

## 1 My Code

The code for this solution is at: <https://github.com/anichols13/CS4400FinalProject>

## 2 References

- <https://pypi.org/project/py-entitymatching/>
- [http://pradap-www.cs.wisc.edu/magellan/how-to-guide/how\\_to\\_guide\\_magellan.pdf](http://pradap-www.cs.wisc.edu/magellan/how-to-guide/how_to_guide_magellan.pdf)
- [https://nbviewer.jupyter.org/github/anhaidgroup/py\\_entitymatching/blob/master/notebooks/guides/end\\_to\\_end\\_em\\_guides/Basic%20EM%20Workflow%20Restaurants%20-%202021.ipynb](https://nbviewer.jupyter.org/github/anhaidgroup/py_entitymatching/blob/master/notebooks/guides/end_to_end_em_guides/Basic%20EM%20Workflow%20Restaurants%20-%202021.ipynb)
- [https://nbviewer.jupyter.org/github/anhaidgroup/py\\_entitymatching/blob/master/notebooks/guides/end\\_to\\_end\\_em\\_guides/Basic%20EM%20Workflow%20Restaurants%20-%202021.ipynb](https://nbviewer.jupyter.org/github/anhaidgroup/py_entitymatching/blob/master/notebooks/guides/end_to_end_em_guides/Basic%20EM%20Workflow%20Restaurants%20-%202021.ipynb)
- <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

## 3 Solution Outline

The solution includes five steps: (1) Data reading and EDA, (2) Blocking, (3) Feature engineering, (4) Model training and (5) Generating output.

### 3.1 Data reading and EDA

In this step, we read the left table and right table, as well as the training set. We explore the dataset to get some ideas of designing the solution. For example, we found that the left table has 2554 rows and the right table has 22074 rows, so there are  $2554 \times 22074 = 56376996$  pairs. Examining every pair is very inefficient, so we will need a blocking step to reduced the number of pairs that we will work on.

### 3.2 Blocking

We perform blocking on the attribute "brand", generating a candidate set of id pairs where the two ids in each pair share the same brand. This is based on the intuition that two products with different brand are unlikely to be the same entity. Our blocking method reduces the number of pairs from 56376996 to 280723. To conduct my blocking, I used the `OverlapBlocker()` function from the entitymatching module to block on the attribute "brand."

### 3.3 Feature engineering

Through importation of the entitymatching module, I was able to utilize features when looking at the product's title. This import of py-entitymatching granted me access to features such as `get_features()` and `extract_feature_vectors()` to eventually discover the tuple pairs between two tables that refer to the same entities. Please refer to the citation section to see the entitymatching user guide.

### 3.4 Model training

Similar to the sample solution, I used the random forest classifier. I then established a prediction for the candidate set matches by training the model through the dataframe training that was implemented. Since the number of non-matches is much more than the number of matches in the training set, we set `class_weight="balanced"` in random forest to handle this training data imbalance problem. We perform prediction on  $X_c$  to get predicted labels  $y_c$  for the candidate set.

### 3.5 Generating output

The pairs with  $y_c = 1$  are our predicted matching pairs  $M$ . We remove the matching pairs already in the training set from  $M$  to obtain  $M^-$ . Finally, we save  $M^-$  to output.csv.