# Restaurant Review Rating Prediction

## Multiclass Classification Using Spark MLlib

Boris Korotkov
SCS-3252-008

# Objective
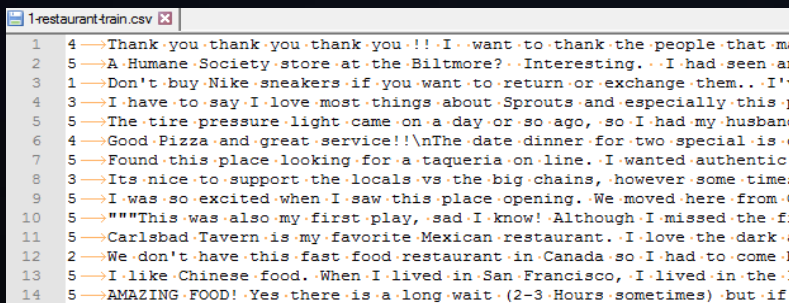
- Predict Review Rating

# Data



82, 065
records

# Feature Engineering

- Tokenizer

- StopWordsRemover

- HashingTF

- IDF (Inverse document frequency)

# Model selection and training

- 80% Training / 20% Test

- Random Forest Classifier

- Cross Validation training
  - 5 folds
  - 2 parameters: Max Depth & NumTrees

# Model evaluation

- MulticlassClassificationEvaluator
- The accuracy on Test set is 35.08%

# Improvement opportunities

- Additional feature engineering
- Different models
- Advanced hyperparameter tuning