



MIDDLE EAST TECHNICAL UNIVERSITY

DEPARTMENT OF COMPUTER ENGINEERING

CENG 300

Summer Practice Report

METU Data Mining Research Group

Start Date: End Date:

Total Working Dates:

October 4, 2020

Student:
Onat ÖZDEMİR

Instructors:
Prof.Dr.Pınar
KARAGÖZ

Prof.Dr.İsmail Hakkı
TOROSLU

Student's Signature

Organization Approval

Contents

1	Introduction	2
2	Project	2
2.1	Analysis Phase	2
2.2	Design Phase	2
2.3	Implementation Phase	2
2.3.1	Neo4j	3
2.3.2	Numpy	4
2.3.3	Scipy	4
2.4	Eigentrust Weighted Trust Based Recommender	5
2.4.1	About Eigentrust	5
2.5	Inverse Distance Weighted Trust Based Recommender	6
2.5.1	Graph Module	6
2.5.1.1	Constructing Adjacency Matrix and Distance Matrix	6
2.5.1.2	Trust Calculation	8
2.5.2	Filterer Module	8
2.5.3	Recommender Module	8
2.6	Testing Phase	9
2.6.1	Methods	9
2.6.1.1	Leave-one-out Cross Validation	9
2.6.1.2	K-Fold Cross Validation	9
2.6.2	Results	9
2.6.3	Libraries I Used	9
2.6.3.1	Surprise	9
2.6.3.2	Matplotlib	10
3	Organization	10
3.1	METU Data Mining Research Group	10
4	Conclusion	10

1 Introduction

I have done my summer internship at METU Data Mining Research Group under the supervision of Prof.Dr.Pınar KARAGÖZ and Prof.Dr.İsmail Hakkı TOROSLU. The task I have worked on was implementing a Trust Based Recommender using the collaborative filtering method and testing it on provided dataset. The dataset contains information about customers and the products they have bought. In addition to dataset, I was able to use eigentrust calculation and community detection modules provided by the TACOREC.

2 Project

During the internship, I implemented two trust based recommenders with different weightening methods:

1. Eigentrust Weighted Recommender
2. Inverse Distance Weighted Recommender

The details of these two recommenders can be found in section 2.4 and 2.5, respectively.

2.1 Analysis Phase

There were two problems I need to solve:

- 1 - Dataset was very sparse
- 2 - There was no explicit trust information

In the both implementations I have made, first case handled by filtering the customers and products which purchased and were bought more than filtering threshold times. To gain better understanding on the second problem, I studied implicit trust calculation methods and looked into lots of research papers.

2.2 Design Phase

2.3 Implementation Phase

Since there are two different implementations, I have divided the implementational details of the recommenders into two subsections: section 2.4 and 2.5.

Under this subsection, libraries and technologies used in implementations are explained.

2.3.1 Neo4j

Driver Installation :

```
1 pip install neo4j
2
```

Configuration :

```
1 import neo4j
2 ...
3
4 uri = self._config["database"]["neo4j"]["uri"]
5 user = self._config["database"]["neo4j"]["user"]
6 password = self._config["database"]["neo4j"]["password"]
7
8 self._driver = neo4j.Driver(uri, auth=(user, password))
9
10
```

Sample Usage :

```
1 import neo4j
2 ...
3
4 def get_customer_trust(self, customer_id):
5
6     query = (
7         f"MATCH (u:Customer)-[r:BELONGS_IN]->(c:Community) "
8         f"WHERE u.id = {repr(customer_id)} "
9         f"RETURN r.eigentrust"
10    )
11
12    with self._driver.session() as session:
13        return tuple(session.run(query).single())
14
15
```

Listing 1: Neo4j driver example

2.3.2 Numpy

Installation :

```
1 pip install numpy
2
```

Sample Usage :

```
1 import numpy as np
2
3 class TrustBasedFilterer(object):
4     ...
5
6     def _create_customers_versus_products_table(self):
7
8         self._customers_versus_products_table = np.zeros(
9             (self._unique_customers.shape[0],
10              self._unique_products.shape[0]),
11             dtype=np.bool,
12             )
13
14         self._customers_versus_products_table[
15             self._sales[:, 0],
16             self._sales[:, 1],
17             ] = True
18
```

Listing 2: Numpy example

2.3.3 Scipy

Installation :

```
1 pip install scipy
2
```

Sample Usage :

```
1 from scipy.sparse import csr_matrix
2 from scipy.sparse.csgraph import dijkstra
3
4 class Graph(object):
5     ...
6
```

```

7  def _create_distance_matrix(self):
8
9  self._create_adjacency_matrix()
10
11 self._adjacency_matrix = \
12 csr_matrix(self._adjacency_matrix)
13
14 self._distance_matrix = dijkstra(
15 csgraph=self._adjacency_matrix,
16 directed=False,
17 return_predecessors=False,
18 unweighted=True,
19 limit=self._max_distance)
20
21 self._distance_matrix\
22 [~np.isfinite(self._distance_matrix)] = 0
23

```

Listing 3: Scipy example

2.4 Eigentrust Weighted Trust Based Recommender

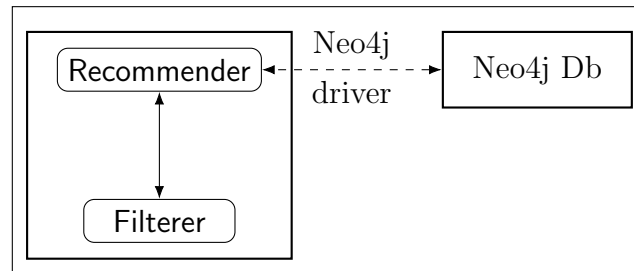


Figure 1: Recommender Structure

2.4.1 About Eigentrust

Eigentrust[4] is a reputation calculation algorithm mainly designed for peer-to-peer networks. In our case, Eigentrust represents how strongly connected the customers are to their communities. Eigentrust values calculated by the eigentrust module provided by TACoRec[1] and stored in Neo4j database as a property of the relationship between a customer and his/her community.

Problem encountered with Eigentrust: Especially for the customers connected to communities with small size and low densities, eigentrust values stored in the database are either very small or equal to zero. Most of the customers with zero eigentrust values are eliminated after filtering the network from customers with a small number of products. Unfortunately, eigentrust values are still quite small.

2.5 Inverse Distance Weighted Trust Based Recommender

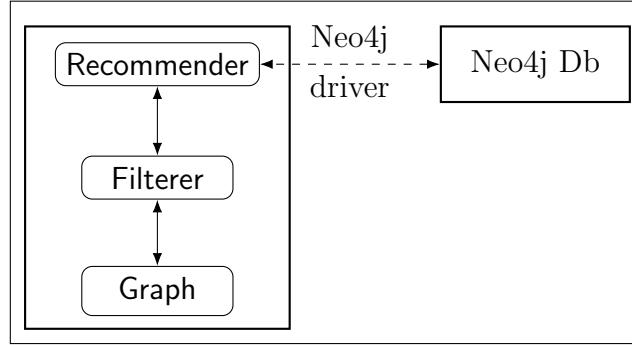


Figure 2: Recommender Structure

Inverse Distance Weighted Trust Based Recommender consists of three modules:

2.5.1 Graph Module

Responsible for creating adjacency matrix from given customer versus product matrix, calculating the distances between customers applying djikstra algorithm to the adjacency matrix using scipy library, and calculating the trust by taking the reciprococal of the distance matrix.

2.5.1.1 Constructing Adjacency Matrix and Distance Matrix

To construct the "adjacency matrix" from customer versus products table, I propose two methods:

Proposed Method 1: Unweighted Graph

In this method, the "adjacency matrix" is constructed based on whether customers purchased a joint product or not. In other words, edge between two customers can exist if and only if the intersection of the set of products they purchased is not the empty set.

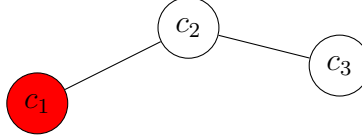


Figure 3: c_1 and c_2 have purchased at least 1 joint product, but c_1 and c_3 do not have a joint product

Proposed Method 2: Euclidean Distance Weighted Graph

In this method, the "adjacency matrix" is constructed based on the "euclidean distances" between customers.

$$adj[c_1][c_2] = \sqrt{\sum_{i \in I_1 \cap I_2} (r1_i - r2_i)^2} \quad (1)$$

where $r1_i$ and $r2_i$ represents ratings given by c_1 and c_2 for product i . Unlike the commonly used "euclidean distance" calculation, in this method, only ratings given to joint products are included in the calculation.

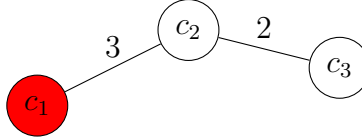


Figure 4: euclidean distance between c_1 and c_2 equals 3, and c_1 and c_3 do not have a joint product

Problem encountered with euclidean distance method: As the module is generally tested on sparse datasets,

Dijkstra's Algorithm

To construct the "distance matrix", the graph module uses "Dijkstra's Algorithm" [2] which takes the adjacency matrix as parameter and returns the distance matrix.

2.5.1.2 Trust Calculation

After calculating the shortest distance between each pair of customers using "Dijkstra's Algorithm", to calculate the trust score between customers Graph module uses

$$T(c_1, c_2) = \begin{cases} \frac{1}{d(c_1, c_2)} & d(c_1, c_2) \neq np.inf \\ 0 & d(c_1, c_2) = np.inf \end{cases}$$

function where $d(c_1, c_2)$ represents the shortest distance between the *customer*₁ and *customer*₂.

A benefit of the method: Especially for excessively sparse datasets, recommenders using euclidean distance-based similarity fails since they cannot calculate similarity score for the the customer pairs with no common products. Since the "Dijkstra's Algorithm" propagates weights even for the customer pairs with no common products, we are able to calculate trust scores between the customers.

2.5.2 Filterer Module

The module initially takes list of unsorted and unfiltered transaction list and filters it to create more denser customer versus product matrix. Then creates a graph object by giving the customer versus product matrix as parameter and use the trust matrix created by the graph object to calculate recommendation coefficient for each product. Finally, for each user, the module sorts all products with respect to their recommendation coefficients in descending order and recommends top k of them.

$$w(c_1, c_2) = \frac{2 * sim(c_1, c_2) * trust(c_2)}{sim(c_1, c_2) + trust(c_2)} \quad (2)$$

Since the dataset contains implicit feedbacks, I didn't use a predictor algorithm for it. However, for the Movielens 100k dataset[3], (2) was used to predict the rating given by the c_{target} for item i .

$$p(i) = \frac{\sum_{c \in C} w(c_{target}, c) * r_c}{\sum_{c \in C} w(c_{target}, c)} \quad (3)$$

2.5.3 Recommender Module

This module sends the list of transactions to the `trust_based_filterer` and takes the list of recommendations for each customer in the database as a

response. The main task of this module is writing the recommendations to the database using neo4j driver.

2.6 Testing Phase

2.6.1 Methods

2.6.1.1 Leave-one-out Cross Validation

2.6.1.2 K-Fold Cross Validation

2.6.2 Results

2.6.3 Libraries I Used

2.6.3.1 Surprise

Installation :

```
1 pip install scikit-surprise
2
```

Sample Usage :

```
1 from surprise import AlgoBase, PredictionImpossible,
   Dataset
2 from surprise.model_selection import cross_validate
3
4 class Inverse_distance_weighted_tbr(AlgoBase):
5     ...
6
7     reader = reader = Reader(line_format='user item rating
   timestamp', sep=';', rating_scale=(1, 5))
8
9     data = Dataset.load_from_file('./dataset.csv', reader=
   reader)
10    algo = Inverse_distance_weighted_tbr()
11
12    cross_validate(algo, data, cv=5, verbose=True)
13
```

Listing 4: Surprise example

2.6.3.2 Matplotlib

Installation :

```
1 pip install matplotlib
2
```

Sample Usage :

```
1 import matplotlib.pyplot as plt
2 ...
3 plt.hist(eigentrust_list,
4 color = 'blue',
5 edgecolor = 'black',
6 bins = bins)
7 plt.title('Distribution of the Eigentrust Values')
8 plt.xlabel('Range')
9 plt.ylabel('Number of Eigentrust Values in the range')
10 plt.show()
11
```

Listing 5: Matplotlib example

3 Organization

3.1 METU Data Mining Research Group

4 Conclusion

References

- [1] AKSOY, K., ODABAS, M., BOZDOGAN, I., AND TEMUR, A. Tacorec. <https://senior.ceng.metu.edu.tr/2020/tacorec/>, 2020.
- [2] DIJKSTRA, E. W. Dijkstra’s algorithm. https://en.wikipedia.org/wiki/Dijkstra%27s_algorithm.
- [3] HARPER, F. M., AND KONSTAN, J. A. The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.* 5, 4 (Dec. 2015).

- [4] KAMVAR, S. D., SCHLOSSER, M. T., AND GARCIA-MOLINA, H. The eigentrust algorithm for reputation management in p2p networks. In *Proceedings of the 12th International Conference on World Wide Web* (New York, NY, USA, 2003), WWW '03, Association for Computing Machinery, p. 640–651.