

Predicting Pulsar Stars

Feature Importance and Reduction

Duc Ngo - CS4300

April 22, 2020

Contents

1	Introduction	3
2	Dataset	3
2.1	Visualization of the distribution of each input features	4
2.2	Distribution of the output labels	6
3	Data Processing	6
3.1	Data Splitting	6
3.2	Data Normalization	6
3.3	Normalized Data	7
4	Modelling	9
4.1	Selected Neural Network Architecture	9
4.1.1	Baseline Model Performance	9
4.2	Learning Curve of Baseline Model	9
4.3	Modifying Activation Function	10
4.4	Learning Curve of Model using Linear Activation (last neuron)	10
4.5	Learning Curve of Model using Linear Activation (all neuron)	10
4.6	Learning Curve of Model using Sigmoid Activation (last neuron)	11
4.7	Learning Curve of Model using Sigmoid Activation (all neuron)	11
4.8	Overfitting	12
5	Model Evaluation	13
5.1	Test Accuracy	13
5.2	Custom Function and Keras Function	15
6	Feature Importance Analysis	16
6.1	Significance of individual features	16
6.2	Performance after removing less important features	17
6.3	Performance of two features	18
7	TBA...	19

1 Introduction

What is a Pulsar? Whenever you look at the beautiful night sky, you may find yourself looking at a stationary shiny object and there is a slim chance that you are looking at a pulsar. Pulsars frequently look like glinting stars, and they appear to sparkle with a normal cadence. Nonetheless, the light from pulsars doesn't sparkle or beat, and these celestial bodies are not stars. Pulsars are round, minimal objects that are about the size of a huge city however contain more mass than our sun or other average stars. [1]

Why are they so fascinating? Researchers are utilizing pulsars to analyze the extreme states of matter, look for planets past Earth's nearby planetary group and measure cosmic distances. Pulsars likewise could assist researchers with finding gravitational waves, which could guide the way to energetic cosmic events like collisions between supermassive black holes. Pulsar is mysterious and still new to the field of science which is why we are studying them. These celestial objects act so differently on an extreme scale that we never experience here on Earth. They are by definition is an extreme science. [1]

Motivation Artificial intelligence could be the perfect tool for exploring the Universe since finding a specific celestial object is a slow and tedious job. With this in mind, the motivation of this project is to try to write supervised learning algorithms to predict the class of pulsars with binary classification. Binary classification is a machine learning technique that categorizes whether the element is either true or false OR 1 or 0.

Google Colab <https://colab.research.google.com/drive/1Y4umEtZC9rIDyOMKkOtXH6FponMZICV6>

2 Dataset

The "Predicting a Pulsar Star" dataset was obtained from the Kaggle Data Science [4][5]. The pulsar candidate data were obtained during the High Time Resolution Universe (HTRU) survey. It contains 16,259 apocryphal (negative) examples caused by RFI/noise and 1,639 real pulsars (positive) examples which totals to 17,898 samples in the dataset. These samples have all been checked by human annotators. Each row lists the variables first that described by 8 continuous variables, and a single class label as the final entry. The class labels used are 0 (negative) and 1 (positive). The input features are for the following fields:

1. Mean of the integrated profile.
2. The standard deviation of the integrated profile.
3. Excess kurtosis of the integrated profile.
4. The skewness of the integrated profile.
5. Mean of the DM-SNR curve.
6. The standard deviation of the DM-SNR curve.
7. Excess kurtosis of the DM-SNR curve.
8. The skewness of the DM-SNR curve.
9. Class

2.1 Visualization of the distribution of each input features

The histogram plot of each info highlights indicating their most extreme and least value as well as how they are distributed can be found in the pictures given underneath.

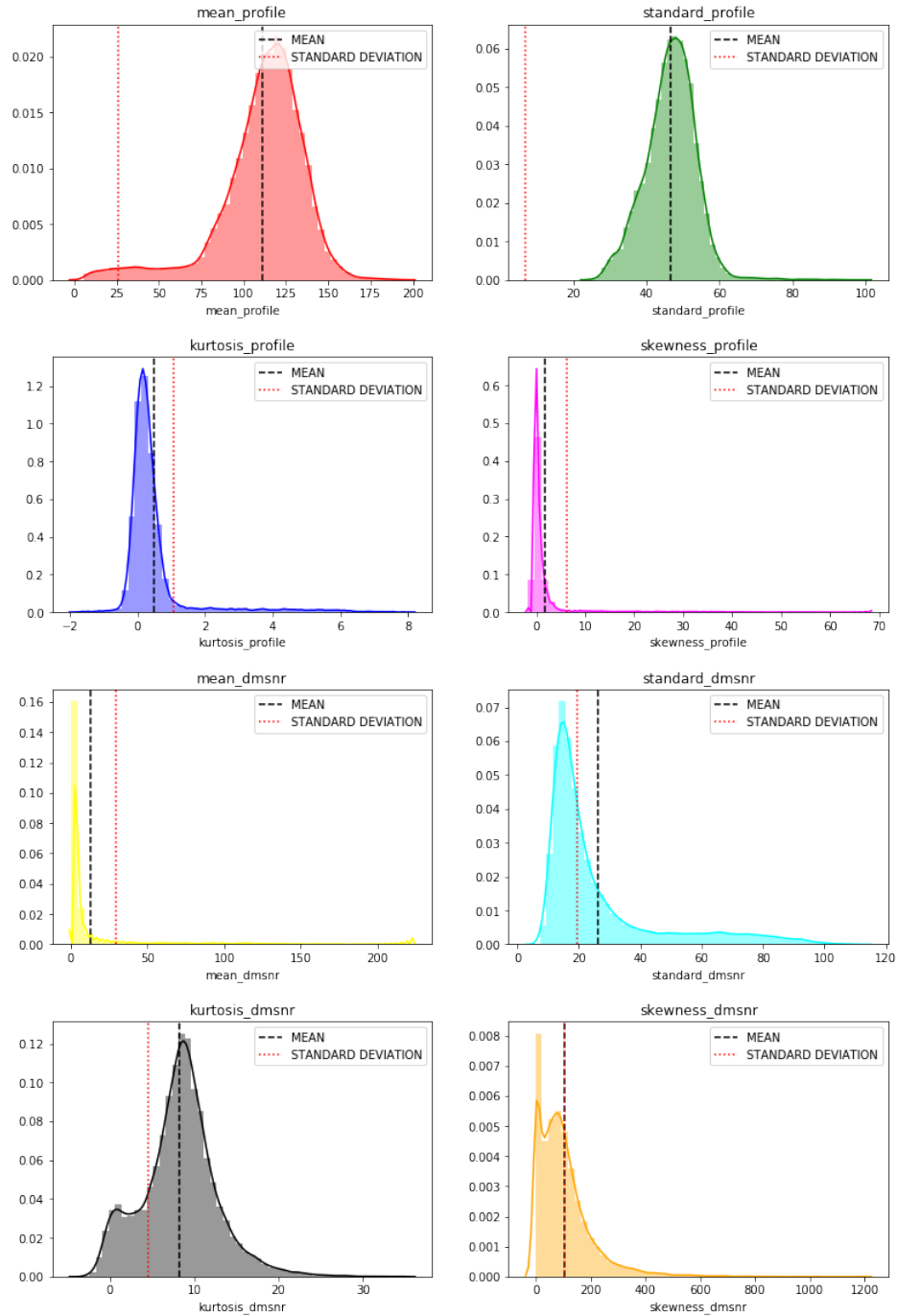


Figure 1: Input Data Distribution Histograms

	mean_profile	standard_profile	kurtosis_profile	skewness_profile	mean_dmsnr	standard_dmsnr	kurtosis_dmsnr	skewness_dmsnr	target
count	17898.000000	17898.000000	17898.000000	17898.000000	17898.000000	17898.000000	17898.000000	17898.000000	17898.000000
mean	111.079968	46.549532	0.477857	1.770279	12.614400	26.326515	8.303556	104.857709	0.091574
std	25.652935	6.843189	1.064040	6.167913	29.472897	19.470572	4.506092	106.514540	0.288432
min	5.812500	24.772042	-1.876011	-1.791886	0.213211	7.370432	-3.139270	-1.976976	0.000000
25%	100.929688	42.376018	0.027098	-0.188572	1.923077	14.437332	5.781506	34.960504	0.000000
50%	115.078125	46.947479	0.223240	0.198710	2.801839	18.461316	8.433515	83.064556	0.000000
75%	127.085938	51.023202	0.473325	0.927783	5.464256	28.428104	10.702959	139.309331	0.000000
max	192.617188	98.778911	8.069522	68.101622	223.392140	110.642211	34.539844	1191.000837	1.000000

Figure 2: Input Feature Statistics

2.2 Distribution of the output labels

Notice that the data is imbalanced and may need to be resampled (such as oversampling, undersampling, or generate synthetic samples), but for now it is acceptable.

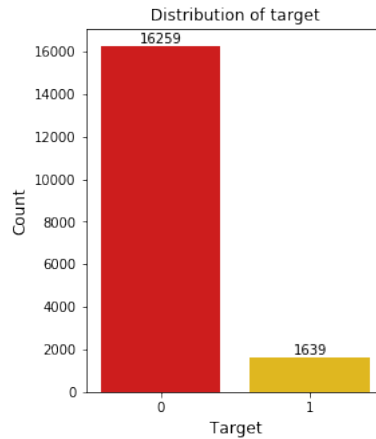


Figure 3: Output Data Distribution Histogram

3 Data Processing

3.1 Data Splitting

The data was randomly shuffled and then the dataset was split into training and validation, where 80% of the dataset was allocated for training and 20% was allocated for validation.

3.2 Data Normalization

Before data mining itself, data preprocessing plays a crucial role. As can be seen, the data was not distributed uniformly. In this manner, we have to pre-process the data with normalization techniques. Normalization makes the optimization problem more numerically stable and makes training less sensitive to the scale of features, so we can better solve for coefficients. When applying normalization, all the values are all now between 0 and 1, and the outliers are eliminate but remain visible within our normalized data. There are two normalization techniques that can be utilized and each of them has its own consequences, but for now, any of them is sufficient.

Mean Normalization Formula

$$X_{new} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Min-Max Normalization Formula

$$X_{new} = \frac{X - X_{mean}}{X_{max} - X_{min}}$$

3.3 Normalized Data

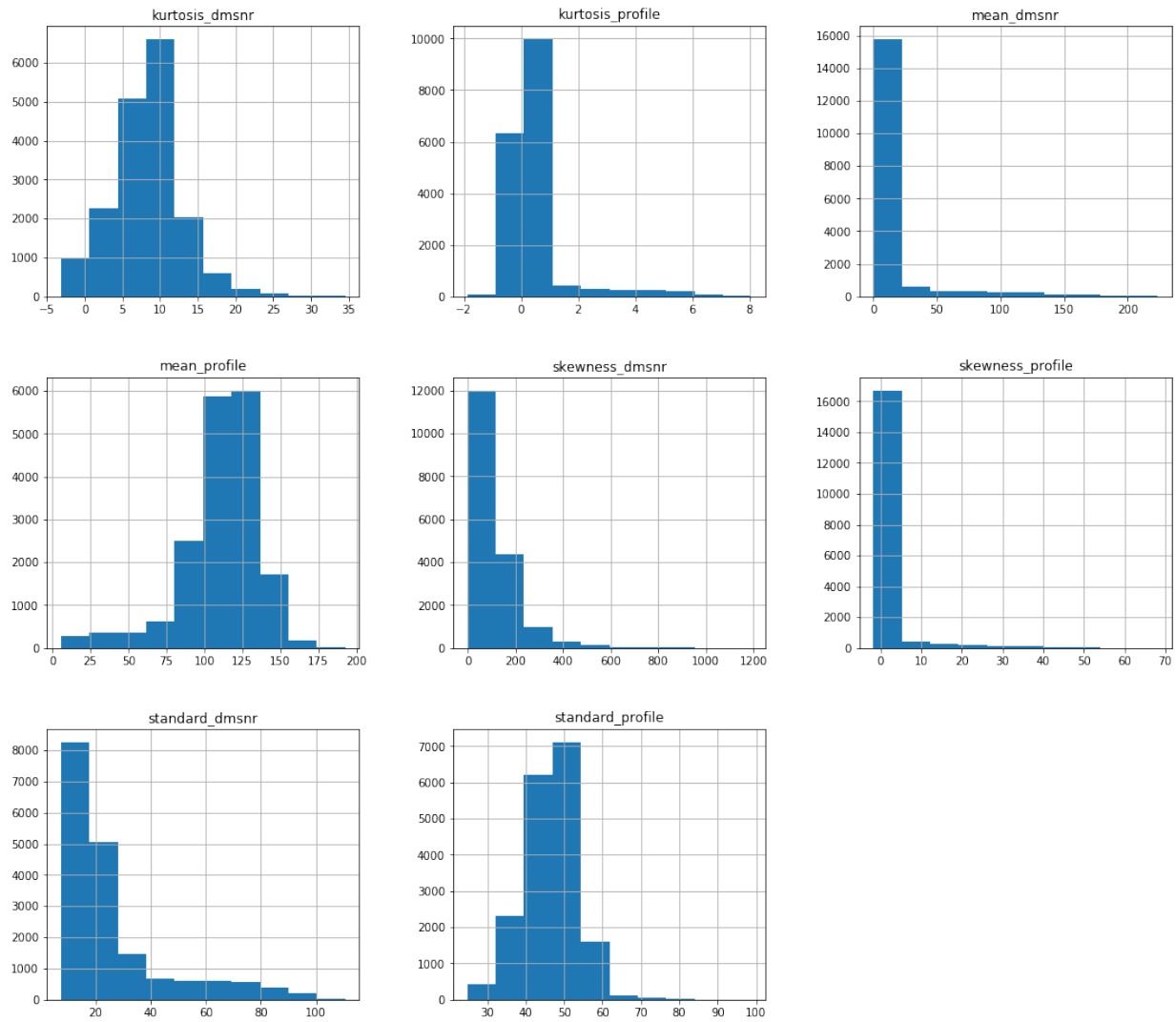


Figure 4: Input Data Before Normalization (Min-Max)

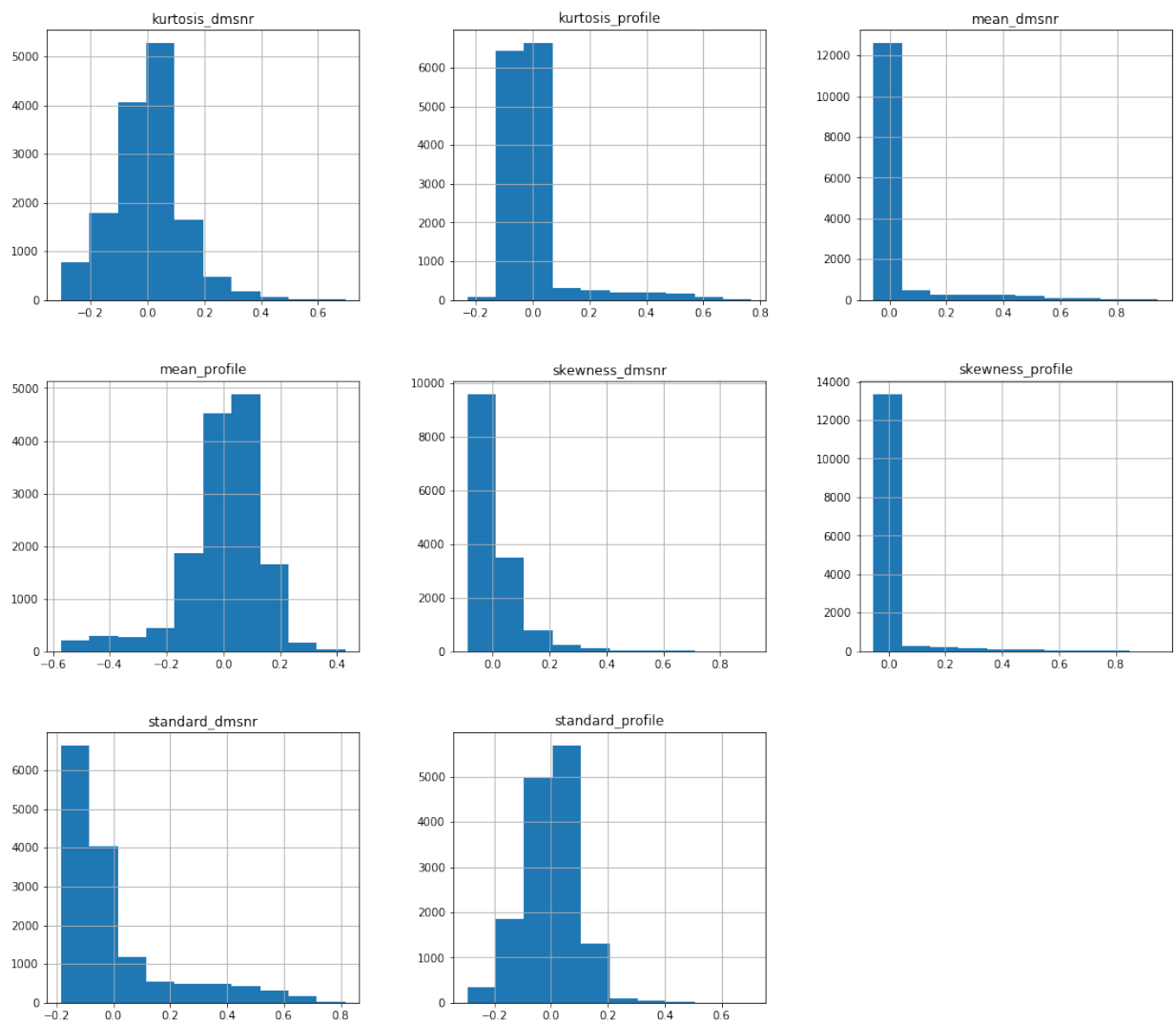


Figure 5: Input Data After Normalization (Min-Max)

4 Modelling

A feed forward artificial neural network architectures was used to create the model.

NOTE: Data is shuffle. Thus, the result will vary every time. All models were compiled and fit on March 17, 2020.

4.1 Selected Neural Network Architecture

The first model is a baseline model (act as a control) that has a basic architecture of one input and one output layer. It will then gradually increase by one hidden layer and up to 2 hidden layers at the end.

4.1.1 Baseline Model Performance

Hidden	Accuracy	Loss
0 Layer	98.07	6.68
1 Layer	98.16	6.46
2 Layer	98.16	6.32

Table 1: performance comparison for different hidden layers

As can be seen from the above table, the basic architecture with just one hidden layer overall is performing better than other architecture with zero or multiple hidden layers for all datasets. Thus, only one hidden layer will be applied to other architecture.

4.2 Learning Curve of Baseline Model

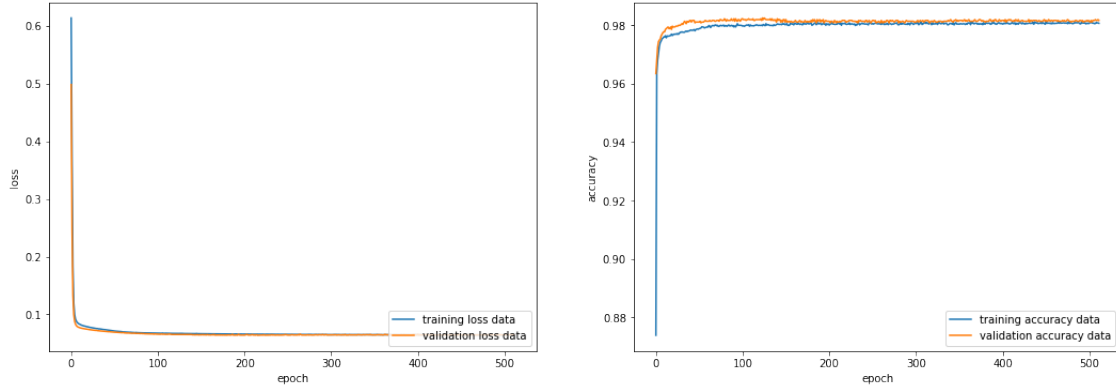


Figure 6: curve showing change in loss/accuracy vs epoch

4.3 Modifying Activation Function

The linear activation function performed poorer in comparison with a sigmoid activation function. In general, linear activation is configured with mean squared error (MSE) or mean absolute error (MAE) loss function and these functions are usually a bad choice (depend on a dataset) for binary classification problems. When using MSE, it is assumed that the underlying data has been generated from a normal distribution or a bell-shaped curve. However, in reality, the dataset that classified into two categories is not from a normal distribution. Lastly, the MSE function is non-convex for binary classification which means it is not guaranteed to minimize the loss function.[3] The performance comparison can be seen in the table shown below.

Activation Function	Accuracy	MAE	Loss
Linear (Output)	-	3.39	1.48
Linear (All)	-	9.28	2.52
Sigmoid (Output)	98.10	-	6.74
Sigmoid (All)	98.13	-	6.73

Table 2: performance comparison for different activation function

4.4 Learning Curve of Model using Linear Activation (last neuron)

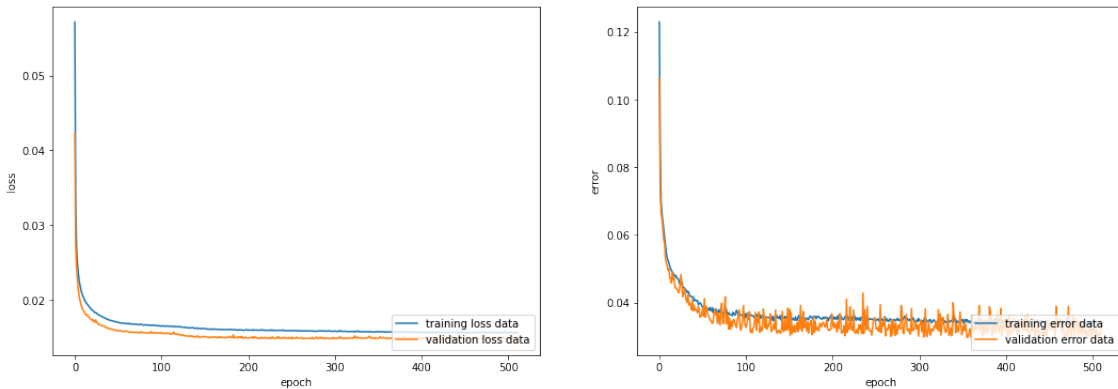


Figure 7: curve showing change in loss/error vs epoch

4.5 Learning Curve of Model using Linear Activation (all neuron)

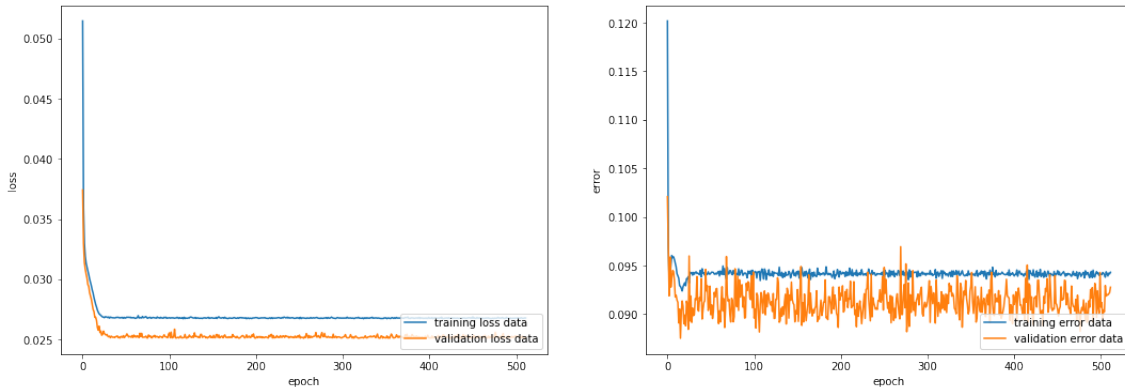


Figure 8: curve showing change in loss/error vs epoch

4.6 Learning Curve of Model using Sigmoid Activation (last neuron)

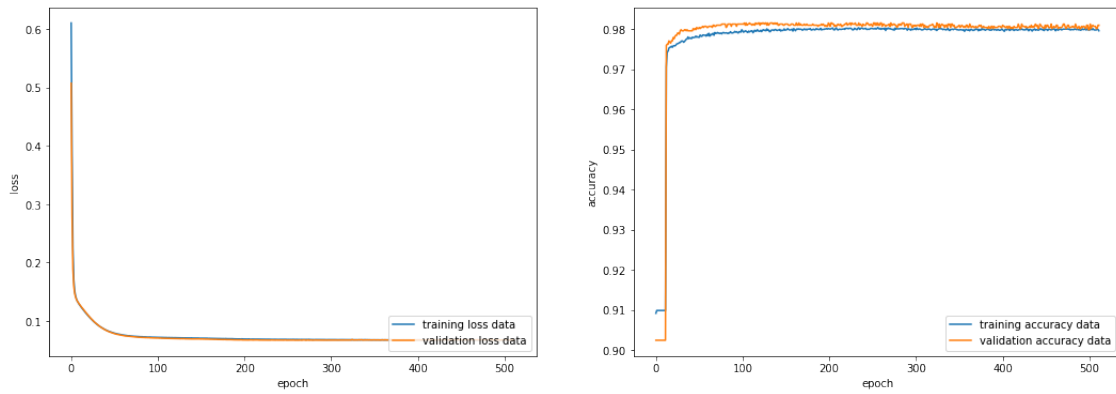


Figure 9: curve showing change in loss/accuracy vs epoch

4.7 Learning Curve of Model using Sigmoid Activation (all neuron)

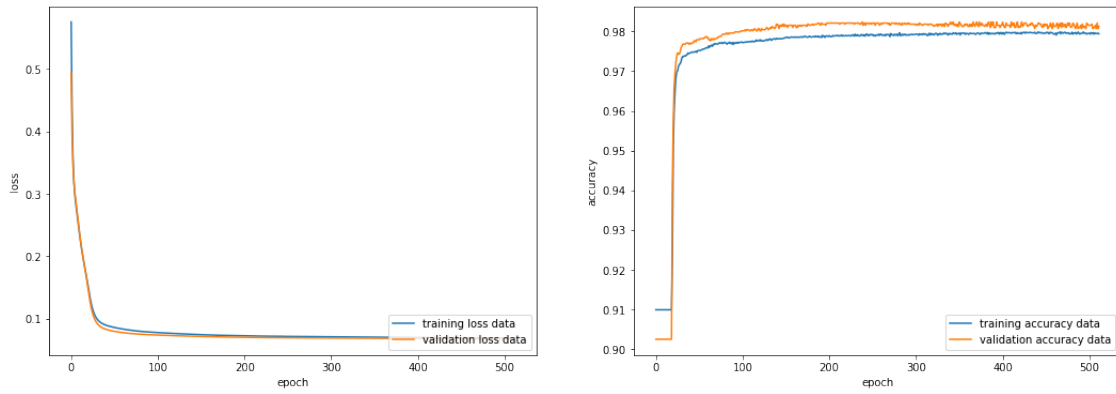


Figure 10: curve showing change in loss/accuracy vs epoch

4.8 Overfitting

The number of neurons in each layer was increased by a factor of 10 to overfit the model.

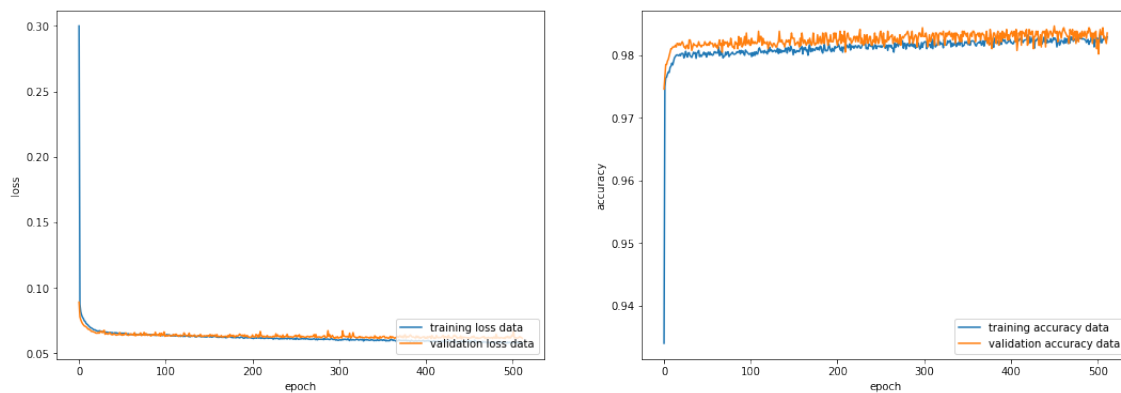


Figure 11: curve showing change in loss/accuracy vs epoch

	Accuracy	Loss
Overfitting	98.35	6.14

Table 3: performance of overfitting model

5 Model Evaluation

Three essential classification model metrics to evaluate. The table given below shows the precision, recall and f1 score for the neural network model.

1. Precision: what proportion of positive identifications was actually correct?
2. Recall: what proportion of actual positives was identified correctly?
3. F1-Score: evaluation metric for classification algorithms, where the best value is at 1 and the worst is at 0.

Model	Precision	Recall	F1-Score
Baseline	94.08	86.53	0.90
Linear (last)	94.14	87.39	0.91
Linear (all)	96.31	74.79	0.84
Sigmoid (last)	93.50	86.53	0.90
Sigmoid (all)	93.79	86.53	0.90
Overfitting	94.48	88.25	0.91

Table 4: predictions evaluation of all model

A useful tool when predicting the probability of a binary outcome is the Receiver Operating Characteristic curve or ROC curve. The area covered by the curve is the area between the red line and the axis. This area covered is AUC. The bigger the area covered, the better the machine learning models are at distinguishing the given classes. In other words, the AUC can be used as a summary of the model skill. The ideal value for AUC is 1.[2]

5.1 Test Accuracy

The closer the graph is to the top and left-hand borders, the more accurate the test. Likewise, the closer the graph to the diagonal, the less accurate the test. In a perfect test, it would go straight from zero up the top-left corner and then straight across the horizontal. The figure is given below shows the receiver operating characteristic curve for the baseline model, the linear model, and the sigmoid model.

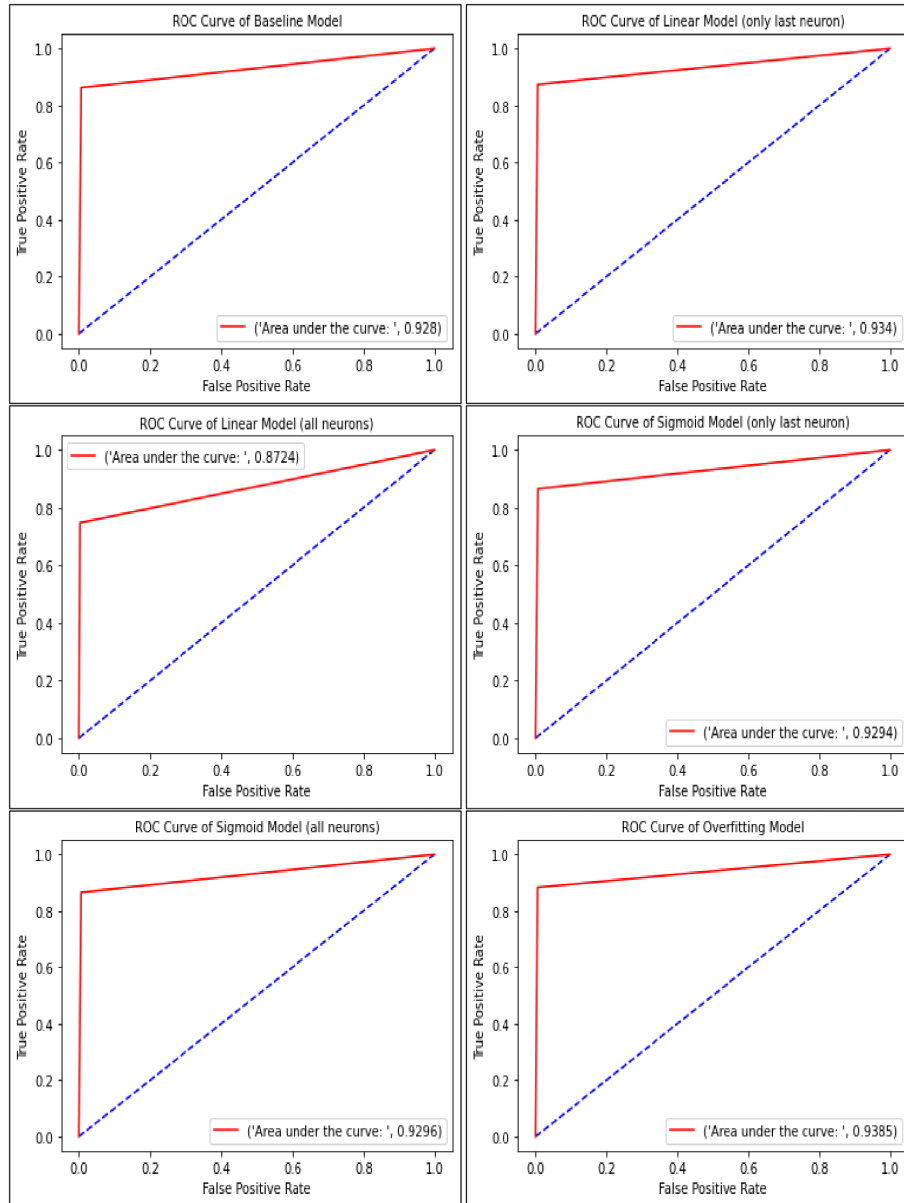


Figure 12: Receiver Operating Characteristic Curve for all model

5.2 Custom Function and Keras Function

After the model was trained, all the weights were extracted. A custom function/method was build to serves as the model. The extracted weights were then applied to the custom function to predict the outputs. The goals are to verify if the custom predictions that yield are the same as the trained model. The comparison will be on the sigmoid model (last neuron).

```

Layer #0
Weights:
[[ 0.0  0.3 -2.0  1.5 -0.3  0.8 -0.8 -0.4]
 [-0.6  0.8 -2.6  0.7  0.7  0.4 -1.1  0.5]
 [ 0.7  0.1  3.4 -1.1  0.2  0.3  0.9  0.1]
 [ 0.2  0.1  0.0 -0.1 -0.2  0.1  0.3 -0.4]
 [ 0.3  1.1 -1.2  0.5  0.4 -2.2 -0.1  1.3]
 [-1.5  0.0  0.3  1.3  0.6 -2.1  0.2  0.3]
 [ 0.7  0.3  2.8 -0.4  0.5 -0.1  0.2  0.1]
 [ 0.6 -1.1 -3.8  0.4  0.6 -1.1 -1.3  1.1]]
Bias:
[ 0.5  0.7  0.3 -0.1  0.4  0.1  0.2  0.4]

Layer #1
Weights:
[[-0.5 -0.1 -0.7 -0.0  0.0 -0.4  0.3 -0.0]
 [ 0.9  1.5 -1.5 -0.5  1.1  1.1 -1.0  2.6]
 [ 1.0  0.2 -1.5 -0.4  1.6  1.5 -1.3  1.5]
 [-8.5 -0.4  0.4 -0.4 -1.7  0.4  0.2 -0.6]]
Bias:
[-0.1  0.4  0.4  0.4]

Layer #2
Weights:
[[ 0.4 -1.9 -1.5 -2.6]]
Bias:
[ 4.6]

Accuracy: 98.10%
Precision: 93.50%
Recall: 86.53%
F1-score: 0.90

X=[ 0.2  0.1 -0.1 -0.0 -0.0 -0.1 -0.0 -0.0], Predicted=[False]
X=[ 0.2  0.1 -0.1 -0.0  0.1  0.2 -0.2 -0.1], Predicted=[False]
X=[-0.3 -0.2  0.2  0.1  0.1  0.4 -0.2 -0.1], Predicted=[ True]
X=[ 0.0 -0.0 -0.0 -0.0 -0.0 -0.1  0.1  0.0], Predicted=[False]
X=[ 0.0  0.1 -0.0 -0.0 -0.0 -0.1  0.0  0.0], Predicted=[False]
X=[ 0.1 -0.0 -0.0 -0.0 -0.0 -0.1  0.0 -0.0], Predicted=[False]
X=[ 0.1  0.1 -0.1 -0.0 -0.0 -0.1 -0.0 -0.0], Predicted=[False]
X=[ 0.1  0.2 -0.0 -0.0 -0.0 -0.1  0.1  0.0], Predicted=[False]
X=[ 0.1 -0.0 -0.1 -0.0 -0.0 -0.1  0.1  0.0], Predicted=[False]
X=[-0.1 -0.0 -0.0 -0.0 -0.0  0.1 -0.1 -0.1], Predicted=[False]

```

Figure 13: custom function result

```

3579/3579 [=====] - 0s 76us/sample - loss: 0.0674 - acc: 0.9810
loss: 6.74%

acc: 98.10%

[[3209  21]
 [  47 302]]
98.10002794076557%

Accuracy: 98.10%
Precision: 93.50%
Recall: 86.53%
F1-score: 0.90

X=[ 0.2  0.1 -0.1 -0.0 -0.0 -0.1 -0.0 -0.0], Predicted=[False]
X=[ 0.2  0.1 -0.1 -0.0  0.1  0.2 -0.2 -0.1], Predicted=[False]
X=[-0.3 -0.2  0.2  0.1  0.1  0.4 -0.2 -0.1], Predicted=[ True]
X=[ 0.0 -0.0 -0.0 -0.0 -0.0 -0.1  0.1  0.0], Predicted=[False]
X=[ 0.0  0.1 -0.0 -0.0 -0.0 -0.1  0.0  0.0], Predicted=[False]
X=[ 0.1 -0.0 -0.0 -0.0 -0.0 -0.1  0.0 -0.0], Predicted=[False]
X=[ 0.1  0.1 -0.1 -0.0 -0.0 -0.1 -0.0 -0.0], Predicted=[False]
X=[ 0.1  0.2 -0.0 -0.0 -0.0 -0.1  0.1  0.0], Predicted=[False]
X=[ 0.1 -0.0 -0.1 -0.0 -0.0 -0.1  0.1  0.0], Predicted=[False]
X=[-0.1 -0.0 -0.0 -0.0 -0.0  0.1 -0.1 -0.1], Predicted=[False]

```

Figure 14: keras function result

6 Feature Importance Analysis

As of now, we have a pretty good idea of which model to use, the number of epochs, and a reasonable validation set to use for the network architecture. The next step is to find out which input features is redundant or insignificant.

6.1 Significance of individual features

Five input features (standard profile, mean dmsnr, standard dmsnr, kurtosis dmsnr, and skewness dmsnr) slightly impact the overall accuracy of the model. We can see the graph below for the performance.

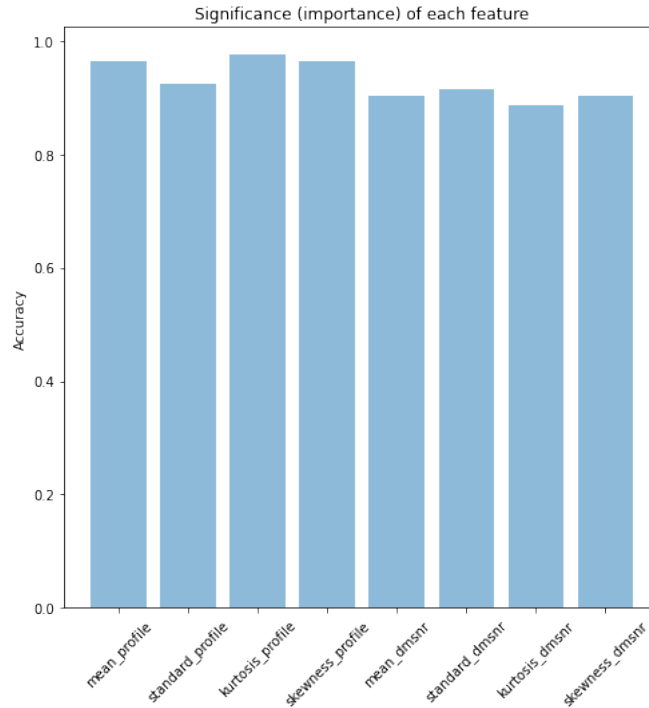


Figure 15: significance of individual features

```
2020/04/21 08:04:42 PM | mean_profile | 0.9638666418327435
2020/04/21 08:05:31 PM | standard_profile | 0.924380704041721
2020/04/21 08:05:56 PM | kurtosis_profile | 0.9759731793630099
2020/04/21 08:06:10 PM | skewness_profile | 0.9655429316446266
2020/04/21 08:07:13 PM | mean_dmsnr | 0.9027751909107842
2020/04/21 08:08:39 PM | standard_dmsnr | 0.914695474017508
2020/04/21 08:08:52 PM | kurtosis_dmsnr | 0.8869435649096666
2020/04/21 08:10:45 PM | skewness_dmsnr | 0.9027751909107842
```

Figure 16: significance of individual features data

6.2 Performance after removing less important features

There was no significant drop (at all) in the performance after removing less important features one at a time. We can see the graph below for the performance.

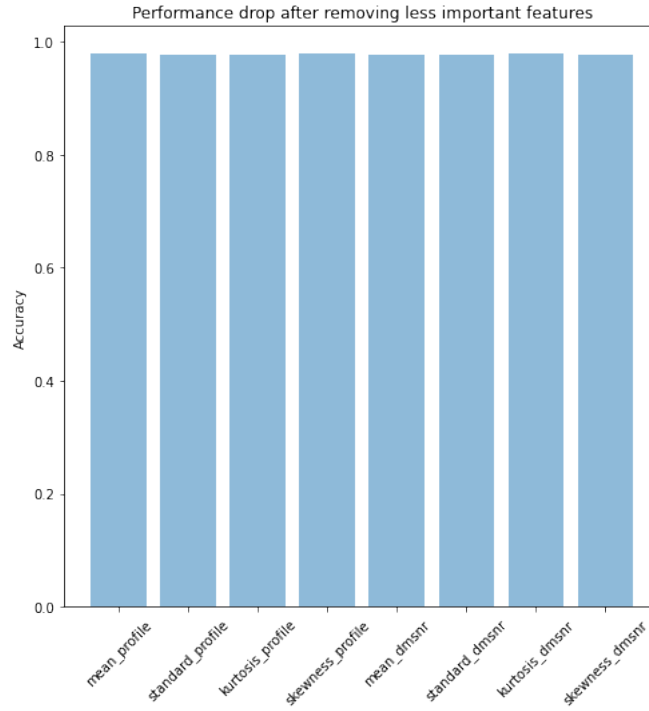


Figure 17: performance after removing less important features

```
2020/04/21 08:13:22 PM | mean_profile | 0.9778357235984355
2020/04/21 08:14:46 PM | standard_profile | 0.9774632147513503
2020/04/21 08:16:24 PM | kurtosis_profile | 0.9770907059042653
2020/04/21 08:17:22 PM | skewness_profile | 0.9782082324455206
2020/04/21 08:19:15 PM | mean_dmsnr | 0.9769044514807227
2020/04/21 08:21:26 PM | standard_dmsnr | 0.9776494691748929
2020/04/21 08:23:04 PM | kurtosis_dmsnr | 0.9787669957161482
2020/04/21 08:24:40 PM | skewness_dmsnr | 0.9774632147513503
```

Figure 18: performance after removing less important features data

6.3 Performance of two features

Even after removing two input features, there was no significant drop (at all) in the performance after removing less important features one at a time. We can see the graph below for the performance.

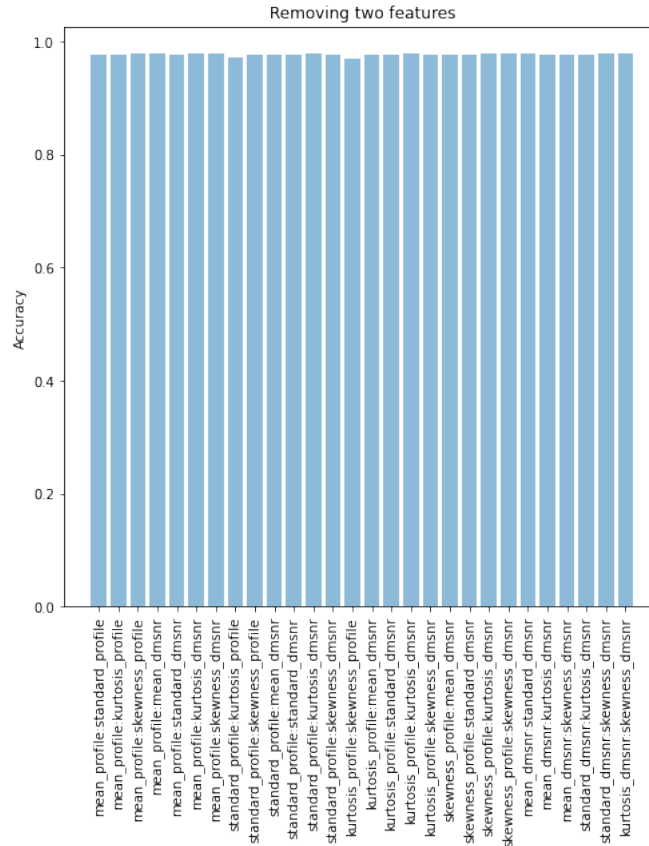


Figure 19: performance of two features

```

2020/04/21 08:26:22 PM | mean_profile:standard_profile | 0.9770907059042653
2020/04/21 08:28:14 PM | mean_profile:kurtosis_profile | 0.976345688210095
2020/04/21 08:29:51 PM | mean_profile:skewness_profile | 0.9778357235984355
2020/04/21 08:31:24 PM | mean_profile:mean_dmsnr | 0.9783944868690632
2020/04/21 08:33:30 PM | mean_profile:standard_dmsnr | 0.9774632147513503
2020/04/21 08:35:10 PM | mean_profile:kurtosis_dmsnr | 0.9782082324455206
2020/04/21 08:37:12 PM | mean_profile:skewness_dmsnr | 0.9778357235984355
2020/04/21 08:38:08 PM | standard_profile:kurtosis_profile | 0.9718755820450735
2020/04/21 08:40:10 PM | standard_profile:skewness_profile | 0.9774632147513503
2020/04/21 08:41:29 PM | standard_profile:mean_dmsnr | 0.9776494691748929
2020/04/21 08:42:43 PM | standard_profile:standard_dmsnr | 0.9770907059042653
2020/04/21 08:44:19 PM | standard_profile:kurtosis_dmsnr | 0.9782082324455206
2020/04/21 08:46:14 PM | standard_profile:skewness_dmsnr | 0.9774632147513503
2020/04/21 08:47:48 PM | kurtosis_profile:skewness_profile | 0.9696405289625628
2020/04/21 08:49:29 PM | kurtosis_profile:mean_dmsnr | 0.9772769603278078
2020/04/21 08:53:32 PM | kurtosis_profile:standard_dmsnr | 0.9772769603278078
2020/04/21 08:55:08 PM | kurtosis_profile:kurtosis_dmsnr | 0.978021978021978
2020/04/21 08:56:38 PM | kurtosis_profile:skewness_dmsnr | 0.9774632147513503
2020/04/21 08:58:04 PM | skewness_profile:mean_dmsnr | 0.9776494691748929
2020/04/21 08:59:33 PM | skewness_profile:standard_dmsnr | 0.9774632147513503
2020/04/21 09:01:27 PM | skewness_profile:kurtosis_dmsnr | 0.978021978021978
2020/04/21 09:03:44 PM | skewness_profile:skewness_dmsnr | 0.9782082324455206
2020/04/21 09:06:01 PM | mean_dmsnr:standard_dmsnr | 0.978021978021978
2020/04/21 09:07:22 PM | mean_dmsnr:kurtosis_dmsnr | 0.9767181970571801
2020/04/21 09:08:52 PM | mean_dmsnr:skewness_dmsnr | 0.9772769603278078
2020/04/21 09:10:05 PM | standard_dmsnr:kurtosis_dmsnr | 0.9772769603278078
2020/04/21 09:12:38 PM | standard_dmsnr:skewness_dmsnr | 0.9782082324455206
2020/04/21 09:13:54 PM | kurtosis_dmsnr:skewness_dmsnr | 0.978021978021978

```

Figure 20: performance of two features data

7 TBA...

References

- [1] Calla Cofield. *What Are Pulsars?* Apr. 2016. URL: <https://www.space.com/32661-pulsars.html>. (accessed: 2.20.2020).
- [2] Guest Contributor. *Understanding ROC Curves with Python*. Feb. 2019. URL: <https://stackabuse.com/understanding-roc-curves-with-python/>. (accessed: 3.16.2020).
- [3] Hong Jing. *Why Linear Regression is not suitable for Classification*. May 2019. URL: <https://jinglescode.github.io/datascience/2019/05/07/why-linear-regression-is-not-suitable-for-classification/>. (accessed: 3.16.2020).
- [4] Dr Robert Lyon. *HTRU2 Data Set*. Feb. 2017. URL: <https://archive.ics.uci.edu/ml/datasets/HTRU2>. (accessed: 2.20.2020).
- [5] Pavan Raj. *Predicting a Pulsar Star*. May 2018. URL: <https://www.kaggle.com/pavanraj159/predicting-a-pulsar-star>. (accessed: 2.20.2020).