

# MTA-KDD'19: A Dataset for Malware Traffic Detection

University of L'Aquila



The logo for ITA SEC 2020, the Italian Conference on Cybersecurity. It features the text "ITA SEC 20" in large red letters, with "ITA" in white and "SEC 20" in red. Below it, in smaller red text, is "ITALIAN CONFERENCE ON CYBERSECURITY". Underneath that, in larger red text, is "Ancona, 4-7 February 2020". At the bottom is the logo for "cini Cybersecurity National Lab", which includes a stylized map of Italy and the word "cini" in blue.

Authors:

Ivan Letteri

Giuseppe Della Penna

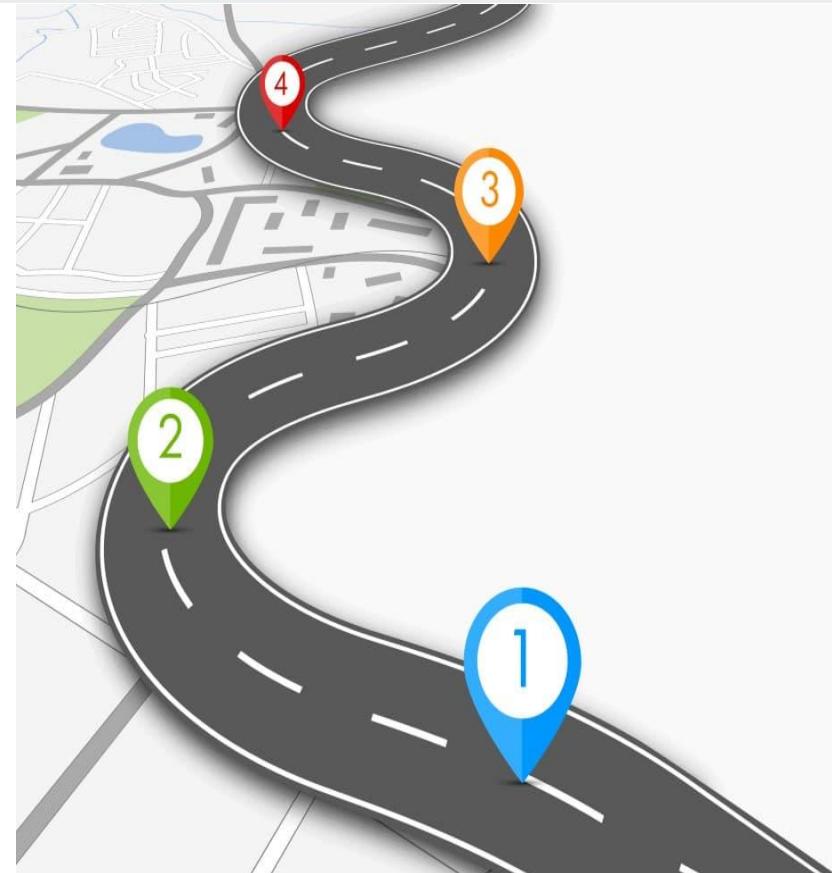
Luca Di Vita

Maria Teresa Grifa



# Road Map

- **1st - MTA-KDD'19 Dataset**
  - legitimate traffic
  - malware traffic
- **2nd - Dataset Features**
  - from 50 to 33 features
- **3rd - Dataset Generation Process**
  - *phase I.* Feature Extraction
  - *phase II.* Dataset Profiling
  - *phase III.* Imputation
  - *phase IV.* Scaling
- **4th - Dataset Evaluation**
  - Outlier Detection
  - Classification Experiment





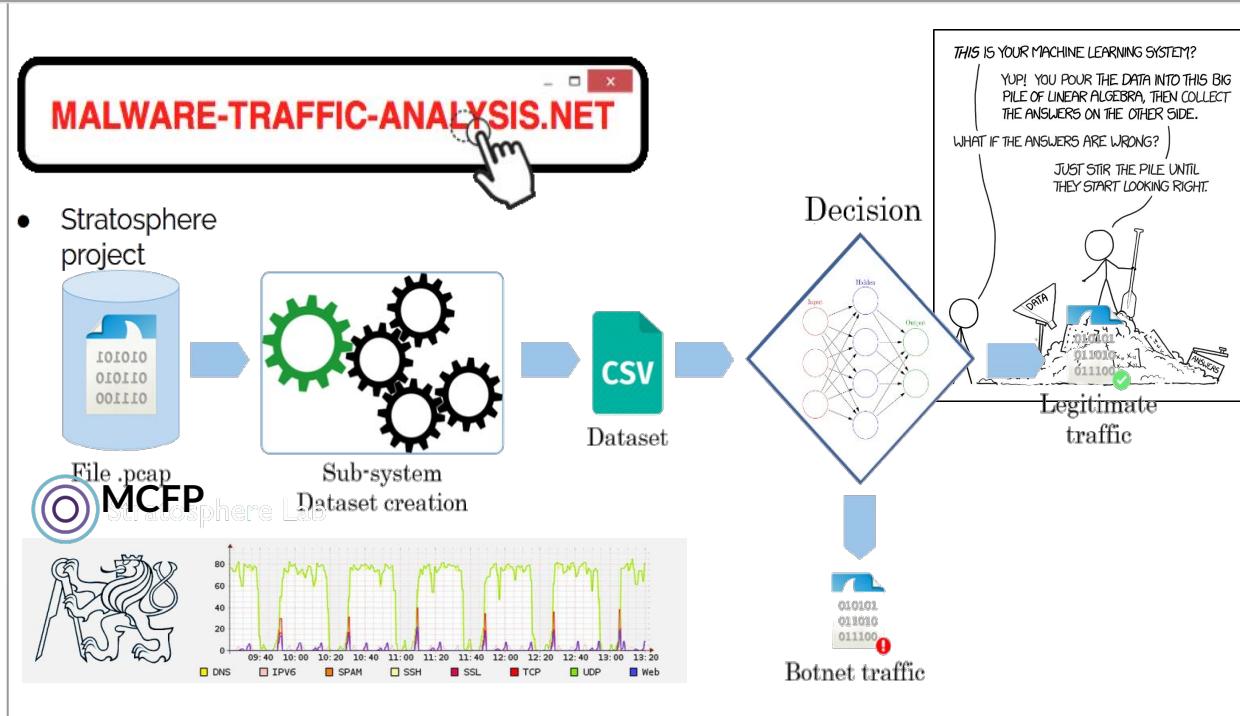
# Generation Process

- Network logs

- Feature Extraction phase from pcap files

- Legitimate traffic

- Malware traffic



- Legitimate traffic come from pcap files marked as Normal in MCFP composed of 15 pcap files with a total size greater than 7 GBytes which produced  $\approx$  45 thousand samples
- Malware traffic come from the MTA repository, is made up of  $\approx$  2112 pcap files, with a total size of more than 4.8 GB. These observations cover a time span from June 2013 to August 2019.



# Generation Process ( 50 Dataset Features )

1. {Ack,Syn,Fin,Psh,Urg,Rst}FlagDist
2. {TCP,UDP,DNS}OverIP
3. MaxLen, MinLen, AvgLen,  
StdDevLen, MaxIAT, MinIAT, ...
4. PktIORatio
5. FirstPktLen:
6. DNSQD, DNSAD, DNSRD, ...
7. {RepeatedPkt, SmallPkt} Ratio
8. AvgDomain{Char, Dot, Hyph, Digit}
9. AvgTTL
10. NumConnections, SynAckSyn
11. NumDstAddr, NumPorts
12. DistinctUA, AvgDistinctUALen, ...

feature	formula	notes	ref	rem
$f\text{FlagDist}$	$\frac{ \mathcal{S}^{TCP} }{ \mathcal{S} } \cdot  \mathcal{S}^f $	with $f \in \{\text{Ack, Syn, Fin, Psh, Rst}\}$ available if not TCP	(1)	-
$p\text{OverIP}$	$\frac{ \mathcal{S}^P }{ \mathcal{S} }$	with $P \in \{\text{TCP, UDP, DNS}\}$	(2)	-
$\text{AvgDeltaTime}$	$\frac{t(p_n) - t(p_0)}{ \mathcal{S} }$	zero if $ \mathcal{S}  = 1$	(3)	(1)
$\text{AvgDistinctUA}$	$\frac{\sum_{p_j \in \mathcal{S}} nchar(s)}$	NaN if $ \mathcal{S}^{HTTP}  = 0$	(12)	(3)
$\text{AvgDomainChar}$	$\frac{\sum_{p_j \in \mathcal{S}} nchar(dom(p_j))}{\sum_{p_j \in \mathcal{S}}  dom(p_j) }$	NaN if $ \mathcal{S}^{DQ}  = 0$	(8)	(3)
$\text{AvgDomainDigit}$	$\frac{\sum_{p_j \in \mathcal{S}} ndigit(dom(p_j))}{\sum_{p_j \in \mathcal{S}}  dom(p_j) }$	NaN if $ \mathcal{S}^{DQ}  = 0$	(8)	(3)
$\text{AvgDomainDot}$	$\frac{\sum_{p_j \in \mathcal{S}} ndot(dom(p_j))}{\sum_{p_j \in \mathcal{S}}  dom(p_j) }$	NaN if $ \mathcal{S}^{DQ}  = 0$	(8)	(3)
$\text{AvgDomainHyph}$	$\frac{\sum_{p_j \in \mathcal{S}} nhyp(dom(p_j))}{\sum_{p_j \in \mathcal{S}}  dom(p_j) }$	NaN if $ \mathcal{S}^{DQ}  = 0$	(8)	(3)
$\text{AvgTTL}$	$\frac{\sum_{p_j \in \mathcal{S}} TTL(p_j)}{ \mathcal{S} }$	-	(9)	(3)
$\text{DeltaTime}$	$t(p_n) - t(p_0)$	zero if $ \mathcal{S}  = 1$	(3)	-
$\text{DistinctUA}$	$ \mathcal{U} $	-	(12)	(3)
$\text{DNSAD}$	$\frac{ \mathcal{S}^{DNS} }{ \mathcal{S} } \cdot \sum_{p_j \in \mathcal{S}} ndnsans(p_j)$	-	(6)	(1)
$\text{DNSQD}$	$\frac{ \mathcal{S}^{DNS} }{ \mathcal{S} } \cdot \sum_{p_j \in \mathcal{S}} ndnsque(p_j)$	-	(6)	-
$\text{DNSRD}$	$\frac{ \mathcal{S}^{DNS} }{ \mathcal{S} } \cdot \sum_{p_j \in \mathcal{S}} ndnsadd(p_j)$	-	(6)	(1)
$\text{DNSSD}$	$\frac{ \mathcal{S}^{DNS} }{ \mathcal{S} } \cdot \sum_{p_j \in \mathcal{S}} ndnsaut(p_j)$	-	(6)	(1)
$\text{EndFlow}$	$t(p_n)$	-	(3)	(1)
$\text{FirstPktLen}$	$len(p_0)$	-	(5)	-
$\text{FlowLen}$	$\sum_{p_j \in \mathcal{S}} len(p_j)$	-	(3)	-
$\text{FlowLenRx}$	$\sum_{p_j \in \mathcal{R}} len(p_j)$	-	(3)	-
$\text{HTTPPkts}$	$ \mathcal{S}^H $	-	(12)	-
$\text{MaxLATRx}$	$\min\{ (t(p_i) - t(p_{i-1}))   p_i \in \mathcal{R} \}$	zero if $ \mathcal{R}  < 2$	(3)	(1)
$\text{MinLAT}$ , $\text{MaxIAT}$ , $\text{AvgIAT}$	$\min, \max, \text{avg of } \{ (t(p_j) - t(p_{j-1}))   p_j \in \mathcal{S} \}$	zero if $ \mathcal{S}  < 2$	(3)	-
$\text{MinIATRx}$ , $\text{AvgIATRx}$	$\min, \text{avg of } \{ (t(p_j) - t(p_{j-1}))   p_j \in \mathcal{R} \}$	zero if $ \mathcal{R}  < 2$	(3)	-
$\text{MinLen}$ , $\text{MaxLen}$ , $\text{AvgLen}$ , $\text{StdDevLen}$	$\min, \max, \text{avg, stddev of } \{ len(p_j)   p_j \in \mathcal{S} \}$	StdDevLen is NaN if $ \mathcal{S}  < 2$	(3)	-
$\text{MinLenRx}$ , $\text{MaxLenRx}$ , $\text{AvgLenRx}$ , $\text{StdDevLenRx}$	$\min, \max, \text{avg, stddev of } \{ len(p_j)   p_j \in \mathcal{R} \}$	StdDevLenRx is NaN if $ \mathcal{R}  < 2$	(3)	-
$\text{NumConnections}$	$ \mathcal{S}^{ACKSYN} $	unavailable if not TCP	(10)	-
$\text{NumDstAddr}$	$ \mathcal{D}_a $	-	(11)	-
$\text{NumPorts}$	$\max\{  \mathcal{D}_a  \}$	unavailable not TCP or UDP	(11)	-
$\text{PktIORatio}$	$\frac{ \mathcal{S} }{ \mathcal{S} }$	-	(4)	-
$\text{RepeatedPktLenRatio}$	$\frac{\max occur(\{len(p_j)   p_j \in \mathcal{S}\})}{ \mathcal{S} }$	-	(7)	-
$\text{SmallPktRatio}$	$\frac{ \mathcal{S}^{small} }{ \mathcal{S} }$	-	(7)	-
$\text{StartFlow}$	$t(p_0)$	-	(3)	-
$\text{SynAckSynRatio}$	$\frac{ \mathcal{S}^S }{ \mathcal{S}^{ACKSYN} }$	unavailable if not TCP, NaN if $ \mathcal{S}^{ACKSYN}  = 0$	(10)	(3)

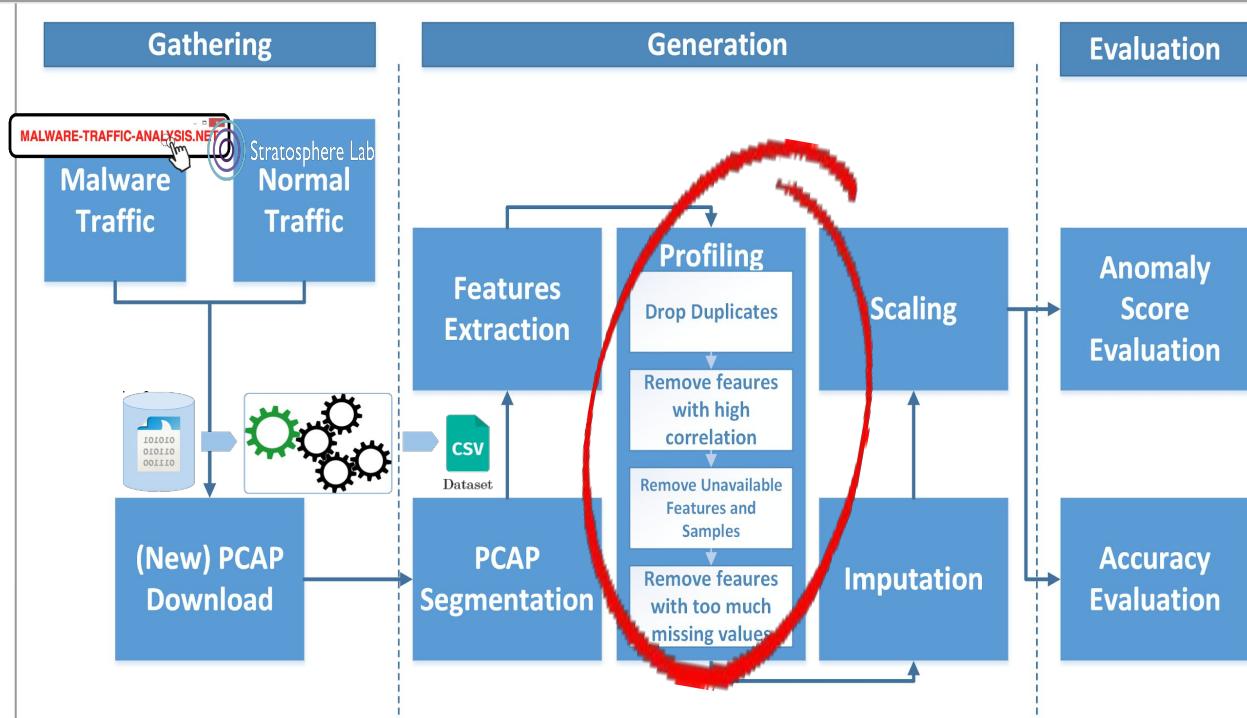
Value	Count	Frequency (%)
Malware	39544	55.3%
Legitimate	31926	44.7%



# Generation Process

## • Gathering

- where came from Dataset
  - MTA
  - MCFP
- 50 features extracted
  - dataset in CSV file format
- Next: *profiling activity*



## Dataset Profiling

- drop duplicates, missing values ...
- remove feature with high correlation, and so on ...
- ... idea to apply a minimal preprocessing on the data



# Generation Process (*Profiling*)

- Issues

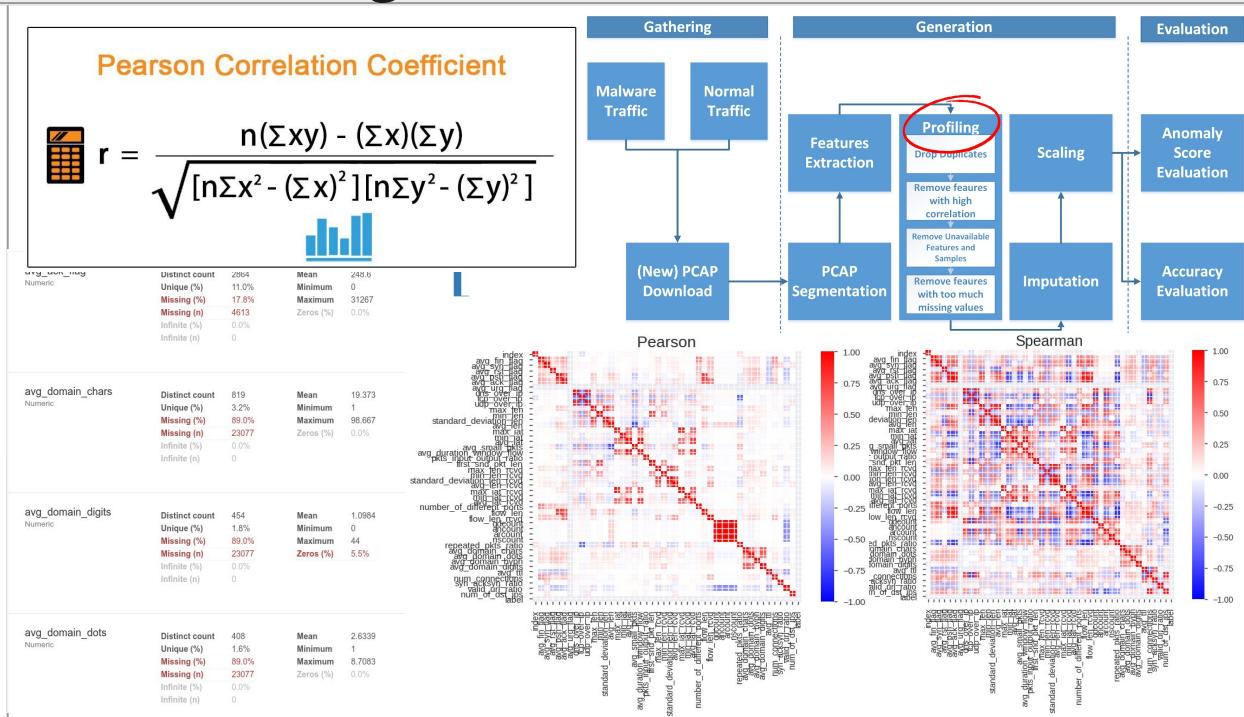
- detect low-complexity
- remove low-complexity

- High-correlation

- Pearson coefficient
- High correlation > 0.95

- Not a Number (NaN)

- ‘Real’ Zeros
- Unavailable (*incomplete*)
- Missing Values (*noisy*)

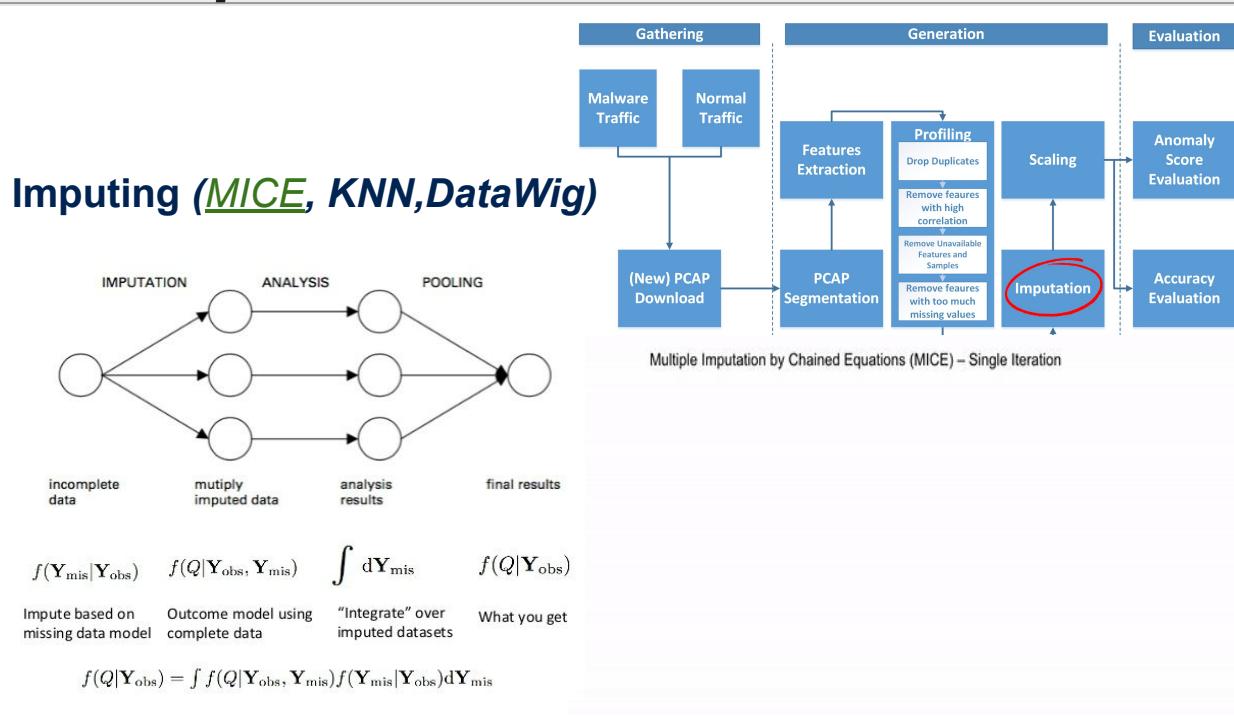


- To detect and remove certain low-complexity issues that may affect the data
- High-correlated features measured using the Pearson coefficient (from 50 features to 33 features)
- some features may have zero values for different reasons: *Zeros*, *Unavailable*, and *Missing* values



# Generation Process (*Imputation*)

- remaining NaN
  - StdDevLen feature
  - StdDevLenRx feature
  
- replacing NaN
  - with real(istic) values
  - MICE ✓
  - Datawig, KNN ✗



- MICE imputation works by filling the missing data multiple times
- Multiple imputations are better than a single imputation as they measure the uncertainty of the missing values more precisely
- chained equations approach is flexible and handle data types (i.e., continuous or binary)



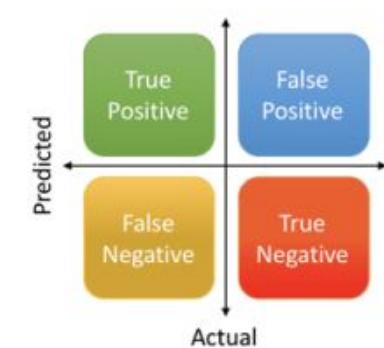
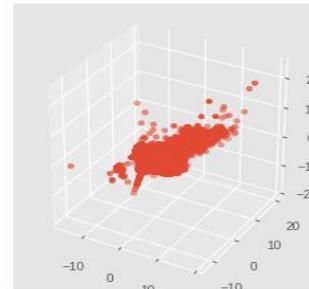
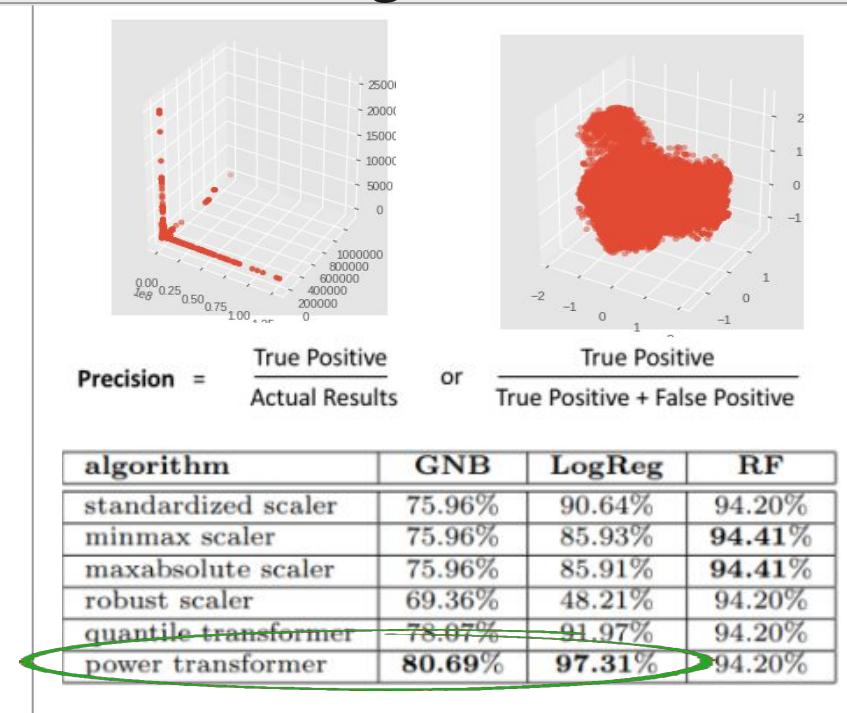
# Generation Process (Scaling)

- **4 Scalers and 2 Transformers**

- MinMax *scaler*
- MaxAbs *scaler*
- Standard *scaler*
- Robust *scaler*
- Quantile *transformer*
- Power *transformer*

- **Precision metric**

- Gaussian Naive Bayes (GNB)
- Logistic Regression (LogReg)
- Random Forest (RF)



- Dataset contains features highly varying in magnitudes, units, and range
- feature scaling may heavily influence the results of some algorithms
- to have an idea of the impact of such scaling we trained 3 “fast” classifiers (GNB, LogReg, RF)

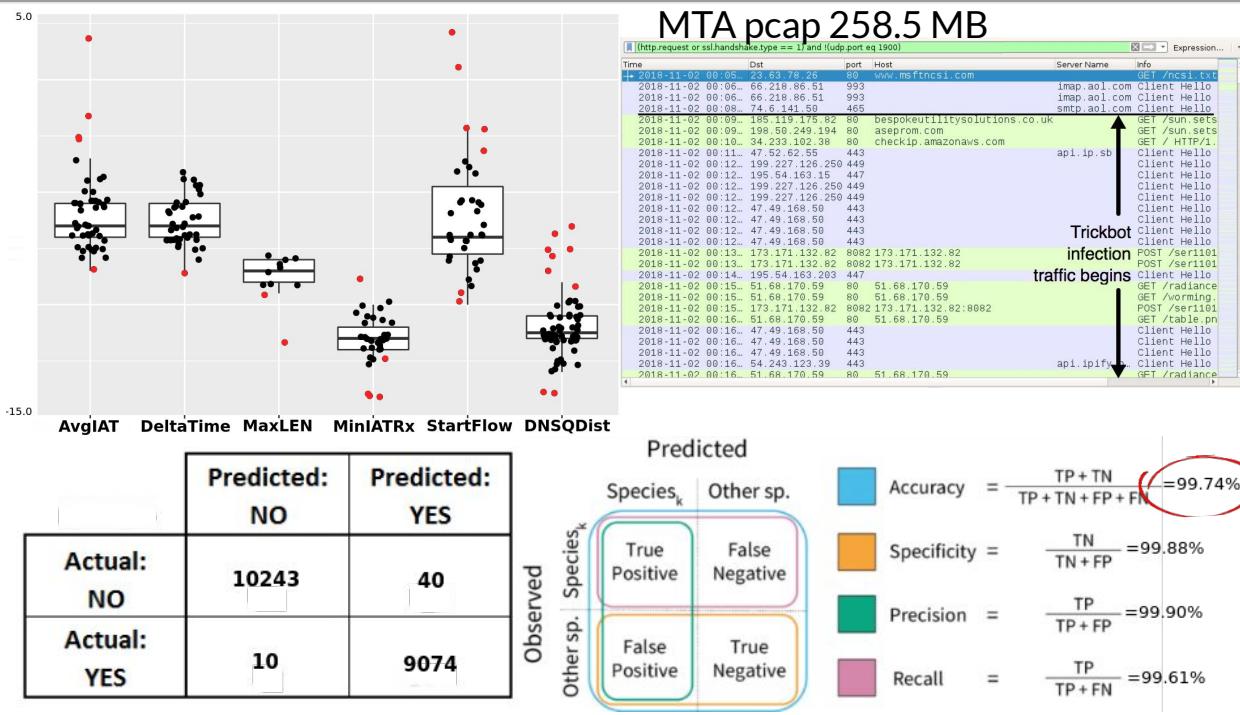


# Dataset Evaluation

## • Outlier Detection

- feature importance with Random Forest (*Gini index*)
- 6 most important features
- Matching between:

*Outlier in BoxPlot <-> pkt in PCAP*



- Further evaluation of quality looking for the matching between *Outliers* in the CSV dataset and the traffic in *PCAP* files on a subset of 6 “important” features
- Figure to the bottom shows the confusion matrix taken from one of the experiments, with an accuracy of 99.74%



# Conclusion

- MTA-KDD'19

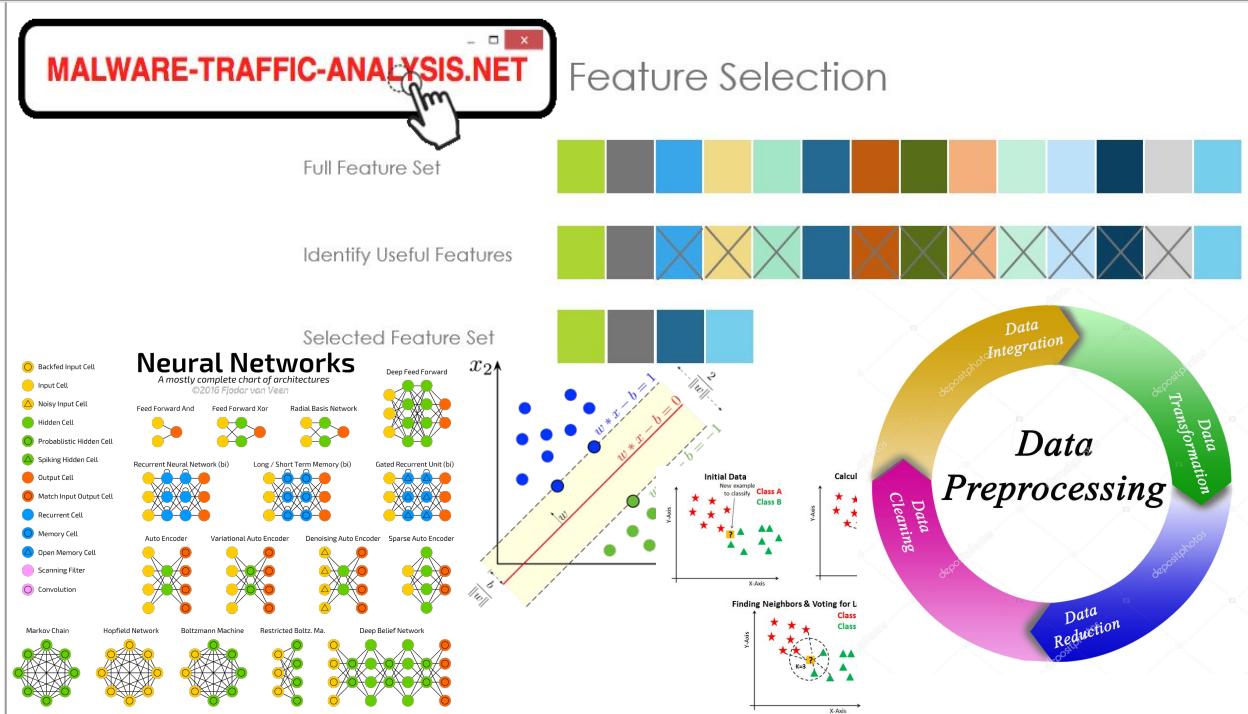
- malware traffic analysis
- Data source up-to-date
- Open source

- More Accurate

- Feature Selection
- Data preprocessing

- Future Work

- more complex features
- different NN architectures
- other ML models
- handle Unavailable values



- The data sources are continuously updated, making the dataset quite realistic
- **Accurate feature selection and data preprocessing** to make the dataset as small as possible
- **More complex feature selection strategies** to further reduce the number of features to the most informative ones.
- Validate the dataset **different neural network architectures** and **other machine learning models**.

- Reduced Dataset

- 33 features

[github.com/IvanLetteri/MTA-KDD-19](https://github.com/IvanLetteri/MTA-KDD-19)



- Source Code

- PCAP downloader
- Feature Extractor



- Entire Dataset

- Dataset 4 pkts tracking

Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data									
Dataset Overview Data</									