

Personal Coverage Notes on Machine Learning

Authored by Dun-Ming Huang; Department of Computer Science, UC Berkeley

Preface(s)

0.1 Preface

In this personal coverage note, I try to record some thought processes and learning results I obtained from watching lectures and videos about numerous machine learning literature or applications. For now, the personal coverage note covers the following sections:

0.1.1 COMPSCI 189: Introduction to Machine Learning

This part covers the basic concepts of machine learning, including supervised learning, unsupervised learning, and focuses specifically on statistical learning techniques.

0.1.2 COMPSCI 285: Deep Reinforcement Learning

This part covers fundamental knowledge regarding deep reinforcement learning, sampled from Sergey Levine's Deep Reinforcement Learning lecture.

0.1.3 COMPSCI 294-158: Deep Unsupervised Learning

This part covers fundamental knowledge regarding deep unsupervised learning, sampled from Pieter Abbeel's Deep Unsupervised Learning lecture.

0.1.4 Mu Li's Videos on YouTube

This part covers the notes I took from watching Mu Li's videos on YouTube, which covers a wide range of topics in machine learning.

Contents

0.1	Preface	3
I	COMPSCI 189: Introduction to Machine Learning	9
1	Fundamentals of Machine Learning	10
1.1	The Framework of Machine Learning	10
1.2	Classification as Example Machine Learning Task	10
1.3	The Train-Validate-Test Framework	11
2	Linear Classifiers	13
2.1	section title	13
3	Gradient Descent	14
3.1	section title	14
4	Support Vector Machine	15
4.1	section title	15
5	It Is All About The Layers of Abstraction	16
5.1	section title	16
6	Decision Theory, Bayesian Decision Rule	17
6.1	section title	17
7	Gaussian Discriminant Analysis and Maximum Likelihood Estimation	18
7.1	section title	18
8	Eigendecomposition of Symmetric Matrices	19
8.1	section title	19
9	Abstractions of a Regression Problem	20
9.1	section title	20
10	Newton's Method and Logistic Regression	21

<i>CONTENTS</i>	5
10.1 section title	21
11 Statistical Justifications for Regressions	22
11.1 section title	22
12 Ridge Regression and Regularization	23
12.1 section title	23
13 Decision Trees	24
13.1 section title	24
14 The Kernel Trick	25
14.1 section title	25
15 Introduction to Neural Networks	26
15.1 section title	26
16 Tricks and Heuristics for Neural Networks	27
16.1 section title	27
17 Convolutional Neural Network	28
17.1 section title	28
18 Unsupervised Learning: PCA	29
18.1 section title	29
19 Unsupervised Learning: Clustering Algorithms	30
19.1 section title	30
20 Unsupervised Learning: Clustering Algorithms	31
20.1 section title	31
21 Geometry of High-Dimensional Space	32
21.1 section title	32
22 Learning Theory	33
22.1 section title	33
23 AdaBoost	34
23.1 section title	34
24 k-Nearest Neighbor Approaches	35
24.1 section title	35

II COMPSCI 285: Deep Reinforcement Learning	36
25 Introduction to Deep Reinforcement Learning	37
25.1 Motivation to Reinforcement Learning	37
25.2 Introduction to Reinforcement Learning	37
25.3 Motivation Towards Deep Reinforcement Learning	38
26 Supervised Learning of Behaviors	39
26.1 Terminology and Notation in DRL	39
26.2 Imitation Learning	40
26.3 Theoretical Analysis of Failure in Behavioral Cloning	40
26.4 Addressing Problem of Imitation Learning	41
26.5 A Hint at the Need of Reward and Cost Signals	42
27 Introduction to Reinforcement Learning	43
27.1 Foundations, in A Comprehensive Manner	43
27.2 Value Functions	44
27.3 Algorithms in Reinforcement Learning	45
28 Policy Gradients	47
28.1 Foundations of Policy Gradient	47
28.2 Reducing Variance in Policy Gradient	48
28.3 Off-Policy Policy Gradient	49
28.4 Covariant Policy Gradient	50
29 Actor-Critic Algorithms	51
29.1 section title	51
30 Value Function Methods	52
30.1 section title	52
31 Deep RL with Q-Functions	53
31.1 section title	53
32 Advanced Policy Gradients	54
32.1 section title	54
33 Optimal Control and Planning	55
33.1 section title	55
34 Model-Based Reinforcement Learning	56
34.1 section title	56

<i>CONTENTS</i>	7
35 Model-Based Policy Learning	57
35.1 section title	57
36 Exploration (Part 1)	58
36.1 section title	58
37 Exploration (Part 2)	59
37.1 section title	59
38 Offline Reinforcement Learning (Part 1)	60
38.1 section title	60
39 Offline Reinforcement Learning (Part 2)	61
39.1 section title	61
40 Reinforcement Learning Theory	62
40.1 section title	62
41 Variational Inference and Generative Models	63
41.1 section title	63
42 Connection between Inference and Control	64
42.1 section title	64
43 Inverse Reinforcement Learning	65
43.1 section title	65
44 Reinforcement Learning with Sequence Models	66
44.1 section title	66
45 Meta-Learning and Transfer Learning	67
45.1 section title	67
46 Challenges and Open Problems	68
46.1 section title	68
 III COMPSCI 294-158: Deep Unsupervised Learning	 69
47 Introduction to Deep Unsupervised Learning	70
47.1 Introduction to the Subject	70
48 Autoregressive Models	71
48.1 section title	71

49 Flow Models	72
49.1 section title	72
50 Latent Variable Models	73
50.1 section title	73
51 GAN and Implicit Models	74
51.1 section title	74
52 Diffusion Models	75
52.1 section title	75
53 Self-Supervised Learning, Non-Generative Representation Learning	76
53.1 section title	76
54 Leage Language Models	77
54.1 section title	77
55 Video Generation	78
55.1 section title	78
56 Semi-Supervised Learning and Unsupervised Distribution Alignment	79
56.1 section title	79
57 Compression	80
57.1 section title	80
58 Multimodal Models	81
58.1 section title	81
59 Parallelization	82
59.1 section title	82
60 AI for Science (Gues Instructor)	83
60.1 section title	83
61 Neural Radiance Fields (Guest Instructor)	84
61.1 section title	84
IV Mu Li's Videos on YouTube	85
62 Note Name	86
62.1 section title	86

Part I

COMPSCI 189: Introduction to Machine Learning

Chapter 1

Fundamentals of Machine Learning

In this chapter, we cover the introductory lecture of COMPSCI 189.

Learning Goals:

- Understand the fundamental workings of machine learning.
- Learn about classification as an example machine learning task.

Machine learning is a popular topic over the recent centuries. It is a subset of artificial intelligence that focuses on the development of algorithms that allow computers to learn from and make predictions based on data. The study of machine learning per se is a long journey, even disregarding the participation of research activities. In this section of the note, we concentrate on statistical learning, which involve fundamental techniques of machine learning that are popularized before the current pulvinar of deep learning. Deep learning techniques will be addressed in later sections of the entire note.

1.1 The Framework of Machine Learning

Machine learning is the use and development of computer systems that can learn without explicit instructions; that is, they learn a specific pattern of the provided data via statistical measures, in an autonomous and algorithmic manner. The learning process is largely valuable on the ability of machine learning algorithms to draw insights, or **inferences**, upon the provided training data. Fundamentally, statistical learning is all about finding patterns in data, and using them to make predictions. An abstraction of this will be issued in Lecture 5 of the section.

Machine learning is a data-driven approach. As mentioned before, all that a machine learning can learn from is what the distribution of a provided training data provides. This is an important insight in the future. Just as how humans cannot learn what a cat is if they have never learned anything about a cat, a machine cannot learn about cat if the data we provide to its algorithm never describes what a cat is. In summary, what an algorithm learns is largely dependent on the data we provide to it.

1.2 Classification as Example Machine Learning Task

Classification, as you may have learned in highschool biology, is the process of categorizing things based on their properties. In machine learning, classification is a task that involves predicting the category of a given data point. For example, provided a picture that may entail a cat or a dog, a machine learning algorithm would be asked to classify it as either a cat or a dog. This is convenient in that humans do not have to process this judgment manually, and can instead automate this task with a fairly accurate algorithm.

How do we really decide if a given data point is a cat or a dog? For example, suppose the datapoint I am provided is an image, how do I transform this image into a decision's label (a cat, versus a dog)? In classification, we usually use

numbers to denote the label of a category (hereby we call it a “class”). A classifier h , therefore, is a function that is provided a datapoint \vec{x} and outputs a numeric label for the representing class:

$$h(\vec{x}) = \begin{cases} 1 & \text{if the algorithm considers } \vec{x} \text{ is a cat} \\ 0 & \text{if the algorithm considers } \vec{x} \text{ is a dog} \end{cases}$$

The question now comes down to:

1. \vec{x} : How is the colorful image we see in human eyes represented as a vector?
2. “if the algorithm considers \vec{x} as a something”: how is this rule implemented programmatically?

For question (1), the image’s pixels can be converted into color-representing numeric values, then flattened into a vector based on the spatial ordering of pixels. For question (2), the algorithm learns a function h that outputs the above mapping. The takeaways of above questions are as follows:

1. The machine learning algorithm receives data in numeric form, such as a list of numbers (vectors), but not qualitatively.
2. The machine learning algorithm learns a function that is tailored to our need.

1.3 The Train-Validate-Test Framework

In machine learning, an algorithm usually follows the framework of train, validate, test. These aspects of the paradigm are summarized as follows.

1.3.1 Aspects of the T V T Framework

Training a Model. Recall that any machine learning algorithm produces a function h , which we also call a model, by having the algorithm detect patterns in our datapoints \vec{x} ’s. The act of learning an appropriate function h that behaves well on our given datapoints is called **training** a model. That is, we are training a machine learning model on a provided dataset, and the resulting model should learn a function h that accurately predicts the class label (dog vs. cat) of a given datapoint \vec{x} (an image). In this phase, models are provided a labeled dataset; that is, a set of images that are labeled either as a cat or a dog. Such dataset we use to train the model is otherwise known as a **training set**. Usually, we continue with the training phase until the algorithm’s model has reached a satisfying accuracy for the training set.

Validating a Model. We have trained a model with images of cats and dogs, and now it’s time to evaluate the model. More precisely, it’s time to evaluate the model on datapoints it has not seen yet. After all, when a model is deployed into the real world, it is expected for the model to be able to classify cats and dogs from pictures immediately before us, mostly unseen to anyone, rather than just known images that are already labeled in a dataset. The dataset that we use to validate the model, entirely unseen during the training phase, is known as the **validation set**. We will stay in this phase until our model has reached a satisfying accuracy for the validation set.

Testing a Model. At last, we evaluate our model again using another unseen dataset, called the **testing set**. The testing set is used to evaluate the model’s performance on a dataset that is entirely unseen during the training and validation phases. This is the final phase of the train-validate-test framework, and the model’s performance on the testing set is the final metric of the model’s performance.

1.3.2 Justification: Overfitting and Underfitting

Hi

1.3.3 A Summary of Questions up to This Point

Hi

Chapter 2

Linear Classifiers

Chapter Description.

2.1 section title

Section.

Theorem 2.1.1. Tested Theorem

I am the bone of my sword.

Definition 2.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 2.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 3

Gradient Descent

Chapter Description.

3.1 section title

Section.

Theorem 3.1.1. Tested Theorem

I am the bone of my sword.

Definition 3.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 3.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 4

Support Vector Machine

Chapter Description.

4.1 section title

Section.

Theorem 4.1.1. Tested Theorem

I am the bone of my sword.

Definition 4.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 4.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 5

It Is All About The Layers of Abstraction

Chapter Description.

5.1 section title

Section.

Theorem 5.1.1. Tested Theorem

I am the bone of my sword.

Definition 5.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 5.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 6

Decision Theory, Bayesian Decision Rule

Chapter Description.

6.1 section title

Section.

Theorem 6.1.1. Tested Theorem

I am the bone of my sword.

Definition 6.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 6.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 7

Gaussian Discriminant Analysis and Maximum Likelihood Estimation

Chapter Description.

7.1 section title

Section.

Theorem 7.1.1. Tested Theorem

I am the bone of my sword.

Definition 7.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 7.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 8

Eigendecomposition of Symmetric Matrices

Chapter Description.

8.1 section title

Section.

Theorem 8.1.1. Tested Theorem

I am the bone of my sword.

Definition 8.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 8.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 9

Abstractions of a Regression Problem

Chapter Description.

9.1 section title

Section.

Theorem 9.1.1. Tested Theorem

I am the bone of my sword.

Definition 9.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 9.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 10

Newton's Method and Logistic Regression

Chapter Description.

10.1 section title

Section.

Theorem 10.1.1. Tested Theorem

I am the bone of my sword.

Definition 10.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 10.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 11

Statistical Justifications for Regressions

Chapter Description.

11.1 section title

Section.

Theorem 11.1.1. Tested Theorem

I am the bone of my sword.

Definition 11.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 11.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 12

Ridge Regression and Regularization

Chapter Description.

12.1 section title

Section.

Theorem 12.1.1. Tested Theorem

I am the bone of my sword.

Definition 12.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 12.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 13

Decision Trees

Chapter Description.

13.1 section title

Section.

Theorem 13.1.1. Tested Theorem

I am the bone of my sword.

Definition 13.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 13.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 14

The Kernel Trick

Chapter Description.

14.1 section title

Section.

Theorem 14.1.1. Tested Theorem

I am the bone of my sword.

Definition 14.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 14.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 15

Introduction to Neural Networks

Chapter Description.

15.1 section title

Section.

Theorem 15.1.1. Tested Theorem

I am the bone of my sword.

Definition 15.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 15.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 16

Tricks and Heuristics for Neural Networks

Chapter Description.

16.1 section title

Section.

Theorem 16.1.1. Tested Theorem

I am the bone of my sword.

Definition 16.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 16.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 17

Convolutional Neural Network

Chapter Description.

17.1 section title

Section.

Theorem 17.1.1. Tested Theorem

I am the bone of my sword.

Definition 17.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 17.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 18

Unsupervised Learning: PCA

Chapter Description.

18.1 section title

Section.

Theorem 18.1.1. Tested Theorem

I am the bone of my sword.

Definition 18.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 18.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 19

Unsupervised Learning: Clustering Algorithms

Chapter Description.

19.1 section title

Section.

Theorem 19.1.1. Tested Theorem

I am the bone of my sword.

Definition 19.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 19.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 20

Unsupervised Learning: Clustering Algorithms

Chapter Description.

20.1 section title

Section.

Theorem 20.1.1. Tested Theorem

I am the bone of my sword.

Definition 20.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 20.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 21

Geometry of High-Dimensional Space

Chapter Description.

21.1 section title

Section.

Theorem 21.1.1. Tested Theorem

I am the bone of my sword.

Definition 21.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 21.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 22

Learning Theory

Chapter Description.

22.1 section title

Section.

Theorem 22.1.1. Tested Theorem

I am the bone of my sword.

Definition 22.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 22.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 23

AdaBoost

Chapter Description.

23.1 section title

Section.

Theorem 23.1.1. Tested Theorem

I am the bone of my sword.

Definition 23.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 23.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 24

k-Nearest Neighbor Approaches

Chapter Description.

24.1 section title

Section.

Theorem 24.1.1. Tested Theorem

I am the bone of my sword.

Definition 24.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 24.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Part II

COMPSCI 285: Deep Reinforcement Learning

Chapter 25

Introduction to Deep Reinforcement Learning

Chapter Description.

25.1 Motivation to Reinforcement Learning

Let us begin a motivating problem: how can we let a robot hand pick up something?

In classical robotics, the problemsolving process is to: (1) Define the problem in modeling perspectives, (2) Model the problem using mathematical equations, and (3) Solve the problem via a designed algorithm. However, as we accumulate technological knowledge, now we have a second option, which is to set it up as a machine learning problem. With the knowledge we have currently learned in statistical learning, we are inclined to use supervised learning; that is, provided some data of the robot and the environment, we train some model that can provide the robot an action to comply with. However, this approach is not well-informed from human experience, and the crafting of such data is still difficult. So, instead, we follow the line of thought of letting robots earn their own experiences via trial and error. This approach develops into a **reinforcement learning** setting.

In a reinforcement learning setting, robots collect examples of their own behavior, and label its success (significance as well) based on a state-derived reward signal (function). The robot then learns to maximize the reward signal by adjusting its behavior. Eventually, we obtain a **policy** (a function $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that provides an action in response to a seen state) that the robot can follow to pick an object up.

So, reinforcement learning is really an experience-collecting framework: it allows for a freedom of trial, a disregard for a pre-existing dataset (although in most situations we will still have one), and inherits from machine learning approaches the waive of need to manually design solutions for each specific problem. Like statistical learning, reinforcement learning is also a massively scaled process of density estimations for underlying distributions (say, $p_\theta(x)$, or $p_\theta(y|x)$) of the training data. However, reinforcement learning is different in that it enables learning via agent-environment interactions, which provides a new source of information; and, it allows for new applications like evolutionary algorithms, controls, and optimizations. Reinforcement learning is mainly an approach for a design of behavior that does not require human intervention. These behavior are impressive because it is unthought of, as well as because of its delicate mimicry for human results (say, artistic outcomes).

25.2 Introduction to Reinforcement Learning

Reinforcement learning is both a mathematical formalism for learning-based decision making, and an approach for learning decision making and control from experience.

In supervised learning, the framework follows as a provided dataset $\mathcal{D} = \{(x_i, y_i)\}$ that contributes to the learning of a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, which itself is the ultimate fruit of a supervised learning paradigm. Recall that supervised learning makes several assumptions regarding its paradigm. One, the data provided to us are i.i.d. samples from an underlying distribution. Two, we have known ground truth outputs in training.

In reinforcement learning, however, we ignore both of these assumptions. One, data is not i.i.d., because previous outputs influence future inputs (in a Markovian fashion). Two, ground truth answer is not known, and we are only provided a reward signal that notifies us whether a demonstration from the agent is successful or not. Reinforcement learning is therefore ran on reward-labeled data, rather than ground-truth-labeled data.

To summarize the paradigm of reinforcement learning, it is comprised of the following aspects:

1. An **agent** that interacts with the world to achieve a specific task.
2. An **environment** that the agent interacts with. This can be a Minecraft flat world for an agent that tries to learn walking on. This environment can be both in-real-life and simulated. In most occasions, it is simulated.
3. The input of learning is a **state** s_t that represents the current situation of the agent.
4. The output of learning, generally, is a **action** a_t that the agent takes in response to the state.
5. A **reward** r_t that the agent receives after taking an action.
6. A **policy** π that the agent follows to take an action.

The data we receive for reinforcement learning is therefore a sequence of (s_t, a_t, r_t) tuples that the agent collects from the environment as it interacts with it.

25.3 Motivation Towards Deep Reinforcement Learning

The fusion of data-driven AI and reinforcement learning provides us a complementary approach. In data-driven AI (deep learning), while we extract valuable inferences about the real world from data, we don't actively attempt to perform better than the data. Meanwhile, in reinforcement learning, while we extract emergent behavior to do better than existing data, we are not prepared with a way to extract inferences regarding the environment, and are not provided a means of using data at scale. That is, Data-Driven AI is about using data, while reinforcement learning is about using optimizations. Therefore, deep reinforcement learning is expected to excel at both learning and searching: learning from data and searching for (discovering) better ways to interact with the environment provided the data.

Noted, we have deep neural network architectures that extracts inference well, and RL algorithms that are compatible with these approaches. However, at the current stage, learning-based control in truly real-world settings remains a major open problem. We will discuss these topics at lengths with later sections of the note.

In the current state, we face the following open challenges:

1. We don't yet have amazing methods that both use data and Reinforcement Learning
2. Humans can learn incredibly quickly, but deep RL methods are usually slow, even in simulators.
3. Humans can reuse past knowledge, but domain transfer is a problem to deep reinforcement learning.
4. The role of prediction and design of reward functions are still not very clear in reinforcement learning.

Chapter 26

Supervised Learning of Behaviors

Chapter Description.

26.1 Terminology and Notation in DRL

In this section, we will cover several terminologies and notations that the community uses regarding learning situations.

In our paradigm, we concern an input, a system that processes the input, and the output. Decision-making problems consider these as respectively observations o_t , policy $\pi_\theta(a_t|o_t)$, and actions a_t . These symbols are subscripted by time because the context of a decision-making problem is usually a chronological sequence of events. Note that, the production of next observation, o_{t+1} , is based on the impacted state o_t upon the transpiration of a_t . Note that the format of a_t , for example, is not limited to a vector; it can also be a continuous distribution, as we would sometimes like to sample actions to take rather than using an almost deterministic policy. The policy π_θ is parameterized by θ , which is a vector of parameters that the policy uses to make decisions (that is, to provide action provided observation). They, therefore, assign the probability to all possible actions provided a specific state.

We also introduce the notion of a state, s_t , which is a partial observation. This notion is introduced by the fact that the entirety of an environment is not always observable (and in most situations, we do not observe the entirety of an environment). Therefore, realistically, the policy we learn is $\pi_\theta(a_t|s_t)$, which we call a partially observed policy. A very appropriate analogy is perhaps our visual-neural system: our eyes guide our decisions based on what objects are posed in our environment, particularly what is in front of our eyes. But, we do not gain information from our eyes regarding what is behind us, simply because the environment we are situated in is partially observable: the sensors we have (eyes) simply do not detect what is behind. Therefore, in reinforcement learning paradigms, we work with state-action pairs rather than observation-action pairs.

We develop this into what we call a Markov Decision Paradigm. In this paradigm, we assume that the state s_t is sufficient to make decisions, and that the future state s_{t+1} is independent of the past states s_{t-1}, s_{t-2}, \dots given the current state s_t and the action a_t . The paradigm per se contains Markov-ness; that is, s_t only depends on s_{t-1} , and the connection of states is only nonzero for $p(s_{t+1}|s_t, a_t)$. The involvement of only states and actions is a reflection of the partial observability we suffer in environments, which produces a Partially Observable Markov Decision Process (POMDP).

Note: in some literature, states will be issued as x_t , and action u_t instead, due to the involvement of background in control theory from some influential figures.

26.2 Imitation Learning

Imitation learning concerns the learning of a policy π_θ from a dataset of expert demonstrations. **Behavioral cloning**, an approach of imitation learning, is the idea that, via supervised learning, we learn a policy that regresses an underlying function $f : \mathcal{S} \rightarrow \mathcal{A}$ that maps the expert-induced states to expert-performed actions. These expert demonstrations are, as assumed, provided to the algorithm. However, behavioral cloning is usually not a good solution to general problems. Therefore, we expect the agent to practically clone the behavior (state-action reaction) of the expert.

The reason why is because, although expert demonstrations provide us many trajectories, once we receive a state outside of the provided trajectories, the actions our agent provides are not well-defined in nature. It is moreover that, once the agent is exposed to an unseen state, the action it provides via its policy does not guarantee a continued cloning of the expert-demonstrated trajectory. This is because the nature of a trajectory destructs the i.i.d. assumption of supervised learning; that is, because of temporal dependencies, we encounter problems in behavioral cloning that does not appear in conventional supervised learning problems. (A spam classifier generalizes well to unseen emails, but a behavioral cloning agent does not generalize well to unseen states). We can propose ad-hoc solutions, such as data augmentation that increases familiarity of the agent to unseen states, but these solutions are not always guaranteed to work.

To make behavioral cloning work, we must make modifications to the existing paradigm. Let us begin with lessons of the story:

1. Different from conventional supervised studying problems, imitation learning via behavioral cloning is not guaranteed to work.
2. To work against it, we can use data augmentation as well as modify data collection methods.
3. An exotic solution, like a multi-task learning formulation, may help to generalize to unseen states and perform good imitation learning.
4. The most intuitive approach is perhaps modifying the algorithm along which we do imitation learning.

26.3 Theoretical Analysis of Failure in Behavioral Cloning

The main culprit of behavioral cloning's failure is the distributional shift problem. Suppose that we have some policy $\pi_\theta(a_t|o_t)$ trained on a dataset of expert demonstrations. Here, let us say that the distribution provided by expert dataset is $p_{data}(o(t))$, but the distribution of observation the policy faces is $p_{\pi_\theta}(o_t)$. Note that, since our policy is trained under $p_{data}o(t)$, the objective of that training is $\max_\theta \mathbb{E}_{o_t \sim p_{data}(o(t))} \log \pi_\theta(a_t|o_t)$. However, the policy is evaluated under $p_{\pi_\theta}(o_t)$, which is not the same as the training distribution. This problem is otherwise known as **distributional shift**, and this occurs due to the policy's own deviation from the training distribution.

The lesson of such problem is that we should perhaps define more precisely what we want to define as “well-learned”. A policy that is good would perhaps not be a point-estimate for actions, since it can easily lead to deviations in policy behavior. Perhaps we may define a cost otherwise. Suppose that π^* is the optimal policy for us to clone:

$$c(s_t, a_t) = \begin{cases} 0 & \text{if } a_t = \pi^*(s_t) \\ 1 & \text{otherwise} \end{cases}$$

Now, then, our training objective becomes:

$$\min \mathbb{E}_{s_t \sim p_{\pi_\theta}(s_t)} [c(s_t, \pi_\theta(s_t))]$$

Assume that we have an upper bound of the policy mistake, $\pi_\theta(a \neq \pi^*(s)|s) \leq \epsilon$. Then, we observe an incurred cost of:

$$\begin{aligned} \mathbb{E}[\sum_t c(s_t, \pi_\theta(s_t))] \\ \leq \epsilon T + (1 - \epsilon)(\epsilon(T - 1) + (1 - \epsilon)(\dots)) \in O(\epsilon T^2) \end{aligned}$$

Therefore, the error of behavioral cloning is quadratic to the length of its trajectory (T).

It turns out that an analysis composed by Ross et al. (2011) shows that the error of behavioral cloning is quadratic to the length of the trajectory. That is, imitation learning is prone to failure at any timestep, and lacks a means of recovery from a policy's small failures.

26.4 Addressing Problem of Imitation Learning

To address the problem of imitation learning, we can consider the following approaches:

1. Data Augmentation and collection
2. Powerful models that make very few mistakes
3. Multi-task learning formulation of imitation learning
4. Changing the algorithm of use (where we discuss DAgger)

Data Augmentation and Collection. Behavioral cloning is difficult because the model doesn't generalize well to mistakes. What if we involve data regarding mistakes instead? If the dataset involves mistakes, due to additional steps in data collection and augmentation process, although the training set will be diluted, we can now access corrections in our BC paradigm.

Powerful Models. Failures from fitting the expert can stem from several reasons. First, the expert's behavior is non-Markovian. While the policy we train is Markovian-ly conditioned, the expert's behavior are not formulated on a Markovian approach. That is, provided exposure to states $s_t = s_{t'}$, it is likely that the action elicited from these states are different. Therefore, human demonstrators post a very unnatural circumstance for a Markovian policy. Perhaps one remedy is to use the entire history of a trajectory, and a sequence model is capable of processing it as a temporal sequence of frames. However, the exploitation of entire history may still work poorly, simply because including the history may still be harmed by incomplete information, which can lead to causal confusion. Second, the expert's behavior may be multimodal. That is, the expert may have multiple ways of solving a problem, and the policy we train may not be able to capture all of these modes. This is specifically unhelpful for a continuous distribution of actions, which rely on the use of mean and variance that is difficultly characterizing for bimodal distributions. Although we can instead choose expressive continuous distributions for actions (namely, use other classes of distributions), or use an autoregressive discretization with high-dimensional action spaces, they do pose higher computational costs to this procedure.

Multi-task Learning Formulation. Training a policy that has one sole destination of trajectory can be difficult, but perhaps a multi-task learning formulation that attempts to let policies reach multiple different destinations. This approach can be summarized as "goal-conditioned behavioral cloning", which can provide more opportunities to learn corrections despite distributional shift; that is, to maximize $\log \pi_\theta(a_t^i | s_t^i, g = s_T^i)$. However, this approach actually introduces a secondary source for distributional shift, making the approach theoretically worse.

Changing the Algorithm of Use: DAgger. The idea of DAgger, Dataset Aggregation, is to make $p_{data}(o_t) = p_{\pi_\theta}(o_t)$. That is, we collect training data from $p_{\pi_\theta}(o_t)$ by running the policy $\pi_\theta(a_t | o_t)$ in the environment. Then, we aggregate the dataset of expert demonstrations and the dataset of the policy's own data, and train the policy on the aggregated dataset. The algorithm is as follows:

1. Collect a dataset of expert demonstrations $\mathcal{D} = \{(o_t^i, a_t^i)\}_{i=1}^N$.
2. For $k = 1, 2, \dots, K$:
 - (a) Train the policy π_θ on \mathcal{D} .
 - (b) Collect a dataset of the policy's own data $\mathcal{D}_{\pi_\theta} = \{(o_t^i, \pi_\theta(o_t^i))\}_{i=1}^N$.
 - (c) Aggregate the datasets $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_{\pi_\theta}$.

26.5 A Hint at the Need of Reward and Cost Signals

Deep learning works best when data is plentiful, but humans are finite sources of data (that is, humans cannot provide all data and label all data). Therefore, to resolve the disruptive demand of large ground truths, we expect our reinforcement learning algorithm to learn autonomously. To enable algorithms to evaluate their experiences, and exceeding their own performances, we offer a reward signal called **reward function**, $r(s, a)$, which provides a numeric evaluation for an observed pair of state and action. For example, for imitation learning, we can propose the reward function $r(s, a) = \log p(a = \pi^*(s)|s)$.

Chapter 27

Introduction to Reinforcement Learning

Chapter Description.

27.1 Foundations, in A Comprehensive Manner

In prior, we have learned that reinforcement learning focuses on learning a policy $\pi_\theta(a_t|o_t)$ that provides an action provided an observation. Here, a_t is an action at timestep t , while o_t is an observation at timestep t . Note that, in a not fully observed context, we instead discuss s_t as a state at timestep t in place of the observation. And, that, the decision paradigm we adopt here is Markovian: the only determining factor of the current timestep is the previous timestep.

Different from imitation learning, a supervised learning setting, we use the reward function to notate the success of a performance. A reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a function that provides a scalar reward for a state-action pair; as the name suggests, a higher reward resembles a better performance. However, reinforcement learning is not about short-term maximizations of reward functions, but long-term maintenance of high reward values. A temporal abstraction is therefore also imposed upon reinforcement learning formulations.

27.1.1 Mathematical Objects

A Markov Chain is a stochastic object defined as $\mathcal{M} = \{\mathcal{S}, \mathcal{P}\}$ where \mathcal{S} is the set of all states (state-space), and \mathcal{T} is a transition operator that encodes $p(s_{t+1}|s_t)$ for all possible pairs (s_t, s_{t+1}) . The transition operator is expressible as a matrix, where each row sums to 1, and each element is a probability of transitioning from one state to another. That is, $\mathcal{T}_{i,j} = p(s_{t+1} = i | s_t = j)$.

A reinforcement learning paradigm, then, can be framed as a variation of Markov Chain, called a Markov Decision Process. This process is formulated as $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{T}, r\}$, where:

- \mathcal{S} is the state-space, as before.
- \mathcal{A} is the action-space, a set of all possible actions.
- \mathcal{T} is the transition operator, as before.
- r is the reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, as before.

A partially observed Markov decision process is just an augmented MDP: $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{E}, r\}$, where \mathcal{E} is the emission probability $p(o_t|s_t)$.

27.1.2 Objective of Reinforcement Learning

In reinforcement learning, we learn a policy-representing object $\pi_\theta(a|s)$. In the chain rule of states within some trajectory, we may discover that:

$$p_\theta(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_\theta(a_t|s_t) p(s_{t+1}|s_t, a_t)$$

We denote the trajectory $s_1, a_1, \dots, s_T, a_T$ as τ . The objective of reinforcement learning is to maximize the expected return, or the expected sum of rewards:

$$J(\theta) = \mathbb{E}_{\tau \sim p_\theta} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

and we would like a policy equipped with the parameter θ that maximizes this return:

$$\theta^* = \arg \max_{\theta} J(\theta)$$

Let us attempt to define the finite horizon case of the expected return using an augmented state-space for markov chain on state-action pairs. Therefore, by linearity of expectation,

$$\begin{aligned} \theta^* &= \arg \max_{\theta} J(\theta) \\ &= \arg \max_{\theta} \mathbb{E}_{\tau \sim p_\theta} \left[\sum_{t=1}^T r(s_t, a_t) \right] \\ &= \arg \max_{\theta} \sum_{t=1}^T \mathbb{E}_{\tau \sim p_\theta} [r(s_t, a_t)] \end{aligned}$$

That also means, in an infinite-horizon case, where $T = \infty$, our objective may become ill-defined unless we obtain a finite mean of reward.

27.2 Value Functions

We can also express our expected reward objective in a recursive manner:

$$\mathbb{E}_{s_1 \sim p(s_1)} \left[\mathbb{E}_{a_1 \sim p(a_1)} \left[r(s_1, a_1) + \mathbb{E}_{s_2 \sim p(s_2|s_1, a_1)} \left[\mathbb{E}_{a_2 \sim p(a_2|s_2)} [r(s_2, a_2) + \dots] | s_1, a_1 \right] | s_1 \right] \right]$$

To simplify this expression, we express:

$$Q(s_1) = r(s_1, a_1) + \mathbb{E}_{s_2 \sim p(s_2|s_1, a_1)} \left[\mathbb{E}_{a_2 \sim p(a_2|s_2)} [r(s_2, a_2) + \dots] | s_1, a_1 \right]$$

such that,

$$\mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right] = \mathbb{E}_{s_1 \sim p(s_1)} \left[\mathbb{E}_{a_1 \sim p(a_1)} [Q(s_1) | s_1] \right]$$

Now, with this concise notation, we may easily modify π_θ based on the function Q . The question, then, is how do we know the function Q , otherwise called the Q-function?

So, as we have defined in prior, the Q-function is expressed as:

$$Q^\pi(s_t, a_t) = \sum_{t=t'}^T \mathbb{E}_{\pi_\theta} [r(s_t, a_t) | s_t, a_t]$$

The value function, then, is defined as:

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi_\theta(a_t|s_t)} [Q^\pi(s_t, a_t)]$$

Now, we can also express the reinforcement learning objective as $\mathbb{E}_{s_1 \sim p(s_1)} [V^\pi(s_1)]$.

How can we use Q-functions and value functions? One immediate application of them is to formulate a policy that only outputs the maximum-Q-inducing action. Another idea would be to compute the gradient of Q^π to increase the probability of a good action a , building on the precondition that $Q^\pi(s, a) > V^\pi(s)$ implies a as better than an average action, so that we modify $\pi(a|s)$ to increase the probability of the aforementioned condition.

27.3 Algorithms in Reinforcement Learning

Reinforcement learning algorithms share the following high-level anatomy:

- Collect data from the environment (via running the policy).
- Update the policy using the collected data (to fit the policy's model).
- Improving the policy (this is some form of optimization).

Here, generating samples has an expense that depends on the task of reinforcement learning; for example, under a robotics setting, the sample generation process is very costly without a simulator. Fitting a model and improving policies also has an expense dependent on the approach.

27.3.1 Diversity in RL Algorithms

Here is a general categorization of reinforcement learning algorithms.

Policy Gradient algorithms directly differentiate the expected-reward objective. While the gradient can be estimated, the improvement of policy relies on a gradient descent step that will be introduced in the next lecture.

Value-based algorithms estimate a value function or Q function of the optimal policy, rather than having an explicit policy. In this family of algorithms, we use a model to fit $V(s)$ or $Q(s, a)$, and use this model to plan actions. The policy is simply $\pi(s) = \arg \max_a Q(s, a)$.

Actor-Critic algorithms combine the two, where the critic is an estimator of the value function, and the improvement of policy relies on a gradient descent step.

Model-based algorithms learn a model of the environment, and use this model to plan actions. For model-based algorithm, fitting a model also involves learning a transition operator \mathcal{T} , while the policy improvement aspect has many options: (1) using the learned transition to plan eventual actions; (2) backpropagating gradients into the policy network; (3) using the transition model to learn a value function, and use dynamic programming to accelerate the improvement progress.

Different families of algorithms have tradeoffs of performance on several fields, ranging across the following examples:

Sample Efficiency. Sample efficiency is the matter of “how many samples do we need to get a good policy”, and one important factor of this is whether the algorithm is on-policy (online, new samples are needed upon any update of policy) or off-policy (the algorithm improves the policy without generating new samples from that policy). On the spectrum of sample efficiency, on-policy algorithms generally host less sample efficiency. However, having a less sample efficiency is acceptable because of the tradeoff it can bring, such as the need of online-ness that is saved by on-policy algorithms’ online-ness in exchange of smaller sample efficiency.

Stability. Stability and ease of use discusses whether the algorithm converges, how often does it converge, and what does the algorithm converge to (if at all). This is a question because reinforcement learning might not be using gradient-based optimizations. For example, Q-learning is a fixed-point iteration scheme, while model-based RL concentrates on a transition model not optimized for expected reward.

Assumptions from Algorithms. Algorithms pose assumptions about our reinforcement learning paradigm. Some common assumptions include: full observability, episodic learning, and continuity or smoothness of specific objectives.

Chapter 28

Policy Gradients

Chapter Description.

28.1 Foundations of Policy Gradient

Remember that the objective of reinforcement learning is to maximize the expected return of our policy π_θ over time. In deep reinforcement learning, we would formulate that the policy is a neural network taking in some state s_t and outputting some action a_t . And, we may also define a trajectory distribution:

$$p_\theta(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

In this note, the shorthand of such distribution is $p_\theta(\tau)$.

Crucially, in development of model-free algorithms, we do not assume that we know the transition probabilities $p(s_{t+1} | s_t, a_t)$ or the initial state probability $p(s_1)$, but treat real-world interaction as an act of sampling per se. And, the objective of our algorithm is to find parameter θ where:

$$\theta^* = \arg \max_{\theta} J(\theta)$$

where

$$J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

for which the finite horizon case can be simplified into a sum of expectations via the linearity of expectation.

Let us attempt to evaluate the objective of our algorithm before discussing its optimization. We may find an unbiased estimator of $J(\theta)$ by sampling trajectories via policy-world interactions:

$$J(\theta) \sim \frac{1}{N} \sum_i \sum_t r(s_{i,t}, a_{i,t})$$

However, we do want to improve the objective beyond estimating it. The estimate of the derivative of objective also needs to be feasible without knowledge of transition probability and state prior. To satisfy the above demands, let us

propose as follows:

$$\begin{aligned}
J(\theta) &= \mathbb{E}_{\tau \sim p_\theta(\tau)}[r(\tau)] \\
&= \int p_\theta(\tau) r(\tau) d\tau \\
\nabla_\theta J(\theta) &= \int \nabla_\theta p_\theta(\tau) r(\tau) d\tau \\
&= \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) r(\tau) d\tau \\
&= \mathbb{E}_{\tau \sim p_\theta(\tau)}[\nabla_\theta \log p_\theta(\tau) r(\tau)] \\
&= \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]
\end{aligned}$$

To evaluate the policy gradient, we can run the policy to generate samples from $p_\theta(\tau)$, and then multiply them by the gradients. Then, we improve the policy by taking a step in the direction of the gradient:

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

The sketch of the algorithm, otherwise known as REINFORCE, is as follows:

- Sample $\{\tau^i\}$ from running the policy π_θ .
- Compute the policy gradient: $\nabla_\theta J(\theta) = \sum_i (\sum_t \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t})) (\sum_t r(s_{i,t}, a_{i,t}))$.
- Update the policy: $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$.
- Repeat.

Policy gradients are then reward-weighted versions of maximum likelihood objective. This may be observed from their expressions:

$$\begin{cases} \nabla_\theta J(\theta) & \sim \frac{1}{N} \sum_i (\sum_t \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t})) (\sum_t r(s_{i,t}, a_{i,t})) \\ \nabla_\theta J_{ML}(\theta) & \sim \frac{1}{N} \sum_i (\sum_t \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t})) \end{cases}$$

That is, in policy gradient, high reward trajectories have increased log probabilities, while low reward trajectories have decreased log probabilities.

Under the constraint of continuous actions, on the other hand, we select a representation of policy π that outputs a representation of continuous distribution, such as a Gaussian policy: $\pi_\theta(a_t | s_t) = \mathcal{N}(f(s_t), \Sigma)$. Then, the log probability of an action can be written in terms of an anisotropic Gaussian distribution's, or some other continuous distribution's PDF.

Note that, the implementation of policy gradient demands automatic differentiation, such that the policy gradient can be computed by a library within reasonable computation time. Typically, this is inefficiency, because neural networks have more parameters within itself than samples that it uses. To calculate the gradient in a neural network, we would often instead employ back-propagation, which the automatic differentiation will set up a computational graph for us as well when it comes to computation of policy gradient. We may, fortunately, construct a computational graph for which the gradient of it is the policy gradient. This is often by the involvement of Q values, and having an alternative objective $\tilde{J}(\theta)$ which grants us the policy gradient on its computational graph of gradients.

28.2 Reducing Variance in Policy Gradient

However, the property of policy gradients to skew away from bad trajectories and towards good ones, and this can introduce pathological changes led by high variance of its estimator. The policy gradient estimator we described before has very high variance depending on the acquired samples.

28.2.1 Causality Trick

The first trick we can use to reduce variance is causality, which is that a policy at time t' cannot affect the reward at time t when $t < t'$. This allows us to rewrite the policy gradient as:

$$\nabla_{\theta} J(\theta) \sim \frac{1}{N} \sum_i \sum_t \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \left(\sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) \right)$$

Which allows us to discard rewards from the past that our reward policy cannot impact. Then, we obtain a lower-variance estimate, because the total sum is a smaller number (which accompanies a smaller variance). This is otherwise known as a reward-to-go formulation.

28.2.2 Baseline Trick

The second trick we can use to reduce variance is the baseline trick. Because subtracting a baseline from the reward aspect in policy gradient does not change its expectation (a proof is provided in lecture), but changes its variance, for any baseline b , the following trick reduces the variance of the policy gradient while maintaining the estimator as unbiased. Let $b = \frac{1}{N} \sum_{i=1}^N r(\tau_i)$. Then, we use the baselined policy gradient instead:

$$\nabla_{\theta} J(\theta) \sim \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b]$$

Furthermore, via derivation, we find the optimal variance-reducing baseline to have value:

$$b^* = \frac{\mathbb{E}[g(\tau)^2 r(\tau)]}{\mathbb{E}[g(\tau)^2]}$$

28.3 Off-Policy Policy Gradient

Because the policy gradient is a sample-based estimator, the requirement of sampling from the policy upon every policy update becomes a large operational cost. This is a disadvantage of the REINFORCE algorithm's online-ness (it being on-policy). This immediately attracts us to convert the algorithm into an off-policy one, via an access to some $\bar{p}(\tau)$ that is not the policy we are optimizing for, so that we can update the policy while being offline (that is to be off-policy, then).

This may be operated along an importance sampling trick:

$$\begin{aligned} \mathbb{E}_{x \sim p(x)} [f(x)] &= \int p(x) f(x) dx \\ &= \int \frac{p(x)}{q(x)} q(x) f(x) dx \\ &= \mathbb{E}_{x \sim q(x)} \left[\frac{p(x)}{q(x)} f(x) \right] \end{aligned}$$

Then, we obtain our reinforcement learning objective as:

$$J(\theta) = \mathbb{E}_{\tau \sim \bar{p}(\tau)} \left[\frac{p_{\theta}(\tau)}{\bar{p}(\tau)} r(\tau) \right]$$

where

$$\begin{aligned} \frac{p_{\theta}(\tau)}{\bar{p}(\tau)} &= \frac{p(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)}{p(s_1) \prod_{t=1}^T \bar{\pi}(a_t | s_t) \bar{p}(s_{t+1} | s_t, a_t)} \\ &= \prod_{t=1}^T \frac{\pi_{\theta}(a_t | s_t)}{\bar{\pi}(a_t | s_t)} \end{aligned}$$

The policy gradient is then re-expressed to be:

$$\begin{aligned}\nabla_{\theta'} J(\theta') &= \mathbb{E}_{\tau \sim \bar{p}(\tau)} \left[\frac{p_{\theta}(\tau)}{\bar{p}(\tau)} \nabla_{\theta'} \log \pi_{\theta'}(\tau) r(\tau) \right] \\ &= \frac{1}{N} \sum_i \left(\sum_t \prod_{t=1}^T \frac{p_{\theta}(\tau)}{\bar{p}(\tau)} \right) \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left(\sum_t r(s_{i,t}, a_{i,t}) \right)\end{aligned}$$

This expression can be further advanced using the causality trick.

28.4 Covariant Policy Gradient

There exist other problems with policy gradient algorithms. For instance, suppose we have a one-dimensional action space, where the policy is Gaussian. Then, as the gradient with respect to the variance becomes larger due to the reduction of variance over time development, we will find the policy gradient not quite pointing the agent towards an optimum (destination). Therefore, following the policy gradient takes a long time until the optimal parameter is reached.

One approach to prevent this problem is known as covariant/natural policy gradient. To choose the learning rate α of our policy gradient step, we take into account that some parameters change probabilities a lot more than others, and therefore give learning rates corresponding to the sensitivity of each parameter. That is, we instead follow the step:

$$\theta' \leftarrow \arg \max_{\|\theta' - \theta\|_2^2 \leq \epsilon} (\theta' - \theta)^T \nabla_{\theta'} J(\theta)$$

We are free to replace the constraint $\|\theta' - \theta\|_2^2 \leq \epsilon$ with a different parameterization-independent divergence measure, such as $D_{KL}(\pi'_{\theta}, \pi_{\theta}) \leq \epsilon$. The quadratic Taylor expansion of divergence constraints further enable different optimization program formulations. Specifically, the quadratic Taylor expansion of KL divergence stated above grants us a policy gradient step:

$$\theta \leftarrow \theta + \alpha F^{-1} \nabla_{\theta} J(\theta)$$

where F is resulted from the use of KL divergence and known as the Fisher Information matrix:

$$F = \mathbb{E}_{\pi_{\theta}} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^T \right]$$

And we may solve for an optimal learning rate α while solving $F^{-1} \nabla_{\theta} J(\theta)$.

Chapter 29

Actor-Critic Algorithms

Chapter Description.

29.1 section title

Section.

Theorem 29.1.1. Tested Theorem

I am the bone of my sword.

Definition 29.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 29.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 30

Value Function Methods

Chapter Description.

30.1 section title

Section.

Theorem 30.1.1. Tested Theorem

I am the bone of my sword.

Definition 30.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 30.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 31

Deep RL with Q-Functions

Chapter Description.

31.1 section title

Section.

Theorem 31.1.1. Tested Theorem

I am the bone of my sword.

Definition 31.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 31.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 32

Advanced Policy Gradients

Chapter Description.

32.1 section title

Section.

Theorem 32.1.1. Tested Theorem

I am the bone of my sword.

Definition 32.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 32.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 33

Optimal Control and Planning

Chapter Description.

33.1 section title

Section.

Theorem 33.1.1. Tested Theorem

I am the bone of my sword.

Definition 33.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 33.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 34

Model-Based Reinforcement Learning

Chapter Description.

34.1 section title

Section.

Theorem 34.1.1. Tested Theorem

I am the bone of my sword.

Definition 34.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 34.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 35

Model-Based Policy Learning

Chapter Description.

35.1 section title

Section.

Theorem 35.1.1. Tested Theorem

I am the bone of my sword.

Definition 35.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 35.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 36

Exploration (Part 1)

Chapter Description.

36.1 section title

Section.

Theorem 36.1.1. Tested Theorem

I am the bone of my sword.

Definition 36.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 36.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 37

Exploration (Part 2)

Chapter Description.

37.1 section title

Section.

Theorem 37.1.1. Tested Theorem

I am the bone of my sword.

Definition 37.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 37.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 38

Offline Reinforcement Learning (Part 1)

Chapter Description.

38.1 section title

Section.

Theorem 38.1.1. Tested Theorem

I am the bone of my sword.

Definition 38.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 38.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 39

Offline Reinforcement Learning (Part 2)

Chapter Description.

39.1 section title

Section.

Theorem 39.1.1. Tested Theorem

I am the bone of my sword.

Definition 39.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 39.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 40

Reinforcement Learning Theory

Chapter Description.

40.1 section title

Section.

Theorem 40.1.1. Tested Theorem

I am the bone of my sword.

Definition 40.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 40.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 41

Variational Inference and Generative Models

Chapter Description.

41.1 section title

Section.

Theorem 41.1.1. Tested Theorem

I am the bone of my sword.

Definition 41.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 41.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 42

Connection between Inference and Control

Chapter Description.

42.1 section title

Section.

Theorem 42.1.1. Tested Theorem

I am the bone of my sword.

Definition 42.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 42.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 43

Inverse Reinforcement Learning

Chapter Description.

43.1 section title

Section.

Theorem 43.1.1. Tested Theorem

I am the bone of my sword.

Definition 43.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 43.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 44

Reinforcement Learning with Sequence Models

Chapter Description.

44.1 section title

Section.

Theorem 44.1.1. Tested Theorem

I am the bone of my sword.

Definition 44.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 44.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 45

Meta-Learning and Transfer Learning

Chapter Description.

45.1 section title

Section.

Theorem 45.1.1. Tested Theorem

I am the bone of my sword.

Definition 45.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 45.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 46

Challenges and Open Problems

Chapter Description.

46.1 section title

Section.

Theorem 46.1.1. Tested Theorem

I am the bone of my sword.

Definition 46.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 46.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Part III

COMPSCI 294-158: Deep Unsupervised Learning

Chapter 47

Introduction to Deep Unsupervised Learning

47.1 Introduction to the Subject

Deep unsupervised learning concerns capturing rich patterns in raw data with deep networks in a label-free way. This involves generative models, which recreate a raw data distribution, as well as self-supervised learning, which are puzzle tasks that require semantic understanding.

The expectation of required unsupervised learning (and a massive amount of it) comes from the large number of “synapses” that a brain must use to learn (which we equate to parameters). As expected, then, foundational models need a tremendous amount of information to build a generalization of semantic understanding, “common sense”. A notable model “LeCake” was proposed in this perspective, stating that a “cake” of machine learning is constructed as follows:

- Unsupervised/Predictive Learning: The body of cake, the most important part of the cake
- Supervised Learning: The icing of cake, for a generally smaller bit per sample.
- Reinforcement Learning: The cherry on top, for small bit per sample that finishes up an application.

An Ideal Intelligence that we refer to in this discipline discusses compression of dataset into a simple expression, and finding a pattern as a short description of raw data (which we call a low Kolmogorov Complexity), such as a summary vector. That is, we aim to present datasets in a compressed, efficient form. And, we also aim for optimal inference, which we know by Solomonoff Induction and demands induction via shortest code-length (with the model being referred to as a program of codes). We also want something that is extensible to optimal action making agents (AIXI). The ultimate demand, and a major assumption, follows as:

Assume we pretrain unsupervised on distribution \mathcal{D}_1 and then finetune on \mathcal{D}_2 . If \mathcal{D}_1 and \mathcal{D}_2 are related, then compressing \mathcal{D}_2 conditioned on \mathcal{D}_1 should be more efficient than compressing \mathcal{D}_2 directly. Therefore, pretraining accelerates learning.

Aside from theoretical interests, DUL has made powerful applications across generating novel data, conditional synthesis technology, compression of data, improving downstream tasks via pretraining, as well as being flexible building blocks.

Chapter 48

Autoregressive Models

Chapter Description.

48.1 section title

Section.

Theorem 48.1.1. Tested Theorem

I am the bone of my sword.

Definition 48.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 48.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 49

Flow Models

Chapter Description.

49.1 section title

Section.

Theorem 49.1.1. Tested Theorem

I am the bone of my sword.

Definition 49.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 49.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 50

Latent Variable Models

Chapter Description.

50.1 section title

Section.

Theorem 50.1.1. Tested Theorem

I am the bone of my sword.

Definition 50.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 50.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 51

GAN and Implicit Models

Chapter Description.

51.1 section title

Section.

Theorem 51.1.1. Tested Theorem

I am the bone of my sword.

Definition 51.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 51.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 52

Diffusion Models

Chapter Description.

52.1 section title

Section.

Theorem 52.1.1. Tested Theorem

I am the bone of my sword.

Definition 52.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 52.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 53

Self-Supervised Learning, Non-Generative Representation Learning

Chapter Description.

53.1 section title

Section.

Theorem 53.1.1. Tested Theorem

I am the bone of my sword.

Definition 53.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 53.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 54

League Language Models

Chapter Description.

54.1 section title

Section.

Theorem 54.1.1. Tested Theorem

I am the bone of my sword.

Definition 54.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 54.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 55

Video Generation

Chapter Description.

55.1 section title

Section.

Theorem 55.1.1. Tested Theorem

I am the bone of my sword.

Definition 55.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 55.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 56

Semi-Supervised Learning and Unsupervised Distribution Alignment

Chapter Description.

56.1 section title

Section.

Theorem 56.1.1. Tested Theorem

I am the bone of my sword.

Definition 56.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 56.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 57

Compression

Chapter Description.

57.1 section title

Section.

Theorem 57.1.1. Tested Theorem

I am the bone of my sword.

Definition 57.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 57.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 58

Multimodal Models

Chapter Description.

58.1 section title

Section.

Theorem 58.1.1. Tested Theorem

I am the bone of my sword.

Definition 58.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 58.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 59

Parallelization

Chapter Description.

59.1 section title

Section.

Theorem 59.1.1. Tested Theorem

I am the bone of my sword.

Definition 59.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 59.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 60

AI for Science (Gues Instructor)

Chapter Description.

60.1 section title

Section.

Theorem 60.1.1. Tested Theorem

I am the bone of my sword.

Definition 60.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 60.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Chapter 61

Neural Radiance Fields (Guest Instructor)

Chapter Description.

61.1 section title

Section.

Theorem 61.1.1. Tested Theorem

I am the bone of my sword.

Definition 61.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 61.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Part IV

Mu Li's Videos on YouTube

Chapter 62

Note Name

Chapter Description.

62.1 section title

Section.

Theorem 62.1.1. Tested Theorem

I am the bone of my sword.

Definition 62.1.1. Tested Theorem

Steel is my body and fire is my blood.

Example Question 62.1.1: Rules?

I have created over a thousand blades.

Unknown to death, nor known to life.

Revision Log

- August 22nd: Note is created