

ESE 650: Learning in Robotics

Spring 2023

Instructor

Pratik Chaudhari pratikac@seas.upenn.edu

Teaching Assistants

Jianning Cui (cuijn)

Swati Gupta (gswati)

Chris Hsu (chsu8)

Gaurav Kuppaa (gakuppaa)

Alice Kate Li (alicekl)

Pankti Parekh (pankti81)

Aditya Singh (adiprs)

Haoxiang You (youhaox)

April 19, 2023

Contents

1	What is Robotics?	5
1.1	Perception-Learning-Control	6
1.2	Goals of this course	7
1.3	Some of my favorite robots	7
2	Introduction to State Estimation	9
2.1	A review of probability	9
2.1.1	Random variables	12
2.2	Using Bayes rule for combining evidence	16
2.2.1	Coherence of Bayes rule	18
2.3	Markov Chains	19
2.4	Hidden Markov Models (HMMs)	21
2.4.1	The forward algorithm	24
2.4.2	The backward algorithm	25
2.4.3	Bayes filter	26
2.4.4	Smoothing	27
2.4.5	Prediction	28
2.4.6	Decoding: Viterbi's Algorithm	28
2.4.7	Shortest path on a Trellis graph	32
2.5	Learning an HMM from observations	33
3	Kalman Filter and its variants	37
3.1	Background	37
3.2	Linear state estimation	39
3.2.1	One-dimensional Gaussian random variables	40
3.2.2	General case	41
3.2.3	Incorporating Gaussian observations of a state	42
3.2.4	An example	45
3.3	Background on linear and nonlinear dynamical systems	45
3.3.1	Linear systems	46
3.3.2	Linear Time-Invariant (LTI) systems	47
3.3.3	Nonlinear systems	48
3.4	Markov Decision Processes (MDPs)	48
3.4.1	Back to Hidden Markov Models	51
3.5	Kalman Filter (KF)	52
3.5.1	Step 0: Observing that the state estimate at any timestep should be a Gaussian	53
3.5.2	Step 1: Propagating the dynamics by one timestep	53

3.5.3	Step 2: Incorporating the observation	54
3.5.4	Discussion	55
3.6	Extended-Kalman Filter (EKF)	56
3.6.1	Propagation of statistics through a nonlinear transformation	57
3.6.2	Extended Kalman Filter	59
3.7	Unscented Kalman Filter (UKF)	61
3.7.1	Unscented Transform	63
3.7.2	The UT with tuning parameters	65
3.7.3	Unscented Kalman Filter (UKF)	66
3.7.4	UKF vs. EKF	67
3.8	Particle Filters (PFs)	68
3.8.1	Importance sampling	68
3.8.2	Resampling particles to make the weights equal	70
3.8.3	Particle filtering: the algorithm	72
3.8.4	Example: Localization using particle filter	73
3.8.5	Theoretical insight into particle filtering	74
3.9	Discussion	77
4	Rigid-body transforms and mapping	78
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	97
4.6	Discussion	98
5	Dynamic Programming	100
5.1	Formulating the optimal control problem	101
5.2	Dijkstra's algorithm	102
5.2.1	Dijkstra's algorithm in the backwards direction	104
5.3	Principle of Dynamic Programming	104
5.3.1	Q-factor	106
5.4	Stochastic dynamic programming: Value Iteration	107
5.4.1	Infinite-horizon problems	110
5.4.2	Dynamic programming for infinite-horizon problems	113
5.4.3	An example	114
5.4.4	Some theoretical results on value iteration	115
5.5	Stochastic dynamic programming: Policy Iteration	117
5.5.1	An example	119

6	Linear Quadratic Regulator (LQR)	121
6.1	Discrete-time LQR	121
6.1.1	Solution of the discrete-time LQR problem	124
6.2	Hamilton-Jacobi-Bellman equation	126
6.2.1	Infinite-horizon HJB	128
6.2.2	Solving the HJB equation	129
6.2.3	Continuous-time LQR	130
6.3	Stochastic LQR	132
6.4	Linear Quadratic Gaussian (LQG)	133
6.4.1	(Optional material) The duality between the Kalman Filter and LQR	135
6.5	Iterative LQR (iLQR)	136
6.5.1	Iterative LQR (iLQR)	138
7	Imitation Learning	141
7.1	A crash course in supervised learning	143
7.1.1	Fitting a machine learning model	145
7.1.2	Deep Neural Networks	146
7.2	Behavior Cloning	149
7.2.1	Behavior cloning with a stochastic controller	151
7.2.2	KL-divergence form of Behavior Cloning	152
7.2.3	Some remarks on Behavior Cloning	153
7.3	DAGger: Dataset Aggregation	154
8	Policy Gradient Methods	157
8.1	Standard problem setup in RL	158
8.2	Cross-Entropy Method (CEM)	159
8.2.1	Some remarks on sample complexity of simulation-based methods	161
8.3	The Policy Gradient	162
8.3.1	Reducing the variance of the policy gradient	164
8.4	An alternative expression for the policy gradient	166
8.4.1	Implementing the new expression	168
8.5	Actor-Critic methods	168
8.5.1	Advantage function	170
8.6	Discussion	171
9	Q-Learning	173
9.1	Tabular Q-Learning	173
9.1.1	How to perform exploration in Q-Learning	176
9.2	Function approximation (Deep Q Networks)	177
9.2.1	Embellishments to Q-Learning	179
9.3	Q-Learning for continuous control spaces	183
10	Model-based Reinforcement Learning	186
10.1	Learning a model of the dynamics	187
10.2	Some model-based methods	188
10.2.1	Bagging multiple models of the dynamics	188
10.2.2	Model-based RL in the latent space	190

11 Offline Reinforcement Learning	191
11.1 Why is offline reinforcement learning difficult?	192
11.2 Regularized Bellman iteration	194
11.2.1 Changing the fixed point of the Bellman iteration to be more conservative	194
11.2.2 Estimating the uncertainty of the value function	195
12 Meta-Learning	196
12.1 Problem formulation for image classification	198
12.1.1 Fine-tuning	199
12.1.2 Prototypical networks	199
12.1.3 Model-agnostic meta-learning (MAML)	201
12.2 Problem formulation for meta-RL	203
12.2.1 A context variable	204
12.2.2 Discussion	205
Bibliography	207

Chapter 1

What is Robotics?

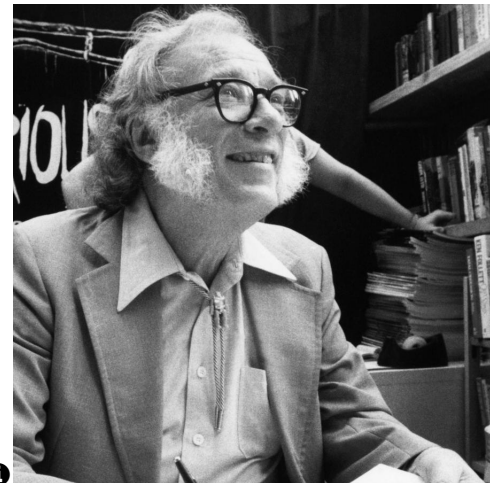
Reading

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

The word *robotics* was first used by a Czech writer Karel Capek in a play named “Rossum’s Universal Robots” where the owner of this company Mr. Rossum builds robots, i.e., agents who do forced labor, and effectively an artificial man. The word was popularized by Isaac 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. 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 Collins that one of her co-workers is infatuated with her. However, when



1 she confronts RB-34 later by pointing out that lying to people can end up
2 hurting them, the robot experiences a logical conflict within its laws and
3 becomes unresponsive.

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

6 “Robotics is the science of perceiving and manipulating
7 the physical world through computer-controlled mechanical
8 devices.” — Sebastian Thrun in *Probabilistic Robotics*

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

17 I find the Third Law really insightful to understand intelligence as well.
18 Let us define intelligence as the ability of an organism to survive¹. We
19 will all agree that trees are less intelligent than animals, an ant is less than
20 intelligent than a dog, which is less intelligent than a human. A program
21 like AlphaGo is not very intelligent because you can disable it by simply
22 switching it off. A key indicator of intelligence is the ability to sense
23 possible harm and take actions to change the outcome.

24 Robotics is Embodied Artificial Intelligence.

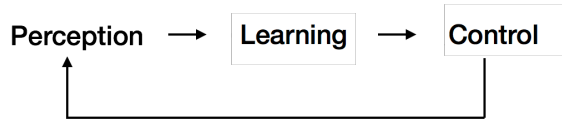
25 A robot is a machine that senses its environment using sensors,
26 interacts with this environment using actuators to perform a
27 given task and does so efficiently using previous experience
28 of performing similar tasks.

29 We will cover the fundamentals of these three aspects of robotics:
30 perception, planning and learning.

31 1.1 Perception-Learning-Control

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

¹feel free to come up with another definition



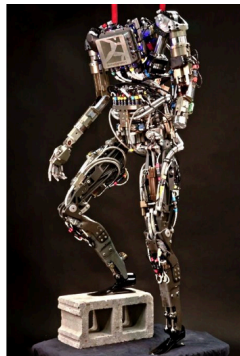
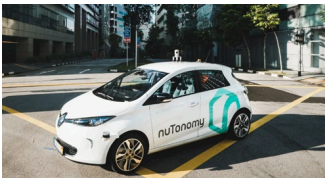
1.2 Goals of this course

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 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



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

Chapter 2

Introduction to State Estimation

Reading

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

2.1 A review of probability

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

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

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

We want to pick this set to be right granularity to answer relevant questions, e.g., it is correct but not very useful for Ω to be the position of all the

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

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

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

$$P(A)$$

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

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

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

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

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

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

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

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

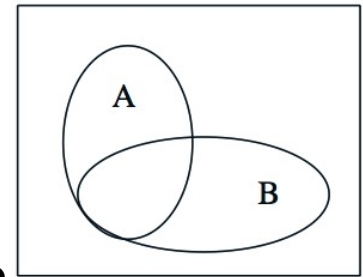
$$P(A | B) = \frac{P(A \cap B)}{P(B)}; \quad (2.1)$$

22 the probability is undefined if $P(B) = 0$. Using this definition, we can
 23 compute the probability of events like "what is the probability of rolling a
 24 2 on a die given that an even number was rolled".

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

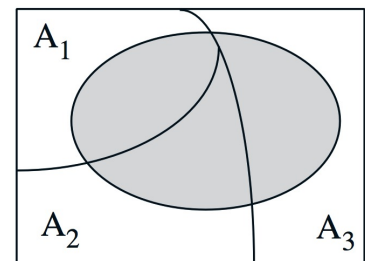
$$A_i \cap A_j = \emptyset \forall i, j, \text{ and } \bigcup_i A_i = \Omega$$

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



❶

❶ Partitioning the sample space



1 **Bayes' rule** Imagine that instead of someone telling us that the condi-
 2 tioning 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
 4 compute the conditional probability $P(B | A_i)$ using (2.1). Say we run
 5 our experiment and observe that B occurred, how would our belief on the
 6 events A_i change? In other words, we wish to compute

$$P(A_i | B).$$

7 This is the subject of Bayes' rule.

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

8 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

$$P(A \cap B) = P(A) P(B). \quad (2.4)$$

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.
 15 The sample space is easy enough to understand, e.g., $\Omega = [0, 1]$ for your
 16 score at the end of this course. We however run into difficulties if we
 17 define the probability of general subsets of Ω in terms of the probabilities
 18 of elementary outcomes (*elements* of Ω). For instance, if we wish to
 19 model all elements $\omega \in \Omega$ to be equally likely, we are forced to assign each
 20 element ω a probability of zero (to be consistent with the second axiom of
 21 probability). This is not very helpful in determining the probability of the
 22 score being 0.9. If you instead assigned some small non-zero number to
 23 $P(\omega_i)$, then we have undesirable 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
 26 sets. While we would ideally like to be able to specify the probability of
 27 every subset of Ω , it turns out that we cannot do so in a mathematically
 28 consistent way. The trick then is to work with a smaller object known as a
 29 σ -algebra, that is the set of "nice" subsets of Ω .

Given a sample space Ω , a σ -algebra \mathcal{F} (also called a σ -field) is a collection of subsets of Ω 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, σ -algebra is a collection of subsets of Ω that is closed under complement and countable unions. The pair (Ω, \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 σ -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 σ -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 σ -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 be $\omega \in \Omega$ when the realized value of the random variable is a particular number x equal to 5.

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 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

❗ 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.

❓ 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} .

❗ 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 subsets of Ω belong to \mathcal{F} .

1 number of values. For instance, if X is the number of coin tosses until the
 2 first head, if we assume that our tosses are independent $P(H) = p > 0$,
 3 then we have

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

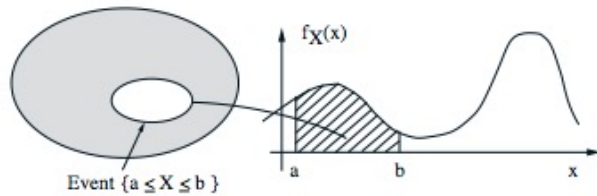
4 for all $k = 1, 2, \dots$. This is what is called a geometric probability mass
 5 function.

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

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

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

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



13

14 If $F_X(x)$ is the CDF of an r.v. X and X takes values in \mathbb{R} , the
 15 probability density function (PDF) $f_X(x)$ (or sometimes also denoted by
 16 $p_X(x)$) is defined to be

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

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

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

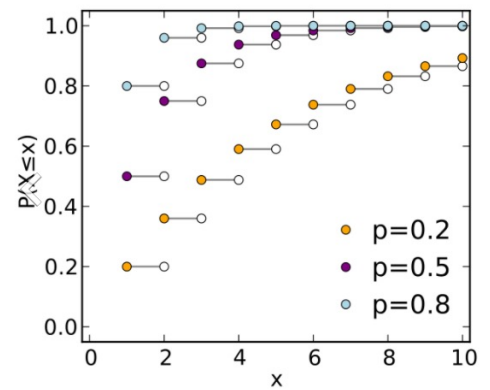
19 This leads to the following interpretation of the probability density func-
 20 tion:

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

21 **Expectation and Variance** The expected value of a random variable X
 22 is

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

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 random variable, since the values that X takes belong to the set of integers, the CDF is constant between any two integers.

1 and denotes the center of gravity of the probability mass function. Roughly
 2 speaking, it is the average of a large number of repetitions of the same
 3 experiment. Expectation is a linear, i.e.,

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

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

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

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

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

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

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

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

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

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

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

13 the definition of variance remains the same.

14 **Joint distributions** We often wish to think of the joint probability distri-
 15 bution of multiple random variables, say the location of an autonomous car
 16 in all three dimensions. The cumulative distribution function associated
 17 with this is therefore

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

18 Just like we have the probability density of a single random variable, we
 19 can also write the joint probability density of multiple random variables
 20 $f_{X,Y,Z}(x, y, z)$. In this case we have

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

- 1 The joint probability density factorizes if two random variables are
2 independent:

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

- 3 Two random variables are uncorrelated if and only if

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

- 4 Note that independence implies uncorrelatedness, they are not equivalent.
5 The covariance is defined as

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

- 6 **Conditioning** As we saw before, for a single random variable X we
7 have

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

- 8 For two random variables, by analogy we would like

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

- 9 The conditional probability density function of X given Y is defined to be

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

- 10 For any given y , the conditional PDF is a normalized section of the joint
11 PDF, as shown below.

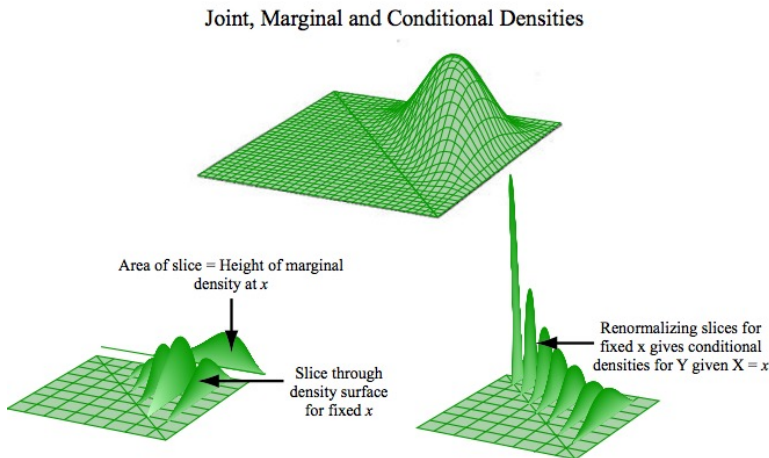


Image by MIT OpenCourseWare, adapted from
Probability, by J. Pittman, 1999.

1 **Continuous form of Bayes rule** We can show using the definition of
 2 conditional probability that

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

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

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

4 2.2 Using Bayes rule for combining evidence

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



7

8 We will abstract each observation by the sensors of the robot as a
 9 random variable Y . This could be the image from its camera after running
 10 some algorithm to check the state of the door, the reading from a laser
 11 sensor (if the time-of-flight of the laser is very large then the door is open),
 12 or any other mechanism. We have two kinds of conditional probabilities
 13 in this problem

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

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

14 The second one is called a causal quantity because the specific Y we
 15 observe depends upon whether the door is open or not. The first one is
 16 called a diagnostic quantity because using this observation Y we can infer
 17 the state of the environment, i.e., whether the door is open or not. Next
 18 imagine how you would calibrate the sensor in a lab: for each value of
 19 the state of the door open, not open you would record all the different
 20 observations received Y and calculate the conditional probabilities. The
 21 causal probability is much easier to calculate in this context, one may
 22 even use some knowledge of elementary physics to model the probability
 23 $P(Y | \text{open})$, or one may count the number of times the observation is
 24 $Y = y$ for a given state of the door.

25 Bayes rule allows us to transform causal knowledge into diagnostic

1 knowledge

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

2 Remember that the left hand side (diagnostic) is typically something that
 3 we desire to calculate. Let us put some numbers in this formula. Let
 4 $P(Y | \text{open}) = 0.6$ and $P(Y | \text{not open}) = 0.3$. We will imagine that the
 5 door is open or closed with equal probability: $P(\text{open}) = P(\text{not open}) =$
 6 0.5 . We then have

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

7 Notice something very important, the original (prior) probability of the
 8 state of the door was 0.5. If we have a sensor that fires with higher
 9 likelihood if the door is open, i.e., if

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

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

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

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

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

21 We now have

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

22 where

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

23 is the denominator. We can now expand the diagnostic probability on the

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

1 right-hand side recursively to get

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

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

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

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

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

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

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

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

8 2.2.1 Coherence of Bayes rule

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

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

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

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

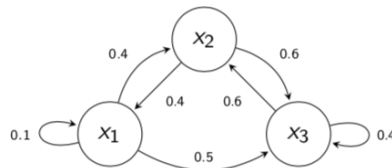
2.3 Markov Chains

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

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

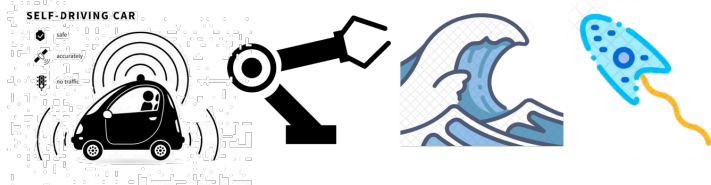
You can see the animations at <https://setosa.io/ev/markov-chains> to build more intuition.

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?

1 Consider the paramecium above. Its position depends upon a large
 2 number of factors: its own motion from the previous time-step but also
 3 the viscosity of the material in which it is floating around. One may
 4 model the state of the environment around the paramecium as a liquid
 5 whose molecules hit thousands of times a second, essentially randomly,
 6 and cause disturbances in how the paramecium moves. Let us call this
 7 disturbance “noise in the dynamics”. If the motion of the molecules of the
 8 liquid has some correlations (does it, usually?), this induces correlations in
 9 the position of the paramecium. The position of the organism is no longer
 10 Markov. This example is important to remember, the Markov property
 11 defined above also implies that the noise in the state transition matrix is
 12 independent.

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

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

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

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

19 is such that

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

20 We can now write the above formula as

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

21 ¹ and compute the distribution $\pi^{(t)}$ for all times

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

22 The numbers $P(X_k = x_i)$ stop changing with time k . Under certain tech-
 23 nical conditions, the distribution $\pi^{(\infty)}$ is unique (single communicating

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

1 class for a Markov chain with a finite number states). We can compute
2 this invariant distribution by writing

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

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

6 **Example 2.1.** Consider a Markov chain on two states where the transition
7 matrix is given by

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

8 The invariant distribution is

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

9 Note that the constraint for π being a probability distribution, i.e., $\pi^{(1)} +$
10 $\pi^{(2)} = 1$ is automatically satisfied by the two equations. We can solve for
11 $\pi^{(1)}, \pi^{(2)}$ to get

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

12 2.4 Hidden Markov Models (HMMs)

13 2

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

19 Hidden Markov Models (HMMs) are an abstraction to reason about
20 observations of the state of a Markov chain. An HMM is a sequence
21 of random variables Y_1, Y_2, \dots, Y_n such that the distribution of Y_k only
22 depends upon the hidden state X_k of the associated Markov chain.

²Parts 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
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

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

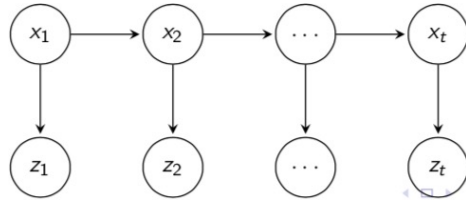


Figure 2.1: A Hidden Markov Model with the underlying Markov chain, the observation at time k only depends upon the hidden state at that time instant. Ignore the notation Z_1, \dots, Z_t we will denote the observations by Y_k .

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

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

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

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

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

16 The state transition probabilities of the associated Markov chain are

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

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

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

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

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

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

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

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

- 10 3. **Prediction:** Given observations up to time k , compute the distribu-
 11 tion of the state at a time $j > k$

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

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

🔗 Why is this true?

- 17 4. **Decoding:** Find the most likely state trajectory X_1, \dots, X_k that
 18 maximizes the probability

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

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

- 28 5. **Likelihood of observations:** Given the observation trajectory,
 29 Y_1, \dots, Y_k , compute the probability

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

30 As you may recall, this is the denominator that we need for the
 31 recursive application of Bayes rule. It is made difficult by the fact
 32 that we do not know the state trajectory X_1, \dots, X_k corresponding
 33 to these observations.

34 These problems are closely related with each other and we will next dig
 35 deeper into them. We will first discuss two building blocks, called the

1 forward and backward algorithm that together help solve all the above
2 problems.

3 2.4.1 The forward algorithm

4 Consider the problem of computing the likelihood of observations. We
5 can certainly write

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

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

8

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

Forward algorithm We can simplify the above computation using the 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

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

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 .

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

2.4.2 The backward algorithm

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

$$Y_1, \dots, Y_t$$

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

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

Notice that the backward variables β_k with the conditioning on $X_k = x$ are slightly different than the forward variables α_k which are the joint probability of the observation trajectory and $X_k = x$.

Backward algorithm We can compute the variables $\beta_k(x)$ recursively again as follows.

1. Initialize

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

This simply indicates that since we are at 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

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

Implementing the forward and backward algorithms in practice The update equations for both α_k and β_k can be written using a matrix vector multiplication. We maintain the vectors

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

What is the computational complexity of the Forward algorithm?

What is the computational complexity of running the backward algorithm?

1 and can write the updates as

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

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

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

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

11 2.4.3 Bayes filter

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

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

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

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

15 where since $\mathbb{P}(X_k = x \mid Y_1, \dots, Y_k)$ is a legitimate probability distribu-
16 tion on x , we have

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

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

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

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

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

26 This is an extremely important concept to remember, in state-estimation we
27 are always interested in computing the state given all available observations.

1 In the same context, is the following statement true?

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

2 2.4.4 Smoothing

3 Given observations till time t , we would like to compute

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

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

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

5 Study the first step carefully, the numerator is *not* equal to $\alpha_k(x)$ because
6 observations go all the way till time t . The final step uses both the Markov
7 and the HMM properties: future observations Y_{k+1}, \dots, Y_t depend only
8 upon future states X_{k+1}, \dots, X_t (HMM property) which are independent
9 of the past observations and states give the current state $X_k = x$ (Markov
10 property).

11 Smoothing can therefore be implemented by running the forward
12 algorithm to update α_k from $k = 1, \dots, t$ and the backward algorithm to
13 update β_k from time $k = t, \dots, 1$.

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

🔗 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?

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

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

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

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

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

1 Run the forward backward algorithm and see that

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

2 and

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

3 Using these, we can now compute the filtering and the smoothing state
4 distributions, let us denote them by π^f and π^s respectively.

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

5 and

$$\pi_1^s = (0.999, 0, 0), \pi_2^s = (0.0529, 0.2328, 0.7143), \pi_3^s = (0.045, 0.2487, 0.7063)$$

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

6

7 2.4.5 Prediction

8 We would like to compute the future probability of the state give observa-
9 tions up to some time

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

10 Here is a typical scenario when you would need this estimate. Imagine
11 that you are tracking the position of a car using images from your camera.
12 You are using a deep network to detect the car in each image Y_k and since
13 the neural network is quite slow, the car moves multiple time steps forward
14 before you get the next observation. As you can appreciate, it would help
15 us compute a more accurate estimate of the conditional probability of
16 $X_k = x$ if we propagated the position of the car in between successive
17 observations using our Markov chain. This is easy to do.

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

$$\pi_{i+1} = T' \pi_i.$$

22 2.4.6 Decoding: Viterbi's Algorithm

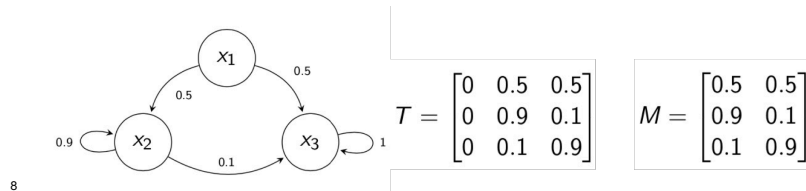
23 Both filtering and smoothing calculate the probability distribution of the
24 state at time k . For instance, after recording a few observations, we can
25 compute the probability distribution of the position of the car at each time
26 instant. How do we get the most likely trajectory of the car?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

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

1 at each instant and output

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

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



9 There are three states x_1, x_2, x_3 with known initial distribution $\pi =$
 10 $(1, 0, 0)$ and transition probabilities and observations given by matrices
 11 T, M respectively. Let us say that we only have two observations $\{y_2, y_3\}$
 12 this time and get the observation sequence

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

13 from our sensor. The filtering estimates are as follows.

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

14
 15 The most likely state at each instant is marked in blue. The point-wise
 16 most likely sequence of states is

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

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

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

1

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

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

3 Because the smoothing estimate at time k also takes into account the
 4 observations from the future $t > k$, it effectively eliminates the impossible
 5 transition from $x_3 \rightarrow x_2$. This is still not however the most likely
 6 trajectory.

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

$$\delta_k(x) = \max_{(x_1, \dots, x_{k-1})} \text{P}(X_1 = x_1, \dots, X_{k-1} = x_{k-1}, X_k = x, Y_1, \dots, Y_k); \quad (2.13)$$

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

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

12 the joint probability that the most likely trajectory ends up at state x at
 13 time $k + 1$ is the maximum of among the joint probabilities that end up
 14 at any state x' at time k multiplied by the one-step state transition $T_{x'x}$
 15 and observation $M_{x,y_{k+1}}$ probabilities. We would like to iterate upon this
 16 identity to find the most likely path. The key idea is to maintain a pointer
 17 to the parent state $\text{parent}_k(x)$ of the most likely trajectory, i.e., the state
 18 from which you could have reached $X_k = x$ given observations. Let us
 19 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

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

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

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

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

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

10 The most likely path is the one that ends in 3 with joint probability 0.0432.
 11 This path is (1, 3, 3).

12 Let us also compute Viterbi's algorithm for a longer observation
 13 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

1
2 The most likely trajectory is

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

3 Notice that if we had only 8 observations, the most likely trajectory would
4 be

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

5
6 What is the computational complexity of Viterbi's algorithm? It is
7 linear in the time-horizon t and quadratic in the number of states in the
8 Markov chain. We are plucking out the most likely trajectory out of
9 $\text{card}(X)^t$ possible trajectories using the δ_k variables. Does this remind
10 you of some other problem that you may have seen before?

11 2.4.7 Shortest path on a Trellis graph

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

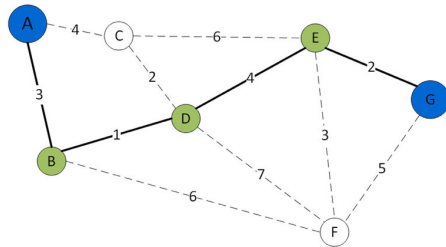


Figure 2.2: A graph with costs assigned to every edge. Dijkstra's algorithm finds the shortest path in this graph between nodes A and B using dynamic programming.

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

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 we see in this example).

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 $-\log P(X_1, X_2, X_3 |$
 3 $Y_1, Y_2, Y_3)$. The term $\log P(Y_1, Y_2, Y_3)$ does not depend on X_1, X_2, X_3
 4 and is a constant as far as the most likely path given observations is
 5 concerned. We can now write down the “Trellis” graph as shown below.

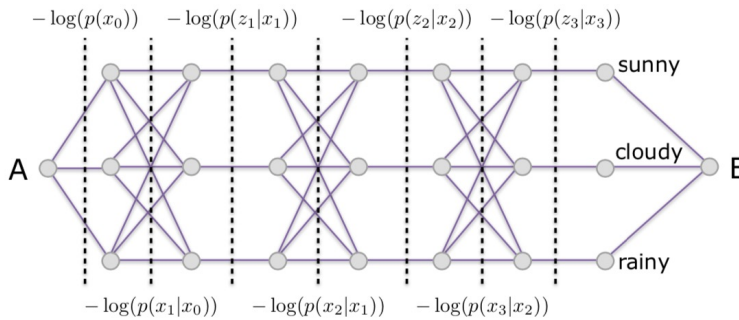


Figure 2.3: A Trellis graph for a 3-state HMM for a sequence of three observations. Disregard the subscript x_0 .

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.
 9 The edge-costs of the final three states, in this case sunny/cloudy/rainy,
 10 are zero. The costs from node A to the respective states are the log-
 11 probabilities of the initial state distribution. Dijkstra’s algorithm, which
 12 we will study in Module 2 in more detail, now gives the shortest path on the
 13 Trellis graph. This approach is the same as that of the Viterbi’s algorithm:
 14 our parent pointers $\text{parent}_k(x)$ are the parent nodes in Dijkstra’s algorithm
 15 and our delta variables $\delta_k(x)$ is the cost of each node in the Trellis graph
 16 maintained by the Dijkstra’s algorithm.

17 2.5 Learning an HMM from observations

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

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

21 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
 3 observation sequence more likely. This is the simplest instance of *learning*
 4 an HMM. The prototypical problem to imagine that our original HMM λ
 5 comes from is our knowledge of the original problem (say a physics model
 6 of the dynamics of a robot and its sensors). Given more data, namely
 7 the observations, we want to update this model. The most natural way to
 8 update the model is to maximize the likelihood of observations given our
 9 model, i.e.,

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

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

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

14 Let us consider a simple problem. We are going to imagine that the
 15 FBI is trying to catch the dangerous criminal Keyser Soze who is known
 16 to travel between two cities Los Angeles (LA) which will be state x_1 and
 17 New York City (NY) which will be state x_2 . The FBI initially have no clue
 18 about his whereabouts, so their initial belief on his location is uniform
 19 $\pi = [0.5, 0.5]$. His movements are modeled using a Markov chain

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

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

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

23 Say that they received an observation sequence of 20 periods

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

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

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

26 the smoothing probability. We can also compute the most likely state
 27 trajectory he could have taken given our observations using decoding. Let
 28 us focus on the 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

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
 5 how should we update the model λ given this data. We are going to learn
 6 the entries 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,
 9 i.e.,

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

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

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

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

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

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

13 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
6 the HMM? Recall that we are trying to compute the best HMM given the
7 observations, 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
11 expected 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
14 gone down because the Markov chain does not have those states.

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

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

Chapter 3

Kalman Filter and its 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

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

3.1 Background

Multi-variate random variables and linear algebra For d -dimensional 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
 2 whether X, Y are correlated. The covariance matrix of a random variable
 3 is defined as

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

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

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

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

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

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

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

- 17 • For matrices A, B we have

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

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

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

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

21 This operation can be thought of as taking the inner product between
 22 two matrices.

23 **Gaussian/Normal distribution** We will spend a lot of time working
 24 with the Gaussian/Normal distribution. The multi-variate d -dimensional
 25 Normal distribution has the probability density

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

26 where $\mu \in \mathbb{R}^d$, $\Sigma \in \mathbb{R}^{d \times d}$ denote the mean and covariance respectively.
 27 You 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
2 integrates to 1.

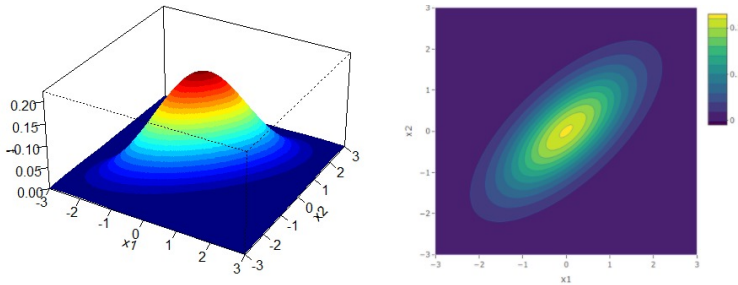


Figure 3.1: Probability density (left) and iso-probability contours (right) of a bi-variate Normal distribution. Warm colors denote regions of high probability.

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

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

4 with covariance

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

5 where

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

6 the matrix Σ_{YX} is defined similarly. If X, Y are independent (or uncorre-
7 lated) 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
10 mean

$$E[Y] = E[AX] = AE[X] = A\mu_X$$

11 and covariance

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

12 This is an important result that you should remember.

13 3.2 Linear state estimation

14 With that background, let us now look at the basic estimation problem.

15 Let $X \in \mathbb{R}^d$ denote the true state of a system. We would like to build an

1 estimator for this state, this is denote by

$$\hat{X}.$$

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

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

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

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

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

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

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

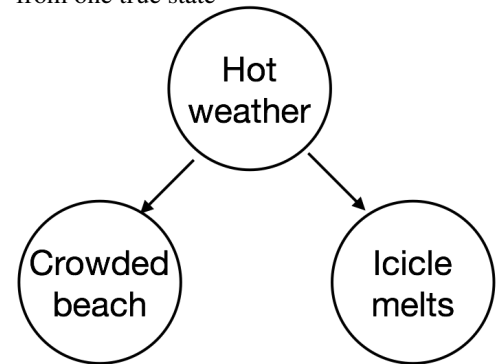
18 We would like to combine the two to obtain a *better* estimate of what the
19 state could be. *Better* can mean many different quantities depending upon
20 the problem but in general in this course we are interested in improving
21 the error covariance. Our goal is then

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 $\text{tr}(\Sigma_{\tilde{X}})$.

① Conditionally independent observations from one true state



22 3.2.1 One-dimensional Gaussian random variables

23 Consider the case when $\hat{X}_1, \hat{X}_2 \in \mathbb{R}$ are Gaussian random variables with
24 means μ_1, μ_2 and variances σ_1^2, σ_2^2 respectively. Assume that both are
25 unbiased estimators of $X \in \mathbb{R}$. Let us combine them linearly to obtain a
26 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
2 new estimator to be unbiased, so

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

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

23 **3.2.2 General case**

24 Let us now perform the same exercise for multi-variate Gaussian random
25 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
 2 for 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
 5 the trace of this covariance matrix. We know that the original covariances
 6 Σ_1 and Σ_2 are symmetric. We will use the following identity for the
 7 partial derivative 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
 9 to 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
 23 $C \in \mathbb{R}^{p \times d}$ being something that is unique to the particular sensor. This
 24 observation is not precise and we will model the sensor as having zero-
 25 mean Gaussian noise

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

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

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

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

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

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

11 to get that

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

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

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

15 Let us now optimize K as before to compute the estimator with minimal
 16 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 .

17 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](#)
- 2 who 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 Σ' 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.

3.2.4 An example

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

$$y_i = x + \nu_i$$

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

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

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

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

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

The updated mean using (3.9) again

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

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

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

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

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

3.3 Background on linear and nonlinear dynamical systems

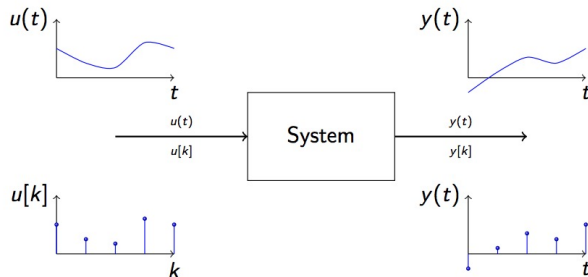
The true state X need not be static. We will next talk about models for how the state of the world evolves using ideas in dynamical systems.

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

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

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

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



8

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.

9 **State of a system** We know that if the system is causal, in order to
 10 compute its output at a time t_0 , we only need to know all the input from
 11 time $t = (-\infty, t_0]$. This is a lot of information. The concept of a state,
 12 about which we have been cavalier until now helps with this. The state
 13 $x(t_1)$ of a causal system at time t_1 is the information needed, together
 14 with the input u between times t_1 and t_2 to *uniquely compute* the output
 15 $y(t_2)$ at time t_2 , for all times $t_2 \geq t_1$. In other words, the state of a system
 16 summarizes the whole history of what happened between $(-\infty, t_1)$.

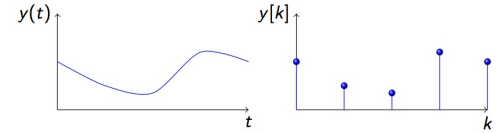
17 Typically the state of a system is a d -dimensional vector in \mathbb{R}^d . The
 18 dimension of a system is the minimum d required to define a state.

19 3.3.1 Linear systems

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

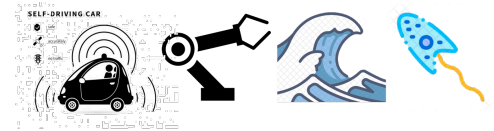
$$\begin{aligned} u_1 &\rightarrow y_1 \\ u_2 &\rightarrow y_2 \\ au_1 + bu_2 &\rightarrow ay_1 + by_2. \end{aligned}$$

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



❓ 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 Linearity is a very powerful property. For instance, it suggests that if
 2 we can decompose a complicated input into the sum of simple signals,
 3 then the output of the system is also a sum of the outputs of these simple
 4 signals. For example, if we can write the input as a Fourier series
 5 $u(t) = \sum_{i=0}^{\infty} a_i \cos(it) + b_i \sin(it)$ we can pass each of the terms in
 6 this summation to system and get the output of $u(t)$ by summing up the
 7 individual outputs.

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

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

or $m\ddot{z} + c\dot{z} + kz = u$

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

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

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

14 We can now rewrite the dynamics as

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

15 3.3.2 Linear Time-Invariant (LTI) systems

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

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

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

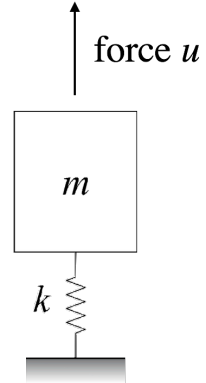
24 Since the state at time t encapsulates everything that happened to
 25 the system due to the inputs $\{u(-\infty), u(t)\}$, we can say that the system
 26 computes its output $y(t)$ as a function of the state $x(t)$ and the latest input
 27 $u(t)$

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

28 If this function is linear we have

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

ⓘ A second-order spring mass system



1 The pair of equations (3.10) and (3.11) together are the so-called state-
 2 space model of an LTI system. The development for discrete-time systems
 3 is completely analogous, we will have

$$\begin{aligned}x_{k+1} &= A_{\Delta t}x_k + B_{\Delta t}u_k \\y_k &= Cx_k + Du_k.\end{aligned}\tag{3.12}$$

We have used the subscript Δt to denote that these are discrete-time matrices and are different from the continuous-time ones in (3.10) and (3.11). This is an important point to keep in mind.

4 If the dynamics matrices A, B, C, D change with time, we have a
 5 time-varying system.

6 3.3.3 Nonlinear systems

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

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

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

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

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

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

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

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

18 Again the discrete-time nonlinear dynamics has a different equation than
 19 the corresponding one in (3.13).

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

20 3.4 Markov Decision Processes (MDPs)

21 Let us now introduce a concept called MDPs which is very close to
 22 Markov chains that we saw in the previous chapter. In fact, you are already
 23 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 $f(x_k, u_k)$.

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

3 1. We did not do a good job of identifying the function $f : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{X}$.

4 This may happen when you are driving a car on an icy road, if you
5 undertake the same control as you do on a clean road, you might
6 reach a different future state x_{k+1} .

7 2. We did not use the correct state-space \mathcal{X} . You could write down
8 the state of the car as given by $(x, y, \theta, \dot{x}, \dot{y}, \dot{\theta})$ where x, y are the
9 Euclidean co-ordinates of the car and θ is its orientation. This is
10 not a good model for studying high-speed turns, which are affected by
11 other quantities like wheel slip, the quality of the suspension etc.

12 We may not even know the full state sometimes. This occurs when
13 you are modeling how users interact with an online website like
14 Amazon.com, you'd like to model the change in state of the user
15 from "perusing stuff" to "looking stuff to buy it" to "buying it"
16 but there are certainly many other variables that affect the user's
17 behavior. As another example, consider the path that an airplane
18 takes to go from Philadelphia to Los Angeles. This path is affected
19 by the weather at all places along the route, it'd be cumbersome
20 to incorporate the weather to find the shortest-time path for the
21 airplane.

22 3. We did not use the correct control-space \mathcal{U} for the controller. This
23 is akin to the second point above. The gas pedal which one may
24 think of as the control input to a car is only one out of the large
25 number of 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” ϵ_k is not under our control. The quantity ϵ_k is not arbitrary however, we will assume that

1. noise ϵ_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 ϵ_k . Note that the distribution of ϵ_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

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

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

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

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.

🔗 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=ym64NFCWO> costs upwards of \$100,000. Would you model it as a stochastic dynamical system?

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

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

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

1 infinitesimal time Δt is a Gaussian random variable

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

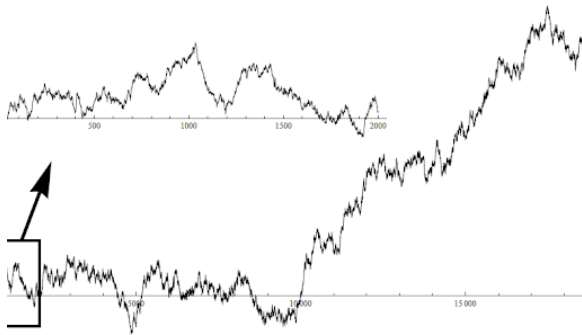


Figure 3.2: A typical Brownian motion signal $\epsilon(t)$. You can also see an animation at https://en.wikipedia.org/wiki/File:Brownian_Motion.ogv

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

6 3.4.1 Back to Hidden Markov Models

7 Since our sensors measure the state x of the world, it will be useful to think
8 of the output y of a dynamical system as the observations from Chapter
9 2. This idea neatly ties back our development of dynamical systems
10 to observations. Just like we considered an HMM with observation
11 probability

$$P(Y_k = y \mid X_k = x)$$

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

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

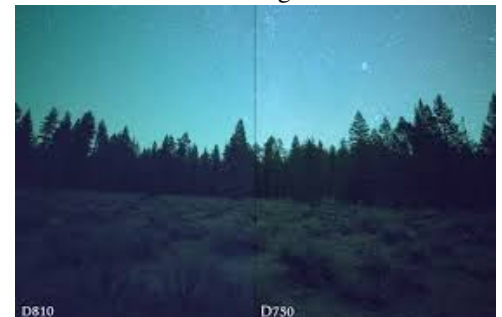
16 The noise ν_k is similar to the noise in the dynamics ϵ_k in (3.14). Analo-
17 gously, we can also have noise in the observations of a linear system

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

Hidden Markov Models with underlying MDPs/Markov chains

❓ 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?

❗ 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?

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

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

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

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 .

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.

3.5 Kalman Filter (KF)

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

Consider a linear dynamical system with linear observations

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

where the noise vectors

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

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

are both zero-mean and Gaussian with covariances R and Q respectively.

We have also assumed that $D = 0$ because typically the observations do not depend on the control.

i We will assume that the distribution of noise ϵ_k, ν_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.

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

observations

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

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.

3.5.1 Step 0: Observing that the state estimate at any timestep should be a Gaussian

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

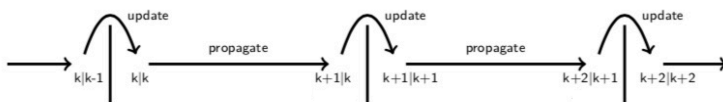
$$\hat{x}_{k|k} \sim P(x_k | y_1, \dots, y_k) \equiv N(\mu_{k|k}, \Sigma_{k|k}).$$

The subscript

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

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

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



3.5.2 Step 1: Propagating the dynamics by one timestep

Suppose we had an estimate $\hat{x}_{k|k}$ after k observations/time-steps. Since the dynamics is linear, we can use the prediction problem to compute the estimate of the state at time $k+1$ before the next observation arrives

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

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

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

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

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

5 using our calculation in (3.1).

6 3.5.3 Step 2: Incorporating the observation

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

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

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

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

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

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

15 The covariance is given by our same calculation again

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

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

The new estimate of the state is

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

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

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 $\Sigma_{k+1|k}$.

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

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

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 ϵ_k and the observation noise ν_{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

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

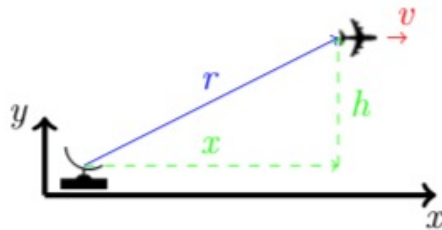
The Wikipedia webpage at https://en.wikipedia.org/wiki/Kalman_filter#Example_application,_technical gives a simple example of a Kalman Filter.

🔗 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. For most robots, both of these are 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 ν . 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

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

Since we do not really know how the plane might change its altitude, let's assume that it maintains a constant altitude

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

The above equations are our model for how the state of the airplane evolves and could of course be wrong. As we discussed, we will model the

1 discrepancy as noise.

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

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

2 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
3 Gaussian noise. The dynamics in this case is linear but the observations
4 are a nonlinear function of the state.

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

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

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

11 3.6.1 Propagation of statistics through a nonlinear trans- 12 formation

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

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

15 If we had a nonlinear function of x

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

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

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

18 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, θ are the XY-coordinates and the angle of the steering wheel respectively. This model is known as a Dubins car.

1 This gives

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

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

A simple example

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

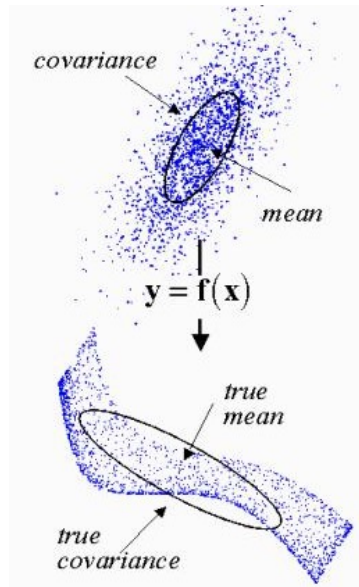
5 We have

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

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

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

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
4 of 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
6 can therefore use the Kalman Filter equations to incorporate this fake
7 observation

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

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

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

10 to get the EKF equations for incorporating one observation

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

The Extended Kalman Filter estimates the state of a nonlinear system by linearizing the dynamics and observation equations at each timestep.

1. Say we have the current estimate $\mu_{k|k}$ and $\Sigma_{k|k}$.
2. After a control input u_k the new estimate is

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

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

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

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

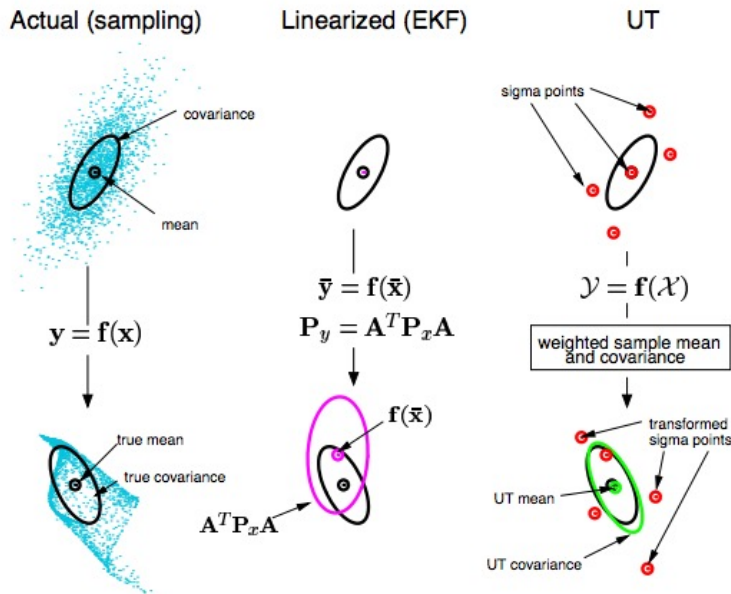
$$\begin{aligned} x_{t+\Delta t} &= x(t) + f(x(t), u(t)) \Delta t + \epsilon \\ &\equiv f^{\text{discrete-time}}(x(t), u(t)) + \epsilon. \end{aligned}$$

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

31 3.7 Unscented Kalman Filter (UKF)

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

37 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
- 4 in top right).
- 5 2. Transform *each of the points* using the nonlinear dynamics f (red
- 6 points in bottom right).
- 7 3. Compute their mean and covariance to get $\mu_{k+1|k}$ and $\Sigma_{k+1|k}$.
- 8 Notice how the green ellipse is slightly different than the black
- 9 ellipse (which is the true mean and covariance). Both of these would
- 10 be different from the mean and covariance obtained by linearization
- 11 of f (middle column) 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.

3.7.1 Unscented Transform

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

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

$$\Sigma = SS.$$

We can compute this via diagonalization as follows.

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

We can therefore define

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

Notice that

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

We can also define the matrix square root using the Cholesky decomposition $\Sigma = LL^T$ which is numerically more stable than computing the square root using the above expression. Recall that matrices L and Σ have the same eigenvectors. Typical applications of the Unscented Transform will use this 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^{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
7 be 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 $\Sigma_x =$
9 $\begin{bmatrix} \sigma_r^2 & 0 \\ 0 & \sigma_\theta^2 \end{bmatrix}$. We would like to compute the probability distribution of
10 $y = f(x) = \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix}$ which is a polar transformation. Since x is two-
11 dimensional, we will have 4 sigma points with equal weights $w^{(i)} = 0.25$.
12 The square root in the sigma 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

$$y^{(1)} = \begin{bmatrix} r^{(1)} \cos \theta^{(1)} \\ r^{(1)} \sin \theta^{(1)} \end{bmatrix} = \begin{bmatrix} 0 \\ 1 + \sqrt{2}\sigma_r \end{bmatrix}$$

$$y^{(2)} = \begin{bmatrix} r^{(2)} \cos \theta^{(2)} \\ r^{(2)} \sin \theta^{(2)} \end{bmatrix} = \begin{bmatrix} \cos(\pi/2 + \sqrt{2}\sigma_\theta) \\ \sin(\pi/2 + \sqrt{2}\sigma_\theta) \end{bmatrix}$$

$$y^{(3)} = \begin{bmatrix} 0 \\ 1 - \sqrt{2}\sigma_r \end{bmatrix}$$

$$y^{(4)} = \begin{bmatrix} \cos(\pi/2 - \sqrt{2}\sigma_\theta) \\ \sin(\pi/2 - \sqrt{2}\sigma_\theta) \end{bmatrix}.$$

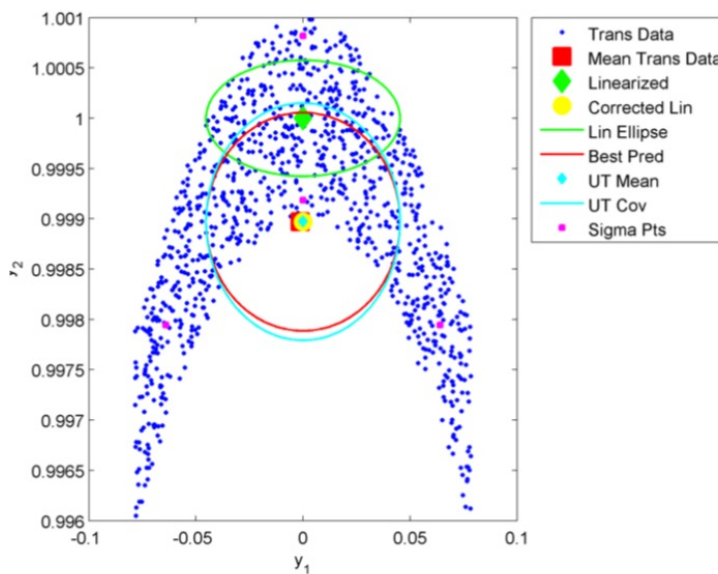


Figure 3.3: Note that the true mean is being predicted very well by the UT and is clearly a better estimate than the linearized mean.

2 3.7.2 The UT with tuning parameters

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

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

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

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

and covariances

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

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 μ_x, Σ_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

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

1 and $\mu_z = [\mu_x \quad \mu_y]$ and

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

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

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

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

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

7 and this gain leads to the new covariance

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

8 The nice thing about the Kalman gain in (3.35) is that we can compute it
9 now using expressions of Σ_{xy} and Σ_{yy} in terms of the sigma points. This
10 goes 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).

11 3.7.4 UKF vs. EKF

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

15 In practice, we typically use the UKF with some tuning parameters in

1 the Unscented Transform as discussed in Section 3.7.2. In practice, the
 2 EKF also has tuning parameters where we may wish to perform multiple
 3 updates of the dynamics equations with a smaller time-discretization
 4 before the next observation comes in to alleviate the effect of linearizing
 5 the dynamics. A well-tuned EKF is often only marginally worse than an
 6 UKF: the former requires us to compute Jacobians at each step which the
 7 latter does not, but 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.

8 3.8 Particle Filters (PFs)

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

13 3.8.1 Importance sampling

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

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

17 Say all weights are equal $1/n$. Depending upon how we pick the samples
 18 $x^{(i)}$, we can get very different approximations

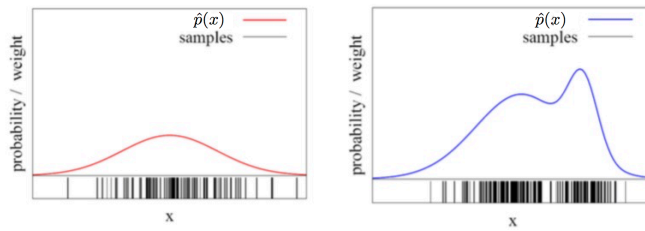
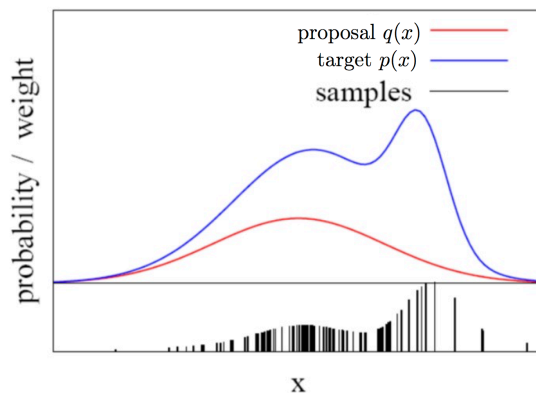


Figure 3.4: Black lines denote particles $x^{(i)}$, while red and blue curves denote the approximations obtained using them. If there are a large number of particles in a given region, the approximated probability density of that region is higher.

We see in Figure 3.4 that depending upon the samples, the approximated 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

$$\begin{aligned} \text{For } i = 1, \dots, n, \\ x^{(i)} &\sim q \\ w^{(i)} &= \frac{p(x^{(i)})}{q(x^{(i)})}. \end{aligned}$$

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 $p(x)$.



3.8.2 Resampling particles to make the weights equal

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

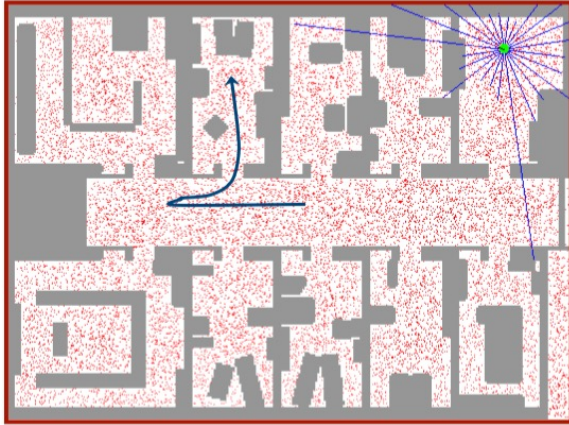


Figure 3.5: An example run of a particle filter. The robot is shown by the green dot in the top right. Observations from a laser sensor (blue rays) attached to the robot measure its distance in a 360-degree field of view around it. Red dots are particles, i.e., possible locations of the robot that we 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 ap-

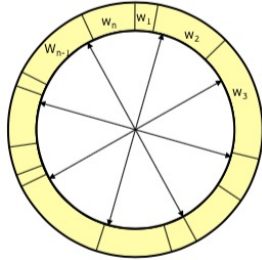
proximate a probability density $p(x)$

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

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

- 11 1. Let r be a uniform random variable in interval $[0, 1/n]$. Pick
 12 $c = w^{(1)}$ and initialize $i = 1$.
- 13 2. For each $m = 1, \dots, n$, let $u = r + (m - 1)/n$. Increment
 14 $i \leftarrow i + 1$ and $c \leftarrow c + w^{(i)}$ while $u > c$ and set new particle
 15 location $x'^{(m)} = x^{(i)}$.

16 **It is important to notice that** the resampling procedure does not actually
 17 change the locations of particles. Particles with weights much lower than
 18 $1/n$ will be eliminated while particles with weights much higher than $1/n$
 19 will be “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. Fancyer resampling methods also change the locations of the particles. The goal remains the same, namely to eliminate particles with low weights.

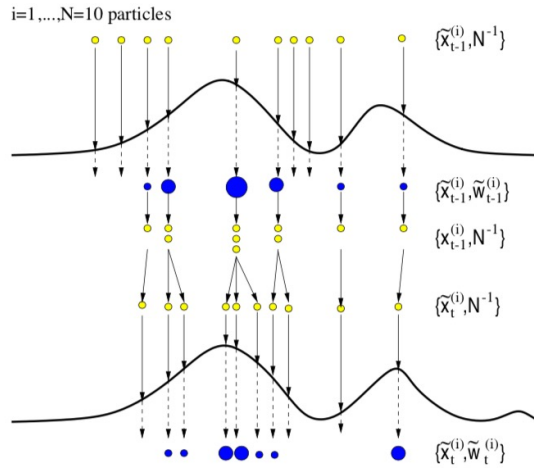


Figure 3.6: A cartoon depicting resampling. Disregard the different notation in this cartoon. Resampling does not change the probability distribution that we wish to approximate; it simply changes the particles and their weights.

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

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

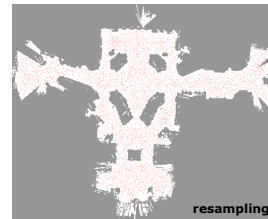
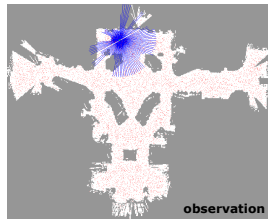
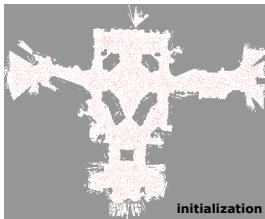
Note that $P(y_{k+1} | x_{k+1|k}^{(i)})$ is a Gaussian and depends upon the Gaussian observation noise ν_k . The mean of this Gaussian is $g(x_{k+1|k}^{(i)})$ and its variance is equal to Q , i.e.,

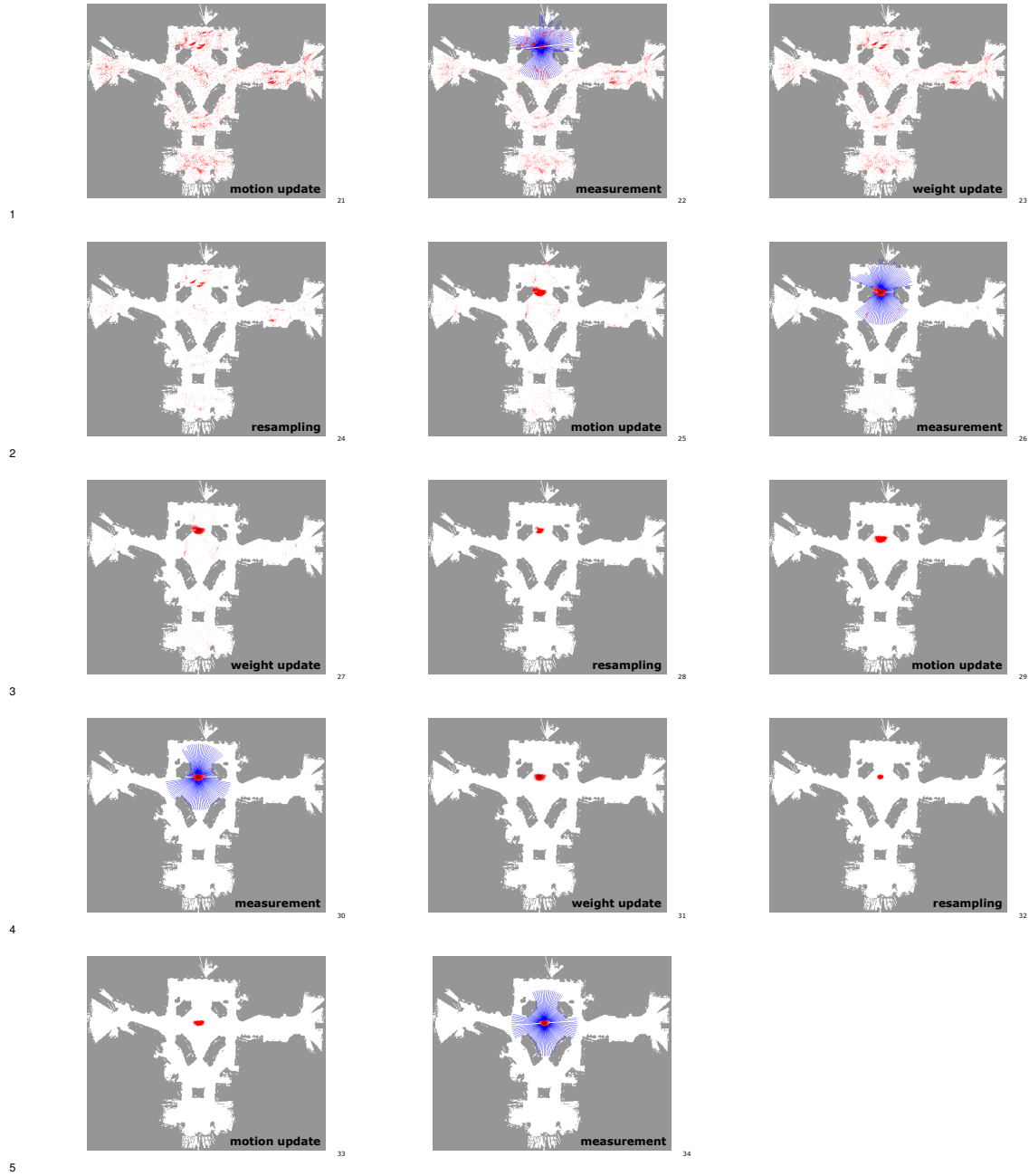
$$\begin{aligned} P(y_{k+1} | x_{k+1|k}^{(i)}) &= P(\nu_{k+1} \equiv y_{k+1} - g(x_{k+1|k}^{(i)})) \\ &= \frac{1}{\sqrt{(2\pi)^p \det(Q)}} \exp\left(-\frac{\nu_{k+1}^\top Q^{-1} \nu_{k+1}}{2}\right). \end{aligned}$$

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





6 3.8.5 Theoretical insight into particle filtering

7 **Step 1: Propagating the dynamics** As we introduced in the section on
 8 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

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

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
2 Suppose we had performed resampling on our particle set from the
3 distribution $x_{k|k}$ and have a set of n particles $\{x_{k|k}^{(i)}\}$ with equal weights
4 $1/n$

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

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

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

10 Let us focus on the numerator. We have

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

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

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

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

17 Say we keep the particle locations unchanged, i.e., $x_{k+1|k}^{(i)} = x_{k|k}^{(i)}$.

1 We then have

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

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

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

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

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

7 with weight of the i^{th} particle given by

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

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

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

12 **Step 2: Incorporating the observation** The target and proposal distri-
13 butions in this case are

$$\begin{aligned} \text{target} &: \mathbf{P}(x_{k+1} \mid y_1, \dots, y_k, y_{k+1}) \\ \text{proposal} &: \mathbf{P}(x_{k+1} \mid y_1, \dots, y_k). \end{aligned}$$

14 Since we have particles $x_{k+1|k}^{(i)}$ with weights $w_{k+1|k}^{(i)}$ for the proposal
15 distribution obtained from the propagation step, we now like to update
16 them to incorporate the latest observation y_{k+1} . Let us imagine for a
17 moment that the weights $w_{k+1|k}^{(i)}$ are uniform. We would then set weights

$$\begin{aligned} w(x) &= \frac{\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k, y_{k+1})}{\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k)} \\ &\propto \frac{\mathbf{P}(y_{k+1} \mid x_{k+1} = x) \mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k)}{\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_k)} \quad (\text{by Bayes rule}) \\ &= \mathbf{P}(y_{k+1} \mid x_{k+1} = x). \end{aligned}$$

❗ Draw a picture of how this approximation looks.

1 for each particle $x = x_{k+1|k}^{(i)}$ to get the approximated distribution as

$$\mathbf{P}(x_{k+1} = x \mid y_1, \dots, y_{k+1}) \approx \sum_{i=1}^n \mathbf{P}(y_{k+1} \mid x_{k+1|k}^{(i)}) w_{k+1|k}^{(i)} \delta_{x_{k+1|k}^{(i)}}(x) \quad (3.41)$$

2 You will notice that the right hand side is not normalized and the distribution
3 does not integrate to 1 (why? because we did not write the proportionality
4 constant in the Bayes rule above). This is easily fixed by normalizing the
5 coefficients $\mathbf{P}(y_{k+1} \mid x_{k+1|k}^{(i)}) w_{k+1|k}^{(i)}$ to sum to 1 as follows

$$w_{k+1|k+1}^{(i)} := \frac{\mathbf{P}(y_{k+1} \mid x_{k+1|k}^{(i)}) w_{k+1|k}^{(i)}}{\sum_j \mathbf{P}(y_{k+1} \mid x_{k+1|k}^{(j)}) w_{k+1|k}^{(j)}}.$$

6 **Step 3: Resampling step** As we discussed in the previous section, after
7 incorporating the observation, some particles may have very small weights.
8 The resampling procedure resamples particles so that all of them have
9 equal weights $1/n$.

$$\left\{ x_{k+1|k+1}^{(i)}, 1/n \right\}_{i=1}^n = \text{resample} \left(\left\{ x_{k+1|k+1}^{(i)}, w_{k+1|k+1}^{(i)} \right\}_{i=1}^n \right).$$

10 3.9 Discussion

11 This brings our study of filtering to a close. We have looked at some of
12 the most important algorithms for a variety of dynamical systems, both
13 linear and nonlinear. Although, we focused on filtering in this chapter,
14 all these algorithms have their corresponding “smoothing” variants, e.g.,
15 you can read about how a typical Kalman smoother is implemented at
16 https://en.wikipedia.org/wiki/Kalman_filter#Fixed-lag_smoother. Filter-
17 ing, and state estimation, is a very wide area of research even today and
18 you will find variants of these algorithms in almost every device which
19 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

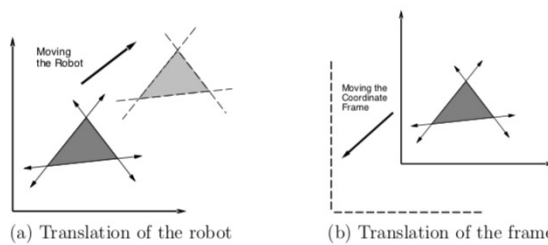
Let us imagine that the robot has a rigid body, we think of this as a subset $A \subset \mathbb{R}^2$. Say the robot is a disc

$$A = \{(x, y) \in \mathbb{R}^2 : x^2 + y^2 \leq 1\}.$$

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 if it were 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 2×2 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
5 using a 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 cor-
8 responding to $(x, y) \in \mathbb{R}^2$ and the matrix T is called a *homogeneous*
9 *transformation matrix*. The peculiar names comes from the fact that even
10 if the matrix T maps rotations and translations of rigid bodies $A \subset \mathbb{R}^2$, it
11 is just a linear transformation of the point $(x, y, 1)$ if viewed in the larger
12 space \mathbb{R}^3 .

13 **Rigid-body transformations** The transformations $R \in \mathbb{R}^{2 \times 2}$ or $T \in$
14 $\mathbb{R}^{3 \times 3}$ are called rigid-body transformations. Mathematically, it means
15 that they do not cause the distance between any two points inside the set A
16 to change. Rigid-body transformations are what are called an orthogonal
17 group in mathematics.

18 **A group** is a mathematical object which imposes certain conditions
19 upon how two operations, e.g., rotations, can be composed together. For
20 instance, if G is the group of rotations, then (i) the composition of two
21 rotations is a rotation, we say that it satisfies closure $R(\theta_1)R(\theta_2) \in G$,
22 (ii) rotations are associative

$$R(\theta_1) \{R(\theta_2)R(\theta_3)\} = \{R(\theta_1)R(\theta_2)\} R(\theta_3),$$

23 and, (iii) there exists an identity and inverse rotation

$$R(0), R(-\theta) \in G.$$

❗ It is important to remember that T represents rotation *followed* by a translation, not the other way around.

1 **An orthogonal group** is a group whose operations preserve distances
 2 in Euclidean space, i.e., $g \in G$ is an element of the group that acts on two
 3 points $x, y \in \mathbb{R}^d$ then

$$\|g(x) - g(y)\| = \|x - y\|.$$

4 If we identify the basis in Euclidean space to be the set of orthonormal
 5 vectors $\{e_1, \dots, e_d\}$, then equivalently, the orthogonal group $O(d)$ is the
 6 set of orthogonal matrices

$$O(d) := \{O \in \mathbb{R}^{d \times d} : OO^\top = O^\top O = I\}.$$

7 This implies that the square of the determinant of any element $a \in O(d)$
 8 is 1, i.e., $\det(a) = \pm 1$.

9 **The Special Orthogonal Group** is a sub-group of the orthogonal group
 10 where the determinant of each element is +1. You can see that rotations
 11 are a special orthogonal group. We denote rotations of objects in \mathbb{R}^2 as

$$SO(2) := \{R \in \mathbb{R}^{2 \times 2} : R^\top R = RR^\top = I, \det(R) = 1\}. \quad (4.3)$$

12 Each group element $g \in SO(2)$ denotes a rotation of the XY -plane about
 13 the Z -axis. The group of 3D rotations is called the Special Orthogonal
 14 Group $SO(3)$ and is defined similarly

$$SO(3) := \{R \in \mathbb{R}^{3 \times 3} : R^\top R = RR^\top = I, \det(R) = 1\}. \quad (4.4)$$

15 **The Special Euclidean Group** $SE(2)$ is simply a composition of a 2D
 16 rotation $R \in SO(2)$ and a 2D translation $\mathbb{R}^2 \ni v \equiv (x_t, y_t)$

$$SE(2) = \left\{ \begin{bmatrix} R & v \\ 0 & 1 \end{bmatrix} : R \in SO(2), v \in \mathbb{R}^2 \right\} \subset \mathbb{R}^{3 \times 3}. \quad (4.5)$$

17 The Special Euclidean Group $SE(3)$ is defined similarly as

$$SE(3) = \left\{ \begin{bmatrix} R & v \\ 0 & 1 \end{bmatrix} : R \in SO(3), v \in \mathbb{R}^3 \right\} \subset \mathbb{R}^{4 \times 4}; \quad (4.6)$$

18 again, remember that it is rotation *followed* by a translation.

19 4.1.1 3D transformations

20 Translations and rotations in 3D are conceptually similar to the two-
 21 dimensional case; however the details appear a bit more difficult because
 22 rotations in 3D are more complicated.

🔗 Check that any rotation matrix R belongs to an orthogonal group.

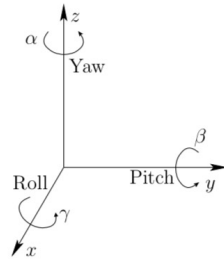


Figure 4.1: Any three-dimensional rotation can be described as a sequence of rotations about each of the cardinal axes. We usually give these specific names: rotation about the Z -axis is called yaw, rotation about the X -axis is called roll and rotation about Y -axis is called pitch. You should commit 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
 2 of the axes is written in terms of a matrix, say yaw of α -radians about the
 3 Z -axis

$$R_z(\alpha) = \begin{bmatrix} \cos \alpha & -\sin \alpha & 0 \\ \sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

- 4 Notice that this is a 3×3 matrix that keeps the Z -coordinate unchanged
 5 and only affects the other two coordinates. Similarly we have for pitch (β
 6 about the Y -axis) and roll (γ about the X -axis)

$$R_y(\beta) = \begin{bmatrix} \cos \beta & 0 & \sin \beta \\ 0 & 1 & 0 \\ -\sin \beta & 0 & \cos \beta \end{bmatrix}, R_x(\gamma) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \gamma & -\sin \gamma \\ 0 & \sin \gamma & \cos \gamma \end{bmatrix}.$$

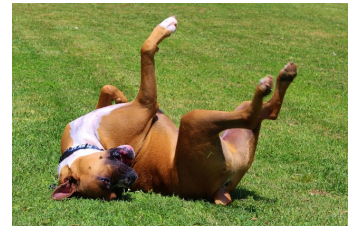
- 7 By convention, a rotation matrix in three dimensions is understood as a
 8 sequential application of rotations, *first roll, then pitch, and then yaw*

$$\mathbb{R}^{3 \times 3} = R(\gamma, \beta, \alpha) = R_z(\alpha)R_y(\beta)R_x(\gamma). \quad (4.7)$$

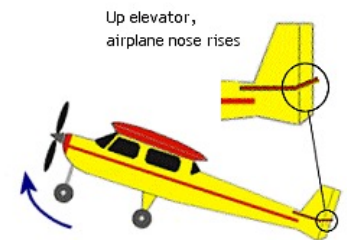
- 9 The angles (γ, β, α) (in order: roll, pitch, yaw) are called Euler angles.
 10 Imagine how the body frame of the robot changes as successive rotations
 11 are applied. If you were sitting in a car, a pure yaw would be similar to
 12 the car turning left; the Z -axis corresponding to this yaw would however
 13 only be pointing straight up perpendicular to the ground if you had not
 14 performed a roll/pitch before. If you had done so, the Z -axis of the body
 15 frame with respect to the world will be tilted.

- 16 Another important thing to note is that one single parameter determines
 17 all possible rotations about one axis, i.e., $SO(2)$. But three Euler angles
 18 are used to parameterize general rotations in three-dimensions. You
 19 can watch <https://www.youtube.com/watch?v=3Zjf95Jw2UE> to get more
 20 intuition about 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 which is not these two.

1 **Rotation matrices to Euler angles** We can back-calculate the Euler
2 angles 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 de-
5 termined using the signs of the numerator and the denominator. So
6 you should use the function atan2 in Python/C++ to implement these
7 expressions correctly. Notice that some of the expressions have r_{11} and
8 r_{33} in the denominator, this means that we need $r_{11} = \cos \alpha \cos \beta \neq 0$
9 and $r_{33} = \cos \beta \cos \gamma \neq 0$. A particular physical rotation can be parame-
10 terized in many different ways using Euler angles (depending upon the
11 order in which roll, pitch and yaw are applied), so the map from rotation
12 matrices to Euler angles is not unique.

13 **Homogeneous coordinates in three dimensions** Just like the 2D case,
14 we can define a 4×4 matrix that transforms points $(x, y, z) \in \mathbb{R}^3$ to their
15 new locations after a rotation by Euler angles (γ, β, α) and a translation
16 by a vector $v = (x_t, y_t, z_t) \in \mathbb{R}^3$

$$T = \begin{bmatrix} R(\gamma, \beta, \alpha) & v \\ 0 & 1 \end{bmatrix}.$$

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
19 by a unit vector $\omega \in \mathbb{R}^3$ with an angular velocity of 1 radian/sec. The
20 instantaneous 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 $a, b \in \mathbb{R}^3$

$$a \times b = \begin{bmatrix} a_2b_3 - a_3b_2 \\ a_3b_1 - a_1b_3 \\ a_1b_2 - a_2b_1 \end{bmatrix}$$

❗ In practice, e.g., if we run a Kalman filter to estimate the Euler angles, we need be careful in cases when α, β or $\gamma \approx \pi/2$.

Consider when β crosses $\pi/2$, i.e., it goes from $\pi/2 - \epsilon$ to $\pi/2 + \epsilon$ for some small value of ϵ . In this case, $\alpha = \tan^{-1}(r_{21}/r_{11})$, assuming $r_{21} > 0$, will jump from $\tan^{-1}(\infty) = \pi/2$ to $\tan^{-1}(-\infty) = -\pi/2$ —a jump of 180 degrees.

Another classic problem when using Euler angles occurs in what is called “Gimbal lock”. This refers to the situation when one of the angles, say $\beta = \pi/2$ (pitch up by 90 degrees). In this case, the $SO(3)$ rotation matrix is

$$R = \begin{bmatrix} 0 & 0 & 1 \\ \sin(\alpha + \gamma) & \cos(\alpha + \gamma) & 0 \\ -\cos(\alpha + \gamma) & \sin(\alpha + \gamma) & 0 \end{bmatrix}.$$

Notice here that changing α (yaw) and γ (roll) have the same effect on the rotation. We cannot distinguish the effect of yaw from roll if the pitch is 90 degrees. And such a “lock” persists until $\beta = \pi/2$. Such a gimbal lock happened on Apollo 11; the mechanism that the engineers had designed to flip the orientation by 180 degrees and escape this degeneracy did not work.

These kind of things make it very cumbersome to work with Euler angles in computer code. They are best used for visualization.

1 which we can equivalently denote as a matrix vector multiplication
 2 $a \times b = \hat{a}b$ where

$$\hat{a} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix}. \quad (4.10)$$

3 is a skew-symmetric matrix. The solution of the differential equation (4.9)
 4 at time $t = \theta$ is

$$r(\theta) = \exp(\hat{\omega}\theta) r(0)$$

5 where the matrix exponential of a matrix A is defined as

$$\exp(A) = I + A + \frac{A^2}{2!} + \frac{A^3}{3!} + \dots = \sum_{k=0}^{\infty} \frac{A^k}{k!}.$$

6 This is an interesting observation: a rotation about a fixed axis ω by an
 7 angle θ can be represented by the matrix

$$R = \exp(\hat{\omega}\theta). \quad (4.11)$$

8 You can check that this matrix is indeed a rotation by showing that
 9 $R^T R = I$ and that $\det(R) = +1$. We can expand the matrix exponential
 10 and collect odd and even powers of $\hat{\omega}$ to get

$$R = I + \sin\theta \hat{\omega} + (1 - \cos\theta)\hat{\omega}^2. \quad (4.12)$$

11 which is the Rodrigues' formula that relates the angle θ and the axis ω to
 12 the rotation matrix. We can also go in the opposite direction, i.e., given a
 13 matrix R calculate what angle θ and axis ω it corresponds to using

$$\begin{aligned} \cos\theta &= \frac{\text{tr}(R) - 1}{2} \\ \hat{\omega} &= \frac{R - R^T}{2\sin\theta}. \end{aligned} \quad (4.13)$$

14 Note that both the above formulae make sense only for $\theta \neq 0$.

15 4.2 Quaternions

16 We know two ways to think about rotations: we can either think in terms
 17 of the three Euler angles (γ, β, α) , or we can consider a rotation matrix
 18 $R \in \mathbb{R}^{3 \times 3}$. We also know ways to go back and forth between these two
 19 forms with the caveat that solving for Euler angles using (4.8) may be
 20 degenerate in some cases. While rotation matrices are the most general
 21 representation of rotations, using them in computer code is cumbersome
 22 (it is, after all, a matrix of 9 elements). So while we can build an EKF
 23 where the state is a rotation matrix, it would be a bit more expensive to
 24 run. We can also implement the same filter using Euler angles but doing
 25 so will require special care due to the degeneracies.

i Groups such as $SO(2)$ and $SO(3)$ are topological spaces (viewed as a subset of \mathbb{R}^{n^2}) and operations such as multiplication and inverses are continuous functions on these groups. These groups are also smooth manifolds (a manifold is a generalization of a curved surface) and that is why they are called *Lie groups* (after Sophus Lie). Associated to each Lie group is a Lie algebra which is the tangent space of the manifold at identity. The Lie algebra of $SO(3)$ is denoted by $so(3)$ and likewise we have $so(2)$. In a sense, the Lie algebra achieves a “linearization” of the Lie group. And the exponential map unlinearizes, i.e., it takes objects in the Lie algebra to objects in the Lie group

$$\exp : so(3) \mapsto SO(3).$$

What we have written in (4.11) is really just this map:

$$\begin{aligned} so(n) \ni \omega &\equiv \hat{\omega}\theta \\ SO(n) \ni R &= \exp(\omega) = \exp(\hat{\omega}\theta). \end{aligned}$$

Therefore, if an object whose frame had a rotation matrix R with respect to the origin were rotating with an angular velocity ω (remember that angular velocity is a vector whose magnitude is the rate of rotation and direction is axis about which the object is rotation), then the rate of change of R would be given by

$$\dot{R} = \hat{\omega}R.$$

If we were to implement a Kalman filter whose state is the rotation matrix R , then this would be the dynamics equation and one would typically have an observation for the velocity ω using a gyroscope.

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

Quaternions are a neat way to avoid the problems with both the rotation matrix and Euler angles, they parametrize the space of rotations using 4 numbers. The central idea behind quaternions is Euler's theorem which says that any 3D rotation can be considered as a pure rotation by an angle $\theta \in \mathbb{R}$ about an axis given by the unit vector ω . This is the result that we also exploited in Rodrigues' formula.

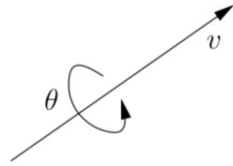


Figure 4.2: Any rotation in 3D can be represented using a unit vector ω and an angle $\theta \in \mathbb{R}$. Notice that there are two ways to encode the same rotation, the unit vector $-\omega$ and angle $2\pi - \theta$ would give the same rotation. Mathematicians say this as quaternions being a double-cover of $SO(3)$.

i As you see in the adjoining figure, quaternions also have degeneracies but they are rather easy ones.

1 A quaternion q as a four-dimensional vector $q \equiv (u_0, u_1, u_2, u_3)$ and
2 we write it as

$$q \equiv (u_0, u), \text{ or} \quad (4.14)$$

$$q = u_0 + u_1 i + u_2 j + u_3 k,$$

3 with i, j, k being three “imaginary” components of the quaternion with
4 “complex-numbers like” relationships

$$i^2 = j^2 = k^2 = ijk = -1. \quad (4.15)$$

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
7 components i, j, k have no relationship with the square roots of negative
8 unity used to define standard complex numbers. You should simply think
9 of the quaternion as a four-dimensional vector. A unit quaternion, i.e., one
10 with

$$u_0^2 + u_1^2 + u_2^2 + u_3^2 = 1,$$

11 is special: unit quaternions can be used to represent rotations in 3D.

12 **Quaternion to axis-angle representation** The quaternion $q = (u_0, u)$
13 corresponds to a counterclockwise rotation of angle θ about a unit
14 vector ω where θ and ω are such that

$$u_0 = \cos(\theta/2), \text{ and } u = \omega \sin(\theta/2). \quad (4.16)$$

1 So given an axis-angle representation of rotation like in Rodrigues' formula
2 (θ, ω) we can write the quaternion as

$$q = (\cos(\theta/2), \omega \sin(\theta/2)).$$

3 Using this, we can also compute the inverse of a quaternion (rotation of
4 angle θ about the opposite axis $-\omega$) as

$$q^{-1} := (\cos(\theta/2), -\omega \sin(\theta/2)).$$

5 The inverse quaternion is therefore the quaternion where all entries except
6 the first have their signs flipped.

7 **Multiplication of quaternions** Just like two rotation matrices multiply
8 together to give a new rotation, quaternions are also a representation
9 for the group of rotations and we can also multiply two quaternions
10 $q_1 = (u_0, u), q_2 = (v_0, v)$ together using the quaternion identities for
11 i, j, k in (4.15) to get a new quaternion

$$q_1 q_2 \equiv (u_0, u) \cdot (v_0, v) = (u_0 v_0 - u^\top v, u_0 v + v_0 u + u \times v).$$

12

13 **Pure quaternions** A pure quaternion is a quaternion with a zero scalar
14 value $u_0 = 0$. This is very useful to simply store a standard 3D vector
15 $u \in \mathbb{R}^3$ as a quaternion $(0, u)$. We can then rotate points easily between
16 different frames as follows. Given a vector $x \in \mathbb{R}^3$ we can form a
17 quaternion $(0, x)$. It turns out that

$$q \cdot (0, x) \cdot q^* = (0, R(q)x). \quad (4.17)$$

18 where $q^* = (u_0, -u)$ is the conjugate quaternion of $q = (u_0, u)$; the
19 conjugate is the same as the inverse for unit quaternions. Notice how
20 the right-hand side is the vector $R(q)x$ corresponding to the vector x
21 rotation by a matrix $R(q)$. This is a very useful trick to transform points
22 across coordinate frames instead of multiplying each point $x \in \mathbb{R}^3$ by the
23 corresponding $SE(3)$ matrix element.

24 **Quaternions to rotation matrix** The rotation matrix corresponding to
25 a quaternion is

$$\begin{aligned} R(q) &= (u_0^2 - u^\top u)I_{3 \times 3} + 2\frac{u_0 u}{\|u\|} + 2uu^\top \\ &= \begin{bmatrix} 2(u_0^2 + u_1^2) - 1 & 2(u_1 u_2 - u_0 u_3) & 2(u_1 u_3 + u_0 u_2) \\ 2(u_1 u_2 + u_0 u_3) & 2(u_0^2 + u_2^2) - 1 & 2(u_2 u_3 - u_0 u_1) \\ 2(u_1 u_3 - u_0 u_2) & 2(u_2 u_3 + u_0 u_1) & 2(u_0^2 + u_3^2) - 1 \end{bmatrix}. \end{aligned} \quad (4.18)$$

26 Using this you can show the identity that rotation matrix corresponding to
27 the product of two quaternions is the product of the individual rotation

i Quaternions belong to a larger group than rotations called the Symplectic Group $Sp(1)$.

1 matrices

$$R(q_1 q_2) = R(q_1)R(q_2).$$

2 **Rotation matrix to quaternion** We can also go in the reverse direction.

3 Given a rotation matrix R , the quaternion is

$$\begin{aligned}
 u_0 &= \frac{1}{2} \sqrt{r_{11} + r_{22} + r_{33} + 1} \\
 \text{if } u_0 \neq 0, u_1 &= \frac{r_{32} - r_{23}}{4u_0} \\
 u_2 &= \frac{r_{13} - r_{31}}{4u_0} \\
 u_3 &= \frac{r_{21} - r_{12}}{4u_0} \tag{4.19} \\
 \text{if } u_0 = 0, u_1 &= \frac{r_{13}r_{12}}{\sqrt{r_{12}^2 r_{13}^2 + r_{12}^2 r_{23}^2 + r_{13}^2 r_{23}^2}} \\
 u_2 &= \frac{r_{12}r_{23}}{\sqrt{r_{12}^2 r_{13}^2 + r_{12}^2 r_{23}^2 + r_{13}^2 r_{23}^2}} \\
 u_3 &= \frac{r_{13}r_{23}}{\sqrt{r_{12}^2 r_{13}^2 + r_{12}^2 r_{23}^2 + r_{13}^2 r_{23}^2}}.
 \end{aligned}$$

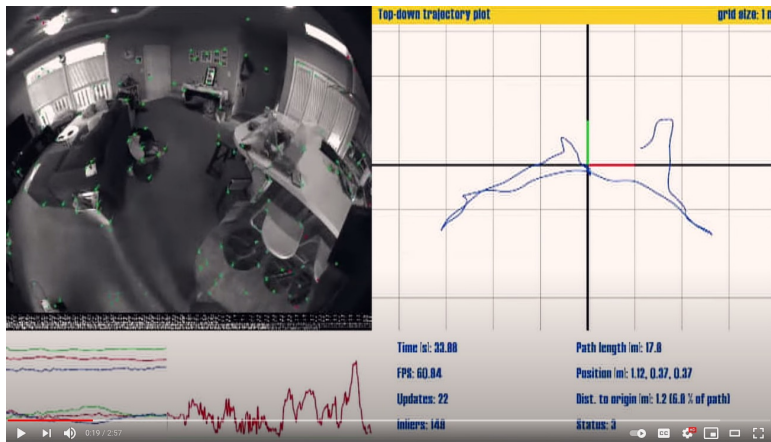
4

5 4.3 Occupancy Grids

6 Rotation matrices and quaternions let us capture the dynamics of a rigid
7 robot body. We will next look at how to better understand observations.

8 **What is location and what is mapping?** Imagine a robot that is
9 moving around in a house. A natural representation of the state of this
10 robot is the 3D location of all the interesting objects in the room, e.g.,
11 <https://www.youtube.com/watch?v=Qe10ExwzCqk>. At each time-instant,
12 we record an observation from our sensor (in this case, a camera) that
13 indicates how far an object is from the robot. This helps us discover the
14 location of the objects in the room. After gathering enough observations,
15 we would have created a *map* of the entire house. This map is the set of
16 positions of all interesting objects in the room. Such a map is called a
17 “feature map”, these are all the green points in the image below

❗ There is little need to memorize these expressions or trying to understand patterns between them. While building a new code base for your robot, you will usually code up these formulae once and all your code will use them again and again.



2 The main point to understand about feature map is that we can hand
 3 over this map to another robot that comes to the same house. The robot
 4 compares images from its camera and if it finds one of the objects inside
 5 the map, it can get an estimate of its location/orientation in the room with
 6 respect to the known location of the object in the map. The map is just
 7 a set of “features” that help identify salient objects in the room (objects
 8 which can be easily detected in images and relatively uniquely determine
 9 the location inside the room). The second robot using this map to estimate
 10 its position/orientation in the room is called the *localization problem*. We
 11 already know how to solve the localization problem using filtering.

12 The first robot was solving a harder problem called *Simultaneous*
 13 *Localization And Mapping (SLAM)*: namely that of discovering the location
 14 of both itself and the objects in the house. This is a very important and
 15 challenging problem in robots but we will not discuss it further. MEAM
 16 620 digs deeper 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.

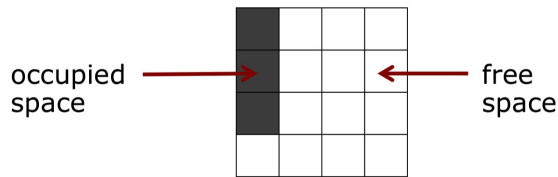
17 **Grid maps** We will first discuss two-dimensional grid maps, they look
 18 as follows.



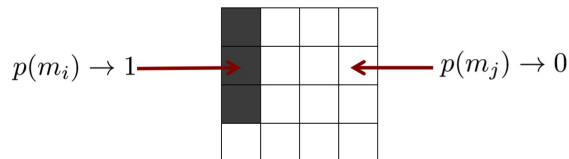
Figure 4.3: A grid map (also called an occupancy grid) is a large gray-scale image, each pixel represents a cell in the physical world. In this picture, cells that are occupied are colored black and empty cells represent free space. A grid map is a useful representation for a robot to localize in this house using observations from its sensors and comparing those to the map.

1 To get a quick idea of what we want to do, you can watch the mapping
 2 being performed in <https://www.youtube.com/watch?v=JJhEKIA1xSE>.
 3 We are interested in learning such maps from the observations that a
 4 robot collects as it moves around the physical space. Let us make two
 5 simplifying assumptions.

6 **Assumption 1: each cell is either free or occupied**



7
 8 This is neat: we can now model each cell as a binary random variable that
 9 indicates occupancy. Let the probability that the cell m_i be occupied be
 10 $p(m_i)$



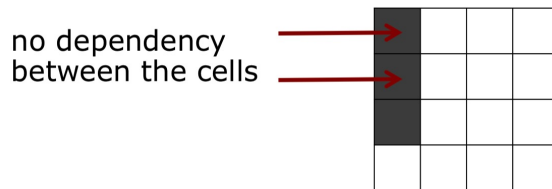
11

12 If we have $p(m_i) = 0$, then the cell is not occupied and if we have
 13 $p(m_i) = 1$, then the cell is occupied. A priori, we do not know the state
 14 of the cell so we will set the prior probability to be $p(m_i) = 0.5$.

15 **Assumption 2: the world is static** Objects in the world do not move.
 16 This is reasonable if we are interested in estimating in building a map of
 17 the walls inside the room. Note that it is not a reasonable assumption if

1 there are moving people inside the room. We will see a clever hack where
 2 the Bayes rule helps automatically disregard such moving objects in this
 3 section.

4 **Assumption 3: cells are independent of each other** This is another
 5 drastic simplification. The state of our system is the occupancy of each cell
 6 in the grid map. We assume that before receiving any observations, the
 7 occupancy of each individual cell is independent; it is a Bernoulli variable
 8 with probability 1/2 since we have assumed the prior to be uniform in
 9 Assumption 1.



10

11 This means that if cells in the map are denoted by a vector $m = (m_1, \dots)$,
 12 then the probability of the cells being occupied/not-occupied can be written
 13 as

$$p(m) = \prod_i p(m_i). \quad (4.20)$$



example map
(4-dim vector)



4 individual cells

14

15 4.3.1 Estimating the map from the data

16 Say that the robot pose (position and orientation) is given by the sequence
 17 x_1, \dots, x_k . While proceeding along this sequence, the robot receives
 18 observations y_1, \dots, y_k . Our goal is to estimate the state of each cell
 19 $m_i \in \{0, 1\}$ (aka “the map” $m = (m_1, m_2, \dots)$)

$$P(m \mid x_1, \dots, x_k, y_1, \dots, y_k) = \prod_i P(m_i \mid x_1, \dots, x_k, y_1, \dots, y_k). \quad (4.21)$$

20 This is called the “static state” Bayes filter and is conceptually exactly the
 21 same as the recursive application of Bayes rule in Chapter 2 for detecting
 22 whether the door was open or closed.

23 We will use a short form to keep the notation clear

$$y_{1:k} = (y_1, y_2, \dots, y_k);$$

24 the quantity $x_{1:k}$ is defined similarly. As usual we will use a recursive

1 Bayes filter to compute this probability as follows.

$$\begin{aligned}
 \mathbb{P}(m_i \mid x_{1:k}, y_{1:k}) &\stackrel{\text{Bayes rule}}{=} \frac{\mathbb{P}(y_k \mid m_i, y_{1:k-1}, x_{1:k}) \mathbb{P}(m_i \mid y_{1:k-1}, x_{1:k})}{\mathbb{P}(y_k \mid y_{1:k-1}, x_{1:k})} \\
 &\stackrel{\text{Markov}}{=} \frac{\mathbb{P}(y_k \mid m_i, x_k) \mathbb{P}(m_i \mid y_{1:k-1}, x_{1:k-1})}{\mathbb{P}(y_k \mid y_{1:k-1}, x_{1:k})} \\
 &\stackrel{\text{Bayes rule}}{=} \frac{\mathbb{P}(m_i \mid y_k, x_k) \mathbb{P}(y_k \mid x_k) \mathbb{P}(m_i \mid y_{1:k-1}, x_{1:k-1})}{\mathbb{P}(m_i \mid x_k) \mathbb{P}(y_k \mid y_{1:k-1}, x_{1:k})} \\
 &\stackrel{\text{Independence}}{=} \frac{\mathbb{P}(m_i \mid y_k, x_k) \mathbb{P}(y_k \mid x_k) \mathbb{P}(m_i \mid y_{1:k-1}, x_{1:k-1})}{\mathbb{P}(m_i) \mathbb{P}(y_k \mid y_{1:k-1}, x_{1:k})}.
 \end{aligned} \tag{4.22}$$

2 We have a similar expression for the opposite probability

$$\mathbb{P}(\neg m_i \mid x_{1:k}, y_{1:k}) = \frac{\mathbb{P}(\neg m_i \mid y_k, x_k) \mathbb{P}(y_k \mid x_k) \mathbb{P}(\neg m_i \mid y_{1:k-1}, x_{1:k-1})}{\mathbb{P}(\neg m_i) \mathbb{P}(y_k \mid y_{1:k-1}, x_{1:k})}.$$

3 Let us take the ratio of the two to get

$$\begin{aligned}
 \frac{\mathbb{P}(m_i \mid x_{1:k}, y_{1:k})}{\mathbb{P}(\neg m_i \mid x_{1:k}, y_{1:k})} &= \frac{\mathbb{P}(m_i \mid y_k, x_k)}{\mathbb{P}(\neg m_i \mid y_k, x_k)} \frac{\mathbb{P}(m_i \mid y_{1:k-1}, x_{1:k-1})}{\mathbb{P}(\neg m_i \mid y_{1:k-1}, x_{1:k-1})} \frac{\mathbb{P}(\neg m_i)}{\mathbb{P}(m_i)} \\
 &= \underbrace{\frac{\mathbb{P}(m_i \mid y_k, x_k)}{1 - \mathbb{P}(m_i \mid y_k, x_k)}}_{\text{uses observation } y_k} \underbrace{\frac{\mathbb{P}(m_i \mid y_{1:k-1}, x_{1:k-1})}{1 - \mathbb{P}(m_i \mid y_{1:k-1}, x_{1:k-1})}}_{\text{recursive term}} \underbrace{\frac{1 - \mathbb{P}(m_i)}{\mathbb{P}(m_i)}}_{\text{prior}}.
 \end{aligned} \tag{4.23}$$

4 This is called the odds ratio. Notice that the first term uses the latest
 5 observation y_k , the second term can be updated recursively because it
 6 is a similar expression as the left-hand side and the third term is a prior
 7 probability of the cell being occupied/non-occupied. Let us rewrite this
 8 formula using the log-odds-ratio that makes implementing it particularly
 9 easy. The log-odds-ratio of the probability $p(x)$ of a binary variable x is
 10 defined as

$$l(x) = \log \frac{p(x)}{1 - p(x)}, \text{ and } p(x) = 1 - \frac{1}{1 + e^{l(x)}}.$$

11 The product in (4.23) now turns into a sum as

$$l(m_i \mid y_{1:k}, x_{1:k}) = l(m_i \mid y_k, x_k) + l(m_i \mid y_{1:k-1}, x_{1:k-1}) - l(m_i). \tag{4.24}$$

12 This expression is used to update the occupancy of each cell. The term

$$\text{sensor model} = l(m_i \mid y_k, x_k)$$

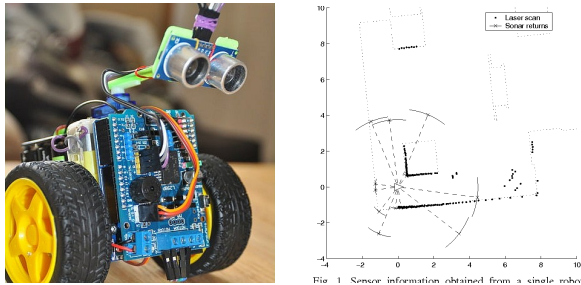
13 is different for different sensors and we will investigate it next.

14 4.3.2 Sensor models

15 **Sonar** This works by sending out an ultrasonic chirp and measuring the
 16 time between emission and reception of the signal. The time gives an

❓ We assumed that the map was static. Can you think of why (4.24) automatically lets us handle some moving objects? Think of what the prior odds $l(m_i)$ does to the log-odds-ratio $l(m_i \mid y_{1:k}, x_{1:k})$.

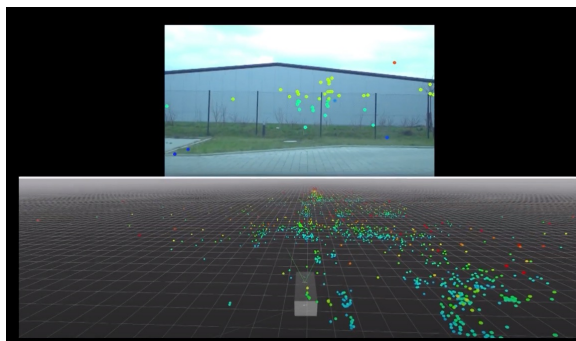
1 estimate of the distance of an object to the robot.



2 Fig. 1 Sense information obtained from a simple robot

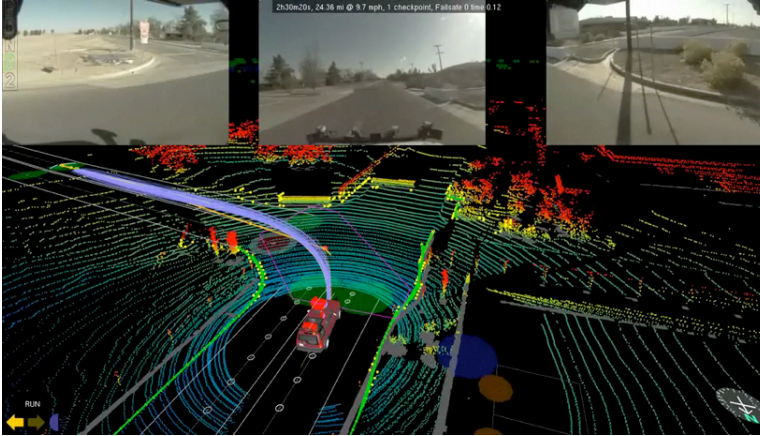
3 The figure above shows a typical sonar sensor (the two “eyes”) on a
 4 low-cost robot. Data from the sensor is shown on the right, a sonar is a
 5 very low resolution sensor and has a wide field of view, say 15 degrees,
 6 i.e., it cannot differentiate between objects that are within 15 degrees
 7 of each other and registers them as the same point. Sophisticated sonar
 8 technology is used today in marine environments (submarines, fish finders,
 9 detecting mines etc.).

10 **Radar** works in much the same way as a sonar except that it uses
 11 pulses of radio waves and measures the phase difference between the
 12 transmitted and the received signal. This is a very versatile sensor
 13 (it was invented by the US army to track planes and missiles during
 14 World War II) but is typically noisy and requires sophisticated process-
 15 ing to be used for mainstream robotics. Autonomous cars, collision
 16 warning systems on human-driven cars, weather sensing, and certainly
 17 the military use the radar today. The following picture and the video
 18 https://www.youtube.com/watch?v=hwKUcu_7F9E will give you an ap-
 19 preciation of the kind of data that a radar records. Radar is a very long
 20 range sensor (typically 150 m) and works primarily to detect metallic
 21 objects.



22
 23 **LiDAR** LiDAR, which is short for Light Detection and Ranging,
 24 (<https://en.wikipedia.org/wiki/Lidar>) is a portmanteau of light and radar.
 25 It is a sensor that uses a pulsed laser as the source of illumination and
 26 records the time it takes (nanoseconds typically) for the signal to return
 27 to the source. See <https://www.youtube.com/watch?v=NZKvf1cXe8s>
 28 for how the data from a typical LiDAR (Velodyne) looks like. While a

1 Velodyne contains an intricate system of rotating mirrors and circuitry to
 2 measure time elapsed, there are new solid state LiDARs that are rapidly
 3 evolving to match the needs of the autonomous driving industry. Most
 4 LiDARs have a usable range of about 100 m.



6 **A typical autonomous car** This is a picture of MIT's entry named Talos
 7 to the DARPA Urban Challenge ([https://en.wikipedia.org/wiki/DARPA_Grand_Challenge_\(2007\)](https://en.wikipedia.org/wiki/DARPA_Grand_Challenge_(2007)))
 8 which was a competition where teams had to traverse a 60 mile urban
 9 route within 6 hours, while obeying traffic laws, understanding incoming
 10 vehicles etc. Successful demonstrations by multiple teams led by (CMU,
 11 Stanford, Odin College, MIT, Penn and Cornell) in this competition jump-
 12 started the wave of autonomous driving. While the number of sensors
 13 necessary to drive well has come down (Tesla famously does not like
 14 to use LiDARs and rely exclusively on cameras and radars), the type of
 15 sensors and challenges associated with them remain essentially the same.



16 **Waymo's autonomous car**



17 4.3.3 Back to sensor modeling

18 Let us go back to understanding our sensor model $l(m_i | y_k, x_k)$ where m_i
 19 is a particular cell of the occupancy grid, y_k and x_k are the observations
 20 and robot position/orientation at time k .

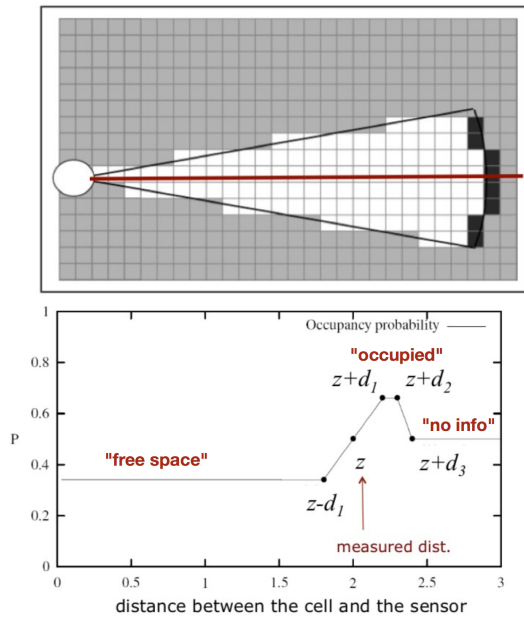


Figure 4.4: **Model for sonar data.** (Top) A sonar gives one real-valued reading corresponding to the distance measured along the red axis. (Bottom) if we travel along the optical axis, the occupancy probability $P(m_i | y_k = z, x_k)$ can be modeled as a spike around the measured value z . It is very 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 cell along the axis.

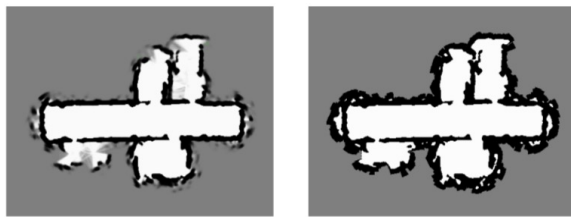
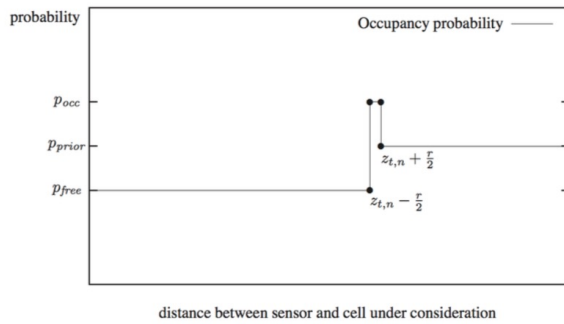


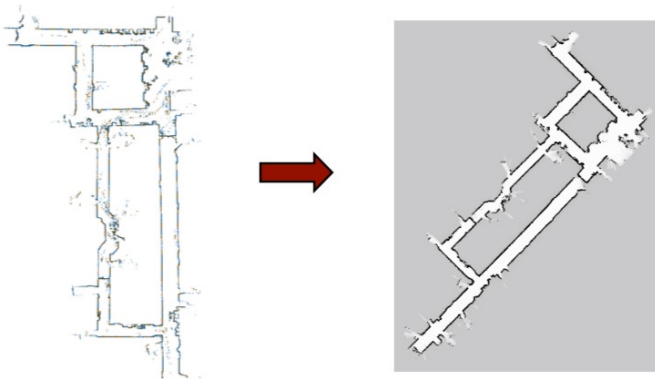
Figure 4.5: (Left) A typical occupancy grid created using a sonar sensor by updating the log-odds-ratio $l(m_i | x_{1:k}, y_{1:k})$ for all cells i for multiple time-steps k . At the end of the map building process, if $l(m_i | x_{1:k}, y_{1:k}) > 0$ for a particular cell, we set its occupancy to 1 and 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
- 2 than the sonar, what we really mean is that the sensor model $P(m_i | y_k, x_k)$
- 3 looks as follows.



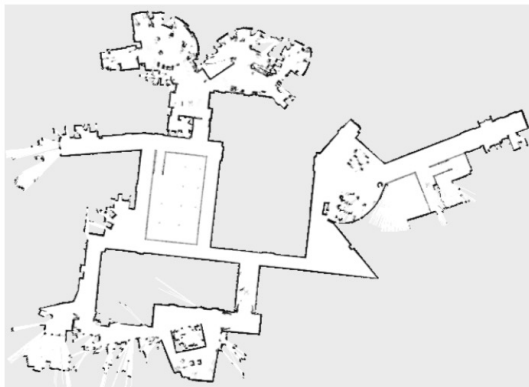
1

2 As a result, we can create high-resolution occupancy grids using a LiDAR.



3

Example: MIT CSAIL 3rd Floor



4

❓ 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?

5 4.4 3D occupancy grids

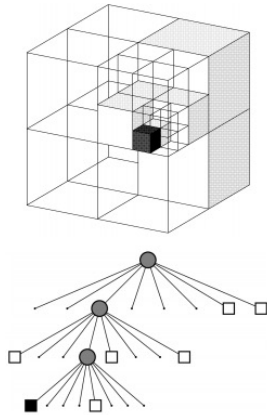
6 Two-dimensional occupancy grids are a fine representation for toy problems
 7 but they run into some obvious issues. Since the occupancy grid is a “top
 8 view” of the world, we cannot represent non-trivial objects in it correctly
 9 (a large tree with a thin trunk eats up all the free space). We often desire a
 10 fundamentally three-dimensional representation of the physical world.



1

2 We could simply create cells in 3D space and our method for occupancy
 3 grid would work but this is no longer computationally cheap. For instance,
 4 if we want to build a map of Levine Hall (say $100\text{ m} \times 100\text{ m}$ area and
 5 height of 25 m), a 3D grid map with a resolution of $5\text{ cm} \times 5\text{ cm} \times \text{cm}$
 6 would have about 2 billion cells (if we store a float in each cell this map will
 7 require about 8 GB memory). It would be cumbersome to carry around
 8 so many cells and update their probabilities after each sensor reading (a
 9 Velodyne gives data at about 30 Hz). More importantly, observe that most
 10 of the volume inside Levine is free space (inside of offices, inner courtyard
 11 etc.) so we do not really need fine resolution in those regions.

12 **Octrees** We would ideally have an occupancy grid whose resolution
 13 adapts with the kind of objects that are detected by the sensors. If nearby
 14 cells are empty we want to collapse them together to save on memory and
 15 computation, on the other hand, if nearby cells are all occupied, we want
 16 to *refine* the resolution in that area so has to more accurately discern the
 17 shape of the underlying objects. Octrees are an efficient representation for
 18 3D volumes.

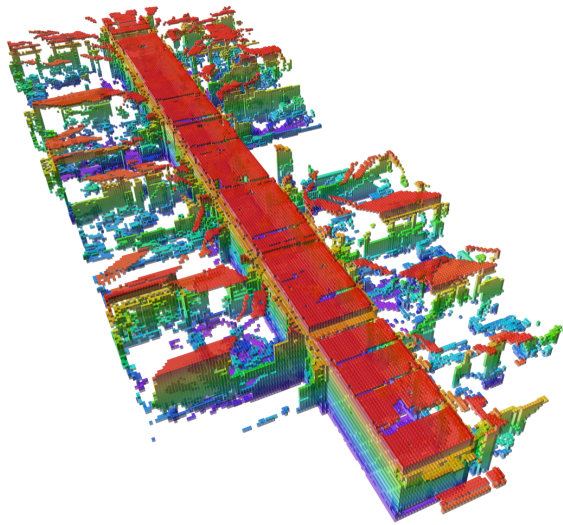


19

20 An octree is a hierarchical data structure that recursively sub-divides the
 21 3D space into octants and allocates volumes as needed for a particular data
 22 point observed by a range sensor. It is analogous to a kd-tree. Imagine if
 23 the entire space in the above picture were empty (the tree only has a root
 24 node), and we receive a reading corresponding to the dark shaded region.
 25 An octree would sub-divide the space starting from the root (each node
 26 in the tree populates is the parent of its eight child octants) recursively
 27 until some pre-determined minimum resolution is reached. This leaf

1 node is grid cell; notice how different cells in the octree have different
 2 resolutions. Occupancy probabilities of each leaf node are updated using
 3 the same formula as that of (4.24). A key point here is that octrees are
 4 designed for accurate sensors (LiDARs) where there is not much noise
 5 in the observations returned by the sensor (and thereby we do not refine
 6 unnecessary parts of the space)

7 Octrees are very efficient at storing large map, I expect you can store
 8 the entire campus of Penn in about a gigabyte. Ray tracing (following all
 9 the cells m_i in tree along the axis of the sensor in Figure 4.4) is harder
 10 in this case but there are efficient algorithms devised for this purpose.
 11 An example OctoMap (an occupancy map created using an Octree) of a
 building on the campus of 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>

13 4.5 Local Map

14 In this chapter, we primarily discussed occupancy grids of static environ-
 15 ments as the robot moves around in the environment. The purpose of doing
 16 so is localization, namely, finding the pose of the robot by comparing
 17 the observations of the sensors with the map (think of the particle filter
 18 localization example in Chapter 3). In typical problems, we often maintain
 19 two kinds of maps, (i) a large occupancy grid for localization (say as big
 20 as city), and (ii) another smaller map, called the local map, that is used to
 21 maintain the locations of objects (typically objects that can move) in the
 22 vicinity of the robot, say a $100\text{ m} \times 100\text{ m}$ area.

23 The local map is used for planning and control purposes, e.g., to check
 24 if the planned trajectory of the robot does not collide with any known
 25 obstacles. See an example of the local map at the 1:42 min mark at
 26 <https://www.youtube.com/watch?v=2va15BE-7lQ>. Some people also call
 27 the local map a “cost map” because occupied cells in the local map indicate
 28 a high collision cost of moving through that cell. The local map is typically

1 constructed in the body frame and evolves as the robot moves around
 2 (objects appear in the front of the robot and are spawned in the local map
 3 and disappear 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.



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.

4 4.6 Discussion

5 Occupancy grids are a very popular approach to represent the environment
 6 given the poses of the robot as it travels in this environment. We can also
 7 use occupancy grids to localize the robot in a future run (which is usually
 8 the purpose of creating them). Each cell in an occupancy grid stores the
 9 posterior probability of the cell being occupied on the basis of multiple
 10 observations $\{y_1, \dots, y_k\}$ from respective poses $\{x_1, \dots, x_k\}$. This is
 11 a very efficient representation of the 3D world around us with the one
 12 caveat that each cell is updated independently of the others. But since
 13 one gets a large amount of data from typical range sensors (a 64 beam
 14 Velodyne (<https://velodynelidar.com/products/hdl-64e>) returns about a
 15 2 million points/sec and cheaper versions of this sensor will cost about
 16 \$100), this caveat does not hurt us much in practice. You can watch this talk

- 1 (https://www.youtube.com/watch?v=V8JMwE_L5s0) by the head of Uber's
- 2 autonomous driving group to get more perspective about localization and
- 3 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 terminology for quantities, e.g., the value function, that you might see being written very succinctly in other places. We will mostly follow the notation of Dmitri Bertsekas' book on "Reinforcement Learning and Optimal Control" (<http://www.mit.edu/~dimitrib/RLbook.html>). You will get used to the extra notation and it will become second nature once you become more familiar with the concepts.

5.1 Formulating the optimal control problem

Let us denote the state of a robot (and the world) by $x_k \in X \subset \mathbb{R}^n$ at the k^{th} timestep. We can change this state using a control input $u_k \in U \subset \mathbb{R}^p$ and this change is written as

$$x_{k+1} = f_k(x_k, u_k) \quad (5.1)$$

for $k = 0, 1, \dots, T - 1$ starting from some initial given state x_0 . This is deterministic nonlinear dynamical system (no noise ϵ in the dynamics). We will let the dynamics f_k also be a function of time k . The time T is some time-horizon up to which we care about running the system. The state-space is X (which we will assume does not change with time k) and the control-space is U .

Recall, that we can safely assume that the system is Markov. The reason for it is as follows. If it is not, and say if x_{k+1} depends upon both x_k and the previous step x_{k-1} , then we can expand the state-space to write a new dynamics in the expanded state-space. We will follow a similar program as that of Module 1: we first describe very general algorithms (dynamic programming) for general systems (Markov Decision Processes), then specialize our methods to a restricted class of systems that are useful in practice (linear dynamical systems) and then finally discuss a very general class of systems again with 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.

1 **Run-time cost and terminal cost** We will take a very general view of
 2 the above problem and formalize it as follows. Consider a cost function

$$q_k(x_k, u_k) \in \mathbb{R}$$

3 which gives a scalar real-valued output for every pair (x_k, u_k) . This
 4 models the fact that you do not want to walk more than you need to get to
 5 School, i.e., we would like to minimize q_k . You also want to make sure
 6 the trajectory actually reaches the lecture venue, we write this down as
 7 another cost $q_f(x_T)$. We want to pick control inputs $(u_0, u_1, \dots, u_{T-1})$
 8 such that

$$J(x_0; u_0, u_1, \dots, u_{T-1}) = q_f(x_T) + \sum_{k=0}^{T-1} q_k(x_k, u_k) \quad (5.2)$$

9 is minimized. The cost $q_f(x_T)$ is called the terminal cost, it is high if
 10 x_T is not the lecture room and small otherwise. The cost q_k is called the
 11 run-time cost, it is high for instance if you have to use large control inputs,
 12 e.g., x_k is a 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 .

13 5.2 Dijkstra's algorithm

14 If the state-space X and control-space U are discrete and finite sets, we
 15 can solve (5.3) as a shortest path problem using very fast algorithms.
 16 Consider the following picture. This is what would be called a transition
 17 graph for a deterministic finite-state dynamics.

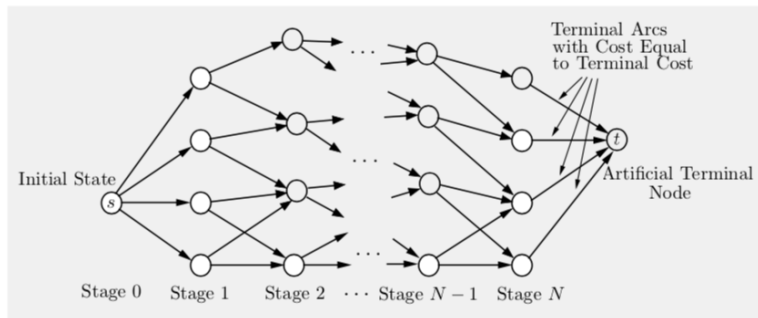


Figure 5.1: Transition graph for Dijkstra's algorithm

1 The graph has one source node x_0 . Each node in the graph is x_k , each
 2 edge depicts taking a certain control u_k . Depending on which control we
 3 pick, we move to some other node x_{k+1} given by the dynamics $f(x_k, u_k)$.
 4 Note that this is *not* a transition like that of a Markov chain, everything is
 5 deterministic in this graph. On each edge we write down the cost

$$\text{cost}(x_k, x_{k+1}) := q_k(x_k, u_k)$$

6 where $x_{k+1} = f_k(x_k, u_k)$ and “close” the graph with a dummy terminal
 7 node with the cost $q_f(x_T)$ on every edge leading to an artificial terminal
 8 node (sink).

9 Minimizing the cost in (5.3) is now the same as finding the shortest
 10 path in this graph from the source to the sink. The algorithm to do so is
 11 quite simple and is called Dijkstra's algorithm after Edsgar Dijkstra who
 12 used it around 1956 as a test program for a new computer named ARMAC
 13 (<http://www-set.win.tue.nl/UnsungHeroes/machines/armac.html>).

- 14 1. Let Q be the set of nodes that are currently unvisited; all nodes in
 15 the graph are added to it at the beginning. S is an empty set. An
 16 array called `dist` maintains the distance of every node in the graph
 17 from the source node x_0 . Initialize $\text{dist}(x_0) = 0$ and $\text{dist} = \infty$ for
 18 all other nodes.
- 19 2. At each step, if Q is not empty, pop a node $v \in Q$ such that $v \notin S$
 20 with the smallest $\text{dist}(v)$. Add v to S . Update the `dist` of all nodes
 21 u connected to v . For each u , if

$$\text{dist}(u) > \text{dist}(v) + \text{cost}(u, v)$$

22 update the distance of u to be $\text{dist}(v) + \text{cost}(u, v)$. If the above
 23 condition is not true do nothing.

24 The algorithm terminates when the set Q is empty.

25 You might know that there are many other variants of Dijkstra's
 26 algorithm, e.g., the A^* algorithm that are quicker to find shortest paths.
 27 We will look at some of these in the next chapter.

🔍 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?

The quantity dist is quite special: observe that after Dijkstra's algorithm finishes running and the set Q is empty, the dist function gives the optimal cost to go from each node in the graph to the source 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 source.

5.2.1 Dijkstra's algorithm in the backwards direction

We can run Dijkstra's algorithm in the backwards direction to get the same answer as well. The sets Q and S are initialized as before. In this case we will let $\text{dist}(v)$ denote the distance of a node v to the sink node. The algorithm proceeds in the same fashion, it pops a node $v \in Q, v \notin S$ and updates the dist of all nodes u connected to v . For each u , if

$$\text{dist}(u) > \text{dist}(v) + \text{cost}(u, v)$$

then we update $\text{dist}(u)$ to be the right-hand side of this inequality. Running Dijkstra's algorithm in reverse (from sink to the source) is completely equivalent to running it in the forward direction (from source to the sink).

❗ 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 dist . The number of edges in the transition graph in Figure 5.1 is $m = \mathcal{O}(T|X|)$.

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.

We can prove this as follows. Suppose that we find the optimal control sequence $(u_0^*, u_1^*, \dots, u_{T-1}^*)$ for the problem in (5.3). Our system is deterministic, so this control sequence results in a *unique* sequence of states $(x_0, x_1^*, \dots, x_T^*)$. Each successive state is given by $x_{k+1}^* = f_k(x_k^*, u_k^*)$ with $x_0^* = x_0$. The principle of optimality, or the principle of dynamic programming, states that if one starts from a state x_k^* at time k and wishes to minimize the "cost-to-go"

$$q_f(x_T) + q_k(x_k^*, u_k) + \sum_{i=k+1}^{T-1} q_i(x_i, u_i)$$

over the (now assumed unknown) sequence of controls $(u_k, u_{k+1}, \dots, u_{T-1})$, then the optimal control sequence for this truncated problem is exactly $(u_k^*, \dots, u_{T-1}^*)$.

The proof of the above assertion is an easy case of proof by contradiction: if the truncated sequence were not optimal starting from x_k^* there

1 exists some other optimal sequence of controls for the truncated problem,
 2 say $(v_k^*, \dots, v_{T-1}^*)$. If so, the solution of the original problem where one
 3 takes controls v_k^* from this new sequence for time-steps $k, k+1, \dots, T-1$
 4 would have a lower cost. Hence the original sequence of controls would
 5 not have been optimal.

Principle of dynamic programming. The essence of dynamic programming 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 $x \in X$.

6 After running the above algorithm we have the optimal cost-to-go $J_0^*(x)$
 7 for each state $x \in X$, in particular, we have the cost-to-go for the initial
 8 state $J_0^*(x_0)$. If we remember the minimizer u_k^* in (5.4) while running the
 9 algorithm, we also have the optimal sequence $(u_0^*, u_1^*, \dots, u_{T-1}^*)$. The
 10 function $J_0^*(x)$ (often shortened to simply $J^*(x)$) is the optimal cost-to-go
 11 from the state $x \in X$.

12 Again, we really only wanted to calculate $J_0^*(x_0)$ but had to do all this
 13 extra work of computing J_k^* for all the states.

14 **Curse of dimensionality** What is the complexity of running dynamic
 15 programming? The cost of the minimization over U is $\mathcal{O}(|U|)$, it is a
 16 bunch of comparisons between floats. The number of operations at each
 17 iteration for setting the values $J_k^*(x)$ for all $x \in X$ is $|X|$. So the total
 18 complexity is $\mathcal{O}(T |X| |U|)$.

19 The terms $|X|$ and $|U|$ are often the hurdle in implementing dynamic
 20 programming or any variant of it. Think of the grid-world in Problem 1 in
 21 HW 1, it had 200×200 cells which amounts to $|X| = 40,000$. This may
 22 seem a reasonable number but it explodes quickly as the dimensionality
 23 of the state-space grows. For a robot manipulator with six degrees-of-
 24 freedom, if we discretize each joint angle into 5 degree cells, the number
 25 of states is $|X| \approx 140$ billion. The number of states $|X|$ is exponential in

1 the dimensionality of the state-space and dynamic programming quickly
 2 becomes prohibitive beyond 4 dimensions or so. Bellman called this the
 3 *curse of dimensionality*.

4 **Cost of dynamic programming is linear in the time-horizon** Notice a
 5 very important difference between (5.4)

$$J_k^*(x) = \min_{u_k \in U} \{q_k(x, u_k) + J_{k+1}^*(f_k(x, u_k))\}$$

6 for iterations $i = T - 1, \dots, 0$ and (5.3)

$$J^*(x_0) = \min_{u_k \in U, k=0, \dots, T-1} J(x_0; u_0, \dots, u_{T-1}).$$

7 The latter has a minimization over a sequence of controls $(u_0, u_1, \dots, u_{T-1})$
 8 while the former has a minimization over only the control at time k , u_k
 9 over T iterations. The former is much much easier to solve because it is
 10 a sequence of $\mathcal{O}(T)$ smaller optimization problems: it is really easy to
 11 compute $\min_{u_k \in U}$ for each state x separately than to solve the gigantic
 12 minimization problem in (5.3) because in the latter case, the variable of
 13 optimization is the entire control trajectory and has size $|U|^T$.

14 **Dynamic programming and Viterbi's algorithm** We have seen the
 15 principle of dynamic programming in action before in Viterbi's algorithm
 16 in Chapter 2. The transition graph in Figure 5.1 is the same as the Trellis
 17 graph for Viterbi's algorithm, the run-time cost was

$$q_k(x_k, u_k) := -\log P(Y_k | X_k) - \log P(X_{k+1} | X_k)$$

18 and instead of a terminal cost q_f , we had an initial cost $-\log P(X_1)$.
 19 Viterbi's algorithm computed the most likely path given observations
 20 of the HMM, i.e., the path (X_1, \dots, X_T) that maximizes the proba-
 21 bility $P(X_1, \dots, X_T | Y_1, \dots, Y_T)$ is simply the solution of dynamic
 22 programming for the Trellis graph.

23 5.3.1 Q-factor

24 The quantity

$$Q_k^*(x, u) := q_k(x, u) + J_{k+1}^*(f_k(x, u))$$

25 is called the Q-factor. It is simply the expression that is minimized in the
 26 right-hand side of (5.4) and denotes the cost-to-go if control u was picked
 27 at state x (corresponding to cost $q_k(x, u)$) and the the optimal control
 28 trajectory was followed after that (corresponding to cost $J_{k+1}^*(f_k(x, u))$)
 29 from state $x' = f_k(x, u)$. This nomenclature was introduced by Watkins
 30 in his thesis.

🔗 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 Q-factors and the cost-to-go are equivalent ways of thinking about
 2 dynamic programming. Given the Q-factor, we can obtain the cost-to-go
 3 J_k^* as

$$J_k^*(x) = \min_{u_k \in U} Q_k^*(x, u_k). \quad (5.5)$$

4 which is precisely the dynamic programming update (by definition) in (5.4).
 5 We can also write dynamic programming completely in terms of Q-factors
 6 as follows.

Dynamic programming written in terms of the Q-factor

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 be 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

$$\begin{aligned} \text{value function} &\equiv \text{cost-to-go } J^*(x) \\ \text{action-value function} &\equiv \text{Q-factor } Q^*(x, u). \end{aligned}$$

Since the two functions are equivalent, we will call both as “value functions”. The difference will be clear from context.

7 5.4 Stochastic dynamic programming: Value 8 Iteration

9 Let us now see how dynamic programming looks for a Markov Decision
 10 Process (MDP). As we saw in Chapter 3, we can think of MDPs as
 11 stochastic dynamical systems denoted by

$$x_{k+1} = f_k(x_k, u_k) + \epsilon_k; \quad x_0 \text{ is given.}$$

12 We will assume that we know the statistics of the noise ϵ_k at each time-step
 13 (say it is a Gaussian). Stochastic dynamical systems are very different from
 14 deterministic dynamical systems, given the same sequence of controls
 15 (u_0, \dots, u_{T-1}) , we may get different state trajectories (x_0, x_1, \dots, x_T)

1 depending upon the realization of noise $(\epsilon_0, \dots, \epsilon_{T-1})$. How should
 2 we find a good control trajectory then? One idea is to modify (5.3) to
 3 minimize the expected value of the cost over all possible state-trajectories

$$J(x_0; u_0, \dots, u_{T-1}) = \mathbb{E}_{(\epsilon_0, \dots, \epsilon_{T-1})} \left[q_f(x_T) + \sum_{k=0}^{T-1} q_k(x_k, u_k) \right] \quad (5.7)$$

4 Suppose we minimized the above expectation and obtained the value
 5 function $J^*(x_0)$ and the optimal control trajectory $(u_0^*, \dots, u_{T-1}^*)$. As
 6 the robot starts executing this trajectory, the realized versions of the noise
 7 ϵ_k might differ a lot from their expected value, and the robot may find
 8 itself in very different states x_k than the average-case states considered
 9 in (5.10).

10 **Feedback controls** The concept of feedback control is a powerful way
 11 to resolve this issue. Instead of seeking $u_k^* \in U$ as the solutions of (5.10),
 12 we instead seek a *function*

$$u_k(x) : X \mapsto U \quad (5.8)$$

13 that maps the state-space X to a control U . Effectively, given a feedback
 14 control $u_k(x)$ the robot knows what control to apply at its current realized
 15 state $x_k \in X$, namely $u_k(x_k)$, even if the realized state x_k is very different
 16 from the average-case state. Feedback controls are everywhere and are
 17 critical to using controls in the real world. For instance, when you tune
 18 the shower faucet to give you a comfortable water temperature, you
 19 are constantly estimating the state (feedback using the temperature) and
 20 turning the faucet accordingly. Doing this without feedback would leave
 21 you terribly cold or scalded. We will denote the space of all feedback
 22 controls $u_k(\cdot)$ that depend on the state $x \in X$ by

$$u_k(\cdot) \in \mathcal{U}(X).$$

23 **Control policy** A sequence of feedback controls

$$\pi = (u_0(\cdot), u_1(\cdot), \dots, u_{T-1}(\cdot)). \quad (5.9)$$

24 is called a control policy. This is an object that we will talk about often. It
 25 is important to remember that a control policy is set of controllers (usually
 26 feedback controls) that are executed at each time-step of a dynamic
 27 programming problem.

The **stochastic optimal control problem** finds a sequence of

❗ 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.

feedback controls $(u_0(\cdot), u_1(\cdot), \dots, u_T(\cdot))$ that minimizes

$$J(x_0; u_0(\cdot), \dots, u_{T-1}(\cdot)) = \mathbb{E}_{(\epsilon_0, \dots, \epsilon_{T-1})} \left[q_f(x_T) + \sum_{k=0}^{T-1} q_k(x_k, u_k(x_k)) \right].$$

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.

❗ 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 Dijkstra's algorithm no longer works, as is, if the edges in the graph
2 are stochastic but we can use the principal of dynamic programming
3 to write the solution for the stochastic optimal control problem. The
4 idea remains the same, we compute a sequence of cost-to-go functions
5 $J_T^*(x), J_{T-1}^*(x), \dots, J_0^*(x)$, and in particular $J_0^*(x_0)$, proceeding *back-*
6 *wards*.

Finite-horizon dynamic programming for stochastic systems.

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 $x \in X$.

7 Just like (5.4), we solve a sub-problem for one time-instant at each
8 iteration. But observe a few importance differences in (5.11) compared
9 to (5.4).

- 10 1. There is an expectation over the noise ϵ_k in the second term in the
11 curly brackets. The second term in the curly brackets is the average of
12 the cost-to-go of the truncated sub-problems from time $k + 1, \dots, T$
13 over all possible starting states $x' = f_k(x_k, u_k(x_k)) + \epsilon_k$. This
14 makes sense, after taking the control $u_k(x_k)$, we may find the robot
15 at any of the possible states $x' \in X$ depending upon different
16 realizations of noise ϵ_k and the cost-to-go from x_k is therefore the
17 average of the cost-to-go from each of those states (according to the
18 principal of dynamic programming).
- 19 2. The minimization in (5.11) is performed over a function

$$\mathcal{U}(X) \ni u_k(\cdot) : X \mapsto U.$$

Since our set of states and controls is finite, this involves finding a table of size $|X| \times |U|$ for each iteration. In (5.4), we only had to search over a set of values $u_k \in U$ of size $|U|$. At the end of dynamic programming, we have a sequence of feedback controls

$$(u_0^*(\cdot), u_1^*(\cdot), \dots, u_{T-1}^*(\cdot)).$$

Each feedback control $u_k^*(x)$ tells us what control the robot should pick if it finds itself at a state x at time k .

3. If we know the dynamical system, not in its functional form $x_{k+1} = f_k(x_k, u_k) + \epsilon_k$ but rather as a transition matrix $P(x_{k+1} | x_k, u_k)$ (like we had in Chapter 2) then the expression in (5.11) simply becomes

$$J_k^*(x) = \min_{u_k(\cdot) \in \mathcal{U}(X)} \left\{ q_k(x, u_k(x)) + \mathbb{E}_{x' \sim P(\cdot | x_k, u_k(x_k))} [J_{k+1}^*(x')] \right\} \quad (5.12)$$

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?

Computational complexity The form in (5.12) helps us understand the computational complexity, each sub-problem performs $|X| \times |X| \times |U|$ amount of work and therefore the total complexity of stochastic dynamic programming is

$$\mathcal{O}(T|X|^2|U|).$$

Naturally, the quadratic dependence on the size of the state-space is an even bigger hurdle while implementing dynamic programming for stochastic systems.

5.4.1 Infinite-horizon problems

In the previous section, we put a lot of importance on the horizon T for dynamic programming. This is natural: if the horizon T changes, say you are in a hurry to get to school, the optimal trajectory may take control inputs that incur a lot of runtime cost simply to reach closer to the goal state (something that keeps the terminal cost small). In most, real-world problems, it is not very clear what value of T we should pick. We therefore formulate the dynamic programming problem as something that also allows a trajectory of infinite steps but also encourages the length of the trajectory to be small enough in order to be meaningful. Such problems are called infinite-horizon problems ($T \rightarrow \infty$).

Stationary dynamics and run-time cost We think of infinite-horizon problems in the following way: at any time-step, the length of the trajectory remaining for the robot to traverse is infinite. It helps in this case to solve a restricted set of problems where the system dynamics and run-time cost do not change as a function of time (they only change as a function of the

1 state and the control). We will set

$$q(x, u) \equiv q_k(x, u),$$

$$f(x, u) \equiv f_k(x, u)$$

2 for all $x \in X$ and $u \in U$. Such a condition is called stationarity. If the
 3 system is stochastic, we also require that the distribution of noise ϵ_k does
 4 not change as a function of time (it could change in (5.11) but we did not
 5 write it so). The infinite-horizon setting is never quite satisfied in practice
 6 but it is a reasonable formulation for problems that run for a long length
 7 of time.

8 **Infinite-horizon objective** The objective that we desire be minimized
 9 by an infinite-horizon control policy

$$\pi = (u_0(\cdot), u_1(\cdot), \dots, u_T(\cdot), u_{T+1}(\cdot), \dots)$$

10 is defined in terms of an asymptotic limit

$$J(x_0; \pi) = \lim_{T \rightarrow \infty} \mathbf{E}_{(\epsilon_0, \dots, \epsilon_{T-1})} \left[\sum_{k=0}^{T-1} \gamma^k q(x_k, u_k(x_k)) \right]. \quad (5.13)$$

11 and we again wish to solve for the optimal cost-to-go

$$J^*(x_0) = \operatorname{argmin}_{\pi} J^*(x_0; \pi). \quad (5.14)$$

12 Thus the infinite horizon costs of a policy is the limit of its finite horizon
 13 costs as the horizon tends to infinity. Notice a few important differences
 14 when compared to (5.7).

- 15 1. The objective is a limit, it is effectively the cost of the trajectory as
 16 it is allowed to stretch for a larger and larger time-horizon.
- 17 2. There is no terminal cost in the objective function; this makes sense
 18 because an explicit terminal state x_T does not exist anymore. In
 19 infinite-horizon problems, you should think of the terminal cost
 20 as being incorporated inside the run-time cost $q(x, u)$ itself, e.g.,
 21 move the robot to minimize the fuel used at *this* time instant but also
 22 move it in a way that it reaches the goal *at some time in the future*.
- 23 3. **Discount factor**— Depending upon what controls we pick, the
 24 summation

$$\sum_{k=0}^T q(x_k, u_k(x_k))$$

25 can diverge to infinity as $T \rightarrow \infty$ and thereby a meaningful solution
 26 to the infinite-horizon problem may not exist. In order to avoid this,
 27 we use a scalar

$$\gamma \in (0, 1)$$

28 known as the discount factor in the formulation. It puts more

1 emphasis on costs incurred earlier in the trajectory than later ones
 2 and thereby encourages the length of the trajectory to be small.
 3 Notice that $\sum_{k=0}^{\infty} \alpha^k = 1/(1 - \alpha)$ if $|\alpha| < 1$, so if the cost
 4 $|q(x_k, u_k(x_k))| < 1$, then we know that the objective in (5.13)
 5 always converges.

6 **Stochastic shortest path problems** It is important to remember that
 7 the discount factor is chosen by the user, no one prescribes it. There is
 8 also a class of problems where we may choose $\gamma = 1$ but in these cases,
 9 there should exist some *essentially terminal* state in the state space where
 10 we can keep taking a control such that the runtime cost $q(x, u)$ is zero.
 11 Otherwise, the objective will diverge. The goal region in the grid-world
 12 problem could be an example of such state. Such problems are called
 13 stochastic shortest path problems because the time-horizon is not *actually*
 14 infinite, we just do not know how many time-steps it will take for the robot
 15 to go to the goal location. Naturally, stochastic shortest path problems
 16 are a generalization of the shortest path problem solved by Dijkstra's
 17 algorithm. The algorithms we discuss next will work for such problems.

18 **Stationary policy** It seems a bit cumbersome to carry around an infinitely
 19 long sequence of feedback controls in infinite-horizon problems. Since
 20 there is an infinitely-long trajectory yet to be traveled *at any given time-*
 21 *step*, the optimal control action that we take should only depend upon the
 22 current state. This is indeed true mathematically. If $J^*(x)$ is the optimal
 23 cost-to-go in the infinite-horizon problem starting from a state x , using
 24 the principle of dynamic programming, we should also have that we can
 25 split this cost as the best one-step cost of the current state x added to the
 26 optimal cost-to-go from the state $f(x, u)$ realized after taking the optimal
 27 control u :

$$J^*(x) = \min_{u(x) \in \mathcal{U}(X)} \mathbb{E}_{\epsilon} [q(x, u(x)) + J^*(f(x, u(x)) + \epsilon)]. \quad (5.15)$$

28 We will study this equation in depth soon. But if we find the minimum at
 29 $u^*(x)$ for this equation, then we can run the policy

$$\pi^* = (u^*(\cdot), u^*(\cdot), \dots, u^*(\cdot), \dots)$$

30 for the entire infinite horizon. Such a policy is called a stationary
 31 policy. Intuitively, since the future optimization problem (tail of dynamic
 32 programming) from a given state x looks the same regardless of the time
 33 at which we start, optimal policies for the infinite-horizon problem can
 34 be found even inside the restricted class of policies where the feedback
 35 control does not change with time k .

36 We will almost exclusively deal with stationary policies in this course.

5.4.2 Dynamic programming for infinite-horizon problems

We wish to compute the optimal cost-to-go of starting from a state x and taking an infinitely long trajectory that minimizes the objective (5.13). We will exploit the equation in (5.15) and develop an iterative algorithm to compute 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.,

$$\forall x \in X, \quad \left| J^{(i)}(x) - J^{(i+1)}(x) \right| < \text{small tolerance.}$$

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}_ϵ 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

Let us observe a few important things in the above sequence of updates. First, at each iteration, we are updating the values of all $|X|$ states. This involves $|X|^2|U|$ amount of work per iteration. How many such iterations N do we need until the value function converges? We will see in a bit, that

$$\forall x \in X, \quad J^*(x) = \lim_{N \rightarrow \infty} J^{(N)}(x).$$

Again, we really only wanted to compute the cost-to-go $J^*(x_0)$ from some initial state x_0 but computed the value function at all states $x \in X$.

1 **Q-Iteration** Just like we wrote dynamic programming in terms of the
 2 Q-factor, we can also write value iteration to find the optimal Q-factor
 3 $Q^*(x, u)$, i.e., the optimal cost-to-go of the infinitely-long trajectory that
 4 starts at state x , takes a control u at the first time-step and therefore follows
 5 the optimal policy.

6 1. We can again initialize $Q^{(0)}(x, u) = q(x, u)$ for all $x \in X$ and
 7 $u \in U$.

8 2. The Bellman update in terms of the Q-factor becomes

$$Q^{(i+1)}(x, u) = E_{\epsilon} \left[q(x, u) + \gamma \min_{u' \in U} Q^{(i)}(f(x, u) + \epsilon, u') \right] \tag{5.18}$$

9 and update this for all $x \in X$ and all $u \in U$.

10 3. The feedback control is the control at that state that minimizes the
 11 Q-factor

$$\forall x \in X, \quad u^*(x) = \operatorname{argmin}_{u' \in U} Q^{(N)}(x, u') \tag{5.19}$$

12 and the control policy is

$$\pi^* = (u^*, u^*, \dots)$$

13 Notice how we can directly find the u' that has the smallest value of
 14 $Q^{(N)}$ and set it to be our feedback control.

15 5.4.3 An example

16 Let us consider a grid-world example. A robot would like to reach a
 17 goal region (marked in green) and we are interested in computing the
 18 cost-to-go from different parts of the domain. Gray cells are obstacles that
 19 the robot cannot enter. At each step the robot can move in four directions
 20 (north, east, west, south) with a small dynamics noise which keeps it at
 21 the original cell in spite of taking the control. These pictures are when
 22 the run-time-cost is *negative*, i.e., the robot gets a certain reward $q(x, u)$
 23 for taking the control u at cell x . Dynamic programming (and value
 24 iteration) also works in this case and we simply replace all minimizations
 25 by maximizations in the equations.

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	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

26

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	0.75	0
0	0	0	0	0	0	0	0.75	1	0
0	0	0	0	0	0	0	0	0	0

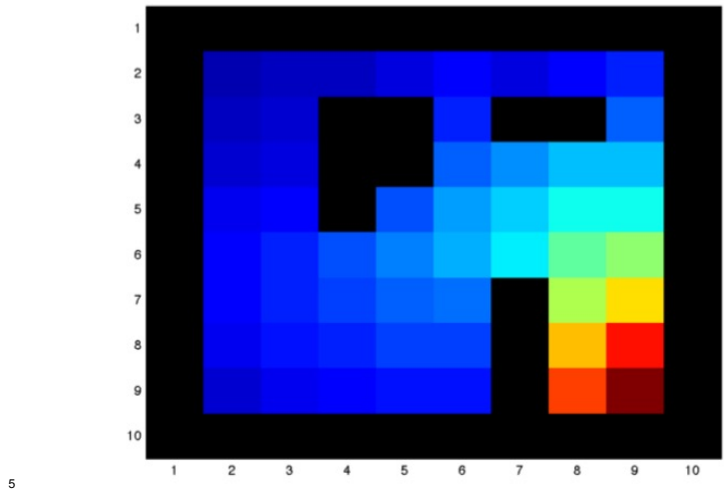
2

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	0.51	0
0	0	0	0	0	0	0	0.56	1.43	0
0	0	0	0	0	0	0	1.43	1.9	0
0	0	0	0	0	0	0	0	0	0

3

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	0.59	0.69	0	0	1.52	0	0	2.13	0
0	0.75	0.90	0	0	2.12	2.55	2.98	3.00	0
0	0.95	1.18	0	2.00	2.70	3.22	3.80	3.88	0
0	1.20	1.55	1.87	2.41	2.92	3.51	4.52	5.00	0
0	1.15	1.47	1.74	2.05	2.25	0	5.34	6.47	0
0	0.99	1.26	1.49	1.72	1.74	0	6.69	8.44	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

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.
 8 These results are valid under a very general set of conditions and make
 9 value iteration work for a large number of real-world problems; they are
 10 at the the heart of all modern algorithms. We will not derive them (it is

1 easy but cumbersome) but you should commit them to memory and try to
2 understand them intuitively.

3 **Value iteration converges.** Given *any* initialization $J^{(0)}(x)$ for all
4 $x \in X$, the sequence of value iteration estimates $J^{(i)}(x)$ converges to the
5 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
7 the 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
9 words, if we find some other function $J'(x)$ that satisfies the Bellman
10 equation, we are 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
14 short. 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
16 have a policy in hand, and wish to find the cost-to-go of this policy, i.e.,
17 “evaluate the policy” we can initialize $J^{(0)}(x) = 0$ for all $x \in X$ and
18 perform the sequence 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 Observe that the corresponding equation for policy equation (5.20) does
22 not have the minimization over controls. This allows us to write the
23 updates in (5.21) as the solution of a linear system of equations. Since we
24 are in a finite state-space, we can write the cost-to-go as a large vector

$$J^\pi := [J^\pi(x_1), J^\pi(x_2), \dots, J^\pi(x_n)]^\top$$

25 where n is the number of total states in the state-space. We create a similar
26 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 ϵ is equivalent to an expecta-
 2 tion over the next state of the system, let us rewrite the dynamics part
 3 $f(x, u(x)) + \epsilon$ in terms of the Markov transition matrix

$$T_{x,x'} = P(x' | x, u(x))$$

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
 7 of a particular control policy π .

8 5.5 Stochastic dynamic programming: Policy 9 Iteration

10 Value iteration converges exponentially quickly, but asymptotically. Note
 11 that the number of states $|X| = n$ is finite and so is the number of controls
 12 $|U|$. So this should seem funny, one would expect that we should be able
 13 to find the optimal cost $J^*(x)$ in finite time if the problem is finite. After
 14 all we need to find $|X|$ numbers $J^*(x_1), \dots, J^*(x_n)$. This intuition is
 15 correct and, in this section, we will discuss an algorithm called policy
 16 iteration which is a more efficient version of value iteration.

17 The idea behind policy iteration is quite simple: given a stationary
 18 policy for an infinite-horizon problem $\pi = (u(\cdot), \dots, u(\cdot))$, we can
 19 evaluate this policy to obtain its cost-to-go $J^{\pi}(x)$. If we now set the
 20 feedback control to be

$$\tilde{u}(x) = \underset{u \in U}{\operatorname{argmin}} \mathbb{E}_{\epsilon} [q(x, u) + \gamma J^{\pi}(f(x, u) + \epsilon)], \quad (5.23)$$

21 i.e., we construct a *new control policy* that finds the best control to execute
 22 in the first step $\tilde{u}(\cdot)$ and thereafter it executes the old feedback control $u(\cdot)$

$$\pi^{(1)} = (\tilde{u}(\cdot), u(\cdot), \dots),$$

23 then the cost-to-go of the new policy $\pi^{(1)}$ has to be better:

$$\forall x \in X, \quad J^{\pi^{(1)}}(x) \leq J^{\pi}(x).$$

24 We don't have to stop at one time-step, we can patch the old policy at the
 25 first two time-steps to get

$$\pi^{(2)} = (\tilde{u}(\cdot), \tilde{u}(\cdot), \dots),$$

26 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
4 π^* starting from some initial stationary policy.

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 π^* .

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 $\pi^{(k)}$. In other words, 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 can 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

$$u^{(k+1)}(x) = \operatorname{argmin}_{u \in U} \mathbb{E}_{\epsilon} \left[q(x, u) + \gamma J^{\pi^{(k)}}(f(x, u) + \epsilon) \right]$$

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

$$\forall x \in X, \quad u^{(k+1)}(x) = u^{(k)}(x).$$

5 Just like value iteration converges to the optimal value function, it can
6 be shown that policy iteration produces a sequence of improved policies

$$\forall x \in X, \quad J^{\pi^{(k+1)}}(x) \leq J^{\pi^{(k)}}(x)$$

i For large problems, we use methods for solving large linear systems such as Lanczos iteration. Typical policy evaluation problems are also sparse (why?) so we can use things like the Kaczmarz method to solve the linear system.

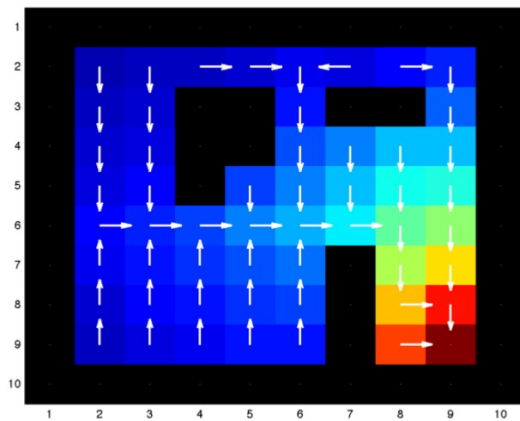
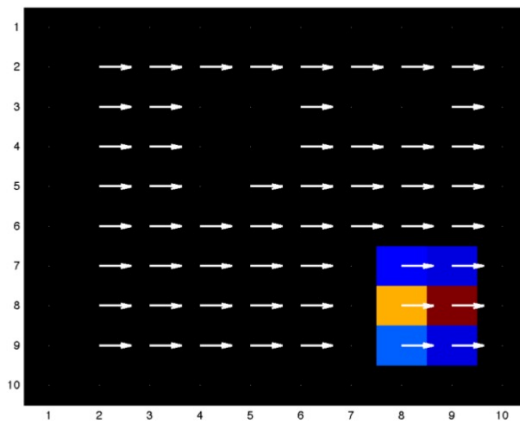
1 and converges to the optimal cost-to-go

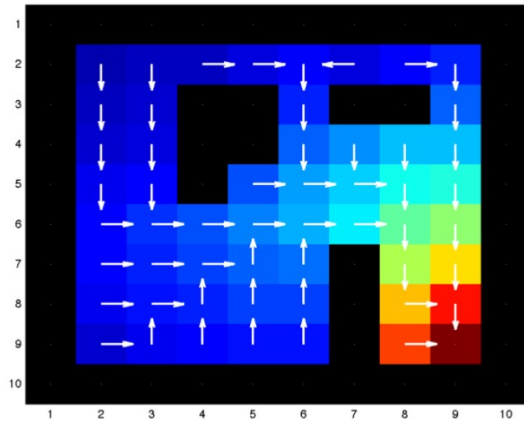
$$\forall x \in X, \quad J^*(x) = \lim_{N \rightarrow \infty} J^{\pi(N)}(x).$$

2 The key property of policy iteration is that we need a finite number of
 3 updates to the policy to find the optimal policy. Notice that this does
 4 not mean always mean that we are doing less work than value iteration
 5 in policy iteration. Observe that the policy evaluation step in the policy
 6 iteration algorithm performs a number of Bellman equation updates. But
 7 typically, it is observed in practice that policy iteration is much cheaper
 8 computationally than value iteration.

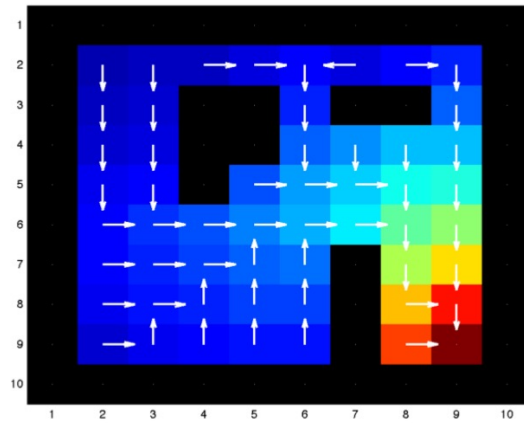
9 5.5.1 An example

10 Let us go back to our example for value iteration. In this case, we will
 11 visualize the controller $u^{(k)}(x)$ at each cell x as arrows pointing to some
 12 other cell. The cells are colored by the value function for that particular
 13 stationary policy.





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

1 in dynamic programming, we want to pick a feedback control which leads
 2 to a trajectory that achieves a minimum of some run-time cost and a
 3 terminal cost. We will assume that both the run-time and terminal costs
 4 are *quadratic* in the state and control input, i.e.,

$$q(x, u) = \frac{1}{2}x^\top Qx + \frac{1}{2}u^\top Ru \quad (6.1)$$

5 where $Q \in \mathbb{R}^{d \times d}$ and $R \in \mathbb{R}^{m \times m}$ are symmetric, positive semi-definite
 6 matrices

$$Q = Q^\top \succeq 0, \quad R = R^\top \succeq 0.$$

7 Effectively, if Q were a diagonal matrix, a large diagonal entry would Q_{ii}
 8 models our desire that the trajectory of the system should not have a large
 9 value of the state x_i along its trajectories. We want these matrices to be
 10 positive semi-definitive to prevent dynamic programming from picking
 11 a trajectory which drives down the run-time cost to negative infinity by
 12 picking.

13 **Example** Consider the discrete-time equivalent of the so-called double
 14 integrator $\ddot{z}(t) = u(t)$. The linear system in this case (obtained by creating
 15 two states $x := [z(t), \dot{z}(t)]$ is

$$x_{k+1} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} x_k + \begin{bmatrix} 0 \\ \Delta t \end{bmatrix} u_k.$$

16
 17 First, note that a continuous-time linear dynamical system $\dot{x} = Ax$ is
 18 asymptotically stable, i.e., from any initial condition $x(0)$ its trajectories go
 19 to the equilibrium point $x = 0$ ($x(t) \rightarrow 0$ as $t \rightarrow \infty$). Asymptotic stability
 20 for continuous-time dynamical systems occurs if all eigenvalues of A are
 21 strictly negative. A discrete-time linear dynamical system $x_{k+1} = Ax_k$
 22 is asymptotically stable if all eigenvalues of A have magnitude strictly
 23 smaller than 1, $|\lambda(A)| < 1$.

24 A typical trajectory of the double integrator will look as follows.
 25 Suppose we would like to pick a different controller that more quickly
 26 brings the system to its equilibrium. One way of doing so is to minimize

$$J = \sum_{k=0}^T \|x_k\|^2$$

27 which represents how far away both the position and velocity are from zero
 28 over all times k . The following figure shows the trajectory that achieves a
 29 small value of J .

i 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.

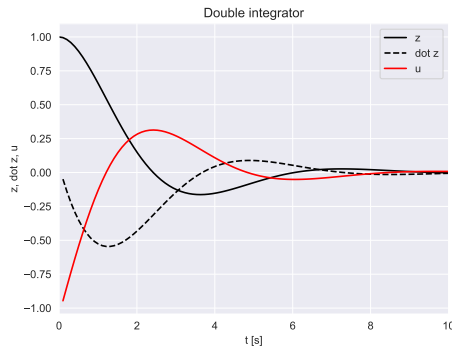


Figure 6.1: The trajectory of $z(t)$ as a function of time t for a double integrator $\ddot{z}(t) = u$ where we have chosen a stabilizing (i.e., one that makes the system asymptotically stable) controller $u = -z(t) - \dot{z}(t)$. Notice how the trajectory starts from some initial condition (in this case $z(0) = 1$ and $\dot{z}(0) = 0$) and moves towards its equilibrium point $z = \dot{z} = 0$.

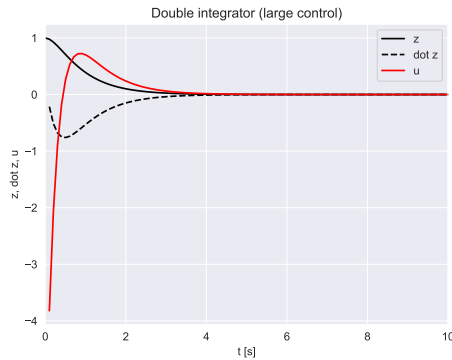


Figure 6.2: The trajectory of $z(t)$ as a function of time t for a double integrator $\ddot{z}(t) = u$ where we have chosen a large stabilizing control at each time $u = -5z(t) - 5\dot{z}(t)$. Notice how quickly the state trajectory converges to the equilibrium without much oscillation as compared to Figure 6.1 but how large the control input is at certain times.

- 1 This is obviously undesirable for real systems where we may want the
- 2 control input to be bounded between some reasonable values (a car cannot
- 3 accelerate by more than a certain threshold). A natural way of enforcing
- 4 this is to modify our our desired cost of the trajectory to be

$$J = \sum_{k=0}^T \left(\|x_k\|^2 + \rho \|u_k\|^2 \right)$$

- 5 where the value of the parameter ρ is something chosen by the user to
- 6 give a good balance of how quickly the trajectory reaches the equilibrium
- 7 point and how much control is exerted while doing so. Linear-Quadratic-
- 8 Regulator (LQR) is a generalization of this idea, notice that the above
- 9 example is equivalent to setting $Q = I_{d \times d}$ and $R = \rho I_{m \times m}$ for the
- 10 run-time cost in (6.1).

1 **Back to LQR** With this background, we are now ready to formulate
 2 the Linear-Quadratic-Regulator (LQR) problem which is simply dynamic
 3 programming for a linear dynamical system with quadratic run-time cost.
 4 In order to enable the system to reach the equilibrium state even if we have
 5 only a finite time-horizon, we also include a quadratic cost

$$q_f(x) = \frac{1}{2}x^\top Q_f x. \quad (6.2)$$

6 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.

7 6.1.1 Solution of the discrete-time LQR problem

8 We know the principle of dynamic programming and can apply it to solve
 9 the LQR problem. As usual, we will compute the cost-to-go of a trajectory
 10 that starts at some state x and goes further by $T - k$ time-steps, $J_k(x)$
 11 backwards. Set

$$J_T^*(x) = \frac{1}{2}x^\top Q_f x \quad \text{for all } x.$$

12 Using the principle of dynamic programming, the cost-to-go J_{T-1} is
 13 given by

$$\begin{aligned} J_{T-1}^*(x_{T-1}) &= \min_u \left\{ \frac{1}{2} (x_{T-1}^\top Q x_{T-1} + u^\top R u) + J_T^*(Ax_{T-1} + Bu) \right\} \\ &= \min_u \left\{ \frac{1}{2} (x_{T-1}^\top Q x_{T-1} + u^\top R u + (Ax + Bu)^\top Q_f (Ax_{T-1} + Bu)) \right\}. \end{aligned}$$

14 We can now take the derivative of the right-hand side with respect to u to
 15 get

$$\begin{aligned} 0 &= \frac{d\text{RHS}}{du} \\ &= Ru + B^\top Q_f (Ax_{T-1} + Bu) \\ \Rightarrow u_{T-1}^* &= -(R + B^\top Q_f B)^{-1} B^\top Q_f A x_{T-1} \\ &\equiv -K_{T-1} x_{T-1}. \end{aligned} \quad (6.4)$$

16 where

$$K_{T-1} = (R + B^\top Q_f B)^{-1} B^\top Q_f A$$

1 is (surprisingly) also called the Kalman gain. The second derivative is
2 positive semi-definite

$$\frac{d^2\text{RHS}}{du^2} = R + B^\top Q_f B \succeq 0$$

3 so we know that u_{T-1}^* is a minimum of the convex quantity on the right-
4 hand side. Notice that the optimal control u_{T-1}^* is a linear function of the
5 state x_{T-1} . Let us now expand the cost-to-go J_{T-1} using this optimal
6 value (the subscript $T - 1$ on the curly bracket simply means that all
7 quantities are at time $T - 1$)

$$\begin{aligned} J_{T-1}^*(x_{T-1}) &= \frac{1}{2} \left\{ x^\top Q x + u^{*\top} R u^* + (Ax + Bu^*)^\top Q_f (Ax + Bu^*) \right\}_{T-1} \\ &= \frac{1}{2} x_{T-1}^\top \left\{ Q + K^\top R K + (A - BK)^\top Q_f (A - BK) \right\}_{T-1} x_{T-1} \\ &\equiv \frac{1}{2} x_{T-1}^\top P_{T-1} x_{T-1} \end{aligned}$$

8 where we set the stuff inside the curly brackets to the matrix P which is
9 also positive semi-definite. This is great, the cost-to-go is also a quadratic
10 function of the state x_{T-1} . Let us assume that this pattern holds for all
11 time steps and the cost-to-go of the optimal LQR trajectory starting from
12 a state x and proceeding forwards for $T - k$ time-steps is

$$J_k^*(x) = \frac{1}{2} x^\top P_k x.$$

13 We can now repeat the same exercise to get a recursive formula for P_k in
14 terms of P_{k+1} . This is the *solution* of dynamic programming for the LQR
15 problem and it looks as follows.

$$\begin{aligned} P_T &= Q_f \\ K_k &= (R + B^\top P_{k+1} B)^{-1} B^\top P_{k+1} A \\ P_k &= Q + K_k^\top R K_k + (A - BK_k)^\top P_{k+1} (A - BK_k), \end{aligned} \quad (6.5)$$

16 for $k = T - 1, T - 2, \dots, 0$. There are a number of important observations
17 to be made from this calculation:

- 18 1. The optimal controller $u_k^* = -K_k x_k$ is a linear function of the state
19 x_k . This is only true for linear dynamical systems with quadratic
20 costs. Notice that both the state and control space are infinite sets
21 but we have managed to solve the dynamic programming problem
22 to get the optimal controller. We could not have done it if the run-
23 time-terminal costs were not quadratic or if the dynamical system
24 were not linear. Can you say why?
- 25 2. The cost-to-go matrix P_k and the Kalman gain K_k do not depend
26 upon the state and can be computed ahead of time if we know what
27 the time horizon T is going to be.
- 28 3. The Kalman gain changes with time k . Effectively, the LQR

1 controller picks a large control input to quickly reduce the run-time
 2 cost at the beginning (if the initial condition were such that the
 3 run-time cost of the trajectory would be very large) and then gets
 4 into a balancing act where it balances the control effort and the state-
 5 dependent part of the run-time cost. LQR is an optimal way to strike
 6 a balance between the two examples in Figure 6.1 and Figure 6.2.

7 The careful reader will notice how the equations in (6.5) and our
 8 remarks about them are similar to the update equations of the Kalman filter
 9 and our remarks there. In fact we will see shortly how spookily similar the
 10 two are. The key difference is that Kalman filter updates run forwards in
 11 time and update the covariance while LQR updates run backwards in time
 12 and update the cost-to-go matrix P . This is not surprising because LQR
 13 is an optimal control problem, its update equations should run backward
 in time like the Dijkstra's algorithm.

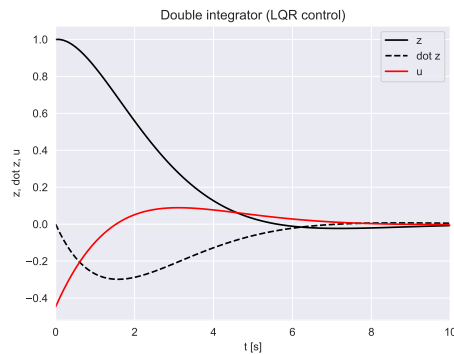


Figure 6.3: The trajectory of $z(t)$ as a function of time t for a double integrator $\ddot{z}(t) = u$ where we have chosen a controller obtained from LQR with $Q = I$ and $R = 5$. This gives the controller to be about $u = -0.45z(t) - 1.05\dot{z}(t)$. Notice how we still get stabilization but the control 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 .

❗ 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

15 6.2 Hamilton-Jacobi-Bellman equation

16 This section will show how the principle of dynamic programming looks
 17 for continuous-time deterministic dynamical systems

$$\dot{x} = f(x, u), \quad \text{with } x(0) = x_0.$$

18 As we discussed in Chapter 3, we can think of this as the limit of discrete-
 19 time dynamical system $x_{k+1} = f^{\text{discrete}}(x_k, u_k)$ as the time discretization
 20 goes to zero. Just like we have a sequence of controls in the discrete-time
 21 case, we have a continuous curve that determines the control (let us also
 22 call it the control sequence)

$$\{u(t) : t \in \mathbb{R}_+\}$$

1 which gives rise to a trajectory of the states

$$\{x(t) : t \in \mathbb{R}_+\}$$

2 for the dynamical system. Let us consider the case when we want to
3 find control sequences that minimize the integral of the cost along the
4 trajectory that stops at some fixed, finite time-horizon T :

$$q_f(x(T)) + \int_0^T q(x(t), u(t)) dt.$$

5 This cost is again a function of the run-time cost and a terminal cost.

❗ 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.

6 The principle of dynamic programming principle is still valid: if we
7 have an optimal control trajectory $\{u^*(t) : t \in [0, T]\}$ we can chop it up
8 into two parts at some intermediate time $t \in [0, T]$ and claim that the tail
9 is optimal. In preparation for this, let us define the cost-to-go of going
10 forward by $T - t$ 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

$$\begin{aligned} J^*(x(t), t) &= \min_{u(s), t \leq s \leq T} \left\{ q_f(x(T)) + \int_t^T q(x(s), u(s)) ds \right\} \\ &= \min_{u(s), t \leq s \leq T} \left\{ q_f(x(T)) + \int_t^{t+\Delta t} q(x(s), u(s)) ds + \int_{t+\Delta t}^T q(x(s), u(s)) ds \right\} \\ &= \min_{u(s), t \leq s \leq t+\Delta t} \left\{ J^*(x(t + \Delta t), t + \Delta t) + \int_t^{t+\Delta t} q(x(s), u(s)) ds \right\}. \end{aligned}$$

11 We now take the Taylor approximation of the term $J^*(x(t + \Delta t), t + \Delta t)$

1 as follows

$$\begin{aligned} & J^*(x(t + \Delta t), t + \Delta t) - J^*(x(t), t) \\ & \approx \partial_x J^*(x(t), t) (x(t + \Delta t) - x(t)) + \partial_t J^*(x(t), t) \Delta t \\ & \approx \partial_x J^*(x(t), t) f(x(t), u(t)) \Delta t + \partial_t J^*(x(t), t) \Delta t \end{aligned}$$

2 where $\partial_x J^*$ and $\partial_t J^*$ denote the derivative of J^* with respect to its first
3 and second argument respectively. We substitute this into the minimization
4 and collect terms of Δt to get

$$0 = \partial_t J^*(x(t), t) + \min_{u(t) \in U} \{q(x(t), u(t)) + f(x(t), u(t)) \partial_x J^*(x(t), t)\}. \quad (6.7)$$

5 Notice that the minimization in (6.7) is only over *one* control input
6 $u(t) \in U$, this is the control that we should take at time t . (6.7) is called
7 the Hamilton-Jacobi-Bellman (HJB) equation. Just like the Bellman
8 equation

$$J_k^*(x) = \min_{u \in U} \{q_k(x, u) + J_{k+1}^*(f(x, u))\}.$$

9 has two quantities x and the time k , the Hamilton-Jacobi-Bellman equation
10 also has two quantities x and continuous time t . Just like the Bellman
11 equation is solved backwards in time starting from T with $J_k^*(x) = q_f(x)$,
12 the HJB 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.

13 6.2.1 Infinite-horizon HJB

14 The infinite-horizon problem with the HJB equation is easy: since we
15 know that the optimal cost-to-go is not a function of time, we have

$$\partial_t J^*(x, t) = 0$$

16 and therefore $J^*(x)$ satisfies

$$0 = \min_{u \in U} \{q(x, u) + f(x, u) \partial_x J^*(x)\}. \quad (6.8)$$

17 In this case, the above equation makes sense only if the integral of the run-
18 time cost with the optimal controller $\int_0^\infty q(x(t), u^*(x(t))) dt$ remains
19 bounded and does not diverge to infinity. Therefore typically in this
20 problem we will set $q(0, 0) = 0$, i.e., there is no cost for the system being
21 at the origin with zero control, otherwise the integral of the run-time cost
22 will never be finite. This also gives the boundary condition $J^*(0) = 0$ for
23 the HJB equation.

6.2.2 Solving the HJB equation

The HJB equation is a partial differential equation (PDE) because there is one cost-to-go from every state $x \in X$ and for every time $t \in [0, T]$. It belongs to a large and important class of PDEs, collectively known as Hamilton-Jacobi-type equations. As you can imagine, since dynamic programming is so pervasive and solutions of DP are very useful in practice for a number of problems, there have been many tools invented to solve the HJB equation. These tools have applications to a wide variety of problems, from understanding how sound travels in crowded rooms to how light diffuses in an animated movie scene, to even obtaining better algorithms to train deep networks (<https://arxiv.org/abs/1704.04932>). HJB equations are usually never exactly solvable and a number of approximations need to be made in order to 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.

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

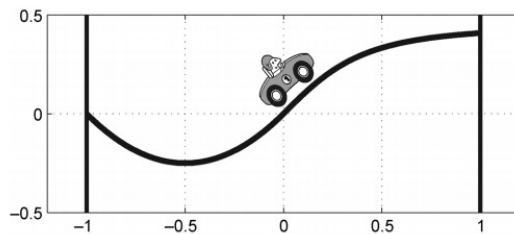


Figure 6.4: A car whose position is given by $z(t)$ would like to climb the hill to its right and reach the top with minimal velocity. The car rolls on the hill without friction. The run-time cost is zero everywhere inside the state-space. Terminal cost is -1 for hitting the left boundary ($z = -1$) 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.

and velocity (z, \dot{z}) and we can solve a two-dimensional HJB equation to obtain the optimal cost-to-go from any state, as done by the authors Yuval Tassa and Tom Erez in “Least Squares Solutions of the HJB Equation With Neural Network Value-Function Approximators” (<https://homes.cs.washington.edu/fodorov/courses/amath579/reading/NeuralNet.pdf>). In practice, while solving the HJB PDE, one discretizes the state-space at given set of states and solves the HJB equation (6.7) on this grid using

- 1 numerical methods (these authors used neural networks to solve it). The
 2 end result looks as follows.

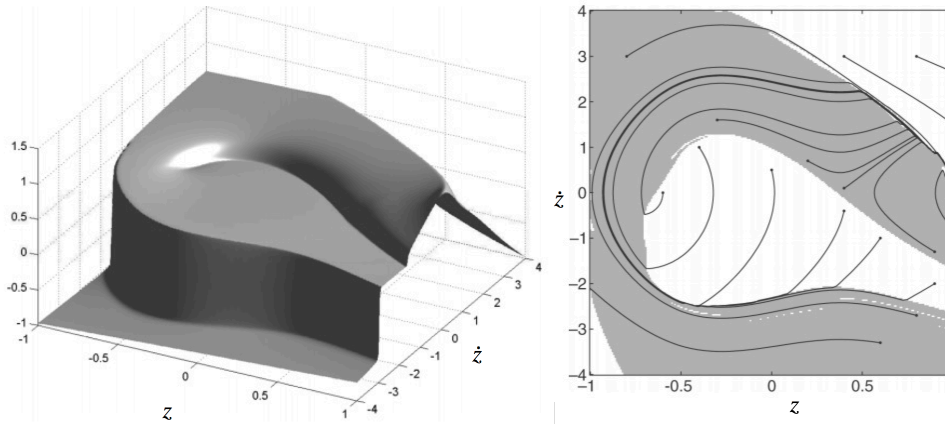


Figure 6.5: The left-hand side picture shows the infinite-horizon cost-to-go $J^*(z, \dot{z})$ for the car-on-the-hill problem. Notice how the value function is non-smooth at various places. This is quite typical of difficult dynamic programming problems. The right-hand side picture shows the optimal 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.

2

3 6.2.3 Continuous-time LQR

- 4 Consider a linear continuous-time dynamical system given by

$$\dot{x} = A x + B u; \quad x(0) = x_0.$$

- 5 In the LQR problem, we are interested in finding a control trajectory that
 6 minimizes, as usual, a cost function that is quadratic in states and controls,
 7 except that we have an integral of the run-time cost because our system is
 8 a continuous-time system

$$\frac{1}{2} x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T x(t)^\top Q x(t) + u(t)^\top R u(t) dt.$$

- 9 This is a very nice setup for using the HJB equation from the previous
 10 section.

- 11 Let us use our intuition from the discrete-time LQR problem and say
 12 that the optimal cost is quadratic in the states, namely,

$$J^*(x, t) = \frac{1}{2} x(t)^\top P(t) x(t);$$

- 13 notice that as usual the optimal cost-to-go is a function of the states x

1 and the time t because is the optimal cost of the continuous-time LQR
 2 problem if the system starts at a state x at time t and goes on until time
 3 $T \geq t$. We will now check if this J^* satisfies the HJB equation (we don't
 4 write the arguments $x(t)$, $u(t)$ etc. to keep the notation clear)

$$-\partial_t J^*(x, t) = \min_{u \in U} \left\{ \frac{1}{2} (x^\top Q x + u^\top R u) + (A x + B u)^\top \partial_x J^*(x, t) \right\} \quad (6.9)$$

5 from (6.7). The minimization is over the control input that we take *at time*
 6 t . Also notice the partial derivatives

$$\begin{aligned} \partial_x J^*(x, t) &= P(t) x. \\ \partial_t J^*(x, t) &= \frac{1}{2} x^\top \dot{P}(t) x. \end{aligned}$$

7 It is convenient in this case to see that the minimization can be per-
 8 formed using basic calculus (just like the discrete-time LQR problem), we
 9 differentiate with respect to u and set it to zero.

$$\begin{aligned} 0 &= \frac{\text{RHS of HJB}}{du} \\ \Rightarrow u^*(t) &= -R^{-1} B^\top P(t) x(t) \\ &\equiv -K(t) x(t). \end{aligned} \quad (6.10)$$

10 where $K(t) = R^{-1} B^\top P(t)$ is the Kalman gain. The controller is again
 11 linear in the states $x(t)$ and the expression for the gain is very simple in
 12 this case, much simpler than discrete-time LQR. Since $R \succ 0$, we also
 13 know that $u^*(t)$ computed here is the global minimum. If we substitute
 14 this value of $u^*(t)$ back into the HJB equation we have

$$\left. \left\{ \right\} \right|_{u^*(t)} = \frac{1}{2} x^\top \{ P A + A^\top P + Q - P B R^{-1} B^\top P \} x.$$

15 In order to satisfy the HJB equation, we must have that the expression
 16 above is equal to $-\partial_t J^*(x, t)$. We therefore have, what is called the
 17 Continuous-time Algebraic Riccati Equation (CARE), for the matrix
 18 $P(t) \in \mathbb{R}^{d \times d}$

$$-\dot{P} = P A + A^\top P + Q - P B R^{-1} B^\top P. \quad (6.11)$$

19 This is an ordinary differential equation for the matrix P . The derivative
 20 $\dot{P} = \frac{dP}{dt}$ stands for differentiating every entry of P individually with
 21 time t . The terminal cost is $\frac{1}{2} x(T)^\top Q_f x(T)$ which gives the boundary
 22 condition for the ODE as

$$P(T) = Q_f.$$

23 Notice that the ODE for the $P(t)$ travels backwards in time.

24 Continuous-time LQR has particularly easy equations, as you can see
 25 in (6.10) and (6.11) compared to those for discrete-time ((6.4) and (6.5)).
 26 Special techniques have been invented for solving the Riccati equation. I

1 used the function `scipy.linalg.solve_continuous_are` to obtain Figure 6.3
 2 using the continuous-time equations; the corresponding function for
 3 solving Discrete-time Algebraic Riccati Equation (DARE) which is given
 4 in (6.5) is `scipy.linalg.solve_discrete_are`. The continuous-time point-of-
 5 view also gives powerful connections to the Kalman filter, where you can
 6 show that the Kalman filter and LQR are duals of each other: in fact the
 7 equations for the Kalman filter (in continuous-time) and continuous-time
 8 LQR turn out to be exactly the same after you interchange appropriate
 9 quantities (!).

10 **Infinite-horizon LQR** Just like the infinite-horizon HJB equation has
 11 $\partial_t J^*(x, t) = 0$, if we have an infinite-horizon LQR problem, the cost
 12 matrix P should not be a function of time

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

13 The continuous-time algebraic Riccati equation in (6.11) now becomes

$$PA + A^\top P + Q - PBR^{-1}B^\top P.$$

14 with the cost-to-go being given by $J^*(x) = \frac{1}{2}x^\top Px$.

15 6.3 Stochastic LQR

16 We will next look at a very powerful result. Say we have a stochastic linear
 17 dynamical system

$$\dot{x}(t) = Ax(t) + Bu(t) + B_\epsilon \epsilon(t); x(0) \text{ is given}$$

18 where $\epsilon(t)$ is standard Gaussian noise $\epsilon(t) \sim N(0, I)$ that is uncorrelated
 19 in time and would like to find a control sequence $\{u(t) : t \in [0, T]\}$ that
 20 minimizes a quadratic run-time and terminal cost

$$\mathbb{E}_{\epsilon(t):t \in [0, T]} \left[\frac{1}{2}x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T x(t)^\top Q x(t) + u(t)^\top R u(t) dt \right].$$

21 over a finite-horizon T . Notice that since the system is stochastic now,
 22 we should minimize the expected value of the cost over all possible
 23 realizations of the noise $\{\epsilon(t) : t \in [0, T]\}$. This is a very challenging
 24 problem, conceptually it is the equivalent of dynamic programming for an
 25 MDP with an infinite number of states $x(t) \in \mathbb{R}^d$ and an infinite number
 26 of controls $u(t) \in \mathbb{R}^m$.

27 However, it turns out that the optimal controller that we should pick in
 28 this case is also given by the standard LQR problem

$$u^*(t) = -R^{-1}B^\top P(t) x(t)$$

with $-\dot{P} = PA + A^\top P + Q - PBR^{-1}B^\top P; P(T) = Q_f.$

29 We will not do the proof (it is easy but tedious, you can try to show it

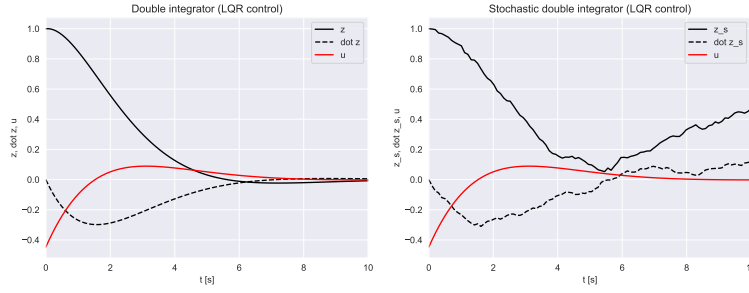


Figure 6.6: Comparison of the state trajectories of deterministic LQR and stochastic LQR problem with $B_\epsilon = [0.1, 0.1]$. The left panel is the same as that in Figure 6.3. The control input is the same in both cases but notice that the states in the plot on the right need not converge to the equilibrium 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 by writing the HJB equation for the stochastic LQR problem). This is a
 2 very surprising result because it says that even if the dynamical system
 3 had noise, the optimal control we should pick is exactly the same as the
 4 control we would have picked had the system been deterministic. It is a
 5 special property of the LQR problem and not true for other dynamical
 6 systems (nonlinear ones, or ones with non-Gaussian noise) or other costs.
 7 We know that the control $u^*(t)$ is the same as the deterministic case.
 8 Is the cost-to-go $J^*(x, t)$ also the same? If you think about this, the
 9 cost-to-go in the stochastic case has to be a bit larger than the deterministic
 10 case because the noise $\epsilon(t)$ is always going to non-zero when we run the
 11 system, the LQR cost $J^*(x_0, 0) = \frac{1}{2}x_0^\top P(0)x_0$ is, after all, only the cost
 12 of the deterministic problem. It turns out that the cost for the stochastic
 13 LQR case for an initial state x_0 is

$$\begin{aligned} J^*(x_0, 0) &= \mathbb{E}_{\epsilon(t):t \in [0, T]} \left[\frac{1}{2}x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T \dots dt \right] \\ &= \frac{1}{2}x_0^\top P(0)x_0 + \frac{1}{2} \int_0^T \text{tr}(P(t)B_\epsilon B_\epsilon^\top) dt. \end{aligned}$$

14 The first term is the same as that of the deterministic LQR problem. The
 15 second term is the penalty we incur for having a stochastic dynamical
 16 system. This is the minimal cost achievable for stochastic LQR but it is
 17 not the same as that of the deterministic LQR.

18 6.4 Linear Quadratic Gaussian (LQG)

19 Our development in the previous sections and the previous chapter was
 20 based on a Markov Decision Process, i.e., we know the state $x(t)$ at each
 21 instant in time t even if this state $x(t)$ changes stochastically. We said that
 22 the optimal control for the linear dynamics is still $u^*(t) = -K(t)x(t)$.

1 What should one do if we cannot observe the state exactly?

2 Imagine a “continuous-time” form the observation equation in the
3 Kalman filter where we receive observations of the form

$$y(t) = Cx(t) + D\nu.$$

4 where $\nu \sim N(0, I)$ is standard Gaussian noise that corrupts our observa-
5 tions y . If we extrapolate the definitions of the Kalman filter mean and
6 covariance to this continuous-time setting, we can write the KF as follows.
7 We know that the Kalman filter is the optimal estimate of the state given
8 all past observations, so it computes

$$\mu(t) = \mathbb{E}_{\epsilon(s), \nu(s): s \in [0, t]} [x(t) \mid y(s) : s \in [0, t]].$$

9 There exists a “continuous-time version” of the Kalman filter (which was
10 actually invented first), called the Kalman-Bucy filter. If the covariance of
11 the estimate is

$$\Sigma(t) = \mathbb{E}_{\epsilon(s), \nu(s): s \in [0, t]} \left[x(t) x(t)^\top \mid y(s) : s \in [0, t] \right],$$

12 the Kalman-Bucy filter updates $\mu(t), \Sigma(t)$ using the differential equation

$$\begin{aligned} \frac{d}{dt} \mu(t) &= Ax(t) + Bu(t) + K(t) (y(t) - C\mu(t)) \\ \frac{d}{dt} \Sigma(t) &= A\Sigma(t) + \Sigma(t)A^\top + B_\epsilon B_\epsilon^\top - K(t)DD^\top K(t)^\top \quad (6.12) \end{aligned}$$

where $K(t) = \Sigma(t) C^\top (DD^\top)^{-1}$.

13 This equation is very close to the Kalman filter equations you saw in
14 Chapter 3. In particular, notice the close similarity of the expression for
15 the Kalman gain $K(t)$ with the Kalman gain of the LQR problem. You
16 can read more at https://en.wikipedia.org/wiki/Kalman_filter.

Linear Quadratic Gaussian (LQG) It turns out that we can plug

i 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 σ -fields constructed using increasing subsets of the set $\{y(s) : s \in [0, t]\}$. Let us not worry about this here.

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.

1 **Certainty equivalence** For instance, even if we are using a particle
 2 filter to estimate the state of the system, we usually use the mean of the
 3 state estimate at time t given by $\mu(t)$ “as if” it were the true state of the
 4 system. Even if we were using some other feedback control $u(x)$ different
 5 than the LQR control (say feedback linearization), we usually plug in this
 6 estimate $\mu(t)$ in place of $x(t)$. Doing so is called “certainty equivalence”
 7 in control theory/robotics, which is a word borrowed from finance where
 8 one takes decisions (controls) directly using the estimate of the state (say
 9 stock price) while fully knowing the the stock price will change in the
 10 future stochastically.

11 **6.4.1 (Optional material) The duality between the Kalman** 12 **Filter and LQR**

13 We can re-write the covariance in (6.12) using the identity

$$\frac{d}{dt} (\Sigma(t)^{-1}) = \Sigma(t)^{-1} \dot{\Sigma}(t) \Sigma(t)^{-1}$$

14 to get

$$\dot{S} = C^\top (DD^\top)^{-1} C - A^\top S - SA - SB_w B_w^\top S \quad (6.14)$$

15 where we have defined $S := \Sigma^{-1}$.

16 Notice that the two equations, updates to the LQR cost matrix in (6.11)

$$-\dot{P} = PA + A^\top P + Q - PBR^{-1}B^\top P$$

17 look quite similar to this equation. In fact, they are identical and you can
 18 substitute the following.

LQR	Kalman-Bucy filter
P	Σ^{-1}
A	$-A$
$BR^{-1}B$	$B_w B_w^\top$
Q	$C^\top (DD^\top)^{-1} C$
t	$T - t$

Let us analyze this equivalence. Notice that the inverse of the Kalman filter covariance is like the cost matrix of LQR. This is conceptually easy to understand, our figure of merit for filtering is the inverse covariance matrix (smaller the better) and our figure of merit for the LQR problem is the cost matrix P (smaller the better). Similarly, smaller the LQR cost, better the controller. The “dynamics” of the Kalman filter is the reverse of the dynamics of the LQR problem, this shows that the P matrix is updated backwards in time while the covariance Σ is updated forwards in time. The next identity

$$BR^{-1}B^\top = B_w B_w^\top$$

is very interesting. Imagine a situation where we have a fully-actuated system with $B = I$ and B_w being a diagonal matrix. This identity suggests that larger the control cost R_{ii} of a particular actuator i , lower is the noise of using that actuator $(B_w)_{ii}$, and vice-versa. This is how muscles in your body have evolved: muscles that are cheap to use (low R) are also very noisy in what they do whereas muscles that are expensive to use (large R) which are typically the biggest muscles in the body are also the least noisy and most precise. You can read more about this in the paper titled “General duality between optimal 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

1 dynamical system

$$\dot{x} = f(x, u); x(0) = x_0 \text{ is given.}$$

2 We will consider a deterministic continuous-time dynamical system, the
3 modifications to following section that one would make if the system
4 is discrete-time, or stochastic, are straightforward and follow the same
5 strategy. First consider the problem where the run-time and terminal costs
6 are quadratic

$$\frac{1}{2}x(T)^\top Q_f x(T) + \frac{1}{2} \int_0^T x(t)^\top Q x(t) + u(t)^\top R u(t) dt.$$

7 **Receding horizon control and Model Predictive Control (MPC)** One
8 easy way to solve the dynamic programming problem, i.e., find a control
9 trajectory of the *nonlinear* system that minimizes this cost functional,
10 approximately, is by linearizing the system about the initial state x_0 and
11 some reference control u_0 (this can usually be zero). Let the linear system
12 be

$$\dot{z} = A_{x_0, u_0} z + B_{x_0, u_0} v; z(0) = 0; \quad (6.15)$$

13 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
14 Jacobians of the nonlinear function $f(x, u)$ with respect to the state and
15 control respectively. The state of the linearized dynamics is

$$z := x - x_0, \text{ and } v := u - u_0,$$

16 We have emphasized the fact that the matrices $A_{x_0, u_0}, B_{x_0, u_0}$ depend
17 upon the reference state and control using the subscript. Given the above
18 linear system, we can find a control sequence $u^*(\cdot)$ that minimizes the
19 cost functional using the standard LQR formulation. Notice now that even
20 we computed this control trajectory using the approximate linear system,
21 it can certainly be *executed* on the nonlinear system, i.e., at run-time we
22 will simply set $u \equiv u^*(z)$.

23 The linearized dynamics in (6.15) is potentially going to be very
24 different from the nonlinear system. The two are close in the neighborhood
25 of x_0 (and u_0) but as the system evolves using our control input to
26 move further away from x_0 , the linearized model no longer is a faithful
27 approximation of the nonlinear model. A reasonable way to fix matters
28 is to linearize about another point, say the state and control after $t = 1$
29 seconds, x_1, u_1 to get a new system

$$\dot{z} = A_{x_1, u_1} z + B_{x_1, u_1} v; z(0) = 0$$

30 and take the LQR-optimal control corresponding to this system for the
31 next second.

32 The above methodology is called “receding horizon control”. The
33 idea is that we compute the optimal control trajectory $u^*(\cdot)$ using an
34 approximation of the original system and recompute this control every few
35 seconds when our approximation is unlikely to be accurate. This is a very

1 popular technique to implement optimal controllers in typical applications.
 2 The concept of using an approximate model (almost invariably, a linear
 3 model with LQR cost) to plan for the near-term future and resolving the
 4 problem in receding horizon fashion once the system is at the end of this
 5 short time-horizon is called “Model Predictive Control”.

6 MPC is, perhaps, the second most common control algorithm im-
 7 plemented in the world. It is responsible for running most complex
 8 engineering systems that you can think of—power grids, oil refineries,
 9 chemical plants, rockets, aircrafts etc. Essentially, one never implements
 10 LQR directly, it is always implemented inside an MPC. For instance, in
 11 autonomous driving, the trajectory that the vehicle plans for traveling
 12 between two points A and B depends upon the current locations of the
 13 other cars/pedestrians in its vicinity, and potentially some prediction model
 14 of where they will be in the future. As the vehicle starts moving along
 15 this trajectory, the rest of the world evolves around it and we recompute
 16 the optimal trajectory to take into account the actual locations of the
 17 cars/pedestrians in the future.

🔗 Can you guess what is *the* most common control algorithm in the world?

18 6.5.1 Iterative LQR (iLQR)

19 Now let us consider the situation when in addition to a nonlinear system,

$$\dot{x} = f(x, u); x(0) = x_0,$$

20 the run-time and terminal cost is also nonlinear

$$q_f(x(T)) + \int_0^T q(x(t), u(t)) dt.$$

21 We can solve the dynamic programming problem in this case approximately
 22 using the following iterative algorithm.

23 Assume that we are given an initial control trajectory $u^{(0)}(\cdot) =$
 24 $\{u^{(0)}(t) : t \in [0, T]\}$. Let $x^{(0)}(\cdot)$ be the state trajectory that corresponds
 25 to taking this control on the nonlinear system, with of course $x^{(0)}(0) = x_0$.
 26 At each iteration k , the Iterative LQR algorithm performs the following
 27 steps.

28 **Step 1** Linearize the nonlinear system about the state trajectory $x^{(k)}(\cdot)$
 29 and $u^{(k)}(\cdot)$ using

$$z(t) := x(t) - x^{(k)}(t), \text{ and } v(t) := u(t) - u^{(k)}(t)$$

30 to get a new system

$$\dot{z} = A^{(k)}(t)z + B^{(k)}(t)v; z(0) = 0$$

31 where

$$A^{(k)}(t) = \left. \frac{df}{dx} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}$$

$$B^{(k)}(t) = \left. \frac{df}{du} \right|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}$$

1 and compute the Taylor series approximation of the nonlinear cost up to
2 the second order

$$q_f(x(T)) \approx \text{constant} + z(T)^\top \frac{dq_f}{dx} \Big|_{x(T)=x^{(k)}(T)} \\ + z(t)^\top \frac{d^2q_f}{dx^2} \Big|_{x(T)=x^{(k)}(T)} z(t),$$

3

$$q(x, u, t) \approx \text{constant} + z(t)^\top \underbrace{\frac{dq}{dx} \Big|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\text{affine term}} \\ + v(t)^\top \underbrace{\frac{dq}{du} \Big|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\text{affine term}} \\ + z(t)^\top \underbrace{\frac{d^2q}{dx^2} \Big|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\equiv Q} z(t) \\ + v(t)^\top \underbrace{\frac{d^2q}{du^2} \Big|_{x(t)=x^{(k)}(t), u(t)=u^{(k)}(t)}}_{\equiv R} v(t).$$

4 This is an LQR problem with run-time cost that depends on time (like our
5 discrete-time LQR formulation, the continuous-time formulation simply
6 has Q, R to be functions of time t in the Riccati equation) and which also
7 has terms that are affine in the state and control in addition to the usual
8 quadratic cost terms.

9 **Step 2** Solve the above linearized problem using standard LQR formula-
10 tion to get the new control trajectory

$$u^{(k+1)}(t) := u^{(k)}(t) - Kz(t).$$

11 Simulate the *nonlinear* system using the control $u^{(k+1)}(\cdot)$ to get the new
12 state trajectory $x^{(k+1)}(\cdot)$.

13 Some important comments to remember about the iLQR algorithm.

14 1. There are many ways to pick the initial control trajectory $u^{(0)}(\cdot)$, e.g.,
15 using a spline to get an arbitrary control sequence, using a spline
16 to interpolate the states to get a trajectory $x^{(0)}(\cdot)$ and then back-
17 calculate the control trajectory, using the LQR solution based on the
18 linearization about the initial state, feedback linearization/differen-
19 tial flatness (https://en.wikipedia.org/wiki/Feedback_linearization)
20 etc.

21 2. The iLQR algorithm is an approximate solution to dynamic pro-
22 gramming for nonlinear system with general, nonlinear run-time and
23 terminal costs. This is because the the algorithm uses a linearization
24 about the previous state and control trajectory to compute the new

❓ 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 control trajectory. iLQR is not guaranteed to find the optimal
2 solution of dynamic programming, although in practice with good
3 implementations, it works excellently.

- 4 3. We can think of iLQR as an algorithm to track a given state trajectory
5 $x^g(t)$ by setting

$$q_f = 0, \text{ and } q(x, u) = \|x^g(t) - x(t)\|^2.$$

6 This is often how iLQR is typically used in practice, e.g., to make
7 an autonomous race car closely follow the racing line (see the paper
8 “BayesRace: Learning to race autonomously using prior experience”
9 <https://arxiv.org/abs/2005.04755> and <https://www.youtube.com/watch?v=dgIpf0Lg8Ek>
10 for a clever application of using MPC to track a challenging race
11 line), or to make a drone follow a given desired trajectory
12 (<https://www.youtube.com/watch?v=QREeZvHg0IQ>).

13 **Differential Dynamic Programming (DDP)** is a suite of techniques
14 that is a more powerful version of iterated LQR. Instead of linearizing
15 the dynamics and taking a second order Taylor approximation of the cost,
16 DDP takes a second order approximation of the Bellman equation directly.
17 The two are not the same; DDP is the more correct version of iLQR but is
18 much more challenging computationally.

19 Broadly speaking, iLQR and DDP are used to perform control for some
20 of the most sophisticated robots today, you can see an interesting discussion
21 of the trajectory planning of some of the DARPA Humanoid Robotics
22 Challenge at <https://www.cs.cmu.edu/~cga/drc/atlas-control>. Techniques
23 like feedback linearization work excellently for drones where we do not
24 really care for optimal cost (see “Minimum snap trajectory generation and
25 control for quadrotors” <https://ieeexplore.ieee.org/document/5980409>)
26 while LQR and its variants are 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
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

$$J_\varphi(x, t) = \frac{1}{2}x(t)^\top \underbrace{(\text{some function of } A, B, Q, R)}_{\text{function of } \varphi} x(t)$$

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

1 of parameters φ . We know that any cost-to-go that satisfies the
2 Bellman equation is the optimal cost-to-go, so we can now “fit”
3 the candidate function $J_\varphi(x, t)$ to satisfy the Bellman equation.
4 Similarly, one may also express the optimal feedback control $u(\cdot)$
5 using some parameters θ as

$$u_\theta(\cdot).$$

6 We will see how to fit such functions in this chapter.

7 **2. Learning from data.** It may happen that we do not know very
8 much about the dynamical system, e.g., we do not know a good
9 model for what drives customers as they buy items in an online
10 merchandise platform, or a robot traveling in a crowded area may
11 not have a good model for how large crowds of people walk around
12 it. One may collect data from these systems fit some model of the
13 form $\dot{x} = f(x, u)$ to the data and then go back to the techniques of
14 Module 2. It is typically not clear how much data one should collect.
15 RL gives a suite of techniques to learn the cost-to-go in these
16 situations by collecting and assimilating the data *automatically*.
17 These techniques go under the umbrella of policy gradients, on-
18 policy methods etc. One may also simply “memorize” the data
19 provided by an expert operator, this is called Imitation Learning
20 and we will discuss it next.

21 **Some motivation** Imitation Learning is also called “learning from
22 demonstrations”. This is in fact one of the earliest successful examples of
23 using a neural network for driving. The ALVINN project at CMU by Dean
24 Pomerleau in 1988 (<https://www.youtube.com/watch?v=2KMAAmkz9go>)
25 used a two-layer neural network with 5 hidden neurons, about 1000 inputs
26 from the pixels of a camera and 30 outputs. It successfully drove in
27 different parts of the United States and Germany. Imitation learning has
28 also been responsible for numerous other early-successes of RL, e.g.,
29 acrobatic maneuvers on an RC helicopter ([http://ai.stanford.edu/~acoates/-
30 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_{train} consists of $n = 50$ RGB images of size 100×100 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

🔗 How many such binary classifiers are there at most?

1 images to their labels (orange: +1 or no orange: -1). If x_{ij}^k is the $(ij)^{\text{th}}$
 2 pixel for image x^k , then we use the function

$$f(x) = \begin{cases} y^k & \text{if } x_{ij}^k = x_{ij} \text{ for some } k = 1, \dots, n \\ -1 & \text{otherwise.} \end{cases}$$

3 This predictor certainly solves the task. It correctly works for all images
 4 in the training set. Does it work for images outside the training set?

5 Our task in machine learning is to learn a predictor that works *outside*
 6 the training set. The training set is only a source of information that Nature
 7 gives us to find such a predictor.

Designing a predictor that is accurate on D_{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_{train} .

8 **Generalization** If we never see data from outside D_{train} why should we
 9 hope to do well on it? The key is the distribution P . Machine learning is
 10 formalized as constructing a predictor that works well on new data that is
 11 also drawn independently from the distribution P . We will call this set of
 12 data the “test set”

$$D_{\text{test}}$$

13 and it is constructed similarly. This assumption is important. It provides
 14 coherence between past and future samples: past samples that were used
 15 to train and future samples that we will wish to predict upon. How to find
 16 such predictors that work well on new data? The central idea in machine
 17 learning is 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.

18 The *right* class of functions f cannot be too large, otherwise we will
 19 find our binary classifier above as the solution, and that is not very useful.
 20 The class of functions cannot be too small either, otherwise we won't be
 21 able to predict difficult images. If the predictor does not even work well
 22 on the training set, there is no reason why we should expect it to work on
 23 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?

7.1.1 Fitting a machine learning model

Let us now solve a classification problem. We will again go around the model selection problem and consider the class of linear classifiers. Assume binary labels $Y \in \{-1, 1\}$. To keep the notation clear, we will use the trick of appending a 1 to the data x and hide the bias term b in the linear classifier. The predictor is now given by

$$\begin{aligned} f(x; w) &= \text{sign}(w^\top x) \\ &= \begin{cases} +1 & \text{if } w^\top x \geq 0 \\ -1 & \text{else.} \end{cases} \end{aligned} \quad (7.1)$$

We have used the sign function denoted as sign to get binary $\{-1, +1\}$ outputs from our real-valued prediction $w^\top x$. This is the famous perceptron model of Frank Rosenblatt.

We want the predictions of the model to match those in the training data and devise an objective to fit/train the perceptron.

$$\ell_{\text{zero-one}}(w) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{y^i \neq f(x^i; w)\}}. \quad (7.2)$$

The indicator function inside the summation measures the number of mistakes the perceptron makes on the training dataset. The objective here is designed to find weights w that minimizes the average number of mistakes, also known as the training error. Such a loss that measures the mistakes is called the zero-one loss, it incurs a penalty of 1 for a mistake and zero otherwise.

Surrogate losses The zero-one loss is the clearest indication of whether the perceptron is working well. It is however non-differentiable, so we cannot use powerful ideas from optimization theory to minimize it. This is why surrogate losses are constructed in machine learning. These are proxies for the loss function, typically for the classification problems and look as follows. The exponential loss is

$$\ell_{\text{exp}}(w) = e^{-y (w^\top x)}$$

or the logistic loss is

$$\ell_{\text{logistic}}(w) = \log(1 + e^{-y w^\top x}).$$

Stochastic Gradient Descent (SGD) SGD is a very general algorithm to optimize objectives typically found in machine learning. We can use it so long as we have a dataset and an objective that is differentiable. Consider an optimization problem where we want to solve for

$$w^* = \underset{w}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \ell^i(w)$$

🔗 Can you think of some quantity other than the zero-one error that we may wish to optimize?

1 where the function ℓ^i denotes the loss on the sample (x^i, y^i) and $w \in \mathbb{R}^p$
 2 denotes the weights of the classifier. Solving this problem using SGD
 3 corresponds to iteratively updating the weights using

$$w^{(t+1)} = w^{(t)} - \eta \frac{d\ell^{\omega_t}(w)}{dw} \Big|_{w=w^{(t)}},$$

4 i.e., we compute the gradient one sample with index ω_t in the dataset. The
 5 index ω_t is chosen uniformly randomly from

$$\omega_t \in \{1, \dots, n\}.$$

6 In practice, at each time-step t , we typically select a few (not just one) input
 7 data ω_t from the training dataset and average the gradient $\frac{d\ell^{\omega_t}(w)}{dw} \Big|_{w=w^{(t)}}$
 8 across them; this is known as a “mini-batch”. The gradient of the loss
 9 $\ell^{\omega_t}(w)$ with respect to w is denoted by

$$\nabla \ell^{\omega_t}(w^{(t)}) := \frac{d\ell^{\omega_t}(w)}{dw} \Big|_{w=w^{(t)}} = \begin{bmatrix} \nabla_{w_1} \ell^{\omega_t}(w^{(t)}) \\ \nabla_{w_2} \ell^{\omega_t}(w^{(t)}) \\ \vdots \\ \nabla_{w_p} \ell^{\omega_t}(w^{(t)}) \end{bmatrix} \in \mathbb{R}^p.$$

10 The gradient $\nabla \ell^{\omega_t}(w^{(t)})$ is therefore a vector in \mathbb{R}^p . We have written

$$\nabla_{w_1} \ell^{\omega_t}(w^{(t)}) = \frac{d\ell^{\omega_t}(w)}{dw_1} \Big|_{w=w^{(t)}}$$

11 for the scalar-valued derivative of the objective $\ell^{\omega_t}(w^{(t)})$ with respect to
 12 the first weight $w_1 \in \mathbb{R}$. We can therefore write SGD as

$$w^{(t+1)} = w^{(t)} - \eta \nabla \ell^{\omega_t}(w^{(t)}). \quad (7.3)$$

13 The non-negative scalar $\eta \in \mathbb{R}_+$ is called the step-size or the learning rate.
 14 It governs the distance traveled along the negative gradient $-\nabla \ell^{\omega_t}(w^{(t)})$
 15 at each iteration.

16 7.1.2 Deep Neural Networks

17 The Perceptron in (7.1) is a linear model: it computes a linear function
 18 of the weights $w^\top x$ and uses this function to make the predictions
 19 $f(x; w) = \text{sign}(w^\top x)$. Linear models try to split the data (say we have
 20 binary labels $Y = \{-1, 1\}$) using a hyper-plane with w denoting the
 21 normal to this hyper-plane. This does not work for all situations of course,
 22 as the figure below shows, there is no hyper-plane that cleanly separates
 23 the two classes (i.e., achieves zero mis-prediction error) but there *is* a
 24 nonlinear function that can do the job.

25 A deep neural network is one such nonlinear function. First consider
 26 a “two-layer” network

$$f(x; v, S) = \text{sign}(v^\top \sigma(S^\top x))$$

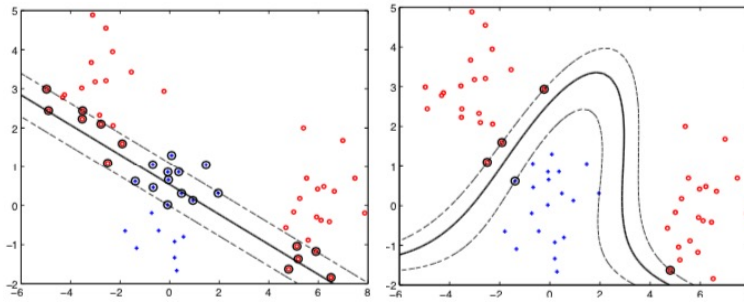


Figure 7.1

1 where the matrix $S \in \mathbb{R}^{d \times p}$ and a vector $v \in \mathbb{R}^p$ are the parameters or
 2 “weights” of the classifier. The “nonlinearity” σ is usually set to be what
 3 is called a Rectified Linear Unit (ReLU)

$$\begin{aligned} \sigma(x) &:= \text{ReLU}(x) = |x|_+ \\ &= \max(0, x). \end{aligned} \quad (7.4)$$

4 Just like the case of a Perceptron, we can use an objective $\frac{1}{n} \sum_{i=1}^n \ell^i(v, S)$
 5 that depends on both v, S to fit this classifier on training data. A deep
 6 neural network takes the idea of a two-layer network to the next step and
 7 has multiple “layers”, each with a different weight matrix S_1, \dots, S_L .
 8 The classifier is therefore given by

$$f(x; v, S_1, \dots, S_L) = \text{sign}(v^\top \sigma(S_L^\top \dots \sigma(S_2^\top \sigma(S_1^\top x)) \dots)). \quad (7.5)$$

9 We call each operation of the form $\sigma(S_k^\top \dots)$, as a *layer*. Consider the
 10 second layer: it takes the features generated by the first layer, namely
 11 $\sigma(S_1^\top x)$, multiplies these features using its feature matrix S_2^\top and applies
 12 a nonlinear function $\sigma(\cdot)$ to this result element-wise before passing it on
 13 to the third layer.

A deep network creates new features by composing older features.

14 This composition is very powerful. Not only do we not have to
 15 pick a particular feature vector, we can create very complex features by
 16 sequentially combining simpler ones. For example Figure 7.2 shows the
 17 features (more precisely, the kernel) learnt by a deep neural network.
 18 The first layer of features are called Gabor-like, and incidentally they are
 19 similar to the features learned by the human brain in the first part of the
 20 visual cortex (the one closest to the eyes). These features are *combined*
 21 linearly along with a nonlinear operation to give richer features (spirals,
 22 right angles) in the middle panel. The third layer combines the lower
 23 features to get even more complex features, these look like patterns (notice
 24 a soccer ball in the bottom left), a box on the bottom right etc.

1 **Deep networks are universal function approximators** The multi-layer
 2 neural network is a powerful class of classifiers: depending upon how many
 3 layers we have and what is the dimensionality of the the weight matrices
 4 S_k at each layer, we can fit *any* training data. In fact, this statement,
 5 which is called the *universal approximation property* holds even for a
 6 two-layer neural network $v^\top \sigma(S^\top x)$ if the number of columns in S is big
 7 enough. This property is the central reason why deep networks are so
 8 widely applicable, we can model complex machine learning problems if
 9 we choose a big enough deep network.

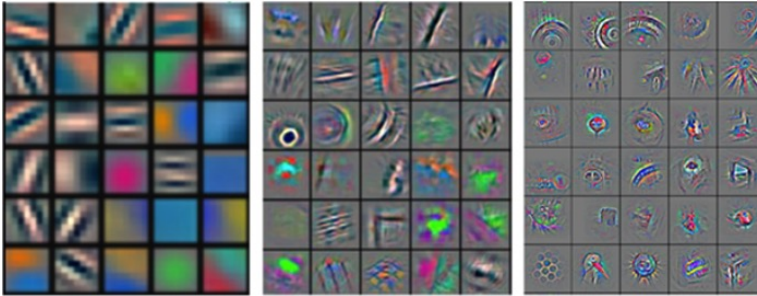


Figure 7.2

10 **Logits for multi-class classification.** The output

$$\hat{y} = v^\top \sigma(S_L^\top \dots \sigma(S_2^\top \sigma(S_1^\top x)) \dots)$$

11 is called the logits corresponding to the different classes. This name
 12 comes from logistic regression where logits are the log-probabilities of an
 13 input datum belonging to one of the two classes. A deep network provides
 14 an easy way to solve a multi-class classification problem, we simply set

$$v \in \mathbb{R}^{p \times C}$$

15 where C is the total number of classes in the data. Just like logistic
 16 regression predicts the logits of the two classes, we would like to *interpret*
 17 the vector \hat{y} as the log-probabilities of an input belonging to one of the
 18 classes.

19 **Weights** It is customary to not differentiate between the parameters of
 20 different layers of a deep network and simply say *weights* when we want
 21 to refer to all parameters. The set

$$w := \{v, S_1, S_2, \dots, S_L\}$$

22 is the set of *weights*. This set is typically stored in PyTorch as a set of
 23 matrices, one for each layer. Using this new notation, we will write down
 24 a deep neural network classifier as simply

$$f(x, w) \tag{7.6}$$

🔗 What would the shape of w be if you were performing regression using a deep network?

1 and fitting the deep network to a dataset involves the optimization problem

$$w^* = \operatorname{argmin}_w \frac{1}{n} \sum_{i=1}^n \ell(y^i, \hat{y}^i). \quad (7.7)$$

2 We will also sometimes denote the loss of the i^{th} sample as

$$\ell^i(w) := \ell(y^i, \hat{y}^i).$$

3 **Backpropagation** The Backpropagation algorithm is a method to com-
4 pute the gradient of the objective while fitting a deep network using SGD,
5 i.e., it computes $\nabla_w \ell^i(w)$. For the purposes of this course, the details of
6 how this is done are not essential, so we will skip them. You can read more
7 in the notes of ESE 546 at https://pratikac.github.io/pub/20_ese546.pdf.

8 **PyTorch** We will use a library called PyTorch (<https://pytorch.org>) to
9 code up deep neural networks for the reinforcement learning part of this
10 course. You can find some excellent tutorials for it at
11 <https://pytorch.org/tutorials/beginner/basics/intro.html>. We have also
12 uploaded two recitations from the Fall 2020 offering of ESE 546 on
13 Canvas which guide you through various typical use-cases of PyTorch.
14 You are advised to go through, at least, the first recitation if you are
15 not familiar with PyTorch. For the purposes of this course, you do not
16 need to know the intricacies of PyTorch, we will give you enough code
17 to work with deep networks so that you can focus on implementing the
18 reinforcement learning-specific parts.

19 7.2 Behavior Cloning

20 With that background, we are ready to tackle what is potentially the simplest
21 problem in RL. We will almost exclusively deal with discrete-time systems
22 for RL. Let us imagine that we are given access to n trajectories each of
23 length $T + 1$ time-steps from an expert demonstrator for our system. We
24 write this as a training dataset

$$D = \{(x_t^i, u_t^i)_{t=0,1,\dots,T}\}_{i=1,\dots,n}$$

25 At each step, we record the state $x_t^i \in \mathbb{R}^d$ and the control that the expert
26 took at that state u_t^i . We would like to learn a deterministic feedback
27 control for the robot that is parametrized by parameters θ

$$u_\theta(x) : X \mapsto U \subset \mathbb{R}^m.$$

28 using the training data. The idea is that if $u_\theta(x^i(t)) \approx u^i(t)$ for all i
29 and all times t , then we can simply run our learned controller $u_\theta(x)$ on
30 the robot instead of having the expert. A simple example is a baby deer
31 learning to imitate how its mother in how to run.

1 **Parameterizing the controller** Our function u_θ may represent many
 2 different families of controllers. For example, $u_\theta(x) = \theta x$ where $\theta \in$
 3 $\mathbb{R}^{d \times p}$ is a linear controller; this is much like the control for LQR except
 4 that we can fit θ to the expert's data instead of solving the LQR problem
 5 to find the Kalman gain. We could also think of some other complicated
 6 function, e.g., a two-layer neural network,

$$u_\theta(x) = v \sigma(S^\top x)$$

7 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,
 8 say ReLU. As we did above, we will use

$$\theta := (v, S)$$

9 to denote all the weights of this two-layer neural network. Multi-layer
 10 neural networks are also another possible avenue. In general, we want
 11 to the parameterization of the controller to be rich enough to fit some
 12 complex controller that the expert may have used on the system.

13 **How to fit the controller?** Given our chosen model for $u_\theta(x)$, say a
 14 two-layer neural network with weights θ , fitting the controller involves
 15 finding the best value for the parameters θ such that $u_\theta(x_t^i) \approx u_t^i$ for data
 16 in our dataset. There are many ways to do this, e.g., we can solve the
 17 following optimization problem

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \ell(\theta) := \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{T+1} \sum_{t=0}^T \|u_t^i - u_\theta(x_t^i)\|_2^2}_{\ell^i(\theta)} \quad (7.8)$$

18 The difficulty of solving the above problem depends upon how difficult the
 19 model $u_\theta(x)$ is, for instance, if the model is linear θx , we can solve (7.8)
 20 using ordinary least squares. If the model is a neural network, one would
 21 have to use SGD to solve the optimization problem above. After fitting
 22 this model, we have a new controller

$$u_{\hat{\theta}}(x) \in \mathbb{R}^m$$

23 that we can use *anywhere* in the domain $X \subset \mathbb{R}^d$, even at places where
 24 we had no expert data. This is known as Behavior Cloning, i.e., cloning
 25 the controls of the expert into a parametric model.

26 **Generalization performance of behavior cloning** Note that the data
 27 provided by the expert is not iid, of course the state x_{t+1}^i in the expert's
 28 trajectory depends upon the previous state x_t^i . Standard supervised
 29 learning makes the assumption that Nature gives training data that is
 30 independent and identically distributed from the distribution P . While
 31 it is still reasonable to fit the regression loss in (7.8) for such correlated
 32 data, one should remember that if the expert trajectories do not go to all
 33 parts of the state-space, the learned controller fitted on the training data

1 may not work outside these parts. Of course, if we behavior clone the
 2 controls taken by a generic driver, they are unlikely to be competitive for
 3 racing, and vice-versa. It is very important to realize that this does *not*
 4 mean that BC does not generalize. Generalization in machine learning is
 5 a concept that suggests that the model should work well on data *from the*
 6 *same distribution*. What does the the “distribution” of the expert mean, in
 7 this case, it simply refers to the distribution of the states that the expert’s
 8 trajectories typically visit, e.g. a race driver typically drives at the limits
 9 of tire friction and throttle, this is different from a usual city-driver who
 10 would rather maximize the longevity of their tires and engine-life.

❗ Discuss generalization performance in behavior cloning.

11 7.2.1 Behavior cloning with a stochastic controller

12 In this case, we have always chosen feedback feedback controllers that
 13 are deterministic, i.e., there is a single value of control u that is taken at
 14 the state x . Going forward, we will also talk about stochastic controllers,
 15 i.e., controllers which sample a control from a distribution. There can
 16 be a few reasons of using such a controller. First, we will see in later
 17 lectures how this may help in training a reinforcement learning algorithm;
 18 this is because in situations where you do not know the system dynamics
 19 precisely, it helps to “hedge” the feedback to take a few different control
 20 actions instead of simply the one that the value function deems as the
 21 maximizing one. This is not very different from having a few different
 22 stocks in your portfolio. Second, we benefit from this hedging even at
 23 test-time when we run a stochastic feedback control, e.g., in situations
 24 where the limited training data may not want to always pick the best
 25 control (because the best control was computed using an imprecise model
 26 of the system dynamics and could be wrong), but rather hedge our bets by
 27 choosing between a few different controls.

28 A stochastic feedback control is denoted by

$$u \sim u_\theta(\cdot | x) = \mathbf{P}(\cdot | x)$$

29 notice that $u_\theta(\cdot | x)$ is a probability distribution on the control space U
 30 that depends on the state x , and in this case the parameters θ . The control
 31 taken at a state x is a sample drawn from this probability distribution. The
 32 deterministic controller is a special case of this setup where

$$u_\theta(u | x) = \delta_{u_\theta(x)}(u) \equiv u_\theta(x)$$

33 is a Dirac-delta distribution at $u_\theta(x)$. If the control space U is discrete,
 34 then $u_\theta(\cdot | x)$ could be a categorical distribution. If the control space U
 35 is continuous, then you may wish to think of the controls being sampled
 36 from a Gaussian distribution with some mean $\mu_\theta(x)$ and variance $\sigma_\theta^2(x)$

$$\mathbb{R}^m \ni u \sim u_\theta(\cdot | x) = N(\mu_\theta(x), \Sigma_\theta(x)).$$

1 **Maximum likelihood estimation** Let's pick a particular stochastic
 2 controller, say a Gaussian. How should we fit the parameters θ for this?
 3 We would like to find parameters θ that make the expert's data in our
 4 dataset very likely. The log-likelihood of each datum is

$$\log u_\theta(u_t^i | x_t^i)$$

5 and maximizing the log-likelihood of the entire dataset amounts to solving

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{T+1} \sum_{t=0}^T -\log u_\theta(u_t^i | x_t^i)}_{\ell^i(\theta)}. \quad (7.9)$$

6 **Fitting BC with a Gaussian controller** Notice that if we use a Gaussian
 7 distribution

$$u_\theta(\cdot | x) = N(\mu_\theta(x), I)$$

8 as our stochastic controller, the objective in (7.9) is the same as that
 9 in (7.8).

$$u_\theta(\cdot | x) = N(\mu_\theta(x), \sigma_\theta^2(x)I)$$

10 we have that

$$-\log u_\theta(u | x) = \frac{\|\mu_\theta(x) - u\|_2^2}{\sigma_\theta^2(x)} + 2cp \log \sigma_\theta(x).$$

11 where c is a constant.

12 7.2.2 KL-divergence form of Behavior Cloning

13 **Background on KL divergence** The Kullback-Leibler (KL) divergence
 14 is a quantity to measure the distance between two probability distributions.
 15 There are many similar distances, for example, given two probability
 16 distributions $p(x)$ and $q(x)$ supported on a discrete set X , the total
 17 variation distance between them is

$$\operatorname{TV}(p, q) = \frac{1}{2} \sum_{x \in X} |p(x) - q(x)|.$$

18 Hellinger distance (https://en.wikipedia.org/wiki/Hellinger_distance), f -
 19 divergences (<https://en.wikipedia.org/wiki/F-divergence>) and the Wasser-
 20 stein metric
 21 (https://en.wikipedia.org/wiki/Wasserstein_metric) are a few other exam-
 22 ples of ways to measure how different two probability distributions are
 23 from each other.

24 The Kullback-Leibler divergence (KL) between two distributions is
 25 given by

$$\operatorname{KL}(p || q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}. \quad (7.10)$$

26 This is a distance and not a metric, i.e., it is always non-zero and zero

1 if and only if the two distributions are equal, but the KL-divergence
 2 is not symmetric (like a metric has to be). Also, the above formula is
 3 well-defined only if for all x where $q(x) = 0$, we also have $p(x) = 0$.
 4 Notice that it is not symmetric

$$\text{KL}(q \parallel p) = \sum_{x \in X} q(x) \log \frac{q(x)}{p(x)} \neq \text{KL}(p \parallel q).$$

5 The funny notation $\text{KL}(p \parallel q)$ was invented by Shun-ichi Amari
 6 (https://en.wikipedia.org/wiki/Shun%27ichi_Amari) to emphasize the fact
 7 that the KL-divergence is asymmetric. The KL-divergence is always
 8 positive: you can show this using an application of Jensen's inequality.
 9 For distributions with continuous support, we integrate over the entire
 10 space X and define KL divergence as

$$\text{KL}(p \parallel q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx.$$

11 **Behavior Cloning** Let us now imagine the expert is also a parametric
 12 stochastic feedback controller $u_{\theta^*}(\cdot \mid x)$. Our data is therefore drawn by
 13 running this controller for n trajectories, T time-steps on the system. This
 14 dataset now consists of samples from

$$p_{u_{\theta^*}}(x, u)$$

15 which is the joint distribution on the state-space X and the control-space U .
 16 We have denoted the parameters of the feedback controller which creates
 17 this distribution as the subscript u_{θ^*} . Our behavior cloning controller
 18 creates a similar distribution $p_{u_{\theta}}(x, u)$ and the general version of the
 19 objective in (7.9) is therefore

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \text{KL}(p_{u_{\theta^*}} \parallel p_{u_{\theta}}); \quad (7.11)$$

20 The objective in (7.9) corresponds to this for Gaussian stochastic con-
 21 trollers, but we can just as easily imagine some other distribution for the
 22 stochastic 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.

23 7.2.3 Some remarks on Behavior Cloning

24 **Worst-case performance** Performance of Behavior Cloning can be
 25 quite bad in the worst case. The authors in “Efficient reductions for
 26 imitation learning” ([https://www.cs.cmu.edu/~ross1/publications/Ross-](https://www.cs.cmu.edu/~ross1/publications/Ross-AIStats11-NoRegret.pdf)
 27 [AIStats11-NoRegret.pdf](https://www.cs.cmu.edu/~ross1/publications/Ross-AIStats11-NoRegret.pdf)) show that if the learned controller $u_{\hat{\theta}}$ differs
 28 from the control taken by the expert controller u_{θ^*} with a probability ϵ at

1 each time-step, over a horizon of length T time-steps, it can be $\mathcal{O}(T^2\epsilon)$ off
 2 from the cost-to-go of the expert *as averaged over states that the learned*
 3 *controller visits*. This is because once the robot makes a mistake and goes
 4 away from the expert's part in the state-space, future states of the robot
 5 and the expert can be very different.

6 **Model-free nature of BC** Observe that our learned controller $u_{\hat{\theta}}(\cdot | x)$
 7 is a feedback controller and works for entire state-space X . We did not
 8 need to know the dynamics of the system to build this controller. The
 9 data from the expert is conceptually the same as the model $\dot{x} = f(x, u)$ of
 10 the dynamics, and you can learn controllers from both. Do you however
 11 notice a catch?

12 7.3 DAgger: Dataset Aggregation

13 The expert's dataset in Behavior Cloning determines the quality of the
 14 controller learned. If we collected very few trajectories from the expert,
 15 they may not cover all parts of the state-space and the behavior cloned
 16 controller has no data to fit the model in those parts.

17 Let us design a simple algorithm, of the same spirit as iterative-LQR,
 18 to mitigate this. We start with a candidate controller, say $u_{\theta^{(0)}}(x)$; one
 19 may also 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_{θ^*} , then the robot controller at a state x is

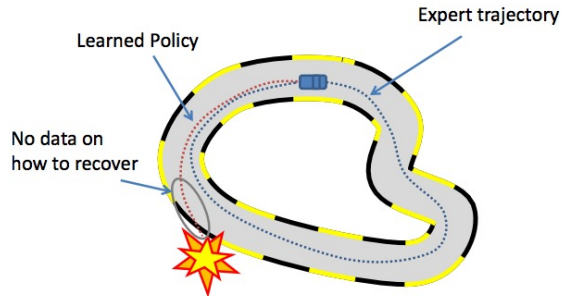
$$u \sim p \delta_{u_{\theta^*}(x)} + (1 - p) \delta_{u_{\theta^{(k-1)}}(x)}.$$

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 $D^{(k)}$.

20 The above algorithm iteratively updates the BC controller $u_{\hat{\theta}}$ by
 21 drawing new data from the expert. The robot first bootstraps off the
 22 expert's data, this simply means that it uses the expert's data to fit its
 23 controller $u_{\theta^{(0)}}(x)$. As we discussed above, this controller may veer off
 24 the expert's trajectory if the robot starts at states that are different from

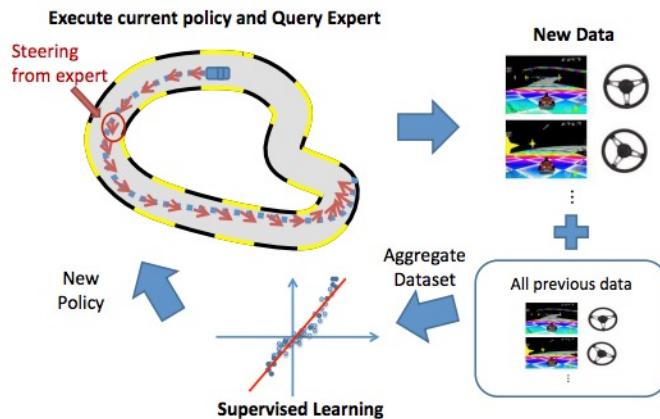
📌 Draw a picture of the amplifying errors of running behavior cloning in real-time.

1 the dataset, or even if it takes a slightly different control than the expert
2 midway through a trajectory.



3

4 To fix this, the robot collects more data at each iteration. It uses a
5 combination of the expert and its controller to collect such data. This,
6 allows *collecting a dataset of expert's controls in states that the robot*
7 *visits* and iteratively expands the dataset $D^{(k)}$.



8

9 In the beginning we may wish to be close to the expert's data and use
10 a large value of p , as the fitted controller $u_{\theta_{k+1}}$ becomes good, we can
11 reduce the value of p and rely less on the expert.

12 DAgger is an iterative algorithm which expands the controller to handle
13 larger and larger parts of the state-space. Therefore, the cost-to-go of
14 the controller learned via DAgger is $\mathcal{O}(T)$ off from the cost-to-go of the
15 expert *as averaged over states that the learned controller visits*.

16 **DAgger with expert annotations at each step** DAgger is a conceptual
17 framework where the expert is queried repeatedly for new control actions.
18 This is obviously problematic because we need to expert on hand at each
19 iteration. We can also cook up a slightly version of DAgger where we
20 start with the BC controller $u_{\theta^{(k)}} = u_{\hat{\theta}}$ and at each step, we run the
21 controller on the real system and ask the expert to relabel the data after

🔗 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 $\theta^{(k+1)} \approx \theta^{(k)}$.

1 that run. The dataset $D^{(k)}$ collected by the algorithm expands at each
2 iteration and although the states x_t^i are those visited by our controller, their
3 annotations are those given by the expert. This is a much more natural
4 way of implementing DAgger.

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/mniha16.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_{θ} 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 controller with parameters (weights) θ . We denote them as follows

$$\begin{aligned} x_{k+1} &\sim p(\cdot \mid x_k, u_k) \text{ with noise denoted by } \epsilon_k \\ u_k &\sim u_\theta(x_k). \end{aligned}$$

We will also change perspective and instead of minimizing the infinite-horizon sum of a runtime cost, maximize the sum of a runtime reward

$$r(x, u) := -q(x, u).$$

We do so simply to conform to tradition and standard notation in reinforcement learning; the two are mathematically completely equivalent. We are interested in maximizing the expected value of the cumulative rewards over infinite-horizon trajectories of the system

$$J(\theta; x_0) = \mathbb{E}_{x_1, x_2, \dots} \left[\underbrace{\sum_{k=0}^{\infty} \gamma^k r(x_k, u_k)}_{\text{discounted return}} \mid x_0 \right]; \quad (8.1)$$

where each $u_k \sim u_\theta(\cdot \mid x_k)$ and each $x_{k+1} \sim p(\cdot \mid x_k, u_k)$.

Trajectory space Let us write out one trajectory of such a system a bit more explicitly. We know that the probability of the next state x_{k+1} given x_k is $p(x_{k+1} \mid x_k, u_k)$. The probability of taking a control u_k at state x_k is $u_\theta(u_k \mid x_k)$. We denote an infinite trajectory by

$$\tau = x_0, u_0, x_1, u_1, \dots$$

The probability of this entire trajectory occurring is

$$p_\theta(\tau) = \prod_{k=0}^{\infty} p(x_{k+1} \mid x_k, u_k) u_\theta(u_k \mid x_k);$$

we have emphasized that the distribution of trajectories depends on the weights of controller θ . If we take the logarithm,

$$\log p_\theta(\tau) = \sum_{k=0}^{\infty} \log p(x_{k+1} \mid x_k, u_k) + \log u_\theta(u_k \mid x_k).$$

Given a trajectory $\tau = x_0, u_0, x_1, u_1, \dots$, the sum

$$R(\tau) = \sum_{k=0}^{\infty} \gamma^k r(x_k, u_k) \quad (8.2)$$

1 is called the discounted return of the trajectory τ . Sometimes we will
 2 also talk of the undiscounted return of the trajectory which is the sum of
 3 the rewards up to some fixed finite horizon T without the discount factor
 4 pre-multiplier. Using this notation, we can write out objective from (8.1)
 5 as

$$J(\theta; x_0) = \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) | x_0] \quad (8.3)$$

6 where $p(\tau)$ is the probability distribution of an infinitely long trajectory τ .
 7 Observe what is probably the most important point in policy-gradient
 8 based reinforcement learning: the probability of trajectory is an infinite
 9 product of terms. All terms are smaller than 1 (they are probabilities),
 10 so it is essentially zero even if the state-space and the control-space are
 11 finite (even if they are small). Any given infinite (or long) trajectory is
 12 quite rare under the probability distribution of the stochastic controller.
 13 Policy-gradient methods sample lots of trajectories from the system and
 14 average the returns across these trajectories. Since the set of trajectories
 15 of even a small MDP is so large, sampling lots of trajectories, or even the
 16 most likely ones, is also very hard. This is a key challenge in getting RL
 17 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

$$\hat{\theta} = \operatorname{argmax}_{\theta} J(\theta; x_0) := \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) | x_0]. \quad (8.4)$$

The objective $J(\theta)$ is called the average return of the controller u_θ .

18 **Computing the average return $J(\theta)$** Before we move on to optimizing
 19 $J(\theta)$, let us discuss how to compute it for given weights θ of the stochastic
 20 controller. We can sample n trajectories from the system and compute the
 21 an estimate of the expectation

$$\hat{J}(\theta) \approx \frac{1}{n} \sum_{i=0}^n \sum_{k=0}^T \gamma^k r(x_k^i, u_k^i) \quad (8.5)$$

22 for some large time-horizon T and where each $u_k^i \sim u_\theta(\cdot | x_k^i)$.

23 8.2 Cross-Entropy Method (CEM)

24 Let us first consider a simple method to compute the best controller. The
 25 basic idea is to solve the problem

$$\hat{\theta} = \operatorname{argmax}_{\theta} J(\theta)$$

🔗 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 using gradient descent. We would like to update weights θ iteratively

$$\theta^{(k+1)} = \theta^{(k)} + \eta \nabla J(\theta).$$

2 where the step-size is $\eta > 0$ and $\nabla J(\theta)$ is the gradient of the objective
3 $J(\theta)$ with respect to weights θ . Instead of computing the exact $\nabla J(\theta)$
4 which we will do in the next section, let us simply compute the gradient
5 using a finite-difference approximation. The i^{th} entry of the gradient is

$$(\widehat{\nabla} J(\theta))_i = \frac{J(\theta + \epsilon e_i) - J(\theta - \epsilon e_i)}{2\epsilon} \approx \frac{\widehat{J}(\theta + \epsilon e_i) - \widehat{J}(\theta - \epsilon e_i)}{2\epsilon}.$$

6 where $e_i = [0, 0, \dots, 0, 1, 0, \dots]$ is a vector with 1 on the i^{th} entry. Each
7 quantity \widehat{J} is computed as the empirical average return of n trajectories
8 from the system. We compute all entries of the objective using this
9 approximation and update the parameters using

$$\theta^{(k+1)} = \theta^{(k)} + \eta \widehat{\nabla} J(\theta^{(k)}).$$

10 **A more efficient way to compute the gradient using finite-differences**

11 Instead of picking perturbations e_i along the cardinal directions, let us
12 sample them from a Gaussian distribution

$$\xi^i \sim N(0, \sigma^2 I)$$

13 for some user-chosen covariance σ^2 . We can however no longer use
14 the finite-difference formula to compute the derivative because the noise
15 e is not aligned with the axes. We can however use a Taylor series
16 approximation as follows. Observe that

$$J(\theta + \xi) \approx J(\theta) + \langle \nabla J(\theta), \xi \rangle$$

17 where $\langle \cdot, \cdot \rangle$ is the inner product. Given m samples ξ^1, \dots, ξ^m observe
18 that

$$\begin{aligned} \widehat{J}(\theta + \xi^1) &= \widehat{J}(\theta) + \langle \nabla J(\theta), \xi^1 \rangle \\ \widehat{J}(\theta + \xi^2) &= \widehat{J}(\theta) + \langle \nabla J(\theta), \xi^2 \rangle \\ &\vdots \\ \widehat{J}(\theta + \xi^m) &= \widehat{J}(\theta) + \langle \nabla J(\theta), \xi^m \rangle. \end{aligned} \tag{8.6}$$

19 is a linear system of equations in $\nabla J(\theta) \in \mathbb{R}^p$ where $\theta \in \mathbb{R}^p$. All
20 quantities \widehat{J} are estimated as before using trajectories drawn from the
21 system. We solve this linear system, e.g., using least-squares if $m > p$, to
22 get an estimate of the gradient $\widehat{\nabla} J(\theta) \in \mathbb{R}^p$.

The Cross-Entropy Method is a more crude but simpler way to

implement the above least-squares formulation. At each iteration it updates the parameters using the formula

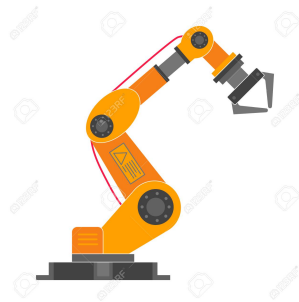
$$\theta^{(k+1)} = \mathbb{E}_{\theta \sim N(\theta^{(k)}, \sigma^2 I)} \left[\theta \mathbf{1}_{\{\hat{J}(\theta) > \hat{J}(\theta^{(k)})\}} \right]. \quad (8.7)$$

In simple words, the CEM samples a few stochastic controllers u_θ from a Gaussian (or any other distribution) centered around the current controller $u_{\theta^{(k)}}$ and updates the weights θ^k in a direction that leads to an increase in $\hat{J}(\theta) > \hat{J}(\theta^{(k)})$.

8.2.1 Some remarks on sample complexity of simulation-based methods

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 θ in high-dimensions. Typical applications are actually the latter, remember that we are interested in using a deep network as a stochastic controller and θ 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

$$\begin{aligned} p &= \underbrace{(12 \times 64 + 64)}_{\text{for } \mu_\theta(x)} + \underbrace{(64 \times 6 + 6)}_{\text{for } \sigma_\theta^2(x)} + \underbrace{(12 \times 64 + 64)}_{\text{for } \mu_\theta(x)} + \underbrace{(64 \times 6 + 6)}_{\text{for } \sigma_\theta^2(x)} \\ &= 2,444. \end{aligned}$$

This is a very high-dimensional space to sample exhaustively. Note that 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.

1 We need two estimates $\hat{J}(\theta - \epsilon e_i)$ and $\hat{J}(\theta + \epsilon e_i)$ for every dimension
 2 $i \in \{1, \dots, p\}$. Each estimate requires us to obtain n trajectories from
 3 the system. Since the number of trajectories that a robot can take is quite
 4 diverse, we should use a large n , so let's pick $n = 100$. The total number
 5 of trajectories required to update the parameters $\theta^{(k)}$ at each iteration is

$$2 p n \approx 10^6.$$

6 *This is an absurdly large number*, and things are even more daunting
 7 when we realize that each update of the weights requires us to sample
 8 these many trajectories from the system. It is not reasonable to sample
 9 such a large number of trajectories from an actual robot, that too for each
 10 update of the weights.

11 **Using fast simulators for RL** If we expand our horizon and think
 12 of learning controllers in simulation, things feel much more reasonable.
 13 While running a large number of trajectories may degrade a real robot
 14 beyond use, doing so requires just computation time in a robot sim-
 15 ulator. There is a large number of simulators that are available with
 16 various capabilities, e.g., Gazebo (<http://gazebo.org>) is a sophisticated
 17 simulator inside ROS that uses a number of Physics engines such as
 18 Bullet (<https://pybullet.org/wordpress>), MuJoCo (<http://www.mujoco.org>)
 19 is incredibly fast although not very good modeling contact, Unity is
 20 a popular platform to simulate driving and more complicated scenes
 21 (<https://docs.nvidia.com/isaac/isaac/doc/simulation/unity3d.html>), Drake
 22 (<https://drake.mit.edu>) is better at contact modeling but more complex and
 23 slower. Most robotics companies develop their own driving simulators
 24 in-house. The assigned reading (#2) for this chapter is a paper which
 25 develops a very fast implementation of the CEM for use in simulators.

26 **Working well in simulation does not mean that a controller works
 27 well on the real robot** It is important to realize that a simulator is not
 28 equivalent to the physical robot. Each simulator makes certain trade-offs
 29 in capturing the dynamics of the real system and it is not a given that a
 30 controller that was learned using data from a simulator will work well on
 31 a real robot. For instance, OpenAI had to develop a large number of tricks
 32 (which took about a year) to modify the simulator to enable the learned
 33 policy to work well on a robot (<https://openai.com/blog/learning-dexterity>)
 34 for a fairly narrow set of tasks.

35 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.

❗ 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 We would like to solve the optimization problem

$$\max_{\theta} J(\theta) := \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau) \mid x_0]$$

2 We will suppress the dependence on x_0 to keep the notation clear. The
3 expectation is taken over all trajectories starting at state x_0 realized using
4 the stochastic controller $u_{\theta}(\cdot \mid x)$. We to update weights θ using gradient
5 descent which amounts to

$$\theta^{(k+1)} = \theta^{(k)} + \eta \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau)].$$

6 First let us note that the distribution p_{θ} using which we compute the
7 expectation also depends on the weights θ . This is why we cannot simply
8 move the derivative ∇_{θ} inside the expectation

$$\nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau)] \neq \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} R(\tau)].$$

9 We need to think of a new technique to compute the gradient above.
10 Essentially, we would like to do the chain rule of calculus but where
11 one of the functions in the chain is an expectation. The likelihood-ratio
12 trick described next allows us to take such derivatives. Here is how the
13 computation goes

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}} [R(\tau)] &= \nabla_{\theta} \int R(\tau) p_{\theta}(\tau) d\tau \\ &= \int R(\tau) \nabla_{\theta} p_{\theta}(\tau) d\tau \\ &\quad \text{(move the gradient inside, integral is over trajectories } \tau \text{ which do not depend on } \theta \text{ themselves)} \\ &= \int R(\tau) p_{\theta}(\tau) \frac{\nabla p_{\theta}(\tau)}{p_{\theta}(\tau)} d\tau \\ &= \int R(\tau) p_{\theta}(\tau) \nabla \log p_{\theta}(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau) \nabla \log p_{\theta}(\tau)] \\ &\approx \frac{1}{n} \sum_{i=1}^n R(\tau^i) \nabla \log p_{\theta}(\tau^i) \end{aligned} \tag{8.8}$$

14 This is called the likelihood-ratio trick to compute the policy gradient. It
15 simply multiplies and divides by the term $p_{\theta}(\tau)$ and rewrites the term
16 $\frac{\nabla p_{\theta}}{p_{\theta}} = \nabla \log p_{\theta}$. It gives us a neat way to compute the gradient: we
17 sample n trajectories τ^1, \dots, τ^n from the system and average the return
18 of each trajectory $R(\tau^i)$ weighted by the gradient of the likelihood of
19 taking each trajectory $\log p_{\theta}(\tau^i)$. The central point to remember here is

1 that the gradient

$$\begin{aligned} \nabla_{\theta} \log p_{\theta}(\tau^i) &= \nabla_{\theta} \sum_{k=0}^T \log p(x_{k+1}^i | x_k^i, u_k^i) + \log u_{\theta}(u_k^i | x_k^i) \\ &= \sum_{k=0}^T \nabla_{\theta} \log u_{\theta}(u_k^i | x_k^i) \end{aligned} \quad (8.9)$$

2 is computed using backpropagation for a neural network. This expression
3 is called the policy gradient because it is the gradient of the objective $J(\theta)$
4 with respect to the parameters of the controller/policy θ .

5 **Variance of policy gradient** The expression for the policy gradient may
6 seem like a sleight of hand. It is a clean expression to get the gradient of
7 the objective but also comes with a number of problems. Observe that

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [R(\tau)] &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[R(\tau) \frac{\nabla p_{\theta}(\tau)}{p_{\theta}(\tau)} \right] \\ &\approx \frac{1}{n} \sum_{i=1}^n R(\tau^i) \frac{\nabla p_{\theta}(\tau^i)}{p_{\theta}(\tau^i)}. \end{aligned}$$

8 If we sample trajectories τ^i that are not very likely under the distribution
9 $p_{\theta}(\tau)$, the denominator in some of the summands can be very small.
10 For trajectories that are likely, the denominator is large. The empirical
11 estimate of the expectation using n trajectories where some terms are
12 very small and some others very large, therefore has a large variance. So
13 one does need lots of trajectories from the system/simulator to compute a
14 reasonable approximation of the policy gradient.

15 8.3.1 Reducing the variance of the policy gradient

16 **Control variates** You will perhaps appreciate that computing the accu-
17 rate policy gradient is very hard in practice. Control variates is a general
18 concept from the literature on Monte Carlo integration and is typically
19 introduced as follows. Say we have a random variable X and we would
20 like to guess its expected value $\mu = \mathbb{E}[X]$. Note that X is an unbiased
21 estimator of μ but it may have a large variance. If we have another random
22 variable Y with known expected value $\mathbb{E}[Y]$, then

$$\hat{X} = X + c(Y - \mathbb{E}[Y]) \quad (8.10)$$

23 is also an unbiased estimator for μ for any value of c . The variance of \hat{X} is

$$\text{Var}(\hat{X}) = \text{Var}(X) + c^2 \text{Var}(Y) + 2c \text{Cov}(X, Y).$$

24 which is minimized for

$$c^* = -\frac{\text{Cov}(X, Y)}{\text{Var}(Y)}$$

1 for which we have

$$\text{Var}(\hat{X}) = \text{Var}(X) - c^{*2} \text{Var}(Y) = \left(1 - \left(\frac{\text{Cov}(X, Y)}{\text{Var}(Y)}\right)^2\right) \text{Var}(X).$$

2 By subtracting $Y - \mathbb{E}[Y]$ from our observed random variable X , we have
 3 reduced the variance of X if the correlation between X and Y is non-zero.
 4 Most importantly, note that no matter what Y we plug into the above
 5 expression, we can never increase the variance of X ; the worst that can
 6 happen is that we pick a Y that is completely uncorrelated with X and
 7 end up achieving nothing.

8 **Baseline** We will now use the concept of a control variate to reduce the
 9 variance of the policy gradient. This is known as “building a baseline”.
 10 The simplest baseline one can build is to subtract a constant value from
 11 the return. Consider the PG given by

$$\begin{aligned} \nabla J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}} [R(\tau) \nabla \log p_{\theta}(\tau)] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [(R(\tau) - b) \nabla \log p_{\theta}(\tau)]. \end{aligned}$$

12 Observe that

$$\begin{aligned} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [b \nabla \log p_{\theta}(\tau)] &= \int d\tau b p_{\theta}(\tau) \nabla \log p_{\theta}(\tau) \\ &= \int d\tau b \nabla p_{\theta}(\tau) = b \nabla \int d\tau p_{\theta}(\tau) = b \nabla 1 = 0. \end{aligned}$$

13 **Example 1: Using the average returns of a mini-batch as the baseline**

14 What is the simplest baseline b we can cook up? Let us write the mini-batch
 15 version of the policy gradient

$$\hat{\nabla} J(\theta) := \frac{1}{\ell} \sum_{i=1}^{\ell} [R(\tau^i) \nabla \log p_{\theta}(\tau^i)].$$

16 where $\tau^1, \dots, \tau^{\ell}$ are trajectories that are a part of our mini-batch. We can
 17 set

$$b = \frac{1}{\ell} \sum_{i=1}^{\ell} R(\tau^i)$$

18 can use the variance-reduced gradient

$$\hat{\nabla} J(\theta) = \frac{1}{\ell} \sum_{i=1}^{\ell} [(R(\tau^i) - b) \nabla \log p_{\theta}(\tau^i)].$$

19 This is a one-line change in your code for policy gradient so there is no
 20 reason not to do it.

Example 2: A weighted averaged of the returns using the log-likelihood of the trajectory The previous example showed how we can use one constant baseline, namely the average of the discounted returns of all trajectories in a mini-batch. What is the **best** constant b we can use? We can perform a similar computation as done in the control variate introduction to minimize the variance of the policy gradient to get the following.

$$\begin{aligned} \delta \left(\widehat{\nabla}_{\theta_i} J(\theta) \right) &= \mathbb{E}_{\tau} \left[((R(\tau) - b_i) \nabla_{\theta_i} \log p_{\theta}(\tau))^2 \right] - \left(\mathbb{E}_{\tau} [((R(\tau) - b_i)) \nabla_{\theta_i} \log p_{\theta}(\tau)] \right)^2 \\ &= \mathbb{E}_{\tau} \left[((R(\tau) - b_i) \nabla_{\theta_i} \log p_{\theta}(\tau))^2 \right] - \left(\widehat{\nabla}_{\theta_i} J(\theta) \right)^2. \end{aligned}$$

Set

$$\frac{\delta \left(\widehat{\nabla}_{\theta_i} J(\theta) \right)}{db_i} = 0$$

in the above expression to get

$$b_i = \frac{\mathbb{E}_{\tau} \left[(\nabla_{\theta_i} \log p_{\theta}(\tau))^2 R(\tau) \right]}{\mathbb{E}_{\tau} \left[(\nabla_{\theta_i} \log p_{\theta}(\tau))^2 \right]}$$

which is the baseline you should subtract from the gradient of the i^{th} parameter θ_i to result in the largest variance reduction. This expression is just the expected return but it is weighted by the magnitude of the gradient, this again 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

$$d^{\theta}(x) = \sum_{k=0}^{\infty} \gamma^k \mathbb{P}(x_k = x | x_0, u_k \sim u_{\theta}(\cdot | x_k)).$$

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

1 empirically and we will not worry about it here; RL algorithms do not
2 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.

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\tau \sim p_\theta} [R(\tau) \nabla \log p_\theta(\tau)] \\ &= \mathbb{E}_{x \sim d^\theta} \mathbb{E}_{u \sim u_\theta(\cdot | x)} [q^\theta(x, u) \nabla_\theta \log u_\theta(u | x)].\end{aligned}\quad (8.11)$$

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

$$q^\theta(x, u) = \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) | x_0 = x, u_0 = u].\quad (8.12)$$

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^θ on the parameters θ of the controller explicit.

3 Compare the above formula for the policy gradient with the one we
4 had before in (8.8)

$$\begin{aligned}\widehat{\nabla} J(\theta) &= \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)] \\ &= \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\left(\sum_{k=0}^T \gamma^k r(x_k, u_k) \right) \left(\sum_{k=0}^T \nabla \log u_\theta(u_k | x_k) \right) \right].\end{aligned}$$

5 It is important to notice that this is an expectation over trajectories;
6 whereas (8.11) is an expectation over states x sampled from the discounted
7 state visitation frequency. The control u_k for both is sampled from the
8 stochastic controller at each time-step k . The most important distinction
9 is that (8.11) involves the expectation of the Q-function q^θ weighted by
10 the gradient of the log-likelihood of picking each control action. There
11 are numerous hacky ways of deriving (8.11) from (8.8) but remember that
12 they are fundamentally different expressions of the *same policy gradient*.

13 This expression allows understanding of a number of properties of
14 reinforcement learning.

- 15 1. While the algorithm collects the data, states that are unlikely under
16 the distribution d^θ contribute little to (8.11). In other words, the
17 policy gradient is insensitive to such states. The policy update will
18 not consider these unlikely states that the system is prone to visit
19 infrequently using the controller u_θ .
- 20 2. The opposite happens for states which are very likely. For two
21 controls u_1, u_2 at the same state x , the policy increases the log-
22 likelihood of taking the controls weighted by their values $q^\theta(x, u_1)$

❗ 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>.

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 σ 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_θ 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^θ .

8.5 Actor-Critic methods

We can of course do more sophisticated things to evaluate the Q-function q^θ 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.

1 But we do have the ability to get trajectories τ^i from the system.

2 Let's say (x_k^i, u_k^i) lie on τ^i at time-step k . We can then estimate the
3 expectation over $P(\cdot | x_k^i, u_k^i)$ using simply x_{k+1}^i and the expectation over
4 the controls using simply u_{k+1}^i to write

$$q^\theta(x_k^i, u_k^i) \approx r(x_k^i, u_k^i) + \gamma q^\theta(x_{k+1}^i, u_{k+1}^i) \quad \text{for all } i \leq n, k \leq T.$$

5 This is a nice constraint on the Q-function. If this were a discrete-state,
6 discrete-control MDP, it is a set of linear equations for the q-values. These
7 constraints would be akin to our linear equations for evaluating a policy in
8 dynamic programming except that instead of using the dynamics model
9 (the transition matrix), we are using trajectories sampled from the system.

10 **Parameterizing the Q-function using a neural network** If we are deal-
11 ing with a continuous state/control-space, we can think of parameterizing
12 the q -function using parameters φ

$$q_\varphi^\theta(x, u) : X \times U \rightarrow \mathbb{R}.$$

13 The parameterization is similar to the parameterization of the controller,
14 e.g., just like we would write a deterministic controller as

$$u_\theta(x) = \theta^\top x$$

15 we can think of a linear Q-function of the form

$$q_\varphi^\theta(x, u) = \varphi^\top \begin{bmatrix} x \\ u \end{bmatrix}, \quad \varphi \in \mathbb{R}^{m+d}$$

16 which is a linear function in the states and controls. You can also think of
17 using something like

$$q_\varphi^\theta(x, u) = \begin{bmatrix} 1 & x & u \end{bmatrix} \varphi \begin{bmatrix} 1 \\ x \\ u \end{bmatrix} \quad \varphi \in \mathbb{R}^{(m+d+1) \times (m+d+1)}.$$

18 which is quadratic in the states and controls, or in general a deep network
19 with weights φ as the Q-function.

20 **Fitting the Q-function** We can now “fit” the parameters of the Q-
21 function by solving the problem

$$\hat{\varphi} = \underset{\varphi}{\operatorname{argmin}} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{k=0}^T \left\| q_\varphi^\theta(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma q_\varphi^\theta(x_{k+1}^i, u_{k+1}^i) \right\|^2. \quad (8.14)$$

22 which is nothing other than enforcing the Bellman equation in (8.13).
23 If the Q-function is linear in $[x, u]$ this is a least squares problem, if it
24 is quadratic the problem is a quadratic optimization problem which can
25 also be solved efficiently, in general we would solve this problem using
26 stochastic gradient descent if we are parameterizing the Q-function using a

1 deep network. Such a Q-function is called the “critic” because it *evaluates*
 2 the controller u_θ , or the “actor”. This version of the policy gradient where
 3 one fits the parameters of both the controller and the Q-function are called
 4 Actor-Critic methods.

Actor-Critic Methods We fit a deep network with weights θ to parameterize a stochastic controller $u_\theta(\cdot | x)$ and another deep network with weights φ to parameterize the Q-function of this controller, $q_\varphi^\theta(x, u)$. Let the controller weights at the k^{th} iteration be $\theta^{(k)}$ and the Q-function weights be φ^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 φ^k .
3. Compute the policy gradient using the alternative expression in (8.11) and update parameters of the policy to $\theta^{(k+1)}$.

5 8.5.1 Advantage function

6 The new expression for the policy gradient in (8.11) also has a large
 7 variance; this should be no surprise, it is after all equal to the old
 8 expression. We can however perform variance reduction on this using the
 9 value function.

10 Our goal as before would be construct a relevant baseline to subtract
 11 from the Q-function. It turns out that any function that depends only upon
 12 the state x is a valid baseline. This gives a powerful baseline for us to to
 13 use in policy gradients. We can use the value function as the baseline. The
 14 value function for controls taken by the controller $u_\theta(\cdot | x)$ (notice that
 15 this is not the optimal value function, it is simply the policy evaluation) is
 16 given by

$$v^\theta(x) = \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau) | x_0 = x]$$

17 where $u_k \sim u_\theta(\cdot | x_k)$ at each timestep. We also know that the value
 18 function is the expected value of the Q-function across different controls
 19 sampled by the controller

$$v^\theta(x) = \mathbb{E}_{u \sim u_\theta(\cdot | x)} [q^\theta(x, u)]. \quad (8.15)$$

20 The value function again satisfies the dynamic programming principle/-

i We will be pedantic and always write the q -function as q_φ^θ . The superscript θ denotes that this is the q -function corresponding to the θ -parameterized controller u_θ . The subscript denotes that the q -function is parameterized by parameters φ .

1 Bellman equation

$$v^\theta(x) = \mathbb{E}_{u \sim u_\theta(\cdot|x)} \left[r(x, u) + \gamma \mathbb{E}_{x' \sim \mathcal{P}(\cdot|x, u)} [v^\theta(x')] \right].$$

2 We again parameterize the value function

$$v_\psi^\theta(x) : X \rightarrow \mathbb{R}$$

3 using parameters ψ and fit it to the data in the same way as (8.14) to get

$$\hat{\psi} = \underset{\psi}{\operatorname{argmin}} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{k=0}^T \|v_\psi^\theta(x_k^i) - r(x_k^i, u_k^i) - \gamma v_\psi^\theta(x_{k+1}^i)\|^2. \quad (8.16)$$

Using this baseline can modify the policy gradient to be

$$\nabla J(\theta) = \mathbb{E}_{x \sim d^\theta} \mathbb{E}_{u \sim u_\theta(\cdot|x)} \left[\left(\underbrace{q_\varphi^\theta(x, u) - v_\psi^\theta(x)}_{a_{\varphi, \psi}^\theta(x, u)} \right) \nabla_\theta \log u_\theta(u | x) \right]. \quad (8.17)$$

where each of the functions q_φ^θ and v_ψ^θ are themselves fitted using (8.14) and (8.16) respectively. The difference

$$\begin{aligned} a_{\varphi, \psi}^\theta(x, u) &= q_\varphi^\theta(x, u) - v_\psi^\theta(x) \\ &\approx q_\varphi^\theta(x, u) - \mathbb{E}_{u \sim u_\theta(\cdot|x)} [q_\varphi^\theta(x, u)] \end{aligned} \quad (8.18)$$

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_φ^θ and the estimate v_ψ^θ . 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?

4 **8.6 Discussion**

5 This brings to an end the discussion of policy gradients. They are, in
6 general, a complicated suite of algorithms to implement. You will see
7 some of this complexity when you implement the controller for something
8 as simple as the simple pendulum. The key challenges with implementing
9 policy gradients come from the following.

- 10 1. Need lots of data, each parameter update requires fresh data from
- 11 the systems. Typical problems may need a million trajectories, most
- 12 robots would break before one gets this much data from them if one
- 13 implements these algorithms naively.

- 1 2. The log-likelihood ratio trick has a high variance due to $u_{\theta}(\cdot | x)$
2 being in the denominator of the expression, so we need to implement
3 complex variance reduction techniques such as actor-critic methods.

- 4 3. Fitting the Q-function and the value function is not easy, each
5 parameter update of the policy ideally requires you to solve the
6 entire problems (8.14) and (8.16). In practice, we only perform a
7 few steps of SGD to solve the two problems and reuse the solution
8 of k^{th} iteration update as an initialization of the $k + 1^{\text{th}}$ update. This
9 is known as “warm start” in the optimization literature and reduces
10 the computational cost of fitting the Q/value-functions from scratch
11 each time.

- 12 4. The Q/value-function fitted in iteration k may be poor estimates of
13 the Q/value at iteration $k + 1$ for the new controller $u_{\theta^{(k+1)}}(\cdot | x)$.
14 If the controller parameters change quickly, $\theta^{(k+1)}$ is very dif-
15 ferent from $\theta^{(k)}$, then so are $q^{\theta^{(k+1)}}$ and $v^{\theta^{(k+1)}}$. There is a
16 very fine balance between training quickly and retaining the ef-
17 ficiency of warm start; and tuning this in practice is quite dif-
18 ficult. A large number of policy gradient algorithms like TRPO
19 (<https://arxiv.org/abs/1502.05477>) and PPO (<https://arxiv.org/abs/1707.06347>)
20 try to alleviate this with varying degrees of success.

- 21 5. The latter, PPO, is a good policy-gradient-based algorithm to try
22 on a new problem. For instance, in a very impressive demon-
23 stration, it was used to build an RL agent to play StarCraft
24 (<https://openai.com/blog/openai-five>). We will see better RL meth-
25 ods 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

- 1 number of controls, with each entry in this table being the value $q(x, u)$.
 2 Value iteration when written using the Q-function at the k^{th} iteration for
 3 the tabular setting looks like

$$\begin{aligned} q^{(k+1)}(x, u) &= \sum_{x' \in X} P(x' | x, u) \left(r(x, u) + \gamma \max_{u'} q^{(k)}(x', u') \right) \\ &= \mathbb{E}_{x' \sim P(\cdot | x, u)} \left[r(x, u) + \gamma \max_{u'} q^{(k)}(x', u') \right]. \end{aligned}$$

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.

- 4 We will imagine the robot as using an *arbitrary* controller

$$u_e(\cdot | x)$$

- 5 that has a fairly large degree of randomness in how it picks actions. We
 6 call such a controller an “exploratory controller”. Conceptually, its goal
 7 is to lead the robot to diverse states in the state-space so that we get a
 8 faithful estimate of the expectation in value iteration. We maintain the
 9 value $q^{(k)}(x, u)$ for all states $x \in X$ and controls $u \in U$ and update these
 10 values to get $q^{(k+1)}$ after *each step* of the robot.

- 11 From the results on Bellman equation, we know that any Q-function
 12 that satisfies the above equation is the optimal Q-function; we would
 13 therefore like our Q-function to satisfy

$$q^*(x_k, u_k) \approx r(x_k, u_k) + \gamma \max_{u'} q^*(x_{k+1}, u').$$

- 14 over samples (x_k, u_k, x_{k+1}) collected as the robot explores the environ-
 15 ment.

Tabular Q-Learning Let us imagine the robot travels for n trajectories of T time-steps each. We can now solve for q^* by minimizing the objective

$$\min_q \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{k=0}^T \left(q(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma \max_{u' \in U} q(x_{k+1}^i, u') \right)^2. \quad (9.1)$$

on the data collected by the robot. The variable of optimization here are all values $q^*(x, u)$ for $x \in X$ and $u \in U$.

- 16 Notice a few important things about the above optimization problem.
 17 First, the last term is a maximization over $u' \in U$, it is $\max_{u' \in U} q(x_{k+1}^i, u')$
 18 and not $q(x_{k+1}^i, u_{k+1}^i)$. In practice, you should imagine a robot performing

1 Q-Learning in a grid-world setting where it seeks to find the optimal
 2 trajectory to go from a source location to a target location. If at each
 3 step, the robot has 4 controls to choose from, computing this last term
 4 involves taking the maximum of 4 different values (4 columns in the
 5 tabular Q-function).

6 Notice that for finite-horizon dynamic programming we initialized
 7 the Q-function at the terminal time to a known value (the terminal cost).
 8 Similarly, for infinite-horizon value iteration, we discussed how we can
 9 converge to the optimal Q-function with any initialization. In the above
 10 case, we do not impose any such constraint upon the Q-function, but there
 11 is an implicit constraint. All values $q(x, u)$ have to be consistent with
 12 each other and ideally, the residual

$$q(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma \max_{u' \in U} q(x_{k+1}^i, u') = 0$$

13 for all trajectories i and all timesteps T .

14 **Solving tabular Q-Learning** How should we solve the optimization
 15 problem in (9.1)? This is easy, every entry $q(x, u)$ for $x \in U$ and $u \in U$
 16 is a variable of this objective and each $(\cdot)^2$ term in the objective simply
 17 represents a constraint that ties these different values of the Q-function
 18 together. We can solve for all $q(x, u)$ iteratively as

$$\begin{aligned} q(x, u) &\leftarrow q(x, u) - \eta \nabla_{q(x, u)} \ell(q) \\ &= (1 - \eta) q(x, u) - \eta \left(r(x, u) + \gamma \max_{u'} q(x', u') \right) \end{aligned} \quad (9.2)$$

19 where $\ell(q)$ is the entire objective $\frac{1}{n(T+1)} \sum_i \sum_k \dots$ above and $(x, u, x') \equiv$
 20 $(x_k^i, u_k^i, x_{k+1}^i)$ in the second equation. An important point to note here is
 21 that although the robot collects a finite number of data

$$D = \{(x_k^i, u_k^i)_{k=0,1,\dots,T}\}_{i=1}^n$$

22 we have an estimate for the value $q(x, u)$ at all states $x \in X$. As an
 23 intuition, tabular Q-learning looks at the returns obtained by the robot
 24 after starting from a state x (the reward-to-come $J(x)$) and patches the
 25 returns from nearby states x, x' using the constraints in the objective (9.1).

26 **Terminal state** One must be very careful about the terminal state in such
 27 implementations of Q-learning. Typically, most research papers imagine
 28 that they are solving an infinite horizon problem but use simulators that
 29 have an explicit terminal state, i.e., the simulator does not proceed to the
 30 next timestep after the robot reaches the goal. A workaround for using
 31 such simulators (essentially all simulators) is to modify (9.2) as

$$q(x, u) = (1 - \eta) q(x, u) - \eta \left(r(x, u) + \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \max_{u'} q(x', u') \right).$$

1 Effectively, we are setting $q(x', u) = 0$ for all $u \in U$ if x' is a terminal state
 2 of problem. This is a very important point to remember and Q-Learning
 3 will never work if you forget to include the term $\mathbf{1}_{\{x' \text{ is terminal}\}}$ in your
 4 expression.

5 **What is the controller in tabular Q-Learning?** The controller in
 6 tabular Q-Learning is easy to get after we solve (9.1). At test time, we use
 7 a deterministic controller given by

$$u^*(x) = \operatorname{argmax}_{u'} q^*(x, u').$$

8 9.1.1 How to perform exploration in Q-Learning

9 The exploratory controller used by the robot $u_e(\cdot | x)$ is critical to perform
 10 Q-Learning well. If the exploratory controller does not explore much,
 11 we do not get states from all parts of the state-space. This is quite bad,
 12 because in this case the estimates of Q-function at *all states* will be bad,
 13 not just at the states that the robot did not visit. To make this intuitive,
 14 imagine if we cordoned off some nodes in the graph for the backward
 15 version of Dijkstra's algorithm and never used them to update the dist
 16 variable. We would never get to the optimal cost-to-go for *all* states in this
 17 case because there could be trajectories that go through these cordoned
 18 off states that lead to a smaller cost-to-go. So it is quite important to pick
 19 the right exploratory controller.

20 It turns out that a random exploratory controller, e.g., a controller
 21 $u_e(\cdot | x)$ that picks controls uniformly randomly is pretty good. We can
 22 show that our tabular Q-Learning will converge to the optimal Q-function
 23 $q^*(x, u)$ as the amount of data drawn from the random controller goes to
 24 infinity, even if we initialize the table to arbitrary values. In other words,
 25 if we are guaranteed that the robot visits each state in the finite MDP
 26 infinitely often, it is a classical result that updates of the form (9.2) for
 27 minimizing the objective in (9.1) converge to the optimal Q-function.

28 **Epsilon-greedy exploration** Instead of the robot using a arbitrary
 29 controller $u_e(\cdot | x)$ to gather data, we can use the current estimate of the
 30 Q-function with some added randomness to ensure that the robot visits all
 31 states in the state-space. This is a key idea in Q-Learning and is known as
 32 "epsilon-greedy" exploration. We set

$$u_e(u | x) = \begin{cases} \operatorname{argmax}_u q(x, u) & \text{with probability } 1 - \epsilon \\ \operatorname{uniform}(U) & \text{with probability } \epsilon. \end{cases} \quad (9.3)$$

33 for some user-chosen value of ϵ . Effectively, the robot repeats the controls
 34 it took in the past with probability $1 - \epsilon$ and uniformly samples from
 35 the entire control space with probability ϵ . The former ensures that the
 36 robot moves towards the parts of the state-space where states have a high
 37 return-to-come (after all, that is the what the Q-function $q(x, u)$ indicates).

❗ 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 The latter ensures that even if the robot's estimate of the Q-function is bad,
2 it is still visiting every state in the state-space infinitely often.

3 **A different perspective on Q-Learning** Conceptually, we can think of
4 tabular Q-learning as happening in two stages. In the first stage, the robot
5 gathers a large amount of data

$$D = \{(x_k^i, u_k^i)_{k=0,1,\dots,T}\}_{i=1}^n$$

6 using the exploratory controller $u_e(\cdot | x)$; let us consider the case
7 when we are using an arbitrary exploratory controller, not epsilon-greedy
8 exploration. Using this data, the robot fits a model for the system, i.e., it
9 learns the underlying MDP

$$P(x' | x, u);$$

10 this is very similar to the step in the Baum-Welch algorithm that we saw
11 for learning the Markov state transition matrix of the HMM in Chapter 2.
12 We simply take frequency counts to estimate this probability

$$P(x' | x, u) \approx \frac{1}{N} \sum_i \mathbf{1}_{\{x' \text{ was reached from } x \text{ using control } u\}}$$

13 where N is the number of the times the robot took control u at state x .
14 Given this transition matrix, we can now perform value iteration on the
15 MDP to learn the Q-function

$$q^{(k+1)}(x, u) = \mathbf{E}_{x' \sim P(\cdot | x, u)} \left[r(x, u) + \gamma \max_{u'} q^{(k)}(x', u) \right].$$

16 The success of this two-stage approach depends upon how accurate our
17 estimate of $P(x' | x, u)$ is. This in turn depends on how much the robot
18 explored the domain and the size of the dataset it collected, both of these
19 need to be large. We can therefore think of Q-learning as interleaving
20 these two stages in a single algorithm, it learns the dynamics of the system
21 and the Q-function for that dynamics simultaneously. But the Q-Learning
22 algorithm does not really maintain a representation of the dynamics, i.e.,
23 at the end of running Q-Learning, we do not know what $P(x' | x, u)$ is.

24 9.2 Function approximation (Deep Q Networks)

25 Tabular methods are really nice but they do not scale to large problems.
26 The grid-world in the homework problem on policy iteration had 100
27 states, a typical game of Tetris has about 10^{60} states. For comparison, the
28 number of atoms in the known universe is about 10^{80} . The number of
29 different states in a typical Atari game is more than 10^{300} . These are all
30 problems with a discrete number of states and controls, for continuous
31 state/control-space, the number of distinct states/controls is infinite. So
32 it is essentially impossible to run the tabular Q-Learning method from

the previous section for most real-world problems. In this section, we will look at a powerful set of algorithms that parameterize the Q-function using a neural network to work around this problem.

We use the same idea from the previous chapter, that of parameterizing the Q-function using a deep network. We will denote

$$q_\varphi(x, u) : X \times U \mapsto \mathbb{R}$$

as the Q-function and our goal is to fit the deep network to obtain the weights $\hat{\varphi}$, instead of maintaining a very large table of size $|X| \times |U|$ for the Q-function. Fitting the Q-function is quite similar to the tabular case: given a dataset $D = \{(x_t^i, u_t^i)\}_{t=0,1,\dots,T}^n$ from the system, we want to enforce

$$q_\varphi(x_t^i, u_t^i) = r(x_t^i, u_t^i) + \gamma \max_{u'} q_\varphi(x_{t+1}^i, u')$$

for all tuples $(x_t^i, u_t^i, x_{t+1}^i)$ in the dataset. Just like the previous section, we will solve

$$\begin{aligned} \hat{\varphi} &= \operatorname{argmin}_\varphi \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=1}^T \left(q_\varphi(x_t^i, u_t^i) - r(x_t^i, u_t^i) - \gamma \underbrace{\left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right)}_{\text{target}(x_{t+1}^i; \varphi)} \max_{u'} q_\varphi(x_{t+1}^i, u') \right)^2 \\ &\equiv \operatorname{argmin}_\varphi \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=1}^T (q_\varphi(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi))^2 \end{aligned} \quad (9.4)$$

The last two terms in this expression above are together called the “target” because the problem is very similar to least squares regression, except that the targets also depend on the weights φ . This is what makes it challenging to solve.

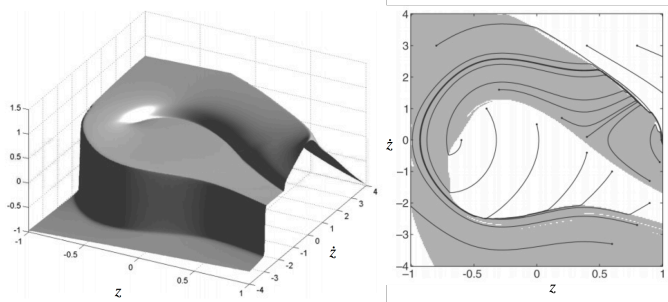
As discussed above, Q-Learning with function approximation is known as “Fitted Q Iteration”. Remember that very important point that the robot collects data using the exploratory controller $u_e(\cdot | x)$ but the Q-function that we fit is the *optimal* Q-function.

Fitted Q-Iteration with function approximation may not converge to the optimal Q-function

It turns out that (9.4) has certain mathematical intricacies that prevent it from converging to the optimal Q-function. We will give the intuitive reason here. In the tabular Q-Learning setting, if we modify some entry $q(x, u)$ for an $x \in X$ and $u \in U$, the other entries (which are tied together using the Bellman equation) are all modified. This is akin to you changing the dist value of one node in Dijkstra’s algorithm; the dist values of *all* other nodes will have to change to satisfy the Bellman equation. This is what (9.2) achieves if implemented with a decaying step-size η ; see <http://users.isr.ist.utl.pt/~mtjspaan/readingGroup/ProofQlearning.pdf> for the proof. This does not hold for (9.4). Even if the objective in (9.4) is zero on our collected dataset, i.e., the Q-function fits data collected by the robot perfectly, the Q-function may not be the *optimal* Q-function. An

❗ 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.

1 intuitive way of understanding this problem is that even if the Bellman
 2 error is zero on samples in the dataset, the optimization objective says
 3 nothing about states that are not present in the dataset; the Bellman error
 4 on them is completely dependent upon the smoothness properties of the
 5 function expressed by the neural architecture. Contrast this comment with
 6 the solution of the HJB equation in Chapter 6 where the value function
 7 was quite non-smooth at some places. If our sampled dataset does not
 8 contain those places, there is no way the neural network can know the
 9 optimal form of the value function.



11 9.2.1 Embellishments to Q-Learning

12 We next discuss a few practical aspects of implementing Q-Learning.
 13 Each of the following points is extremely important to understand how to
 14 get RL to work on real-world problems, so you should internalize these.

15 **Pick mini-batches from different trajectories in SGD** . In practice,
 16 we fit the Q-function using stochastic gradient descent. At each iteration
 17 we sample a mini-batch of inputs $(x_t^i, u_t^i, x_{t+1}^i)$ from different trajectories
 18 $i \in \{1, \dots, n\}$ and update the weights φ in the direction of the negative
 19 gradient.

$$\varphi^{k+1} = \varphi^k - \eta \nabla_{\varphi} (q_{\varphi^k}(x, u) - \text{target}(x'; \varphi^k))^2.$$

20 The mini-batch is picked to have samples from different trajectories
 21 because samples from the same trajectory are correlated to each other
 22 (after all, the robot obtains the next tuple (x', u', x'') from the previous
 23 tuple (x, u, x')).

24 **Replay buffer** The dataset D is known as the replay buffer.

25 **Off-policy learning** The replay buffer is typically not fixed during
 26 training. Instead of drawing data from the exploratory controller u_e , we
 27 can think of the following algorithm. Initialize the Q-function weights to
 28 φ^0 and the dataset to $D = \emptyset$. At the k^{th} iteration,

- 29 • Draw a dataset D^k of n trajectories from the ϵ -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 θ^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 ϵ -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 ϵ -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 ϵ for exploration is critical Towards the beginning of training, we want a large value for ϵ to gather diverse data from the environment. As training progresses, we want to reduce ϵ 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

$$\left| q_{\varphi}(x_t^i, u_t^i) - r(x_t^i, u_t^i) - \gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}} \right) \max_{u'} q_{\varphi}(x_{t+1}^i, u') \right|$$

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.

1 **Using robust regression to fit the Q-function** There may be states in
 2 the replay buffer with very high Bellman error, e.g., the kinks in the value
 3 function for the mountain car obtained from HJB above, if we happen to
 4 sample those. For instance, these are states where the controller “switches”
 5 and is discontinuous function of state x . In these cases, instead of these
 6 few states dominating the gradient for the entire dataset, we can use ideas
 7 in robust regression to reduce their effect on the gradient. A popular way
 8 to do so is to use a Huber-loss in place of the quadratic loss in (9.4)

$$\text{huber}_\delta(a) = \begin{cases} \frac{a^2}{2} & \text{for } |a| \leq \delta \\ \delta(|a| - \frac{\delta}{2}) & \text{otherwise.} \end{cases} \quad (9.5)$$

9 **Delayed target** Notice that the target also depends upon the weights φ :

$$\text{target}(x'; \varphi) := r(x, u) + \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \max_{u'} q_\varphi(x', u').$$

10 This creates a very big problem when we fit the Q-function. Effectively,
 11 both the covariate and the target in (9.4) depend upon the weights of the
 12 Q-function. Minimizing the objective in (9.4) is akin to performing least
 13 squares regression where the targets keep changing every time you solve
 14 for the solution. This is the root cause of why Q-Learning is difficult to
 15 use in practice. A popular hack to get around this problem is to use some
 16 old weights to compute the target, i.e., use the loss

$$\frac{1}{n(T+1)} \sum_{i,t} \left(q_{\varphi^k}(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi^{k'}) \right)^2. \quad (9.6)$$

17 in place of (9.4). Here k' is an iterate much older than k , say $k' = k - 100$.
 18 This trick is called “delayed target”.

19 **Exponential averaging to update the target** Notice that in order to
 20 implement delayed targets as discussed above we will have to save all
 21 weights $\varphi^k, \varphi^{k-1}, \dots, \varphi^{k-100}$, which can be cumbersome. We can
 22 however do yet another clever hack and initialize two copies of the weights,
 23 one for the actual Q-function φ^k and another for the target, let us call it
 24 φ'^k . We set the target equal to the Q-function at initialization. The target
 25 copy is updated at each iteration to be

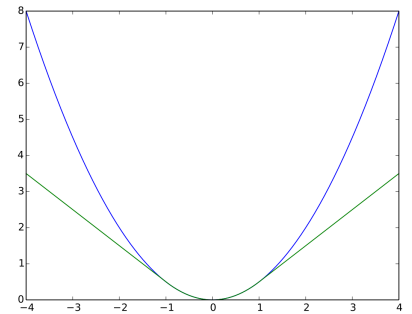
$$\varphi'^{k+1} = (1 - \alpha)\varphi'^k + \alpha\varphi^{k+1} \quad (9.7)$$

26 with some small value, say $\alpha = 0.05$. The target’s weights are therefore
 27 an exponentially averaged version of the weights of the Q-function.

28 **Why are delayed targets essential for Q-Learning to work?** There
 29 are many explanations given why delayed targets are essential in practice
 30 but the correct one is not really known yet.

31 1. For example, one reason could be that since $q_{\varphi^k}(x, u)$ for a given
 32 state typically increases as we train for more iterations in Q-Learning,

❶ Huber loss for $\delta = 1$ (green) compared to the squared error loss (blue).



1 the old weights inside a delayed target are an underestimate of the
 2 true target. This might lead to some stability in situations when the
 3 Q-function's weights φ^k change too quickly when we fit (9.4) or we
 4 do not have enough data in the replay buffer yet.

- 5 2. Another reason one could hypothesize is related to concepts like
 6 self-distillation. For example, we may write a new objective for
 7 Q-Learning that looks like

$$(q_{\varphi^k}(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi^k))^2 + \frac{1}{2\lambda} \|\varphi^k - \varphi^{k'}\|_2^2$$

8 where the second term is known as proximal term that prevents the
 9 weights φ^k from change too much from their old values $\varphi^{k'}$. Proxi-
 10 mal objectives are more stable versions of the standard quadratic
 11 objective in (9.4) and help in cases when one is solving Q-Learning
 12 using SGD updates.

13 **Double Q-Learning** Even a delayed target may not be sufficient to get
 14 Q-Learning to lead to good returns in practice. Focus on one state x . One
 15 problem arise from the max operator in (9.4). Suppose that the Q-function
 16 q_{φ^k} corresponds to a particularly bad controller, say a controller that picks
 17 a control

$$\operatorname{argmax}_u q_{\varphi^k}(x, u)$$

18 that is very different from the optimal control

$$\operatorname{argmax}_u q^*(x, u)$$

19 then, even the delayed target $q_{\varphi^{k'}}$ may be a similarly poor controller. The
 20 ideal target is of course the return-to-come, or the value of the optimal
 21 Q-function $\max_{u'} q^*(x', u')$, but we do not know it while fitting the Q-
 22 function. The same problem also occurs if our Q-function (or its delayed
 23 version, the target) is too optimistic about the values of certain control
 24 inputs, it will consistently pick those controls in the max operator. One
 25 hack to get around this problem is to pick the maximizing control input
 26 using the non-delayed Q-function but use the value of the delayed target

$$\text{target}_{\text{DDQN}}(x_{t+1}^i; \varphi^{k'}) = r(x, u) + \gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right) q_{\varphi^{k'}}(x_{t+1}^i, u'). \quad (9.8)$$

27 where

$$u' = \underbrace{\operatorname{argmax}_u q_{\varphi^k}(x_{t+1}^i, u)}_{\text{control chosen by the Q-function}} .$$

28 **Training two Q-functions** We can also train two copies of the Q-function
 29 simultaneously, each with its own delayed target and mix-and-match their
 30 targets. Let $\varphi^{(1)k}$ and $\varphi'^{(1)k}$ be one Q-function and target pair and $\varphi^{(2)k}$
 31 and $\varphi'^{(2)k}$ be another pair. We update both of them using the following

1 objective.

$$\begin{aligned} \text{For } \varphi^{(1)} : & \left(q^{(1)k}(x, u) - r(x, u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \text{target}_{\text{DDQN}}(x', \varphi'^{(2)k}) \right)^2 \\ \text{For } \varphi^{(2)} : & \left(q^{(2)k}(x, u) - r(x, u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \text{target}_{\text{DDQN}}(x', \varphi'^{(1)k}) \right)^2 \end{aligned} \quad (9.9)$$

2 Sometimes we also use only one target that is the minimum of the two
3 targets (this helps because it is more pessimistic estimate of the true target)

$$\text{target}(x') := \min \left\{ \text{target}_{\text{DDQN}}(x', \varphi'^{(1)k}), \text{target}_{\text{DDQN}}(x', \varphi'^{(2)k}) \right\}.$$

4 You will also see many papers train multiple Q-functions, many more than
5 2. In such cases, it is a good idea to pick the control for evaluation using
6 all the Q-functions:

$$u^*(x) := \operatorname{argmax}_u \sum_k q_{\varphi^{(k)}}(x, u).$$

7 rather than only one of them, as is often done in research papers.

8 **A remark on the various tricks used to compute the target** It may
9 seem that a lot of these tricks are about being pessimistic while computing
10 the target. This is our current understanding in RL and it is born out of
11 the following observation: typically in practice, you will observe that the
12 Q-function estimates can become very large. Even if the TD error is small,
13 the values $q_{\varphi}(x, u)$ can be arbitrarily large; see Figure 1 in [Continuous](#)
14 [Doubly Constrained Batch Reinforcement Learning](#) for an example in a
15 slightly different setting. This occurs because we pick the control that
16 maximizes the Q-value of a particular state x in (9.8). Effectively, if the
17 Q-value $q_{\varphi}(x', u)$ of a particular control $u \in U$ is an over-estimate, the
18 target will keep selecting this control as the maximizing control, which
19 drives up the value of the Q-function at $q_{\varphi}(x, u)$ as well. This problem is
20 a bit more drastic in the next section on continuous-valued controls. It is
21 however unclear how to best address this issue and design mathematically
22 sound methods that do not use arbitrary heuristics such as “pessimism”.

23 9.3 Q-Learning for continuous control spaces

24 All the methods we have looked at in this chapter are for discrete control
25 spaces, i.e., the set of controls that the robot can take is a finite set. In this
26 case we can easily compute the maximizing control of the Q-function.

$$u^*(x) = \operatorname{argmax}_u q_{\varphi}(x, u).$$

27 Certainly a lot of real-world problems have continuous-valued controls
28 and we therefore need Q-Learning-based methods to handle this.

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 **Deterministic policy gradient** A natural way, although non-rigorous,
 2 to think about this is to assume that we are given a Q-function $q_\varphi(x, u)$
 3 (we will leave the controller for which this is the Q-function vague for
 4 now) and a dataset $D = \{(x_t^i, u_t^i)_{t=0}^T\}_{i=1}^n$. We can find a deterministic
 5 feedback controller that takes controls that lead to good values as

$$\theta^* = \max_{\theta} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=0}^T q_\varphi(x_t^i, u_\theta(x_t^i)). \quad (9.10)$$

6 Effectively we are fitting a feedback controller that takes controls $u_{\theta^*}(x)$
 7 that are the maximizers of the Q-function. This is a natural analogue
 8 of the argmax over controls for discrete/finite control spaces. Again we
 9 should think of having a deep network that parametrizes the deterministic
 10 controller and fitting its parameters θ using stochastic gradient descent
 11 on (9.10)

$$\begin{aligned} \theta^{k+1} &= \theta^k + \eta \nabla_{\theta} q_\varphi(x^\omega, u_{\theta^k}(x^\omega)) \\ &= \theta^k + \eta (\nabla_u q_\varphi(x^\omega, u)) (\nabla_{\theta} u_{\theta^k}(x^\omega)) \end{aligned} \quad (9.11)$$

12 where ω is the index of the datum in the dataset D . The equality was
 13 obtained by applying the chain rule. This result is called the “deterministic
 14 policy gradient” and we should think of it as the limit of the policy gradient
 15 for a stochastic controller as the stochasticity goes to zero. Also notice
 16 that the term

$$\nabla_u q_\varphi(x^\omega, u)$$

17 is the gradient of the output of the Q-function $q_\varphi : X \times U \mapsto \mathbb{R}$ with
 18 respect to its second input u . Such gradients can also be easily computed
 19 using backpropagation in PyTorch. It is different than the gradient of the
 20 output with respect to its weights

$$\nabla_{\varphi} q_\varphi(x^\omega, u).$$

21 **On-policy deterministic actor-critic** Let us now construct an analogue
 22 of the policy gradient method for the case of a deterministic controller. The
 23 algorithm would proceed as follows. We initialize weights of a Q-function
 24 φ^0 and weights of the deterministic controller θ^0 .

- 25 1. At the k^{th} iteration, we collect a dataset from the robot using the
 26 latest controller u_{θ^k} . Let this dataset be D^k that consists of tuples
 27 (x, u, x', u') .
- 28 2. Fit a Q-function q^{θ^k} to this dataset by minimizing the temporal
 29 difference error

$$\varphi^{k+1} = \underset{\varphi}{\operatorname{argmin}} \sum_{(x, u, x', u') \in D^k} (q_\varphi(x, u) - r(x, u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) q_{\varphi'}(x', u'))^2. \quad (9.12)$$

30 Notice an important difference in the expression above, instead of
 31 using \max_u in the target, we are using the control that the current

1 controller, namely u_{θ^k} has taken. This is because we want to
 2 evaluate the controller u_{θ^k} and simply parameterize the Q-function
 3 using weights φ^{k+1} . More precisely, we hope that we have

$$q_{\varphi^{k+1}}(x, u_{\theta^k}(x)) \approx \max_u q^{\theta^k}(x, u).$$

4 3. We can now update the controller using this Q-function:

$$\theta^{k+1} = \theta^k + \eta \nabla_{\theta} q_{\varphi^{k+1}}(x^{\omega}, u_{\theta^k}(x^{\omega})) \quad (9.13)$$

5 This algorithm is called “on-policy SARSA” because at each iteration we
 6 draw fresh data D^k from the environment; this is the direct analogue of
 7 actor-critic methods that we studied in the previous chapter for deterministic
 8 controllers.

9 **Off-policy deterministic actor-critic methods** We can also run the
 10 above algorithm using data from an exploratory controller. The only
 11 difference is that the we now do not throw away the data D^k from older
 12 iterations

$$D = D^1 \cup \dots \cup D^k$$

13 and therefore have to change (9.12) to be

$$\varphi^{k+1} = \underset{\varphi}{\operatorname{argmin}} \sum_{(x, u, x', u') \in D} \left(q_{\varphi}(x, u) - r(x, u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) q_{\varphi'}(x', \underbrace{u_{\theta^k}(x')}_{\text{notice the difference}}) \right)^2. \quad (9.14)$$

14 Effectively, we are fitting the optimal Q-function using the data D but
 15 since we can no longer take the maximum over controls directly, we plug
 16 in the controller in the computation of the target. This is natural; we
 17 think of the controller as the one that maximizes the Q-function when we
 18 update (9.13). When used with deep networks, this is called the “deep
 19 deterministic policy gradient” algorithm, it is popular by the name DDPG.

❗ 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 . . .

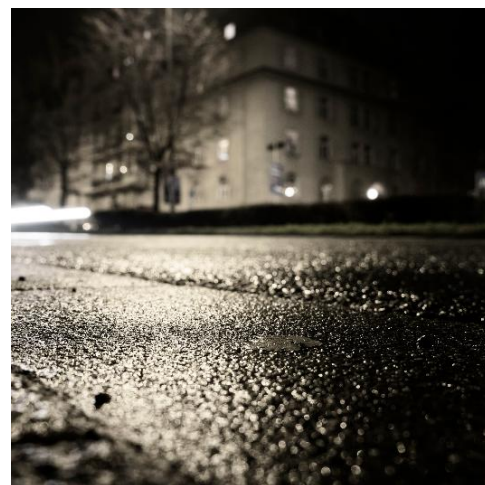
Chapter 10

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-and-intelligent-locomotion-via.html>. Very rarely you will see RL methods being used to train robots directly.

It makes sense to combine model-based and model-free methods if we want to reduce the number of data required from the system to learn a controller. Such methods are typically called model-based RL methods.

10.1 Learning a model of the dynamics

Imagine that we have a robot with dynamics

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

and obtained some data from this robot using an exploratory controller $u_e(\cdot | x)$. Let us call this dataset $D = \{(x_t^i, u_t^i)_{t=0}^T\}_{i=1}^n$; it consists of n trajectories each of length T timesteps. We can fit a deep network to learn the dynamics. This involves parameterizing the unknown dynamics using a deep network with weights w

$$f_w : X \times U \mapsto \mathbb{R}$$

and minimizing a regression error of the form

$$w^* = \operatorname{argmin}_w \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=0}^T \|x_{t+1}^i - f_w(x_t^i, u_t^i)\|_2^2. \quad (10.1)$$

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 we know that given some new control $u' \neq u_t^i$, we can, for instance, estimate the new future state $x' = f_w(x, u')$. In principle, we can use this model now to execute algorithms like iterated LQR to find an optimal controller. We could also imagine using this as our own simulator for the robot, i.e., instead of drawing new trajectories in model-free RL from the actual robot, we use our learned model of the dynamics to obtain more data.

An inverse model of the dynamics We can also learn what is called the inverse model of the system that maps the current state x_t^i and the next state x_{t+1}^i to the control that takes the system from the former to the latter:

$$f_w^{\text{inv}} : X \times X \mapsto U.$$

The regression error for one sample in this case would be $\|u_t^i - f_w^{\text{inv}}(x_t^i, x_{t+1}^i)\|_2^2$.

This is often a more useful model to learn from the data, e.g., if we want to use this model in a Rapidly Exploring Random Tree (RRT), we can sample states in the configuration space of the robot and have the learned dynamics guess the control between two states. Also see the paper on contact-invariant optimization (<https://homes.cs.washington.edu/~todorov/papers/-MordatchSIGGRAPH12.pdf>) and a video at https://www.youtube.com/watch?v=mhr_jtQrhVA

1 for an impressive demonstration of using an inverse model.

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^{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

$$D \leftarrow D \cup D^1.$$

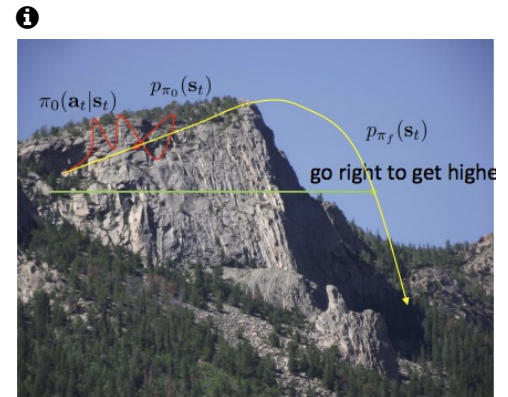
This is a simple mechanism that ensures that we can collect more data from the system. If the controller goes to parts of the state-space that the model is incorrect at, we get samples from such regions and iteratively improve both the learned dynamics model $f_w^k(x, u)$ and the controller $u^k(x)$ using this model.

2 10.2 Some model-based methods

3 10.2.1 Bagging multiple models of the dynamics

4 Let us look at bagging, which is a method to estimate the uncertainty of
5 a learning-based predictor. Bagging is short for *bootstrap aggregation*,

6 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 from Emilio Frazzoli (ETH/MIT).



1 and can be explained using a simple experiment. Suppose we wanted to
 2 estimate the average height μ of people in the world. We can measure the
 3 height of N individuals and obtain *one* estimate of the mean μ . This is of
 4 course unsatisfying because we know that our answer is unlikely to be the
 5 mean of the entire population. Bootstrapping computes multiple estimates
 6 of the mean μ_k over many *subsets* of the data N and reports the answer as

$$\mu := \text{mean}(\mu_k) + \text{stddev}(\mu_k).$$

7 Each subset of the data is created by sampling the original data with
 8 N samples *with replacement*. This is among the most influential ideas
 9 in statistics see “Bootstrap Methods: Another Look at the Jackknife”
 10 [https://projecteuclid.org/journals/annals-of-statistics/volume-7/issue-1/Bootstrap-](https://projecteuclid.org/journals/annals-of-statistics/volume-7/issue-1/Bootstrap-Methods-Another-Look-at-the-Jackknife/10.1214/aos/1176344552.full)
 11 [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) be-
 12 cause it is a very simple and general procedure to obtain the uncertainty
 13 of the estimate. Also see a very famous paper “Bagging predictors”
 14 at <https://link.springer.com/article/10.1007/BF00058655> that invented
 15 random forests based on this idea.

16 **Training an ensemble** We are going to train multiple dynamics models
 17 $f_{w^k}(x, u)$ for $k \in \{1, \dots, M\}$, one each for bootstrapped versions of the
 18 training dataset $\{D^1, \dots, D^M\}$. Each subset D^k is built by sampling a
 19 fraction, say 60% of the data uniformly randomly from our training dataset
 20 D . In other words, the bootstrapped versions of the data are not disjoint
 21 but their union is likely to be the entire training dataset. We are going to
 22 use this ensemble as follows. For each pair (x, u) , we run all models and
 23 set

$$\hat{x}' = \frac{1}{M} \sum_{k=1}^M f_{w^k}(x, u);$$

24 i.e., the ensemble predicts the next state of the robot using the mean. The
 25 important benefit of using an ensemble is that we can also get an estimate
 26 of the error in these predictions

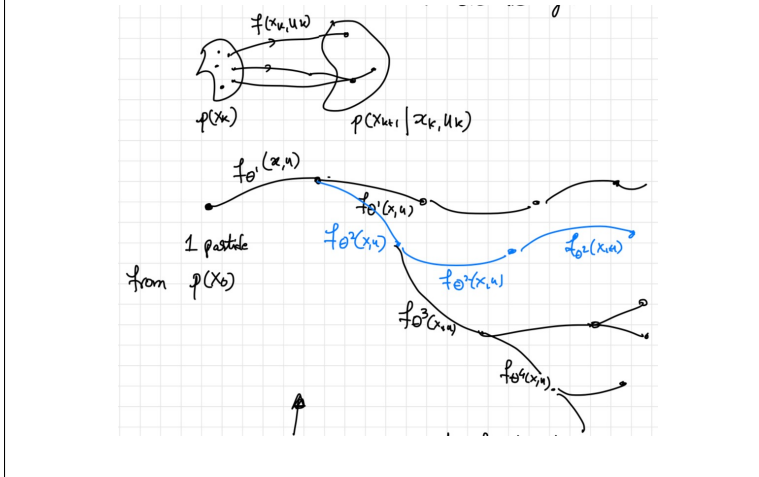
$$\text{error in } \hat{x}' = (\delta_k f_{w^k}(x, u))^{1/2}.$$

27 Different members of the ensemble are training on different datasets
 28 and make different predictions as to what the next state could be. The
 29 mis-match between them is an indicator of the error in our dynamics model.
 30 This need not be an accurate estimate of the error (i.e., the difference
 31 between the predicted \hat{x}' and the actual next state x' of the true dynamics)
 32 but is often a good proxy to use if we do bootstrapping.

33 Bagging is perhaps the most useful idea in machine learning (by far).
 34 It is always good to keep it in your mind. The winners of most high-profile
 35 machine learning competitions, e.g., the Netflix Prize
 36 (https://en.wikipedia.org/wiki/Netflix_Prize) or the ImageNet challenge,
 37 have been bagged classifiers created by fitting multiple architectures on
 38 the same dataset.

❗ 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.

How to use an ensemble for hallucinating new data The procedure is 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** A powerful Gaussian Process-based algorithm to incorporate
- 2 uncertainty in the predictions of a learned dynamics model is “PILCO: A
- 3 Model-Based and Data-Efficient Approach to Policy Search”
- 4 <http://mlg.eng.cam.ac.uk/pub/pdf/DeiRas11.pdf>. Instead of using boot-
- 5 strapping of an ensemble to estimate the uncertainty, this algorithm
- 6 explicitly models the uncertainty as

$$p(x_{k+1} | x_k, u_k) = N(x_k + E[\Delta_k], \text{Var}(\Delta_k))$$

- 7 where $\Delta_k = x_{k+1} - x_k$ in the training data, using a Gaussian Process.
- 8 See a preliminary but great tutorial on Gaussian Processes at [https://distill.pub/2019/visual-](https://distill.pub/2019/visual-exploration-gaussian-processes)
- 9 [exploration-gaussian-processes](https://distill.pub/2019/visual-exploration-gaussian-processes). PILCO is a complicated algorithm to
- 10 implement but you can see the source code by the original authors at
- 11 <https://mloss.org/software/view/508>. You can also look at the paper ti-
- 12 tled “BayesRace: Learning to race autonomously using prior experience”
- 13 <https://arxiv.org/abs/2005.04755> which uses model-based RL in

14 10.2.2 Model-based RL in the latent space

Chapter 11

Offline Reinforcement Learning

Reading

1. Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems by [Levine et al. \(2020\)](#)
2. Continuous doubly constrained batch reinforcement learning by [Fakoor et al. \(2021\)](#)

So far, we have imagined that we have access to a model of a robot, a simulator for it, or access to the actual robot that allows us to obtain data from the robot. But there are many problems in which we cannot get any of these. For example, when Amazon sells merchandise to its customers, it is quite difficult for them to model or simulate each customer, or even a canonical one. It is possible to do rollouts using actual customer data but that would not be wide because an exploratory policy, by the time it learns, would lead to huge losses. This would also not be desirable for the customers. Amazon also needs reinforcement learning for this problem because there is not many other ways to explore what customers like and do not like. One could take another example from a very different domain. Suppose a hospital is trying to develop a new protocol to handle incoming patients. E.g., a patient comes into ER, there is a fixed set of checks that are performed quickly on them called as “triage”. The number and kind of checks are performed directly affects patient care. If there are too many checks, then the patient loses precious time before they are treated. If there are too few, then there is a large bottleneck when these patients are referred to doctors. It is not easy to model or simulate this problem. It is possible to do rollouts to discover a policy but that would not be easy, or wise.

These problems have a few commonalities.

- We do not have models or simulators for such systems *and that is why* we need to use ideas from reinforcement learning to build controllers for them;
- There are both existing systems, i.e., there exists a controller that is currently deployed. This controller may be very complex, e.g., in Amazon’s case it is the result of many scientists building this system over a decade (and thereby even if one could look at it, there is no way one would be able to model/simulate it). In the case of the hospital, the current triage policy was likely created by the doctors over experience and refined by the triage nurses by looking at actual cases. Because the system evolved over a long time, a lot of this knowledge is not accessible in a codified form.

Offline reinforcement learning are a suite of techniques that allow us to learn the optimal value function from data that need not be coming from an optimal policy—without drawing any new data from the environment. Note that we do not just want to evaluate an existing policy, we want to learn the optimal value function, or at the very least improve upon the current policy.

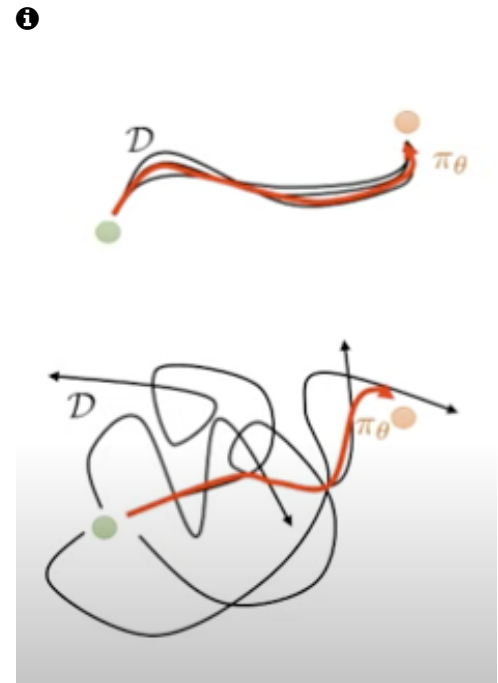
There are many other problems of this kind: data is plentiful, just that we cannot get more.

Technically speaking, offline learning is a very clean problem in that we are close to the supervised learning setting (we do not know the true targets). A meaningful theoretical analysis of typical reinforcement learning algorithms is difficult because there are a lot of moving parts in the problem definition: exploratory controllers, the fact that we are adding correlated samples into our dataset as we draw more trajectories, function approximation properties of the neural network that does not allow Bellman iteration to remain a contraction etc. Some of these hurdles are absent in the analysis of offline learning.

11.1 Why is offline reinforcement learning difficult?

Suppose that we wanted to use behavior cloning for this problem. After all, we could build a state vector x and do behavior cloning with a neural network to learn a policy $u_\theta(x)$ from the existing dataset $D = \{(x_t^i, u_t^i)_{t=0, \dots, T}\}_{i=1}^n$. In principle, we can also learn the value function corresponding to this policy, $q_\varphi^\theta(x, u)$ where u is the control input. Note that this is not the optimal value function. We will call this the behavior cloning solution to offline reinforcement learning.

We know further that value iteration converges (for tabular MDPs) from any initial condition. In lieu of the actual model of the Markov Decision Process, we can imagine using the entire dataset D to calculate

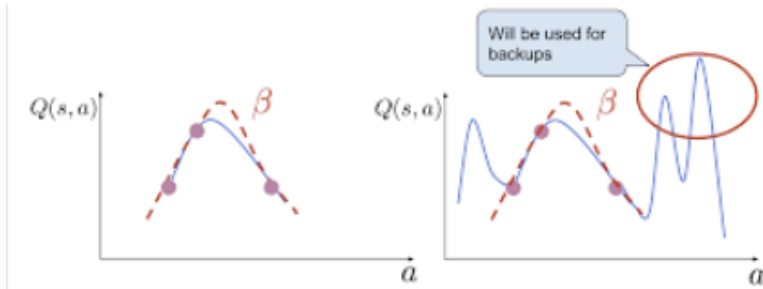


A good intuition for offline reinforcement learning is given by this picture. In imitation learning, or behavior cloning, we simply want to learn a policy that mimics the expert (or recorded data). In offline reinforcement learning, we can “stitch together” multiple sub-optimal policies in different parts of the domain.

1 Bellman updates for a value function parameterized by a neural network:

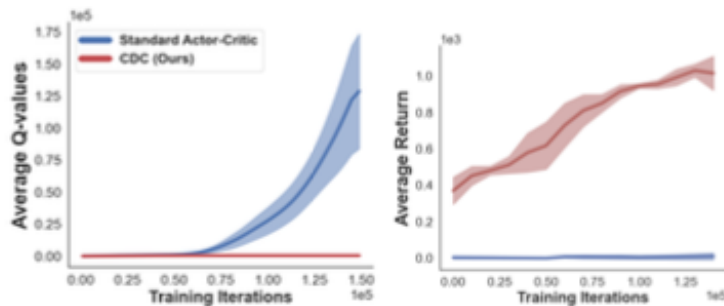
$$\min_{\varphi} \sum_{i,t} \left(q_{\varphi}^{\theta^{(k+1)}}(x_t^i, u_t^i) - r(x_t^i, u) - \gamma \max_{u' \in U} q_{\varphi}^{\theta^{(k)}}(x_{t+1}^i, u') \right)^2.$$

2 This approach is unlikely to work very well. Observe the following figure.



3

4 We do not really know whether the initial value function q_{φ}^{θ} assigns large
 5 returns to controls that are outside of the ones in the dataset. If it does, then
 6 these controls will be chosen in the maximization step while calculating
 7 the target. If there are states where the value function over different
 8 controls looks like this picture, then their targets will cause the value at
 9 all other states to grow unbounded during training. This is exactly what
 10 happens in practice. For example,



11

12 In offline reinforcement learning, this phenomenon is often called the
 13 “extrapolation error”. It arises because we do not know a natural way to
 14 force the network to avoid predicting large values for controls that are not
 15 a part of the training dataset. It is instructive to ask why off-policy or
 16 on-policy reinforcement learning does not suffer from extrapolation error.
 17 Both of these algorithms explicitly draw more data from the simulator. If
 18 the value function were to over-estimate the value of certain control actions
 19 at a state, then the controller would take those controls during exploration
 20 and discover that the value was in fact an over-estimate. Methods that
 21 draw more data have this natural self-correcting behavior that we cannot
 22 get in offline learning.

23 There is a second problem associated with computing the maximum
 24 in value iteration. For problems with continuous controls, we do not know
 25 of computationally effective ways to compute the maximum $\max_{u \in U}$.

i We have discussed how Bellman updates are a set of consistent constraints on the optimal value function. This entails that if we over-estimate the value at a given state, then the estimated return to come (i.e., the value) of *all* other states becomes incorrect.

1 Typical implementations of offline learning fit a controller that maximizes
2 the value function

$$\theta = \operatorname{argmax} \frac{1}{nT} \sum_{i,t} q_{\varphi}^{\theta}(x_t^i, u_{\theta}(x))$$

3 to make it easy to calculate this maximum. But this forces us to use
4 another network (and it is not a given that the parameterized controller
5 can calculate the correct maximum).

6 11.2 Regularized Bellman iteration

7 There are two broad class of techniques that are believed to give reasonable
8 results for offline learning. These are both quite new and ad hoc and
9 effective offline learning is essentially an open problem today. To wit,
10 current offline learning methods not only fail to learn optimal value
11 functions or policies from sub-optimal data, but they often do any better
12 than behavior cloning.

13 11.2.1 Changing the fixed point of the Bellman iteration 14 to be more conservative

15 Since extrapolation error is fundamentally caused by the value function
16 taking large values, it is reasonable strategy to modify the Bellman updates
17 to regularize the value function in some way. A basic version would look
18 like a coupled system of updates

$$\begin{aligned} \varphi^* &= \min_{\varphi} \frac{1}{nT} \sum_{i,t} \left(q_{\varphi}^{\theta^{(k+1)}}(x_t^i, u_t^i) - r(x_t^i, u) - \gamma q_{\varphi}^{\theta^{(k)}}(x_{t+1}^i, u_{\theta}(x_{t+1}^i)) \right)^2 \\ &\quad - \lambda \Omega(q_p^{\theta}) \\ \theta^* &= \min_{\theta} \frac{1}{nT} \sum_{i,t} q_{\varphi}^{\theta}(x_t^i, u_{\theta}(x_t^i)). \end{aligned} \tag{11.1}$$

19 We will use the regularizer

$$\Omega = \frac{1}{nT} \sum_{i,t} (q_{\varphi}^{\theta}(x_t^i, u_{\theta}(x_t^i)) - q_{\varphi}^{\theta}(x_t^i, u_t^i))_+^2$$

20 The notation $(\cdot)_+$ denotes rectification, i.e., $(x)_+ = x$ if $x > 0$ and zero
21 otherwise. Notice that second term in the objective for fitting the value
22 function forces the value of the control $u_{\theta}(x_t^i)$ to be smaller than the
23 value of the control u_t^i taken at the same state in the dataset. This is a
24 conservative choice. While it does not fix the issue of extrapolation error,
25 it forces the value network to predict smaller values and prevents it from
26 blowing up.

❗ One can use a linear value function, e.g., $q_{\varphi}^{\theta}(x, u) = \langle [1, x, u], \varphi[1, x, u] \rangle$ or any other set of basis functions. But this is not enough to resolve the issue in theory. Bellman updates are not a contraction when the TD-error is projected onto a set of bases. In practice, this approach works reasonably well.

1 A second, similar, strategy looks like

$$\Omega = \frac{1}{nT} \sum_{i,t} \max_{u' \sim u_\theta(x_t^i)} (q_\varphi^\theta(x_t^i, u') - q_\varphi^\theta(x_t^i, u_t^i))_+^2$$

2 where the maximum of actions is computed over a large number of samples.

3 A yet another strategy looks like

$$\Omega = \frac{1}{nT} \sum_{i,t} \log \int_u \exp(q_\varphi^\theta(x_t^i, u)).$$

4 These strategies have been found to be somewhat useful in the sense that
 5 they prevent the value function from taking large values. But it is also
 6 clear that the solution of the problems is not the optimal value function.

7 **11.2.2 Estimating the uncertainty of the value function**

8 It is reasonable to ask the question whether initializing the value function

i Practically speaking, it is reasonable to expect that even if we do not find the optimal value function, if we can obtain a better policy than the existing policy that is running on the system, then offline reinforcement learning is a viable approach.

Chapter 12

Meta-Learning

Reading

1. Learning To Learn: Introduction (1996),
<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.21.3140>
2. Prototypical Networks for Few-shot Learning
<http://papers.nips.cc/paper/6996-prototypical-networks-for-few-shot-learning>
3. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks <https://arxiv.org/abs/1703.03400>
4. A Baseline for Few-Shot Image Classification
<https://arxiv.org/abs/1909.02729>
5. Meta-Q-Learning <https://arxiv.org/abs/1910.00125>

The human visual system is proof that we do not need lots of images to learn to identify objects or lots of experiences to learn about concepts. Consider the mushrooms shown in the image below.



The one on the left is called *Muscaria* and you'd be able to identify the bright spots on this mushroom very easily after looking at one image. The differences between an edible one in the center (*Russula*) and the one on the right (*Phalloides*) may sometimes be subtle but a few samples from each of them are enough for humans to learn this important distinction.

1 There are also more everyday examples of this phenomenon. You
 2 touched a hot stove once as a child and have forever learnt not to do it.
 3 You learnt to ride a bike as a child and only need a few minutes on a
 4 completely new bike to be able to ride it these days. At the same time, you
 5 could start learning to juggle today and will be able to juggle 3 objects
 6 with a couple of days of practice.

7 The hallmark of human perception and control is the ability to gen-
 8 eralize. This generalization comes in two forms. The first is the ability
 9 to do a task better if you see more samples from the same task; this is
 10 what machine learning calls generalization. The second is the ability to
 11 mix-and-match concepts from previously seen tasks to do well on new
 12 tasks; doing well means obtaining a lower error/higher reward as well
 13 as learning the new task quickly with few days. This second kind is the
 14 subject of what is called “learning to learn” or meta-learning.

15 Standard machine learning \Rightarrow generalization across samples

16 Meta-Learning \Rightarrow generalization across tasks

What is a task? If we are going to formalize meta-learning, we better define what a task is. This is harder than it sounds. Say we are doing image classification, classifying cats vs. dogs could be considered Task 1; Task 2 could be classifying apples vs. oranges. It is reasonable to expect that learning low-level features such as texture, colors, shapes etc. while learning Task 1 could help us to do well on Task 2.

17 This is not always the case, two tasks can also fight each other. Say,
 18 you design a system to classify ethnicities of people using two kinds
 19 of features. Task 1 uses the shape of the nose to classify Caucasians
 20 (long nose) vs. Africans (wide nose). Task 2 uses the kind of hair to
 21 classify Caucasians (straight hair) vs. Africans (curly hair). An image of
 22 a Caucasian person with curly hair clearly results in two tasks fighting
 23 each other.

24 The difficulty in meta-learning begins with defining what a task is.

25 While understanding what a task is may seem cumbersome but doable
 26 for image classification, it is even harder for robotics systems.



27
 28 We can think of two different kinds of tasks.

- 1 1. The first, on the left, is picking up different objects like a soccer
2 ball, a cup of coffee, a bottle of water etc. using a robot manipulator.
3 We may wish to learn how to pick up a soccer ball quickly given
4 data about how to pick up the bottle.
- 5 2. The second kind of tasks is shown on the right. You can imagine
6 that after building/training a model for the robot in your lab you
7 want to install it in a factory. The factory robot might have 6 degrees-
8 of-freedom whereas the one in your lab had only 5; your policy
9 better adapt to this slightly different state-space. An autonomous
10 car typically has about 10 cameras and 5 LIDARs, any learning
11 system on the car better adapt itself to handle the situation when
12 one of these sensors breaks down. The gears on a robot manipulator
13 will degrade over the course of its life, our policies should adapt to
14 this degrading robot.

15 Almost all the current meta-learning/meta-reinforcement learning
16 literature focuses on developing methods to do the first set of tasks. The
17 second suite of tasks are however more important in practice. In the
18 remainder of this chapter, we will discuss two canonical algorithms to
19 tackle these two kinds of tasks.

20 12.1 Problem formulation for image classifica- 21 tion

22 The image classification formulation thinks of each class/category as a
23 “task”.

24 Consider a supervised learning problem with a dataset $D = \{(x^i, y^i)\}_{i=1, \dots, N}$.
25 The labels $y^i \in \{1, \dots, C\}$ for some large K and there are $\frac{N}{C}$ samples in
26 the dataset for each class. Think of this as a large dataset of cars, cats, dogs,
27 airplanes etc. all objects that are very frequent in nature and for which
28 we can get lots of images. This training set is called the meta-training set
29 with C “base tasks”.

30 Say, we were simply interested in obtaining a machine learning model
31 to classify this data. Let us denote the parameters of this model by w . If
32 this model predicts the probability of the input x belonging to each of
33 these known classes we can think of maximizing the log-likelihood of the
34 data under the model (or minimizing the cross-entropy loss)

$$\hat{w} = \operatorname{argmax}_w \frac{1}{N} \sum_{i=1}^N \log p_w(y^i | x^i). \quad (12.1)$$

35 This is the standard multi-class image classification setup. Since we like
36 to think of one task as one category, this is also the multi-task learning
37 setup. The model

$$p_{\hat{w}}(\cdot | x)$$

🔗 Meta-learning vs. multi-task learning

We have talked about adaptation as the way to handle new tasks in the previous remark.

Consider the following situation: in standard machine learning, we know that larger the size of the training data we collect better the performance on the test data; a large number of images help capture lots of variability in the data, e.g., dogs of different shapes, sizes and colors. You can imagine then that in order to do well on lots of different tasks, i.e., meta-learn, we should simply collect data from lots of different tasks. Can you think as to why mere multi-task learning may not work well for meta-learning?

1 after fitting on the training data will be good at classifying some new input
 2 image x as to whether it belongs to one of the C training classes. Note that
 3 we have written the model as providing the probability distribution $p_{\hat{w}}(\cdot |$
 4 $x)$ as the output, one real-valued scalar per candidate class $\{1, \dots, C\}$.

5 12.1.1 Fine-tuning

6 Let us now consider the following setup. In addition to our original dataset
 7 of the base tasks, we are given a “few-shot dataset” that has c new classes
 8 and s labeled samples per class, a total of $n = cs$ new samples

$$D' = \{(x^i, y^i)\}_{i=N+1, \dots, N+n}; \quad y^i \in \{C+1, \dots, C+c\}.$$

9 The words “few-shot” simply mean that s is small, in particular we are
 10 given much fewer images per class than the meta-training dataset,

$$\frac{n}{c} = s \ll \frac{N}{C}.$$

11 This models the situation where the model is forced to classify images
 12 from rare classes, e.g., the three kinds of strawberries grown on a farm in
 13 California after being trained on data of cars/cats/dogs/planes etc.

14 We would like to adapt the parameters \hat{w} using this labeled few-shot
 15 data. Here is one solution, we simply train the model again on the new
 16 data. This looks like solving another optimization problem of the form

$$w^* = \operatorname{argmin}_w -\frac{1}{n} \sum_{i=N+1}^{N+n} \log p_w(y^i | x^i) + \frac{\lambda}{2} \|w - \hat{w}\|_2^2. \quad (12.2)$$

17 The new parameters w can potentially do well on the new classes even if
 18 the shot s is small because training is initialized using the parameters \hat{w} .
 19 We write down this initialization using the second term

$$\frac{\lambda}{2} \|w - \hat{w}\|_2^2$$

20 which keeps the parameters being optimized w closed to their initialization
 21 using a quadratic spring-force controlled by the parameter λ . We can
 22 expect the new model p_{w^*} to perform well on the new classes if the
 23 initialization \hat{w} was good, i.e., if the new tasks and the base tasks were
 24 close to each other. This method is called fine-tuning, it is the easiest trick
 25 to implement to handle new classes.

26 12.1.2 Prototypical networks

27 The cross-entropy objective used in (12.1) to train the model $p_{\hat{w}}$ simply
 28 maximizes the log-likelihood of the labels given the data. It is reasonable
 29 to think that since the base classes are not going to show up as the few-shot
 30 classes, we should not be fitting to this objective.

🔗 Say, we are interested in classifying images from classes that are different than those in the training set. The model has only C outputs, effectively the universe is partitioned into C categories as far as the model is concerned and it does not know about any other classes. How should one formalize the problem of meta-learning then?

🔗 Think of a multi-layer neural network \hat{w} that has K outputs. The new network should now produce m outputs, how should we modify this network?

The idea behind a prototypical loss is to train the model to be a good discriminator among different classes.

1 Let us imagine the features of the model, e.g., the activations of the
2 last layer in the neural network,

$$z = \varphi_w(x)$$

3 for a particular image x . Note that the features z depend on the parameters
4 w . During standard cross-entropy training, there is a linear layer on top of
5 these features and our prediction probability for class y is

$$p_w(y | x) = \frac{e^{w_y^\top z}}{\sum_{y'} e^{w_{y'}^\top z}}$$

6 where $w_y \in \mathbb{R}^{\dim(z)}$. This is the softmax operation and the vectors w are
7 the weights of the last layer of the network; when we wrote (12.1) we
8 implicitly folded those parameters into the notation w .

9 Prototypical networks train the model to be a discriminator as follows.

10 1. Each mini-batch during training time consists of a few-shot task
11 created out of the meta-training set by sub-sampling.

$$\begin{aligned} D^{\text{episode}} &= D^{\text{support}} \cup D^{\text{query}} \\ &= \{(x^i, y^i); y^i \in \{1, \dots, C\}\}_{i \in \{1, \dots, Cs\}}; \\ &\quad \cup \{(x^i, y^i); y^i \in \{1, \dots, C\}\}_{i \in \{1, \dots, Cq\}}. \end{aligned}$$

12 with $|D^{\text{support}}| = Cs$ and $|D^{\text{query}}| = Cq$. This is called an “episode”
13 by researchers in this literature. Each episode comes with some
14 more data from the same classes called the “query-shot” in this
15 literature. The query-shot is akin to the data from the new classes
16 that the model is forced to predict during adaptation time. Let us
17 have q query-shot per class in each episode.

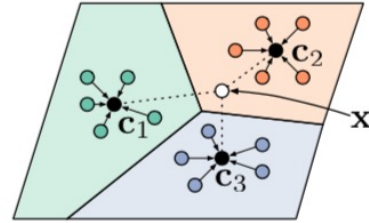
18 2. We know the labels of the $N = Cs$ labeled data and can compute
19 the prototypes, which are simply centroids of the features,

$$\mu_y = \frac{1}{s} \sum_{(x^i, y^i) \in D^{\text{episode}}} \mathbf{1}_{\{y^i=y\}} \varphi_w(x^i).$$

20 3. You can now impose a clustering loss to force the query samples to
21 be labeled correctly, i.e., maximize

$$p_{w, \mu_y}(y | x) = \frac{e^{-\|\varphi_w(x) - \mu_y\|_2}}{\sum_{y'} e^{-\|\varphi_w(x) - \mu_{y'}\|_2}}$$

22 where $y = y^i$ and $x = x^i$ for each of the samples (x^i, y^i) in the
23 query-set of the episode.



(a) Few-shot

- 1
- 2 4. The objective maximized at each mini-batch is

$$\frac{1}{Cq} \sum_{(x,y) \in D^{\text{query}}} \log p_{w, \mu_{y^i}}(y^i | x^i).$$

3 Note that the gradient of the above expression is both on the weights
4 w of the top layer and the weights w of the lower layers.

- 5 5. We can now use the trained model for classifying new classes by
6 simply feeding the new images through the network, computing the
7 prototypes using the few labeled data and computing the clustering
8 loss on the unlabeled query data at test time to see which prototype
9 the feature of a particular query datum is closest to.

10 **Discussion** Prototypical loss falls into the general category of metric-
11 based approaches to few-shot learning. We make a few remarks next.

- 12 1. It is a very natural setting for learning representations of the data for
13 classification that can be transferred easily. If the model is going to
14 be used for new classes, it seems reasonable that the prototypes of
15 the new classes should be far away from each other and the z s of the
16 query samples should be clustered around their correct prototypes.
- 17 2. Prototypical networks perform well if you can estimate the proto-
18 types accurately. In practice, this requires that you have about 10
19 labeled data per new class.
- 20 3. We used the ℓ_2 metric $\|\cdot\|_2$ in the z -space to compute the affinities
21 of the query samples. This may not be a reasonable metric to use for
22 some problems, so a large number of approaches try to devise/learn
23 new metrics.
- 24 4. A key point of prototypical networks is that there is no gradient-based
25 learning going on upon the new categories; we simply compute
26 the prototypes and the affinities and use those to classify the new
27 samples.

28 12.1.3 Model-agnostic meta-learning (MAML)

29 We will next look at a simple algorithm for gradient-based adaptation of
30 the model on the new categories. The key idea is to update the model

1 using the same objective in (12.1) but avoid overfitting the model on the
 2 meta-training data so that the model can be quickly adapted using the
 3 few-shot data via gradient-updates.

4 Here we consider an episode $D^{\text{episode}} = D^{\text{support}}$ and $D^{\text{query}} = \emptyset$, i.e.,
 5 there are no query shots. Let us define

$$\ell(w; D^{\text{support}}) = \frac{1}{Cs} \sum_{(x^i, y^i) \in D^{\text{support}}} \log p_w(y^i | x^i);$$

6 this is the same objective as that in (12.1) so if we maximized the objective

$$\ell(w; D^{\text{support}})$$

7 we will perform standard cross-entropy training. At each mini-batch/episode,
 8 the MAML algorithm instead maximizes the objective

$$\ell_{\text{maml}}(w; D^{\text{support}}) = \ell(w + \alpha \nabla \ell(w; D^{\text{support}}); D^{\text{support}}). \quad (12.3)$$

9 In other words, MAML uses a “look ahead” gradient: the gradient of
 10 $\ell(w; D^{\text{support}})$ is not in the steepest ascent direction of $\ell(w; D^{\text{support}})$
 11 but in the steepest ascent direction after one update of the parameters
 12 $w + \alpha \nabla \ell(w; D^{\text{support}})$.

13 **Adaptation on the few-shot data** One we have a model trained using
 14 MAML

$$\hat{w} = \underset{w}{\operatorname{argmax}} \ell_{\text{maml}}(w; D)$$

15 we can update it on new data simply by maximizing the standard cross-
 16 entropy objective again, i.e.,

$$w^* = \underset{w}{\operatorname{argmax}} \frac{1}{n} \sum_{i=N+1}^{N+n} \log p_w(y^i | x^i) - \frac{1}{2\lambda} \|w - \hat{w}\|_2^2. \quad (12.4)$$

🔗 How does look-ahead in MAML help?

17 The adaptation phase is exactly the same as standard cross-entropy training.

18 **MAML as an approximation of a second order optimization method**

19 MAML is not specific to few-shot learning. We can use the MAML
 20 objective for any other standard supervised learning problem, is this going
 21 to help? Indeed it will, gradient descent/stochastic gradient descent are
 22 myopic algorithms because they update parameters only in the direction
 23 of the steepest gradient, you can potentially do better by computing the
 24 lookahead gradient. The caveat is that is computationally difficult to
 25 compute the lookahead gradient. Observe that

$$\begin{aligned} \ell_{\text{maml}}(w) &= \ell(w + \alpha \nabla \ell(w)) \\ &\approx \ell(w) + \alpha (\nabla \ell(w))^\top \nabla \ell(w) \\ \Rightarrow \nabla \ell_{\text{maml}}(w) &= \nabla \ell(w) + 2\alpha \nabla^2 \ell(w) \nabla \ell(w). \end{aligned}$$

1 So MAML is secretly a second-order optimization method, computing the
 2 gradient of the MAML objective requires having access to the Hessian of
 3 the objective $\nabla^2 \ell(w)$. For large models such as neural networks this is
 4 very expensive to compute.

5 **Remark 12.1.** Let us consider a meta-training set with two mini-batches/episodes/tasks,
 6 $D = D^1 \cup D^2$. The MAML algorithm uses the gradient

$$\begin{aligned} \nabla \ell_{\text{maml}}(w; D) &= \nabla \frac{1}{2} \sum_{i=1}^2 \ell_{\text{maml}}(w; D^i) \\ &= \frac{1}{2} \sum_{i=1}^2 \nabla \ell(w; D^i) + 2\alpha \nabla^2 \ell(w; D^i) \nabla \ell(w; D^i) \end{aligned}$$

7 Observe now that if there exist parameters w that have $\nabla \ell(w; D^i)$ for all
 8 the episodes D^i then the MAML gradient is also zero. In other words, if
 9 there exist parameters w that work well for all tasks then MAML may find
 10 such parameters. However, in this case, the simple objective

$$\ell_{\text{multi-task}}(w; D) = \frac{1}{2} \sum_{i=1}^2 \ell(w; D^i) \quad (12.5)$$

11 that sums up the losses of all the mini-batches/episodes/tasks will also
 12 find these parameters. This objective known as the multi-task learning
 13 objective is much simpler than MAML's because it requires only the
 14 first-order gradient.

🔗 What happens if the two tasks are different (as is likely to be the case) in which case there don't exist parameters that work well for all the tasks?

15 12.2 Problem formulation for meta-RL

16 One mathematical formulation of meta-RL is as follows. Let k denote a
 17 task and there is an underlying (unknown) dynamics for this task given by

$$x_{t+1}^k = f^k(x_t^k, u_t^k, \xi_t)$$

18 We will assume that all the tasks have a shared state-space $x_t^k \in X$ and a
 19 shared control-space $u_t^k \in U$. The reward function of each task is different
 20 $r^k(x, u)$ but we are maximizing the same infinite-horizon discounted
 21 objective for each task. The q -function is then defined to be

$$q^{k, \theta^k}(x, u) = \mathbb{E}_{\xi(\cdot)} \left[\sum_{t=0}^{\infty} \gamma^t r^k(x_t, u_t) \mid x_0 = x, u_0 = u, u_t = u_{\theta^k}(x_t) \right].$$

22 where $u_{\theta^k}(x_t)$ is a deterministic controller for task k . Given all these
 23 meta-training tasks, our objective is to learn a controller that can solve a
 24 new task $k \notin \{1, \dots, K\}$ upon being presented a few trajectories from
 25 the new task. Think of you learning to pick up different objects during
 26 training time and then adapting to picking up a new object not in the
 27 training set.

1 Let us consider the off-policy Q-learning setting and learn separate
 2 controllers for all the tasks for now. As usual, we want the q -function
 3 to satisfy the Bellman equation, i.e., if we are using parameters φ^k to
 4 approximate the q -function, we want to find parameters φ^k such that

$$\operatorname{argmin}_{\varphi^k} \mathbb{E}_{(x,u,x') \in D^k} \left[\left(q_{\varphi^k}^{k,\theta^k}(x,u) - r^k(x,u) - \gamma q_{\varphi^k}^{k,\theta^k}(x',u_{\theta^k}(x')) \right)^2 \right] \quad (12.6)$$

5 where the dataset D^k is created using some exploratory policy for the task
 6 k . The controllers u_{θ^k} are trained to behave like the greedy policy for the
 7 particular q -function $q_{\varphi^k}^{k,\theta^k}$

$$\operatorname{argmax}_{\theta^k} \mathbb{E}_{(x,u) \in D^k} \left[q_{\varphi^k}^{k,\theta^k}(x, u_{\theta^k}(x)) \right]. \quad (12.7)$$

8 The above development is standard off-policy Q-learning and we
 9 have seen it in earlier lectures. The different controllers u_{θ^k} do not
 10 learn anything from each other, they are trained independently on their
 11 own datasets. We can now construct a multi-task learning objective for
 12 meta-RL, in this we will learn a single q -function and a single controller
 13 for all tasks. We modify (12.6) and (12.7) to simply work on all the
 14 datasets together

$$\begin{aligned} \operatorname{argmin}_{\varphi} \mathbb{E}_{(x,u,x') \in D^1 \cup D^2 \dots} \left[\left(q_{\varphi}^{\theta}(x,u) - r^k(x,u) - \gamma q_{\varphi}^w(x',u_w(x')) \right)^2 \right] \\ \operatorname{argmax}_w \mathbb{E}_{(x,u) \in D^1 \cup D^2 \dots} \left[q_{\varphi}^w(x, u_w(x)) \right] \end{aligned} \quad (12.8)$$

15 This is the multi-task learning objective for RL. This is unlikely to work
 16 well because depending upon the task, the controllers for the different
 17 tasks will conflict with each other, it is unlikely that there is a single set
 18 of parameters for the controller and the q -function that works well for all
 19 tasks.

20 12.2.1 A context variable

21 Reinforcement Learning offers a very interesting way to solve the few-
 22 shot/meta-learning problem. We can append the state-space to include
 23 a context variable that is a representation of the particular task. Let us
 24 construct features for a trajectory using a set of basis functions

$$\{\phi_1(x, u, r), \phi_2(x, u, r), \dots, \phi_m(x, u, r)\}$$

25 We will now construct a variable $\mu^k(\tau)$ for a trajectory $\tau_{0:t} = (x_0, u_0, \dots, x_t, u_t, \dots)$
 26 from task k as

$$\mu(\tau_{0:t}) = \sum_{s=0}^t \sum_{i=1}^m \gamma^t \alpha_i \phi_i(x_s^k, u_s^k, r^k(x_s^k, u_s^k)).$$

🔗 Imagine a planning task with multiple goals. The optimal trajectory that goes to one goal location will bifurcate from the optimal trajectory that goes to some other goal. We will never be able to learn a controller that goes to both goals using one neural network. How to fix this? You can use MAML certainly to under-fit the controller and the q -function to all the tasks and then adapt them using some data from the new task using gradient updates.

1 We will call this variable a “context” because we can use it to guess which
 2 task a particular trajectory is coming from. **It is important to note that**
 3 **the mixing coefficients α_i are shared across all the tasks.** We would
 4 like to think of this feature vector $\mu(\tau)$ as a kind of indicator of whether a
 5 trajectory τ belongs to the task k or not. We now learn a q -function and
 6 controller that also depend on $\mu(\tau)$

$$\begin{aligned} q_\varphi^\theta(x_t, u_t, \mu(\tau_{0:t})) \\ u_\theta(x_t, \mu(\tau_{0:t})). \end{aligned} \quad (12.9)$$

7 Including a context variable like $\mu(\tau)$ allows the q -function to detect the
 8 particular task that it is being executed for using the past t time-steps of
 9 the trajectory $\tau_{0:t}$. This is similar to learning independent q -functions
 10 $q_{\varphi_k}^{k, \theta^k}$ and controllers u_{θ^k} but there is parameter sharing going on in (12.9).
 11 We will still solve the multi-task learning problem like (12.8) but also
 12 optimize the parameters α_i s that combine the basis functions.

$$\begin{aligned} \operatorname{argmin}_{\varphi, \alpha_i} \sum_{i=1}^K \mathbb{E}_{(x, u, x') \in D^k} \left[\left(q_\varphi^\theta(x, u, \mu(\tau)) - r^k(x, u) - \gamma q_\varphi^\theta(x', u_\theta(x', \mu(\tau))) \right)^2 \right] \\ \operatorname{argmax}_{\theta, \alpha_i} \sum_{i=1}^K \mathbb{E}_{(x, u) \in D^k} \left[q_\varphi^\theta(x, u_\theta(x, \mu(\tau)), \mu(\tau)) \right]. \end{aligned} \quad (12.10)$$

13 The parameters α_i of the context join the q -functions and the controllers
 14 of the different tasks together but also allow the controller the freedom to
 15 take different controls depending on which task it is being trained for.

16 **Adapting the meta-learned controller to a new task** Suppose we
 17 trained on K tasks using the above setup and have the parameters
 18 $\hat{\theta}, \hat{\varphi}, \{\hat{\alpha}_i\}_{i=1, \dots, m}$ in our hands. How should we adapt to a new task? This
 19 is easy, we can run an exploration policy on the new task (the current policy
 20 u_θ will work just fine if the control space U is the same) to collect some
 21 data and update our off-policy Q-learning parameters $\hat{\theta}, \hat{\varphi}$ on this data
 22 using (12.6) and (12.7) while keeping the results close to our meta-trained
 23 parameters using penalties like

$$\frac{1}{2\lambda} \|w - \hat{w}\|_2^2 \quad \text{and} \quad \frac{1}{2\lambda} \|\varphi - \hat{\varphi}\|_2^2.$$

24 We don't update the context parameters α_i s during such adaptation.

❓ Does adaptation always improve performance on the new task?

25 12.2.2 Discussion

26 This brings an end to the chapter on meta-learning and Module 4. We
 27 focused on adapting learning-based models for robotics to new tasks. This
 28 adaptation can take the form of learning a reward (inverse RL), learning
 29 the dynamics (model-based RL) or learning to adapt (meta-learning).
 30 Adaptation to new data/tasks with few samples is a very pertinent problem
 31 because we want learning-based methods to generalize a variety of different

1 tasks than the ones they have been trained for. Such adaptation also comes
2 with certain caveats, adaptation may not always improve the performance
3 on new tasks; understanding when one can/cannot adapt forms the bulk of
4 the research on meta-learning.

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 Fakoor, R., Mueller, J. W., Asadi, K., Chaudhari, P., and Smola, A. J. (2021). Continuous doubly constrained
5 batch reinforcement learning. *Advances in Neural Information Processing Systems*, 34:11260–11273.
- 6 Levine, S., Kumar, A., Tucker, G., and Fu, J. (2020). Offline reinforcement learning: Tutorial, review, and
7 perspectives on open problems. *arXiv preprint arXiv:2005.01643*.
- 8 Turing, A. M. (2009). Computing machinery and intelligence. In *Parsing the Turing Test*, pages 23–65.
9 Springer.