

# The recommendation system of H&M products

Based on purchase history

Xuan zhang  
the college of sciences  
northeastern university  
Boston US  
zhang.xuan8@northeastern.edu

## ABSTRACT

This recommender system is based on content. In this article, I make an Item file and a User file to do cosine similarity. In general, the recommender system will calculate the cosine similarity between the User file and Item file. But, to H&M data, there are lots of products that are in this data. Thus, I use t-sne and a hierarchical cluster to observe the data and use k-meaning to do good classification. This way helps the system to find recommended products quicker. Finally, each customer can get 10 recommended products. And the average of RMSE from 100 customers is not over 0.12. RMSE is very low.

## CCS CONCEPTS

• Applied computing • Electronic commerce • Online shopping

## KEYWORDS

User file, Item file, hierarchical cluster, T-sne, k-mean cluster, the cosine similarity, RMSE, Precision and Accuracy

ACM Reference format:

## 1 Introduce the data

H&M is a big company, and it provides a huge dataset that includes products' pictures, customers' purchase history, product detailed metadata, etc. In this project, I only use two CVS files data.

There are two CSV files Transactions\_train.csv containing the purchases of each customer to date.

	t_dat	customer_id	article_id	price	sales_channel_id
0	2018-09-20	000058a12d5b43e67d225668fa1f8d618c13dc232df0ca...	663713001	0.050831	2
1	2018-09-20	000058a12d5b43e67d225668fa1f8d618c13dc232df0ca...	541518023	0.030492	2
2	2018-09-20	00007d2de826758b65a93dd24ce629ed96842531df6699...	505221004	0.015237	2
3	2018-09-20	00007d2de826758b65a93dd24ce629ed96842531df6699...	685687003	0.016932	2
4	2018-09-20	00007d2de826758b65a93dd24ce629ed96842531df6699...	685687004	0.016932	2

Articles.csv contains detailed metadata for each article\_id available. In the articles.csv, there are nine numerical attributes that were selected to do the K-mean cluster and calculate cosine similarity. There are product\_type\_no, graphical \_ appearance \_ no, colour \_group \_ code,perceived\_colour\_value\_id,perceived\_colour\_master\_id,d epartment\_no,index\_group\_no,section\_no,garment\_group\_no. These nine attributes will distinct products of H&M from different aspects which include colour, type, graphs and garment etc.

article_id	product_code	product_type_no	graphical_appearance_no	colour_group_code	perceived_colour_value_id	perceived_colour_master_id	department
0	108775015	108775	253	1010016	9	4	5
1	108775044	108775	253	1010016	10	3	9
2	108775051	108775	253	1010017	11	1	9
3	110065001	110065	306	1010016	9	4	5
4	110065002	110065	306	1010016	10	3	9
5	110065011	110065	306	1010016	12	1	11
6	111565001	111565	304	1010016	9	4	5
7	111565003	111565	302	1010016	13	2	11
8	111565001	111565	273	1010016	9	4	5

## 2 Methods and algorithm

### 2.1 Introduction about assumption and definition

In this project, I assume that some potential customers in the same area (like Boston) once bought H&M group's products. They may browse H&M group's website and APP. When they browse its website and APP, we will provide some recommendations about H&M group's products. These products are being stored in warehouses which are around Boston. All products are seasonal clothes. Thus, I just use a part of customers' information and a part of the products' information from this data which provides all customers' purchase history and product information.

Making a good Item file and User file is very important for us to do a good recommendation to customers. I wash the data and analysis its attributes of the data. For each attribute of customers, I will do a weight depending on purchase history.

## 2.4 Cosine similarity

$$\text{Formula: } \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

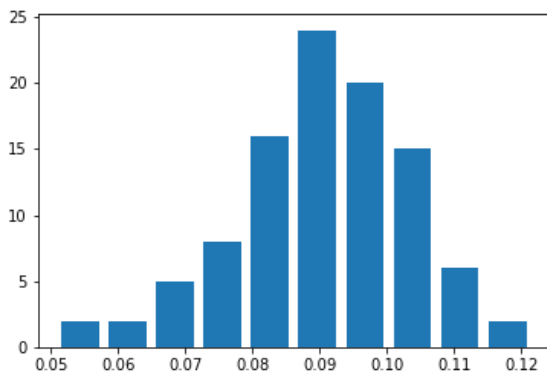
I use K-mean to group one customer from the User file. Then Use the customer data to calculate cosine similarity with the Item file and find ten similar items and return these items' id.

## 2.5 analysis result about RMSE

To result in the similarity of the User file and Item file, we use RMSE to calculate the difference. I select 100 customers; each customer is recommended ten products. Then I calculate the average of RMSE by the value of the customer and the value of the ten products. I choose the mean of each customer's RMSE.

$$\text{Formula: } RMES\_mean\_eachcustomer = \sqrt{\sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{100n}}$$

I use python to do a histogram to show the distribution of the value of RMES\_mean\_eachcustomer.



From this table, we can find lots of the average of RMSE is about 0.09. It means the difference between attributes of recommendations and attributes of customers is very low. I get a good result from validation.

## 3 Procedure

1. Observe data and select useful attributes from the data
2. Making an Item file and a User file
3. Using T-sne and the hierarchical cluster to observe the distribution of the data, then using the K-mean to do classification.

To a special customer, I select this customer's data from the User file and classify this customer into one special group, then calculate cosine similarity between items in this group and this customer and recommend products that are the biggest similarity ten.

4. Calculate RMES

## 4 Further research

We can try to use collaborative filtering in this project. The

collaborative filtering will be based on the User file and Item file to fill out empty of these two files. I think it will make this model of this project better. In addition, there are some parts of the H&m data I did not use in this project. It is possible to extract useful features from the products' pictures. It will help this model of this project to do a better classification.

## ACKNOWLEDGMENTS

Thank Professor Eliassi-Rad Tina for her help of the validation of the model and guidance of the recommender system.

## REFERENCES

- [1] Jure Lekovec, Anand Rajaraman, Jeffrey D.Ullman,2014. Mining of Massive datasets. DOI: <http://infolab.stanford.edu/~ullman/mmds/ch9.pdf>.