# Classification and visualization of web attacks
# using HTTP headers and machine learning techniques

Nicolas Ricardo Enciso

UNSECURE

# ACNS 2019

Bogotá, Colombia,
June 5-7 2019

**Call for papers**

17th International
Conference on Applied
Cryptography and
Network Security

The conference seeks submissions presenting novel research on all technical aspects of applied cryptography, cyber security (including network and computer security) and privacy. Topics of interest include but are not limited to:

Access control
Applied cryptography
Automated security analysis
Biometric security/privacy
Block chain and cryptocurrencies
Cloud security/privacy
Complex systems security
Critical infrastructure
Cryptographic primitives
Cryptographic protocols
Data protection
Database/system security
Digital rights management
Email, app and web security
Future Internet security
Human factors in security
Identity management

IP protection
Internet fraud, cybercrime
Internet-of-Things security
Intrusion detection
Key management
Malware
Mobile/wireless/5G security
Network security protocols
Privacy/anonymity, PETs
Security in e-commerce
Security in grid systems
Security in P2P systems
Security/privacy metrics
Trust management
Ubiquitous security/privacy
Usability in security/privacy

**Submission**
22 January 2019

**Notification**
22 March 2019

**Final Version**
5 April 2019

*This year there will be a 1000 EUR prize for the Best Student Paper Award sponsored by Springer*

PC co-chairs:
**Robert Deng**
Singapore Management University
**Moti Yung**
Google and Columbia University

General co-chairs:
**Valérie Gauthier**
Universidad del Rosario
**Martín Ochoa**
Universidad del Rosario

Springer

Universidad del Rosario

MACC
Aplied Mathematics
and Computer Science

www.acns19.com

---

# Conference on Applied Cryptography and Network Security ACNS 2019

# GENERAL OVERVIEW

**Abstract.** This paper presents a methodology to identify web attacks such as XSS, CRLF and SQL injection using a data set that contains normal and anomalous items. The proposed methodology uses dimensional reduction techniques for visualization (PCA, t-SNE) and machine learning algorithms (SVM, Naive Bayes, random forest, logistic regression) to perform classification of URLs contained in HTTP headers. Results show that visualization is useful to present a general overview of attacks and classification experiments show an accuracy of 83% to detect attacks.

# Data set

- 25065 attacks (anomalous) marked as 1
- 36000 normal (normal) marked as 0
- Total 61065 cases in the dataset.

UNSECURE

# Splitted data

- 70% for training
- 30% for testing
- Random election of the cases from the data
- Dataset :
    - Original : 8 features with no changes

# Feature extraction

**Sample HTTP header**

GET http://localhost:8080/tienda1/publico/anadir.jsp?id=2&nombre=Jam
User-Agent: Mozilla/5.0 (compatible; Konqueror/3.5; Linux) KHTML/3.5.8 (like Gecko)
Pragma: no-cache
Cache-control: no-cache
Accept: text/xml,application/xml,application/xhtml+xml,text/html
Accept-Encoding: x-gzip, x-deflate, gzip, deflate
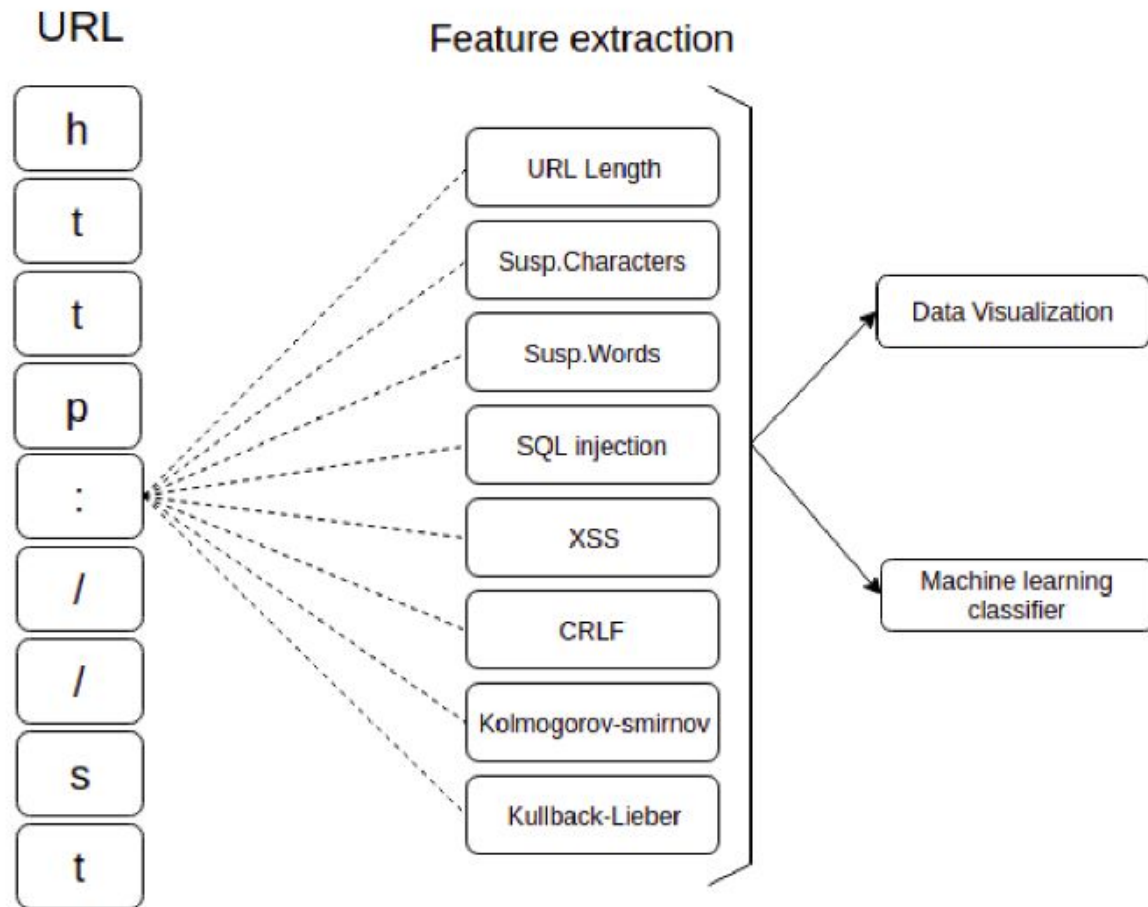Accept-Charset: utf-8, utf-8;q=0.5, *;q=0.5
Accept-Language: en
Host: localhost:8080
Cookie: JSESSIONID=B92A8B48B9008CD29F622A994E0F650D
Connection: close

# Feature extraction

# Classifier algorithms

- Naive Bayes classifier:
    - Gaussian classifier
    - Multinomial classifier
- Support Vector Machine:
    - Linear
    - Gaussian (C = 1.11 and gamma = 0.09)
    - Sigmoid
- Logistic Regression
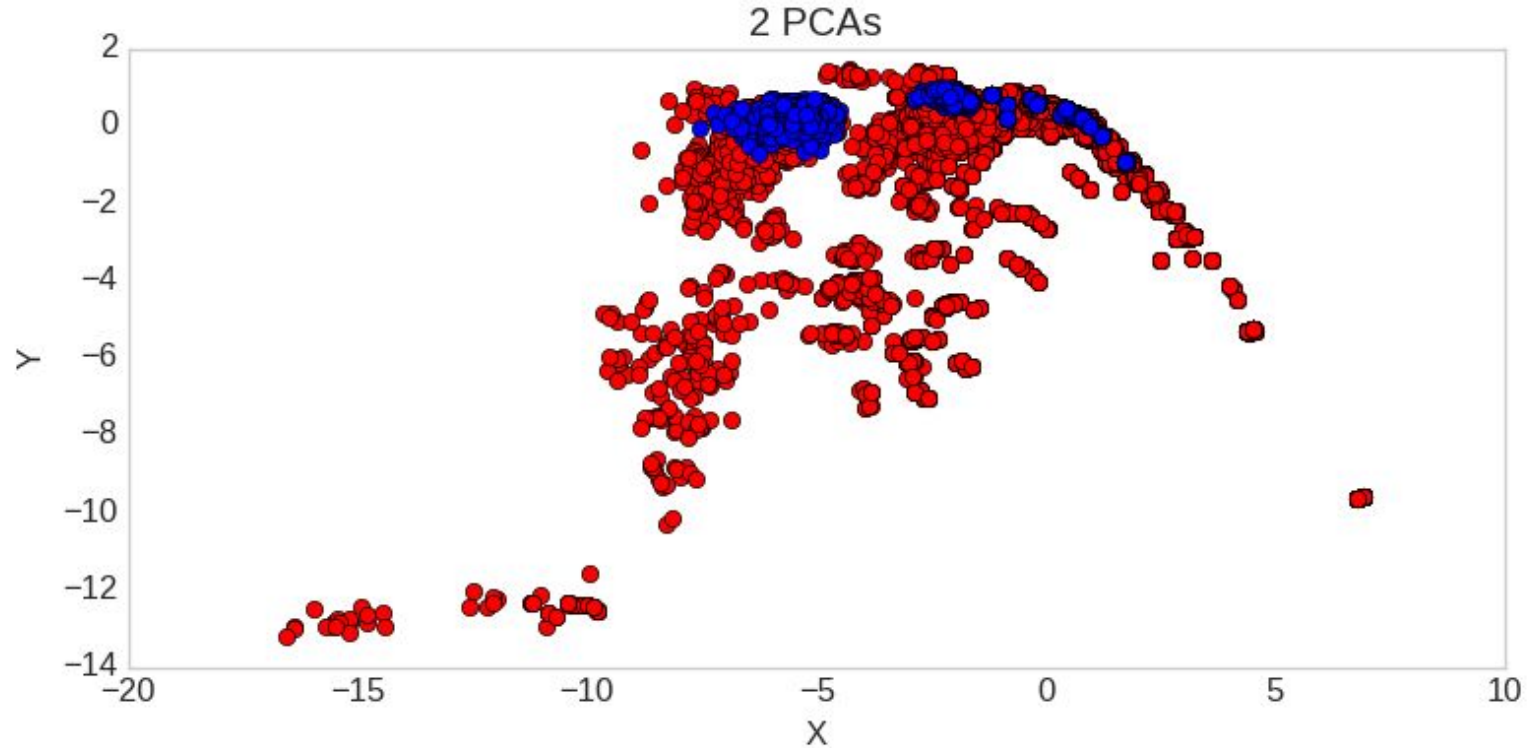- Random Forest ( estimators = 100)

# Visualization and dimensionality reduction methods

- PCA : principal component analysis


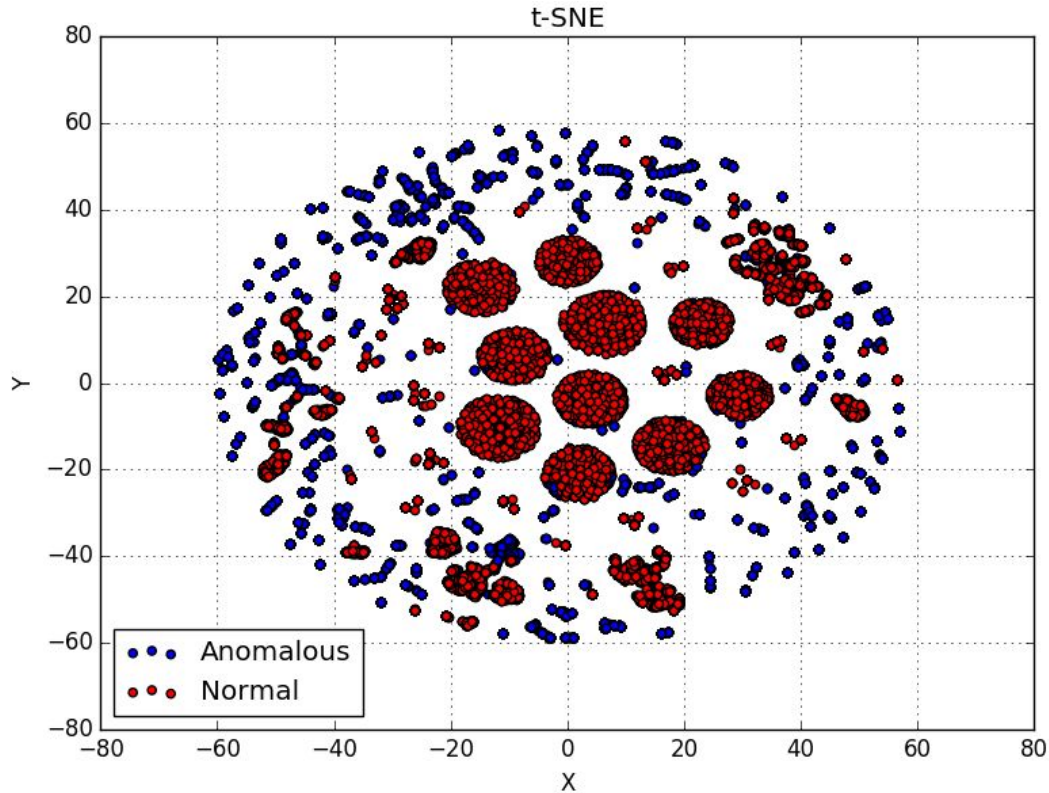- t-SNE : t-Distributed Stochastic Neighbor Embedding

# Types of attacks

- XSS cross site scripting
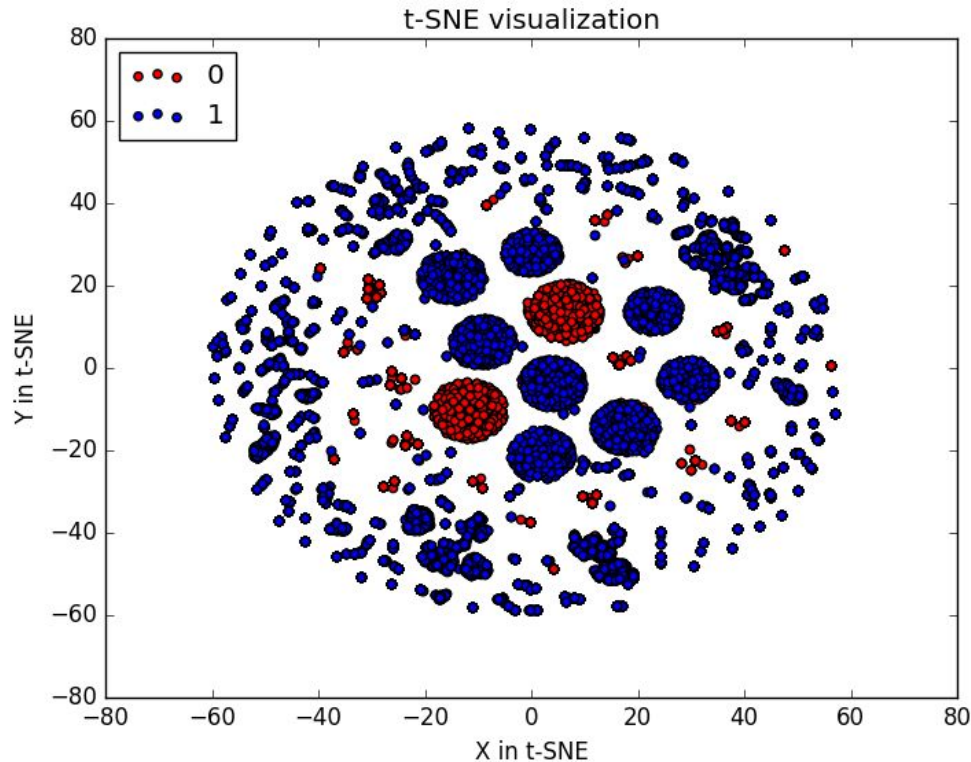- CRLF Carriage return line feed
- SQL injection
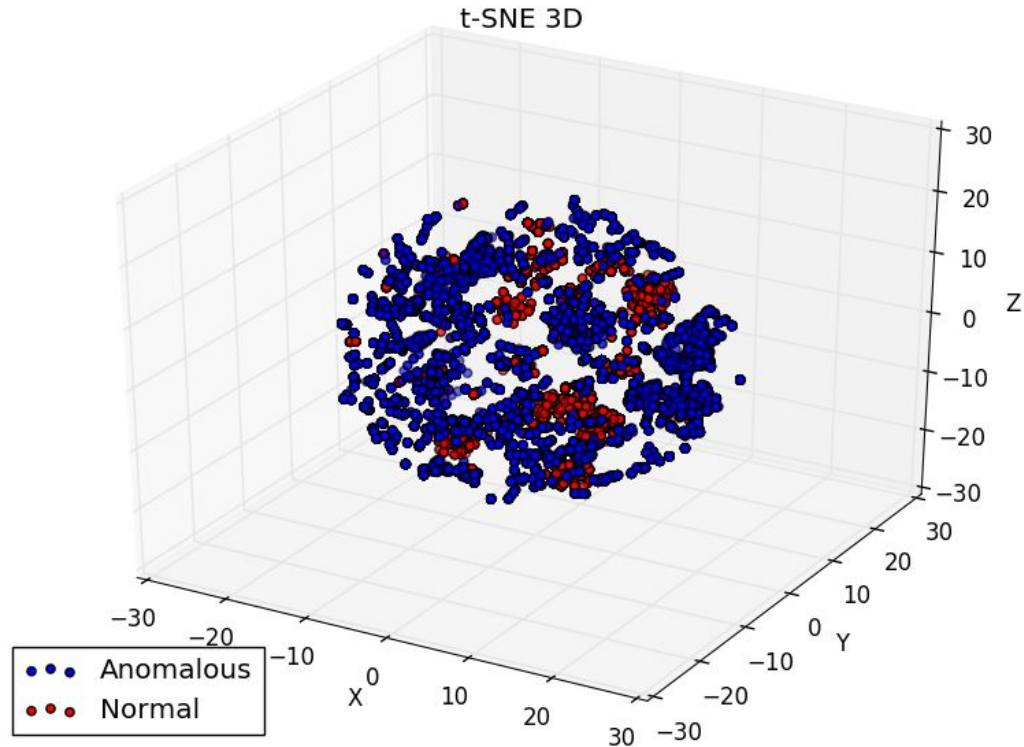
# Visualization data: PCA (Anomalous in blue)

# Visualization data: t-SNE 2D

# Visualization data: t-SNE (1 anomalous - 0 normal

# Visualization data: t-SNE 3D



t-SNE 3D

# Results

# Results performance classification

| | Precision | Recall | F1-Score | Accuracy | AUC |
|---|---|---|---|---|---|
| | | Anomalous predictions Original data | | | |
| **NB + Gaussian** | 0,720 | 0,200 | 0,320 | 0,643504366812 | 0,657966634589 |
| **NB + Multinomial** | 0,750 | 0,230 | 0,350 | 0,653056768559 | 0,650886405340 |
| **SVM + Sigmoid** | 0,000 | 0,000 | 0,000 | 0,592139737991 | 0,390236405688 |
| **SVM + linear** | 0,810 | 0,260 | 0,390 | 0,674672489083 | 0,693580346245 |
| **SVM + Gaussian** | **0,830** | 0,620 | 0,710 | 0,794596069869 | 0,870937300964 |
| **Logistic Regression** | 0,700 | 0,410 | 0,520 | 0,688373362445 | 0,707206599566 |
| **Random Forest** | 0,750 | **0,880** | **0,810** | **0,832860262009** | **0,933937629906** |

# Results performance classification

| | Precision | Recall | F1-Score | Accuracy | AUC |
|---|---|---|---|---|---|
| | | Normal cases predictions Original data | | | |
| NB + Gaussian | 0,630 | 0,950 | 0,760 | 0,643504366812 | 0,657966634589 |
| NB + Multinomial | 0,640 | 0,950 | 0,760 | 0,653056768559 | 0,650886405340 |
| SVM + Sigmoid | 0,590 | **1,000** | 0,740 | 0,592139737991 | 0,390236405688 |
| SVM + linear | 0,650 | 0,960 | 0,780 | 0,674672489083 | 0,693580346245 |
| SVM + Gaussian | **0,780** | 0,910 | 0,840 | 0,794596069869 | 0,870937300964 |
| Logistic Regression | 0,680 | 0,880 | 0,770 | 0,688373362445 | 0,707206599566 |
| Random Forest | 0,910 | 0,800 | **0,850** | **0,832860262009** | **0,933937629906** |

# ROC



ROC curve ML classifier