
Otto Group Product Classification Challenge

Classify products into the correct category

By

ZHUFENG LI, HAO LIU



Comprendre le monde,
construire l'avenir



Computer Science Department of the Faculty of Science

UNIVERSITY PARIS SUD

NOVEMBER 5, 2018

OTTO GROUP PRODUCT CLASSIFICATION

1. Introduction

For this competition, we are provided a dataset with 93 features for more than 200,000 products. The objective is to build a predictive model which is able to distinguish between our main product categories. Submissions are evaluated using the multi-class logarithmic loss. Each product has been labeled with one true category. For each product, you must submit a set of predicted probabilities (one for every category). The formula is then,

$$\text{logloss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij})$$

where N is the number of products in the test set, M is the number of class labels, \log is the natural logarithm, y_{ij} is 1 if observation i is in class j and 0 otherwise, and p_{ij} is the predicted probability that observation i belongs to class j . More description on the following link

2. Data analysis

We will study the correlation between all these variables. We plotted the correlation matrix which is a table showing correlation coefficients between sets of variables. In this matrix, +1 is the case of a perfect direct (increasing) linear relationship (correlation), -1 is the case of a perfect decreasing (inverse). We also did principal component analysis on all the features. A distribution of variance of all the principal components is given below. The total variance is the sum of variances of all individual principal components.

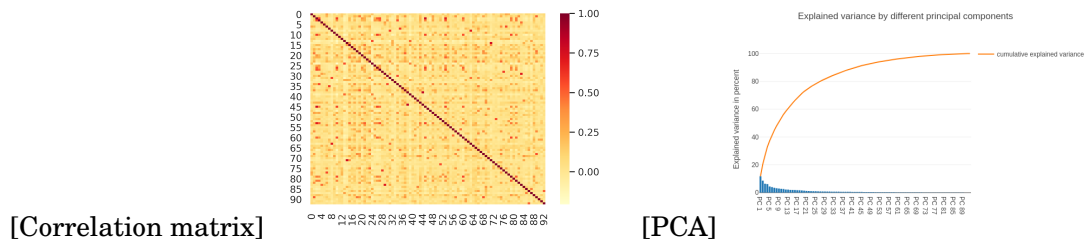


Figure 1: Figures side by side

3. Machine learning models

Several different methods of machine learning are used in this project. They are decision tree, random forest, naive bayes, k nearest neighbors, support vector machines, and neural network. After testing and fine tuning hyperparameters, we benchmarked their performances with kaggle submission score. As it shows in the table below TABLE 1, a finely tuned neural network gave us the best score w.r.t the others. Submissions made by k nearest neighbors are worse than others in general. The results of SVM are acceptable. However due to the complexity of SVM, it took quite a while to predict p_{ij} . Decision tree performed very poorly on this challenge. Because, it can only provide a categorical prediction instead of a probability. Thus, it has a huge log loss. Random forest gave a very good score, especially when it comes to computational efficiency.

Comparison of performance of models			
Model	Score	Model	Score
Neural network (Complex ver.)	0.46422	KNN(100)	0.69295
Random forest	0.47818	SVM(poly)	0.78819
Neural Network (Simple ver.)	0.48763	SVM(linear)	0.78844
SVM(rbf)	0.57218	KNN(1000)	0.81102

TABLE 1. Comparison of performance

4. Results

11 submissions for Hao LIU		Sort by Private Score	
All	Successful	Selected	
Submission and Description	Private Score	Public Score	Use for Final Score
submission.csv a few seconds ago by Hao LIU Neural Nets	0.46958	0.46422	<input type="checkbox"/>
RFsub.csv 5 days ago by Hao LIU 300 trees RF	0.48018	0.47818	<input type="checkbox"/>
RFNormsub.csv 5 days ago by Hao LIU RF with norm	0.48264	0.48055	<input type="checkbox"/>
benchmark.csv 13 days ago by Hao LIU With random forest model	0.48288	0.47988	<input type="checkbox"/>
SVM_submission.csv 5 days ago by Hao LIU SVM with proba	0.57124	0.57218	<input type="checkbox"/>

FIGURE 2. Results from submission board