

# Evaluation of Classification algorithms for Distributed Denial of Service Attack Detection

Maulik Gohil  
Sathish Kumar Ph.D.  
(Cleveland State University)

# Outline



Introduction



Term paper objective



Dataset Explanation



Detection Approach - Data Analysis Methods



Experimentation



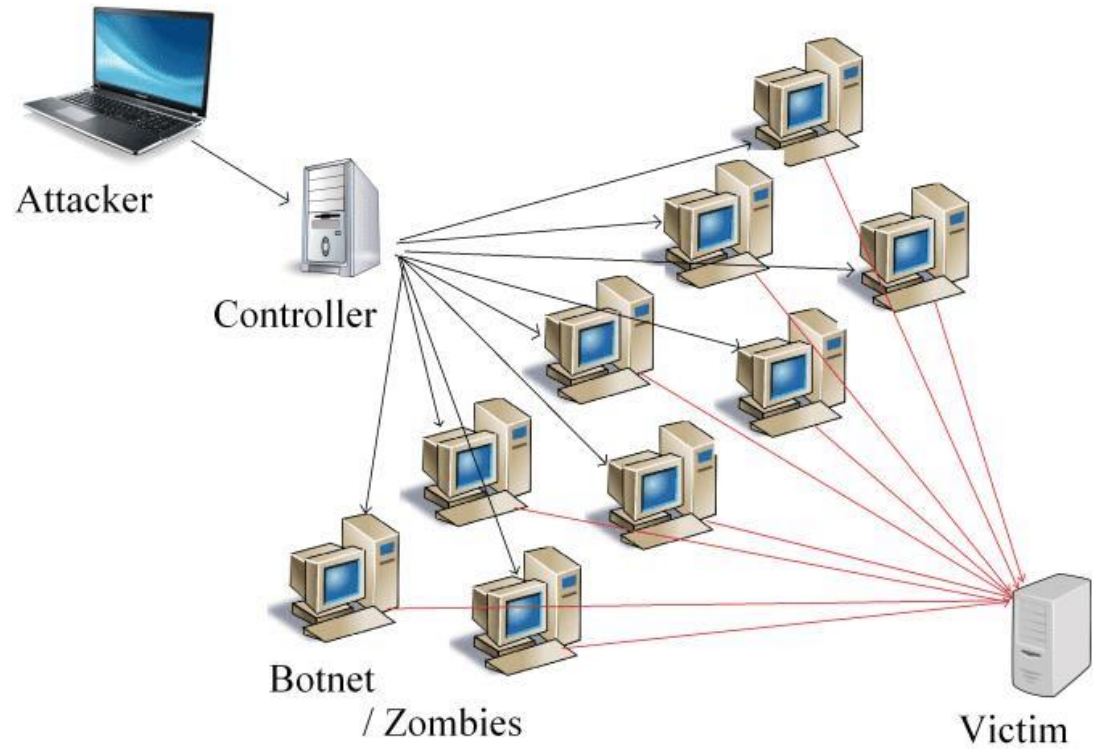
Analysis Results



Conclusion and Future work

# What is DDoS Attack?

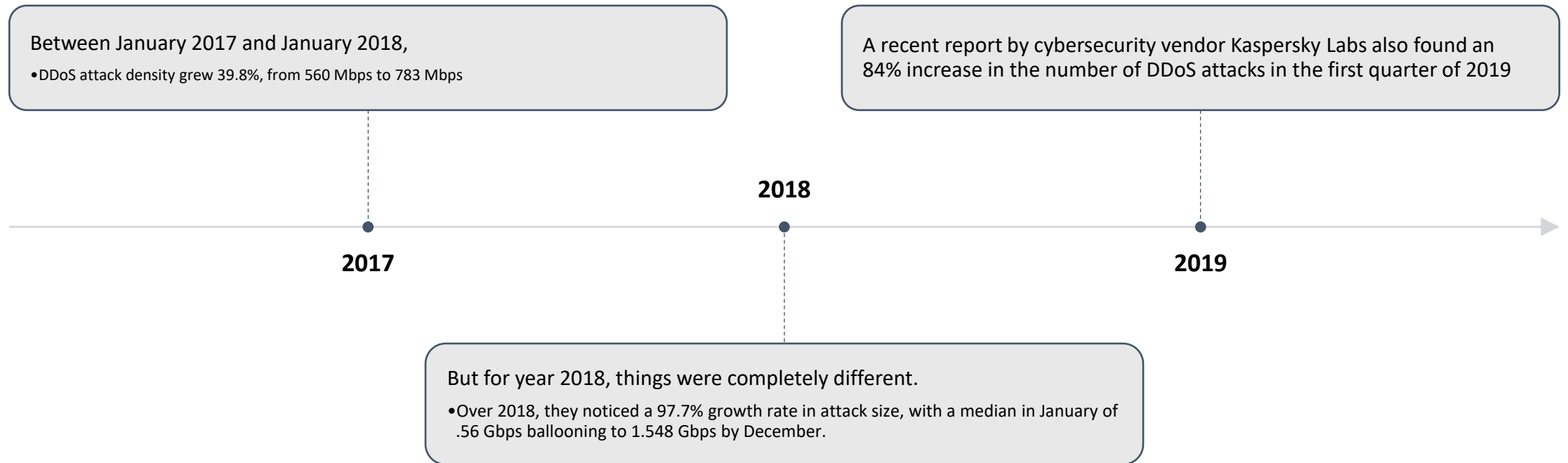
Source: Google Images



# How has it started?

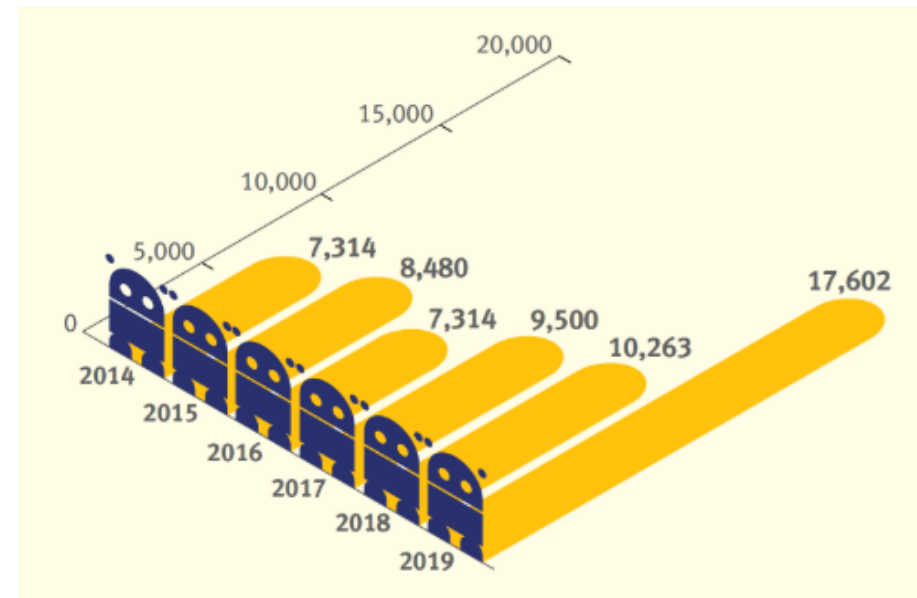
- In July 1999, a set of computers infected with the Trin00 malware **attacked** and took down the network of the University of Minnesota.
- The episode marked the **first** recorded case of a distributed-denial-of-service (**DDoS**) **attack**.

# Facts provided by Security Giant Akamai



# Spamhaus Report

- Spamhaus Malware Labs identified and blocked 17,602 botnet Command & Control servers hosted on 1,210 different networks in the year of 2019.
  - That is an enormous 71.5% increase from the number of botnet C&Cs seen in 2018.



<https://www.spamhaus.org/news/article/793/spamhaus-botnet-threat-report-2019>

# Impact of DDoS Attack

- A new study conducted by Corero confirmed that the erosion of customer trust and confidence is the single most damaging consequence of DDoS attacks for businesses today
- It is ranked that, the loss of customer trust and confidence as the worst effect of a DDoS attack (42%), followed by data theft (26%), potential revenue losses (13%) and intellectual property theft (10%).

# Objective



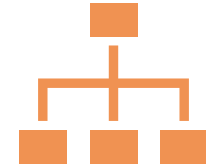
Most of the research papers published in conferences and journals have used old dataset (KDD Cup '99, DARPA) for doing their analysis which is less impactful because as the time passes, the cybercrimes and attacks are taking place in an artful way to intrude the target environment.



So, doing analysis on recent dataset which has all the variety of novel attack signatures, is much better when the security is concerned. Therefore, here I will be using the CICDDoS2019 dataset to do the analysis.



The core objective of this paper is doing data analysis with the most recent dataset specifically of DDoS attack and comparison between different classification algorithms which will be used in conducting the analysis so that it helps us to reduce the False Positives with highest accuracy.



And this will help any security Administrator/Engineer to get notified in real time that eventually helps in betterment of organization's availability-production system Uptime, as well as the reputation.



# Dataset CICDDoS2019

- It contains benign and the most up-to-date common DDoS attacks, which resembles the true real-world data (PCAPs)
  - Dataset includes 12 DDoS attacks includes NTP, DNS, LDAP, MSSQL, NetBIOS, SNMP, SSDP, UDP, UDP-Lag, WebDDoS, SYN and TFTP on the training day and 7 attacks including PortScan, NetBIOS, LDAP, MSSQL, UDP, UDP-Lag and SYN in the testing day
- The dataset has been organized per day. For each day, they recorded the raw data including the network traffic (Pcaps) and event logs (windows and Ubuntu event Logs) per machine.
- In features extraction process from the raw data, we used the [CICFlowMeter-V3](#) and extracted more than 80 traffic features and saved them as a CSV file per machine.

# More on Dataset

- The dataset package contains 7 csv files with different DDoS signatures.
- The dataset has more than 20 million of rows in total when we merge all the csv files.
- Each CSV has 88 features with label, and total size of dataset is 8+ Gigabytes.
- Dataset is created by University of New Brunswick and publicly available on the Canadian Institute of Cyber security website.

CSV File Name	Total Rows	Benign Rows
LDAP	2113234	5124
MSSQL	5775786	2794
NetBIOS	3455899	1321
Syn	4320541	35790
UDP	3782206	3134
UDPLag	725165	4068
Portmap	191694	4734
<b>Total</b>	<b>20364525</b>	<b>56965</b>

# Feature Selection

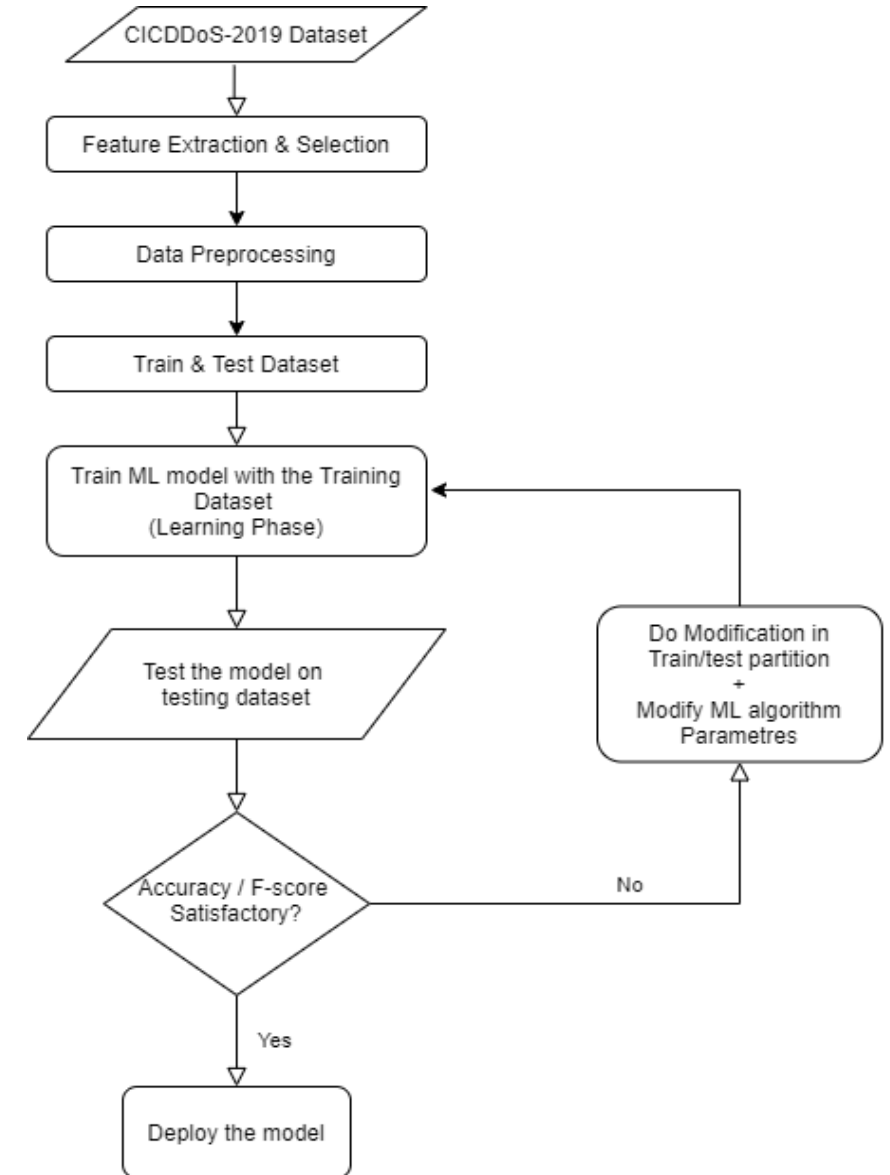
- We have used **SelectKBest** class that helped us to get most correlated features with the Class label.
- For our experiment we used top 25 correlated features as shown.

Feature Name	Description
<i>Source_IP</i>	Source IP address from where attack has been initiated
<i>Source_Port</i>	Source Port number
<i>Destination_IP</i>	IP Address of Target Machines
<i>Destination_Port</i>	Port number of Target Machine
<i>Timestamp</i>	Timestamp of the Packet has been received from Target machine
<i>Protocol</i>	Which protocol has been exploited for the DDoS Attack
<i>Flow_Duration</i>	Duration of the flow in Microsecond
<i>Total_Fwd_Packets</i>	Total packets in the forward direction
<i>Total_Backward_Packets</i>	Total packets in the backward direction
<i>Total_Length_of_Fwd_Packets</i>	Total size of packet in forward direction
<i>Total_Length_of_Bwd_Packets</i>	Total size of packet in backward direction
<i>Fwd_Packet_Length_Mean</i>	Mean size of packet in forward direction
<i>Bwd_Packet_Length_Mean</i>	Mean size of packet in backward direction
<i>Flow_IAT_Mean</i>	Mean time between two packets sent in the flow
<i>Fwd_IAT_Mean</i>	Mean time between two packets sent in the forward direction
<i>Bwd_IAT_Mean</i>	Mean time between two packets sent in the backward direction
<i>Fwd_Header_Length</i>	Total bytes used for headers in the forward direction
<i>Bwd_Header_Length</i>	Total bytes used for headers in the backward direction
<i>Packet_Length_Mean</i>	Mean length of a packet
<i>Fwd_Avg_Bulk_Rate</i>	Average number of bulk rate in the forward direction
<i>Bwd_Avg_Bulk_Rate</i>	Average number of bulk rate in the backward direction
<i>Down_Up_Ratio</i>	Download and upload ratio
<i>Average_Packet_Size</i>	Average size of packet
<i>Subflow_Fwd_Packets</i>	The average number of packets in a sub flow in the forward direction
<i>Subflow_Bwd_Packets</i>	The average number of packets in a sub flow in the backward direction
<i>Inbound Label</i>	Traffic is inbound or not Class Label

# Detection Approaches

- We are analyzing the dataset with 6 different classification algorithms
  - Decision Tree
  - Naïve Bayes
  - Logistic Regression
  - SVM
  - K-NN
  - Random Forest

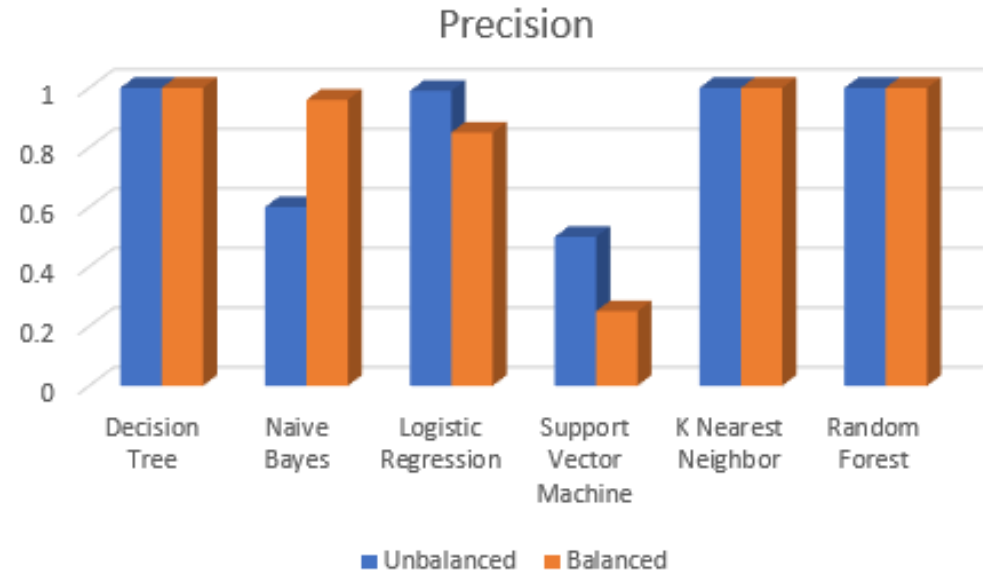
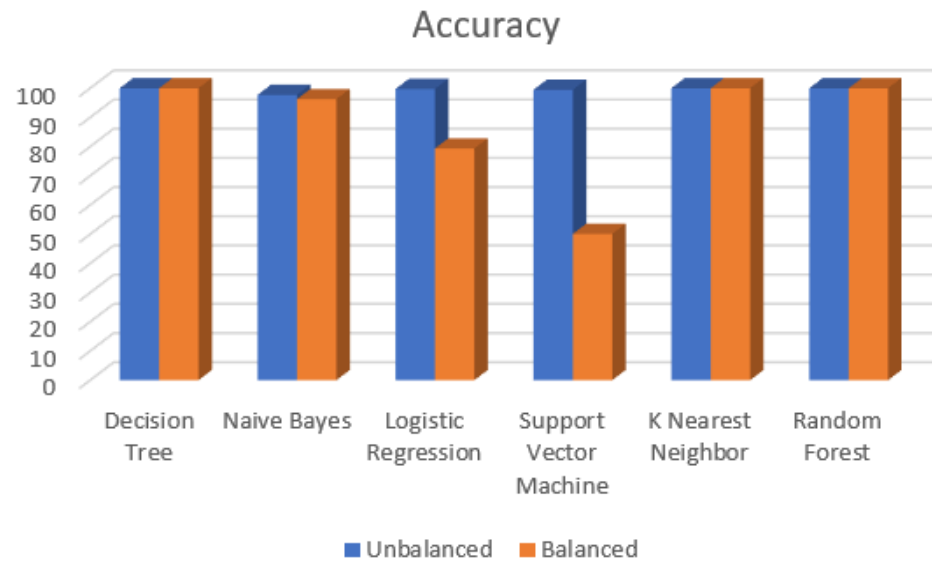
# Classification Process



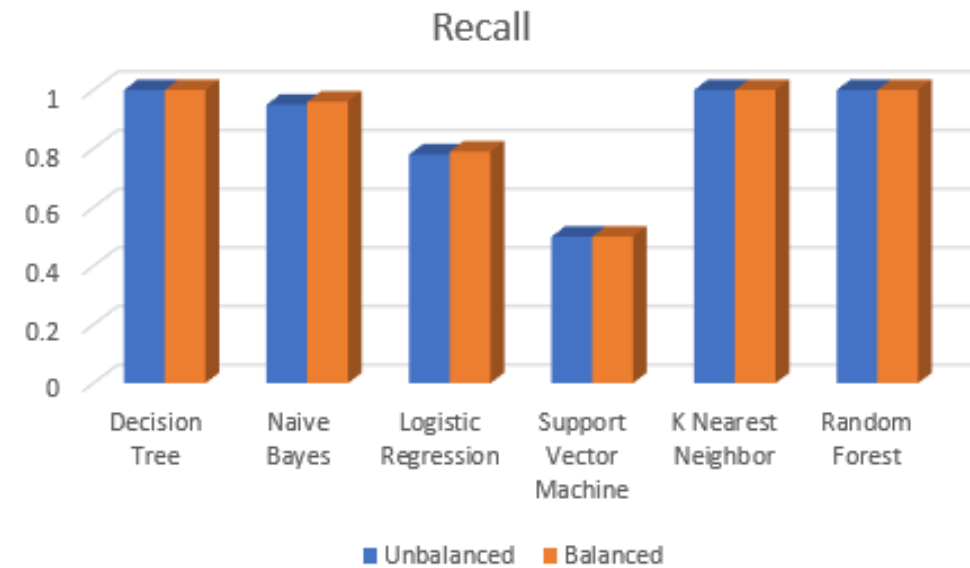
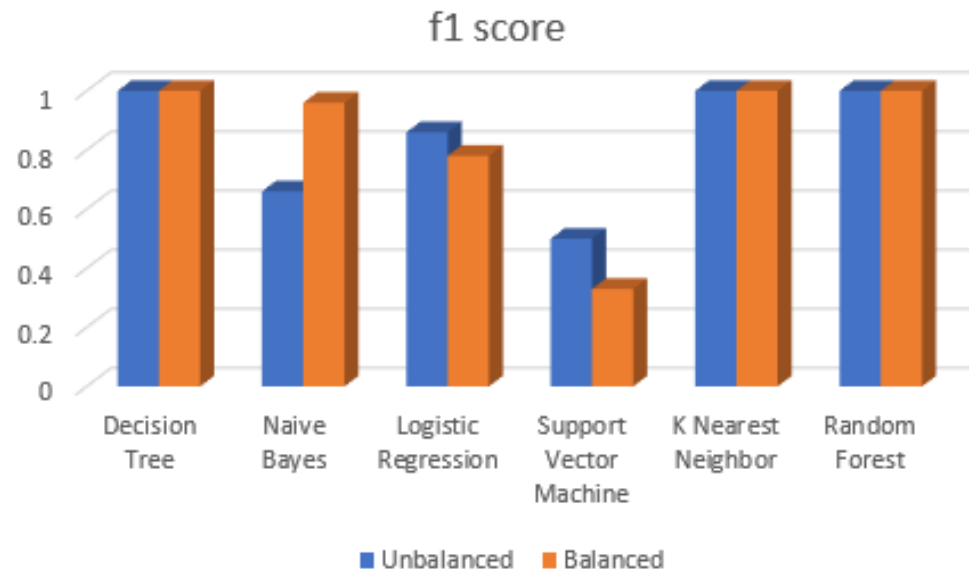
# Experimentation

- We have done two experiment one with Balanced Dataset and another with unbalanced Dataset.
- Both dataset consist 200K rows for the experiment.
- We kept train test split ratio to 70:30
- For the balanced dataset, we only had 56965 rows of normal traffic, so We randomly chose 50K attack traffic and 50K benign traffic and then appended the same data to reach 200K rows
- But for unbalanced dataset, we just randomly chosen 200K rows from the dataset. But this was quite biased.

# Result -Charts



# Result -Charts





# Result – Unbalanced Dataset

						macro avg		
Unbalanced Dataset	TP	TN	FN	FP	Accuracy	Precession	Recall	F1 Score
Decision Tree	62599	398	3	0	99.99523	1	1	1
Naive Bayes	61199	370	31	1400	97.72857	0.6	0.95	0.66
Logistic Regression	62619	213	164	4	99.73333	0.99	0.78	0.86
Support Vector Machine	62663	0	337	0	99.46507	0.5	0.5	0.5
K Nearest Neighbor	62598	401	0	1	99.99841	1	1	1
Random Forest	62602	397	0	1	99.99841	1	1	1

# Result – Balanced Dataset

						macro avg		
Balanced Dataset	TP	TN	FN	FP	Accuracy	Precession	Recall	F1 Score
Decision Tree	31577	31449	0	0	100	1	1	1
Naive Bayes	31387	29278	2290	71	96.25392	0.96	0.96	0.96
Logistic Regression	31276	8730	12819	201	79.34185	0.85	0.79	0.78
Support Vector Machine	31577	0	31449	0	50.10154	0.25	0.5	0.33
K Nearest Neighbor	31477	31549	0	0	100	1	1	1
Random Forest	31477	31549	0	0	100	1	1	1

# Conclusion

- We have used CICDDoS2019 dataset which is fairly a recent dataset that includes most recent attack signatures for DDoS.
- The experimentation has been carried out using major supervised classification algorithms to classify the attack accurately from the legitimate flows
- When the results are compared with other algorithms among all classifiers, decision tree, random forest & K-NN performed the best.

# Future work

- One can so the same experiment with bigger scale in terms of data.
- We can perform the same experiment with different feature selection techniques.
- Rather than applying classification algorithms directly on mixed data having all the attack signatures, we can keep it separate and then can-do data analysis applying all the six-classification algorithms separately.
  - It might be possible that some attack can be efficiently identified by specific ML algorithm.