# Product Matching between Zalando and Aboutyou: CA684 Machine Learning Assignment 2022

Khizer Ahmed Biyabani
*School of Computing*
*Dublin City University*
Dublin, Ireland
khizer.biyabani2@mail.dcu.ie

*Abstract*—Consumers may now buy things from tens of thousands of online stores. The completeness of product descriptions and the taxonomies used to organize the items, on the other hand, varies amongst e-shops. Approaches for product integration on the Web are needed to improve the user experience, such as enabling for easy comparison of offers from multiple suppliers. The following report presents a method for implementing a product matching algorithm for matching products for the fashion company Zalando with that of the Aboutyou using different classification approach in machine learning.

*Index Terms*—Product Matching, Text pre-processing, Text Similarity Distance Measures, Levenshtein distance, Damerau–Levenshtein distance, Hamming Distance, Token Set Ratio, Fuzzy string matching, Jaro and Jaro-Winkler, DummyClassifier, KNeighborsClassifier, XGBClassifier, DecisionTreeClassifier, RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier, Perceptron, MLP, Confusion Matrix, F1 Score.

## I. INTRODUCTION

The e-commerce sector, such as Amazon, Flipkart, and many other leading online purchasing sites, has recently experienced a huge upswing. With the rise in popularity of online selling platforms, each and every one of these online retailers require a technique for comparing their items to those of their competitors.

Zalando is a well-known European online fashion and leisure platform that brings together customers, businesses, and partners. Many of Zalando's rivals sell the same brand of items; for example, 'Aboutyou' is a fast-growing fashion platform and one of Zalando's competitors.

The following article presents a method for comparing various items between Zalando and Aboutyou, and it matches the common products between the two stores. Detailed information regarding the product data for the two stores for which the match must be discovered may be found throughout the text. To anticipate the match, a number of approaches have been developed, one of which is the text-similarity measure between the product color, title, and description of each of the two shops. We may determine the proximity or extent of each product's relationship using these similarities, as discussed in the next sections of this report.

We then use the similarity metrics to determine whether there is a match and train various machine learning classification models such as Dummy Classifier, KNearest Neighbours, and many more. The performance of each model is then evaluated in order to determine which model is the best fit. Finally, we make predictions about the matches in the provided test set, which is our final result.

### A. About the Dataset

The dataset was given by Zalando and it was a licensed dataset, the dataset contains the following files:

1) offer_train.parquet : This file is the training set. It cotains the product offers of each of the zalando and aboutyou products. This dataset is neede fo the training of the model. It contains 102884 product offers and 10 attributes or columns of each. Out of which 61980 are zalando offers and 40904 are anoutyou offers.

2) offer_test.parquet : The final test set contains the offers of which matches has to be predicted using our trained model. It is same as the training set, it has 106741 values in which 70105 and 36636 are zalando and aboutyou offers respectively.

Each of the above datasets contain the following attributes:

- offer_id : unique identifier for an offer of a product (i.e. a product x shop combination, where we don't know the product component).
- shop: "zalando", "aboutyou".
- lang: "de" (German).
- brand: brand of each of the shop for example 'Nike' it is different for each of the shop.
- color: Color of the product e.g. "blue" - note that there could be more than one and different 'shop's might have different 'brand' nomenclature ("ocean", "light blue", "...") and may have more than one color (ordering matters).
- title: The actual title of the products e.g. "White Nike tennis top".
- description: a long product description that can may contain material composition, size and different attributes for the. products, cleaning instructions, etc.
- price: price in euro without any discount.
- url: url of the product description page.

- image_urls: list of product images such as stock photo, with model, lifestyle photo, or close up.

3) match_training.parquet : This file contains the ground truth matches present in the training dataset.
   The dataset has the following attributes:
   - zalando: offer_id from "zalando" shop.
   - aboutyou: offer_id from "aboutyou" shop.
   - brand: unique identifier for the brand representing the match.

## II. PROPOSED METHODOLOGY

### A. System Architecture

The following diagram gives the summary of the entire process adopted :



Fig. 1. System Architecture: The process flow diagram

The further sections discuss in detail about each and every step in the above flow diagram.

### B. Loading the Dataset:

With the help of pandas library of python we can read the parquet file of our datasets. The following snippet in figure 2 shows the first five rows of the training set.



Fig. 2. First five rows of the training data.

### C. Pre-processing of Data

Data pre-processing is the process of converting raw data into a format that can be understood. Data in the real world is frequently partial, inconsistent, redundant, and loud. Data preparation entails a number of stages that aid in the conversion of unprocessed data into a usable state.

This mainly includes the data cleaning, the raw data that we have is usually corrupt and this data needs to processed before giving as an input to the machine learning algorithm. Both the training and test data contains certain anomalies, such as NoneTypes or there are certain values missing in various columns of the dataset.

*1) Removing NoneTypes::* In Python, the data type of an object is NoneType when the object has no value. The NoneType object can be started with the keyword None. There is a function created to clean the NoneTypes from the daatset as shown below:



Fig. 3. Python function for NoneType Removal

*2) Treating Null or Missing Values::* As you can see in the image below there certain missing or null values present in the dataset, we can find that the color column contains 2 missing values, these values can be replaced with an empty string (""). The missing values in the image_urls, are also there but as we are working on the text so we can ignore these values.



Fig. 4. Treating Missing Values

*3) Changing the text to lowercase::* For fast processing of text for color, title and description we standardise the case of the text to lower, this make the text uniform for the shop offers.

## D. Preparing the Training data using Ground Truth Matches

As we have the ground-truth matches for the training from that we can prepare a dataset that contains the matched offers of both zalando and aboutyou along-with the attributes of each of these product offers, i.e. the color, title and description. These matched offers we can combine with the unmatched offers. Also, there is column called 'match' with value 1 for a matched offer and 0 for an unmatched offer. This dataset looks like the following in the image below:

Fig. 5. The matched and unmatched offers training set

## E. Measuring text similarities and preparing the Training set

Now, the dataset shown above in figure 5 contains only the textual information but the machine learning model needs some numerical measures with help of which it can classify the match as 0 or 1. Hence we calculate the text similarity score each title, color and description of the offers of both the shops. There are several similarity measures [4] that has been used in this assignment some of which include:

*1) Levenshtein distance:* The Levenshtein distance is a measure of word similarity. The distance between two words indicates how many modifications are required to change one word into another. Editing can be accomplished using one of three methods: Insertion, deletion and substitution. [1]

Another version of this is the 'Damerau–Levenshtein distance' which is the smallest number of operations (consisting of insertions, deletions, or replacements of a single character, or transposition of two adjacent characters) necessary to turn one word into the other.

*2) Hamming Distance:* It is the distance is a measurement of how different two bitstrings of the same dimension are. The classic definition of a Hamming distance is as follows: given two bitstrings (more broadly, strings of symbols) of the same dimension, the Hamming distance is the smallest number of symbol modifications required to turn one string into the other. [2]

*3) Token Set Ratio:* The strings are tokenized and pre-processed in token sort ratio by converting to lower case and removing punctuation. After that, the strings are alphabetically sorted and put together. Following that, the Levenshtein distance similarity ratio between the strings is determined.

There are other Fuzzy string matching techniques like matching the sub-string of the two strings which is the 'Partial Ratio' and various other methods have been used to find the similarity score. [3] Apart from that, other similarity measures such as Jaro and Jaro-Winkler similarity, and even matching numbers in both string is also used. Python's 'JellyFish' and 'Fuzzywuzzy' has inbuilt libraries for all these similarity scores.

Finally, the dataset created with the text similarities and dropping all the unnecessary columns is shown in the figure below:

Fig. 6. The Final Training set: Contains different text similarity measures.

The correlation with each of the similarity measures can be seen in the following figure. It is evident that maximum similarity measures are positively correlated with the match.

Fig. 7. The Correlation between Similarity Measures and the Match

## F. Building the Classification Model

After preparing the training dataset with similarity scores of each of the offers from zalando and aboutyou. The training data is then splitted up into train and test set. The train set for training the model and the test set for the validation purpose and results evaluation.

Fig. 8. Splitting the Training set into train and validation or test set

There several classification algorithms in machine learning that can be used. This assignment uses the following classifiers to classify the matches:

1) DummyClassifier
2) KNeighborsClassifier
3) XGBClassifier
4) DecisionTreeClassifier
5) RandomForestClassifier
6) AdaBoostClassifier
7) GradientBoostingClassifier

8) Perceptron
9) MLP

Each of the above models have been fitted with the training and the testing or validation set and the best model is chosen among them.

## III. EVALUATION AND RESULTS

There are several evaluation metrics that has been used and based on which each classification model has been compared. [4]

1) The Confusion Matrix is a performance metric for machine learning classification tasks with two or more classes as output. As illustrated in the figure below, it's a table with a mix of expected and actual data.

Fig. 9. Confusion Matrix: General Form

We can determine Recall, Precision from the above matrix.

2) Recall : Recall is the proportion of relevant documents that has been retrieved.
$$Recall = \frac{\text{True Positive retrieved}}{\text{True Positive + False Negative}}$$

3) Precision: It is the proportion of documents retrieved that is relevant.
$$Recall = \frac{\text{True Positive retrieved}}{\text{True Positive + False Positive}}$$
One must note that both precision and recall are inversely proportional to each other.

4) F1 Score: It gives the combined idea about precision and recall.
$$F1 - Score = 2 * \frac{\text{Precision * Recall}}{\text{Precision + Recall}}$$

Now that we have a better understanding of the evaluation metrics, we can use them to assess the prediction outcomes of each of our categorization models. The evaluation of each of the above-mentioned trained classification models is summarized in the table shown in figure 10.

It can be observed that XGBooster Classifier outperforms the other algorithms, with F1 score of 96.59% than the other algorithms. Hence, XGB Classifier can be used to predict the matches in the given testing dataset. The confusion matrix shown below explains the predictions done by the XGB Classifier, it clearly shows that only 400 products have been predicted wrong, which us quite a good result.

Fig. 10. Confusion Matrix resembling evaluation of each classification model.

Fig. 11. Confusion Matrix for XGBoost Classifier

Also, we can note that the the F1 score for match is 96% and that for no match is 97% percent from the classification report below.

Fig. 12. Classification Report for XGBoost Classifier

Thus, we can use XGBooster to predict the matches for the given test set.

## IV. WORKING ON THE TEST SET: PREDICTING THE UNSEEN MATCHES

As we mentioned in section I-A we have a offers_test.parquet file for which the matches has to be predicted. The dataset is loaded the same as above for the training set and it is pre-processed and the similarity for each of the offers in the test set.

The following is the snippet of the test dataset we loaded:

Fig. 13. The Test data

The final predictions is done on the above test set and the matches have been predicted and even verified as shown in the images 14 and 15 below:
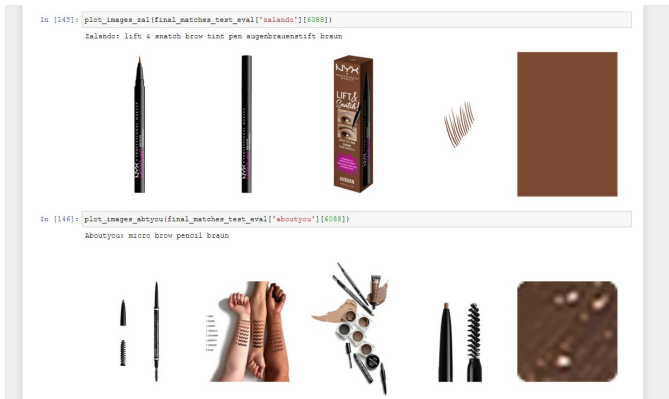
Fig. 14. Predictions done on Test data



Fig. 15. An example of a match in test data

## V. CONCLUSION

To summarize, we performed a supervised machine learning job for both zalando and aboutyou by categorizing matched and unmatched products.We started with training the model with the similarity scores of the text strings from the title, color and description of each of the products of the two shops, and classified them using different classification algorithms. Finally, we then predicted the unseen matches present in the given test data.

Apart from this, there were several difficulties encountered during the procedure, such as memory concerns and a time-consuming difficulty due to the big datasets. But still there were predictions done successfully. There is still more that can be done when matching priducts, we can consider the images also, while matching but this assignement was mainly text focussed,

## REFERENCES

[1] S. Zhang, Y. Hu and G. Bian, "Research on string similarity algorithm based on Levenshtein Distance," 2017 IEEE 2nd Advanced Information Technology, Electronic and Automation Control Conference (IAEAC), 2017, pp. 2247-2251, doi: 10.1109/IAEAC.2017.8054419.
[2] Bookstein, Abraham Kulyukin, Vladimir Raita, Timo. (2002). Generalized Hamming Distance. Information Retrieval. 5. 10.1023/A:1020499411651.
[3] Fuzzy String Matching – A Hands-on Guide by Mimi Dutta — July 15, 2021. https://www.analyticsvidhya.com/blog/2021/07/fuzzy-string-matching-a-hands-on-guide/
[4] Metrics to Evaluate your Classification Model to take the right decisions by Sumeet Kumar Agrawal — July 20, 2021 https://www.analyticsvidhya.com/blog/2021/07/metrics-to-evaluate-your-classification-model-to-take-the-right-decisions/
[5] Product matching algorithms find identical products on ecommerce sites so users can compare products and retailers can compare prices, by Matt Clarke, Saturday, March 13, 2021. https://practicaldatascience.co.uk/machine-learning/how-to-create-a-product-matching-model-using-xgboost