**6.867 Machine Learning (Fall 2018)**
Prof.: Suvrit Sra
TA: Alex Turner

**Youssef Berrada** (ID: 911844154)
**Georgy Guryev** (ID: 924667514)
**Sheng Yao** (ID: 916855968)

# Project Milestone 4: Initial Results

In this milestone, we will get an initial version of your system running end-to-end and produce initial results. Basically, after cleaning the data and pre-processing it as describe below, we tried different algorithm such as a neural network, and some decision tree models (CART, Random forest, boosting) with arbitrary hyper-parameters in order to have the full workflow.

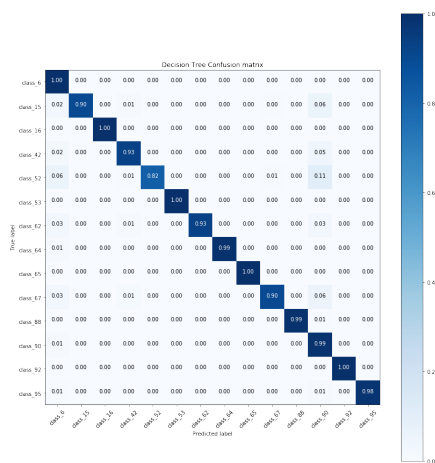## Data Description and Feature Engineering

The datasets that we used to predict astronomical bodies are from two types. The first source includes meta-data about the objects such as its location in the sky, red shift, etc. The other source of data includes 6 different light pass-band measurements observed across 2 years (each object is observed during different periods within two years). Our main challenge with feature engineering is to incorporate the time-series data from different frequency band measurements to classification. As a first-pass, we simply take the average of each passband's time series, for each object, thereby reducing the time dimension of the data. Our next steps for future milestones are to apply different methods to preserve the time-series information and incorporate them into classification.
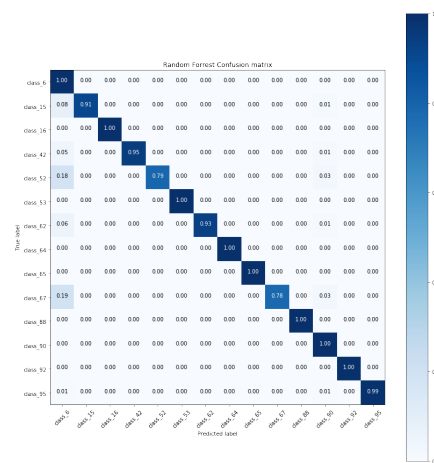
## Initial Results

### Tree-based Models

#### CART

The CART Model overfits the training set perfectly (100% precision) with max-depth of 20. It is undesirable but we include it as a first try of the method.
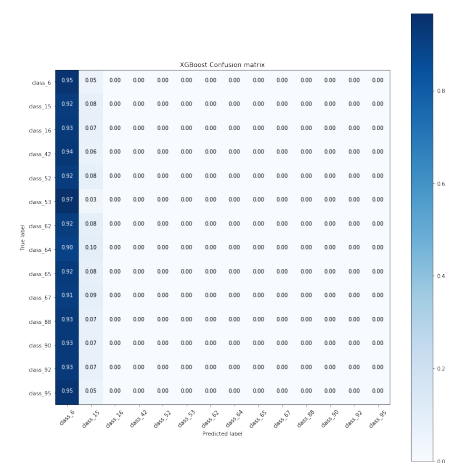
#### Random Forest

The Random Forest model is less prone to overfit than CART Model theoretically, but with same max-depth of 20. We can see that it also tends towards overfitting.

#### Boosting Tree

Based on the previous two results, we would expect the Boosting method to exhibit similar behavior. However, the method underfits the training set and produces only 1%precision



### Neural Network

We implemented a feed-forward Neural Network with 4 hidden layers and a softmax output layer. To prevent overfitting, we add dropout and batch normalization at each layer. The activation function at each hidden layer is ReLu and we train the model with 5 epochs and batch size of 100.

| dense_21_input: InputLayer | input: | (None, 31) |
|---|---|---|
| | output: | (None, 31) |

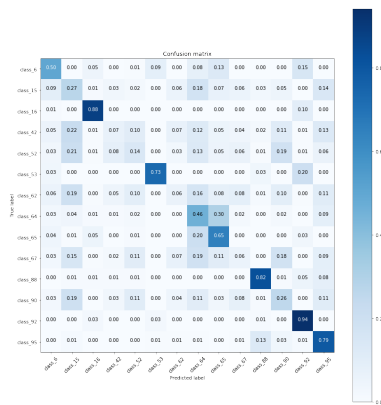| dense_21: Dense | input: | (None, 31) |
|---|---|---|
| | output: | (None, 512) |

| h_normalization_17: BatchNormalization | input: | (None, 512) |
|---|---|---|
| | output: | (None, 512) |

| dropout_17: Dropout | input: | (None, 512) |
|---|---|---|
| | output: | (None, 512) |

| dense_22: Dense | input: | (None, 512) |
|---|---|---|
| | output: | (None, 256) |

| h_normalization_18: BatchNormalization | input: | (None, 256) |
|---|---|---|
| | output: | (None, 256) |

| dropout_18: Dropout | input: | (None, 256) |
|---|---|---|
| | output: | (None, 256) |

| dense_23: Dense | input: | (None, 256) |
|---|---|---|
| | output: | (None, 128) |

| h_normalization_19: BatchNormalization | input: | (None, 128) |
|---|---|---|
| | output: | (None, 128) |

| dropout_19: Dropout | input: | (None, 128) |
|---|---|---|
| | output: | (None, 128) |

| dense_24: Dense | input: | (None, 128) |
|---|---|---|
| | output: | (None, 64) |

| h_normalization_20: BatchNormalization | input: | (None, 64) |
|---|---|---|
| | output: | (None, 64) |

| dropout_20: Dropout | input: | (None, 64) |
|---|---|---|
| | output: | (None, 64) |

| dense_25: Dense | input: | (None, 64) |
|---|---|---|
| | output: | (None, 14) |

Figure 1: The results of astronomical object classification with the feed-forward Neural Net

## Conclusion and Next Steps

Running the above neural network led us to be ranked 230 in Kaggle, which is a strong improvement from the previous iteration with a loss going from 6.6 to 1.70 where the leader is around 0.8. This first approach forces us to have a complete work flow from the data pre-processing to the data visualization. This helped us get to know the data and produce some initial (bad) results... For the next step, we would like to:

- Refine the approach with aggregated data to get finer results across different methodology.
- Once we are confident we our models with this approach, we would like to switch to a time series analysis and use the RNN framework to see if the results will be better. The time-series approach offers a wide possibility of features engineering, as we can identify pattern using technical indicators ( Bollinger bands...)

## Division of Labor

All three team members contributed significantly to the milestone. We actually all brainstormed at the same time to understand each algorithm and discussed what methods should we explore and package we should use. We then all went over all algorithms and libraries we plan to use.