

# Machine Learning Nanodegree

## Udacity Connect Intensive

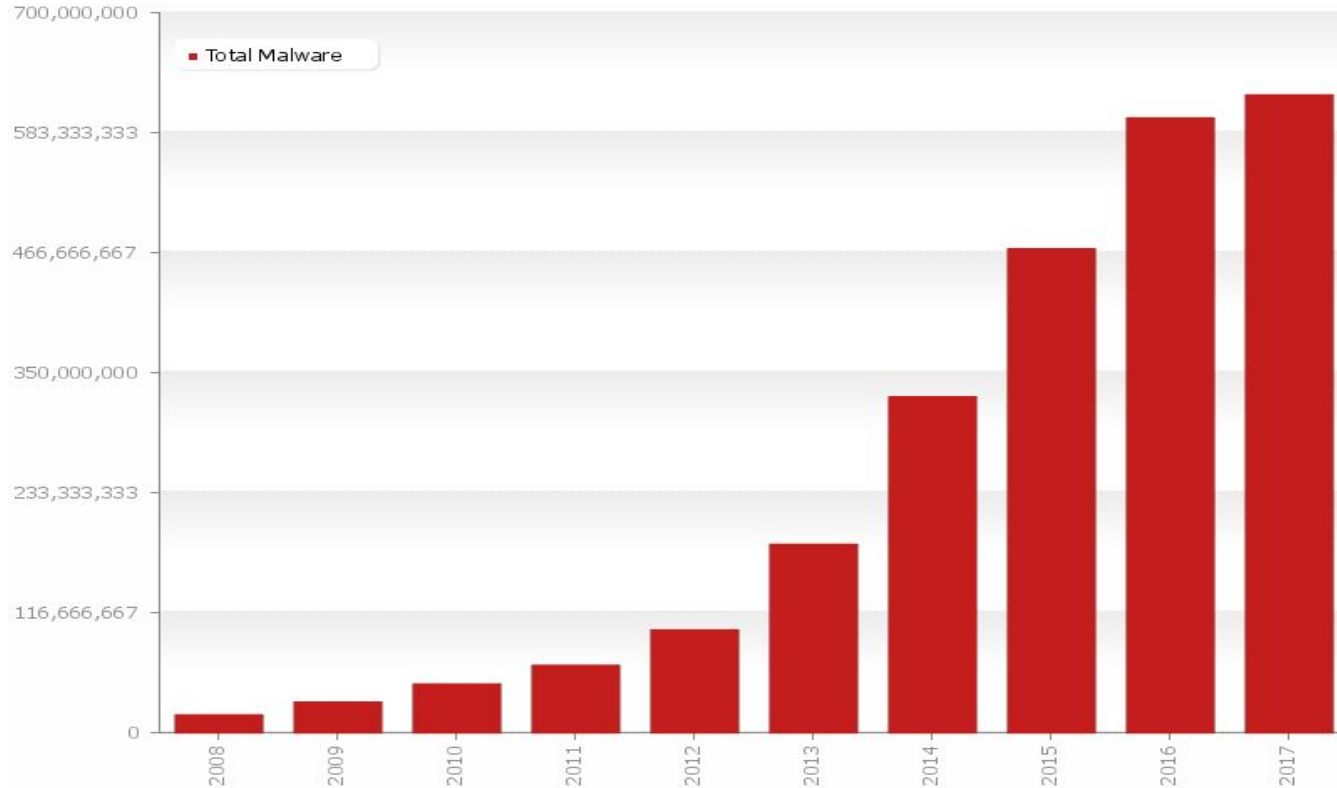
### Capstone Proposal

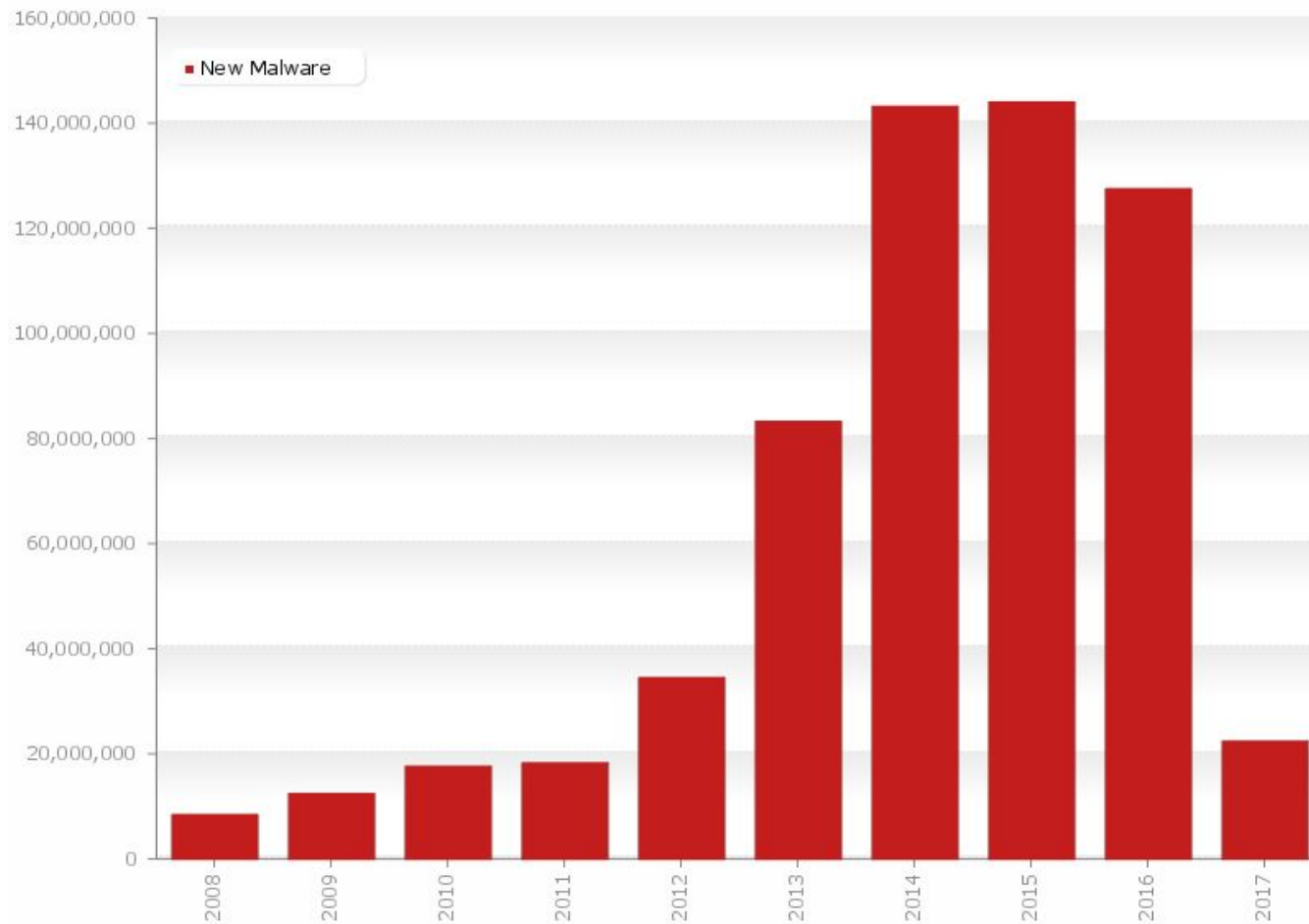
Using Machine Learning for Malware detection

Dileep Tiku



# Malware-the most common threat to IT security





# Problem Statement

**Develop a **behavior analysis** model that can identify whether a software is **malware or not**.**

- Allow organizations to use the malware detection model, to assist in advanced threat protection not possible only by traditional methods of malware detection, which are signature based.
- Malware poses the biggest threat to security in organization.
- With the advent of Industrial internet of things besides compromising information technology assets, malware can impact operational technology which has potential of causing harm to society at large.

# Datasets and Inputs

The dataset used for this project is a collection 766 malware traces

The malware traces were obtained from the [CSMining website](#).

## ***Training Data***

<i>Diagnosis</i>	<i>No. of Malware Traces</i>	<i>Percent</i>
Infected*	320	82.47%
Benign	68	17.53%

## **Test Data**

<i>No Of Malware Traces</i>	378
-----------------------------	-----

\*Worms, Trojans, Viruses are labeled as infected

# Utilize behavior analysis to diagnose malware

Malware is primarily **diagnosed by using signatures of the known types**. A classification model will be developed based on the behavior data of the software in a sandbox environment, where monitoring tools will record the behavior data. The software will be treated as a black box for generating the behavior data.

The training data set from CSMining web site contains a list of API Calls, which has been generated by executing malware and benign software in a controlled environment. The sequence of API Calls are labeled as malware or benign software.

The first step of the solution is to extract features from the list of API calls. The features will be extracted by generating n-grams for each API Call sequence. The n-grams with maximum number of frequencies will then be selected for the SVM classification model.

## Benchmark: model

The benchmark model will be the detection accuracy achieved by using the SVM classifier with kernel = 'linear'.

## Evaluation Metrics

tp: true positive  
fp: false positive

fn: false negative  
tn: true negative

Actual Class	Predicted Class	
	Malware	Benign
Malware	tp	fn
Benign	fp	tn

$$\text{Accuracy} = \frac{tp + tn}{tp + tn + fp + fn} \quad \text{Precision} = \frac{tp}{tp + fp} \quad \text{Recall} = \frac{tp}{tp + fn}$$

$$F_1 = 2 \cdot \frac{1}{\frac{1}{\text{recall}} + \frac{1}{\text{precision}}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

**Type II errors** are associated with false negatives. In this project it is when a software is **diagnosed** as **benign** but **actually** is malware. These errors are measured by **recall**; the more type II errors the lower the recall will be and vice versa.

### Malware removal

Infected files should be quarantined to keep the IT assets safe.

Type I errors are not generally an issue

If malware is **not detected**, it can potentially cause huge damage.

$$\text{Recall} = \frac{tp}{tp + fn}$$

# Algorithms and Techniques

1. Extract the features from CSMining dataset.
  - The feature set will be extracted from the list of API calls by using a text mining technique utilizing bag of words.
  - The technique involves computing the tf-idf (term frequency times the inverse document frequency)
2. Train a linear kernel based SVM classifier with the extracted features from the list of API calls.
3. Test the classifier with the test data set and compute the performance metrics.
4. Using the extracted feature set, train the SVM classifier with Stochastic Gradient Descent with loss = 'hinge' to notice any difference in test accuracy.
5. Test the new model with test set and compare the performance metrics with the results from step 3.