

# BDMA - Machine Learning

Jose Antonio Lorenzo Abril

Fall 2023



Professor: Tom Dupuis

Student e-mail: [jose-antonio.lorenco-abril@student-cs.fr](mailto:jose-antonio.lorenco-abril@student-cs.fr)

This is a summary of the course *Machine Learning* taught at the Université Paris Saclay - CentraleSupélec by Professor Tom Dupuis in the academic year 23/24. Most of the content of this document is adapted from the course notes by Dupuis, [1], so I won't be citing it all the time. Other references will be provided when used.

## Contents

<b>I</b>	<b>Deep Learning</b>	<b>4</b>
<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	AI History . . . . .	5
<b>2</b>	<b>Machine Learning Basics</b>	<b>6</b>
2.1	Linear Algebra Basics . . . . .	6
2.2	Probability Basics . . . . .	7
2.3	Machine Learning Basics . . . . .	9
<b>3</b>	<b>Deep Neural Networks</b>	<b>11</b>
3.1	Perceptron . . . . .	11
3.2	Multi-layer perceptron . . . . .	11
3.3	Cost Functions . . . . .	13
3.3.1	Choice of cost function . . . . .	13
3.3.2	Cross-entropy . . . . .	14
3.4	Why deep NN? . . . . .	15
3.5	Gradient-based Learning . . . . .	16
3.5.1	Back-propagation . . . . .	17
<b>II</b>	<b>Reinforcement Learning</b>	<b>19</b>
<b>4</b>	<b>Introduction</b>	<b>19</b>
<b>5</b>	<b>Definition and components</b>	<b>20</b>
5.1	Maximising the value by taking actions . . . . .	20
5.2	Markov Decision Processes . . . . .	21
5.3	Policies . . . . .	22
5.4	Value Functions . . . . .	22
5.4.1	Value Function Approximations . . . . .	22
5.5	Model . . . . .	23
5.6	Agent categories . . . . .	24
5.7	Subproblems of RL . . . . .	24

---

<b>6</b>	<b>Markov Decision Processes</b>	<b>24</b>
6.1	Solving Reinforcement Learning Problems with Bellman Equations . . . . .	28
6.2	Dynamic Programming . . . . .	30
6.2.1	Policy Evaluation . . . . .	30
6.2.2	Policy Improvement . . . . .	32
6.3	Extensions to Dynamic Programming . . . . .	33
6.3.1	Asynchronous Dynamic Programming . . . . .	33

## Part I

# Deep Learning

## 1 Introduction

**Artificial Intelligence** is a wide concept, encompassing different aspects and fields. We can understand the term AI as the multidisciplinary field of study that aims at recreating human intelligence using artificial means. This is a bit abstract, and, in fact, there is no single definition for what this means. Intelligence is not fully understood, and thus it is hard to assess whether an artificial invention has achieved intelligence, further than intuitively thinking so.

For instance, AI involves a whole variety of fields:

- Perception
- Knowledge
- Cognitive System
- Planning
- Robotics
- Machine Learning (Neural Networks)
- Natural Language Processing

Leveraging all of these, people try to recreate or even surpass human performance in different tasks. For example, a computer program that can play chess better than any human could ever possibly play, such as Stockfish, or a system that is able to understand our messages and reply, based on the knowledge that it has learnt in the past, such as ChatGPT and similar tools. Other examples are self-driving cars, auto-controlled robots, etc.

Therefore, AI is a very wide term, which merges many different scientific fields. **Machine Learning**, on the other side, is a narrower term, which deals with the study of the techniques that we can use to make a computer learn to perform some task. It takes concepts from Statistics, Optimization Theory, Computer Science, Algorithms, etc. A relevant subclass of Machine Learning, which has come to be one of the most prominent fields of research in the recent years, is **Neural Networks** or **Deep Learning**, which consists on an ML technique based on the human brain. Many amazing use cases that we see everywhere, like Siri (Apple assistant), Cortana (Windows assistant), Amazon recommender system, Dall-E (OpenAI image generation system), etc. Not only this, but the trend is growing, and the interest in DL is continuously increasing.

This is partly also due to the increase in computing resources, and the continuous optimization that different techniques are constantly experiencing. For instance, for a model trained on one trillion data points, in 2021 the training process required around 16500x less compute than a model trained in 2012.

But not everything is sweet and roses when using DL. Since these systems are being involved in decision making processes, there are some questions that arise, like whose responsibility is it when a model fails? Moreover, data is needed to train the models, so it is relevant to address how datasets should be collected, and to respect the privacy of the people that produce data. In addition, the recent technologies that are able to generate new content and to modify real content, make it a new issue that AI can create false information, mistrust, and even violence or paranoia.

Nonetheless, let's not focus on the negative, there are lots of nice applications of DL, and it is a key component to deal with data, achieving higher performance than traditional ML techniques for huge amount of data.

## 1.1 AI History

In 1950, Alan Turing aimed to answer the question '*Can machines think?*' through a test, which came to be named the **Turing Test**, and consists in a 3 players game. First, a similar game is the following: 2 talkers, a man and a female, and 1 interrogator. The interrogator asks questions to the talkers, with the aim of determining who is the man and who is the female. The man tries to trick the interrogator, while the woman tries to help him to identify her.

Then, the Turing Test consists in replacing the man by an artificial machine. Turing thought that a machine that could trick a human interrogator, should be considered intelligent.

Later, in 1956, in the Dartmouth Workshop organized by IBM, the term **Artificial Intelligence** was first used to describe *every aspect of learning or any other feature of intelligence can be so precisely described that a machine can be made to simulate it*.

From this year on, there was a focus on researching about **Symbolic AI**, specially in three areas of research:

- Reasoning as search: a different set of actions leads to a certain goal, so we can try to find the best choice of action to obtain the best possible outcome.
- Natural Language: different tools were developed, following grammar and language rules.
- Micro world: small block based worlds, that the system can identify and move.

In 1958, the **Perceptron** was conceived, giving birth to what is called the connectionism, an approach to AI based on the human brain, and a big hype that encouraged funding to support AI research. At this era, scientists experience a bit of lack of perspective, thinking that the power of AI was much higher than it was. For instance, H. A. Simon stated in 1965 that '*machines will be capable, within twenty years, of doing any work a man can do.*' We can relate to our time, with the huge hype that AI is experiencing, as well as the many apocalyptic theories that some people are making. Maybe we are again overestimating the power of AI.

The time from 1974 to 1980 is seen as the first winter of AI, in which research was slowed down and funding was reduced. This was due to several problems found at the time:

- There were few computational resources.
- The models at the time were not scalable.
- The Moravec's paradox: it is comparatively easy to make computers exhibit adult level performance on intelligence test or playing checkers, and difficult or impossible to give them the skills of a one-year-old when it comes to perception and mobility.
- Marvin Minsky made some devastating critics to connectionism, compared to symbolic, rule-based models:
  - Limited capacity: Minsky showed that single-layer perceptrons (a simple kind of neural network) could not solve certain classes of problems, like the XOR problem. While it was later shown that multi-layer perceptrons could solve these problems, Minsky's work resulted in a shift away from neural networks for a time.
  - Lack of clear symbols: Minsky believed that human cognition operates at a higher level with symbols and structures (like frames and scripts), rather than just distributed patterns of activation. He often argued that connectionist models lacked a clear way to represent these symbolic structures.
  - Generalization and Abstraction: Minsky was concerned that connectionist models struggled with generalizing beyond specific training examples or abstracting high-level concepts from raw data.
  - Inefficiency: Minsky pointed out that many problems which seemed simple for symbolic models could be extremely computationally intensive for connectionist models.
  - Lack of explanation: Connectionist models, especially when they become complex, can be seen as "black boxes", making it difficult to interpret how they arrive at specific conclusions.

- Over-reliance on learning: Minsky believed that not all knowledge comes from learning from scratch, and some of it might be innate or structured in advance. He felt connectionism put too much emphasis on learning from raw data.

In 1980, there was a boom in expert knowledge systems that made AI recover interest. An **expert system** solves specific tasks following an ensemble of rules based on knowledge facilitated by experts. A remarkable use case was the XCON sorting system, developed for the Digital Equipment Corporation, which helped them save 40M\$ per year. In addition, connectionism also came again on scene, thanks to the development of **backpropagation** applied to neurons, by Geoffrey Hinton. All these achievement made funding to come back to the field.

Nonetheless, there came a second winter of AI, from 1987 to 1994, mainly because several companies were disappointed and AI was seen as a technology that couldn't solve wide varieties of tasks. The funding was withdrawn from the field and a lot AI companies went bankrupt.

Luckily, from 1995 there started a new return of AI in the industry. The Moore's Law states that speed and memory of computer doubles every two years, and so computing power and memory was rapidly increasing, making the use of AI systems more feasible each year. During this time, many new concepts were introduced, such as **intelligent agents** as systems that perceive their environment and take actions which maximize their chances of success; or different **probabilistic reasoning tools** such as Bayesian networks, hidden Markov models, information theory, SVM,... In addition, AI researchers started to reframe their work in terms of mathematics, computer science, physics, etc., making the field more attractive for funding. A remarkable milestone during this time was the victory of Deep Blue against Garry Kasparov.

The last era of AI comes from 2011 to today, with the advent and popularization of **Deep Learning** (DL), which are deep graph processing layers mimicking human neurons interactions. This happened thanks to the advances of hardware technologies, that have enabled the enormous computing requirements needed for DL. The huge hype comes from the spectacular results shown by this kind of systems in a huge variety of tasks, such as computer vision, natural language processing, anomaly detection,...

In summary, we can see how the history of AI has been a succession of hype and dissapointment cycles, with many actors involved and the industry as a very important part of the process.

## 2 Machine Learning Basics

In this section, we review some notation, and basic knowledge of Linear Algebra, Probability and Machine Learning.

### 2.1 Linear Algebra Basics

A **scalar** is a number, either real and usually denoted  $x \in \mathbb{R}$ , or natural and denoted  $n \in \mathbb{N}$ . A **vector** is an array of numbers, usually real,  $x \in \mathbb{R}^n$ , or

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}.$$

A **matrix** is a 2-dimensional array of numbers,  $A \in \mathbb{R}^{n \times m}$ , or

$$A = \begin{bmatrix} A_{11} & \dots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{m1} & \dots & A_{mn} \end{bmatrix}.$$

A **tensor** is an  $n$ -dimensional array of numbers, for example  $A \in \mathbb{R}^{m \times k \times p}$  is a 3-dimensional tensor.

Usually, we will be working with matrices, which can be operated in different ways:

- Transposition:  $A^T$  is the transposed of  $A$ , defined as  $(A^T)_{ij} = A_{j,i}$ .
- Multiplication: Let  $A \in \mathbb{R}^{m \times k}$ ,  $B \in \mathbb{R}^{k \times n}$ , their multiplication,  $C \in \mathbb{R}^{m \times n}$  is defined as

$$C = A \cdot B = AB = (C_{ij})_{i \leq m, j \leq n} = \left( \sum_k A_{ik} B_{kj} \right)_{i \leq m, j \leq n}.$$

Note that the following holds for every matrix  $A, B$ :

$$(AB)^T = B^T A^T.$$

- Point-wise operations: if we have two matrices of the same size,  $A, B \in \mathbb{R}^{m \times n}$ , we can use apply scalar operator point-wise to each pair of elements in the same position in the two matrices. For example, the sum or the subtraction of matrices.

There are also special matrices:

- Identity matrix: the identity matrix is a square matrix that preserves any vector it is multiplied with. For vectors of size  $n$ , the identity matrix  $I_n$  verifies

$$I_n x = x, \forall x \in \mathbb{R}^n.$$

- Inverse matrix: the inverse of a square matrix,  $A \in \mathbb{R}^{n \times n}$ , when it exists, is defined as the only matrix  $A^{-1}$  such that

$$A^{-1} A = A A^{-1} = I_n.$$

Another important concept is that of the norm, which is basically measuring how far a point is from the origin of the space and can be used to measure distances:

**Definition 2.1.** A **norm** is a function  $f$  that measures the size of vectors, and must have the following properties:

- $f(x) = 0 \iff x = 0$ ,
- $f(x + y) \leq f(x) + f(y)$ , and
- $\forall \alpha \in \mathbb{R}, f(\alpha x) = |\alpha| f(x)$ .

A very important family of norms is the  $L^p$  norm, defined as

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}.$$

The **Euclidean norm** is the  $L^2$  norm, noted  $\|x\|$  and equivalent to computing  $\sqrt{x^T x}$ . In Machine Learning, it is not uncommon to find the use of the squared Euclidean norm, since it maintains the ordinals and is easier to operate with. The **Manhattan norm** is the  $L^1$  norm, and it is used when the difference between zero and nonzero elements is important. Finally, the **Max norm** is the  $L^\infty$ , or  $\|x\|_\infty = \max_i |x_i|$ .

## 2.2 Probability Basics

A **random variable**,  $X$ , is a variable that can take different values,  $x$ , randomly. They can be **discrete**, like the number drawn from a dice, or **continuous**, like the humidity in the air.

A probability distribution,  $p$ , is a **Probability Mass Function (PMF)** for discrete variables, and a **Probability Density Function (PDF)** for continuous random variables. It must satisfy:

- The domain of  $p$  describe all possible states of  $X$ .
- $\forall x \in X, p(x) \geq 0$ .
- $\int_{x \in X} p(x) dx = 1$ .

It is usual to have two (or more) random variables,  $X$  and  $Y$ , and to be interested in the probability distribution of their combination,  $p(x, y)$ . In this context, we define the **marginal probability** of the variable  $X$  as

$$p(X = x) = \int_{y \in Y} p(x, y) dy, \forall x \in X.$$

The **conditional probability** of the variable  $Y$  conditioned to  $X = x$  is

$$p(Y = y|X = x) = \frac{p(Y = y, X = x)}{P(X = x)}.$$

Finally, there is the **chain rule of conditional probabilities**, in which we start with  $n$  random variables,  $X_1, \dots, X_n$ , and it follows:

$$p(X_1 = x_1, \dots, X_n = x_n) = p(X_1 = x_1) \prod_{i=2}^n p(X_i = x_i | X_1 = x_1, \dots, X_{i-1} = x_{i-1}).$$

**Example 2.1.** For example, let's say  $X = \{1, 2, 3\}$ ,  $Y = \{1, 2\}$  and  $Z = \{1, 2\}$  with the following probabilities:

$X$	$Y$	$Z$	$p(x, y, z)$
1	1	1	$\frac{1}{6}$
1	1	2	$\frac{1}{6}$
1	2	1	$\frac{1}{12}$
1	2	2	$\frac{1}{24}$
2	1	1	$\frac{1}{24}$
2	1	2	$\frac{1}{20}$
2	2	1	$\frac{1}{20}$
2	2	2	$\frac{1}{20}$
3	1	1	$\frac{1}{6}$
3	1	2	$\frac{1}{12}$
3	2	1	$\frac{1}{20}$
3	2	2	$\frac{1}{20}$

Then, the marginal probabilities for the variable  $X$  are

$$\begin{aligned}
 P(X = 1) &= \frac{1}{6} + \frac{1}{6} + \frac{1}{12} + \frac{1}{24} = \frac{11}{24}, \\
 P(X = 2) &= \frac{1}{24} + \frac{1}{20} + \frac{1}{20} + \frac{1}{20} = \frac{23}{120}, \\
 P(X = 3) &= \frac{1}{6} + \frac{1}{12} + \frac{1}{20} + \frac{1}{20} = \frac{21}{60} = \frac{7}{20}.
 \end{aligned}$$

The conditional probability for the event  $\{Y = 1|X = 3\}$  is:

$$P(Y = 1|X = 3) = \frac{P(Y = 1, X = 3)}{P(X = 3)} = \frac{\frac{1}{6} + \frac{1}{12}}{\frac{7}{20}} = \frac{\frac{4}{12}}{\frac{7}{20}} = \frac{5}{7}.$$

The conditional probability for the event  $\{Z = 1|X = 3, Y = 1\}$  is:

$$P(Z = 1|X = 3, Y = 1) = \frac{P(X = 3, Y = 1, Z = 1)}{P(X = 3, Y = 1)} = \frac{\frac{1}{6}}{\frac{4}{12}} = \frac{2}{3}.$$



The probability of the event  $\{X = 3, Y = 1, Z = 1\}$  could be computed from the conditional probabilities as follows, in case we only knew these:

$$\begin{aligned} P(X = 3, Y = 1, Z = 1) &= P(X = 3) \cdot P(Y = 1|X = 3) \cdot P(Z = 1|X = 3, Y = 1) \\ &= \frac{7}{20} \cdot \frac{5}{7} \cdot \frac{2}{3} = \frac{10}{60} = \frac{1}{6}. \end{aligned}$$

When there are several variables, it is possible that the value of one of them is dependant, somehow, on the values that the other variables take; or that it is not:

**Definition 2.2.** Two random variables  $X$  and  $Y$  are **independent**, denoted  $X \perp Y$ , if  $\forall x \in X, y \in Y, p(X = x, Y = y) = p(X = x) \cdot p(Y = y)$ .

$X$  and  $Y$  are **conditionally independent** given the random variable  $Z$ , written  $X \perp_Z Y$  if  $\forall x \in X, y \in Y, z \in Z$ ,

$$p(X = x, Y = y|Z = z) = p(X = x|Z = z) \cdot p(Y = y|Z = z).$$

In Statistics and Machine Learning, there are some measures that summarize information about random variables, and that hold great importance.

**Definition 2.3.** The **expectation** of a function  $f(x)$  where  $x \sim p(x)$  is the average value of  $f$  over  $x$ :

$$\mathbb{E}_{x \sim p}[f(x)] = \int_{x \in X} p(x) f(x) dx.$$

The **variance** of  $f(x)$  measures how the values of  $f$  varies from its average:

$$\text{Var}[f(x)] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2],$$

and the **standard deviation** is the square root of the variance.

The **covariance** of two random variables provides information about how much two values are linearly related. More generally, if we apply two functions  $f(x)$ , where  $x \sim p(x)$ , and  $g(y)$ , where  $y \sim p(y)$ , the covariance between them is:

$$\text{Cov}[f(x), g(y)] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])].$$

## 2.3 Machine Learning Basics

To finalize with this review chapter, we are going to remember some basic concepts of Machine Learning.

First, let's give a definition of the concept:

**Definition 2.4.** A computer program is said to **learn** from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .

- The **task**  $T$  can be classification, regression, translation, generation, anomaly detection,...
- The **performance measure**  $P$  is specific to the tasks involved, and can be accuracy for classification, for example. It is measured on a **test set**.
- The **experience**  $E$  is divided into two main categories:
  - **Supervised learning**: a dataset of points associated with a label or a target determines the expected outcome of each event.

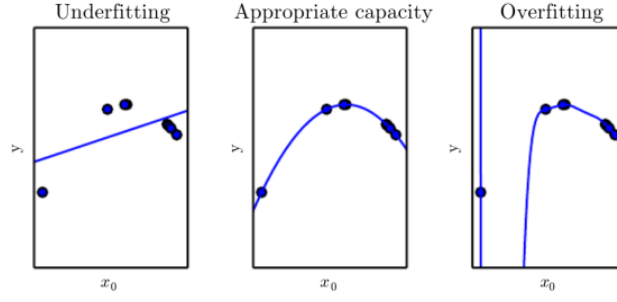


Figure 1: Appropriate capacity, overfitting and underfitting visualization.

- **Unsupervised learning:** a dataset of points without labels or targets, in which the desirable outcome needs to be defined in some different way.

Mathematically, we can formalize this as having a dataset of  $m$  points and  $k$  features, which can be represented as a matrix  $X \in \mathbb{R}^{m \times k}$ . In the case of supervised learning,  $X$  is associated with a vector of labels,  $y$ , and we aim to learn a joint distribution,  $p(X, y)$  to infer

$$p(Y = y | X = x) = \frac{p(x, y)}{\sum_{y'} p(x, y')}.$$

The goal is then to find a function  $\hat{f}$  that associates each  $x$  to the best approximation of  $y$ , and that is capable of generalizing to unseen data. Usually,  $\hat{f}$  is parameterized by a set of parameters,  $\theta$ , which are learnt during training.

The main challenge of an ML model is **generalization** to unseen data estimated on test data after the training on training data. **Overfitting** occurs when the gap between training error and test error is too large, while **underfitting** occurs when the training error is too large. The **capacity** of a model is the range of functions that it is able to learn and control how likely the model can overfit or underfit. This is visualized in Figure 1.

When we want to train a model, we will define the parameters that characterize it, and then we need to obtain the best possible of the parameters, according to the data. For this, we use estimators:

**Definition 2.5.** Given an unknown parameter  $\theta$ , we estimate it through an **estimator**,  $\hat{\theta}$ . A **point estimator** is a function of the data,  $X$ ,

$$\hat{\theta} = g(X).$$

The **bias** of an estimator is

$$\text{bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta.$$

An estimator is **unbiased** if  $\text{bias}(\hat{\theta}) = 0$ .

The **variance** of an estimator is  $\text{Var}(\hat{\theta})$ .

There are different ways to construct estimators, but one that is frequently used and that has solid mathematical foundations is the **maximum likelihood estimator**. Consider a dataset  $X = \{x_1, \dots, x_n\}$  and  $p(x; \theta)$  a parametric family of probability distribution that maps for each  $x$  the probability  $p_{\text{data}}(x)$ . This is, for each  $\theta$ ,  $p(x; \theta)$  is a probability density function. The maximum likelihood estimator is then

$$\begin{aligned} \theta_{ML} &= \arg \max_{\theta} p_{\text{model}}(X; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n p_{\text{model}}(x_i; \theta), \end{aligned}$$

considering that all instances of data are independent and identically distributed (iid). It is also a common practice to use the maximum **log**-likelihood instead, removing the product and avoiding floating point issues, since when the dataset is large, the product will rapidly go to 0. In addition, the logarithm does not modify the ordinals of the function. Therefore, we can use:

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^n \log(p_{model}(x_i; \theta)).$$

## 3 Deep Neural Networks

### 3.1 Perceptron

A deeper explanation of the perceptron can be read in my notes from another course, [https://lorenccio.github.io/BDMA\\_Notes/universities/UPC/Machine\\_Learning\\_summary.pdf](https://lorenccio.github.io/BDMA_Notes/universities/UPC/Machine_Learning_summary.pdf).

A perceptron is an algorithm for supervised learning of binary classifiers. That is, we have a dataset  $X \in \mathbb{R}^{n \times m}$  associated with a vector of labels  $y \in \{0, 1\}^n$ . Then, the perceptron learns a function  $\hat{f}$  parametrized by a vector of weights  $w \in \mathbb{R}^m$  and a bias  $b$ , such that:

$$\hat{f}(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ 0 & \sim \end{cases}.$$

Therefore, it is a linear classifier, which divides the input space into two regions separated by a hyperplane. This means that a perceptron cannot separate non-linear data.

### 3.2 Multi-layer perceptron

A deeper explanation of the MLP can be read in my notes from another course, [https://lorenccio.github.io/BDMA\\_Notes/universities/UPC/Machine\\_Learning\\_summary.pdf](https://lorenccio.github.io/BDMA_Notes/universities/UPC/Machine_Learning_summary.pdf).

When we say 'Deep' neural network, we refer to a series of stacked perceptrons. However, just like this, the model is still linear. This is why activation functions are introduced. An **activation function** is a function that is applied to the output of a perceptron, to make it non linear.

For example, ReLU is a piecewise-linear function defined as

$$ReLU(z) = \begin{cases} z & \text{if } z \geq 0 \\ 0 & \sim \end{cases} = \max\{z, 0\}.$$

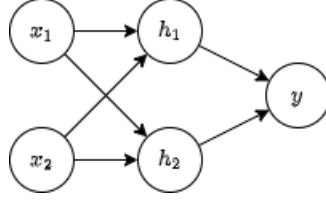
This function preserves much of the good optimization properties of a linear function, i.e., it is differentiable (apart from one point), and its derivative is constant.

**Example 3.1.** Learn the XOR function with a 2-layer MLP.

The XOR function is represented with the table:

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0

We want to use a 2-layer MLP to learn this function:



In  $h_1$ , it will be

$$w_{11}x_1 + w_{12}x_2 + b_1,$$

and in  $h_2$

$$w_{21}x_1 + w_{22}x_2 + b_2.$$

This can be represented as

$$h = W_h^T X + b_h.$$

Then, we apply ReLU

$$\max(0, W_h^T X + b_h),$$

and finally the output layer

$$y = W_y^T \max(0, W_h^T X + b_h) + b_y.$$

Let's see the different inputs:

$x_1$	$x_2$	$h$	$y$
0	0	$\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$	$\begin{pmatrix} w_1^y & w_2^y \end{pmatrix} \begin{pmatrix} \max(0, b_1) \\ \max(0, b_2) \end{pmatrix} + b_y \leq 0$
0	1	$\begin{pmatrix} w_{12} + b_1 \\ w_{22} + b_2 \end{pmatrix}$	$\begin{pmatrix} w_1^y & w_2^y \end{pmatrix} \begin{pmatrix} \max(0, w_{12} + b_1) \\ \max(0, w_{22} + b_2) \end{pmatrix} + b_y > 0$
1	0	$\begin{pmatrix} w_{11} + b_1 \\ w_{21} + b_2 \end{pmatrix}$	$\begin{pmatrix} w_1^y & w_2^y \end{pmatrix} \begin{pmatrix} \max(0, w_{11} + b_1) \\ \max(0, w_{21} + b_2) \end{pmatrix} + b_y > 0$
1	1	$\begin{pmatrix} w_{11} + w_{12} + b_1 \\ w_{21} + w_{22} + b_2 \end{pmatrix}$	$\begin{pmatrix} w_1^y & w_2^y \end{pmatrix} \begin{pmatrix} \max(0, w_{11} + w_{12} + b_1) \\ \max(0, w_{21} + w_{22} + b_2) \end{pmatrix} + b_y \leq 0$

A solution is:

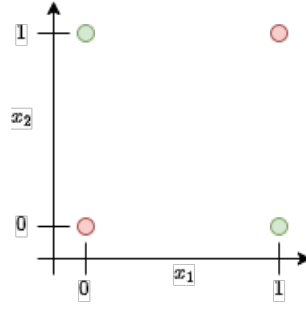
$$W_h = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, b_h = \begin{pmatrix} 0 \\ -1 \end{pmatrix}, W_y = \begin{pmatrix} 1 \\ -2 \end{pmatrix}, b_y = 0.$$

Let's check:

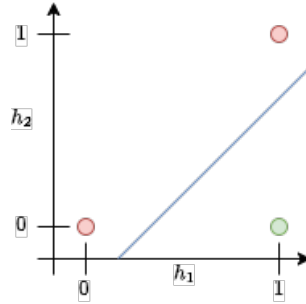
$x_1$	$x_2$	$h$	$y$
0	0	$\begin{pmatrix} 0 \\ -1 \end{pmatrix}$	$\begin{pmatrix} 1 & -2 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} = 0 \leq 0$
0	1	$\begin{pmatrix} 1 \\ 1 - 1 \end{pmatrix}$	$\begin{pmatrix} 1 & -2 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = 1 > 0$
1	0	$\begin{pmatrix} 1 \\ 1 - 1 \end{pmatrix}$	$\begin{pmatrix} 1 & -2 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = 1 > 0$
1	1	$\begin{pmatrix} 1 + 1 \\ 1 + 1 - 1 \end{pmatrix}$	$\begin{pmatrix} 1 & -2 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = 0 \leq 0$

So, it works! Note that this solution is not unique!

What happens is actually that the solution for the XOR problem is not linearly separable:



But, the hidden layer transforms this space, making the problem linearly separable, and therefore solvable in the last layer:



### 3.3 Cost Functions

The cost function is important when working with neural networks, because our goal is to ultimately train the model to solve some problem, and the cost functions will be the function that our model will aim at minimizing, thus guiding the training process.

Usually, we will need to choose a cost function  $\mathcal{L}$  that is suitable for our problem. Then, we will minimize  $\mathcal{L}$  with stochastic gradient descent, by:

- Training on a training dataset.
- Estimating error on an evaluation dataset.
- Computing the gradients using backpropagation.

In this process, we will aim to find good local minima, instead of global minimum. This is related to overfitting (learning only the training data, losing generalization capabilities), and to the empirical fact that deep neural network have surprisingly good local and non-global optima.

#### 3.3.1 Choice of cost function

In the general case, we use the maximum likelihood principle, taking the output types of the network into account. This means that we assume our dataset  $\{x_1, \dots, x_n\}$  to be independently and identically distributed (i.i.d.) from an unknown distribution,  $p_{data}(x)$ . We choose a parametric model family  $p_{model}(x; \theta)$  represented as a neural network, which we use to estimate an approximation of the true distribution. For this, we utilize the **maximum likelihood estimator**, defined as

$$\theta_{ML} = \arg \max_{\theta} \prod_{i=1}^n p_{model}(x_i; \theta).$$

Usually, to avoid floating point errors, the log-likelihood is used instead:

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^n \log p_{model}(x_i; \theta).$$

In the maximum likelihood estimation framework, we might apply activation functions to the output layer to get a desired structure for our distribution. This choice will also influence the mathematical form of the cost function. For example, we can use linear units for regression or for Gaussian distributions, sigmoid units for binary classification or softmax units for multi-class classification.

**Linear units for regression** A **linear output layer** is such that, given the features  $h$ , the output is

$$\hat{y} = W^T h + b,$$

where  $W$  is the weights vector and  $b$  the bias.

We can use this to predict real or vector valued variables, such as prices, biometrics,...

**Linear unit for Gaussian distribution** A **Gaussian output unit** is such that, given features  $h$ , a linear layer produces a vector  $\hat{y}$  representing the mean and the covariance matrix of a conditional Gaussian distribution:

$$p(y|x) = \mathcal{N}(y; \hat{y}, I).$$

Covariance is usually not modelled or simplified to be diagonal (in which case we need to ensure that the output is non-negative).

**Binary classification** In this case, the objective is to predict a binary variable,  $y$ : the neural network must predict  $P(y = 1|x)$ . Thus, we must ensure that the output is a probability, in the interval  $[0, 1]$ . For this, we can take

$$P(y = 1|x) = \max\{0, \min\{1, W^T h + b\}\}.$$

The problem with this approach is that if  $W^T h + b \notin [0, 1]$  then the gradient is 0 and the training will stop. To solve this issue, we can use a **sigmoid unit**, which is

$$\hat{y} = \sigma(W^T h + b) = \frac{1}{1 + e^{-W^T h + b}}.$$

**Softmax unit for multi-class classification** Now our objective is to classify the input data into one among  $N > 2$  classes. We want to predict  $\hat{y}$  with  $\hat{y}_i P(y = i|x)$ , subject to  $\hat{y}_i \in [0, 1], \forall i$  and  $\sum_i \hat{y}_i = 1$ .

In the output layer we can have  $N$  perceptrons, each of them computing  $z_i = \log P(y = i|x)$ , i.e., the **logits**. With this, we can apply the **softmax output unit** to all of them, obtaining our vector of probabilities, as

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}.$$

### 3.3.2 Cross-entropy

In classification problems we want to estimate the probability of different outcomes. Let the estimated probability of outcome  $i$  be  $p_{model}(x = i)$  with to-be-optimized parameters  $\theta$  and let the frequency of outcome  $i$  in the training set be  $p(x = i)$ . Given  $N$  conditionally independent samples in the training set, then the likelihood of the parameters  $\theta$  of the model  $p_{model}(x = i)$  on the training set is:

$$\mathcal{L}(\theta) = \prod_{i \in X} p_{model}(x = i)^{N \cdot p(x=i)}.$$

Therefore, the log-likelihood, divided by  $N$ , is

$$\frac{1}{N} \log (\mathcal{L}(\theta)) = \frac{1}{N} N \cdot \sum_{i \in X} p(x=i) \log p_{\text{model}}(x=i) = \sum_{i \in X} p(x=i) \log p_{\text{model}}(x=i).$$

Cross-entropy minimization is frequently used in optimization and rare-event probability estimation. When comparing a distribution  $q$  against a fixed reference distribution  $p$ , cross-entropy and KL divergence are identical up to an additive constant (since  $p$  is fixed): According to the Gibbs' inequality, both take on their minimal values when  $p = q$ , which is 0 for KL divergence, and  $H(p)$  for cross-entropy. In the engineering literature, the principle of minimizing KL divergence (Kullback's "Principle of Minimum Discrimination Information") is often called the Principle of Minimum Cross-Entropy (MCE), or Minxent.

### 3.4 Why deep NN?

Depth is the longest data path data can take from input to output. For a deep feed forward NN, depth is the number of hidden layers plus the output layer. State-of-the-art architectures used in practice have dozens to hundreds of layers.

**Theorem 3.1. Universal Approximation Theorem**

Let  $\varphi(\cdot)$  be a nonconstant, bounded, and monotonically increasing continuous function. Let  $I_{m_0}$  denote the  $m_0$ -dimensional unit hypercube,  $[0, 1]^{m_0}$ . The space of continuous functions on  $I_{m_0}$  is denoted by  $C(I_{m_0})$ .

Then, given any function  $f \in C(I_{m_0})$  and  $\varepsilon > 0$ , there exists an integer  $m_1$  and sets of real constants  $\alpha_i, b_i$  and  $w_{ij} \in \mathbb{R}$  where  $i = 1, \dots, m_1$  and  $j = 1, \dots, m_0$  such that we may define

$$F(x) = \sum_{i=1}^{m_1} \alpha_i \cdot \varphi \left( \sum_{j=1}^{m_0} w_{ij} \cdot x_j + b_i \right)$$

as an approximate realization of the function  $f$ , that is

$$|F(x) - f(x)| < \varepsilon, \forall x \in I_m.$$

This theorem is very relevant, because it says that for any mapping function  $f$  in supervised learning, there exists a MLP with  $m_1$  neurons in the hidden layer which is able to approximate it with a desired precision.

However, it only proves the existence of a shallow (just one hidden layer) MLP with  $m_1$  neurons in the hidden layer that can approximate the function, but it does not tell how to find this number.

As a rule of thumb for the generalization error, it is

$$\varepsilon = \frac{VC_{\text{dim}}(MLP)}{N},$$

where  $VC_{\text{dim}}$  is the Vapnik-Chervonenkis dimension, a measure of the capacity of a model. It refers to the largest set of points that the model can shatter. It is not easy to compute, but a rough upper bound for a FFNN is  $O(W \log W)$ , with  $W$  being the total number of weight in the network.

Also, this theorem hints us that having more neurons in the hidden layers will give us better training error, but worse generalization error: overfitting.

However, for most functions  $m_1$  is very high, and becomes quickly computationally intractable: so we need to go deeper.

**Theorem 3.2. No Free Lunch Theorem**

Multiple informal formulations:

- For every learning algorithm  $A$  and  $B$ , there are as many problems where  $A$  has a better generalization error than problems where  $B$  has a better one.
- All learning algorithms have the same generalization error if we average over all learning problems.
- There is no universally better learning algorithm.

**Depth Property**

The number of polygonal regions generated by a MLP with a ReLU function,  $d$  inputs,  $n$  neurons per hidden layer and  $l$  layers is

$$O\left(\binom{n}{d}^{d(l-1)} n^d\right).$$

This number grows exponentially with depth. This means that adding depth basically allows for more transformations of the input space.

**3.5 Gradient-based Learning**

The **gradient** of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , is  $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , defined at the point  $p = (x_1, \dots, x_n)$  as

$$\nabla f(p) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{pmatrix}.$$

This is, it's the local derivative or slope of each dimension at a certain point.

Going in the opposite direction of the gradient is a naïve but practical guess of the direction of the local minimum. This is the base for the gradient descent method.

**Gradient-descent method** (Cauchy, 1847)

A parametric function  $f(\theta)$  can be iteratively minimized by following the opposite direction of the gradient:

$$\theta_{t+1} = \theta_t - \varepsilon \nabla_{\theta} f(\theta),$$

where  $\varepsilon > 0$  is the **learning rate**.

We stop iterating when the gradient is near to 0.

Notice that this is useless if we have a close form for the gradient! In that case it is easier to just minimize it. This is useful when this is not the case, which always happens for neural networks.

In addition, there are variations to the method, for example, we can vary  $\varepsilon$  during training.



**Stochastic gradient descent**

Given a cost function  $f(\theta)$ , parameters of the network are updated with

$$\theta \leftarrow \theta - \varepsilon \nabla_{\theta} f(\theta).$$

For the negative log-likelihood (MLE), the function is:

$$f(\theta) = \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)}, \theta),$$

so the estimated gradient is

$$\nabla_{\theta} f(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(x^{(i)}, y^{(i)}, \theta).$$

The **problem** with this approach is that to take a single step of gradient descent, we must compute the loss over the whole dataset everytime, making the method not scalable at all. This is called **batch gradient descent**.

One **solution** is to compute the gradient with 1 sample only at each step, which is very noisy and inefficient, but works. This is the **stochastic gradient descent**.

In the middle ground, we find the **mini-batch gradient descent**, which divides the dataset into subsets, and updates the parameters after processing each of these subsets. A batch is a collection of samples used at each iteration for performing SDG in DL. A bigger batch provides a better gradient estimation, and therefore a faster learning, but also implies more device memory and slower descent.

Therefore, there is a tradeoff between money and performance at companies. In practice, a batch is set between 1 to 256 on one GPU.

But there is an even **greater problem**, the computation of the gradient is computationally very costly. To go around this problem, **back-propagation** was invented, as an efficient technique for gradient computation.

**3.5.1 Back-propagation**

The back-propagation algorithm is based on the chain rule for the derivative of composite functions: if we have  $y = g(x)$  and  $z = f(y) = f(g(x)) = (f \circ g)(x)$ , then

$$\frac{df}{dx}(x) = f'(g(x)) g'(x),$$

or, abusing notation,

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}.$$

This is generalized to multivariate functions as follows: let  $x \in \mathbb{R}^m, y \in \mathbb{R}^n, g: \mathbb{R}^m \rightarrow \mathbb{R}^n$  and  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ . If  $z = f(y) = f(g(x))$ , then

$$\frac{\partial z}{\partial x_i}(x) = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i},$$

or,

$$\nabla_x z = \left( \frac{\partial y}{\partial x} \right)^T \nabla_y z,$$

where  $\left( \frac{\partial y}{\partial x} \right)$  is the Jacobian of  $g$ .

Now, back-propagation is a recursive application of the chain rule, starting from the cost function. The algorithm works as follows:

1. **Forward pass**: a feedforward network takes as input  $x$  and produces the output  $\hat{y}$ . The information flows from layer to layer.

2. **Cost function:** compute the error between expected output and actual output.
3. **Back-propagate:** evaluate the individual gradient of each parameter and propagates them backwards to update them. For this, we use the concept of **local derivative**: the derivative of connected nodes are computed locally on the edges of the graph. For non-connected nodes, we multiply the edges connected between the nodes, and we sum over all incoming edges.

If we do this in a forward way, summing over all paths becomes intractable pretty quickly, while when doing it in a backwards way, it allows to obtain the derivative of the output with respect to every node directly in one pass. This leads to massive parallelization.

I did a more detailed explanation, with visualizations in my previous notes, [https://lorenccio.github.io/BDMA\\_Notes/universities/UPC/Machine\\_Learning\\_summary.pdf](https://lorenccio.github.io/BDMA_Notes/universities/UPC/Machine_Learning_summary.pdf).

## Part II

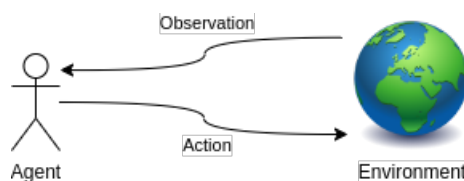
# Reinforcement Learning

## 4 Introduction

People and animals learn by interacting with the environment that surround them, differing from certain other types of learning. This process is active, rather than passive: the subject needs to perform interactions with the environment, to obtain knowledge. Also, the interactions are usually sequential, with future interactions possibly depending on earlier ones.

Not only this, but we are goal oriented: we act towards an objective. And, more importantly, we can learn without examples of optimal behavior! Instead, we optimise some reward signal obtained from the outcome of our actions.

It is in these observation that **Reinforcement Learning (RL)** arises as a learning paradigm, based on the **interaction loop**: there is an agent in an environment; the agent can make actions in the environment, and get observations from it.



RL relies on the reward hypothesis:

**Conjecture 4.1.** *Reward Hypothesis*

*Any goal can be formalized as the outcome of maximizing a cumulative reward.*

This hypothesis basically says that every objective that an agent can have, can be stated in terms of maximizing a reward associated to the actions of the agent with respect to this objective.

For example, if the objective of the agent is to fly a helicopter from point A to point B, then the reward could be negatively affected by the distance to point B, by the time taken to reach B,...

Now, it is important to realize that there exist different reasons to learn:

- Find solutions to problems.
- Adapt online to unforeseen circumstances.

Well, RL can provide algorithm for both cases! Note that the second point is not just about generalization, but also to cope with the so-called data shift, efficiently, during operation.

With all this, now we can define RL:

**Definition 4.1.** Reinforcement Learning is the science and framework of learning to make decisions from interaction.

This requires us to think about time, consequences of actions, experience gathering, future prediction, uncertainty,...

It has a huge potential scope and is a formalisation of the AI problem.

## 5 Definition and components

**Definition 5.1.** The **environment** is the world of the problem at hand, with **agents** in it that can perform actions, over time.

At each time step  $t$ , the agent:

- Receives observation  $O_t$  and reward  $R_t$  from the environment.
- Executes action  $A_t$ .

And the environment:

- Receives action  $A_t$ .
- Emits observation  $O_{t+1}$  and reward  $R_{t+1}$ .

But, what is a reward?

A **reward**,  $R_t$ , is a scalar feedback signal which indicates how well the agent is doing at step  $t$ : it defines how well the goal is being accomplished!

Therefore, the agent's job is to maximize the cumulative reward of the future steps, i.e.,

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$$

This is called the **return**. But, when one thinks about it carefully, one realizes that it is hard to know the future rewards with such precision. Therefore, it is also usual to use the **value**, which is the expected return, taking into account the current state,  $s$ :

$$v(s) = \mathbb{E}[G_t | S_t = s].$$

This depends on the actions the agents takes, and the goal is to **maximize it**! To achieve this, the agent must pick suitable actions.

Therefore, rewards and values define the utility of states and actions, and in this setup there is no supervised feedback.

Note, also, that this values can be defined recursively as

$$G_t = R_{t+1} + G_{t+1},$$

$$v(s) = \mathbb{E}[R_{t+1} + v(S_{t+1}) | S_t = s].$$

The **environment state** is the environment's internal state, which is usually invisible or partially visible to the agent. It is very important, but it can also contain lots of irrelevant information.

An environment is **fully observable** when the agent can see the full environment state, so every observation reveals the whole environment state. That is, the agent state could just be the observation:

$$S_t = O_t.$$

Note that  $S_t$  is the agent state, not the environment state!

### 5.1 Maximising the value by taking actions

As we have outlined, the goal is to select the actions that maximise the value. For this, we may need to take into account that actions may have long term consequences, delaying rewards. Thus, it may be better to sacrifice immediate reward to gain long-term reward.

The decision making process that for a given state chooses which action to take is called a **policy**.

To decide which action to take, we can also condition the value on actions:

$$q(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a],$$

so, for a given state  $s$ , a possible set of actions  $A_t^s$ , we could decide which action to take as

$$a_t = \max_{a \in A_t^s} q(s, a).$$

Then, the **history** is the full sequence of observation, actions and rewards:

$$\mathcal{H}_t = O_0, A_0, R_1, O_1, \dots, O_{t-1}, A_{t-1}, R_t, O_t.$$

## 5.2 Markov Decision Processes

Markov Decision Processes (MDPs) are a useful mathematical framework, defined as:

**Definition 5.2.** A decision process is Markov if

$$p(r, s | S_t, A_t) = p(r, s | \mathcal{H}_t, A_t).$$

This means that the current state is the only information needed to make a decision, we don't need the full story. For example, think in a chess game: there are many ways to arrive to a certain position, but it really does not matter how to got to the position, the past does not affect your choice now.

In order for a process to be Markov, full observability is required. When the situation is of **partial observability**, the observations are not Markovian, so using the observation as state is not enough to make the decision. This is called a **partially observable Markov decision process (POMDP)**. Note that the environment state can still be Markov, but the agent does not know it. In this case, we might be able to construct a Markov agent state.

In the general case, the agent state is a function of the history:

$$S_{t+1} = u(S_t, A_t, R_{t+1}, O_{t+1}),$$

where  $u$  is a **state update function**.

Usually, the agent state is much smaller than the environment state.

**Example 5.1.** A not Markov process:

Consider the following maze to be the full environment:



And consider the following observations:



They are indistinguishable! This process is not Markov, because only taking into account the current state, we cannot identify where we are.

To deal with partial observability, agent can construct suitable state representations. Some examples of agent states are:

- Last observation:  $S_t = O_t$  (might not be enough).
- Complete history:  $S_t = \mathcal{H}_t$  (might be too large).
- A generic update:  $S_t = u(S_{t-1}, A_{t-1}, R_t, O_t)$  (but how to design  $u$ ?)

Constructing a fully Markovian agent state is often not feasible and, more importantly, the state should allow for good policies and value predictions.

### 5.3 Policies

As we saw, a **policy** defines the agent's behavior: it is a map from the agent state to an action. Policies can be deterministic,

$$A = \pi(S),$$

or stochastic,

$$\pi(A|S) = p(A|S).$$

### 5.4 Value Functions

We saw the value before, which is the expected return. However, it is usual to introduce a **discount factor**,  $\gamma \in [0, 1]$ , which trades off importance of immediate and long-term rewards. This way, the value becomes

$$v_\pi(s) = \mathbb{E}[G_t | S_t = s, \pi] = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s, \pi].$$

The value depends on the policy,  $\pi$ , and can be used to evaluate the desirability of states, as well as to select between actions.

Note the role of the discount factor: the higher it is, the higher the focus on long term outcomes.

Now, using the recursive expression of the return,  $G_t = R_{t+1} + \gamma G_{t+1}$ , we can rewrite the value as

$$v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t \sim \pi(s)],$$

where  $A \sim \pi(s)$  means  $A$  is chosen by policy  $\pi$  in state  $s$ . This is known as a **Bellman equation**. A similar equation holds for the optimal value, i.e., the highest possible value:

$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a].$$

Note how this does not depend on a policy, it is just the maximum achievable value from the current state.

#### 5.4.1 Value Function Approximations

Agents often approximate value functions, and with an accurate value function approximation, the agent can behave optimally, or very well, even in intractably big domains.

## 5.5 Model

A **model** predicts what the environment will do next. For example,  $\mathcal{P}$  predicts the next state, given the current state and an action:

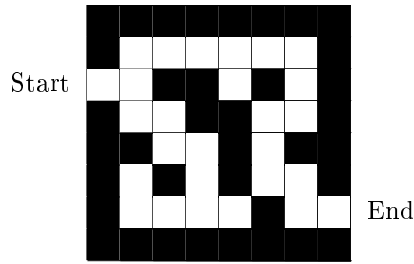
$$\mathcal{P}(s, a, s') \approx p(S_{t+1} = s' | S_t = s, A_t = a).$$

Or  $\mathcal{R}$  predicts the next immediate reward:

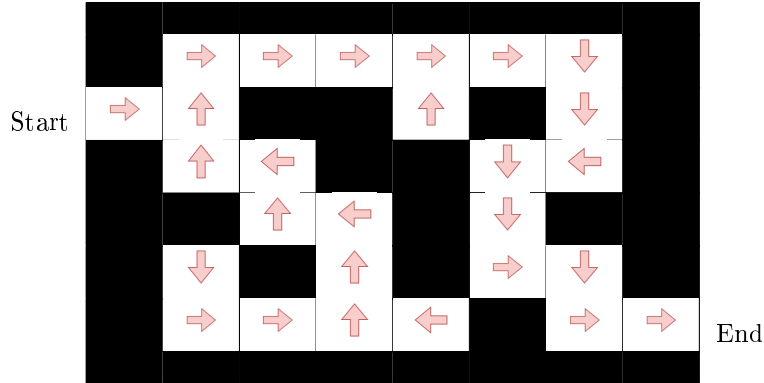
$$\mathcal{R}(s, a) \approx \mathbb{E}[R_{t+1} | S_t = s, A_t = a].$$

Note that a model does not immediately give us a good policy! We still need to plan and see how actions and states are related.

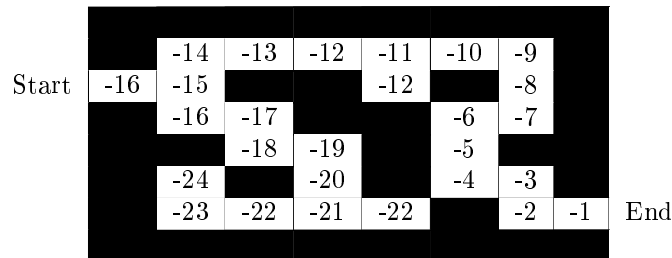
**Example 5.2.** Consider the following maze, where the rewards are -1 per time-step, the actions are to go N, E, S and W, and the states are the agent's location:



The following arrows represent the policy,  $\pi(s)$ , for each state  $s$ :



In the following one, the numbers represent the value  $v_\pi(s)$  of each state  $s$ :



The grid layout represents the partial transition model  $\mathcal{P}_{ss'}^a$ , and numbers represent the immediate reward,  $\mathcal{R}_{ss'}^a$ , from each state  $s$ , which is -1 for all  $a$  and  $s'$  in this case.

## 5.6 Agent categories

An agent is **model free** when the behavior of the environment is not known. The agent needs a policy or a value function to operate and there is no model. On the other hand, it is **model based** when the environment is known by means of a model. In this case, a policy and a value function might be optional, since it is possible that the agent can operate just knowing the model.

Model free agents are simpler, while model based agents are more sample efficient.

Another categorization is the following:

- Value based: there is no policy, it is implicit in the value function.
- Policy based: there is no value function, the model operates only by means of the policy.
- Actor critic: they have both a policy and a value function.

## 5.7 Subproblems of RL

**Prediction** consists in evaluating the future, for a given policy, i.e., what are the values in each state?

**Control** refers to the problem of optimising the future to find the best policy, i.e., which actions to take?

These two problems are strongly related, because the best actions to take will be decided using our predictions about the future:

$$\pi_*(s) = \arg \max_{\pi} v_{\pi}(s).$$

Two fundamental problems in RL are:

- Learning: the environment is initially unknown and the agent interacts with it to learn.
- Planning/search: a model of the environment is given or learnt, and the agent plans in this model.

In order to learn, we need to define all components of the problem as functions:

- Policy:  $\pi : S \rightarrow A$  (or probabilities over  $A$ ).
- Value functions:  $v : S \rightarrow \mathbb{R}$ .
- Models:  $p : S \rightarrow S$  or  $r : S \rightarrow \mathbb{R}$ .
- State update:  $u : S \times O \rightarrow S$ .

Then, we can use, for example, neural networks and deep learning techniques to learn. But we do need to be careful, because in RL it is usual to violate assumptions made in supervised learning, such as having i.i.d. samples, or stationarity.

## 6 Markov Decision Processes

We saw the notion of MDP, and now we formalize it:

**Definition 6.1.** A **Markov Decision Process (MDP)** is a tuple  $(S, A, p, \gamma)$  where:

- $S$  is the set of all possible states with the Markov Property.
- $A$  is the set of all possible actions.
- $p(r, s'|s, a)$  is the joint probability of a reward,  $r$ , and next state,  $s'$ , given a state  $s$  and an action  $a$ .
- $\gamma \in [0, 1]$  is a discount factor that trades off later rewards to earlier ones.



*Remark 6.1.*  $p$  defines the dynamics of the problem.

Sometimes, it is useful to marginalise out the state transitions or expected rewards:

$$p(s'|s, a) = \sum_r p(s', r|s, a),$$

to obtain the probability of arriving to a certain state.

Also, the expected reward:

$$\mathbb{E}[R|s, a] = \sum_r r \sum_{s'} p(r, s'|s, a).$$

There is an alternative equivalent definition, which introduces the notion of the expected reward into the concept, and takes it out of the probability function:

**Definition 6.2.** A MDP is a tuple  $(S, A, p, r, \gamma)$  where:

- $S$  is the set of all possible states with the Markov Property.
- $A$  is the set of all possible actions.
- $p(s'|s, a)$  is the probability of transitioning to  $s'$ , given a state  $s$  and an action  $a$ .
- $r : S \times A \rightarrow \mathbb{R}$  is the expected reward, achieved on a transition starting in  $(s, a)$ ,

$$r = \mathbb{E}[R|s, a].$$

- $\gamma \in [0, 1]$  is a discount factor that trades off later rewards to earlier ones.

Now, we have to clarify what is the Markov Property:

**Definition 6.3.** Consider a sequence of random variables,  $\{S_t\}_{t \in \mathbb{N}}$ , indexed by time and taken from a set of states  $S$ . Consider also the set of actions  $A$  and rewards in  $\mathbb{R}$ .

A state  $s$  has the **Markov Property** when, for all  $s' \in S$ ,

$$p(S_{t+1} = s' | S_t = s) = p(S_{t+1} = s' | h_{t-1}, S_t = s),$$

for all possible histories  $h_{t-1} = \{S_1, \dots, S_{t-1}, A_1, \dots, A_{t-1}, R_1, \dots, R_{t-1}\}$ .

Therefore, in an MDP, the current state captures all relevant information from the history, it is a sufficient statistic of the past. So, once the state is known, the history may be thrown away.

**Exercise 6.1.** In an MDP, which of the following statements are true?

1.  $p(S_{t+1} = s' | S_t = s, A_t = a) = p(S_{t+1} = s' | S_1, \dots, S_{t-1}, A_1, \dots, A_t, S_t = s)$ : false, the RHS does not condition on  $A_t = a$ .
2.  $p(S_{t+1} = s' | S_t = s, A_t = a) = p(S_{t+1} = s' | S_1, \dots, S_{t-1}, S_t = s, A_t = a)$ : true.
3.  $p(S_{t+1} = s' | S_t = s, A_t = a) = p(S_{t+1} = s' | S_1, \dots, S_{t-1}, S_t = s)$ : false, the RHS does not condition on  $A_t = a$ .
4.  $p(R_{t+1} = r, S_{t+1} = s' | S_t = s) = p(R_{t+1} = r, S_{t+1} = s' | S_1, \dots, S_{t-1}, S_t = s)$ : true.

It is also worth noting that most MDPs are discounted, and there are several reasons for these:

- Problem specification: immediate rewards may actually be more valuable. For instance, animal/human behavior shows preference for immediate reward.

- Solution side: it is mathematically convenient to discount rewards, because it allows for easier proofs of convergence, and avoids infinite returns in cyclic Markov processes.

As we outlined previously, the **goal of an RL agent** is to find a behavior policy that maximises the expected return  $G_t$ . Recall our definition for value function

$$v_\pi(s) = \mathbb{E}[G_t | S_t = s, \pi].$$

Similarly, we can define the **state-action values**, as

$$q_\pi(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a, \pi],$$

and there is the following connection between them:

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a) = \mathbb{E}[q_\pi(S_t, A_t) | S_t = s, \pi], \forall s \in S.$$

Also, we can define the maximum possible value functions:

**Definition 6.4.** The **optimal state value function**,  $v^*(s)$ , is the maximum value function over all policies,

$$v^*(s) = \max_{\pi} v_\pi(s).$$

The **optimal state-action value function**,  $q^*(s, a)$ , is the maximum state-action value function over all policies,

$$q^*(s, a) = \max_{\pi} q_\pi(s, a).$$

The optimal value function specifies the best possible performance in the MDP. We can consider the MDP to be solved when we know the optimal value function.

In addition, value functions allow us to define a partial ordering over policies, having

$$\pi \geq \pi' \iff v_\pi(s) \geq v_{\pi'}(s), \forall s \in S.$$

With this partial ordering, the following theorem state that optimal policies exist for every MDP:

**Theorem 6.1. Optimal Policies Theorem**

For any MDP:

- There exists an optimal policy  $\pi^*$  that is better than or equal to all other policies,

$$\pi^* \geq \pi, \forall \pi.$$

- All optimal policies achieve the optimal value function,

$$v^{\pi^*}(s) = v^*(s).$$

- All optimal policies achieve the optimal state-action value function,

$$q^{\pi^*}(s, a) = q^*(s, a).$$

To find an optimal policy, we can maximise over  $q^*(s, a)$ :

$$\pi^*(s, a) = \begin{cases} 1 & \text{if } a = \arg \max_{a \in A} q^*(s, a) \\ 0 & \text{otherwise} \end{cases}.$$

That is, the optimal policy is to take action  $a$  in state  $s$  if  $a$  is the action that gives the highest state-action value given state  $s$ .

*Remark 6.2.* There is always a deterministic optimal policy for any MDP, and we know  $q^*(s, a)$ , we know the optimal policy immediately.

Also, there can be multiple optimal policies, and if multiple actions maximize  $q_*(s, \cdot)$ , we can pick any of them.

Now, recall the Bellman Equations we saw previously. The following theorem explains how to express the value functions by means of these equations:

**Theorem 6.2. Bellman Expectation Equations**

Given an MDP,  $M = (S, A, p, r, \gamma)$ , for any policy  $\pi$ , the value functions obey the following expectation equations:

$$v_\pi(s) = \sum_a \pi(s, a) \left[ r(s, a) + \gamma \sum_{s'} p(s'|a, s) v_\pi(s') \right],$$

and

$$q_\pi(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|a, s) \sum_{a' \in A} \pi(a'|s') q_\pi(s', a').$$

**Theorem 6.3. Bellman Optimality Equations**

Given an MDP,  $M = (S, A, p, r, \gamma)$ , the optimal value functions obey the following expectation equations:

$$v^*(s) = \max_{a \in A} \left[ r(s, a) + \gamma \sum_{s'} p(s'|a, s) v^*(s') \right],$$

$$q^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|a, s) \max_{a' \in A} q^*(s', a').$$

*Remark 6.3.* There can not be a policy with a higher value than  $v^*(s) = \max_\pi v_\pi(s)$ ,  $\forall s$ .

*Intuition on the proof for the Bellman Optimality Equations:*

An optimal policy can be found by maximising over  $q^*(s, a)$ ,

$$\pi^*(s, a) = \begin{cases} 1 & \text{if } a = \arg \max_{a \in A} q^*(s, a), \\ 0 & \text{otherwise.} \end{cases}$$

Applyin the Bellman Expectation Equation:

$$\begin{aligned} q_{\pi^*}(s, a) &= r(s, a) + \gamma \sum_{s'} p(s'|a, s) \sum_{a' \in A} \pi^*(a'|s') q_{\pi^*}(s', a') \\ &= r(s, a) + \gamma \sum_{s'} p(s'|a, s) \max_{a' \in A} q^*(s', a'). \end{aligned}$$

We can also express the Bellman equations in matrix form, as

$$V = R^\pi + \gamma P^\pi V,$$

where  $v_i = v(s_i)$ ,  $R_i^\pi = \mathbb{E}[R_{t+1}|S_t = s_i, A_t \sim \pi(S_t)]$  and  $P_{ij}^\pi = p(s_j|s_i) = \sum_a \pi(a|s_i) p(s_j|s_i, a)$ .

This is a linear equations, that can be solved directly:

$$\begin{aligned} (I - \gamma P^\pi) V &= R^\pi \\ V &= (I - \gamma P^\pi)^{-1} R^\pi. \end{aligned}$$

The computational complexity is  $O(|S|^3)$ , making this only feasible for small problems. This makes it helpful to design other methods for larger problems. For example, there are iterative methods such as dynamic programming, monte-calro evaluation and temporal-difference learning.

## 6.1 Solving Reinforcement Learning Problems with Bellman Equations

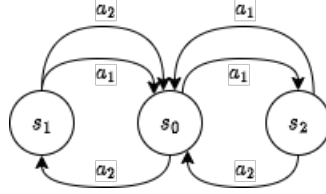
Pb1. Estimating  $v_\pi$  or  $q_\pi$  is called **policy evaluation**, or **prediction**:

- (a) Given a policy, what is my expected return under that behavior?
- (b) Given this treatment protocol/trading strategy, what is my expected return?

Pb2. Estimating  $v_*$  or  $q_*$  is sometimes called **control**, because these can be used for **policy optimisation**:

- (a) What is the optimal way of behaving? What is the optimal value function?
- (b) What is the optimal treatment? What is the optimal control policy to minimise time, fuel consumption, etc?

**Exercise 6.2.** Consider the following MDP:



The actions have a 0.9 probability of success and with 0.1 probability we remain in the same state.

$R_t = 0$  for all transitions that end up in  $S_0$  and  $R_t = -1$  for all other transitions.

Discount factor:  $\gamma = 0.9$ .

Questions:

- What is  $v_\pi$  for  $\pi(s) = a_1(\rightarrow), \forall s$ ?

According to the Bellman Equations:

$$v_\pi(s) = \sum_a \pi(s, a) \left[ r(s, a) + \gamma \sum_{s'} p(s'|a, s) v_\pi(s') \right],$$

$$\begin{aligned} v_\pi(s_0) &= \pi(s_0, a_1 \checkmark) [r(s_0, a_1 \checkmark) + 0.9 \cdot [p(s_1|a_1, s_0) v_\pi(s_1) + p(s_2|a_1, s_0) v_\pi(s_2)]] + \pi(s_0, a_1 \times) \cdot [r(s_0, a_1 \times) + 0.9 \cdot p(s_0|a_1, s_0) v_\pi(s_0)] \\ &= 0.9 \cdot [-1 + 0.9 \cdot [0 + v_\pi(s_2)]] + 0.09 \cdot v_\pi(s_0) \\ &= -0.9 + 0.81 \cdot v_\pi(s_2) + 0.09 \cdot v_\pi(s_0). \end{aligned}$$

Isolating,  $v_\pi(s_0)$ , we get

$$0.91 v_\pi(s_0) = -0.9 + 0.81 \cdot v_\pi(s_2) \implies v_\pi(s_0) = -0.99 + 0.89 \cdot v_\pi(s_2).$$

Now, let's go for  $v_\pi(s_2)$ :

$$\begin{aligned} v_\pi(s_2) &= \pi(s_2, a_1 \checkmark) [r(s_2, a_1 \checkmark) + 0.9 \cdot [p(s_0|a_1, s_2) v_\pi(s_0) + p(s_1|a_1, s_2) v_\pi(s_1)]] + \pi(s_2, a_1 \times) \cdot [-1 + 0.9 \cdot v_\pi(s_2)] \\ &= 0.9 \cdot [0 + 0.9 \cdot [v_\pi(s_0) + 0]] + 0.1 [-1 + 0.9 \cdot v_\pi(s_2)] \\ &= 0.81 \cdot v_\pi(s_0) - 0.1 + 0.09 \cdot v_\pi(s_2) \\ &= 0.81 \cdot (-1 + 0.9 \cdot v_\pi(s_2)) - 0.1 + 0.09 \cdot v_\pi(s_2) \\ &= -0.91 + 0.82 \cdot v_\pi(s_2). \end{aligned}$$

Therefore,

$$v_\pi(s_2) = \frac{-0.91}{0.18} = -5.06.$$

This means that

$$v_{\pi}(s_0) = -5.49.$$

Finally,

$$\begin{aligned} v_{\pi}(s_1) &= \pi(s_1, a_1) [r(s_1, a_1) + 0.9 \cdot [p(s_0|a_1, s_1) v_{\pi}(s_0) + p(s_2|a_1, s_1) v_{\pi}(s_2)]] + \pi(s_1, a_1 \times) \cdot [-1 + 0.9 \cdot v_{\pi}(s_1)] \\ &= 0.9 \cdot [0 + 0.9 \cdot [-5.49 + 0]] + 0.1 \cdot [-1 + 0.9 \cdot v_{\pi}(s_1)] \\ &= -4.45 - 0.1 + 0.09 \cdot v_{\pi}(s_1) . \\ &= -4.55 + 0.09 \cdot v_{\pi}(s_1) . \end{aligned}$$

So

$$v_{\pi}(s_1) = \frac{-4.55}{0.91} = -5.$$

That is,  $v_{\pi}(s_0) = -5.49$ ,  $v_{\pi}(s_1) = -5$  and  $v_{\pi}(s_2) = -5.06$ .

- What is  $v_{\pi}$  for the uniformly random policy?

This is: with probability 0.1, stay in the same state, with probability 0.45 choose  $a_1$  and with probability 0.45 choose  $a_2$ . For  $s_0$ :

$$\begin{aligned} v_{\pi}(s_0) &= 0.1 \cdot 0.9 \cdot v_{\pi}(s_0) + 0.45 \cdot [-1 + 0.9 \cdot v_{\pi}(s_2)] + 0.45 \cdot [-1 + 0.9 \cdot v_{\pi}(s_1)] \\ &= 0.09 \cdot v_{\pi}(s_0) - 0.9 + 0.405 \cdot v_{\pi}(s_2) + 0.405 \cdot v_{\pi}(s_1) , \end{aligned}$$

so

$$v_{\pi}(s_0) = -0.99 + 0.445 \cdot v_{\pi}(s_2) + 0.445 \cdot v_{\pi}(s_1) .$$

For  $s_1$ :

$$\begin{aligned} v_{\pi}(s_1) &= 0.1 \cdot [-1 + 0.9 \cdot v_{\pi}(s_1)] + 0.45 \cdot [0 + 0.9 \cdot v_{\pi}(s_0)] + 0.45 \cdot [0 + 0.9 \cdot v_{\pi}(s_0)] \\ &= -0.1 + 0.09 \cdot v_{\pi}(s_1) + 0.81 \cdot v_{\pi}(s_0) , \end{aligned}$$

so

$$v_{\pi}(s_1) = -0.1 + 0.9 \cdot v_{\pi}(s_0) .$$

For  $s_2$ :

$$\begin{aligned} v_{\pi}(s_2) &= 0.1 \cdot [-1 + 0.9 \cdot v_{\pi}(s_2)] + 0.45 \cdot [0 + 0.9 \cdot v_{\pi}(s_0)] + 0.45 \cdot [0 + 0.9 \cdot v_{\pi}(s_0)] \\ &= -0.1 + 0.09 \cdot v_{\pi}(s_2) + 0.81 \cdot v_{\pi}(s_0) , \end{aligned}$$

so

$$v_{\pi}(s_2) = -0.1 + 0.9 \cdot v_{\pi}(s_0) .$$

Therefore

$$\begin{aligned} v_{\pi}(s_0) &= -0.99 + 0.89 \cdot (-0.1 + 0.9 \cdot v_{\pi}(s_0)) \\ &= -0.99 - 0.089 + 0.8 \cdot v_{\pi}(s_0) \\ &= -1.08 + 0.8 \cdot v_{\pi}(s_0) . \end{aligned}$$

That is,

$$v_{\pi}(s_0) = -5.4.$$

And,

$$v_{\pi}(s_1) = v_{\pi}(s_2) = -4.96.$$

- Same policy evaluation problems for  $\gamma = 0$ ? What do you notice?

First,  $\pi(s) = a_1 (\rightarrow), \forall s$ :

$$\begin{aligned} v_{\pi}(s_0) &= 0.1 \cdot 0 + 0.9 \cdot (-1) \\ &= -0.9. \end{aligned}$$

$$\begin{aligned} v_{\pi}(s_1) &= 0.1 \cdot (-1) + 0.9 \cdot (0) \\ &= -0.1. \end{aligned}$$

$$v_{\pi}(s_2) = -0.1.$$

Second,  $\pi(s)$  the random policy:

$$\begin{aligned} v_{\pi}(s_0) &= 0.1 \cdot 0 + 0.45 \cdot (-1) + 0.45 \cdot (-1) \\ &= -0.9. \end{aligned}$$

$$\begin{aligned} v_{\pi}(s_1) &= 0.1 \cdot (-1) + 0.9 \cdot (0) \\ &= -0.1. \end{aligned}$$

$$v_{\pi}(s_2) = -0.1.$$

We can observe that if we don't take the discount factor into account, two very different policies can give us the same (short term) values.

## 6.2 Dynamic Programming

Dynamic programming refers to a collection of algorithms that can be used to compute optimal policies given a perfect model of the environment as a Markov Decision Process (MDP).

-Sutton & Barto, 2018

We will discuss several dynamic programming methods to solve MDPs, all of which consist of two important parts:

- Policy evaluation.
- Policy improvement.

### 6.2.1 Policy Evaluation

We start by discussing how to estimate

$$v_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | s, \pi].$$

The idea is to turn the equality into an update rule. The process is described in the following algorithm:

```

1 def policy_eval(S, R, pi, gamma):
2
3     v = [0 for s in S]
4
5     repeat until converge:
6         for s in S:
7             v_new(s) = E[R + gamma * v(S_t+1) | S_t = s, pi]
8
9         v = v_new
10
11     return v

```

Note that this algorithm always converge under appropriate conditions, like  $\gamma < 1$ . We will delve into this later.

**Example 6.1.** Policy Evaluation example.

Take the following MDP:

	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

The possible actions are to go up, down, left, right and when you reach the red cells, you finish. Each transitions costs -1 point.

Let's evaluate the random policy with  $\gamma = 1$ .

Initialization:

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

Step 1: we do

$$v(s) = \mathbb{E}[R + \gamma v(S) | S, \pi],$$

so for example, for cell (1), it is:

$$v(1) = \frac{1}{3}(-1 + 1 \cdot 0) \cdot 3 = -1.$$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

Step 2:

$$v(1) = \frac{1}{3}(-1 + 1 \cdot 0) + \frac{2}{3}(-1 - 1) = -1.7.$$

0	-1.7	-2	-2
-1.7	-2	-2	-2
-2	-2	-2	-1.7
-2	-2	-1.7	0

Step 3:

$$v(1) = \frac{1}{3}(-1 + 0) + \frac{2}{3}(-1 - 2) = -2.3.$$

0	-2.3	-2.9	-3
-2.3	-2.9	-3	-2.9
-2.9	-3	-2.9	-2.3
-3	-2.9	-2.3	0

And so on... It would converge at

0	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	0

### 6.2.2 Policy Improvement

We can use the values computed with policy evaluation to improve the policy. The simplest way to achieve this is with a **greedy policy improvement** approach, which is as follows:

```

1 def policy_improvement(S, R, pi, gamma)
2     pi_new = {}
3
4     v = [0 for s in S]
5
6     for s in S:
7         v(s) = policy_eval(S, R, pi, gamma)
8         pi_new[v] = argmax_a E[R + gamma*v(S_t+1) | S_t = s, A_t = a]
9         pi = pi_new
10
11     return pi_new

```

*Claim 6.1.* It is possible to show that

$$v_{\pi_{new}}(s) \geq v_{\pi}(s), \forall s.$$

**Example 6.2.** We can use this greedy approach combined with the previous example:

Initialization, with random policy:

0	0	0	0	0	↓, ←, →	↓, ←, →	↓, ←
0	0	0	0	↑, ↓, →	↑, ↓, ←, →	↑, ↓, ←, →	↑, ↓, ←
0	0	0	0	↑, ↓, →	↑, ↓, ←, →	↑, ↓, ←, →	↑, ↓, ←
0	0	0	0	↑, →	↑, ←, →	↑, ←, →	

Step 1:

0	-1	-1	-1	0	←	↓, ←, →	↓, ←
-1	-1	-1	-1	↑	↑, ↓, ←, →	↑, ↓, ←, →	↑, ↓, ←
-1	-1	-1	-1	↑, ↓, →	↑, ↓, ←, →	↑, ↓, ←, →	↓
-1	-1	-1	0	↑, →	↑, ←, →	→	

Step 2:

0	-1.7	-2	-2	0	←	←	↓, ←
-1.7	-2	-2	-2	↑	↑, ←	↑, ↓, ←, →	↓
-2	-2	-2	-1.7	↑	↑, ↓, ←, →	↓, →	↓
-2	-2	-1.7	0	↑, →	→	→	

Step 3:

0	-2.3	-2.9	-3	0	←	←	↓, ←
-2.3	-2.9	-3	-2.9	↑	↑, ←	↑, ↓, ←, →	↓
-2.9	-3	-2.9	-2.3	↑	↑, ↓, ←, →	↓, →	↓
-3	-2.9	-2.3	0	↑, →	→	→	

Step converged:

0	-14	-20	-22	0	←	←	↓, ←
-14	-18	-20	-20	↑	↑, ←	↑, ↓, ←, →	↓
-20	-20	-18	-14	↑	↑, ↓, ←, →	↓, →	↓
-22	-20	-14	0	↑, →	→	→	

Observe how in the second iteration we already found the optimal policy!

In this example, we showed how we can use evaluation to improve our policy, and in fact we obtained the optimal policy. However, the greedy approach does not always ensure reaching the optimal policy.



This approach is called **policy iteration**:

- Policy evaluation: estimate  $v^\pi$ .
- Policy improvement: generate  $\pi' \geq \pi$ .

It is natural to ask if policy evaluation need to converge to  $v^\pi$  or if we should stop the evaluation at some point. Ways to stop it are to put a threshold of minimum change between iterations, or simply after  $k$  iterations. One extreme, which is in fact quite usual in practice, is to stop after  $k = 1$ , which is equivalent to **value iteration**:

```

1 def value_iter(S, R, pi, gamma):
2
3     v = [0 for s in S]
4
5     repeat until converge
6         for s in S:
7             v_new(s) = max_a E[R + gamma * v(S_t+1) | S_t = s, A_t = a]
8
9         v = v_new
10
11     return v

```

In the following table, we sum up the different approaches:

Problem	Bellman Equation	Algorithm
Prediction	Bellman Expectation Eq	Iterative Policy Eval
Control	Bellman Expectation Eq + (Greedy) Policy Improvement	Policy Iteration
Control	Bellman Optimality Eq	Value Iteration

Observations

- The algorithms are based on state-value functions  $v_\pi(s)$  or  $v^*(s)$ , with complexity  $O(|A||S|^2)$  per iteration.
- It could also be applied to action-value functions  $q_\pi(s, a)$  or  $q^*(s, a)$ , with complexity  $O(|A|^2|S|^2)$  per iteration.

## 6.3 Extensions to Dynamic Programming

### 6.3.1 Asynchronous Dynamic Programming

DP methods described so far used **synchronous** updates, meaning all states are updated in parallel. In contrast, **asynchronous DP** backs up states individually, in any order. This can significantly reduce computation, and it is guaranteed to converge if all states continue to be selected.

We are going to see three approaches for ADP:

**In-Place DP** Before, with synchronous value iteration, we stored two copies of the value function:

$$v_{new}(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s, A_t = a],$$

$$v = v_{new}.$$

Now, in-place value iteration stores only one copy of the value function, doing:

$$v(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s, A_t = a].$$

**Prioritised Sweeping** We can use the magnitude of the Bellman error to guide the state selection. For example:

$$E = \left| \max_a \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s, A_t = a] - v(s) \right|.$$

We then backup the state with the largest remaining Bellman error, and update the Bellman error of affected states after each backup. This requires knowledge of reverse dynamics (which are the predecessor states). It can be implemented efficiently with a priority queue.

**Real-Time DP** The idea in this case is to only update states that are relevant to the agent. For example, if the agent is in state  $S_t$ , we update that state value, or the states that it expects to be in soon.

**Full-Width backups** Standard DP uses full-width backups. This means that, for each backup, it being synchronous or asynchronous:

- Every successor state and action is considered.
- Using the true model of transitions and reward function.

DP is effective for medium-sized problems (with millions of states). For large problems DP suffers from the curse of dimensionality, since the number of states grows exponentially with the number of state variables, and even one full backup can be too expensive.

**Sample Backups** This approach consists in using sample rewards and sample transitions  $(s, a, r, s')$  instead of the reward function  $R$ , and the transition dynamics,  $P$ .

It presents some advantages:

- It's model free: knowledge about the MDP is not required.
- Breaks the curse of dimensionality through sampling.
- Cost of backup is constant, independent of  $n = |S|$ .

## References

- [1] Tom Dupuis. Machine learning. Lecture Notes.