

# Learning to Recommend using a Deep Reinforcement Agent

*Hernan Santiago Gonzalez Toral*

*Department of Computer Science*

*University College London*

**Supervisor**

*Dr. Jun Wang*

*jun.wang.l@cs.ucl.ac.uk*

This report is submitted as part requirement for the

**MSc in Web Science & Big Data Analytics**

at

**University College London.**

It is substantially the result of my own work except  
where explicitly indicated in the text.

Department of Computer Science

University College London

September 5, 2016

The report may be freely copied and distributed provided the source is explicitly  
acknowledged.

**This page is purposely left blank.**

# Abstract

Nowadays, an ever-increasing volume of information on the Internet has made many websites offer recommendation systems to help customers to quickly find items that satisfy their needs. Such accurate systems should not act in a greedy manner and must balance the need to explore the current user preferences and the exploitation of the rating information for recommendation.

Reinforcement learning has demonstrated to be a potentially effective approach to solve the exploration/exploitation dilemma by generating recommendations as a sequential optimization problem; nevertheless, it has received relatively little attention and found only limited application as its computational complexity made the problem intractable when there is a large number of items to recommend.

This dissertation project aims is to demonstrate that recent advances in deep reinforcement learning are now useful to build a model that learns how to generate good recommendations by reducing the complexity of acting over a large set of items. A comparison between a reinforcement learning agent based on the Deep Deterministic Policy Gradient algorithm with k-nearest neighbors and two different item representations, and a new policy approach based on Factorization Machines demonstrate

**This side is purposely left blank.**

# Acknowledgements

Firstly, I would like to express my sincere gratitude to my supervisor and programme director Dr. Jun Wang for his advice and comments on this work, and for being our mentor at University College London.

I would also like to thank to the National Secretariat of Higher Education, Science, Technology and Innovation (SENESCYT), and especially, to Dr. Rafael Correa Delgado, president of the Republic of Ecuador, for believing in the skills of all Ecuadorian citizens studying in the best Universities around the world. Your trust and our sacrifice will be transformed into innovation that will turn our Country into a republic of opportunities.

Finally, I would like to thank to my mom, friends and family for their support and encouragement throughout my studies.

**This side is purposely left blank.**

# Contents

<b>1</b>	<b>Introduction</b>	<b>13</b>
1.1	Problem Overview . . . . .	13
1.2	Dissertation Objectives and Contributions . . . . .	15
1.3	Dissertation Objectives and Structure . . . . .	15
<b>2</b>	<b>Background</b>	<b>17</b>
2.1	Recommender Systems . . . . .	17
2.1.1	Collaborative Filtering . . . . .	18
2.1.2	Content-based RS . . . . .	19
2.1.3	Hybrid approaches . . . . .	20
2.1.4	The cold-start problem . . . . .	20
2.1.5	Trends in Recommendation Systems . . . . .	20
2.2	Reinforcement Learning . . . . .	22
2.2.1	Markov Decision Processes . . . . .	24
2.2.2	A model-free Reinforcement Learning Setup . . . . .	25
2.2.3	Exploration/Exploitation Dilemma . . . . .	26
2.3	Recommender Systems and Reinforcement Learning . . . . .	27
2.4	Deep Reinforcement Learning . . . . .	34
<b>3</b>	<b>Models</b>	<b>41</b>
3.1	Motivation . . . . .	41
3.2	Baseline Agent Definition . . . . .	42
3.3	FM Agent Model . . . . .	45
3.3.1	Proposed Approach . . . . .	45

<b>4 Experiments</b>	<b>47</b>
4.1 Reinforcement Learning Environment . . . . .	47
4.1.1 Recommender System Environment . . . . .	48
4.1.2 Simulation algorithm . . . . .	49
4.2 Data Set . . . . .	49
4.3 Evaluation Scheme . . . . .	50
4.4 Baselines and Experimental Settings . . . . .	50
4.5 Results . . . . .	52
<b>5 Conclusion and Discussion</b>	<b>53</b>
5.1 Summary . . . . .	53
5.2 Future Work . . . . .	53
<b>Appendices</b>	<b>54</b>
<b>A DDPG Algorithm using the FM-policy Architecture</b>	<b>55</b>
<b>B Transition probability Matrix algorithm</b>	<b>57</b>
<b>C Recommender System Simulated Environment</b>	<b>59</b>
<b>D Hyperparameter configuration for models</b>	<b>61</b>
<b>Bibliography</b>	<b>62</b>



# List of Figures

2.1	Typical user rating matrix . . . . .	18
2.2	Recommender Systems taxonomy . . . . .	21
2.3	Agent-Environment interface in a Reinforcement Learning problem	23
4.1	Feature vector representation for the FM model. Source: <i>Factor- ization machines</i> [17] . . . . .	52



# List of Tables

4.1	Properties of the 100-k and 1M Movielens datasets . . . . .	49
D.1	Hyperparameter and network settings for DRL models . . . . .	61



## Chapter 1

# Introduction

The volume of information available on the internet have been increasing rapidly with the explosive growth of the World Wide Web, e-Commerce and other online services, making difficult for users to find relevant products in a short period of time. To avoid this problem, many websites use recommendation systems to help customers to find items that satisfy their needs [1]. Suggestions for books on Amazon, or movies on Netflix, are real-world examples of Recommender Systems (RS) that have demonstrated that the evolution of RS and the web should go hand-in-hand to achieve the common goal of efficiently managing an ever growing set of items and thus generate meaningful recommendations to a collection of users with common interests.

### 1.1 Problem Overview

The design of such recommendation engines depends on the domain and the particular characteristics of the data available. While *Collaborative Filtering* systems analyze historical interactions alone, *Content-based* RS systems are based on user profile attributes; and finally, *hybrid* techniques combine both of these designs. Latest studies in CF [2] showed that combining conceptual and usage information can improve the quality of web recommendation.

Current RS typically act in a *greedy* manner by recommending items with the highest user ratings. However, greedy recommendations are suboptimal over the long term. This approach does not actively gather information on

user preferences and fails to recommend novel songs that are potentially interesting. A successful recommender system must balance the need to explore user preferences and to exploit this information for recommendation.

Therefore, the architecture of recommender systems and their performance on real-world problems is an active area of research. where approaches using a diverse Information Retrieval, Web Mining and Machine Learning techniques have been proposed [2]. Reinforcement learning[3] for example, can be used to try to solve the recommendation task.

In general, once a user selects a item based on a list of recommendations given by a RS, a new list of recommended items is then presented. This recommendation process can be then thought as a sequential process, where users choices are sequential by nature (e.g. a user buy a book by the author of a recent book she liked).

Indeed, reinforcement learning has demonstrated in the past to be a potentially effective approach in the domain and particularly good to solve the cold-start problem and in some extent, the exploration/exploitation trade-off [4]. Nevertheless, it has received relatively little attention and found only limited application as its computational complexity made the problem intractable when there is a high number of items to recommend.

On the other hand, most recent research [5] have focused on making a generalization of recent advances in Deep Learning [6] from Identically and Independently Distributed (i.i.d.) input to the non-i.i.d. (CF-based) input, and propose hierarchical Bayesian models that jointly performs deep representation learning for the content information and CF for estimating the ratings matrix. Experiments carried out shown that Deep Learning approaches can significantly outperform the state of the art. As a result, current research has been getting involved on using deep representational approaches that can lead to better recommendations.

Furthermore, a recent work presented by Dulac-Arnold et al. [7] use deep neural neural networks and reinforcement learning in large discrete action

spaces and propose a new policy architecture that not only allows RL methods to be applied to large-scale learning problems, and to operate efficiently with a large number of actions, but also is able to generalize over the action set in logarithmic time. This algorithm showed to a promising convergence of the model among different continuous environments.

## 1.2 Dissertation Objectives and Contributions

The main objective of this project is to show that a deep reinforcement learning model is able to perform well under the RS domain by learning how to generate good recommendations to users over a large set of items. After defining the recommendation problem under the reinforcement learning domain, we describe a baseline model using the Deep Deterministic Policy Gradient algorithm with a k-nearest neighbours (kNN) policy for action space reduction, and trained it using a content-based and a collaborative-filtering representation of items. Additionally, we propose a new policy architecture based on a Factorization Machine model for ranking recommendation, and compare the performance of the reinforcement agent under different policy definitions and network configurations.

A quantitative analysis on the effectiveness of the proposed policies and the average precision gathered by each reinforcement agent over the learning period showed that [\*\*\*](Missing explanation of experiments and conclusions)

## 1.3 Dissertation Objectives and Structure

The remaining chapters of this report are organized as follows. Chapter 2 defines the recommender problem and introduces Reinforcement Learning. Following, it reviews some related work that has been made for modeling recommender systems using reinforcement learning approaches and deep learning. The proposed model and its implementation details are presented in Chapter 3, while the experimental evaluation and results are analyzed in Chapter 4. Finally, we present a summary of lessons learned and conclusions, and propose some ideas for future work.





## **Chapter 2**

# **Background**

In this chapter, we will present all the necessary background concepts as well as some related work on recommendation systems and reinforcement learning methods that help to understand the work of this dissertation project. First, we give an overview of the existing recommendation systems techniques and approaches that have been used to solve the recommendation problem. Then, we introduce the basic concepts to define a reinforcement learning task and expose the first attempts made on the field to create reinforcement agents that learn to give recommendation to users. Finally, recent work on reinforcement learning algorithms using deep learning is reviewed.

## **2.1 Recommender Systems**

In recent years, there has been growing focus on the study of Recommender Systems (RSs), where different techniques have been published in order to address the problem of information overload on the Internet. Consequently, the evolution of RS and the web usually go hand-in hand. Resnick and Varian[1] defined RS as systems that help users limit their search by supplying a list of items that might interest them. Such systems has found an active application area for diverse Information Retrieval, Web Mining and Machine Learning techniques that help to solve typical recommendation problems. Additionally, different RS mechanisms have been categorized depending on the way they analyze the data sources to find affinities between users and items.

The most general setting in which recommender systems are studied is presented in Figure 2.1. Having a rating matrix of  $n$  users and  $m$  items representing the current user preferences, the task of a RS is to predict all missing ratings  $r_{a,i}$  for the active user  $a$ , and then recommend the item(s) with the highest rate[8]. However the user ratings matrix is typically sparse, as most users do not rate most items. Different approaches to solve this task can be categorized into the following types.

**Figure 2.1:** Typical user ratings matrix, where each cell  $r_{u,i}$  corresponds to the rating of user  $u$  for item  $i$ .

	Items					
	1	2	...	$i$	...	$m$
Users	1	5	3		2	1
	2	3	4	3	4	
	:		2			3
	$u$	5			1	5
	:		4	2		1
	$n$	4		1	1	5
	$a$	3	5		?	1

### 2.1.1 Collaborative Filtering

In Collaborative Filtering (CF) systems, items are recommended by exploiting the similarities amongst several users based on the feedback of previously consumed items. Usually, databases that fit into a CF approach, store the past users interactions in the form of explicit ratings (e.g., 1 to 5 range), or implicit ratings (e.g. a song played by a user, or an item bought by her). In general, CF methods are subdivided into: memory-based and model-based approaches.

- **Memory-based CF:** items are recommended using a weighted combination of ratings in a subset of users that are chosen based on their rating similarity to the target user. The most commonly used measure is the Pearson correlation coefficient[9]. Alternatively, the similarity between the ratings of two users (represented as a vector in an  $m$ -dimensional space) can be and computed using the Cosine similarity or the adjusted Cosine similarity measure which overcomes the issue of users with different rating behavior schemes[10]. Previous studies found that correlation[11] and adjusted cosine similarity[10] performs slightly

better.

Memory-based CF methods have been extended and improved over the years. Linden, Smith, and York[12] proposed an *Item-based* (or item-item) CF method to overcome the problem of dimensionality found when conventional memory-based algorithms cannot scale well when finding similarities between millions of users and items. This approach, which matches a user's rated items using Pearson correlation, lead to faster online systems and improved recommendations.

- **Model-based CF:** provides recommendations by using statistical models for predicting user ratings. The most widely used models are Bayesian classifiers, neural networks, fuzzy systems, genetic algorithms, latent features and matrix factorization[2]. For instance, CF methods can be turned into a classification problem, where a classifier is built for each active user representing items as features over users and available ratings as labels. However, latent factor and matrix factorization models have emerged as a state-of-the-art methodology in this class of techniques[13].

### 2.1.2 Content-based RS

Content-based (CB) RS (compared to pure CF that only utilizes the user rating matrix) tries to make a better personalized recommendation by exploiting the knowledge about a user (e.g. demographic information), or the properties of items (e.g. genre of a movie). Several approaches have treated this problem as an information retrieval (IR) task, where the content associated with the user's preferences is treated as a query, and the unrated items are scored with relevance/similarity to this query[14]. Alternatively, CB-RS has also been treated as a classification task using algorithms such as k-Nearest Neighbors (k-NN), decision trees, and neural networks[15]

### 2.1.3 Hybrid approaches

In order to leverage the strengths of both CB-RS and CF methods, several hybrid approaches have been proposed. For instance, Melville et al. presented in [16] a general framework for content-boosted collaborative filtering, where content-based predictions are applied to convert a sparse user ratings matrix into a full ratings matrix, and then a CF method is used to provide recommendations. In essence, this approach has been shown to perform better than pure CF, pure CB-RS, and a linear combination of the two.

### 2.1.4 The cold-start problem

The cold-start problem is a common issue that happens in recommender systems when it is not possible to make reliable recommendations due to an initial lack of ratings. There are three kinds of known cold-start problems: new community, new item and new user[2]. However, the latter represents one of the greatest difficulties faced by an RS in operation. Since new users have not yet provided any rating in the RS, they cannot receive any personalized recommendations when using memory-based CF. Usually, when the users enter their first ratings they expect the RS to offer them personalized recommendations, but there are not enough ratings yet to be able to make reliable predictions. As a consequence, new users feel that the RS does not offer the service they expected.

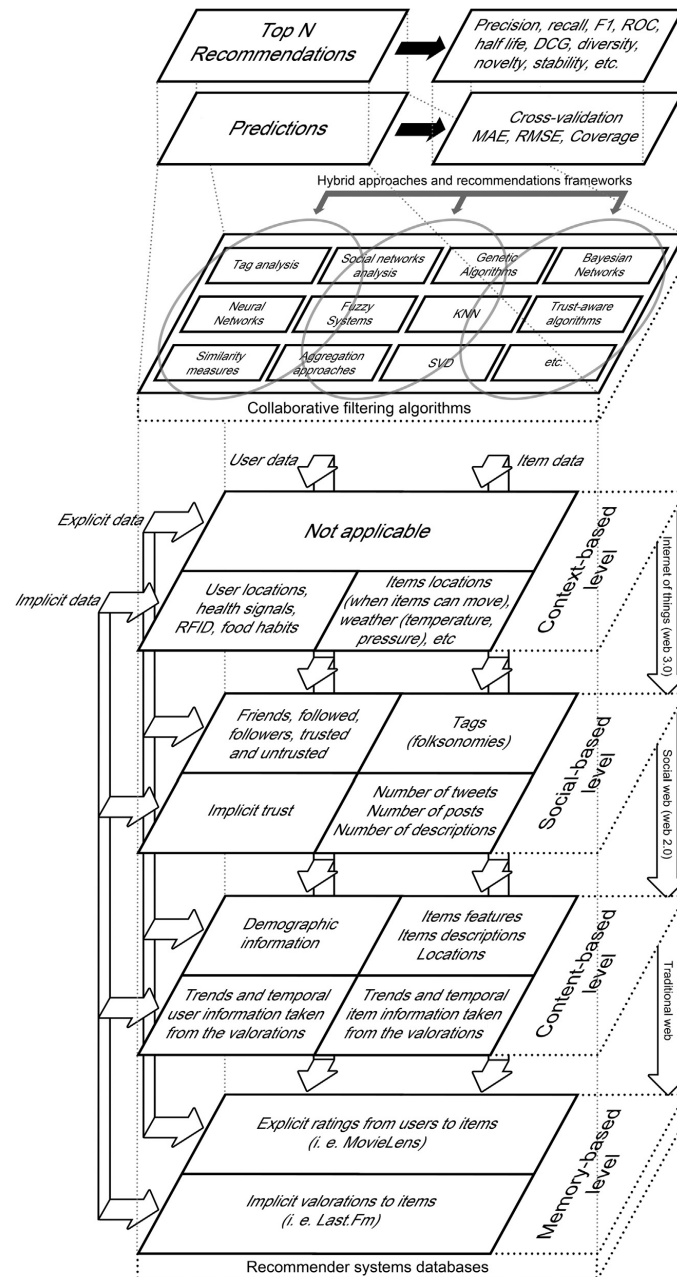
This problem is often faced using hybrid approaches such as CF-CB RS, CF-demographic based RS, or CF-social based RS[2]. However, these methods can be combined with clustering techniques over items in order to improve the prediction quality.

### 2.1.5 Trends in Recommendation Systems

Latest studies in CF showed that combining conceptual (explicit) and usage (implicit) information can improve the quality of web recommendations. Bobadilla et al.[2] presented a taxonomy for RS which unifies the current recommender methods and algorithms that can be applied to incorporate

memory-based, social and content-based information into their CF system depending on the type of information available. The taxonomy depicted in Figure 2.2, also detail at its higher levels the current evaluation methods for RS in terms of quality measures, diversity and novelty.

**Figure 2.2:** Recommender System taxonomy. Source: Recommender systems survey[2]



Moreover, Bobadilla et al. argue that one the most widely used algorithm

for CF is the k-nearest neighbors (kNN) model. In general, kNN generates recommendations by executing the following tasks: (1) determine k users neighbors; (2) implement an aggregation approach with the ratings for the neighborhood in items not rated by a; and (3) select the top N recommendations from predictions obtained in step 2.

On the other hand, a more robust model that combines the advantages of Support Vector Machines with factorization models was introduced by Rendle S. in [17]. The Factorization Machine (FM) model is able to compute all interactions between variables using factorized parameters even in scenarios with huge sparsity like recommender systems. Later, Rendle S. presented in [18] the *LibFM* software package that implements three learning methods have been proposed for FMs: Stochastic Gradient Descent (SGD)[17], Alternating Least-Squares (ALS)[19] and the Markov Chain Monte Carlo (MCMC)[20].

Most recent research like [5], made a generalization of recent advances in Deep Learning [6] of the Identically and Independently Distributed (i.i.d.) input and applied it to the non-i.i.d. (e.g. CF-based) input by presenting hierarchical Bayesian models that jointly performs deep representation learning for the content information and CF for the ratings matrix. Additionally, Wu Y. et al. in [21] proposed a Collaborative Deep Auto Encoders (CDAE) model which formulates the top-N recommendation problem using Denoising Auto-Encoders to learn the latent factors from corrupted inputs. Experiments carried out over different domains have shown that the two deep learning approaches can significantly outperform the state of the art. Therefore, recent research has started to focused on defining deep representational approaches that can lead to better recommendations.

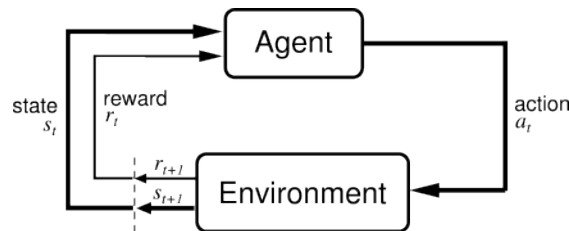
## 2.2 Reinforcement Learning

Reinforcement learning (RL) [22] is a family of machine learning algorithms that optimize sequential decision making processes based on scalar evaluations or rewards. Similarly to an n-armed bandit model[23], RL considers the

problem as a goal-directed agent interacting with an uncertain environment, where the main objective is to perform actions that maximize the expected sum of future reward for each state in the long term. The most important feature distinguishing RL from other types of learning is that it uses training information that evaluates the actions taken rather than instructs by giving correct actions (like supervised learning algorithms). In other words, an agent must be able to learn from its own experience, like a trial-and-error search.

RL uses a formal framework which defines the continuous interaction in terms of states, actions, and rewards, between the *agent* and an *environment*[3]. At each time step  $t$ , the agent receives some state representation of the environment  $s_t \in \mathcal{S}$ , where  $\mathcal{S}$  is the set of possible states. Then it selects an action  $a_t \in \mathcal{A}(s_t)$ , where  $\mathcal{A}(s_t)$  is the set of actions available in the observed state. One time step later, the agent receives a numerical reward  $r_{t+1} \in \mathcal{R}$ , and the environment turns into a new state  $s_{t+1}$ . Figure ?? shows the whole interaction in an agent-environment interface.

**Figure 2.3:** The Agent-Environment interface in a Reinforcement Learning problem.  
Source: Reinforcement Learning: a Survey[3]



This framework is intended to be a simple way of representing four essential features of the artificial intelligence problem: (1) a *policy* (commonly stochastic) defines a mapping from the perceived states of the environment to actions to be taken when in those states; (2) a *reward* function maps each state-action pair to a single number that represents how good or bad the new state is for the environment; (3) a *value* function specifies the total amount of reward an agent can expect to accumulate over the future and starting from a given state. This function is considered the most important as it is used by the agent during decision-making and planning; and optionally (4) a model that

mimics the behavior of the environment (transitions and rewards) and provides a way of deciding which action to perform considering possible future situations before they are actually experienced.

In general, agents are categorized depending on how the reinforcement learning problem needs to be modeled. *Value-based* and *Policy-based* agents are solely based on the value and policy functions respectively. *Actor-critic* agents use both policy and value functions; *Model-free* agents use policy and/or value function but no model; and *Model-based* agents have all properties mentioned above.

On the other hand, The *return*  $R_t$  is defined as the sum of the discounted future rewards over an episode of  $T$  time steps that an agent actually seeks to maximize:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \quad (2.1)$$

where  $\gamma, 0 \leq \gamma \leq 1$  is a discount factor that determines the present value of future rewards

### 2.2.1 Markov Decision Processes

A Markov Decision Process (MDP) is a model for sequential stochastic decision problems. As such, it is widely used in applications where an autonomous agent is influencing its surrounding environment through actions[4]. It is defined by a tuple  $\langle S, A, R, Pr \rangle$  representing a Markov Chain with values and decisions that follows the Markov property: *a state  $S$  is Markov if and only if  $\mathbb{P}[S_{t+1}|S_t] = \mathbb{P}[S_{t+1}|S_1, \dots, S_t]$ . In other words, in a MDP the future is independent of the past given the present.* Therefore, if a reinforcement learning task satisfies the Markov property and the state and action spaces are finite, then it can be modeled as a finite Markov Decision Process.

A particular finite MDP is defined by its state and action sets, a reward function  $R$  (equation 2.2) that assigns a real value to each state/action pair, and a state-transition function  $Pr$  (equation 2.3) which provides the probability of transitioning between every pair of states given each action. Altogether,



these quantities completely specify the most important aspects of the dynamics of a finite MDP.

$$\mathcal{R}_{ss'}^a = \mathbb{E}r_{t+1} | s_t = s, a_t = a, s_{t+1} = s' \quad (2.2)$$

$$\mathcal{P}_{ss'}^a = \text{Pr} s_{t+1} = s' | s_t = s, a_t = a \quad (2.3)$$

Therefore, the decision-maker's goal is to find an optimal policy  $\pi$ , such that at each stage of the decision process, the agent needs only to obtain the current state  $s$  and execute the action  $a = \pi(s)$ . Various exact and approximate algorithms have been proposed for estimating the optimal policy  $\pi^*$ , such as policy iteration, value iteration, Monte-Carlo learning, temporal difference, Q-learning, SARSA, etc[22][3].

### 2.2.2 A model-free Reinforcement Learning Setup

A reinforcement learning setup consists of an agent interacting with an environment  $E$  at discrete time steps. At each time step  $t$ , the agent receives an observation  $s_t$ , takes an action  $a_t$  and receives a reward  $r_t$ . Additionally, the environment  $E$  may be stochastic so it can be modeled as an MDP with a state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , an initial state distribution  $\rho(s_1)$ , transition dynamics  $\rho(s_{t+1} | s_t, a_t)$ , and reward function  $r(s_t, a_t)$ . On the other hand, the agent's behavior is defined by a policy  $\pi$ , which maps states to a probability distribution over the actions  $\pi : \mathcal{S} \rightarrow P(\mathcal{A})$ . Finally, the return from a state is defined as the sum of the discounted future reward  $R_t = \sum_{i=t}^T \gamma^{(i-t)} r(s_i, a_i)$ . As the return depends on the actions chosen and therefore on  $\pi$ , it may be stochastic.

The goal in reinforcement learning is to learn a policy which maximizes the expected return from a start distribution  $J = \mathbb{E}_{r_i, s_i \sim E, a_i \sim \pi} [R_1]$ . The action-value function is used in many RL algorithms and describes the expected return after taking an action  $a_t$  in state  $s_t$  and thereafter following policy  $\pi$ :

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_{i \geq t}, s_{i > t} \sim E, a_{i > t} \sim \pi} [R_t | s_t, a_t] \quad (2.4)$$

Many approaches in RL opt to represent the action-value function as a recursive relationship using the Bellman equation (eq. 2.5). Nevertheless, If the target policy is deterministic we can describe it as a function  $\mu : \mathcal{S} \leftarrow \mathcal{A}$  and avoid the inner expectation as shown in equation 2.6. As the expectation depends only on the environment, it is possible to learn  $Q^\mu$  off-policy, using transitions which are generated from a different stochastic behavior policy  $\mu$ .

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} [r(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1} \sim \pi} [Q^\pi(s_{t+1}, a_{t+1})]] \quad (2.5)$$

$$Q^\mu(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} [r(s_t, a_t) + \gamma Q^\mu(s_{t+1}, \mu_{t+1})] \quad (2.6)$$

Q-learning[24], is a commonly used off-policy algorithm which uses a *greedy* policy  $\mu(s) = \operatorname{argmax}_a Q(s, a)$  in order to estimate the action that gives the maximum reward. As the algorithm considers function approximators parameterized by  $\theta^Q$ , it can be used as an optimizer that minimize the loss:

$$L(\theta^Q) = \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \beta, r_t \sim E} [(Q(s_t, a_t | \theta^Q) - y_t)^2] \quad (2.7)$$

where  $y_t = \mathbb{E}_{s' \sim \mathcal{E}} [r(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, \mu(s_{t+1}) | \theta^Q)]$ .

### 2.2.3 Exploration/Exploitation Dilemma

One of the challenges that arise in reinforcement learning is the trade-off between exploration and exploitation[22][3]. To obtain a lot of reward, a reinforcement learning agent must prefer actions that has already tried in the past and found to be effective in terms of reward. In order to discover such patterns, it needs to try actions that it has not been selected before. so the agent has to be able to exploit actions that are already known as effective, but it also needs to explore in order to make better action selections in the future.

Naïve approaches use a greedy policy  $\pi(s) = \operatorname{argmax}_a Q(s, a)$  to get the action with maximum reward as a pure exploitative model. However, an  $\epsilon$ -greedy policy can be considered instead to allow exploration in the environment and improve the estimation of the non-greedy action values. Therefore, the agent

will select a random action with a probability  $\epsilon$ , and will act greedily with probability  $1 - \epsilon$ . The advantage of  $\epsilon$ -greedy policy over greedy policy is that the former continuously explore and improve the chances of recognizing possible actions that can lead to better rewards.

Nevertheless, this has still been an issue to solve in some reinforcement learning tasks, so research has continuously focused on trying different techniques to solve the exploration/exploitation trade-off that can lead to the discovery of better policies[22].

## 2.3 Recommender Systems and Reinforcement Learning

Although reinforcement learning seems to have great potential for improving the recommendation problem, it has received relatively little attention and found only limited application. The first experiments on using reinforcement learning techniques were focused on developing recommendation systems for web content using implicit data from server files which stored the navigational logs of the users throughout their contents.

Ten Hagen et al. [25] considered RL approach with an exploration/exploitation trade-off for discovering unknown areas and automatically improve their recommendation policy to helping users navigate through an adaptive web site. They stated that a model-free algorithm like Q-Learning can become infeasible if the number of actions are relatively high and showed that a greedy policy can indeed lead to a suboptimal recommendation model as the use of such kind of policy function, where the best actions (e.g with higher probability) are taken, does not always improve the policy. They finally demonstrated that this problem is solved by starting the algorithm with a lot of exploration and gradually reduce the parameter  $\epsilon$  as the policy is getting closer to an optimum. Results concluded that a recommender system without exploration potentially can get trapped into a local maxima.

Rojanavasut et al. in [26] presented a general framework for web recom-

mendation that learns directly from customer's past behavior by applying an RL process based on the SARSA method and a  $\epsilon$ -greedy policy. The system was composed of two models: a global model to keep track of customers trends as a whole, and a local model to record the user's individual browsing history. In general, the model treated pages as states of the system and links within a page as actions. To predict the next state, it used a ranking system that is also separated into 2 parts. i) a global ranking system using the data from a  $Q_{global}$ -matrix obtained from the action-value function and an  $\epsilon$ -greedy policy that gave the chance for rank new items that have few clicks but may match a user's interest; and ii) a local ranking  $Q_{local}$  using an inverse  $\epsilon$ -greedy policy. The system then finds the final total reward using  $Q_{total} = Q_{local} + wQ_{global}$  where  $w \in (0 - 1]$  is a weight hyper parameter of the model.

Overall, the system provided customers the chance to explore other products than the ones they already visited. Experimental results showed that a value of  $\epsilon = 0.2$  preserves the balance between exploration and exploitation, meanwhile If  $\epsilon < 0.2$  the system will gave small chances to users to explore new items, or otherwise may include products that do not match to the customer's interest ( $\epsilon > 0.2$ ). On the other hand, results also prove that even if the purpose of  $Q_{local}$  was to discover new products, it appeared to be less effective than  $Q_{global}$ .

Shani et al. in [4] argued that it is more appropriate to formulate the problem of generating recommendations as a sequential optimization problem so they described a novel approach for a commercial web site based on MDPs together with a predictive model. The MDP-based recommender system took into account the expected utility and the long-time effect of a particular recommendation. Then it suggested items whose immediate reward is lower, but leads to more *profitable* rewards in the future. However, the benefits of this model are offset by the fact that the model parameters are unknown, randomly initialized and that they take considerable time to converge. So, they defined a strong initial model-based for collaborative filtering, which solves quickly,

and that does not consume too much memory.

Before implementing the recommender, they initialize a predictive model of user behavior using data extracted from the web site, and then used it to provide the initial parameters for the MDP. They proposed to use maximum-likelihood estimation with three major enhancements: skipping, clustering, mixture modeling, so the predictive model can be thought of as a first-order Markov chain (MC) of user dynamics in which states correspond to sequences of  $k$  events (in this case: previous selections) representing relevant information about the users interaction. Then, the transition function described the probability that a user ,whose  $k$  recent selections were  $x_1, \dots, x_k$  will select the item  $x'$  next.

Authors experimented with small values of  $k$  ( $k \in [1 - 5]$ ) in order to overcome data sparsity and MDP solution complexity issues. Moreover, enhancements on the maximum-likelihood n-gram formulation showed to be useful to find a good estimate of the transition function for the initial predictive model without suffering the problem of data sparsity and bad performance presented on other model variations.

On the other hand, the components of the reinforcement learning model were defined as follows: a) the states are the  $k$ -tuples of items purchased; b) actions correspond to a recommendation of an item; c) rewards depends only on the last item defining the current state, where its net profit was used; and d) the transition function as the stochastic element of the model that estimates the user's actual choice.

In order to solve the MDP problem, a policy iteration algorithm was used as the defined state space presented certain characteristics that lead to fast convergence. The transition function was carefully initialized in order to be fairly accurate when the system was first deployed to avoid the cold-start problem, causing states that are infrequently observed to be updated faster than already observed states. Moreover, when the first transition to a state is observed, its probability is initialized to 0.9. In this way, the system balances

the need to explore unobserved items in order to improve its model by recommending non-optimal items occasionally until getting their actual counts. Finally, to update the model, the system used an off-line approach which keeps track of the recommendations and the user selections and build a new model at fixed time intervals (e.g. once a week).

Experiments demonstrated that a deployed system using three mixture components (combined using an equal weight), with history length  $k \in [1, 3]$  was able to generate 28% higher average user reward compared to pure MC model, while in terms of computational costs, the MDP-based model was built quickly and provided fastest recommendations at the price of more memory use. All in all, the off-line predictive model approach provided an adequate initial performance that overcomes the cold-start problem, however, authors avoid using some form of reinforcement learning technique as they argued that at that time, its implementation requires many calls and computations by a recommender system online, which lead to slower responses and undesirable results for the web site owner.

A machine learning perspective, introduced by Taghipour et al. in [27], used an off-line reinforcement learning model to solve the recommendation problem using the Q-Learning algorithm while employing concepts and techniques commonly applied in the web usage mining domain. They argued that this method is appropriate for the nature of web page recommendation problem as it provides a mechanism which is constantly learning, does not need periodic updates, can be easily adapted to changes in the website structure and new trends in users behavior.

The RL problem formulation was considered as as a competition between different recommender systems to gather more points, together with a stochastic model with the following properties:

- *states*: showing the history of pages visited by the user so far. For this case, a notion of N-Grams was adopted and a sliding window of size  $w$  was set to limit the page visit sequences to a constant number and avoid

large state spaces.

- *actions*: consisting of a single page recommendation at each state.
- *policy*: rewards actions positively if it recommends a page that will be visited in one of the consequent states
- *reward function*: defined as  $r(s, a) = \text{reward}(\text{Dist}(Rs', p), \text{Time}(p_w^v))$ , where  $\text{Dist}(Rs', p)$  is the distance of page  $p$  from the end of the recommended pages list to state  $s'$ , and  $\text{Time}(p_w^v)$  indicates the time user has spent on the last page of the state. Finally, the  $\text{reward}(\text{Dist}, \text{Time})$  function was defined as a linear combination of both values:  $\text{reward}(\text{Dist}, \text{Time}) = \alpha \times \text{dist} + \beta \times \text{Time}$  with  $\alpha + \beta = 1$ .

The proposed off-line RL method used a simplified formulation of the recommendation problem but it obtained much better results compared with two baselines methods: association rules and item-based collaborative filtering (with probabilistic similarity measure). As the coverage increase, naturally accuracy decrease in all systems, but the RL approach outperformed the other two systems displaying a lower rate in which its accuracy decreased. Authors concluded that the algorithm is a good candidate for solving the recommendation problem as it does not rely on any previous assumptions regarding the probability distribution of visiting a page after having visited a sequence of pages, and that the nature of the problem matches perfectly with the notion of delayed reward (known as temporal difference). Additionally, they suggested that in order to produce better results it could be possible to use a more complicated formulation of the reward function such as a neural networks, rather than a linear combination of factors.

Later in [28] they exploited the previous RL framework and present a hybrid web recommendation method enriched with semantic knowledge about the usage behavior and thus obtain a more generalized solution regarding to the usage data it has. The new system used the incremental Document Conceptual Clustering[29] algorithm to map pages to higher level concepts, and

thus exploit the hierarchical and conceptual document clustering to provide a semantic relationship between them in combination with the usage data in a user session.

Therefore, new states consist of a sequence of concepts visited by the user, while actions are recommendation of pages that belong to a specific concept. This definition resulted in a much smaller state-action space as the size of the state space is now dependent on the number of distinct page clusters. Finally, the new reward function takes into account the content similarity of the recommended and visited pages along with the usage-based reward defined in the previous approach.

The modified algorithm do not make predictions based on weak usage patterns as the new states represented a generalized view of many single visit sequences. Furthermore, evaluation results showed the flexibility of the new RL approach to incorporate different sources of information in order to improve the quality of recommendations.

A model-based RL agent for music playlist recommendation proposed by Liebman et al. [30] modeled the preferences of both songs and song transitions as MDPs, and demonstrated to be potentially effective in the domain and particularly good to solve the cold-start problem as it is able to generate personalized song sequences within a single listening session and with no prior knowledge of the new user's preferences.

The adaptive playlist generation problem was defined as an episodic MDP  $\langle S, A, P, R, T \rangle$  with the following components:

- a state space  $S$  composed by the ordered sequence of songs played,  $S = \{(a_1, a_2, \dots, a_i) | 1 \leq i \leq k; \forall j \leq i, a_j \in \mathcal{M}\}$
- an actions space  $A$  formed by all possible candidates for being the next song to play,  $a_k \in A$  which means  $A = \mathcal{M}$
- a deterministic transition function  $P(s, a) = s'$  which indicates the existing transitions from state  $s$  to  $s'$  by taking action  $a$



- an utility function  $R(s, a)$  derived from hearing song  $a$  when in state  $s$
- $T = \{(a_1, a_2, \dots, m_k)\}$ : the set of playlists of length  $k$ .

In order to find an optimal policy  $\pi^*$  (and thus obtain the most pleasing sequence of songs to the listener), a model-based approach was chosen arguing that even though model-free approaches learn the value of taking an action  $a$  from state  $s$  directly, it requires a lot of data to converge, and it is considerably scarce in the domain music recommendation (as well as in other kind of applications). On the other hand, model-based formulations often requires a lot of computation to find an approximate solution to the MDP, but the trade-off of computational expense for data efficiency makes the model-based approach a good option for this problem.

Since  $P$  is deterministic, only the listener's reward function  $R$  was required to be modeled. Therefore, the reward function defined as  $R(s, a) = R_s(a) + R_t(s, a)$  represented the sum of two distinct components: (1) the listener's preference over songs,  $R_s : A \rightarrow \mathcal{R}$ , and (2) the preference over transitions from songs played to a new song,  $R_t : S \times A \rightarrow \mathcal{R}$ .

In essence, the model represents each song as a compacted vector of spectral auditory descriptors that capture meaningful differences in user's preferences, so the system is able to leverage knowledge using only a few transition examples to plan a future sequence of songs. On the other hand, the agent architecture was composed by a module for learning the listener parameters during initialization and learning on the fly processes; and an additional module for planning a sequence of songs in the playlist.

The initialization is divided in two parts: (1) initialization of song preferences polls the listener for her  $k_s$  favorite songs in the database and updates the user's preference vector using the feature descriptors of each of the selected items; and 2) initialization of the transition preferences by presenting different possible transitions that encapsulate the variety in the dataset, and directly asking which of a possible set of options the listener would prefer. Updates are carried out in the same manner as the initialization of song pref-

erences.

After initialization, the agent begins to play songs for the listener, waiting for her feedback (reward), and updating the user’s preferences accordingly. This learning of the fly procedure is perceived as a temporal-difference update with an attenuating learning rate that balances the trust between the previous history of observations and the newly obtained signal. Then, a Monte Carlo Tree Search (MCTS) heuristic[31] is used for planning. The iterative process chooses a subset of 50% of the songs in the database with highest  $R_s$  score, and then at each point, it simulates a trajectory of future songs selected at random and calculate the expected payoff of the song trajectory using  $R_s$  and  $R_t$ . The process ends when it finds the trajectory which yields to the highest expected payoff, and finally, the first item in the trajectory is selected to be the next song to recommend.

Additionally, to mitigate the problem of high complexity during re-planning under large song spaces, the agent used the canonical K-means algorithm for clustering songs according to song types and reduce the search complexity drastically.

An evaluation of the cumulative reward distribution of the proposed approach showed that overall, the DJ-MC agent outperforms two baselines models (a random agent and a greedy agent), while in terms of transition reward preferences, the new system presented a small but significant boost in performance compared to a model based only on reasoning about song preferences. All in all, this approach demonstrated to be a good improvement as a reinforcement learning model for music recommendation, as it enables learning from relatively few examples, and provides quality on recommendation sequences.

## 2.4 Deep Reinforcement Learning

So far, the RL agents presented above have shown to achieve some success under the recommendation domain, nevertheless, their performance has been conditioned to the quality of the hand-crafted feature engineering made, as

well as the custom definition of value functions or policy representations. Moreover, to use RL successfully in situations approaching to real-world complexity, agents usually need to derive efficient representations of the environment from high-dimensional sensory inputs, and use these to generalise past experiences to new situations.

Recent advances in deep learning have made it possible to extract high-level features from raw sensory inputs[6]. Even if it seems that this approach could fit in a RL task, its applicability to the RL domain may present several challenges from a deep learning perspective, For instance, successful deep learning applications have required large amounts of hand-labelled training data to obtain good predictions results. Additionally, most deep learning algorithms assume the data samples to be i.i.d., while in RL an agent typically encounters with non-i.i.d. sequences of highly correlated states.

On the other hand, agents must be able to learn from a scalar reward signal that is frequently sparse, noisy and delayed, while the RL data distribution changes as the algorithm learns new behaviour. However, thanks to the increasing success of adapting different deep learning techniques to different domains, Deep Reinforcement Learning (DRL) models have started to be proposed, getting promising results.

Mnih et al. presented in [32] [33] the first deep learning model which uses an end-to-end reinforcement learning approach to learn control policies directly from a high-dimensional raw input space. The model, applied to a range of Atari 2600 games, consists of a convolutional neural network and an experience replay mechanism[34] to alleviate the problems of correlated data and non-stationary distributions in typical RL problems.

Moreover, the model's architecture uses only the state representation as input, and a separate output unit for each possible action (target network), corresponding to the predicted Q-values of each individual action that can be applied to the input state. This configuration allows the computation of all possible actions reward in a given state with only a single forward pass.

Finally, the resulting Deep Q neural network (DQN) is trained with a variant of Q-learning which learns from raw pixels input and the underlying environment properties, and outputs a value function estimating the future rewards.

On the other hand, the RL formulation models the environment  $\mathcal{E}$  as a MDP composed by a state space  $\mathcal{S}$ , a discrete action space  $\mathcal{A} = \{1 \dots K\}$ , a scalar reward function  $r(s_t, a_t)$  and a transition dynamics function  $\rho(s_{t+1}|s_t, a_t)$ . Then, given the agent's behaviour defined by a policy  $\pi$ , they estimate the action-value function by using an approximator function of the Bellman equation (presented in eq. 2.5) as an iterative update. A Q-network function approximator with weights  $\theta$  is then trained by minimising the sequence of loss functions  $L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot)}[(y_i - Q(s, a; \theta_i))^2]$  where  $\rho(s, a)$  is the behaviour probability distribution over sequences of states  $s$  and actions  $a$ .

Finally, to compute the full expectations in the gradient of the loss function, they considered stochastic gradient mini-batch updates of the weight parameters, and a learning algorithm which is: a) *model-free*: replaces the expectations by directly using uniform samples from the behaviour distribution  $\rho$ , and b) *off-policy*: follows an  $\epsilon$ -greedy strategy of the behaviour distribution that ensures an adequate exploration of the state space.

Evaluation results showed that this method is able to learn how the value function evolves for a reasonably complex sequence of events. Furthermore, the proposed DRL algorithm demonstrated to perform better compared to the SARSA and Contingency[35] methods as they need to incorporate significant (handcrafted) prior knowledge about the visual problem to get considerable performance, in comparison to the raw inputs that DRL use. Also, DRL achieved better performance than an expert human player in 29 out of 49 games.

However, this model is only able to handle discrete low-dimensional action spaces and cannot be directly applied to tasks under the continuous domain as the iterative process of finding the action that maximises the action-value function becomes intractable. Additionally, a naïve discretization of the action

space would not be a good solution as it may throw away information about the structure of the action domain.

Lillicrap et al. [36] adapted the ideas presented above to the continuous action space and defined an actor-critic, model-free approach, based on the Deterministic Policy Gradient[37] (DPG) algorithm, that is able to find policies end-to-end and whose performance is competitive compared to approaches based on a planning algorithm with access to the dynamics of the domain.

While the model's architecture uses similar features of DQN (network architecture, experience replay memory, and target outputs), along with batch normalization[38] to overcome the internal covariate shift problem in the network input layers during mini-batches, the Deep Deterministic Policy Gradient (DDPG) algorithm maintains two functions: (1) an actor function  $\mu(s|\theta_\mu)$  that obtains the current policy by deterministically mapping states to specific actions; and (2) the critic function  $Q(s, a)$  learned by applying the Bellman equation. During the learning process, DDPG updates the actor by applying the policy gradient to the expected return from a start distribution  $J = \mathbb{E}_{r_i, s_i \sim E, a_i \sim \pi}[R_1]$  and with respect to the actor parameters. Even if convergence is no longer guaranteed due to the added non-linearity of the proposed approximation function, it produces a good generalization over large state spaces.

The DDPG algorithm also implements certain improvements to the RL task. First, it introduces stability in the learning process by applying soft target updates to the target networks instead of directly copying the weights. This is performed by exponentially moving their average values using  $\theta' \leftarrow \tau\theta + (1 - \tau)\theta'$  with  $\tau \ll 1$ . Second, batch normalization allows DDPG to be able to generalize across environments using the same hyper parameters for learning. Finally, the exploration-exploitation trade-off is managed independently from the learning algorithm by adding noise sampled from a noise process  $\mathcal{N}$  to the actor policy  $\mu$ .

To evaluate the algorithm performance, authors used simulated phys-

ical environments with different levels of difficulty, and the Ornstein-Uhlenbeck[39] sampling process for adding exploration to the learning policy. Results demonstrated that even when harder tasks obtain poor Q estimates, DDPG is able to learn good policies across a variety of domains with continuous action spaces and using both low-dimensional feature vector and high-dimensional pixel inputs. Additionally, all the experiments were solved using fewer steps of experience than the DQN algorithm, showing that DDPG may be able to solve even more difficult problems.

However, the DDPG algorithm requires of a large number of training episodes to find solutions and only works under environments with a not so big number of actions. Therefore, a more robust model-free approach is needed in order to tackle these limitation.

Later, Dulac-Arnold et al. [7] presented a new policy architecture, under the same actor-critic framework used above, that not only allows reinforcement learning methods to be applied to large-scale learning problems, and to operate efficiently with a large number of actions, but also can generalize over the action set in logarithmic time. The full policy is then trained and optimized using DDPG. As a result, they obtained a more efficient algorithm that makes both learning and acting tractable in time and allows value-based policies to use the action features to reason about previously unseen actions.

The so-called *Wolpertinger* architecture leverages prior information about the actions and embed them in a continuous space upon which the actor can generalize using a smooth function approximator. After the policy produces a continuous action within this space, it uses an approximate nearest neighbour search to find the set of closest discrete actions in logarithmic time to finally select the action with the highest reward. Consequently, the algorithm avoids the heavy cost of evaluating all actions mainly by defining an efficient action-generating actor, and then using the critic to refine the actor's choices for the full policy.

Having the definition of the full policy as  $\pi_{\theta\pi} = g \circ f_{\theta\pi}(s)$ , where  $f_{\theta\pi}(s) = \hat{\mathbf{a}}$

is the proto-action generation function and  $g$  the actual action function that return the set of  $k$  closest actions to  $\hat{a}$ , the final goal of this approach is to perform policy iteration by alternatively performing policy evaluation on the current policy with Q-learning, and then improving upon the current policy by following the DDPG over  $f_{\theta^\pi}$ , considering that the effects of  $g$  are a deterministic aspect of the environment.

The Wolpertinger agent, evaluated on three environment classes with large actions spaces (including a simulated recommender system), showed that it can scale to real-world MDPs with large number of actions as only a subset of the full set of actions was sufficient in many tasks to converge and provide significant speedups. Nevertheless, there is still an existing issue of exploration when the agent needs to learn from scratch.





## **Chapter 3**

# **Models**

This chapter will develop a reinforcement learning agent definition to solve the recommendation problem using deep learning techniques. After describing the main motivation behind this approach, details of the baseline model derived from previous work on the DDPG algorithm is presented, and finally, as the contribution to this MSc. project, a bayesian personalised ranking model based on factorization machines is introduced into the deep learning algorithm to try to enhance the policy architecture and therefore the system recommendations.

### **3.1 Motivation**

Previously, reinforcement learning agents have received little attention to solve the recommendation tasks due to their action domain becomes intractable in current large-scale recommender systems. However, the great performance reached in several control tasks by end-to-end RL agents based on deep neural networks[33] have demonstrated the capacity to continuously adapt their behaviour without any human intervention and solve challenging human tasks. More interestingly, a single trained general-purpose deep reinforcement agent (using the same hyperparameters) was able to converge and perform well on several task environments with different action space definitions. Consequently, a new opened research area has been using recent breakthroughs in training deep neural networks to derive RL agents models

that can adapt to other kind of existing tasks under both discrete and continuous domains and solve tasks that were previously considered intractable.

Another reason that motivates our approach is the sequential nature of the recommendation process. In general, recommender systems usually work in a sequential manner: once the RS suggest items to the user, she can accept one of the recommendations. At the next stage a new list of recommended items is calculated based on the user feedback and then presented again to the user. This sequential nature allows us to reformulate the recommendation process as a sequential optimization process where optimal recommendations can depend either on the previous items purchased, and in the order in which those items were purchased. For example, Zimdars et al. [40] suggested the use of a k-order Markov chain model (with  $k=3$ ) to represent this behaviour by dividing a sequence of transactions  $X_1, \dots, X_T$  into cases  $(X_{t-k}, \dots, X_{t-1}, X_t)$  for  $t = 1, \dots, T$ . and then, build a model to predict the item at the time step  $t$  given the other columns.

As now deep reinforcement learning approaches can converge to good estimations under tasks considered intractable before, it motivates us to evaluate and improve these models under recommendation tasks, which this thesis project is focused.

## 3.2 Baseline Agent Definition

Therefore, the recommender problem can be considered as a sequence of Markov Decision Processes  $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma \rangle$ , where the action space  $\mathcal{A}$  is defined as the set of items to recommend,  $\mathcal{S}$  represent the current state or item a user is currently consuming, thus both  $\mathcal{S}$  and  $\mathcal{A}$  represent items as feature vectors  $\mathbf{a}$  and  $\mathbf{s} \in \mathbb{R}^n$  respectively.  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the transition probability distribution to move from state  $s_t$  to  $s_{t+1}$ ,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward value function  $\mathcal{R}$  associated with an item accepted by the user, and  $\gamma \in [0, 1]$  is the discount factor for future rewards.

The reinforcement learning agent is then defined as a model-free approach

with the main objective to find the optimal policy  $\pi$  that maximizes the expected total future reward  $\mathbb{E}[R_1]$  over all episodes  $R_t = \sum_{i=t}^T \gamma^{i-t} r(s_i, a_i)$ . The expectation can be obtained using the Bellman equation (presented in section 2.2.2):

$$Q^\pi(s, \mathbf{a}) = r(s, \mathbf{a}) + \gamma \sum_{s'} P(s'|s, \mathbf{a}) Q^\pi(s', \pi(s')) \quad (3.1)$$

In order to approximate  $\pi$  and  $Q$  under large-scale action spaces, standard actor-critic approaches need to define a policy by a parametrized actor function approximator  $\pi_\theta = \arg \max_{\mathbf{a} \in \mathcal{A}} Q(s, \mathbf{a}) : \mathcal{S} \rightarrow \mathcal{A}$  that would be able to generalize and that does not scales linearly in relation to the number of actions (turning the task into an intractable problem). Based on recent work presented by Lillicrap et al. [11] and Dulac-Arnold et al. [7], it is now possible to define generalizable reinforcement agents using a policy architecture with a sub-linear complexity relative to the action space that avoids the heavy cost of evaluating all actions. This actor-critic framework uses a multi-layer neural network as function approximators and is trained using Deep Deterministic Policy Gradient (DDPG).

The action generation function is based on the *Wolpertinger* policy defined in [7]. First, an initial function  $f_{\theta^\pi}(s) = \hat{\mathbf{a}}$  provides a proto-action in  $\mathbb{R}^n$  for a given state  $s$ .  $\hat{\mathbf{a}}$  is then mapped to the discrete action set  $\mathcal{A}$  by applying the k-nearest-neighbor<sup>1</sup> function from equation 3.2, which returns the k closest actions in  $\mathcal{A}$  by computing the  $L_2$  distance between them.

$$g_k(\hat{\mathbf{a}}) = \arg \min_{\mathbf{a} \in \mathcal{A}}^k |\mathbf{a} - \hat{\mathbf{a}}|_2 \quad (3.2)$$

Even if  $g_k$  has the same complexity as the *argmax* in the Q action-value function, each step of evaluation of the  $L_2$  distance is faster than a full value-function evaluation as the lookup is performed in logarithmic time. Finally, the choice of the actual action to be applied to the environment is set according to the highest-scoring action obtained according  $Q_{\theta^Q}$  defined in equation 3.3. As

---

<sup>1</sup>The size of the generated action set k is task specific, and allows for an explicit trade-off between policy quality and speed

a result, the full Wolpertinger policy makes the algorithm significantly more robust to imperfections in the choice of action representation.

$$\pi_\theta(s) = \arg \max_{a \in g_k \circ f_{\theta^\pi}(s)} Q_{\theta^Q}(s, a) \quad (3.3)$$

On the other hand, the DDPG implementation maintains two functions: (1) the actor function  $\mu(s|\theta_\mu)$  that obtains the current policy by deterministically mapping states to specific actions; and (2) the critic function  $Q(s, a)$  learned by applying the Bellman equation. It then updates the actor by applying the chain rule to the expected return from a start distribution  $J = \mathbb{E}_{r_i, s_i \sim E, a_i \sim \pi}[R_1]$  and with respect to the actor parameters:

$$\begin{aligned} \nabla_{\theta^\mu} J &\approx \mathbb{E}_{s_t \sim \rho^\beta} [\nabla_{\theta^\mu} Q(s, a|\theta^Q)|_{s=s_t, a=\mu(s_t|\theta^\mu)}] \\ &= \mathbb{E}_{s_t \sim \rho^\beta} [\nabla_a Q(s, a|\theta^Q)|_{s=s_t, a=\mu(s_t)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s=s_t}] \end{aligned} \quad (3.4)$$

In order to maintain stability in the learning algorithm, a separate target networks  $\mu'(s|\theta_{\mu'})$  and  $Q'(s, a|\theta_{Q'})$  for the actor and critic are used for calculating the target values, where their weight parameters are updated by using an smooth function update  $\theta' \leftarrow \tau\theta + (1 - \tau)\theta'$  with  $\tau \ll 1$  to main an exponential moving average of the trained parameters  $\theta^\mu$  and  $\theta^Q$ .

Among the main advantages of using this deep learning approach, DDPG introduces an experience replay buffer of size N that helps to reduce the correlation among the sequence of samples used by the RL agent. The mechanism allows the behaviour distribution to be averaged over many of its previous states, smoothing out learning and avoiding instability or divergence in the parameters.

When executing the DDPG algorithm, the critic is trained from samples generated and stored in the replay buffer by the full policy, while the policy gradient  $\nabla_a Q_{\theta^Q}(s, a)$  is taken at the true output  $\hat{a} = f_{\theta^\pi}(s)$ . The target action in the Q-update is generated by the full policy.

### 3.3 FM Agent Model

Our implementation attempts to replace the k-nearest neighbour computation in the Wolpertinger policy for a Bayesian Personalized Ranking (BPR) model using Factorization Machines to find the top-K items from a given state to then select the item that yields to the highest Q value in the long term. The deep reinforcement agent will finally learn using the same DDPG architecture as the baseline model.

The Bayesian formulation presented in Rendle et al. in [41] find the correct personalized ranking for all items  $i \in I$  by maximizing the posterior probability  $p(\theta | >_u) \approx p(>_u | \theta)p(\theta)$  where  $\theta$  represents the parameter vector of an arbitrary model class (e.g. matrix factorization) and  $>_u$  is the desired but latent preference structure for user  $u$ . In order to guarantee a personalized order, totality, antisymmetry and transitivity properties are fulfilled by defining the individual probability that a user really prefers item  $i$  to  $j$  as  $p(i >_u j | \theta) := \sigma(\hat{x}_{uij}(\theta))$  where  $\sigma$  is the sigmoid function and  $\hat{x}$  is an arbitrary real-valued function of the model parameter vector which captures the special relationship between user  $u$ , item  $i$  and item  $j$ , such as Matrix Factorization of adaptive kNN.

Rendle S. showed in [18] that factorization models, such as the BPR model for top-K recommendation, can be expressed using Factorization Machines, where vectors  $x$  are ordered by the score of  $\hat{y}(x)$  and optimization is done over pairs of instance vectors  $(x_i, x_j) \in \mathcal{D}$ .

#### 3.3.1 Proposed Approach

The new policy architecture called *FM-policy* tries to improve how an agent selects the action with the highest Q value by evaluating the top-k returned items from a given state (i.e. the item returned by the actor function and last item consumed by the user) by a pre-trained Factorization Machine model, rather than obtaining the k-nearest neighbors using the euclidean distance. As a result, the FM-policy maintains the sublinear complexity of the Wolpertinger policy but it introduces an stronger (factorization) model. The FM-policy algo-

**Algorithm 1:** FM policy

---

State  $s$  previously received from environment.  
 $\hat{\mathbf{a}} = f_{\theta^\pi}(s)$  Receive proto-action from actor.  
 $\mathcal{A}_k = g_k(\hat{\mathbf{a}}, s)$  Retrieve Top-k recommendation using BPR-FM-model  
 $\mathbf{a}_{TOP} = \arg \max_{\mathbf{a}_j \in \mathcal{A}_k} Q_{\theta^Q}(s, \mathbf{a}_j)$   
 Apply  $\mathbf{a}_{TOP}$  to environment; receive  $r, s'$

---

rithm is described in Algorithm 1.

After the actor function provides a proto-action  $\hat{\mathbf{a}}$  in  $\mathbb{R}$  given the current state or last consumed item by the user, a pre-trained BPR-FM model obtains the the Top-K recommendations based on the current observation:  $(\hat{\mathbf{a}}, s)$ . Then, the item in  $\mathcal{A}_k$  that yields to the highest reward is applied to the environment.

Similarly to the baseline model, when executing the DDPG algorithm, the critic is trained from samples generated and stored in the replay buffer by the FM-policy, while the policy gradient is taken at the true output in the actor function  $f_{\theta^\pi}(s)$ . As a reference, the pseudocode of the DDPG algorithm is included in Appendix A.

## Chapter 4

# Experiments

In the previous chapter, deep reinforcement learning agents were described to solve the recommendation task. Now, in order to test and evaluate their performance, a simulated environment for Recommender Systems is going to be defined. Then, after defining the baseline models and the experimental settings, two sets of experiments are carried out: (i) a performance evaluation of the proposed deep learning models and a pure Collaborative Filtering baseline model; and (ii) an analysis of the convergence and speed-up between the deep reinforcement learning agent based on previous work in the area, and an DRL agent with the new policy approach. Final conclusions on the experimental results obtained will be presented in the next chapter.

### 4.1 Reinforcement Learning Environment

For the purpose to demonstrate how the deep reinforcement learning agent perform on the recommendation task, a simulated environment is modelled under the OpenAI Gym<sup>1</sup> reinforcement learning toolkit. It that allow us to define Partially Observed MDP environments (POMDPs) under an episodic setting where the agent's experience is broken down into a series of episodes [42]. In general, OpenAI Gym manages two core concepts: (1) a free-to-code *agent* that maximizes the convenience for users allowing them to implement different styles of agent interface; and (2) a common *environment* and

---

<sup>1</sup><https://gym.openai.com/>

action/observation interfaces that ease the implementation and testing of different reinforcement learning problems under the same framework. The final goal in OpenAI Gym is to maximize the expectation of total reward per episode, and to achieve a high level of performance in as few episodes as possible.

During each episode, the agent’s initial state is randomly sampled from a distribution, and the interaction proceeds until the environment reaches a terminal state. After performing an action step, the environment interface returns the observation of current state, the reward achieved by the action performed, a done flag indicating if an episode ends and an information object containing useful information for debugging and learning. Finally, the framework also allow us to measure the performance of a RL algorithm under an environment along two axes: *final performance* or average reward per episode, after learning is complete; and, *sample complexity* or the amount of time it takes to learn.

### 4.1.1 Recommender System Environment

The recommender system environment, characterized by a set of  $n$  items to recommend to  $m$  users, is defined (based on a previous definition in [7]) as a sequence of MDPs  $\langle \mathcal{A}, \mathcal{S}, r, \mathcal{P} \rangle$  composed by an action set  $\mathcal{A}$  that correspond to set of items to recommend, a state space  $\mathcal{S}$  holding the item the user is currently consuming, a reward  $r$  defined as the rating value given by a user to an item if she accepts it, and a transition probability matrix  $\mathcal{W}$  which defines the probability of that a user will accept item  $j$  given that the last item she accepted was item  $i$ .

The transition probability matrix is generated using the ideas exposed by Yildirim et al. in [10] to build a random walk recommender system using a Markov Chain model, where the probability of being in a state only depends on the previous step  $Pr(X_{u,k+1} = i | X_{u,k})$ . The algorithm presented in Appendix B details the underlying process to set up the environment.



### 4.1.2 Simulation algorithm

The implemented algorithm is presented in Appendix C and describes the complete simulation process of the recommendation task. At each time-step, the agent presents an item  $i$  to the user with action  $\mathcal{A}_i$ . The user then accepts the item according to the transition probability matrix  $\mathcal{P}$  or selects a random item. To simulate the user patience in a session as in [7], an episode of learning ends with probability 0.1 if the user accepts the item, and with probability 0.2 if a random item is selected instead. Finally, after each episode the environment is reset by selecting a random item from a likely subset for an specific user provided by a single item-based collaborative filtering model based on the similarity between items.

## 4.2 Data Set

The 100-k and 1M Movielens datasets [43] from the GroupLens Research Lab<sup>2</sup> are used to train and test the implemented deep reinforcement agents. Table 4.1 summarizes the statistical properties of each dataset. The ratings files were pre-processed in order to generate the user rating matrix and the transition probability matrix. Additionally, train and test sets were generated to apply a five-fold cross validation to the baseline CF ratings predictor in the simulated environment.

	100K Movielens	1M Movielens
Total ratings	100,000	1,000,029
Users	943	6,040
Movies	1682	3,883
Density	0.63	0.42
Mean ratings/user	106	165
Min. ratings/user	20	20
Max. ratings/user	737	2314

**Table 4.1:** Properties of the 100-k and 1M Movielens datasets

---

<sup>2</sup><http://grouplens.org/datasets/movielens/>

### 4.3 Evaluation Scheme

For each model proposed below, we measure the average return obtained by the learning agent over steps, as well as the average precision along user sessions or episodes of learning. Finally, a quantitative and qualitative analysis of the results of each model are compared and contrasted against the performance of a random policy which acts randomly by sampling the next action to perform.

### 4.4 Baselines and Experimental Settings

To evaluate the strength and accuracy of the deep reinforcement learning model under the recommender system environment, we use two variants of the baseline model described in section 3.2, and a model using the new policy architecture. These models are:

- **DRL-kNN-CB:** the *DRL using Content-based item similarity* model uses the DDPG algorithm with a bag-of-words feature vectors to represent actions and states. The k-nearest neighbours index is then created using the KDtree algorithm [?] (with 30 leaf nodes checked) to find the k-closest set of similar items  $\mathcal{A}_k$  from the action estimated by the actor function. Finally, the Wolpertinger policy applies the action in  $\mathcal{A}_k$  that yields to the maximum return.
- **DRL-kNN-CF:** the *DRL with Item-based similarity* model uses the same policy and learning algorithm as the model described above but it implements an item-to item version of the k-NN algorithm by using the item vector representation from the user rating matrix.
- **DRL-FM:** the *DRL with Factorization Machines policy* model replaces the k-NN algorithm from the models previously described with the BPR-FM model to return the Top-K items that can be recommended from a given state. Then, the FM-policy will select the action from the top-K recommendations that yields to the highest Q value in the long term.

The hyperparameters and network settings used in the experiments are detailed in Appendix D. The network configuration for the DDPG algorithm in the DRL-kNN-CB and DRL-FM models consist of a fully connected network with two hidden layers of 400 and 300 neurons each. while the DRL-kNN-CF model consists of two hidden layers of 1000 neurons. The input and output of the network are the current state (item previously recommended) and the item that the agents predicts to be the next recommendation, respectively. Therefore, their size depends on the vector representation of items in each model.

For the DRL-kNN-CB model, feature vectors for each item were created using a combination of word embeddings extracted using *word2vec*<sup>3</sup>, and genre categorization as a binary vector. The final feature vector consists of 119 elements. Whereas, for the DRL-kNN-CF, items representation were extracted from the user rating matrix, so their size are 943 and 6040 elements for the 100-k and 1M datasets respectively.

On the other hand, the network input in the DRL-FM model consists of 119 elements (as it uses the same item representation as the DRL-kNN-CB model). Nevertheless, the output size depends on the binary representation of the action set identifiers. For example, for the 100-k dataset, the output size is 11 as the maximum identifier number (1682: the total number of actions) can be represented using eleven binary digits ( $1682_{10} = 11010010010_2$ ).

To build the proposed policy architectures, we use the *NearestNeighbors* model in the *scikit-learn* library<sup>4</sup> to train a k-NN model for the Wolpertinger policy using the KD-Tree algorithm, while for the FM-policy, a Ranking FM model is trained using the algorithm and default hyperparameter configuration provided in the *Graphlab Create* API<sup>5</sup>. Finally, the DRL models are trained using a value of  $k$  that corresponds to the 5% of the action set size (which yielded to good performance in [7]).

---

<sup>3</sup><https://radimrehurek.com/gensim/index.html>

<sup>4</sup><http://scikit-learn.org>

<sup>5</sup><https://turi.com/products/create/>

In order to train and use the FM model, items and the current state respectively are formatted using the  $SVM^{light}$  representation introduced in [45]. Figure 4.1 presents an example of the feature vector used by the FM model, where the first  $|\mathcal{U}|$  binary indicators variables (blue) represent the active user of a transaction, and the next  $|\mathcal{I}|$  binary indicator variables (red) hold the active item. Each feature vector also contains normalized indicator variables (yellow) for all other items the user has ever rated. Finally, the vector contains a variable (green) holding the time when the rating was registered, and another variable containing information of the last item the user has rated before.

**Figure 4.1:** Feature vector representation for the FM model. Source: *Factorization machines* [17]

Feature vector $\mathbf{x}$																	Target $y$					
$x_1$	1	0	0	...	1	0	0	0	...	0.3	0.3	0.3	0	...	13	0	0	0	0	...	5	$y_1$
$x_2$	1	0	0	...	0	1	0	0	...	0.3	0.3	0.3	0	...	14	1	0	0	0	...	3	$y_2$
$x_3$	1	0	0	...	0	0	1	0	...	0.3	0.3	0.3	0	...	16	0	1	0	0	...	1	$y_3$
$x_4$	0	1	0	...	0	0	1	0	...	0	0	0.5	0.5	...	5	0	0	0	0	...	4	$y_4$
$x_5$	0	1	0	...	0	0	0	1	...	0	0	0.5	0.5	...	8	0	0	1	0	...	5	$y_5$
$x_6$	0	0	1	...	1	0	0	0	...	0.5	0	0.5	0	...	9	0	0	0	0	...	1	$y_6$
$x_7$	0	0	1	...	0	0	1	0	...	0.5	0	0.5	0	...	12	1	0	0	0	...	5	$y_7$
	A	B	C	...	TI	NH	SW	ST	...	TI	NH	SW	ST	...	Time	TI	NH	SW	ST	...		
	User				Movie					Other Movies rated						Last Movie rated						

## 4.5 Results

## **Chapter 5**

# **Conclusion and Discussion**

## **5.1 Summary**

## **5.2 Future Work**

- More robust definition of the simulated environment

  - robust model of user behaviour

  - using a live recommender system

  - monte carlo tree search

- Try to solve the problem using the Q-learning Model-based acceleration algorithm: Normalized Advantage Functions (NAF) to apply Q-learning with experience replay to continuous tasks. it outperforms actor-critic based algorithm (DDPG)

  - test the model using dataset with higher sparsity level



## **Appendix A**

# **DDPG Algorithm using the FM-policy Architecture**





## Appendix B

# Transition probability Matrix algorithm

---

**Algorithm 2:** Transition Probability Matrix(R)

---

**Input:**  $R$ : Initial User Rating Matrix

**Result:**  $\mathcal{P}$ : Transition Probability Matrix

$S \leftarrow \text{ComputeCosineSimilarityMatrix}(R)$

**for**  $i \leftarrow 1$  **to**  $m$  **do**

$sum \leftarrow \sum_{j=1}^m S_{ij}$

**for**  $j \leftarrow 1$  **to**  $m$  **do**

        /\* user walks through a direct neighborith probability  $\beta$  or  
        jumps to an arbitrary item with uniform probability \*/

$\mathcal{P}_{ij} \leftarrow \beta S_{ij}/sum + (1 - \beta)/m$

**end**

**end**

**return**  $\mathcal{P}$

---



## Appendix C

# Recommender System Simulated Environment

---

**Algorithm 3:** Simulated environment

---

```
Input:  $R$ : Ratings Matrix;  $\mathcal{P}$ : Transition Prob. Matrix;  $\mathcal{I}$ : Items set  
 $u_{random} \leftarrow \text{selectRandomUser}(R)$  # selects an active user  
 $state \leftarrow \text{selectLikelyItem}(\mathcal{I}, u_{random})$  # initial guided exploration  
 $term \leftarrow \text{False}$   
while not  $term$  do  
   $action \leftarrow \text{waitForAgentAction}()$   
  /* probability of choosing action given current state */  
   $\rho_{action} \leftarrow \text{computeTransition}(\mathcal{P}, state, action)$   
   $rating \leftarrow \text{getUserRating}(u_{random}, action)$   
  /* get the mean and std dev of the transition probability  
    distr. for a given state */  
   $\mu, \sigma \leftarrow \text{transProbDistribution}(\mathcal{P}, state)$   
   $threshold \leftarrow \mu$   
   $is\_chosen \leftarrow \rho_{action} > threshold$  and  $rating > 0$   
  if  $is\_chosen$  then  
     $\rho_{end} \leftarrow 0.1$   
     $reward \leftarrow rating$   
  else  
     $\rho_{end} \leftarrow 0.2$   
     $reward \leftarrow 0$   
     $action \leftarrow \text{selectRandomItem}(u_{random})$   
     $state \leftarrow action$   
  if  $\text{random.uniform}() < \rho_{end}$  then  
     $term \leftarrow \text{True}$   
   $\text{sendSignalToAgent}(state, reward, term)$ 
```

---



## Appendix D

# Hyperparameter configuration for models

Hyperparameter	DRL-kNN-CB	DRL-kNN-CF	DRL-FM
Hidden Layer 1 Size	400	1,000	400
Hidden Layer 2 Size	300	1,000	400
Replay Memory size	50,000		
$k - NN/Top - K$	5% of action set size		
$\mu$ $L_2$ reg.	0		
$Q$ $L_2$ reg.	0.01		
$\nabla \mu$ learning rate	0.0001		
$\nabla Q$ learning rate	0.001		
$\tau$	0.001		
$\gamma$	0.99		

**Table D.1:** Hyperparameter and network settings for DRL models



# Bibliography

- [1] Paul Resnick and Hal R Varian. Recommender systems. *Communications of the ACM*, 40(3):56–58, 1997.
- [2] Jesús Bobadilla, Fernando Ortega, Antonio Hernando, and Abraham Gutiérrez. Recommender systems survey. *Knowledge-Based Systems*, 46:109–132, 2013.
- [3] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [4] Guy Shani, David Heckerman, and Ronen I Brafman. An mdp-based recommender system. *Journal of Machine Learning Research*, 6(Sep):1265–1295, 2005.
- [5] Hao Wang, Naiyan Wang, and Dit-Yan Yeung. Collaborative deep learning for recommender systems. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1235–1244. ACM, 2015.
- [6] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [7] Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep Reinforcement Learning in Large Discrete Action Spaces. 2015.

- [8] Prem Melville and Vikas Sindhwani. Recommender systems. In *Encyclopedia of machine learning*, pages 829–838. Springer, 2011.
- [9] Paul Resnick, Neophytos Iacovou, Mitesh Suchak, Peter Bergstrom, and John Riedl. Grouplens: an open architecture for collaborative filtering of netnews. In *Proceedings of the 1994 ACM conference on Computer supported cooperative work*, pages 175–186. ACM, 1994.
- [10] Hilmi Yildirim and Mukkai S Krishnamoorthy. A random walk method for alleviating the sparsity problem in collaborative filtering. In *Proceedings of the 2008 ACM conference on Recommender systems*, pages 131–138. ACM, 2008.
- [11] John S Breese, David Heckerman, and Carl Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence*, pages 43–52. Morgan Kaufmann Publishers Inc., 1998.
- [12] Greg Linden, Brent Smith, and Jeremy York. Amazon. com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing*, 7(1):76–80, 2003.
- [13] Yehuda Koren, Robert Bell, Chris Volinsky, et al. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- [14] Marko Balabanović and Yoav Shoham. Fab: content-based, collaborative recommendation. *Communications of the ACM*, 40(3):66–72, 1997.
- [15] Michael Pazzani and Daniel Billsus. Learning and revising user profiles: The identification of interesting web sites. *Machine learning*, 27(3):313–331, 1997.
- [16] Prem Melville, Raymond J Mooney, and Ramadass Nagarajan. Content-boosted collaborative filtering for improved recommendations. In *Aaai/iaai*, pages 187–192, 2002.



- [17] Steffen Rendle. Factorization machines. In *2010 IEEE International Conference on Data Mining*, pages 995–1000. IEEE, 2010.
- [18] Steffen Rendle. Factorization machines with libfm. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 3(3):57, 2012.
- [19] Steffen Rendle, Zeno Gantner, Christoph Freudenthaler, and Lars Schmidt-Thieme. Fast context-aware recommendations with factorization machines. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*, pages 635–644. ACM, 2011.
- [20] Christoph Freudenthaler, Lars Schmidt-Thieme, and Steffen Rendle. Bayesian factorization machines. 2011.
- [21] Yao Wu, Christopher DuBois, Alice X Zheng, and Martin Ester. Collaborative denoising auto-encoders for top-n recommender systems. In *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*, pages 153–162. ACM, 2016.
- [22] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4:237–285, 1996.
- [23] Michael N Katehakis and Arthur F Veinott Jr. The multi-armed bandit problem: decomposition and computation. *Mathematics of Operations Research*, 12(2):262–268, 1987.
- [24] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- [25] Stephan Ten Hagen, Maarten Van Someren, and Vera Hollink. Exploration/exploitation in adaptive recommender systems. *proceedings of Eunit 2003*, 2003.

- [26] Pornthep Rojanavas, Phaitoon Srinil, and Ouen Pinngern. New recommendation system using reinforcement learning. *Special Issue of the Intl. J. Computer, the Internet and Management*, 13, 2005.
- [27] Nima Taghipour, Ahmad Kardan, and Saeed Shiry Ghidary. Usage-based web recommendations: a reinforcement learning approach. In *Proceedings of the 2007 ACM conference on Recommender systems*, pages 113–120. ACM, 2007.
- [28] Nima Taghipour and Ahmad Kardan. A hybrid web recommender system based on q-learning. In *Proceedings of the 2008 ACM symposium on Applied computing*, pages 1164–1168. ACM, 2008.
- [29] Daniela Godoy and Analía Amandi. Modeling user interests by conceptual clustering. *Information Systems*, 31(4):247–265, 2006.
- [30] Elad Liebman, Maytal Saar-Tsechansky, and Peter Stone. Dj-mc: A reinforcement-learning agent for music playlist recommendation. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pages 591–599. International Foundation for Autonomous Agents and Multiagent Systems, 2015.
- [31] Guillaume Chaslot. Monte-carlo tree search. *Maastricht: Universiteit Maastricht*, 2010.
- [32] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [33] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

- [34] Sander Adam, Lucian Busoniu, and Robert Babuska. Experience replay for real-time reinforcement learning control. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 42(2):201–212, 2012.
- [35] Marc G Bellemare, Joel Veness, and Michael Bowling. Investigating contingency awareness using atari 2600 games. In *AAAI*, 2012.
- [36] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [37] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic Policy Gradient Algorithms. *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 387–395, 2014.
- [38] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [39] Enrico Bibbona, Gianna Panfilò, and Patrizia Tavella. The ornstein-uhlenbeck process as a model of a low pass filtered white noise. *Metrologia*, 45(6):S117, 2008.
- [40] Andrew Zimdars, David Maxwell Chickering, and Christopher Meek. Using temporal data for making recommendations. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pages 580–588. Morgan Kaufmann Publishers Inc., 2001.
- [41] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. Bpr: Bayesian personalized ranking from implicit feedback. In *Proceedings of the twenty-fifth conference on uncertainty in artificial intelligence*, pages 452–461. AUAI Press, 2009.

- [42] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [43] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 5(4):19, 2016.
- [44] Immanuel Bayer. Fastfm: a library for factorization machines. *arXiv preprint arXiv:1505.00641*, 2015.
- [45] Thorsten Joachims. Making large scale svm learning practical. Technical report, Universität Dortmund, 1999.