

UNIVERSITÀ DEGLI STUDI DI MILANO
FACULTY OF SCIENCE AND TECHNOLOGY

DEPARTMENT OF COMPUTER SCIENCE



Master in
Computer Science

ALGORITHMS FOR MASSIVE DATASETS
FINAL REPORT ABOUT RECOMMENDER SYSTEM

Teacher: Prof. Dario Malchiodi

Final report written by:
Samuele Simone
Matr. Nr. 11910A

ACADEMIC YEAR 2022-2023

I/We declare that this material, which I/We now submit for assessment, is entirely my/our own work and has not been taken from the work of others, save and to the extent that such work has been cited and acknowledged within the text of my/our work. I/We understand that plagiarism, collusion, and copying are grave and serious offences in the university and accept the penalties that would be imposed should I engage in plagiarism, collusion or copying. This assignment, or any part of it, has not been previously submitted by me/us or any other person for assessment on this or any other course of study.

Contents

Index

1	Introduction	2
2	Environment setup	3
2.1	Pyspark setup	4
2.2	Kaggle setup	4
2.3	Pandas setup	5
2.4	Scikit-learn setup	5
2.5	Numpy setup	6
3	Dataset: A look inside	7
3.1	Data loading	7
3.2	Visualizing the data	8
3.3	Data Filtering	8
3.4	Data preprocessing	9
4	Recommender system	11
4.1	Content-Based Recommendations	11
4.2	Collaborative Filtering	12
5	Creation of a Content-based Recommender system	13
5.1	Business similarity suggestion with Pandas and cosine similarity . .	13
5.1.1	Experiment	14
5.2	Business similarity suggestion with k-NN algorithm	15
5.3	Business similarity suggestion with PySpark and cosine similarity .	16
5.4	Collaborative filtering with ALS ml library	17
5.5	Content based recommendation from scratch in PySpark	17
6	Scalability and complexity	18
6.1	System specification	18

CONTENTS

7	Results and experiments	19
7.1	NDCG	19
8	Conclusion	20

This report is made up of 8 chapters. In the **Chapter 1** we give at the reader a general view of what is a recommender system and where we can find it. **Chapter 2**, we start to discuss about the development setup for the project. Then, in the **Chapter 3** we focus on the dataset, how is composed and we explore the data. Later, in the **Chapter 4** we explain the recommender system. the different approaches and the mechanisms behind. In order to evaluate the project developed there are some notion about scalability and complexity in the **Chapter 6**. Consequently , we summarize the aspects and the results obtained during the various experiments in the **Chapter 7**. Finally there is the conclusion in last **Chapter 8**.

Chapter 1

Introduction

A recommender system is used everywhere nowadays. Indeed all big companies are pushing in these systems because they can increase the sells about their product, e.g., when we are scrolling a product on Amazon, then they show us a list of recommendation based on the item selected.

Recommendation system use a number of different technologies. We can split these system into two broad groups: [1]

- *Content-based systems*: examine properties of the items recommended. For instance, if a Netflix user has watched many cowboy movies, then recommend a movie classified in the database as having the "cowboy" genre.
- *Collaborative filtering* systems recommend items based on similarity measures between users and/or items.

Chapter 2

Environment setup

Before going deeper into the project we must discuss about the entire environment was setup. Indeed I used different libraries in order to create the recommender system. Here the snippet of all libraries used in the project:

```
1 #importing the required pyspark library
2 import pyspark
3 from pyspark import SparkConf
4 from pyspark.sql import SparkSession
5 from pyspark.ml.evaluation import RegressionEvaluator
6 from pyspark.ml.recommendation import ALS
7 from pyspark.sql.functions import split
8 from pyspark.sql.functions import array_contains
9 from pyspark.sql.functions import col
10 from pyspark.sql.functions import from_json
11 from pyspark.sql.functions import when
12 from pyspark.sql import functions as F
13 from pyspark.ml.feature import VectorAssembler
14 from pyspark.ml.feature import StringIndexer
15 from pyspark.ml.feature import MinMaxScaler
16 from pyspark.ml.feature import Normalizer
17 from pyspark.ml.linalg import Vectors
18 from pyspark.sql import Row
19 from pyspark.sql.window import Window
20 #import for manage global var
21 import os
22 #import for graphics
23 import matplotlib.pyplot as plt
24 import pandas as pd
25 #import for regular expression
26 import re
27 import numpy as np
28 from sklearn.metrics.pairwise import cosine_similarity
29 from sklearn.model_selection import train_test_split
30 from sklearn.neighbors import KNeighborsClassifier
```

```
31 from sklearn.model_selection import cross_val_score
32 from sklearn.metrics import accuracy_score
```

Listing 2.1: Loading all the python libraries

2.1 Pyspark setup

PySpark is the Python API for Apache Spark. It enables you to perform real-time, large-scale data processing in a distributed environment using Python. It also provides a PySpark shell for interactively analyzing your data.

PySpark combines Python's learnability and ease of use with the power of Apache Spark to enable processing and analysis of data at any size for everyone familiar with Python.

PySpark supports all of Spark's features such as Spark SQL, DataFrames, Structured Streaming, Machine Learning (MLlib) and Spark Core. [2]

Referring to the Listing 2.1 we can see that there are several libraries about PySpark. I will discuss about the most important:

- **SparkConf**: Configuration for a Spark application. Used to set various Spark parameters as key-value pairs. [3]
- **SparkSession**: The entry point to programming Spark with the Dataset and DataFrame API. [4]
- **pyspark.ml**: DataFrame-based machine learning APIs to let users quickly assemble and configure practical machine learning pipelines. [5]
- **pyspark.sql.functions**: A list of function that allow the user to explore dataframe [6]

In the project I used PySpark Dataframe. They are distributed collections of data that can be run on multiple machines and organize data into named columns. These dataframes can pull from external databases, structured data files or existing resilient distributed datasets (RDDs).

2.2 Kaggle setup

The data were hosted on the Kaggle platform. Kaggle is a subsidiary of Google, it is an online community of data scientists and machine learning engineers.

Kaggle allows users to find datasets they want to use in building AI models, publish datasets, work with other data scientists and machine learning engineers, and enter competitions to solve data science challenges.

Kaggle got its start in 2010 by offering machine learning and data science competitions as well as offering a public data and cloud-based business platform for data science and AI education. [7]

In order to access to the data Kaggle platform gives to every register account an API credential that is necessary to download the datasets.

```
1 os.environ['KAGGLE_USERNAME'] = 'your_kaggle_username'
2 os.environ['KAGGLE_KEY'] = 'your_kaggle_key'
```

Listing 2.2: API Credential for accessing the data

2.3 Pandas setup

Pandas is a Python package providing fast, flexible, and expressive data structures designed to make working with “relational” or “labeled” data both easy and intuitive. It aims to be the fundamental high-level building block for doing practical, real-world data analysis in Python. Additionally, it has the broader goal of becoming the most powerful and flexible open source data analysis/manipulation tool available in any language. It is already well on its way toward this goal. [8]

I used it to explore the data, as we will see in the chapter 3, and to build the content-based recommendation system. My goal was to implement the recommendation system with different types of approaches, and Pandas is one of them. However, in the course of the paper I will emphasize the different ways in which I tried my hand at it.

2.4 Scikit-learn setup

Scikit-learn is a library in Python that provides many unsupervised and supervised learning algorithms. It’s built upon some of the technology like NumPy, pandas, and Matplotlib. The functionality that scikit-learn provides include [9]:

- Regression, including Linear and Logistic Regression
- Classification, including K-Nearest Neighbors
- Clustering, including K-Means and K-Means++
- Model selection
- Preprocessing, including Min-Max Normalization

Indeed some of these functionality are used in this project such as the K-NN, the second different approach to build the Content based recommender system.

2.5 Numpy setup

NumPy is the fundamental package for scientific computing in Python. It is a Python library that provides a multidimensional array object, various derived objects (such as masked arrays and matrices), and an assortment of routines for fast operations on arrays, including mathematical, logical, shape manipulation, sorting, selecting, I/O, discrete Fourier transforms, basic linear algebra, basic statistical operations, random simulation and much more. [\[10\]](#)

It was useful in constructing the function `cosine_similarity_scratch` to perform multiplication and normalization of vectors.

Chapter 3

Dataset: A look inside

3.1 Data loading

We start to look how the data is loaded inside the project. First of all the data is downloaded from Kaggle with this code:

```
1 !kaggle datasets download -d yelp-dataset/yelp-dataset
```

Listing 3.1: Download dataset

and proceeded to unzip the dataset obtaining:

```
1 Archive:  /content/yelp-dataset.zip
2   inflating: yelp-dataset/Dataset_User_Agreement.pdf
3   inflating: yelp-dataset/yelp_academic_dataset_business.json
4   inflating: yelp-dataset/yelp_academic_dataset_checkin.json
5   inflating: yelp-dataset/yelp_academic_dataset_review.json
6   inflating: yelp-dataset/yelp_academic_dataset_tip.json
7   inflating: yelp-dataset/yelp_academic_dataset_user.json
```

Listing 3.2: Unzip dataset

These json files are huge and to manage it we load the data into a Pyspark df with this command:

```
1 df_review = spark.read.json('/content/yelp-dataset/
   yelp_academic_dataset_review.json')
```

Listing 3.3: Loading json data into Pyspark df

We repeat this operation many time as the number of the json files. So we obtain these dataframe:

- **df_review:** (review_id,user_id,business_id,stars_review). These are the attributes I selected from the df. However, there is one field that I think should be mentioned which is the text field, very useful for doing further

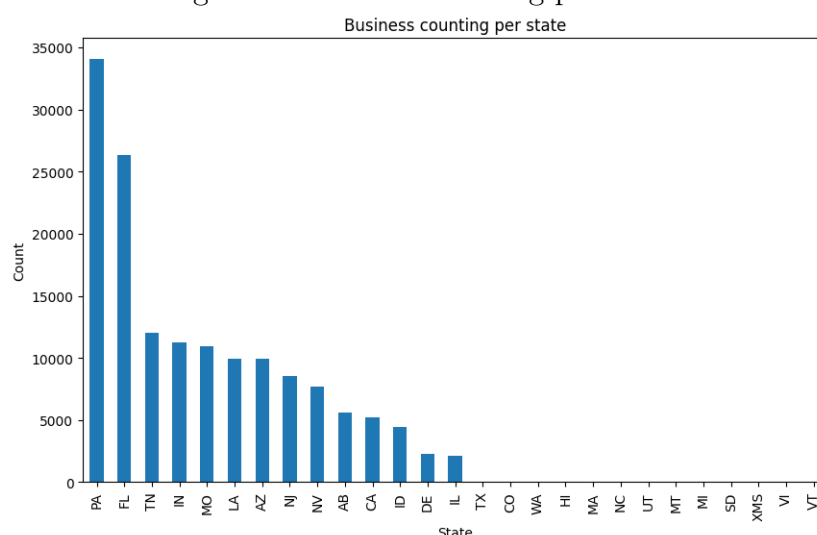
analysis. In my case it was not selected as I believe that the quoted data is sufficient for our purposes.

- **df_users:**(user_id,username,average_stars)
- **df_business:** (business_id,address,attributes,categories,city,is_open,name,stars,state).
There are others column that I prefer to remove because is not adding values for our analysis.

3.2 Visualizing the data

To give a main idea of the data I create a chart where on the x-axes there are all the states and on the y-axes the count of the total business for that particular state. From the Figure 1 we can see that in "PA", we refer to Pennsylvania, a

Figure 1: Business counting per state



state located in the northeastern part of the United States. So I decided to conduct my analysis over the "PA" state due the amount of businesses present over there. Furthermore due the high dimension of the df it's difficult to create others type of graphics.

3.3 Data Filtering

As we discussed in the previous section 3.2 we want to targeting our recommender system in a specific area for a specific type of business. So I decided to take "PA" as

state and "Restaurant" as type of business. Moreover I filtered the new dataframe called `restaurant_df` with another parameter called `is_open` that when is equal to 1 means that this business is actually a running business.

```
1 # Filtering data
2 restaurant_df = df_business_pandas[(df_business_pandas['state'] ==
    state) & (df_business_pandas['is_open'] == 1) &
    df_business_pandas['categories'].str.contains(type_business)==
    True].reset_index()
```

Listing 3.4: Filtering Pyspark df

Here an example of `restaurant_df` printed: Note that the attribute and category

index	address	attributes	business_id	categories	city	is_open	name	stars	state
3	935 Race St	(None, None, 'full_bar', {'touristy': False, '...'}	MTSW4McQd7CbVtyjqoe9uww	Restaurants, Food, Bubble Tea, Coffee & Tea, B...	Philadelphia	1	St Honore Pastries	4.0	"PA"

Table 1: Example of a `restaurant_df`'s row

columns must be preprocessed before being used.

3.4 Data preprocessing

Lets start to analyze the attribute column. As we can see they are in this form:

```
1 Row(AcceptsInsurance=None, AgesAllowed=None, Alcohol="u'none'",
    Ambience=None, BYOB=None,...)
```

Listing 3.5: Attributes item example

So we need a function that allow us to extract key-value pairs in a format understandable to Python. Here the snippet:

```
1 # Function for extracting key-value from attributes
2 def extract_values(row):
3     if row is None or row == "{}":
4         return {}
5
6     attributes = {}
7     pattern = r"(\w+)\s*=\s*([^\s,]+)"
8     matches = re.findall(pattern, row)
9
10    for key, value in matches:
11        value = value.strip('"')
12        attributes[key] = value
13
14    return attributes
```

Listing 3.6: Extracting key-value from attributes

Table 2: Example of `restaurant_df_attr` elements

BestNights	RestaurantsCounterService	WiFi	...	Smoking	CoatCheck
0	1	0	...	0	0
1	0	1	...	1	1

To achieve the result I used the regular expression. More specifically the line 7 defines a regular expression pattern. The pattern captures two groups:

- the key (composed of one or more word characters `\w+`)
- the value (composed of one or more characters that are not a comma `\[,]+`), separated by an equal sign with optional whitespace around it.

Then we apply the function over all the rows of column attributes. Now that the format seems to be good we proceed to create an entire `restaurant_df_attr` in which the attributes are extracted making a new column of the df and populating it with 0 or 1 based on their original value. An example are in the Table 2.

We apply this reasoning also for categories obtaining the `df_categories_dm`. It will be useful to create the `df_total`, a fundamental component for the recommender system. So in the Chapter 4 we are going to understand how recommendation systems work in theory.

Chapter 4

Recommender system

As we mentioned in the Chapter 1, there are two main basic architectures for a recommendation system. In my project I developed from scratch more on a *Content-Based* recommendation system. They are based on properties of items and the similarity of these are determined by measuring the similarity in their properties. However I also worked with the ALS as *Collaborative Filtering* using the PySpark library to compare mine with the one developed with the library. In this Chapter we give to the reader the general notion behind recommender system based on the reference book [1].

4.1 Content-Based Recommendations

We are starting to understand what is an item profiles. It's a record or collection of records representing important characteristics of the item. For example if we consider the restaurants as type of business of this project, we can say that the features of a restaurants can be all the attributes such as *Good-ForKids, BusinessAcceptsBitcoin, DogsAllowed...*. In order to be a recommender system we need also the user. Indeed also the users has their profiles that describes their preferences.

Another fundamental aspect is the utility matrix. It's giving for each user-item pair, a value that represents what is known about the degree of preference of that user for that item. Without a utility matrix, it is almost impossible to recommend items. Usually there are two main approaches:

- We can ask users to evaluate items. We fall in this case because the from the dataset **review** we can access to all user-business pair reviews with the respective rate.

- We can make inferences from user's behavior. Indeed if one user buy something online or just watched a video we can say that he liked it.

Our ultimate goal for content-based recommendation is to create both a item profile consisting of feature-value pairs and a user profile summarizing the preferences of the user, based of their row of the utility matrix. If we consider profiles as vectors we can compute the cosine similarity in the space of the features and understand if the vectors appears closer or not. In this way we can recommend business to the user based on his preferences.

4.2 Collaborative Filtering

When we are referring to Collaborative Filtering we are focusing on the similarity of the user ratings for two items. In place of the item-profile vector for an item, we use its column in the utility matrix. User are similar if their vectors are close according to some notion of distance measures such as cosine distance. The process of identifying similar user and recommending what similar users like is called collaborative filtering.

Chapter 5

Creation of a Content-based Recommender system

Before I talk about the content-based recommendation system that I created from scratch, I want to highlight some of the steps that allowed me to arrive at the solution or at least allowed me to learn more about the data I had at my disposal.

5.1 Business similarity suggestion with Pandas and cosine similarity

First of all I tried to create a suggestion system by looking the business similarity. For achieve this goal I used in this step Pandas. I performed this steps:

1. Extract the reference restaurant index
2. Mapping the rate (1-3 stars will be mapped in 0 and 4-5 in 1)
3. Create the `cosine_similarity_scratch` function

```
1 def cosine_similarity_scratch(vector1, vector2):
2     dot_product = np.dot(vector1, vector2)
3     norm_vector1 = np.linalg.norm(vector1)
4     norm_vector2 = np.linalg.norm(vector2)
5
6     if norm_vector1 == 0 or norm_vector2 == 0:
7         return 0.0
8
9     similarity = dot_product / (norm_vector1 * norm_vector2)
10    return similarity
```

Listing 5.1: Cosine similarity function

4. Calculate business similarity respect the one selected. Order the similar restaurants in a descending way and then select n (= 5 in this case) top similar restaurant and print it.

```

1 restaurant_features = features_df.iloc[restaurant_index,
    6:-2].values # Reference restaurant feature vector
2
3 similarities = []
4 for i in range(len(features_df)):
5     if i != restaurant_index:
6         other_restaurant_features = features_df.iloc[i,
    6:-2].values
7         similarity = cosine_similarity_scratch(
    restaurant_features, other_restaurant_features)
8         similarities.append(similarity)
9
10 similarities = np.array(similarities)
11 similar_restaurants = similarities.argsort()[::-1] # Similar
    restaurant ordered in a desc way
12
13 # Visualize 5 top suggested restaurants with cosine sim
    column
14 top_similar_restaurants = similar_restaurants[:5]
15 recommended_restaurants = df_total.iloc[
    top_similar_restaurants][['name', 'stars']]
16 recommended_restaurants['cosine_similarity'] = similarities[
    top_similar_restaurants]
17
18 print(recommended_restaurants)

```

Listing 5.2: Suggestion system in Pandas

So now I will show the results obtained.

5.1.1 Experiment

I choose the 8071 restaurant index called Adelita Taqueria & Restaurant. It's a mexican restaurant located in 1108 S 9th St Philadelphia, PA 19147. In the note the link of the restaurant on Yelp.¹ The restaurant counts 39 reviews with an average rating of 4.5 stars.

The similar business that the system recommend are these reported in the Table 3 where there are linked their Yelp page. They are all mexican restaurant so it's perfectly matching our selected restaurant. I tried also to use a cosine similarity library from `sklearn.metrics.pairwise` import `cosine_similarity` and the results are quite similar as shown in Table 4.

¹<https://www.yelp.com/biz/adelita-taqueria-and-restaurant-philadelphia?osq=Adelita+Taqueria%26+Restaurant>

Table 3: Top similar business recommended in Pandas

	name	stars	cosine_similarity
7109	Los Taquitos de Puebla	4	1.000000
5990	San Antonio Mexican cousine	5	1.000000
3534	Los Cuatro Soles	5	0.942809
3292	El Limon - Norristown	4	0.942809
1525	La Hacienda Mexican Restaurant	4	0.942809

Table 4: Top similar business recommended in Pandas with sklearn cosine similarity

name	cosine similarity	stars
San Antonio Mexican cousine	0.9999999999999999	5
Los Taquitos de Puebla	0.9486832980505138	4
El Charro Negro	0.9486832980505138	4
El Limon - Norristown	0.9128709291752769	4
Plaza Garibaldi	0.8999999999999999	4

5.2 Business similarity suggestion with k-NN algorithm

I continued my search for similar restaurants also based on the k-nn algorithm. In brief it's a supervised machine learning approach called k-nearest neighbors (k-NN) and it is utilized for both classification and regression problems. It is a non-parametric algorithm that bases its predictions on how closely new input data and training examples resemble each other.

The "k" in the k-NN method denotes how many nearest neighbors are taken into account while producing predictions. The steps done are the following:

- Take features excluding name and stars from the `df_total` as x and stars as y
- Splitting the data into training set and test set
- To figure out which k to choose, I applied 5-fold cross-validation and the best value for k is 20.
- Then I fitted the model and printed the results

Score on training set: 0.5624903205823137 ; Score on test set:
0.52198142486068

- Testing the model and use the last row as the validation test
- Getting the similar restaurants compared with the one taken in the validation test.

In the end, the results obtained are comparable with the other two approaches (cosine similarity from scratch and sklearn).

5.3 Business similarity suggestion with PySpark and cosine similarity

In this section I wanted to focus more on the scalability of the system. Therefore, I converted the approach used in section 5.1 to PySpark, performing all the preprocessing by obtaining `df_final_spark`. Specifically, I summarize the operations:

- Transform the categories into a list of categories for the business
- Perform filtering by state and business type.
- Extract the attributes as columns and then having some columns of nested attributes I re-run the operation
- Transform the values to 0-1 obtaining the `df_business_attr_spark`
- Same for categories creating the `df_categories_dm`
- Finally I get the `df_final_spark` by the join between the categories, attributes and the `restaurant_df_sel_spark` which contains the name and ratings of the various businesses.

Once I had completed all the preprocessing phase and obtained the various df's I could proceed with the recommendation. As seen in the code shown I created the `df_matrix` that contains *business_id, features*. The features were generated using `VectorAssembler`. It is a feature transformer that combines multiple input columns into a single vector column.

```

1 # Exclude the 'business_id' column from df_final
2 columns = [c for c in df_final_spark.columns if c != 'business_id'
3             ]
4 # Convert binary values to integers
5 df_numeric = df_final_spark.select(['business_id'] + [col(c).cast(
6     'integer') for c in columns])
7 # Fill NaN values with 0

```

```

8 df_filled = df_numeric.fillna(0, subset=columns)
9
10 # Use VectorAssembler to create a feature vector column
11 assembler = VectorAssembler(inputCols=columns, outputCol='features'
12                             ')
12 df_matrix = assembler.transform(df_filled).select('business_id', '
           features')

```

Listing 5.3: Suggestion system in PySpark

After that I proceeded to run the same test on the *Adelita Taqueria & Restaurant* using `cosine_similarity_spark` slightly modified to make it work properly with the PySpark pipeline. In this case I wanted to select 15 as the number of recommendations. This is the table obtained:

Table 5: Top similar business recommended in PySpark

business_id	cosine_similarity	name	stars
F2QwLwzS3vF9os0...	[[1.0]]	San Antonio Mexican cousine	4.5
DwOhLOd9Say...	[[0.9803789354850793]]	Los Taquitos de Puebla	3.5
EoQiJ5D-pyWc...	[[0.9714285714285714]]	El Primo Taqueria	4.5
XQewVfTaosZ3U-4g...	[[0.9714285714285714]]	Los Cuatro Soles	4.5
mKJ_WV7TvrjyDjm...	[[0.9710083124552245]]	El Limon - Bensalem	4.5
SVf23pjKERkedqCdWl6ECA	[[0.9660917830792959]]	Teresa's Mesa	4.0
5ItgryJvadUrKVIjjJ8l4g	[[0.9660917830792959]]	Que Chula Es Puebla	4.0
rNg75hKR0UIB5-jX5WiVQg	[[0.9660917830792959]]	Los Mariachis	3.5
1MVMKUvZfWwIqkhxP3rYvQ	[[0.9613406389911041]]	Indian Garden	4.0
sKUs4ISUgn3j6SeLYvSByg	[[0.9613406389911041]]	El Charro Negro	4.0
x1CtK2qnlCr_1DJFNo1_vw	[[0.9583148474999098]]	El Limon - Norristown	4.0
venWfi69QVylbyrvaPt0nQ	[[0.9578414886923188]]	Cafe Ynez	4.5
t0Qyogb4x-K9i5b0AoDCg	[[0.9578414886923188]]	Taco Maya	4.5
fcAaBcYFo1YqetIjs9l8Iw	[[0.9566222088265397]]	Sungate Diner	4.5
uh-387pbyipzjCMcj5Dy6w	[[0.9566222088265397]]	Las Palmas Del Sur	4.5

Also you can see that they are all restaurants that cook Mexican style.

5.4 Collaborative filtering with ALS ml library

5.5 Content based recommendation from scratch in PySpark

Chapter 6

Scalability and complexity

6.1 System specification

Chapter 7

Results and experiments

7.1 NDCG

Chapter 8

Conclusion

Listings

2.1	Loading all the python libraries	3
2.2	API Credential for accessing the data	5
3.1	Download dataset	7
3.2	Unzip dataset	7
3.3	Loading json data into Pyspark df	7
3.4	Filtering Pyspark df	9
3.5	Attributes item example	9
3.6	Extracting key-value from attributes	9
5.1	Cosine similarity function	13
5.2	Suggestion system in Pandas	14
5.3	Suggestion system in PySpark	16

Bibliography

- [1] Anand Rajaraman, Jure Leskovec, and Jeffrey D. Ullman. *Mining Massive Datasets*. Cambridge University Press, 2014.
- [2] Apache Spark. Pyspark overview, 2023.
- [3] Apache Spark. pyspark.sparkconf, 2023.
- [4] Apache Spark. pyspark.sparksession, 2023.
- [5] Apache Spark. pyspark.ml, 2023.
- [6] Apache Spark. pyspark.sql.function, 2023.
- [7] Mahmoud Hassan Mahmoud. What is kaggle?, 2022.
- [8] Pandas. Pandas package overview, 2023.
- [9] codecademy. What is scikit-learn?, 2023.
- [10] codecademy. What is numpy?, 2023.