

Collaborative Filtering for the Prediction of Drug-Target Associations

Mirco Perna

A.A. 2024/2025

1 Introduction

The aim of the project is to study and develop Collaborative Filtering, a technique used in Recommender Systems to predict a user's interests based on the preferences of other users with similar tastes, to predict new targets for a given drug, based on the protein-protein interactions (PPI).

The main idea is to identify new targets for existing drugs by searching for among 'neighbors' of their known targets within the associated protein-protein interaction network. It is important to note that not every protein that physically interacts with a drug's known targets should necessarily be considered a potential target itself. However, if a set of known drug-target associations is used to train the recommendation system, only those proteins that exhibit behavior similar to that of actual targets are likely to be identified.

2 Collaborative Filtering

2.1 Recommendation systems

Recommendation systems are made up of filtering algorithms that aim to predict a rating or preference a user would assign to a given item, which have become increasingly important across a variety of commercial domain.

These systems generally produce recommendations through one of two methods:

1. Content based filtering
2. Collaborative based filtering

Content-based filtering techniques use attributes of an item in order to recommend future items with similar attributes. Collaborative filtering builds a model from a user's past behavior, activities, or preferences and makes recommendations to the user based on similarities with other users. However, collaborative filtering methods generally offer higher accuracy than content-based approaches.

2.2 Approaches

Collaborative filtering is broken down into two primary approaches: model-based and memory-based.

The **memory-based** approaches take user rating data to compute similarities between users and items in order to make a recommendation. The most famous memory-based approach are neighborhood-based algorithms. One of the most commonly used algorithms for this approach is the K-nearest neighbor (KNN) algorithm.

The **model-based** (or matrix factorization based) approaches build models based on modern machine learning algorithms discovering patterns in the training data. The models are then used to make predictions on real data.

2.3 Matrix Factorization based

Matrix factorization is a simple embedding model. Given the feedback matrix

$$A \in \mathbb{R}^{M \times N}$$

where M is the number of users and N is the number of items. The model learns:

- A user embedding matrix

$$U \in \mathbb{R}^{M \times D}$$

- A item embedding matrix

$$V \in \mathbb{R}^{D \times N}$$

The value of the dimension D is a system parameter and it is called latent factors. The embeddings are learned such that the product of matrix U and the transpose of matrix V is a good approximation of the feedback matrix A.

$$A \approx UV^T$$

Observe that (i, j) the entry of UV^T is simply the dot product $\langle U_i \cdot V_j \rangle$ of the embeddings of user i and item j. Ideally $a_{i,j} = \langle u_i, v_j \rangle \forall i, j$ but in practice we need to minimize the loss function, which in this case is the root mean square error (RMSE)

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (p_{u,v} - a_{u,v})^2}$$

where $p_{u,v}$ and $r_{u,v}$ are respectively predicted and observed ratings.

2.4 Alternating Least Square

The chosen algorithm is Alternating Least Squares (ALS), which performs well on large and sparse datasets. The ALS rotates between fixing one of the unknowns u_i or v_j . When one is fixed the other can be computed by solving the least-squares problem. This approach is useful because it turns the previous non-convex problem into a quadratic that can be solved optimally

Algorithm 1 ALS

Step 1: Initialize matrix V by assigning the average rating for that item as the first row, and small random numbers for the remaining entries.

Step 2: Fix V, solve U by minimizing the RMSE function.

Step 3: Fix U, solve V by minimizing the RMSE function similarly.

Step 4: Repeat Steps 2 and 3 until convergence.

3 Experiment

3.1 Dataset

The Data Set was built manually and is composed by two sets of pairs:

- Protein-Protein interactions (PPI)

$$\langle Protein, Protein \rangle$$

- Drug-Target associations

$$\langle Drug, Protein \rangle$$

After the two sets were built, they were inserted in an instance of MongoDB, a document database used to build highly available and scalable internet applications, and it uses flexible schema approach. The document used to create the database is shown below.

```

1 {
2   "_id": "collaborative_filtering",
3   "PPI": [
4     { "ID_Protein_A": "PT0001", "ID_Protein_B": "PT1002" },
5     { "ID_Protein_A": "PT0001", "ID_Protein_B": "PT1004" },
6     { "ID_Protein_A": "PT0001", "ID_Protein_B": "PT1005" },
7     { "ID_Protein_A": "PT0002", "ID_Protein_B": "PT1004" },
8     { "ID_Protein_A": "PT0003", "ID_Protein_B": "PT1005" },
9     { "ID_Protein_A": "PT0003", "ID_Protein_B": "PT1004" },
10    { "ID_Protein_A": "PT0004", "ID_Protein_B": "PT1001" },
11    { "ID_Protein_A": "PT0005", "ID_Protein_B": "PT1006" },
12    { "ID_Protein_A": "PT0005", "ID_Protein_B": "PT1000" },
13    { "ID_Protein_A": "PT0006", "ID_Protein_B": "PT1001" },
14    { "ID_Protein_A": "PT0007", "ID_Protein_B": "PT1005" },
15    { "ID_Protein_A": "PT0007", "ID_Protein_B": "PT1003" },
16    { "ID_Protein_A": "PT0008", "ID_Protein_B": "PT1003" },
17    { "ID_Protein_A": "PT0009", "ID_Protein_B": "PT1001" },
18    { "ID_Protein_A": "PT0009", "ID_Protein_B": "PT1005" },
19    { "ID_Protein_A": "PT0010", "ID_Protein_B": "PT1003" },
20    { "ID_Protein_A": "PT0011", "ID_Protein_B": "PT1005" },
21    { "ID_Protein_A": "PT0011", "ID_Protein_B": "PT1003" },
22    { "ID_Protein_A": "PT0011", "ID_Protein_B": "PT1002" },
23    { "ID_Protein_A": "PT0012", "ID_Protein_B": "PT1000" },
24    { "ID_Protein_A": "PT0013", "ID_Protein_B": "PT1002" },
25    { "ID_Protein_A": "PT0013", "ID_Protein_B": "PT1003" },
26    { "ID_Protein_A": "PT0014", "ID_Protein_B": "PT1005" },
27    { "ID_Protein_A": "PT0015", "ID_Protein_B": "PT1006" },
28    { "ID_Protein_A": "PT0015", "ID_Protein_B": "PT1005" },
29    { "ID_Protein_A": "PT0019", "ID_Protein_B": "PT1000" },
30    { "ID_Protein_A": "PT0019", "ID_Protein_B": "PT1001" },
31    { "ID_Protein_A": "PT0020", "ID_Protein_B": "PT1001" }
32  ],
33  "Drug_Target": [
34    { "ID_DrugBank": "DB0001", "ID_Target": "PT0001" },
35    { "ID_DrugBank": "DB0001", "ID_Target": "PT0002" },
36    { "ID_DrugBank": "DB0001", "ID_Target": "PT0003" },
37    { "ID_DrugBank": "DB0001", "ID_Target": "PT0004" },
38    { "ID_DrugBank": "DB0001", "ID_Target": "PT0005" },
39    { "ID_DrugBank": "DB0002", "ID_Target": "PT0006" },
40    { "ID_DrugBank": "DB0002", "ID_Target": "PT0007" },
41    { "ID_DrugBank": "DB0002", "ID_Target": "PT0008" },
42    { "ID_DrugBank": "DB0002", "ID_Target": "PT0009" },
43    { "ID_DrugBank": "DB0002", "ID_Target": "PT0010" },
44    { "ID_DrugBank": "DB0003", "ID_Target": "PT0011" },
45    { "ID_DrugBank": "DB0003", "ID_Target": "PT0012" },
46    { "ID_DrugBank": "DB0003", "ID_Target": "PT0013" },
47    { "ID_DrugBank": "DB0003", "ID_Target": "PT0014" },
48    { "ID_DrugBank": "DB0003", "ID_Target": "PT0015" },
49    { "ID_DrugBank": "DB0004", "ID_Target": "PT0016" },
50    { "ID_DrugBank": "DB0004", "ID_Target": "PT0017" },
51    { "ID_DrugBank": "DB0004", "ID_Target": "PT0018" },
52    { "ID_DrugBank": "DB0004", "ID_Target": "PT0019" },
53    { "ID_DrugBank": "DB0004", "ID_Target": "PT0020" },
54    { "ID_DrugBank": "DB0005", "ID_Target": "PT0001" },
55    { "ID_DrugBank": "DB0005", "ID_Target": "PT0006" },
56    { "ID_DrugBank": "DB0005", "ID_Target": "PT0011" },
57    { "ID_DrugBank": "DB0005", "ID_Target": "PT0016" },
58    { "ID_DrugBank": "DB0005", "ID_Target": "PT0020" }
59  ]
60 }
61 }

```

3.2 Code

Python was used as a programming language and the following libraries:

- PySpark
- Numpy
- PyMongo

The main steps of the experiment are as follows, and they are described in detail in the following paragraphs:

1. Obtain the two sets by MongoDB
2. Process the sets to obtain the dataset for training the model
3. Training the model
4. Evaluate the recommended protein targets for the drugs

3.3 DataAccess

In order to obtain the sets from the database, the DataAccess class was implemented. It contains the methods get_PPIS and get_drug_target_interactions, which are used to extract the sets from the database and return a list for each of them, respectively.

```
1 from pymongo import MongoClient
2 from PPI import PPI
3 from drug_target import DrugTarget
4
5 class DataAccess():
6     def __init__(self):
7         self.client = MongoClient("mongodb://localhost:27017/")
8         self.database = self.client["Biological_Network"]
9         self.collection = self.database["Collaborative_Filtering"]
10        self.id_collection = "collaborative_filtering"
11
12
13    def get_PPIS(self):
14        data = self.collection.find_one({"_id":self.id_collection})
15
16        PPIS = []
17        for interaction in data["PPI"]:
18            new_ppi = PPI()
19            new_ppi.set_properties(interaction["ID_Protein_A"], interaction["
                ID_Protein_B"])
20            PPIS.append(new_ppi)
21
22        return PPIS
23
24    def get_drug_target_interactions(self):
25        data = self.collection.find_one({"_id":self.id_collection})
26
27        drug_target_interactions = []
28        for interaction in data["Drug_Target"]:
29            new_drug_target = DrugTarget()
30            new_drug_target.set_properties(interaction["ID_DrugBank"],
                interaction["ID_Target"])
31            drug_target_interactions.append(new_drug_target)
32
33        return drug_target_interactions
34
35    def close_connection(self):
36        self.client.close()
```

3.4 DataFrame

In order to obtain the dataset required to train and test our model, the class DataFrame was implemented.

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import collect_list
3 from pyspark.sql.functions import lit
4 from pyspark.sql.functions import expr
5 from pyspark.ml.feature import StringIndexer
6 from pyspark.ml import Pipeline
7
8 class DataFrame():
9
10     def __init__(self, PPIs, drug_target_interactions):
11         self.spark = SparkSession.builder \
12             .appName("Collaborative_Filtering") \
13             .config("spark.driver.host", "localhost") \
14             .config("spark.driver.bindAddress", "127.0.0.1") \
15             .getOrCreate()
16
17         self.spark.sparkContext.setLogLevel("ERROR")
18
19         self.create_ppi_dataframe(PPIs)
20
21         self.create_drug_target_interactions_dataframe(drug_target_interactions)
22
23
24     def create_ppi_dataframe(self, PPIs):
25         PPI_list = [PPI.to_list() for PPI in PPIs]
26         self.PPIs_df = self.spark.createDataFrame(PPI_list, ["ID_Protein_A", "
27             ID_Protein_B"])
28
29     def create_drug_target_interactions_dataframe(self, drug_target_interactions)
30     :
31         drug_target_interactions_list = [interaction.to_list() for interaction
32             in drug_target_interactions]
33         self.drug_target_interactions_df = self.spark.createDataFrame(
34             drug_target_interactions_list, ["ID_DrugBank", "ID_Target"])
35
36     def create_date_frame_for_als(self):
37         def create_drug_interactions():
38             self.drug_interactions_df = self.drug_target_interactions_df
39                 .groupBy("ID_DrugBank")
40                 .agg(collect_list("ID_Target")
41                     .alias("Proteins"))
42                 .orderBy("ID_DrugBank")
43
44         def join_drug_target_with_ppi():
45             self.joined_df = self.drug_target_interactions_df.join(self.PPIs_df
46                 , self.drug_target_interactions_df.ID_Target == self.PPIs_df.
47                 ID_Protein_A)
48             self.joined_df = self.joined_df.withColumnRenamed("ID_DrugBank", "
49                 ID_Drug")
50
51         def filter_interactions():
52             self.joined_df = self.joined_df.join(self.drug_interactions_df,
53                 self.joined_df.ID_Drug == self.drug_interactions_df.
54                 ID_DrugBank)
55             self.joined_df = self.joined_df.select(self.joined_df["ID_Drug"],
56                 self.joined_df["ID_Target"], self.joined_df["ID_Protein_A"],
57                 self.joined_df["ID_Protein_B"], self.joined_df["Proteins"])
```

```

47         self.joined_df = self.joined_df.withColumn("Interactor_drug_target
48             ", expr("array_contains(Proteins, ID_Protein_B)"))
49
50         self.joined_df = self.joined_df.filter(self.joined_df.
51             Interactor_drug_target == False)
52
53     def calculate_interaction():
54         columns_to_group_by = ["ID_Drug", "ID_Protein_B"]
55         self.joined_df = self.joined_df
56             .groupBy(columns_to_group_by)
57             .count()
58
59         self.joined_df = self.joined_df
60             .withColumnRenamed("count", "Interactions")
61         self.joined_df = self.joined_df
62             .withColumnRenamed("ID_Protein_B", "ID_Protein")
63
64         create_drug_interactions()
65
66         join_drug_target_with_ppi()
67
68         filter_interactions()
69
70         calculate_interaction()

```

The role of each method will not be explained in detail; instead, the focus will be on the general process that leads to the construction of the final dataset which consists of the following steps:

1. Create the PPIs_df and drug_target_interactions_df DataFrames.
2. Create the drug_interactions_df DataFrame, where the first column contains the drug and in the other column the list of proteins associated with that drug.
3. Create the joined_df DataFrame which is generated by the join of PPIs_df and drug_target_interactions_df DataFrames, where the ID_Target matches ID_Protein_A
4. Join the joined_df and drug_interactions_df DataFrames, where the ID_Drug matches ID_DrugBank. Then, use a select operation to extract the columns: ID_Drug, ID_Target, ID_Protein_A, ID_Protein_B and Proteins
5. Remove all the records where ID_Protein_B is in Proteins
6. Using a group-by operation on the ID_Target and ID_Protein_B columns, we compute the value of the Interactions column, which represents the number of target proteins associated with the drug ID_Target that interact with the protein ID_Protein_B. This value will be used as a rating for the recommendation system.
7. Finally, perform a select operation to extract the columns ID_Drug , ID_Protein_B, which will be renamed to ID_Protein, and Interactions in order to obtain the final DataFrame for the model.

3.5 ALSModel

The ALSModel class includes methods to train the model and to evaluate new potential targets. The ALSModel class includes methods to train the model and to evaluate new potential targets. Since the training values must be of type Integer or Float, the StringIndexer function from PySpark was used to encode all string-type fields. However, as this function can only be applied to one column at a time, and the resulting indexes must remain distinct across columns, the Pipeline function from PySpark was used to manage and apply the transformations efficiently. The indexing_name method is responsible for all these operations.

During training, a evaluation is performed to find the optimal hyperparameters that minimize the RMSE. The model requires the following hyperparameters:

Hyperparameter	Description
maxIter	Maximum number of iterations to run
regParam	Specifies the regularization parameter in ALS
rank	Number of latent factors in the model
alpha	Parameter applicable to the implicit feedback variant of ALS that governs the baseline confidence in preference observations
coldStartStrategy	Strategy to manage new or unknown items/users. Setting it to 'drop' excludes these items from the results

To evaluate the RMSE value for the chosen regParam, rank and alpha parameters the RegressionEvaluator function is used.

Finally, the method calculate_recommended_proteins evaluates new potential targets.

```

1 from pyspark.ml.recommendation import ALS
2 from pyspark.ml.evaluation import RegressionEvaluator
3 from pyspark.sql.functions import explode
4 from pyspark.sql.functions import first
5 from pyspark.ml.feature import StringIndexer
6 from pyspark.ml.feature import IndexToString
7 from pyspark.ml import Pipeline
8
9 class ALSModel():
10     def __init__(self, data):
11         self.indexing_name(data)
12
13     def indexing_name(self, data):
14         self.drug_indexer = StringIndexer(inputCol = "ID_Drug", outputCol = "
15             ID_Drug_Index").fit(data)
16         self.protein_indexer = StringIndexer(inputCol = "ID_Protein",
17             outputCol = "ID_Protein_Index").fit(data)
18         pipeline = Pipeline(stages = [self.drug_indexer, self.protein_indexer
19             ])
20         self.data = pipeline.fit(data).transform(data)
21         self.data.show()
22
23     def train(self):
24         (training, test) = self.data.randomSplit([0.8, 0.2], seed=42)
25
26         regParams = [0.01, 0.1]
27         ranks = [25,30,35]
28         alphas = [10.0, 20.0, 40.0, 60.0, 80.0, 100.0]
29
30         self.aus_regParam = 0.0
31         self.aus_rank = 0
32         self.aus_alpha = 0.0

```

```

30     self.aus_rmse = 0.0
31
32     for regParam in regParams:
33         for rank in ranks:
34             for alpha in alphas:
35                 aus_als = ALS(maxIter = 10, regParam = regParam, rank =
36                     rank, alpha = alpha, userCol = "ID_Drug_Index",
37                     itemCol = "ID_Protein_Index", ratingCol = "
38                         Interactions", coldStartStrategy = "drop"
39                     )
40
41                 aus_model = aus_als.fit(training)
42                 predictions = aus_model.transform(test)
43                 evaluator = RegressionEvaluator(metricName = "rmse",
44                     labelCol = "Interactions", predictionCol = "prediction
45                     ")
46                 rmse = evaluator.evaluate(predictions)
47
48                 if(self.aus_rmse == 0.0 or rmse < self.aus_rmse):
49                     self.aus_regParam = regParam
50                     self.aus_rank = rank
51                     self.aus_alpha = alpha
52                     self.aus_rmse = rmse
53                     self.model = aus_model
54
55                 print("For regParam: {0}, rank:{1}, alpha:{2}, RMSE:{3}".
56                     format(regParam, rank, alpha, rmse))
57
58     print("Chosen parameters: regParam: {0}, rank:{1}, alpha:{2}, RMSE:{3}
59         ".format(self.aus_regParam, self.aus_rank, self.aus_alpha, self.
60             aus_rmse))
61
62     def from_index_to_name(self, proteins_recommended):
63         drug_name = IndexToString(inputCol = "ID_Drug_Index", outputCol =
64             "ID_Drug", labels = self.drug_indexer.labels)
65
66         prontein_name= IndexToString(inputCol = "ID_Protein_Index",
67             outputCol = "ID_Protein",
68
69             labels = self.prontein_indexer.labels)
70
71         pipeline = Pipeline(stages = [drug_name, prontein_name])
72
73         self.drug_proteins_recommended = pipeline.fit(self.data)
74         .transform(proteins_recommended)
75
76         self.drug_proteins_recommended = self.drug_proteins_recommended.
77             select("ID_Drug", "ID_Protein", "rating").orderBy("ID_Drug", "
78                 rating")
79
80     def calculate_recommended_proteins(self):
81         self.train()
82
83         amount_proteins_for_drug = 4
84
85         proteins_recommended = self.model.recommendForAllUsers(
86             amount_proteins_for_drug)
87
88         proteins_recommended = proteins_recommended.withColumn("
89             proteinAndRating", explode(proteins_recommended.
90                 recommendations))
91         .select("ID_Drug_Index", "proteinAndRating.*")

```


77
78

```
self.from_index_to_name(proteins_recommended)
```

3.6 Results

The `calculate_recommended_proteins` function evaluates the top protein recommendations for each drug. The following figure shows the predicted proteins for the DB0003 drug are actually those with the highest Interactions values.

ID_Drug	ID_Protein	rating
DB0003	PT1005	2.9748125
DB0003	PT1004	2.027082
DB0003	PT1002	2.0130541
DB0003	PT1003	2.0003428
DB0003	PT1001	1.6661369
DB0003	PT1000	1.0055522
DB0003	PT1006	0.6756939

ID_Drug	ID_Protein	Interactions
DB0003	PT1005	3
DB0003	PT1003	2
DB0003	PT1002	2
DB0003	PT1006	1
DB0003	PT1000	1

Finally, the following image shows both the table of actual values and the table of predicted values, highlighting the small differences between them.

ID_Drug	PT1000	PT1001	PT1002	PT1003	PT1004	PT1005	PT1006
DB0001	1	1	1	NULL	3	2	1
DB0002	NULL	2	NULL	3	NULL	2	NULL
DB0003	1	NULL	2	2	NULL	3	1
DB0004	1	2	NULL	NULL	NULL	NULL	NULL
DB0005	NULL	2	2	1	1	2	NULL

ID_Drug	PT1000	PT1001	PT1002	PT1003	PT1004	PT1005	PT1006
DB0001	0.9926122	1.0018413	1.1318059	0.9785331	2.984873	2.000412	0.9949578
DB0002	0.660511	1.2651672	1.5605621	2.996085	1.0428447	2.2388773	0.3476149
DB0003	1.0055522	1.6661369	2.0130541	2.0003428	2.027082	2.9748125	0.6756939
DB0004	0.99036026	1.9968805	2.082445	1.3171701	1.7327752	2.4978905	0.5775917
DB0005	0.7375037	1.9904449	1.9897145	1.0000297	0.69525766	2.0055344	0.23175254