

AUTONOMOUS UNIVERSITY OF BARCELONA

MASTER THESIS

Solving a Complex Optimization Problem for Product Delivery with Reinforcement Learning and Artificial Neural Networks

Autor:

Daniel Salgado Rojo

Director:

Toni Lozano Bagén¹

Master's degree in Modelling for Science and Engineering
Mathematics Department¹

Wednesday 30th May, 2018



“Artificial Intelligence is the new electricity.”

Andrew Ng

Data Scientist and co-founder of Coursera

Abstract

To do.

[1]

- Mention tensorflow library somewhere [2]
- Open AI gym environment (explicar que hemos hecho uno para los algoritmos de reinforcement learning que simulamos, todos menos los del Q-learning que hice un código a parte) [3]

Acknowledgements

To do.

Nomenclature

The units we specify here are the ones that will be used mostly for simulations.

Symbol	Units	Description
t	day	Time,
s	—	System's state
\mathcal{S}	—	Set of all possible states
a	—	Action
\mathcal{A}	—	Set of all actions
$\mathcal{A}(s)$	—	Set of possible (not forbidden) actions from state s .
r	—	Reward after some action (scalar)
R	—	Function of rewards (function)
T	—	Transition probability function
π	—	A policy
Q^π	—	Q-value (or state-action value) function with regards to some policy π
V^π	—	V-value (or state value) function with regards to some policy π
τ	day	Length (number of days) of an episode
\mathcal{G}	—	Weighted and directed graph that defines the connections between shops and the depot
$W_{\mathcal{G}}$	—	Weights of the model graph edges \mathcal{G}
$A_{\mathcal{G}}$	—	Adjacency matrix of the model graph \mathcal{G}
\mathcal{N}	—	The set of shops (modelled as tanks) in the system
n	—	Number of shops in the system
C_i	—	Maximum stock capacity of the i -th shop
c_i	—	Current stock of the i -th shop
$x^{(i)} = \frac{c_i}{C_i}$	—	Fraction of available stock in the i -th shop

q_i	—	Expected total consumption of stock (after the end of the current day) for the i -th shop
\mathcal{K}	—	The set of trucks in the system
k	—	Number of trucks in the system
L_i	—	Maximum load capacity of the i -th truck
l_i	—	Current load of the i -th truck
p_i	—	Id (number) of the shop or depot where the i -th truck is located
$\lambda^{(i)}$	—	Quantity of product that the i -th truck is going to unload in a given shop
V	—	Number of shops that a given truck can at most visit every day
V_{\max}	—	Maximum number of shops (among all trucks) that a truck can visit every day
E	—	An episode
R_t	—	???
γ	—	Discount factor (for the discounted rewards)

Shortcut	Description
MDP	Markov Decision Process
RL	Reinforcement Learning
PG	Policy Gradient
ANN, NN	Artificial Neural Network
DNN	Deep Neural Network
GD	Gradient Descent
SGD	Stochastic Gradient Descent

Contents

Abstract	iii
Acknowledgements	v
Nomenclature	vii
1 Introduction	1
1.1 Goal optimization problem	2
1.2 Overview	2
2 Markov Decision Processes and Reinforcement Learning	3
2.1 Mathematical formalism of Markov Decision Processes	3
2.1.1 Markov Decision Processes	4
2.1.2 Policies	6
2.1.3 Optimality Criteria and Discounting	6
2.1.4 Value Functions	7
2.2 Reinforcement Learning	8
3 Problem approach	11
3.1 Problem description and constraints	11
3.2 Mathematical model	13
3.2.1 States: simple approach	14
3.2.2 States: complex approach	15
3.2.3 Actions	15
3.2.4 Rewards function	15
3.3 A simplified (toy) model	17
3.3.1 States, actions and rewards	18
4 Classical Reinforcement Learning	19

4.1	The Q-learning algorithm	19
4.2	Simulations	21
4.3	Results	23
4.3.1	Deterministic simulations	23
4.3.2	Stochastic simulations	27
5	Neural Networks and Deep Learning	31
5.1	Introduction to Artificial Neural Networks	31
5.2	Training of Deep Neural Networks	33
5.2.1	Gradient descent	34
5.2.2	The Backpropagation algorithm	36
5.3	Improving the way neural networks learn	42
5.3.1	Activation functions	43
5.3.2	Weight and biases initialization	44
5.3.3	Batch normalization	45
5.3.4	Regularization techniques	46
5.4	Learning a simple policy with a Deep Neural Network	47
5.4.1	Deep Neural Networks and Classification	47
5.4.2	Simulation framework for learning a hard-coded policy	48
5.4.3	Simulations and Results (TO DO)	49
6	Deep Reinforcement Learning	55
6.1	Policy Gradient algorithms	55
6.2	Simulations	55
6.3	Results	55
7	(? En caso que hiciésemos algo más como lo de separar en k softmaxes y aprendiese, o bien rutas, etc.)	57
8	Conclusions and Future Work	59
A	A Product Delivery gym environment	65
B	Derivation of $J_{\text{levels}}^{(i)}$	67
C	Default Q-learning simulation parameter values	69
C.1	Model system parameters	69

Chapter 1

Introduction

Introducción/ co

- Comentar que hay 3 campos en machine learning para distinguir en qual nos centramos y dar una idea de que va cada uno?

- Motivation artificial intelligence? AlphaGO,...

1.1 Goal optimization problem

Imagine we are a petrol company which owns some number of petrol stations and we are paying a transport company to bring our product to the stations with their trucks in order to refuel them (see Figure 1.1).

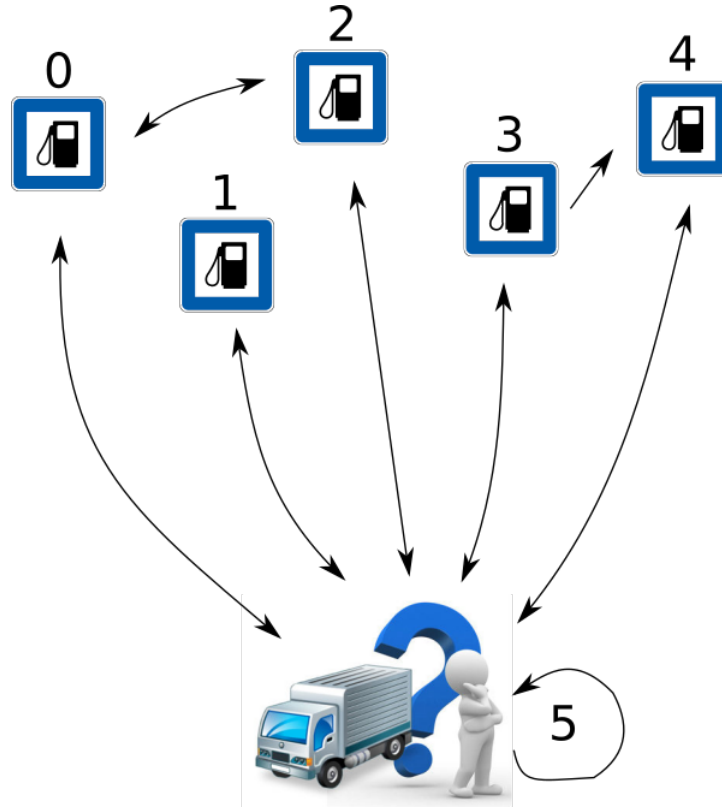


Figure 1.1: Product-delivery problem representation. (Truck image from: <https://www.pinterest.es/pin/521925044288845297/>)

The goal is to minimize the total amount of money (lets say per month or per year) that the petrol company has to pay to the transport company, ensuring that petrol is always available for costumers (among other constraints). Note that this problem can be generalized to the situation in which we have a company that is selling some product, and we have to pay to a transport company to distribute it from the loading dock so that it is available in all of our shops at any time.

solo hay que min-
imos de cosas
nó también qual-
ie no mueran las

1.2 Overview

Overview de las secciones para ubicar al lector.

Chapter 2

Markov Decision Processes and Reinforcement Learning

Markov Decision Processes (MDP) [4] are the fundamental mathematical formalism for *decision-theoretic planning* (DTP) [5], reinforcement learning (RL) [6, 7, 8] and other learning problems in stochastic domains [9].

In this chapter we introduce the basic concepts about Markov Decision Processes and Reinforcement Learning, focusing on the approaches that are adaptable to solve our goal optimization problem of product delivery. We mainly follow the first chapter in [9].

2.1 Mathematical formalism of Markov Decision Processes

In MDP models an *environment* is modelled as a set of *states* and *actions* that can be performed to control the system's state. The goal is to control the system by means of some *policy*, in such a way that some performance criteria is maximized (or minimized). MDP are commonly used to model problems such as stochastic planning problems, learning robot control and game playing.

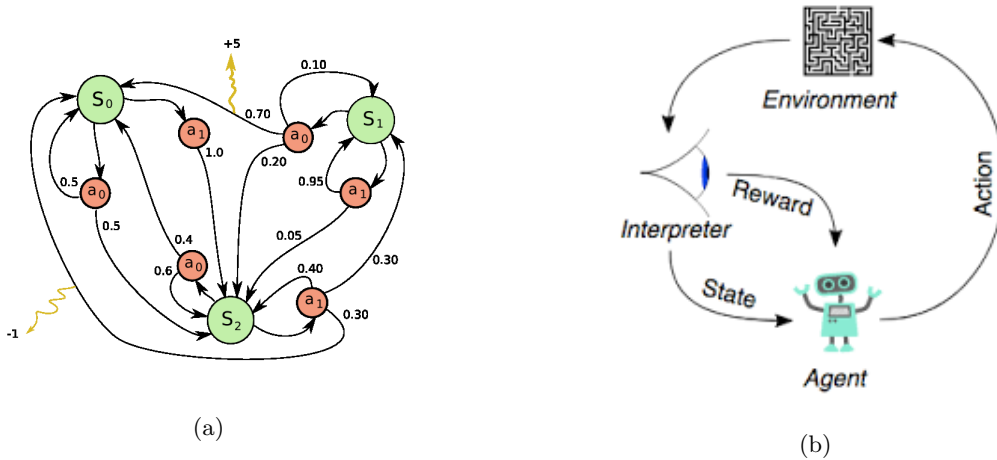


Figure 2.1: (a) “Example of a simple MDP with three states (green circles) and two actions (orange circles), with two rewards (orange arrows)” [10]. (b) “The typical framing of a Reinforcement Learning (RL) scenario: an agent takes actions in an environment, which is interpreted into a reward and a representation of the state, which are fed back into the agent” [11].

RL is a general class of algorithms whose goal is to make an *agent* learn how to behave in an environment, where the only feedback consists of a scalar *reward* signal [9]. The goal of the agent is to perform actions that maximize rewards in the long run.

In Figure 2.1a we see the directed graph representation of a simple MDP with three states (green nodes) and some actions (orange nodes). In a given state, only some actions can be performed, and each can lead to one or more states according to different probabilities. Transitions between states and the intermediate actions are represented by directed edges. In some state transitions, a positive or negative reward can be obtained (orange arrows) if we think of a MDP in the context of Reinforcement Learning (see Figure 2.1b).

2.1.1 Markov Decision Processes

The definition of MDPs we present at the end of this section consist of states, actions, transitions between states and a reward function. Although general MDPs may have infinite state and action spaces, we focus on discrete MDP problems with finite state and action spaces.

States

We denote $\mathcal{S} = \{s^1, \dots, s^N\}$ the set of possible states that can define the environment in a particular instant, and $|\mathcal{S}| = N$. Each $s \in \mathcal{S}$ is going to be considered as a tuple of *features* or properties that uniquely determine the environment state in a certain situation. For instance, in the *tic-tac-toe* game (see Figure 2.2 on the right), a state could be a tuple of 9 components corresponding to each cell of the game's board and it may contain either 0, 1 or 2 (or any other three symbols) depending on if the cell is empty or occupied by one of the players.

Actions

We denote $\mathcal{A} = \{a^1, \dots, a^\Lambda\}$ the set of actions that can be applied to control the environment by changing its current state; $|\mathcal{A}| = \Lambda$. Usually not all actions are going to be applicable when the system is in a particular state $s \in \mathcal{S}$. Thus, we define $\mathcal{A}(s) \subseteq \mathcal{A}$ the set of actions that are applicable in state s .

The transition operator

If an action $a \in \mathcal{A}$ is applied in a state $s \in \mathcal{S}$, the system makes a transition from s to a new state $s' \in \mathcal{S}$ based on a probability distribution over the set of possible transitions [9]. We define the transition function T as

$$T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \longrightarrow [0, 1]$$

$$(s, a, s') \longmapsto P(s'|s, a)$$

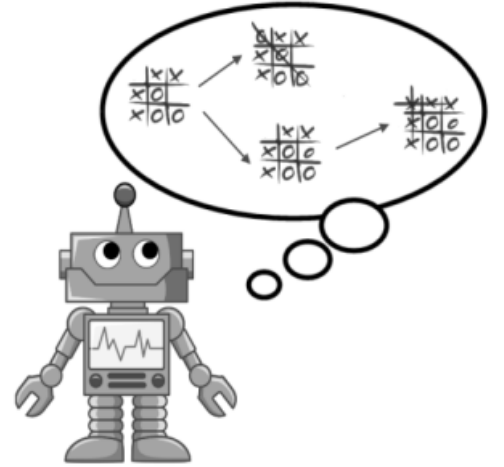


Figure 2.2: <https://mostafa-samir.github.io/Tic-Tac-Toe-AI/>

Hence, $T(s, a, s')$ is the probability of a system in a state s to make a transition to a new state s' after being applied action a .

There are some conditions to be satisfied so that T defines a proper probability distribution over possible next states:

- i) For all $s, s' \in \mathcal{S}, a \in \mathcal{A}, 0 \leq T(s, a, s') \leq 1$,
- ii) To model the fact that some actions are not applicable when being in some states, one sets $T(s, a, s') = 0$ for all triples (s, a, s') with $s', s \in \mathcal{S}$ so that $a \notin \mathcal{A}(s)$.
- iii) For all $s \in \mathcal{S}, a \in \mathcal{A}, \sum_{s' \in \mathcal{S}} T(s, a, s') = 1$.

For talking about the *order* in which actions occur one defines a discrete *global clock*, $t \in \{0, 1, 2, \dots\}$, so that s_t, a_t denote the state and action at time t , respectively.

Definition 2.1. *The system being controlled is Markovian if the result of an action does not depend on the previous actions and visited states, but only depends on the current state, i.e. [9]*

$$P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots) = P(s_{t+1}|s_t, a_t) = T(s_t, a_t, s_{t+1}). \quad (2.1)$$

In our context one assumes that the system being controlled is Markovian.

The rewards function

In classical optimization problems one seeks a *cost function* and aims to either maximize or minimize it under some set of conditions. In the context of RL, one talks about *rewards*, and the goal is to make an agent learn how to maximize the rewards obtained. For the type of problems we are interested to solve, we define the reward function R as

$$\begin{aligned} R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} &\longrightarrow \mathbb{R} \\ (s, a, s') &\longmapsto R(s, a, s') \end{aligned}$$

Thus, R is a scalar feedback signal which can be interpreted as a *punishment*, if negative, or a *reward*, if positive.

For the *tic-tac-toe* example, one could assign either high positive or negative rewards, or zero rewards, to actions that (when being in a given state) lead to win or lose the game, or having a draw, respectively. In this situation, the goal of the agent is to reach positive valued states, which means winning the game. Sometimes the problem is more complex and it is common to assign non-zero reward to combinations of states and actions that are good or bad under some criteria. For instance, having two aligned pieces in *tic-tac-toe* is good, so we may assign a positive reward for these situations.

In conclusion, rewards are used to give a *direction* in which way the MDP system should be controlled [9].

elaborar un poco l

The Markov Decision Process

With all this stuff, we can now give the definition of a *Markov decision process*.

Definition 2.2. A *Markov decision process* (MDP) is a tuple $(\mathcal{S}, \mathcal{A}, T, R)$ in which \mathcal{S} is a finite set of states, \mathcal{A} a finite set of actions, T a transition probability function $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ and R a reward function, $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$. One says that the pair T, R define the model of the MDP.

This definition of MDP allows us to model *episodic tasks* and *continuing tasks*:

- In episodic tasks, there is the notion of *episodes* of some length τ where the goal is to take the agent from some starting state to a *goal* state. In these cases, one can distinguish between *fixed horizon tasks* in which each episode consists of a fixed number of steps, or *indefinite horizon tasks* in which each episode has an end but episodes can have arbitrary length [9]. An example for the former type would be *tic-tac-toe* and for the latter, *chess* board game. Another example for indefinite horizon tasks would be a robot trying to learn how to run or walk; at the beginning it will fall down early (so that episodes will be short), but if it is learning properly, at some point it will be able to run or walk forever.

In each episode the initial state of the system is initialized with some *initial state distribution* $I : \mathcal{S} \rightarrow [0, 1]$.

- In continuing or *infinite horizon* tasks the system does not end unless done it in purpose (theoretically it never ends). A possible example would be a software dedicated to sustain the energy supply for a city by controlling the power plants.

2.1.2 Policies

Given an MDP $(\mathcal{S}, \mathcal{A}, T, R)$, for us a (deterministic) policy is a function $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maps states to actions that are applicable in that state, i.e., $\pi(s) = a \in \mathcal{A}(s)$ for all $s \in \mathcal{S}$ ¹.

A policy π can be used to make evolve, i.e. to control, a MDP system in the following way:

- Starting from a initial state $s_0 \in \mathcal{S}$, the next action the agent will do is taken as $a_0 = \pi(s_0)$.
- After the action is performed by the agent, according to the transition probability function T and the reward function R , a transition is made from s_0 to some state s_1 , with probability $T(s_0, a, s_1)$ and a obtained reward $r_0 = R(s_0, a_0, s_1)$.
- By iterating this process, one obtains a sequence $s_0, a_0, r_0, s_1, a_1, r_1, \dots$ of state-action-reward triples over $\mathcal{S} \times \mathcal{A} \times \mathbb{R}$ which constitute a trajectory (or path) of the MDP.

If the task is episodic, the sequence of state-action-reward triples ends in a finite number of iterations τ , the system is restarted and the process starts again with a new sampled initial state. If the task is continuing, the sequence can be extended indefinitely ($\tau = \infty$) [9].

2.1.3 Optimality Criteria and Discounting

The core problem of MDPs is to find an optimal policy π^* to control the system “optimally”. Therefore, we need a model for optimality with regard to policies.

¹In a general framework, policies are considered to be stochastic functions $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$, where $\mathcal{P}(\mathcal{A})$ is the set of probability measures on \mathcal{A} .

Although there are several models of optimality for a MDP, we focus here on the so called *discounted average reward* criteria, which is applicable in both finite (episodic) and infinite (continuing) horizon tasks.

The discounted sum of rewards received or *return* by the agent starting from state s_t is defined as

$$R_t^\gamma = \sum_{k=0}^{\tau} \gamma^k r_{t+k} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots + \gamma^\tau r_{t+\tau} \quad (2.2)$$

where $\tau < \infty$ if the task is episodic, and $\tau = \infty$ if it is continuing, and $\gamma \in [0, 1]$ is a discount factor for future rewards ($\gamma < 1$ if the task is continuing).

En episodic tasks t
sentido que la sum
 $\tau - t$!!!!!!

The factor γ determines the importance of future rewards:

- For $\gamma = 0$ (and taking $\gamma^0 = 1$) the only term that survives in equation (2.2) is the one corresponding to the present reward. In this case the agent will be *myopic* (short-sighted) by only considering current rewards.
- On the contrary, for $\gamma \approx 1$ the discount factor will make the agent strive for a long-term high reward [12]. For instance, if $\gamma = 0.95$, rewards 13 steps far into the future count approximately half as much as immediate rewards ($0.95^{13} \approx 0.51$), while for $\gamma = 0.99$, rewards start counting half as much immediate ones from the 69th step ($0.99^{69} \approx 0.499$) [13].
- For episodic tasks and the case $\gamma = 1$, R_t^γ is the sum of rewards at each step of an episode.

The goal of the discounted average reward criteria in the context of MDP is to find a policy π^* that maximizes the expected return $\mathbb{E}_\pi[R_t^\gamma]$.

Poner la suma par
tienda como se hac
anza?

Esto es lo que quiero comentar (estas expresiones son las que creo que saldrian cuando la policy es determinista)

$$\begin{aligned} \mathbb{E}_\pi[R_t^\gamma | s_t = s] &= \sum_{k=0}^{\tau} \gamma^k E_\pi[r_{t+k}(s_{t+k}, \pi(s_{t+k}), s_{t+k+1})] \\ &= \sum_{k=0}^{\tau} \gamma^k \sum_{s' \in \mathcal{S}} T(s_{t+k}, \pi(s_{t+k}), s') R(s_{t+k}, \pi(s_{t+k}), s') \end{aligned}$$

$$\begin{aligned} \mathbb{E}_\pi[R_t^\gamma] &= \sum_{s \in \mathcal{S}} I(s) E_\pi[R_t | s_t = s] \\ &= \sum_{s \in \mathcal{S}} I(s) \left\{ \sum_{k=0}^{\tau} \gamma^k \sum_{s' \in \mathcal{S}} T(s_{t+k}, \pi(s_{t+k}), s') R(s_{t+k}, \pi(s_{t+k}), s') \right\}_{s_t=s} \end{aligned}$$

This idea is formalized with the definition of the so-called value functions.

2.1.4 Value Functions

In order to link the criteria of optimality introduced in the precedent subsection to the policies, one considers the so-called *value functions*, which are a way to quantify “how good” it is for the agent

being in a certain state (*state-value function* V^π) or making a certain action when being in some particular state (*action-value function* Q^π).

Definition 2.3. The *value of a state s under policy π* , denoted $V^\pi(s)$ is the expected return when starting in state s and following π thereafter [9]. If we consider the discounted average reward criteria, we have the following expression:

$$V^\pi(s) = \mathbb{E}_\pi [R_t | s_t = s] = \mathbb{E}_\pi \left[\sum_{k=0}^{\tau} \gamma^k r_{t+k} | s_t = s \right] \quad (2.3)$$

where $\tau < \infty$ if the task is episodic, and $\tau = \infty$ if it is continuing.

Definition 2.4. The *state-action value of s, a under policy π* , denoted $Q^\pi(s, a)$ is the expected return when starting in state s , taking action a and thereafter following π [9]. For the discounted average reward criteria we have the following expression:

$$Q^\pi(s, a) = \mathbb{E}_\pi [R_t | s_t = s, a_t = a] = \mathbb{E}_\pi \left[\sum_{k=0}^{\tau} \gamma^k r_{t+k} | s_t = s, a_t = a \right] \quad (2.4)$$

where $\tau < \infty$ if the task is episodic, and $\tau = \infty$ if it is continuing.

The value functions are known to satisfy the so-called *Bellman equations* [14] that we do not introduce in this work (see for example [9, 14]).

The main goal of a given MDP is to find the policy that receives the most reward. This is equivalent to find the policy that maximizes the value function for each possible state.

Definition 2.5. A policy π^* is said to be *optimal* if it is such that $V^{\pi^*}(s) \geq V^\pi(s)$ for all $s \in \mathcal{S}$ and all policies π . $V^* := V^{\pi^*}$ is called the *optimal value function*. One may write

$$V^*(s) = \max_{\pi} V^\pi(s). \quad (2.5)$$

Similarly one defines the *optimal Q-value*, Q^* .

To illustrate how value functions are used in practice; assume we know the optimal Q-values for each state-action pair (i.e., $Q^*(s, a)$ for all $s \in \mathcal{S}, a \in \mathcal{A}$), or an algorithm able to estimate them. Then, one can greedily select an optimal action using the greedy Q-policy π_Q defined as

$$\pi_Q(s) = \arg \max_{a \in \mathcal{A}(s)} Q^*(s, a), \quad \forall s \in \mathcal{S}, \quad (2.6)$$

and it is an optimal policy according to definition 2.5.

2.2 Reinforcement Learning

Once we have defined MDPs, policies, optimality criteria and value functions, the next step is to consider the question of how to compute optimal policies.

Model-free and model-based algorithms

Although there are several approaches, there are usually two distinguished types of algorithms with regards to the available information about the MDP model: *model-based* and *model-free*.

Model-based algorithms assume that a model of the MDP is given, i.e. that a transition and a reward function pair (T, R) is known. This class of algorithms is also known as Dynamic Programming (DP).

On the contrary, model-free algorithms, which are under the general name of Reinforcement Learning, do not rely on the availability of a perfect model. “Instead, they rely on *interaction* with the environment, i.e. a *simulation* of the policy thereby generating *samples* of state transitions and rewards” [9]. Then, for instance, these samples can be used to estimate the Q-values (or the V-values) for each visited state-action pair (s, a) , $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$, by means of some iterative algorithm. Since a model of the MDP is not known a priori, the agent has to learn how to behave by experience. Thus, the agent has to *explore* the MDP to learn about its environment and how to behave in order to maximize the rewards obtained. “This naturally induces a *exploration-exploitation* trade-off which has to be balanced to obtain an optimal policy” [9].

As we are going to see during the following chapters, for the product delivery problem we introduced in section 1.1, we will have a model for the rewards function but we do not know anything about the transition function of our system modelled as an MDP. For this reason, from now on we assume to be in the model-free, RL, framework where the MDP model (either T , R or both) is initially unknown.

Value and Policy estimation algorithms

So far, for the sake of simplicity we have assumed MDPs with finite state and action spaces. There are some problems where the state or/and the action spaces are continuous so that one can consider that $\mathcal{S} \subseteq \mathbb{R}^d$ and $\mathcal{A} \subseteq \mathbb{R}^m$, for some integers $d, m \geq 1$. In these cases the concepts introduced in section 2.1 are practically the same but with another formalism (see for example Chapter 7 in [9]).

Generally one can distinguish between two RL methods for finding optimal policies: value-approximation and policy-approximation algorithms.

On the one hand, in value-approximation algorithms, experience samples are used to update a value function (such as Q or V) that gives an approximation of the current or the optimal policy, and from these estimations one can define an (approximated) optimal policy. For instance, usually one considers the greedy policy defined in (2.6). The difference between having finite or continuous state-action spaces is that in the first case, value functions are tabular functions², and in the continuous case³ one would consider some parametrized functions Q_θ, V_φ with respect to some set of parameters θ, φ , and from those functions define a policy.

On the other hand, policy-approximation algorithms focus directly on finding optimal policies by improving them as the system interacts with the environment. This approach usually corresponds to the context of continuous state-action spaces where it is considered a parametrized function π_θ , and the goal is to find the set of parameters θ^* such that $\pi^* := \pi_{\theta^*}$ is an optimal policy.

When considering parametrized functions one assumes enough differentiability with respect to the

²For instance, the Q -value function would be a table with one row per state and one column per possible action. And the V function would just be an array with as many components as states.

³Also in the finite case but with very large state and/or action spaces.

parameters so that one can perform the desired optimization procedures on them.

In this work we are going to focus mainly on two RL approaches:

- In chapter 4, we consider a product delivery framework where states and actions are finite so that we can consider tabular-based value-approximation algorithms such as the so-called Q-learning, the one we will focus on.
- In chapter 6, we consider a more general framework where the state space is continuous and we resort to policy-approximation algorithms where neural networks are used as parametrized policy functions π_θ , where the parameters θ are the weights of the Neural Network (see chapter 5 for more details on Neural Networks). The algorithm considered there is the so called Policy-Gradient, which in general words consists in updating parameters θ taking the direction of the gradient of π_θ scaled in proportion to the discounted average reward obtained (in this way, one achieves to make better states more probable, and worse states less probable).

Chapter 3

Problem approach

Poner un párrafo de introducción.

3.1 Problem description and constraints

A single product is to be delivered to several shops from a depot using several trucks (see Figure ??). We call *unload* the fact that a truck delivers (unloads) some quantity of product to a shop; unloads can be either *simple*, if a truck only delivers to one shop and then returns to the depot, or *shared* if a truck delivers to more than one shop before returning to the depot. Trucks can only be loaded in the depot.

Moreover, some initial restrictions and assumptions are considered:

1. Trucks leave the depot fully loaded (lets say at 8 in the morning) and return completely empty (lets say at the end of the day, i.e., 23:59:59).
2. Trucks perform at most one *unload* (simple or shared) every day. Hence, it is possible that some day some of the trucks do not go to any shop.
3. A truck can not visit the same shop to perform a *unload* more than one time during the same day.
4. Trucks can perform a shared unloads of at most $V \geq 1$ shops. In general V may depend on the truck.
5. The price that the company owning the shops has to pay to the transport company is given by some cost function J_{costs} which may depend on the distance travelled by the

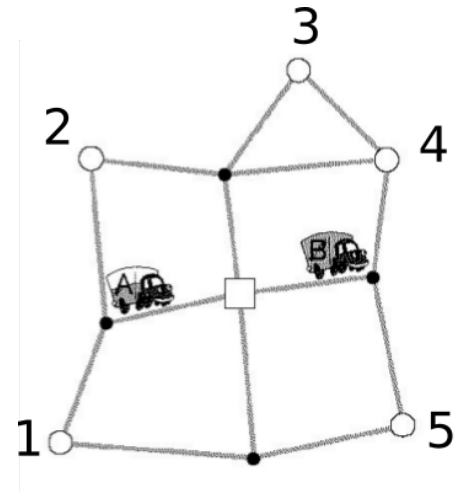


Figure 3.1: Product-delivery domain: the square in the centre is the depot (load doc) and white circles are the shops.. Adapted from [15].

Poner que basamos el planteamiento de bastante des de una vista de un problema

En algun sitio hay que poner que estamos intentando hacer aprender al algoritmo que las plantas no mueven optimally en un raro caso (o 'X' dias vis

trucks from the depot to the corresponding shops, the quantity of product delivered and other contributions imposed by the transport company (for example, impose an additional cost if that day is a holiday, an additional cost depending on the number of shops visited for *shared* unloads, etc.).

Another assumption in our model is that the company that owns the shops has some criteria to decide if the quantity of stock available in a shop is good or bad. For instance, it is clearly bad having zero stock in some shop.

Lets consider a particular shop and imagine that it stores the product in a cilindric container we will refer to as *tank* (see Figure 3.2). There are considered several levels of stock:

1. **Minimum capacity or zero level:** when there is no product in stock at all.
2. **Minimum level:** below this level it is considered as having a break of stock (high risk). Thus, we consider that below this level is like having no product in stock.
3. **Danger level:** from this level to the minimum level it is considered that we are taking a moderate risk, so that the shop should receive a *unload* of product soon. The smaller the value of stock below this level, the higher the risk of having a break of stock.
4. **Maximum level:** the desired maximum stock. This level can be overpassed in some quantity but never surpassing the maximum capacity level. The higher the value of stock above this level, the higher the risk of having a break of stock.
5. **Maximum capacity level:** when it is physically impossible to put more product in the container (i.e., however much we compress the products inside the container, we cannot fit more product).

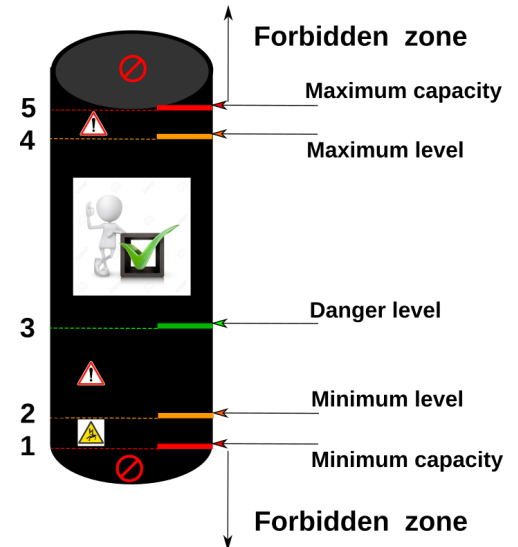


Figure 3.2: Stock levels for an abstract storage container (*tank*) of a given shop.

The ideal and desired state of stock is one between the danger level (3.) and the maximum level (4.).

For us, the minimum level and the danger level will correspond to the expected consumption of product in 12 and 36 hours, respectively.

3.2 Mathematical model

To formalize the model we use the Markov Decision Processes (MDP) nomenclature introduced in chapter 2.

Our system of shops and trucks can be modelled as a weighted and directed graph \mathcal{G} , a set of Trucks \mathcal{K} and a set of Tanks \mathcal{N} (the abstract storage containers of each shop, as in Figure 3.2). \mathcal{G} is determined by an adjacency matrix $A_{\mathcal{G}}$ and a matrix of weights $W_{\mathcal{G}}$, which are defined as follows:

$$(A_{\mathcal{G}})_{ij} = \begin{cases} 1 & \text{if there is an edge from shop } i \text{ to shop } j \\ 0 & \text{otherwise} \end{cases} \quad (3.1)$$

$$(W_{\mathcal{G}})_{ij} = \begin{cases} w_{ij} \in \mathbb{R}^+ & \text{if there is an edge from shop } i \text{ to shop } j \\ \infty & \text{otherwise} \end{cases} \quad (3.2)$$

Hence, weights are used to quantify the cost of going from a shop to another, and this cost is infinity when it is not possible to go from some shop to some other. Note that a component of $W_{\mathcal{G}}$ is ∞ if and only if the corresponding component of matrix $A_{\mathcal{G}}$ is zero.

As an example, the following matrix is the adjacency matrix of the system of 5 shops from Figure 1.1. The depot is considered to be the 6th node.

$$A_{\mathcal{G}} = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix}$$

We consider n tanks and k trucks. Thus, $|\mathcal{K}| = k$ and $|\mathcal{N}| = n$.

If V_{\max} is the maximum number of shops that any of the k trucks of the system can visit each day; then, we obtain a sequence of time labels that define a clock for our product delivery system

$$t \in \{t_{0,0}, \dots, t_{0,V_{\max}}, \dots, t_{\tau-1,0}, \dots, t_{\tau-1,V_{\max}}\},$$

such that the first index indicates the day and the second is aimed to divide a day in $V_{\max} + 1$ parts, and τ is some “last day” to consider. When the second index is zero, we assume that all trucks are fully loaded in the depot. We assume that it is only at the end of the day when the total amount of product in each tank is reduced by some amount according to what was consumed during all that day.

Each $N_i \in \mathcal{N}$ have the following features:

- An integer id, $i \in \{0, \dots, n-1\}$,
- A maximum load (maximum capacity) C_i ,
- The current load, $c_i = c_i(t)$,

- A discrete partition of the interval $[0, C_i]$ with d_i parts: $h_\beta = [\beta \cdot C_i/d_i, (\beta + 1) \cdot C_i/d_i]$, $\beta = 0, \dots, d_i - 1$.

Each $K_i \in \mathcal{K}$ has the following features:

- An integer id, $i \in \{0, \dots, k - 1\}$,
- A maximum load (maximum capacity) L_i ,
- The current load, $l_i = l_i(t)$,
- Its current position: characterized by the tank's id where the truck is located, we call it $p_i = p_i(t)$. If $i = n$, the truck is at the depot.
- A discrete partition of the interval $[0, L_i]$ with m_i parts: $f_\alpha = [\alpha \cdot L_i/m_i, (\alpha + 1) \cdot L_i/m_i]$, $\alpha = 0, \dots, m_i - 1$.

This leads to a discrete set of possible “deliveries” that a truck can do when it delivers product to some tank: $\{\lambda_0^{(i)}, \dots, \lambda_{m_i-1}^{(i)}\}$, where $\lambda_j^{(i)} = (j + 1) \cdot L_i/m_i$, $j = 0, \dots, m_i - 1$.

3.2.1 States: simple approach

Following the notation introduced above, we define the state of the system s_t in a given instant t by a tuple of three tuples consisting of the positions (p_0, \dots, p_{k-1}) of each truck, the current load of each truck (l_0, \dots, l_{n-1}) and the current load of each tank, (c_0, \dots, c_{n-1}) ¹. We assume that the consumption rate of each tank is deterministic and is the same every day.

In equation (3.6) we have the real (continuous) states of the system, and in equation (3.4) its discretized version.

Real state of the system

$$s = ((p_0, \dots, p_{k-1}), (l_0, \dots, l_{n-1}), (c_0, \dots, c_{n-1})) \in \mathcal{S} \subseteq \mathbb{N}^k \times \mathbb{R}^k \times \mathbb{R}^n \quad (3.3)$$

Discrete state of the system

$$s = ((p_0, \dots, p_{k-1}), (f_{\alpha_0}, \dots, f_{\alpha_{k-1}}), (h_{\beta_0}, \dots, h_{\beta_{n-1}})) \quad (3.4)$$

with $0 \leq \alpha_0, \dots, \alpha_{k-1} < m_i$, $0 \leq \beta_0, \dots, \beta_{n-1} < d_i$; and f_{α_j} , h_{β_j} are such that $l_j \in f_{\alpha_j}$ and $c_j \in h_{\beta_j}$, where j varies between 0 and $k - 1$ and between 0 and $n - 1$ respectively.

For the discrete case, the total number of possible states is:

$$|\mathcal{S}| = (n + 1)^k \times \prod_{i=1}^k m_i \times \prod_{i=1}^n d_i \quad (3.5)$$

¹For simplicity we omit the dependencies on t .

3.2.2 States: complex approach

Now we consider a more general case where the consumption rates of the tanks are added as state variables, so that they could model the fact that consumption rates are not the same every day and allow them to be the predicted consumptions coming from another model.

Thus, the real state of the system in this case would be like

$$s = ((p_0, \dots, p_{k-1}), (l_0, \dots, l_{k-1}), (c_0, \dots, c_{n-1}), (q_0, \dots, q_{n-1})), \quad (3.6)$$

where q_i is the consumption of product for the next time (day) and for tank i . Concretely, it is the consumption after transitioning from state s at the end of some day d (corresponding to time $t_{d, V_{\max}}$) to some other state s' (corresponding to time $t_{d+1, 0}$) through an action a .

Similarly, we could generalize even more this approach by considering not only the predictions of consumption for the next first day, but also the predictions for the coming $D > 1$ days. Therefore, the tuple vector (q_0, \dots, q_{n-1}) will become a set of D tuples $(q_0^j, \dots, q_{n-1}^j)$, $j = 1, \dots, D$, that we can think as a $D \times n$ matrix.

3.2.3 Actions

In a given state, an action consists in deciding the next location, either a shop or the depot, that each truck has to visit, and in the case of being a shop, how much product has to be delivered. Therefore, if p'_i denotes the new location where the i -th truck has to go and $\lambda^{(i)}$ the quantity of product to be unload.

$$a = ((p'_0, \dots, p'_{k-1}), (\lambda^{(0)}, \dots, \lambda^{(k-1)})) \in \mathcal{A} \subseteq \mathbb{N}^k \times \mathbb{R}^k \quad (3.7)$$

The dimension of the actions space \mathcal{A} when the truck loads are discretized satisfies:

$$|\mathcal{A}| \leq (n+1)^k \times \prod_{i=1}^k m_i \quad (3.8)$$

The inequality comes from the fact that not all actions need to be valid actions in a given state.

3.2.4 Rewards function

As we commented in section 2.1.1, the rewards function is a scalar score that allows the agent to know how good or bad it is to take an action in a given state.

In the product delivery problem we want to solve, we split the function of rewards in three main contributions (see eq. (3.13)):

$$R(s, a, s') = \mu_1 J_{\text{costs}} + \mu_2 J_{\text{levels}} + \sum_j \mu_{3,j} J_{\text{extra},j}, \quad \mu_i \geq 0, \quad (3.9)$$

- The first term, $J_{\text{costs}}(s, a, s')$, would be the one representing the economical costs to pay to the transport company. In particular, $J_{\text{costs}} \leq 0$.

- The second term, $J_{\text{levels}}(s')$, should be a reward function that penalizes negatively in dangerous and forbidden regions of stock for each tank, and positively in the optimal region of stock according to Figure 3.2 and the qualitative importance of the stock levels explained in section 3.1.
- Finally, a third term of extra penalties $J_{\text{extra},j}(s, a, s')$ that the agent may need to learn additional rules or prohibitions (such as forbidden actions).

Since the range of values taken by the different J contributions can be very different, the coefficients μ_i allow to balance the importance of each contribution term and deal with different scales.

Transport and unloads costs contribution

If we assume that costs are proportional to the distance travelled by trucks and to the total amount of product unloaded, we can consider a function of the form

$$J_{\text{costs}}(s, a, s') = C_{\text{costs}} \sum_{i=0}^{k-1} w_{p_i, p'_i} \cdot \lambda^{(i)}, \quad (3.10)$$

for some constant C_{costs} that makes J_{costs} be dimensionless.

Levels of stock contribution

If $x^{(i)}(s')$ denotes the fraction between the current load c_i of tank i with respect to its maximum capacity C_i , then we consider the following decomposition for the levels of stock contribution to the total reward function:

$$J_{\text{levels}}(s') = \sum_{i=0}^{n-1} J_{\text{levels}}^{(i)}(x^{(i)}(s')), \quad (3.11)$$

where

$$J_{\text{levels}}^{(i)}(x) = \begin{cases} P_2 & \text{if } x < a \text{ or } x > f \\ C_{ab} \exp\left(\frac{\alpha_{ab}}{x}\right) & \text{if } a \leq x \leq b \\ m_{bc}x + n_{bc} & \text{if } b < x \leq c \\ m_{cd}x + n_{cd} & \text{if } c < x \leq d \\ m_{de}x + n_{de} & \text{if } d < x \leq e \\ C_{ef} \exp\left(\frac{\alpha_{ef}}{1-x}\right) & \text{if } e < x \leq f, \end{cases}$$

and the coefficients m_{ij} , n_{ij} and C_{ij} are determined such that the functions $J_{\text{levels}}^{(i)}(x)$ look like the one in Figure 3.3. As it is detailed in Appendix B, for the suggested function there are 6 free tunable parameters: b, c, e and M, P_1, P_2 .

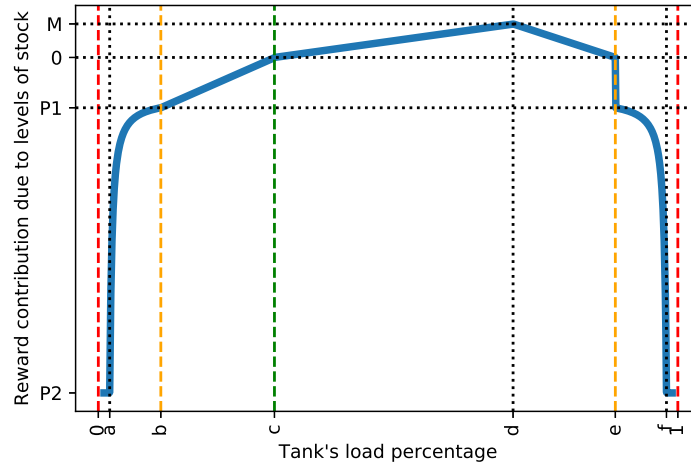


Figure 3.3: TO DO

Additional contributions

Some examples for additional contributions to the rewards function could be a $J_{\text{extra, costs}}$ term that could be included in the J_{costs} term, to account for extra costs due to different prices for holidays, for the fact that shared unloads are cheaper or more expensive than simple unloads,...

Moreover, since during the simulations the agent does not know which actions are valid or not in a given state so that it may choose non-performable actions, we need to add an extra (negative) penalty to make the agent know that this action is “bad” in some sense.

For instance, the agent may at some point try to fill a tank with a quantity of product that would make the shop surpass its maximum stock capacity. For this reason, we may add a contribution to the reward of the form:

$$J_{\text{extra, forbidden}}(s, a, s') = -C_{\text{extra, forbidden}} \sum_{i=0}^{k-1} \mathbb{1}(\text{the action performed by truck } i \text{ is forbidden}) \quad (3.12)$$

3.3 A simplified (toy) model

For the simulations and computations we will present in the next chapters, otherwise noted we are going to restrict ourselves to a simplified version of the model we have introduced so far in this chapter. We are going to focus on the first states approach described in section 3.2.1.

As it can be deduced from Figure 3.4, the first important assumption we make is that the only connections available are the ones between shops and the depot. Thus, each truck can at most visit one shop every day so that we are restricting to the case where $V_{\text{max}} = 1$ (the maximum number of visited shops by a truck).

In this frame, the adjacency matrix that model our system’s graph is the following:

$$A_{\mathcal{G}} = \begin{pmatrix} 0_{n \times n} & 1_n \\ 1_n^T & 1 \end{pmatrix}$$

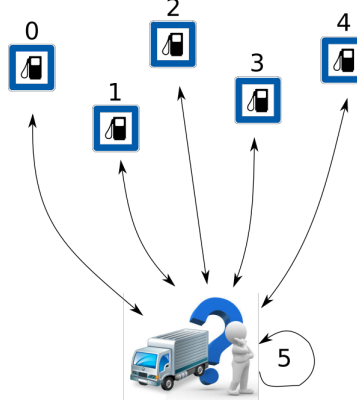


Figure 3.4: Simplified Product-delivery problem representation. Case $n = 5$ shops.

where $0_{n \times n}$ is a matrix with n rows and n columns of zeros, and 1_n is a column vector with n components equal to 1.

3.3.1 States, actions and rewards

Since at most only one shop is visited by each truck, there is no need to include the current positions of the trucks in the system's state. We know that at the beginning of the day, each truck will visit one shop or stay at the depot and at the end of the day it will return to the depot to be filled up to be ready for the next day.

The first assumption we considered in section 3.1 requires that if a truck visits a shop, the current load of the truck - which under the assumptions we are making corresponds to the whole capacity of the truck - must fit in that shop. For this reason, we should consider a term on the reward function such as the one from equation (3.12) to help the agent avoid these kind of forbidden actions. Moreover, now there is no need to include the current load of each truck in the system's state since it will be always the same; the same applies for the actions, where now only the new positions (shop or depot) for each truck are needed. In short,

$$s = (c_0, \dots, c_{n-1}) \in \mathcal{S} \subseteq \mathbb{R}^n$$

$$a = (p'_0, \dots, p'_{k-1}) \in \mathcal{A} \subseteq \mathbb{N}^k$$

Moreover, for now we forget about the transport and unload costs contributions and just focus on the contribution made by the levels of stock criteria to see if the reinforcement learning algorithms we will present are able to learn how to maintain the shops alive, or even better, maintain their stocks around the optimal regions (remember Figure 3.2).

In short, the function of rewards for the simplified model version that we will use to perform most of the simulations has the following form:

$$R(s, a, s') = \mu_2 J_{\text{levels}}(s') + \mu_{3, \text{forbidden}} J_{\text{extra, forbidden}}(s, a, s') \quad (3.13)$$

where the two terms are given in equations (3.11) and (3.12), respectively.

After this, we will (probably) consider $\mu_1 \neq 0$ to study the balance between the levels criteria and the economical costs due to transport and unloads

Chapter 4

Classical Reinforcement Learning

1 Párrafo de intro

4.1 The Q-learning algorithm

In practice, the value functions we introduced in section 2.1.4 are learned through the agent's interaction with the environment. If we focus on the Q-values, the popular algorithms that have been used more are Q-learning [16, 17] and SARSA¹ [18] from the so-called *temporal difference approach*², which are typically used when the model of the environment, i.e. the pair (T, R) , is unknown.

Assuming we can simulate the system in `n_episodes` episodes of length τ by means of some exploitation-exploration strategy for selecting actions in a given state, this leads to a set of simulated episodes E_j that we can write as

$$E_j = (s_0^j, a_0^j, r_0^j, s_1^j, a_1^j, r_1^j, \dots, s_{\tau-1}^j, a_{\tau-1}^j, r_{\tau-1}^j, s_\tau^j)$$

for $j \in \{1, \dots, \text{n_episodes}\}$.

The most exploitative action-selection criteria would be the greedy one, which would consist on using equation (2.6) for the current learned Q -values if the current state is known (and otherwise chose an action randomly). The most exploratory criteria would be the one which selects an action completely randomly in every step.

An interesting compromise between the two action-selection extremes is the ε -greedy policy, we denote π_ε , which is widely used in practice [21], and it is defined as follows:

$$\pi_\varepsilon(s) = \begin{cases} \text{random action from } \mathcal{A}(s) & \text{if } p < \varepsilon \\ \arg \max_{a \in \mathcal{A}(s)} Q(s, a) & \text{otherwise,} \end{cases} \quad (4.1)$$

where $p \in [0, 1]$ is a uniform random number drawn at each time step (of each episode).

Policy (4.1) executes the greedy policy (2.6) with probability $1 - \varepsilon$ and the random policy with probability ε . The use of this combination of policies gives a balance between exploration and exploitation that both guarantees convergence and often good performance [22].

¹Which stands for State-Action-Reward-State-Action)

²See for example [9, 13] for an academic point of view, or the original papers [19, 20].

In Algorithm 1 we present the pseudocode for Q-learning which is hopeful self-explanatory. There is only one parameter we have not talked about yet which is the learning rate α :

ONER / DE-
CO DE DONDE
DATE RULE DE

The learning rate or step size determines to what extent newly acquired information overrides old information. A factor of 0 makes the agent learn nothing, while a factor of 1 makes the agent consider only the most recent information. In fully deterministic environments, a learning rate of $\alpha = 1$ is optimal. When the problem is stochastic, the algorithm converges under some technical conditions on the learning rate that required it to decrease to zero (see Theorem 4.1) [12]. In general, the learning rate can be considered a function time step, the current state and the current chosen action.

Algorithm 1 Q-learning (train) algorithm

```

1: procedure Q-LEARNING(  $\gamma, \alpha(a, s), \tau, \mathbf{n\_episodes}$ )
2:    $Q \leftarrow 0$  ▷ Initialise  $Q(s, a)$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .
3:   for  $j \in \{1, \dots, \mathbf{n\_episodes}\}$  do
4:      $s_0^j \leftarrow \text{Random\_choice}(s \in \mathcal{S})$  ▷ The system is initialized to some initial state randomly
5:     for  $t \in \{0, \dots, \tau - 1\}$  do
6:       Choose  $a_t^j \in \mathcal{A}(s_t^j)$  based on the current  $Q$  and an exploration strategy (e.g.  $\pi_\epsilon$ ).
7:       Perform action  $a_t^j$ .
8:       Observe the new state  $s_{t+1}^j$  and the received reward  $r_t^j$ .
9:       Update  $Q$  with the following rule:


$$Q(s_t^j, a_t^j) \leftarrow Q(s_t^j, a_t^j) + \alpha_t^j(s_t^j, a_t^j) \left[ r_t^j + \gamma \max_{a \in \mathcal{A}(s_{t+1}^j)} Q(s_{t+1}^j, a) - Q(s_t^j, a_t^j) \right]$$


10:    end for
11:  end for
12: end procedure

```

For when the problem we are dealing with is not completely deterministic, the following theorem gives us a criteria to ensure that the Q-learning algorithm will converge as long as the learning rate satisfies some conditions.

Theorem 4.1. *Given a finite MDP $(\mathcal{S}, \mathcal{A}, T, R)$ such that \mathcal{S}, \mathcal{A} are finite, $\gamma \in (0, 1)$ and R is deterministic and bounded, the Q-learning algorithm given by the update rule*

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t(s_t, a_t) \left[r_t + \gamma \max_{a \in \mathcal{A}(s_{t+1})} Q(s_{t+1}, a) - Q(s_t, a_t) \right] \quad (4.2)$$

converges with probability 1 to the optimal Q-value function as long as

$$\sum_{t=0}^{\infty} \alpha_t(s, a) = +\infty, \quad \sum_{t=0}^{\infty} \alpha_t(s, a)^2 < +\infty, \quad 0 \leq \alpha_t(s, a) < 1 \quad \text{for all } (s, a) \in \mathcal{S} \times \mathcal{A} \quad (4.3)$$

Condition (4.3) requires that all state-action pairs are visited infinitely often.

Proof. The proof of this theorem can be found for example in [23]. □

Example 4.1. The conditions for the learning rates that appear in Theorem 4.1; for the case of an episodic task they can be simply rewritten as

$$\sum_{j=0}^{\infty} \sum_{t=0}^{\tau-1} \alpha_t^j(s, a) = +\infty, \quad \sum_{j=0}^{\infty} \sum_{t=0}^{\tau-1} \alpha_t^j(s, a)^2 < +\infty, \quad 0 \leq \alpha_t^j(s, a) < 1 \quad \text{for all } (s, a) \in \mathcal{S} \times \mathcal{A} \quad (4.4)$$

where the index k denotes the episode.

As an example, consider a learning rate that decreases with the number of episodes but does not depend on t (when fixing an episode), neither on the state nor the action. For example, take $\alpha_t^j = \frac{\alpha_0}{1+\beta j}$, where $\alpha_0, \beta \in (0, 1)$. We call α_0 the initial learning rate, and β the learning rate decay.

Note that α_t^j satisfies the conditions in (4.4) since:

$$\sum_{j=0}^{\infty} \sum_{t=0}^{\tau-1} \frac{\alpha_0}{1+\beta j} \sim \sum_{j=0}^{\infty} \frac{1}{1+j} \sim \sum_{n=1}^{\infty} \frac{1}{n} = +\infty,$$

and

$$\sum_{j=0}^{\infty} \sum_{t=0}^{\tau-1} \left(\frac{\alpha_0}{1+\beta j} \right)^2 \sim \sum_{j=0}^{\infty} \frac{1}{(1+j)^2} \sim \sum_{n=1}^{\infty} \frac{1}{n^2} < +\infty.$$

4.2 Simulations

We have implemented the Q-learning algorithm 1 in Python [1], which has been tested for the simplified case considered in section 3.3. To make the algorithm computationally feasible we should resort to a discretization of the states and actions as we explained in section 3.2, since the tabular function $Q(s, a)$ that we encode as a python's dictionary structure³ should not be arbitrarily large; otherwise, the Q-learning algorithm will not perform well in a reasonable amount of execution time.

At the beginning of each episode the state of the system, i.e. the initial stocks, are sampled via $c_i \sim \text{Uniform}(C_i/d_i, \frac{d_i-1}{d_i}C_i)$ for $i = 0, 1, \dots, n-1$.

As exploration strategy we use π_ε as defined in equation (4.1) with a parameter ε that starts from ε_0 and decreases until a value ε_{\min} after each episode E_j as follows:

$$\varepsilon(E_j) = \max \left(\varepsilon_{\min}, \frac{\varepsilon_0}{1 + \varepsilon_{\text{decay}}(j-1)} \right), j \geq 1, \quad (4.5)$$

where $\varepsilon_{\text{decay}}$ is a decay parameter that determines the rate at which ε decreases as the number of episodes increases.

The π_ε strategy taking the ε parameter according to equation (4.5) allows to be very exploratory at the beginning, and start being more and more exploitative as the number of episodes increases. Moreover, in order to be always a bit exploratory, we perform a maximum operation to cut ε to a minimum threshold value ε_{\min} .

In Table 4.1 we summarize the default parameter values that will be used in the simulations of this chapter.

³Whose *keys* are the state-action pairs encoded as strings with *values* the corresponding Q-values.

Table 4.1: Some constant parameters for the Q-learning simulations.

Parameter	Value	Parameter	Value
n	5	d_i	4
k	2	m_i	1
ε_0	1.0	$\varepsilon_{\text{decay}}$	The inverse of some fraction of <code>n_episodes</code> (e.g. $(0.05 \cdot \text{n_episodes})^{-1}$)
ε_{\min}	0.05	γ	0.9
α_0	1	β	0
τ	30	<code>n_episodes</code>	100.000

Note that we restrict to the case where we have 5 shops and 2 trucks; each shop divides its stock capacity in 4 discrete levels, and each truck does it in only 1 level (the truck’s capacity). During the action selection via the policy π_ε we start being purely exploratory ($\varepsilon_0 = 1$) and end being a bit exploratory 5% of the times ($\varepsilon_{\min} = 0.05$) (see Figure 4.1).

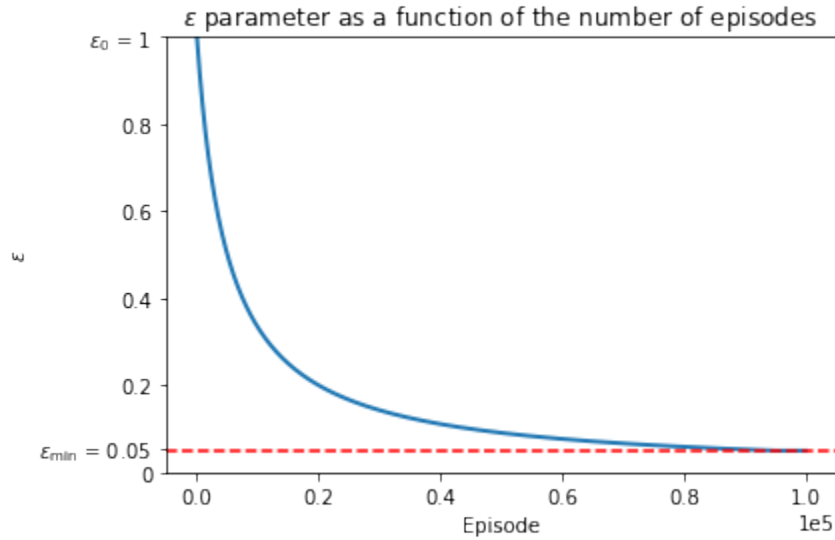


Figure 4.1: TODO

at we are using a
ng rate equal to 1
that our system

Although our implementation of Q-learning follows almost literally Algorithm1, there is a subtle detail that we must point out: just after applying the action returned by π_ε and before both “observing” the next state and computing the reward in that step, the current stocks of each shop is reduced in some amount equal to a “daily consumption” that may depend on the shop and also on the day if we consider stochasticity.

When we consider a deterministic system such that all shops have a constant daily consumption rate, it will be given by the expected consumption for the next 24 hours. For us it is approximated by $C_i \frac{b+c}{2}$ (i.e. the mean between the expected consumption in the next 12 and 36 hours), where i denotes the shop.

If we want to add some stochasticity to the consumption of stock in shops, we will consider con-

sumption rates sampled via

$$C_i \frac{b+c}{2} + \delta \cdot \text{Uniform}(-1, 1), \quad i \in \{0, \dots, n-1\}, \quad (4.6)$$

where typically $\delta \in [0, 1]$. For $\delta = 0$ we are in the deterministic case and for $\delta > 0$ we say that we have a noise of the $\delta \cdot 100$ %.

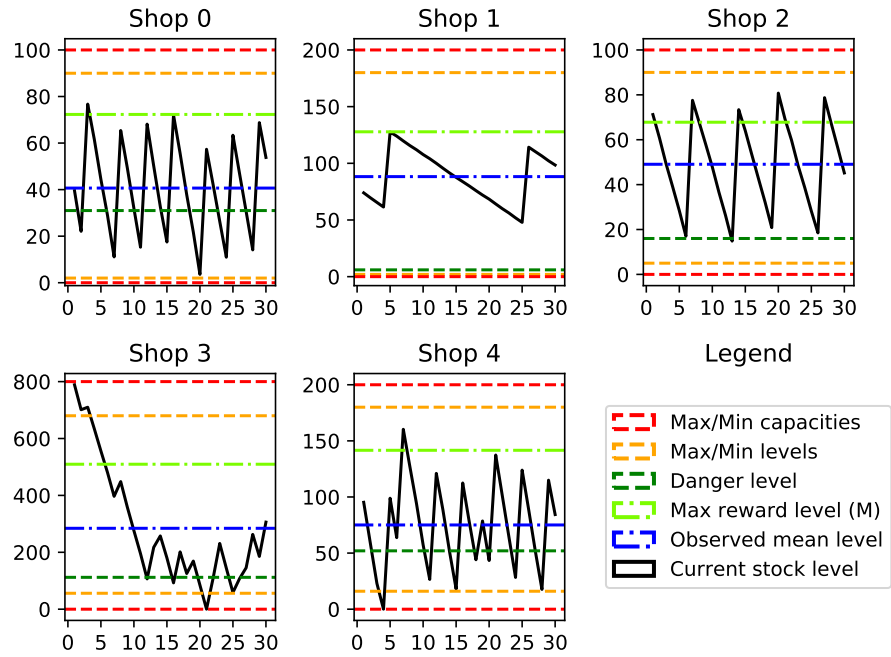
Table 4.2: Simulation parameters (that differ from the default ones in table ??)

Id	Simulation	μ_1	Stochastic consumption
1)	Deterministic without transport/unload costs	0	No
2.1)	Deterministic with transport/unload costs	10^{-6}	No
2.2)	Deterministic with transport/unload costs	10^{-4}	No
2.3)	Deterministic with transport/unload costs	10^{-3}	No
3.1)	Stochastic without transport/unload costs	0	Yes, 25%
3.2)	Stochastic without transport/unload costs	0	Yes, 10%
4.1)	Stochastic with transport/unload costs	10^{-6}	Yes, 25%
4.2)	Stochastic with transport/unload costs	10^{-6}	Yes, 10%

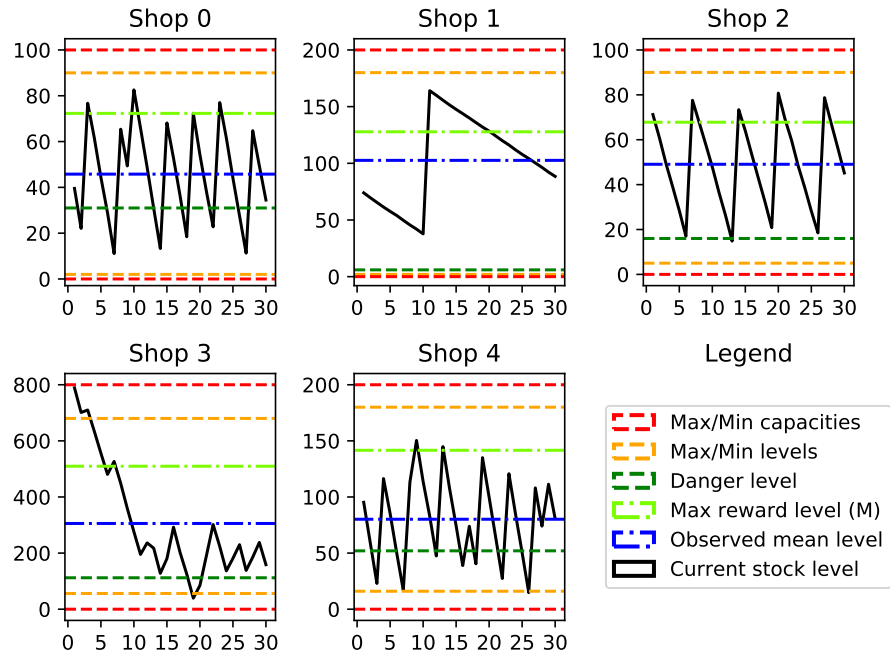
4.3 Results

4.3.1 Deterministic simulations

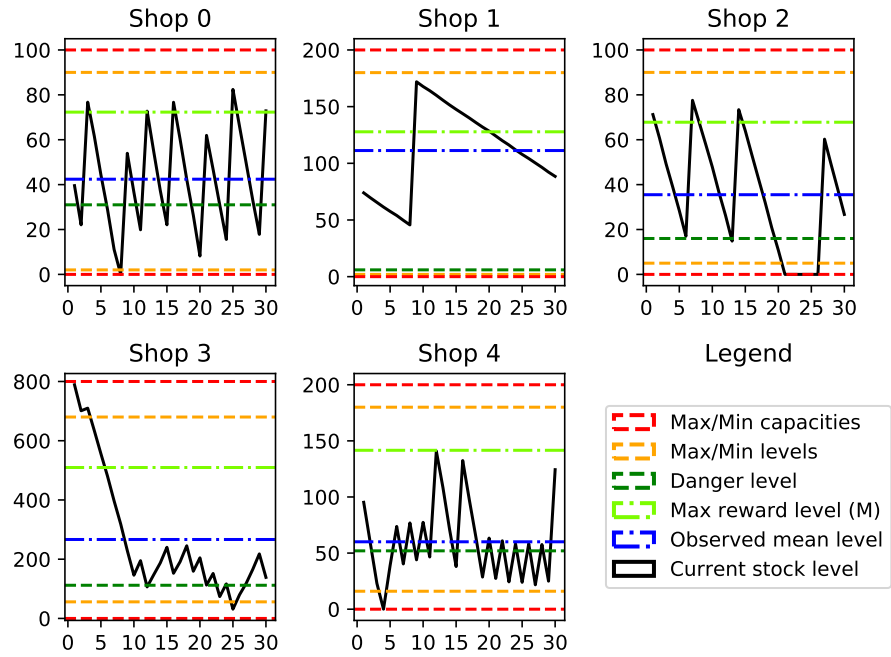
Simulation 1



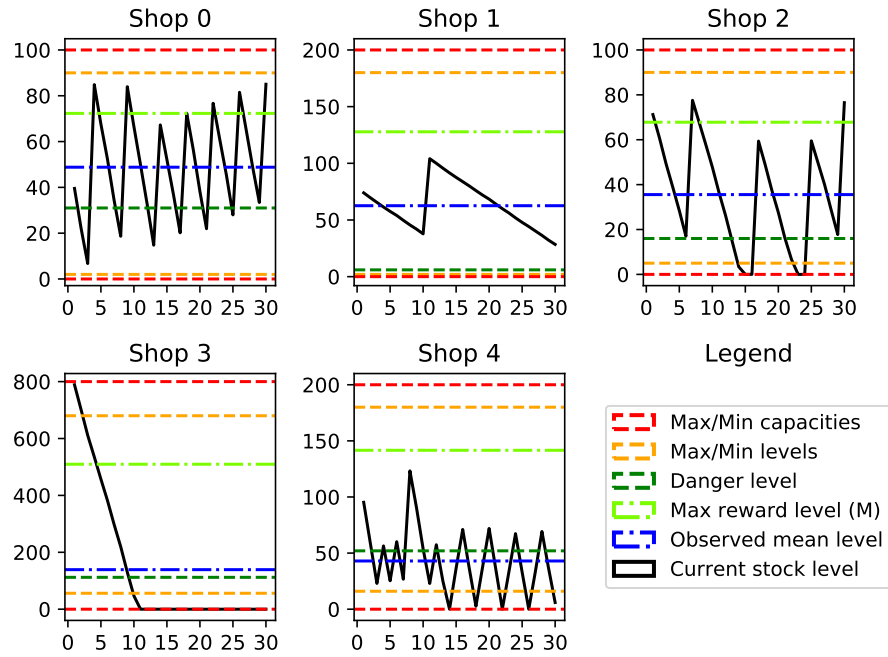
Simulation 2.1



Simulation 2.2

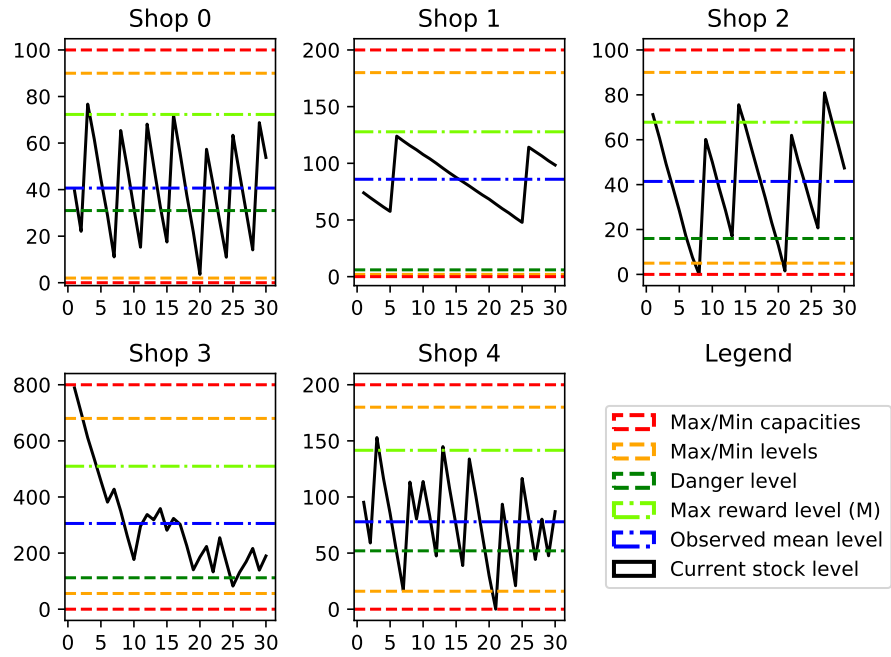


Simulation 2.3



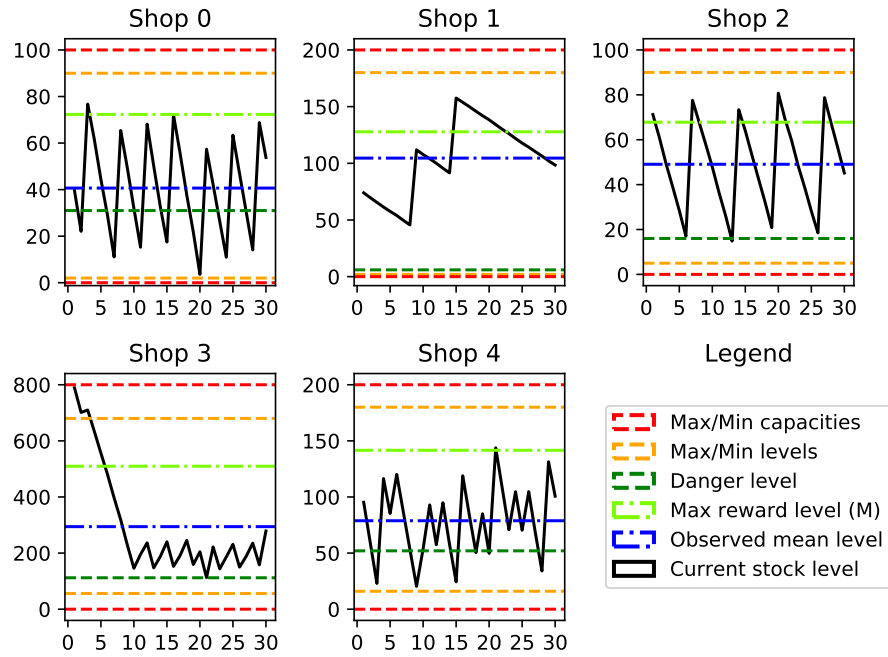
4.3.2 Stochastic simulations

Simulation 3.1



Simulation 3.2

Simulation 4.1



Simulation 4.2

Chapter 5

Neural Networks and Deep Learning

5.1 Introduction to Artificial Neural Networks

Artificial Neural Networks (ANN)¹ are computational structures inspired on the biological neural networks of the (human) brain with the aim of learning and adapting to a wide variety of tasks. Application areas for NN include system identification and control (vehicle control, trajectory prediction, process control), game-playing and decision making, pattern recognition (radar systems, face identification, object recognition), sequence recognition (gesture, speech, hand written and printed text recognition), medical diagnosis, finance (automated trading systems), data mining, visualization, social network filtering and e-mail spam filtering [24].

In 1943, Warren McCulloch and Walter Pitts [25] introduced the first mathematical and algorithmic models for NN that paved the way for NN research to split into two approaches: one focused on biological processes in the brain (Neuroscience), and another focused on the application of NN to Artificial Intelligence (e.g. machine learning). However, their real powerfulness had not been able to be exploited until the 1990s, thanks to the tremendous increase in computing power, the huge quantity of data available to train neural networks and the improvement of algorithms.

In this work we introduce NN with a functional point of view, which personally is more intuitive and understandable. In this way, we consider a NN to be a parametrized function $F_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ whose inputs x are observations coming from the *train data* $X \subseteq \mathcal{X}$ that one wants to use to make the NN learn to predict some output $y \in \mathcal{Y}$. The final goal is to tune the $\theta \subseteq \mathbb{R}^m$ parameters so that F_θ is able to perform the desired task (e.g. a regression task).

Concretely, a NN can be represented as a weighted, directed and acyclic graph (DAG), whose nodes are called *neurons*. These neurons are organized in ordered *layers* in such a way that the only edges that arrive to neurons in the l -th layer are those coming from the $l - 1$ -th layer. In particular, we focus on fully-connected *feed-forward* NN, where the output from one layer is the input of every neuron in the next layer and there are no loops (see Figure 5.1). At least it makes sense to consider NN with $L \geq 3$ layers. The first one (left-most) is called *input layer*, the L -th layer (right-most) is the *output layer* and all other layers in the middle are called *hidden layers*. A NN with more than one hidden layer is referred to as Deep Neural Network (DNN).

¹We will refer to them simply as Neural Networks (NN). The “Artificial” term is used to distinguish them from “biological” neural networks.

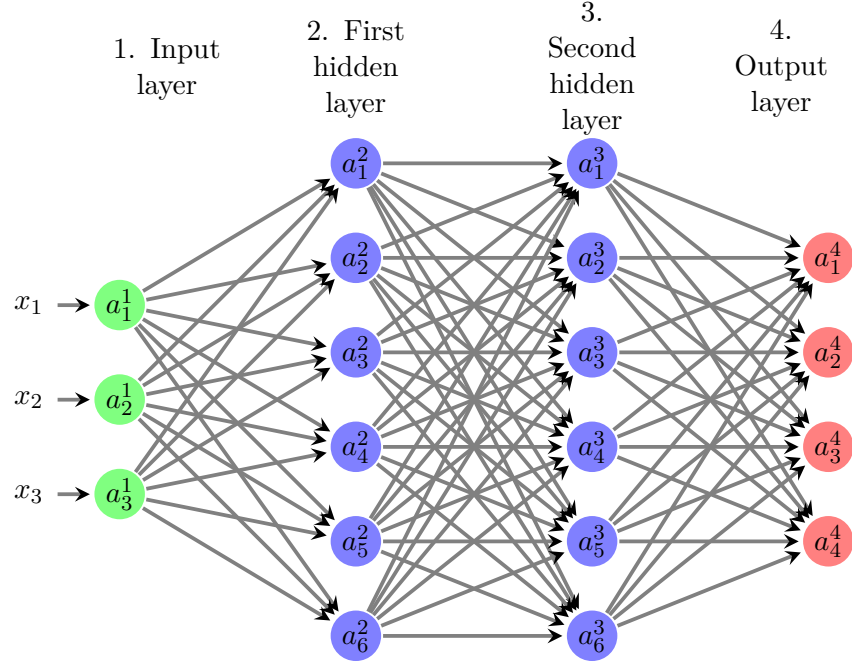


Figure 5.1: Representation of a feed-forward, fully connected Deep Neural Network with two hidden layers of six neurons each, an input layer of three neurons and an output layer with four neurons. With this idea one could generalize to a (feed-forward, fully connected) DNN with arbitrary number of hidden layers and neurons.

In order to define a NN function properly, we should focus on individual neurons first. For the sake of simplicity, let's consider that every layer has the same number of neurons n .

If we consider the j -th neuron in layer l , it receives some input values $a^{l-1} = (a_1^{l-1}, \dots, a_n^{l-1})^T$, with associated weights $w_j^l = (w_{j1}^l, \dots, w_{jn}^l)$ (w_{jk}^l being the weight of the edge that connects the k -th neuron from the $(l-1)$ -th layer to the j -th neuron from the l -th layer), and a constant value b_j^l coming from an auxiliary *bias* node (see Figure 5.2).

As illustrated in Figure 5.2, after the neuron receives the inputs, it computes an intermediate value we call *weighted input*, which is “the weighted sum of inputs plus the bias”:

$$z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l = w_j^l \cdot a^{l-1} + b_j^l, \quad l = 1, \dots, L, \quad (5.1)$$

and then this value is passed to a so-called *activation function* σ ² to produce a final output value $a_j^l = \sigma(z_j^l)$. For simplicity, we are assuming that each neuron uses the same function σ , but in general we could use a different activation function for each layer or even for each neuron.

In this way the j -th neuron in the l -th layer can be thought as a function f_j^l that takes a^{l-1} , w_{jk}^l (for each neuron k in the previous layer) and b_j^l as input, and outputs a value a_j^l , i.e. $f_j^l(a^{l-1}, w_{jk}^l, b_j^l) = \sigma(z_j^l) = a_j^l$. In this way, it is obtained an overall output value $a^l = (a_1^l, \dots, a_n^l)^T$ from the l -th layer, one component per neuron.

Similarly, for each layer we can define a vector of weighted inputs $z^l = (z_1^l, \dots, z_n^l)^T$, a vector of biases $b^l = (b_1^l, \dots, b_n^l)^T$, and a matrix of weights $w^l = (w_{jk}^l)_{j,k}$, so that with the convention of thinking the activation σ as acting to a vector component by component (vectorized form of σ), we can write

²Usually the symbol σ is used for the so-called *sigmoid* function. Here it is an arbitrary function.

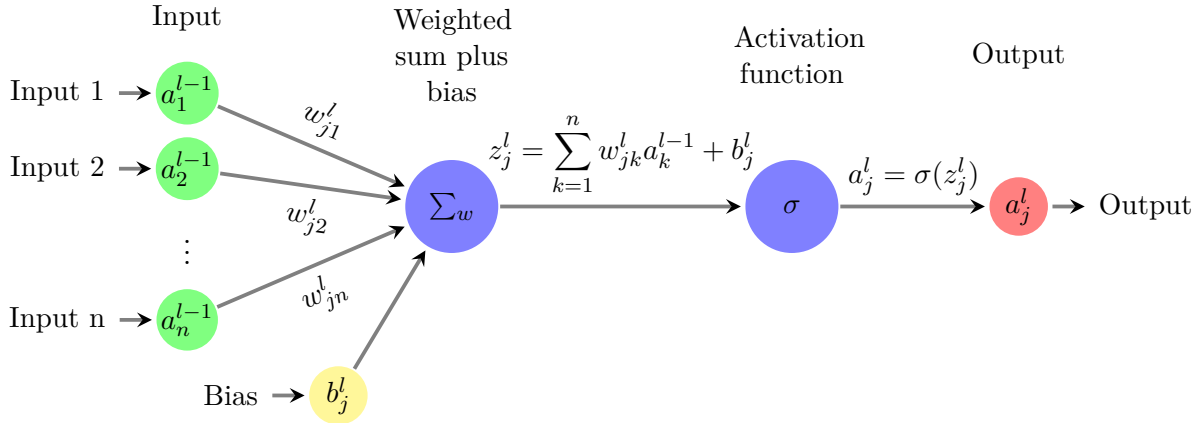


Figure 5.2: Schematic representation of how the inputs of a single neuron are converted to a final output which will go to the next layer of the network to be one of the new inputs in that layer.

the following relation in a compact form:

$$a^l = \sigma(z^l) = \sigma(w^l a^{l-1} + b^l), \quad l \in \{1, \dots, L\}. \quad (5.2)$$

$$a^1 = \sigma(z^1), \quad z^1 := x \quad (5.3)$$

With this notation in mind, we remark that for the input neurons ($l = 1$) there are neither input weights nor bias, so the “weighted input” for the first layer is $z^1 = x$ and the outputs $a^1 = (a_1^1, \dots, a_n^1)^T$ are obtained by applying some activation function (**usually the identity**) to the input values $x = (x_1, \dots, x_n)^T$.

To sum up, given a weighted DAG, it leads to a NN function F_θ which depends on linear combinations and compositions of the neural functions f_j^l , and the vector of parameters $\theta = (w_{jk}^l, b_j^l)$ for $l = 1, \dots, L$ and j, k varying according to the number of neurons in each of the layers. In particular, F_θ takes $x \in \mathcal{X}$ as input, and outputs $a^L \in \mathcal{Y}$, i.e. $F_\theta(x) = a^L$.

5.2 Training of Deep Neural Networks

In the last section we introduced neural networks as parametrized functions F_θ but we do not know yet how to tune the parameters θ so that the neural network performs the task we want. The process to tune NN parameters is also known as *training*.

Assuming we have a train dataset X and some *loss* function $J_\theta(X) := J(\{F_\theta(x) | x \in X\})$ that quantifies how well (and how bad) a neural network F_θ is performing according to its learning goal, and also enough regularity on F_θ and J_θ ; as we will see more concretely in the coming sections, then the way of tuning the parameters θ is by minimizing J_θ with respect to θ . Note that the J_θ is a function of the whole training set X , and in fact we are saying that it can be written as a function of the outputs $a^L = F_\theta(x)$.

Usual loss functions can be decomposed as a sum of the form,

$$J_\theta(X) = \frac{1}{n} \sum_{x \in X} \tilde{J}_\theta(x). \quad (5.4)$$

where $n = |X|$ is the number of training samples and \tilde{J}_θ a function that depends on θ and a single input observation x .

5.2.1 Gradient descent

The idea of Gradient Descent (GD) ³ algorithm in the context of optimization ⁴ problems is to tweak the parameters θ of the function to optimize iteratively, by following the opposite direction of the gradient. Once the gradient is zero, hopefully we would have reached a minimum. GD is a particular case of the so-called descent methods [26] ⁵.

In equation (5.6) we show the simple expression used to take a GD step for a given training set X ,

$$\theta^{k+1} \leftarrow \theta^k - \alpha \nabla_{\theta} J_{\theta}(X)|_{\theta=\theta^k}, \quad (5.5)$$

where $\alpha > 0$ is the *learning rate* (or *stepsize*) hyperparameter that determines how “big” are the jumps when going from one parameter value to the next. In particular, α indirectly determines the speed of convergence (if any) of the algorithm.

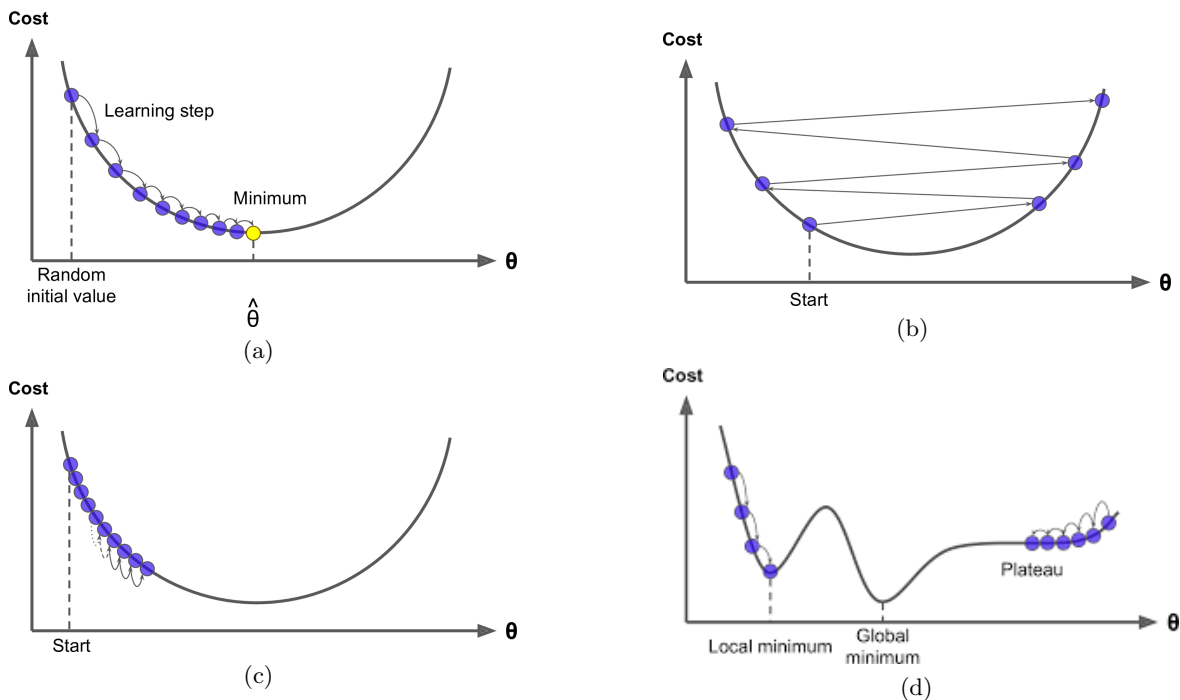


Figure 5.3: Gradient descent performance with (a) appropriate learning rate (b) big learning rate (c) small learning rate. (d) Illustration of possible problems for convergence to a global minimum (the presence of plateaus and local minima). Figures taken from [13].

In the pictures from Figure 5.5 we illustrate the role of α assuming that the computed gradients are very close the exact ones (since this would not always be the case for SGD). On the one hand, very big learning rates could make the GD algorithm to diverge (Figure 5.3b). On the other hand, very small learning rates will probably make the algorithm converge, but the number of steps required to do so will probably be very large (Figure 5.3c). Hence, we see that one should try to find some appropriate values for α so that both convergence and relatively fast convergence is achieved (Figure 5.3a).

We should be aware of the fact that GD algorithm for general optimization problems is usually the worst algorithm. As an example, consider the non-convex function depicted in Figure (5.3d) for

³Also known as *steepest descent* method.

⁴With the goal of minimizing some function.

⁵In case we wanted to maximize, we would follow the direction of the gradient, and not the opposite.

the case of having a single parameter θ . If the initialized parameter is in the very left, then it will converge to a *local minimum* instead of the *global minimum*. If it starts on the very right, then it will take very long time to surpass the plateau so that if we stop the iterations too early we will never reach the global minimum [13].

However, in the context of NN, loss functions are usually convex, so that in these cases GD is the best algorithm (in terms of speed and performance) to be used assuming we use good approximations for the gradients. Otherwise, we take the risk of finding local minima instead of global, stay in a plateau for a long time or even having divergence of the method if the learning rate is not carefully chosen at each step.

citation?

Finally, it is worth commenting about *feature scaling*, a technique that is commonly used when gradient descent or its variants are used for minimization ⁶. To understand why feature scaling usually improves convergence, we refer to Figure 5.4 and quote a very nice answer from the Q&A platform quora ⁷:

Essentially, scaling the inputs gives the error surface a more spherical shape, where it would otherwise be a very high curvature ellipse. Since gradient descent is curvature-ignorant, having an error surface with high curvature will mean that we take many steps which aren't necessarily in the optimal direction. When we scale the inputs, we reduce the curvature, which makes methods that ignore curvature (like gradient descent) work much better. When the error surface is circular (spherical), the gradient points right at the minimum, so learning is easy.

where “error” refers to the loss function.

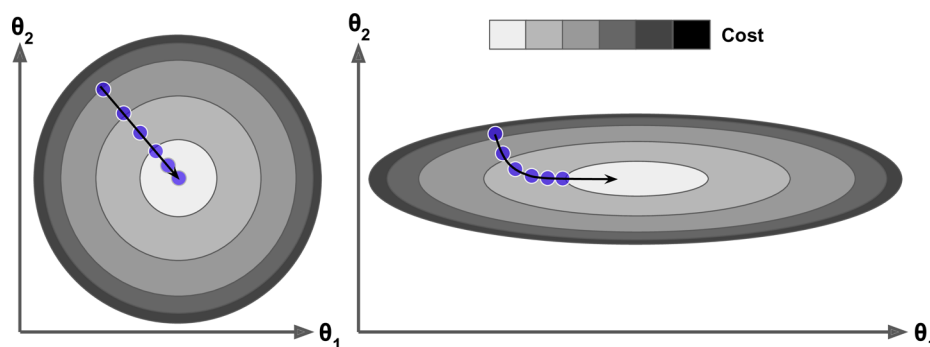


Figure 5.4: Representation of a loss function dependent on two parameters $\theta = (\theta_1, \theta_2)^T$ after feature scaling (on the left) and before feature scaling (on the right). Figures taken from [13].

In short, the idea is that since the learning rate is usually a constant value that multiplies all the components of the gradient vector in every GD step; if the loss function is more elongated in some direction than in others, a learning rate able to improve parameters faster in less elongated directions, it will not be that faster in for more elongated directions.

However, we have to be aware of the fact that feature scaling will not always improve convergence, perhaps in some cases this technique may worsen it.

⁶More usually in the context of classification. A case example would be the one we present in the next section 5.4.

⁷See: <https://www.quora.com/Why-does-mean-normalization-help-in-gradient-descent>

Mini-batch Gradient Descent

Note that if the cost function decomposes as in (5.4), the GD step can be written as

$$\theta^{k+1} \leftarrow \theta^k - \frac{\alpha}{n} \sum_{x \in X} \nabla_{\theta} \tilde{J}_{\theta}(x) \big|_{\theta=\theta^k}. \quad (5.6)$$

In practice, to compute ∇J_{θ} , we need to compute the gradients $\nabla_{\theta} \tilde{J}_{\theta}(x)$ for each training input x and then average them. However, if the number of training inputs is very large, this computations can take a long time and the process of learning becomes very slow.

A way to overcome this computational problem, what is done is to consider a random sample of $m < n$ (usually $m \ll n$) observations $x_1, \dots, x_m \in X$, we call it a *mini-batch*, and approximate the loss function gradient from (5.6) as

$$\nabla J_{\theta} = \frac{1}{n} \sum_{x \in X} \nabla_{\theta} \tilde{J}_{\theta}(x) \approx \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \tilde{J}_{\theta}(x_i). \quad (5.7)$$

The particular case with $m = 1$ is the so-called Stochastic Gradient Descent (SGD).

As one can intuitively deduce, the *batch-size* hyperparameter m leads to a trade-off between being very fast but perhaps less accurate ($m = 1$) and being very slow but probably very accurate ($m = n$) when computing the estimations of the gradient. Another point of view would be the exploration-exploitation trade-off for small or big values of m , respectively.

Adam optimizer

In this work, when training neural networks we will be using one of the most recent and accepted optimization methods by the time of writing this thesis. This method is known as *Adam* (adaptive moment estimation) [27].

Although the majority of methods used in practice, such as *Adam*, technically differ from standard gradient descent; they have the common property of needing the gradients of the loss function with respect to the parameters and are inspired by equation (5.6). In particular, for *Adam* there are computed estimations of the first and second moments (mean and variance) of the gradients and combined to play the same role as the gradient in equation (5.6).

Moreover, it is an adaptive learning rate algorithm, which means that starting from an initial learning rate α_0 neither too big nor too small, the algorithm itself is able to tweak the learning rate in every iteration to achieve better convergence and performance.

5.2.2 The Backpropagation algorithm

So far we have seen how to update the parameters $\theta = (w_{jk}^l, b_j^l)_{j,k,l}$ ⁸ of a NN using the gradients of a loss function, but we have not explained how to compute them yet. Note that the GD step from equation (5.6) can be written in terms of w_{jk}^l and b_j^l (instead of the whole parameter vector θ) in

⁸To simplify the notation in the coming explanations, we omit the lower and upper indices and just write w or b to refer to some weight or bias coefficient.

one equation for each parameter,

$$w_{jk}^l \leftarrow w_{jk}^l - \alpha \frac{\partial J_\theta}{\partial w_{jk}^l}, \quad (5.8)$$

$$b_j^l \leftarrow b_j^l - \alpha \frac{\partial J_\theta}{\partial b_j^l}. \quad (5.9)$$

Thus, compute the gradient of the loss function is equivalent to compute all those derivatives with respect to individual parameters.

The most famous and fast algorithm for computing gradients is known as *backpropagation* [28] and is “the workhorse of learning in Neural Networks” [29].

In section 5.2 we indirectly considered two assumptions about the loss J_θ that should be taken so that backpropagation works [29]:

- 1) First assumption is that J_θ can be decomposed as in equation (5.4). To simplify notation, we denote $\tilde{J}_\theta(x)$ simply by J_x , so that

$$J_\theta \equiv J_\theta(X) = \frac{1}{n} \sum_{x \in X} J_x \quad (5.10)$$

- 2) The second assumption is that given an input $x \in X$, the loss function J_θ can be written as a function of the outputs $a^L(x) := F_\theta(x)$, i.e.,

$$J_\theta(X) \equiv J(\{a^L(x) \mid x \in X\}). \quad (5.11)$$

We are going to abbreviate $a^L(x)$ simply by a , so that $J_\theta(x) \equiv J(a)$.

The first assumption is made because what is computed in backpropagation are the partial derivatives $\frac{\partial J_x}{\partial w}, \frac{\partial J_x}{\partial b}$ for a single training input $x \in X$, and then are computed the “total” derivatives $\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}$ by averaging over all x (due to equation (5.10) and the linearity of the derivatives). From now, to simplify notation we omit the subindex x for the cost function J_x of a single observation.

The second assumption is made so that one can consider the gradients $\nabla_a J$ of J with respect to the output “ a ”⁹ of the neural network and be able to derive the backpropagation equations in a very simple way.

Before going on with backpropagation we need to introduce some definitions: the hadamard product of two vectors (or matrices) and what we will call here the δ -error measure.

Definition 1 (Hadamard (or Schur) product). *Let K be a field and E a K -vector space of dimension n . Let $u, v \in E$ and $A, B \in \mathcal{M}_{r \times c}(K)$. The Hadamard product of u and v , and of A and B is just the element-wise product denoted \circ and defined as:*

- $(u \circ v)_j = u_j \cdot v_j, j \in \{1, \dots, n\},$
- $(A \circ B)_{ij} = A_{ij} \cdot B_{ij}, i \in \{1, \dots, r\}, j \in \{1, \dots, c\},$

where \cdot denotes the product in K .

⁹Abbreviation for a^L .

Definition 2 (δ -error). The δ -error in layer l , we denote by δ^l , is the vector whose components are defined as

$$\delta_j^l = \frac{\partial J}{\partial z_j^l}, \quad (5.12)$$

and we say that δ_j^l is the error of neuron j in layer l .

To understand the intuitive meaning of δ_j^l , imagine that for some reason the j -th neuron in layer l “mistakenly” outputs $\sigma(z_j^l + \Delta z_j^l)$ instead of $\sigma(z_j^l)$. Then, the error committed will be given by

$$J(z_j^l + \Delta z_j^l) - J(z_j^l) = \frac{\partial J}{\partial z_j^l} \Delta z_j^l + \mathcal{O}(\Delta z_j^l)^2 = \delta_j^l \Delta z_j^l + \mathcal{O}(\Delta z_j^l)^2 \approx \delta_j^l \Delta z_j^l, \quad (5.13)$$

for Δz_j^l sufficiently small.

Hence, δ_j^l is an amplification factor of how much varies the loss function J given an “error” Δz_j^l on the weighted input of the j -th neuron in layer l .

The backpropagation equations

We are now ready to derive the fundamental equations for the backpropagation algorithm. In Theorem 5.1 we formalize the four equations that are needed, and give a proof of them that only requires basic multi-variable calculus knowledge and the application of the chain rule.

The first equation (BP1) is an expression for the δ -error δ^L in the output layer, which is then used by the second equation (BP2) to recursively compute the δ -errors in the other layers. Finally, the last two equations (BP3) and (BP4) are the ones that allow us to compute the gradients with respect to the parameters of the neural network assuming that the δ -errors given by the other two equations have been computed and the activations (the outputs) in each neuron are available.

Theorem 5.1 (The backpropagation equations). *Consider a loss function J satisfying the two assumptions presented at the beginning of this section, and an activation function σ that is applied by each neuron of a Neural Network, both differentiable in every point. Then, the following relations are satisfied:*

$$\delta^L = \nabla_a J \circ \sigma'(z^L), \quad (\text{BP1})$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \circ \sigma'(z^l) \quad (\text{BP2})$$

$$\frac{\partial J}{\partial b_j^l} = \delta_j^l \quad (\text{BP3})$$

$$\frac{\partial J}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \quad (\text{BP4})$$

Proof.

(BP1) First of all note that since $\nabla_a J = \left(\frac{\partial J}{\partial a_1^L}, \dots, \frac{\partial J}{\partial a_n^L} \right)$ by definition, the j -th component of (BP1) is

$$\delta_j^L = \frac{\partial J}{\partial a_j^L} \sigma'(z_j^L). \quad (5.14)$$

If we start from the definition 5.12 of δ_j^l for $l = L$ and apply the chain rule thinking J as a function of $a = a^L$ with components a_k^L that can depend on z_j^L (among other things), then

$$\delta_j^L = \frac{\partial J}{\partial z_j^L} = \sum_{k=1}^n \frac{\partial J}{\partial a_k^L} \frac{\partial a_k^L}{\partial z_j^L} = \frac{\partial J}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L}, \quad (5.15)$$

where in the last equality we have used the fact that only the activation a_j^L , i.e. the output of the j -th neuron in the L -th layer, depends on z_j^L (the other neurons depend on z_i^L for some $i \neq j$), so that $\frac{\partial a_k^L}{\partial z_j^L} = 0$ for $k \neq j$. Finally, remember that by definition $a_j^L = \sigma(z_j^L)$, so that $\frac{\partial a_j^L}{\partial z_j^L}$ can be written as $\sigma'(z_j^L)$. Thus, from (5.15) we have arrived to equation (5.14). This ends the proof for (BP1).

(BP2) Given that w^{l+1} is the matrix with components w_{jk}^{l+1} , its transpose will have components $(w^{l+1})_{jk}^T = w_{kj}^{l+1}$. Hence, the j -th component of $((w^{l+1})^T \delta^{l+1})$ will be given by the product of the j -th row of $(w^{l+1})^T$ and δ^{l+1} :

$$((w^{l+1})^T \delta^{l+1})_j = \sum_{k=1}^n (w^{l+1})_{jk}^T \delta_k^{l+1} = \sum_{k=1}^n w_{kj}^{l+1} \delta_k^{l+1}.$$

With this, the j -th component of the equation (BP1) we want to prove is

$$\delta_j^l = \left(\sum_{k=1}^n w_{kj}^{l+1} \delta_k^{l+1} \right) \sigma'(z_j^l). \quad (5.16)$$

So it is enough to prove (5.16) for an arbitrary j . With the same idea we used to prove equation (BP1), we use the definition (5.12) of δ_j^l and apply the chain rule but now with respect to z_k^{l+1} :

$$\delta_j^l = \frac{\partial J}{\partial z_j^l} = \sum_{k=1}^n \frac{\partial J}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_{k=1}^n \delta_k^{l+1} \frac{\partial z_k^{l+1}}{\partial z_j^l} \quad (5.17)$$

In order to compute $\frac{\partial z_k^{l+1}}{\partial z_j^l}$, remember that by definition

$$z_k^{l+1} = \sum_i w_{ki}^{l+1} a_i^l + b_k^{l+1} = \sum_i w_{ki}^{l+1} \sigma(z_i^l) + b_k^{l+1}$$

Now if we take derivatives with respect to z_j^l , the bias term cancels out and, in the summation, the only term that survives is that with $i = j$. Thus,

$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \frac{\partial \sigma(z_j^l)}{\partial z_j^l} = w_{kj}^{l+1} \sigma'(z_j^l) \quad (5.18)$$

Putting (5.16) in (5.17) we obtain

$$\delta_j^l = \sum_{k=1}^n \delta_k^{l+1} w_{kj}^{l+1} \sigma'(z_j^l), \quad (5.19)$$

which is exactly (5.16) since $\sigma'(z_j^l)$ does not depend on k and it can be taken out of the summation.

(BP3) Starting from the definition of δ_j^l and applying the chain rule with respect to b_j^l we arrive to

$$\delta_j^l = \frac{\partial J}{\partial z_j^l} = \frac{\partial J}{\partial b_j^l} \frac{\partial b_j^l}{\partial z_j^l} \quad (5.20)$$

Now given that $z_j^l = \sum_{k=1}^n w_{jk}^l a_k^l + b_j^l$, we also have that $b_j^l = z_j^l - \sum_{k=1}^n w_{jk}^l a_k^l$, and taking the derivative with respect to z_j^l we trivially obtain that $\frac{\partial b_j^l}{\partial z_j^l} = 1$. And this proves (BP3).

(BP4) To find the expression for $\frac{\partial J}{\partial w_{jk}^l}$ we apply the chain rule with respect to z_j^l :

$$\frac{\partial J}{\partial w_{jk}^l} = \frac{\partial J}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} = \delta_j^l \frac{\partial z_j^l}{\partial w_{jk}^l}.$$

Now since $z_j^l = \sum_s w_{js}^l a_s^{l-1} + b_j$, if we derivate with respect to w_{kj}^l , the only term of the sum that survives is that corresponding to $s = k$. Hence,

$$\frac{\partial z_j^l}{\partial w_{jk}^l} = a_k^{l-1},$$

and this proves that

$$\frac{\partial J}{\partial w_{jk}^l} = \delta_j^l a_k^{l-1}.$$

□

To illustrate how the implementation of backpropagation looks like algorithmically, in Algorithm 2 we present a simplified pseudocode.

Algorithm 2 Backward-propagation algorithm [29]

```

1: procedure BACKPROPAGATION( $x \in X$ )
2:    $a^1 \leftarrow \sigma(x)$  ▷ Set the corresponding activation for the input layer.
3:   for  $l \in \{2, \dots, L\}$  do
4:     Compute  $z^l = w^l a^{l-1} + b^l$  and  $a^l = \sigma(z^l)$  ▷ Feedforward
5:   end for
6:
7:   Compute  $\delta^L$  using (BP1) ▷ Output  $\delta$ -error
8:
9:   for  $l \in \{L-1, \dots, 2\}$  do
10:    Compute  $\delta^l$  using (BP2) ▷ Backpropagation of the  $\delta$ -error
11:  end for
12:
13:  Use the computed  $\delta$ -errors and equations (BP3), (BP4) to obtain estimations of the derivatives of  $J$  with respect to the parameters of the neural network (weights and biases).
14:
15:  return:  $\frac{\partial J_x}{\partial w_{jk}^l}, \frac{\partial J_x}{\partial b_j^l}$  for all  $l, j, k$ .
16: end procedure

```

Once we know how to compute gradients, we can combine backpropagation with mini-batch gradient descent to train our neural network.

To end up with this section we need to introduce some standard notation in the NN framework. Imagine we have a training set X with n observations and consider mini-batches randomly sampled from X of size $m_B < n$. Assume for simplicity that n is some multiple s of n , so that we can split X in s mini-batches of size m_B . Then, we call a training epoch the fact of applying backpropagation plus mini-batch gradient descent (or any other optimization algorithm, such as *Adam*) to every mini-batch one time. Thus, if the number of epochs in our training algorithm is M , it means that we have used the training set X exactly M times to update the parameters of the NN.

In Algorithm (3) we present pseudocode for training a NN by updating parameters with backpropagation and mini-batch gradient descent. Note that we could use another algorithm such as *Adam* instead of mini-batch gradient descent which also uses mini-batches to estimate gradients, but the parameter's update rules are slightly different and more sophisticated.

It is very important to notice that the loop over minibatches (**Que loop es la clave para paralelizar? ahora no lo veo claro**)

Que loop es la clave para paralelizar? ahora no lo veo claro

Algorithm 3 NN training.

```

1: procedure TRAIN(A training dataset  $X$ ,  $M$ ,  $m_B$ ,  $s$ )
2:   for epoch  $\in \{1, \dots, M\}$  do ▷ Begin loop over epochs
3:     Split  $X$  in  $s$  mini-batches  $X_1, \dots, X_s$  of size  $m_B$  randomly.
4:     for  $i \in \{1, \dots, s\}$  do ▷ Begin loop over mini-batches  $X_i$ 
5:       for  $x \in X_i$  do
6:          $\frac{\partial J_x}{\partial w_{jk}^l}, \frac{\partial J_x}{\partial b_j^l} \leftarrow \text{BACKPROPAGATION}(x)$  ▷ Backpropagation
7:       end for
8:
9:       for  $l \in \{L, L-1, \dots, 2\}$  do
10:        Update the NN parameters as ▷ Mini-batch Gradient Descent step

$$w_{jk}^l \leftarrow w_{jk}^l - \frac{\alpha}{m_B} \sum_{x \in X_i} \frac{\partial J_x}{\partial w_{jk}^l}, \tag{5.21}$$


$$b_j^l \leftarrow b_j^l - \frac{\alpha}{m_B} \sum_{x \in X_i} \frac{\partial J_x}{\partial b_j^l}. \tag{5.22}$$

11:      end for
12:
13:    end for ▷ End loop over mini-batches  $X_i$ 
14:  end for ▷ End loop over epochs
15: end procedure

```

5.3 Improving the way neural networks learn

In the context of training NN there are two main problems that one has to deal with in order to obtain a NN model with good performance (e.g. the accuracy of the predictions or desired outputs):

- The first is known as the vanishing or exploding gradients problem.
- The second, which is more general in the context of Machine Learning, is known as *overfitting*.

Vanishing gradients problem appears when in some layers the derivatives with respect to some parameters become very small so that the updates of the parameters given by the equations (5.21), (5.22) in algorithm 3 - for instance - become insignificant (we say that the parameters *learn slowly*). Similarly, *exploding gradients* appear when the derivatives become very big in value so that the training of the parameters usually becomes unstable. This processes of the gradients becoming small or big, is propagated and accentuated more and more as backpropagation goes from the output layer to lower layers.

On the other hand, we can think of *overfitting* as the fact that the NN have memorized all the train dataset, in the sense that it outputs correctly the desired values for every train observation in that dataset. However, generally it would imply that the weights and biases have been adapted to the train dataset so much that when trying to use the NN with new observations not used during training, the performance is reduced considerably. Thus, it is preferred to avoid overfitting in the sense of training the NN in a way that it is capable of learning from the train dataset but also generalize properly to data that has never seen before.

In this section we present some techniques and options that can help to avoid the vanishing (exploiting gradients) gradients and overfitting, as well as improve the speed of convergence during training. We follow mainly [13].

5.3.1 Activation functions

One of the main reasons for vanishing gradient problems can be attributed to activation functions.

Traditionally, due to the activation functions used in models for describe biological neurons, there have been used “S”-shaped functions such as the sigmoid (or *logistic*, see Figure 5.5a) and the hyperbolic tangent functions as activation functions to define neural networks. We say that these type of activation functions “saturate” (to either 0 or 1 in the case of the sigmoid) in the extremes if the input values are quite small or quite big. In these cases the derivatives σ' are nearly zero, and if we remember the backpropagation equations from Theorem (5.1), the derivatives of the loss with respect to the parameters are proportional to σ' so that they will become small if the neurons “saturate”.

Another activation that was popular and also inspired on biological NN was the so-called ReLU function defined as

$$\text{ReLU}(z) = \max(0, z). \quad (5.23)$$

Then, variants of this function have been proposed by several authors and are reported to be better for training in terms of avoiding vanishing (or exploiting) gradient problems due to “saturation”.

In Figures 5.5b, 5.5c and 5.5d we show the currently most used variants of the ReLU activation function.

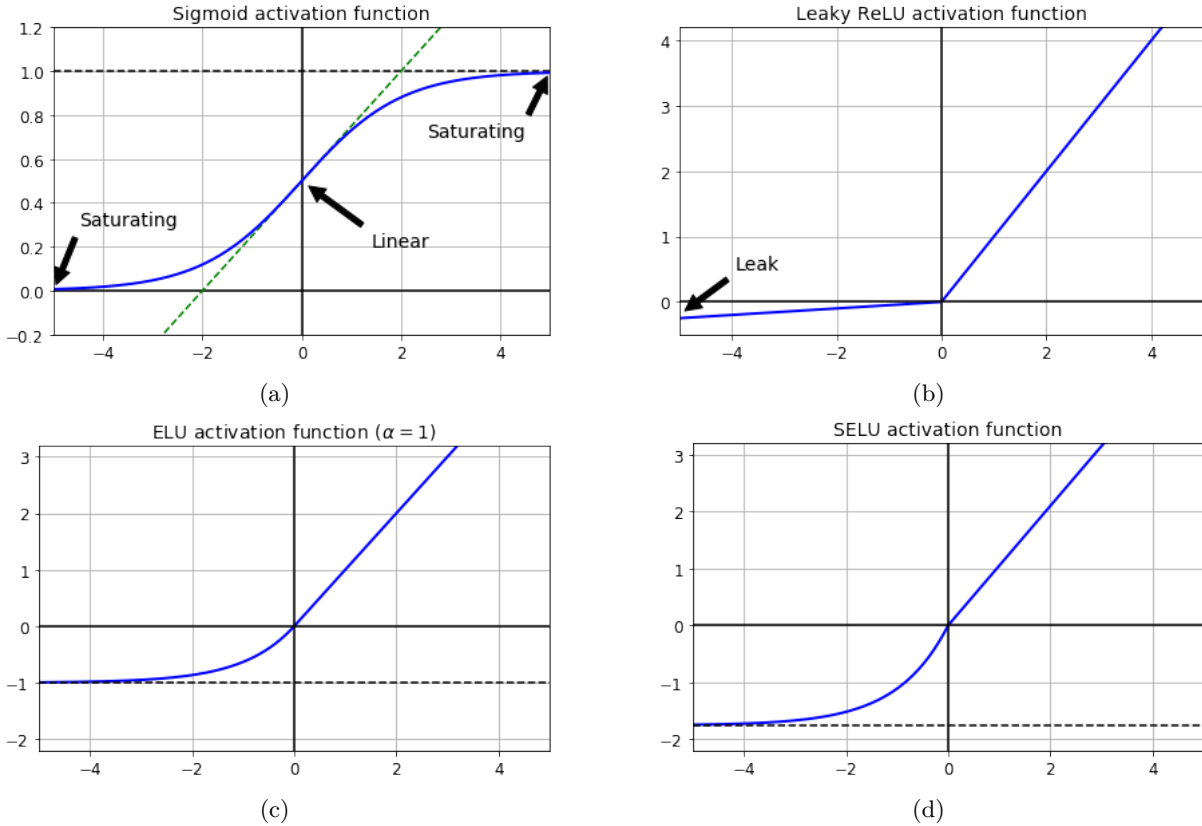


Figure 5.5: (a) Sigmoid (b) Leaky-ReLU [?] (c) ELU with $\alpha = 1$ [30] (d) SELU with $\beta \approx 1.05$, $\alpha \approx 1.67$, [31]. Figures taken from [13].

Mathematically, these functions are defined as follows:

$$\text{Leaky-ReLU}(z) = \max(0, \alpha z), \quad \alpha > 0 \quad (5.24)$$

$$\text{ELU}_{\alpha}(z) = \begin{cases} \alpha(\exp(z) - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases} \quad (5.25)$$

$$\text{SELU}_{\alpha}(z) = \beta \text{ELU}_{\alpha}(z), \quad \beta > 1 \quad (5.26)$$

Note that for the activation functions we have considered, apart from the sigmoid, all are non differentiable in any neighbourhood of zero, so technically they do strictly not satisfy the conditions of the backpropagation theorem 5.1. However, in practice we will not have values that small.....

5.3.2 Weight and biases initialization

When we explained the backpropagation algorithm to compute the gradients of the loss function with respect to the parameters of a NN, we talked about how to update them, but we did not talk about how to initialize them. ¹⁰

It has been seen that if the weights (and biases) are initialized randomly (for example with a normal distribution with mean 0 and variance 1) with the same distribution, the variance of the outputs a^l to a given layer are much greater than the variance of the inputs to this layer. Thus, going forward

¹⁰talk quizás no es el verbo adecuado

in the network, the variance keeps increasing after passing each layer until at some point it is very probable that from some layer on the activation functions saturate (if we use sigmoid or hyperbolic tangent, for instance).

In 2010 Xavier Glorot and Yoshua Bengio proposed a way to initialize weights to avoid these phenomena [32] for activation functions that can saturate. The idea was to initialize the weights of each layer based on the number of input neurons in the previous layer (n_{in}) and the number of neurons in the current layer (n_{out}). Inspired by these ideas, some authors [33] have proposed similar initialization criteria for different activation functions. In table 5.1 we summarize these criteria for initializing with uniform and normal random distributions.

Table 5.1: Initialization parameters for each type of activation function [13]. The normal distributions always have zero mean.

Activation function	Uniform distribution $[-r, r]$	Normal distribution
Logistic (sigmoid)	$r = \sqrt{\frac{6}{n_{\text{in}} + n_{\text{out}}}}$	$\sigma = \sqrt{\frac{2}{n_{\text{in}} + n_{\text{out}}}}$
Hyperbolic tangent	$r = 4\sqrt{\frac{6}{n_{\text{in}} + n_{\text{out}}}}$	$\sigma = 4\sqrt{\frac{2}{n_{\text{in}} + n_{\text{out}}}}$
ReLU (and its variants)	$r = \sqrt{2}\sqrt{\frac{6}{n_{\text{in}} + n_{\text{out}}}}$	$\sigma = \sqrt{2}\sqrt{\frac{2}{n_{\text{in}} + n_{\text{out}}}}$

These criteria to initialize weights is also known as *variance scaling initialization* and it is the one we will use for initializing NN.

5.3.3 Batch normalization

Although by combining novel activation functions with variance scaling initialization can significantly reduce vanishing/exploding gradients problems at the beginning of training [13], it does not ensure that these problems will not come back at some point during training.

In 2015 Sergey Ioffe and Christian Szegedy proposed a technique called *Batch Normalization* (BN) [34], and it has become very popular for being able to improve the majority of DNN. This technique consists of adding an operation BN just before applying the activation function of each layer.

Imagine we wanted to apply the BN operator in the l -th layer. Then, instead of having the output a^l of this layer defined as in equation (5.2), it would now be defined as

$$a^l = \sigma(\text{BN}(z^l)). \quad (5.27)$$

Basically, BN zero-centres and normalizes the weighted inputs z^l , and then scales and shifts the result using two additional parameters (per layer) ν^l (scale), η^l (offset) that can be learnt during training.

In Algorithm 4 it is described how the BN operator works. In order to zero-centring and normalize the inputs, the algorithm estimates the inputs' mean and standard deviation using the empirical mean and standard deviation over the observations of the current mini-batch. Note that when working with a mini-batch, equation (5.27) has to be thought as a component by component equality (one for each element of the mini-batch). Moreover, it is considered a tiny number ϵ to avoid division by zero when normalizing inputs (*smoothing term*).

Algorithm 4 Batch Normalization operator [34] for a given input mini-batch X_j of size m_B .

```

1: procedure BN( $z_x = z(x), \forall x \in X_j, \nu, \eta, \epsilon$ )
2:    $\mu_B \leftarrow \frac{1}{m_B} \sum_{x \in X_j} z_x$  ▷ Estimation of the mean
3:    $\sigma_B^2 \leftarrow \frac{1}{m_B} \sum_{x \in X_j} (z_x - \mu_B)^2$  ▷ Estimation of the variance
4:   for  $x \in X_j$  do
5:      $\hat{z}_x \leftarrow \frac{z_x - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$  ▷ Zero-centred and normalized input
6:      $\hat{z}_x \leftarrow \nu \hat{z}_x + \eta$  ▷ Scaled and shifted input
7:   end for
8:
9:   return:  $\hat{z}_x$  for all  $x \in X_j$ .
10: end procedure

```

A drawback of BN is the slowdown of training, so it is worth checking first if the DNN works well enough without applying BN.

5.3.4 Regularization techniques

Comentar solo la idea, breve, de cada una. Aunque luego no apliquemos algunas.

Early stopping

ℓ_1, ℓ_2 - regularization

Dropout

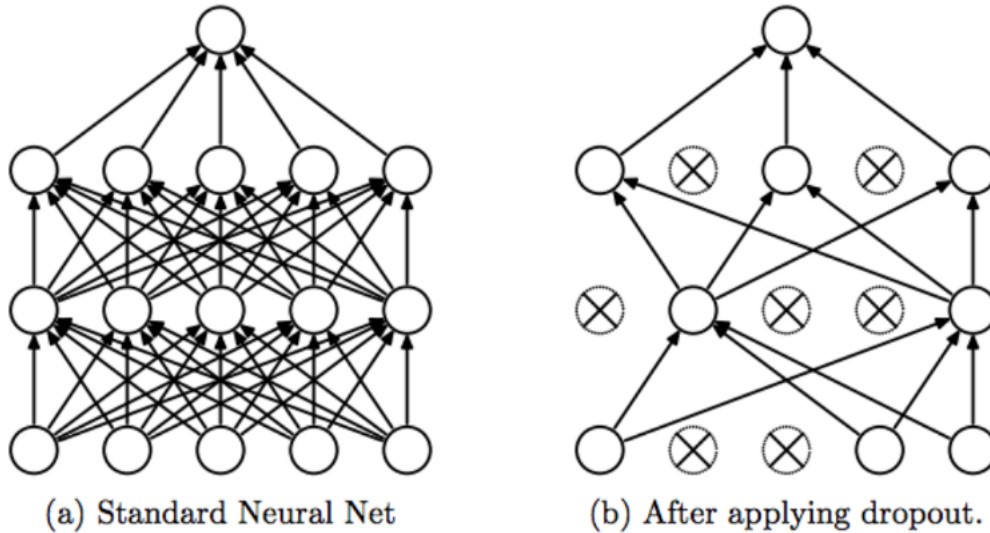


Figure 5.6: [35] **TO DO**

5.4 Learning a simple policy with a Deep Neural Network

In order to put in practice what we have learnt about Deep Neural Networks, in this section we consider one of the classical applications of Supervised Learning: classification.

First in section 5.4.1 we explain briefly how to use DNN for classification problems; then in 5.4.2 we set up a classification problem related to our product delivery problem, and finally in section 5.4.3 we show the results of solving that classification problem with DNN.

5.4.1 Deep Neural Networks and Classification

Imagine that we had data about the volume, the diameter and the weight of a sample of fruits, each of them being either an apple, an orange, a banana or a kiwi. In other words, we have four classes of fruits and data about three features (volume, diameter and weight) that can be easily measured from representatives to these classes.

Our goal is to have a model such that given a set of values for the volume, the diameter and the weight, it is able to predict from which class of fruit (either an apple, an orange, a banana or a kiwi) these values come from. To train the model we need a train set X such that each $x \in X$ is a triple of feature values, and a target set Y containing the labels $y(x)$ for each observation in X (e.g., the name of the fruit that corresponds to each observation in the train set).

The easiest approach would be to consider a DNN with one output neuron per class, and consider the output value a_i^L as a measure or score for the i -th class. Then, for instance, given a training instance $x \in X$, corresponding to class $y(x) = i$, we would expect a higher score in the i -th output value; i.e., $a_i^L > a_j^L$, $i \neq j$. In this framework a^L are called the *logits*.

Remember that in order to train a DNN (for classification, in this case), we need some loss function to minimize. As it is shown in Figure 5.7, what is done in practise is to apply a *softmax* activation function to the usual output layer which we now call “logits layer”. In this way, fixed $x \in X$, the final outputs p_i of the DNN are normalized to sum 1, so that they can be considered as a probability distribution $p(x)$ (see eq. (5.28)) over classes.

$$p_i = \frac{\exp(a_i^L)}{\sum_j \exp(a_j^L)} \quad (5.28)$$

Then we can define a “true probability distribution” $q = q(x)$ as follows:

$$q_i = \begin{cases} 1 & \text{if } y(x) = i \\ 0 & \text{if } y(x) \neq i \end{cases} \quad (5.29)$$

With these probability distributions p and q defined, one can consider the loss function to be a distance between them. For instance, the usual distance chosen is the *cross entropy*, which can be written as in equation (5.30).

$$H(p, q) = \sum_i q_i \log p_i \quad (5.30)$$

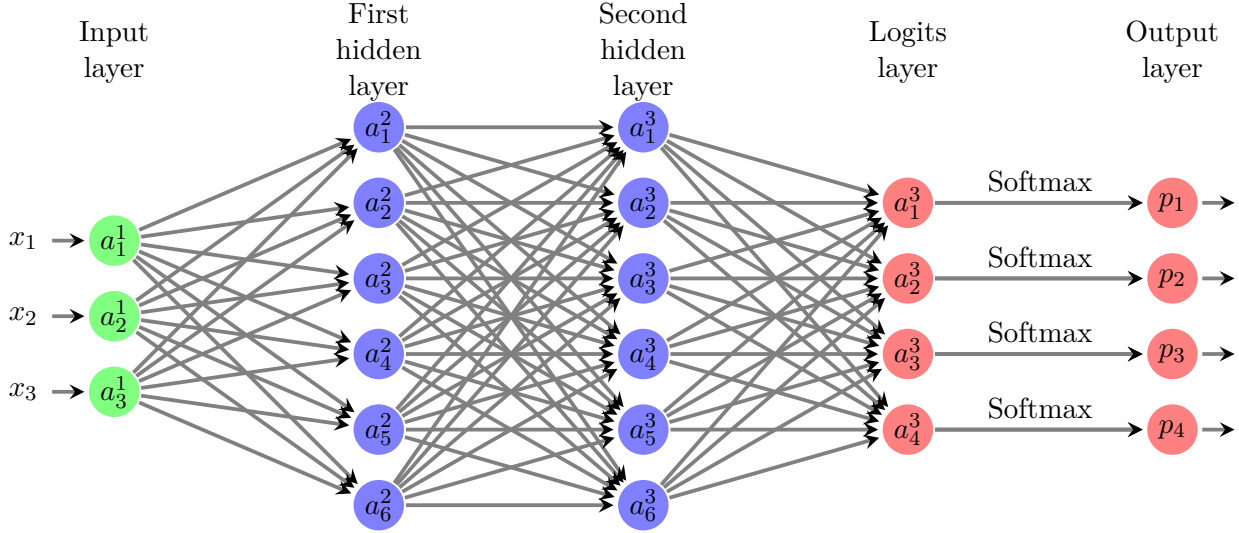


Figure 5.7: Representation of a feed-forward fully connected DNN for a classification problem where data has three features and there are four classes. We consider a simplified case with two hidden layers with six neuron each. The logits layer behaves in the same way as a hidden layer but then to the outputs of this layer, the *logits*, we apply a softmax activation function to normalize outputs so that they sum 1.

Since one usually uses mini-batches, the final loss function used for training a DNN such as the one in Figure 5.7 becomes

$$J_{\theta} = -\frac{1}{m_B} \sum_{x \in X_j} H(p(x), q(x)), \quad (5.31)$$

where X_j is a mini-batch of size m_B .

5.4.2 Simulation framework for learning a hard-coded policy

In order to get started with DNN and the product delivery problem we modelled with the point of view of Reinforcement Learning in chapter 3, we study the capability of a DNN of learning a simple hard-coded policy.

First, remember that we are restricting to the toy model we explained in section 3.3, and we consider a system with $n = 3$ shops and a single truck ($k = 1$). In this case, the state of the system will be given by the current stock of each shop, i.e., $s_t = (c_0, c_1, c_2)$, and the $n + 1 = 4$ possible actions for the truck would be going to either one of the shops or staying at the depot, i.e. $a_t = i$ for $i \in \{0, 1, 2, 3\}$ ¹¹.

Next, let's consider the policy π_m that sends the truck to the shop with minimum stock, if it is possible (i.e., if the load L of the truck fits all in the shop without surpassing the maximum level of stock), or makes it stay at the depot otherwise. Mathematically, speaking, we define π_m as follows:

$$\pi_m(s) = \begin{cases} \arg \min(s) & \text{if } \min(s) + L \leq C_{\arg \min(s)} \\ n & \text{otherwise.} \end{cases} \quad (5.32)$$

Finally, note that if we now generate a train dataset X containing possible states of the system using some distribution $I(s)$ (which is equivalent to generate a table with a row for each state and

¹¹Remember that $a_t = n$ is the action of staying at the depot.

columns contain the level of stock of each shop), and a target set Y given by $\pi_m(X)$ (where here it means to apply π_m to each row of X); then we are in the framework of a classification problem like the one we exposed in the previous section. In short, our train dataset consists of observations with three features (the stock of each shop) and an associated target value equal to the action to take according to π_m .

With these ideas, if we consider a DNN F_θ with the same structure as in Figure 5.7, that is, n input neurons, some hidden layers and $n + 1$ output neurons with a final softmax activation, if it learns properly after training, $\hat{\pi}_m := \arg \max F_\theta$ will act almost or exactly equally as the policy π_m , i.e., $\hat{\pi}_m(s) = \pi_m(s)$ for most of the possible states s we may use as input to de DNN.

Habría que comentar que $\arg \max$ se hace numérico, es decir, se hace un $\arg \max$ sobre los valores de salida de la red, que son números reales. Sería conveniente hacer una normalización para que sea un poco más estable, pero eso ya es un detalle de implementación.

5.4.3 Simulations and Results (TO DO)

In this section we present results for different DNN simulations prepared in the framework from the previous section.

To define DNN in PYTHON code and train them, we use the popular library Tensorflow from Google [2]. Moreover, to work with a product delivery environment adapted to our definitions and model of states and actions, we have developed an Open AI gym environment [3]. All the code to reproduce the results we show in this section can be found in a Jupyter notebook in the *Imitation-Learning* folder from [1].

First of all, we need to generate train, test and validation datasets containing state and action pairs of the form $(s, \pi_m(s))$. To choose states, we use a uniform random distribution for each of the components of an state $s = (c_0, c_1, c_2)$ ranging from 0 to the maximum capacities C_i , i.e., we sample $s \sim I(s)$ where I is such that $c_i \sim \text{Uniform}([0, C_i])$, $i = 0, 1, 2$.

For all simulations we have generated a train data set with 10^4 observations and a test and validation sets of 10^3 observations.

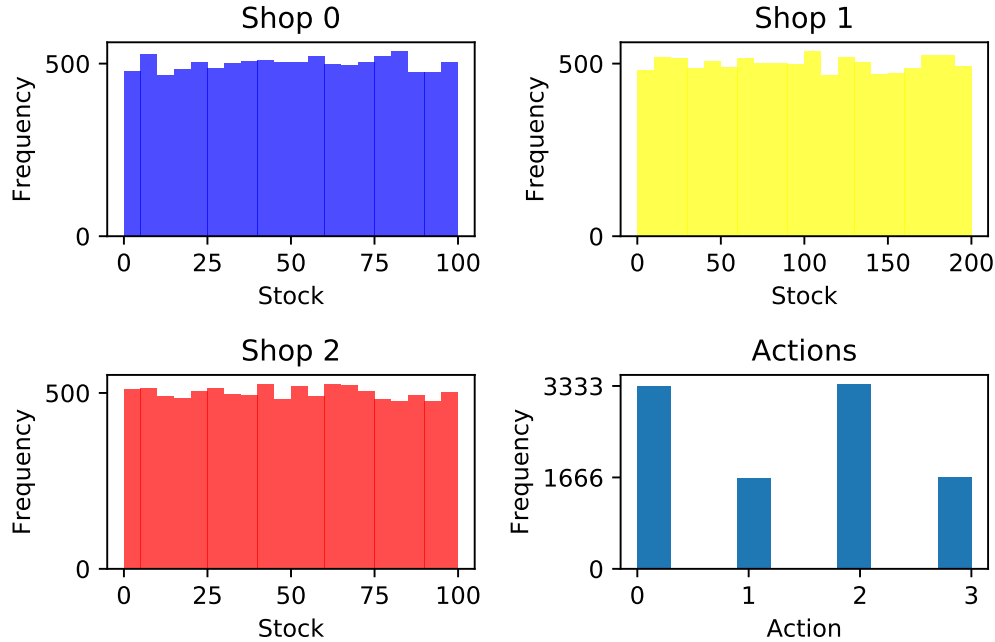


Figure 5.8: Histograms for the stocks of each and the corresponding actions for a train dataset of 10^4 observations. The bin size for the stocks is 20.

In Figure 5.8 we show the histograms of the train dataset, and in Table 5.2 the main features of the system used to generate the data.

Table 5.2: TO DO

	Maximum capacity	Daily consumption rate
Shop 1	100	16.5
Shop 2	200	4.0
Shop 3	100	10.5
Truck	50	either 0 or 50

We can see that the empirical distributions for the stock of each shop approximates very well a uniform between 0 and the corresponding maximum capacity. About the distribution of actions, we can see that the times that the action of sending the truck to shop 1 is about half the times that a truck is sent to shops 0 and 2. This makes sense since shop 1 has the double maximum capacity when compared to shops 0 and 2, so it is half likely that shop 1 has the minimum current load when initializing stocks uniformly. As a consequence, this leads the frequency of staying at the depot to be the same as sending a truck to shop 1.

- Explain simulation's set-up:

Table 5.3: Hyperparameter values

Hyperparameter	Value
Learning rate	0.01
<code>batch_size</code>	50
τ (episode length)	30
epochs	1000

Table 5.4

Simulation id	Classifier	Train accuracy	Test Accuracy
1	DNN ($n_1 = 100, n_2 = 50$)	0.9964	0.995
2	DNN ($n_1 = 20, n_2 = 10$)	0.9921	0.993
3	DNN ($n_1 = 10, n_2 = 5$)	0.9954	0.996
4	DNN ($n_1 = 4, n_2 = 4$)	0.9889	0.985
7	DNN ($n_1 = 500, n_2 = 200$)	0.9974	0.994

- Explain results of each simulation in a comparative way, figure 5.9:

TO DO

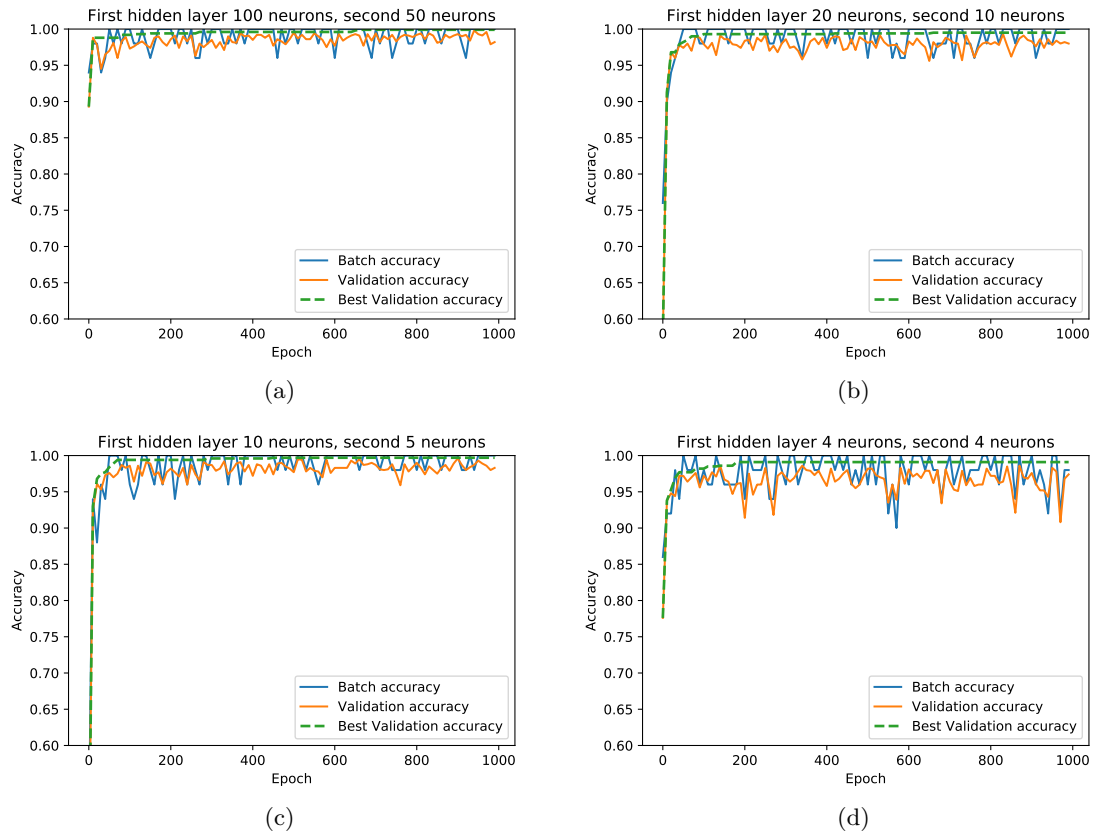


Figure 5.9

- Comment a bit figure 5.10

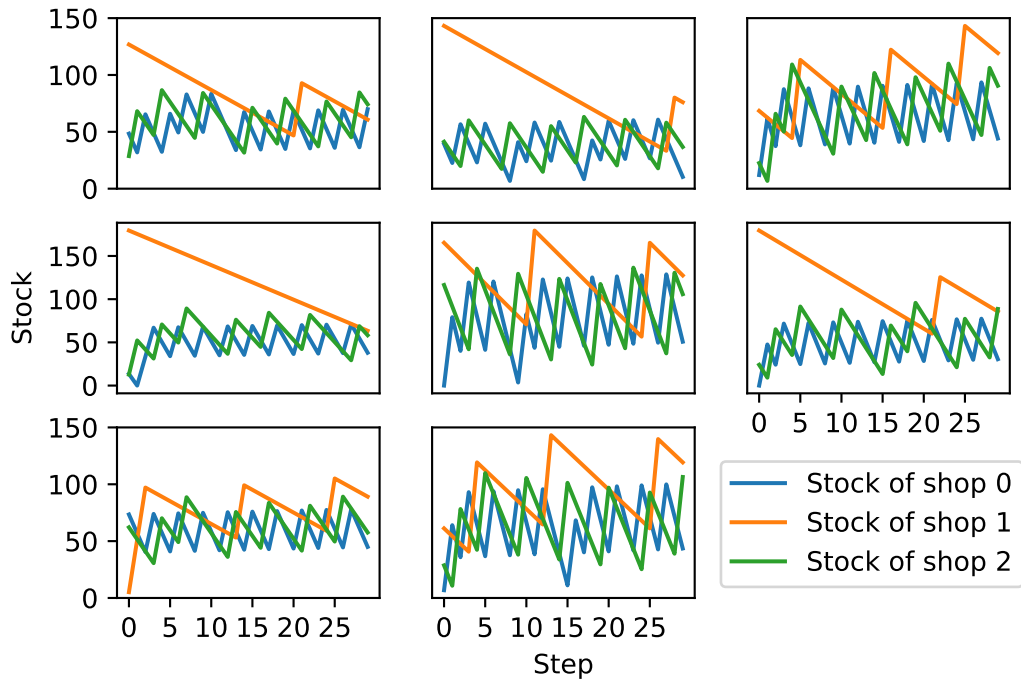


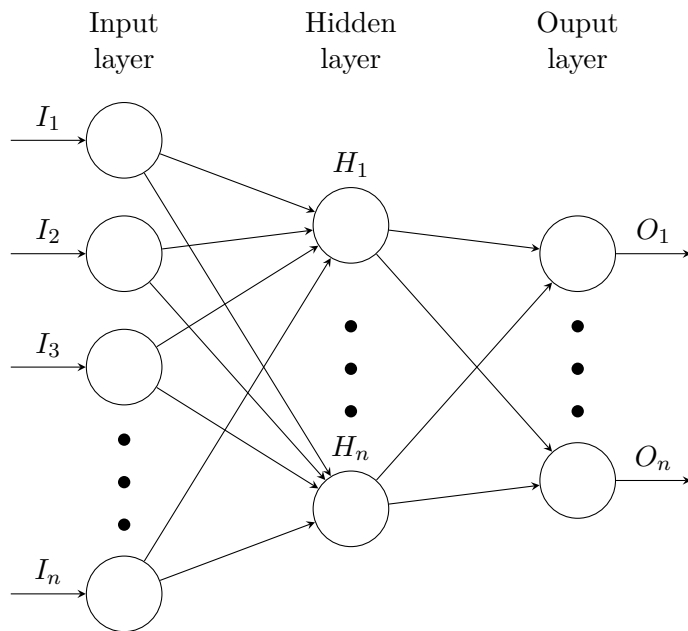
Figure 5.10: 8 test sims for sim 3

Some further work:

- Incorporar el cálculo de la confusion matrix para las DNN.,
- guardar los stados que clasifica mal para ver si es algun motivo senzillo'?

Chapter 6

Deep Reinforcement Learning



6.1 Policy Gradient algorithms

6.2 Simulations

- Deterministic simple approach (with different n and k)
- Stochastic simple approach (with noise) (future work)
- Deterministic complex approach (future work)
- Stochastic complex approach (future work)

6.3 Results

Chapter 7

(? En caso que hiciésemos algo más como lo de separar en k softmaxes y aprendiese, o bien rutas, etc.)

Chapter 8

Conclusions and Future Work

Bibliography

- [1] Daniel Salgado Rojo. Master’s thesis. <https://github.com/dsalgador/master-thesis>, 2018.
- [2] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.
- [4] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1st edition, 1994.
- [5] Craig Boutilier, Thomas L. Dean, and Steve Hanks. Decision-theoretic planning: Structural assumptions and computational leverage. *CoRR*, abs/1105.5460, 2011.
- [6] J. N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5):674–690, May 1997.
- [7] Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Proceedings of the 12th International Conference on Neural Information Processing Systems*, NIPS’99, pages 1057–1063, Cambridge, MA, USA, 1999. MIT Press.
- [8] Leslie Pack Kaelbling, Michael L. Littman, and Andrew W. Moore. Reinforcement learning: A survey. *CoRR*, cs.AI/9605103, 1996.
- [9] Martijn van Otterlo and Marco Wiering. *Reinforcement Learning and Markov Decision Processes*, pages 3–42. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.
- [10] Wikipedia contributors. Markov decision process — wikipedia, the free encyclopedia, 2018. [Online; accessed 18-March-2018].
- [11] Wikipedia contributors. Reinforcement learning — wikipedia, the free encyclopedia, 2018. [Online; accessed 18-March-2018].

- [12] Wikipedia contributors. Q-learning — wikipedia, the free encyclopedia, 2018. [Online; accessed 20-March-2018].
- [13] Aurélien Géron. *Hands-on machine learning with Scikit-Learn and TensorFlow : concepts, tools, and techniques to build intelligent systems*. O'Reilly Media, Sebastopol, CA, 6 edition, 2017.
- [14] Richard Bellman. *Dynamic Programming*. Princeton University Press, Princeton, NJ, USA, 1 edition, 1957.
- [15] Scott Proper and Prasad Tadepalli. Scaling model-based average-reward reinforcement learning for product delivery. In Johannes Fürnkranz, Tobias Scheffer, and Myra Spiliopoulou, editors, *Machine Learning: ECML 2006*, pages 735–742, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- [16] Christopher J. C. H. Watkins and Peter Dayan. Q-learning. In *Machine Learning*, pages 279–292, 1992.
- [17] Christopher John Cornish Hellaby Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Cambridge, UK, May 1989.
- [18] G. A. Rummery and M. Niranjan. On-line q-learning using connectionist systems. Technical report, 1994.
- [19] R. S. Sutton and A. G. Barto. Time-derivative models of Pavlovian reinforcement. In J. W. Moore and M. Gabriel, editors, *Learning and Computational Neuroscience*, pages 497–537. MIT Press, Cambridge, MA, 1990.
- [20] Gerald Tesauro. Temporal difference learning and td-gammon. *Commun. ACM*, 38(3):58–68, March 1995.
- [21] Tokic M and Palm G. Value-difference based exploration: Adaptive control between epsilon-greedy and softmax. *Lecture Notes in Computer Science*, 7006, Springer, Berlin, Heidelberg, KI 2011.
- [22] Eyal Even-Dar and Yishay Mansour. Convergence of optimistic and incremental q-learning. In *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*, NIPS’01, pages 1499–1506, Cambridge, MA, USA, 2001. MIT Press.
- [23] Francisco S. Melo. Convergence of q-learning: a simple proof. <http://users.isr.ist.utl.pt/~mtjspaam/readingGroup/ProofQlearning.pdf>.
- [24] Wikipedia contributors. Artificial neural network — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=Artificial_neural_network&oldid=840077868, 2018. [Online; accessed 9-May-2018].
- [25] Warren S. McCulloch and Walter Pitts. Neurocomputing: Foundations of research. chapter A Logical Calculus of the Ideas Immanent in Nervous Activity, pages 15–27. MIT Press, Cambridge, MA, USA, 1988.
- [26] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical Mathematics*. Texts in applied mathematics. Springer, 2000.
- [27] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.

- [28] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1. chapter Learning Internal Representations by Error Propagation, pages 318–362. MIT Press, Cambridge, MA, USA, 1986.
- [29] Michael A. Nielsen. Neural networks and deep learning, 2018.
- [30] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *CoRR*, abs/1511.07289, 2015.
- [31] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. *CoRR*, abs/1706.02515, 2017.
- [32] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.
- [33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *CoRR*, abs/1502.01852, 2015.
- [34] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR*, abs/1502.03167, 2015.
- [35] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.

Appendix A

A Product Delivery gym environment

Appendix B

Derivation of $J_{\text{levels}}^{(i)}$

Appendix C

Default Q-learning simulation parameter values

C.1 Model system parameters

In this appendix we present Table C.1 with the default parameters used to define the system of $n = 5$ shops and $k = 2$ trucks.

Table C.1

Parameter(s)	Description	Value(s)
C_i	Maximum capacity of the shops	[100, 200, 100, 800, 200]
L_i	Maximum capacity of the trucks	[70, 130]
b	Minimum level (fraction of C_i)	[0.02, 0.01, 0.05, 0.07, 0.08]
c	Danger level (fraction of C_i)	[0.31, 0.03, 0.16, 0.14, 0.26]
e	Maximum level (fraction of C_i)	[0.9, 0.9, 0.9, 0.85, 0.9]
a	(see Appendix B)	$b/10$
f	(see Appendix B)	$1 - (1 - e)/10$
d	(see Appendix B)	$p_0e + (1 - p_0)c$, $p_0 = 0.7$
w_{5i}	Weights from the depot to shop i (or depot if $i = 5$)	[32, 159, 162, 156, 156, 0]
M	(see Appendix B)	10
P_1	(see Appendix B)	-10^3
P_2	(see Appendix B)	-10^6
μ_1	Reward coefficient for the costs contribution	0
C_{costs}	(see equation (5.4))	0.051591
$\mu_{\text{extra},1}$	Reward coefficient a “a truck not delivering” contribution	10^{-6}
μ_2	Reward coefficient for the levels of stock contribution	10^{-6}