

# MINI-PROJECT 3: NETWORK INTRUSION DETECTION

Version 1.0

CSC215, Fall 2018

Oct 22<sup>nd</sup>, 2018

Prepared by: Chandini Nagendra

Siddharth Chittora

## Contents

1. Problem Statement.....	3
2. Methodology.....	3
3. Experimental Results and Analysis.....	3
4. Task Division.....	5
4.1. Chandini Nagendra: .....	5
4.2. Siddharth Chittora.....	5
5. Project Reflection.....	5
6. Additional Features.....	6
7. References .....	8

## 1. Problem Statement

This project aims to build a network intrusion detector, a predictive model capable of distinguishing between bad connections, called intrusions or attacks, and good normal connections.

## 2. Methodology

Here we are using dataset containing a wide variety of intrusions simulated in a military network environment. We have implemented Sklearn models for Logistic Regression, SVM, Gaussian Naïve Bayes, KNN, TensorFlow models and CNN.

- We are using the entire dataset.
- The dataset lacked column headers, so we manually added them.
- We normalized numeric features and encoded categorical features.
- We dropped rows with missing values.
- We encode the Outcome column with 0 for normal connection and 1 for all the other intrusion.
- We then split the data for training and testing.
- We implemented the Sklearn models for Logistic Regression, SVM, Gaussian Naïve Bayes, KNN.
- We created three Tensorflow models with activation function ReLU, Sigmoid and Tanh.
- We also used Early Stopping and checkpointing and 4 hidden layers in these models.
- We implemented CNN model to handle numeric data.
- Here we transformed the shape of the data to make it look like an 2D image to be used for Conv2D.
- We plotted the Confusion Matrix and ROC curve for each model to analyze the performance of the models for the given data.
- We obtained the accuracy, recall, precision and F1-score of all the models, and the comparison can be seen in the table below.

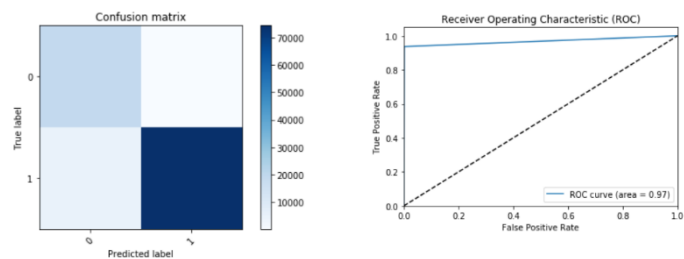
## 3. Experimental Results and Analysis

- We observed that when we consider the full data set, for every model that we run, we got 99% accuracy.
- We tried to analyze whether it is because of the model tuning or the redundancy in the data set.
- We are showing below the confusion matrix and the ROC curve of the Gaussian Naïve Bayes and CNN.
- The confusion matrix and ROC curve of all the other models are like CNN.

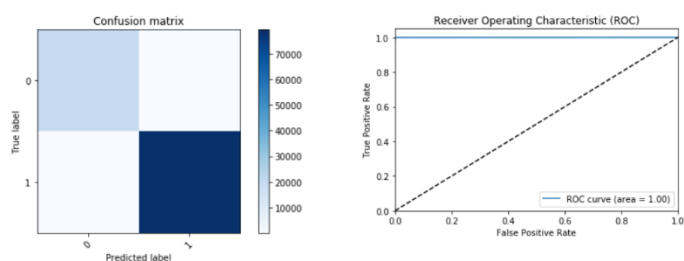
Model & Tuning	Accuracy	Precision	Recall	F1 Score
Logistic Regression	0.99462572365 49128	0.994682832303 6672	0.99462572365 49128	0.9946408290 044061
KNN	0.99955467389 98421	0.999554935399 6092	0.99955467389 98421	0.9995547520 894394
SVM	0.99833002712 44079	0.998331562815 7213	0.99833002712 44079	0.9983305969 065448
Gaussian NB	0.94890895105 46131	0.959301763742 0566	0.94890895105 46131	0.9509965159 983609
<b>Fully Connected NN</b>				
ReLU + adam + early stopping and checkpoint + 4 hidden layers	0.99928140561 11089	0.999281343572 176	0.99928140561 11089	0.9992813705 348386
Sigmoid + sgd + early stopping and checkpoint + 4 hidden layers	0.99493947613 45695	0.995049473076 8633	0.99493947613 45695	0.9949616324 501922
Tanh + adagrad + early stopping and checkpoint + 4 hidden layers	0.99925104246 79163	0.999251410520 4905	0.99925104246 79163	0.9992511739 686027
<b>CNN</b> conv2d_1 (41 , ReLU, (1,3), (1,1))+ conv2d_2 (82 , ReLU, (1,3))+ max_pooling2d_1(1,2) + dropout_1 (0.25)+ flatten_8 + dense_1 (164, ReLU) + dropout_2 (0.5) + dense_2 (2) output.	0.99941297923 16101	0.999413496087 5988	0.99941297923 16101	0.9994131280 64702

## Confusion Matrix and ROC

- Gaussian Naïve Bayes



- CNN



## 4. Task Division

### 4.1. Chandini Nagendra:

- Logistic Regression
- Gaussian NB
- Fully Connected Neural Networks
- Report

### 4.2. Siddharth Chittora

- SVM
- KNN
- CNN
- Report

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

## 5. Project Reflection

- All the models have an outstanding accuracy of 99%.
- We tried to remove null values and found out that there were no null values in the dataset.
- We removed redundancy which reduced the dataset significantly.
- Even after the removing redundancy the accuracy stayed at 99%.
- Though the accuracy remained the same we observed that the dataset with redundant record are highly skewed towards intrusion or attacks as an outcome.
- After the removal of redundant records, we observed that maximum of the connections was normal.
- Since we did not have complete domain knowledge, we researched existing papers on data mining for network intrusion detection.
- We were surprised to find that features suggested in the papers were almost similar to the features selected after Feature Importance Analysis using tree-based estimators.
- It was challenging to perceive the data as an image, whether to use it as greyscale image of size 1X41 or color image of size 1X1 with 41 channels.
- Therefore, we chose to go with the feature selected by the tree-based estimator.
- We experimented with a new CNN architecture to see how it would impact the accuracy of the model, because every other model that we worked with gave 99% accuracy.

## 6. Additional Features

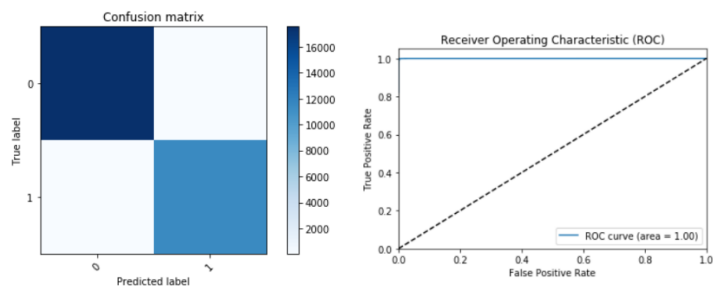
- The data set had a lot redundancy which skewed the results.
- We removed redundancy from the entire dataset and ran all the models.
- We performed feature importance analysis using tree-based estimators which computes the feature importance
- We coupled it with the Sklearn feature selection meta-transformer to discard the irrelevant features and reduced the number of features to train the models.
- We also researched to see if there was an existing proven CNN architecture for this type of problem.
- We found an architecture which we have implemented in the additional features section.
- The implemented model has sixteen layers, which is a mix of eight Conv2d layers, three Max-pooling layers, one flatten layer, two hidden layers and one dropout layer.

Model & Tuning	Accuracy	Precision	Recall	F1 Score
<b>Redundancy Removal</b>				
Logistic Regression	0.98763608888 27832	0.987655519814 2505	0.98763608888 27832	0.9876232148 457863
<b>KNN</b>	<b>0.99835147851 77044</b>	<b>0.998351438697 4491</b>	<b>0.99835147851 77044</b>	<b>0.9983513798 219226</b>
SVM	0.99220386715 66439	0.992202586638 6067	0.99220386715 66439	0.9922020510 243128
Gaussian NB	0.91829515403 3726	0.926633232685 079	0.91829515403 3726	0.9163503763 232373
<b>Fully Connected NN</b>				
<b>ReLU + adam + early stopping and checkpoint + 4 hidden layers</b>	<b>0.99852319950 54435</b>	<b>0.998523328125 6618</b>	<b>0.99852319950 54435</b>	<b>0.9985230557 019816</b>
Sigmoid + sgd + early stopping and checkpoint + 4 hidden layers	0.98633100937 5966	0.986336941032 662	0.98633100937 5966	0.9863197626 089032
Tanh + adagrad + early stopping and checkpoint + 4 hidden layers	0.99807672493 73218	0.998076683117 1078	0.99807672493 73218	0.9980765809 54231
<b>CNN</b> conv2d_1 (41 , ReLU, (1,3), (1,1))+ conv2d_2 (82 , ReLU, (1,3))+ max_pooling2d_1(1,2) + dropout_1 (0.25)+ flatten_8 + dense_1 (164, ReLU) + dropout_2 (0.5) + dense_2 (2) output	<b>0.99941297923 16101</b>	<b>0.999413496087 5988</b>	<b>0.99941297923 16101</b>	<b>0.9994131280 64702</b>
<b>Feature Selection based on Feature Importance Analysis</b>				
Logistic Regression	0.98372085036 23313	0.983778415803 2781	0.98372085036 23313	0.9836938095 343001
<b>KNN</b>	<b>0.99811106913 48697</b>	<b>0.998111899331 3256</b>	<b>0.99811106913 48697</b>	<b>0.9981107146 517584</b>

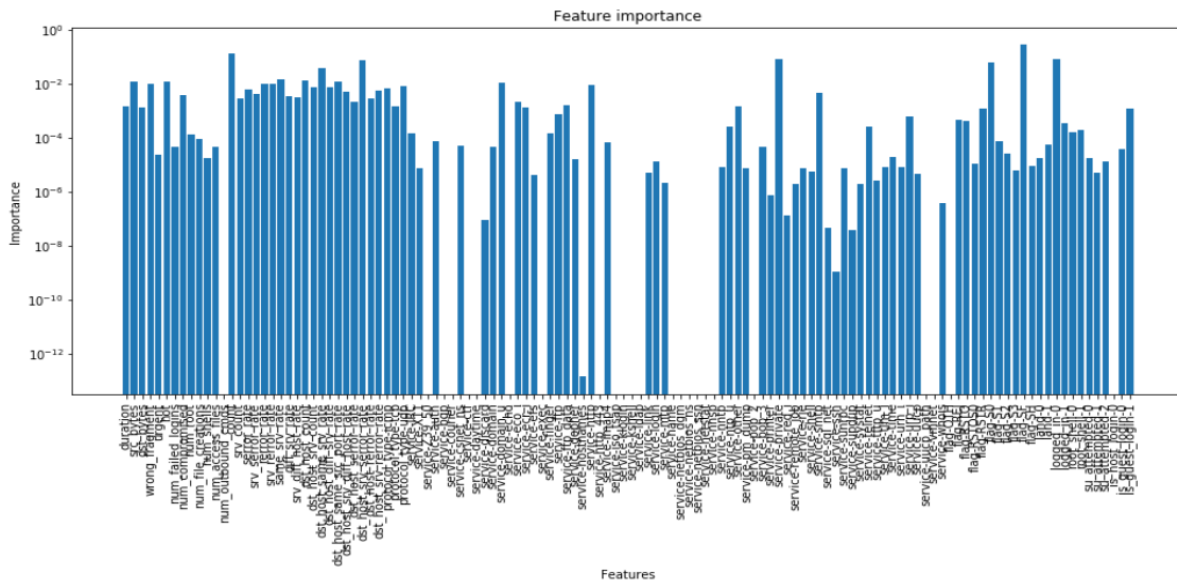
SVM	0.98052683999 03836	0.980974768826 7506	0.98052683999 03836	0.9804479756 98628
Gaussian NB	0.97829446714 97751	0.978595264164 6312	0.97829446714 97751	0.9782192710 782204
<b>Fully Connected NN</b>				
ReLU + adam + early stopping and checkpoint + 4 hidden layers	0.99745852938 1461	0.997460551139 1783	0.99745852938 1461	0.9974589455 556856
Sigmoid + sgd + early stopping and checkpoint + 4 hidden layers	0.98691486073 42789	0.986922766611 9975	0.98691486073 42789	0.9869037894 03367
Tanh + adagrad + early stopping and checkpoint + 4 hidden layers	0.99704639901 08871	0.997051160023 4251	0.99704639901 08871	0.9970471884 585685
<b>CNN</b> conv2d_1 (41 , ReLU, (1,3), (1,1))+ conv2d_2 (82 , ReLU, (1,3))+ max_pooling2d_1(1,2) + dropout_1 (0.25)+ flatten_8 + dense_1 (164, ReLU) + dropout_2 (0.5) + dense_2 (2) output	0.99773328296 18436	0.997734941418 7415	0.99773328296 18436	0.9977336205 22209
<b>Experimental CNN Architecture</b>				
<b>CNN</b> 8 Conv2d layers + 3 Max-pooling layers + 1 flatten layer + 2 Dense hidden layers + 1 dropout layer	0.99866057629 56349	0.998660731732 4747	0.99866057629 56349	0.9986604458 692393

## Confusion Matrix and ROC

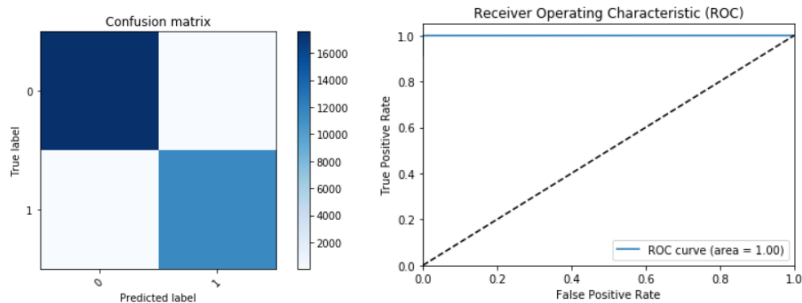
- CNN after Redundancy Removal



- Feature Importance Analysis Plot



- CNN after Feature Selection



## 7. References

[1] Md Moin Uddin Chowdhury and *et. al.* "A Few-shot Deep Learning Approach for Improved Intrusion Detection", IEEE UEMCOM 2017, October 2017