# ML based network traffic analysis on UNSW-NB15 public dataset

Francesco Bocchi
*Department of Computer Engineering - Cybersecurity*

# Overview

- Introduction to NB15

- Preprocessing

- Multi class classification
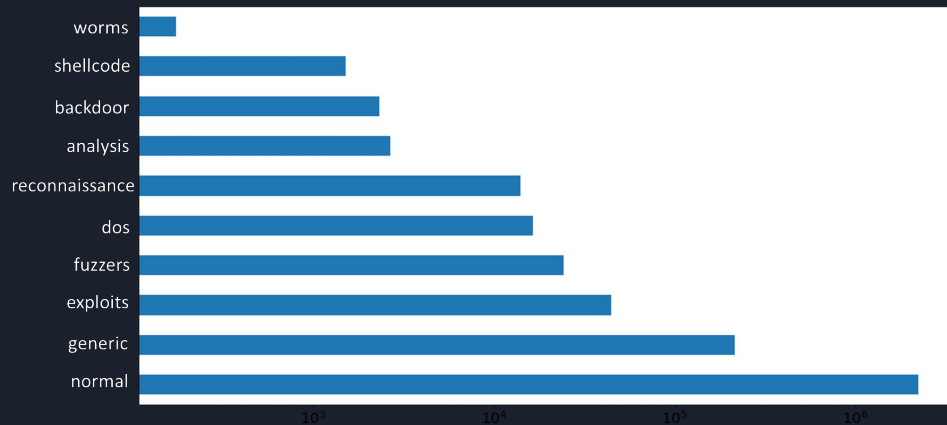
- Binary classification

- Conclusion

# UNSW-NB15: An unbalanced dataset

The UNSW-NB 15 dataset contains an hybrid of realistic normal traffic activities and synthetic attack behaviors simulated in a laboratory environment with three interconnected virtual machines.

It contains 9 different attack categories:
*Analysis, Backdoor, Dos, Exploits, Fuzzers, Generic, Reconnaissance, Shellcode, Worms*



**Distribution of network traffic:**

| | |
|---|---|
| **normal** | **2218761** |
| **generic** | **215481** |
| **exploits** | **44525** |
| **fuzzers** | **24246** |
| **dos** | **16353** |
| **reconnaissance** | **13987** |
| **analysis** | **2677** |
| **backdoor** | **2329** |
| **sheelcode** | **1511** |
| **worms** | **174** |

# UNSW-NB15: A packet-based analysis

It contains a total of 49 features divided in **flow-based** and **packet-based** features.

**Flow-based:** Source IP, Source Port, Destination IP, Destination Port of the three virtual machines in the laboratory environment

**Packet-based**: Divided in four sub-features
1. *Basic features* (i.e bytes, type of the service, bps…)
2. *Content features* (i.e body length in case of HTTP, TCP window of Source and Destination…)
3. *Time features* (i.e jitter, time between SYN/SYN-ACK, time between SYN-ACK/ACK…)
4. *Additional features* (i.e if the FTP session is accessed with a password…)

**Considered features:** 49-4-2 (start timestamp, final timestamp) = **43**

Two errors in *attack_cat* feature: "normal" category is missing and values are incorrect

| | |
|---|---|
| Generic | 215481 |
| Exploits | 44525 |
| Fuzzers | 19195 |
| DoS | 16353 |
| Reconnaissance | 12228 |
| Fuzzers | 5051 |
| Analysis | 2677 |
| Backdoor | 1795 |
| Reconnaissance | 1759 |
| Shellcode | 1288 |
| Backdoors | 534 |
| Shellcode | 223 |
| Worms | 174 |

Adjust values, fix missing 'normal' traffic

| | |
|---|---|
| normal | 2218761 |
| generic | 215481 |
| exploits | 44525 |
| fuzzers | 24246 |
| dos | 16353 |
| reconnaissance | 13987 |
| analysis | 2677 |
| backdoor | 2329 |
| shellcode | 1511 |
| worms | 174 |

The attributes *ct_flw_http_mthd* and *is_ftp_login* and *ct_ftp_cmd* contains null values:

1. *No. of flows that has methods such as Get and Post in HTTP service (ct_flw_http_mthd, numerical variable)*
2. *If the FTP session is accessed by user and password then 1 else 0 (is_ftp_login, binary variable)*
3. *No. of flows that has a command in FTP session (ct_ftp_cmd , numerical variable)*

*Moreover,*

| Name | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| is_sm_ips_ports | 2540044.0 | 0.001651 | 0.040596 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| is_ftp_login | 2540044.0 | 0.017351 | 0.133457 | 0.0 | 0.0 | 0.0 | 0.0 | 4.0 |

```
df.is_ftp_login.value_counts()
✓ 0.4s
0.0    1066591
1.0      43389
4.0        156
2.0         30
```

**Solution adopted for null values:**
→ **0** where the service is not HTTP (for *ct_flw_http_mthd* ) or FTP (for *is_ftp_login*)
→ **The attribute mean** for all sample belonging to the same class where service is HTTP/FTP

**Solution adopted for binary out of range:**
→ Substituting with **the most probable value** (that is 0)

Three types of attributes:

1. **Numerical attributes**

   - Z-score normalization has been applied
   - sklearn.preprocessing.StandardScaler

2. **Categorical attributes** (*proto, state, service*):

   - converted into binary: from **43 up to 202 features** after conversion
   - sklearn.preprocessing.OneHotEncoder

3. **Binary attributes**

# UNSW-NB15: Classification

Starting from the entire dataset, a **Stratified K-Fold cross validation** is applied

- it provides train/test indices to split data in train/test sets

- K=5

- stratified folds: the folds are made by preserving the percentage of samples for each class.

Different types of experiment: classification with a

- unbalanced training set

- a 60/40 rebalanced training set with undersampling of the majority class

- a 50/50 rebalanced training set with a combination of undersampling and SMOTE of the minority class

# UNSW-NB15: Classification

## Classification algorithm

✔ **Naive Bayes**

✔ **Decision Tree**

- experimented with GINI index and min_samples_leaf=[3,5,7]

✔ **Random Forest**

- experimented with GINI index and n_estimators=[10,20]

## Additional

✔ **Random Forest with PCA (n_components=0.99)**

- experimented with GINI index and n_estimators=20

# UNSW-NB15: Multi-class classification

From left to right results for:
**Naive Bayes**, **Decision Tree**, **Random Forest**

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| analysis | 0.0614 | 0.3313 | 0.1747 | 535.4 |
| backdoor | 0.0465 | 0.0965 | 0.0743 | 465.8 |
| dos | 0.2387 | 0.0098 | 0.0121 | 3270.6 |
| exploits | 0.7941 | 0.0395 | 0.0488 | 8905.0 |
| fuzzers | 0.0925 | 0.0663 | 0.0702 | 4849.2 |
| generic | 0.5783 | 0.0922 | 0.1085 | 43096.2 |
| normal | 0.8804 | 0.6844 | 0.7161 | 443752.2 |
| reconnaissance | 0.1637 | 0.0088 | 0.0109 | 2797.4 |
| shellcode | 0.0407 | 0.9960 | 0.1734 | 302.2 |
| worms | 0.0002 | 0.8336 | 0.0010 | 34.8 |
| avg_acc | 0.6081 | | | |

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| analysis | 0.0435 | 0.3549 | 0.1405 | 535.4 |
| backdoor | 0.1333 | 0.1498 | 0.1315 | 465.8 |
| dos | 0.3028 | 0.1091 | 0.1245 | 3270.6 |
| exploits | 0.7839 | 0.5401 | 0.5751 | 8905.0 |
| fuzzers | 0.5426 | 0.4680 | 0.4698 | 4849.2 |
| generic | 0.9411 | 0.9859 | 0.9763 | 43096.2 |
| normal | 0.9951 | 0.9940 | 0.9942 | 443752.2 |
| reconnaissance | 0.8989 | 0.7556 | 0.7797 | 2797.4 |
| shellcode | 0.5597 | 0.6081 | 0.5856 | 302.2 |
| worms | 0.5494 | 0.4478 | 0.4579 | 34.8 |
| avg_acc | 0.9716 | | | |

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| analysis | 0.0299 | 0.1277 | 0.0760 | 535.4 |
| backdoor | 0.0675 | 0.3541 | 0.1713 | 465.8 |
| dos | 0.3423 | 0.0875 | 0.1023 | 3270.6 |
| exploits | 0.7968 | 0.5708 | 0.6045 | 8905.0 |
| fuzzers | 0.5654 | 0.4819 | 0.4840 | 4849.2 |
| generic | 0.9484 | 0.9848 | 0.9770 | 43096.2 |
| normal | 0.9946 | 0.9948 | 0.9947 | 443752.2 |
| reconnaissance | 0.8629 | 0.7511 | 0.7687 | 2797.4 |
| shellcode | 0.6150 | 0.5916 | 0.5924 | 302.2 |
| worms | 0.7201 | 0.3334 | 0.3710 | 34.8 |
| avg_acc | 0.9726 | | | |

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| analysis | 0.0615 | 0.2797 | 0.1597 | 535.4 |
| backdoor | 0.0463 | 0.1601 | 0.1051 | 465.8 |
| dos | 0.3203 | 0.0111 | 0.0137 | 3270.6 |
| exploits | 0.8017 | 0.0393 | 0.0486 | 8905.0 |
| fuzzers | 0.0915 | 0.0650 | 0.0689 | 4849.2 |
| generic | 0.5728 | 0.0874 | 0.1034 | 43096.2 |
| normal | 0.8991 | 0.8429 | 0.8533 | 443752.2 |
| reconnaissance | 0.5633 | 0.0057 | 0.0071 | 2797.4 |
| shellcode | 0.0409 | 0.9933 | 0.1739 | 302.2 |
| worms | 0.0004 | 0.8218 | 0.0023 | 34.8 |
| avg_acc | 0.7462 | | | |

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| analysis | 0.0575 | 0.3122 | 0.1628 | 535.4 |
| backdoor | 0.0608 | 0.3361 | 0.1683 | 465.8 |
| dos | 0.2340 | 0.1462 | 0.1575 | 3270.6 |
| exploits | 0.7625 | 0.5005 | 0.5367 | 8905.0 |
| fuzzers | 0.4999 | 0.7899 | 0.6804 | 4849.2 |
| generic | 0.9547 | 0.9842 | 0.9780 | 43096.2 |
| normal | 0.9997 | 0.9855 | 0.9883 | 443752.2 |
| reconnaissance | 0.7831 | 0.7717 | 0.7717 | 2797.4 |
| shellcode | 0.3494 | 0.7478 | 0.5893 | 302.2 |
| worms | 0.2323 | 0.6260 | 0.4660 | 34.8 |
| avg_acc | 0.9670 | | | |

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| analysis | 0.0657 | 0.2887 | 0.1655 | 535.4 |
| backdoor | 0.0774 | 0.3520 | 0.1834 | 465.8 |
| dos | 0.3233 | 0.1267 | 0.1437 | 3270.6 |
| exploits | 0.8062 | 0.5347 | 0.5728 | 8905.0 |
| fuzzers | 0.4796 | 0.8402 | 0.7056 | 4849.2 |
| generic | 0.9464 | 0.9840 | 0.9759 | 43096.2 |
| normal | 0.9999 | 0.9860 | 0.9888 | 443752.2 |
| reconnaissance | 0.7938 | 0.7847 | 0.7844 | 2797.4 |
| shellcode | 0.3736 | 0.8312 | 0.6458 | 302.2 |
| worms | 0.2925 | 0.4650 | 0.4113 | 34.8 |
| avg_acc | 0.9685 | | | |

# UNSW-NB15: Binary classification

From left to right results for:
**Naive Bayes**, **Decision Tree**, **Random Forest**

| type | precision | recall | f-score | support |
|------|-----------|--------|---------|---------|
| attack | 0.6943 | 0.1799 | 0.2101 | 64256.6 |
| normal | 0.8922 | 0.9838 | 0.9640 | 443752.2 |
| avg_acc | 0.8821 | | | |

| type | precision | recall | f-score | support |
|------|-----------|--------|---------|---------|
| attack | 0.9574 | 0.9656 | 0.9638 | 64256.6 |
| normal | 0.9950 | 0.9936 | 0.9939 | 443752.2 |
| avg_acc | 0.9901 | | | |

| type | precision | recall | f-score | support |
|------|-----------|--------|---------|---------|
| attack | 0.9596 | 0.9685 | 0.9665 | 64256.6 |
| normal | 0.9954 | 0.9939 | 0.9942 | 443752.2 |
| avg_acc | 0.9907 | | | |

| type | precision | recall | f-score | support |
|------|-----------|--------|---------|---------|
| attack | 0.7025 | 0.1754 | 0.2059 | 64256.6 |
| normal | 0.8919 | 0.9862 | 0.9657 | 443752.2 |
| avg_acc | 0.8836 | | | |

| type | precision | recall | f-score | support |
|------|-----------|--------|---------|---------|
| attack | 0.9138 | 0.9971 | 0.9787 | 64256.6 |
| normal | 0.9995 | 0.9858 | 0.9885 | 443752.2 |
| avg_acc | 0.9873 | | | |

| type | precision | recall | f-score | support |
|------|-----------|--------|---------|---------|
| attack | 0.9142 | 0.9996 | 0.9808 | 64256.6 |
| normal | 0.9999 | 0.9859 | 0.9886 | 443752.2 |
| avg_acc | 0.9876 | | | |

# UNSW-NB15: PCA and Random Forest in multi-class classification

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| analysis | 0.0659 | 0.2689 | 0.1604 | 535.4 |
| backdoor | 0.0487 | 0.2365 | 0.1236 | 465.8 |
| dos | 0.2073 | 0.1037 | 0.1146 | 3270.6 |
| exploits | 0.7641 | 0.4787 | 0.5168 | 8905.0 |
| fuzzers | 0.4363 | 0.8048 | 0.6690 | 4849.2 |
| generic | 0.9430 | 0.9762 | 0.9689 | 43096.2 |
| normal | 0.9999 | 0.9855 | 0.9883 | 443752.2 |
| reconnaissance | 0.6269 | 0.7230 | 0.6986 | 2797.4 |
| shellcode | 0.3193 | 0.6201 | 0.5153 | 302.2 |
| worms | 0.1167 | 0.1954 | 0.1708 | 34.8 |
| avg_acc | 0.9652 | | | |

| type | precision | recall | f-score | support |
|---|---|---|---|---|
| attack | 0.9101 | 0.9996 | 0.9798 | 64256.6 |
| normal | 0.9999 | 0.9852 | 0.9881 | 443752.2 |
| avg_acc | 0.9870 | | | |

# UNSW-NB15: Conclusion

Possible future analysis:
- Approach an analysis on attributes by studying correlation among them (for ex. Pearson coefficient)
- Try an NLS-KDD approach by leave out from training set some attack categories and evaluate performance

**THANK YOU FOR YOUR ATTENTION**