

# **ESE 650: Learning in Robotics**

## **Spring 2021**

### **Instructor**

Pratik Chaudhari    [pratikac@seas.upenn.edu](mailto:pratikac@seas.upenn.edu)

### **Teaching Assistants**

Shuo Li    [lishuo1@seas](mailto:lishuo1@seas)  
Lakshay Sharma    [lakshays@seas](mailto:lakshays@seas)  
Yihang Xu    [yihangxu@seas](mailto:yihangxu@seas)

May 13, 2021

# Contents

<b>1</b>	<b>What is Robotics?</b>	<b>4</b>
1.1	Perception-Learning-Control	5
1.2	Goals of this course	6
1.3	Some of my favorite robots	6
<b>2</b>	<b>Introduction to State Estimation</b>	<b>7</b>
2.1	A review of probability	7
2.1.1	Random variables	10
2.2	Using Bayes rule for combining evidence	13
2.2.1	Coherence of Bayes rule	16
2.3	Markov Chains	17
2.4	Hidden Markov Models (HMMs)	19
2.4.1	The forward algorithm	21
2.4.2	The backward algorithm	22
2.4.3	Bayes filter	23
2.4.4	Smoothing	24
2.4.5	Prediction	25
2.4.6	Decoding: Viterbi's Algorithm	26
2.4.7	Shortest path on a Trellis graph	29
2.5	Learning an HMM from observations	30
<b>3</b>	<b>Kalman Filter and its variants</b>	<b>34</b>
3.1	Background	34
3.2	Linear state estimation	37
3.2.1	One-dimensional Gaussian random variables	37
3.2.2	General case	38
3.2.3	Incorporating Gaussian observations of a state	39
3.2.4	An example	41
3.3	Background on linear and nonlinear dynamical systems	42
3.3.1	Linear systems	43
3.3.2	Linear Time-Invariant (LTI) systems	44
3.3.3	Nonlinear systems	45
3.4	Markov Decision Processes (MDPs)	45
3.4.1	Back to Hidden Markov Models	47
3.5	Kalman Filter (KF)	48
3.5.1	Step 0: Observing that the state estimate at any timestep is a Gaussian	49
3.5.2	Step 1: Propagating the dynamics by one timestep	49

---

3.5.3	Step 2: Incorporating the observation	50
3.5.4	Discussion	51
3.6	Extended-Kalman Filter (EKF)	52
3.6.1	Propagation of statistics through a nonlinear transformation	53
3.6.2	Extended Kalman Filter	55
3.7	Unscented Kalman Filter (UKF)	57
3.7.1	Unscented Transform	59
3.7.2	The UT with tuning parameters	61
3.7.3	Unscented Kalman Filter (UKF)	62
3.7.4	UKF vs. EKF	63
3.8	Particle Filters (PFs)	64
3.8.1	Importance sampling	64
3.8.2	Resampling particles to make the weights equal	65
3.8.3	Particle filtering: the algorithm	68
3.8.4	Example: Localization using particle filter	69
3.8.5	Theoretical insight into particle filtering	74
3.9	Discussion	77
<b>4</b>	<b>Rigid-body transforms and mapping</b>	<b>78</b>
4.1	Rigid-Body Transformations	79
4.1.1	3D transformations	81
4.1.2	Rodrigues' formula: an alternate view of rotations	83
4.2	Quaternions	84
4.3	Occupancy Grids	87
4.3.1	Estimating the map from the data	90
4.3.2	Sensor models	91
4.3.3	Back to sensor modeling	93
4.4	3D occupancy grids	95
4.5	Local Map	96
4.6	Discussion	97
<b>5</b>	<b>Dynamic Programming</b>	<b>99</b>
5.1	Formulating the optimal control problem	100
5.2	Dijkstra's algorithm	101
5.2.1	Dijkstra's algorithm in the backwards direction	102
5.3	Principle of Dynamic Programming	103
5.3.1	Q-factor	105
5.4	Stochastic dynamic programming: Value Iteration	106
5.4.1	Infinite-horizon problems	109
5.4.2	Dynamic programming for infinite-horizon problems	111
5.4.3	An example	112
5.4.4	Some theoretical results on value iteration	113
5.5	Stochastic dynamic programming: Policy Iteration	115
5.5.1	An example	117

---

<b>6</b>	<b>Linear Quadratic Regulator (LQR)</b>	<b>119</b>
6.1	Discrete-time LQR	119
6.1.1	Solution of the discrete-time LQR problem	122
6.2	Hamilton-Jacobi-Bellman equation	124
6.2.1	Infinite-horizon HJB	126
6.2.2	Solving the HJB equation	126
6.2.3	Continuous-time LQR	127
6.3	Stochastic LQR	130
6.4	Linear Quadratic Gaussian (LQG)	131
6.4.1	(Optional material) The duality between the Kalman Filter and LQR	133
6.5	Iterative LQR (iLQR)	134
6.5.1	Iterative LQR (iLQR)	135
<b>7</b>	<b>Imitation Learning</b>	<b>138</b>
7.1	A crash course in supervised learning	140
7.1.1	Fitting a machine learning model	141
7.1.2	Deep Neural Networks	143
7.2	Behavior Cloning	146
7.2.1	Behavior cloning with a stochastic controller	147
7.2.2	KL-divergence form of Behavior Cloning	149
7.2.3	Some remarks on Behavior Cloning	150
7.3	DAgger: Dataset Aggregation	150
<b>8</b>	<b>Policy Gradient Methods</b>	<b>153</b>
8.1	Standard problem setup in RL	153
8.2	Cross-Entropy Method (CEM)	155
8.2.1	Some remarks on sample complexity of simulation-based methods	157
8.3	The Policy Gradient	158
8.3.1	Reducing the variance of the policy gradient	160
8.4	An alternative expression for the policy gradient	162
8.4.1	Implementing the new expression	163
8.5	Actor-Critic methods	164
8.5.1	Advantage function	165
8.6	Discussion	167
<b>9</b>	<b>Q-Learning</b>	<b>168</b>
9.1	Tabular Q-Learning	168
9.1.1	How to perform exploration in Q-Learning	171
9.2	Function approximation (Deep Q Networks)	172
9.2.1	Embellishments to Q-Learning	174
9.3	Q-Learning for continuous control spaces	178
<b>10</b>	<b>Model-based Reinforcement Learning</b>	<b>181</b>
10.1	Learning a model of the dynamics	182
10.2	Some model-based methods	183
10.2.1	Bagging multiple models of the dynamics	183
10.2.2	Model-based RL in the latent space	185
	<b>Bibliography</b>	<b>186</b>

# Chapter 1

## What is Robotics?

### Reading

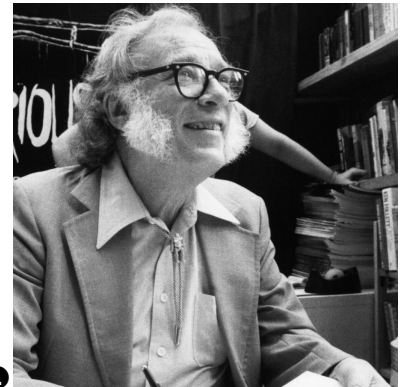
1. Computing machinery and intelligence, Turing (2009)
2. Thrun Chapter 1
3. Barfoot Chapter 1

The word *robotics* was essentially coined by Issac Asimov in one of his short stories named *Liar!*. This is about a *robot* named RB-34 which, through a manufacturing fault, happens to be able to read the minds of humans around it. Around 1942, Isaac Asimov started using the word robot in his writings, is among its first references. This is also when he introduced the Three Laws of Robotics as the theme for how robots would interact with others in his stories/books. These are as follows.

1. A robot may not injure a human being or, through inaction, allow a human being to come to harm.
2. A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.
3. A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.

Asimov would go on to base his stories on the counter-intuitive ways in which robots could apply these laws. In this case, RB-34 adheres to the First Law and in order to not hurt the feelings of humans and make them happy, it deliberately lies to them. It tells the *robopsychologist* Susan Colins that one of her co-workers is infatuated with her. However, when she confronts RB-34 later by pointing out that lying to people can end up hurting them, the robot experiences a logical conflict within its laws and becomes unresponsive.

This is, after all, science fiction but these laws give us insight into what robots are. Let's see what modern *roboticists* have to say.



1 “Robotics is the science of perceiving and manipulating the phys-  
2 ical world through computer-controlled mechanical devices.” —  
3 Sebastian Thrun in Probabilistic Robotics

4 “EVERYTHING comes together in the field of robotics. The de-  
5 sign of an autonomous robots involves: the choice of the mechan-  
6 ical platform, the choice of actuators, the choice of sensors, the  
7 choice of the energy source, the choices of algorithms (perception,  
8 planning, and control). Each of these subproblems corresponds  
9 to a discipline in itself, with its design trade-offs of achievable  
10 performance vs limited resources.” — Andrea Censi in [Censi](#)  
11 [\(2016\)](#).

12 I find the Third Law really insightful to understand intelligence as well. Let  
13 us define intelligence as the ability of an organism to survive<sup>1</sup>. We will all agree  
14 that trees are less intelligent than animals, an ant is less than intelligent than a  
15 dog, which is less intelligent than a human. A program like AlphaGo is not  
16 very intelligent, you can disable it by simply switching it off. A key indicator  
17 of intelligence is the ability to move around in the physical environment, to  
18 sense possible harm and take actions to change the outcome.

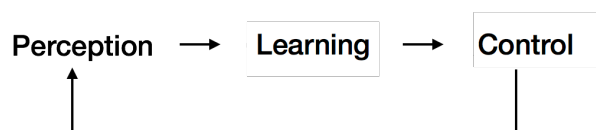
19 Robotics is Embodied Artificial Intelligence.

20 A robot is a machine that senses its environment using sensors,  
21 interacts with this environment using actuators to perform a given  
22 task and does so efficiently using previous experience of perform-  
23 ing similar tasks.

24 We will cover the fundamentals of these three aspects of robotics: percep-  
25 tion, planning and learning.

## 26 1.1 Perception-Learning-Control

27 Perception refers to the sensory mechanisms to gain information about the  
28 environment (eyes, ears, tactile input etc.). Action refers to your hands, legs, or  
29 motors/engines in machines that help you move on the basis of this information.  
30 Learning is kind of the glue in between. It helps crunch information of your  
31 sensors quickly, compare it with past data, guesses what future data may  
32 look like and computes actions that are likely to succeed. The three facets  
33 of intelligence are not sequential and robotics is not merely a feed-forward  
34 process. Your sensory inputs depend on the previous action you took.



---

<sup>1</sup>feel free to come up with another definition

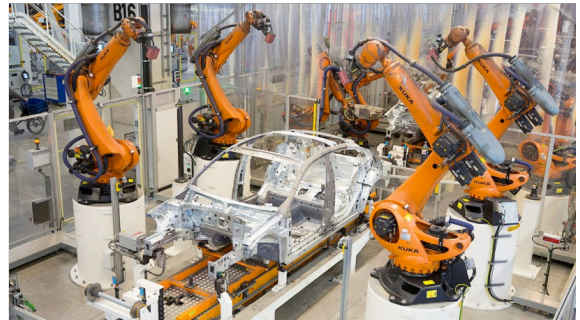
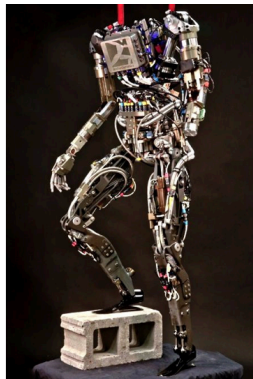
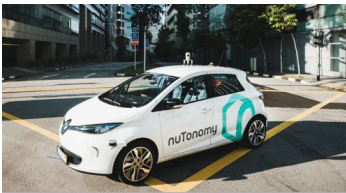
## 1.2 Goals of this course

The goals of this course is to develop the main ideas in robotic perception, learning and control. Robotics *is* everything, so we will focus on understanding how they are *combined together* to build a typical robot. After this course, we expect you to be able to choose one among the different robotics algorithms to perform a particular task, think critically about these algorithms and build new ones.

**Other courses** Some examples of other courses at Penn that address various aspects of this picture above are

- Perception: CIS 580, CIS 581, CIS 680
- Learning: CIS 520, CIS 521, CIS 522, CIS 620, CIS 700, ESE 545, ESE 546
- Control: ESE 650, MEAM 520, MEAM 620, ESE 500, ESE 505, ESE 619

## 1.3 Some of my favorite robots



These videos should give you an idea of how the everyday life of a roboticist looks like: [Kiva's robots](#), [Waymo's 360 experience](#), [Boston Dynamics' Spot](#), [JPL-MIT team at the DARPA Sub-T Challenge](#), [Romeo and Juliet at Ferrari's factory](#), [Anki's Vector](#), and the [DARPA Humanoid Challenge](#).

# 1 Chapter 2

## 2 Introduction to State 3 Estimation

### Reading

1. Barfoot, Chapter 2.1-2.2
2. Thrun, Chapter 2
3. Russell Chapter 15.1-15.3

### 4 2.1 A review of probability

5 Probability is a very useful construct to reason about real systems which we  
6 cannot model fully at all scales. It is a fundamental part of robotics. No matter  
7 how sophisticated your camera, it will have noise in how it measures the real  
8 world around it. No matter how good your model for a motor is, there will be  
9 modeled effects which make it move a little differently than how you would  
10 expect. We begin with a quick review of probability, you can read more at  
11 many sources, e.g., [MIT's OCW](#).

12 An experiment is a procedure which can be repeated infinitely and has a  
13 well-defined set of possible outcomes, e.g., the toss of a coin or the roll of dice.  
14 The outcome itself need not always be deterministic, e.g., depending upon  
15 your experiment, the coin may come up heads or tails. We call the set  $\Omega$  the  
16 *sample space*, it is the set of all possible outcomes of an experiment. For two  
17 coins, this set would be

$$\Omega = \{HH, HT, TH, TT\}.$$

18 We want to pick this set to be right granularity to answer relevant questions,  
19 e.g., it is correct but not useful for  $\Omega$  to be position of all the molecules in the



1 coin. After every experiment, in this case tossing the two coins once each, we  
2 obtain an event, it is a subset *event*  $A \subset \Omega$  from the sample space.

$$A = \{HH\}.$$

3 Probability theory is a mathematical framework that allows us to reason  
4 about phenomena or experiments whose outcome is uncertain. Probability of  
5 an event

$$P(A)$$

6 is a function that maps each event  $A$  to a number between 0 and 1: closer to 1  
7 this number, stronger our belief that the outcome of the experiment is going to  
8 be  $A$ .

9 **Axioms** Probability is formalized using a set of three basic axioms that are  
10 intuitive and yet very powerful, they are known a Kolmogorov's axioms. They  
11 are

- 12 • Non-negativity:  $P(A) \geq 0$
- 13 • Normalization:  $P(\Omega) = 1$
- 14 • Additivity: If two events  $A, B$  are such that  $A \cap B = \emptyset$ , then

$$P(A \cup B) = P(A) + P(B).$$

15 You can use these axioms to say things like  $P(\emptyset) = 0$ ,  $P(A^c) = 1 - P(A)$ , or  
16 if  $A \subseteq B$  then  $P(A) \leq P(B)$ .

17 **Conditioning on events** Conditioning helps us answer questions like

$$P(A | B) := \text{probability of } A \text{ given that } B \text{ occurred.}$$

18 Effectively, the sample space has now shrunk from  $\Omega$  to the event  $B$ . It would  
19 be silly to have a null sample-space, so let's say that  $P(B) \neq 0$ . We define  
20 conditional probability as

$$P(A | B) = \frac{P(A \cap B)}{P(B)}; \quad (2.1)$$

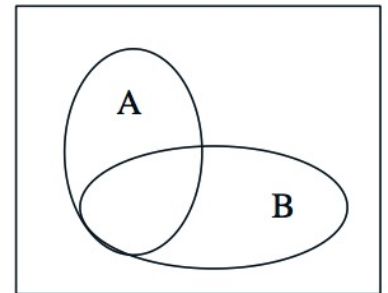
21 the probability is undefined if  $P(B) = 0$ . Using this definition, we can  
22 compute the probability of events like "what is the probability of rolling a 2  
23 on a die given that an even number was rolled".

24 We can use this trick to get the law of total probability: if a finite number  
25 of events  $\{A_i\}$  form a partition of  $\Omega$ , i.e.,

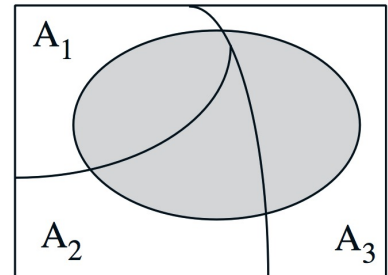
$$A_i \cap A_j = \emptyset \quad \forall i, j, \text{ and } \bigcup_i A_i = \Omega$$

$$26 \quad P(B) = \sum_i P(B | A_i) P(A_i). \quad (2.2)$$

❗



❗ Partitioning the sample space



1 **Bayes' rule** Imagine that instead of someone telling us that the conditioning  
 2 event actually happened, we simply had a belief

$$P(A_i)$$

3 about the possibility of such events  $\{A_i\}$ . For each of  $A_i$ , we can compute the  
 4 conditional probability  $P(B | A_i)$  using (2.1). Say we run our experiment and  
 5 observe that  $B$  occurred, how would our belief on the events  $A_i$  change? In  
 6 other words, we wish to compute

$$P(A_i | B).$$

7 This is the subject of Bayes' rule.

$$\begin{aligned} P(A_i | B) &= \frac{P(A_i \cap B)}{P(B)} \\ &= \frac{P(A_i) P(B|A_i)}{P(B)} \\ &= \frac{P(A_i) P(B|A_i)}{\sum_j P(A_j) P(B | A_j)}. \end{aligned} \tag{2.3}$$

8 The Bayes' rule naturally leads to the concept of independent events. Two  
 9 events  $A, B \subset \Omega$  are independent if observing one does not give us any  
 10 information about the other

$$P(A \cap B) = P(A) P(B). \tag{2.4}$$

11 This is different from disjoint events. Disjoint events never co-occur, i.e.,  
 12 observing one tells us that the other one *did not* occur.

13 **Probability for experiments with real-valued outcomes** We need some  
 14 more work in defining probability for events with real-valued outcomes. The  
 15 sample space is easy enough to understand, e.g.,  $\Omega = [0, 1]$  for your score  
 16 at the end of this course. We however run into difficulties if we define the  
 17 probability of general subsets of  $\Omega$  in terms of the probabilities of elementary  
 18 outcomes (*elements* of  $\Omega$ ). For instance, if we wish to model all elements  
 19  $\omega \in \Omega$  to be equally likely, we are forced to assign each element  $\omega$  a probability  
 20 of zero (to be consistent with the second axiom of probability). This is not very  
 21 helpful in determining the probability of the score being 0.9. If you instead  
 22 assigned some small non-zero number to  $P(\omega_i)$ , then we have undesirable  
 23 conclusions such as

$$P(\{1, 1/2, 1/3, \dots\}) = \infty.$$

24 The way to fix this is to avoid defining the probability of a set in terms  
 25 of the probability of elementary outcomes and work with more general sets.  
 26 While we would ideally like to be able to specify the probability of every  
 27 subset of  $\Omega$ , it turns out that we cannot do so in a mathematically consistent  
 28 way. The trick then is to work with a smaller object known as a  $\sigma$ -algebra, that  
 29 is the set of "nice" subsets of  $\Omega$ .

Given a sample space  $\Omega$ , a  $\sigma$ -algebra  $\mathcal{F}$  (also called a  $\sigma$ -field) is a collection of subsets of  $\Omega$  such that

- $\emptyset \in \mathcal{F}$
- If  $A \in \mathcal{F}$ , then  $A^c \in \mathcal{F}$ .
- If  $A_i \in \mathcal{F}$  for every  $i \in \mathbb{N}$ , then  $\cup_{i=1}^{\infty} A_i \in \mathcal{F}$ .

In short,  $\sigma$ -algebra is a collection of subsets of  $\Omega$  that is closed under complement and countable unions. The pair  $(\Omega, \mathcal{F})$ , also called a measurable space, is now used to define probability of events. A set  $A$  that belongs to  $\mathcal{F}$  is called an event. The probability measure

$$P : \mathcal{F} \rightarrow [0, 1].$$

assigns a probability to events in  $\mathcal{F}$ . We cannot take  $\mathcal{F}$  to be too small, e.g., elements of  $\mathcal{F} = \{\emptyset, \Omega\}$  are easy to construct our  $P$  but are not very useful. For technical reasons, the  $\sigma$ -algebra cannot be too large; notice that we used this concept to *avoid* considering every subset of the sample space  $\mathcal{F} = 2^\Omega$ . Modern probability is defined using a Borel  $\sigma$ -algebra. Roughly speaking, this is an  $\mathcal{F}$  that is just large enough to do interesting things but small enough that mathematical technicalities do not occur.

### 2.1.1 Random variables

A random variable is an assignment of a value to every possible outcome. Mathematically, in our new language of a measurable space, a random variable is a function

$$X : \Omega \rightarrow \mathbb{R}$$

if the set  $\{\omega : X(\omega) \leq c\}$  is  $\mathcal{F}$ -measurable for every number  $c \in \mathbb{R}$ . This is equivalent to saying that every preimage of the Borel  $\sigma$ -algebra on reals  $\mathcal{B}(\mathbb{R})$  is in  $\mathcal{F}$ . A statement  $X(\omega) = x = 5$  means that the outcome of our experiment happens to  $\omega \in \Omega$  when the realized value of the random variable is a particular number  $x$  equal to 5 in our case.

We can now define functions of random variables, e.g., if  $X$  is a random variable, the function  $Y = X^3(\omega)$  for every  $\omega \in \Omega$ , or  $Y = X^3$  for short, is a new random variable. An indicator random variable is a special. If  $A \subset \Omega$ , let  $I_A : \Omega \rightarrow \{0, 1\}$  be the indicator function of this set  $A$ , i.e.,  $I_A(\omega) = 1$  if  $\omega \in A$  and zero otherwise. If our set  $A \in \mathcal{F}$ , then  $I_A$  is an indicator random variable.

**Probability mass functions** The probability law, or a probability distribution, of a random variable  $X$  is denoted by

$$p_X(x) := P(X = x) = P(\{\omega \in \Omega : X(\omega) = x\}).$$

We denote probability distribution using a lower-case  $p$ . It is a function of the realized value  $x$  in the range of a random variable, and  $p_X(x) \geq 0$  (the probability is non-zero) and  $\sum_x p_X(x) = 1$  if  $X$  takes on a discrete number of values. For instance, if  $X$  is the number of coin tosses until the first head, if

❶ Random variables are typically denoted using capital letters,  $X, Y, Z$  although we will be sloppy and not always do so in this course to avoid complicated notation. The distinction between a random variable and the value that it takes will be clear from context.

❷ Let us check that  $Y$  satisfies our definition of a random variable. If  $\{\omega : X(\omega) \leq c\}$  lies in  $\mathcal{F}$  then the set  $\{\omega : Y(\omega) \leq c^{1/3}\}$  also lies in  $\mathcal{F}$ .

❸ The function  $I_A$  is not a random variable if  $A \notin \mathcal{F}$ , but this is, as we said in the previous section, a mathematical corner case. Most subset of  $\Omega$  belong to  $\mathcal{F}$ .

1 we assume that our tosses are independent  $P(H) = p > 0$ , then we have

$$p_X(k) = P(X = k) = P(TT \cdots TH) = (1 - p)^{k-1} p$$

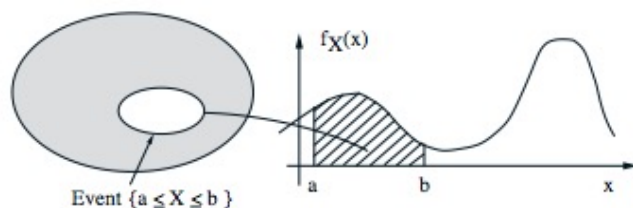
2 for all  $k = 1, 2, \dots$ . This is what is called a geometric probability mass  
3 function.

4 **Cummulative distribution function** A cummulative distribution function  
5 (CDF) is the probability of a random variable  $X$  taking a value less than an  
6 particular  $x \in \mathbb{R}$ , i.e.,

$$F_X(x) = P(X \leq x).$$

7 The CDF  $F_X(x)$  is a non-decreasing function of  $x$ . It converges to zero as  
8  $x \rightarrow -\infty$  and goes to 1 as  $x \rightarrow \infty$ .

9 **Probability density functions** A continuous random variable, i.e., one that  
10 takes values in  $\mathbb{R}$  is described by a probability density function.



11

12 If  $F_X(x)$  is the CDF of an r.v.  $X$  and  $X$  takes values in  $\mathbb{R}$ , the probability  
13 density function (PDF)  $f_X(x)$  (or sometimes also denoted by  $p_X(x)$ ) is defined  
14 to be

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx.$$

15 We also have the following relationship between the CDF and the PDF, the  
16 former is the integral of the latter:

$$P(-\infty \leq X \leq x) = F_X(x) = \int_{-\infty}^x f_X(x) dx.$$

17 This leads to the following interpretation of the probability density function:

$$P(x \leq X \leq x + \delta) \approx f_X(x) \delta.$$

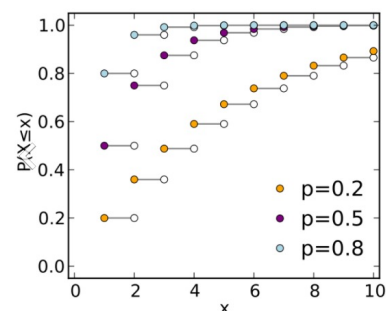
18 **Expectation and Variance** The expected value of a random variable  $X$  is

$$E[X] = \sum_x x p_X(x)$$

19 and denotes the center of gravity of the probability mass function. Roughly  
20 speaking, it is the average of a large number of repetitions of the same experi-  
21 ment. Expectation is a linear, i.e.,

$$E[aX + b] = aE[X] + b$$

❶ The CDF of a geometric random variable for different values of  $p$



Note that CDFs need not be continuous, in the case of a geometric r.v. since the values that  $X$  takes belong to the set of integers, the CDF is constant between any two integers.

1 for any constants  $a, b$ . For two independent random variables  $X, Y$  we have

$$E[XY] = E[X] E[Y].$$

2 We can also compute the expected value of any function  $g(X)$  using the  
3 same formula

$$E[g(X)] = \sum_x g(x) p_X(x).$$

4 In particular, if  $g(x) = x^2$  we have the second moment  $E[X^2]$ . The variance  
5 is defined to be

$$\begin{aligned} \text{var}(X) &= E[(X - E[X])^2] \\ &= \sum_x (x - E[X])^2 p_X(x) \\ &= E[X^2] - (E[X])^2. \end{aligned}$$

6 The variance is always non-negative  $\text{var}(X) \geq 0$ . For an affine function of  $X$ ,  
7 we have

$$\text{var}(aX + b) = a^2 \text{var}(X).$$

8 For continuous-valued random variables, the expectation is defined as

$$E[X] = \int_{-\infty}^{\infty} x p_X(x) dx;$$

9 the definition of the variance remains the same.

10 **Joint distributions** We often wish to think of the joint probability distribu-  
11 tion of multiple random variables, say the location of an autonomous car in all  
12 three dimensions. The cumulative distribution function associated with this  
13 is therefore

$$F_{X,Y,Z}(x, y, z) = P(X \leq x, Y \leq y, Z \leq z).$$

14 Just like we have the probability density of a single random variable, we  
15 can also write the joint probability density of multiple random variables  
16  $f_{X,Y,Z}(x, y, z)$ . In this case we have

$$F_{X,Y,Z}(x, y, z) = \int_{-\infty}^x \int_{-\infty}^y \int_{-\infty}^z f_{X,Y,Z}(x, y, z) dz dy dx.$$

17 The joint probability density factorizes if two random variables are independ-  
18 ent:

$$f_{X,Y}(x, y) = f_X(x) f_Y(y) \quad \text{for all } x, y.$$

19 Two random variables are uncorrelated if and only if

$$E[XY] = E[X] E[Y].$$

20 Note that independence implies uncorrelatedness, they are not equivalent. The  
21 covariance is defined as

$$\text{cov}(X, Y) = E[XY] - E[X] E[Y].$$

1 **Conditioning** As we saw before, for a single random variable  $X$  we have

$$P(x \leq X \leq x + \delta) \approx f_X(x) \delta.$$

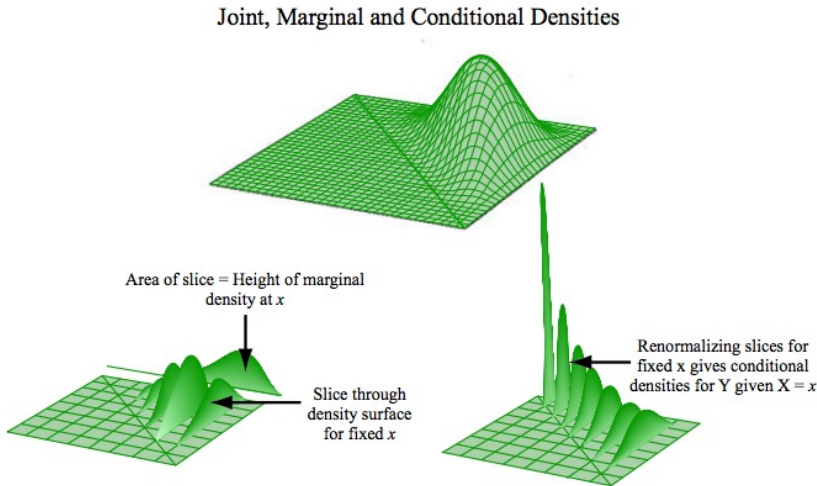
2 For two random variables, by analogy we would like

$$P(x \leq X \leq x + \delta | Y \approx y) \approx f_{X|Y}(x | y) \delta.$$

3 The conditional probability density function of  $X$  given  $Y$  is defined to be

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)} \quad \text{if } f_Y(y) > 0.$$

4 For any given  $y$ , the conditional PDF is a normalized section of the joint PDF,  
5 as shown below.



6  
7 **Continuous form of Bayes rule** We can show using the definition of condi-  
8 tional probability that

$$f_{Y|X}(y | x) = \frac{f_{X|Y}(x | y)}{f_X(x)}. \quad (2.5)$$

9 Similarly we also have the law of total probability in the continuous form

$$f_X(x) = \int_{-\infty}^{\infty} f_{X|Y}(x | y) f_Y(y) dy.$$

## 10 2.2 Using Bayes rule for combining evidence

11 We now study a prototypical state estimation problem. Let us consider a  
12 robot that is trying to check whether the door to a room is open or not.



1

2 We will abstract each observation by the sensors of the robot as a random  
 3 variable  $Y$ . This could be the image from its camera after running some  
 4 algorithm to check the state of the door, the reading from a laser sensor (if  
 5 the time-of-flight of the laser is very large then the door is open), or any  
 6 other mechanism. Let us first distinguish between the cause and effect of the  
 7 observation: if the door is open, we get an observation of a certain kind (laser  
 8 sensor says that the time of flight is large), of course the observation is not  
 9 *caused* by the open door. In other words, we have two kinds of conditional  
 10 probabilities in this problem

$P(\text{open} | Y)$  is a diagnostic quantity, while

$P(Y | \text{open})$  is a causal quantity.

11 Imagine how you would calibrate the sensor in a lab: for each value of the  
 12 state of the door open, not open you would record all the different observations  
 13 received  $Y$  and calculate the conditional probabilities. The causal probability  
 14 is much easier to calculate in this context, one may even use some knowledge  
 15 of elementary physics to model the probability  $P(Y | \text{open})$ , or one may count  
 16 the number of times the observation is  $Y = y$  for a given state of the door.

17 The Bayes rule allows us to transform causal knowledge into diagnostic  
 18 knowledge

$$P(\text{open} | Y) = \frac{P(Y | \text{open}) P(\text{open})}{P(Y)}.$$

19 Remember that the left hand side (diagnostic) is typically something that  
 20 we desire to calculate. Let us put some numbers in this formula. Let  
 21  $P(Y | \text{open}) = 0.6$  and  $P(Y | \text{not open}) = 0.3$ . We will imagine that the  
 22 door is open or closed with equal probability:  $P(\text{open}) = P(\text{not open}) = 0.5$ .  
 23 We then have

$$\begin{aligned} P(\text{open} | Y) &= \frac{P(Y | \text{open}) P(\text{open})}{P(Y)} \\ &= \frac{P(Y | \text{open}) P(\text{open})}{P(Y | \text{open}) P(\text{open}) + P(Y | \text{not open}) P(\text{not open})} \\ &= \frac{0.6 \times 0.5}{0.6 \times 0.5 + 0.3 \times 0.5} = \frac{2}{3}. \end{aligned}$$

24 Notice something very important, the original (prior) probability of the state  
 25 of the door is was 0.5. If we have a sensor that fires with higher likelihood if

1 the door is open, i.e., if

$$\frac{P(Y \mid \text{open})}{P(Y \mid \text{not open})} > 1$$

2 then the probability of the door being open after receiving an observation  
3 *increases*. If the likelihood were less than 1, then observing a realization of  $Y$   
4 would reduce our estimate of the probability of the door being open.

5 **Combining evidence for Markov observations** Say we updated the prior  
6 probability using our first observation  $Y_1$ , let us take another observation  $Y_2$ .  
7 How can we integrate this new observation? It is again an application of Bayes  
8 rule using two observations, or in general multiple observations  $Y_1, \dots, Y_n$ .  
9 Let us imagine this time that  $X = \text{open}$ .

$$P(X \mid Y_1, \dots, Y_n) = \frac{P(Y_n \mid X, Y_1, \dots, Y_{n-1}) P(X \mid Y_1, \dots, Y_{n-1})}{P(Y_n \mid Y_1, \dots, Y_{n-1})}.$$

10 Let us make the very natural assumption that says that our observations from  
11 the sensor  $Y_1, \dots, Y_n$  are independent given the state of the door  $X$ . This is  
12 known as the Markov assumption.

13 We now have

$$\begin{aligned} P(X \mid Y_1, \dots, Y_n) &= \frac{P(Y_n \mid X) P(X \mid Y_1, \dots, Y_{n-1})}{P(Y_n \mid Y_1, \dots, Y_{n-1})} \\ &= \eta P(Y_n \mid X) P(X \mid Y_1, \dots, Y_{n-1}) \end{aligned}$$

14 where

$$\eta^{-1} = P(Y_n \mid Y_1, \dots, Y_{n-1})$$

15 is the denominator. We can now expand the diagnostic probability on the  
16 right-hand side recursively to get

$$P(X \mid Y_1, \dots, Y_n) = \prod_{i=1}^n \eta_i P(Y_i \mid X) P(X). \quad (2.6)$$

17 where  $\eta_i^{-1} = P(Y_i \mid Y_1, \dots, Y_{i-1})$ .

The calculation in (2.6) is very neat and you should always remember

❗ The denominator in Bayes rule, i.e.,  $P(Y)$  is called the evidence in statistics.



it. Given multiple observations  $Y_1, \dots, Y_n$  of the same quantity  $X$ , we can compute the conditional probability  $P(X | Y_1, \dots, Y_n)$  if we code up two functions to compute

- the causal probability (also called the likelihood of an observation)  $P(Y_i | X)$ , and
- the denominator  $\eta_i^{-1}$ .

Given these two functions, we can use the recursion to update multiple observations. The same basic idea also holds if you have two quantities to estimate, e.g.,  $X_1 = \text{open door}$  and  $X_2 = \text{color of the door}$ . The recursive application of Bayes rule lies at the heart of all state estimation methods.

1 Let us again put some numbers into these formulae, imagine that the  
2 observation  $Y_2$  was taken using a different sensor which now has

$$P(Y_2 | \text{open}) = 0.5 \text{ and } P(Y_2 | \text{not open}) = 0.6.$$

3 We have from our previous calculation that  $P(\text{open} | Y_1) = 2/3$  and

$$\begin{aligned} P(\text{open} | Y_1, Y_2) &= \frac{P(Y_2 | \text{open}) P(\text{open} | Y_1)}{P(Y_2 | \text{open}) P(\text{open} | Y_1) + P(Y_2 | \text{not open}) P(\text{not open} | Y_1)} \\ &= \frac{0.5 \times 2/3}{0.5 \times 2/3 + 0.6 \times 1/3} = \frac{5}{8} = 0.625. \end{aligned}$$

4 Notice in this case that the probability that the door is open has reduced from  
5  $P(\text{open} | Y_1) = 2/3$ .

### 6 2.2.1 Coherence of Bayes rule

7 Would the probability change if we used sensor  $Y_2$  before using  $Y_1$ ? In this  
8 case, the answer to this question is no and you are encouraged to perform this  
9 computation for yourselves. Bayes rule is coherent, it will give the same result  
10 regardless of the order of observations.

11 The order of incorporating observation matters if the state of the world  
12 changes while we make observations, e.g., if we have a sensor that tracks  
13 the location of a car, the car presumably moves in between two observations  
14 and we would get the wrong answer if our question was “is there a car at this  
15 location”.

16 As we motivated in the previous chapter, movement is quite fundamental to  
17 robotics and we are typically concerned with estimating the state of a dynamic  
18 world around us using our observations. We will next study the concept of a  
19 Markov Chain with is a mathematical abstraction for the evolution of the state  
20 of the world.

🔗 Can you think of a situation when the order of incorporating observations matters?

## 2.3 Markov Chains

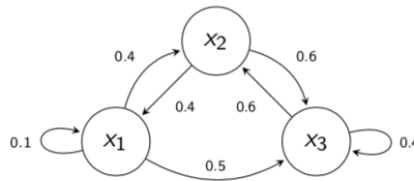
Consider the Whack-The-Mole game: a mole has burrowed a network of three holes  $x_1, x_2, x_3$  into the ground. It keeps going in and out of the holes and we are interested in finding which hole it will show up next so that we can give it a nice whack.

- Three holes:

$$X = \{x_1, x_2, x_3\}.$$

- Transition probabilities:

$$T = \begin{bmatrix} 0.1 & 0.4 & 0.5 \\ 0.4 & 0 & 0.6 \\ 0 & 0.6 & 0.4 \end{bmatrix}$$



This is an example of a Markov chain. There is a transition matrix  $T$  which determines the probability  $T_{ij}$  of the mole resurfacing on a given hole  $x_j$  given that it resurfaced at hole  $x_i$  the last time. The matrix  $T^k$  is the  $k$ -step transition matrix

$$T_{ij}^k = P(X_k = x_j \mid X_0 = x_i).$$

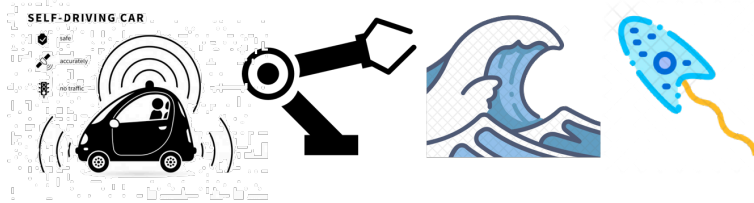
You can see the animations at <https://setosa.io/ev/markov-chains> to build more intuition.

The key property of a Markov chain is that the next state  $X_{k+1}$  is independent of all the past states  $X_1, \dots, X_{k-1}$  given the current state  $X_k$ .

$$X_{k+1} \perp\!\!\!\perp X_1, \dots, X_{k-1} \mid X_k$$

This is known as the Markov property and all systems where we can define a “state” which governs their evolution have this property. Markov chains form a very broad class of systems. For example, all of Newtonian physics fits this assumption.

What is the state of the following systems?



❓ Does a deterministic dynamical system, e.g., a simple pendulum, also satisfy the Markov assumption? What is the transition matrix in this case?

❓ Can you think of a system which does not have the Markov property?

Consider the paramecium above. Its position depends upon a large number of factors: its own motion from the previous time-step but also the viscosity of the material in which it is floating around. One may model the state of the environment around the paramecium as a liquid whose molecules hit thousands of times a second, essentially randomly, and cause disturbances in how the paramecium moves. Let us call this disturbance “noise in the dynamics”. If the

1 motion of the molecules of the liquid has some correlations (does it, usually?),  
 2 this induces correlations in the position of the paramecium. The position of  
 3 the organism is no longer Markov. This example is important to remember,  
 4 the Markov property defined above also implies that the noise in the state  
 5 transition matrix is independent.

6 **Evolution of a Markov chain** The probability of being in a state  $x^i$  at time  
 7  $k + 1$  can be written as

$$P(X_{k+1} = x_i) = \sum_{j=1}^N P(X_{k+1} = x_i \mid X_k = x_j) P(X_k = x_j).$$

8 This equation governs how the probabilities  $P(X_k = x_i)$  change with time  $k$ .  
 9 Let's do the calculations for the Whack-The-Mole example. Say the mole was  
 10 at hole  $x_1$  at the beginning. So the probability distribution of its presence

$$\pi^{(k)} = \begin{bmatrix} P(X_k = x_1) \\ P(X_k = x_2) \\ P(X_k = x_3) \end{bmatrix}$$

11 is such that

$$\pi^{(1)} = [1, 0, 0]^T.$$

12 We can now write the above formula as

$$\pi^{(k+1)} = T' \pi^{(k)} \quad (2.7)$$

13 <sup>1</sup> and compute the distribution  $\pi^{(t)}$  for all times

$$\begin{aligned} \pi^{(2)} &= T' \pi^{(1)} = [0.1, 0.4, 0.5]^T; \\ \pi^{(3)} &= T' \pi^{(2)} = [0.17, 0.34, 0.49]^T; \\ \pi^{(4)} &= T' \pi^{(3)} = [0.153, 0.362, 0.485]^T; \\ &\vdots \\ \pi^{(\infty)} &= \lim_{k \rightarrow \infty} T'^k \pi^{(1)} \\ &= [0.158, 0.355, 0.487]^T. \end{aligned}$$

14 The numbers  $P(X_k = x_i)$  stop changing with time  $k$ . Under certain technical  
 15 conditions, the distribution  $\pi^{(\infty)}$  is unique (single communicating class for a  
 16 Markov chain finite number states). We can compute this invariant distribution  
 17 by writing

$$\pi^{(\infty)} = T' \pi^{(\infty)}.$$

18 We can also compute the distribution  $\pi^{(\infty)}$  directly: the invariant distribution  
 19 is the right-eigenvector of the matrix  $T'$  corresponding to the eigenvalue 1.

20 **Example 2.1.** Consider a Markov chain on two states where the transition

🔍 Do we always know that the transition matrix has an eigenvalue that is 1?

<sup>1</sup>Let us denote the transpose of the matrix  $T$  using the Matlab notation  $T'$  instead of  $T^\top$  for clarity.

1 matrix is given by

$$T = \begin{bmatrix} 0.5 & 0.5 \\ 0.4 & 0.6 \end{bmatrix}.$$

2 The invariant distribution is

$$\begin{aligned} \pi^{(1)} &= 0.5\pi^{(1)} + 0.4\pi^{(2)} \\ \pi^{(2)} &= 0.5\pi^{(1)} + 0.6\pi^{(2)}. \end{aligned}$$

3 Note that the constraint for  $\pi$  being a probability distribution, i.e.,  $\pi^{(1)} + \pi^{(2)} =$   
 4 1 is automatically satisfied by the two equations. We can solve for  $\pi^{(1)}, \pi^{(2)}$   
 5 to get

$$\pi^{(1)} = 4/9 \quad \pi^{(2)} = 5/9.$$

## 6 2.4 Hidden Markov Models (HMMs)

7 2

8 Markov chains are a good model for how the state of the world evolves  
 9 with time. We may not always know the exact of these systems and only have  
 10 sensors, e.g., cameras, LiDARs, and radars, to record observations. These  
 11 sensors are typically noisy. So we model the observations as random variables.

12 Hidden Markov Models (HMMs) are an abstraction to reason about ob-  
 13 servations of the state of a Markov chain. An HMM is a sequence of random  
 14 variables  $Y_1, Y_2, \dots, Y_n$  such that the distribution of  $Y_k$  only depends upon on  
 15 the hidden state  $X_k$  of the associated Markov chain.

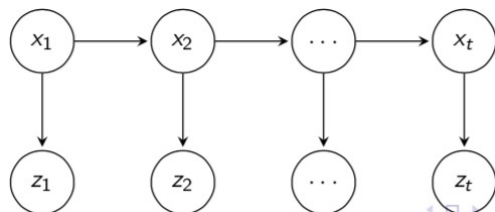


Figure 2.1: A Hidden Markov Model with the underlying Markov chain, the observation at time  $k$  only depends upon the hidden state at that time instant. Ignore the notation  $Z_1, \dots, Z_t$  we will denote the observations by  $Y_k$ .

16 Notice that an HMM always has an underlying Markov chain behind it.  
 17 For example, if we model the position of a car  $X_k$  as a Markov chain, our  
 18 observation of the position at time  $k$  would be  $Y_k$ . In our example of the  
 19 robot sensing whether the door is open or closed using multiple observations  
 20 across time, the Markov chain is trivial, it is simply the transition matrix  
 21  $P(\text{not open} \mid \text{not open}) = P(\text{open} \mid \text{open}) = 1$ . Just like Markov chains,  
 22 HMMs are a very general class of mathematical models that allow us to think  
 23 about multiple observations across time of a Markov chain.

<sup>2</sup>Parts of this section closely follow Emilio Frazzoli's course notes at [https://ocw.mit.edu/courses/aeronautics-and-astronautics/16-410-principles-of-autonomy-and-decision-making-fall-2010/lecture-notes/MIT16\\_410F10\\_lec20.pdf](https://ocw.mit.edu/courses/aeronautics-and-astronautics/16-410-principles-of-autonomy-and-decision-making-fall-2010/lecture-notes/MIT16_410F10_lec20.pdf) and [https://ocw.mit.edu/courses/aeronautics-and-astronautics/16-410-principles-of-autonomy-and-decision-making-fall-2010/lecture-notes/MIT16\\_410F10\\_lec21.pdf](https://ocw.mit.edu/courses/aeronautics-and-astronautics/16-410-principles-of-autonomy-and-decision-making-fall-2010/lecture-notes/MIT16_410F10_lec21.pdf)

Let us imagine that the observations of our HMM are also finite in number, e.g., your score in this course  $\in [0, 100]$  where the associated state of the Markov chain is your expertise in the subject matter. We will write a matrix of observation probabilities

$$M_{ij} = P(Y_k = y_j | X_k = x_i). \quad (2.8)$$

The matrix  $M$  has non-negative entries, after all, each entry is a probability. Since each state has to result in *some* observation, we also have

$$\sum_j M_{ij} = 1.$$

The state transition probabilities of the associated Markov chain are

$$T_{ij} = P(X_{k+1} = x_j | X_k = x_i).$$

Given the abstraction of an HMM, we may be interested in solving a number of problems. We will consider the problem where the state  $X_k$  is the position of a car (which could be stationary or moving) and observations  $Y_k$  give us some estimate of the this position.

1. **Filtering:** Given observations up to time  $k$ , compute the distribution of the state at time  $k$

$$P(X_k | Y_1, \dots, Y_k).$$

This is the most natural problem to understand: we want to find the probability of the car being at a location at time  $k$  given all previous observations. This is a temporally causal prediction, i.e., we are not using any information from the future to reason about the present.

2. **Smoothing:** Given observations up to time  $k$ , compute the distribution of the state at any time  $j < k$

$$P(X_j | Y_1, \dots, Y_k) \quad \text{for } j < k.$$

The observation at a future time  $Y_{k+1}$  gives us some indication of where the car might have been at time  $k$ . In this case we are interested in using the entire set of observations from the past  $Y_1, \dots, Y_j$ , the future  $Y_{j+1}, \dots, Y_k$  to estimate the position of the car. Of course, this problem can only be solved *ex post facto*, i.e., after the time instant  $j$ . An important thing to remember is that we are interested in the position of the car for all  $j < k$  in smoothing.

3. **Prediction:** Given observations up to time  $k$ , compute the distribution of the state at a time  $j > k$

$$P(X_j | Y_1, \dots, Y_k) \quad \text{for } j > k.$$

This is the case when we wish to make predictions about the state of the car  $j > k$  given only observations until time  $k$ . If we knew the underlying Markov chain for the HMM and its transition matrix  $T$ , this would amount to running (2.7) forward using the output of the filtering problem as the initial distribution of the state.

🔍 Why is this true?

- 1 4. **Decoding:** Find the most likely state trajectory  $X_1, \dots, X_k$  that maxi-  
2 mizes the probability

$$P(X_1, \dots, X_k | Y_1, \dots, Y_k)$$

3 given observations  $Y_1, \dots, Y_k$ . Observe that the smoothing problem is  
4 essentially solved independently for all time-steps  $j < k$ . It stands to  
5 reason that if we knew a certain state (say car made a right turn) was  
6 likely given observations at time  $k+1$  and that the traffic light was green  
7 at time  $k$  (given our observations of the traffic light), then we know that  
8 the car did not stop at the intersection at time  $k$ . The decoding problem  
9 allows us to reason about the joint probability of the states and outputs  
10 the most likely trajectory given all observations.

- 11 5. **Likelihood of observations:** Given the observation trajectory,  $Y_1, \dots, Y_k$ ,  
12 compute the probability

$$P(Y_1, \dots, Y_k).$$

13 As you may recall, this is the denominator that we need for the recur-  
14 sive application of Bayes rule. It is made difficult by the fact that we  
15 do not know the state trajectory  $X_1, \dots, X_k$  corresponding to these  
16 observations.

17 These problems are closely related with each other and we will next dig deeper  
18 into them. We will first discuss two building blocks, called the forward and  
19 backward algorithm that together help solve all the above problems.

## 20 2.4.1 The forward algorithm

21 Consider the problem of computing the likelihood of observations. We can  
22 certainly write

$$\begin{aligned} & P(Y_1, \dots, Y_k) \\ &= \sum_{\text{all } (x_1, \dots, x_k)} P(Y_1, \dots, Y_k | X_1, \dots, X_k) P(X_1, \dots, X_k) \\ &= \sum_{\text{all } (x_1, \dots, x_k)} \prod_{i=1}^k P(Y_i = y_i | X_i = x_i) P(X_1 = x_1) \prod_{i=2}^k P(X_i = x_i | X_{i-1} = x_{i-1}) \\ &= \sum_{\text{all } (x_1, \dots, x_k)} M_{x_1 y_1} M_{x_2 y_2} \dots M_{x_k y_k} \pi_{x_1} T_{x_1 x_2} \dots T_{x_{k-1} x_k}. \end{aligned}$$

23 But this is a very large computation, for each possible trajectory  $(x_1, \dots, x_k)$   
24 the states could have taken, we need to perform  $2k$  matrix multiplications.

**Forward algorithm** We can simplify the above computation using the

🔗 How many possible state trajectories are there? What is the total cost of computing the likelihood of observations?

Markov property of the HMM as follows. We will define a quantity known as the forward variable

$$\alpha_k(x) = P(Y_1, \dots, Y_k, X_k = x) \quad (2.9)$$

where  $Y_1, \dots, Y_k$  is our observation sequence up to time  $k$ . Observe now that

1. We can initialize

$$\alpha_1(x) = \pi_x M_{x,y_1} \quad \text{for all } x.$$

2. For each time  $i = 1, \dots, k - 1$ , for all states  $x$ , we can compute

$$\alpha_{k+1}(x) = M_{xy_{k+1}} \sum_{x'} \alpha_k(x') T_{x'x}.$$

using the law of total probability.

3. Finally, we have

$$P(Y_1, \dots, Y_k) = \sum_x \alpha_k(x)$$

by marginalizing over the state variables  $X_k$ .

- 1 This recursion in the forward algorithm is a powerful idea and is much
- 2 faster than our naive summation above.

### 3 2.4.2 The backward algorithm

- 4 Just like the forward algorithm performs the computation recursively in the
- 5 forward direction, we can also perform a backward recursion to obtain the
- 6 probability of the observations. Let us imagine that we have an observation
- 7 trajectory

$$Y_1, \dots, Y_t$$

- 8 up to some time  $t$ . We first define the so-called backward variables which are
- 9 the probability of a future trajectory given the state of the Markov chain at a
- 10 particular time instant

$$\beta_k(x) = P(Y_{k+1}, Y_{k+2}, \dots, Y_t | X_k = x). \quad (2.10)$$

- 11 Notice that the backward variables  $\beta_k$  with the conditioning on  $X_k = x$  are
- 12 slightly different than the forward variables  $\alpha_k$  which are the joint probability
- 13 of the observation trajectory and  $X_k = x$ .

**The Backward algorithm** We can compute the variables  $\beta_k(x)$  recur-

🔗 What is the computational complexity of the Forward algorithm?

sively again as follows.

1. Initialize

$$\beta_t(x) = 1 \quad \text{for all } x.$$

This simply indicates that since we are the end of the trajectory, the future trajectory  $Y_{t+1}, \dots$  does not exist.

2. For all  $k = t - 1, t - 2, \dots, 1$ , for all  $x$ , update

$$\beta_k(x) = \sum_{x'} \beta_{k+1}(x') T_{xx'} M_{x'y_{k+1}}.$$

3. We can now compute

$$P(Y_1, \dots, Y_t) = \sum_x \beta_1(x) \pi_x M_{xy_1}.$$

- 1 **Implementing the forward and backward algorithms in practice** The update equations for both  $\alpha_k$  and  $\beta_k$  can be written using a matrix vector multiplication. We maintain the vectors

$$\begin{aligned} \alpha_k &:= [\alpha_k(x_1), \alpha_k(x_2), \dots, \alpha_k(x_N)] \\ \beta_k &:= [\beta_k(x_1), \beta_k(x_2), \dots, \beta_k(x_N)] \end{aligned}$$

- 4 and can write the updates as

$$\alpha_{k+1}^\top = M_{\cdot, y_{k+1}} \odot (\alpha_k^\top T)$$

- 5 where  $\odot$  denotes the element-wise product and  $M_{\cdot, y_{k+1}}$  is the  $y_{k+1}^{\text{th}}$  column of the matrix  $M$ . The update equation for the backward variables is

$$\beta_k = (\beta_{k+1} \odot M_{\cdot, y_{k+1}}) T.$$

- 7 You must be careful about directly implement these recursions however, because we are iteratively multiplying by matrices  $T, M$  whose entries are all smaller than 1 (they are all probabilities after all), we can quickly run into difficulties where  $\alpha_k, \beta_k$  become too small for some states and we get numerical underflow. You can implement these algorithms in the log-space by writing similar update equations for  $\log \alpha_k$  and  $\log \beta_k$  to avoid such numerical issues.

### 13 2.4.3 Bayes filter

- 14 Let us now use the forward and backward algorithms to solve the filtering problem. We want to compute

$$P(X_k = x \mid Y_1, \dots, Y_k)$$

🔗 What is the computational complexity of running the backward algorithm?



1 for all states  $x$  in the Markov chain. We have that

$$2 \quad \mathbb{P}(X_k = x \mid Y_1, \dots, Y_k) = \frac{\mathbb{P}(X_k = x, Y_1, \dots, Y_k)}{\mathbb{P}(Y_1, \dots, Y_k)} = \eta \alpha_k(x) \quad (2.11)$$

2 where since  $\mathbb{P}(X_k = x \mid Y_1, \dots, Y_k)$  is a legitimate probability distribution on  
3  $x$ , we have

$$\eta = \left( \sum_x \alpha_k(x) \right)^{-1} .$$

4 As simple as that. In order to estimate the state at time  $k$ , we run the forward  
5 algorithm to update variables  $\alpha_i(x)$  from  $i = 1, \dots, k$ . We can implement  
6 this using the matrix-vector multiplication in the previous section.

7 This is a commonly used algorithm known as the Bayes filter and is our  
8 first insight into state estimation.

9 **An important fact** Even if the filtering estimate is computed recursively  
10 using each observation as it arrives, the estimate is actually the probability of  
11 the current state given *all* past observations.

$$\mathbb{P}(X_k = x \mid Y_1, \dots, Y_k) \neq \mathbb{P}(X_k = x \mid Y_k)$$

12 This is an extremely important concept to remember, in state-estimation we  
13 are always interested in computing the state given all available observations.  
14 In the same context, is the following statement true?

$$\mathbb{P}(X_k = x \mid Y_1, \dots, Y_k) = \mathbb{P}(X_k = x \mid Y_k, X_{k-1})$$

## 15 2.4.4 Smoothing

16 For smoothing given observations till time  $t$ , we would like to compute

$$\mathbb{P}(X_k = x \mid Y_1, \dots, Y_t)$$

17 for all time instants  $k = 1, \dots, t$ . Observe the filtering

$$\begin{aligned} \mathbb{P}(X_k = x \mid Y_1, \dots, Y_t) &= \frac{\mathbb{P}(X_k = x, Y_1, \dots, Y_t)}{\mathbb{P}(Y_1, \dots, Y_t)} \\ &= \frac{\mathbb{P}(X_k = x, Y_1, \dots, Y_k, Y_{k+1}, \dots, Y_t)}{\mathbb{P}(Y_1, \dots, Y_t)} \\ &= \frac{\mathbb{P}(Y_{k+1}, \dots, Y_t \mid X_k = x, Y_1, \dots, Y_k) \mathbb{P}(X_k = x, Y_1, \dots, Y_k)}{\mathbb{P}(Y_1, \dots, Y_t)} \\ &= \frac{\mathbb{P}(Y_{k+1}, \dots, Y_t \mid X_k = x) \mathbb{P}(X_k = x, Y_1, \dots, Y_k)}{\mathbb{P}(Y_1, \dots, Y_t)} \\ &= \frac{\beta_k(x) \alpha_k(x)}{\mathbb{P}(Y_1, \dots, Y_t)} \end{aligned} \quad (2.12)$$

18 Study the first step carefully, the numerator is *not* equal to  $\alpha_k(x)$  because  
19 observations go all the way till time  $t$ . The final step uses both the Markov  
20 and the HMM properties: future observations  $Y_{k+1}, \dots, Y_t$  depend only upon

1 future states  $X_{k+1}, \dots, X_t$  (HMM property) which are independent of the  
 2 past observations and states give the current state  $X_k = x$  (Markov property).

3 Smoothing can therefore be implemented by running the forward algorithm  
 4 to update  $\alpha_k$  from  $k = 1, \dots, t$  and the backward algorithm to update  $\beta_k$  from  
 5 time  $k = t, \dots, 1$ .

6 To see an example of smoothing in action, see [ORB-SLAM 2](#). What do  
 7 you think is the state of the Markov chain in this video?

8 **Example for the Whack-the-mole problem** Let us assume that we do not  
 9 see which hole the mole surfaces from (say it is dark outside) but we can hear  
 10 it. Our hearing is not very precise so we have an observation probabilities

$$M = \begin{bmatrix} 0.6 & 0.2 & 0.2 \\ 0.2 & 0.6 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}.$$

11 Assume that the mole surfaces three times and we make the measurements

$$Y_1 = 1, Y_2 = 3, Y_3 = 3.$$

12 We want to compute the distribution of the states the mole could be in at each  
 13 time. Assume that the we know the mole was in hole 1 at the first step, i.e.,  
 14  $\pi_1 = (1, 0, 0)$  for the Markov chain, like we had in Section 2.3.

15 Run the forward backward algorithm and see that

$$\alpha_1 = (0.6, 0, 0), \alpha_2 = (0.012, 0.048, 0.18), \alpha_3 = (0.0041, 0.0226, 0.0641),$$

16 and

$$\beta_3 = (1, 1, 1), \beta_2 = (0.4, 0.44, 0.36), \beta_1 = (0.1512, 0.1616, 0.1392).$$

17 Using these, we can now compute the filtering and the smoothing state distri-  
 18 butions, let us denote them by  $\pi^f$  and  $\pi^s$  respectively.

$$\pi_1^f = (1, 0, 0), \pi_2^f = (0.05, 0.2, 0.75), \pi_3^f = (0.045, 0.2487, 0.7063)$$

19 and

$$\pi_1^s = (1, 0, 0), \pi_2^s = (0.0529, 0.2328, 0.7143), \pi_3^s = (0.045, 0.2487, 0.7063).$$

20

## 21 2.4.5 Prediction

22 We would like to compute the future probability of the state give observations  
 23 up to some time

$$P(X_k = x | Y_1, \dots, Y_t) \quad \text{for } t < k.$$

24 Here is a typical scenario when you would need this estimate. Imagine that  
 25 you are tracking the position of a car using images from your camera. You are

🔗 Both the filtering problem and the smoothing problem give us the probability of the state given observations. Discuss which one should we should use in practice and why?

🔗 Do you notice any pattern in the solution returned by the filtering and the smoothing problem? Explain why that is the case.

using a deep network to detect the car in each image  $Y_k$  and since the neural network is quite slow, the car moves multiple time steps forward before you get the next observation. As you can appreciate, it would help us compute a more accurate estimate of the conditional probability of  $X_k = x$  if we propagated the position of the car in between successive observations using our Markov chain. This is easy to do.

1. We compute the filtering estimate  $\pi_t^f = P(X_t = x | Y_1, \dots, Y_t)$ , using the forward algorithm.
2. Propagate the Markov chain forward for  $k - t$  time-steps using  $\pi_t^f$  as the initial condition using

$$\pi_{i+1} = T^t \pi_i.$$

## 2.4.6 Decoding: Viterbi's Algorithm

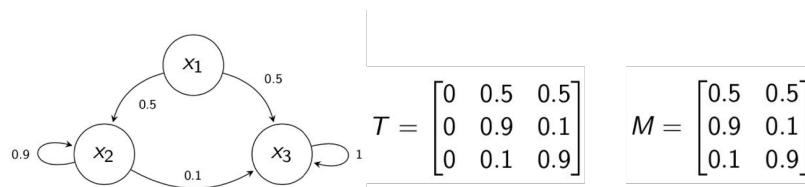
Both filtering and smoothing calculate the probability distribution of the state at time  $k$ . For instance, after recording a few observations, we can compute the probability distribution of the position of the car at each time instant. How do we get most likely trajectory of the car? One option is to choose

$$\hat{X}_k = \underset{x}{\operatorname{argmax}} P(X_k = x | Y_1, \dots, Y_t)$$

at each instant and output

$$(\hat{X}_1, \dots, \hat{X}_t)$$

as the answer. This is however only the point-wise best estimate of the state. This sequence may not be the most likely trajectory of Markov chain underlying our HMM. In the decoding problem, we are interested in computing the most likely state trajectory, not the point-wise most likely sequence of states. Let us take an example of the Whack-the-mole again. We will use a slightly different Markov chain shown below.



There are three states  $x_1, x_2, x_3$  with known initial distribution  $\pi = (1, 0, 0)$  and transition probabilities and observations given by matrices  $T, M$  respectively. Let us say that we only have two observations  $\{y_2, y_3\}$  this time and get the observation sequence

$$(2, 3, 3, 2, 2, 2, 3, 2, 3)$$

from our sensor. The filtering estimates are as follows.

$t$	$x_1$	$x_2$	$x_3$
1	<b>1.0000</b>	0	0
2	0	0.1000	<b>0.9000</b>
3	0	0.0109	<b>0.9891</b>
4	0	0.0817	<b>0.9183</b>
5	0	0.4165	<b>0.5835</b>
6	0	<b>0.8437</b>	0.1563
7	0	0.2595	<b>0.7405</b>
8	0	<b>0.7328</b>	0.2672
9	0	0.1771	<b>0.8229</b>

1

2 The most likely state at each instant is marked in blue. The point-wise most  
3 likely sequence of states is

$$(1, 3, 3, 3, 3, 2, 3, 2, 3).$$

4 Observe that this is not even feasible for the Markov chain. The transition  
5 from  $x_3 \rightarrow x_2$  is not even possible so this answer is clearly wrong. Let us  
6 look at the smoothing estimates.

$t$	$x_1$	$x_2$	$x_3$
1	<b>1.0000</b>	0	0
2	0	<b>0.6297</b>	0.3703
3	0	<b>0.6255</b>	0.3745
4	0	<b>0.6251</b>	0.3749
5	0	<b>0.6218</b>	0.3782
6	0	<b>0.5948</b>	0.4052
7	0	0.3761	<b>0.6239</b>
8	0	0.3543	<b>0.6457</b>
9	0	0.1771	<b>0.8229</b>

7

8 The point-wise most likely states in this case are feasible

$$(1, 2, 2, 2, 2, 2, 3, 3, 3).$$

9 Because the smoothing estimate at time  $k$  also takes into account the observa-  
10 tions from the future  $t > k$ , it effectively eliminates the impossible transition  
11  $x_3 \rightarrow x_2$ . This is still not however the most likely trajectory.

12 We will exploit the Markov property again to calculate the most likely  
13 state trajectory recursively. Let us define the “decoding variables” as

$$\delta_k(x) = \max_{(X_1, \dots, X_{k-1})} \mathbb{P}(X_1 = X_1, \dots, X_{k-1} = X_{k-1}, X_k = x, Y_1, \dots, Y_k); \quad (2.13)$$

14 this is the joint probability of the most likely state trajectory that ends at the  
15 state  $x$  at time  $k$  while generating observations  $Y_1, \dots, Y_k$ . We can now see  
16 that

$$\delta_{k+1}(x) = \max_{x'} \delta_k(x') T_{x'x} M_{x, y_{k+1}}; \quad (2.14)$$

17 the joint probability that the most likely trajectory ends up at state  $x$  at time  
18  $k + 1$  is the maximum of among the joint probabilities that end up at any state

1  $x'$  at time  $k$  multiplied by the one-step state transition  $T_{x'x}$  and observation  
 2  $M_{xy_{k+1}}$  probabilities. We would like to iterate upon this identity to find the  
 3 most likely path. The key idea is to maintain a pointer to the parent state  
 4  $\text{parent}_k(x)$  of the most likely trajectory, i.e., the state from which you could  
 5 have reached  $X_k = x$  given observations. Let us see how.

**Viterbi's algorithm** First initialize

$$\delta_1(x) = \pi_x M_{xy_1}$$

$$\text{parent}_k(x) = \text{null}.$$

for all states  $x$ . For all times  $k = 1, \dots, t - 1$ , for all states  $x$ , update

$$\delta_{k+1}(x) = \max_{x'} \delta_k(x') T_{x'x} M_{x,y_{k+1}}$$

$$\text{parent}_{k+1}(x) = \operatorname{argmax}_{x'} (\delta_k(x') T_{x'x}).$$

The most likely final state is

$$\hat{x}_t = \operatorname{argmax}_{x'} \delta_t(x')$$

and we can now backtrack using our parent pointers to find the most likely trajectory that leads to this state

$$\hat{x}_k = \text{parent}_{k+1}(\hat{x}_{k+1}).$$

The most likely trajectory given observations is

$$\hat{x}_1, \hat{x}_2, \dots, \hat{x}_t$$

and the joint probability of this trajectory and all observations is

$$P(X_1 = \hat{x}_1, \dots, X_t = \hat{x}_t, Y_1 = y_1, \dots, Y_t = y_t) = \delta_t(\hat{x}_t).$$

6 This is a very widely used algorithm, both in robotics and in other areas  
 7 such as speech recognition (given audio, find the most likely sentence spoken  
 8 by the person), wireless transmission and reception, DNA analysis (e.g., the  
 9 state of the Markov chain is the sequence ACTG... and our observations are  
 10 functions of these states at periodic intervals). Its name comes from Andrew  
 11 Viterbi who developed the algorithm in the late 60s, he is one of the founders  
 12 of Qualcomm Inc.

13 Here is how Viterbi's algorithm would look like for our whack-the-model  
 14 example.

$$\delta_1 = (0.6, 0, 0), \delta_2 = (0.012, 0.048, 0.18), \delta_3 = (0.0038, 0.0216, 0.0432)$$

$$\text{parent}_1 = (\text{null}, \text{null}, \text{null}), \text{parent}_2 = (1, 1, 1), \text{parent}_3 = (2, 3, 3).$$

15 The most likely path is the one that ends in 3 with joint probability 0.0432.

- 1 This path is (1, 3, 3).
- 2 Let us also compute Viterbi's algorithm for a longer observation sequence.

$t$	$x_1$	$x_2$	$x_3$
1	0.5/0	0	0
2	0/1	0.025/1	0.225/1
3	0/1	0.00225/2	0.2025/3
4	0/1	0.0018225/2	0.02025/3
5	0/1	0.0014762/2	0.002025/3
6	0/1	0.0011957/2	0.0002025/3
7	0/1	0.00010762/2	0.00018225/3
8	0/1	8.717e-05/2	1.8225e-05/3
9	0/1	7.8453e-06/2	1.6403e-05/3

- 3
- 4 The most likely trajectory is
 
$$(1, 3, 3, 3, 3, 3, 3, 3, 3).$$
- 5 Notice that if we had only 8 observations, the most likely trajectory would be
 
$$(1, 2, 2, 2, 2, 2, 2, 2).$$

6  
 7 What is the computational complexity of Viterbi's algorithm? It is line  
 8 in the time-horizon  $t$  and quadratic in the number of states in the Markov  
 9 chain. We are plucking out the most likely trajectory out of  $\text{card}(X)^t$  possible  
 10 trajectories using the  $\delta_k$  variables. Does this remind you of some other problem  
 11 that you may have seen before?

### 12 2.4.7 Shortest path on a Trellis graph

13 You may have seen Dijkstra's algorithm before that computes the shortest path  
 14 to reach a node in the graph given costs of traversing every edge.

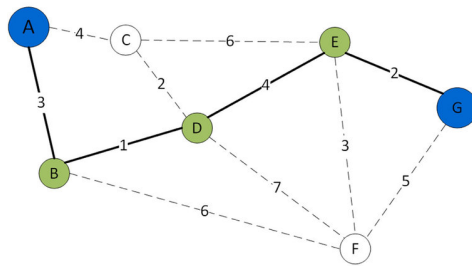


Figure 2.2: A graph with costs assigned to every edge. Dijkstra's algorithm finds the shortest path in this graph between nodes A and B using dynamic programming.

15 In the case of Viterbi's algorithm, we are also interested in finding the

**i** Just like the Bayes filter, Viterbi's algorithm is typically implemented using  $\log \delta_k(x)$  to avoid numerical underflows. This is particularly important for Viterbi's algorithm: since  $\delta_k(x)$  is the probability of an entire state and observation trajectory it can get small very quickly for unlikely states (as seen in the figure).

1 most likely path. For example we can write our joint probabilities as

$$P(X_1, X_2, X_3 | Y_1, Y_2, Y_3) = \frac{P(Y_1 | X_1) P(Y_2 | X_2) P(Y_3 | X_3) P(X_1) P(X_2 | X_1) P(X_3 | X_2)}{P(Y_1, Y_2, Y_3)}.$$

$$\Rightarrow \log P(X_1, X_2, X_3 | Y_1, Y_2, Y_3) = \log P(Y_1 | X_1) + \log P(Y_2 | X_2) + \log P(Y_3 | X_3) \\ + \log P(X_1) + \log P(X_2 | X_1) + \log P(X_3 | X_2) - \log P(Y_1, Y_2, Y_3).$$

2 To find the most likely trajectory, we want to minimize  $-\log P(X_1, X_2, X_3 | Y_1, Y_2, Y_3)$ .

3 The term  $\log P(Y_1, Y_2, Y_3)$  does not depend on  $X_1, X_2, X_3$  and is a constant

4 as far as the most likely path given observations is concerned. We can now

5 write down the “Trellis” graph as shown below.

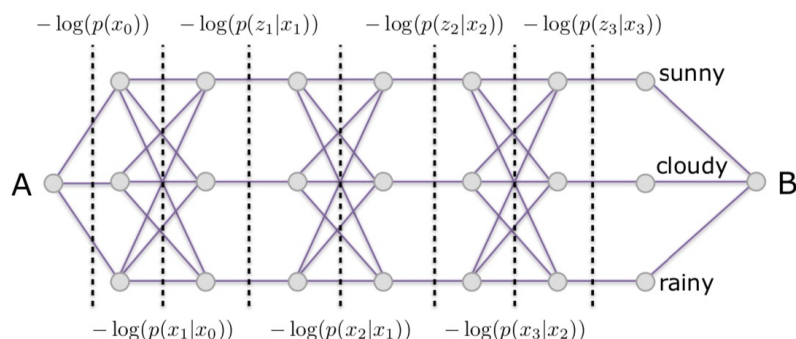


Figure 2.3: A Trellis graph for a 3-state HMM for a sequence of three observations. Disregard the subscript  $x_0$ .

6 Each edge is either the log-probability of the transition of the Markov  
 7 chain, or it is the log-probability of the receiving the observation given a  
 8 state. We create a dummy initial node A and a dummy terminal node B. The  
 9 edge-costs of the final three states, in this case sunny/cloudy/rainy, are zero.  
 10 The costs from node A to the respective states are the log-probabilities of the  
 11 initial state distribution. Dijkstra’s algorithm, which we will study in Module 2  
 12 in more details, now gives the shortest path on the Trellis graph. This approach  
 13 is the same as that of the Viterbi’s algorithm: our parent pointers  $\text{parent}_k(x)$   
 14 are the parent nodes in Dijkstra’s algorithm and our delta variables  $\delta_k(x)$  is the  
 15 cost of each node in the Trellis graph maintained by the Dijkstra’s algorithm.

## 16 2.5 Learning an HMM from observations

17 In the previous sections, given an HMM that had an initial distribution  $\pi$  for the  
 18 Markov chain, a transition matrix  $T$  for the Markov chain and an observation  
 19 matrix  $M$

$$\lambda = (\pi, T, M)$$

20 we computed various quantities such as

$$P(Y_1, \dots, Y_t; \lambda)$$

1 for an observation sequence  $Y_1, \dots, Y_t$  of the HMM. Given an observation  
 2 sequence, we can also go back and update our HMM to make this observation  
 3 sequence more likely. This is the simplest instance of *learning* an HMM. The  
 4 prototypical problem to imagine that our original HMM  $\lambda$  comes from our  
 5 knowledge of the original problem (say a physics model of the dynamics of a  
 6 robot and its sensors). Given more data, namely the observations, we want to  
 7 update this model. The most natural way to update the model is to maximize  
 8 the likelihood of observations given our model, i.e.,

$$\lambda^* = \underset{\lambda}{\operatorname{argmax}} P(Y_1, \dots, Y_t; \lambda).$$

9 This is known as maximum-likelihood estimation (MLE). In this section  
 10 we will look at the Baum-Welch algorithm which solves the MLE problem  
 11 iteratively. Given  $\lambda$ , it finds a new HMM  $\lambda' = (\pi', T', M')$  (the  $'$  denotes a  
 12 new matrix, not the transpose here) such that

$$P(Y_1, \dots, Y_t; \lambda') > P(Y_1, \dots, Y_t; \lambda).$$

13 Let us consider a simple problem. We are going to imagine that the FBI  
 14 is trying to catch the dangerous criminal Keyser Soze who is known to travel  
 15 between two cities Los Angeles (LA) which will be state  $x_1$  and New York  
 16 City (NY) which will be state  $x_2$ . The FBI initially have no clue about his  
 17 whereabouts, so their initial belief on his location is uniform  $\pi = [0.5, 0.5]$ .  
 18 His movements are modeled using a Markov chain

$$T = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix},$$

19 e.g., if Soze is in LA, he is likely to stay in LA or go to NY with equal  
 20 probability. The FBI can make observations about him, they either observe  
 21 him to be in LA ( $y_1$ ), NY ( $y_2$ ) or do not observe anything at all (null,  $y_3$ ).

$$M = \begin{bmatrix} 0.4 & 0.1 & 0.5 \\ 0.1 & 0.5 & 0.4 \end{bmatrix}.$$

22 Say that they received an observation sequence of 20 periods

(null, LA, LA, null, NY, null, NY, NY, NY, null, NY, NY, NY, NY, NY, null, null, LA, LA, NY).

23 Can we say something about the probability of Soze's movements? At each  
 24 time  $k$  we can compute

$$\gamma_k(x) := P(X_k = x \mid Y_1, \dots, Y_t)$$

25 the smoothing probability. We can also compute the most likely state trajectory  
 26 he could have take given our observations using decoding. Let us focus on the  
 27 smoothing probabilities  $\gamma_k(x)$  as shown below.



$t$	LA	NY
1	0.5556	0.4444
2	0.8000	0.2000
3	0.8000	0.2000
...	...	...
18	0.8000	0.2000
19	0.8000	0.2000
20	0.1667	0.8333

1

2 The point-wise most likely sequence of states after doing so turns out to be

(LA, LA, LA, LA, NY, LA, NY, NY, NY, LA, NY, NY, NY, NY, NY, LA, LA, LA, LA, NY).

3 Notice how smoothing fills in the missing observations above.

4 **Expected state visitation counts** The next question we should ask is how  
 5 should we update the model  $\lambda$  given this data. We are going to learn the entries  
 6 of the state-transition using

$$T'_{x,x'} = \frac{E[\text{number of transitions from } x \text{ to } x']}{E[\text{number of times the Markov chain was in state } x]}.$$

7 What is the denominator, it is simply the sum of the probabilities that the  
 8 Markov chain was at state  $x$  at time  $1, 2, \dots, t-1$  given our observations, i.e.,

$$E[\text{number of times the Markov chain was in state } x] = \sum_{k=1}^{t-1} \gamma_k(x).$$

9 The numerator is given in a similar fashion. We will define a quantity

$$\begin{aligned} \xi_k(x, x') &:= P(X_k = x, X_{k+1} = x' | Y_1, \dots, Y_t) \\ &= \eta \alpha_k(x) T_{x,x'} M_{x',y_{k+1}} \beta_{k+1}(x'); \end{aligned} \quad (2.15)$$

🔗 Derive the expression for  $\xi_k(x, x')$  for yourself.

10 where  $\eta$  is a normalizing constant such that  $\sum_{x,x'} \xi_k(x, x') = 1$ . Observe that  
 11  $\xi_k$  is the joint probability of  $X_k$  and  $X_{k+1}$

$$\begin{aligned} \xi_k(x, x') &= P(X_{k+1} = x' | X_k = x, Y_1, \dots, Y_t) \gamma_k(x) \\ &\neq T_{x,x'} \gamma_k(x) \\ &= P(X_{k+1} = x' | X_k = x) P(X_k = x | Y_1, \dots, Y_t). \end{aligned}$$

12 The expected value of transitioning between states  $x$  and  $x'$  is

$$E[\text{number of transitions from } x \text{ to } x'] = \sum_{k=1}^{t-1} \xi_k(x, x').$$

1 This gives us our new state transition matrix, you will see in the homework  
2 that it comes to be

$$T' = \begin{bmatrix} 0.47023 & 0.52976 \\ 0.35260 & 0.64739 \end{bmatrix}.$$

3 This is a much better informed FBI than the other we had before beginning  
4 the problem where the transition matrix was all 0.5s.

5 **The new initial distribution** What is the new initial distribution for the  
6 HMM? Recall that we are trying to compute the best HMM given the observa-  
7 tions, so if the initial distribution was

$$\pi = P(X_1)$$

8 before receiving any observations from the HMM, it is now

$$\pi' = P(X_1 | Y_1, \dots, Y_t) = \gamma_1(x);$$

9 the smoothing estimate at the first time-step.

10 **Updating the observation matrix** We can use a similar logic at the ex-  
11 pected state visitation counts to write

$$\begin{aligned} M'_{x,y} &= \frac{E[\text{number of times in state } x, \text{ when observation was } y]}{E[\text{number of times the Markov chain was in state } x]} \\ &= \frac{\sum_{k=1}^t \gamma_k(x) \mathbf{1}_{\{y_k=y\}}}{\sum_{k=1}^t \gamma_k(x)}. \end{aligned}$$

12 You will see in your homework problem that this matrix comes up to be

$$M' = \begin{bmatrix} 0.39024 & 0.20325 & 0.40650 \\ 0.06779 & 0.706214 & 0.2259 \end{bmatrix}.$$

13 Notice how the observation probabilities for the unknown state  $y_3$  have gone  
14 down because the Markov chain does not have those states.

15 The ability to start with a rudimentary model of the HMM and update it  
16 using observations is quite revolutionary. Baum et al. proved in the paper  
17 Baum, Leonard E., et al. "A maximization technique occurring in the stati-  
18 statistical analysis of probabilistic functions of Markov chains." The annals of  
19 mathematical statistics 41.1 (1970): 164-171. Discuss the following questions:

- 20 • When do we stop in our iterated application of the Baum-Welch algo-  
21 rithm?
- 22 • Are we always guaranteed to find the same HMM irrespective of our  
23 initial HMM?
- 24 • If our initial HMM  $\lambda$  is the same, are guaranteed to find the same HMM  
25  $\lambda'$  across two different iterations of the Baum-Welch algorithm?
- 26 • How many observations should we use to update the HMM?

# 1 Chapter 3

## 2 Kalman Filter and its 3 variants

### Reading

1. Barfoot, Chapter 3, 4 for Kalman filter
2. Thrun, Chapter 3 for Kalman filter, Chapter 4 for particle filters
3. Russell Chapter 15.4 for Kalman filter

4 Hidden Markov Models (HMMs) which we discussed in the previous  
5 chapter were a very general class of models. As a consequence algorithms for  
6 filtering, smoothing and decoding that we prescribed for the HMM are also  
7 very general. In this chapter we will consider the situation when we have a little  
8 more information about our system. Instead of writing the state transition and  
9 observation matrices as arbitrary matrices, we will use the framework of linear  
10 dynamical systems to model them better. Since we know the system a bit better,  
11 algorithms that we prescribe for these models for solving filtering, smoothing  
12 and decoding will also be more efficient. We will almost exclusively focus on  
13 the filtering problem in this chapter. The other two, namely smoothing and  
14 decoding, can also be solved easily using these ideas but are less commonly  
15 used for these systems.

### 16 3.1 Background

17 **Multi-variate random variables and linear algebra** For  $d$ -dimensional  
18 random variables  $X, Y \in \mathbb{R}^d$  we have

$$E[X + Y] = E[X] + E[Y];$$

1 this is actually more surprising than it looks, it is true regardless of whether  
2  $X, Y$  are correlated. The covariance matrix of a random variable is defined as

$$\text{cov}(X) = \mathbb{E}[(X - \mathbb{E}[X]) (X - \mathbb{E}[X])^\top];$$

3 we will usually denote this by  $\Sigma \in \mathbb{R}^{d \times d}$ . Note that the covariance matrix is,  
4 by construction, symmetric and positive semi-definite. This means it can be  
5 factorized as

$$\Sigma = U \Lambda U^\top$$

6 where  $U \in \mathbb{R}^{d \times d}$  is an orthonormal matrix (i.e.,  $UU^\top = I$ ) and  $\Lambda$  is a  
7 diagonal matrix with non-negative entries. The trace of a matrix is the sum of  
8 its diagonal entries. It is also equal to the sum of its eigenvalues, i.e.,

$$\text{tr}(\Sigma) = \sum_{i=1}^d \Sigma_{ii} = \sum_{i=1}^d \lambda_i(\Sigma)$$

9 where  $\lambda_i(S) \geq 0$  is the  $i^{\text{th}}$  eigenvalue of the covariance matrix  $S$ . The trace  
10 is a measure of the uncertainty in the multi-variate random variable  $X$ , if  $X$   
11 is a scalar and takes values in the reals then the covariance matrix is also, of  
12 course, a scalar  $\Sigma = \sigma^2$ .

13 A few more identities about the matrix trace that we will often use in this  
14 chapter are as follows.

- 15 • For matrices  $A, B$  we have

$$\text{tr}(AB) = \text{tr}(BA);$$

16 the two matrices need not be square themselves, only their product does.

- 17 • For  $A, B \in \mathbb{R}^{m \times n}$

$$\text{tr}(A^\top B) = \text{tr}(B^\top A) = \sum_{i=1}^m \sum_{j=1}^n B_{ij} A_{ij}.$$

18 This operation can be thought of as taking the inner product between  
19 two matrices.

20 **Gaussian/Normal distribution** We will spend a lot of time working with  
21 the Gaussian/Normal distribution. The multi-variate  $d$ -dimensional Normal  
22 distribution has the probability density

$$f(x) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left\{-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right\}$$

23 where  $\mu \in \mathbb{R}^d$ ,  $\Sigma \in \mathbb{R}^{d \times d}$  denote the mean and covariance respectively. You  
24 should commit this formula to memory. In particular remember that

$$\int_{x \in \mathbb{R}^d} \exp\left\{-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right\} dx = \sqrt{\det(2\pi\Sigma)}$$

🔗 Why is it so ubiquitous?

- 1 which is simply expressing the fact that the probability density function inte-  
 2 grates to 1.

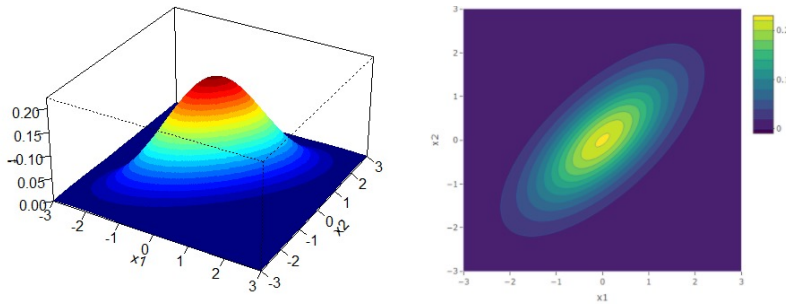


Figure 3.1: Probability density (left) and iso-probability contours (right) of a bi-variate Normal distribution. Warm colors denote regions of high probability.

- 3 Given two Gaussian rvs.  $X, Y \in \mathbb{R}^d$  and  $Z = X + Y$  we have

$$\mathbb{E}[Z] = \mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

- 4 with covariance

$$\text{cov}(Z) = \Sigma_Z = \Sigma_X + \Sigma_Y + \Sigma_{XY} + \Sigma_{YX}$$

- 5 where

$$\mathbb{R}^{d \times d} \ni \Sigma_{XY} = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])^\top];$$

- 6 the matrix  $\Sigma_{YX}$  is defined similarly. If  $X, Y$  are independent (or uncorrelated)  
 7 the covariance simplifies to

$$\Sigma_Z = \Sigma_X + \Sigma_Y.$$

- 8 If we have a linear function of a Gaussian random variable  $X$  given by  
 9  $Y = AX$  for some *deterministic* matrix  $A$  then  $Y$  is also Gaussian with mean

$$\mathbb{E}[Y] = \mathbb{E}[AX] = A\mathbb{E}[X] = A\mu_X$$

- 10 and covariance

$$\begin{aligned} \text{cov}(Y) &= \mathbb{E}[(AX - A\mu_X)(AX - A\mu_X)^\top] \\ &= \mathbb{E}[A(X - \mu_X)(X - \mu_X)^\top A^\top] \\ &= A\mathbb{E}[(X - \mu_X)(X - \mu_X)^\top]A^\top \\ &= A\Sigma_X A^\top. \end{aligned} \tag{3.1}$$

- 11 This is an important result that you should remember.

## 3.2 Linear state estimation

With that background, let us now look at the basic estimation problem. Let  $X \in \mathbb{R}^d$  denote the true state of a system. We would like to build an estimator for this state, this is denote by

$$\hat{X}.$$

An estimator is any quantity that indicates our belief of what  $X$  is. The estimator is created on the basis of observations and we will therefore model it as a random variable. We would like the estimator to be unbiased, i.e.,

$$E[\hat{X}] = X;$$

this expresses the concept that if we were to measure the state of the system many times, say using many sensors or multiple observations from the same sensor, the resultant estimator  $\hat{X}$  is correct on average. The error in our belief is

$$\tilde{X} = \hat{X} - X.$$

The error is zero-mean  $E[\tilde{X}] = 0$  and its covariance  $\Sigma_{\tilde{X}}$  is called the covariance of the estimator.

**Optimally combining two estimators** Let us now imagine that we have two estimators  $\hat{X}_1$  and  $\hat{X}_2$  for the same true state  $X$ . We will assume that the two estimators were created independently (say different sensors) and therefore are conditionally independent random variables given the true state  $X$ . Say both of them are unbiased but each of them have a certain covariance of the error

$$\Sigma_{\tilde{X}_1} \text{ and } \Sigma_{\tilde{X}_2}.$$

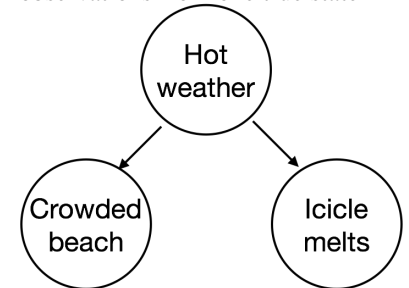
We would like to combine the two to obtain a *better* estimate of what the state could be. *Better* can mean many different quantities depending upon the problem but in general in this course we are interested in improving the error covariance. Our goal is then

Given two estimators  $\hat{X}_1$  and  $\hat{X}_2$  of the true state  $X$  combine them to obtain a new estimator

$$\hat{X} = \text{some function}(\hat{X}_1, \hat{X}_2)$$

which has the best error covariance  $\text{tr}(\Sigma_{\tilde{X}})$ .

📌 Conditionally independent observations from one true state



### 3.2.1 One-dimensional Gaussian random variables

Consider the case when  $\hat{X}_1, \hat{X}_2 \in \mathbb{R}$  are Gaussian random variables with means  $\mu_1, \mu_2$  and variances  $\sigma_1^2, \sigma_2^2$  respectively. Assume that both are unbiased estimators of  $X \in \mathbb{R}$ . Let us combine them linearly to obtain a new estimator

$$\hat{X} = k_1 \hat{X}_1 + k_2 \hat{X}_2.$$

- 1 How should we pick the coefficients  $k_1, k_2$ ? We would of course like the new  
2 estimator to be unbiased, so

$$\begin{aligned} \mathbb{E}[\hat{X}] &= \mathbb{E}[k_1\hat{X}_1 + k_2\hat{X}_2] = (k_1 + k_2)X = X \\ \Rightarrow k_1 + k_2 &= 1. \end{aligned}$$

- 3 The variance of the  $\hat{X}$  is

$$\text{var}(\hat{X}) = k_1^2\sigma_1^2 + k_2^2\sigma_2^2 = k_1^2\sigma_1^2 + (1 - k_1)^2\sigma_2^2.$$

- 4 The optimal  $k_1$  that leads to the smallest variance is thus given by

$$k_1 = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}.$$

- 5 We set the derivative of  $\text{var}(\hat{X})$  with respect to  $k_1$  to zero to get this. The final  
6 estimator is

$$\hat{X} = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2} \hat{X}_1 + \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2} \hat{X}_2. \quad (3.2)$$

- 7 It is unbiased of course and has variance

$$\sigma_{\hat{X}}^2 = \frac{\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2}.$$

- 8 Notice that since  $\sigma_2^2/(\sigma_1^2 + \sigma_2^2) < 1$ , the variance of the new estimator is  
9 smaller than either of the original estimators. This is an important fact to  
10 remember, combining two estimators *always* results in a better estimator.

### 11 **Some comments about the optimal combination.**

- 12 • It is easy to see that if  $\sigma_2 \gg \sigma_1$  then the corresponding estimator, namely  
13  $\hat{X}_2$  gets less weight in the combination. This is easy to understand, if  
14 one of our estimates is very noisy, we should rely less upon it to obtain  
15 the new estimate. In the limit that  $\sigma_2 \rightarrow \infty$ , the second estimator is not  
16 considered at all in the combination.
- 17 • If  $\sigma_1 = \sigma_2$ , the two estimators are weighted equally and since  $\sigma_{\hat{X}}^2 =$   
18  $\sigma_1^2/2$  the variance reduces by half after combination.
- 19 • The minimal variance of the combined estimator is not zero. This is easy  
20 to see because if we have two noisy estimates of the state, combining  
21 them need not lead to us knowing the true state with certainty.

### 22 **3.2.2 General case**

- 23 Let us now perform the same exercise for multi-variate Gaussian random  
24 variables. We will again combine the two estimators linearly to get

$$\hat{X} = K_1\hat{X}_1 + K_2\hat{X}_2$$

1 where  $K_1, K_2 \in \mathbb{R}^{d \times d}$  are matrices that we would like to choose. In order for  
2 the estimator to be unbiased we again have the condition

$$\begin{aligned} \mathbb{E}[\hat{X}] &= \mathbb{E}[K_1 \hat{X}_1 + K_2 \hat{X}_2] = (K_1 + K_2)X = X \\ &\Rightarrow K_1 + K_2 = I_{d \times d}. \end{aligned}$$

3 The covariance of  $\hat{X}$  is

$$\begin{aligned} \Sigma_{\hat{X}} &= K_1 \Sigma_1 K_1^\top + K_2 \Sigma_2 K_2^\top \\ &= K_1 \Sigma_1 K_1^\top + (I - K_1) \Sigma_2 (I - K_1)^\top. \end{aligned}$$

4 Just like the minimized the variance in the scalar case, we will minimize the  
5 trace of this covariance matrix. We know that the original covariances  $\Sigma_1$  and  
6  $\Sigma_2$  are symmetric. We will use the following identity for the partial derivative  
7 of a matrix product

$$\frac{\partial}{\partial A} \text{tr}(ABA^\top) = 2AB \quad (3.3)$$

8 for a symmetric matrix  $B$ . Minimizing  $\text{tr}(\Sigma_{\hat{X}})$  with respect to  $K_1$  amounts to  
9 setting

$$\frac{\partial}{\partial K_1} \text{tr}(\Sigma_{\hat{X}}) = 0$$

10 which yields

$$\begin{aligned} 0 &= K_1 \Sigma_1 - (I - K_1) \Sigma_2 \\ \Rightarrow K_1 &= \Sigma_2 (\Sigma_1 + \Sigma_2)^{-1} \text{ and } K_2 = \Sigma_1 (\Sigma_1 + \Sigma_2)^{-1}. \end{aligned}$$

11 The optimal way to combine the two estimators is thus

$$\hat{X} = \Sigma_2 (\Sigma_1 + \Sigma_2)^{-1} \hat{X}_1 + \Sigma_1 (\Sigma_1 + \Sigma_2)^{-1} \hat{X}_2. \quad (3.4)$$

12 You should consider the similarities of this expression with the one for the  
13 scalar case in (3.2). The same broad comments hold, i.e., if one of the  
14 estimators has a very large variance, that estimator is weighted less in the  
15 combination.

### 16 3.2.3 Incorporating Gaussian observations of a state

17 Let us now imagine that we have a sensor that can give us observations of  
18 the state. The development in this section is analogous to our calculations  
19 in Chapter 2 with the recursive application of Bayes rule or the observation  
20 matrix of the HMM. We will consider a special type of sensor that gives  
21 observations

$$\mathbb{R}^p \ni Y = CX + \nu \quad (3.5)$$

22 which is a linear function of the true state  $X \in \mathbb{R}^d$  with the matrix  $C \in \mathbb{R}^{p \times d}$   
23 being something that is unique to the particular sensor. This observation is not  
24 precise and we will model the sensor as having zero-mean Gaussian noise

$$\nu \sim N(0, Q)$$



1 of covariance  $Q \in \mathbb{R}^{p \times p}$ . Notice something important here, the dimensionality  
 2 of the observations need not be the same as the dimensionality of the state.  
 3 This should not be surprising, after all the the number of observations in the  
 4 HMM need not be the same as the number of the states in the Markov chain.

5 We will solve the following problem. Given an existing estimator  $\hat{X}'$  we  
 6 want to combine it with the observation  $Y$  to update the estimator to  $\hat{X}$ , in the  
 7 best way, i.e., in a way that gives the minimal variance. We will again use a  
 8 linear combination

$$\hat{X} = K' \hat{X}' + KY.$$

9 Again we want the estimator to be unbiased, so we set

$$\begin{aligned} \mathbb{E}[\hat{X}] &= \mathbb{E}[K' \hat{X}' + KY] \\ &= K' X + K \mathbb{E}[Y] \\ &= K' X + K \mathbb{E}[CX + \nu] \\ &= K' X + KCX \\ &= X. \end{aligned}$$

10 to get that

$$\begin{aligned} I &= K' + KC. \\ \Rightarrow \hat{X} &= (I - KC) \hat{X}' + KY \\ &= \hat{X}' + K(Y - C \hat{X}'). \end{aligned} \tag{3.6}$$

11 This is special form which you will do well to remember. The old estimator  
 12  $\hat{X}'$  gets an additive term  $K(Y - C \hat{X}')$ . For reasons that will soon become  
 13 clear, we call this term

$$\text{innovation} = Y - C \hat{X}'.$$

14 Let us now optimize  $K$  as before to compute the estimator with minimal  
 15 variance. We will make the following important assumption in this case.

We will assume that the observation  $Y$  is independent of the estimator  $\hat{X}'$  given  $X$ . This is a natural assumption because presumably our original estimator  $\hat{X}'$  was created using past observations and the present observation  $Y$  is therefore independent of it given the state  $X$ .

16 The covariance of  $\hat{X}$  is

$$\Sigma_{\hat{X}} = (I - KC) \Sigma_{\hat{X}'} (I - KC)^\top + KQK^\top.$$

We optimize the trace of  $\Sigma_{\hat{X}}$  with respect to  $K$  to get

$$\begin{aligned} 0 &= \frac{\partial}{\partial K} \text{tr}(\Sigma_{\hat{X}}) \\ 0 &= -2(I - KC) \Sigma_{\hat{X}'} C^\top + 2KQ \\ \Rightarrow \Sigma_{\hat{X}'} C^\top &= K(C \Sigma_{\hat{X}'} C^\top + Q) \\ \Rightarrow K &= \Sigma_{\hat{X}'} C^\top (C \Sigma_{\hat{X}'} C^\top + Q)^{-1}. \end{aligned}$$

- 1 The matrix  $K \in \mathbb{R}^{d \times p}$  is called the “Kalman gain” after **Rudolph Kalman** who  
 2 developed this method in the 1960s.

**Kalman gain** This is an important formula and it helps to have a mnemonic and a slightly simpler notation to remember it by. If  $\Sigma'$  is the covariance of the previous estimator,  $Q$  is the covariance of the zero-mean observation and  $C$  is the matrix that gives the observation from the state, then the Kalman gain is

$$K = \Sigma_{\hat{X}} C^\top (C \Sigma_{\hat{X}} C^\top + Q)^{-1}. \quad (3.7)$$

and the new estimator for the state is

$$\hat{X} = \hat{X}' + K(Y - C\hat{X}').$$

The covariance of the updated estimator  $\hat{X}$  is given by

$$\begin{aligned} \Sigma_{\hat{X}} &= (I - KC)\Sigma_{\hat{X}'}(I - KC)^\top + KQK^\top \\ &= \left( \Sigma_{\hat{X}'}^{-1} + C^\top Q^{-1} C \right)^{-1}. \end{aligned} \quad (3.8)$$

If  $C = I$ , the Kalman gain is the same expression as the optimal coefficient in (3.4). This should not be surprising because the observation is an estimator for the state.

The second expression for  $\Sigma_{\hat{X}}$  follows by substituting the value of the Kalman gain  $K$ . Yet another way of remembering this equation is to notice that

$$\begin{aligned} \Sigma_{\hat{X}}^{-1} &= \Sigma_{\hat{X}'}^{-1} + C^\top Q^{-1} C \\ K &= \Sigma_{\hat{X}}^{-1} C^\top Q^{-1} \\ \hat{X} &= \hat{X}' + \Sigma_{\hat{X}}^{-1} C^\top Q^{-1} (Y - C\hat{X}'). \end{aligned} \quad (3.9)$$

**i** Derive these expressions for the Kalman gain and the covariance yourself.

### 3 3.2.4 An example

- 4 Consider the scalar case when we have multiple measurements of some scalar  
 5 quantity  $x \in \mathbb{R}$  corrupted by noise.

$$y_i = x + \nu_i$$

- 6 where  $y_i \in \mathbb{R}$  and the scalar noise  $\nu_i \sim N(0, 1)$  is zero-mean and standard  
 7 Gaussian. Find the updated estimate of the state  $x$  after  $k$  such measurements;  
 8 this means both the mean and the covariance of the state.

- 9 You can solve this in two ways, you can either use the measurement matrix  
 10  $C$  to be  $\mathbb{1}_k = [1, \dots, 1]$  to be a vector of all ones and apply the formula  
 11 in (3.7) and (3.8) Show that the estimate  $\hat{x}_k$  after  $k$  measurements has mean

1 and covariance

$$\begin{aligned} \mathbb{E}[\hat{x}_k] &= \frac{1}{k} \sum_{i=1}^k y_i. \\ \text{cov}(\hat{x}_k) &= C^\top C^{-1} = \frac{1}{k}. \end{aligned}$$

2 If we take one more measurement  $y_{k+1} = x + \nu_{k+1}$  with noise  $\nu_{k+1} \sim$   
3  $N(0, \sigma^2)$ , show using (3.9) that

$$\begin{aligned} \text{cov}(\hat{x}_{k+1})^{-1} &= \text{cov}(\hat{x}_k)^{-1} + \frac{1}{\sigma^2} \\ \Rightarrow \text{cov}(\hat{x}_{k+1}) &= \frac{\sigma^2}{\sigma^2 k + 1}. \end{aligned}$$

4 The updated mean using (3.9) again

$$\begin{aligned} \mathbb{E}[\hat{x}_{k+1}] &= \hat{x}_k + \text{cov}(\hat{x}_{k+1}) \frac{1}{\sigma^2} (y_{k+1} - \hat{x}_k) \\ &= \hat{x}_k + \frac{y_{k+1} - \hat{x}_k}{\sigma^2 k + 1}. \end{aligned}$$

5 You will notice that if the noise on the  $k + 1^{\text{th}}$  observation is very small, even  
6 after  $k$  observations, the new estimate fixates on the latest observation

$$\sigma \rightarrow 0 \Rightarrow \hat{x}_{k+1} \rightarrow y_{k+1}.$$

7 Similarly, if the latest observation is very noisy, the estimate does not change  
8 much

$$\sigma \rightarrow \infty \Rightarrow \hat{x}_{k+1} \rightarrow \hat{x}_k.$$

### 9 3.3 Background on linear and nonlinear dynamical systems

The true state  $X$  need not static. We will next talk about models for how the state of the world evolves using ideas in dynamical systems.

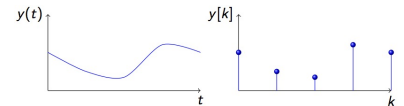
11 A continuous-time signal is a function that associates to each time  $t \in \mathbb{R}$  a real  
12 number  $y(t)$ . We denote signals by

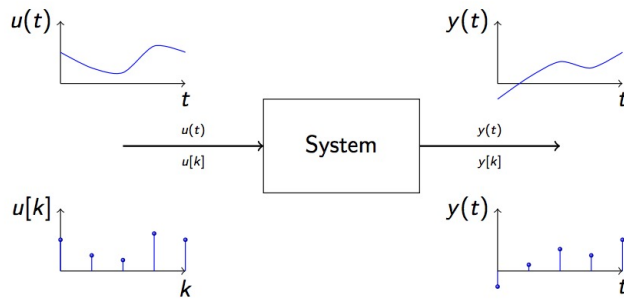
$$y : t \mapsto y(t).$$

13 Similarly a discrete-time signal is a function that associates to each integer  $k$  a  
14 real number  $y(k)$ , we have been denoting quantities like this by  $y_k$ .

15 A dynamical system is an operator (a box) that transforms an input signal  
16  $u(t)$  or  $u_k$  to an output  $y(t)$  or  $y_k$  respectively. We call the former a continuous-  
17 time system and the latter a discrete-time system.

**i** A continuous-time signal  $y(t)$  and discrete-time signal  $y_k$ .





Almost always in robotics, we will be interested in systems that are temporally *causal*, i.e., the output at time  $t_0$  is only a function the input *up to* time  $t_0$ . Analogously, the output at time  $k_0$  for a discrete-time system is dependent only on the input up to time  $k_0$ . Most systems in the physical world are temporally causal.

**State of a system** We know that if the system is causal, in order to compute its output at a time  $t_0$ , we only need to know all the input from time  $t = (-\infty, t_0]$ . This is a lot of information. The concept of a state, about which we have been cavalier until now helps with this. The state  $x(t_1)$  of a causal system at time  $t_1$  is the information needed, together with the input  $u$  between times  $t_1$  and  $t_2$  to *uniquely compute* the output  $y(t_2)$  at time  $t_2$ , for all times  $t_2 \geq t_1$ . In other words, the state of a system summarizes the whole history of what happened between  $(-\infty, t_1)$ .

Typically the state of a system is a  $d$ -dimensional vector in  $\mathbb{R}^d$ . The dimension of a system is the minimum  $d$  required to define a state.

### 3.3.1 Linear systems

A system is called a linear system if for any two input *signals*  $u_1$  and  $u_2$  and any two real numbers  $a, b$

$$\begin{aligned} u_1 &\rightarrow y_1 \\ u_2 &\rightarrow y_2 \\ au_1 + bu_2 &\rightarrow ay_1 + by_2. \end{aligned}$$

Linearity is a very powerful property. For instance, it suggests that if we can decompose a complicated input into the sum of simple signals, then the output of the system is also a sum of the outputs of these simple signals. For example, if we can write the input as a Fourier series  $u(t) = \sum_{i=0}^{\infty} a_i \cos(it) + b_i \sin(it)$  we can pass each of the terms in this summation to system and get the output of  $u(t)$  by summing up the individual outputs.

Finite-dimensional systems can be written using a set of differential equations as follows. Consider the spring-mass system. If  $z(t)$  denotes the position of the mass at time  $t$  and  $u(t)$  is the force that we are applying upon it at time

❓ Can you give an example of a dynamical system that is non-causal? Think of how a DVD Ripper, or a pre-programmed acrobatic maneuver on a plane works.

❓ Discuss some examples of the state.



Is the state of a system uniquely defined?

1  $t$ , the position of the mass satisfies the differential equation

$$m \frac{d^2 z(t)}{dt^2} + c \frac{dz(t)}{dt} + kx(t) = u(t)$$

$$\text{or } m\ddot{x} + c\dot{x} + kx = u$$

2 in short. Here  $m$  is the mass of the block,  $c$  is the damping coefficient of the  
3 spring and  $k$  is the spring force constant. Let us define

$$z_1(t) := z(t)$$

$$z_2(t) := \frac{dz(t)}{dt}.$$

4 We can now rewrite the dynamics as

$$\begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -k/m & -c/m \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u$$

### 5 3.3.2 Linear Time-Invariant (LTI) systems

6 If we define the state  $x(t) = \begin{bmatrix} z_1(t) \\ z_2(t) \end{bmatrix}$ , then the above equation can be written  
7 as

$$\dot{x}(t) = Ax(t) + Bu(t). \quad (3.10)$$

8 This is a linear system that takes in the input  $u(t)$  and has a state  $x(t)$ . You can  
9 check the conditions for linearity to be sure. It is also a linear time-invariant  
10 (LTI) system because the matrices  $A, B$  do not change with time  $t$ . **The input**  
11  **$u(t)$  is also typically called the control (or action, or the control input)**  
12 **and essentially the second half of the course is about computing good**  
13 **controls.**

14 Since the state at time  $t$  encapsulates everything that happened to the  
15 system due to the inputs  $\{u(-\infty), u(t)\}$ , we can say that the system computes  
16 its output  $y(t)$  as a function of the state  $x(t)$  and the latest input  $u(t)$

$$y(t) = \text{function}(x(t), u(t))$$

17 If this function is linear we have

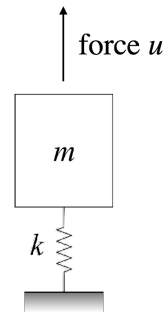
$$y(t) = Cx(t) + Du(t). \quad (3.11)$$

18 The pair of equations (3.10) and (3.11) together are the so-called state-space  
19 model of an LTI system. The development for discrete-time systems is com-  
20 pletely analogous, we will have

$$x_{k+1} = Ax_k + Bu_k$$

$$y_k = Cx_k + Du_k. \quad (3.12)$$

21 If the matrices  $A, B, C, D$  change with time, we have a time-varying system.



### 3.3.3 Nonlinear systems

Nonlinear systems are defined entirely analogously as linear systems. Imagine if we had a non-linear spring in the spring-mass system whereby the dynamics of the block was given by

$$m\ddot{z} + c\dot{z} + (k_1z + k_2z^2) = u.$$

The state of the system is still  $x = [z_1, z_2]^\top$ . But we cannot write this second-order differential equation as two first-order linear differential equations. We are forced to write

$$\dot{x} = \begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -k_1/m - k_2z_1/m & -c/m \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u.$$

Such systems are called nonlinear systems. We will write them succinctly as

$$\begin{aligned} \dot{x} &= f(x, u) \\ y &= g(x, u). \end{aligned} \tag{3.13}$$

The function  $f : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{X}$  that maps the state-space and the input space to the state-space is called the dynamics of the system. Analogously, for discrete-time nonlinear systems we will have

$$\begin{aligned} x_{k+1} &= f(x_k, u_k) \\ y_k &= g(x_k, u_k). \end{aligned}$$

🔗 Is the nonlinear spring-mass system time-invariant?

## 3.4 Markov Decision Processes (MDPs)

Let us now introduce a concept called MDPs which is very close to Markov chains that we saw in the previous chapter. In fact, you are already implementing an MDP in your HW 1 problem on the Bayes filter.

MDPs are a model for the scenario when we do not completely know the dynamics  $f(x_k, u_k)$ .

This may happen for a number of reasons and it is important to appreciate them in order to understand the widespread usage of MDPs.

1. We did not do a good job of identifying the function  $f : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{X}$ . This may happen when you are driving a car on an icy road, if you undertake the same control as you do on a clean road, you might reach a different future state  $x_{k+1}$ .
2. We did not use the correct state-space  $\mathcal{X}$ . You could write down the state of the car as given by  $(x, y, \theta, \dot{x}, \dot{y}, \dot{\theta})$  where  $x, y$  are the Euclidean co-ordinates of the car and  $\theta$  is its orientation. This is not a good model

1 for studying high-speed turns, which are affected by other quantities like  
2 wheel slip, the quality of the suspension etc.

3 We may not even know the full state sometimes. This occurs when  
4 you are modeling how users interact with an online website like Ama-  
5 zon.com, you'd like to model the change in state of the user from  
6 "perusing stuff" to "looking stuff to buy it" to "buying it" but there are  
7 certainly many other variables that affect the user's behavior. As another  
8 example, consider the path that an airplane takes to go from Philadelphia  
9 to Los Angeles. This path is affected by the weather at all places along  
10 the route, it'd be cumbersome to incorporate the weather to find the  
11 shortest-time path for the airplane.

- 12 3. We did not use the correct control-space  $U$  for the controller. This is  
13 akin to the second point above. The gas pedal which one may think  
14 of as the control input to a car is only one out of the large number of  
15 variables that affect the running of the car's engine.

**MDPs are a drastic abstraction of all the above situations.** We write

$$x_{k+1} = f(x_k, u_k) + \epsilon_k \quad (3.14)$$

where the "noise"  $\epsilon_k$  is not under our control. The quantity  $\epsilon_k$  is not arbitrary however, we will assume that

1. noise  $\epsilon_k$  is a random variable and we know its distribution. For example, you ran your car lots of times on icy road and measured how the state  $x_{k+1}$  deviates from similar runs on a clean road. The difference between the two is modeled as  $\epsilon_k$ . Note that the distribution of  $\epsilon_k$  may be a function of time  $k$ .
2. noise at different timesteps  $\epsilon_1, \epsilon_2, \dots$ , is independent.

Instead of a deterministic transition for our system from  $x_k$  to  $x_{k+1}$ , we now have

$$x_{k+1} \sim P(x_{k+1} | x_k, u_k).$$

which is just another way of writing (3.14).

The latter is a probability table of size  $|\mathcal{X}| \times |\mathcal{U}| \times |\mathcal{X}|$  akin to the transition matrix of a Markov chain except that there is a different transition matrix for every control  $u \in \mathcal{U}$ . The former version (3.14) is more amenable to analysis. MDPs can be alternatively called stochastic dynamical systems, we will use either names for them in this course. For completeness, let us note down that linear stochastic systems will be written as

$$x_{k+1} = Ax_k + Bu_k + \epsilon_k.$$

🔗 You should think about the state-space, control-space and the noise in the MDP for the Bayes filter problem in HW 1. Where do we find MDPs in real-life? There are lots of expensive robots in GRASP, e.g., a Kuka manipulator such as this <https://www.youtube.com/watch?v=ym64NFC> costs upwards of \$100,000. Would you model it as a stochastic dynamical system?

The moral of this section is to remember that as pervasive as noise

seems in all problem formulations in this course, it models different situations depending upon the specific problem. Understanding where noise comes from is important for real-world applications.

1 **Noise in continuous-time systems** You will notice that we only talked about  
 2 discrete-time systems with noise in (3.14). We can also certainly talk about  
 3 continuous-time systems whose dynamics  $f$  we do not know precisely

$$\dot{x}(t) = f(x(t), u(t)) + \epsilon(t) \quad (3.15)$$

4 and model the gap in our knowledge as noise  $\epsilon(t)$ . While this may seem quite  
 5 natural, it is mathematically very problematic. The hurdle stems from the fact  
 6 that if we want  $\epsilon(t)$  to be a random variable *at each time instant*, then the  
 7 signal  $\epsilon(t)$  may not actually exist, e.g., it may not even be continuous. Signals  
 8 like  $\epsilon(t)$  exist only in very special cases, one of them is called “Brownian  
 9 motion” where the increment of the signal after infinitesimal time  $\Delta t$  is a  
 10 Gaussian random variable

$$\epsilon(t + \Delta t) - \epsilon(t) = N(0, \Delta t).$$

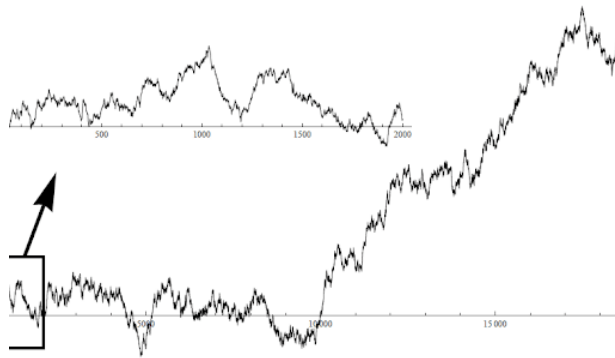


Figure 3.2: A typical Brownian motion signal  $\epsilon(t)$ . You can also see an animation at [https://en.wikipedia.org/wiki/File:Brownian\\_Motion.ogv](https://en.wikipedia.org/wiki/File:Brownian_Motion.ogv)

11 We will not worry about this technicality in this course. We will talk about  
 12 continuous-time systems with noise but with the implicit understanding that  
 13 there is some underlying real-world discrete-time system and the continuous-  
 14 time system is only an abstraction of it.

### 15 3.4.1 Back to Hidden Markov Models

16 Since our sensors measure the state  $x$  of the world, it will be useful to think of  
 17 the output  $y$  of a dynamical system as the observations from Chapter 2. This  
 18 idea neatly ties back our development of dynamical systems to observations.

❓ Do continuous-time systems, stochastic or non-stochastic, exist in the real world? Consider the Kuka manipulator again, do you think the dynamics of this robot is a continuous-time system? Would you model it so?



1 Just like we considered an HMM with observation probability

$$P(Y_k = y | X_k = x)$$

2 we will consider dynamical systems for which we do not precisely know how  
3 the output computation. We will model the gap in our knowledge of the exact  
4 observation mechanism as the output being a noisy function of the state. This  
5 is denoted as

$$y_k = g(x_k) + \nu_k. \quad (3.16)$$

6 The noise  $\nu_k$  is similar to the noise in the dynamics  $\epsilon_k$  in (3.14). Analogously,  
7 we can also have noise in the observations of a linear system

$$y_k = Cx_k + Du_k + \nu_k.$$

Hidden Markov Models with underlying MDPs/Markov chains and stochastic dynamical systems with noisy observations are two different ways to think of the same concept, namely getting observations across time about the true state of a dynamic world.

In the former we have

$$\text{(state transition matrix)} \quad P(X_{k+1} = x' | X_k = x, u_k = u)$$

$$\text{(observation matrix)} \quad P(Y_k = y | X_k = x),$$

while in the latter we have

$$\text{(nonlinear dynamics)} \quad x_{k+1} = f(x_k, u_k) + \epsilon_k$$

$$\text{(nonlinear observation model)} \quad y_k = g(x_k) + \nu_k,$$

or

$$\text{(linear dynamics)} \quad x_{k+1} = Ax_k + Bu_k + \epsilon_k$$

$$\text{(linear observation model)} \quad y_k = Cx_k + Du_k + \nu_k.$$

HMMs are easy to use for certain kinds of problems, e.g., speech-to-text, or a robot wandering in a grid world (like the Bayes filter problem in HW 1). Dynamical systems are more useful for certain other kinds of problems, e.g., a Kuka manipulator where you can use Newton's laws to simply write down the functions  $f, g$ .

## 8 3.5 Kalman Filter (KF)

9 We will now introduce the Kalman Filter. It is the analog of the Bayes filter  
10 from the previous chapter. This is by far the most important algorithm in  
11 robotics and it is hard to imagine running any robot without the Kalman filter  
12 or some variant of it.

13 Consider a linear dynamical system with linear observations

$$\begin{aligned} x_{k+1} &= Ax_k + Bu_k + \epsilon_k \\ y_k &= Cx_k + \nu_k. \end{aligned} \quad (3.17)$$

**i** Observation noise and dynamics noise are different in subtle ways. The former may not always be due to our poor modeling. For instance, the process by which a camera acquires its images has some inherent noise. You may have seen a side-by-side comparison of different cameras using their ISOs



An image taken from a camera with low lighting has a lot of “noise”. What causes this noise?

**i** You will agree that creating the state-transition matrix for the Bayes filter problem in HW 1 was really the hardest part of the problem. If the state-space were continuous and not a discrete cell-based world, you could have written the dynamics very easily in one line of code.

1 where the noise vectors

$$\epsilon_k \sim N(0, R)$$

$$\nu_k \sim N(0, Q)$$

2 are both zero-mean and Gaussian with covariances  $R$  and  $Q$  respectively. We  
3 have also assumed that  $D = 0$  because typically the observations do not  
4 depend on the control.

Our goal is to compute the best estimate of the state after multiple observations

$$P(x_k | y_1, \dots, y_k).$$

Notice that this is the same as the filtering problem that we solved for Hidden Markov Models. Just like we used the forward algorithm to compute the filtering estimate recursively, we are going to use our development of the Kalman gain to incorporate a new observation recursively.

**i** We will assume that the distribution of noise  $\epsilon_k, \nu_k$  does not change with time  $k$ . If it does change in your problem, you will see that following equations are quite easy to modify.

### 5 3.5.1 Step 0: Observing that the state estimate at any timestep 6 is a Gaussian

7 Maintaining the entire probability distribution  $P(x_k | y_1, \dots, y_k)$  is difficult  
8 now, as opposed to the HMM with a finite number of states. We will exploit  
9 the following important fact. If we assume that the initial distribution of  $x_0$   
10 was a Gaussian, since all operations in (3.17) are linear, our new estimate of  
11 the state  $\hat{x}_k$  at time  $k$  is also a Gaussian

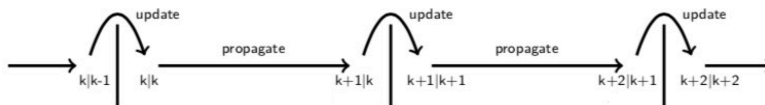
$$\hat{x}_{k|k} \sim P(x_k | y_1, \dots, y_k) \equiv N(\mu_{k|k}, \Sigma_{k|k}).$$

12 The subscript

$$\hat{x}_{k+1|k}$$

13 denotes that the quantity being talked about, i.e.,  $\hat{x}_{k+1|k}$ , or others like  $\mu_{k+1|k}$ ,  
14 is of the  $(k+1)$ <sup>th</sup> timestep and was calculated on the basis of observations  
15 up to (and including) the  $k$ <sup>th</sup> timestep. We will therefore devise recursive  
16 updates to obtain  $\mu_{k+1|k+1}, \Sigma_{k+1|k+1}$  using their old values  $\mu_{k|k}, \Sigma_{k|k}$ . We  
17 will imagine that our initial estimate for the state  $\hat{x}_{0|0}$  has a known distribution

$$\hat{x}_{0|0} \sim N(\mu_{0|0}, \Sigma_{0|0}).$$



18

### 19 3.5.2 Step 1: Propagating the dynamics by one timestep

20 Suppose we had an estimate  $\hat{x}_{k|k}$  after  $k$  observations/time-steps. Since the  
21 dynamics is linear, we can use the prediction problem to compute the estimate

1 of the state at time  $k + 1$  before the next observation arrives

$$P(x_{k+1} \mid y_1, \dots, y_k).$$

2 From the first equation of (3.17), this is given by

$$\hat{x}_{k+1|k} = A\hat{x}_{k|k} + Bu_k + \epsilon_{k+1}$$

3 Notice that the subscript on the left-hand side is  $k + 1|k$  because we did  
4 not take into account the observation at timestep  $k + 1$  yet. The mean and  
5 covariance of this estimate are given by

$$\begin{aligned} \mu_{k+1|k} &= \mathbb{E}[\hat{x}_{k+1|k}] = \mathbb{E}[A\hat{x}_{k|k} + Bu_k + \epsilon_{k+1}] \\ &= A\mu_{k|k} + Bu_k. \end{aligned} \quad (3.18)$$

6 We can also calculate the covariance of the estimate  $\hat{x}_{k+1|k}$  to see that

$$\begin{aligned} \Sigma_{k+1|k} &= \text{cov}(\hat{x}_{k+1|k}) = \text{cov}(A\hat{x}_{k|k} + Bu_k + \epsilon_{k+1}) \\ &= A\Sigma_{k|k}A^\top + R, \end{aligned} \quad (3.19)$$

7 using our calculation in (3.1).

### 8 3.5.3 Step 2: Incorporating the observation

9 After the dynamics propagation step, our estimate of the state is  $\hat{x}_{k+1|k}$ , this is  
10 the state of the system that we believe is true after  $k$  observations. We should  
11 now incorporate the latest observation  $y_{k+1}$  to update this estimate to get

$$P(x_{k+1} \mid y_1, \dots, y_k, y_{k+1}).$$

12 This is exactly the same problem that we saw in Section 3.2.3. Given the  
13 measurement

$$y_{k+1} = Cx_{k+1} + \nu_{k+1}$$

14 we first compute the Kalman gain  $K_{k+1}$  and the updated mean of the estimate  
15 as

$$\begin{aligned} K_{k+1} &= \Sigma_{k+1|k}C^\top (C\Sigma_{k+1|k}C^\top + Q)^{-1} \\ \mu_{k+1|k+1} &= \mu_{k+1|k} + K_{k+1}(y_{k+1} - C\mu_{k+1|k}). \end{aligned} \quad (3.20)$$

16 The covariance is given by our same calculation again

$$\begin{aligned} \Sigma_{k+1|k+1} &= (I - K_{k+1}C)\Sigma_{k+1|k}, \text{ or} \\ &= (I - K_{k+1}C)\Sigma_{k+1|k}(I - K_{k+1}C)^\top + K_{k+1}QK_{k+1}^\top, \text{ or} \\ &= \left(\Sigma_{k+1|k}^{-1} + C^\top Q^{-1}C\right)^{-1}. \end{aligned} \quad (3.21)$$

17 The second expression is known as Joseph's form and is numerically more  
18 stable than the other expressions.

**i** Observe that even if we knew the state dynamics precisely, i.e., if  $R = 0$ , we still have a non-trivial propagation equation for  $\Sigma_{k+1|k}$ .

The new estimate of the state is

$$\hat{x}_{k+1|k+1} \sim \mathbf{P}(x_{k+1} | y_1, \dots, y_{k+1}) \equiv N(\mu_{k+1|k+1}, \Sigma_{k+1|k+1}).$$

and we can again proceed to Step 1 for the next timestep.

### 3.5.4 Discussion

There are several important observations to make and remember about the Kalman Filter (KF).

- **Recursive updates to compute the best estimate given all past observations.** The KF is a recursive filter (just like the forward algorithm for HMMs) and incorporates observations one by one. The estimate that it maintains, namely  $\hat{x}_{k+1|k+1}$ , depends upon all past observations

$$\hat{x}_{k+1|k+1} \sim \mathbf{P}(x_{k+1} | y_1, \dots, y_{k+1}).$$

We have simply *computed* the estimate recursively.

- **Optimality of the KF for linear systems with Gaussian noise.** The KF is optimal in the following sense. Imagine if we had access to all the observations  $y_1, \dots, y_k$  beforehand and computed some other estimate

$$\hat{x}_{k|k}^{\text{fancy filter}} = \text{some function}(\hat{x}_{0|0}, y_1, \dots, y_k).$$

We use some other fancy method to design this estimator, e.g., nonlinear combination of the observations or incorporating observations across multiple timesteps together etc. to obtain something that has the smallest error with respect to the true state  $x_k$

$$\text{tr} \left( \mathbb{E}_{\epsilon_1, \dots, \epsilon_k, \nu_1, \dots, \nu_k} \left[ (\hat{x}_{k|k}^{\text{fancy filter}} - x_k)(\hat{x}_{k|k}^{\text{fancy filter}} - x_k)^\top \right] \right). \quad (3.22)$$

Then this estimate would be exactly the same as that of the KF

$$\hat{x}_{k|k}^{\text{fancy filter}} = \hat{x}_{k|k}^{\text{KF}}.$$

This is a deep fact. First, the KF estimate was created recursively and yet we can do no better than it with our fancy estimator. This is analogous to the fact that the forward algorithm computes the correct filtering estimate even if it incorporates observations one by one recursively. Second, the KF combines the new observation and the old estimate linearly in (3.20). You could imagine that there is some other way to incorporate new observations, but it turns out that for linear dynamical systems with Gaussian noise, the KF is the best solution, we can do no better.

- **The KF is the best linear filter.** If we had a nonlinear dynamical system or a non-Gaussian noise with a linear dynamics/observations, there are other filters that can give a smaller error (3.22) than the KF.

In the next section, we will take a look at one such example. However, even in these cases, the KF is the *best linear filter*.

- **Assumptions that are implicit in the KF.** We assumed that both the dynamics noise  $\epsilon_{k+1}$  and the observation noise  $\nu_{k+1}$  are uncorrelated with the estimate  $\hat{x}_{k+1|k}$  computed prior to them (where did we use these assumptions?). This implicitly assumes that dynamics and observation noise are “white”, i.e., uncorrelated in time

$$\begin{aligned} E[\epsilon_k \epsilon_{k'}^\top] &= 0 \quad \text{for all } k, k' \\ E[\nu_k \nu_{k'}^\top] &= 0 \quad \text{for all } k, k'. \end{aligned}$$

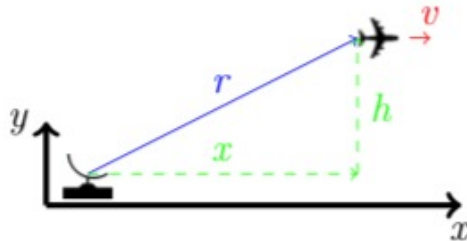
The Wikipedia webpage at [https://en.wikipedia.org/wiki/Kalman\\_filter#Example\\_application,\\_technical](https://en.wikipedia.org/wiki/Kalman_filter#Example_application,_technical) gives a simple example of a Kalman Filter.

🔗 How should one modify the KF equations if we have multiple sensors in a robot, each coming in at different frequencies?

### 3.6 Extended-Kalman Filter (EKF)

The KF heavily exploits the fact that our dynamics/measurements are linear. In robotics problems, either of them or typically both, will be nonlinear. The Extended-Kalman Filter (EKF) is a modification of the KF to handle such situations.

**Example of a nonlinear dynamical system** The state of most real problems evolves as a nonlinear function of their current state and control. This is the same for sensors such as cameras measure a nonlinear function of the state. We will first see how to linearize a given nonlinear system shown below.



We have a radar sensor that measures the distance of the plane  $r$  from the radar trans-receiver up to noise  $\nu$ . We would like to measure its distance  $x$  and height  $h$ . If the plane travels with a constant velocity, we have

$$\dot{x} = v, \text{ and } \dot{h} = 0,$$

and

$$r = \sqrt{x^2 + h^2} + \nu.$$

Since we do not know how the plane might change its altitude, we will assume that it maintains a constant altitude

$$\dot{h} = 0.$$

1 Notice that the above systems are our model for how the state of the airplane  
 2 evolves and could of course be wrong. As we discussed, we will model this as  
 3 noise.

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \epsilon;$$

$$r = \sqrt{x_1^2 + x_3^2} + \nu;$$

4 here  $x_1 \equiv x$ ,  $x_2 \equiv v$  and  $x_3 = h$ , and  $\epsilon \in \mathbb{R}^3$ ,  $\nu \in \mathbb{R}$  are zero-mean Gaussian  
 5 noise. The dynamics in this case is linear but the observations are a nonlinear  
 6 function of the state.

7 One way to use the Kalman Filter for this problem is to linearize the  
 8 observation equation around some state, say  $x_1 = x_2 = x_3 = 0$  using the  
 9 Taylor series

$$\begin{aligned} r_{\text{linearized}} &= r(0, 0, 0) + \left. \frac{\partial r}{\partial x_1} \right|_{x_1=0, x_3=0} (x_1 - 0) + \left. \frac{\partial r}{\partial x_3} \right|_{x_1=0, x_3=0} (x_3 - 0) \\ &= 0 + \left. \frac{2x_1}{2\sqrt{x_1^2 + x_3^2}} \right|_{x_1=0, x_3=0} x_1 + \left. \frac{2x_3}{2\sqrt{x_1^2 + x_3^2}} \right|_{x_1=0, x_3=0} x_3 \\ &= x_1 + x_3. \end{aligned}$$

10 In other words, upto first order in  $x_1, x_3$ , the observations are linear and we  
 11 can therefore run the KF for computing the state estimate after  $k$  observations.

12

### 13 3.6.1 Propagation of statistics through a nonlinear transfor- 14 mation

15 Given a Gaussian random variable  $\mathbb{R}^d \ni x \sim N(\mu_x, \Sigma_x)$ , we saw how to  
 16 compute the mean and covariance after an *affine* transformation  $y = Ax$

$$E[y] = AE[x], \text{ and } \Sigma_y = A\Sigma_x A^\top.$$

17 If we had a nonlinear function of  $x$

$$\mathbb{R}^p \ni y = f(x)$$

18 we can use the Taylor series by linearizing around the mean of  $x$  to approximate  
 19 the first and second moments of  $y$  as follows.

$$\begin{aligned} y = f(x) &\approx f(\mu_x) + \left. \frac{df}{dx} \right|_{x=\mu_x} (x - \mu_x) \\ &= Jx + (f(\mu_x) - J\mu_x). \end{aligned}$$

20 where we have defined the Jacobian matrix

$$\mathbb{R}^{p \times d} \ni J = \left. \frac{df}{dx} \right|_{x=\mu_x}. \quad (3.23)$$

🔗 You can try to perform a similar linearization for a simple model of a car

$$\begin{aligned} \dot{x} &= \cos \theta \\ \dot{y} &= \sin \theta \\ \dot{\theta} &= u. \end{aligned}$$

where  $x, y, \theta$  are the XY-coordinates and the angle of the steering wheel respectively. This model is known as a Dubins car.

1 This gives

$$\begin{aligned} \mathbb{E}[y] &\approx \mathbb{E}[Jx + (f(\mu_x) - J\mu_x)] = f(\mu_x) \\ \Sigma_y &\approx \mathbb{E}[(y - \mathbb{E}[y])(y - \mathbb{E}[y])^\top] \approx J\Sigma_x J^\top. \end{aligned} \quad (3.24)$$

2 Observe how, up to first order, the mean  $\mu_x$  is directly transformed by the  
 3 nonlinear function  $f$  while the covariance  $\Sigma_x$  is transformed as if there were a  
 4 linear operation  $y \approx Jx$ .

### A simple example

$$y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = f \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = \begin{bmatrix} x_1^2 + x_2 x_3 \\ \sin x_2 + \cos x_3 \end{bmatrix}.$$

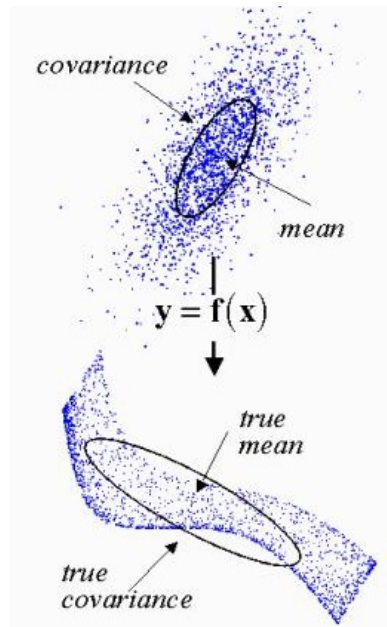
5 We have

$$\frac{df}{dx} = \nabla f(x) = \begin{bmatrix} 2x_1 & x_3 & x_2 \\ 0 \cos x_2 & -\sin x_3 & \end{bmatrix}.$$

6 The Jacobian at  $\mu_x = [\mu_{x_1}, \mu_{x_2}, \mu_{x_3}]^\top$  is

$$J = \nabla f(x) \Big|_{x=\mu_x} = \begin{bmatrix} 2\mu_{x_1} & \mu_{x_3} & \mu_{x_2} \\ 0 \cos \mu_{x_2} & -\sin \mu_{x_3} & \end{bmatrix}.$$

It is very important to remember that we are approximating the distribution of  $P(f(x))$  as a Gaussian. Even if  $x$  is a Gaussian random variable, the distribution of  $y = f(x)$  need not be Gaussian. Indeed  $y$  is only Gaussian if  $f$  is an affine function of  $x$ .



### 3.6.2 Extended Kalman Filter

The above approach of linearizing the observations of the plane around the origin may lead to a lot of errors. This is because the point about which we linearize the system is fixed. We can do better by linearizing the system at each timestep. Let us say that we are given a nonlinear system

$$\begin{aligned}x_{k+1} &= f(x_k, u_k) + \epsilon \\ y_k &= g(x_k) + \nu.\end{aligned}$$

The central idea of the Extended Kalman Filter (EKF) is to linearize a nonlinear system at each timestep  $k$  around the latest state estimate given by the Kalman Filter and use the resultant linearized dynamical system in the KF equations for the next timestep.

🔍 Can you say where will our linearized observation equation incur most error?

#### Step 1: Propagating the dynamics by one timestep

We will linearize the dynamics equation around the mean of the previous state estimate  $\mu_{k|k}$

$$\begin{aligned}x_{k+1} &= f(x_k, u_k) + \epsilon \\ &\approx f(\mu_{k|k}, u_k) + \left. \frac{\partial f}{\partial x} \right|_{x=\mu_{k|k}} (x_k - \mu_{k|k}) + \epsilon_k.\end{aligned}$$

Let the Jacobian be

$$A(\mu_{k|k}) = \left. \frac{\partial f}{\partial x} \right|_{x=\mu_{k|k}}. \quad (3.25)$$

The mean and covariance of the EKF after the dynamics propagation step is therefore given by

$$\begin{aligned}\mu_{k+1|k} &= f(\mu_{k|k}, u_k) \\ \Sigma_{k+1|k} &= A\Sigma_{k|k}A^\top + R.\end{aligned} \quad (3.26)$$

It is worthwhile to notice the similarities of the above set of equations with (3.18) and (3.19). The mean  $\mu_{k|k}$  is propagated using a nonlinear function  $f$  to get  $\mu_{k+1|k}$ , the covariance is propagated using the Jacobian  $A(\mu_{k|k})$  which is recomputed using (3.25) at each timestep.

#### Step 2: Incorporating the observation

We have access to  $\mu_{k+1|k}$  after Step 1, so we can linearize the nonlinear observations at this state.

$$\begin{aligned}y_{k+1} &= g(x_{k+1}) + \nu \\ &\approx g(\mu_{k+1|k}) + \left. \frac{dg}{dx} \right|_{x=\mu_{k+1|k}} (x_{k+1} - \mu_{k+1|k}) + \nu\end{aligned}$$



1 Again define the Jacobian

$$C(\mu_{k+1|k}) = \left. \frac{\partial g}{\partial x} \right|_{x=\mu_{k+1|k}}. \quad (3.27)$$

2 Consider the fake observation which is a transformed version of the actual  
3 observation  $y_{k+1}$  (think of this as a new sensor or a post-processed version of  
4 the original sensor)

$$y'_{k+1} = y_{k+1} - g(\mu_{k+1|k}) + C\mu_{k+1|k} \approx Cx_{k+1}.$$

5 Our fake observation is a nice linear function of the state  $x_{k+1}$  and we can  
6 therefore use the Kalman Filter equations to incorporate this fake observation

$$\begin{aligned} \mu_{k+1|k+1} &= \mu_{k+1|k} + K(y'_{k+1} - C\mu_{k+1|k}) \\ \text{where } K &= \Sigma_{k+1|k} C^\top (C\Sigma_{k+1|k} C^\top + Q)^{-1}. \end{aligned}$$

7 Let us resubstitute our fake observation in terms of the actual observation  
8  $y_{k+1}$ .

$$y'_{k+1} - C\mu_{k+1|k} = y_{k+1} - g(\mu_{k+1|k}),$$

9 to get the EKF equations for incorporating one observation

$$\begin{aligned} \mu_{k+1|k+1} &= \mu_{k+1|k} + K(y_{k+1} - g(\mu_{k+1|k})) \\ \Sigma_{k+1|k+1} &= (I - KC)\Sigma_{k+1|k}. \end{aligned} \quad (3.28)$$

**Summary of the Extended Kalman Filter** The EKF can use a nonlinear model of the system but linearizes the dynamics and observation equations at each timestep before plugging them into the Kalman Filter formulae.

1. Just like the KF, say we have the current state estimate  $\mu_{k|k}$  and  $\Sigma_{k|k}$ .
2. The system applies a control input  $u_k$  and we update our state estimate to be

$$\begin{aligned} \mu_{k+1|k} &= f(\mu_{k|k}, u_k) \\ \Sigma_{k+1|k} &= A\Sigma_{k|k}A^\top + R. \end{aligned}$$

where  $A$  depends on  $\mu_{k|k}$ .

3. We next incorporate an observation by linearizing the observation equations around  $\mu_{k+1|k}$

$$\begin{aligned} K &= \Sigma_{k+1|k} C^\top (C\Sigma_{k+1|k} C^\top + Q)^{-1} \\ \mu_{k+1|k+1} &= \mu_{k+1|k} + K(y_{k+1} - g(\mu_{k+1|k})) \\ \Sigma_{k+1|k+1} &= (I - KC)\Sigma_{k+1|k} \end{aligned}$$

where again  $C$  depends on  $\mu_{k+1|k}$

## 1 Discussion

- 2 1. The EKF dramatically expands the applicability of the Kalman Filter. It  
3 can be used for most real systems, even with very complex models  $f, h$ .  
4 It is very commonly used in robotics and can handle nonlinear observa-  
5 tions from complex sensors such as a LiDAR and camera easily. For in-  
6 stance, sophisticated augmented/virtual reality systems like Google AR-  
7 Core/Snapchat/iPhone etc. ([https://www.youtube.com/watch?v=cape\\_Af9j7w](https://www.youtube.com/watch?v=cape_Af9j7w))  
8 run EKF to track the motion of the phone or of the objects in the image.
- 9 2. The KF was special because it is the optimal linear filter, i.e., KF  
10 estimates have the smallest mean squared error with respect to the true  
11 state for linear dynamical systems with Gaussian. The EKF is a clever  
12 application of KF to nonlinear systems but it no longer has this property.  
13 There do exist filters for nonlinear systems that will have a smaller  
14 mean-squared error than the EKF. We will look at some of them in the  
15 next section.
- 16 3. Linearization is the critical step in the implementation of the EKF and  
17 EKF state estimate can be easily quite bad, especially if the system goes  
18 into states where the linearized matrix  $A$  and the nonlinear dynamics  
19  $f(x_k, u_k)$  differ significantly. A common trick for handling this is to  
20 perform multiple steps of dynamics propagation using a continuous-time  
21 model of the system between successive observations. Say we have a  
22 system

$$\dot{x} = f(x(t), u(t)) + \epsilon(t)$$

23 where  $\epsilon(t + \delta t) - \epsilon(t)$  is a Gaussian random variable  $N(0, R\delta t)$  as  
24  $\delta \rightarrow 0$ ; see the section on Brownian motion for how to interpret noise  
25 in continuous-time systems. We can construct a discrete-time system  
26 from this as

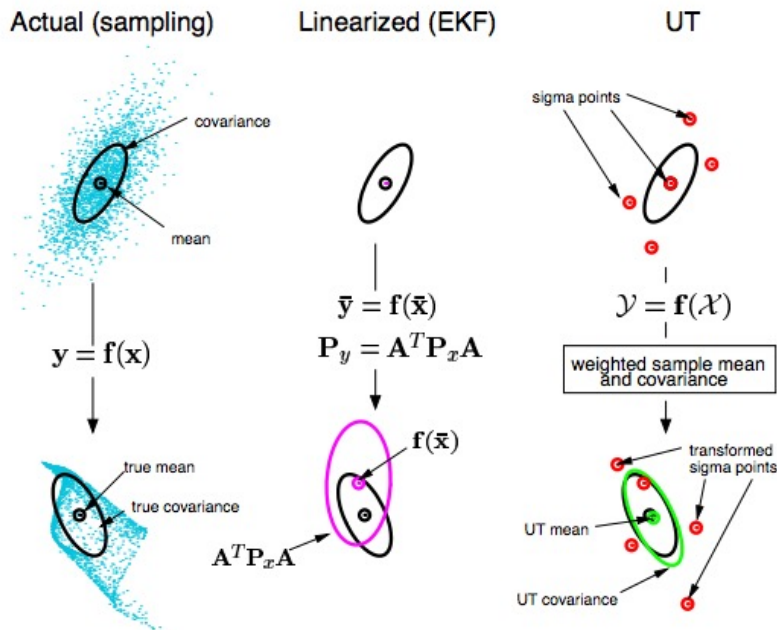
$$\begin{aligned} x_{t+\Delta t} &= x(t) + f(x(t), u(t)) \Delta t + \epsilon \\ &\equiv f^{\text{discrete-time}}(x(t), u(t)) + \epsilon. \end{aligned}$$

27 where  $\epsilon \sim N(0, R\Delta t)$  is noise. This is now a discrete-time dynamics  
28 and we can perform Step 1 of the EKF multiple times to obtain a more  
29 accurate estimate of  $\mu_{k+1|k}$  and  $\Sigma_{k+1|k}$ .

## 30 3.7 Unscented Kalman Filter (UKF)

31 Linearization of the dynamics in the EKF is a neat trick to use the KF equations.  
32 But as we said, this can cause severe issues in problems where the dynamics  
33 is very nonlinear. In this section, we will take a look at a powerful method to  
34 handle nonlinear dynamics that is better than linearization.

35 Let us focus on Step 1 which propagates the dynamics in the EKF.



1

We know that even if  $x$  is Gaussian (faint blue points in top left picture), the transformed variable  $y = f(x)$  need not be Gaussian (faint blue points in bottom left). The EKF is really approximating the probability distribution  $P(x_{k+1} | y_1, \dots, y_k)$  as a Gaussian; this distribution could be very different from a Gaussian. This is really the crux of the issue in filtering for nonlinear systems. This approximation, which happens because we are linearizing about the mean  $\mu_{k|k}$ .

2 Let us instead do the following:

- 3 1. Sample a few points from the Gaussian  $N(\mu_{k|k}, \Sigma_{k|k})$  (red points in
- 4     top right).
- 5 2. Transform *each of the points* using the nonlinear dynamics  $f$  (red points
- 6     in bottom right).
- 7 3. Compute their mean and covariance to get  $\mu_{k+1|k}$  and  $\Sigma_{k+1|k}$ . Notice
- 8     how the green ellipse is slightly different than the black ellipse (which
- 9     is the true mean and covariance). Both of these would be different from
- 10    the mean and covariance obtained by linearization of  $f$  (middle column)
- 11    but the green one is more accurate.

In general, we would need a large number of sample points (red) to

accurately get the mean and covariance of  $y = f(x)$ . The Unscented Transform (UT) uses a special set of points known as “sigma points” (these are the ones actually shown in red above) and transforms those points. Sigma points have the special property that the empirical mean of the transformed distribution (UT mean in the above picture) is close to the true mean up to third order; linearization is only accurate up to first order. The covariance (UT covariance) and true covariance also match up to third order.

### 1 3.7.1 Unscented Transform

2 Given a random variable  $x \sim N(\mu_x, \Sigma_x)$ , the Unscented Transform (UT) uses  
 3 sigma points to compute an approximation of the probability distribution of  
 4 the random variable  $y = f(x)$ .

5 **Preliminaries: matrix square root.** Given a symmetric matrix  $\Sigma \in \mathbb{R}^{n \times n}$ ,  
 6 the matrix square root of  $\Sigma$  is a matrix  $S \in \mathbb{R}^{n \times n}$  such that

$$\Sigma = SS.$$

7 We can compute this via diagonalization as follows.

$$\begin{aligned} \Sigma &= VD V^{-1} \\ &= V \begin{bmatrix} d_{11} & \cdots & 0 \\ 0 & & \\ \cdots & 0 & \\ 0 & \cdots & d_{nn} \end{bmatrix} V^{-1} \\ &= V \begin{bmatrix} \sqrt{d_{11}} & \cdots & 0 \\ 0 & & \\ \cdots & 0 & \\ 0 & \cdots & \sqrt{d_{nn}} \end{bmatrix}^2 V^{-1}. \end{aligned}$$

8 We can therefore define

$$S = V \begin{bmatrix} \sqrt{d_{11}} & \cdots & 0 \\ 0 & & \\ \cdots & 0 & \\ 0 & \cdots & \sqrt{d_{nn}} \end{bmatrix} V^{-1}.$$

9 Notice that

$$SS = (VD^{1/2}V^{-1})(VD^{1/2}V^{-1}) = VD V^{-1} = \Sigma.$$

10 We can also define the matrix square root using the Cholesky decomposition  
 11  $\Sigma = LL^T$  which is numerically more stable than computing the square root  
 12 using the above expression. Recall that matrices  $L$  and  $\Sigma$  have the same  
 13 eigenvectors. Typical applications of the Unscented Transform will use this  
 14 method.

1 Given a random variable  $\mathbb{R}^n \ni x \sim N(\mu, \Sigma)$ , we will use the matrix  
2 square root to compute the sigma points as

$$\begin{aligned} x^{(i)} &= \mu + \sqrt{n\Sigma_i}^\top \\ x^{(n+i)} &= \mu - \sqrt{n\Sigma_i}^\top \end{aligned} \quad (3.29)$$

for  $i = 1, \dots, n$ ,

3 where  $\sqrt{n\Sigma_i}$  is the  $i^{\text{th}}$  row of the matrix  $\sqrt{n\Sigma}$ . There are  $2n$  sigma points

$$\{x^{(1)}, \dots, x^{(2n)}\}$$

4 for an  $n$ -dimensional Gaussian. Each sigma point is assigned a weight

$$w^{(i)} = \frac{1}{2n}. \quad (3.30)$$

5 We then transform each sigma point to get the transformed sigma points

$$y^{(i)} = f(x^{(i)}).$$

6 The mean and covariance of the transformed random variable  $y$  can now be  
7 computed as

$$\begin{aligned} \mu_y &= \sum_{i=1}^{2n} w^{(i)} y^{(i)} \\ \Sigma_y &= \sum_{i=1}^{2n} w^{(i)} (y^{(i)} - \mu_y) (y^{(i)} - \mu_y)^\top. \end{aligned} \quad (3.31)$$

8 **Example** Say we have  $x = \begin{bmatrix} r \\ \theta \end{bmatrix}$  with  $\mu_x = [1, \pi/2]$  and  $\Sigma_x = \begin{bmatrix} \sigma_r^2 & 0 \\ 0 & \sigma_\theta^2 \end{bmatrix}$ .

9 We would like to compute the probability distribution of  $y = f(x) = \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix}$   
10 which is a polar transformation. Since  $x$  is two-dimensional, we will have 4  
11 sigma points with equal weights  $w^{(i)} = 0.25$ . The square root in the sigma  
12 point expression is

$$\sqrt{n\Sigma} = \begin{bmatrix} \sqrt{2}\sigma_r & 0 \\ 0 & \sqrt{2}\sigma_\theta \end{bmatrix}$$

13 and the sigma points are

$$\begin{aligned} x^{(1)} &= \begin{bmatrix} 1 \\ \pi/2 \end{bmatrix} + \begin{bmatrix} \sqrt{2}\sigma_r \\ 0 \end{bmatrix}, & x^{(3)} &= \begin{bmatrix} 1 \\ \pi/2 \end{bmatrix} - \begin{bmatrix} \sqrt{2}\sigma_r \\ 0 \end{bmatrix} \\ x^{(2)} &= \begin{bmatrix} 1 \\ \pi/2 \end{bmatrix} + \begin{bmatrix} 0 \\ \sqrt{2}\sigma_\theta \end{bmatrix}, & x^{(4)} &= \begin{bmatrix} 1 \\ \pi/2 \end{bmatrix} - \begin{bmatrix} 0 \\ \sqrt{2}\sigma_\theta \end{bmatrix}. \end{aligned}$$

🔗 Compute the mean and covariance of  $y$  by linearizing the function  $f(x)$ .

1 The transformed sigma points are

$$\begin{aligned}
 y^{(1)} &= \begin{bmatrix} r^{(1)} \cos \theta^{(1)} \\ r^{(1)} \sin \theta^{(1)} \end{bmatrix} = \begin{bmatrix} 0 \\ 1 + \sqrt{2}\sigma_r \end{bmatrix} \\
 y^{(2)} &= \begin{bmatrix} r^{(2)} \cos \theta^{(2)} \\ r^{(2)} \sin \theta^{(2)} \end{bmatrix} = \begin{bmatrix} \cos(\pi/2 + \sqrt{2}\sigma_\theta) \\ \sin(\pi/2 + \sqrt{2}\sigma_\theta) \end{bmatrix} \\
 y^{(3)} &= \begin{bmatrix} 0 \\ 1 - \sqrt{2}\sigma_r \end{bmatrix} \\
 y^{(4)} &= \begin{bmatrix} \cos(\pi/2 - \sqrt{2}\sigma_\theta) \\ \sin(\pi/2 - \sqrt{2}\sigma_\theta) \end{bmatrix}.
 \end{aligned}$$

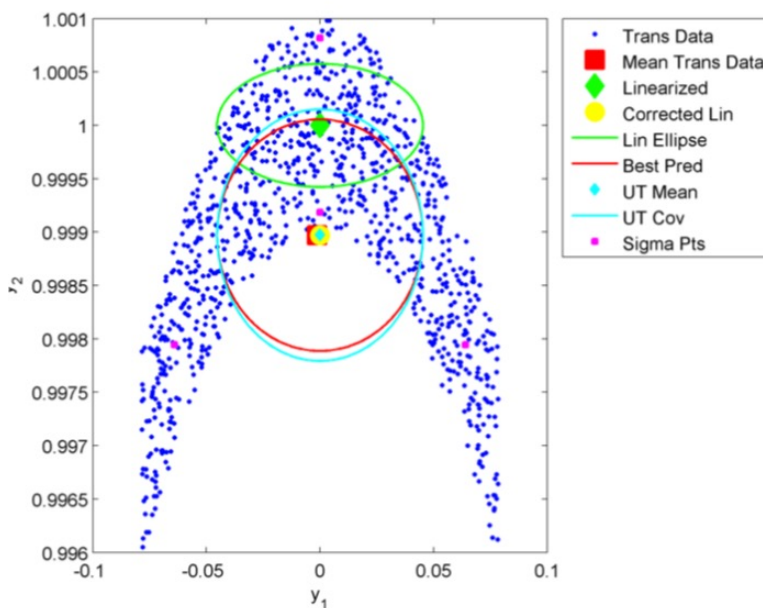


Figure 3.3: Note that the true mean is being predicted very well by the UT and is clearly a better estimate than the linearized mean.

## 2 3.7.2 The UT with tuning parameters

3 The UT is a basic template for a large suite of techniques that capture the  
 4 covariance  $\Sigma_x$  as a set of points and transform those points through the nonlin-  
 5 earity. You will see many alternative implementations of the UT that allow for  
 6 user-tunable parameters. For instance, sometimes the UT is implemented with  
 7 an additional sigma point  $x^{(0)} = \mu$  with weight  $w^{(0)} = \frac{\lambda}{n+\lambda}$  and the weights  
 8 of the other points are adjusted to be  $w^{(i)} = \frac{1}{2(n+\lambda)}$  for a user-chosen parame-  
 9 ter  $\lambda$ . You may also see people using one set of weights  $w^{(i)}$  for computing  
 10 the mean  $\mu_y$  and another set of weights for computing the covariance  $\Sigma_y$ .

🔍 Are the transformed sigma points  $y^{(i)}$  the sigma points of  $P(y) = N(\mu_y, \Sigma_y)$ ?

🔍 We are left with a big lingering question. Why do you think this method is called the “unscented transform”?

### 3.7.3 Unscented Kalman Filter (UKF)

The Unscented Transform gives us a way to accurately estimate the mean and covariance of the transformed distribution through a nonlinearity. We can use the UT to modify the EKF to make it a more accurate state estimator. The resultant algorithm is called the Unscented Kalman Filter (UKF).

**Step 1: Propagating the dynamics by one timestep** Given our current state estimate  $\mu_{k|k}$  and  $\Sigma_{k|k}$ , we use the UT to obtain the updated estimates  $\mu_{k+1|k}$  and  $\Sigma_{k+1|k}$ . If  $x^{(i)}$  are the sigma points with corresponding weights  $w^{(i)}$  for the Gaussian  $N(\mu_{k|k}, \Sigma_{k|k})$ , we set

$$\begin{aligned}\mu_{k+1|k} &:= \sum_{i=1}^{2n} w^{(i)} f(x^{(i)}, u_k) \\ \Sigma_{k+1|k} &:= R + \sum_{i=1}^{2n} w^{(i)} \left( f(x^{(i)}) - \mu_{k+1|k} \right) \left( f(x^{(i)}) - \mu_{k+1|k} \right)^\top\end{aligned}\quad (3.32)$$

**Step 2.1: Incorporating one observation** The observation step is also modified using the UT. The key issue in this case is that we need a way to compute the Kalman gain in terms of the sigma points in the UT. We proceed as follows.

Using *new* sigma points  $x^{(i)}$  for the updated state distribution  $N(\mu_{k+1|k}, \Sigma_{k+1|k})$  with equal weights  $w^{(i)} = 1/2n$ , we first compute their mean after the transformation

$$\hat{y} = \sum_{i=1}^{2n} w^{(i)} g(x^{(i)}) \quad (3.33)$$

and covariances

$$\begin{aligned}\Sigma_{yy} &:= Q + \sum_{i=1}^{2n} w^{(i)} \left( g(x^{(i)}) - \hat{y} \right) \left( g(x^{(i)}) - \hat{y} \right)^\top \\ \Sigma_{xy} &:= \sum_{i=1}^{2n} w^{(i)} \left( f(x^{(i)}) - \mu_{k+1|k} \right) \left( g(x^{(i)}) - \hat{y} \right)^\top.\end{aligned}\quad (3.34)$$

**Step 2.2: Computing the Kalman gain** Until now we have written the Kalman gain using the measurement matrix  $C$ . We will now discuss a more abstract formulation that gives the same expression.

Say we have a random variable  $x$  with known  $\mu_x, \Sigma_x$  and get a new observation  $y$ . We saw how to incorporate this new observation to obtain a better estimator for  $x$  in Section 3.2.3. We will go through a similar analysis as before but in a slightly different fashion, one that does not involve the matrix  $C$ . Let

$$p = \begin{bmatrix} x \\ y \end{bmatrix}$$

and  $\mu_p = [\mu_x \quad \mu_y]$  and

$$\Sigma_p = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}.$$

1 Finding the best (minimum variance estimator)  $\hat{x} = \mu_x + K(y - \mu_y)$  amounts  
 2 to minimizing

$$\min_K \mathbb{E} \left[ \text{tr} (\hat{x} - x) (\hat{x} - x)^\top \right].$$

3 This is called the least squares problem, which you have seen before perhaps  
 4 in slightly different notation. You can solve this problem to see that the best  
 5 gain  $K$  is given by

$$K^* = \Sigma_{xy} \Sigma_{yy}^{-1}. \quad (3.35)$$

6 and this gain leads to an error of

$$(\hat{x} - x) (\hat{x} - x)^\top = \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} = \Sigma_{xx} - K^* Q_{yy} K^{*\top}.$$

7 The nice thing about the Kalman gain in (3.35) is that we can compute it now  
 8 using expressions of  $\Sigma_{xy}$  and  $\Sigma_{yy}$  in terms of the sigma points. This goes as  
 9 as follows:

$$\begin{aligned} K^* &= \Sigma_{xy} \Sigma_{yy}^{-1} \\ \mu_{k+1|k+1} &= \mu_{k+1|k} + K (y_{k+1} - \hat{y}) \\ \Sigma_{k+1|k+1} &= \Sigma_{k+1|k} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \\ &= \Sigma_{k+1|k} - K^* \Sigma_{yy} K^{*\top}. \end{aligned} \quad (3.36)$$

#### Summary of UKF

1. The Unscented Transform (UT) is an alternative to linearization. It gives a better approximation of the mean and covariance of the random variable after being transformed using a nonlinear function than taking the Taylor series approximation.
2. The UKF uses the UT and its sigma points for propagation of uncertainty through the dynamics (3.32) and observation nonlinearities (3.36).

### 3.7.4 UKF vs. EKF

11 As compared to the Extended Kalman Filter, the UKF is a better approximation  
 12 for nonlinear systems. Of course, if the system is linear, both EKF and UKF  
 13 are equivalent to a Kalman Filter.

14 In practice, we typically use the UKF with some tuning parameters in  
 15 the Unscented Transform as discussed in Section 3.7.2. In practice, the EKF  
 16 also has tuning parameters where we may wish to perform multiple updates  
 17 of the dynamics equations with a smaller time-discretization before the next  
 18 observation comes in to alleviate the effect of linearizing the dynamics. A  
 19 well-tuned EKF is often only marginally worse than an UKF: the former  
 20 requires us to compute Jacobians at each step which the latter does not, but  
 21 the latter is often a more involved implementation.



**UKF/EKF approximate filtering distribution as a Gaussian** An important point to remember about both the UKF and EKF is that even if they can handle nonlinear systems, they still approximate the filtering distribution

$$P(x_k | y_1, \dots, y_k)$$

as a Gaussian.

## 3.8 Particle Filters (PFs)

We next look at particle filters (PFs) which are a generalization of the UKF and can handle non-Gaussian filtering distributions. Just like the UT forms the building block of the UKF, the building block of a particle filter is the idea of importance sampling.

### 3.8.1 Importance sampling

Consider the following problem, given a probability distribution  $p(x)$ , we want to approximate it as a sum of Dirac-delta distributions at points  $x^{(i)}$ , also called “particles”, each with weight  $w^{(i)}$

$$p(x) \approx \sum_{i=1}^n w^{(i)} \delta_{x^{(i)}}(x).$$

Say all weights are equal  $1/n$ . Depending upon how we pick the samples  $x^{(i)}$ , we can get very different approximations

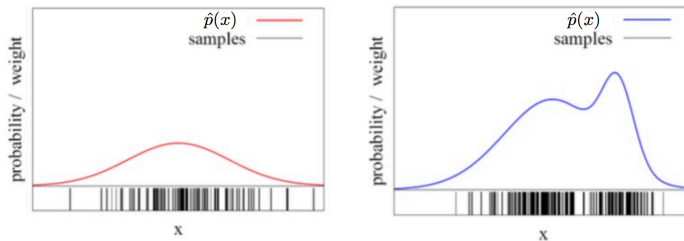


Figure 3.4: Black lines denote particles  $x^{(i)}$ , while red and blue curves denote the approximations obtained using them. If there are a large number of particles in a given region, the approximated probability density of that region is higher.

We see in Figure 3.4 that depending upon the samples, the approxi-

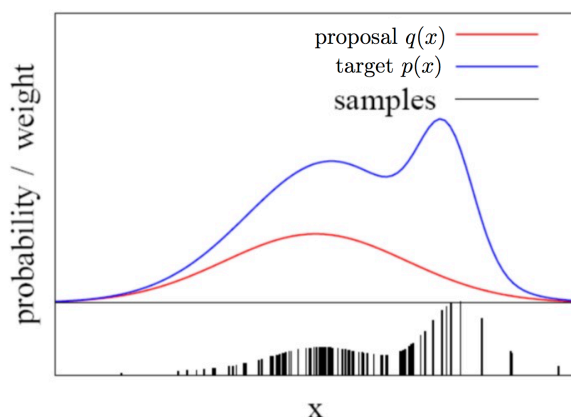
mated probability distributions  $\hat{p}(x)$  can be quite different. Importance sampling is a technique to sample the particles to approximate a given probability distribution  $p(x)$ . The main idea is to use another *known* probability distribution, let us call it  $q(x)$  to *generate particles*  $x^{(i)}$  and account for the differences between the two by assigning weights to each particle

For  $i = 1, \dots, n$ ,

$$x^{(i)} \sim q$$

$$w^{(i)} = \frac{p(x^{(i)})}{q(x^{(i)})}.$$

The original distribution  $p(x)$  is called the “target” and our chosen distribution  $q(x)$  is called the “proposal”. If the number of particles  $n$  is large, we can expect a better approximation of the target density  $p(x)$ .



1

### 2 3.8.2 Resampling particles to make the weights equal

- 3 A particle filter modifies the weights of each particle as it goes through the  
 4 dynamics and observation update steps. This often causes some particles to  
 5 have very low weights and some others to have very high weights.

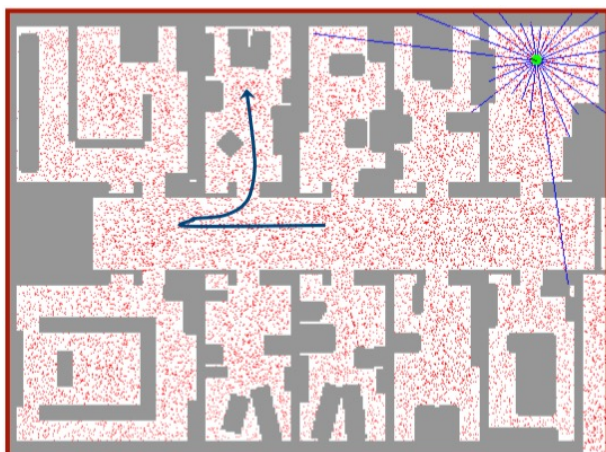


Figure 3.5: An example run of a particle filter. The robot is shown by the green dot in the top right. Observations from a laser sensor (blue rays) attached to the robot measure its distance in a 360-degree field of view around it. Red dots are particles, i.e., possible locations of the robot that we need in order to compute the filtering density  $P(x_k | y_1, \dots, y_k)$ . You should think of this picture as being similar to Problem 1 in Homework 1 where the robot was traveling on a grid. Just like the the filtering density in Problem 1 was essentially zero in some parts of the domain, the particles, say in the bottom left, will have essentially zero weights in a particle filter once we incorporate multiple observations from the robot in top right. Instead of having to carry around these null particles with small weights, the resampling step is used to remove them and sample more particles, say in the top right, where we can benefit from a more accurate approximation of the filtering density.

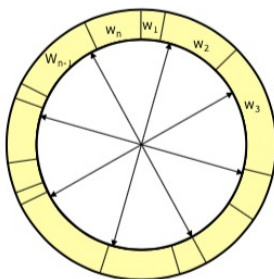
The resampling step takes particles  $\{w^{(i)}, x^{(i)}\}_{i=1}^n$  which approximate a probability density  $p(x)$

$$p(x) = \sum_{i=1}^n w^{(i)} \delta_{x^{(i)}}(x)$$

and returns a new set of particles  $x'^{(i)}$  with equal weights  $w'^{(i)} = 1/n$  that approximate the same probability density

$$p(x) = \frac{1}{n} \sum_{i=1}^n \delta_{x'^{(i)}}(x).$$

The goal of the resampling step is to avoid particle degeneracy, i.e., remove unlikely particles with very low weights and effectively split the particles with very large weights into multiple particles.



1

2 Consider the weights of particles  $\{w^{(i)}\}$  arranged in a roulette wheel as shown  
 3 above. We perform the following procedure: we start at some location, say  
 4  $\theta = 0$ , and move along the wheel in random increments of the angle. After  
 5 each random increment, we add the corresponding particle into our set  $\{x^{(i)}\}$ .  
 6 Since particles with higher weights take up a larger angle in the circle, this  
 7 procedure will often pick those particles and quickly move across particles  
 8 with small weights without picking them too often. We perform this procedure  
 9  $n$  times for  $n$  particles. As an algorithm

- 10 1. Let  $r$  be a uniform random variable in interval  $[0, 1/n]$ . Pick  $c = w^{(1)}$   
 11 and initialize  $i = 1$ .
- 12 2. For each  $m = 1, \dots, n$ , let  $u = r + (m-1)/n$ . Increment  $i \leftarrow i+1$  and  
 13  $c \leftarrow c + w^{(i)}$  while  $u > c$  and set new particle location  $x^{(m)} = x^{(i)}$ .

14 **It is important to notice that** the resampling procedure does not actually  
 15 change the locations of particles. Particles with weights much lower than  $1/n$   
 16 will be eliminated while particles with weights much higher than  $1/n$  will be  
 17 “cloned” into multiple particles each of weight  $1/n$ .

**i** There are many other methods of resampling. We have discussed here, something known as “low variance resampling”, which is easy to remember and code up. Fancier resampling methods also change the locations of the particles. The goal remains the same, namely to eliminate particles with low weights.

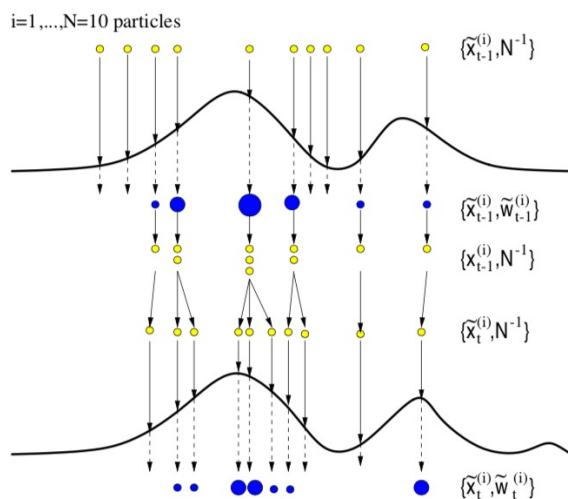


Figure 3.6: A cartoon depicting resampling. Disregard the different notation in this cartoon. Resampling does not change the probability distribution that we wish to approximate; it simply changes the particles and their weights.

### 3.8.3 Particle filtering: the algorithm

The basic template of a PF is similar to that of the UKF and involves two steps, the first where we propagate particles using the dynamics to estimate  $P(x_{k+1} | y_1, \dots, y_k)$  and a second step where we incorporate the observation to compute the updated distribution  $P(x_{k+1} | y_1, \dots, y_{k+1})$ .

Before we look at the theoretical derivation of a particle filter, it will help to go through the algorithm as you would implement on a computer. We assume that we have access to particles  $x_{k|k}^{(i)}$

$$P(x_k | y_1, \dots, y_k) = \frac{1}{n} \sum_{i=1}^n \delta_{x_{k|k}^{(i)}}(x),$$

all with equal weights  $w_{k|k}^{(i)} = 1/n$ .

**1. Step 1: Propagating the dynamics.** Each particle  $i = 1, \dots, n$  is updated by one timestep

$$x_{k+1|k}^{(i)} = f(x_{k|k}^{(i)}, u_k) + \epsilon_k$$

where  $f$  is the system dynamics using Gaussian noise  $\epsilon_k \sim N(0, R)$ . Weights of particles are unchanged  $w_{k+1|k}^{(i)} = w_{k|k}^{(i)} = 1/n$ .

**2. Step 2: Incorporating the observation.** Given a new observation  $y_{k+1}$ , we update the weight of each particle using the likelihood of receiving that observation

$$w_{k+1|k+1}^{(i)} \propto P(y_{k+1} | x_{k+1|k}^{(i)}) w_{k+1|k}^{(i)}.$$

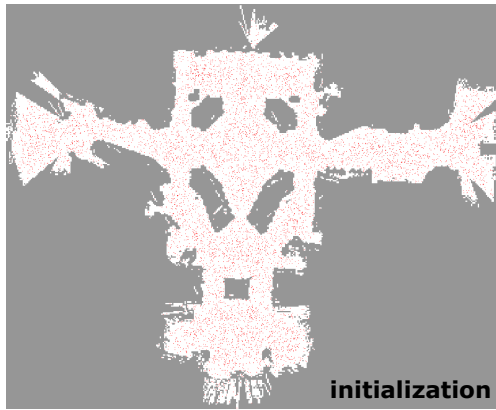
Note that  $P(y_{k+1} | x_{k+1|k}^{(i)})$  is a Gaussian and depends upon the Gaussian observation noise  $\nu_k$ . The mean of this Gaussian is  $g(x_{k+1|k}^{(i)})$  and its variance is equal to  $Q$ , i.e.,

$$\begin{aligned} P(y_{k+1} | x_{k+1|k}^{(i)}) &= P(\nu_{k+1} \equiv y_{k+1} - g(x_{k+1|k}^{(i)})) \\ &= \frac{1}{\sqrt{(2\pi)^p \det(Q)}} \exp\left(-\frac{\nu_{k+1}^\top Q^{-1} \nu_{k+1}}{2}\right). \end{aligned}$$

Normalize the weights  $w_{k+1|k+1}^{(i)}$  to sum up to 1.

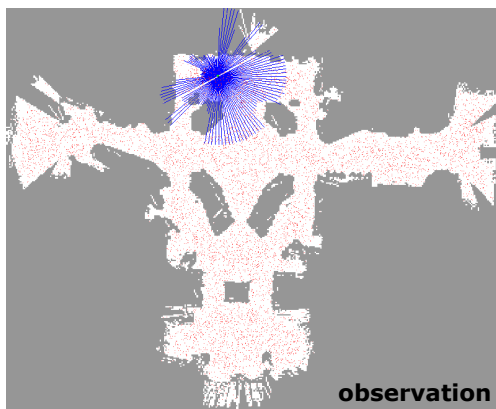
**3. Step 3: Resampling step** Perform the resampling step to obtain new particle locations  $x_{k+1|k+1}^{(i)}$  with uniform weights  $w_{k+1|k+1}^{(i)} = 1/n$ .

### 3.8.4 Example: Localization using particle filter



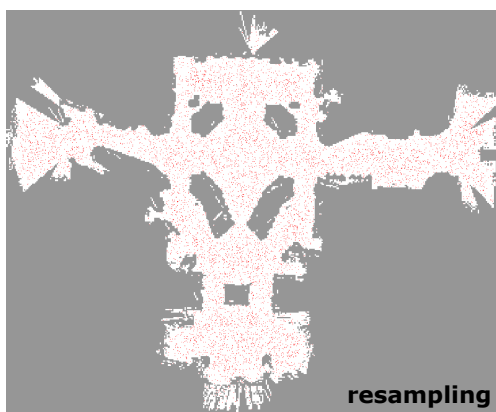
18

2



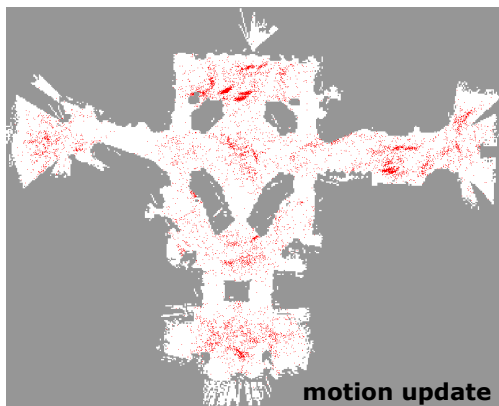
19

3



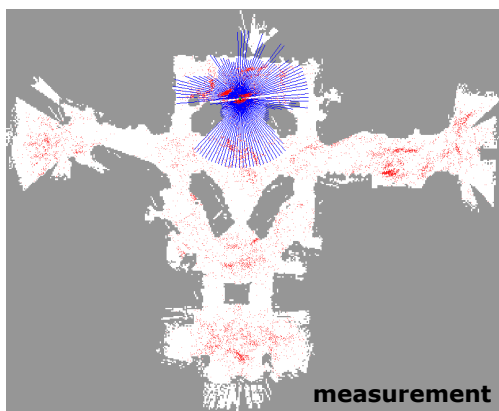
20

4



21

1



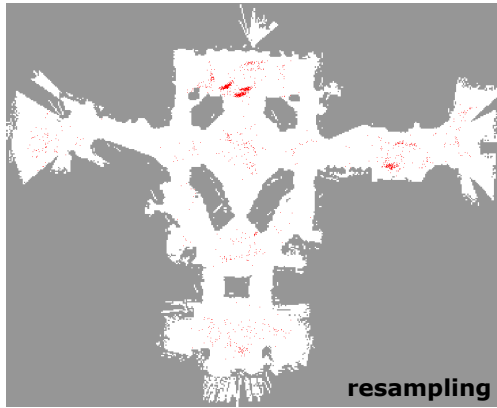
22

2



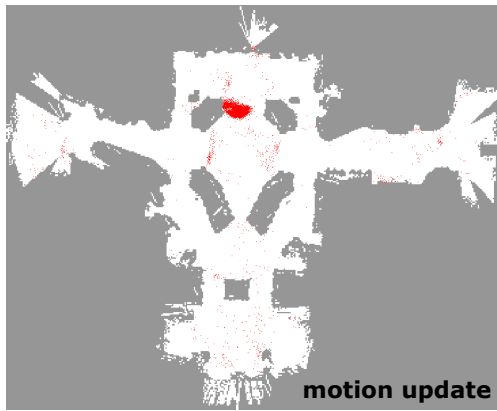
23

3



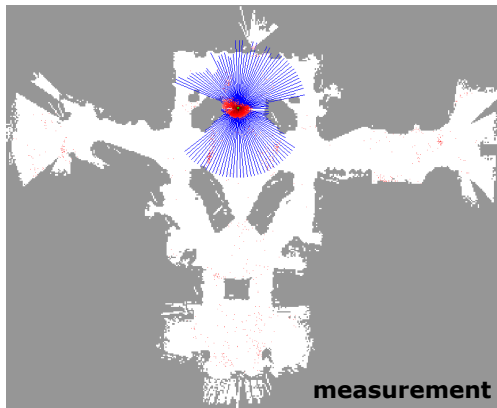
24

1



25

2



26

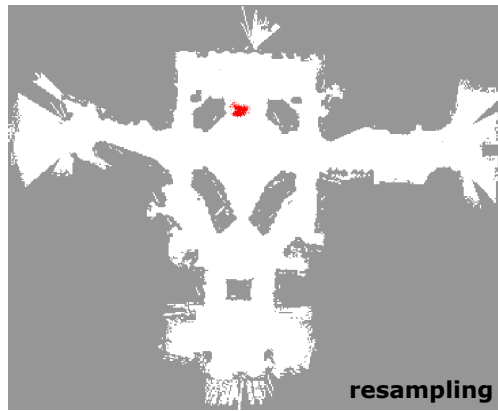
3





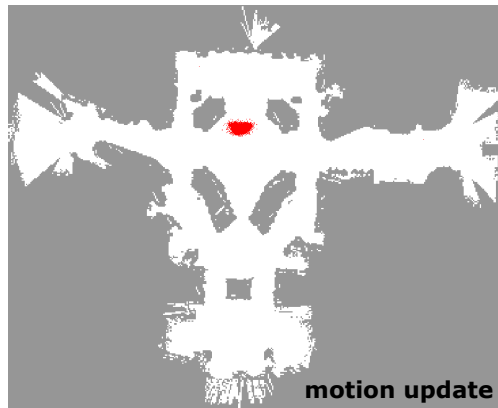
27

1



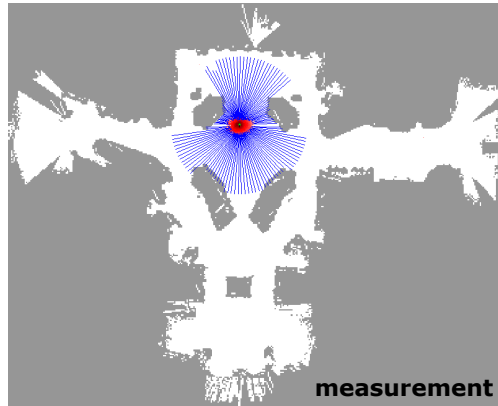
28

2



29

3



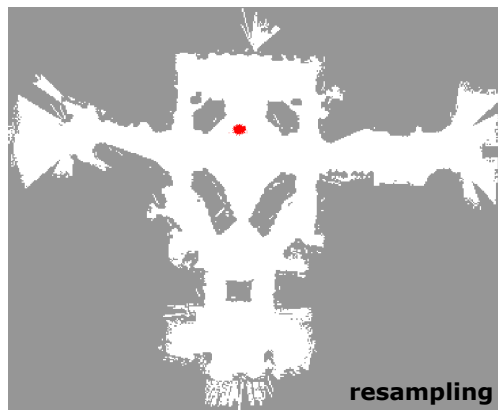
30

1



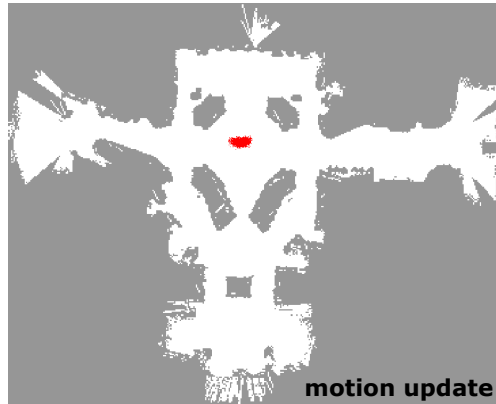
31

2



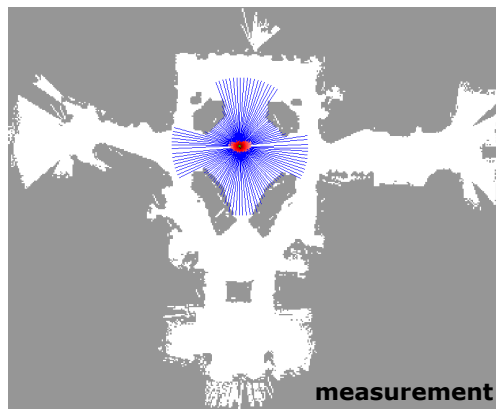
32

3



33

1



34

2

### 3.8.5 Theoretical insight into particle filtering

**Step 1: Propagating the dynamics** As we introduced in the section on Markov Decision Processes (MDPs), a stochastic dynamical system

$$x_{k+1} = f(x_k, u_k) + \epsilon_k$$

is equivalent to a probability transition matrix  $x_{k+1} \sim P(x_{k+1} | x_k, u_k)$ . Our goal is to approximate the distribution of  $x_{k+1|k}$  using particles. What proposal distribution should we choose? The “closest” probability distribution to  $x_{k+1|k}$  that we have available is  $x_{k|k}$ . So we set

$$\begin{aligned} \text{target} &: P(x_{k+1} | y_1, \dots, y_k) \\ \text{proposal} &: P(x_k | y_1, \dots, y_k) \end{aligned}$$

6

Suppose we had performed resampling on our particle set from the distri-

**i** In this sense, picking a proposal distribution to draw particles from is like linearization. Better the match between the proposal and the target, fewer samples we need to approximate the target.

1 bution  $x_{k|k}$  and have a set of  $n$  particles  $\{x_{k|k}^{(i)}\}$  with equal weights  $1/n$

$$P(x_k | y_1, \dots, y_k) \approx \frac{1}{n} \sum_{i=1}^n \delta_{x_{k|k}^{(i)}}(x).$$

2 Propagating the dynamics in a PF involves computing importance sampling  
3 weights. If we had a particle at location  $x$  that was supposed to approximate  
4 the distribution of  $x_{k+1|k}$ , as we saw for importance sampling, its importance  
5 weight is the ratio of the target and proposal densities at that location

$$w_{k+1|k}(x) = \frac{P(x_{k+1} = x | y_1, \dots, y_k)}{P(x_k = x | y_1, \dots, y_k)}.$$

6 Let us focus on the numerator. We have

$$\begin{aligned} P(x_{k+1} = x | y_1, \dots, y_k) &= \int P(x_{k+1} = x, x_k = x' | y_1, \dots, y_k) dx_k \\ &= \int P(x_{k+1} = x | x_k = x', y_1, \dots, y_k) P(x_k = x' | y_1, \dots, y_k) dx_k \\ &= \int P(x_{k+1} = x | x_k = x') P(x_k = x' | y_1, \dots, y_k) dx' \\ &\approx \frac{1}{n} \int P(x_{k+1} = x | x_k = x') \sum_{i=1}^n \delta_{x_{k|k}^{(i)}}(x') dx' \\ &= \frac{1}{n} \sum_{i=1}^n P(x_{k+1} = x | x_k = x_{k|k}^{(i)}, u = u_k), \end{aligned}$$

7 where the system dynamics is  $f(x_k, u_k) + \epsilon_k$  and  $u_k$  is the control at time  $k$ .  
8 The denominator  $P(x_k = x_{k|k}^{(i)} | y_1, \dots, y_k)$  when evaluated at particles  $x_{k|k}^{(i)}$   
9 is simply  $1/n$ . This gives us weights

$$w_{k+1|k}(x) = \sum_{i=1}^n P(x_{k+1} = x | x_k = x_{k|k}^{(i)}, u = u_k). \quad (3.37)$$

10 Let us now think about what particles we should pick for  $x_{k+1|k}$ . We have  
11 from (3.37) a function that lets us compute the correct weight for any particle  
12 we may choose to approximate  $x_{k+1|k}$ .

13 Say we keep the particle locations unchanged, i.e.,  $x_{k+1|k}^{(i)} = x_{k|k}^{(i)}$ . We  
14 then have

$$P(x_{k+1} = x | y_1, \dots, y_k) \approx \sum_{i=1}^n w_{k+1|k}(x_{k|k}^{(i)}) \delta_{x_{k|k}^{(i)}}(x). \quad (3.38)$$

15 You will notice that keeping the particle locations unchanged may be a very  
16 poor approximation. After all, the probability density  $P(x_{k+1} | y_1, \dots, y_k)$  is  
17 large, not at the particles  $x_{k|k}^{(i)}$  (that were a good approximation of  $x_{k|k}$ ), but  
18 rather at the transformed locations of these particles

$$f(x_{k|k}^{(i)}, u_k).$$

**i** Draw a picture of how this approximation looks.

1 We will therefore update the locations of the particles to be

$$x_{k+1|k}^{(i)} = f(x_{k|k}^{(i)}, u_k) \quad (3.39)$$

2 with weight of the  $i^{\text{th}}$  particle given by

$$\begin{aligned} w_{k+1|k}^{(i)} &:= w_{k+1|k}(x_{k+1|k}^{(i)}) = \sum_{j=1}^n \mathbf{P}(x_{k+1} = x_{k+1|k}^{(i)} \mid x_k = x_{k|k}^{(j)}, u = u_k) \\ &\approx \mathbf{P}(x_{k+1} = x_{k+1|k}^{(i)} \mid x_k = x_{k|k}^{(i)}, u = u_k). \end{aligned} \quad (3.40)$$

3 The approximation in the above equation is very crude: we are essentially say-  
4 ing that each particle  $x_{k|k}^{(i)}$  is transformed independently of the other particles  
5 to a new location  $x_{k+1|k}^{(i)} = f(x_{k|k}^{(i)}, u_k)$ . This completes the first step of a  
6 particle filter and we have

$$\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k) \approx \sum_{i=1}^n w_{k+1|k}^{(i)} \delta_{x_{k+1|k}^{(i)}}(x).$$

7 **Step 2: Incorporating the observation** The target and proposal distributions  
8 in this case are

$$\begin{aligned} \text{target} &: \mathbf{P}(x_{k+1} \mid y_1, \dots, y_k, y_{k+1}) \\ \text{proposal} &: \mathbf{P}(x_{k+1} \mid y_1, \dots, y_k). \end{aligned}$$

9 Since we have particles  $x_{k+1|k}^{(i)}$  with weights  $w_{k+1|k}^{(i)}$  for the proposal distri-  
10 bution obtained from the propagation step, we now like to update them to  
11 incorporate the latest observation  $y_{k+1}$ . Let us imagine for a moment that the  
12 weights  $w_{k+1|k}^{(i)}$  are uniform. We would then set weights

$$\begin{aligned} w(x) &= \frac{\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k, y_{k+1})}{\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k)} \\ &\propto \frac{\mathbf{P}(y_{k+1} \mid x_{k+1} = x) \mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k)}{\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k)} \quad (\text{by Bayes rule}) \\ &= \mathbf{P}(y_{k+1} \mid x_{k+1} = x). \end{aligned}$$

13 for each particle  $x = x_{k+1|k}^{(i)}$  to get the approximated distribution as

$$\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_{k+1}) \approx \sum_{i=1}^n \mathbf{P}(y_{k+1} \mid x_{k+1|k}^{(i)}) w_{k+1|k}^{(i)} \delta_{x_{k+1|k}^{(i)}}(x) \quad (3.41)$$

14 You will notice that the right hand side is not normalized and the distribution  
15 does not integrate to 1 (why? because we did not write the proportionality  
16 constant in the Bayes rule above). This is easily fixed by normalizing the

1 coefficients  $\mathbf{P}(y_{k+1} | x_{k+1|k}^{(i)}) w_{k+1|k}^{(i)}$  to sum to 1 as follows

$$w_{k+1|k+1}^{(i)} := \frac{\mathbf{P}(y_{k+1} | x_{k+1|k}^{(i)}) w_{k+1|k}^{(i)}}{\sum_j \mathbf{P}(y_{k+1} | x_{k+1|k}^{(j)}) w_{k+1|k}^{(j)}}.$$

2 **Step 3: Resampling step** As we discussed in the previous section, after  
 3 incorporating the observation, some particles may have very small weights.  
 4 The resampling procedure resamples particles so that all of them have equal  
 5 weights  $1/n$ .

$$\left\{ x_{k+1|k+1}^{(i)}, 1/n \right\}_{i=1}^n = \text{resample} \left( \left\{ x_{k+1|k+1}^{(i)}, w_{k+1|k+1}^{(i)} \right\}_{i=1}^n \right).$$

## 6 3.9 Discussion

7 This brings our study of filtering to a close. We have looked at some of the most  
 8 important algorithms for a variety of dynamical systems, both linear and nonlin-  
 9 ear. Although, we focused on filtering in this chapter, all these algorithms have  
 10 their corresponding “smoothing” variants, e.g., you can read about how a typi-  
 11 cal Kalman smoother is implemented at [https://en.wikipedia.org/wiki/Kalman\\_filter#Fixed-](https://en.wikipedia.org/wiki/Kalman_filter#Fixed-lag_smoother)  
 12 [lag\\_smoother](https://en.wikipedia.org/wiki/Kalman_filter#Fixed-lag_smoother). Filtering, and state estimation, is a very wide area of research  
 13 even today and you will find variants of these algorithms in almost every  
 14 device which senses the environment.

# Chapter 4

## Rigid-body transforms and mapping

### Reading

1. LaValle Chapter 3.2 for rotation matrices, Chapter 4.1-4.2 for quaternions
2. Thrun Chapter 9.1-9.2 for occupancy grids
3. OctoMap: An Efficient Probabilistic 3D Mapping Framework Based on Octrees  
<http://www.arminhornung.de/Research/pub/hornung13auro.pdf>, also see <https://octomap.github.io>.
4. Robot Operating System  
<http://www.willowgarage.com/sites/default/files/icraoss09-ROS.pdf>, Optional: Lightweight Communications and Marshalling (LCM) system  
<https://people.csail.mit.edu/albert/pubs/2010-huang-olson-moore-lcm-iros.pdf>
5. A Perception-Driven Autonomous Urban Vehicle  
<https://april.eecs.umich.edu/media/pdfs/mitduc2009.pdf>
6. Optional reading: Thrun Chapter 10 for simultaneous localization and mapping

In the previous chapter, we looked at ways to estimate the state of the robot in the physical world. We kept our formulation abstract, e.g., the way the robot moves was captured by an abstract expression like  $x_{k+1} = f(x_k, u_k) + \epsilon$  and observations  $y_k = g(x_k) + \nu$  were similarly opaque. In order to actually implement state estimation algorithms on real robots, we need to put concrete functions in place of  $f, g$ .

This is easy to do for some robots, e.g., the robot in Problem 1 in Homework 1 moved across cells. Of course real robots are a bit more complicated, e.g., a car cannot move sideways (which is a huge headache when you parallel park). In the first half of this chapter, we will look at how to model the dynamics  $f$  using rigid-body transforms.

The story of measurement models and sensors is similar. Although we need to write explicit formulae in place of the abstract function  $g$ . In the second half, we will study occupancy grids and dig deeper into a typical state-estimation problem in robotics, namely that of mapping the location of objects in the world around the robot.

## 4.1 Rigid-Body Transformations

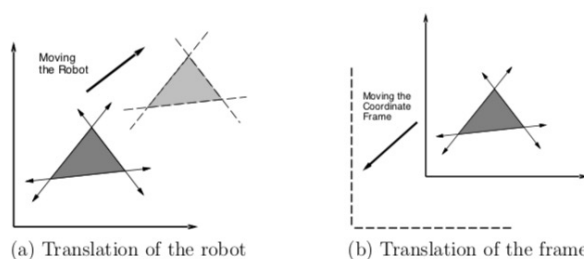
Let us imagine that the robot has a rigid body, we think of this as a subset  $A \subset \mathbb{R}^2$ . Say the robot is a disc

$$A = \{(x, y) \in \mathbb{R}^2 : x^2 + y^2 \leq 1\}.$$

This set  $A$  changes as the robot moves around, e.g., if the center of mass of the robot is translated by  $x_t, y_t \in \mathbb{R}$  the set  $A$  changes to

$$A' = \{(x + x_t, y + y_t) : (x, y) \in A\}.$$

The concept of “degrees of freedom” denotes the maximum number of independent parameters needed to completely characterize the transformation applied to a robot. Since the set of allowed values  $(x_t, y_t)$  is a two-dimensional subset of  $\mathbb{R}^2$ , then the degrees of freedom available to a translating robot is two.



As the above figure shows, there are two ways of thinking about this transformation. We can either think of the robot transforming while the co-ordinate frame of the world is fixed, or we can think of it as the robot remaining stationary and the co-ordinate frame undergoing a translation. The second style is useful if you want to imagine things from the robot’s perspective. But the first one feels much more natural and we will therefore exclusively use the first notion.

If the same robot is where rotated counterclockwise by some angle  $\theta \in [0, 2\pi]$ , we would map

$$(x, y) \mapsto (x \cos \theta - y \sin \theta, x \sin \theta + y \cos \theta).$$



1 Such a map can be written as multiplication by a  $2 \times 2$  rotation matrix

$$R(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}. \quad (4.1)$$

2 to get

$$\begin{bmatrix} x \cos \theta - y \sin \theta \\ x \sin \theta + y \cos \theta \end{bmatrix} = R(\theta) \begin{bmatrix} x \\ y \end{bmatrix}.$$

3 The transformed robot is thus given by

$$A' = \left\{ R \begin{bmatrix} x \\ y \end{bmatrix} : (x, y) \in A \right\}.$$

4 If we perform both rotation and translation, we can the transformation using a  
5 single matrix

$$T = \begin{bmatrix} \cos \theta & -\sin \theta & x_t \\ \sin \theta & \cos \theta & y_t \\ 0 & 0 & 1 \end{bmatrix} \quad (4.2)$$

6 and this transformation looks like

$$\begin{bmatrix} x \cos \theta - y \sin \theta + x_t \\ x \sin \theta + y \cos \theta + y_t \\ 1 \end{bmatrix} = T \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}.$$

7 The point  $(x, y, 1) \in \mathbb{R}^3$  is called homogeneous coordinate space correspond-  
8 ing to  $(x, y) \in \mathbb{R}^2$  and the matrix  $T$  is called a *homogeneous transformation*  
9 *matrix*. The peculiar names comes from the fact that even if the matrix  $T$   
10 maps rotations and translations of rigid bodies  $A \subset \mathbb{R}^2$ , it is just a linear  
11 transformation of the point  $(x, y, 1)$  if viewed in the larger space  $\mathbb{R}^3$ .

12 **Rigid-body transformations** The transformations  $R \in \mathbb{R}^{2 \times 2}$  or  $T \in \mathbb{R}^{3 \times 3}$   
13 are called rigid-body transformations. Mathematically, it means that they do  
14 not cause the distance between any two points inside the set  $A$  to change. Rigid-  
15 body transformations are what are called an orthogonal group in mathematics.

16 **A group** is a mathematical object which imposes certain conditions upon  
17 how two operations, e.g., rotations, can be composed together. For instance,  
18 if  $G$  is the group of rotations, then (i) the composition of two rotations is a  
19 rotation, we say that it satisfies closure  $R(\theta_1)R(\theta_2) \in G$ , (ii) rotations are  
20 associative

$$R(\theta_1) \{R(\theta_2)R(\theta_3)\} = \{R(\theta_1)R(\theta_2)\} R(\theta_3),$$

21 and, (iii) there exists an identity and inverse rotation

$$R(0), R(-\theta) \in G.$$

22 **An orthogonal group** is a group whose operations preserve distances in  
23 Euclidean space, i.e.,  $g \in G$  is an element of the group that acts on two points

**i** It is important to remember that  $T$  represents rotation *followed* by a translation, not the other way around.

1  $x, y \in \mathbb{R}^d$  then

$$\|g(x) - g(y)\| = \|x - y\|.$$

2 If we identify the basis in Euclidean space to be the set of orthonormal vectors  $\{e_1, \dots, e_d\}$ , then equivalently, the orthogonal group  $O(d)$  is the set of  
3 orthogonal matrices  
4

$$O(d) := \{O \in \mathbb{R}^{d \times d} : OO^\top = O^\top O = I\}.$$

5 This implies that the square of the determinant of any element  $a \in O(d)$  is 1,  
6 i.e.,  $\det(a) = \pm 1$ .

🔗 Check that any rotation matrix  $R$  belongs to an orthogonal group.

7 **The Special Orthogonal Group** is a sub-group of the orthogonal group  
8 where the determinant of each element is +1. You can see that rotations are a  
9 special orthogonal group. We denote rotations of objects in  $\mathbb{R}^2$  as

$$SO(2) := \{R \in \mathbb{R}^{2 \times 2} : R^\top R = RR^\top = I, \det(R) = 1\}. \quad (4.3)$$

10 Each group element  $g \in SO(2)$  denotes a rotation of the  $XY$ -plane about the  
11  $Z$ -axis. The group of 3D rotations is called the Special Orthogonal Group  
12  $SO(3)$  and is defined similarly

$$SO(3) := \{R \in \mathbb{R}^{3 \times 3} : R^\top R = RR^\top = I, \det(R) = 1\}. \quad (4.4)$$

13 **The Special Euclidean Group**  $SE(2)$  is simply a composition of a 2D  
14 rotation  $R \in SO(2)$  and a 2D translation  $\mathbb{R}^2 \ni v \equiv (x_t, y_t)$

$$SE(2) = \left\{ \begin{bmatrix} R & v \\ 0 & 1 \end{bmatrix} : R \in SO(2), v \in \mathbb{R}^2 \right\} \subset \mathbb{R}^{3 \times 3}. \quad (4.5)$$

15 The Special Euclidean Group  $SE(3)$  is defined similarly as

$$SE(3) = \left\{ \begin{bmatrix} R & v \\ 0 & 1 \end{bmatrix} : R \in SO(3), v \in \mathbb{R}^3 \right\} \subset \mathbb{R}^{4 \times 4}; \quad (4.6)$$

16 again, remember that it is rotation *followed* by a translation.

### 17 4.1.1 3D transformations

18 Translations and rotations in 3D are conceptually similar to the two-dimensional  
19 case; however the details appear a bit more difficult because rotations in 3D  
20 are more complicated.

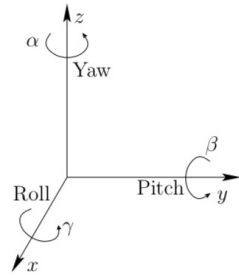


Figure 4.1: Any three-dimensional rotation can be described as a sequence of rotations about each of the cardinal axes. We usually give these specific names: rotation about the  $Z$ -axis is called yaw, rotation about the  $X$ -axis is called roll and rotation about  $Y$ -axis is called pitch. You should commit this picture and these names to memory because it will be of enormous to think about these rotations intuitively.

- 1 **Euler angles** We know that a pure counter-clockwise rotation about one of  
 2 the axes is written in terms of a matrix, say yaw of  $\alpha$ -radians about the  $Z$ -axis

$$R_z(\alpha) = \begin{bmatrix} \cos \alpha & -\sin \alpha & 0 \\ \sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

- 3 Notice that this is a  $3 \times 3$  matrix that keeps the  $Z$ -coordinate unchanged and  
 4 only affects the other two coordinates. Similarly we have for pitch ( $\beta$  about  
 5 the  $Y$ -axis) and roll ( $\gamma$  about the  $X$ -axis)

$$R_y(\beta) = \begin{bmatrix} \cos \beta & 0 & \sin \beta \\ 0 & 1 & 0 \\ -\sin \beta & 0 & \cos \beta \end{bmatrix}, R_x(\gamma) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \gamma & -\sin \gamma \\ 0 & \sin \gamma & \cos \gamma \end{bmatrix}.$$

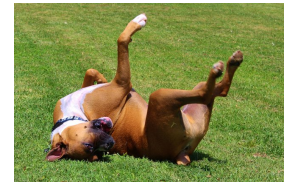
- 6 By convention, a rotation matrix in three dimensions is understood as a se-  
 7 quential application of rotations, *first roll, then pitch, and then yaw*

$$\mathbb{R}^{3 \times 3} = R(\alpha, \beta, \gamma) = R_z(\alpha)R_y(\beta)R_x(\gamma). \quad (4.7)$$

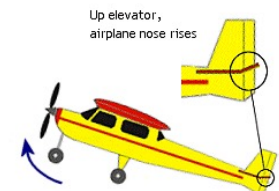
- 8 The angles  $(\alpha, \gamma, \gamma)$  are called Euler angles. Imagine how the body frame of  
 9 the robot changes as successive rotations are applied. If you were sitting in a  
 10 car, a pure yaw would be similar to the car turning left; the  $Z$ -axis correspond-  
 11 ing to this yaw would however only be pointing straight up perpendicular to  
 12 the ground if you had not performed a roll/pitch before. If you had, the  $Z$ -axis  
 13 of the body frame with respect to the world will be tilted.

- 14 Another important thing to note is that one single parameter determines  
 15 all possible rotations about one axis, i.e.,  $SO(2)$ . But three Euler angles are  
 16 used to parameterize general rotations in three-dimensions. You can watch  
 17 <https://www.youtube.com/watch?v=3Zjf95Jw2UE> to get more intuition about  
 18 Euler angles.

❗ Here is how I remember these names. Say you are driving a car, usually in robotics we take the  $X$ -axis to be longitudinally forward, the  $Y$ -axis is your left hand if you are in the driver's seat and the  $Z$ -axis points up by the right-hand thumb rule. Roll



is what a dog does when it rolls, it rotates about the  $X$ -axis. Pitch is what a plane



does when it takes off, its nose lifts up and it rotates about the  $Y$ -axis. Yaw is the one leftover.

1 **Rotation matrices to Euler angles** We can back-calculate the Euler angles  
 2 from a rotation matrix as follows. Given an arbitrary matrix

$$R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix},$$

3 we set

$$\begin{aligned} \alpha &= \tan^{-1}(r_{21}/r_{11}) \\ \beta &= \tan^{-1}\left(-r_{31}/\sqrt{r_{32}^2 + r_{33}^2}\right) \\ \gamma &= \tan^{-1}(r_{32}/r_{33}). \end{aligned} \quad (4.8)$$

4 For each angle, the corresponding quadrant for the Euler angle is determined  
 5 using the signs of the numerator and the denominator. So you should use  
 6 the function atan2 in Python/C++ to implement these expressions correctly.  
 7 Notice that some of the expressions have  $r_{11}$  and  $r_{33}$  in the denominator, this  
 8 means that we need  $r_{11} = \cos \alpha \cos \beta = 0$  and  $r_{33} = \cos \beta \cos \gamma = 0$ . Euler  
 9 angles do not completely parametrize the group of rotations  $SO(3)$ . This is  
 10 really an indicator that this particular physical rotation can be parameterized  
 11 in two different ways using Euler angles, so the map from rotation matrices to  
 12 Euler angles is not unique.

13 **Homogeneous coordinates in three dimensions** Just like the 2D case, we  
 14 can define a  $4 \times 4$  matrix that transforms points  $(x, y, z) \in \mathbb{R}^3$  to their new  
 15 locations after a rotation by Euler angles  $(\alpha, \beta, \gamma)$  and a translation by a vector  
 16  $v = (x_t, y_t, z_t) \in \mathbb{R}^3$

$$T = \begin{bmatrix} R(\alpha, \beta, \gamma) & v \\ 0 & 1 \end{bmatrix}.$$

#### 17 4.1.2 Rodrigues' formula: an alternate view of rotations

18 Consider a point  $r(t) \in \mathbb{R}^3$  that is being rotated about an axis denoted by a  
 19 unit vector  $\omega \in \mathbb{R}^3$  with an angular velocity of 1 radian/sec. The instantaneous  
 20 linear velocity of the head of the vector is

$$\dot{r}(t) = \omega \times r(t) \equiv \hat{\omega}r(t) \quad (4.9)$$

21 where the  $\times$  denotes the cross-product of the two vectors  $a, b \in \mathbb{R}^3$

$$a \times b = \begin{bmatrix} a_2b_3 - a_3b_2 \\ a_3b_1 - a_1b_3 \\ a_1b_2 - a_2b_1 \end{bmatrix}$$

22 which we can equivalently denote as a matrix vector multiplication  $a \times b = \hat{a}b$   
 23 where

$$\hat{a} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix}. \quad (4.10)$$

1 is a skew-symmetric matrix. The solution of the differential equation (4.9) at  
 2 time  $t = \theta$  is

$$r(\theta) = \exp(\hat{\omega}\theta) r(0)$$

3 where the matrix exponential of a matrix  $A$  is defined as

$$\exp(A) = I + A + \frac{A^2}{2!} + \frac{A^3}{3!} + \dots = \sum_{k=0}^{\infty} \frac{A^k}{k!}.$$

4 This is an interesting observation: a rotation about a fixed axis  $\omega$  by an angle  $\theta$   
 5 can be represented by the matrix

$$R = \exp(\hat{\omega} \theta).$$

6 You can check that this matrix is indeed a rotation by showing that  $R^\top R = I$   
 7 and that  $\det(R) = +1$ . We can expand the matrix exponential and collect odd  
 8 and even powers of  $\hat{\omega}$  to get

$$R = I + \sin \theta \hat{\omega} + (1 - \cos \theta) \hat{\omega}^2. \quad (4.11)$$

9 which is the Rodrigues' formula that relates the angle  $\theta$  and the axis  $\omega$  to the  
 10 rotation matrix. We can also go in the opposite direction, i.e., given a matrix  
 11  $R$  calculate what angle  $\theta$  and axis  $\omega$  it corresponds to using

$$\begin{aligned} \cos \theta &= \frac{\text{tr}(R) - 1}{2} \\ \hat{\omega} &= \frac{R - R^\top}{2 \sin \theta}. \end{aligned} \quad (4.12)$$

## 12 4.2 Quaternions

13 We know two ways to think about rotations: we can either think in terms of the  
 14 three Euler angles  $(\alpha, \beta, \gamma)$ , or we can consider a rotation matrix  $R \in \mathbb{R}^{3 \times 3}$ .  
 15 We also know ways to go to and fro between these two forms with the caveat  
 16 that solving for Euler angles using (4.8) may be degenerate in some cases.  
 17 While rotation matrices are the most general representation of rotations, using  
 18 them in computer code is cumbersome (it is, after all, a matrix of 9 elements),  
 19 imagine an EKF where the state is a rotation matrix. You can certainly  
 20 implement the same filter using Euler angles but doing so requires special  
 21 care due to the degeneracies. Quaternions are a neat way to avoid both these  
 22 problems, they parametrize the space of rotations using 4 numbers and do not  
 23 have degeneracies.

24 The central idea behind quaternions is Euler's theorem which says that  
 25 any 3D rotation can be considered as a pure rotation by an angle  $\theta \in \mathbb{R}$  about  
 26 an axis given by the unit vector  $\omega$ . This is the result that we also exploited in  
 27 Rodrigues' formula.

**i** Quaternions were invented by British mathematician William Rowan Hamilton while walking along a bridge with his wife. He was quite excited by this discovery and promptly graffitied the expression into the stone of the bridge

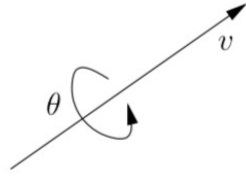


Figure 4.2: Any rotation in 3D can be represented using a unit vector  $\omega$  and an angle  $\theta \in \mathbb{R}$ . Notice that there are two ways to encode the same rotation, the unit vector  $-\omega$  and angle  $2\pi - \theta$  would give the same rotation. Mathematicians say this as quaternions being a double-cover of  $SO(3)$ .

1 A quaternion  $q$  as a four-dimensional vector  $q \equiv (u_0, u_1, u_2, u_3)$  and we  
2 write it as

$$q \equiv (u_0, u), \text{ or} \quad (4.13)$$

$$q = u_0 + u_1 i + u_2 j + u_3 k,$$

3 with  $i, j, k$  being three “imaginary” components of the quaternion with “complex-  
4 numbers like” relationships

$$i^2 = j^2 = k^2 = ijk = -1. \quad (4.14)$$

5 It follows from these relationships that

$$ij = -ji = k, ki = -ik = j, \text{ and } jk = -jk = i.$$

6 Although you may be tempted to think about this, these imaginary components  
7  $i, j, k$  have no relationship with the square roots of negative unity used to  
8 define standard complex numbers. You should simply think of the quaternion  
9 as a four-dimensional vector in some space. A unit quaternion, i.e., one with

$$u_0^2 + u_1^2 + u_2^2 + u_3^2 = 1,$$

10 is special: unit quaternions can be used to represent rotations in 3D.

11 **Quaternion to axis-angle representation** The quaternion  $q = (u_0, u)$  cor-  
12 responds to a counterclockwise rotation of angle  $\theta$  about a unit vector  $\omega$   
13 where  $\theta$  and  $\omega$  are such that

$$u_0 = \cos \frac{\theta}{2}, \text{ and } u = \sin \frac{\theta}{2} \omega. \quad (4.15)$$

14 So given an axis-angle representation of rotation like in Rodrigues’ formula  
15  $(\theta, \omega)$  we can write the quaternion as

$$q = \left( \cos \frac{\theta}{2}, \sin \frac{\theta}{2} \omega \right).$$

1 Using this, we can also compute the inverse of a quaternion (rotation of angle  
2  $\theta$  about the opposite axis  $-\omega$ ) as

$$q^{-1} := \left( \cos \frac{\theta}{2}, -\sin \frac{\theta}{2} \omega \right).$$

3 The inverse quaternion is therefore the quaternion where all entries except the  
4 first have their signs flipped.

5 **Multiplication of quaternions** Just like two rotation matrices multiply to-  
6 gether to give a new rotation (rotations are a group), since quaternions are  
7 also a representation for the group of rotations, we can also multiply two  
8 quaternions  $q_1 = (u_0, u)$ ,  $q_2 = (v_0, v)$  together using the quaternion identities  
9 for  $i, j, k$  in (4.14) to get a new quaternion

$$q_1 q_2 \equiv (u_0, u) \cdot (v_0, v) = (u_0 v_0 - u^\top v, u_0 v + v_0 u + u \times v).$$

10

11 **Pure quaternions** A pure quaternion is a quaternion with a zero scalar value  
12  $u_0 = 0$ . This is very useful to simply store a standard 3D vector  $u \in \mathbb{R}^3$  as a  
13 quaternion  $(0, u)$ . We can then rotate points easily between different frames  
14 as follows. Given a vector  $x \in \mathbb{R}^3$  we can form a quaternion  $(0, x)$  and show  
15 that

$$q \cdot (0, x) \cdot q^* = (0, R(q)x). \quad (4.16)$$

16 where  $q^* = (u_0, -u)$  is the conjugate quaternion of  $q = (u_0, u)$ ; the conjugate  
17 is the same as the inverse for unit quaternions. Notice how the right-hand side  
18 is the vector  $R(q)x$  corresponding to the vector  $x$  rotation by a matrix  $R(q)$ .

19 **Quaternions to rotation matrix** The rotation matrix corresponding to a  
20 quaternion is

$$\begin{aligned} R(q) &= (u_0^2 - u^\top u) I_{3 \times 3} + 2 \frac{u_0 u}{\|u\|} + 2 u u^\top \\ &= \begin{bmatrix} 2(u_0^2 + u_1^2) - 1 & 2(u_1 u_2 - u_0 u_3) & 2(u_1 u_3 + u_0 u_2) \\ 2(u_1 u_2 + u_0 u_3) & 2(u_0^2 + u_2^2) - 1 & 2(u_2 u_3 - u_0 u_1) \\ 2(u_1 u_3 - u_0 u_2) & 2(u_2 u_3 - u_0 u_1) & 2(u_0^2 + u_3^2) - 1 \end{bmatrix}. \end{aligned} \quad (4.17)$$

21 Using this you can show the identity that rotation matrix corresponding to the  
22 product of two quaternions is the product of the individual rotation matrices

$$R(q_1 q_2) = R(q_1) R(q_2).$$

**i** Quaternions belong to a larger group than rotations called the Symplectic Group  $Sp(1)$ .

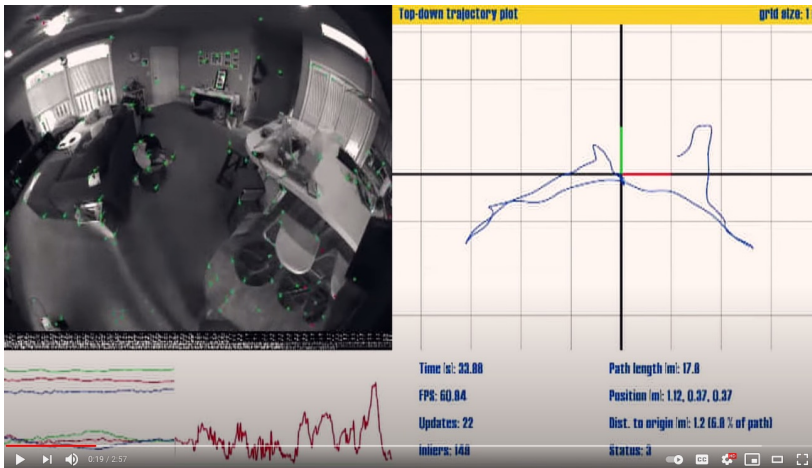
- 1 **Rotation matrix to quaternion** We can also go in the reverse direction.  
 2 Given a rotation matrix  $R$ , the quaternion is

$$\begin{aligned}
 u_0 &= \frac{1}{2} \sqrt{r_{11} + r_{22} + r_{33} + 1} \\
 \text{if } u_0 \neq 0, u_1 &= \frac{r_{32} - r_{23}}{4u_0} \\
 u_2 &= \frac{r_{13} - r_{31}}{4u_0} \\
 u_3 &= \frac{r_{21} - r_{12}}{4u_0} \\
 \text{if } u_0 = 0, u_1 &= \frac{r_{13}r_{12}}{\sqrt{r_{12}^2r_{13}^2 + r_{12}^2r_{23}^2 + r_{13}^2r_{23}^2}} \\
 u_2 &= \frac{r_{12}r_{23}}{\sqrt{r_{12}^2r_{13}^2 + r_{12}^2r_{23}^2 + r_{13}^2r_{23}^2}} \\
 u_3 &= \frac{r_{13}r_{23}}{\sqrt{r_{12}^2r_{13}^2 + r_{12}^2r_{23}^2 + r_{13}^2r_{23}^2}}.
 \end{aligned} \tag{4.18}$$

### 3 4.3 Occupancy Grids

- 4 Rotation matrices and quaternions let us capture the dynamics of a rigid robot  
 5 body. We will next look at how to better understand observations.

- 6 **What is location and what is mapping?** Imagine a robot that is moving  
 7 around in a house. A natural representation of the state of this robot is the 3D lo-  
 8 cation of all the interesting objects in the room, e.g., <https://www.youtube.com/watch?v=Qe10ExwzCqk>.  
 9 At each time-instant, we record an observation from our sensor (in this case,  
 10 a camera) that indicates how far an object is from the robot. This helps us  
 11 discover the location of the objects in the room. After gathering enough ob-  
 12 servations, we would have created a *map* of the entire house. This map is the  
 13 set of positions of all interesting objects in the room. Such a map is called a  
 14 “feature map”, these are all the green points in the image below



- 15  
 16 The main point to understand about feature map is that we can hand over  
 17 this map to another robot that comes to the same house. The robot compares



1 images from its camera and if it finds one of the objects inside the map, it  
 2 can get an estimate of its location/orientation in the room with respect to the  
 3 known location of the object in the map. The map is just a set of “features”  
 4 that help identify salient objects in the room (objects which can be easily  
 5 detected in images and relatively uniquely determine the location inside the  
 6 room). The second robot using this map to estimate its position/orientation in  
 7 the room is called the *localization problem*. We already know how to solve  
 8 the localization problem using filtering.

9 The first robot was solving a harder problem called *Simultaneous Local-*  
 10 *ization And Mapping (SLAM)*: namely that of discovering the location of both  
 11 itself and the objects in the house. This is a very important and challenging  
 12 problem in robots but we will not discuss it further. MEAM 620 digs deeper  
 13 into it.

In this section, we will solve a part of the SLAM problem, namely the mapping problem. We will assume that we know the position/orientation of the robot in the 3D world, and want to build a map of the objects in the world. We will discuss grid maps, which are a more crude way of representing maps than feature maps but can be used easily even if there are lots of objects.

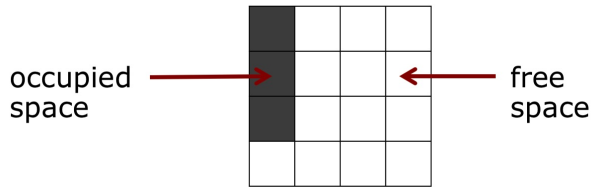
14 **Grid maps** We will first discuss two-dimensional grid maps, they look as  
 15 follows.



Figure 4.3: A grid map (also called an occupancy grid) is a large gray-scale image, each pixel represents a cell in the physical world. In this picture, cells that are occupied are colored black and empty cells represent free space. A grid map is a useful representation for a robot to localize in this house using observations from its sensors and comparing those to the map.

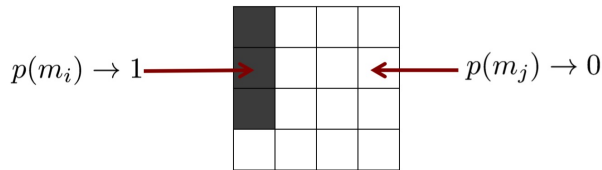
16 To get a quick idea of what we want to do, you can watch the mapping  
 17 being performed in <https://www.youtube.com/watch?v=JJhEkIA1xSE>. We are  
 18 interested in learning such maps from the observations that a robot collects as  
 19 it moves around the physical space. Let us make two simplifying assumptions.

20 **Assumption 1: each cell is either free or occupied**



1

2 This is neat: we can now model each cell as a binary random variable that  
 3 indicates occupancy. Let the probability that the cell  $m_i$  be occupied be  $p(m_i)$

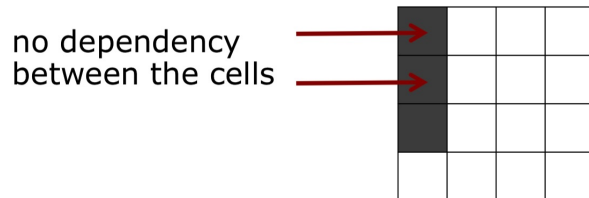


4

5 If we have  $p(m_i) = 0$ , then the cell is not occupied and if we have  $p(m_i) = 1$ ,  
 6 then the cell is occupied. A priori, we do not know the state of the cell so we  
 7 will set the prior probability to be  $p(m_i) = 0.5$ .

8 **Assumption 2: the world is static** Objects in the world do not move. This  
 9 is reasonable if we are interested in estimating in building a map of the walls  
 10 inside the room. Note that it is not a reasonable assumption if there are moving  
 11 people inside the room. We will see a clever hack where the Bayes rule helps  
 12 automatically disregard such moving objects in this section.

13 **Assumption 3: cells are independent of each other** This is another drastic  
 14 simplification. The state of our system is the occupancy of each cell in the  
 15 grid map. We assume that before receiving any observations, the occupancy of  
 16 each individual cell is independent; it is a Bernoulli variable with probability  
 17  $1/2$  since we have assumed the prior to be uniform in Assumption 1.



18

19 This means that if cells in the map are denoted by a vector  $m = (m_1, \dots)$ ,  
 20 then the probability of the cells being occupied/not-occupied can be written as

21

$$p(m) = \prod_i p(m_i). \quad (4.19)$$

example map  
(4-dim vector)

4 individual cells

1

### 2 4.3.1 Estimating the map from the data

3 Say that the robot pose (position and orientation) is given by the sequence  
4  $x_1, \dots, x_k$ . While proceeding along this sequence, the robot receives obser-  
5 vations  $y_1, \dots, y_k$ . Our goal is to estimate the state of each cell  $m_i \in \{0, 1\}$   
6 (aka “the map”  $m = (m_1, m_2, \dots)$ )

$$P(m | x_1, \dots, x_k, y_1, \dots, y_k) = \prod_i P(m_i | x_1, \dots, x_k, y_1, \dots, y_k). \quad (4.20)$$

7 This is called the “static state” Bayes filter and is conceptually exactly the  
8 same as the recursive application of Bayes rule in Chapter 2 for detecting  
9 whether the door was open or closed.

10 We will use a short form to keep the notation clear

$$y_{1:k} = (y_1, y_2, \dots, y_k);$$

11 the quantity  $x_{1:k}$  is defined similarly. As usual we will use a recursive Bayes  
12 filter to compute this probability as follows.

$$\begin{aligned} P(m_i | x_{1:k}, y_{1:k}) &\stackrel{\text{Bayes rule}}{=} \frac{P(y_k | m_i, y_{1:k-1}, x_{1:k}) P(m_i | y_{1:k-1}, x_{1:k})}{P(y_k | y_{1:k-1}, x_{1:k})} \\ &\stackrel{\text{Markov}}{=} \frac{P(y_k | m_i, x_k) P(m_i | y_{1:k-1}, x_{1:k-1})}{P(y_k | y_{1:k-1}, x_{1:k})} \\ &\stackrel{\text{Bayes rule}}{=} \frac{P(m_i | y_k, x_k) P(y_k | x_k) P(m_i | y_{1:k-1}, x_{1:k-1})}{P(m_i | x_k) P(y_k | y_{1:k-1}, x_{1:k})} \\ &\stackrel{\text{Markov}}{=} \frac{P(m_i | y_k, x_k) P(y_k | x_k) P(m_i | y_{1:k-1}, x_{1:k-1})}{P(m_i) P(y_k | y_{1:k-1}, x_{1:k})}. \end{aligned} \quad (4.21)$$

13 We have a similar expression for the opposite probability

$$P(\neg m_i | x_{1:k}, y_{1:k}) = \frac{P(\neg m_i | y_k, x_k) P(y_k | x_k) P(\neg m_i | y_{1:k-1}, x_{1:k-1})}{P(\neg m_i) P(y_k | y_{1:k-1}, x_{1:k})}.$$

14 Let us take the ratio of the two to get

$$\begin{aligned} \frac{P(m_i | x_{1:k}, y_{1:k})}{P(\neg m_i | x_{1:k}, y_{1:k})} &= \frac{P(m_i | y_k, x_k) P(m_i | y_{1:k-1}, x_{1:k-1}) P(\neg m_i)}{P(\neg m_i | y_k, x_k) P(\neg m_i | y_{1:k-1}, x_{1:k-1}) P(m_i)} \\ &= \underbrace{\frac{P(m_i | y_k, x_k)}{1 - P(m_i | y_k, x_k)}}_{\text{uses observation } y_k} \underbrace{\frac{P(m_i | y_{1:k-1}, x_{1:k-1})}{1 - P(m_i | y_{1:k-1}, x_{1:k-1})}}_{\text{recursive term}} \underbrace{\frac{1 - P(m_i)}{P(m_i)}}_{\text{prior}}. \end{aligned} \quad (4.22)$$

15 This is called the odds ratio. Notice that the first term uses the latest obser-  
16 vation  $y_k$ , the second term can be updated recursively because it is a similar

1 expression as the left-hand side and the third term is a prior probability of the  
 2 cell being occupied/non-occupied. Let us rewrite this formula using the log-  
 3 odds-ratio that makes implementing it particularly easy. The log-odds-ratio of  
 4 the probability  $p(x)$  of a binary variable  $x$  is defined as

$$l(x) = \log \frac{p(x)}{1 - p(x)}, \text{ and } p(x) = 1 - \frac{1}{1 + e^{l(x)}}.$$

5 The product in (4.22) now turns into a sum as

$$l(m_i | y_{1:k}, x_{1:k}) = l(m_i | y_k, x_k) + l(m_i | y_{1:k-1}, x_{1:k-1}) - l(m_i). \quad (4.23)$$

6 This expression is used to update the occupancy of each cell. The term

$$\text{sensor model} = l(m_i | y_k, x_k)$$

7 is different for different sensors and we will investigate it next.

### 8 4.3.2 Sensor models

9 **Sonar** This works by sending out an ultrasonic chirp and measuring the time  
 10 between emission and reception of the signal. The time gives an estimate of  
 11 the distance of an object to the robot.

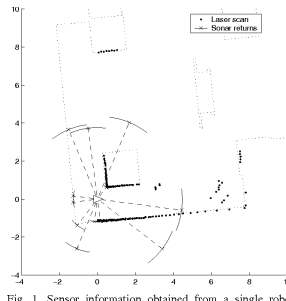
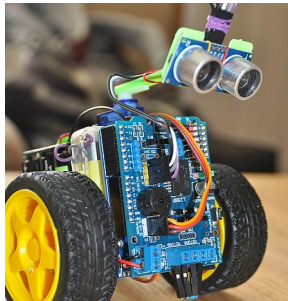


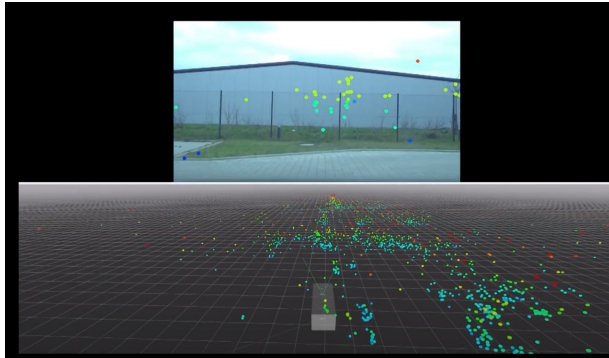
Fig 1 Sensor information obtained from a single robot

13 The figure above shows a typical sonar sensor (the two “eyes”) on a low-cost  
 14 robot. Data from the sensor is shown on the right, a sonar is a very low  
 15 resolution sensor and has a wide field of view, say 15 degrees, i.e., it cannot  
 16 differentiate between objects that are within 15 degrees of each other and  
 17 registers them as the same point. Sophisticated sonar technology is used today  
 18 in marine environments (submarines, fish finders, detecting mines etc.).

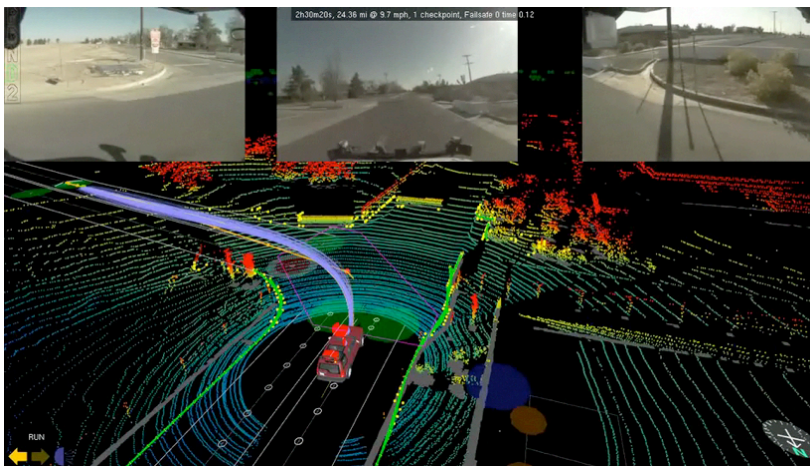
19 **Radar** works in much the same way as a sonar except that it uses pulses  
 20 of radio waves and measures the phase difference between the transmitted  
 21 and the received signal. This is a very versatile sensor (it was invented by the  
 22 US army to track planes and missiles during World War II) but is typically  
 23 noisy and requires sophisticated processing to be used for mainstream robotics.  
 24 Autonomous cars, collision warning systems on human-driven cars, weather  
 25 sensing, and certainly the military use the radar today. The following picture  
 26 and the video [https://www.youtube.com/watch?v=hwKUcu\\_7F9E](https://www.youtube.com/watch?v=hwKUcu_7F9E) will give  
 27 you an appreciation of the kind of data that a radar records. Radar is a very

❓ We assumed that the map was static. Can you think of why (4.23) automatically lets us handle some moving objects? Think of what the prior odds  $l(m_i)$  does to the log-odds-ratio  $l(m_i | y_{1:k}, x_{1:k})$ .

1 long range sensor (typically 150 m) and works primarily to detect metallic  
2 objects.



3  
4 **LiDAR** LiDAR, which is short for Light Detection and Ranging, (<https://en.wikipedia.org/wiki/Lidar>)  
5 is a portmanteau of light and radar. It is a sensor that uses a pulsed laser as the  
6 source of illumination and records the time it takes (nanoseconds typically) for  
7 the signal to return to the source. See <https://www.youtube.com/watch?v=NZKvf1cXe8s>  
8 for how the data from a typical LiDAR (Velodyne) looks like. While a Velo-  
9 dyne contains an intricate system of rotating mirrors and circuitry to measure  
10 time elapsed, there are new solid state LiDARs that are rapidly evolving to  
11 match the needs of the autonomous driving industry. Most LiDARs have a  
12 usable range of about 100 m.



13  
14 **A typical autonomous car** This is a picture of MIT's entry named Talos to  
15 the DARPA Urban Challenge ([https://en.wikipedia.org/wiki/DARPA\\_Grand\\_Challenge\\_\(2007\)](https://en.wikipedia.org/wiki/DARPA_Grand_Challenge_(2007)))  
16 which was a competition where teams had to traverse a 60 mile urban route  
17 within 6 hours, while obeying traffic laws, understanding incoming vehicles  
18 etc. Successful demonstrations by multiple teams led by (CMU, Stanford,  
19 Odin College, MIT, Penn and Cornell) in this competition jump-started the  
20 wave of autonomous driving. While the number of sensors necessary to drive  
21 well has come down (Tesla famously does not like to use LiDARs and rely ex-  
22 clusively on cameras and radars), the type of sensors and challenges associated  
23 with them remain essentially the same.

**i** Waymo's autonomous car





1

### 2 4.3.3 Back to sensor modeling

3 Let us go back to understanding our sensor model  $l(m_i | y_k, x_k)$  where  $m_i$   
 4 is a particular cell of the occupancy grid,  $y_k$  and  $x_k$  are the observations and  
 5 robot position/orientation at time  $k$ .

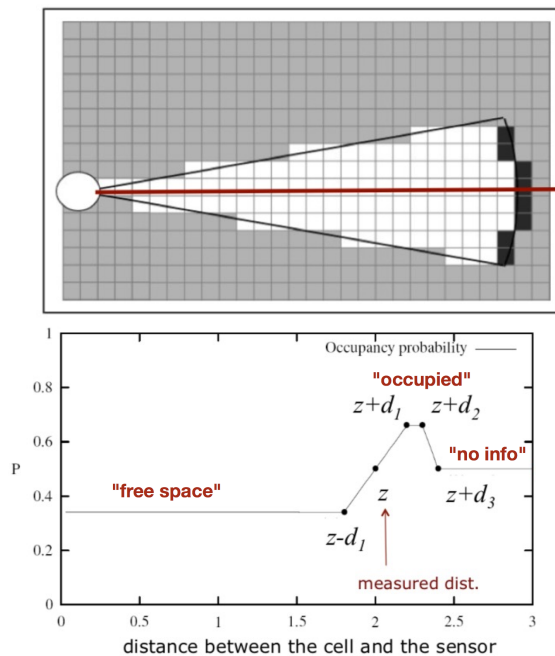


Figure 4.4: **Model for sonar data.** (Top) A sonar gives one real-valued reading corresponding to the distance measured along the red axis. (Bottom) if we travel along the optical axis, the occupancy probability  $P(m_i | y_k = z, x_k)$  can be modeled as a spike around the measured value  $z$ . It is very important to remember that range sensors such as sonar gives us three kinds of information about this ray: (i) all parts of the environment up to  $\approx z$  are *unoccupied* (otherwise we would not have recorded  $z$ ), (ii) there is some object at  $z$  which resulted in the return, (iii) but we do not know anything about what is behind  $z$ . So incorporating a measurement  $y_k$  from a sonar/radar/lidar involves not just updating the cell which corresponds to the return, but also updating the occupancy probabilities of every grid call along the axis.

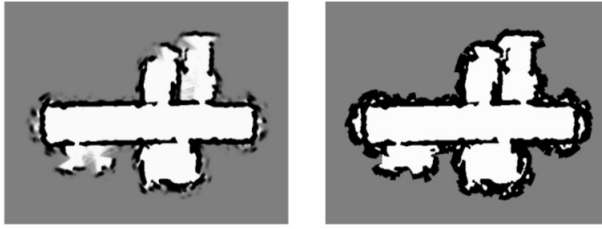
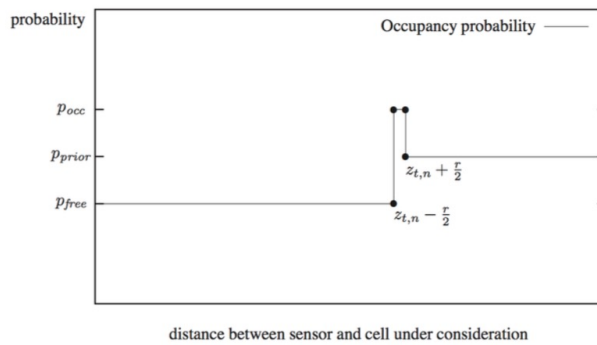
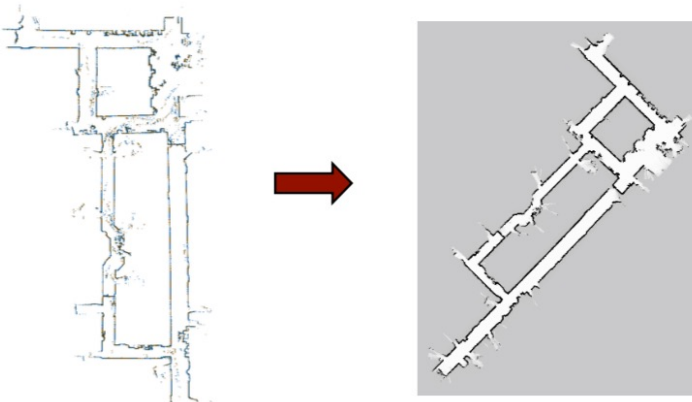


Figure 4.5: (Left) A typical occupancy grid created using a sonar sensor by updating the log-odds-ratio  $l(m_i | x_{1:k}, y_{1:k})$  for all cells  $i$  for multiple time-steps  $k$ . At the end of the map building process, if  $l(m_i | x_{1:k}, y_{1:k}) > 0$  for a particular cell, we set its occupancy to 1 and to zero otherwise, to get the maximum-likelihood estimate of the occupancy grid on the right.

- 1 **LiDAR model** When we say that a LiDAR is a more accurate sensor than
- 2 the sonar, what we really mean is that the sensor model  $P(m_i | y_k, x_k)$  looks
- 3 as follows.

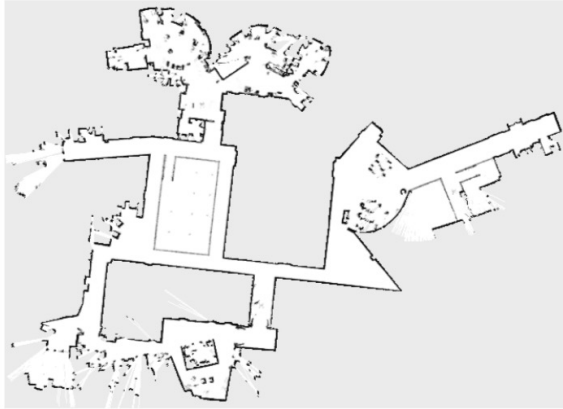


- 4
- 5 As a result, we can create high-resolution occupancy grids using a LiDAR.



6

## Example: MIT CSAIL 3<sup>rd</sup> Floor



❓ How will you solve the localization problem given the map? In other words, if we know the occupancy grid of a building as estimated in a prior run, and we now want to find the position/orientation of the robot traveling in this building, how should we use these sensors?

### 4.4 3D occupancy grids

Two-dimensional occupancy grids are a fine representation for toy problems but they run into some obvious problems: since the occupancy grid is a “top view” of the world, we cannot represent non-trivial objects in it correctly (a large tree with a thin trunk eats up all the free space). We often desire a fundamentally three-dimensional representation of the physical world.

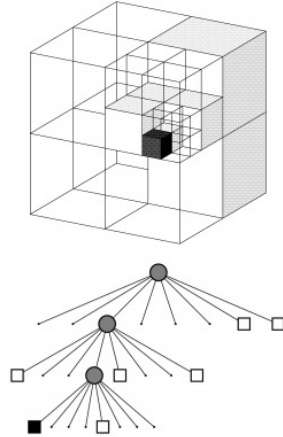


We could simply create cells in 3D space and our method for occupancy grid would work but this is no longer computationally cheap. For instance, if we want to build a map of Levine Hall (say  $50\text{ m} \times 50\text{ m}$  area and height of  $25\text{ m}$ ), a 3D grid map with a resolution of  $5\text{ cm} \times 5\text{ cm} \times \text{cm}$  would have about 500 million cells (if we store a float in each cell this map will require about 2 GB memory). It would be cumbersome to carry around so many cells and update their probabilities after each sensor reading (a Velodyne gives data at about 30 Hz). More importantly, observe that most of the volume inside Levine is free space (inside of offices, inner courtyard etc.) so we do not really need fine resolution in those regions.

**Octrees** We would ideally have an occupancy grid whose resolution adapts with the kind of objects that are detected by the sensors. If nearby cells are empty we want to collapse them together to save on memory and computation,



1 on the other hand, if nearby cells are all occupied, we want to *refine*  
 2 the resolution in that area so has to more accurately discern the shape of the  
 3 underlying objects. Octrees are an efficient representation for 3D volumes.



4

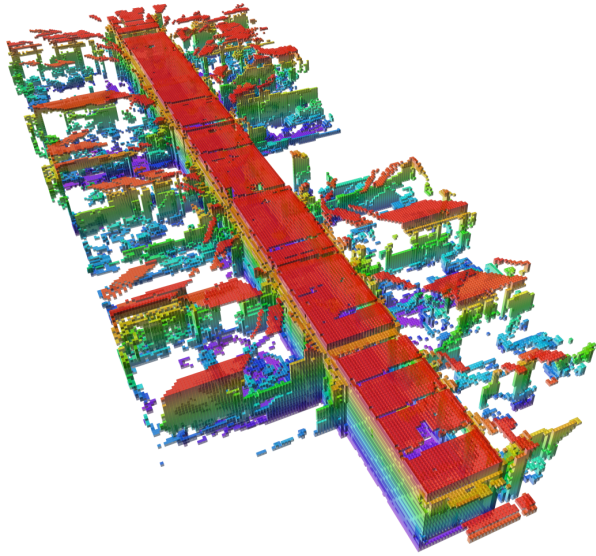
5 An octree is a hierarchical data structure that recursively sub-divides the 3D  
 6 space into octants and allocates volumes as needed for a particular data point  
 7 observed by a range sensor. It is analogous to a kd-tree. Imagine if the entire  
 8 space in the above picture were empty (the tree only has a root node), and we  
 9 receive a reading corresponding to the dark shaded region. An octree would  
 10 sub-divide the space starting from the root (each node in the tree populates  
 11 is the parent of its eight child octants) recursively until some pre-determined  
 12 minimum resolution is reached. This leaf node is grid cell; notice how different  
 13 cells in the octree have different resolutions. Occupancy probabilities of each  
 14 leaf node are updated using the same formula as that of (4.23). A key point  
 15 here is that octrees are designed for accurate sensors (LiDARs) where there is  
 16 not much noise in the observations returned by the sensor (and thereby we do  
 17 not refine unnecessary parts of the space)

18 Octrees are very efficient at storing large map, I expect you can store the  
 19 entire campus of Penn in about a gigabyte. Ray tracing (following all the cells  
 20  $m_i$  in tree along the axis of the sensor in Figure 4.4) is harder in this case but  
 21 there are efficient algorithms devised for this purpose. An example OctoMap  
 22 (an occupancy map created using an Octree) of a building on the campus of  
 23 the University of Freiburg is shown below.

**i** You can find LiDAR maps of the entire United States (taken from a plane) at <https://www.usgs.gov/core-science-systems/ngp/3dep>

## 24 4.5 Local Map

25 In this chapter, we primarily discussed occupancy grids of static environments  
 26 as the robot moves around in the environment. The purpose of doing so is lo-  
 27 calization, namely, finding the pose of the robot by comparing the observations  
 28 of the sensors with the map (think of the particle filter localization example  
 29 in Chapter 3). In typical problems, we often maintain two kinds of maps, (i)  
 30 a large occupancy grid for localization (say as big as city), and (ii) another  
 31 smaller map, called the local map, that is used to maintain the locations of



1 objects (typically objects that can move) in the vicinity of the robot, say a 100  
2 m  $\times$  100 m area.

3 The local map is used for planning and control purposes, e.g., to check if  
4 the planned trajectory of the robot does not collide with any known obstacles.  
5 See an example of the local map at the 1:42 min mark at  
6 <https://www.youtube.com/watch?v=2va15BE-7IQ>. Some people also call the  
7 local map a “cost map” because occupied cells in the local map indicate a  
8 high collision cost of moving through that cell. The local map is typically  
9 constructed in the body frame and evolves as the robot moves around (objects  
10 appear in the front of the robot and are spawned in the local map and disappear  
11 from the map at the back as the robot moves forward).

You should think of the map (and especially the local map) as the filtering estimate of the locations of various objects in the vicinity of the robot computed on the basis of multiple observations received from the robot’s sensors.

## 12 4.6 Discussion

13 Occupancy grids are a very popular approach to represent the environment  
14 given the poses of the robot as it travels in this environment. We can also  
15 use occupancy grids to localize the robot in a future run (which is usu-  
16 ally the purpose of creating them). Each cell in an occupancy grid stores  
17 the posterior probability of the cell being occupied on the basis of mul-  
18 tiple observations  $\{y_1, \dots, y_k\}$  from respective poses  $\{x_1, \dots, x_k\}$ . This  
19 is a very efficient representation of the 3D world around us with the one  
20 caveat that each cell is updated independently of the others. But since one



Figure 4.6: The output of perception modules for a typical autonomous vehicle (taken from <https://www.youtube.com/watch?v=tiwVMrTLUWg>). The global occupancy grid is shown in gray (see the sides of the road). The local map is not shown in this picture but you can imagine that it has occupied voxels at all places where there are vehicles (purple boxes) and other stationary objects such as traffic light, nearby buildings etc. Typically, if we know that so and so voxel corresponds to a vehicle, we run an Extended Kalman Filter for that particular vehicle to estimate the voxels in the local map that it is likely to be in, in the next time-instant. The local map is a highly dynamic data structure that is rich in information necessary for planning trajectories of the robot.

- 1 gets a large amount of data from typical range sensors (a 64 beam Velodyne
- 2 (<https://velodynelidar.com/products/hdl-64e>) returns about a 2 million
- 3 points/sec and cheaper versions of this sensor will cost about \$100), this caveat
- 4 does not hurt us much in practice. You can watch this talk
- 5 ([https://www.youtube.com/watch?v=V8JMwE\\_L5s0](https://www.youtube.com/watch?v=V8JMwE_L5s0)) by the head of Uber's
- 6 autonomous driving group to get more perspective about localization and
- 7 mapping.

# Chapter 5

## Dynamic Programming

### Reading

1. (Thrun) Chapter 15
2. (Sutton & Barto) Chapters 3–4
3. Optional: (Bertsekas) Chapter 1 and 4

This is the beginning of Module 2, this module is about “how to act”. The first module was about “how to sense”. The prototypical problem in the first module was how to assimilate the information gathered by all the sensors into some representation of the world. In the next few lectures, we will assume that this representation is good, that it is accurate in terms of its geometry (small variance of the occupancy grid) and in terms of its information (small innovation in the Kalman filter etc.). Let us also assume that it has all the necessary semantics, e.g., objects are labeled as cars, buses, pedestrians etc (we will talk about how to do this in Module 4).

The prototypical problem investigated in the next few chapters is how to move around in this world, or affect the state of this world to achieve a desired outcome, e.g., drive a car from some place A to another place B.

**Our philosophy about notation** Material on Dynamic Programming and Reinforcement Learning (RL), which we will cover in the following chapters, contains a lot of tiny details (much more than other areas in robotics/machine learning). These details are usually glossed over in most treatments. In the interest of simplicity, other courses or most research papers these days, develop an imprecise notation and terminology to focus on the problem. However, these details of RL matter enormously when you try to apply these techniques to real-world problems. Not knowing all the details or using imprecise terminology to think about RL is unlikely to make us good at real-world applications.

For this reason, the notation and the treatment in this chapter, and the following ones, will be a bit pedantic. We will see complicated notation and

1 terminology for quantities, e.g., the value function, that you might see being  
 2 written very succinctly in other places. We will mostly follow the notation  
 3 of Dmitri Bertsekas' book on "Reinforcement Learning and Optimal Control"  
 4 (<http://www.mit.edu/~dimitrib/RLbook.html>). You will get used to the extra  
 5 notation and it will become second nature once you become more familiar  
 6 with the concepts.

## 7 5.1 Formulating the optimal control problem

8 Let us denote the state of a robot (and the world) by  $x_k \in X \subset \mathbb{R}^n$  at the  $k^{\text{th}}$   
 9 timestep. We can change this state using a control input  $u_k \in U \subset \mathbb{R}^p$  and  
 10 this change is written as

$$x_{k+1} = f_k(x_k, u_k) \quad (5.1)$$

11 for  $k = 0, 1, \dots, T - 1$  starting from some initial given state  $x_0$ . This is  
 12 deterministic nonlinear dynamical system (no noise  $\epsilon$  in the dynamics). We  
 13 will let the dynamics  $f_k$  also be a function of time  $k$ . The time  $T$  is some  
 14 time-horizon up to which we care about running the system. The state-space is  
 15  $X$  (which we will assume does not change with time  $k$ ) and the control-space  
 16 is  $U$ .

17 Recall, that we can safely assume that the system is Markov. The reason  
 18 for it is as follows. If it is not, and say if  $x_{k+1}$  depends upon both  $x_k$  and  
 19 the previous step  $x_{k-1}$ , then we can expand the state-space to write a new  
 20 dynamics in the expanded state-space. We will follow a similar program as that  
 21 of Module 1: we first describe very general algorithms (dynamic programming)  
 22 for general systems (Markov Decision Processes), then specialize our methods  
 23 to a restricted class of systems that are useful in practice (linear dynamical  
 24 systems) and then finally discuss a very general class of systems again with  
 25 more sophisticated algorithms (motion-planning).

The central question in this chapter is how to pick a control  $u_k$ . We want to pick controls that lead to desirable trajectories of the system, e.g., results in a parallel-parked car at time  $T$  and does not collide against any other object for all times  $k \in \{1, 2, \dots, T\}$ . We may also want to minimize some chosen quantity, e.g., when you walk to School, you find a trajectory that avoids a certain street with a steep uphill.

**Finite, discrete state and control-space** In this chapter we will only consider problems with finitely-many states and controls, we will assume that the state-space  $X$  and the control-space  $U$  are finite, discrete sets.

26 **Run-time cost and terminal cost** We will take a very general view of the  
 27 above problem and formalize it as follows. Consider a cost function

$$q_k(x_k, u_k) \in \mathbb{R}$$

1 which gives a scalar real-valued output for every pair  $(x_k, u_k)$ . This models  
 2 the fact that you do not want to walk more than you need to get to School,  
 3 i.e., we would like to minimize  $q_k$ . You also want to make sure the trajectory  
 4 actually reaches the lecture venue, we write this down as another cost  $q_f(x_T)$ .  
 5 We want to pick control inputs  $(u_0, u_1, \dots, u_{T-1})$  such that

$$J(x_0; u_0, u_1, \dots, u_{T-1}) = q_f(x_T) + \sum_{k=0}^{T-1} q_k(x_k, u_k) \quad (5.2)$$

6 is minimized. The cost  $q_f(x_T)$  is called the terminal cost, it is high if  $x_T$  is  
 7 not the lecture room and small otherwise. The cost  $q_k$  is called the run-time  
 8 cost, it is high for instance if you have to use large control inputs, e.g.,  $x_k$  is a  
 9 climb.

**The optimal control problem** Given a system  $x_{k+1} = f_k(x_k, u_k)$ , we want to find control *sequences* that minimize the total cost  $J$  above, i.e., we want to solve

$$J^*(x_0) = \min_{u_k \in U, k=0, \dots, T-1} J(x_0; u_0, \dots, u_{T-1}) \quad (5.3)$$

It is important to realize that the function  $J(x_0; u_0, \dots, u_{T-1})$  depends upon an entire sequence of control inputs and we need to find them all to find the optimal cost  $J^*(x_0)$  of, say reaching the School from your home  $x_0$ .

## 10 5.2 Dijkstra's algorithm

11 If the state-space  $X$  and control-space  $U$  are discrete and finite sets, we can  
 12 solve (5.3) as a shortest path problem using very fast algorithms. Consider  
 13 the following picture. This is what would be called a transition graph for a  
 14 deterministic finite-state dynamics.

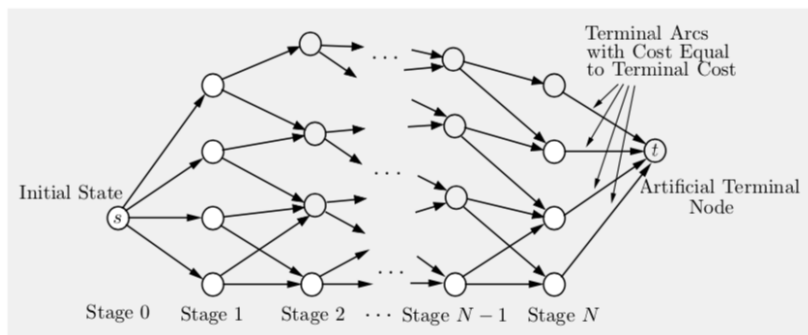


Figure 5.1: Transition graph for Dijkstra's algorithm

15 The graph has one source node  $x_0$ . Each node in the graph is  $x_k$ , each edge

1 depicts taking a certain control  $u_k$ . Depending on which control we picks, we  
 2 move to some other node  $x_{k+1}$  given by the dynamics  $f(x_k, u_k)$ . Note that  
 3 this is *not* a transition like that of a Markov chain, everything is deterministic  
 4 in this graph. On each edge we write down the cost

$$\text{cost}(x_k, x_{k+1}) := q_k(x_k, u_k)$$

5 where  $x_{k+1} = f_k(x_k, u_k)$  and “close” the graph with a dummy terminal node  
 6 with the cost  $q_f(x_T)$  on every edge leading to an artificial terminal node (sink).

7 Minimizing the cost in (5.3) is now the same as finding the shortest path in  
 8 this graph from the source to the sink. The algorithm to do so is quite simple  
 9 and is called Dijkstra’s algorithm after Edsger Dijkstra who used it around  
 10 1956 as a test program for a new computer named ARMAC ([http://www-](http://www-set.win.tue.nl/UnsungHeroes/machines/armac.html)  
 11 [set.win.tue.nl/UnsungHeroes/machines/armac.html](http://www-set.win.tue.nl/UnsungHeroes/machines/armac.html)).

12 1. Let  $Q$  be the set of nodes that are currently unvisited; all nodes in the  
 13 graph are added to it at the beginning.  $S$  is an empty set. An array called  
 14  $\text{dist}$  maintains the distance of every node in the graph from the source  
 15 node  $x_0$ . Initialize  $\text{dist}(x_0) = 0$  and  $\text{dist} = \infty$  for all other nodes.

16 2. At each step, if  $Q$  is not empty, pop a node  $v \in Q$  such that  $v \notin S$   
 17 with the smallest  $\text{dist}(v)$ . Add  $v$  to  $S$ . Update the  $\text{dist}$  of all nodes  $u$   
 18 connected to  $v$ . For each  $u$ , if

$$\text{dist}(u) > \text{dist}(v) + \text{cost}(u, v)$$

19 update the distance of  $u$  to be  $\text{dist}(v) + \text{cost}(u, v)$ . If the above condition  
 20 is not true do nothing.

21 The algorithm terminates when the set  $Q$  is empty.

22 You might know that there are many other variants of Dijkstra’s algorithm,  
 23 e.g., the  $A^*$  algorithm that are quicker to find shortest paths. We will look at  
 24 some of these in the next chapter.

The quantity  $\text{dist}$  is quite special: observe that after Dijkstra’s algorithm finishes running and the set  $Q$  is empty, the  $\text{dist}$  function gives the optimal cost to go from each node in the graph to the sink node. We wanted to only find the cost to go from source  $x_0$  to the sink node but ended up computing the cost from every node in the graph to the sink.

❓ Shortest path algorithms do not work if there are cycles in the graph because the shortest path is not unique. Are there cycles in the above graph?

❓ What should one do if the state/control space is not finite? Can we still use Dijkstra’s algorithm?

## 25 5.2.1 Dijkstra’s algorithm in the backwards direction

26 We can run Dijkstra’s algorithm in the backwards direction to get the same  
 27 answer as well. The sets  $Q$  and  $S$  are initialized as before. In this case we  
 28 will let  $\text{dist}(v)$  denote the distance of a node  $v$  to the sink node. The algorithm  
 29 proceeds in the same fashion, it pops a node  $v \in Q, v \notin S$  and updates the  
 30  $\text{dist}$  of all nodes  $u$  connected to  $v$ . For each  $u$ , if

$$\text{dist}(u) > \text{dist}(v) + \text{cost}(u, v)$$

1 then we update  $\text{dist}(u)$  to be the right-hand side of this inequality. Running Di-  
 2 jkstra's algorithm in reverse (from sink to the source) is completely equivalent  
 3 to running it in the forward direction (from source to the sink).

### 4 5.3 Principle of Dynamic Programming

The principle of dynamic programming is a formalization of the idea behind Dijkstra's algorithm. It was discovered by Richard Bellman in the 1940s. The idea behind dynamic programming is quite intuitive: it says that the remainder of an optimal trajectory is optimal.

5 We can prove this as follows. Suppose that we find the optimal control  
 6 sequence  $(u_0^*, u_1^*, \dots, u_{T-1}^*)$  for the problem in (5.3). Our system is de-  
 7 terministic, so this control sequence results in a *unique* sequence of states  
 8  $(x_0, x_1^*, \dots, x_T^*)$ . Each successive state is given by  $x_{k+1}^* = f_k(x_k^*, u_k^*)$  with  
 9  $x_0^* = x_0$ . The principle of optimality, or the principle of dynamic program-  
 10 ming, states that if one starts from a state  $x_k^*$  at time  $k$  and wishes to minimize  
 11 the "cost-to-go"

$$q_f(x_T) + q_k(x_k^*, u_k) + \sum_{i=k+1}^{T-1} q_i(x_i, u_i)$$

12 over the (now assumed unknown) sequence of controls  $(u_k, u_{k+1}, \dots, u_{T-1})$ ,  
 13 then the optimal control sequence for this truncated problem is exactly  $(u_k^*, \dots, u_{T-1}^*)$ .

14 The proof of the above assertion is an easy case of proof by contradic-  
 15 tion: if the truncated sequence were not optimal starting from  $x_k^*$  there ex-  
 16 ists some other optimal sequence of controls for the truncated problem, say  
 17  $(v_k^*, \dots, v_{T-1}^*)$ . If so, the solution of the original problem where one takes  
 18 controls  $v_k^*$  from this new sequence for time-steps  $k, k+1, \dots, T-1$  would  
 19 have a lower cost. Hence the original sequence of controls would not have  
 20 been optimal.

**Principle of dynamic programming.** The essence of dynamic program-

❗ If Dijkstra's algorithm (forwards or backwards) is run on a graph with  $n$  vertices and  $m$  edges, its computational complexity is  $\mathcal{O}(m + n \log n)$  if we use a priority queue to find the node  $v \in Q, v \notin S$  with the smallest  $\text{dist}$ . The number of edges in the transition graph in Figure 5.1 is  $m = \mathcal{O}(T|X|)$ .



ming is to solve the larger, original problem by sequentially solving the truncated sub-problems. At each iteration, Dijkstra's algorithm constructs the functions

$$J_T^*(x_T), J_{T-1}^*(x_{T-1}), \dots, J_0^*(x_0)$$

starting from  $J_T^*$  and proceeding backwards to  $J_{T-1}^*, J_{T-2}^* \dots$ . The function  $J_{T-k}^*(v)$  is just the array  $\text{dist}(v)$  at iteration  $k$  of the *backwards* implementation of Dijkstra's algorithm. Mathematically, dynamic programming looks as follows.

1. Initialize  $J_T^*(x) = q_f(x)$  for all  $x \in X$ .
2. For iteration  $k = T - 1, \dots, 0$ , set

$$J_k^*(x) = \min_{u_k \in U} \{q_k(x, u_k) + J_{k+1}^*(f_k(x, u_k))\} \quad (5.4)$$

for all  $x \in X$ .

After running the above algorithm we have the optimal cost-to-go  $J_0^*(x)$  for each state  $x \in X$ , in particular, we have the cost-to-go for the initial state  $J_0^*(x_0)$ . If we remember the minimizer  $u_k^*$  in (5.4) while running the algorithm, we also have the optimal sequence  $(u_0^*, u_1^*, \dots, u_{T-1}^*)$ . The function  $J_0^*(x)$  (often shortened to simply  $J^*(x)$ ) is the optimal cost-to-go from the state  $x \in X$ .

Again, we really only wanted to calculate  $J_0^*(x_0)$  but had to do all this extra work of computing  $J_k^*$  for all the states.

**Curse of dimensionality** What is the complexity of running dynamic programming? The cost of the minimization over  $U$  is  $\mathcal{O}(|U|)$ , it is a bunch of comparisons between floats. The number of operations at each iteration for setting the values  $J_k^*(x)$  for all  $x \in X$  is  $|X|$ . So the total complexity is  $\mathcal{O}(T |X| |U|)$ .

The terms  $|X|$  and  $|U|$  are often the hurdle in implementing dynamic programming or any variant of it. Think of the grid-world in Problem 1 in HW 1, it had  $200 \times 200$  cells which amounts to  $|X| = 40,000$ . This may seem a reasonable number but it explodes quickly as the dimensionality of the state-space grows. For a robot manipulator with six degrees-of-freedom, if we discretize each joint angle into 5 degree cells, the number of states is  $|X| \approx 140$  billion. The number of states  $|X|$  is exponential in the dimensionality of the state-space and dynamic programming quickly becomes prohibitive beyond 4 dimensions or so. Bellman called this the *curse of dimensionality*.

**Cost of dynamic programming is linear in the time-horizon** Notice a very important difference between (5.4)

$$J_k^*(x) = \min_{u_k \in U} \{q_k(x, u_k) + J_{k+1}^*(f_k(x, u_k))\}$$

1 for iterations  $i = T - 1, \dots, 0$  and (5.3)

$$J^*(x_0) = \min_{u_k \in U, k=0, \dots, T-1} J(x_0; u_0, \dots, u_{T-1}).$$

2 The latter has a minimization over a sequence of controls  $(u_0, u_1, \dots, u_{T-1})$   
 3 while the former has a minimization over only the control at time  $k$ ,  $u_k$  over  $T$   
 4 iterations. The former is much much easier to solve because it is a sequence of  
 5  $\mathcal{O}(T)$  smaller optimization problems: it is really easy to compute  $\min_{u_k \in U}$   
 6 for each state  $x$  separately than to solve the gigantic minimization problem  
 7 in (5.3) because in the latter case, the variable of optimization is the entire  
 8 control trajectory and has size  $|U|^T$ .

9 **Dynamic programming and Viterbi's algorithm** We have seen the princi-  
 10 ple of dynamic programming in action before in Viterbi's algorithm in Chapter  
 11 2. The transition graph in Figure 5.1 is the same as the Trellis graph for  
 12 Viterbi's algorithm, the run-time cost was

$$q_k(x_k, u_k) := -\log P(Y_k | X_k) - \log P(X_{k+1} | X_k)$$

13 and instead of a terminal cost  $q_f$ , we had an initial cost  $-\log P(X_1)$ . Viterbi's  
 14 algorithm computed the most likely path given observations of the HMM, i.e.,  
 15 the path  $(X_1, \dots, X_T)$  that maximizes the probability  $P(X_1, \dots, X_T | Y_1, \dots, Y_T)$   
 16 is simply the solution of dynamic programming for the Trellis graph.

### 17 5.3.1 Q-factor

18 The quantity

$$Q_k^*(x, u) := q_k(x, u) + J_{k+1}^*(f_k(x, u))$$

19 is called the Q-factor. It is simply the expression that is minimized in the  
 20 right-hand side of (5.4) and denotes the cost-to-go if control  $u$  was picked at  
 21 state  $x$  (corresponding to cost  $q_k(x, u)$ ) and the the optimal control trajectory  
 22 was followed after that (corresponding to cost  $J_k^*(f_k(x, u))$ ) from state  $x' =$   
 23  $f_k(x, u)$ . This nomenclature was introduced by Watkins in his thesis.

24 Q-factors and the cost-to-go are equivalent ways of thinking about dynamic  
 25 programming. Given the Q-factor, we can obtain the cost-to-go  $J_k^*$  as

$$J_k^*(x) = \min_{u_k \in U} Q_k^*(x, u_k). \quad (5.5)$$

26 which is precisely the dynamic programming update (by definition) in (5.4).  
 27 We can also write dynamic programming completely in terms of Q-factors as  
 28 follows.

**Dynamic programming written in terms of the Q-factor**

🔗 The principle of dynamic programming gives us a way to solve an optimization problem (5.3) over a really large space (the space of all control trajectories) using a linear in time-horizon number of optimization problems (5.4). Can we split any optimization problem in sub-problems like this?

🔗 How should one modify dynamic programming if we have a non-additive cost, e.g., the runtime cost at time  $k$  given by  $q_k$  is a function of both  $x_k$  and  $x_{k-1}$ ?

1. Initialize  $Q_T^*(x, u) = q_f(x)$  for all  $x \in X$  and all  $u \in U$ .

2. For iteration  $k = T - 1, \dots, 0$ , set

$$Q_k^*(x, u) = q_k(x, u) + \min_{u' \in U} Q_{k+1}^*(f_k(x, u), u'). \quad (5.6)$$

for all  $x \in X$  and all  $u \in U$ .

As yet, it may seem unnecessary to think of the Q-factor (which is a larger array with  $|X| \times |U|$  entries) instead of the cost-to-go (which only has  $|X|$  entries in the array).

**Value function** The following terminology is commonly used in the literature

value function  $\equiv$  cost-to-go  $J^*(x)$

action-value function  $\equiv$  Q-factor  $Q^*(x, u)$ .

Since the two functions are equivalent, we will call both as “value functions”. The difference will be clear from context.

## 5.4 Stochastic dynamic programming: Value Iteration

Let us now see how dynamic programming looks for a Markov Decision Process (MDP). As we saw in Chapter 3, we can think of MDPs as stochastic dynamical systems denoted by

$$x_{k+1} = f_k(x_k, u_k) + \epsilon_k; \quad x_0 \text{ is given.}$$

We will assume that we know the statistics of the noise  $\epsilon_k$  at each time-step (say it is a Gaussian). Stochastic dynamical systems are very different from deterministic dynamical systems, given the same sequence of controls  $(u_0, \dots, u_{T-1})$ , we may get different state trajectories  $(x_0, x_1, \dots, x_T)$  depending upon the realization of noise  $(\epsilon_0, \dots, \epsilon_{T-1})$ . How should we find a good control trajectory then? One idea is to modify (5.3) to minimize the expected value of the cost over all possible state-trajectories

$$J(x_0; u_0, \dots, u_{T-1}) = \mathbf{E}_{(\epsilon_0, \dots, \epsilon_{T-1})} \left[ q_f(x_T) + \sum_{k=0}^{T-1} q_k(x_k, u_k) \right] \quad (5.7)$$

Suppose we minimized the above expectation and obtained the value function  $J^*(x_0)$  and the optimal control trajectory  $(u_0^*, \dots, u_{T-1}^*)$ . As the robot starts executing this trajectory, the realized versions of the noise  $\epsilon_k$  might differ a lot from their expected value, and the robot may find itself in very different states  $x_k$  than the average-case states considered in (5.10).

**i** Draw the picture of a one-dimensional stochastic dynamical system (random walk on a line) and see that the realized trajectory of the system can be very different from the average trajectory.

1 **Feedback controls** The concept of feedback control is a powerful way to  
 2 resolve this issue. Instead of seeking  $u_k^* \in U$  as the solutions of (5.10), we  
 3 instead seek a *function*

$$u_k(x) : X \mapsto U \quad (5.8)$$

4 that maps the state-space  $X$  to a control  $U$ . Effectively, given a feedback  
 5 control  $u_k(x)$  the robot knows what control to apply at its current realized state  
 6  $x_k \in X$ , namely  $u_k(x_k)$ , even if the realized state  $x_k$  is very different from  
 7 the average-case state. Feedback controls are everywhere and are critical to  
 8 using controls in the real world. For instance, when you tune the shower faucet  
 9 to give you a comfortable water temperature, you are constantly estimating  
 10 the state (feedback using the temperature) and turning the faucet accordingly.  
 11 Doing this without feedback would leave you terribly cold or scalded. We will  
 12 denote the space of all feedback controls  $u_k(\cdot)$  that depend on the state  $x \in X$   
 13 by

$$u_k(\cdot) \in \mathcal{U}(X).$$

14 **Control policy** A sequence of feedback controls

$$\pi = (u_0(\cdot), u_1(\cdot), \dots, u_{T-1}(\cdot)). \quad (5.9)$$

15 is called a control policy. This is an object that we will talk about often. It is  
 16 important to remember that a control policy is set of controllers (usually feed-  
 17 back controls) that are executed at each time-step of a dynamic programming  
 18 problem.

The **stochastic optimal control problem** finds a sequence of feed-  
 back controls  $(u_0(\cdot), u_1(\cdot), \dots, u_T(\cdot))$  that minimizes

$$J(x_0; u_0(\cdot), \dots, u_{T-1}(\cdot)) = \mathbf{E}_{(\epsilon_0, \dots, \epsilon_{T-1})} \left[ q_f(x_T) + \sum_{k=0}^{T-1} q_k(x_k, u_k(x_k)) \right]$$

The value function is given by

$$J^*(x_0) = \min_{u_k(\cdot) \in \mathcal{U}(X), k=0, \dots, T-1} J(x_0; u_0(\cdot), \dots, u_{T-1}(\cdot)) \quad (5.10)$$

The optimal sequence of feedback controls (in short, the optimal control  
 trajectory) is the one that achieves this minimum.

19 Dijkstra's algorithm no longer works, as is, if the edges in the graph are  
 20 stochastic but we can use the principal of dynamic programming to write the  
 21 solution for the stochastic optimal control problem. The idea remains the same,  
 22 we compute a sequence of cost-to-go functions  $J_T^*(x), J_{T-1}^*(x), \dots, J_0^*(x)$ ,  
 23 and in particular  $J_0^*(x_0)$ , proceeding *backwards*.

**Finite-horizon dynamic programming for stochastic systems.**

❗ All this sounds very tricky and abstract but you will quickly get used to the idea of feedback control because it is quite natural. You can think of feedback control as being analogous to the innovation term in the Kalman filter  $K(y_k - C\mu_{k+1|k})$  which corrects the estimate  $\mu_{k+1|k}$  to get a new estimate  $\mu_{k+1|k+1}$  using the *current* observation  $y_k$ . Filtering would not work at all if the innovation term did not depend upon the actual observation  $y_k$  and only depended upon some average observation.

1. Initialize  $J_T^*(x) = q_f(x)$  for all  $x \in X$ .

2. For all times  $k = T - 1, \dots, 0$ , set

$$J_k^*(x) = \min_{u_k(\cdot) \in \mathcal{U}(X)} \left\{ q_k(x, u_k(x)) + \mathbb{E}_{\epsilon_k} [J_{k+1}^*(f_k(x, u_k(x)) + \epsilon_k)] \right\} \quad (5.11)$$

for all  $x \in X$ .

1 Just like (5.4), we solve a sub-problem for one time-instant at each iteration.  
2 But observe a few importance differences in (5.11) compared to (5.4).

- 3 1. There is an expectation over the noise  $\epsilon_k$  in the second term in the curly  
4 brackets. The second term in the curly brackets is the average of the  
5 cost-to-go of the truncated sub-problems from time  $k + 1, \dots, T$  over  
6 all possible starting states  $x' = f_k(x_k, u_k(x_k)) + \epsilon_k$ . This makes sense,  
7 after taking the control  $u_k(x_k)$ , we may find the robot at any of the  
8 possible states  $x' \in X$  depending upon different realizations of noise  
9  $\epsilon_k$  and the cost-to-go from  $x_k$  is therefore the average of the cost-to-  
10 go from each of those states (according to the principal of dynamic  
11 programming).  
12 2. The minimization in (5.11) is performed over a function

$$\mathcal{U}(X) \ni u_k(\cdot) : X \mapsto U.$$

13 Since our set of states and controls is finite, this involves finding a table  
14 of size  $|X| \times |U|$  for each iteration. In (5.4), we only had to search over  
15 a set of values  $u_k \in U$  of size  $|U|$ . At the end of dynamic programming,  
16 we have a sequence of feedback controls

$$(u_0^*(\cdot), u_1^*(\cdot), \dots, u_{T-1}^*(\cdot)).$$

17 Each feedback control  $u_k^*(x)$  tells us what control the robot should pick  
18 if it finds itself at a state  $x$  at time  $k$ .

- 19 3. If we know the dynamical system, not in its functional form  $x_{k+1} =$   
20  $f_k(x_k, u_k) + \epsilon_k$  but rather as a transition matrix  $P(x_{k+1} | x_k, u_k)$  (like  
21 we had in Chapter 2) then the expression in (5.11) simply becomes

$$J_k^*(x) = \min_{u_k(\cdot) \in \mathcal{U}(X)} \left\{ q_k(x, u_k(x)) + \mathbb{E}_{x' \sim P(\cdot | x_k, u_k(x_k))} [J_{k+1}^*(x')] \right\} \quad (5.12)$$

22 **Computational complexity** The form in (5.12) helps us understand the  
23 computational complexity, each sub-problem performs  $|X| \times |X| \times |U|$  amount  
24 of work and therefore the total complexity of stochastic dynamic programming  
25 is

$$\mathcal{O}(T|X|^2|U|).$$

🔗 Why should we only care about minimizing the average cost in the objective in (5.10)? Can you think of any other objective we may wish to use?

1 Naturally, the quadratic dependence on the size of the state-space is an even  
 2 bigger hurdle while implementing dynamic programming for stochastic sys-  
 3 tems.

#### 4 **5.4.1 Infinite-horizon problems**

5 In the previous section, we put a lot of importance on the horizon  $T$  for  
 6 dynamic programming. This is natural: if the horizon  $T$  changes, say you are  
 7 in a hurry to get to school, the optimal trajectory may take control inputs that  
 8 incur a lot of runtime cost simply to reach closer to the goal state (something  
 9 that keeps the terminal cost small). In most, real-world problems, it is not very  
 10 clear what value of  $T$  we should pick. We therefore formulate the dynamic  
 11 programming problem as something that also allows a trajectory of infinite  
 12 steps but also encourages the length of the trajectory to be small enough in  
 13 order to be meaningful. Such problems are called infinite-horizon problems  
 14 ( $T \rightarrow \infty$ ).

15 **Stationary dynamics and run-time cost** We think of infinite-horizon prob-  
 16 lems in the following way: at any time-step, the length of the trajectory  
 17 *remaining* for the robot to traverse is infinite. It helps in this case to solve a  
 18 restricted set of problems where the system dynamics and run-time cost do  
 19 not change as a function of time (they only change as a function of the state  
 20 and the control). We will set

$$q(x, u) \equiv q_k(x, u),$$

$$f(x, u) \equiv f_k(x, u)$$

21 for all  $x \in X$  and  $u \in U$ . Such a condition is called stationarity. If the system  
 22 is stochastic, we also require that the distribution of noise  $\epsilon_k$  does not change  
 23 as a function of time (it could change in (5.11) but we did not write it so). The  
 24 infinite-horizon setting is never quite satisfied in practice but it is a reasonable  
 25 formulation for problems that run for a long length of time.

26 **Infinite-horizon objective** The objective that we desire be minimized by an  
 27 infinite-horizon control policy

$$\pi = (u_0(\cdot), u_1(\cdot), \dots, u_T(\cdot), u_{T+1}(\cdot), \dots)$$

28 is defined in terms of an asymptotic limit

$$J(x_0; \pi) = \lim_{T \rightarrow \infty} \mathbb{E}_{(\epsilon_0, \dots, \epsilon_{T-1})} \left[ \sum_{k=0}^{T-1} \gamma^k q(x_k, u_k(x_k)) \right]. \quad (5.13)$$

29 and we again wish to solve for the optimal cost-to-go

$$J^*(x_0) = \operatorname{argmin}_{\pi} J^*(x_0; \pi). \quad (5.14)$$

30 Thus the infinite horizon costs of a policy is the limit of its finite horizon  
 31 costs as the horizon tends to infinity. Notice a few important differences when

1 compared to (5.7).

- 2 1. The objective is a limit, it is effectively the cost of the trajectory as it is
- 3 allowed to stretch for a larger and larger time-horizon.
- 4 2. There is no terminal cost in the objective function; this makes sense
- 5 because an explicit terminal state  $x_T$  does not exist anymore. In infinite-
- 6 horizon problems, you should think of the terminal cost as being incor-
- 7 porated inside the run-time cost  $q(x, u)$  itself, e.g., move the robot to
- 8 minimize the fuel used at *this* time instant but also move it in a way that
- 9 it reaches the goal *at some time in the future*.
- 10 3. **Discount factor**— Depending upon what controls we pick, the summa-
- 11 tion

$$\sum_{k=0}^T q(x_k, u_k(x_k))$$

12 can diverge to infinity as  $T \rightarrow \infty$  and thereby a meaningful solution to

13 the infinite-horizon problem may not exist. In order to avoid this, we

14 use a scalar

$$\gamma \in (0, 1)$$

15 known as the discount factor in the formulation. It puts more em-

16 phasis on costs incurred earlier in the trajectory than later ones and

17 thereby encourages the length of the trajectory to be small. Notice that

18  $\sum_{k=0}^{\infty} \alpha^k = 1/(1 - \alpha)$  if  $|\alpha| < 1$ , so if the cost  $|q(x_k, u_k(x_k))| < 1$ ,

19 then we know that the objective in (5.13) always converges.

20 **Stochastic shortest path problems** It is important to remember that the

21 discount factor is chosen by the user, no one prescribes it. There is also a

22 class of problems where we may choose  $\gamma = 1$  but in these cases, there should

23 exist some *essentially terminal* state in the state space where we can keep

24 taking a control such that the runtime cost  $q(x, u)$  is zero. Otherwise, the

25 objective will diverge. The goal region in the grid-world problem could be

26 an example of such state. Such problems are called stochastic shortest path

27 problems because the time-horizon is not *actually* infinite, we just do not

28 know how many time-steps it will take for the robot to go to the goal location.

29 Naturally, stochastic shortest path problems are a generalization of the shortest

30 path problem solved by Dijkstra's algorithm. The algorithms we discuss next

31 will work for such problems.

32 **Stationary policy** It seems a bit cumbersome to carry around an infinitely

33 long sequence of feedback controls in infinite-horizon problems. Since there

34 is an infinitely-long trajectory yet to be traveled *at any given time-step*, the

35 optimal control action that we take should only depend upon the current state.

36 This is indeed true mathematically. If  $J^*(x)$  is the optimal cost-to-go in the

37 infinite-horizon problem starting from a state  $x$ , using the principle of dynamic

38 programming, we should also have that we can split this cost as the best

39 one-step cost of the current state  $x$  added to the optimal cost-to-go from the

40 state  $f(x, u)$  realized after taking the optimal control  $u$ :

$$J^*(x) = \min_{u(x) \in \mathcal{U}(X)} \mathbb{E}_{\epsilon} [q(x, u(x)) + J^*(f(x, u(x))) + \epsilon]. \quad (5.15)$$

1 We will study this equation in depth soon. But if we find the minimum at  
2  $u^*(x)$  for this equation, then we can run the policy

$$\pi^* = (u^*(\cdot), u^*(\cdot), \dots, u^*(\cdot), \dots)$$

3 for the entire infinite horizon. Such a policy is called a stationary policy. Intu-  
4 itively, since the future optimization problem (tail of dynamic programming)  
5 from a given state  $x$  looks the same regardless of the time at which we start,  
6 optimal policies for the infinite-horizon problem can be found even inside the  
7 restricted class of policies where the feedback control does not change with  
8 time  $k$ .

9 We will almost exclusively deal with stationary policies in this course.

## 10 5.4.2 Dynamic programming for infinite-horizon problems

11 We wish to compute the optimal cost-to-go of starting from a state  $x$  and  
12 taking an infinitely long trajectory that minimizes the objective (5.13). We will  
13 exploit the equation in (5.15) and develop an iterative algorithm to compute  
14 the optimal cost-to-go  $J^*(x)$ .

**Value Iteration.** The algorithm proceeds iteratively to maintain a sequence of approximations

$$\forall x \in X, \quad J^{(0)}(x), J^{(1)}(x), J^{(2)}(x), \dots,$$

to the optimal value function  $J^*(x)$ . Such an algorithm is called “value iteration”.

1. Initialize  $J^{(0)}(x) = 0$  for all  $x \in X$ .
2. Update using the Bellman equation at each iteration, i.e., for  $i = 1, 2, \dots, N$ , set

$$J^{(i+1)}(x) = \min_{u \in U} \mathbb{E}_\epsilon \left[ q(x, u) + \gamma J^{(i)}(f(x, u) + \epsilon) \right]. \quad (5.16)$$

for all  $x \in X$  until the value function converges at all states, e.g.,

$$\forall x \in X, \quad |J^{(i)}(x) - J^{(i+1)}(x)| < \text{small tolerance.}$$

3. Compute the feedback control and the stationary policy  $\pi^* = (u^*(\cdot), \dots)$  corresponding to the value function estimate  $J^{(N)}$  as

$$u^*(x) = \operatorname{argmin}_{u \in U} \mathbb{E}_\epsilon \left[ q(x, u) + \gamma J^{(N)}(f(x, u) + \epsilon) \right] \quad (5.17)$$

for all  $x \in X$ .

❗ If the dynamics is given as a transition matrix, we can replace the expectation over noise  $\mathbb{E}_\epsilon$  as an expectation over the next state  $x' \sim \mathbb{P}(x' | x, u(x))$  in (5.16) to run value iteration. Everything else remains the same

15 Let us observe a few important things in the above sequence of updates.  
16 First, at each iteration, we are updating the values of all  $|X|$  states. This  
17 involves  $|X|^2|U|$  amount of work per iteration. How many such iterations  $N$





1

0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0.75	0
0	0	0	0	0	0	0	0.75	1	0	0
0	0	0	0	0	0	0	0	0	0	0

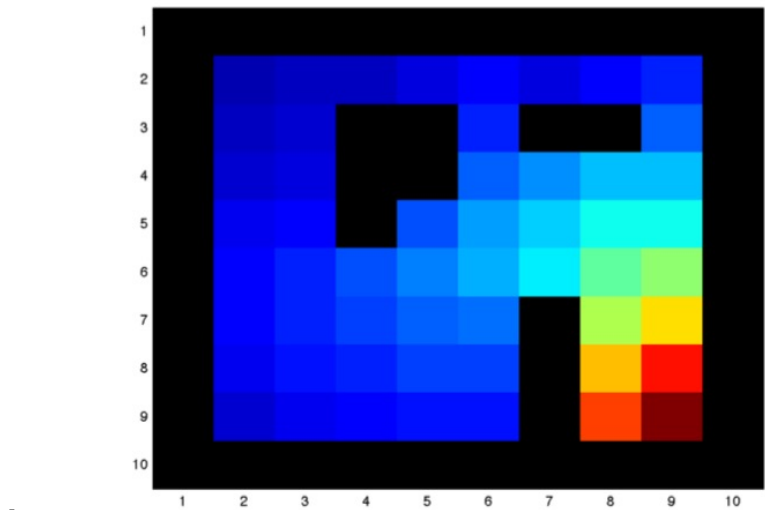
2

0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0.51	0
0	0	0	0	0	0	0	0	0.56	1.43	0
0	0	0	0	0	0	0	0	1.43	1.9	0
0	0	0	0	0	0	0	0	0	0	0

3

0	0	0	0	0	0	0	0	0	0	0
0	0.44	0.54	0.59	0.82	1.15	0.85	1.09	1.52	0	0
0	0.59	0.69	0	0	1.52	0	0	2.13	0	0
0	0.75	0.90	0	0	2.12	2.55	2.98	3.00	0	0
0	0.95	1.18	0	2.00	2.70	3.22	3.80	3.88	0	0
0	1.20	1.55	1.87	2.41	2.92	3.51	4.52	5.00	0	0
0	1.15	1.47	1.74	2.05	2.25	0	5.34	6.47	0	0
0	0.99	1.26	1.49	1.72	1.74	0	6.69	8.44	0	0
0	0.74	0.99	1.17	1.34	1.27	0	7.96	9.94	0	0
0	0	0	0	0	0	0	0	0	0	0

4 The final value function after 50 iterations looks as follows.



6 **5.4.4 Some theoretical results on value iteration**

7 We list down some very powerful theoretical results for value iteration. These  
 8 results are valid under a very general set of conditions and make value iteration

1 work for a large number of real-world problems; they are at the heart of all  
 2 modern algorithms. We will not derive them (it is easy but cumbersome) but  
 3 you should commit them to memory and try to understand them intuitively.

4 **Value iteration converges.** Given *any* initialization  $J^{(0)}(x)$  for all  $x \in X$ ,  
 5 the sequence of value iteration estimates  $J^{(i)}(x)$  converges to the optimal cost

$$\forall x \in X, \quad J^*(x) = \lim_{N \rightarrow \infty} J^{(N)}(x)$$

6 **The solution is unique.** The optimal cost-to-go  $J^*(x)$  of (5.14) satisfies the  
 7 Bellman equation

$$J^*(x) = \min_{u \in U} \mathbb{E}_\epsilon [q(x, u) + \gamma J^*(f(x, u) + \epsilon)].$$

8 The function  $J^*$  is also the *unique* solution of this equation. In other words, if  
 9 we find some other function  $J'(x)$  that satisfies the Bellman equation, we are  
 10 guaranteed that  $J'$  is indeed the optimal cost-to-go.

11 **Policy evaluation: Bellman equation for a particular policy.** Consider  
 12 a stationary policy  $\pi = (u(\cdot), u(\cdot), \dots)$ . The cost of executing this policy  
 13 starting from a state  $x$ , is  $J(x; \pi)$  from (5.13), also denoted by  $J^\pi(x)$  for short.  
 14 It satisfies the equation

$$J^\pi(x) = q(x, u(x)) + \gamma \mathbb{E}_\epsilon [J^\pi(f(x, u(x)) + \epsilon)] \quad (5.20)$$

15 and is also the unique solution of this equation. In other words, if we have a  
 16 policy in hand, and wish to find the cost-to-go of this policy, i.e., “evaluate the  
 17 policy” we can initialize  $J^{(0)}(x) = 0$  for all  $x \in X$  and perform the sequence  
 18 of iterative updates to this initialization

$$J^{(i+1)}(x) = q(x, u(x)) + \gamma \mathbb{E}_\epsilon [J^{(i)}(f(x, u(x)) + \epsilon)]. \quad (5.21)$$

19 As the number of updates goes to infinity, the iterate converges to  $J^\pi(x)$

$$\forall x \in X, \quad J^\pi(x) = \lim_{N \rightarrow \infty} J^{(N)}(x).$$

20 **Policy evaluation is equivalent to solving a linear system of equations.**  
 21 The min operation in (5.16) or (5.18) is particularly problematic in imple-  
 22 menting value iteration. As compared to it, observe that the corresponding  
 23 equation for policy equation (5.20) does not have min operation. This allows  
 24 us to write the updates in (5.21) as the solution of a linear system of equations.  
 25 Since we are in a finite state-space, we can write the cost-to-go as a large  
 26 vector

$$J^\pi := [J^\pi(x_1), J^\pi(x_2), \dots, J^\pi(x_n)]^\top$$

27 where  $n$  is the number of total states in the state-space. We create a similar  
 28 vector for the run-time cost term

$$q^u := [q(x_1, u(x_1)), q(x_2, u(x_2)), \dots, q(x_n, u(x_n))].$$

1 We know that the expectation over noise  $\epsilon$  is equivalent to an expectation over  
 2 the next state of the system, let us rewrite the dynamics part  $f(x, u(x)) + \epsilon$  in  
 3 terms of the Markov transition matrix

$$T_{x,x'} = \mathbb{P}(x' | x, u(x))$$

4 as

$$\gamma \mathbb{E}_{\epsilon} [J^{\pi}(f(x, u(x)) + \epsilon)] = \gamma \sum_{x'} T_{x,x'} J^{\pi}(x') = \gamma T J^{\pi}$$

5 to get a linear system

$$J^{\pi} = q^u + \gamma T J^{\pi} \quad (5.22)$$

6 which can be solved easily for  $J^{\pi} = (I - \gamma T)^{-1} q^u$  to get the cost-to-go of a  
 7 particular control policy  $\pi$ .

## 8 5.5 Stochastic dynamic programming: Policy It- 9 eration

10 Value iteration converges exponentially quickly, but asymptotically. The  
 11 number of states  $|X| = n$  is finite and so is the number of controls  $|U|$ . This  
 12 seems very funny, one would expect that we should be able to find the optimal  
 13 cost  $J^*(x)$  in finite time if the problem is finite, after all we need to find  $|X|$   
 14 numbers  $J^*(x_1), \dots, J^*(x_n)$ . This intuition is rightly placed indeed and in  
 15 this section, we will discuss an algorithm called policy iteration, which is a  
 16 more efficient version of value iteration.

17 The idea behind policy iteration is quite simple: given a stationary policy  
 18 for an infinite-horizon problem  $\pi = (u(\cdot), \dots, u(\cdot))$ , we can evaluate this  
 19 policy to obtain its cost-to-go  $J^{\pi}(x)$ . If we now set the feedback control to be

$$\tilde{u}(x) = \underset{u \in U}{\operatorname{argmin}} \mathbb{E}_{\epsilon} [q(x, u) + \gamma J^{\pi}(f(x, u) + \epsilon)], \quad (5.23)$$

20 i.e., we construct a *new control policy* that executes this new control  $\tilde{u}(\cdot)$  at  
 21 the first step and thereafter executes the old feedback control  $u(\cdot)$

$$\pi^{(1)} = (\tilde{u}(\cdot), u(\cdot), \dots),$$

22 then the cost-to-go of policy  $\pi^{(1)}$  is always better:

$$\forall x \in X, \quad J^{\pi^{(1)}}(x) \leq J^{\pi}(x).$$

23 We don't have to stop at one time-step, we can patch the old policy at the first  
 24 two time-steps to get

$$\pi^{(2)} = (\tilde{u}(\cdot), \tilde{u}(\cdot), \dots),$$

25 and have by the same logic

$$\forall x \in X, \quad J^{\pi^{(2)}}(x) \leq J^{\pi^{(1)}}(x) \leq J^{\pi}(x).$$

🔗 Why? It is simply because (5.23) is at least an improvement upon the feedback control  $u(\cdot)$ . The cost-to-go cannot improve only if the old feedback control  $u(\cdot)$  where optimal to begin with.

1 If we build a new stationary policy

$$\tilde{\pi} = (\tilde{u}(\cdot), \tilde{u}(\cdot), \tilde{u}(\cdot), \dots), \quad (5.24)$$

2 we similarly have

$$\forall x \in X, \quad J^{\tilde{\pi}}(x) \leq J^{\pi}(x).$$

3 This suggests an iterative way to compute the optimal stationary policy  $\pi^*$   
 4 starting from some initial stationary policy (i.e., implicitly a feedback con-  
 5 troller).

**Policy Iteration** The algorithm proceeds to maintain a sequence of stationary policies

$$\pi^{(k)} = (u^{(k)}(\cdot), u^{(k)}(\cdot), u^{(k)}(\cdot), \dots)$$

that converges to the optimal policy  $\pi^*$ .

Initialize  $u^{(0)}(x) = 0$  for all  $x \in X$ . This gives the initial stationary policy  $\pi^{(0)}$ . At each iteration  $k = 1, \dots$ , we do the following two things.

1. **Policy evaluation** Use multiple iterations of (5.21) to evaluate the old policy  $\pi^{(k)}$ . Initialize  $J^{(0)}(x) = 0$  for all  $x \in X$  and iterate upon

$$J^{(i+1)}(x) = q(x, u^{(i)}(x)) + \gamma \mathbb{E}_{\epsilon} \left[ J^{(i)}(f(x, u^{(i)}(x) + \epsilon)) \right]$$

for all  $x \in X$  until convergence. In practice, we always use the linear system of equations in (5.22) to solve for  $J^{\pi^{(k)}}$  directly.

2. **Policy improvement** Update the feedback controller using (5.23) to be

$$u^{(k+1)}(x) = \underset{u \in U}{\operatorname{argmin}}_{\epsilon} \left[ q(x, u) + \gamma J^{\pi^{(k)}}(f(x, u) + \epsilon) \right]$$

for all  $x \in X$  and compute the updated stationary policy

$$\pi^{(k+1)} = (u^{(k+1)}(\cdot), u^{(k+1)}(\cdot), \dots)$$

The algorithm terminates when the controller does not change *at any state*, i.e., when the following condition is satisfied

$$\forall x \in X, \quad u^{(k+1)}(x) = u^{(k)}(x).$$

6 Just like value iteration converges to the optimal value function, it can be  
 7 shown that policy iteration produces a sequence of improved policies

$$\forall x \in X, \quad J^{\pi^{(k+1)}}(x) \leq J^{\pi^{(k)}}(x)$$

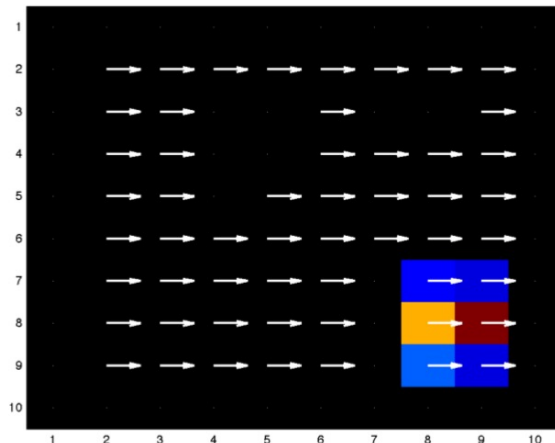
1 and converges to the optimal cost-to-go

$$\forall x \in X, \quad J^*(x) = \lim_{N \rightarrow \infty} J^{\pi(N)}(x).$$

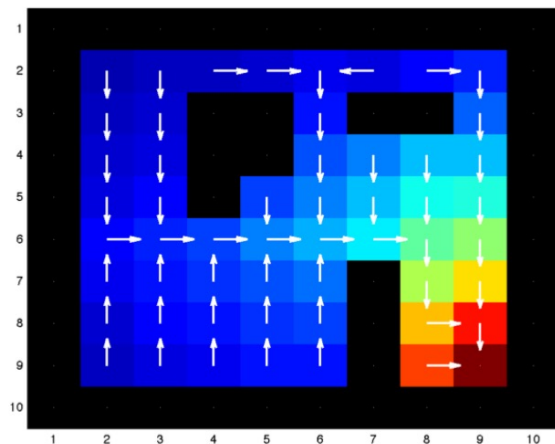
2 The key property of policy iteration is that we need a finite number of updates  
 3 to the policy to find the optimal policy. Notice that this does not mean always  
 4 mean that we are doing less work than value iteration in policy iteration. Ob-  
 5 serve that the policy evaluation step in the policy iteration algorithm performs  
 6 a number of Bellman equation updates. But typically, it is observed in practice  
 7 that policy iteration is much cheaper computationally than value iteration.

### 8 5.5.1 An example

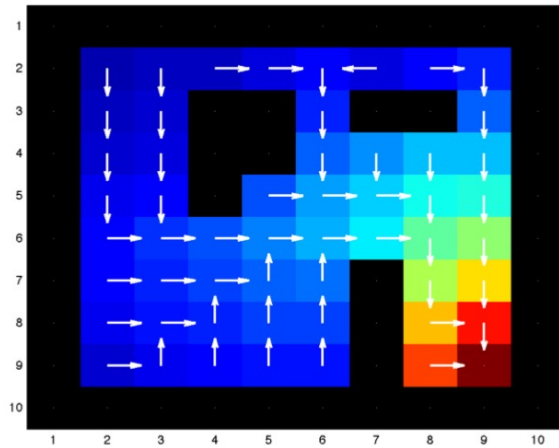
9 Let us go back to our example for value iteration. In this case, we will visualize  
 10 the controller  $u^{(k)}(x)$  at each cell  $x$  as arrows pointing to some other cell. The  
 11 cells are colored by the value function for that particular stationary policy.



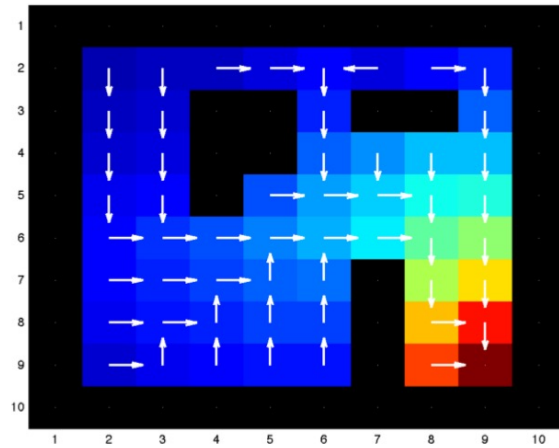
12



13



1



2

3 The evaluated value for the policy after 4 iterations is optimal, compare  
 4 this to the example for value iteration.

0	0	0	0	0	0	0	0	0	0
0	0.45	0.56	0.61	0.84	1.17	0.87	1.11	1.5411	0
0	0.61	0.71	0	0	1.54	0	0	2.16	0
0	0.78	0.93	0	0	2.16	2.59	3.02	3.03	0
0	0.98	1.21	0	2.03	2.74	3.26	3.84	3.91	0
0	1.23	1.58	1.90	2.44	2.95	3.54	4.56	5.03	0
0	1.18	1.50	1.78	2.09	2.28	0	5.38	6.51	0
0	1.02	1.29	1.52	1.76	1.77	0	6.74	8.49	0
0	0.76	1.02	1.20	1.37	1.30	0	8.01	10	0
0	0	0	0	0	0	0	0	0	0

5

# Chapter 6

## Linear Quadratic Regulator (LQR)

### Reading

1. <http://underactuated.csail.mit.edu/lqr.html>, Lecture 3-4 at <https://ocw.mit.edu/courses/aeronautics-and-astronautics/16-323-principles-of-optimal-control-spring-2008/lecture-notes>
2. Optional: Applied Optimal Control by Bryson & Ho, Chapter 4-5

This chapter is the analogue of Chapter 3 on Kalman filtering. Just like Chapter 2, the previous chapter gave us two algorithms, namely value iteration and policy iteration, to solve dynamic programming problems for a finite number of states and a finite number of controls. Solving dynamic programming problems is difficult if the state/control space are infinite. In this chapter, we will look at an important and powerful special case, called the Linear Quadratic Regulator (LQR), when we can solve dynamic programming problems easily. Just like a lot of real-world state-estimation problems can be solved using the Kalman filter and its variants, a lot of real-world control problems can be solved using LQR and its variants.

### 6.1 Discrete-time LQR

Consider a deterministic, *linear* dynamical system given by

$$x_{k+1} = Ax_k + Bu_k; \quad x_0 \text{ is given.}$$

where  $x_k \in \mathbb{R}^d$  and  $u_k \in \mathbb{R}^m$  which implies that  $A \in \mathbb{R}^{d \times d}$  and  $B \in \mathbb{R}^{d \times m}$ . In this chapter, we are interested in calculating a feedback control  $u_k = u(x_k)$  for such a system. Just like we formulated the problem in dynamic programming, we want to pick a feedback control which leads to a trajectory



1 that achieves a minimum of some run-time cost and a terminal cost. We will  
 2 assume that both the run-time and terminal costs are *quadratic* in the state and  
 3 control input, i.e.,

$$q(x, u) = \frac{1}{2}x^\top Qx + \frac{1}{2}u^\top Ru \quad (6.1)$$

4 where  $Q \in \mathbb{R}^{d \times d}$  and  $R \in \mathbb{R}^{m \times m}$  are symmetric, positive semi-definite  
 5 matrices

$$Q = Q^\top \succeq 0, \quad R = R^\top \succeq 0.$$

6 Effectively, if  $Q$  were a diagonal matrix, a large diagonal entry would  $Q_{ii}$   
 7 models our desire that the trajectory of the system should not have a large value  
 8 of the state  $x_i$  along its trajectories. We want these matrices to be positive  
 9 semi-definitive to prevent dynamic programming from picking a trajectory  
 10 which drives down the run-time cost to negative infinity by picking.

11 **Example** Consider the discrete-time equivalent of the so-called double inte-  
 12 grator  $\ddot{z}(t) = u(t)$ . The linear system in this case (obtained by creating two  
 13 states  $x := [z(t), \dot{z}(t)]$  is

$$x_{k+1} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} x_k + \begin{bmatrix} 0 \\ \Delta t \end{bmatrix} u_k.$$

14  
 15 First, note that a continuous-time linear dynamical system  $\dot{x} = Ax$  is  
 16 asymptotically stable, i.e., from any initial condition  $x(0)$  its trajectories go  
 17 to the equilibrium point  $x = 0$  ( $x(t) \rightarrow 0$  as  $t \rightarrow \infty$ ). Asymptotic stability  
 18 occurs if all eigenvalues of  $A$  are strictly negative. A discrete-time linear  
 19 dynamical system  $x_{k+1} = Ax_k$  is asymptotically stable if all eigenvalues of  
 20  $A$  have magnitude strictly smaller than 1,  $|\lambda(A)| < 1$ .

A typical trajectory of the double integrator will look as follows. Suppose

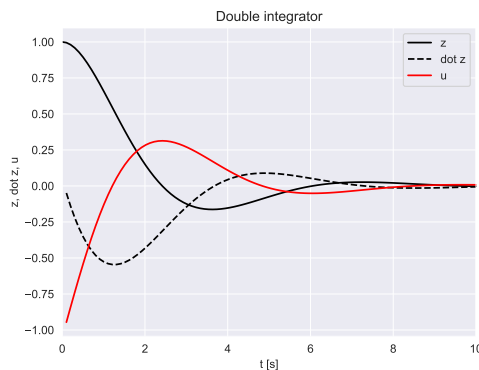


Figure 6.1: The trajectory of  $z(t)$  as a function of time  $t$  for a double integrator  $\ddot{z}(t) = u$  where we have chosen a stabilizing (i.e., one that makes the system asymptotically stable) controller  $u = -z(t) - \dot{z}(t)$ . Notice how the trajectory starts from some initial condition (in this case  $z(0) = 1$  and  $\dot{z}(0) = 0$ ) and moves towards its equilibrium point  $z = \dot{z} = 0$ .

❗ This system is called the double integrator because of the structure  $\ddot{z} = u$ ; if  $z$  denotes the position of an object the equation is simply Newton's law which connects the force applied  $u$  to the acceleration.

1  
2 we would like to pick a different controller that more quickly brings the system  
3 to its equilibrium. One way of doing so is to minimize

$$J = \sum_{k=0}^T \|x_k\|^2$$

4 which represents how far away both the position and velocity are from zero  
5 over all times  $k$ . The following figure shows the trajectory that achieves a  
6 small value of  $J$ .

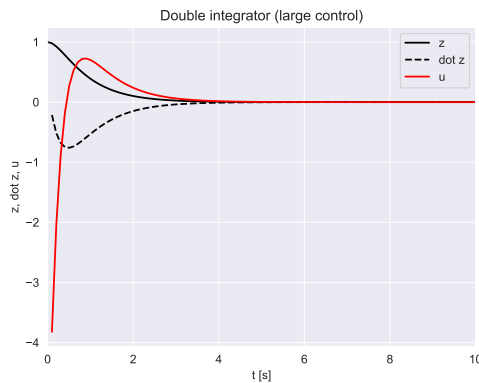


Figure 6.2: The trajectory of  $z(t)$  as a function of time  $t$  for a double integrator  $\ddot{z}(t) = u$  where we have chosen a large stabilizing control at each time  $u = -5z(t) - 5\dot{z}(t)$ . Notice how quickly the state trajectory converges to the equilibrium without much oscillation as compared to Figure 6.1 but how large the control input is at certain times.

7 This is obviously undesirable for real systems where we may want the  
8 control input to be bounded between some reasonable values (a car cannot  
9 accelerate by more than a certain threshold). A natural way of enforcing this  
10 is to modify our our desired cost of the trajectory to be

$$J = \sum_{k=0}^T (\|x_k\|^2 + \rho \|u_k\|^2)$$

11 where the value of the parameter  $\rho$  is something chosen by the user to give  
12 a good balance of how quickly the trajectory reaches the equilibrium point  
13 and how much control is exerted while doing so. Linear-Quadratic-Regulator  
14 (LQR) is a generalization of this idea, notice that the above example is equiva-  
15 lent to setting  $Q = I_{d \times d}$  and  $R = \rho I_{m \times m}$  for the run-time cost in (6.1).

16 **Back to LQR** With this background, we are now ready to formulate the  
17 Linear-Quadratic-Regulator (LQR) problem which is simply dynamic pro-  
18 gramming for a linear dynamical system with quadratic run-time cost. In order  
19 to enable the system to reach the equilibrium state even if we have only a finite  
20 time-horizon, we also include a quadratic cost

$$q_f(x) = \frac{1}{2} x^\top Q_f x. \quad (6.2)$$

1 The dynamic programming problem is now formulated as follows.

**Finite time-horizon LQR problem** Find a sequence of control inputs  $(u_0, u_1, \dots, u_{T-1})$  such that the function

$$J(x_0; u_0, u_1, \dots, u_{T-1}) = \frac{1}{2} x_T^\top Q_f x_T + \frac{1}{2} \sum_{k=0}^{T-1} (x_k^\top Q x_k + u_k^\top R u_k) \quad (6.3)$$

is minimized under the constraint that  $x_{k+1} = Ax_k + Bu_k$  for all times  $k = 0, \dots, T-1$  and  $x_0$  is given.

### 2 6.1.1 Solution of the discrete-time LQR problem

3 We know the principle of dynamic programming and can apply it to solve the  
4 LQR problem. As usual, we will compute the cost-to-go of a trajectory that  
5 starts at some state  $x$  and goes further by  $T - k$  time-steps,  $J_k(x)$  backwards.

6 Set

$$J_T^*(x) = \frac{1}{2} x^\top Q_f x \quad \text{for all } x.$$

7 Using the principle of dynamic programming, the cost-to-go  $J_{T-1}$  is given by

$$\begin{aligned} J_{T-1}^*(x_{T-1}) &= \min_u \left\{ \frac{1}{2} (x_{T-1}^\top Q x_{T-1} + u^\top R u) + J_T^*(Ax_{T-1} + Bu) \right\} \\ &= \min_u \left\{ \frac{1}{2} (x_{T-1}^\top Q x_{T-1} + u^\top R u + (Ax + Bu)^\top Q_f (Ax_{T-1} + Bu)) \right\}. \end{aligned}$$

8 We can now take the derivative of the right-hand side with respect to  $u$  to get

$$\begin{aligned} 0 &= \frac{d\text{RHS}}{du} \\ &= \frac{1}{2} \{ Ru + B^\top Q_f (Ax_{T-1} + Bu) \} \quad (6.4) \\ \Rightarrow u_{T-1}^* &= -(R + B^\top Q_f B)^{-1} B^\top Q_f A x_{T-1} \\ &\equiv -K_{T-1} x_{T-1}. \end{aligned}$$

9 where

$$K_{T-1} = (R + B^\top Q_f B)^{-1} B^\top Q_f A$$

10 is (surprisingly) also called the Kalman gain. The second derivative is positive  
11 semi-definite

$$\frac{d^2\text{RHS}}{du^2} = R + B^\top Q_f B \succeq 0$$

12 so we know that  $u_{T-1}^*$  is a minimum of the convex quantity on the right-hand  
13 side. Notice that the optimal control  $u_{T-1}^*$  is a linear function of the state  
14  $x_{T-1}$ . Let us now expand the cost-to-go  $J_{T-1}$  using this optimal value (the  
15 subscript  $T-1$  on the curly bracket simply means that all quantities are at

1 time  $T - 1$ )

$$\begin{aligned}
 J_{T-1}^*(x_{T-1}) &= \frac{1}{2} \left\{ x^\top Q x + u^{*\top} R u^* + (Ax + Bu^*)^\top Q_f (Ax + Bu^*) \right\}_{T-1} \\
 &= \frac{1}{2} x_{T-1}^\top \left\{ Q + K^\top R K + (A - BK)^\top Q_f (A - BK) \right\}_{T-1} x_{T-1} \\
 &\equiv \frac{1}{2} x_{T-1}^\top P_{T-1} x_{T-1}
 \end{aligned}$$

2 where we set the stuff inside the curly brackets to the matrix  $P$  which is also  
 3 positive semi-definite. This is great, the cost-to-go is also a quadratic function  
 4 of the state  $x_{T-1}$ . Let us assume that this pattern holds for all time steps  
 5 and the cost-to-go of the optimal LQR trajectory starting from a state  $x$  and  
 6 proceeding forwards for  $T - k$  time-steps is

$$J_k^*(x) = \frac{1}{2} x^\top P_k x.$$

7 We can now repeat the same exercise to get a recursive formula for  $P_k$  in terms  
 8 of  $P_{k+1}$ . This is the *solution* of dynamic programming for the LQR problem  
 9 as looks as follows.

$$\begin{aligned}
 P_T &= Q_f \\
 K_k &= (R + B^\top P_{k+1} B)^{-1} B^\top P_{k+1} A \\
 P_k &= Q + K_k^\top R K_k + (A - BK_k)^\top P_{k+1} (A - BK_k),
 \end{aligned} \tag{6.5}$$

10 for  $k = T - 1, T - 2, \dots, 0$ . There are a number of important observations to  
 11 be made from this calculation:

- 12 1. The optimal controller  $u_k^* = -K_k x_k$  is a linear function of the state  
 13  $x_k$ . This is only true for linear dynamical systems with quadratic costs.  
 14 Notice that both the state and control space are infinite sets but we have  
 15 managed to solve the dynamic programming problem to get the optimal  
 16 controller. We could not have done it if the run-time/terminal costs were  
 17 not quadratic or if the dynamical system were not linear. Can you say  
 18 why?
- 19 2. The cost-to-go matrix  $P_k$  and the Kalman gain  $K_k$  do not depend upon  
 20 the state and can be computed ahead of time if we know what the time  
 21 horizon  $T$  is going to be.
- 22 3. The Kalman gain changes with time  $k$ . Effectively, the LQR controller  
 23 picks a large control input to quickly reduce the run-time cost at the  
 24 beginning (if the initial condition were such that the run-time cost of  
 25 the trajectory would be very large) and then gets into a balancing act  
 26 where it balances the control effort and the state-dependent part of the  
 27 run-time cost. LQR is an optimal way to strike a balance between the  
 28 two examples in Figure 6.1 and Figure 6.2.

29 The careful reader will notice how the equations in (6.5) and our remarks  
 30 about them are similar to the update equations of the Kalman filter and our  
 31 remarks there. In fact we will see shortly how spookily similar the two are.  
 32 The key difference is that Kalman filter updates run forwards in time and

- 1 update the covariance while LQR updates run backwards in time and update  
 2 the cost-to-go matrix  $P$ . This is not surprising because LQR is an optimal  
 control problem, its update equations run backward in time.

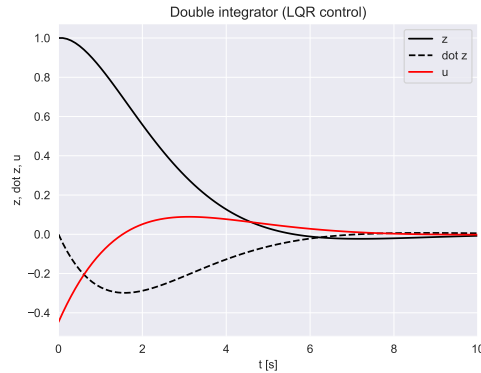


Figure 6.3: The trajectory of  $z(t)$  as a function of time  $t$  for a double integrator  $\ddot{z}(t) = u$  where we have chosen a controller obtained from LQR with  $Q = I$  and  $R = 5$ . This gives the controller to be about  $u = -0.45z(t) - 1.05\dot{z}(t)$ . Notice how we still get stabilization but the control acts more gradually. Using different values of  $R$ , we can get many different behaviors. Another key aspect of LQR as compared to Figure 6.1 where the control was chosen in an ad hoc fashion is to let us prescribe the quality of state trajectories using high-level quantities like  $Q, R$ .

3

## 4 6.2 Hamilton-Jacobi-Bellman equation

- 5 This section will show how the principle of dynamic programming looks for  
 6 continuous-time deterministic dynamical systems

$$\dot{x} = f(x, u), \quad \text{with } x(0) = x_0.$$

- 7 As we discussed in Chapter 3, we can think of this as the limit of discrete-time  
 8 dynamical system  $x_{k+1} = f^{\text{discrete}}(x_k, u_k)$  as the time discretization goes to  
 9 zero. Just like we have a sequence of controls in the discrete-time case, we  
 10 have a continuous curve that determines the control (let us also call it the  
 11 control sequence)

$$\{u(t) : t \in \mathbb{R}_+\}$$

- 12 which gives rise to a trajectory of the states

$$\{x(t) : t \in \mathbb{R}_+\}$$

- 13 for the dynamical system. Let us consider the case when we want to find  
 14 control sequences that minimize the integral of the cost along the trajectory  
 15 that stops at some fixed, finite time-horizon  $T$ :

$$q_f(x(T)) + \int_0^T q(x(t), u(t)) dt.$$

- 16 This cost is again a function of the run-time cost and a terminal cost.

**i** If you are trying this example yourself, I used the formula for continuous-time LQR and then discretized the controller while implementing it. We will see this in Section 6.2

**i** Since  $\{x(t)\}_{t \geq 0}$  and  $\{u(t)\}_{t \geq 0}$  are continuous curves and the cost is now a function of a continuous-curve, mathematicians say that the cost is a “functional” of the state and control trajectory.

**Continuous-time optimal control problem** We again want to solve for

$$J^*(x_0) = \min_{u(t), t \in [0, T]} \left\{ q_f(x(T)) + \int_0^T q(x(t), u(t)) dt \right\} \quad (6.6)$$

with the system satisfying  $\dot{x} = f(x, u)$  at each time instant. Notice that the minimization is over a function of time  $\{u(t) : t \in [0, T]\}$  as opposed to a discrete-time sequence of controls that we had in the discrete-time case. We will next look at the Hamilton-Jacobi-Bellman equation which is a method to solve optimal-control problems of this kind.

1 The principle of dynamic programming principle is still valid: if we have  
 2 an optimal control trajectory  $\{u^*(t) : t \in [0, T]\}$  we can chop it up into two  
 3 parts at some intermediate time  $t \in [0, T]$  and claim that the tail is optimal.  
 4 In preparation for this, let us define the cost-to-go of going forward by  $T - t$   
 5 time as

$$J^*(x, t) = \min_{u(s), s \in [t, T]} \left\{ q_f(x(T)) + \int_t^T q(x(s), u(s)) ds \right\},$$

the cost incurred if the trajectory starts at state  $x$  and goes forward by  $T - t$  time. This is very similar to the cost-to-go  $J_k^*(x)$  we had in discrete-time dynamic programming. Dynamic programming now gives

$$\begin{aligned} J^*(x(t), t) &= \min_{u(s), t \leq s \leq T} \left\{ q_f(x(T)) + \int_t^T q(x(s), u(s)) ds \right\} \\ &= \min_{u(s), t \leq s \leq T} \left\{ q_f(x(T)) + \int_t^{t+\Delta t} q(x(s), u(s)) ds + \int_{t+\Delta t}^T q(x(s), u(s)) ds \right\} \\ &= \min_{u(s), t \leq s \leq T} \left\{ J^*(x(t + \Delta t), t + \Delta t) + \int_t^{t+\Delta t} q(x(s), u(s)) ds \right\}. \end{aligned}$$

6 We now take the Taylor approximation of the term  $J^*(x(t + \Delta t), t + \Delta t)$  as  
 7 follows

$$\begin{aligned} &J^*(x(t + \Delta t), t + \Delta t) - J^*(x(t), t) \\ &\approx \partial_x J^*(x(t), t) (x(t + \Delta t) - x(t)) + \partial_t J^*(x(t), t) \Delta t \\ &\approx \partial_x J^*(x(t), t) f(x(t), u(t)) \Delta t + \partial_t J^*(x(t), t) \Delta t \end{aligned}$$

8 where  $\partial_x J^*$  and  $\partial_t J^*$  denote the derivative of  $J^*$  with respect to its first and  
 9 second argument respectively. We substitute this into the minimization and  
 10 collect terms of  $\Delta t$  to get

$$0 = \partial_t J^*(x(t), t) + \min_{u(t) \in U} \left\{ q(x(t), u(t)) + f(x(t), u(t)) \partial_x J^*(x(t), t) \right\}. \quad (6.7)$$

11 Notice that the minimization in (6.7) is only over *one* control input  $u(t) \in U$ ,  
 12 this is the control that we should take at time  $t$ . (6.7) is called the Hamilton-

1 Jacobi-Bellman (HJB) equation. Just like the Bellman equation

$$J_k^*(x) = \min_{u \in U} \{q_k(x, u) + J_{k+1}^*(f(x, u))\}.$$

2 has two quantities  $x$  and the time  $k$ , the Hamilton-Jacobi-Bellman equation  
 3 also has two quantities  $x$  and continuous time  $t$ . Just like the Bellman equation  
 4 is solved backwards in time starting from  $T$  with  $J_k^*(x) = q_f(x)$ , the HJB  
 5 equation is solved backwards in time by setting

$$J^*(x, T) = q_f(x).$$

You should think of the HJB equation as the continuous-time, continuous-space analogue of Dijkstra's algorithm when the number of nodes in the graph goes to infinity and the length of each edge is also infinitesimally small.

## 6 6.2.1 Infinite-horizon HJB

7 The infinite-horizon problem with the HJB equation is easy: since we know  
 8 that the optimal cost-to-go is not a function of time, we have

$$\partial_t J^*(x, t) = 0$$

9 and therefore  $J^*(x)$  satisfies

$$0 = \min_{u \in U} \{q(x, u) + f(x, u) \partial_x J^*(x)\}. \quad (6.8)$$

10 In this case, the above equation makes sense only if the integral of the run-time  
 11 cost with the optimal controller  $\int_0^\infty q(x(t), u^*(x(t))) dt$  remains bounded and  
 12 does not diverge to infinity. Therefore typically in this problem we will set  
 13  $q(0, 0) = 0$ , i.e., there is no cost for the system being at the origin with zero  
 14 control, otherwise the integral of the run-time cost will never be finite. This  
 15 also gives the boundary condition  $J^*(0) = 0$  for the HJB equation.

## 16 6.2.2 Solving the HJB equation

17 The HJB equation is a partial differential equation (PDE) because there is one  
 18 cost-to-go from every state  $x \in X$  and for every time  $t \in [0, T]$ . It belongs  
 19 to a large and important class of PDEs, collectively known as Hamilton-  
 20 Jacobi-type equations. As you can imagine, since dynamic programming is  
 21 so pervasive and solutions of DP are very useful in practice for a number of  
 22 problems, there have been many tools invented to solve the HJB equation.  
 23 These tools have applications to a wide variety of problems, from under-  
 24 standing how sound travels in crowded rooms to how light diffuses in an  
 25 animated movie scene, to even obtaining better algorithms to train deep net-  
 26 works (<https://arxiv.org/abs/1704.04932>). HJB equations are usually never  
 27 exactly solvable and a number of approximations need to be made in order to  
 28 solve it.

In this course, we will not solve the HJB equation. Rather, we are interested in seeing how the HJB equation looks for continuous-time linear dynamical systems (both deterministic and stochastic ones) and LQR problems for such systems, as done in the following section.

- 1 **An example** We will look at a classical example of the so-called car-on-the-hill problem given below. The state of the problem is the position and

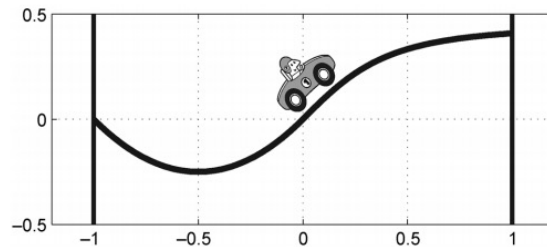


Figure 6.4: A car whose position is given by  $z(t)$  would like to climb the hill to its right and reach the top with minimal velocity. The car rolls on the hill without friction. The run-time cost is zero everywhere inside the state-space. Terminal cost is -1 for hitting the left boundary ( $z = -1$ ) and  $-1 - \dot{z}/2$  for reaching the right boundary ( $z = 1$ ). The car is a single integrator, i.e.,  $\dot{z} = u$  with only two controls ( $u = 4$  and  $u = -4$ ) and cannot exceed a given velocity (in this case  $|\dot{z}| \leq 4$ ). This looks like a simple dynamic programming problem but it is quite hard due to the constraint on the velocity. The car may need to make multiple swing ups before it gains enough velocity (but not too much) to climb up the hill.

- 2  
 3 velocity  $(z, \dot{z})$  and we can solve a two-dimensional HJB equation to obtain the  
 4 optimal cost-to-go from any state, as done by the authors Yuval Tassa and Tom  
 5 Erez in “Least Squares Solutions of the HJB Equation With Neural Network  
 6 Value-Function Approximators”  
 7 (<https://homes.cs.washington.edu/~fodorov/courses/amath579/reading/NeuralNet.pdf>).  
 8 In practice, while solving the HJB PDE, one discretizes the state-space at given  
 9 set of states and solves the HJB equation (6.7) on this grid using numerical  
 10 methods (these authors used neural networks to solve it). The end result looks  
 11 as follows.

### 12 6.2.3 Continuous-time LQR

13 Consider a linear continuous-time dynamical system given by

$$\dot{x} = A x + B u; \quad x(0) = x_0.$$

- 14 In the LQR problem, we are interested in finding a control trajectory that  
 15 minimizes, as usual, a cost function that is quadratic in states and controls,  
 16 except that we have an integral of the run-time cost because our system is a



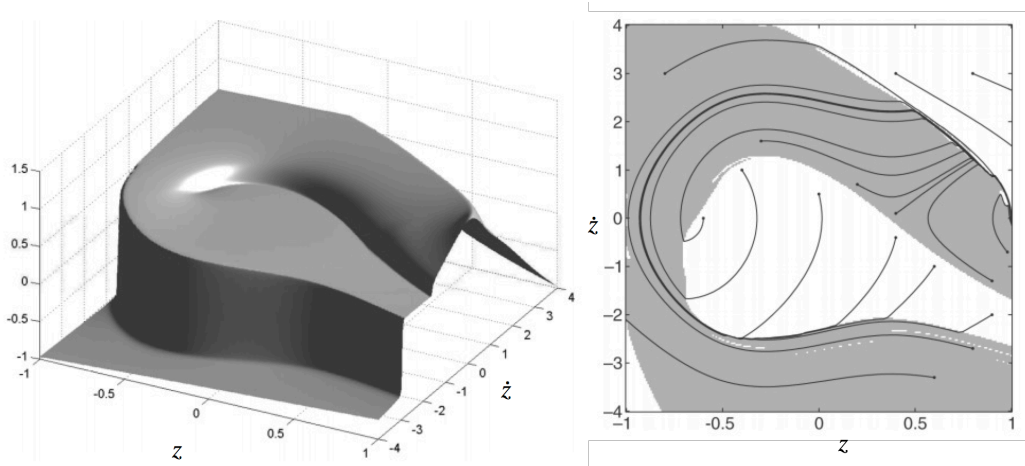


Figure 6.5: The left-hand side picture shows the infinite-horizon cost-to-go  $J^*(z, \dot{z})$  for the car-on-the-hill problem. Notice how the value function is non-smooth at various places. This is quite typical of difficult dynamic programming problems. The right-hand side picture shows the optimal trajectories of the car  $(z(t), \dot{z}(t))$ ; gray areas indicate maximum control and white areas indicate minimum control. The black lines show a few optimal control sequences taken the car starting from various states in the state-space. Notice how the optimal control trajectory can be quite different even if the car starts from nearby states  $(-0.5, 1)$  and  $(-0.4, 1.2)$ . This is also quite typical of difficult dynamic programming problems.

1 continuous-time system

$$\frac{1}{2} x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T x(t)^\top Q x(t) + u(t)^\top R u(t) dt.$$

2 This is a very nice setup for using the HJB equation from the previous section.

3 Let us use our intuition from the discrete-time LQR problem and say that  
4 the optimal cost is quadratic in the states, namely,

$$J^*(x, t) = \frac{1}{2} x(t)^\top P(t) x(t);$$

5 notice that as usual the optimal cost-to-go is a function of the states  $x$  and the  
6 time  $t$  because is the optimal cost of the continuous-time LQR problem if the  
7 system starts at a state  $x$  at time  $t$  and goes on until time  $T \geq t$ . We will now  
8 check if this  $J^*$  satisfies the HJB equation (we don't write the arguments  $x(t)$ ,  
9  $u(t)$  etc. to keep the notation clear)

$$-\partial_t J^*(x, t) = \min_{u \in U} \left\{ \frac{1}{2} (x^\top Q x + u^\top R u) + (A x + B u)^\top \partial_x J^*(x, t) \right\} \quad (6.9)$$

10 from (6.7). The minimization is over the control input that we take at time  $t$ .

1 Also notice the partial derivatives

$$\begin{aligned}\partial_x J^*(x, t) &= P(t) x. \\ \partial_t J^*(x, t) &= \frac{1}{2} x^\top \dot{P}(t) x.\end{aligned}$$

2 It is convenient in this case to see that the minimization can be performed  
3 using basic calculus (just like the discrete-time LQR problem), we differentiate  
4 with respect to  $u$  and set it to zero.

$$\begin{aligned}0 &= \frac{\text{d RHS of HJB}}{\text{d}u} \\ \Rightarrow u^*(t) &= -R^{-1} B^\top P(t) x(t) \\ &\equiv -K(t) x(t).\end{aligned}\tag{6.10}$$

5 where  $K(t) = R^{-1} B^\top P(t)$  is the Kalman gain. The controller is again linear  
6 in the states  $x(t)$  and the expression for the gain is very simple in this case,  
7 much simpler than discrete-time LQR. Since  $R \succ 0$ , we also know that  $u^*(t)$   
8 computed here is the global minimum. If we substitute this value of  $u^*(t)$   
9 back into the HJB equation we have

$$\left. \{ \} \right|_{u^*(t)} = \frac{1}{2} x^\top \{ PA + A^\top P + Q - PBR^{-1}B^\top P \} x.$$

10 In order to satisfy the HJB equation, we must have that the expression above is  
11 equal to  $-\partial_t J^*(x, t)$ . We therefore have, what is called the Continuous-time  
12 Algebraic Riccati Equation (CARE), for the matrix  $P(t) \in \mathbb{R}^{d \times d}$

$$-\dot{P} = PA + A^\top P + Q - PBR^{-1}B^\top P.\tag{6.11}$$

13 This is an ordinary differential equation for the matrix  $P$ . The derivative  
14  $\dot{P} = \frac{dP}{dt}$  stands for differentiating every entry of  $P$  individually with time  $t$ .  
15 The terminal cost is  $\frac{1}{2} x(T)^\top Q_f x(T)$  which gives the boundary condition for  
16 the ODE as

$$P(T) = Q_f.$$

17 Notice that the ODE for the  $P(t)$  travels backwards in time.

18 Continuous-time LQR has particularly easy equations, as you can see  
19 in (6.10) and (6.11) compared to those for discrete-time ((6.4) and (6.5)).  
20 Special techniques have been invented for solving the Riccati equation. I  
21 used the function `scipy.linalg.solve_continuous_are` to obtain Figure 6.3 using  
22 the continuous-time equations; the corresponding function for solving  
23 Discrete-time Algebraic Riccati Equation (DARE) which is given in (6.5)  
24 is `scipy.linalg.solve_discrete_are`. The continuous-time point-of-view also  
25 gives powerful connections to the Kalman filter, where you can show that the  
26 Kalman filter and LQR are duals of each other: in fact the equations for the  
27 Kalman filter (in continuous-time) and continuous-time LQR turn out to be  
28 exactly the same after you interchange appropriate quantities (!).

1 **Infinite-horizon LQR** Just like the infinite-horizon HJB equation has  $\partial_t J^*(x, t) =$   
 2 0, if we have an infinite-horizon LQR problem, the cost matrix  $P$  should not  
 3 be a function of time

$$\dot{P} = 0.$$

4 The continuous-time algebraic Riccati equation in (6.11) now becomes

$$PA + A^\top P + Q - PBR^{-1}B^\top P.$$

5 with the cost-to-go being given by  $J^*(x) = \frac{1}{2}x^\top Px$ .

### 6 6.3 Stochastic LQR

7 We will next look at a very powerful result. Say we have a stochastic linear  
 8 dynamical system

$$\dot{x}(t) = Ax(t) + Bu(t) + B_\epsilon \epsilon(t); x(0) \text{ is given}$$

9 where  $\epsilon(t)$  is standard Gaussian noise  $\epsilon(t) \sim N(0, I)$  that is uncorrelated  
 10 in time and would like to find a control sequence  $\{u(t) : t \in [0, T]\}$  that  
 11 minimizes a quadratic run-time and terminal cost

$$\mathbb{E}_{\epsilon(t):t \in [0, T]} \left[ \frac{1}{2}x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T x(t)^\top Q x(t) + u(t)^\top R u(t) dt \right].$$

12 over a finite-horizon  $T$ . Notice that since the system is stochastic now, we  
 13 should minimize the expected value of the cost over all possible realizations of  
 14 the noise  $\{\epsilon(t) : t \in [0, T]\}$ . This is a very challenging problem, conceptually  
 15 it is the equivalent of dynamic programming for an MDP with an infinite  
 16 number of states  $x(t) \in \mathbb{R}^d$  and an infinite number of controls  $u(t) \in \mathbb{R}^m$ .

17 However, it turns out that the optimal controller that we should pick in this  
 18 case is also given by the standard LQR problem

$$u^*(t) = -R^{-1}B^\top P(t) x(t)$$

$$\text{with } -\dot{P} = PA + A^\top P + Q - PBR^{-1}B^\top P; P(T) = Q_f.$$

19 We will not do the proof (it is easy but tedious, you can try to show it by  
 20 writing the HJB equation for the stochastic LQR problem). This is a very  
 21 surprising result because it says that even if the dynamical system had noise,  
 22 the optimal control we should pick is exactly the same as the control we would  
 23 have picked had the system been deterministic. It is a special property of the  
 24 LQR problem and not true for other dynamical systems (nonlinear ones, or  
 25 ones with non-Gaussian noise) or other costs.

26 We know that the control  $u^*(t)$  is the same as the deterministic case. Is  
 27 the cost-to-go  $J^*(x, t)$  also the same? If you think about this, the cost-to-go  
 28 in the stochastic case has to be a bit larger than the deterministic case because  
 29 the noise  $\epsilon(t)$  is always going to non-zero when we run the system, the LQR  
 30 cost  $J^*(x_0, 0) = \frac{1}{2}x_0^\top P(0)x_0$  is, after all, only the cost of the deterministic  
 31 problem. It turns out that the cost for the stochastic LQR case for an initial

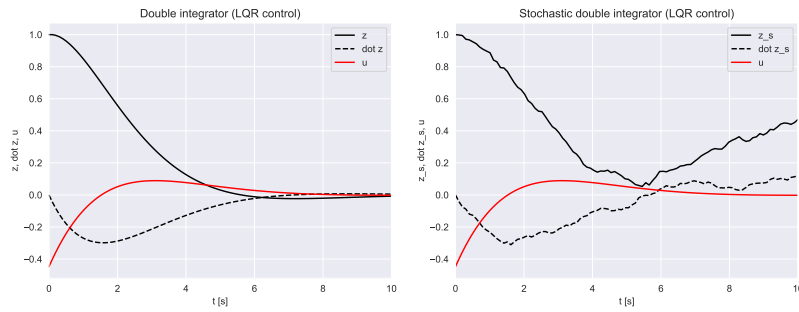


Figure 6.6: Comparison of the state trajectories of deterministic LQR and stochastic LQR problem with  $B_\epsilon = [0.1, 0.1]$ . The left panel is the same as that in Figure 6.3. The control input is the same in both cases but notice that the states in the plot on the right need not converge to the equilibrium due to noise. The cost of the trajectory will also be higher for the stochastic LQR case due to this. The total cost is  $J^*(x_0) = 32.5$  for the deterministic case (32.24 for the quadratic state-cost and 0.26 for the control cost). The total cost  $J^*(x_0)$  is much higher for the stochastic case, it is 81.62 (81.36 for the quadratic state cost and 0.26 for the control cost).

1 state  $x_0$  is

$$\begin{aligned}
 J^*(x_0, 0) &= \mathbf{E}_{\epsilon(t):t \in [0, T]} \left[ \frac{1}{2} x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T \dots dt \right] \\
 &= \frac{1}{2} x_0^\top P(0) x_0 + \frac{1}{2} \int_0^T \text{tr}(P(t) B_\epsilon B_\epsilon^\top) dt.
 \end{aligned}$$

2 The first term is the same as that of the deterministic LQR problem. The  
 3 second term is the penalty we incur for having a stochastic dynamical system.  
 4 This is the minimal cost achievable for stochastic LQR but it is not the same  
 5 as that of the deterministic LQR.

## 6 6.4 Linear Quadratic Gaussian (LQG)

7 Our development in the previous sections and the previous chapter was based  
 8 on a Markov Decision Process, i.e., we know the state  $x(t)$  at each instant in  
 9 time  $t$  even if this state  $x(t)$  changes stochastically. We said that the optimal  
 10 control for the linear dynamics is still  $u^*(t) = -K(t) x(t)$ . What should one  
 11 do if we cannot observe the state exactly?

12 Imagine a “continuous-time” form the observation equation in the Kalman  
 13 filter where we receive observations of the form

$$y(t) = Cx(t) + D\nu.$$

14 where  $\nu \sim N(0, I)$  is standard Gaussian noise that corrupts our observations  
 15  $y$ . If we extrapolate the definitions of the Kalman filter mean and covariance  
 16 to this continuous-time setting, we can write the KF as follows. We know that  
 17 the Kalman filter is the optimal estimate of the state given all past observations,

1 so it computes

$$\mu(t) = \mathbb{E}_{\epsilon(s), \nu(s): s \in [0, t]} [x(t) | y(s) : s \in [0, t]].$$

2 There exists a “continuous-time version” of the Kalman filter (which was  
3 actually invented first), called the Kalman-Bucy filter. If the covariance of the  
4 estimate is

$$\Sigma(t) = \mathbb{E}_{\epsilon(s), \nu(s): s \in [0, t]} [x(t) x(t)^\top | y(s) : s \in [0, t]],$$

5 the Kalman-Bucy filter updates  $\mu(t)$ ,  $\Sigma(t)$  using the differential equation

$$\begin{aligned} \frac{d}{dt} \mu(t) &= Ax(t) + Bu(t) + K(t) (y(t) - C\mu(t)) \\ \frac{d}{dt} \Sigma(t) &= A\Sigma(t) + \Sigma(t)A^\top + B_\epsilon B_\epsilon^\top - K(t)DD^\top K(t)^\top \end{aligned} \quad (6.12)$$

where  $K(t) = \Sigma(t) C^\top (DD^\top)^{-1}$ .

6 This equation is very close to the Kalman filter equations you saw in Chapter  
7 3. In particular, notice the close similarity of the expression for the Kalman  
8 gain  $K(t)$  with the Kalman gain of the LQR problem. You can read more at  
9 [https://en.wikipedia.org/wiki/Kalman\\_filter](https://en.wikipedia.org/wiki/Kalman_filter).

**Linear Quadratic Gaussian (LQG)** It turns out that we can plug in the Kalman filter estimate  $\mu(t)$  of the state  $x(t)$  in order to compute optimal control for LQR if we know the state only through observations  $y(t)$

$$u^*(t) = -K(t) \mu(t). \quad (6.13)$$

It is almost as if, we can blindly run a Kalman Filter in parallel with the deterministic LQR controller and get the optimal control for the stochastic LQR problem even if we did not observe the state of the system exactly. This method is called Linear Quadratic Gaussian (LQG).

This is a very powerful and surprising result. It is only true for linear dynamical systems with linear observations, Gaussian noise in both the dynamics and the observations and quadratic run-time and terminal costs. It is not true in other cases. However, it is so elegant and useful that it inspires essentially all other methods that control a dynamical system using observations from sensors.

10 **Certainty equivalence** For instance, even if we are using a particle filter to  
11 estimate the state of the system, we usually use the mean of the state estimate  
12 at time  $t$  given by  $\mu(t)$  “as if” it were the true state of the system. Even if we  
13 were using some other feedback control  $u(x)$  different than the LQR control  
14 (say feedback linearization), we usually plug in this estimate  $\mu(t)$  in place of  
15  $x(t)$ . Doing so is called “certainty equivalence” in control theory/robotics,  
16 which is a word borrowed from finance where one takes decisions (controls)  
17 directly using the estimate of the state (say stock price) while fully knowing

❗ As we discussed while introducing stochastic dynamical systems, there are various mathematical technicalities associated with conditioning on a continuous-time signal  $\{y(s) : s \in [0, t]\}$ . To be precise mathematicians define what is called a “filtration”  $\mathcal{Y}(t)$  which is the union of the Borel  $\sigma$ -fields constructed using increasing subsets of the set  $\{y(s) : s \in [0, t]\}$ . Let us not worry about this here.

1 the the stock price will change in the future stochastically.

## 2 **6.4.1 (Optional material) The duality between the Kalman** 3 **Filter and LQR**

4 We can re-write the covariance in (6.12) using the identity

$$\frac{d}{dt} (\Sigma(t)^{-1}) = \Sigma(t)^{-1} \dot{\Sigma}(t) \Sigma(t)^{-1}$$

5 to get

$$\dot{S} = C^\top (DD^\top)^{-1} C - A^\top S - SA - SB_w B_w^\top S \quad (6.14)$$

6 where we have defined  $S := \Sigma^{-1}$ .

7 Notice that the two equations, updates to the LQR cost matrix in (6.11)

$$-\dot{P} = PA + A^\top P + Q - PBR^{-1}B^\top P$$

8 look quite similar to this equation. In fact, they are identical and you can  
9 substitute the following.

LQR	Kalman-Bucy filter
$P$	$\Sigma^{-1}$
$A$	$-A$
$BR^{-1}B$	$B_w B_w^\top$
$Q$	$C^\top (DD^\top)^{-1} C$
$t$	$T - t$

11 Let us analyze this equivalence. Notice that the inverse of the Kalman  
12 filter covariance is like the cost matrix of LQR. This is conceptually easy to  
13 understand, our figure of merit for filtering is the inverse covariance matrix  
14 (smaller the better) and our figure of merit for the LQR problem is the cost  
15 matrix  $P$  (smaller the better). Similarly, smaller the LQR cost, better the  
16 controller. The “dynamics” of the Kalman filter is the reverse of the dynamics  
17 of the LQR problem, this shows that the  $P$  matrix is updated backwards in  
18 time while the covariance  $\Sigma$  is updated forwards in time. The next identity

$$BR^{-1}B^\top = B_w B_w^\top$$

19 is very interesting. Imagine a situation where we have a fully-actuated system  
20 with  $B = I$  and  $B_w$  being a diagonal matrix. This identity suggests that  
21 larger the control cost  $R_{ii}$  of a particular actuator  $i$ , lower is the noise of using  
22 that actuator  $(B_w)_{ii}$ , and vice-versa. This is how muscles in your body have  
23 evolved: muscles that are cheap to use (low  $R$ ) are also very noisy in what they  
24 do whereas muscles that are expensive to use (large  $R$ ) which are typically  
25 the biggest muscles in the body are also the least noisy and most precise. You  
26 can read more about this in the paper titled “General duality between optimal  
27 control and estimation” by Emanuel Todorov. The next identity

$$Q = C^\top (DD^\top)^{-1} C$$

is related to the quadratic state-cost in LQR. Imagine the situation where both  $Q, D$  are diagonal matrices. If the noise in the measurements  $D_{ii}$  is large, this is equivalent to the state-cost matrix  $Q_{ii}$  being small; roughly there is no way we can achieve a low state-cost  $x^\top Q x$  in our system that consists of LQR and a Kalman filter (this combination is known as Linear Quadratic Gaussian LQG as saw before) if there is lots of noise in the state measurements. The final identity

$$t = T - t$$

is the observation that we have made many times before: dynamic programming travels backwards in time and the Kalman filter travels forwards in time.

## 6.5 Iterative LQR (iLQR)

This section is analogous to the section on the Extended Kalman Filter. We will study how to solve optimal control problems for a nonlinear dynamical system

$$\dot{x} = f(x, u); x(0) = x_0 \text{ is given.}$$

We will consider a deterministic continuous-time dynamical system, the modifications to following section that one would make if the system is discrete-time, or stochastic, are straightforward and follow the same strategy. First consider the problem where the run-time and terminal costs are quadratic

$$\frac{1}{2}x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T x(t)^\top Q x(t) + u(t)^\top R u(t) dt.$$

**Receding horizon control and Model Predictive Control (MPC)** One easy way to solve the dynamic programming problem, i.e., find a control trajectory of the *nonlinear* system that minimizes this cost functional, approximately, is by linearizing the system about the initial state  $x_0$  and some reference control  $u_0$  (this can usually be zero). Let the linear system be

$$\dot{z} = A_{x_0, u_0} z + B_{x_0, u_0} v; z(0) = 0; \quad (6.15)$$

where  $A_{x_0, u_0} = \left. \frac{df}{dx} \right|_{x=x_0, u=u_0}$  and  $B_{x_0, u_0} = \left. \frac{df}{du} \right|_{x=x_0, u=u_0}$  are the Jacobians of the nonlinear function  $f(x, u)$  with respect to the state and control respectively. The state of the linearized dynamics is

$$z := x - x_0, \text{ and } v := u - u_0,$$

We have emphasized the fact that the matrices  $A_{x_0, u_0}, B_{x_0, u_0}$  depend upon the reference state and control using the subscript. Given the above linear system, we can find a control sequence  $u^*(\cdot)$  that minimizes the cost functional using the standard LQR formulation. Notice now that even we computed this control trajectory using the approximate linear system, it can certainly be *executed* on the nonlinear system, i.e., at run-time we will simply set  $u \equiv u^*(z)$ .

The linearized dynamics in (6.15) is potentially going to be very different from the nonlinear system. The two are close in the neighborhood of  $x_0$  (and

1  $u_0$ ) but as the system evolves using our control input to move further away  
 2 from  $x_0$ , the linearized model no longer is a faithful approximation of the  
 3 nonlinear model. A reasonable way to fix matters is to linearize about another  
 4 point, say the state and control after  $t = 1$  seconds,  $x_1, u_1$  to get a new system

$$\dot{z} = A_{x_1, u_1} z + B_{x_1, u_1} v; z(0) = 0$$

5 and take the LQR-optimal control corresponding to this system for the next  
 6 second.

7 The above methodology is called “receding horizon control”. The idea is  
 8 that we compute the optimal control trajectory  $u^*(\cdot)$  using an approximation  
 9 of the original system and recompute this control every few seconds when our  
 10 approximation is unlikely to be accurate. This is a very popular technique to  
 11 implement optimal controllers in typical applications. The concept of using an  
 12 approximate model (almost invariably, a linear model with LQR cost) to plan  
 13 for the near-term future and resolving the problem in receding horizon fashion  
 14 once the system is at the end of this short time-horizon is called “Model  
 15 Predictive Control”.

16 MPC is, perhaps, the second most common control algorithm implemented  
 17 in the world. It is responsible for running most complex engineering systems  
 18 that you can think of—power grids, oil refineries, chemical plants, rockets,  
 19 aircrafts etc. Essentially, one never implements LQR directly, it is always im-  
 20 plemented inside an MPC. For instance, in autonomous driving, the trajectory  
 21 that the vehicle plans for traveling between two points  $A$  and  $B$  depends upon  
 22 the current locations of the other cars/pedestrians in its vicinity, and potentially  
 23 some prediction model of where they will be in the future. As the vehicle  
 24 starts moving along this trajectory, the rest of the world evolves around it and  
 25 we recompute the optimal trajectory to take into account the actual locations  
 26 of the cars/pedestrians in the future.

🔗 Can you guess what is *the* most common control algorithm in the world?

### 27 6.5.1 Iterative LQR (iLQR)

28 Now let us consider the situation when in addition to a nonlinear system,

$$\dot{x} = f(x, u); x(0) = x_0,$$

29 the run-time and terminal cost is also nonlinear

$$q_f(x(T)) + \int_0^T q(x(t), u(t)) dt.$$

30 We can solve the dynamic programming problem in this case approximately  
 31 using the following iterative algorithm.

32 Assume that we are given an initial control trajectory  $u^{(0)}(\cdot) = \{u^{(0)}(t) : t \in [0, T]\}$ .  
 33 Let  $x^{(0)}(\cdot)$  be the state trajectory that corresponds to taking this control on  
 34 the nonlinear system, with of course  $x^{(0)}(0) = x_0$ . At each iteration  $k$ , the  
 35 Iterative LQR algorithm performs the following steps.

36 **Step 1** Linearize the nonlinear system about the state trajectory  $x^{(k)}(\cdot)$  and



1  $u^{(k)}(\cdot)$  using

$$z(t) := x(t) - x^{(k)}(t), \text{ and } v(t) := u(t) - u^{(k)}(t)$$

2 to get a new system

$$\dot{z} = A^{(k)}(t)z + B^{(k)}(t)v; z(0) = 0$$

3 where

$$A^{(k)}(t) = \left. \frac{df}{dx} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}$$

$$B^{(k)}(t) = \left. \frac{df}{du} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}$$

4 and compute the Taylor series approximation of the nonlinear cost up to the  
5 second order

$$q_f(x(T)) \approx \text{constant} + z(T)^\top \left. \frac{dq_f}{dx} \right|_{x(T)=x^{(k)}(T)}$$

$$+ z(t)^\top \left. \frac{d^2 q_f}{dx^2} \right|_{x(T)=x^{(k)}(T)} z(t),$$

6

$$q(x, u, t) \approx \text{constant} + z(t)^\top \underbrace{\left. \frac{dq}{dx} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\text{affine term}}$$

$$+ v(t)^\top \underbrace{\left. \frac{dq}{du} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\text{affine term}}$$

$$+ z(t)^\top \underbrace{\left. \frac{d^2 q}{dx^2} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\equiv Q} z(t)$$

$$+ v(t)^\top \underbrace{\left. \frac{d^2 q}{du^2} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\equiv R} v(t).$$

7 This is an LQR problem with run-time cost that depends on time (like our  
8 discrete-time LQR formulation, the continuous-time formulation simply has  
9 Q, R to be functions of time  $t$  in the Riccati equation) and which also has  
10 terms that are affine in the state and control in addition to the usual quadratic  
11 cost terms.

12 **Step 2** Solve the above linearized problem using standard LQR formulation to  
13 get the new control trajectory

$$u^{(k+1)}(t) := u^{(k)}(t) - Kz(t).$$

14 Simulate the *nonlinear* system using the control  $u^{(k+1)}(\cdot)$  to get the new state  
15 trajectory  $x^{(k+1)}(\cdot)$ .

16 Some important comments to remember about the iLQR algorithm.

17 1. There are many ways to pick the initial control trajectory  $u^{(0)}(\cdot)$ , e.g.,

🔗 How will you solve for the optimal controller for a linear dynamics for the cost

$$\int_0^T \left( q^\top x + \frac{1}{2} x^\top Q x \right) dt,$$

i.e., when in addition the quadratic cost, we also have an affine term?

1 using a spline to get an arbitrary control sequence, using a spline to  
 2 interpolate the states to get a trajectory  $x^{(0)}(\cdot)$  and then back-calculate  
 3 the control trajectory, using the LQR solution based on the lineariza-  
 4 tion about the initial state, feedback linearization/differential flatness  
 5 ([https://en.wikipedia.org/wiki/Feedback\\_linearization](https://en.wikipedia.org/wiki/Feedback_linearization)) etc.

- 6 2. The iLQR algorithm is an approximate solution to dynamic program-  
 7 ming for nonlinear system with general, nonlinear run-time and terminal  
 8 costs. This is because the the algorithm uses a linearization about the  
 9 previous state and control trajectory to compute the new control trajec-  
 10 tory. iLQR is not guaranteed to find the optimal solution of dynamic  
 11 programming, although in practice with good implementations, it works  
 12 excellently.
- 13 3. We can think of iLQR as an algorithm to track a given state trajectory  
 14  $x^g(t)$  by setting

$$q_f = 0, \text{ and } q(x, u) = \|x^g(t) - x(t)\|^2.$$

15 This is often how iLQR is typically used in practice, e.g., to make  
 16 an autonomous race car closely follow the racing line (see the paper  
 17 “BayesRace: Learning to race autonomously using prior experience”  
 18 <https://arxiv.org/abs/2005.04755> and <https://www.youtube.com/watch?v=dgIpf0Lg8Ek>  
 19 for a clever application of using MPC to track a challenging race line),  
 20 or to make a drone follow a given desired trajectory  
 21 (<https://www.youtube.com/watch?v=QREeZvHg0IQ>).

22 **Differential Dynamic Programming (DDP)** is a suite of techniques that is  
 23 a more powerful version of iterated LQR. Instead of linearizing the dynamics  
 24 and taking a second order Taylor approximation of the cost, DDP takes a  
 25 second order approximation of the Bellman equation directly. The two are  
 26 not the same; DDP is the more correct version of iLQR but is much more  
 27 challenging computationally.

28 Broadly speaking, iLQR and DDP are used to perform control for some of  
 29 the most sophisticated robots today, you can see an interesting discussion of  
 30 the trajectory planning of some of the DARPA Humanoid Robotics Challenge  
 31 at <https://www.cs.cmu.edu/~cga/drc/atlas-control>. Techniques like feedback  
 32 linearization work excellently for drones where we do not really care for opti-  
 33 mal cost (see “Minimum snap trajectory generation and control for quadrotors”  
 34 <https://ieeexplore.ieee.org/document/5980409>) while LQR and its variants are  
 35 still heavily utilized for satellites in space.

# Chapter 7

## Imitation Learning

### Reading

1. The DAGGER algorithm  
(<https://www.cs.cmu.edu/~sross1/publications/Ross-AIStats11-NoRegret.pdf>)
2. [https://www.youtube.com/watch?v=TUBBIgtQL\\_k](https://www.youtube.com/watch?v=TUBBIgtQL_k)
3. An Algorithmic Perspective on Imitation Learning  
(<https://arxiv.org/pdf/1811.06711.pdf>)

This is the beginning of Module 3 of the course. The previous two modules have been about how to estimate the state of the world around the robot (Module 1) and how to move the robot (or the world) to a desired state (Module 2). Both of these required that we maintain a model of the dynamics of the robot; this model may be inaccurate and we fudged over this inaccuracy by modeling the remainder as “noise” in Markov Decision Processes.

The next few lectures introduce different aspects of what is called Reinforcement Learning (RL). This is a very large field and you can think of using techniques from RL in many different ways.

- 1. Dynamic programming with function approximation.** If we are solving a dynamic programming problem, we can think of writing down the optimal cost-to-go  $J^*(x, t)$  as a function of some parameters, e.g., the cost-to-go is

$$J_\varphi(x, t) = \frac{1}{2} x(t)^\top \underbrace{\left( \text{some function of } A, B, Q, R \right)}_{\text{function of } \varphi} x(t)$$

for LQR. We know the stuff inside the brackets to be exactly  $P(t)$  but, if we did not, it could be written down as some generic function of parameters  $\varphi$ . We know that any cost-to-go that satisfies the Bellman

equation is the optimal cost-to-go, so we can now “fit” the candidate function  $J_\varphi(x, t)$  to satisfy the Bellman equation. Similarly, one may also express the optimal feedback control  $u(\cdot)$  using some parameters  $\theta$  as

$$u_\theta(\cdot).$$

We will see how to fit such functions in this chapter.

2. **Learning from data.** It may happen that we do not know very much about the dynamical system, e.g., we do not know a good model for what drives customers as they buy items in an online merchandise platform, or a robot traveling in a crowded area may not have a good model for how large crowds of people walk around it. One may collect data from these systems fit some model of the form  $\dot{x} = f(x, u)$  to the data and then go back to the techniques of Module 2. It is typically not clear how much data one should collect. RL gives a suite of techniques to learn the cost-to-go in these situations by collecting and assimilating the data *automatically*. These techniques go under the umbrella of policy gradients, on-policy methods etc. One may also simply “memorize” the data provided by an expert operator, this is called Imitation Learning and we will discuss it next.

**Some motivation** Imitation Learning is also called “learning from demonstrations”. This is in fact one of the earliest successful examples of using a neural network for driving. The ALVINN project at CMU by Dean Pomerleau in 1988 (<https://www.youtube.com/watch?v=2KMAAmkz9go>) used a two-layer neural network with 5 hidden neurons, about 1000 inputs from the pixels of a camera and 30 outputs. It successfully drove in different parts of the United States and Germany. Imitation learning has also been responsible for numerous other early-successes of RL, e.g., acrobatic maneuvers on an RC helicopter ([http://ai.stanford.edu/~acoates/papers/AbbeelCoatesNg\\_IJRR2010.pdf](http://ai.stanford.edu/~acoates/papers/AbbeelCoatesNg_IJRR2010.pdf)).

Imitation Learning seeks to record data from experts, e.g., humans, and reproduce these desired behaviors on robots. The key questions we should ask, and which we will answer in this chapter, are as follows.

1. Who should demonstrate (experts, amateurs, or novices) and how should we record data (what states, controls etc.)?
2. How should we learn from this data? e.g., fit a supervised regression model for the policy. How should one ignore bad behaviors in non-expert data?
3. And most importantly, what can we do if the robot encounters a situation which was not in the dataset.

## 7.1 A crash course in supervised learning

Nature gives us data  $X$  and targets  $Y$  for this data.

$$X \rightarrow Y.$$

Nature does not usually tell us what property of a datum  $x \in X$  results in a particular prediction  $y \in Y$ . We would like to learn to imitate Nature, namely predict  $y$  given  $x$ .

What does such learning mean? It is simply a notion of being able to identify patterns in the input data without explicitly programming a computer for prediction. We are often happy with a learning process that identifies correlations: if we learn correlations on a few samples  $(x^1, y^1), \dots, (x^n, y^n)$ , we may be able to predict the output for a new datum  $x^{n+1}$ . We may not need to know *why* the label of  $x^{n+1}$  was predicted to be so and so.

Let us say that Nature possesses a probability distribution  $P$  over  $(X, Y)$ . We will formalize the problem of machine learning as Nature drawing  $n$  independent and identically distributed samples from this distribution. This is denoted by

$$D_{\text{train}} = \{(x^i, y^i) \sim P\}_{i=1}^n$$

is called the “training set”. We use this data to identify patterns that help make predictions on some future data.

**What is the task in machine learning?** Suppose  $D_{\text{train}}$  consists of  $n = 50$  RGB images of size  $100 \times 100$  of two kinds, ones with an orange inside them and ones without.  $10^4$  is a large number of pixels, each pixel taking any of the possible  $255^3$  values. Suppose we discover that one particular pixel, say at location  $(25, 45)$ , takes distinct values in all images inside our training set. We can then construct a predictor based on this pixel. This predictor, it is a binary classifier, perfectly maps the training images to their labels (orange: +1 or no orange: -1). If  $x_{ij}^k$  is the  $(ij)^{\text{th}}$  pixel for image  $x^k$ , then we use the function

$$f(x) = \begin{cases} y^k & \text{if } x_{ij}^k = x_{ij} \text{ for some } k = 1, \dots, n \\ -1 & \text{otherwise.} \end{cases}$$

This predictor certainly solves the task. It correctly works for all images in the training set. Does it work for images outside the training set?

Our task in machine learning is to learn a predictor that works *outside* the training set. The training set is only a source of information that Nature gives us to find such a predictor.

Designing a predictor that is accurate on  $D_{\text{train}}$  is trivial. A hash function that memorizes the data is sufficient. This is NOT our task in machine learning. We want predictors that generalize to new data outside  $D_{\text{train}}$ .

🔗 How many such binary classifiers are there at most?

1 **Generalization** If we never see data from outside  $D_{\text{train}}$  why should we hope  
 2 to do well on it? The key is the distribution  $P$ . Machine learning is formalized  
 3 as constructing a predictor that works well on new data that is also drawn  
 4 independently from the distribution  $P$ . We will call this set of data the “test  
 5 set”

$$D_{\text{test}}$$

6 and it is constructed similarly. This assumption is important. It provides  
 7 coherence between past and future samples: past samples that were used to  
 8 train and future samples that we will wish to predict upon. How to find such  
 9 predictors that work well on new data? The central idea in machine learning is  
 10 to restrict the set of possible binary functions that we consider.

We are searching for a predictor that generalizes well but only have the training data to select predictors.

11 The *right* class of functions  $f$  cannot be too large, otherwise we will find  
 12 our binary classifier above as the solution, and that is not very useful. The class  
 13 of functions cannot be too small either, otherwise we won't be able to predict  
 14 difficult images. If the predictor does not even work well on the training set,  
 15 there is no reason why we should expect it to work on the test set.

Finding this correct class of functions with the right balance is what machine learning is all about.

🔗 Can you now think how is machine learning different from other fields you might know such as statistics or optimization?

### 16 7.1.1 Fitting a machine learning model

17 Let us now solve a classification problem. We will again go around the model  
 18 selection problem and consider the class of linear classifiers. Assume binary  
 19 labels  $Y \in \{-1, 1\}$ . To keep the notation clear, we will use the trick of  
 20 appending a 1 to the data  $x$  and hide the bias term  $b$  in the linear classifier. The  
 21 predictor is now given by

$$\begin{aligned} f(x; w) &= \text{sign}(w^\top x) \\ &= \begin{cases} +1 & \text{if } w^\top x \geq 0 \\ -1 & \text{else.} \end{cases} \end{aligned} \quad (7.1)$$

22 We have used the sign function denoted as  $\text{sign}$  to get binary  $\{-1, +1\}$  outputs  
 23 from our real-valued prediction  $w^\top x$ . This is the famous perceptron model of  
 24 Frank Rosenblatt.

25 We want the predictions of the model to match those in the training data  
 26 and devise an objective to fit/train the perceptron.

$$\ell_{\text{zero-one}}(w) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{y^i \neq f(x^i; w)\}}. \quad (7.2)$$

1 The indicator function inside the summation measures the number of mistakes  
 2 the perceptron makes on the training dataset. The objective here is designed to  
 3 find weights  $w$  that minimizes the average number of mistakes, also known as  
 4 the training error. Such a loss that measures the mistakes is called the zero-one  
 5 loss, it incurs a penalty of 1 for a mistake and zero otherwise.

❓ Can you think of some quantity other than the zero-one error that we may wish to optimize?

6 **Surrogate losses** The zero-one loss is the clearest indication of whether the  
 7 perceptron is working well. It is however non-differentiable, so we cannot use  
 8 powerful ideas from optimization theory to minimize it. This is why surrogate  
 9 losses are constructed in machine learning. These are proxies for the loss  
 10 function, typically for the classification problems and look as follows. The  
 11 exponential loss is

$$\ell_{\text{exp}}(w) = e^{-y(w^\top x)}$$

12 or the logistic loss is

$$\ell_{\text{logistic}}(w) = \log(1 + e^{-yw^\top x}).$$

13 **Stochastic Gradient Descent (SGD)** SGD is a very general algorithm to  
 14 optimize objectives typically found in machine learning. We can use it so  
 15 long as we have a dataset and an objective that is differentiable. Consider an  
 16 optimization problem where we want to solve for

$$w^* = \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell^i(w)$$

17 where the function  $\ell^i$  denotes the loss on the sample  $(x^i, y^i)$  and  $w \in \mathbb{R}^p$   
 18 denotes the weights of the classifier. Solving this problem using SGD corre-  
 19 sponds to iteratively updating the weights using

$$w^{t+1} = w^t - \eta \left. \frac{d\ell^{\omega_t}(w)}{dw} \right|_{w=w^t},$$

20 i.e., we compute the gradient one sample with index  $\omega_t$  in the dataset. The  
 21 index  $\omega_t$  is chosen uniformly randomly from

$$\omega_t \in \{1, \dots, n\}.$$

22 In practice, at each time-step  $t$ , we typically select a few (not just one) input  
 23 data  $\omega_t$  from the training dataset and average the gradient  $\left. \frac{d\ell^{\omega_t}(w)}{dw} \right|_{w=w^t}$  across  
 24 them; this is known as a “mini-batch”. The gradient of the loss  $\ell^{\omega_t}(w)$  with  
 25 respect to  $w$  is denoted by

$$\nabla \ell^{\omega_t}(w^t) := \left. \frac{d\ell^{\omega_t}(w)}{dw} \right|_{w=w^t} = \begin{bmatrix} \nabla_{w_1} \ell^{\omega_t}(w^t) \\ \nabla_{w_2} \ell^{\omega_t}(w^t) \\ \vdots \\ \nabla_{w_p} \ell^{\omega_t}(w^t) \end{bmatrix} \in \mathbb{R}^p.$$

1 The gradient  $\nabla \ell^{\omega_t}(w^t)$  is therefore a vector in  $\mathbb{R}^p$ . We have written

$$\nabla_{w_1} \ell^{\omega_t}(w^t) = \left. \frac{d\ell^{\omega_t}(w)}{dw_1} \right|_{w=w^t}$$

2 for the scalar-valued derivative of the objective  $\ell^{\omega_t}(w^t)$  with respect to the  
3 first weight  $w_1 \in \mathbb{R}$ . We can therefore write SGD as

$$w^{t+1} = w^t - \eta \nabla \ell^{\omega_t}(w^t). \quad (7.3)$$

4 The non-negative scalar  $\eta \in \mathbb{R}_+$  is called the step-size or the learning rate. It  
5 governs the distance traveled along the negative gradient  $-\nabla \ell^{\omega_t}(w^t)$  at each  
6 iteration.

## 7 7.1.2 Deep Neural Networks

8 The Perceptron in (7.1) is a linear model: it computes a linear function of  
9 the weights  $w^\top x$  and uses this function to make the predictions  $f(x; w) =$   
10  $\text{sign}(w^\top x)$ . Linear models try to split the data (say we have binary labels  
11  $Y = \{-1, 1\}$ ) using a hyper-plane with  $w$  denoting the normal to this hyper-  
12 plane. This does not work for all situations of course, as the figure below  
13 shows, there is no hyper-plane that cleanly separates the two classes (i.e.,  
14 achieves zero mis-prediction error) but there *is* a nonlinear function that can  
do the job.

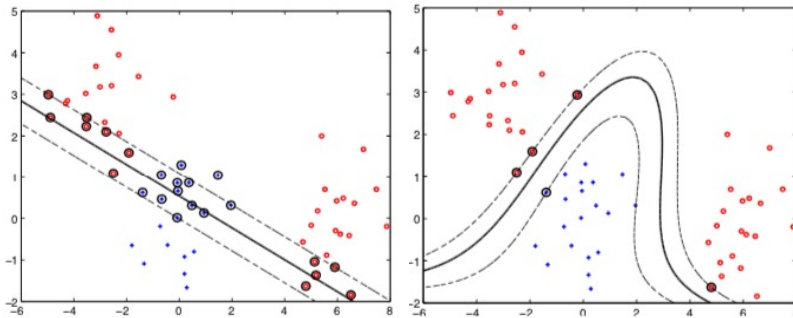


Figure 7.1

15 A deep neural network is one such nonlinear function. First consider a  
16 “two-layer” network  
17

$$f(x; v, S) = \text{sign}(v^\top \sigma(S^\top x))$$

18 where the matrix  $S \in \mathbb{R}^{d \times p}$  and a vector  $v \in \mathbb{R}^p$  are the parameters or  
19 “weights” of the classifier. The “nonlinearity”  $\sigma$  is usually set to be what is  
20 called a Rectified Linear Unit (ReLU)

$$\begin{aligned} \sigma(x) &:= \text{ReLU}(x) = |x|_+ \\ &= \max(0, x). \end{aligned} \quad (7.4)$$

21 Just like the case of a Perceptron, we can use an objective  $\frac{1}{n} \sum_{i=1}^n \ell^i(v, S)$



1 that depends on both  $v, S$  to fit this classifier on training data. A deep neural  
 2 network takes the idea of a two-layer network to the next step and has multiple  
 3 “layers”, each with a different weight matrix  $S_1, \dots, S_L$ . The classifier is  
 4 therefore given by

$$f(x; v, S_1, \dots, S_L) = \text{sign} \left( v^\top \sigma \left( S_L^\top \dots \sigma \left( S_2^\top \sigma \left( S_1^\top x \right) \dots \right) \right) \right). \quad (7.5)$$

5 We call each operation of the form  $\sigma \left( S_k^\top \dots \right)$ , as a *layer*. Consider the  
 6 second layer: it takes the features generated by the first layer, namely  $\sigma \left( S_1^\top x \right)$ ,  
 7 multiplies these features using its feature matrix  $S_2^\top$  and applies a nonlinear  
 8 function  $\sigma(\cdot)$  to this result element-wise before passing it on to the third layer.

A deep network creates new features by composing older features.

9 This composition is very powerful. Not only do we not have to pick a  
 10 particular feature vector, we can create very complex features by sequentially  
 11 combining simpler ones. For example Figure 7.2 shows the features (more  
 12 precisely, the kernel) learnt by a deep neural network. The first layer of features  
 13 are called Gabor-like, and incidentally they are similar to the features learned  
 14 by the human brain in the first part of the visual cortex (the one closest to the  
 15 eyes). These features are *combined* linearly along with a nonlinear operation  
 16 to give richer features (spirals, right angles) in the middle panel. The third  
 17 layer combines the lower features to get even more complex features, these  
 18 look like patterns (notice a soccer ball in the bottom left), a box on the bottom  
 19 right etc.

20 **Deep networks are universal function approximators** The multi-layer  
 21 neural network is a powerful class of classifiers: depending upon how many  
 22 layers we have and what is the dimensionality of the the weight matrices  
 23  $S_k$  at each layer, we can fit *any* training data. In fact, this statement, which  
 24 is called the *universal approximation property* holds even for a two-layer  
 25 neural network  $v^\top \sigma \left( S^\top x \right)$  if the number of columns in  $S$  is big enough. This  
 26 property is the central reason why deep networks are so widely applicable, we  
 27 can model complex machine learning problems if we choose a big enough  
 28 deep network.

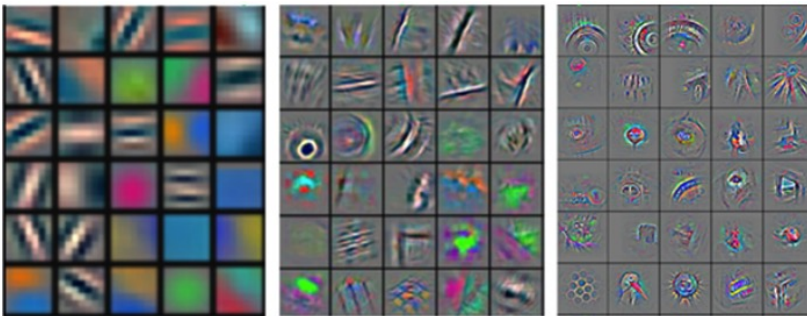


Figure 7.2

1 **Logits for multi-class classification.** The output

$$\hat{y} = v^\top \sigma(S_L^\top \dots \sigma(S_2^\top \sigma(S_1^\top x)) \dots)$$

2 is called the logits corresponding to the different classes. This name comes  
3 from logistic regression where logits are the log-probabilities of an input datum  
4 belonging to one of the two classes. A deep network provides an easy way to  
5 solve a multi-class classification problem, we simply set

$$v \in \mathbb{R}^{p \times C}$$

6 where  $C$  is the total number of classes in the data. Just like logistic regression  
7 predicts the logits of the two classes, we would like to *interpret* the vector  $\hat{y}$  as  
8 the log-probabilities of an input belonging to one of the classes.

9 **Weights** It is customary to not differentiate between the parameters of dif-  
10 ferent layers of a deep network and simply say *weights* when we want to refer  
11 to all parameters. The set

$$w := \{v, S_1, S_2, \dots, S_L\}$$

12 is the set of *weights*. This set is typically stored in PyTorch as a set of matrices,  
13 one for each layer. Using this new notation, we will write down a deep neural  
14 network classifier as simply

$$f(x, w) \tag{7.6}$$

15 and fitting the deep network to a dataset involves the optimization problem

$$w^* = \operatorname{argmin}_w \frac{1}{n} \sum_{i=1}^n \ell(y^i, \hat{y}^i). \tag{7.7}$$

16 We will also sometimes denote the loss of the  $i^{\text{th}}$  sample as

$$\ell^i(w) := \ell(y^i, \hat{y}^i).$$

17 **Backpropagation** The Backpropagation algorithm is a method to compute  
18 the gradient of the objective while fitting a deep network using SGD, i.e., it  
19 computes  $\nabla_w \ell^i(w)$ . For the purposes of this course, the details of how this is  
20 done are not essential, so we will skip them. You can read more in the notes  
21 of ESE 546 at [https://pratikac.github.io/pub/20\\_ese546.pdf](https://pratikac.github.io/pub/20_ese546.pdf).

22 **PyTorch** We will use a library called PyTorch (<https://pytorch.org>) to code  
23 up deep neural networks for the reinforcement learning part of this course.  
24 You can find some excellent tutorials for it at  
25 <https://pytorch.org/tutorials/beginner/basics/intro.html>. We have also uploaded  
26 two recitations from the Fall 2020 offering of ESE 546 on Canvas which guide  
27 you through various typical use-cases of PyTorch. You are advised to go  
28 through, at least, the first recitation if you are not familiar with PyTorch. For  
29 the purposes of this course, you do not need to know the intricacies of PyTorch,

🔗 What would the shape of  $w$  be if you were performing regression using a deep network?

1 we will give you enough code to work with deep networks so that you can  
2 focus on implementing the reinforcement learning-specific parts.

## 3 7.2 Behavior Cloning

4 With that background, we are ready to tackle what is potentially the simplest  
5 problem in RL. We will almost exclusively deal with discrete-time systems for  
6 RL. Let us imagine that we are given access to  $n$  trajectories each of length  
7  $T + 1$  time-steps from an expert demonstrator for our system. We write this  
8 as a training dataset

$$D = \{(x_t^i, u_t^i)_{t=0,1,\dots,T}\}_{i=1,\dots,n}$$

9 At each step, we record the state  $x_t^i \in \mathbb{R}^d$  and the control that the expert took  
10 at that state  $u_t^i$ . We would like to learn a deterministic feedback control for the  
11 robot that is parametrized by parameters  $\theta$

$$u_\theta(x) : X \mapsto U \subset \mathbb{R}^m.$$

12 using the training data. The idea is that if  $u_\theta(x^i(t)) \approx u^i(t)$  for all  $i$  and all  
13 times  $t$ , then we can simply run our learned controller  $u_\theta(x)$  on the robot  
14 instead of having the expert. A simple example is a baby deer learning to  
15 imitate how its mother in how to run.

16 **Parameterizing the controller** Our function  $u_\theta$  may represent many differ-  
17 ent families of controllers. For example,  $u_\theta(x) = \theta x$  where  $\theta \in \mathbb{R}^{d \times p}$  is a  
18 linear controller; this is much like the control for LQR except that we can fit  $\theta$   
19 to the expert's data instead of solving the LQR problem to find the Kalman  
20 gain. We could also think of some other complicated function, e.g., a two-layer  
21 neural network,

$$u_\theta(x) = v \sigma(S^\top x)$$

22 where  $S \in \mathbb{R}^{d \times p}$  and  $v \in \mathbb{R}^{m \times p}$  and  $\sigma : \mathbb{R}^m \mapsto \mathbb{R}^m$  is some nonlinearity, say  
23 ReLU. As we did above, we will use

$$\theta := (v, S)$$

24 to denote all the weights of this two-layer neural network. Multi-layer neural  
25 networks are also another possible avenue. In general, we want to the parame-  
26 terization of the controller to be rich enough to fit some complex controller  
27 that the expert may have used on the system.

28 **How to fit the controller?** Given our chosen model for  $u_\theta(x)$ , say a two-  
29 layer neural network with weights  $\theta$ , fitting the controller involves finding the  
30 best value for the parameters  $\theta$  such that  $u_\theta(x_t^i) \approx u_t^i$  for data in our dataset.  
31 There are many ways to do this, e.g., we can solve the following optimization

1 problem

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \ell(\theta) := \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{T+1} \sum_{t=0}^T \|u_t^i - u_{\theta}(x_t^i)\|_2^2}_{\ell^i(\theta)} \quad (7.8)$$

2 The difficulty of solving the above problem depends upon how difficult the  
 3 model  $u_{\theta}(x)$  is, for instance, if the model is linear  $\theta x$ , we can solve (7.8)  
 4 using ordinary least squares. If the model is a neural network, one would have  
 5 to use SGD to solve the optimization problem above. After fitting this model,  
 6 we have a new controller

$$u_{\hat{\theta}}(x) \in \mathbb{R}^m$$

7 that we can use *anywhere* in the domain  $X \subset \mathbb{R}^d$ , even at places where we had  
 8 no expert data. This is known as Behavior Cloning, i.e., cloning the controls  
 9 of the expert into a parametric model.

10 **Generalization performance of behavior cloning** Note that the data pro-  
 11 vided by the expert is not iid, of course the state  $x_{t+1}^i$  in the expert's trajectory  
 12 depends upon the previous state  $x_t^i$ . Standard supervised learning makes the  
 13 assumption that Nature gives training data that is independent and identically  
 14 distributed from the distribution  $P$ . While it is still reasonable to fit the re-  
 15 gression loss in (7.8) for such correlated data, one should remember that if  
 16 the expert trajectories do not go to all parts of the state-space, the learned  
 17 controller fitted on the training data may not work outside these parts. Of  
 18 course, if we behavior clone the controls taken by a generic driver, they are  
 19 unlikely to be competitive for racing, and vice-versa. It is very important to  
 20 realize that this does *not* mean that BC does not generalize. Generalization in  
 21 machine learning is a concept that suggests that the model should work well  
 22 on data *from the same distribution*. What does the the “distribution” of the  
 23 expert mean, in this case, it simply refers to the distribution of the states that  
 24 the expert's trajectories typically visit, e.g, a race driver typically drives at the  
 25 limits of tire friction and throttle, this is different from a usual city-driver who  
 26 would rather maximize the longevity of their tires and engine-life.

❗ Discuss generalization performance in behavior cloning.

## 27 7.2.1 Behavior cloning with a stochastic controller

28 In this case, we have always chosen feedback feedback controllers that are  
 29 deterministic, i.e., there is a single value of control  $u$  that is taken at the state  $x$ .  
 30 Going forward, we will also talk about stochastic controllers, i.e., controllers  
 31 which sample a control from a distribution. There can be a few reasons of  
 32 using such a controller. First, we will see in later lectures how this may help  
 33 in training a reinforcement learning algorithm; this is because in situations  
 34 where you do not know the system dynamics precisely, it helps to “hedge” the  
 35 feedback to take a few different control actions instead of simply the one that  
 36 the value function deems as the maximizing one. This is not very different  
 37 from having a few different stocks in your portfolio. Second, we benefit from  
 38 this hedging even at test-time when we run a stochastic feedback control, e.g.,

1 in situations where the limited training data may not want to always pick the  
 2 best control (because the best control was computed using an imprecise model  
 3 of the system dynamics and could be wrong), but rather hedge our bets by  
 4 choosing between a few different controls.

5 A stochastic feedback control is denoted by

$$u \sim u_\theta(\cdot | x) = \mathbf{P}(\cdot | x)$$

6 notice that  $u_\theta(\cdot | x)$  is a probability distribution on the control space  $U$  that  
 7 depends on the state  $x$ , and in this case the parameters  $\theta$ . The control taken at a  
 8 state  $x$  is a sample drawn from this probability distribution. The deterministic  
 9 controller is a special case of this setup where

$$u_\theta(u | x) = \delta_{u_\theta(x)}(u) \equiv u_\theta(x)$$

10 is a Dirac-delta distribution at  $u_\theta(x)$ . If the control space  $U$  is discrete,  
 11 then  $u_\theta(\cdot | x)$  could be a categorical distribution. If the control space  $U$  is  
 12 continuous, then you may wish to think of the controls being sampled from a  
 13 Gaussian distribution with some mean  $\mu_\theta(x)$  and variance  $\sigma_\theta^2(x)$

$$\mathbb{R}^m \ni u \sim u_\theta(\cdot | x) = N(\mu_\theta(x), \Sigma_\theta(x)).$$

14 **Maximum likelihood estimation** Let's pick a particular stochastic con-  
 15 troller, say a Gaussian. How should we fit the parameters  $\theta$  for this? We would  
 16 like to find parameters  $\theta$  that make the expert's data in our dataset very likely.  
 17 The log-likelihood of each datum is

$$\log u_\theta(u_t^i | x_t^i)$$

18 and maximizing the log-likelihood of the entire dataset amounts to solving

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{T+1} \sum_{t=0}^T -\log u_\theta(u_t^i | x_t^i)}_{\ell^i(\theta)}. \quad (7.9)$$

19 **Fitting BC with a Gaussian controller** Notice that if we use a Gaussian  
 20 distribution

$$u_\theta(\cdot | x) = N(\mu_\theta(x), I)$$

21 as our stochastic controller, the objective in (7.9) is the same as that in (7.8).

$$u_\theta(\cdot | x) = N(\mu_\theta(x), \sigma_\theta^2(x)I)$$

22 we have that

$$-\log u_\theta(u | x) = \frac{\|\mu_\theta(x) - u\|_2^2}{\sigma_\theta^2(x)} + 2cp \log \sigma_\theta(x).$$

23 where  $c$  is a constant.

## 7.2.2 KL-divergence form of Behavior Cloning

**Background on KL divergence** The Kullback-Leibler (KL) divergence is a quantity to measure the distance between two probability distributions. There are many similar distances, for example, given two probability distributions  $p(x)$  and  $q(x)$  supported on a discrete set  $X$ , the total variation distance between them is

$$\text{TV}(p, q) = \frac{1}{2} \sum_{x \in X} |p(x) - q(x)|.$$

Hellinger distance ([https://en.wikipedia.org/wiki/Hellinger\\_distance](https://en.wikipedia.org/wiki/Hellinger_distance)),  $f$ -divergences (<https://en.wikipedia.org/wiki/F-divergence>) and the Wasserstein metric ([https://en.wikipedia.org/wiki/Wasserstein\\_metric](https://en.wikipedia.org/wiki/Wasserstein_metric)) are a few other examples of ways to measure how different two probability distributions are from each other.

The Kullback-Leibler divergence (KL) between two distributions is given by

$$\text{KL}(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}. \quad (7.10)$$

This is a distance and not a metric, i.e., it is always non-zero and zero if and only if the two distributions are equal, but the KL-divergence is not symmetric (like a metric has to be). Also, the above formula is well-defined only if for all  $x$  where  $q(x) = 0$ , we also have  $p(x) = 0$ . Notice that it is not symmetric

$$\text{KL}(q \parallel p) = \sum_{x \in X} q(x) \log \frac{q(x)}{p(x)} \neq \text{KL}(p \parallel q).$$

The funny notation  $\text{KL}(p \parallel q)$  was invented by Shun-ichi Amari ([https://en.wikipedia.org/wiki/Shun%27ichi\\_Amari](https://en.wikipedia.org/wiki/Shun%27ichi_Amari)) to emphasize the fact that the KL-divergence is asymmetric. The KL-divergence is always positive: you can show this using an application of Jensen's inequality. For distributions with continuous support, we integrate over the entire space  $X$  and define KL divergence as

$$\text{KL}(p \parallel q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx.$$

**Behavior Cloning** Let us now imagine the expert is also a parametric stochastic feedback controller  $u_{\theta^*}(\cdot \mid x)$ . Our data is therefore drawn by running this controller for  $n$  trajectories,  $T$  time-steps on the system. This dataset now consists of samples from

$$p_{u_{\theta^*}}(x, u)$$

which is the joint distribution on the state-space  $X$  and the control-space  $U$ . We have denoted the parameters of the feedback controller which creates this distribution as the subscript  $u_{\theta^*}$ . Our behavior cloning controller creates a similar distribution  $p_{u_{\theta}}(x, u)$  and the general version of the objective in (7.9)

1 is therefore

$$\hat{\theta} = \operatorname{argmin}_{\theta} \operatorname{KL}(p_{u_{\theta^*}} \parallel p_{u_{\theta}}); \quad (7.11)$$

2 The objective in (7.9) corresponds to this for Gaussian stochastic controllers,  
3 but we can just as easily imagine some other distribution for the stochastic  
4 controller of the expert and the robot.

Written this way, BC can be understood as finding a controller  $\hat{\theta}$  whose distribution on the states and controls is close to the distribution of states and controls of the expert.

### 5 7.2.3 Some remarks on Behavior Cloning

6 **Worst-case performance** Performance of Behavior Cloning can be quite  
7 bad in the worst case. The authors in “Efficient reductions for imitation learn-  
8 ing” (<https://www.cs.cmu.edu/~sross1/publications/Ross-AIStats11-NoRegret.pdf>)  
9 show that if the learned controller  $u_{\hat{\theta}}$  differs from the control taken by the  
10 expert controller  $u_{\theta^*}$  with a probability  $\epsilon$  at each time-step, over a horizon of  
11 length  $T$  time-steps, it can be  $\mathcal{O}(T^2\epsilon)$  off from the cost-to-go of the expert *as*  
12 *averaged over states that the learned controller visits*. This is because once the  
13 robot makes a mistake and goes away from the expert’s part in the state-space,  
14 future states of the robot and the expert can be very different.

15 **Model-free nature of BC** Observe that our learned controller  $u_{\hat{\theta}}(\cdot | x)$  is a  
16 feedback controller and works for entire state-space  $X$ . We did not need to  
17 know the dynamics of the system to build this controller. The data from the  
18 expert is conceptually the same as the model  $\dot{x} = f(x, u)$  of the dynamics,  
19 and you can learn controllers from both. Do you however notice a catch?

📌 Draw a picture of the amplifying errors of running behavior cloning in real-time.

## 20 7.3 DAgger: Dataset Aggregation

21 The expert’s dataset in Behavior Cloning determines the quality of the con-  
22 troller learned. If we collected very few trajectories from the expert, they may  
23 not cover all parts of the state-space and the behavior cloned controller has no  
24 data to fit the model in those parts.

25 Let us design a simple algorithm, of the same spirit as iterative-LQR, to  
26 mitigate this. We start with a candidate controller, say  $u_{\theta^{(0)}}(x)$ ; one may also  
27 start with a stochastic controller  $u_{\theta^{(0)}}(\cdot | x)$  instead.

**DAgger:** Let the dataset  $D^{(0)}$  be the data collected from the expert.

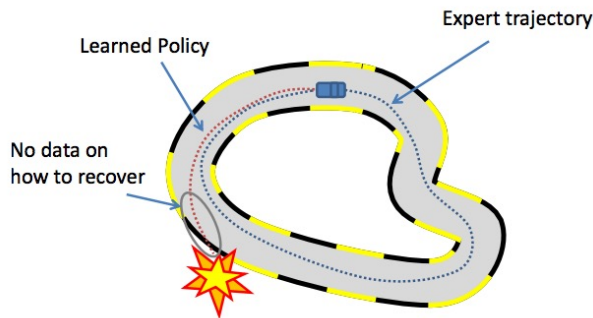
Initialize  $u_{\theta^{(0)}} = u_{\hat{\theta}}$  to be the BC controller learned using data  $D^{(0)}$ . At iteration  $k$

1. The robot queries the expert for a fraction  $p$  of the time-steps and uses its learned controller  $u_{\theta^{(k-1)}}$  for the other time-steps. If the expert corresponds to some controller  $u_{\theta^*}$ , then the robot controller at a state  $x$  is

$$u \sim p \delta_{u_{\theta^*}(x)} + (1 - p) \delta_{u_{\theta^{(k-1)}}(x)}.$$

2. Use  $u(x)$  to collect a dataset  $D = \{(x_t^i, u_t^i)_{t=0, \dots, T}\}_{i=1, \dots, n}$  with  $n$  trajectories.
3. Set the new dataset to be  $D^{(k)} = D^{(k-1)} \cup D$
4. Fit a controller  $u_{\theta^{(k)}}$  using behavior cloning to the new dataset  $D^{(k)}$ .

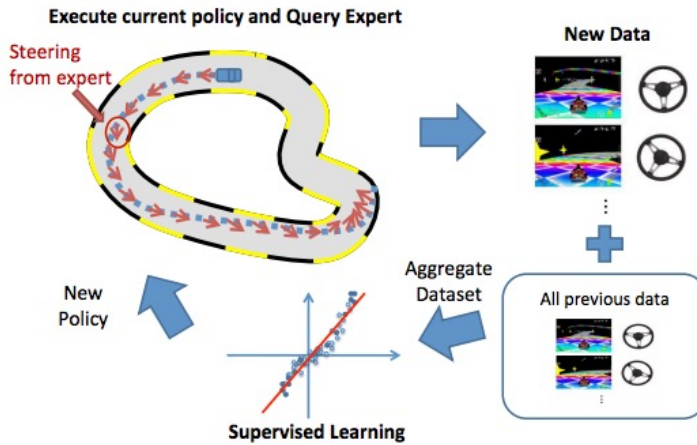
1 The above algorithm iteratively updates the BC controller  $u_{\hat{\theta}}$  by drawing  
 2 new data from the expert. The robot first bootstraps off the expert's data, this  
 3 simply means that it uses the expert's data to fit its controller  $u_{\theta^{(0)}}(x)$ . As we  
 4 discussed above, this controller may veer off the expert's trajectory if the robot  
 5 starts at states that are different from the dataset, or even if it takes a slightly  
 6 different control than the expert midway through a trajectory.



7

8 To fix this, the robot collects more data at each iteration. It uses a combina-  
 9 tion of the expert and its controller to collect such data. This, allows *collecting*  
 10 *a dataset of expert's controls in states that the robot visits* and iteratively  
 11 expands the dataset  $D^{(k)}$ .





1

2 In the beginning we may wish to be close to the expert's data and use a large  
 3 value of  $p$ , as the fitted controller  $u_{\theta_{k+1}}$  becomes good, we can reduce the  
 4 value of  $p$  and rely less on the expert.

5 DAgger is an iterative algorithm which expands the controller to handle  
 6 larger and larger parts of the state-space. Therefore, the cost-to-go of the  
 7 controller learned via DAgger is  $\mathcal{O}(T)$  off from the cost-to-go of the expert as  
 8 averaged over states that the learned controller visits.

9 **DAgger with expert annotations at each step** DAgger is a conceptual  
 10 framework where the expert is queried repeatedly for new control actions.  
 11 This is obviously problematic because we need to expert on hand at each  
 12 iteration. We can also cook up a slightly version of DAgger where we start  
 13 with the BC controller  $u_{\theta^{(k)}} = u_{\hat{\theta}}$  and at each step, we run the controller on  
 14 the real system and ask the expert to relabel the data after that run. The dataset  
 15  $D^{(k)}$  collected by the algorithm expands at each iteration and although the  
 16 states  $x_t^i$  are those visited by our controller, their annotations are those given  
 17 by the expert. This is a much more natural way of implementing DAgger.

❓ What criterion can we use to stop these iterations? We can stop when the incremental dataset collected  $D_k$  is not that different from the cumulative dataset  $D$ , we know that the new controllers are not that different. We can also stop when the parameters of our learned controller are  $\theta^{(k+1)} \approx \theta^{(k)}$ .

# Chapter 8

## Policy Gradient Methods

### Reading

1. Sutton & Barto, Chapter 9–10, 13
2. Simple random search provides a competitive approach to reinforcement learning at <https://arxiv.org/abs/1803.07055>
3. Proximal Policy Optimization Algorithms <https://arxiv.org/abs/1707.06347>
4. Are Deep Policy Gradient Algorithms Truly Policy Gradient Algorithms? <https://arxiv.org/abs/1811.02553>
5. Asynchronous Methods for Deep Reinforcement Learning <http://proceedings.mlr.press/v48/mnih16.pdf>

This chapter discusses methods to learn the controller that minimizes a given cost functional over trajectories of an unknown dynamical system. We will use what is called the “policy gradient” which will be the main section of this chapter.

Recall from the last chapter that we were able to fit stochastic controllers of the form  $u_{\hat{\theta}}(\cdot | x)$  that is a probability distribution on the control-space  $U$  for each  $x \in X$ . We fitted  $u_{\theta}$  using data from the expert in imitation learning. We did not learn the cost-to-go for the fitted controller, like we did in the lectures on dynamic programming. This is a clever choice: it is often easier to learn the controller in a typical problem than to *compute* the optimal cost-to-go as a parametric function  $J^*(x)$ .

🔗 Can you give another instance when we have computed a controller previously in the class without coming up with its cost-to-go?

### 8.1 Standard problem setup in RL

**Dynamics and rewards** In this and the next few chapters we will always consider discrete-time stochastic dynamical systems with a stochastic con-

1 troller with parameters (weights)  $\theta$ . We denote them as follows

$$\begin{aligned} x_{k+1} &\sim p(\cdot | x_k, u_k) \text{ with noise denoted by } \epsilon_k \\ u_k &\sim u_\theta(x_k). \end{aligned}$$

2 We will also change perspective and instead of minimizing the infinite-horizon  
3 sum of a runtime cost, maximize the sum of a runtime reward

$$r(x, u) := -q(x, u).$$

4 We do so simply to conform to tradition and standard notation in reinforcement  
5 learning; the two are mathematically completely equivalent. We are interested  
6 in maximizing the expected value of the cumulative rewards over infinite-  
7 horizon trajectories of the system

$$J(\theta) = \mathbb{E}_{x_1, x_2, \dots} \left[ \underbrace{\sum_{k=0}^{\infty} \gamma^k r(x_k, u_k)}_{\text{discounted return}} \mid x_0 \right]; \quad (8.1)$$

8 where each  $u_k \sim u_\theta(\cdot | x_k)$  and each  $x_{k+1} \sim p(\cdot | x_k, u_k)$ .

9 **Trajectory space** Let us write out one trajectory of such a system a bit  
10 more explicitly. We know that the probability of the next state  $x_{k+1}$  given  
11  $x_k$  is  $p(x_{k+1} | x_k, u_k)$ . The probability of taking a control  $u_k$  at state  $x_k$  is  
12  $u_\theta(u_k | x_k)$ . We denote an infinite trajectory by

$$\tau = x_0, u_0, x_1, u_1, \dots$$

13 The probability this entire trajectory occurring is

$$p_\theta(\tau) = \prod_{k=0}^{\infty} p(x_{k+1} | x_k, u_k) u_\theta(u_k | x_k);$$

14 we have emphasized that the distribution of trajectories depends on the weights  
15 of controller  $\theta$ . If we take the logarithm,

$$\log p_\theta(\tau) = \sum_{k=0}^{\infty} \log p(x_{k+1} | x_k, u_k) + \log u_\theta(u_k | x_k).$$

16 Given a trajectory  $\tau = x_0, u_0, x_1, u_1, \dots$ , the sum

$$R(\tau) = \sum_{k=0}^{\infty} \gamma^k r(x_k, u_k) \quad (8.2)$$

17 is called the discounted return of the trajectory  $\tau$ . Sometimes we will also talk  
18 of the undiscounted return of the trajectory which is the sum of the rewards  
19 up to some fixed finite horizon  $T$  without the discount factor pre-multiplier.

1 Using this notation, we can write out objective from (8.1) as

$$J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) \mid x_0] \quad (8.3)$$

2 where  $p(\tau)$  is the probability distribution of an infinitely long trajectory  $\tau$ .

3 Observe what is probably the most important point in policy-gradient based  
4 reinforcement learning: the probability of trajectory is an infinite product of  
5 terms all of which are smaller than 1 (they are probabilities), so it is essentially  
6 zero even if the state-space and the control-space are finite (even if they are  
7 small). Any given infinite (or long) trajectory is quite rare under the probability  
8 distribution of the stochastic controller. Policy-gradient methods sample lots  
9 of trajectories from the system and average the returns across these trajectories.  
10 Since the set of trajectories of even a small MDP is so large, sampling lots  
11 of trajectories, or even the most likely ones, is also very hard. This is a key  
12 challenge in getting RL algorithms to work.

**Our goal in this chapter** is to compute the best stochastic controller which maximizes the average discounted return. Mathematically, this amounts to finding

$$\hat{\theta} = \operatorname{argmax}_{\theta} J(\theta) := \mathbb{E}_{\tau \sim p(\tau)} [R(\tau) \mid x_0]. \quad (8.4)$$

The objective  $J(\theta)$  is called the average return of the controller  $u_\theta$ .

13 **Computing the average return  $J(\theta)$**  Before we move on to optimizing  
14  $J(\theta)$ , let us discuss how to compute it for given weights  $\theta$  of the stochastic  
15 controller. We can sample  $n$  trajectories from the system and compute the an  
16 estimate of the expectation

$$\hat{J}(\theta) \approx \frac{1}{n} \sum_{i=0}^n \sum_{k=0}^T \gamma^k r(x_k^i, u_k^i) \quad (8.5)$$

17 for some large time-horizon  $T$  and where each  $u_k^i \sim u_\theta(\cdot \mid x_k^i)$ .

## 18 8.2 Cross-Entropy Method (CEM)

19 Let us first consider a simple method to compute the best controller. The basic  
20 idea is to solve the problem

$$\hat{\theta} = \operatorname{argmax}_{\theta} J(\theta)$$

21 using gradient descent. We would like to update weights  $\theta$  iteratively

$$\theta^{k+1} = \theta^k + \eta \nabla J(\theta).$$

22 where the step-size is  $\eta > 0$  and  $\nabla J(\theta)$  is the gradient of the objective  $J(\theta)$   
23 with respect to weights  $\theta$ . Instead of computing the exact  $\nabla J(\theta)$  which

🔗 Contrast (8.5) with the complexity of policy evaluation which was simply a system of linear equations. Evaluating the policy without having access to the dynamical system is harder.

1 we will do in the next section, let us simply compute the gradient using a  
2 finite-difference approximation. The  $i^{\text{th}}$  entry of the gradient is

$$(\widehat{\nabla} J(\theta))_i = \frac{J(\theta + \epsilon e_i) - J(\theta - \epsilon e_i)}{2\epsilon} \approx \frac{\widehat{J}(\theta + \epsilon e_i) - \widehat{J}(\theta - \epsilon e_i)}{2\epsilon}.$$

3 where  $e_i = [0, 0, \dots, 0, 1, 0, \dots]$  is a vector with 1 on the  $i^{\text{th}}$  entry. Each  
4 quantity  $\widehat{J}$  is computed as the empirical average return of  $n$  trajectories from  
5 the system. We compute all entries of the objective using this approximation  
6 and update the parameters using

$$\theta^{k+1} = \theta^k + \eta \widehat{\nabla} J(\theta^k).$$

7 **A more efficient way to compute the gradient using finite-differences** In-  
8 stead of picking perturbations  $e_i$  along the cardinal directions, let us sample  
9 them from a Gaussian distribution

$$\xi^i \sim N(0, \sigma^2 I)$$

10 for some user-chosen covariance  $\sigma^2$ . We can however no longer use the finite-  
11 difference formula to compute the derivative because the noise  $e$  is not aligned  
12 with the axes. We can however use a Taylor series approximation as follows.  
13 Observe that

$$J(\theta + \xi) \approx J(\theta) + \langle \nabla J(\theta), \xi \rangle$$

14 where  $\langle \cdot, \cdot \rangle$  is the inner product. Given  $m$  samples  $\xi^1, \dots, \xi^m$  observe that

$$\begin{aligned} \widehat{J}(\theta + \xi^1) &= \widehat{J}(\theta) + \langle \nabla J(\theta), \xi^1 \rangle \\ \widehat{J}(\theta + \xi^2) &= \widehat{J}(\theta) + \langle \nabla J(\theta), \xi^2 \rangle \\ &\vdots \\ \widehat{J}(\theta + \xi^m) &= \widehat{J}(\theta) + \langle \nabla J(\theta), \xi^m \rangle. \end{aligned} \tag{8.6}$$

15 is a linear system of equations in  $\nabla J(\theta) \in \mathbb{R}^p$  where  $\theta \in \mathbb{R}^p$ . All quantities  
16  $\widehat{J}$  are estimated as before using trajectories drawn from the system. We solve  
17 this linear system, e.g., using least-squares if  $m > p$ , to get an estimate of the  
18 gradient  $\widehat{\nabla} J(\theta)$ .

**The Cross-Entropy Method** is a more crude but simpler to implement version of the above least-squares formulation. At each iteration it updates the parameters using the formula

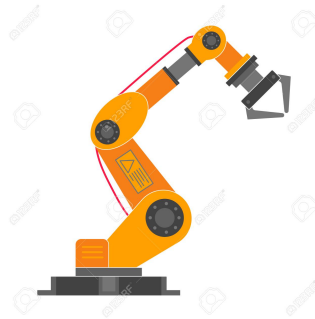
$$\theta^{k+1} = \mathbb{E}_{\theta \sim N(\theta^k, \sigma^2 I)} \left[ \theta \mathbf{1}_{\{\widehat{J}(\theta) > \widehat{J}(\theta^k)\}} \right]. \tag{8.7}$$

In simple words, the CEM samples a few stochastic controllers  $u_\theta$  from a Gaussian (or any other distribution) centered around the current controller  $u_{\theta^k}$  and updates the weights  $\theta^k$  in a direction that leads to an increase in  $\widehat{J}(\theta) > \widehat{J}(\theta^k)$ .

## 8.2.1 Some remarks on sample complexity of simulation-based methods

The CEM may seem to be a particularly bad method to maximize  $J(\theta)$ , after all we are perturbing the weights of the stochastic controller randomly and updating the weights if they result in a better average return  $\hat{J}(\theta)$ . This is likely to work well if the dimensionality of weights  $\theta \in \mathbb{R}^p$ , i.e.,  $p$ , is not too large. But is unlikely to work well if we are sampling  $\theta$  in high-dimensions. Typical applications are actually the latter, remember that we are interested in using a deep network as a stochastic controller and  $\theta$  are the weights of the neural networks.

Let us do a quick computation, if the state is  $x \in \mathbb{R}^d$  and  $u \in \mathbb{R}^m$  with  $d = 12$  (joint angles and velocities) and  $m = 6$  for a six-degree of freedom robot manipulator



and if we use a two-layer neural network with 64 neurons in the hidden layer, the total number of weights  $\theta \in \mathbb{R}^p$  for the function  $u_\theta(\cdot | x) = N(\mu_\theta(x), \sigma_\theta^2(x)I)$  where  $\sigma^2(x)$  is a vector in  $\mathbb{R}^m$ , is

$$p = \underbrace{(12 \times 64 + 64) + (64 \times 6 + 6)}_{\text{for } \mu_\theta(x)} + \underbrace{(12 \times 64 + 64) + (64 \times 6 + 6)}_{\text{for } \sigma_\theta^2(x)}$$

$$= 2,444.$$

This is a very high-dimensional space to sample exhaustively, and it is quite large even if the input and output dimensions of the neural network are not too large. To appreciate the complexity of computing the gradient  $\nabla J(\theta)$ , let us think of how to compute it using finite-differences, we need two estimates  $\hat{J}(\theta - \epsilon e_i)$  and  $\hat{J}(\theta + \epsilon e_i)$  for every dimension  $i \in \{1, \dots, p\}$ . Each estimate requires us to obtain  $n$  trajectories from the system. Since the number of trajectories that a robot can take is quite diverse, we should use a large  $n$ , so let's pick  $n = 100$ . The total number of trajectories required to update the parameters  $\theta^k$  at each iteration is

$$2 p n \approx 10^6.$$

This is an absurdly large number, and things are even more daunting when we realize that each update of the weights requires us to sample these many trajectories from the system. It is not reasonable to sample such a large number of trajectories from a physical robot, that too for each update of the weights.

**i** For comparison, a busy espresso bar in a city makes about 500 shots per day. The espresso machine would have to work for 5 years without breaking down to make  $10^6$  shots.

1 **Using fast simulators for RL** If we expand our horizon and think of learn-  
 2 ing controllers in simulation, things feel much more reasonable. While run-  
 3 ning a large number of trajectories may degrade a robot beyond use, doing  
 4 so requires just computation time in a robot simulator. There is a large num-  
 5 ber of simulators that are available with various capabilities, e.g., Gazebo  
 6 (<http://gazebo.org>) is a sophisticated simulator inside ROS that uses a num-  
 7 ber of Physics engines such as Bullet (<https://pybullet.org/wordpress>), MuJoCo  
 8 (<http://www.mujo.co>) is incredibly fast although not very good modeling  
 9 contact, Unity is a popular platform to simulate driving and complicated scenes  
 10 (<https://docs.nvidia.com/isaac/isaac/doc/simulation/unity3d.html>), Drake  
 11 (<https://drake.mit.edu>) is better at contact modeling but more complex and  
 12 slower; most autonomous driving companies have developed their own driving  
 13 simulators in-house. The assigned reading (#2) for this chapter is a paper  
 14 which develops a very fast implementation of the CEM for use in simulators.

15 **Working well in simulation does not mean that a controller works well**  
 16 **on the real robot** It is important to realize that a simulator is not equivalent  
 17 to the physical robot. Each simulator makes certain trade-offs in capturing  
 18 the dynamics of the real system and it is not a given that a controller that  
 19 was learned using data from a simulator will work well on a real robot. For  
 20 instance, OpenAI had to develop a large number of tricks (which took about  
 21 a year) to modify the simulator to enable the learned policy to work well on  
 22 a robot (<https://openai.com/blog/learning-dexterity>) for a fairly narrow set of  
 23 tasks.

## 24 8.3 The Policy Gradient

In this section, we will study how to take the gradient of the objective  $J(\theta)$ , without using finite-differences.

25 We would like to solve the optimization problem

$$\max_{\theta} J(\theta) := \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau) \mid x_0]$$

26 We will suppress the dependence on  $x_0$  to keep the notation clear. The  
 27 expectation is taken over all trajectories starting at state  $x_0$  realized using  
 28 the stochastic controller  $u_{\theta}(\cdot \mid x)$ . We to update weights  $\theta$  using gradient  
 29 descent which amounts to

$$\theta^{k+1} = \theta^k + \eta \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau)].$$

30 First let us note that the distribution  $p_{\theta}$  using which we compute the expectation  
 31 also depends on the weights  $\theta$ . This is why we cannot simply move the  
 32 derivative  $\nabla_{\theta}$  inside the expectation

$$\nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau)] \neq \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} R(\tau)].$$

1 We need to think of a new technique to compute the gradient above. Essentially,  
 2 we would like to do the chain rule of calculus but where one of the functions  
 3 in the chain is an expectation. The likelihood-ratio trick described next allows  
 4 us to take such derivatives. Here is how the computation goes

$$\begin{aligned}
 \nabla_{\theta} \mathbf{E}_{\tau \sim p_{\theta}} [R(\tau)] &= \nabla_{\theta} \int R(\tau) p_{\theta}(\tau) d\tau \\
 &= \int R(\tau) \nabla_{\theta} p_{\theta}(\tau) d\tau \\
 &\quad \text{(move the gradient inside, integral is over trajectories } \tau \text{ which do not depend on } \theta \text{ themselves)} \\
 &= \int R(\tau) p_{\theta}(\tau) \frac{\nabla p_{\theta}(\tau)}{p_{\theta}(\tau)} d\tau \\
 &= \int R(\tau) p_{\theta}(\tau) \nabla \log p_{\theta}(\tau) d\tau \\
 &= \mathbf{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau) \nabla \log p_{\theta}(\tau)] \\
 &\approx \frac{1}{n} \sum_{i=1}^n R(\tau^i) \nabla \log p_{\theta}(\tau^i)
 \end{aligned} \tag{8.8}$$

5 This is called the likelihood-ratio trick to compute the policy gradient. It simply  
 6 multiplies and divides by the term  $p_{\theta}(\tau)$  and rewrites the term  $\frac{\nabla p_{\theta}}{p_{\theta}} = \nabla \log p_{\theta}$ .  
 7 It gives us a neat way to compute the gradient: we sample  $n$  trajectories  
 8  $\tau^1, \dots, \tau^n$  from the system and average the return of each trajectory  $R(\tau^i)$   
 9 weighted by the gradient of the likelihood of taking each trajectory  $\log p_{\theta}(\tau^i)$ .  
 10 The central point to remember here is that the gradient

$$\begin{aligned}
 \nabla_{\theta} \log p_{\theta}(\tau^i) &= \nabla_{\theta} \sum_{k=0}^T \log p(x_{k+1}^i | x_k^i, u_k^i) + \log u_{\theta}(u_k^i | x_k^i) \\
 &= \sum_{k=0}^T \nabla_{\theta} \log u_{\theta}(u_k^i | x_k^i)
 \end{aligned} \tag{8.9}$$

11 is computed using backpropagation for a neural network. This expression is  
 12 called the policy gradient because it is the gradient of the objective  $J(\theta)$  with  
 13 respect to the parameters of the controller/policy  $\theta$ .

14 **Variance of policy gradient** The expression for the policy gradient may  
 15 seem like a sleight of hand. It is a clean expression to get the gradient of the  
 16 objective but also comes with a number of problems. Observe that

$$\begin{aligned}
 \nabla_{\theta} \mathbf{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau)] &= \mathbf{E}_{\tau \sim p_{\theta}(\tau)} \left[ R(\tau) \frac{\nabla p_{\theta}(\tau)}{p_{\theta}(\tau)} \right] \\
 &\approx \frac{1}{n} \sum_{i=1}^n R(\tau^i) \frac{\nabla p_{\theta}(\tau^i)}{p_{\theta}(\tau^i)}.
 \end{aligned}$$

17 If we sample trajectories  $\tau^i$  that are not very likely under the distribution  
 18  $p_{\theta}(\tau)$ , the denominator in some of the summands can be very small. For  
 19 trajectories that are likely, the denominator is large. The empirical estimate



1 of the expectation using  $n$  trajectories where some terms are very small and  
 2 some others very large, therefore has a large variance. So one does need lots of  
 3 trajectories from the system/simulator to compute a reasonable approximation  
 4 of the policy gradient.

### 5 8.3.1 Reducing the variance of the policy gradient

6 **Control variates** You will perhaps appreciate that computing the accurate  
 7 policy gradient is very hard in practice. Control variates is a general concept  
 8 from the literature on Monte Carlo integration and is typically introduced as  
 9 follows. Say we have a random variable  $X$  and we would like to guess its  
 10 expected value  $\mu = \mathbb{E}[X]$ . Note that  $X$  is an unbiased estimator of  $\mu$  but it  
 11 may have a large variance. If we have another random variable  $Y$  with known  
 12 expected value  $\mathbb{E}[Y]$ , then

$$\hat{X} = X + c(Y - \mathbb{E}[Y]) \quad (8.10)$$

13 is also an unbiased estimator for  $\mu$  for any value of  $c$ . The variance of  $\hat{X}$  is

$$\text{var}(\hat{X}) = \text{var}(X) + c^2 \text{var}(Y) + 2c \text{cov}(X, Y).$$

14 which is minimized for

$$c^* = -\frac{\text{cov}(X, Y)}{\text{var}(Y)}$$

15 for which we have

$$\text{var}(\hat{X}) = \text{var}(X) - c^{*2} \text{var}(Y) = \left(1 - \left(\frac{\text{cov}(X, Y)}{\text{var}(Y)}\right)^2\right) \text{var}(X).$$

16 By subtracting  $Y - \mathbb{E}[Y]$  from our observed random variable  $X$ , we have re-  
 17 duced the variance of  $X$  if the correlation between  $X$  and  $Y$  is non-zero. Most  
 18 importantly, note that no matter what  $Y$  we plug into the above expression, we  
 19 can never increase the variance of  $X$ ; the worst that can happen is that we pick  
 20 a  $Y$  that is completely uncorrelated with  $X$  and end up achieving nothing.

21 **Baseline** We will now use the concept of a control variate to reduce the  
 22 variance of the policy gradient. This is known as “building a baseline”. The  
 23 simplest baseline one can build is to subtract a constant value from the return.  
 24 Consider the PG given by

$$\begin{aligned} \nabla J(\theta) &= \mathbb{E}_{\tau \sim p_\theta} [R(\tau) \nabla \log p_\theta(\tau)] \\ &= \mathbb{E}_{\tau \sim p_\theta(\tau)} [(R(\tau) - b) \nabla \log p_\theta(\tau)]. \end{aligned}$$

25 Observe that

$$\begin{aligned} \mathbb{E}_{\tau \sim p_\theta(\tau)} [b \nabla \log p_\theta(\tau)] &= \int d\tau b p_\theta(\tau) \nabla \log p_\theta(\tau) \\ &= \int d\tau b \nabla p_\theta(\tau) = b \nabla \int d\tau p_\theta(\tau) = b \nabla 1 = 0. \end{aligned}$$

1 **Example 1: Using the average returns of a mini-batch as the baseline**  
 2 What is the simplest baseline  $b$  we can cook up? Let us write the mini-batch  
 3 version of the policy gradient

$$\widehat{\nabla} J(\theta) := \frac{1}{\ell} \sum_{i=1}^{\ell} [R(\tau^i) \nabla \log p_{\theta}(\tau^i)].$$

4 where  $\tau^1, \dots, \tau^{\ell}$  are trajectories that are a part of our mini-batch. We can set

$$b = \frac{1}{\ell} \sum_{i=1}^{\ell} R(\tau^i)$$

5 can use the variance-reduced gradient

$$\widehat{\nabla} J(\theta) = \frac{1}{\ell} \sum_{i=1}^{\ell} [(R(\tau^i) - b) \nabla \log p_{\theta}(\tau^i)].$$

6 This is a one-line change in your code for policy gradient so there is no reason  
 7 not to do it.

8 **Example 2: A weighted averaged of the returns using the log-likelihood**  
 9 **of the trajectory** The previous example showed how we can use one constant  
 10 baseline, namely the average of the discounted returns of all trajectories  
 11 in a mini-batch. What is the **best** constant  $b$  we can use? We can perform a  
 12 similar computation as done in the control variate introduction to minimize  
 13 the variance of the policy gradient to get the following.

$$\begin{aligned} \text{var} \left( \widehat{\nabla}_{\theta_i} J(\theta) \right) &= \mathbb{E}_{\tau} \left[ ((R(\tau) - b_i) \nabla_{\theta_i} \log p_{\theta}(\tau))^2 \right] - \left( \mathbb{E}_{\tau} [((R(\tau) - b_i) \nabla_{\theta_i} \log p_{\theta}(\tau))] \right)^2 \\ &= \mathbb{E}_{\tau} \left[ ((R(\tau) - b_i) \nabla_{\theta_i} \log p_{\theta}(\tau))^2 \right] - \left( \widehat{\nabla}_{\theta_i} J(\theta) \right)^2. \end{aligned}$$

14 Set

$$\frac{\text{var} \left( \widehat{\nabla}_{\theta_i} J(\theta) \right)}{db_i} = 0$$

15 in the above expression to get

$$b_i = \frac{\mathbb{E}_{\tau} \left[ (\nabla_{\theta_i} \log p_{\theta}(\tau))^2 R(\tau) \right]}{\mathbb{E}_{\tau} \left[ (\nabla_{\theta_i} \log p_{\theta}(\tau))^2 \right]}$$

16 which is the baseline you should subtract from the gradient of the  $i^{\text{th}}$  parameter  
 17  $\theta_i$  to result in the largest variance reduction. This expression is just the  
 18 expected return but it is weighted by the magnitude of the gradient, this again  
 19 is 1–2 lines of code.

🔗 Show that any function that only depends on the state  $x$  can be used as a baseline in the policy gradient. This technique is known as reward shaping.

## 8.4 An alternative expression for the policy gradient

We will first define an important quantity that helps us think of RL algorithms.

**Definition 8.1 (Discounted state visitation frequency).** Given a stochastic controller  $u_\theta(\cdot | x)$  the discounted state visitation frequency for a discrete-time dynamical system is given by

$$d^\theta(x) = \sum_{k=0}^{\infty} \gamma^k \mathbf{P}(x_k = x | x_0, u_k \sim u_\theta(\cdot | x_k)).$$

The distribution  $d^\theta(x)$  is the probability of visiting a state  $x$  computed over all trajectories of the system that start at the initial state  $x_0$ . If  $\gamma = 1$ , this is the steady-state distribution of the Markov chain underlying the Markov Decision Process where at each step the MDP chooses the control  $u_k \sim u_\theta(\cdot | x_k)$ . The fact that we have defined the discounted distribution is a technicality; this version is seen in the policy gradient expression. You will also notice that  $d^\theta(x)$  is not a normalized distribution. The normalization constant is difficult to characterize both theoretically and empirically and we will not worry about it here; RL algorithms do not require it.

**Q-function** Using the discounted state visitation frequency, the policy gradient that we saw before can be written in terms of the value function as follows.

$$\begin{aligned} \nabla J(\theta) &= \mathbf{E}_{\tau \sim p_\theta} [R(\tau) \nabla \log p_\theta(\tau)] \\ &= \mathbf{E}_{x \sim d^\theta} \mathbf{E}_{u \sim u_\theta(\cdot | x)} [q^\theta(x, u) \nabla \log u_\theta(u | x)]. \end{aligned} \quad (8.11)$$

The function  $q^\theta(x, u)$  is similar to the cost-to-go that we have studied in dynamic programming and is called the Q-function

$$q^\theta(x, u) = \mathbf{E}_{\tau \sim p_\theta(\tau)} [R(\tau) | x_0 = x, u_0 = u]. \quad (8.12)$$

It is the infinite-horizon discounted cumulative reward (i.e., the return) if the system starts at state  $x$  and takes the control  $u$  in the first step and runs the controller  $u_\theta(\cdot | x)$  for all steps thereafter. We make the dependence of  $q^\theta$  on the parameters  $\theta$  of the controller explicit.

Compare the above formula for the policy gradient with the one we had before in (8.8)

$$\begin{aligned} \widehat{\nabla} J(\theta) &= \mathbf{E}_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)] \\ &= \mathbf{E}_{\tau \sim p_\theta(\tau)} \left[ \left( \sum_{k=0}^T \gamma^k r(x_k, u_k) \right) \left( \sum_{k=0}^T \nabla \log u_\theta(u_k | x_k) \right) \right]. \end{aligned}$$

It is important to notice that this is an expectation over trajectories; whereas

**i** The derivation of this expression is easy although tedious, you can find it in the Appendix of the paper “Policy gradient methods for reinforcement learning with function approximation” at <https://papers.nips.cc/paper/1713-policy-gradient-methods-for-reinforcement-learning-with-function-approximation.pdf>.

(8.11) is an expectation over states  $x$  sampled from the discounted state visitation frequency. The control  $u_k$  for both is sampled from the stochastic controller at each time-step  $k$ . The most important distinction is that (8.11) involves the expectation of the Q-function  $q^\theta$  weighted by the gradient of the log-likelihood of picking each control action. There are numerous hacky ways of deriving (8.11) from (8.8) but remember that they are fundamentally different expressions of the *same policy gradient*.

This expression allows understanding of a number of properties of reinforcement learning.

1. While the algorithm collects the data, states that are unlikely under the distribution  $d^\theta$  contribute little to (8.11). In other words, the policy gradient is insensitive to such states. The policy update will not consider these unlikely states that the system is prone to visit infrequently using the controller  $u_\theta$ .
2. The opposite happens for states which are very likely. For two controls  $u_1, u_2$  at the same state  $x$ , the policy increases the log-likelihood of taking the controls weighted by their values  $q^\theta(x, u_1)$  and  $q^\theta(x, u_2)$ . This is sort of the “definition” of reinforcement learning. In the expression (8.8) the gradient was increasing the likelihood of trajectories with high returns, here it deals with states and controls individually.

### 8.4.1 Implementing the new expression

Suppose we have a stochastic control that is a Gaussian

$$u_\theta(u | x) = \frac{1}{(2\pi\sigma^2)^{p/2}} e^{-\frac{\|u - \theta^\top x\|^2}{2\sigma^2}}$$

where  $\theta \in \mathbb{R}^{d \times p}$  and  $u \in \mathbb{R}^p$ ; the variance  $\sigma$  can be chosen by the user. We can easily compute  $\log u_\theta(u | x)$  in (8.11). How should one compute  $q^\theta(x, u)$  in (8.12)? We can again estimate it using sample trajectories from the system; each of these trajectories would have to start from a state  $x$  and the control at the first step would be  $u$ , with the controller  $u_\theta$  being used thereafter. Note that we have one such trajectory, namely the remainder of the trajectory where we encountered  $(x, u)$  while sampling trajectories for the policy gradient in (8.11). In practice, we do not sample trajectories a second time, we simply take this trajectory, let us call it  $\tau_{x,u}$  and set

$$q^\theta(x, u) = \sum_{k=0}^T \gamma^k r(x_k, u_k)$$

for some large time-horizon  $T$  where  $(x_0, u_0) = (x, u)$  and the summation is evaluated for  $(x_k, u_k)$  that lie on the trajectory  $\tau_{x,u}$ . Effectively, we are evaluating (8.12) using one sample trajectory, a highly erroneous estimate of  $q^\theta$ .

## 8.5 Actor-Critic methods

We can of course do more sophisticated things to evaluate the Q-function  $q^\theta$  in our new expression of the policy gradient.

Actor-Critic methods fit a Q-function to the data collected from the system using the current controller (policy evaluation step) and then use this fitted Q-function in the expression of the policy gradient (8.11) to update the controller. In this sense, Actor-Critic methods are conceptually similar to policy iteration.

In order to understand how to fit the Q-function, first recall that it should satisfy the Bellman equation. This means

$$q^\theta(x, u) = r(x, u) + \gamma \mathbb{E}_{u \sim u_\theta(\cdot | x'), x' \sim P(\cdot | x, u)} [q^\theta(x', u')]. \quad (8.13)$$

We do not know a model for the system so we cannot evaluate the expectation over  $x' \sim P(\cdot | x, u)$  like we used to in dynamic programming. But we do have the ability to get trajectories  $\tau^i$  from the system.

Let's say  $(x_k^i, u_k^i)$  lie on  $\tau^i$  at time-step  $k$ . We can then estimate the expectation over  $P(\cdot | x_k^i, u_k^i)$  using simply  $x_{k+1}^i$  and the expectation over the controls using simply  $u_{k+1}^i$  to write

$$q^\theta(x_k^i, u_k^i) \approx r(x_k^i, u_k^i) + \gamma q^\theta(x_{k+1}^i, u_{k+1}^i) \quad \text{for all } i \leq n, k \leq T.$$

This is a nice constraint on the Q-function. If this were a discrete-state, discrete-control MDP, it is a set of linear equations for the q-values. These constraints would be akin to our linear equations for evaluating a policy in dynamic programming except that instead of using the dynamics model (the transition matrix), we are using trajectories sampled from the system.

**Parameterizing the Q-function using a neural network** If we are dealing with a continuous state/control-space, we can think of parameterizing the q-function using parameters  $\varphi$

$$q_\varphi^\theta(x, u) : X \times U \rightarrow \mathbb{R}.$$

The parameterization is similar to the parameterization of the controller, e.g., just like we would write a deterministic controller as

$$u_\theta(x) = \theta^\top x$$

we can think of a linear Q-function of the form

$$q_\varphi^\theta(x, u) = \varphi^\top \begin{bmatrix} x \\ u \end{bmatrix}, \quad \varphi \in \mathbb{R}^{m+d}$$

1 which is a linear function in the states and controls. You can also think of  
2 using something like

$$q_{\varphi}^{\theta}(x, u) = [1 \quad x \quad u] \varphi \begin{bmatrix} 1 \\ x \\ u \end{bmatrix} \quad \varphi \in \mathbb{R}^{(m+d+1) \times (m+d+1)}.$$

3 which is quadratic in the states and controls, or in general a deep network with  
4 weights  $\varphi$  as the Q-function.

5 **Fitting the Q-function** We can now “fit” the parameters of the Q-function  
6 by solving the problem

$$\hat{\varphi} = \underset{\varphi}{\operatorname{argmin}} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{k=0}^T \|q_{\varphi}^{\theta}(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma q_{\varphi}^{\theta}(x_{k+1}^i, u_{k+1}^i)\|^2. \quad (8.14)$$

7 which is nothing other than enforcing the Bellman equation in (8.13). If the  
8 Q-function is linear in  $[x, u]$  this is a least squares problem, if it is quadratic  
9 the problem is a quadratic optimization problem which can also be solved  
10 efficiently, in general we would solve this problem using stochastic gradient  
11 descent if we are parameterizing the Q-function using a deep network. Such a  
12 Q-function is called the “critic” because it *evaluates* the controller  $u_{\theta}$ , or the  
13 “actor”. This version of the policy gradient where one fits the parameters of  
14 both the controller and the Q-function are called Actor-Critic methods.

❗ We will be pedantic and always write the  $q$ -function as  $q_{\varphi}^{\theta}$ . The superscript  $\theta$  denotes that this is the  $q$ -function corresponding to the  $\theta$ -parameterized controller  $u_{\theta}$ . The subscript denotes that the  $q$ -function is parameterized by parameters  $\varphi$ .

**Actor-Critic Methods** We fit a deep network with weights  $\theta$  to parameterize a stochastic controller  $u_{\theta}(\cdot | x)$  and another deep network with weights  $\varphi$  to parameterize the Q-function of this controller,  $q_{\varphi}^{\theta}(x, u)$ . Let the controller weights at the  $k^{\text{th}}$  iteration be  $\theta^k$  and the Q-function weights be  $\varphi^k$ .

1. Sample  $n$  trajectories, each of  $T$  timesteps, using the current controller  $u_{\theta^k}(\cdot | x)$ .
2. Fit a Q-function  $q_{\varphi^{k+1}}^{\theta^k}$  to this data using (8.14). using stochastic gradient descent. While performing this fitting (although it is not mathematically sound), it is common to use initialize the Q-function weights to their values from the previous iteration  $\varphi^k$ .
3. Compute the policy gradient using the alternative expression in (8.11) and update parameters of the policy to  $\theta^{k+1}$ .

### 15 8.5.1 Advantage function

16 The new expression for the policy gradient in (8.11) also has a large variance;  
17 this should be no surprise, it is after all equal to the old expression. We can  
18 however perform variance reduction on this using the value function.

19 Our goal as before would be construct a relevant baseline to subtract from  
20 the Q-function. It turns out that any function that depends only upon the state

1  $x$  is a valid baseline. This gives a powerful baseline for us to use in policy  
 2 gradients. We can use the value function as the baseline. The value function  
 3 for controls taken by the controller  $u_\theta(\cdot | x)$  (notice that this is not the optimal  
 4 value function, it is simply the policy evaluation) is given by

$$v^\theta(x) = \mathbf{E}_{\tau \sim p_\theta(\tau)} [R(\tau) | x_0 = x]$$

5 where  $u_k \sim u_\theta(\cdot | x_k)$  at each timestep. We also know that the value function  
 6 is the expected value of the Q-function across different controls sampled by  
 7 the controller

$$v^\theta(x) = \mathbf{E}_{u \sim u_\theta(\cdot | x)} [q^\theta(x, u)]. \quad (8.15)$$

8 The value function again satisfies the dynamic programming principle/Bellman  
 9 equation

$$v^\theta(x) = \mathbf{E}_{u \sim u_\theta(\cdot | x)} \left[ r(x, u) + \gamma \mathbf{E}_{x' \sim P(\cdot | x, u)} [v^\theta(x')] \right].$$

10 We again parameterize the value function

$$v_\psi^\theta(x) : X \rightarrow \mathbb{R}$$

11 using parameters  $\psi$  and fit it to the data in the same way as (8.14) to get

$$\hat{\psi} = \operatorname{argmin}_\psi \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{k=0}^T \|v_\psi^\theta(x_k^i) - r(x_k^i, u_k^i) - \gamma v_\psi^\theta(x_{k+1}^i)\|^2. \quad (8.16)$$

Using this baseline can modify the policy gradient to be

$$\nabla J(\theta) = \mathbf{E}_{x \sim d^\theta} \mathbf{E}_{u \sim u_\theta(\cdot | x)} \left[ \underbrace{(q_\varphi^\theta(x, u) - v_\psi^\theta(x))}_{a_{\varphi, \psi}^\theta(x, u)} \nabla_\theta \log u_\theta(u | x) \right]. \quad (8.17)$$

where each of the functions  $q_\varphi^\theta$  and  $v_\psi^\theta$  are themselves fitted using (8.14) and (8.16) respectively. The difference

$$a_{\varphi, \psi}^\theta(x, u) = q_\varphi^\theta(x, u) - v_\psi^\theta(x) \approx q_\varphi^\theta(x, u) - \mathbf{E}_{u \sim u_\theta(\cdot | x)} [q_\varphi^\theta(x, u)] \quad (8.18)$$

is called the advantage function. It measures how much better the particular control  $u$  is for a state  $x$  as compared to the average return of controls sampled from the controller at that state. The form (8.17) is the most commonly implemented form in research papers whenever they say “we use the policy gradient”.

🔍 The advantage function is very useful while doing theoretical work on RL algorithms. But it is also extremely useful in practice. It imposes a constraint upon our estimate  $q_\varphi^\theta$  and the estimate  $v_\psi^\theta$ . If we are not solving (8.14) and (8.16) to completion, we may benefit by imposing this constraint on the advantage function. Can you think of a way?

## 8.6 Discussion

This brings to an end the discussion of policy gradients. They are, in general, a complicated suite of algorithms to implement. You will see some of this complexity when you implement the controller for something as simple as the simple pendulum. The key challenges with implementing policy gradients come from the following.

1. Need lots of data, each parameter update requires fresh data from the systems. Typical problems may need a million trajectories, most robots would break before one gets this much data from them if one implements these algorithms naively.
2. The log-likelihood ratio trick has a high variance due to  $u_{\theta}(\cdot | x)$  being in the denominator of the expression, so we need to implement complex variance reduction techniques such as actor-critic methods.
3. Fitting the Q-function and the value function is not easy, each parameter update of the policy ideally requires you to solve the entire problems (8.14) and (8.16). In practice, we only perform a few steps of SGD to solve the two problems and reuse the solution of  $k^{\text{th}}$  iteration update as an initialization of the  $k + 1^{\text{th}}$  update. This is known as “warm start” in the optimization literature and reduces the computational cost of fitting the Q/value-functions from scratch each time.
4. The Q/value-function fitted in iteration  $k$  may be poor estimates of the Q/value at iteration  $k + 1$  for the new controller  $u_{\theta^{k+1}}(\cdot | x)$ . If the controller parameters change quickly,  $\theta^{k+1}$  is very different from  $\theta^k$ , then so are  $q^{\theta^{k+1}}$  and  $v^{\theta^{k+1}}$ . There is a very fine balance between training quickly and retaining the efficiency of warm start; and tuning this in practice is quite difficult. A large number of policy gradient algorithms like TRPO (<https://arxiv.org/abs/1502.05477>) and PPO (<https://arxiv.org/abs/1707.06347>) try to alleviate this with varying degrees of success.
5. The latter, PPO, is a good policy-gradient-based algorithm to try on a new problem. For instance, in a very impressive demonstration, it was used to build an RL agent to play StarCraft (<https://openai.com/blog/openai-five>). We will see better RL methods in the next chapter.



# Chapter 9

## Q-Learning

### Reading

1. Sutton & Barto, Chapter 6, 11
2. Human-level control through deep reinforcement learning  
<https://www.nature.com/articles/nature14236>
3. Deterministic Policy Gradient Algorithms,  
<http://proceedings.mlr.press/v32/silver14.html>
4. Addressing Function Approximation Error in Actor-Critic Methods  
<https://arxiv.org/abs/1802.09477>
5. An Application of Reinforcement Learning to Aerobatic Helicopter Flight, <https://papers.nips.cc/paper/3151-an-application-of-reinforcement-learning-to-aerobatic-helicopter-flight>

In the previous chapter, we looked at what are called “on-policy” methods, these are methods where the current controller  $u_{\theta^k}$  is used to draw fresh data from the dynamical system and used to update to parameters  $\theta^k$ . *The key inefficiency in on-policy methods is that this data is thrown away in the next iteration.* We need to draw a fresh set of trajectories from the system for  $u_{\theta^{k+1}}$ . This lecture will discuss off-policy methods which are a way to reuse past data. These methods require much fewer data than on-policy methods (in practice, about 10–100× less).

### 9.1 Tabular Q-Learning

Recall the value iteration algorithm for discrete (and finite) state and control spaces; this is also called “tabular” Q-Learning in the RL literature because we can store the Q-function  $q(x, u)$  as a large table with number of rows being the number of states and number of columns being the number of controls, with

1 each entry in this table being the value  $q(x, u)$ . Value iteration when written  
 2 using the Q-function at the  $k^{\text{th}}$  iteration for the tabular setting looks like

$$\begin{aligned} q^{(k+1)}(x, u) &= \sum_{x' \in X} \mathbf{P}(x' | x, u) \left( r(x, u) + \gamma \max_{u'} q^{(k)}(x', u') \right) \\ &= \mathbf{E}_{x' \sim \mathbf{P}(\cdot | x, u)} \left[ r(x, u) + \gamma \max_{u'} q^{(k)}(x', u') \right]. \end{aligned}$$

In the simplest possible instantiation of Q-learning, the expectation in the value iteration above (which we can only compute if we know a model of the dynamics) is replaced by samples drawn from the environment.

3 We will imagine the robot as using an *arbitrary* controller

$$u_e(\cdot | x)$$

4 that has a fairly large degree of randomness in how it picks actions. We call  
 5 such a controller an “exploratory controller”. Conceptually, its goal is to lead  
 6 the robot to diverse states in the state-space so that we get a faithful estimate  
 7 of the expectation in value iteration. We maintain the value  $q^{(k)}(x, u)$  for all  
 8 states  $x \in X$  and controls  $u \in U$  and update these values to get  $q^{(k+1)}$  after  
 9 *each step* of the robot.

10 From the results on Bellman iteration, we know that any Q-function that  
 11 satisfies the above equation is the optimal Q-function; we would therefore like  
 12 our Q-function to satisfy

$$q^*(x_k, u_k) \approx r(x_k, u_k) + \gamma \max_{u'} q^*(x_{k+1}, u').$$

13 over samples  $(x_k, u_k, x_{k+1})$  collected as the robot explores the environment.

**Tabular Q-Learning** Let us imagine the robot travels for  $n$  trajectories each of  $T$  time-steps each. We can now solve for  $q^*$  by minimizing the objective

$$\min_q \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{k=0}^T \|q(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma \max_{u' \in U} q(x_{k+1}^i, u')\|_2^2. \quad (9.1)$$

on the data collected by the robot. The variable of optimization here are all values  $q^*(x, u)$  for  $x \in X$  and  $u \in U$ .

14 Notice a few important things about the above optimization problem. First,  
 15 the last term is a maximization over  $u' \in U$ , it is  $\max_{u' \in U} q(x_{k+1}^i, u')$  and  
 16 not  $q(x_{k+1}^i, u_{k+1}^i)$ . In practice, you should imagine a robot performing Q-  
 17 Learning in a grid-world setting where it seeks to find the optimal trajectory  
 18 to go from a source location to a target location. If at each step, the robot  
 19 has 4 controls to choose from, computing this last term involves taking the  
 20 maximum of 4 different values (4 columns in the tabular Q-function).

1 Notice that for finite-horizon dynamic programming we initialized the  
 2 Q-function at the terminal time to a known value (the terminal cost). Similarly,  
 3 for infinite-horizon value iteration, we discussed how we can converge to the  
 4 optimal Q-function with any initialization. In the above case, we do not impose  
 5 any such constraint upon the Q-function, but there is an implicit constraint. All  
 6 values  $q(x, u)$  have to be consistent with each other and ideally, the residual

$$\|q(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma \max_{u' \in U} q(x_{k+1}^i, u')\|_2^2 = 0$$

7 for all trajectories  $i$  and all timesteps  $T$ .

8 **Solving tabular Q-Learning** How should we solve the optimization prob-  
 9 lem in (9.1)? This is easy, every entry  $q(x, u)$  for  $x \in U$  and  $u \in U$  is a  
 10 variable of this objective and each  $\|\cdot\|_2^2$  term in the objective simply represents  
 11 a constraint that ties these different values of the Q-function together. We can  
 12 solve for all  $q(x, u)$  iteratively as

$$\begin{aligned} q(x, u) &\leftarrow q(x, u) - \eta \nabla_{q(x, u)} \ell(q) \\ &= (1 - \eta) q(x, u) - \eta \left( r(x, u) + \gamma \max_{u'} q(x', u') \right) \end{aligned} \quad (9.2)$$

13 where  $\ell(q)$  is the entire objective  $\frac{1}{n(T+1)} \sum_i \sum_k \dots$  above and  $(x, u, x') \equiv$   
 14  $(x_k^i, u_k^i, x_{k+1}^i)$  in the second equation. An important point to note here is that  
 15 although the robot collects a finite number of data

$$D = \{(x_k^i, u_k^i)_{k=0,1,\dots,T}\}_{i=1}^n$$

16 we have an estimate for the value  $q(x, u)$  at all states  $x \in X$ . Intuitively,  
 17 tabular Q-learning looks at the returns obtained by the robot after starting from  
 18 a state  $x$  (the reward-to-come  $J(x)$ ) and patches the returns from nearby states  
 19  $x, x'$  using the constraints in the objective (9.1).

20 **Terminal state** One must be very careful about the terminal state in such  
 21 implementations of Q-learning. Typically, most research papers imagine that  
 22 they are solving an infinite horizon problem but use simulators that have an  
 23 explicit terminal state, i.e., the simulator does not proceed to the next timestep  
 24 after the robot reaches the goal. A workaround for using such simulators (this  
 25 applies for essentially all simulators) is to modify (9.2) as

$$q(x, u) = (1 - \eta) q(x, u) - \eta \left( r(x, u) + \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \max_{u'} q(x', u') \right).$$

26 Effectively, we are setting  $q(x', u) = 0$  for all  $u \in U$  if  $x'$  is a terminal state  
 27 of problem. This is a very important point to remember and Q-Learning will  
 28 never work if you forget to include the term  $\mathbf{1}_{\{x' \text{ is terminal}\}}$  in your expression.

29 **What is the controller in tabular Q-Learning?** The controller in tabular Q-  
 30 Learning is easy to get after we solve (9.1). At test time, we use a deterministic

1 controller given by

$$u^*(x) = \underset{u'}{\operatorname{argmax}} q^*(x, u').$$

## 2 9.1.1 How to perform exploration in Q-Learning

3 The exploratory controller used by the robot  $u_e(\cdot | x)$  is critical to perform  
 4 Q-Learning well. If the exploratory controller does not explore much, we do  
 5 not get states from all parts of the state-space. This is quite bad, because in this  
 6 case the estimates of Q-function at *all states* will be bad, not just at the states  
 7 that the robot did not visit. To make this intuitive, imagine if we cordoned off  
 8 some nodes in the graph for the backward version of Dijkstra’s algorithm and  
 9 never used them to update the dist variable. We would never get to the optimal  
 10 cost-to-go for *all* states in this case because there could be trajectories that  
 11 go through these cordoned off states that lead to a smaller cost-to-go. So it is  
 12 quite important to pick the right exploratory controller.

13 It turns out that a random exploratory controller, e.g., a controller  $u_e(\cdot | x)$   
 14 that picks controls uniformly randomly is pretty good. We can show that our  
 15 tabular Q-Learning will converge to the optimal Q-function  $q^*(x, u)$  as the  
 16 amount of data drawn from the random controller goes to infinity, even if we  
 17 initialize the table to arbitrary values. In other words, if we are guaranteed that  
 18 the robot visits each state in the finite MDP infinitely often, it is a classical  
 19 result that updates of the form (9.2) for minimizing the objective in (9.1)  
 20 converge to the optimal Q-function.

21 **Epsilon-greedy exploration** Instead of the robot using a arbitrary controller  
 22  $u_e(\cdot | x)$  to gather data, we can use the current estimate of the Q-function  
 23 with some added randomness to ensure that the robot visits all states in the  
 24 state-space. This is a key idea in Q-Learning and is known as “epsilon-greedy”  
 25 exploration. We set

$$u_e(u | x) = \begin{cases} \operatorname{argmax}_u q(x, u) & \text{with probability } 1 - \epsilon \\ \operatorname{uniform}(U) & \text{with probability } \epsilon. \end{cases} \quad (9.3)$$

26 for some user-chosen value of  $\epsilon$ . Effectively, the robot repeats the controls it  
 27 took in the past with probability  $1 - \epsilon$  and uniformly samples from the entire  
 28 control space with probability  $\epsilon$ . The former ensures that the robot moves  
 29 towards the parts of the state-space where states have a high return-to-come  
 30 (after all, that is the what the Q-function  $q(x, u)$  indicates). The latter ensures  
 31 that even if the robot’s estimate of the Q-function is bad, it is still visiting  
 32 every state in the state-space infinitely often.

33 **A different perspective on Q-Learning** Conceptually, we can think of  
 34 tabular Q-learning as happening in two stages. In the first stage, the robot  
 35 gathers a large amount of data

$$D = \{(x_k^i, u_k^i)_{k=0,1,\dots,T}\}_{i=1}^n$$

36 using the exploratory controller  $u_e(\cdot | x)$ ; let us consider the case when we are  
 37 using an arbitrary exploratory controller, not epsilon-greedy exploration. Using

❶ This is again the power of dynamic programming at work. The Bellman equation guarantees the convergence of value iteration provided we compute the expectation exactly. But if the robot does give us lots of data from the environment, then Q-Learning also inherits this property of convergence to the optimal Q-function from any initialization.

1 this data, the robot fits a model for the system, i.e., it learns the underlying  
2 MDP

$$P(x' | x, u);$$

3 this is very similar to the step in the Baum-Welch algorithm that we saw for  
4 learning the Markov state transition matrix of the HMM in Chapter 2. We  
5 simply take frequency counts to estimate this probability

$$P(x' | x, u) \approx \frac{1}{N} \sum_i \mathbf{1}_{\{x' \text{ was reached from } x \text{ using control } u\}}$$

6 where  $N$  is the number of the times the robot took control  $u$  at state  $x$ . Given  
7 this transition matrix, we can now perform value iteration on the MDP to learn  
8 the Q-function

$$q^{(k+1)}(x, u) = \mathbb{E}_{x' \sim P(\cdot | x, u)} \left[ r(x, u) + \gamma \max_{u'} q^{(k)}(x', u) \right].$$

9 The success of this two-stage approach depends upon how accurate our esti-  
10 mate of  $P(x' | x, u)$  is. This in turn depends on how much the robot explored  
11 the domain and the size of the dataset it collected, both of these need to be  
12 large. We can therefore think of Q-learning as interleaving these two stages  
13 in a single algorithm, it learns the dynamics of the system and the Q-function  
14 for that dynamics simultaneously. But the Q-Learning algorithm does not  
15 really maintain a representation of the dynamics, i.e., at the end of running  
16 Q-Learning, we do not know what  $P(x' | x, u)$  is.

## 17 9.2 Function approximation (Deep Q Networks)

18 Tabular methods are really nice but they do not scale to large problems. The  
19 grid-world in the homework problem on policy iteration had 100 states, a  
20 typical game of Tetris has about  $10^{60}$  states. For comparison, the number of  
21 atoms in the known universe is about  $10^{80}$ . The number of different states  
22 in a typical Atari game is more than  $10^{300}$ . These are all problems with a  
23 discrete number of states and controls, for continuous state/control-space, the  
24 number of distinct states/controls is infinite. So it is essentially impossible to  
25 run the tabular Q-Learning method from the previous section for most real-  
26 world problems. In this section, we will look at a powerful set of algorithms  
27 that parameterize the Q-function using a neural network to work around this  
28 problem.

29 We use the same idea from the previous chapter, that of parameterizing the  
30 Q-function using a deep network. We will denote

$$q_\varphi(x, u) : X \times U \mapsto \mathbb{R}$$

31 as the Q-function and our goal is to fit the deep network to obtain the weights  $\hat{\varphi}$ ,  
32 instead of maintaining a very large table of size  $|X| \times |U|$  for the Q-function.  
33 Fitting the Q-function is quite similar to the tabular case: given a dataset

1  $D = \{(x_t^i, u_t^i)_{t=0,1,\dots,T}\}_{i=1}^n$  from the system, we want to enforce

$$q_\varphi(x_t^i, u_t^i) = r(x_t^i, u_t^i) + \gamma \max_{u'} q_\varphi(x_{t+1}^i, u')$$

2 for all tuples  $(x_t^i, u_t^i, x_{t+1}^i)$  in the dataset. Just like the previous section, we  
3 will solve

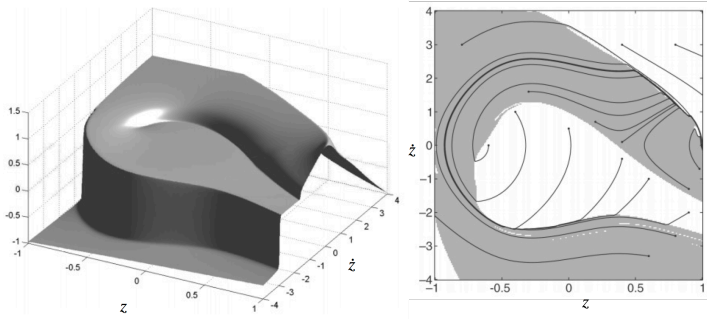
$$\begin{aligned} \hat{\varphi} &= \operatorname{argmin}_{\varphi} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=1}^T \left( q_\varphi(x_t^i, u_t^i) - r(x_t^i, u_t^i) - \underbrace{\gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right)}_{\text{target}(x'; \varphi)} \max_{u'} q_\varphi(x_{t+1}^i, u') \right)^2 \\ &\equiv \operatorname{argmin}_{\varphi} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=1}^T (q_\varphi(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi))^2 \end{aligned} \quad (9.4)$$

4 The last two terms in this expression above are together called the “target”  
5 because the problem is very similar to least squares regression, except that the  
6 targets also depend on the weights  $\varphi$ . This is what makes it challenging to  
7 solve.

8 As discussed above, Q-Learning with function approximation is known  
9 as “Fitted Q Iteration”. Remember that very important point that the robot  
10 collects data using the exploratory controller  $u_e(\cdot | x)$  but the Q-function that  
11 we fit is the *optimal* Q-function.

12 **Fitted Q-Iteration with function approximation may not converge to the**  
13 **optimal Q-function** It turns out that (9.4) has certain mathematical intricacies  
14 that prevent it from converging to the optimal Q-function. We will  
15 first give the intuitive reason. In the tabular Q-Learning setting, if we mod-  
16 ify some entry  $q(x, u)$  for an  $x \in X$  and  $u \in U$ , the other entries (which  
17 are tied together using the Bellman equation) are all modified. This is akin  
18 to you changing the dist value of one node in Dijkstra’s algorithm; the dist  
19 values of *all* other nodes will have to change to satisfy the Bellman equation.  
20 This is what (9.2) achieves if implemented with a decaying step-size  $\eta$ ; see  
21 <http://users.isr.ist.utl.pt/~mtjspaam/readingGroup/ProofQlearning.pdf> for the  
22 proof. This does not hold for (9.4). Even if the objective in (9.4) is zero on our  
23 collected dataset, i.e., the Q-function fits data collected by the robot perfectly,  
24 the Q-function may not be the *optimal* Q-function. An intuitive way of under-  
25 standing this problem is that even if the Bellman error is zero on samples in  
26 the dataset, the optimization objective says nothing about states that are not  
27 present in the dataset; the Bellman error on them is completely dependent upon  
28 the smoothness properties of the function expressed by the neural architecture.  
29 Contrast this comment with the solution of the HJB equation in Chapter 6  
30 where the value function was quite non-smooth at some places. If our sampled  
31 dataset does not contain those places, there is no way the neural network can  
32 know the optimal form of the value function.

**i** The mathematical reason behind this is that the Bellman operator, i.e., the update to the Q/value-function is a contraction for the tabular setting, this is not the case for Fitted Q-Iteration unless the function approximation has some technical conditions imposed upon it.



## 9.2.1 Embellishments to Q-Learning

We next discuss a few practical aspects of implementing Q-Learning. Each of the following points is extremely important to understand how to get RL to work on real-world problems, so you should internalize these.

**Pick mini-batches from different trajectories in SGD** . In practice, we fit the Q-function using stochastic gradient descent. At each iteration we sample a mini-batch of inputs  $(x_t^i, u_t^i, x_{t+1}^i)$  from different trajectories  $i \in \{1, \dots, n\}$  and update the weights  $\varphi$  in the direction of the negative gradient.

$$\varphi^{k+1} = \varphi^k - \eta \nabla_{\varphi} (q_p^k(x, u) - \text{target}(x'; \varphi^k))^2.$$

The mini-batch is picked to have samples from different trajectories because samples from the same trajectory are correlated to each other (after all, the robot obtains the next tuple  $(x', u', x'')$  from the previous tuple  $(x, u, x')$ ).

**Replay buffer** The dataset  $D$  is known as the replay buffer.

**Off-policy learning** The replay buffer is typically not fixed during training. Instead of drawing data from the exploratory controller  $u_e$ , we can think of the following algorithm. Initialize the Q-function weights to  $\varphi^0$  and the dataset to  $D = \emptyset$ . At the  $k^{\text{th}}$  iteration,

- Draw a dataset  $D^k$  of  $n$  trajectories from the  $\epsilon$ -greedy policy

$$u_e(u | x) = \begin{cases} \operatorname{argmax}_u q^k(x, u) & \text{with probability } 1 - \epsilon \\ \operatorname{uniform}(U) & \text{with probability } \epsilon. \end{cases}$$

- Add new trajectories to the dataset

$$D \leftarrow D \cup D^k.$$

- Update weights to  $q^{k+1}$  using all past data  $D$  using (9.4).

Compare this algorithm to policy-gradient-based methods which throw away the data from the previous iteration. Indeed, when we want to compute the gradient  $\nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta^k}} [R(\tau)]$ , we should sample trajectories from current weights  $\theta^k$ , we cannot use trajectories from some old weights. In contrast,

in Q-Learning, we maintain a cumulative dataset  $D$  that contains trajectories from all the past  $\epsilon$ -greedy controllers and use it to find new weights of the Q-function. We can do so because of the powerful Bellman equation, Q-Iteration is learning the *optimal* value function and no matter what dataset (9.4) is evaluated upon, if the error is zero, we are guaranteed that Q-function learned is the optimal one. Policy gradients do not use the Bellman equation and that is why they are so inefficient. This is also the reason Q-Learning with a replay buffer is called “off-policy” learning because it learns the optimal controller even if the data that it uses comes from some other non-optimal controller (the exploratory controller or the  $\epsilon$ -greedy controller).

Using off-policy learning is an old idea, the DQN paper which demonstrated very impressive results on Atari games using RL brought it back into prominence.

**Setting a good value of  $\epsilon$  for exploration is critical** Towards the beginning of training, we want a large value for  $\epsilon$  to gather diverse data from the environment. As training progresses, we want to reduce  $\epsilon$  because presumably we have a few good control trajectories that result in good returns and can focus on searching the neighborhood of these trajectories.

**Prioritized experience replay** is an idea where instead of sampling from the replay buffer  $D$  uniformly randomly when we fit the Q-function in (9.4), we only sample data points  $(x_t^i, u_t^i)$  which have a high Bellman error

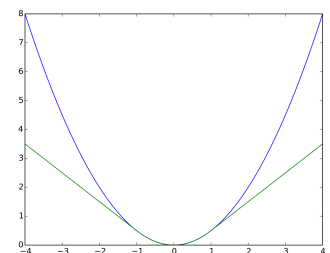
$$|q_\varphi(x_t^i, u_t^i) - r(x_t^i, u_t^i) - \gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right) \max_{u'} q_\varphi(x_{t+1}^i, u')|$$

This is a reasonable idea but is not very useful in practice for two reasons. First, if we use deep networks for parameterizing the Q-function, the network *can* fit even very complex datasets so there is no reason to not use the data points with low Bellman error in (9.4); the gradient using them will be small anyway. Second, there are a lot of hyper-parameters that determine prioritized sampling, e.g., the threshold beyond which we consider the Bellman error to be high. These hyper-parameters are quite difficult to use in practice and therefore it is a good idea to not use prioritized experience replay at the beginning of development of your method on a new problem.

**Using robust regression to fit the Q-function** There may be states in the replay buffer with very high Bellman error, e.g., the kinks in the value function for the mountain car obtained from HJB above, if we happen to sample those. For instance, these are states where the controller “switches” and is discontinuous function of state  $x$ . In these cases, instead of these few states dominating the gradient for the entire dataset, we can use ideas in robust regression to reduce their effect on the gradient. A popular way to do so is to use a Huber-loss in place of the quadratic loss in (9.4)

$$\text{huber}_\delta(a) = \begin{cases} \frac{a^2}{2} & \text{for } |a| \leq \delta \\ \delta \left(|a| - \frac{\delta}{2}\right) & \text{otherwise.} \end{cases} \quad (9.5)$$

**i** Huber loss for  $\delta = 1$  (green) compared to the squared error loss (blue).





1 **Delayed target** Notice that the target also depends upon the weights  $\varphi$ :

$$\text{target}(x'; \varphi) := r(x, u) + \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \max_{u'} q_{\varphi}(x', u').$$

2 This creates a very big problem when we fit the Q-function. Effectively, both  
 3 the covariate and the target in (9.4) depend upon the weights of the Q-function.  
 4 Minimizing the objective in (9.4) is akin to performing least squares regression  
 5 where the targets keep changing every time you solve for the solution. This  
 6 is the root cause of why Q-Learning is difficult to use in practice. A popular  
 7 hack to get around this problem is to use some old weights to compute the  
 8 target, i.e., use the loss

$$\frac{1}{n(T+1)} \sum_{i,t} \left( q_{\varphi^k}(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi^{k'}) \right)^2. \quad (9.6)$$

9 in place of (9.4). Here  $k'$  is an iterate much older than  $k$ , say  $k' = k - 100$ .  
 10 This trick is called “delayed target”.

11 **Exponential averaging to update the target** Notice that in order to im-  
 12 plement delayed targets as discussed above we will have to save all weights  
 13  $\varphi^k, \varphi^{k-1}, \dots, \varphi^{k-100}$ , which can be cumbersome. We can however do yet  
 14 another clever hack and initialize two copies of the weights, one for the actual  
 15 Q-function  $\varphi^k$  and another for the target, let us call it  $\varphi'^k$ . We set the target  
 16 equal to the Q-function at initialization. The target copy is updated at each  
 17 iteration to be

$$\varphi'^{k+1} = (1 - \alpha)\varphi'^k + \alpha\varphi^{k+1} \quad (9.7)$$

18 with some small value, say  $\alpha = 0.05$ . The target’s weights are therefore an  
 19 exponentially averaged version of the weights of the Q-function.

20 **Why are delayed targets essential for Q-Learning to work?** There are  
 21 many explanations given why delayed targets are essential in practice but the  
 22 correct one is not really known yet.

- 23 1. For example, one reason could be that since  $q_{\varphi^k}(x, u)$  for a given state  
 24 typically increases as we train for more iterations in Q-Learning, the old  
 25 weights inside a delayed target are an underestimate of the true target.  
 26 This might lead to some stability in situations when the Q-function’s  
 27 weights  $\varphi^k$  change too quickly when we fit (9.4) or we do not have  
 28 enough data in the replay buffer yet.
- 29 2. Another reason one could hypothesize is related to concepts like self-  
 30 distillation. For example, we may write a new objective for Q-Learning  
 31 that looks like

$$\left( q_{\varphi^k}(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi^k) \right)^2 + \frac{1}{2\lambda} \|\varphi^k - \varphi^{k'}\|_2^2$$

32 where the second term is known as proximal term that prevents the  
 33 weights  $\varphi^k$  from change too much from their old values  $\varphi^{k'}$ . Proximal  
 34 objectives are more stable versions of the standard quadratic objective

in (9.4) and help in cases when one is solving Q-Learning using SGD updates.

**Double Q-Learning** Even a delayed target may not be sufficient to get Q-Learning to lead to good returns in practice. Focus on one state  $x$ . One problem arise from the max operator in (9.4). Suppose that the Q-function  $q_{\varphi^k}$  corresponds to a particularly bad controller, say a controller that picks a control

$$\operatorname{argmax}_u q_{\varphi^k}(x, u)$$

that is very different from the optimal control

$$\operatorname{argmax}_u q^*(x, u)$$

then, even the delayed target  $q_{\varphi^{k'}}$  may be a similarly poor controller. The ideal target is of course the return-to-come, or the value of the optimal Q-function  $\max_{u'} q^*(x', u')$ , but we do not know it while fitting the Q-function. The same problem also occurs if our Q-function (or its delayed version, the target) is too optimistic about the values of certain control inputs, it will consistently pick those controls in the max operator. One hack to get around this problem is to pick the maximizing control input using the non-delayed Q-function but use the value of the delayed target

$$\operatorname{target}_{\text{DDQN}}(x_{t+1}^i; \varphi'^k) = r(x, u) + \gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right) q_{\varphi'^k}(x_{t+1}^i, u'). \quad (9.8)$$

where

$$u' = \underbrace{\operatorname{argmax}_u q_{\varphi^k}(x_{t+1}^i, u)}_{\text{control chosen by the Q-function}} .$$

**Training two Q-functions** We can also train two copies of the Q-function simultaneously, each with its own delayed target and mix-and-match their targets. Let  $\varphi^{(1)k}$  and  $\varphi'^{(1)k}$  be one Q-function and target pair and  $\varphi^{(2)k}$  and  $\varphi'^{(2)k}$  be another pair. We update both of them using the following objective.

$$\begin{aligned} \text{For } \varphi^{(1)} : & \left( q^{(1)k}(x, u) - r(x, u) - \gamma \left(1 - \mathbf{1}_{\{x' \text{ is terminal}\}}\right) \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(2)k}) \right)^2 \\ \text{For } \varphi^{(2)} : & \left( q^{(2)k}(x, u) - r(x, u) - \gamma \left(1 - \mathbf{1}_{\{x' \text{ is terminal}\}}\right) \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(1)k}) \right)^2 \end{aligned} \quad (9.9)$$

Sometimes we also use only one target that is the minimum of the two targets (this helps because it is more pessimistic estimate of the true target)

$$\operatorname{target}(x') := \min \left\{ \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(1)k}), \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(2)k}) \right\} .$$

You will also see many papers train multiple Q-functions, many more than 2. In such cases, it is a good idea to pick the control for evaluation using all the

1 Q-functions:

$$u^*(x) := \operatorname{argmax}_u \sum_k q_{\varphi^{(k)}}(x, u).$$

2 rather than only one of them, as is often done in research papers.

3 **A remark on the various tricks used to compute the target** It may seem  
 4 that a lot of these tricks are about being pessimistic while computing the target.  
 5 This is our current understanding in RL and it is born out of the following ob-  
 6 servation: typically in practice, you will observe that the Q-function estimates  
 7 can become very large. Even if the TD error is small, the values  $q_{\varphi}(x, u)$  can  
 8 be arbitrarily large; see Figure 1 in **Continuous Doubly Constrained Batch**  
 9 **Reinforcement Learning** for an example in a slightly different setting. This  
 10 occurs because we pick the control that maximizes the Q-value of a particular  
 11 state  $x$  in (9.8). Effectively, if the Q-value  $q_{\varphi}(x', u)$  of a particular control  
 12  $u \in U$  is an over-estimate, the target will keep selecting this control as the  
 13 maximizing control, which drives up the value of the Q-function at  $q_{\varphi}(x, u)$   
 14 as well. This problem is a bit more drastic in the next section on continuous-  
 15 valued controls. It is however unclear how to best address this issue and design  
 16 mathematically sound methods that do not use arbitrary heuristics such as  
 17 “pessimism”.

### 18 9.3 Q-Learning for continuous control spaces

19 All the methods we have looked at in this chapter are for discrete control  
 20 spaces, i.e., the set of controls that the robot can take is a finite set. In this case  
 21 we can easily compute the maximizing control of the Q-function.

$$u^*(x) = \operatorname{argmax}_u q_{\varphi}(x, u).$$

22 Certainly a lot of real-world problems have continuous-valued controls and  
 23 we therefore need Q-Learning-based methods to handle this.

24 **Deterministic policy gradient** A natural way, although non-rigorous, to  
 25 think about this is to assume that we are given a Q-function  $q_{\varphi}(x, u)$  (we  
 26 will leave the controller for which this is the Q-function vague for now)  
 27 and a dataset  $D = \{(x_t^i, u_t^i)_{t=0}^T\}_{i=1}^n$ . We can find a deterministic feedback  
 28 controller that takes controls that lead to good values as

$$\theta^* = \max_{\theta} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=0}^T q_{\varphi}(x_t^i, u_{\theta}(x_t^i)). \quad (9.10)$$

29 Effectively we are fitting a feedback controller that takes controls  $u_{\theta^*}(x)$  that  
 30 are the maximizers of the Q-function. This is a natural analogue of the argmax  
 31 over controls for discrete/finite control spaces. Again we should think of  
 32 having a deep network that parametrizes the deterministic controller and fitting

❶ Mathematically, the fundamental problem in function-approximation-based RL is actually clear: even if the Bellman operation is a contraction for tabular RL, it need not be a contraction when we are approximating the Q-function using a neural network. Therefore minimizing TD-error which works quite well for the tabular case need not work well in the function-approximation case. There may exist other, more robust, ways of computing the Bellman fixed point  $q_{\varphi}(x, u) = r(x, u) + \max_{u'} \gamma q_{\varphi}(x', u')$  other than minimizing the the squared TD error but we do not have good candidates yet.

1 its parameters  $\theta$  using stochastic gradient descent on (9.10)

$$\begin{aligned}\theta^{k+1} &= \theta^k + \eta \nabla_{\theta} q_{\varphi}(x^{\omega}, u_{\theta^k}(x^{\omega})) \\ &= \theta^k + \eta (\nabla_u q_{\varphi}(x^{\omega}, u)) (\nabla_{\theta} u_{\theta^k}(x^{\omega}))\end{aligned}\quad (9.11)$$

2 where  $\omega$  is the index of the datum in the dataset  $D$ . The equality was obtained  
3 by applying the chain rule. This result is called the “deterministic policy  
4 gradient” and we should think of it as the limit of the policy gradient for a  
5 stochastic controller as the stochasticity goes to zero. Also notice that the term

$$\nabla_u q_{\varphi}(x^{\omega}, u)$$

6 is the gradient of the output of the Q-function  $q_{\varphi} : X \times U \mapsto \mathbb{R}$  with respect  
7 to its second input  $u$ . Such gradients can also be easily computed using  
8 backpropagation in PyTorch. It is different than the gradient of the output with  
9 respect to its weights

$$\nabla_{\varphi} q_{\varphi}(x^{\omega}, u).$$

10 **On-policy deterministic actor-critic** Let us now construct an analogue of  
11 the policy gradient method for the case of a deterministic controller. The  
12 algorithm would proceed as follows. We initialize weights of a Q-function  $\varphi^0$   
13 and weights of the deterministic controller  $\theta^0$ .

- 14 1. At the  $k^{\text{th}}$  iteration, we collect a dataset from the robot using the  
15 latest controller  $u_{\theta^k}$ . Let this dataset be  $D^k$  that consists of tuples  
16  $(x, u, x', u')$ .
- 17 2. Fit a Q-function  $q^{\theta^k}$  to this dataset by minimizing the temporal differ-  
18 ence error

$$\varphi^{k+1} = \underset{\varphi}{\operatorname{argmin}} \sum_{(x, u, x', u') \in D^k} (q_{\varphi}(x, u) - r(x, u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) q_{\varphi'}(x', u'))^2.$$

(9.12)

19 Notice an important difference in the expression above, instead of using  
20  $\max_u$  in the target, we are using the control that the current controller,  
21 namely  $u_{\theta^k}$  has taken. This is because we want to evaluate the controller  
22  $u_{\theta^k}$  and simply parameterize the Q-function using weights  $\varphi^{k+1}$ . More  
23 precisely, we hope that we have

$$q_{\varphi^{k+1}}(x, u_{\theta^k}(x)) \approx \max_u q^{\theta^k}(x, u).$$

- 24 3. We can now update the controller using this Q-function:

$$\theta^{k+1} = \theta^k + \eta \nabla_{\theta} q_{\varphi^{k+1}}(x^{\omega}, u_{\theta^k}(x^{\omega})) \quad (9.13)$$

25 This algorithm is called “on-policy SARSA” because at each iteration we draw  
26 fresh data  $D^k$  from the environment; this is the direct analogue of actor-critic  
27 methods that we studied in the previous chapter for deterministic controllers.

28 **Off-policy deterministic actor-critic methods** We can also run the above  
29 algorithm using data from an exploratory controller. The only difference is

**i** SARSA is an old algorithm in RL that is the tabular version of what we did here. It stands for state-action-reward-state-action ...

1 that the we now do not throw away the data  $D^k$  from older iterations

$$D = D^1 \cup \dots \cup D^k$$

2 and therefore have to change (9.12) to be

$$\varphi^{k+1} = \underset{\varphi}{\operatorname{argmin}} \sum_{(x,u,x',u') \in D} \left( q_{\varphi}(x,u) - r(x,u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) q_{\varphi'}(x', \underbrace{u_{\theta^k}(x')}_{\text{notice the difference}}) \right)^2. \quad (9.14)$$

3 Effectively, we are fitting the optimal Q-function using the data  $D$  but since  
 4 we can no longer take the maximum over controls directly, we plug in the  
 5 controller in the computation of the target. This is natural; we think of the  
 6 controller as the one that maximizes the Q-function when we update (9.13).  
 7 When used with deep networks, this is called the “deep deterministic policy  
 8 gradient” algorithm, it is popular by the name DDPG.

# Chapter 10

## Model-based Reinforcement Learning

### Reading

1. PILCO: A Model-Based and Data-Efficient Approach to Policy Search, <http://mlg.eng.cam.ac.uk/pub/pdf/DeiRas11.pdf>
2. Embed to Control: A Locally Linear Latent Dynamics Model for Control from Raw Images <https://arxiv.org/abs/1506.07365>
3. Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models <https://arxiv.org/abs/1805.12114>

We have seen a large number of methods which use a known model of the dynamical system to compute the control inputs, these include value/policy iteration, Linear Quadratic Regulator (LQR) and Model Predictive Control (MPC). We also saw a number of methods from the reinforcement learning literature that can work “model-free”, i.e., having access to some data from the environment in lieu of a model. On one hand, model-based methods come with some obvious challenges, if we do not know the model of the system, the controller will not be optimal and worse, it may even be unsafe; think of driving on black ice that is a thin coat of ice which develops after repeated freezing and melting of snow on asphalt. On the other hand, model-free approaches are spectacularly inefficient: policy gradient-based methods require several thousands of trajectories to train a controller and even more efficient ones such as off-policy methods require prohibitive amounts of data; recall the example of an espresso bar in New York City that makes 50 shots a day, it takes more than a month to sample more than 1000 trajectories. This has limited the reach of model-free RL methods primarily to simulation, although there are examples where these policies were run (typically after training) in the real world also; see another example at <https://ai.googleblog.com/2020/05/agile->



1 [and-intelligent-locomotion-via.html](#). Very rarely you will see RL methods  
2 being used to train robots directly.

3 It makes sense to combine model-based and model-free methods if we want  
4 to reduce the number of data required from the system to learn a controller.  
5 Such methods are typically called model-based RL methods.

## 6 10.1 Learning a model of the dynamics

7 Imagine that we have a robot with dynamics

$$x_{k+1} = f(x_k, u_k)$$

8 and obtained some data from this robot using an exploratory controller  $u_e(\cdot |$   
9  $x)$ . Let us call this dataset  $D = \{(x_t^i, u_t^i)_{t=0}^T\}_{i=1}^n$ ; it consists of  $n$  trajectories  
10 each of length  $T$  timesteps. We can fit a deep network to learn the dynamics.  
11 This involves parameterizing the unknown dynamics using a deep network  
12 with weights  $w$

$$f_w : X \times U \mapsto \mathbb{R}$$

13 and minimizing a regression error of the form

$$w^* = \underset{w}{\operatorname{argmin}} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=0}^T \|x_{t+1}^i - f_w(x_t^i, u_t^i)\|_2^2. \quad (10.1)$$

14 If the residual  $\|x_{t+1}^i - f_w(x_t^i, u_t^i)\|_2^2$  is small on average over the dataset, then  
15 we know that given some new control  $u' \neq u_t^i$ , we can, for instance, estimate  
16 the new future state  $x' = f_w(x, u')$ . In principle, we can use this model now  
17 to execute algorithms like iterated LQR to find an optimal controller. We could  
18 also imagine using this as our own simulator for the robot, i.e., instead of  
19 drawing new trajectories in model-free RL from the actual robot, we use our  
20 learned model of the dynamics to obtain more data.

21 **An inverse model of the dynamics** We can also learn what is called the  
22 inverse model of the system that maps the current state  $x_t^i$  and the next state  
23  $x_{t+1}^i$  to the control that takes the system from the former to the latter:

$$f_w^{\operatorname{inv}} : X \times X \mapsto U.$$

24 The regression error for one sample in this case would be  $\|u_t^i - f_w^{\operatorname{inv}}(x_t^i, x_{t+1}^i)\|_2^2$ .  
25 This is often a more useful model to learn from the data, e.g., if we want to  
26 use this model in a Rapidly Exploring Random Tree (RRT), we can sample  
27 states in the configuration space of the robot and have the learned dynamics  
28 guess the control between two states. Also see the paper on contact-invariant  
29 optimization ([https://homes.cs.washington.edu/~todorov/papers/MordatchSIG-](https://homes.cs.washington.edu/~todorov/papers/MordatchSIG-GRAPH12.pdf)  
30 [GRAPH12.pdf](https://homes.cs.washington.edu/~todorov/papers/MordatchSIG-GRAPH12.pdf)) and a video at [https://www.youtube.com/watch?v=mhr\\_jtQrhVA](https://www.youtube.com/watch?v=mhr_jtQrhVA)  
31 for an impressive demonstration of using an inverse model.

**i** For a quick primer on planning using a model, see the notes at [https://ocw.mit.edu/courses/aeronautics-and-astronautics/16-410-principles-of-autonomy-and-decision-making-fall-2010/lecture-notes/MIT16\\_410F10\\_lec15.pdf](https://ocw.mit.edu/courses/aeronautics-and-astronautics/16-410-principles-of-autonomy-and-decision-making-fall-2010/lecture-notes/MIT16_410F10_lec15.pdf) from Emilio Frazzoli (ETH/MIT).

### Models can be wrong at parts of state-space where we have few data

This is really the biggest concern with using models. We have seen in the chapter on deep learning that if we do not have data from some part of the state-space, there are few guarantees of the model  $f_w$  or  $f_w^{\text{inv}}$  working well for those states. A planning algorithm does not however *know* that the model is wrong for a given state. So the central question in learning a model is “how to estimate the uncertainty of the output of the model”, i.e.,

$$P(x_{k+1} \neq f_x(x_k, u_k))$$

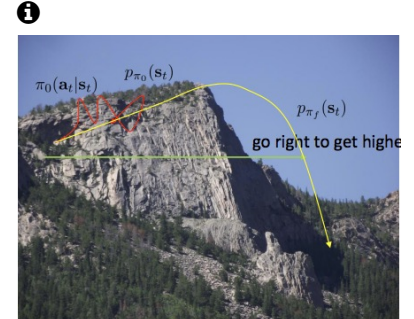
where  $x_{k+1}$  is the true next state of the system and  $f_x(x_k, u_k)$  is our prediction using the model. If we have a good estimate of such uncertainty, we can safely use the model only at parts of the state-space where this uncertainty is small.

**Sequentially querying the environment for data** We can use our ideas from DAgger and off-policy RL to improve our model iteratively by collecting more data using the controller that is being learned. Here is how it would work

1. Draw some data  $D$  from the system, fit a dynamics model  $f_w^0(x, u)$  using (10.1). Learn a feedback controller  $u^0(x)$  using any method we know so far (LQR, MPC, RL-based methods)
2. Run the learned controller  $u^0(x)$  from the real system to collect more data  $D^1$  and add it to the dataset

$$D \leftarrow D \cup D^1.$$

This is a simple mechanism that ensures that we can collect more data from the system. If the controller goes to parts of the state-space that the model is incorrect at, we get samples from such regions and iteratively improve both the learned dynamics model  $f_w^k(x, u)$  and the controller  $u^k(x)$  using this model.



## 10.2 Some model-based methods

### 10.2.1 Bagging multiple models of the dynamics

- 3 Let us look at bagging, which is a method to estimate the uncertainty of a
- 4 learning-based predictor. Bagging is short for *bootstrap aggregation*, and
- 5 can be explained using a simple experiment. Suppose we wanted to estimate
- 6 the average height  $\mu$  of people in the world. We can measure the height
- 7 of  $N$  individuals and obtain *one* estimate of the mean  $\mu$ . This is of course
- 8 unsatisfying because we know that our answer is unlikely to be the mean of
- 9 the entire population. Bootstrapping computes multiple estimates of the mean



1  $\mu_k$  over many *subsets* of the data  $N$  and reports the answer as

$$\mu := \text{mean}(\mu_k) + \text{stddev}(\mu_k).$$

2 Each subset of the data is created by sampling the original data with  $N$  samples  
3 *with replacement*. This is among the most influential ideas in statistics see

4 “Bootstrap Methods: Another Look at the Jackknife” [https://projecteuclid.org/journals/annals-](https://projecteuclid.org/journals/annals-of-statistics/volume-7/issue-1/Bootstrap-Methods-Another-Look-at-the-Jackknife/10.1214/aos/1176344552.full)  
5 [of-statistics/volume-7/issue-1/Bootstrap-Methods-Another-Look-at-the-Jackknife/10.1214/aos/1176344552.full](https://projecteuclid.org/journals/annals-of-statistics/volume-7/issue-1/Bootstrap-Methods-Another-Look-at-the-Jackknife/10.1214/aos/1176344552.full)

6 because it is a very simple and general procedure to obtain the uncertainty  
7 of the estimate. Also see a very famous paper “Bagging predictors” at  
8 <https://link.springer.com/article/10.1007/BF00058655> that invented random  
9 forests based on this idea.

10 **Training an ensemble** We are going to train multiple dynamics models  
11  $f_{w^k}(x, u)$  for  $k \in \{1, \dots, M\}$ , one each for bootstrapped versions of the  
12 training dataset  $\{D^1, \dots, D^M\}$ . Each subset  $D^k$  is built by sampling a  
13 fraction, say 60% of the data uniformly randomly from our training dataset  
14  $D$ . In other words, the bootstrapped versions of the data are not disjoint but  
15 their union is likely to be the entire training dataset. We are going to use this  
16 ensemble as follows. For each pair  $(x, u)$ , we run all models and set

$$\hat{x}' = \frac{1}{M} \sum_{k=1}^M f_{w^k}(x, u);$$

17 i.e., the ensemble predicts the next state of the robot using the mean. The  
18 important benefit of using an ensemble is that we can also get an estimate of  
19 the error in these predictions

$$\text{error in } \hat{x}' = \left( \text{var}_k f_{w^k}(x, u) \right)^{1/2}.$$

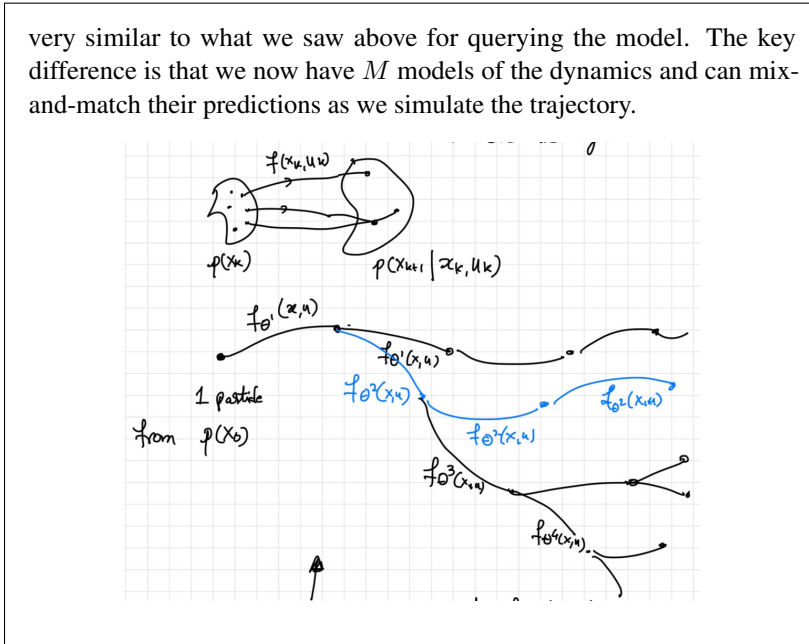
20 Different members of the ensemble are training on different datasets and make  
21 different predictions as to what the next state could be. The mis-match between  
22 them is an indicator of the error in our dynamics model. This need not be an  
23 accurate estimate of the error (i.e., the difference between the predicted  $\hat{x}'$  and  
24 the actual next state  $x'$  of the true dynamics) but is often a good proxy to use  
25 if we do bootstrapping.

26 Bagging is perhaps the most useful idea in machine learning (by far). It  
27 is always good to keep it in your mind. The winners of most high-profile  
28 machine learning competitions, e.g., the Netflix Prize  
29 ([https://en.wikipedia.org/wiki/Netflix\\_Prize](https://en.wikipedia.org/wiki/Netflix_Prize)) or the ImageNet challenge, have  
30 been bagged classifiers created by fitting multiple architectures on the same  
31 dataset.

**How to use an ensemble for hallucinating new data** The procedure is

❗ Typically, while fitting deep networks  $f_{w^k}(x, u)$  using SGD, most RL papers initialize the training process at different weights and do not perform any bootstrapping. The rationale that is usually given is that since the training process is non-convex, two models initialized at different locations train to two different predictors even if they both work on the same data. Although doing this leads to *some* notion of uncertainty, it is not an entirely correct one and performing bootstrapping will always give better estimates.

very similar to what we saw above for querying the model. The key difference is that we now have  $M$  models of the dynamics and can mix-and-match their predictions as we simulate the trajectory.



## 1 PILCO

2 A powerful Gaussian Process-based algorithm to incorporate uncertainty  
3 in the predictions of a learned dynamics model is “PILCO: A Model-Based  
4 and Data-Efficient Approach to Policy Search”

5 <http://mlg.eng.cam.ac.uk/pub/pdf/DeiRas11.pdf>. Instead of using bootstrap-  
6 ping of an ensemble to estimate the uncertainty, this algorithm explicitly  
7 models the uncertainty as

$$p(x_{k+1} | x_k, u_k) = N(x_k + E[\Delta_k], \text{var}(\Delta_k))$$

8 where  $\Delta_k = x_{k+1} - x_k$  in the training data, using a Gaussian Process.

9 See a preliminary but great tutorial on Gaussian Processes at [https://distill.pub/2019/visual-  
10 exploration-gaussian-processes](https://distill.pub/2019/visual-exploration-gaussian-processes). PILCO is a complicated algorithm to imple-  
11 ment but you can see the source code by the original authors at <https://mloss.org/software/view/508>.

12 You can also look at the paper titled “BayesRace: Learning to race au-  
13 tonomously using prior experience” <https://arxiv.org/abs/2005.04755> which  
14 uses model-based RL in

## 15 10.2.2 Model-based RL in the latent space

# 1 Bibliography

- 2 Censi, A. (2016). A class of co-design problems with cyclic constraints and their solution. *IEEE Robotics and*  
3 *Automation Letters*, 2(1):96–103.
- 4 Turing, A. M. (2009). Computing machinery and intelligence. In *Parsing the Turing Test*, pages 23–65.  
5 Springer.