

MINI-PROJECT 2: YELP BUSINESS RATING PREDICTION USING TENSORFLOW

Version 1.0

CSC215, Fall 2018

Oct 8th, 2018

Prepared by: Chandini Nagendra

Siddharth Chittora

Contents

1. Problem Statement.....	3
2. Methodology.....	3
2.1. Regression Problem	4
2.2. Classification Problem.....	4
3. Experimental Results and Analysis.....	5
3.1. Regression.....	5
Analysis and Result:	6
3.2. Classification	7
Analysis and Result:	7
Conclusion:.....	7
4. Task Division.....	8
4.1. Chandini Nagendra:	8
4.2. Siddharth Chittora.....	8
5. Project Reflection.....	8
6. Additional Features.....	9

1. Problem Statement

In this project, we aim to predict a business's stars rating using the reviews of that business and review count based on neural network implementation in Tensorflow. This project is twofold:

Task 1: Consider this problem as a regression problem. Compare the RMSE of the BEST **Tensorflow regression neural network model** you obtained with that of **regression model** you achieved in the last project.

Task 2: Consider this problem as a classification problem. Compare the accuracy of the BEST **Tensorflow classification neural network model** you obtained with that of **each classification model** you achieved in the last project.

2. Methodology

Here we compare Linear and logistic models with Tensor flow models by using early stopping, Model checking and tuning the models with hyperparameters and see how they affect performance.

- we are using data for 10000 businesses.
- For regression problem we worked with linear regression and logistic regression.
- we created tensorflow model with activation function ReLU to compare the best RMSE of the earlier regression models with this model.
- we used Early Stopping and checkpointing with ReLU to see how it affected the model.
- now we tried different optimizers like, adam, SGD, RMSprop, Adagrad, Adamax, Adadelat, Nadam.
- we also experimented with multiple hidden layers and the number of neurons. The experimental results are shown below.
- after learning that the model gave best performance with optimizer Adam and four hidden layers, we then experimented with activation function Sigmoid and Tanh.
- Tanh is also like logistic sigmoid. The range of the tanh function is from (-1 to 1). tanh is also sigmoidal (s - shaped) hence we used the zscore normalized review count for Tanh
- For Classification problem we implemented KNN, SVM, MNB.
- we created tensorflow model for classification model using activation function Softmax.
- we used optimizer adam and four hidden layers as we learnt from our experiment that they give the best result.

2.1. Regression Problem

Model & Tuning	RMSE	R2 Score
Linear Regression	0.56	0.70
Logistic Regression	1.38	0.54
Tensor flow regression neural network models		
ReLU without stopping & checkpoint	0.59	0.64
ReLU with stopping & checkpoint + adam	0.52	0.72
ReLU with stopping & checkpoint + sgd	0.49	0.74
ReLU with stopping & checkpoint + RMSprop	0.51	0.73
ReLU with stopping & checkpoint + Adagrad	0.51	0.73
ReLU with stopping & checkpoint + Adadelta	0.50	0.74
ReLU with stopping & checkpoint + Adamax	0.50	0.74
ReLU with stopping & checkpoint + Nadam	0.52	0.72
ReLU with stopping & checkpoint + Adamax + 2 hidden layers	0.52	0.72
ReLU with stopping & checkpoint + Adamax + 3 hidden layers	0.51	0.72
ReLU with stopping & checkpoint + Adamax + 4 hidden layers	0.51	0.73
ReLU with stopping & checkpoint + Adamax + 5 hidden layers	0.51	0.73
Sigmoid without stopping & checkpoint	0.51	0.73
Sigmoid with stopping & checkpoint + adam	0.49	0.78
Tanh without stopping & checkpoint	0.63	0.60
Tanh with stopping & checkpoint +adam	0.47	0.75

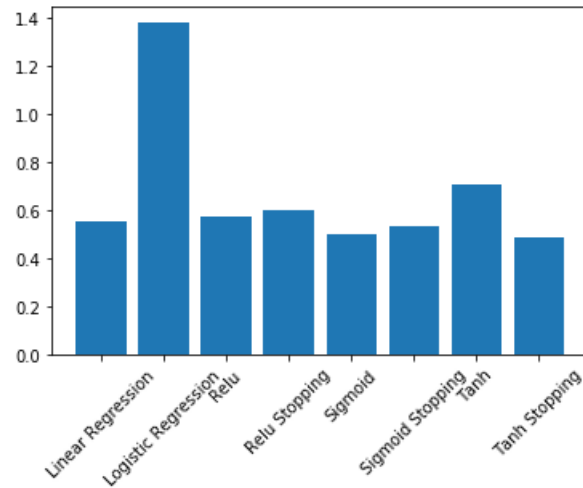
2.2. Classification Problem

Models & Tuning	Accuracy	Precision	Recall	F1 Score
KNN	0.295	0.296	0.295	0.291
SVM	0.436	0.430	0.436	0.427
MNB	0.332	0.299	0.332	0.283
Tensorflow Classification				
ReLU + adam + early stopping and checkpoint + 4 hidden layers	0.486	0.475	0.486	0.474
Sigmoid + adam + early stopping and checkpoint + 4 hidden layers	0.494	0.473	0.491	0.480
Tanh + adam + early stopping and checkpoint + 4 hidden layers	0.491	0.491	0.494	0.485

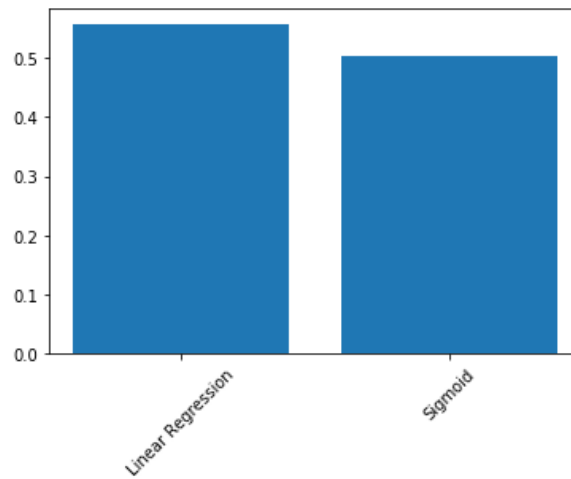
3. Experimental Results and Analysis

3.1. Regression

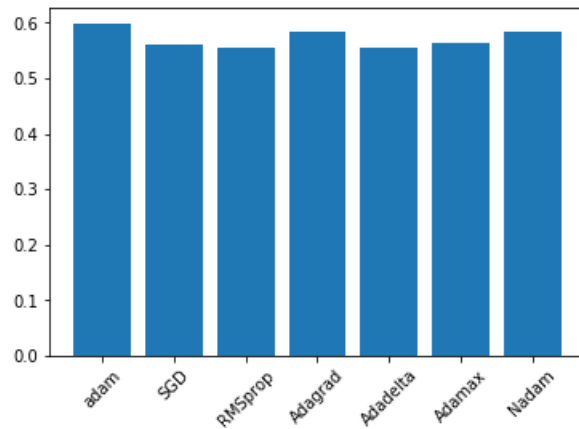
- Comparison of RMSE between all the regression models



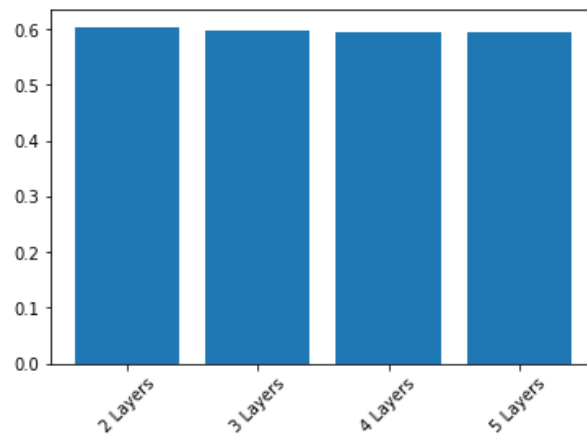
- Comparison of RMSE between the best Tensorflow model with the best Classical Regression model



- Comparison of RMSE between different optimizers



- Comparison of RMSE between models with different number of hidden layers

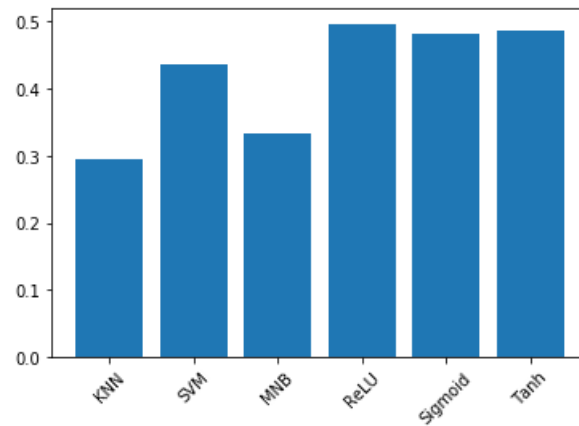


Analysis and Result:

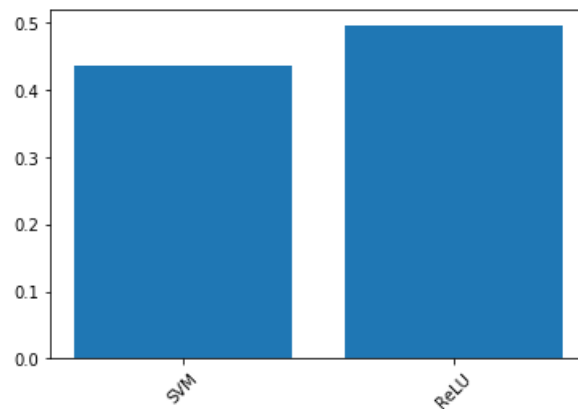
From our experiments for regression problems we observed that Sigmoid with early stopping and checkpointing with optimizer adam and 4 hidden layers had the least RMSE. Linear regression had the next best RMSE score.

3.2. Classification

- Comparison of accuracy between all the classification models



- Comparison of accuracy between best classification models



Analysis and Result:

From our experiments for classification problem we observed that ReLU with early stopping and checkpointing with optimizer adam and 4 hidden layers had the best accuracy. SVM had the next best accuracy.

Conclusion:

From our experiments we observed that regression model best fit the given problem statement.

4. Task Division

4.1. Chandini Nagendra:

- Comparison of Linear regression model with tensorflow model
- Report

4.2. Siddharth Chittora

- Comparison of Classification Model with tensor flow model
- Report

Discussed together on how to improve the model and came up with the solution discussed in the additional features section.

5. Project Reflection

- In Mini project 1, we were extracting only the primary category from each row in the categories column. In this project we are using multilabel binarizer to extract all the categories.
- when using tensorflow for classification, we found that if label encoded stars are used as output directly, throws dimension error. So we one hot code the label encoded stars and use this as our output.
- If we use the same best weight HDF5 file to save the best weight for multiple models it does not give accurate results hence we have to reinitialize and recreate the file each we run a model.
- Tanh is also like logistic sigmoid. The range of the tanh function is from (-1 to 1). tanh is also sigmoidal (s - shaped) hence we used the zscore normalized review count for Tanh

6. Additional Features

- We implemented the L1/L2 regularization for the original dataset.
- We also implemented Dropout to reduce the effect of overfitting.
- We processed postal code, we performed One hot Coding on it to extract features.
- we processed categories using Multilabel Binarizer to extract features.
- we merged it with the original matrix and used this matrix.
- we used activation function ReLU and Sigmoid on the new matrix we created with postal code and categories.
- we also performed regularization and dropout on the new matrix.

Model & Tuning	RMSE	R2
Regularization	0.55	0.71
Dropout	0.52	0.74
Regularization + additional features	0.59	0.66
Dropout + additional features	0.61	0.63
Relu + additional features	0.59	0.65
Sigmoid + additional features	0.59	0.65

