

Fake Review Classification

Anomaly Detection – Challenge 2

Team Pauliguel



Index

- Dataset description
- Undersampling
- Pre-processing
- Models and evaluation
- Tuning
- Kaggle submission

- Data reading problems due to errors in the test reviews file.
- Total of 2908 instances.
- Unbalanced class representation
 - Positive class is 13% of the instances.

- To solve the unbalanced data.
- Random sample of the most common class.
- 784 instances in the training set.

-

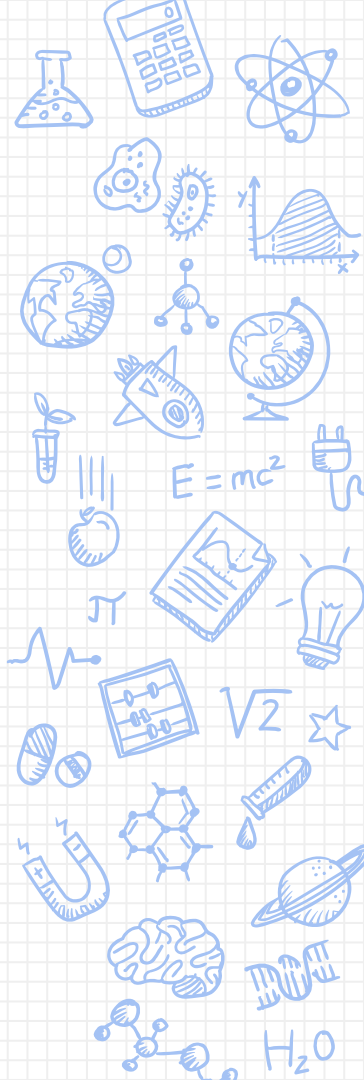
Feature Extraction

- Features from the paper

Text mining

- Length review (characters and and words)
- Bag of words

Personal pronouns	i, we, me, us, my, mine, our and ours.
Associated actions	Went and feel.
Targets	Area, options, price and stay.
Emotion words	Nice, deal, comfort and helpful.



- # Behavioral

- $$MNR(reviewers_n) = \frac{|reviews(reviewers_n)|}{daysBetween(joinDate(reviewers_n), date(mostRecentReview))}$$

- $$PR(reviewer_n) = \frac{|positiveReviews(reviewer_n)|}{|reviews(reviewer_n)|}$$

- ## Behavioral

- $$RD(review_n) = rating(review_n) - rating(hotel(review_n))$$

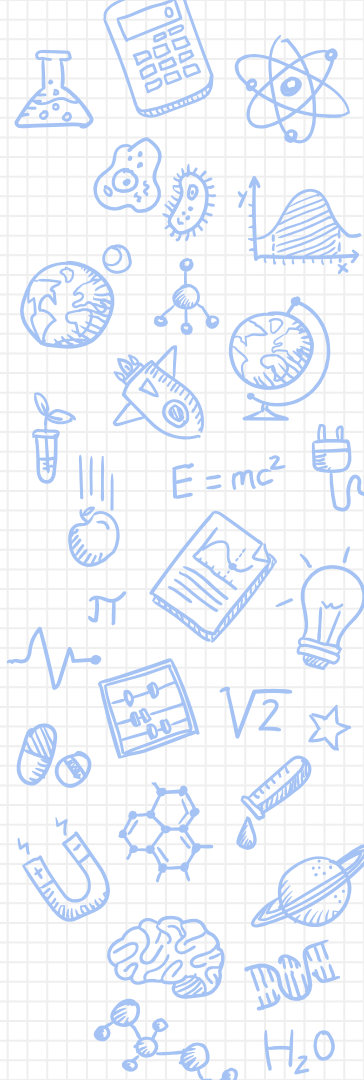
- Maximum Content Similarity (MCS)
 - For each Reviewer gather all the Reviews
 - Compute the average of the cosine similarity between each pair of reviews.

Feature Extraction

- Other features
 - Mark fake users (users with at least one fake review)

- rating_review
- usefulCount_review
- coolCount_review
- funnyCount_review
- friendCount_user
- reviewCount_user
- firstCount_user
- usefulCount_user'

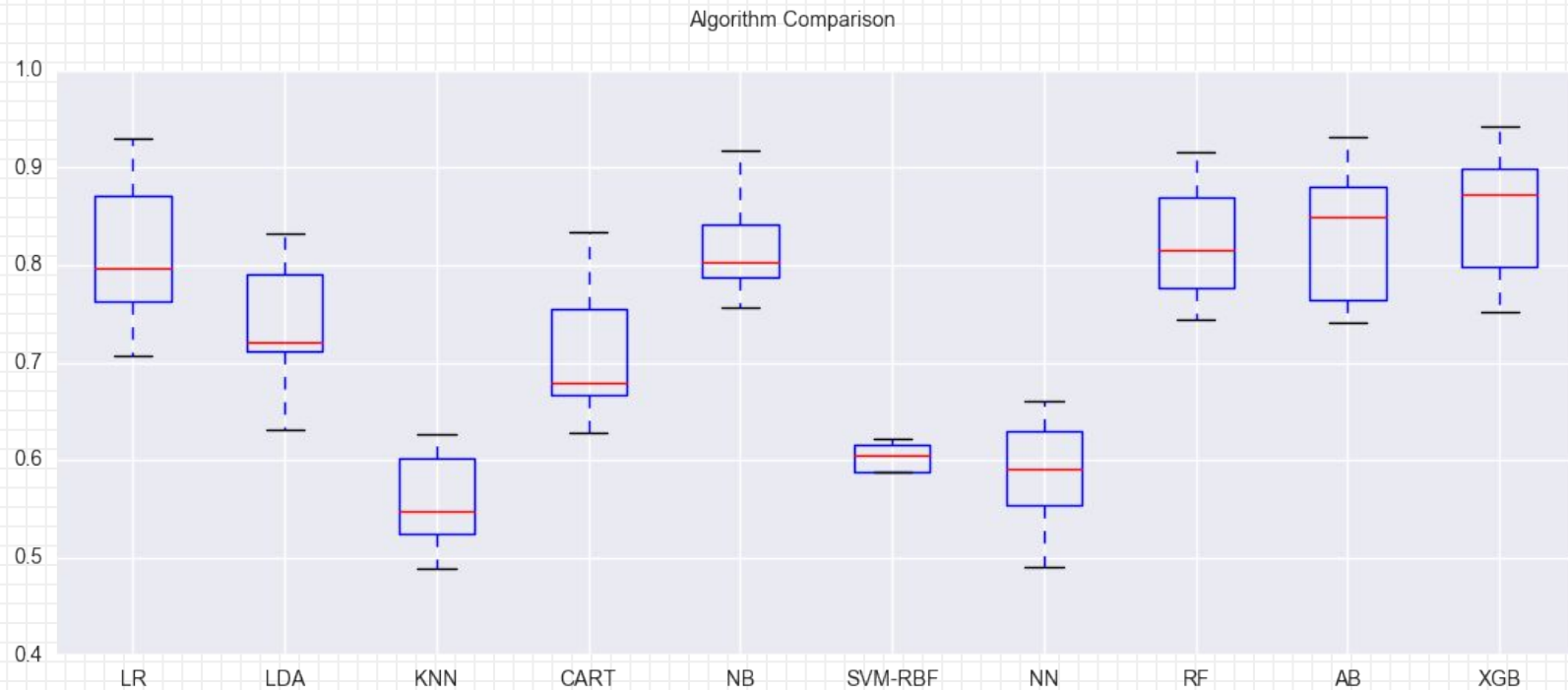
- coolCount_user
- funnyCount_user
- complimentCount_user
- tipCount_user
- fanCount_user
- rating_hotel
- filReviewCount_hotel'



Models used:

- Logistic Regression
- Linear Discriminant Analysis
- K-Nearest Neighbors Classifier
- Decision Tree Classifier (CART)
- Naïve Bayes
- Support Vector Machines (with both “poly” and “rgf” kernels)
- Multilayer Perceptron
- Random Forest
- AdaBoost
- Gradient Boosted Decision Trees (XGBoost)

Results with AUC



- “number of estimators”,
- “criterion”
- “maximum number of features”.

new score = 0.858400

- XGBoost
 - No increase in the ROC-AUC metric

Kaggle submission

0.78924

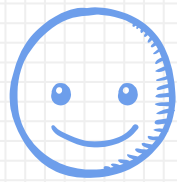
Public score - best submission

0.79813

Private score - best submission

We basically quit social life this week to improve our score.

Maybe more luck next time?



THANK YOU!

Any questions?