Jason Long

Dr. A. Jafari

DATS6203: Machine Learning II

29 April 2019
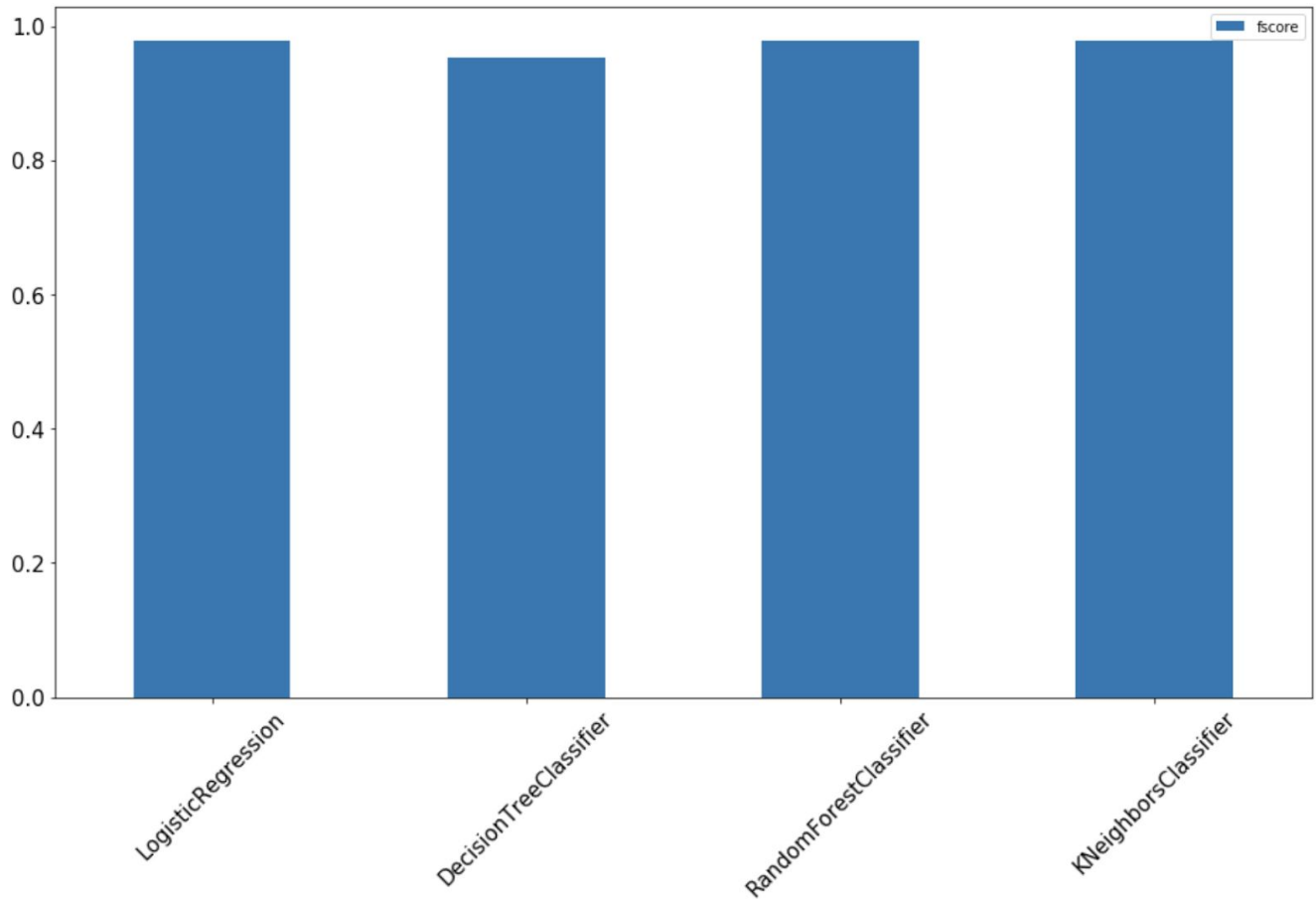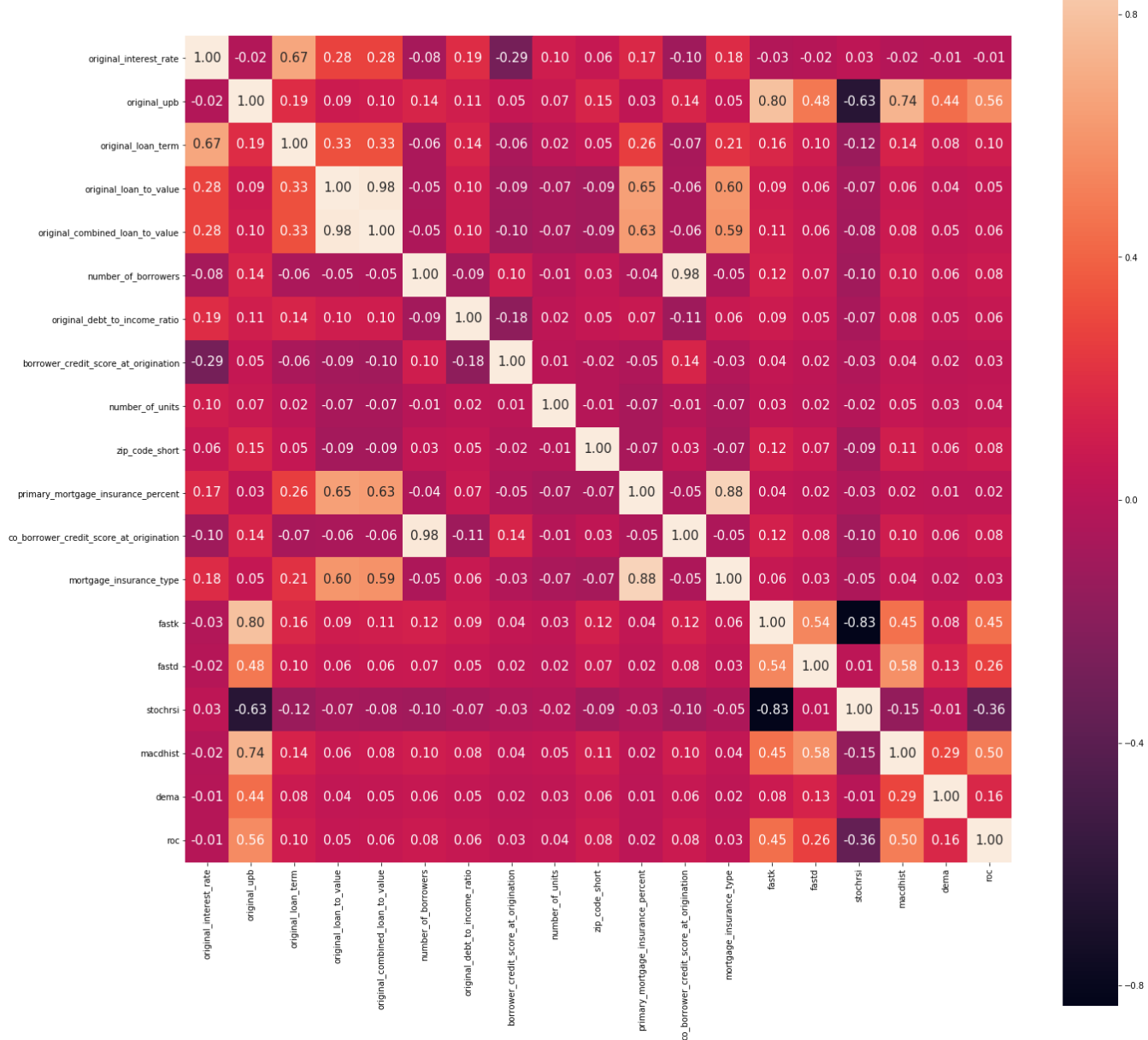
**Individual Final Report**

Using Pytorch and a Multi-Layer Perceptron Network, our goal is to predict whether a mortgage loan for a single-family home will fall delinquent by one or more months. The work was distributed as follows; Luke Bogacz would address getting the data to the model, Vishal Sinha would create and test models, and Jason Long would create visualizations and address project deliverables/management. All would experiment with models as time permits and share high performance models.
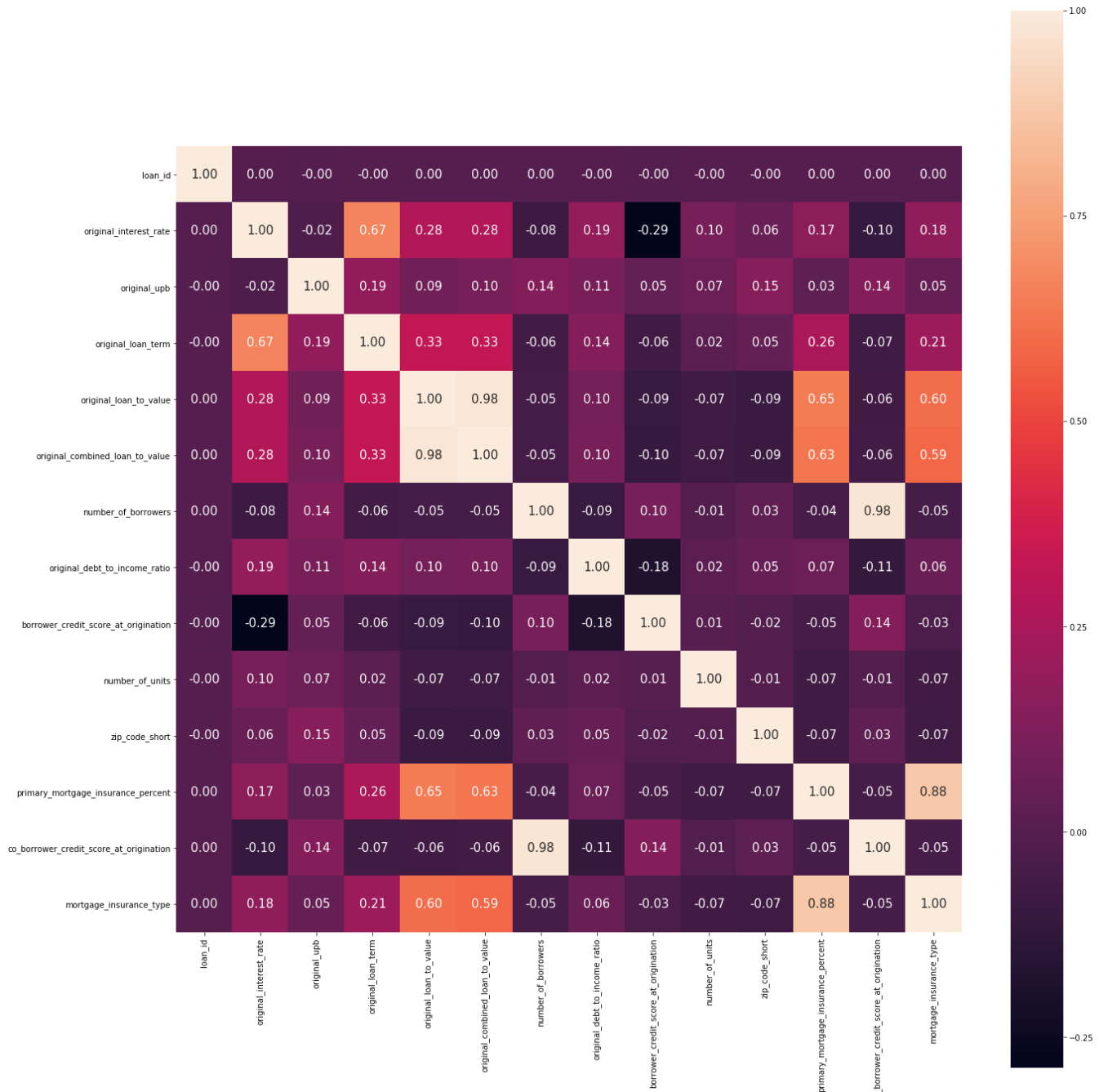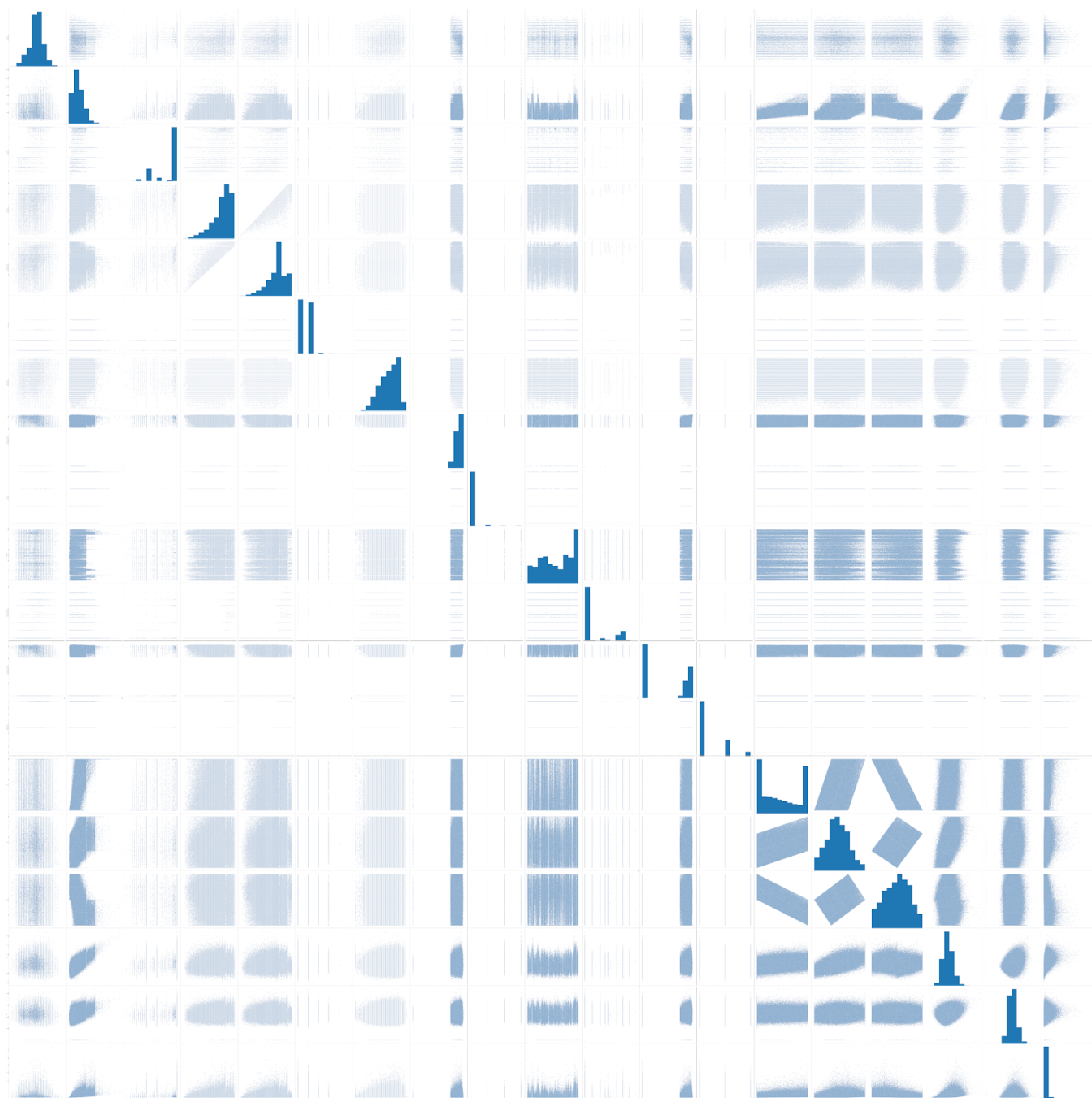
**Individual Work**

While awaiting the model being created so I could create visualizations, I spent time doing analytics on the dataset. Looking for collinearity and multicollinearity within the data to assist in dealing with the size of the data. Early models were having issues with data size as there were over 1,000 features and at times were causing memory issues during testing. In response to this, I upgraded my graphic processing unit to one that had a 50% larger amount of memory. This overcame some of the issues when dealing with the data; however, it exposed that there were some errors within the model. As those were overcome by the team, I continued data analytics.

These analytics were performed with Python within a Anaconda Jupyter Notebook, the notebook and associated code are located within the Code folder in my individual project submission.

Attempts at creating visualizations that operated from a remote server were initially unsuccessful. I initially attempted with HiddenLayer, a package that would create a graph depiction of the model and could assist with the creation of dynamic accuracy and loss. This package is focused to work with Jupyter notebooks. Setting up and operating Jupyter to work on an AWS deep learning instance is initially easy as it is pre-installed along with most

frameworks (pytorch, cafe, etc.) but does present its own issues. By having a complex AMI structure, getting packages to install correctly into the correct anaconda environment took experimentation. Eventually the packages were installed correctly, test code was functionable, but I was unable to get it to properly work with our model. The core issues with deploying the visualizations persisted across several packages. Straight displaying also did not properly function.

## Summary and Conclusion

The dataset is massive but contains many errors that will persist into the future as Fannie Mae relies on external data providers. Much more rigorous data management and transformations would be required to get the most out of the data. With the vast amount of potential training data, it would be possible to build a system that would best determine missing values. One method could be to identify persistent errors from vendors if they exist, and then examining correct records, to determine if there is a variance between them and a generic correct record. Or to examine if using the originator of the loan plays a larger role in the model by using that as a data import lens. That would reduce the size of the data within the model. A real world case would rely on a complex model that would inference on the corresponding originator model to determine probability of default.

My final conclusion is that going forward, it will be important for me to continue to pursue a more robust computer science skills focusing on programming. Luke Bogacz and Vishal Sinah are more advanced in their programming capabilities and it proved challenging keeping up with them in this group assignment. Both provided great code and my personal skills precluded matching their technical contributions to the success of this project.

**Code Calculation**

Percentage of non original code: 91.02%

## Works Cited/References

Facebook, Pytorch, *https://pytorch.org*

Waleed, HiddenLayer, *https://github.com/waleedka/hiddenlayer*

Sergey Zagoruyko, TNT, *https://github.com/pytorch/tnt*