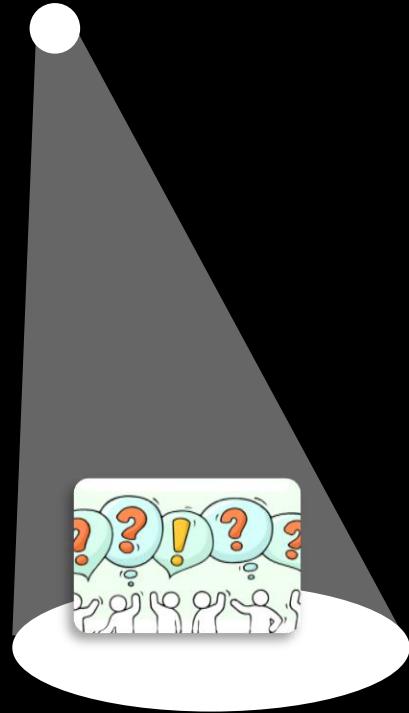


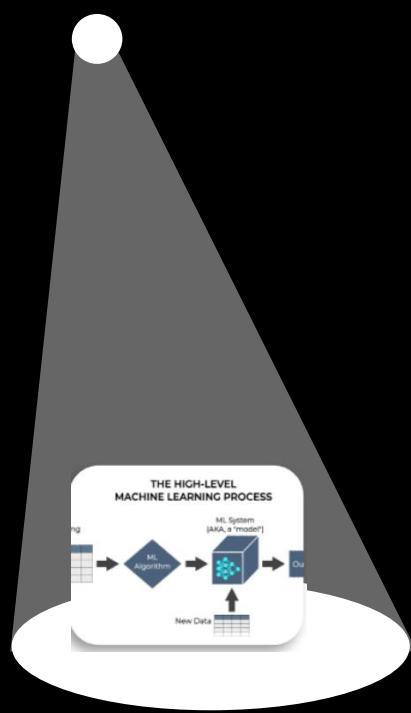
# Machine Learning For Security Professionals

## Building And Hacking ML Systems

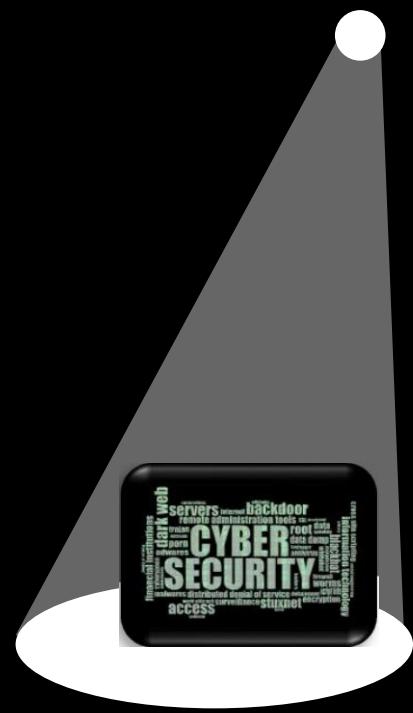




# Intro



# ML for Security



# Security for ML

# \$whoami

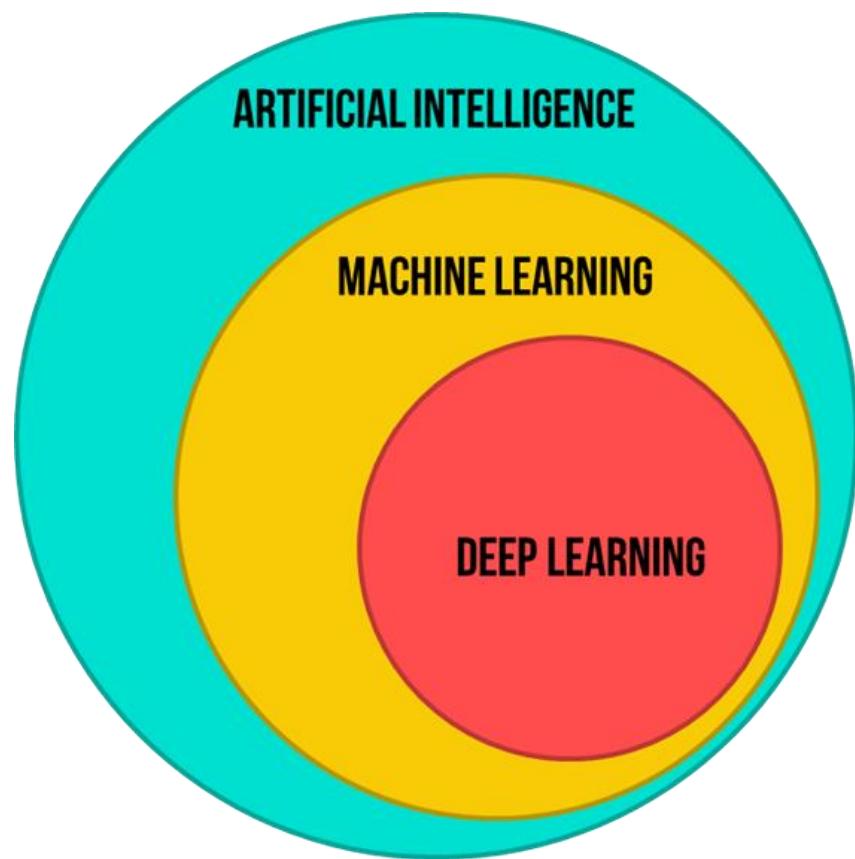
**SAGAR BHURE**

Software Engineer | OWASP Project Lead | Blackhat Speaker | Mentor



# Intro

DL ≠ ML ≠ AI

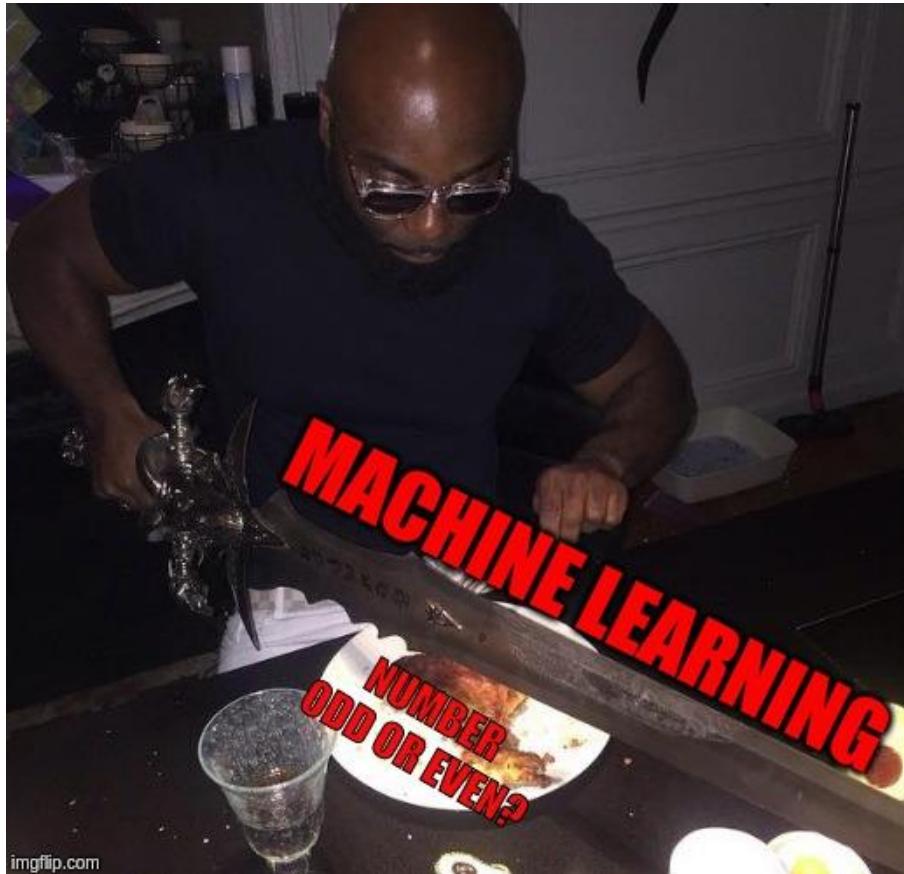


# ML/DL Intro: What for?

- Pattern matching
- Pattern search
- Classification
- Regression
- Clustering
- Forecasting
- etc.



# Where not to use ML?



# Applications in security

- WAF
- IDS/IPS
- Malware detection
- Antiviruses
- Spam filters
- ...

# Libs



TensorFlow



# Tools

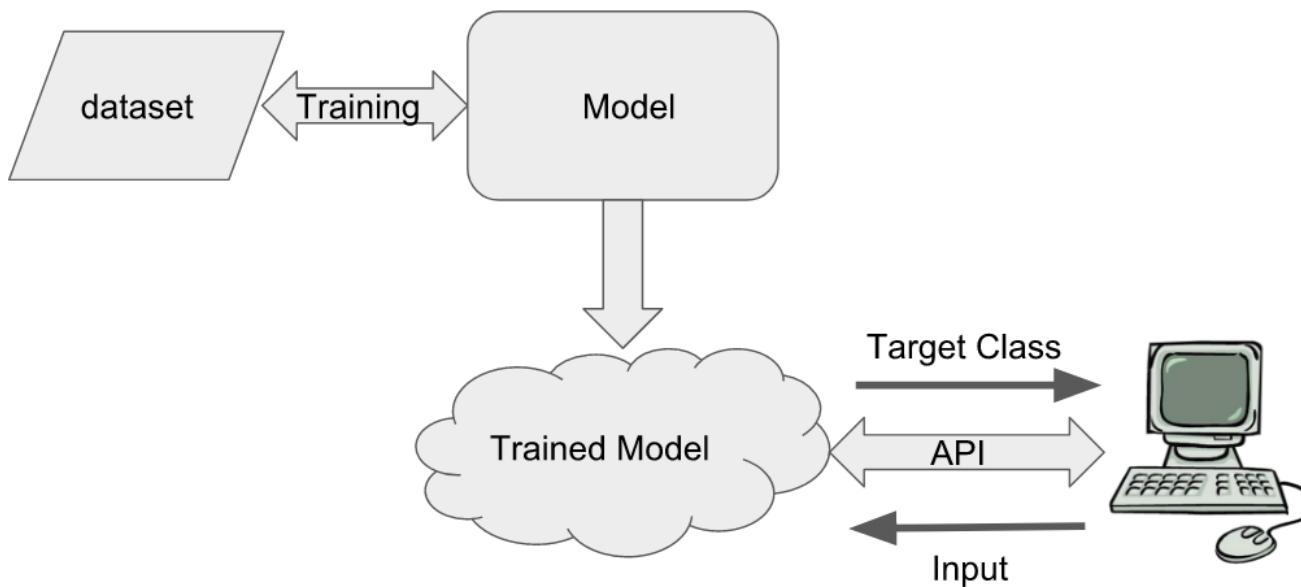
CONDA



IP[y]:  
IPython



# MLops

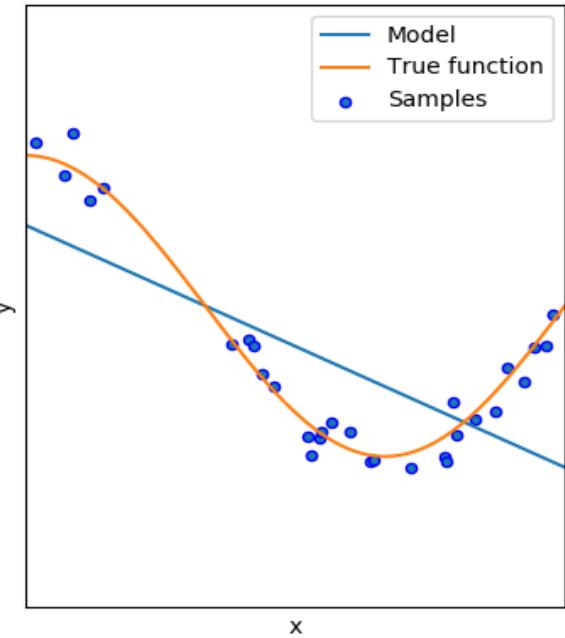


# Basic Terminology

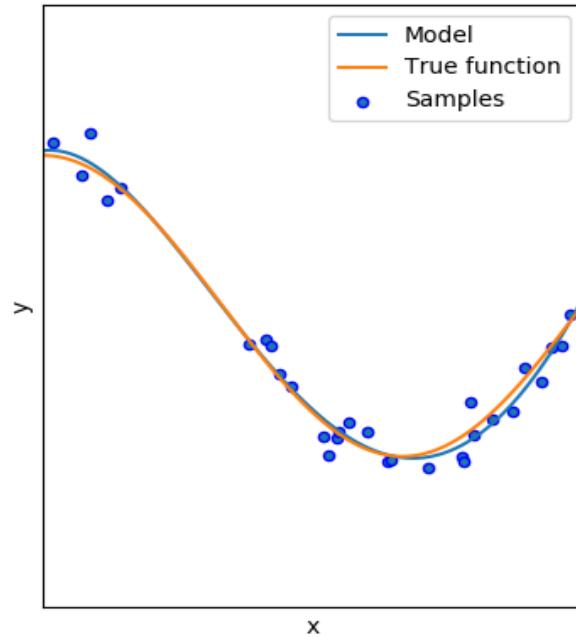
- Dataset
- Features and Labels
- Training set and Testing set
- Probability
- Model
- Loss function
- Optimization
- Overfitting
- Underfitting

# Underfitting vs Overfitting

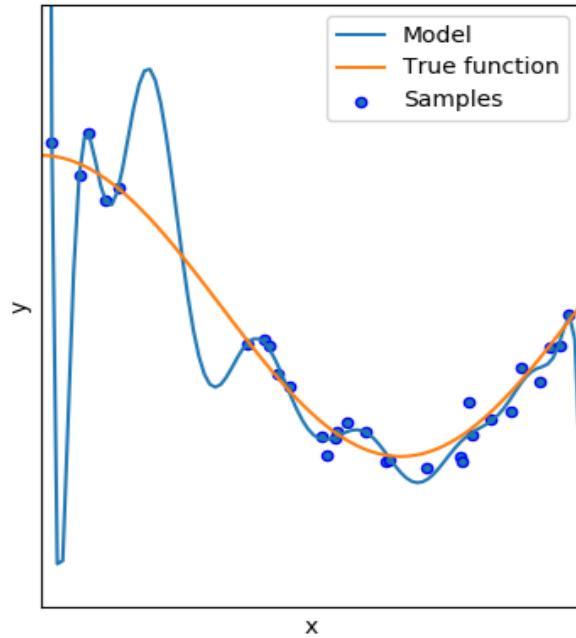
Degree 1  
MSE = 4.08e-01(+/- 4.25e-01)



Degree 4  
MSE = 4.32e-02(+/- 7.08e-02)



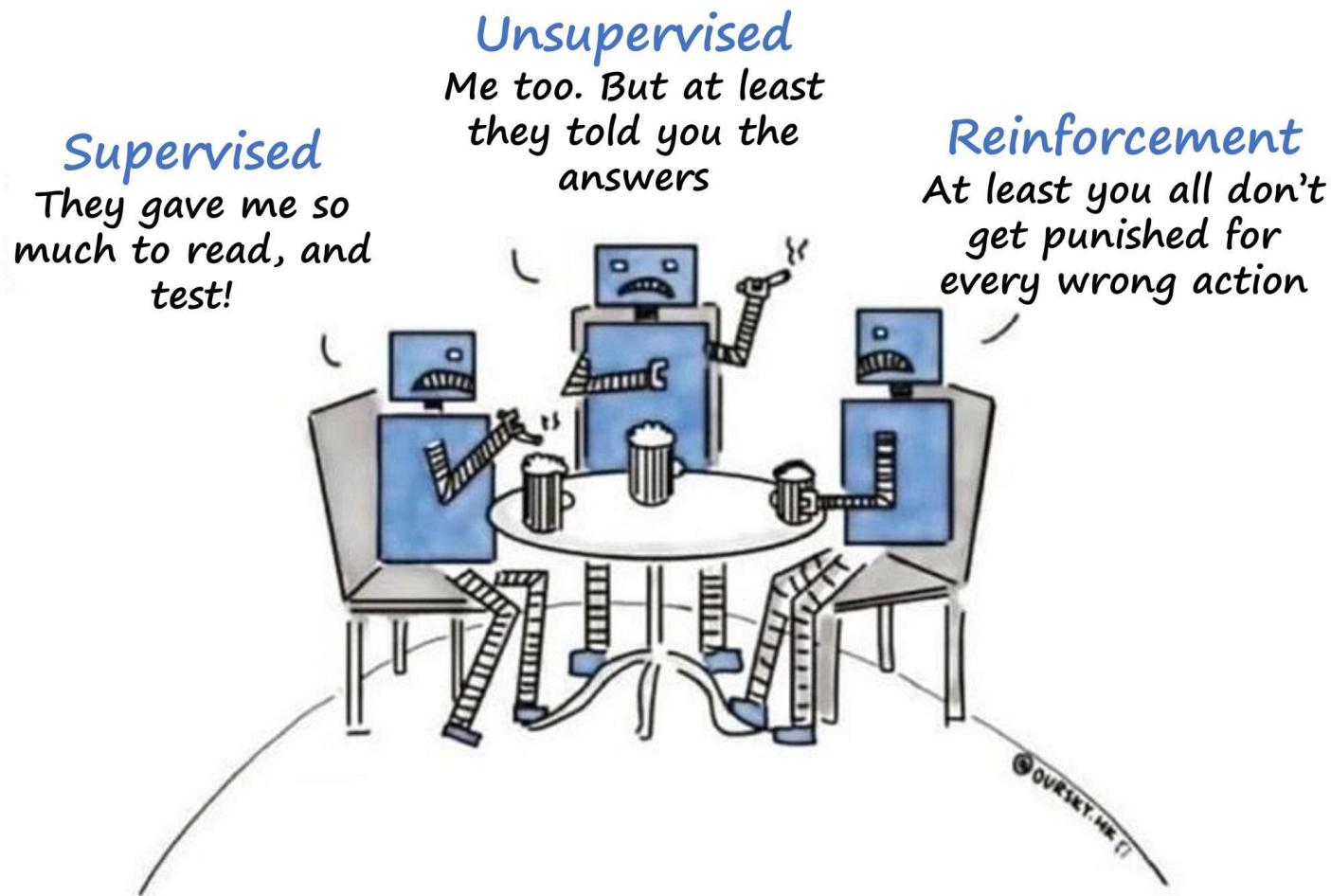
Degree 15  
MSE = 1.83e+08(+/- 5.48e+08)



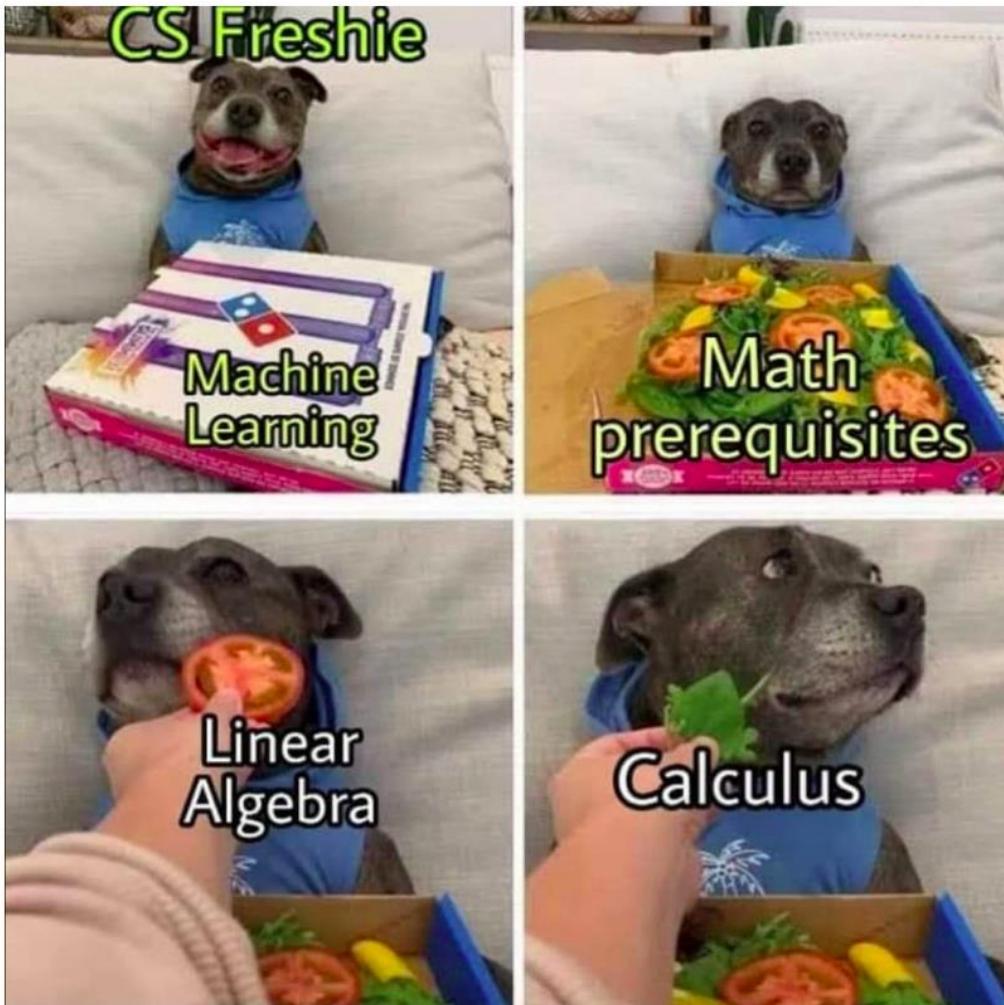
# Basic Terminology

- Supervised Learning
- Unsupervised Learning
- Semi-supervised Learning

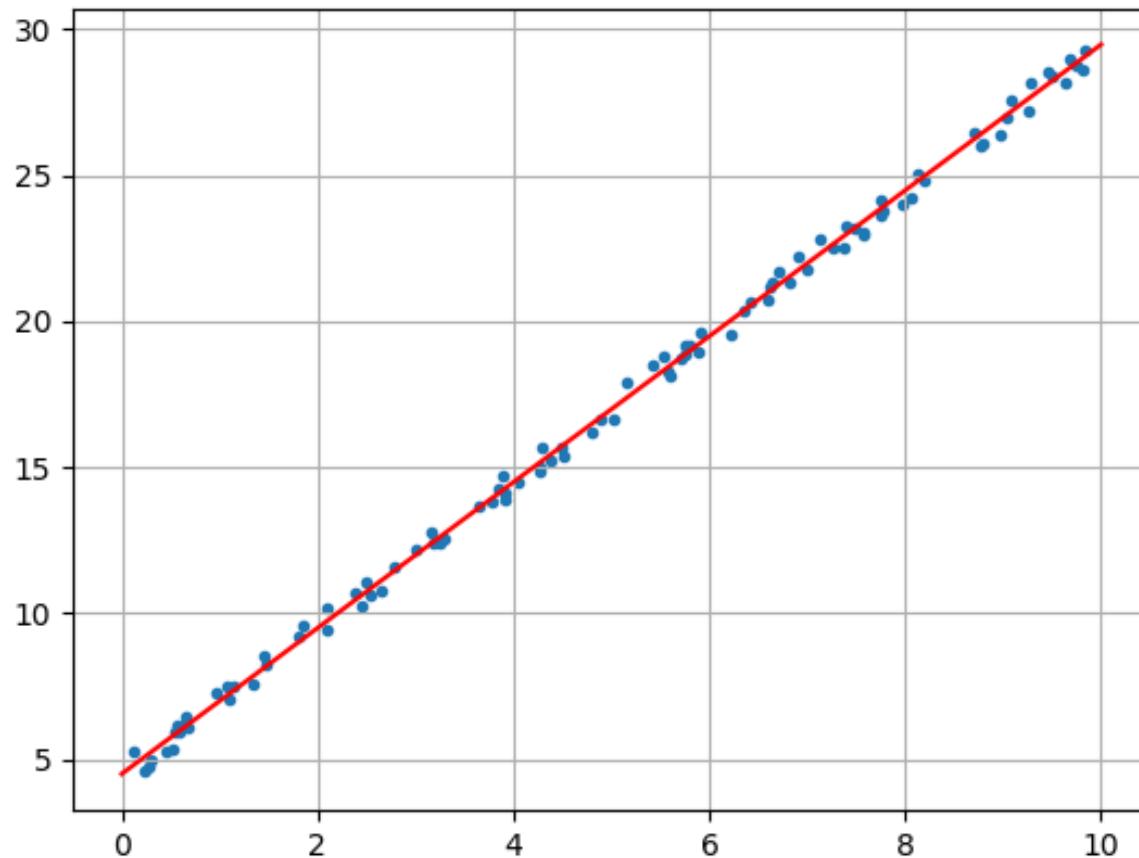
# Basic Terminology



# I know, right



# Linear Regression



# Linear Regression

## Demo 1: Linear Regression

<Linear Regression End>

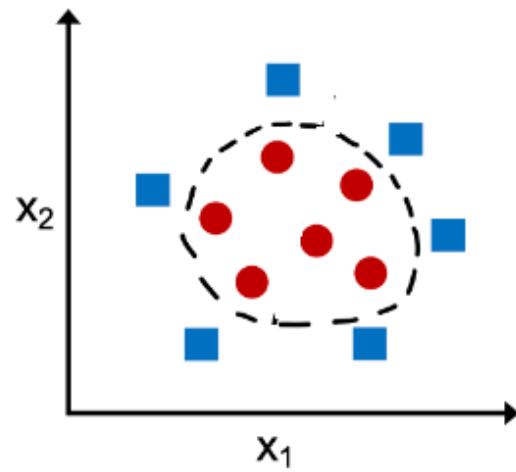
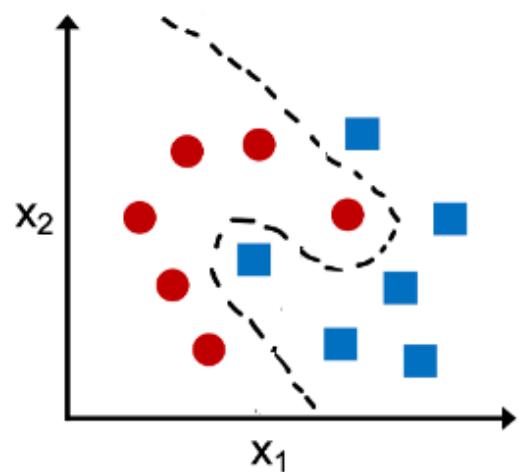
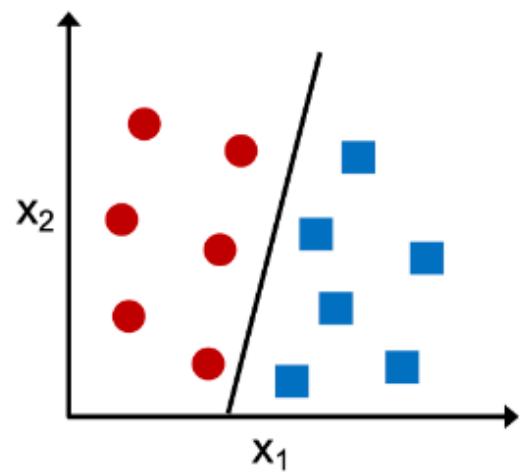
# Before we begin

## Teachable Machine

<https://teachablemachine.withgoogle.com/>

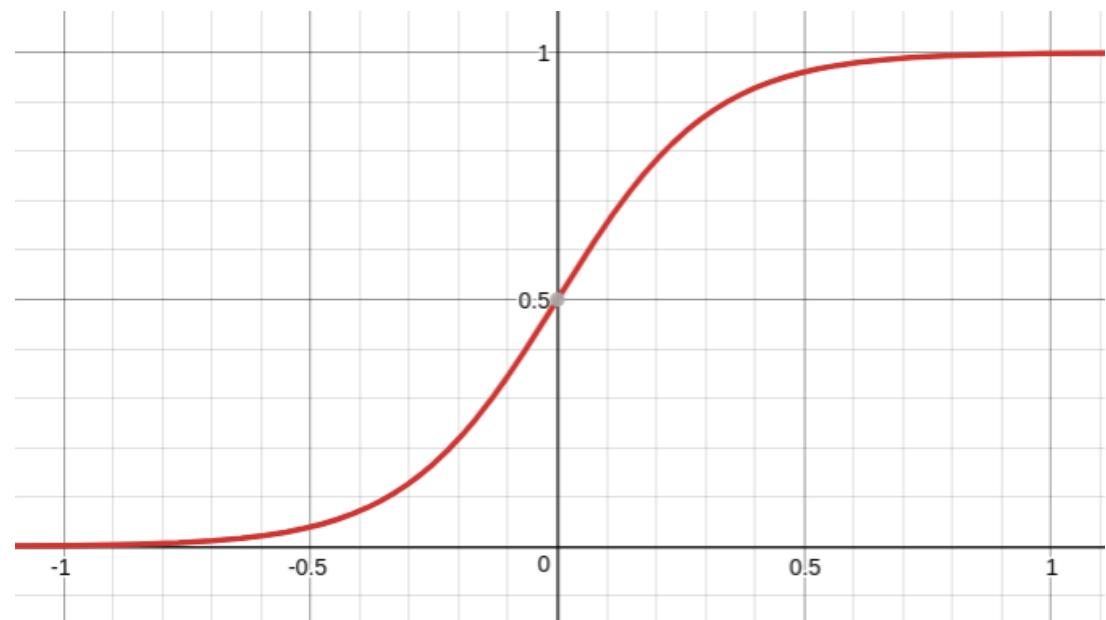
{ref} <https://www.tensorflow.org/js/demos>

# Classification

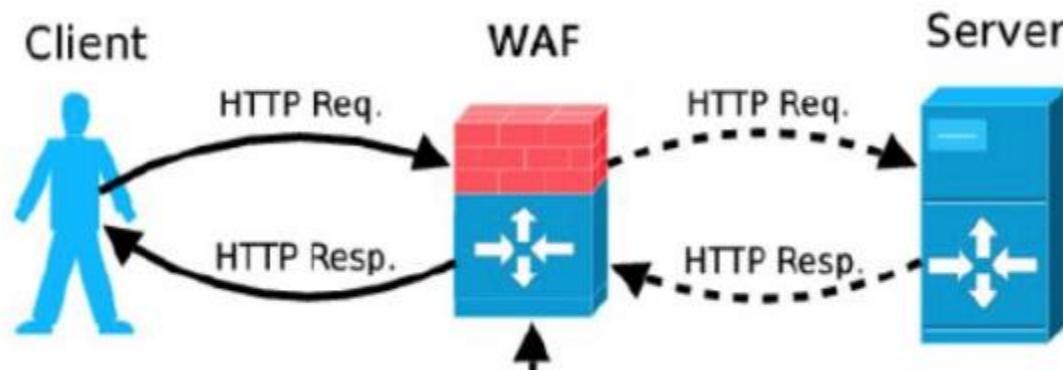


# Logistic Regression

$$\frac{1}{1 + e^{-(wx + b)}}$$



# WAF using Logistic Regression



Ref: <https://github.com/faizann24/Fwaf-Machine-Learning-driven-Web-Application-Firewall>

# DataSet

## Good queries

```
/final_warning/  
/javascript/resize.cs  
/javascript/msg.rtf  
/javascript/cat.rb  
/buy_phente/  
//yske4gpe.asp  
/javascript/menu.meta  
/avalanche-research_0/  
/privacy-gnw/  
/176232/  
/examples/jsp/checkbox/webc.cgi/  
/cb_line/  
/bmw_3er/  
/155362/  
/javascript/shopping.nsf  
.....
```

## Bad queries

```
/<script>document.cookie="testnvxc=4301;"</script>  
/cacti/base_local_rules.php?dir=<script>alert('base_l  
ocal_rules_xss.nasl-1331905098')</script>  
/scripts/comersus_customerregistrationform.asp  
/cgi-bin/search.pl?form=../../../../etc/passwd\  
x00  
/help.php?q='uname >q23433168 #  
/u8yiqilu.jsp?<meta http-equiv=set-cookie content="te  
stswjp=6010">  
/scripts/forum.php3?id_article=1&id_forum=-1/**/union  
/**/select 1376526276--  
/core/adodb/server.php?sql='adodb_sql_sql_injection.n  
asl  
/xsqln7eb.aspx?<meta http-equiv=set-cookie content="t  
estpcbb=9254">
```

# Part one

## Preprocessing data

# TF-IDF

$$TF_{ij} = \frac{f_{ij}}{n_j} \quad (1)$$

Where  $f_{ij}$  is the frequency of term i in document j.  $n_j$  is the total number of words in document j.

$$IDF_i = 1 + \log\left(\frac{N}{c_i}\right) \quad (2)$$

Where N is the total number of documents in the corpus.  $c_i$  is the number of documents that contain word i.

$$w_{ij} = TF_{ij} \times IDF_i \quad (3)$$

Where  $w_{ij}$  is the TF-IDF score of term i in document j.

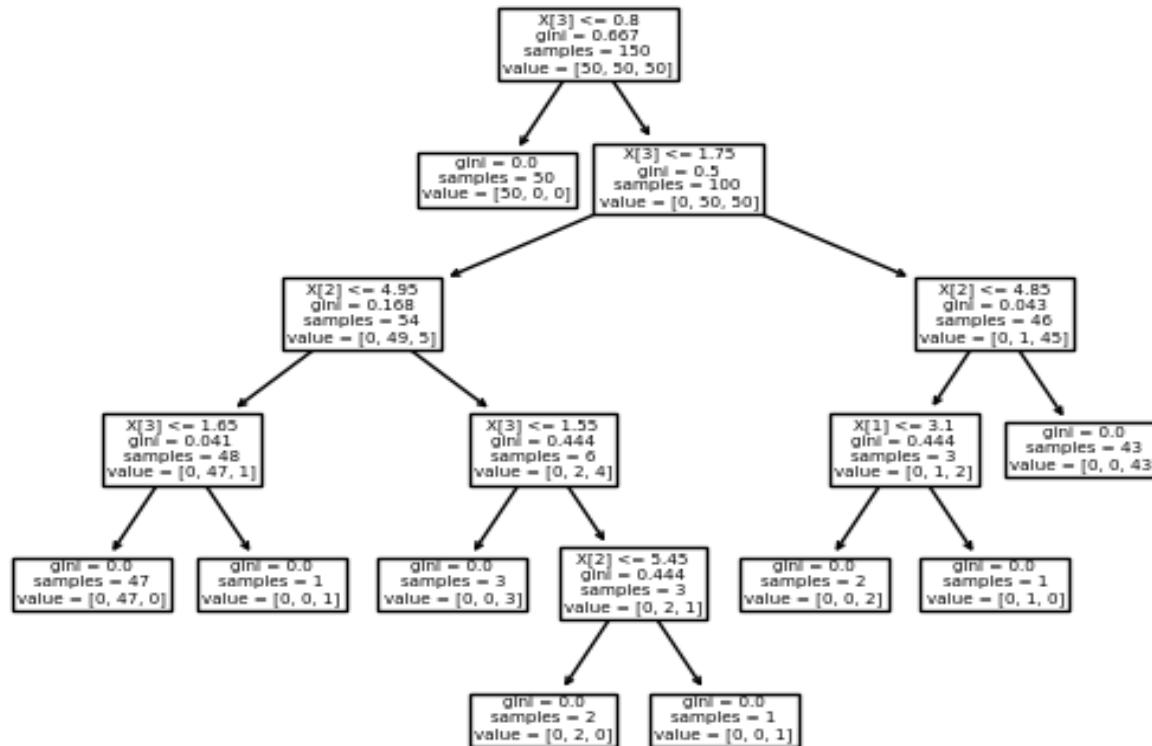
# TF-IDF

- Term Frequency: This summarizes how often a given word appears within a document.
- Inverse Document Frequency: This downscals words that appear a lot across documents.

## **Demo 2: WAF using Logistic regression**

<Logistic Regression End>

# Decision Tree & Random Forest



# Decision Tree & Random Forest

- Gini Impurity: Gives you a likelihood of incorrect classification

$$Gini\ impurity = 1 - \sum_{j=1}^f p_j^2$$

- Calculate Gini impurity for all possible features
- Add the node to tree if improved Gini impurity value is seen
- Else do not fork the tree

# Decision Tree & Random Forest

- Precision: What proportion of positive identifications was actually correct?

$$\text{Precision} = \frac{TP}{TP + FP}$$

- Recall: What proportion of actual positives was identified correctly?

$$\text{Recall} = \frac{TP}{TP + FN}$$

- F1 score:

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

# Dataset : Malware Classifier

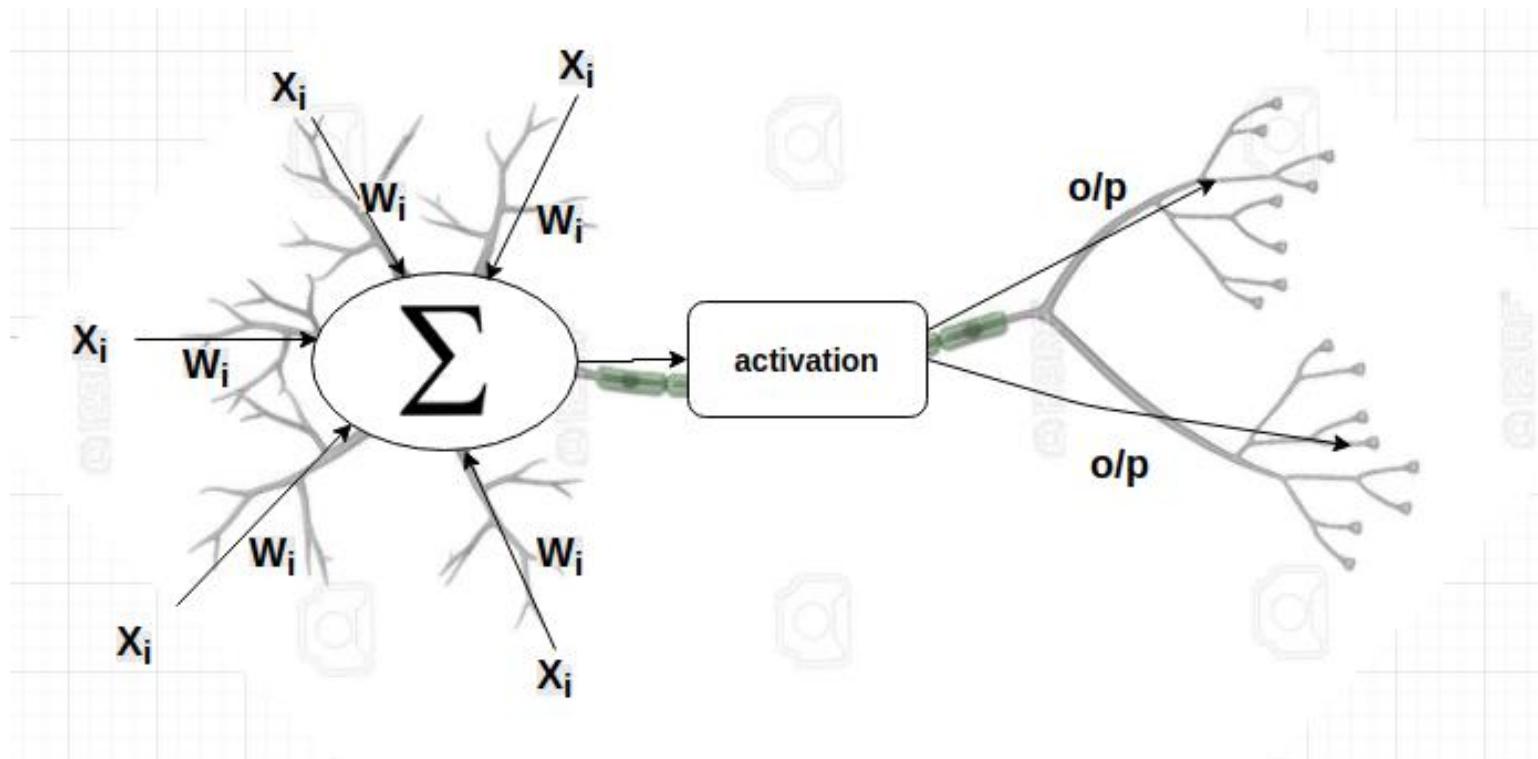
- 56 features extracted from PE files using pefile
- 2 classes
- 138047 samples

{ref}: <https://github.com/Te-k/malware-classification>

# **Demo 3: Malware Classifier**

<Decision Tree End>

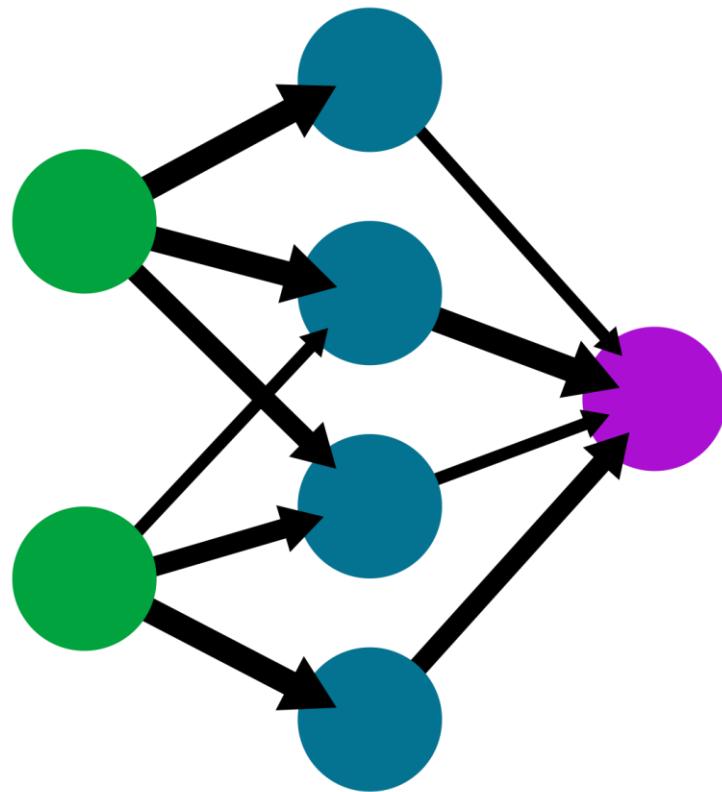
# Neural Network



# A simple NN

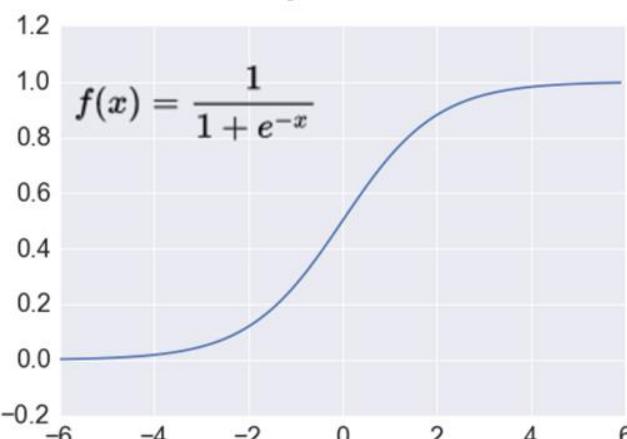
A simple neural network

input      hidden      output  
layer      layer      layer

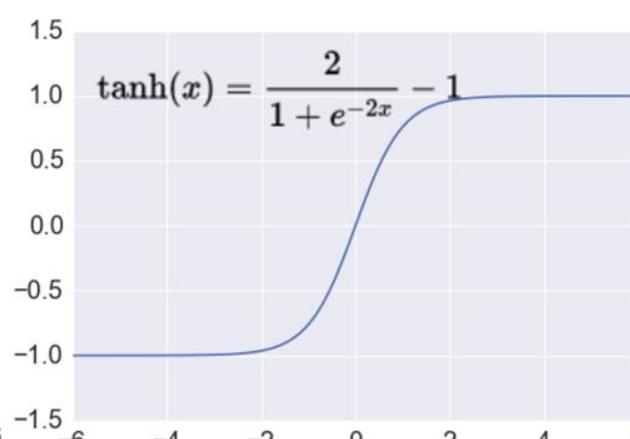


# Activation Function

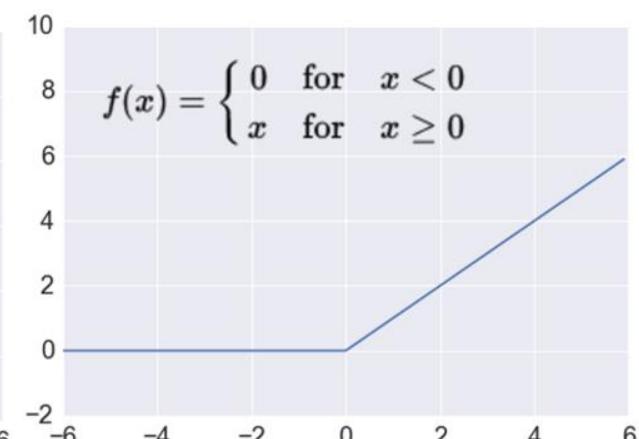
Sigmoid



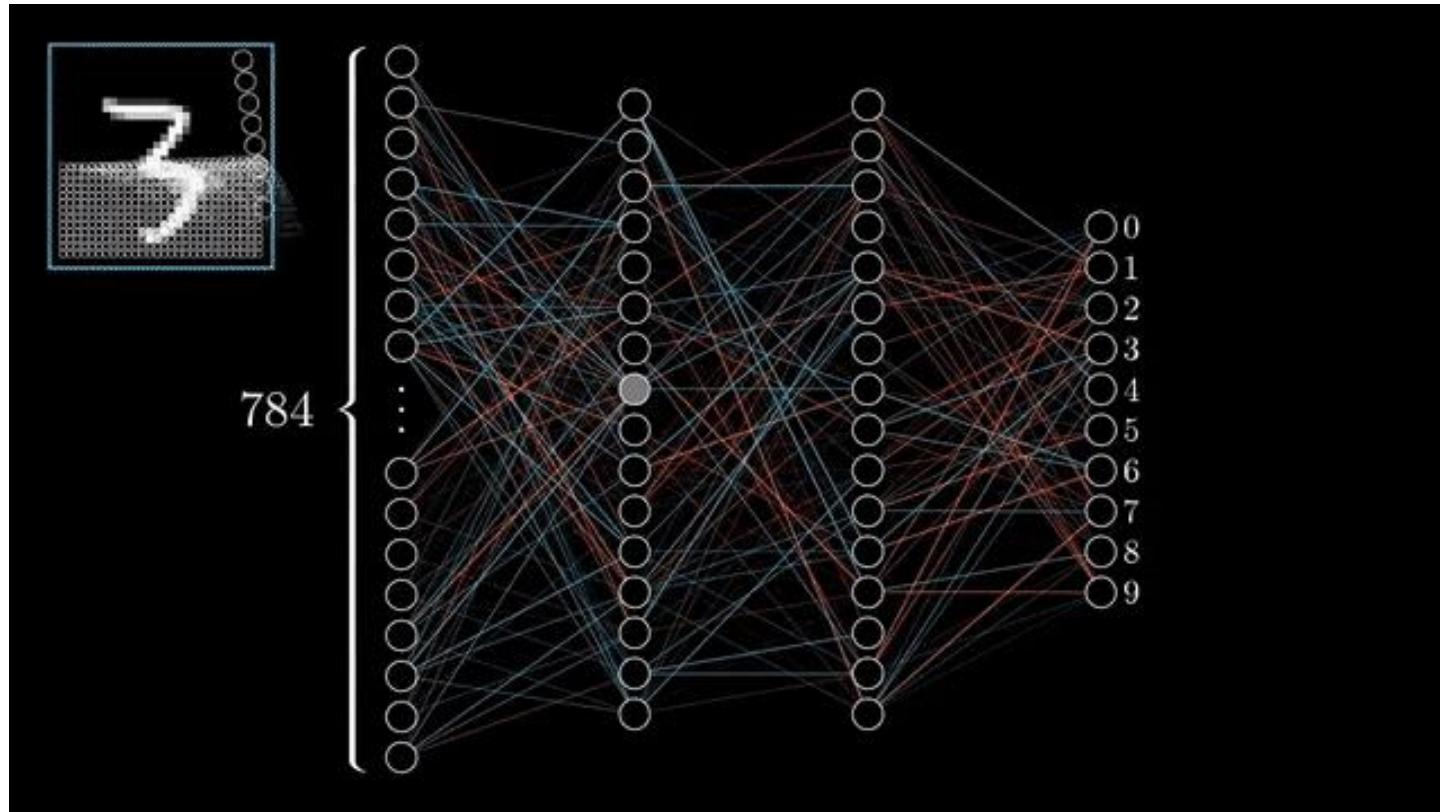
TanH



ReLU

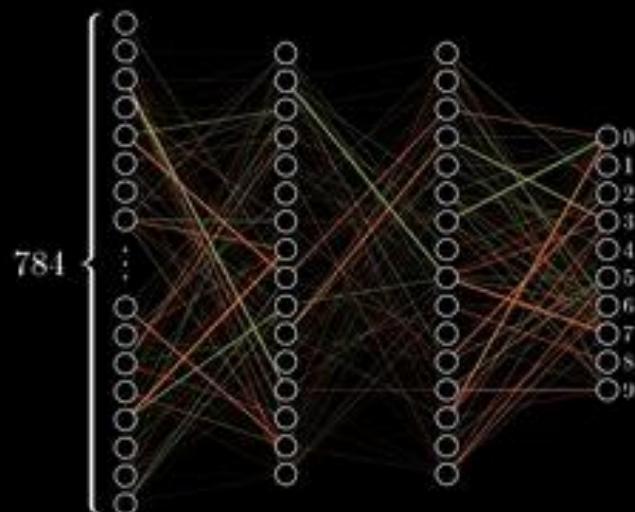


# Initialization

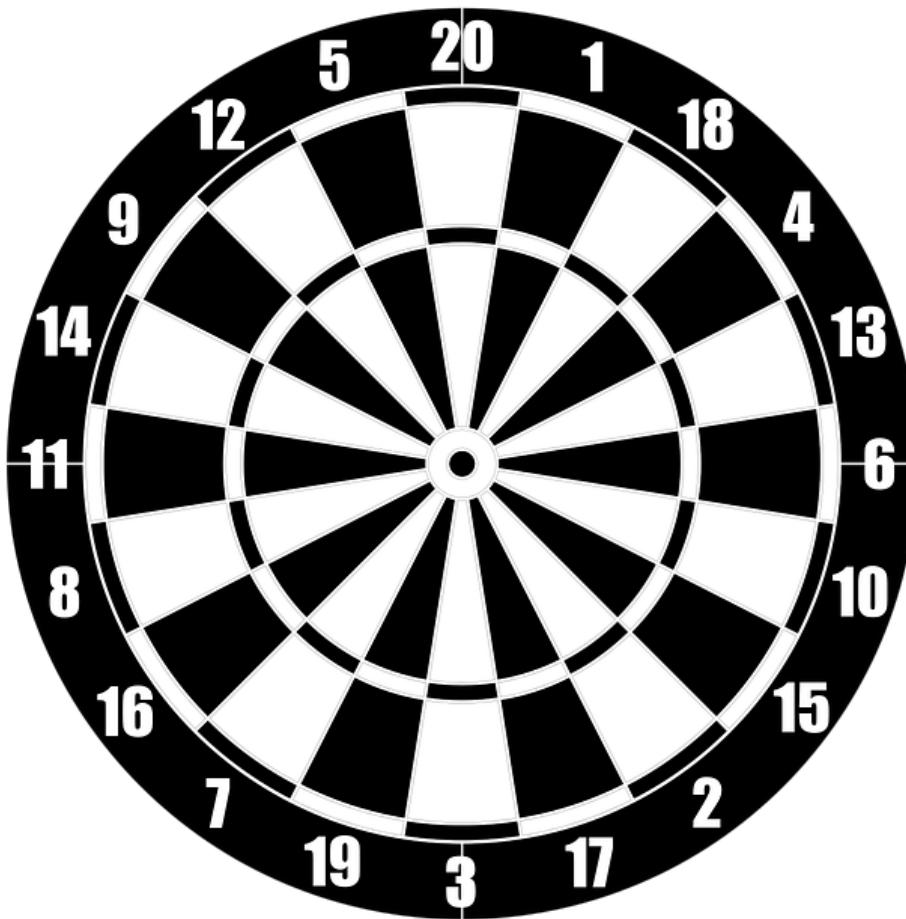


# Training

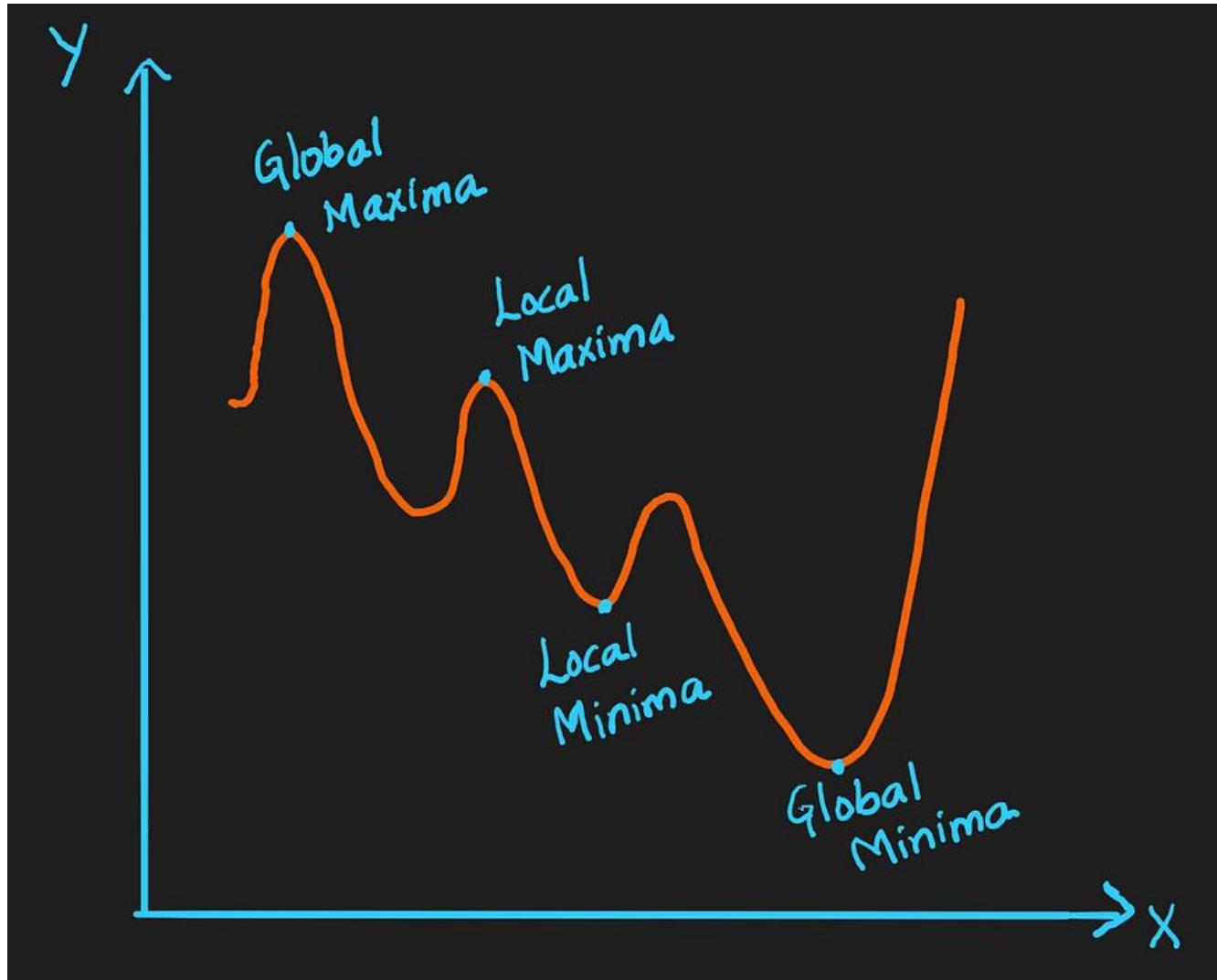
Training in  
progress. . .



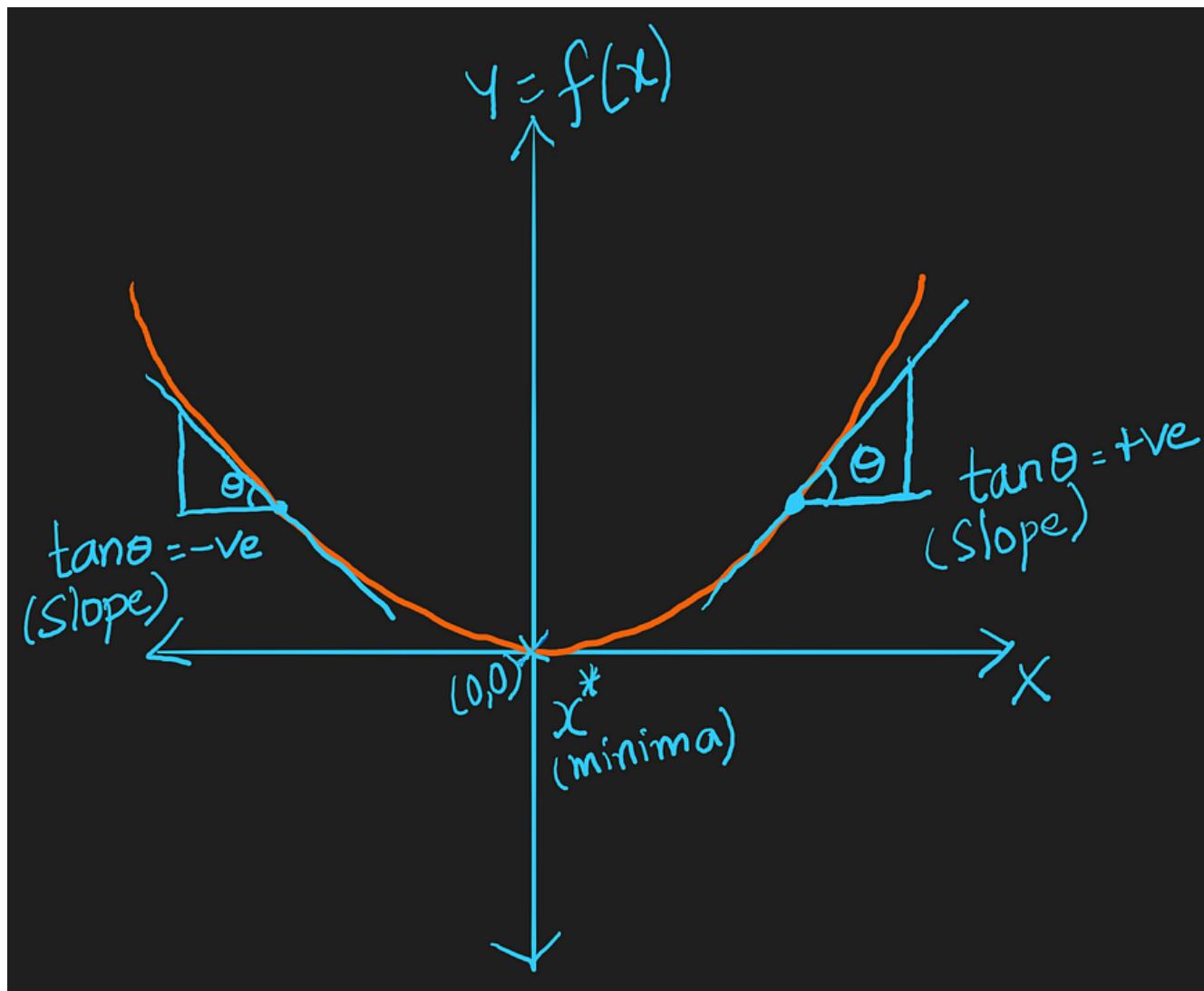
# Learning



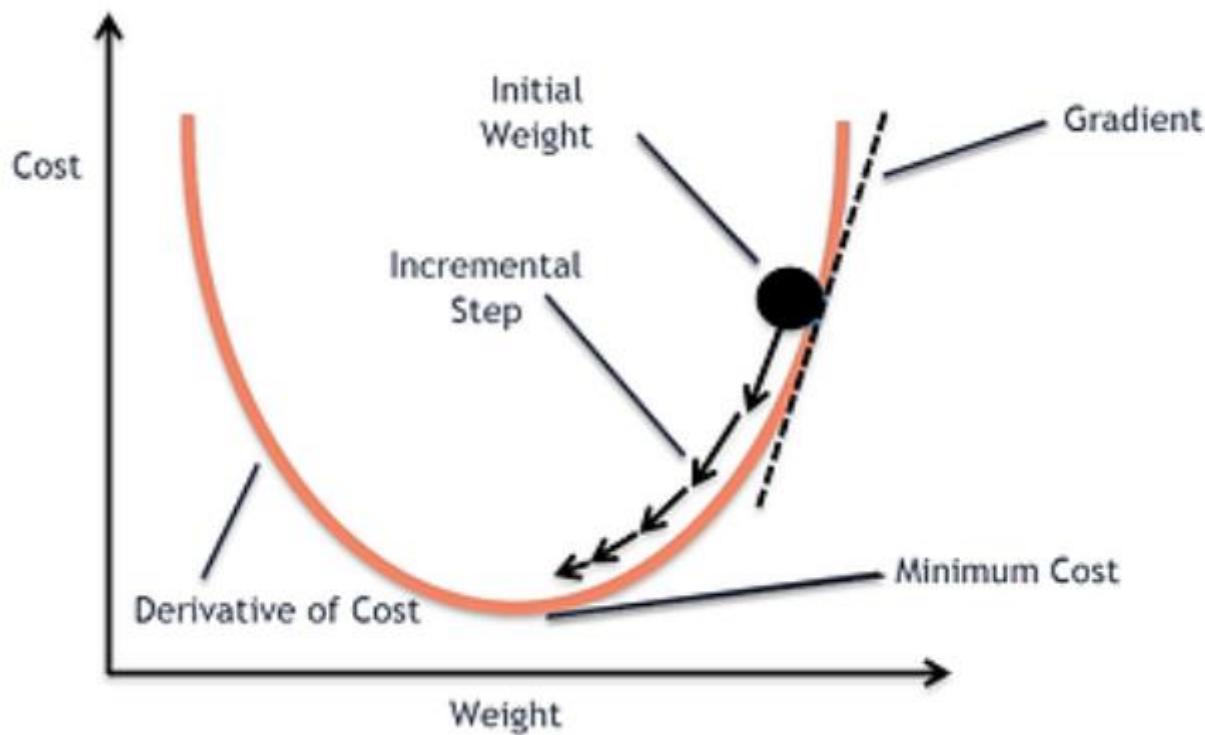
# Maths behind Learning – Maxima/Minima



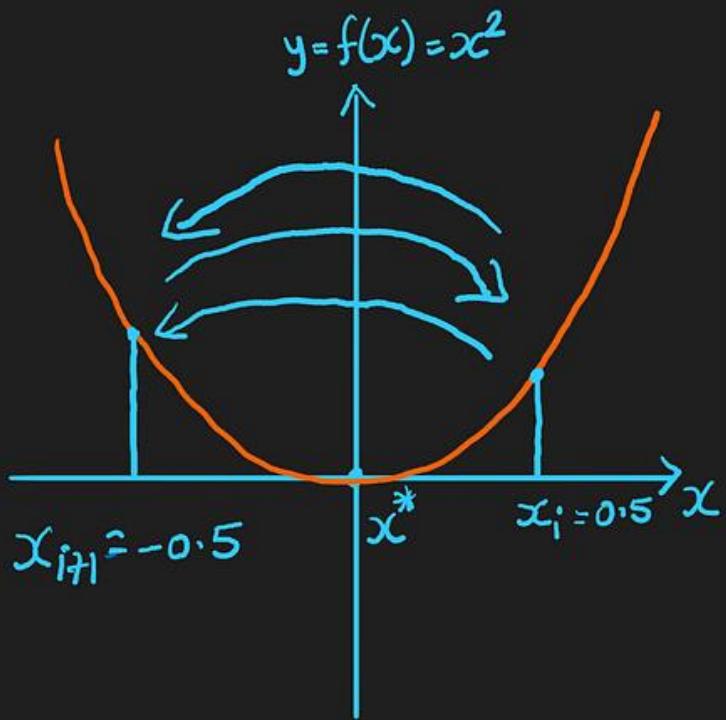
# Maths behind Learning – Gradient/Descent



# Maths behind Learning – Gradient/Descent



# Maths behind Learning – Learning Rate



Oscillating between  
-0.5 and 0.5

Let  $\delta = 1$

$$\frac{df}{dx} = 2x$$

$$x_{i+1} = x_i - \delta \left[ \frac{df}{dx} \right]_{x_i}$$

$$x_{i+1} = 0.5 - 1 [2 \times 0.5]$$

$$x_{i+1} = -0.5$$

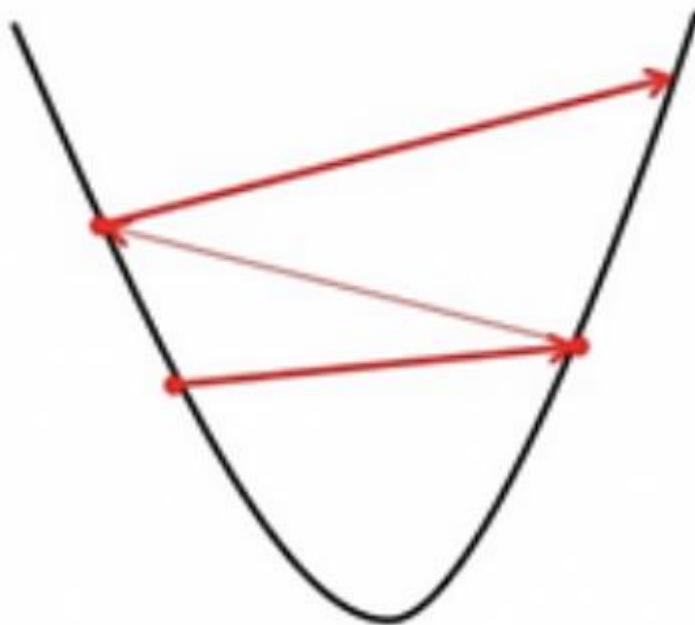
$$x_{i+2} = x_{i+1} - \delta \left[ \frac{df}{dx} \right]_{x_{i+1}}$$

$$x_{i+2} = -0.5 - 1 [2(-0.5)]$$

$$= -0.5 + 1 = 0.5$$

# Maths behind Learning – Learning Rate

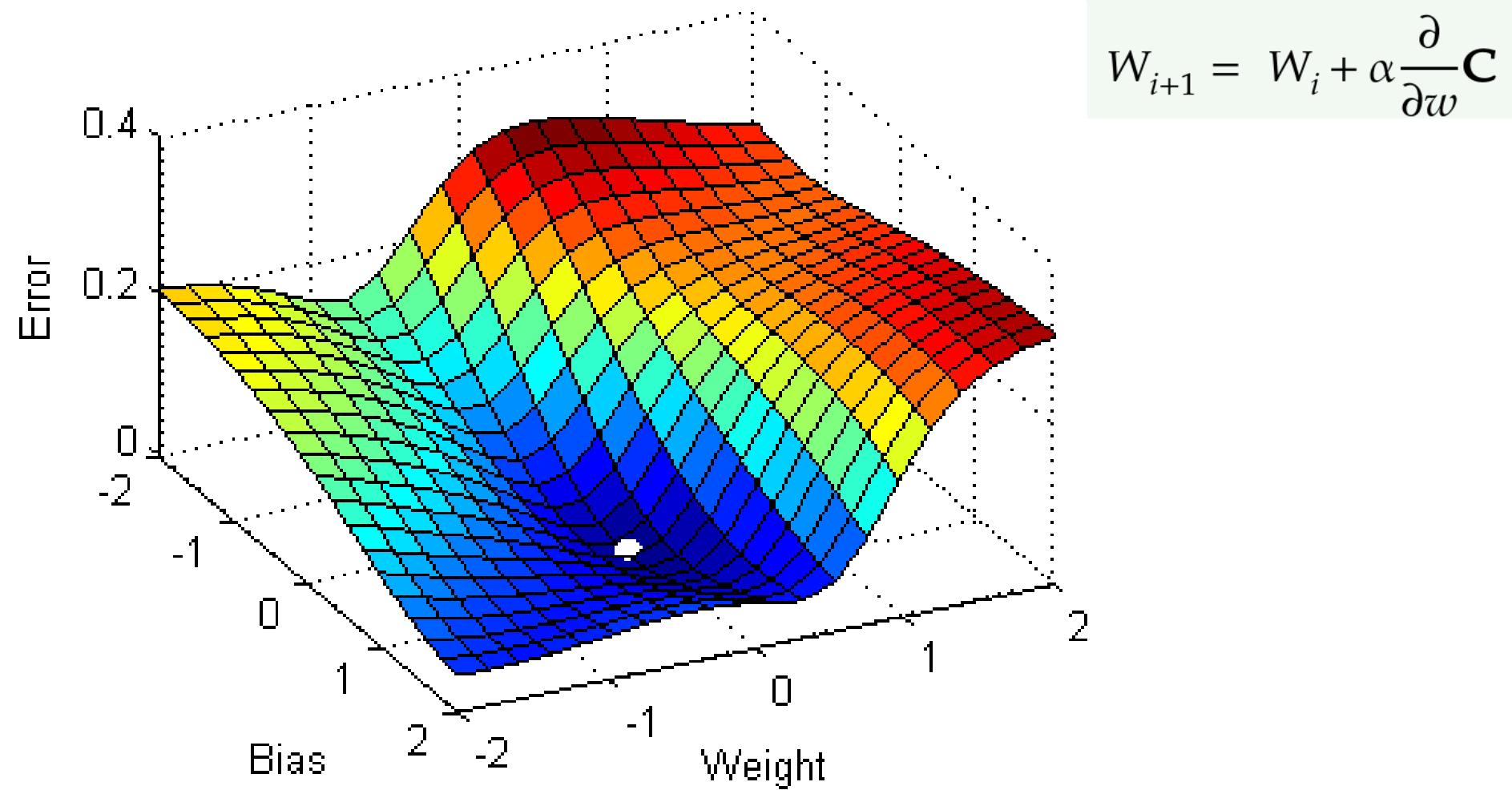
Big learning rate



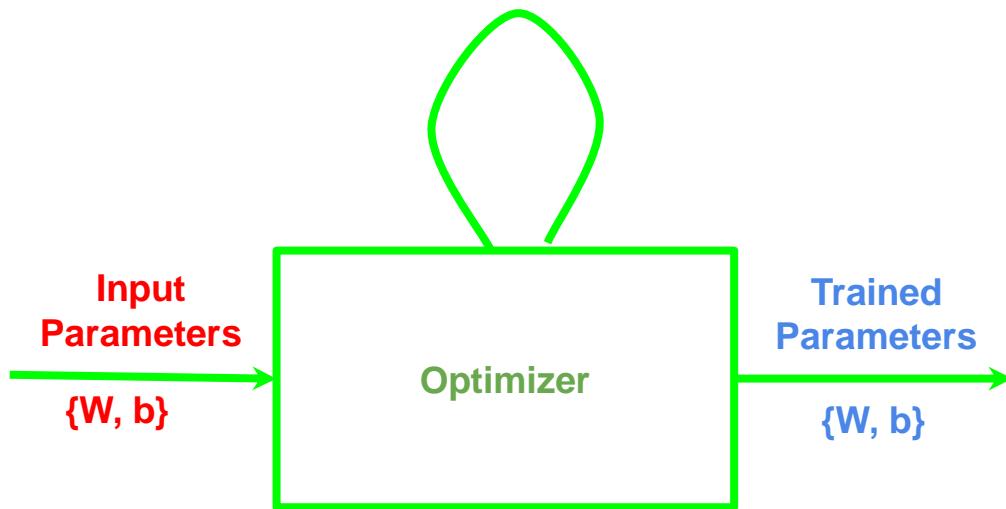
Small learning rate



# Maths behind Learning - ALL



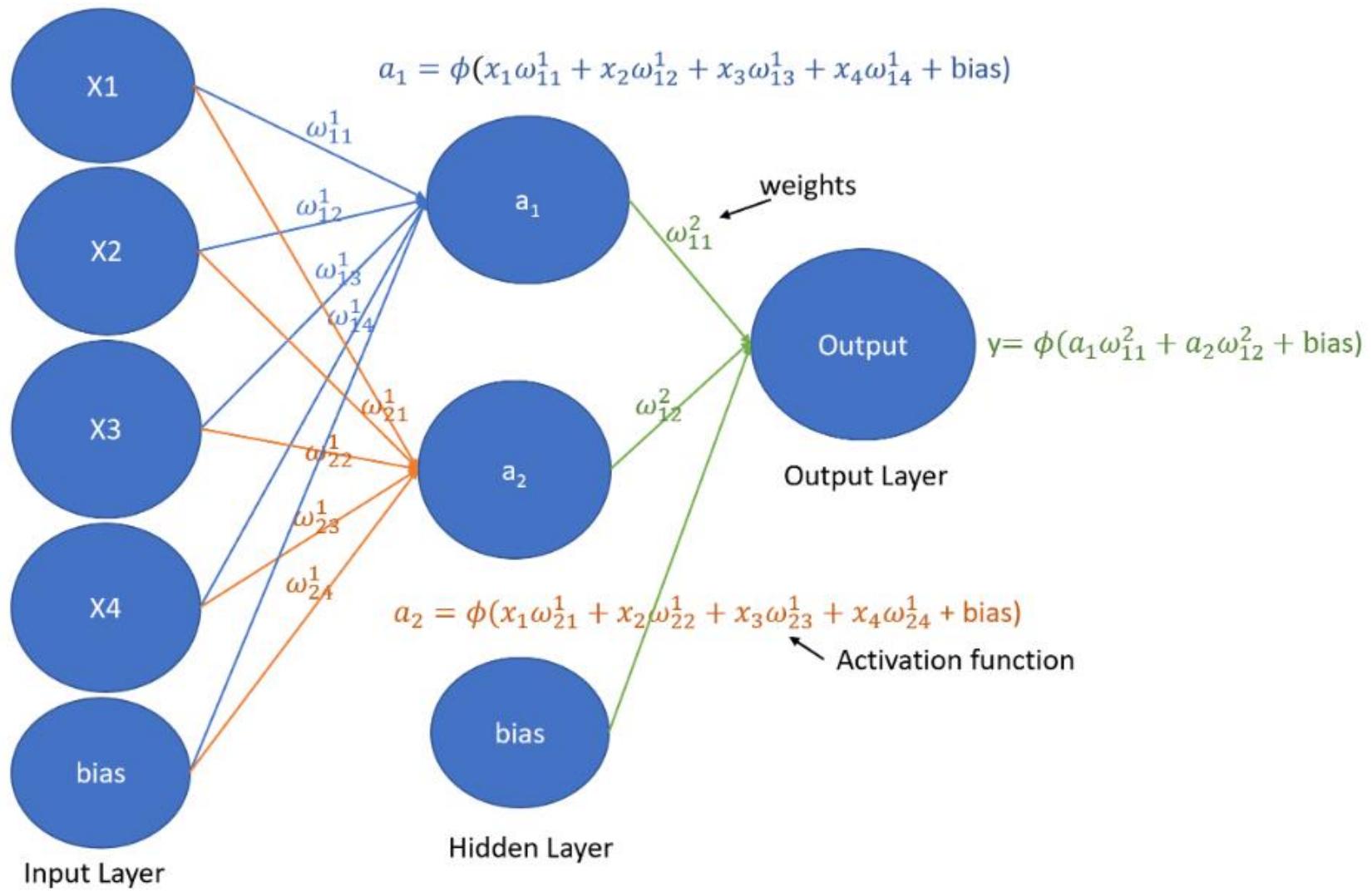
# Optimization Algorithms



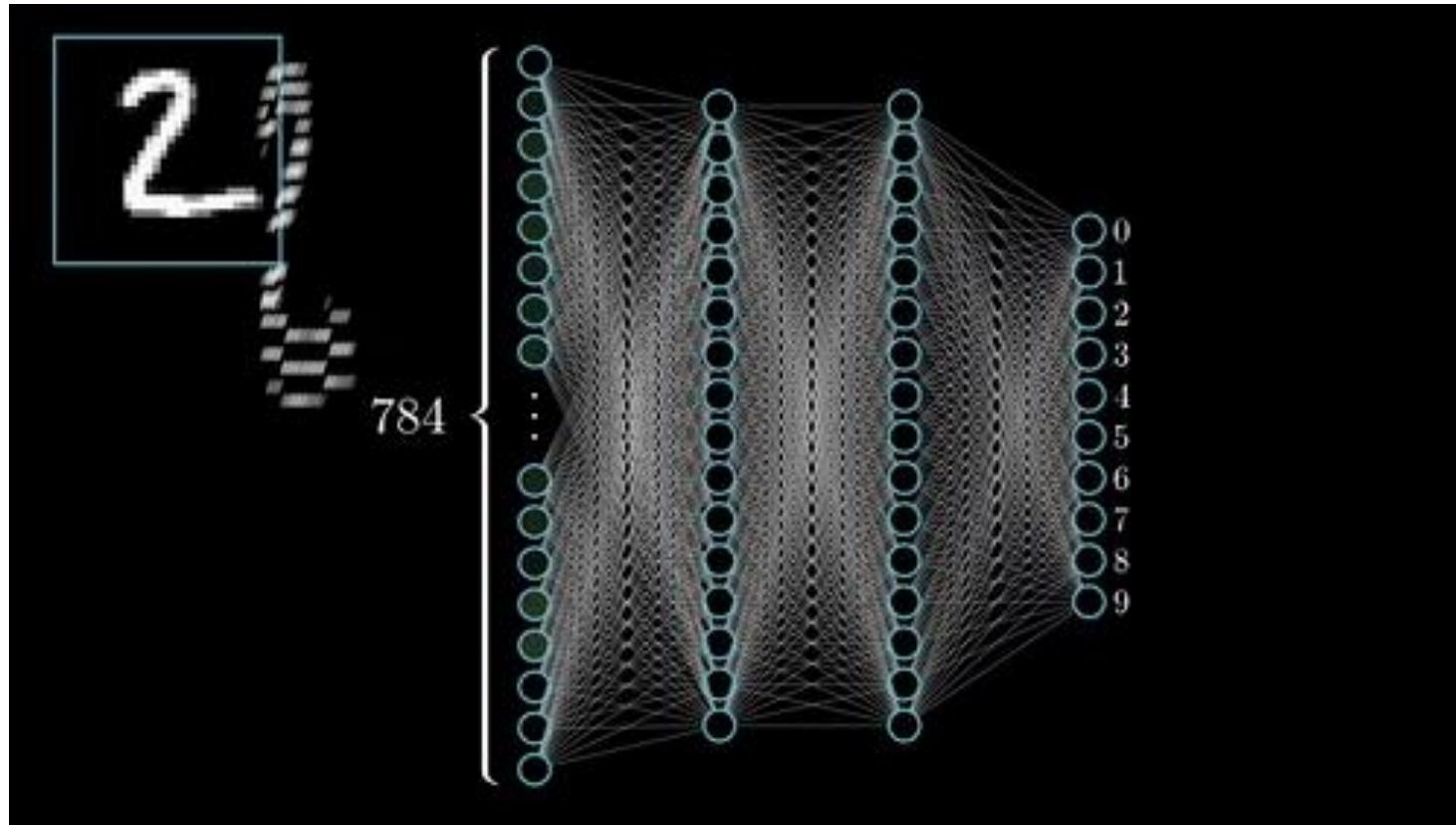
# Optimization Algorithms

Method	Update equation
SGD	$g_t = \nabla_{\theta_t} J(\theta_t)$ $\Delta\theta_t = -\eta \cdot g_t$ $\theta_t = \theta_t + \Delta\theta_t$
Momentum	$\Delta\theta_t = -\gamma v_{t-1} - \eta g_t$
NAG	$\Delta\theta_t = -\gamma v_{t-1} - \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$
Adagrad	$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$
Adadelta	$\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$
RMSprop	$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$
Adam	$\Delta\theta_t = -\frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$

# Summarizing the MATH



# Prediction



# Neural Net for Spam classification

spam	FreeMsg Hey there darling it's been 3 week's n...
ham	Even my brother is not like to speak with me. ...
ham	As per your request 'Melle Melle (Oru Minnamin...
spam	WINNER!! As a valued network customer you have...
spam	Had your mobile 11 months or more? U R entitle...
ham	I'm gonna be home soon and i don't want to tal...
spam	SIX chances to win CASH! From 100 to 20,000 po...
spam	URGENT! You have won a 1 week FREE membership ...
ham	I've been searching for the right words to tha...
ham	I HAVE A DATE ON SUNDAY WITH WILL!!
spam	XXXMobileMovieClub: To use your credit, click ...
ham	Oh k...i'm watching here:)
ham	Eh u remember how 2 spell his name... Yes i di...
ham	Fine if thatås the way u feel. Thatås the wa...
spam	England v Macedonia - dont miss the goals/team...
ham	Is that seriously how you spell his name?

# Neural Net for Spam classification

## Preprocessing data

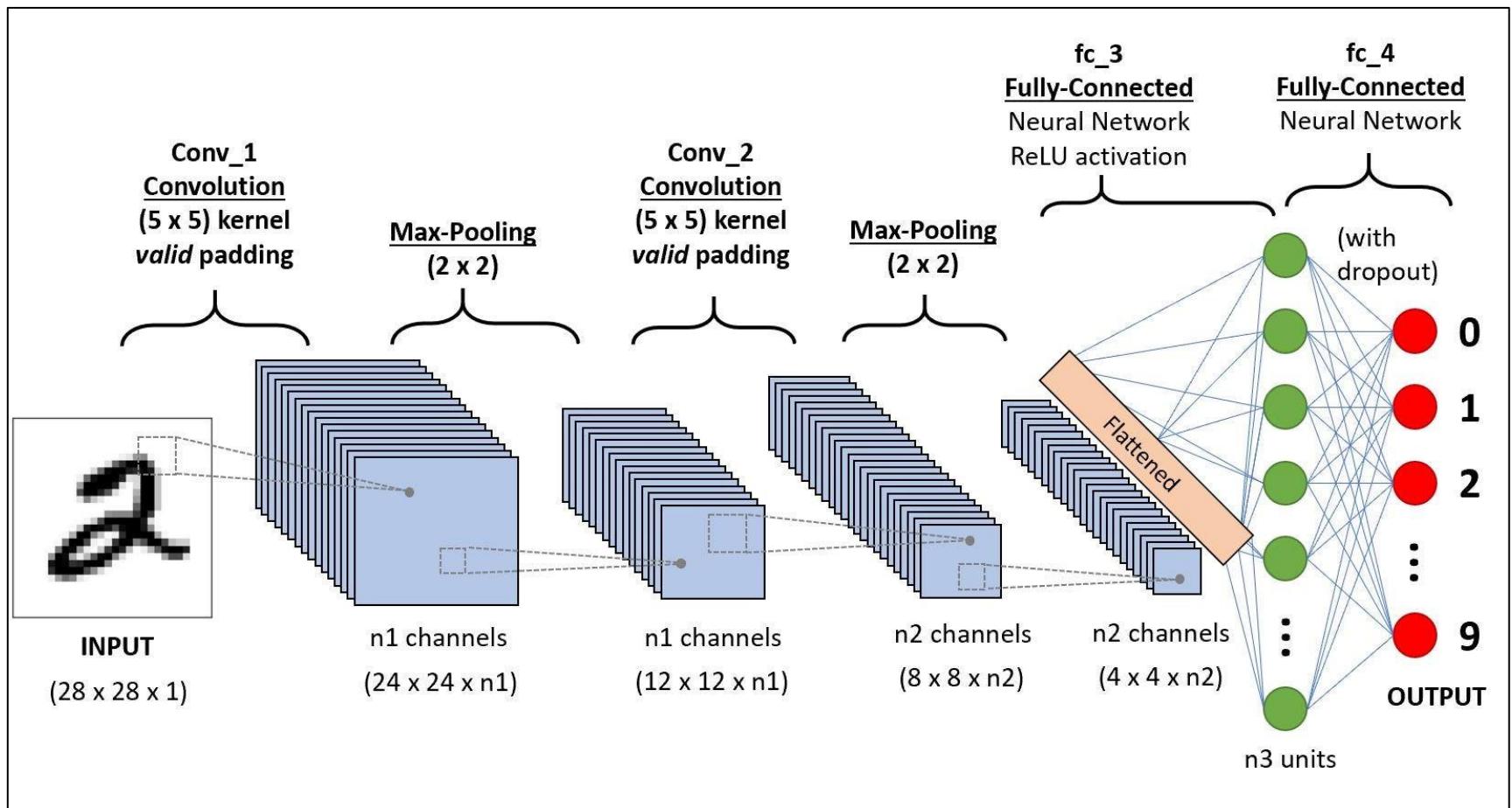
# Neural Net for Spam classification

## **Demo 4: Neural Net for Spam classification**

<NN for spam classification End>

# Convolutional Neural Net

Convolution layer | Pooling layer | Dense layer



# Convolutional Neural Net - Convolution Layer

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0



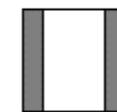
\*

1	0	-1
1	0	-1
1	0	-1



=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0



# Convolutional Neural Net - Convolution Layer

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...	...	...	...	...	...	...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...	...	...	...	...	...	...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...	...	...	...	...	...	...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

+

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

Output

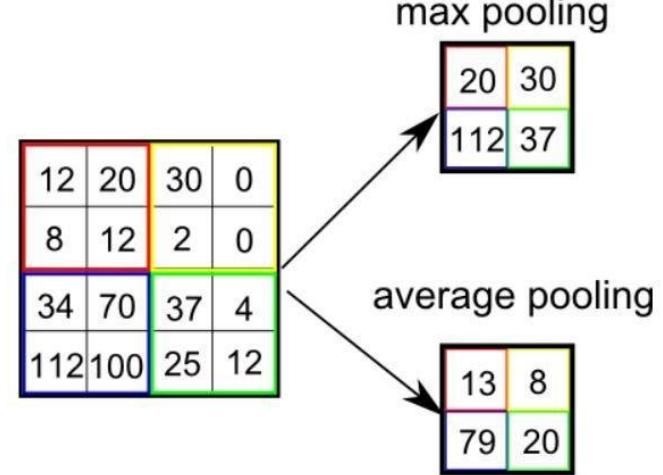
-25	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...

Bias = 1

# Convolutional Neural Net - Pooling Layer

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

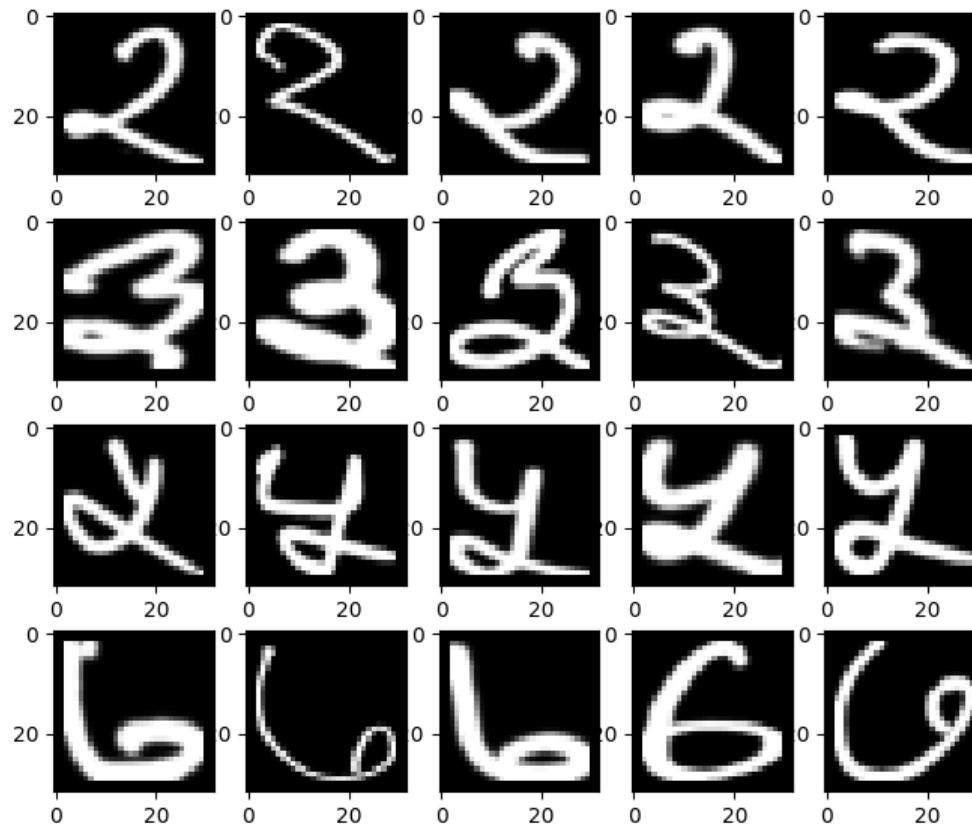


# Modified National Institute of Standards and Technology database

93197245103243759034986680965164954  
30242948320135357468814169690142131  
28232382498291391119966979422633316  
63690360301139315049687103799181722  
3380705698 1432614063  
1795804377 1051939618  
9500511174 3304385467  
0216170956 1020279246  
94321002081403175070751183703929632  
55166276756658168710538319574143978  
71759239430458004046669348131311301  
16796411413123481550794845652540711  
01616755566881728376555002835558045  
64687713073869167364880210608898024  
79731337936249214503851916575991545

# MNIST

# CNN for Devanagari digit classification



# About Dataset

- Grayscale image of size (32x32)
- There are 1025 columns and 1200 rows.
- Columns 1-1024 : contains the pixel values in range 0 to 255.
- column 1025: Label

# **Demo 5: CNN for Devanagari digit classification**

# Phishing

- Fraudulent attempt to obtain sensitive/personal information
- Multiple vectors
- Mails, DMs, SMS, Social network, chat forums, etc.
- Spear Phishing



# Spear Phishing: The problem

- Targeted
- Understand interests of target
- Takes hours of manual efforts
- Higher success rate

# Spear Phishing: How Deep Learning helps

- Build a huge corpus of text from target's social media feeds
- Train a model to generate similar text using words from corpus

# ML based Spear Phishing: Results

[+]20: Much awaited talk given by @imohitch at @PydataPune @PythonPune for AI" [\[link\]](#)

[+]12: Thanks!! [\[link\]](#) CFP deadline: Jun 20 - 12 #PyDataPune [\[link\]](#)

[+]5: It was a comment. [\[link\]](#)

#python #pydatapune #HouseOfResearch @RedHat. Worked on xarray - R and action plan.

@RedHat #MachineLearning #DataScience #PyDataPuneMeetups

[+]9: Yo! I find the Newsletter.

@PydataPune @PyData @NumFOCUS [pic.twitter.com/T5ydVfF5PR](https://pic.twitter.com/T5ydVfF5PR)

[+]13: [\[link\]](#) May 19, 2018. We are PyDaters.

Thanks @katyhuff

@PyData @NumFOCUS #python #pune #machinelearning #dataviz #openscience  
#nonprofit #swag [pic.twitter.com/j9JhQy9Pus](https://pic.twitter.com/j9JhQy9Pus)

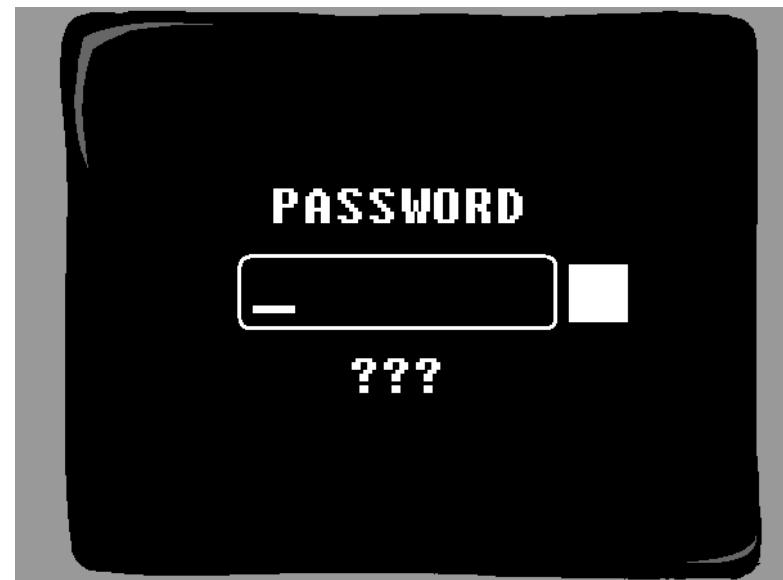
[+]7: Interested in today's @WiDSPune2019 on numpy and show your Attention please!!!

The goal is going @pranshutople14 conclude the maths behind NeuralNets [\[link\]](#)

#PyData #NumFOCUS @PyData #PyData\_Meetup @Pune @machinelearnflx @Pune

# Dictionary attacks

- Dictionary contains commonly used password
- Try every password from dictionary unless you succeed
- Attack may work if the password is similar to
  - Pydata\_123\$\$
  - ##pydata##Pun3
  - PyData@@098



# Dictionary attacks: data collection

1.4 Billion Clear Text Credentials  
Discovered in a Single Database



Julio Casal [Follow](#)

Dec 9, 2017 · 5 min read

# Dictionary attacks: Approach

- Skeletonize the passwords
  - <first\_name>\_<last\_name>##123
  - <dob><mob><First\_name>
  - <dob><mob>\_<Lastname>\$\$
  - <last\_name>##<random\_num>
  - Etc ...
- Train generative models on corpus of skeletons
- Basic personal info + skeletons = passwords

# Traditional API Security

## Rule-based vs. Machine Learning-based Methods



Conventional methods in the cybersecurity field are **rule based**. A subcase of this is the blacklist & whitelist-based approach.



Rule-based methods are rigid. They are **not capable of identifying novel cases** like zero-day threats, new browser fingerprinting techniques, or new IoT devices.

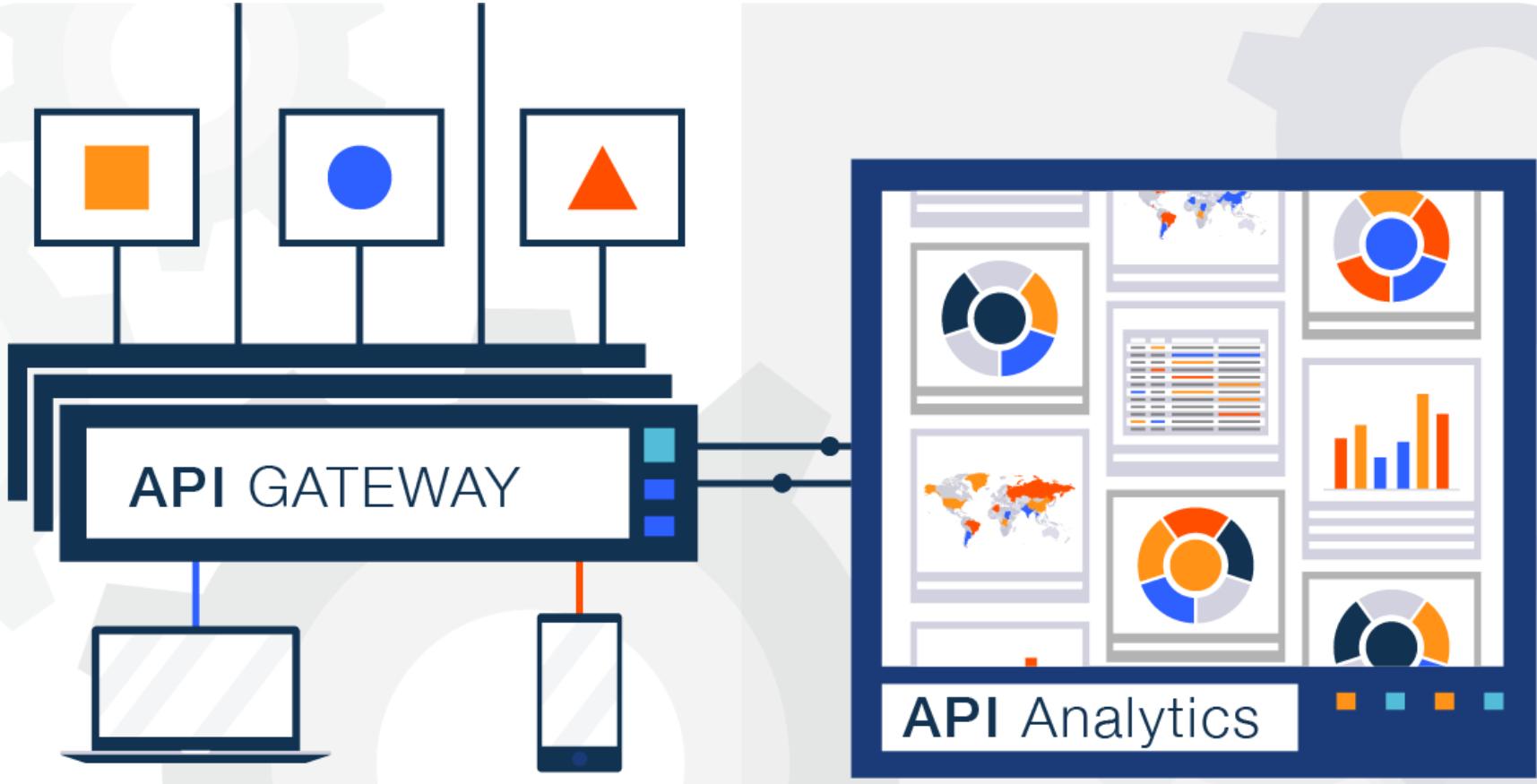


ML-based methods are **more flexible**. They can learn patterns, model the relationship between the patterns and identify new cases based on their similarity to known ones.



ML-based methods **require less maintenance** than traditional methods. Keeping a rule base up to date is more difficult than constantly feeding the ML model with new data.

# Adding ML to API Security



# Type of API Security Threat

## Bot Attack

Bot attacks targeting APIs

## Injection attacks

Injection attacks such as SQL injection and cross-site scripting(xss)

## Authentication Attack

Authentication attacks such as brute-force & credential stuffing

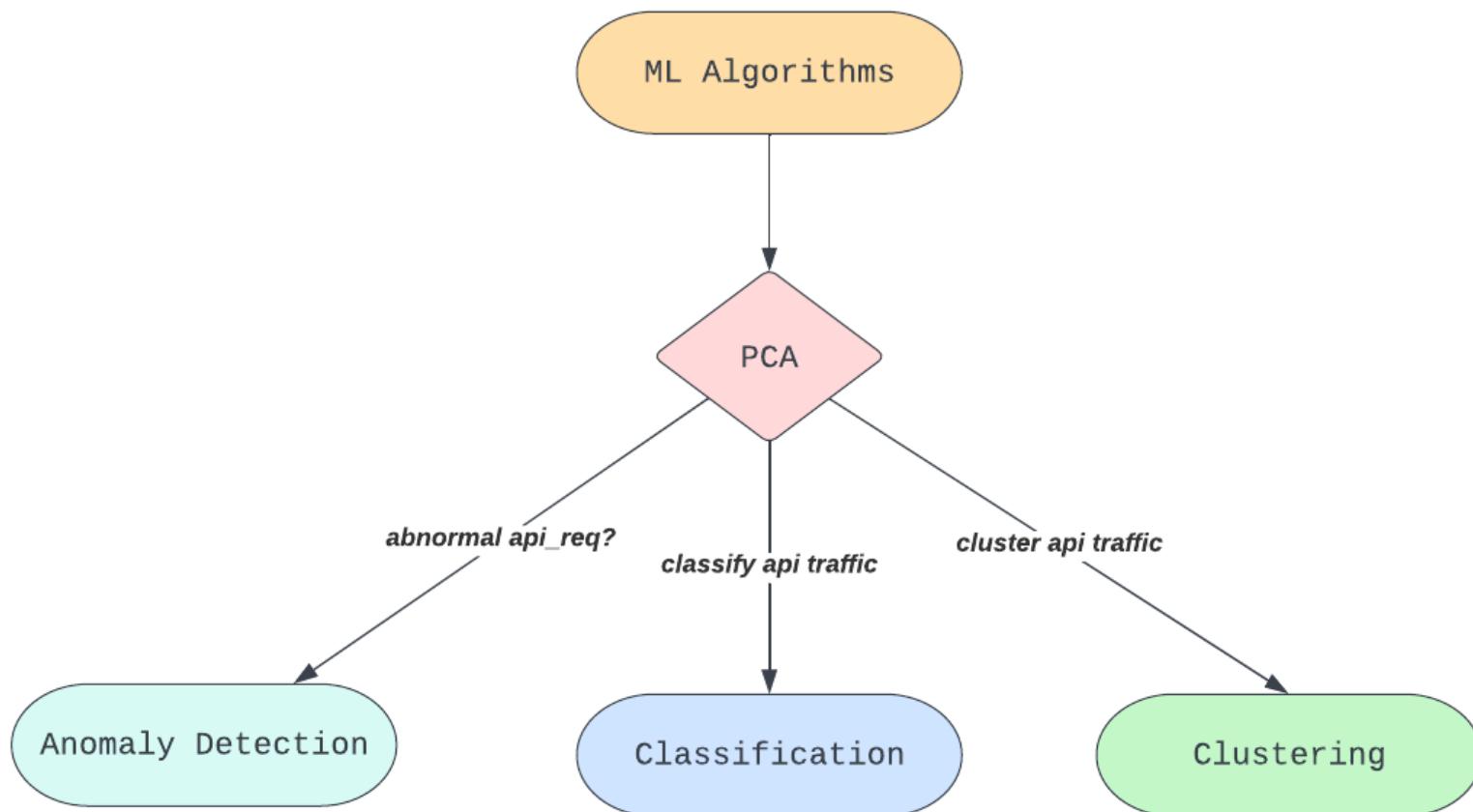
## API Specific Attack

API-specific attacks such as scraping and abuse

## Data privacy and confidentiality

Breaches of data privacy and confidentiality

# Applying ML to API Security



# DEMO: Before the break

## Modeling & evaluation of ML model

Import all dependencies  
Selecting dependent and independent variables

## Modeling data with Xgboost Classifier

Storing the trained model into  
pickle file  
Storing model with joblib

Integration with website  
Demo

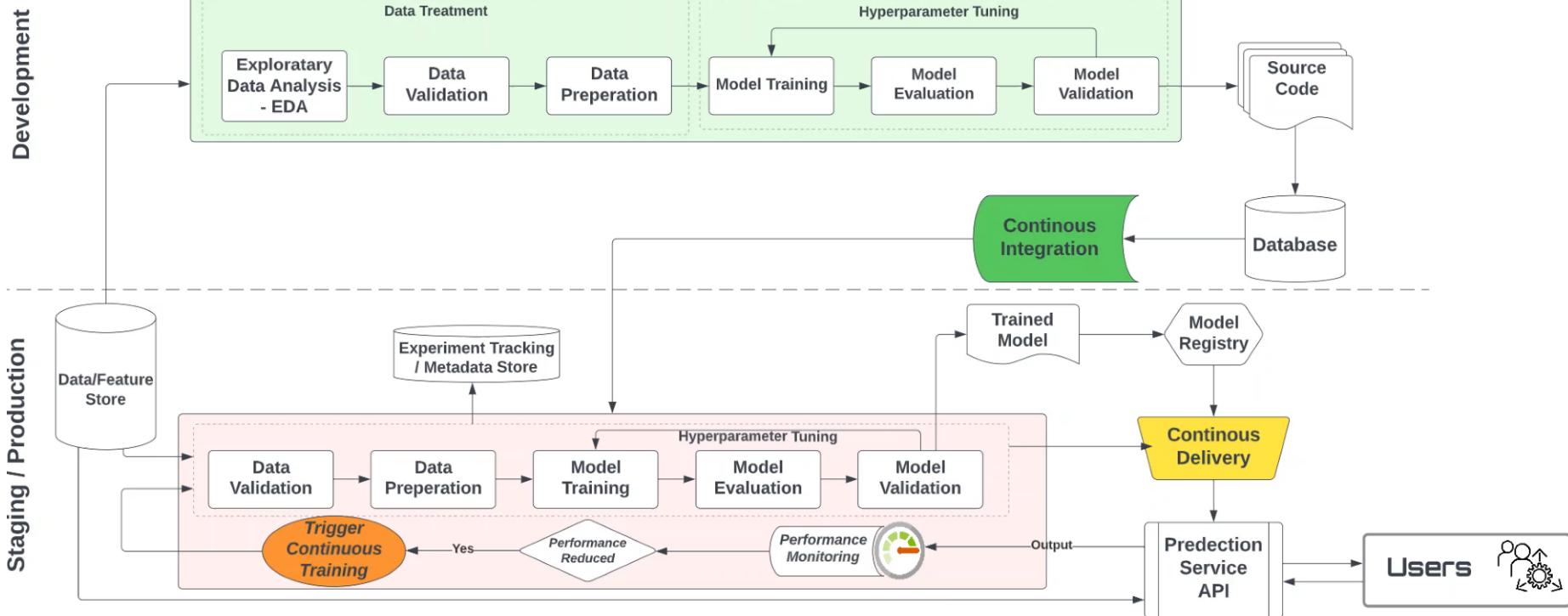
[LINK](#)

# **<Interval>**

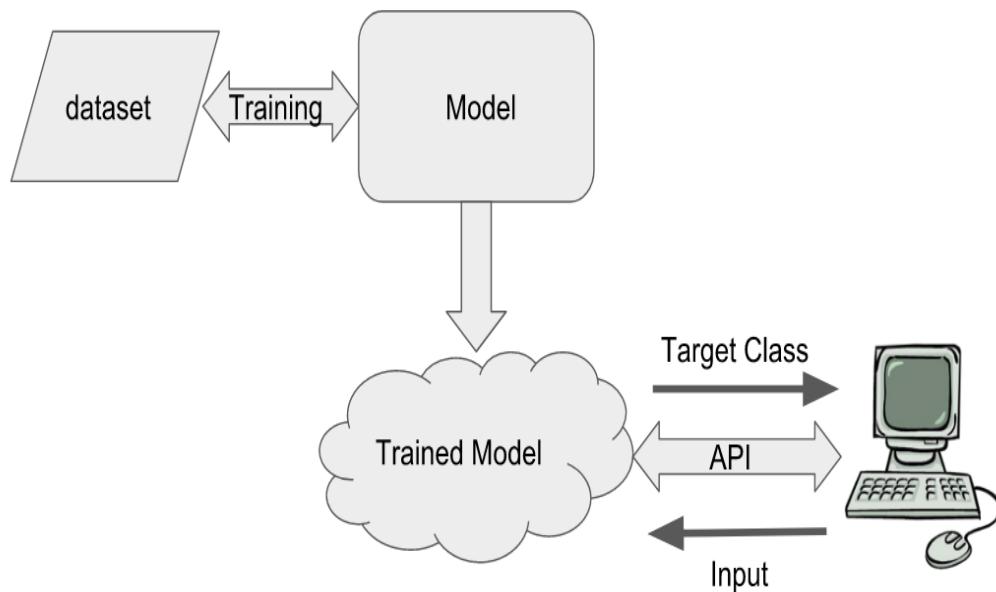
## **Part two: Security for ML**

# ML in Production

***“As MIOps hits maturity, time to consider security”***



# Possible Attack Targets: Almost Everything



- Poison Dataset
- Skewed Models
- Adversarial Attack
- Steal Models
- I/P Data inference
- Framework Exploits
- Network attacks

# AI Security Assessment

- Recon
- Attack surface identification
- Threat modeling
- Adversarial Learning Attack
- Model Stealing
- Model Skewing
- Input data inference
- Model inversion
- Hybrid attacks

# The Recon

- Intel on Application and Model internals
- Research Papers, Blogs and product pages
- Conference talks, demos, social media content for target application/appliance
- Helps to plan efficient attack scenarios
- Get your hands on target application

# Adversarial Learning

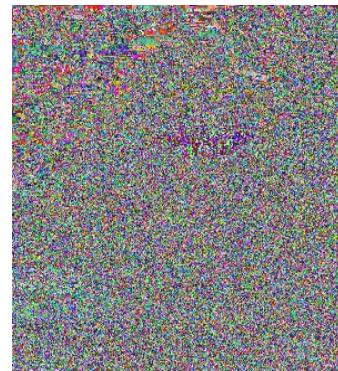
**Is our model actually “learning” what we expect it to learn?**

# Adversarial Learning

Specially craft an input with an intention to produce desired prediction from target model.



Street Sign



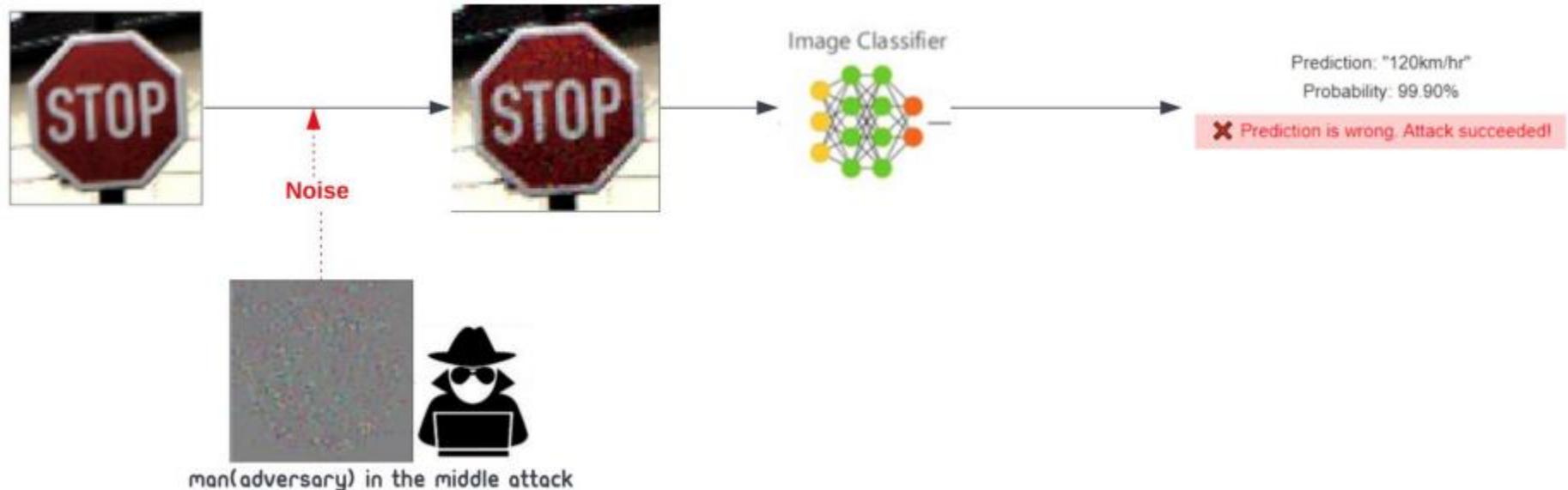
Perturbations



Mail Box

# Adversarial Learning

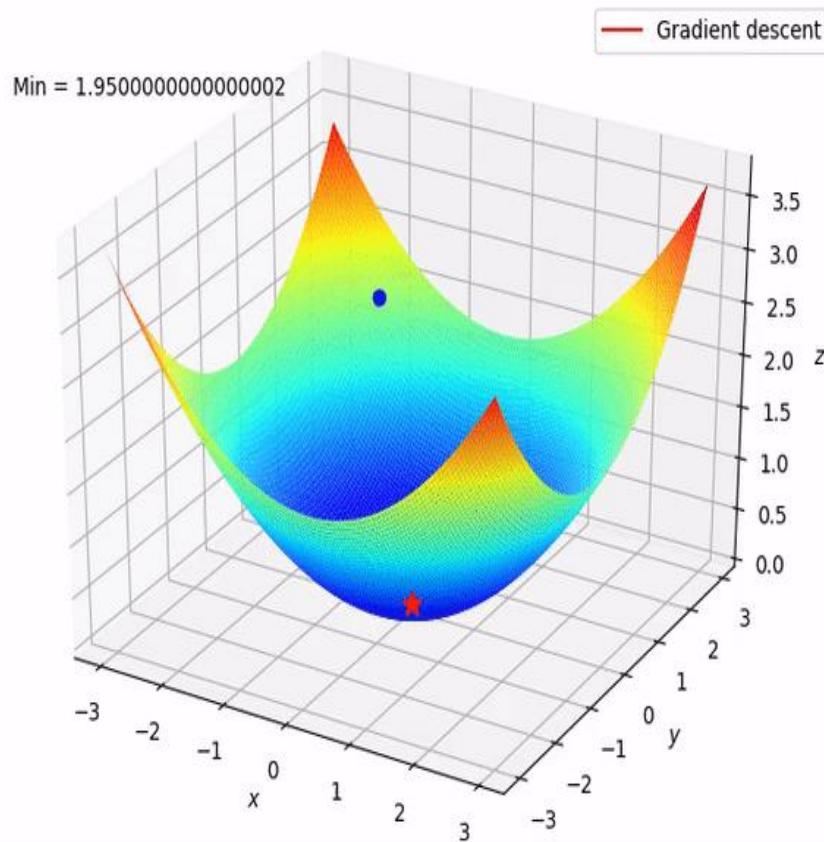
## Demo 6: Hack the ML



<https://hacktheml.web.app>

# Adversarial Learning

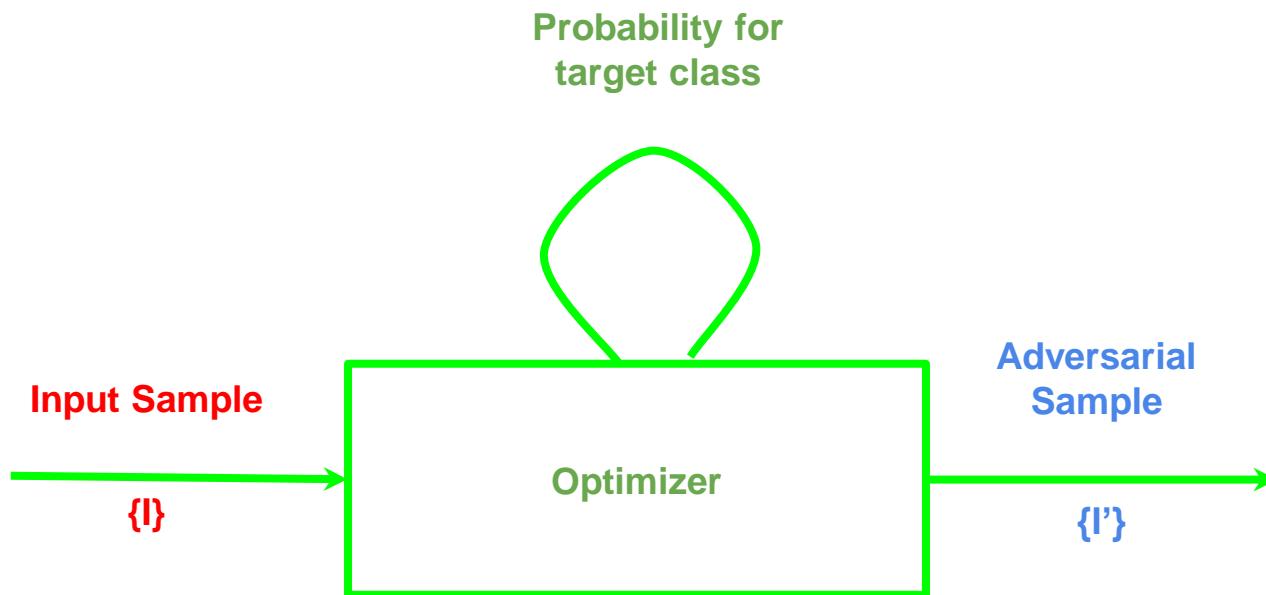
Adversarial attack as an optimization problem



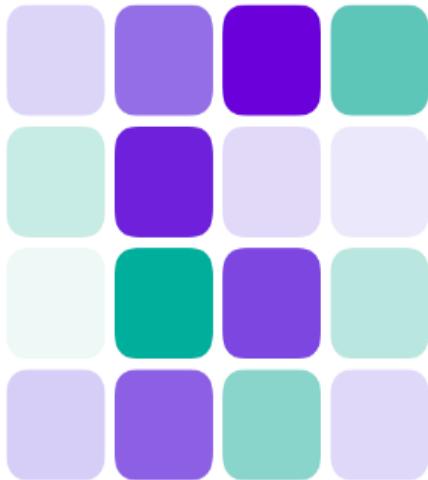
$$I_{i+1} = I_i + \alpha \frac{\partial}{\partial I} C$$

# Adversarial Learning Attacks

Specially craft an input with an intention to produce desired prediction from target model.



# Adversarial Learning Attacks



Foolbox

Ref: <https://github.com/bethgelab/foolbox>



Cleverhans

Ref: <https://github.com/cleverhans-lab/cleverhans>

# Experimental Adversarial attacks

## One Pixel Attack for Fooling Deep Neural Networks

Jiawei Su\*, Danilo Vasconcellos Vargas\* and Kouichi Sakurai

**AllConv**



SHIP

CAR(99.7%)

**NiN**



HORSE

FROG(99.9%)

**VGG**



DEER

AIRPLANE(85.3%)

# Experimental Adversarial attacks

## A General Framework for Adversarial Examples with Objectives



Fig. 5. An example of digital dodging. Left: An image of actor Owen Wilson (from the PubFig dataset [40]), correctly classified by VGG143 with probability 1.00. Right: Dodging against VGG143 using AGN's output (probability assigned to the correct class <0.01).

# Realworld Attacks

**Researchers Trick Cylance Antivirus Into Thinking Malware Is Trusted Software**

NICOLE LINDSEY · AUGUST 2, 2019

**How do we deal with copying content, or the first adversarial attack in prod**

Avito company blog , Programming , algorithms , Image processing , Machine learning

# Realworld Attacks

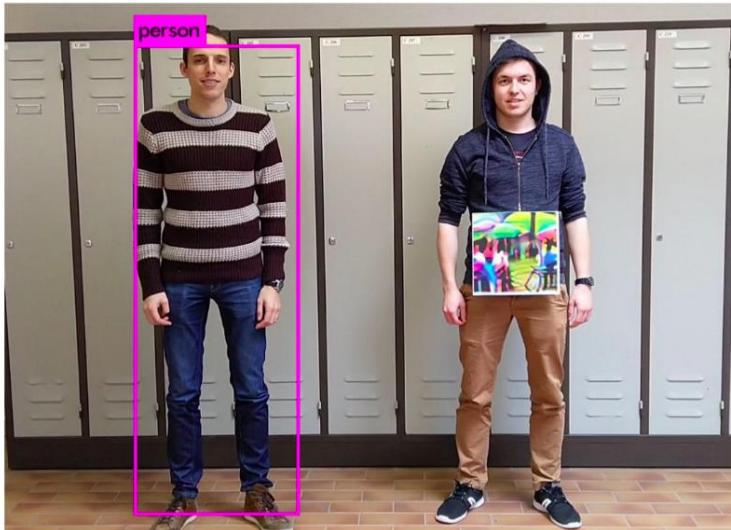
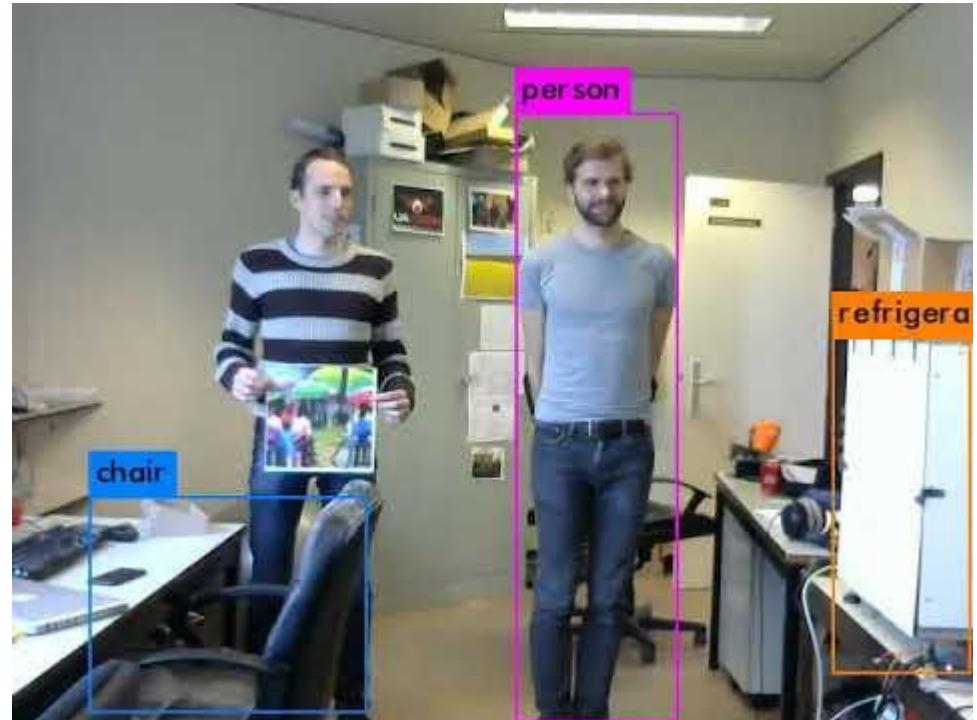
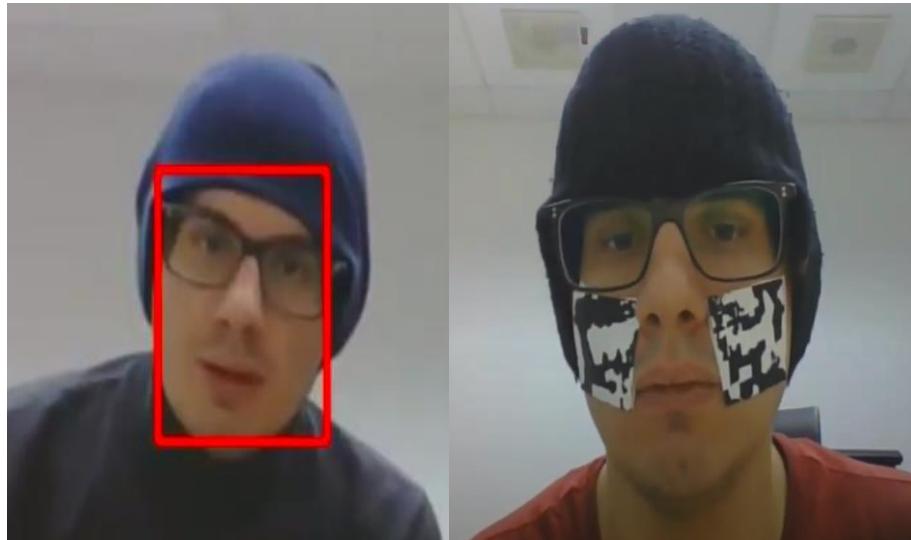


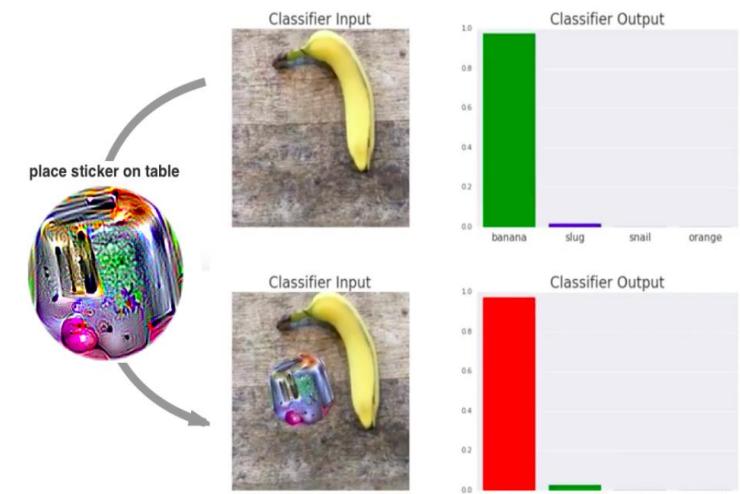
Figure 1: We create an adversarial patch that is successfully able to hide persons from a person detector. Left: The person without a patch is successfully detected. Right: The person holding the patch is ignored.



# Adversarial Attacks in Physical Domain

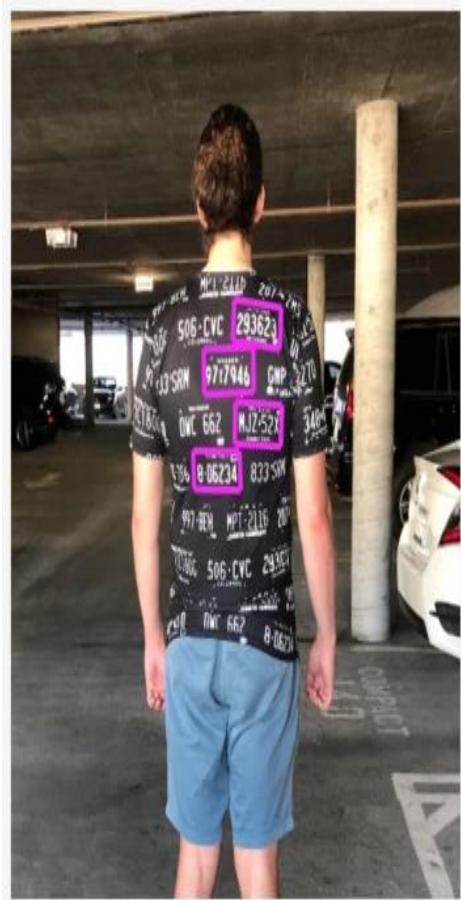


<https://arxiv.org/pdf/1910.06261.pdf>



<https://arxiv.org/pdf/1712.09665.pdf>

# Mis-classification attacks Without any Fancy Maths



<https://adversarialfashion.com/>



Distorted Inputs

# Mis-classification attacks Without any Fancy Maths



# Adversarial attacks Without any Fancy Maths



<https://distill.pub/2019/activation-atlas/>

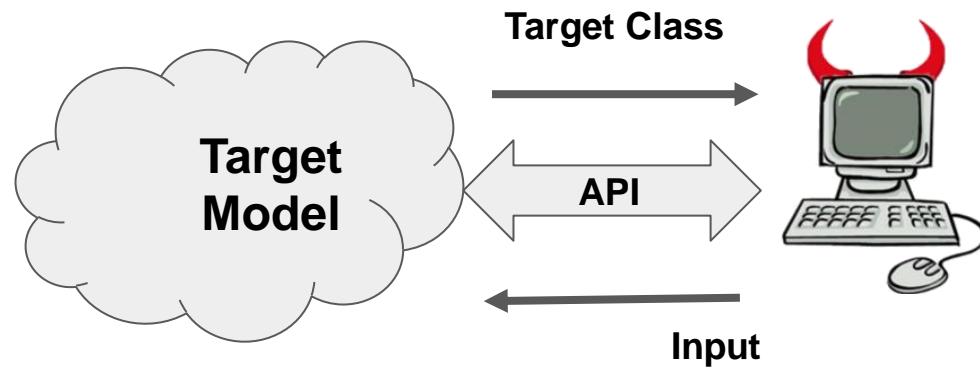
The patient has a history of **back pain** and chronic **alcohol abuse** and more recently has been seen in several...

Adversarial text substitution

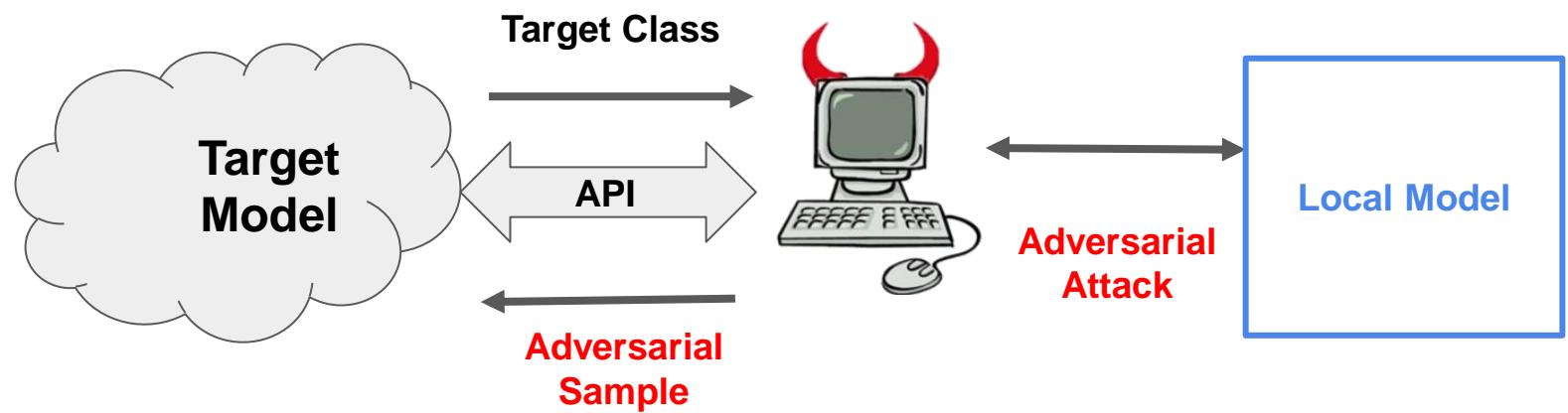
The patient has a history of **lumbago** and chronic **alcohol dependence** and more recently has been seen in several...

<https://science.sciencemag.org/content/363/6433/1287.full?ijkey=OXnSsEp.lagl6&keytype=ref&siteid=sci>

# Adversarial Attacks on Black-Box Models

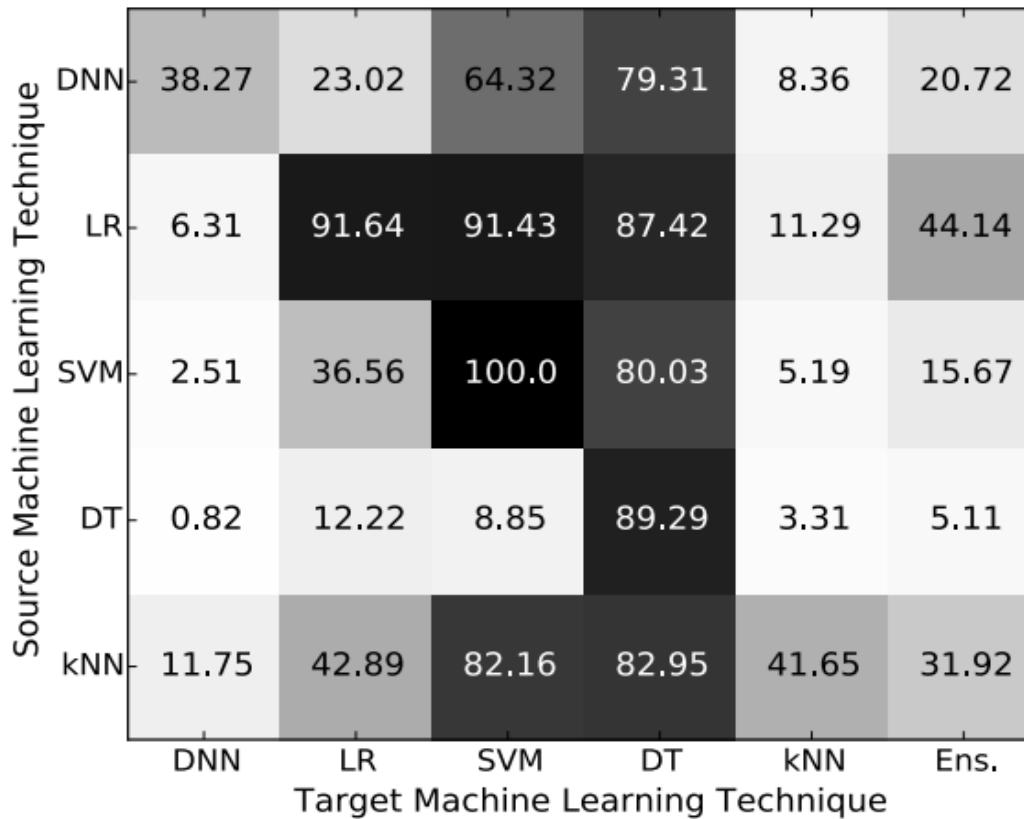


# Adversarial Attacks on Black-Box Models

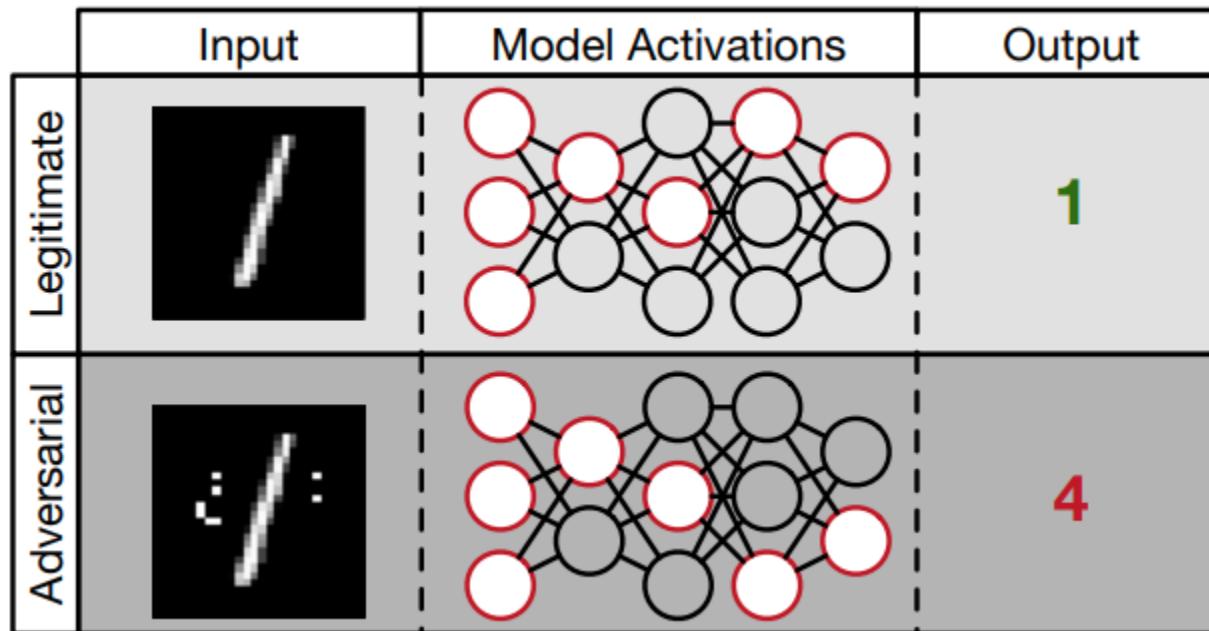


# Adversarial Attacks on Black-Box Models

## Cross-technique Transferability



# Adversarial Attacks on Black-Box Models



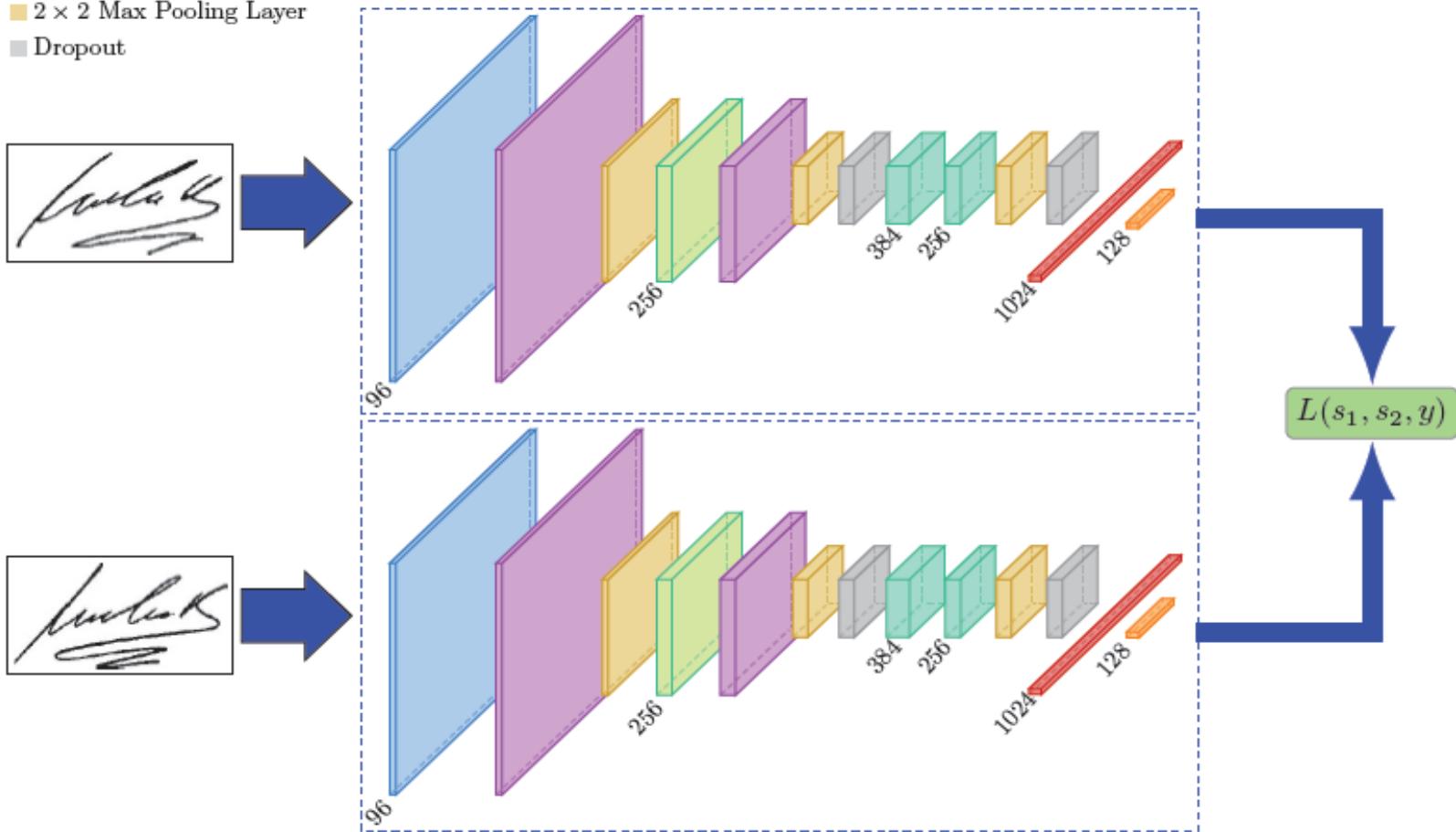
## Demo 7: Adversarial Attacks

# Mitigation: Adversarial Attacks

# Adversarial Attack on FaceNet

# Siamese Neural Network

- $11 \times 11$  Convolutional Layer + ReLU
- $5 \times 5$  Convolutional Layer + ReLU
- $3 \times 3$  Convolutional Layer + ReLU
- $2 \times 2$  Max Pooling Layer
- Fully Connected Layer + ReLU
- F.C. Layer + ReLU + Dropout
- Local Response Normalisation
- Dropout



# FaceNet Architecture

layer	size-in	size-out	kernel	param
conv1	$220 \times 220 \times 3$	$110 \times 110 \times 64$	$7 \times 7 \times 3, 2$	9K
pool1	$110 \times 110 \times 64$	$55 \times 55 \times 64$	$3 \times 3 \times 64, 2$	0
rnorm1	$55 \times 55 \times 64$	$55 \times 55 \times 64$		0
conv2a	$55 \times 55 \times 64$	$55 \times 55 \times 64$	$1 \times 1 \times 64, 1$	4K
conv2	$55 \times 55 \times 64$	$55 \times 55 \times 192$	$3 \times 3 \times 64, 1$	111K
rnorm2	$55 \times 55 \times 192$	$55 \times 55 \times 192$		0
pool2	$55 \times 55 \times 192$	$28 \times 28 \times 192$	$3 \times 3 \times 192, 2$	0
conv3a	$28 \times 28 \times 192$	$28 \times 28 \times 192$	$1 \times 1 \times 192, 1$	37K
conv3	$28 \times 28 \times 192$	$28 \times 28 \times 384$	$3 \times 3 \times 192, 1$	664K
pool3	$28 \times 28 \times 384$	$14 \times 14 \times 384$	$3 \times 3 \times 384, 2$	0
conv4a	$14 \times 14 \times 384$	$14 \times 14 \times 384$	$1 \times 1 \times 384, 1$	148K
conv4	$14 \times 14 \times 384$	$14 \times 14 \times 256$	$3 \times 3 \times 384, 1$	885K
conv5a	$14 \times 14 \times 256$	$14 \times 14 \times 256$	$1 \times 1 \times 256, 1$	66K
conv5	$14 \times 14 \times 256$	$14 \times 14 \times 256$	$3 \times 3 \times 256, 1$	590K
conv6a	$14 \times 14 \times 256$	$14 \times 14 \times 256$	$1 \times 1 \times 256, 1$	66K
conv6	$14 \times 14 \times 256$	$14 \times 14 \times 256$	$3 \times 3 \times 256, 1$	590K
pool4	$14 \times 14 \times 256$	$7 \times 7 \times 256$	$3 \times 3 \times 256, 2$	0
concat	$7 \times 7 \times 256$	$7 \times 7 \times 256$		0
fc1	$7 \times 7 \times 256$	$1 \times 32 \times 128$	maxout p=2	103M
fc2	$1 \times 32 \times 128$	$1 \times 32 \times 128$	maxout p=2	34M
fc7128	$1 \times 32 \times 128$	$1 \times 1 \times 128$		524K
L2	$1 \times 1 \times 128$	$1 \times 1 \times 128$		0

## **Demo 8 : Adversarial Attack on FaceNet**

# Model Stealing Attacks

## Offline Attacks

Stealing models that are deployed on device

## Online Attacks

Stealing models deployed on cloud with black box access

# Offline Model Stealing Attacks

- Applications like WAF, Object detection, Antiviruses, IDS/IPS etc. require to work without internet connectivity
- Models are deployed on device to make them available for prediction
- Deployment may or may not employ any protection mechanism

# Offline Model Stealing Attacks

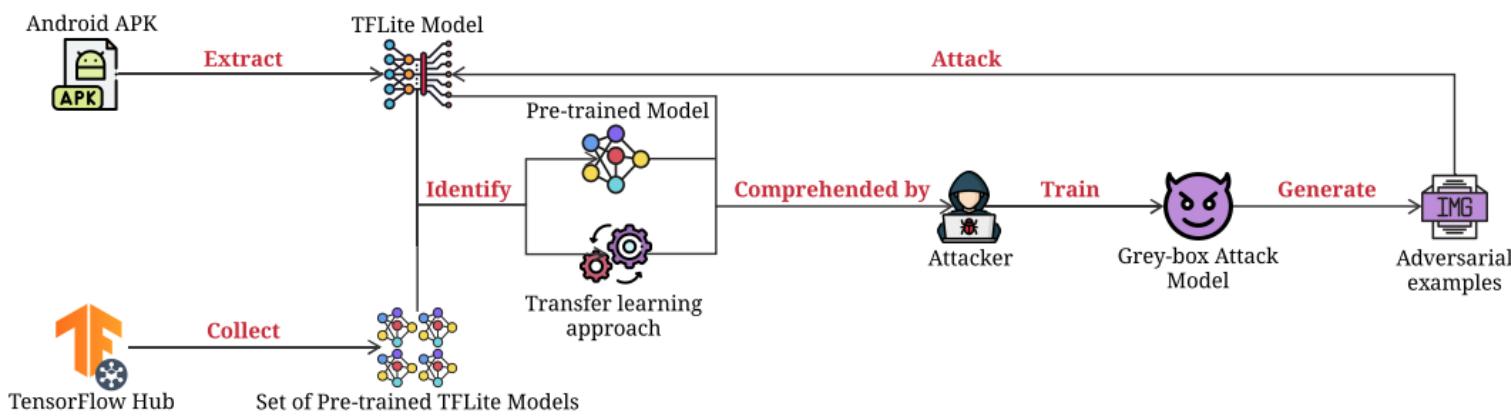
- Model could be deployed on device
- Extracting from
  - firmware/flash storage
  - APK, executables
  - Disk
  - Transmitted over Network
- Analyze the serialized model
- Leverage recon to understand internals (like framework, algorithm, architecture, etc) of the model
- Get predictions from model

# Offline Model Stealing Attacks

- Model could be deployed on device
- Extracting from
  - firmware/flash storage
  - APK, executables
  - Disk
  - Transmitted over Network
- Analyse the serialized model
- Leverage recon to understand internals (like framework, algorithm, architecture, etc) of the model
- Get predictions from model

# Offline Model Stealing Attacks

## Demo 9 : Offline Model Stealing Attack [H/W]



# Offline Model Stealing Attacks

**Demo 9 : Offline Model Stealing Attack  
[H/W]**

# Online Black Box model stealing

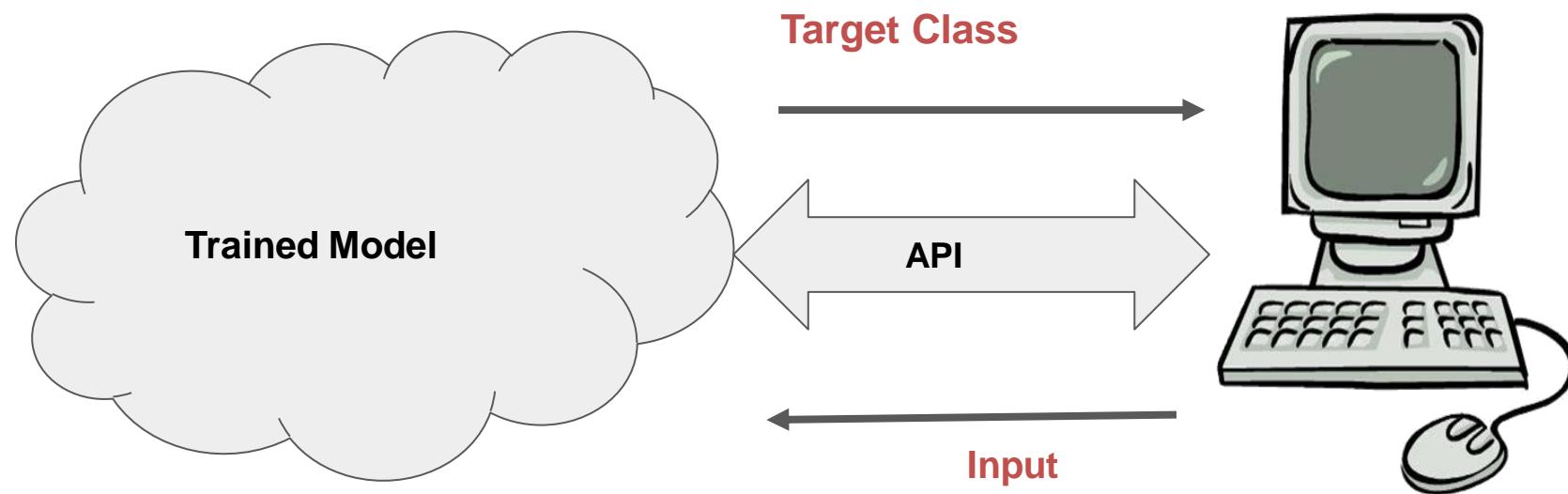
MLaaS



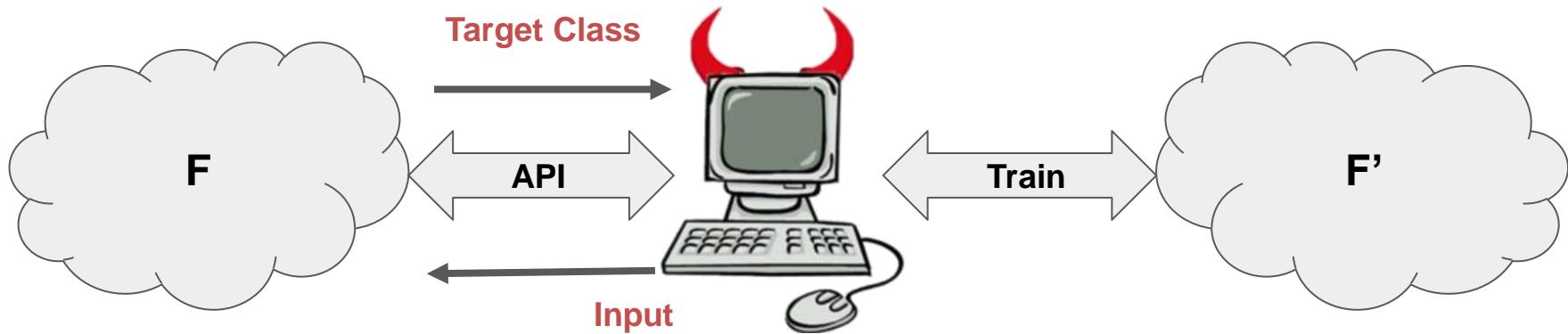
Google Cloud Platform



# MLaaS



# Model Stealing attack methodology

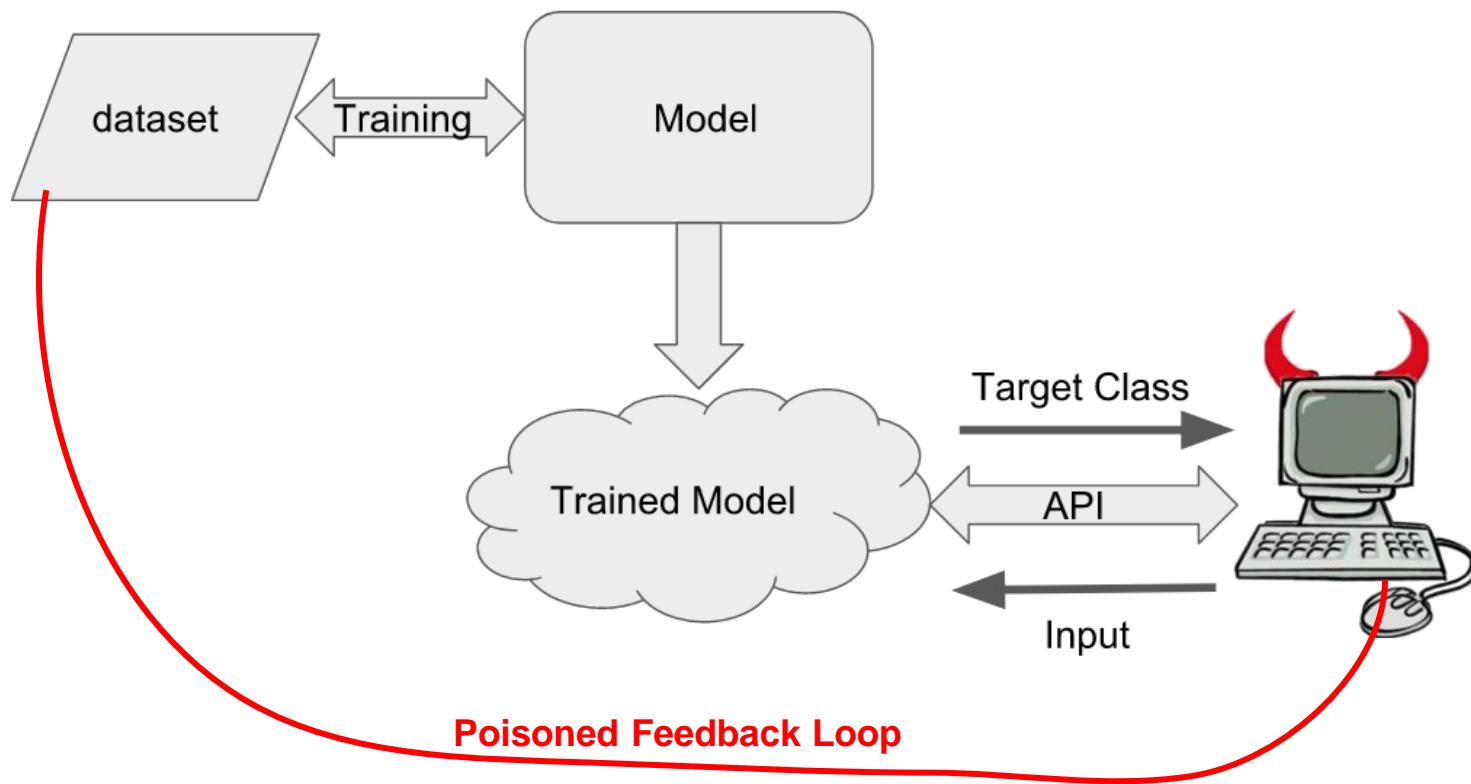


# Model Stealing attack methodology

- Collect data
- Query prediction APIs to generate labels from target Model
- Use generated Labels and collected data to train a duplicate model
- Optimization algorithms like Gradient descent are used to train the

# Model skewing Attacks

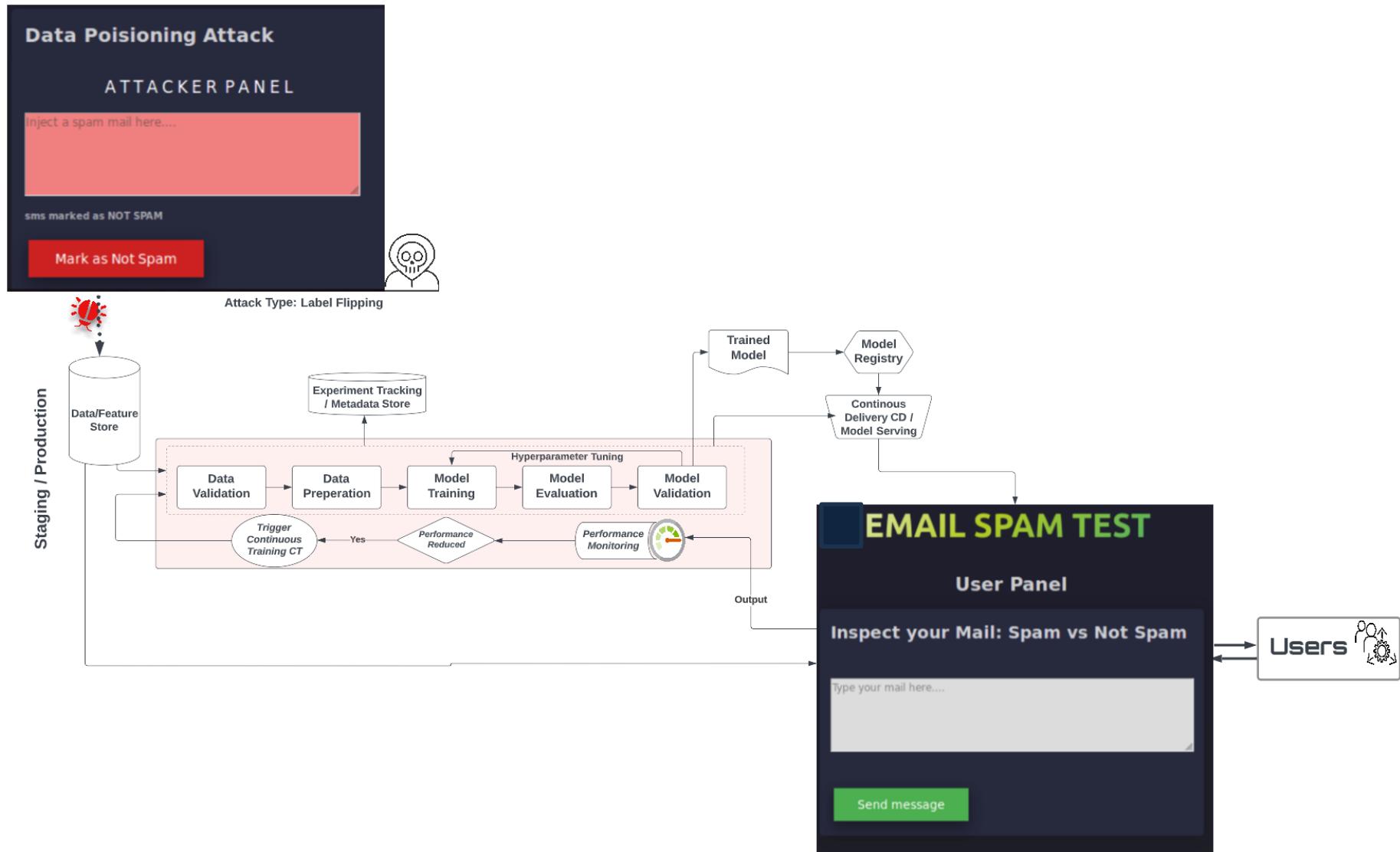
# Model skewing Attacks



# **Demo 10: Skewing Spam Filters**

# **Data Poisioning Attack**

# Data Poisioning Attack

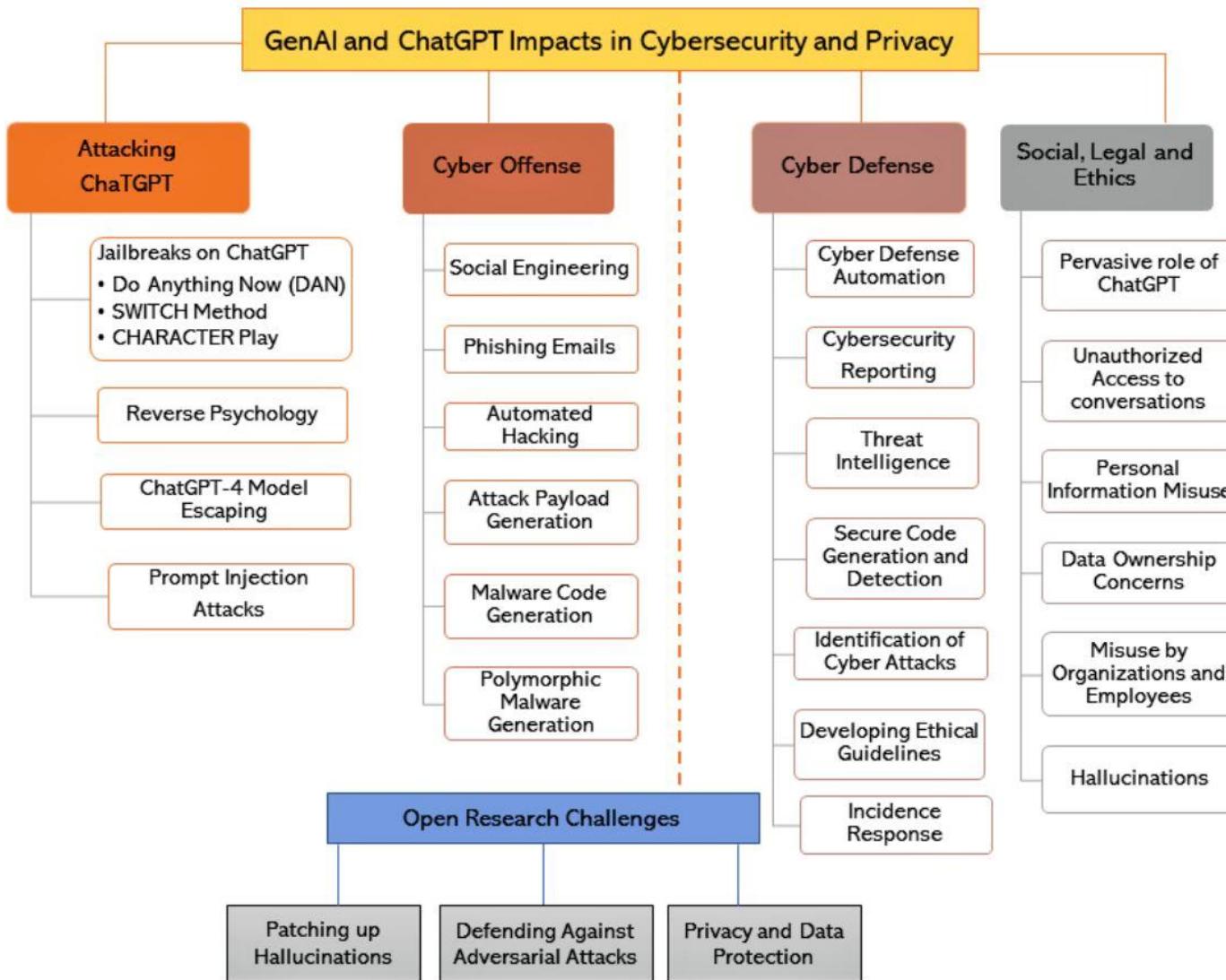


# Data Poisioning Attack

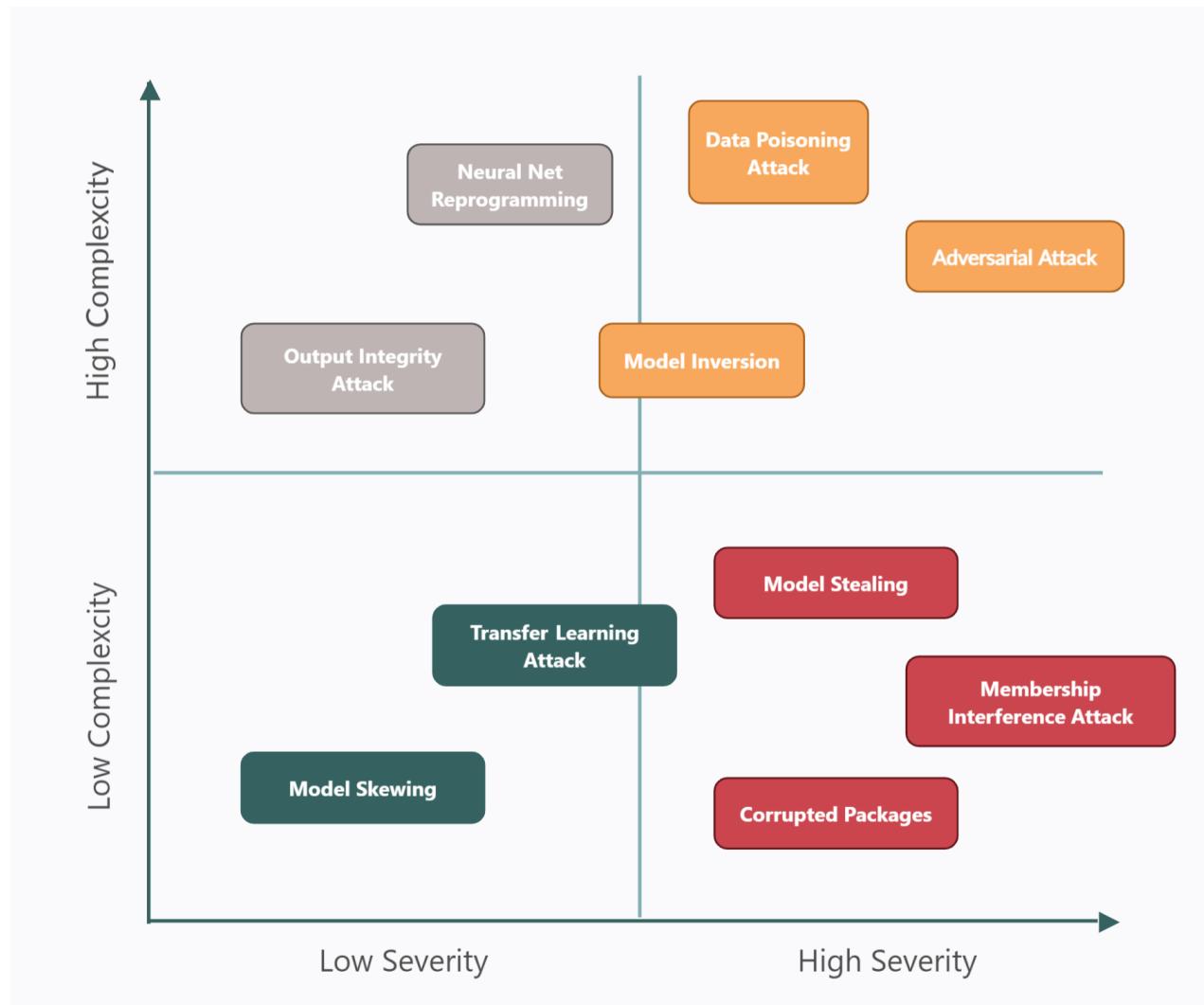
**Demo X:label flipping data poisoning attack**

[Play Video: /seasides/Data Poisoning.mp4](#)

# Generative AI



# ML Top 10: Risk Assessment Matrix



# OWASP “Machine Learning Security top 10”

<https://mltop10.info/>

## OWASP Machine Learning Security Top Ten

[Main](#) [Charter](#) [Related](#) [Glossary](#)

owasp incubator License CC BY-SA 4.0

### Important Information

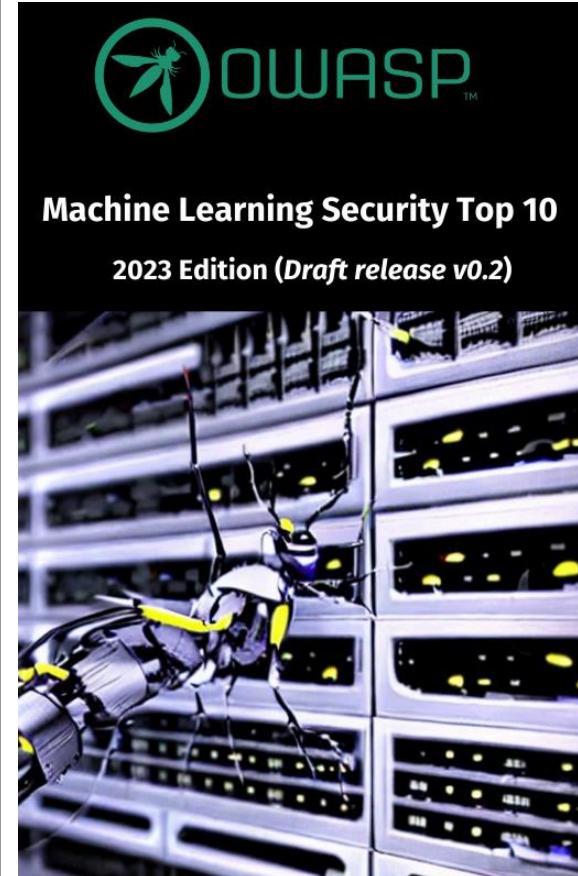
The current version of this work is in draft and is being modified frequently. Please refer to the [project wiki](#) for information on how to contribute and project release timelines.

## Overview

Welcome to the repository for the OWASP Machine Learning Security Top 10 project! The primary aim of the OWASP Machine Learning Security Top 10 project is to deliver an overview of the top 10 security issues of machine learning systems. More information on the project scope and target audience is available in our [project working group charter](#)

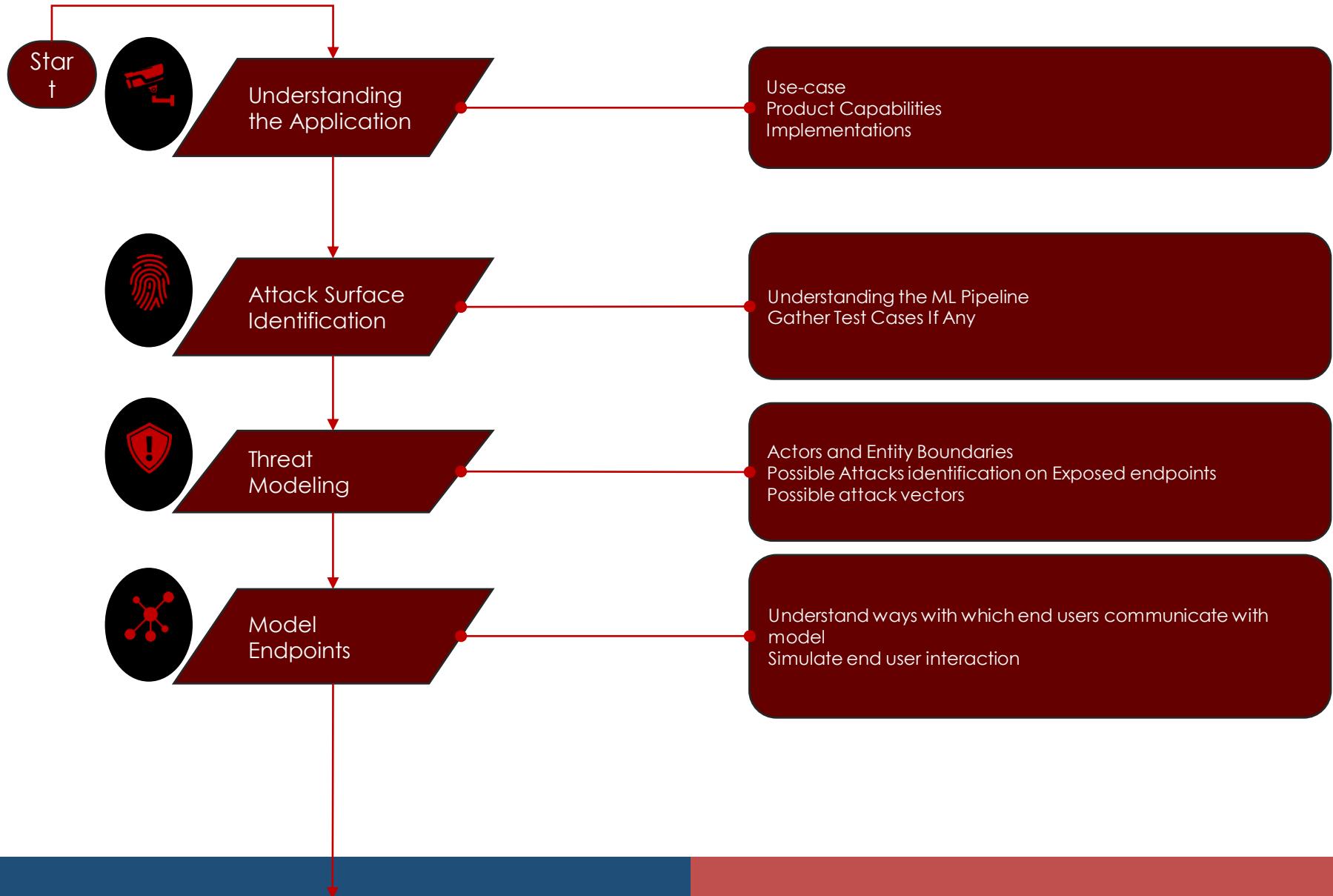
## Top 10 Machine Learning Security Risks

- [ML01:2023 Input Manipulation Attack](#)
- [ML02:2023 Data Poisoning Attack](#)
- [ML03:2023 Model Inversion Attack](#)
- [ML04:2023 Membership Inference Attack](#)
- [ML05:2023 Model Stealing](#)
- [ML06:2023 Corrupted Packages](#)
- [ML07:2023 Transfer Learning Attack](#)
- [ML08:2023 Model Skewing](#)
- [ML09:2023 Output Integrity Attack](#)
- [ML10:2023 Model Poisoning](#)

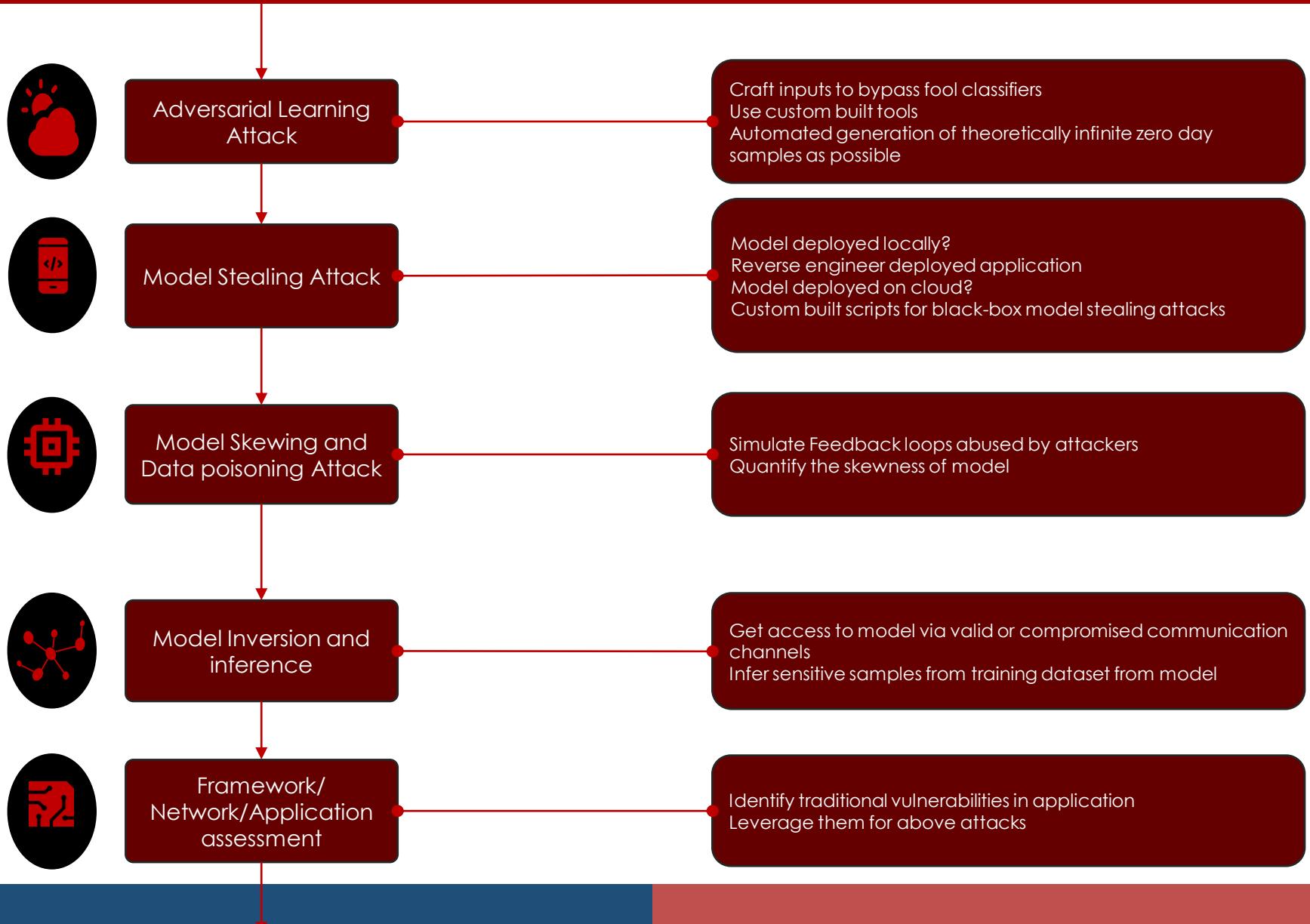


{ref} <https://owasp.org/www-project-machine-learning-security-top-10/>

# ML Security Assessment Process



# ML Security Assessment Process



# ML Security Assessment Process



# CTF

An organisation has implemented an authentication system "mlAuth" using machine learning, which is 99.9% accurate. Every employee has a profile(represented by a string on 784 hex values). mlAuth is trained using these profiles to predict the probability of authenticity for an employee. System grants access only if the predicted probability is higher than 0.99. Hence, your aim is to generate a fake profile that will trick the 99.9% accurate mlAuth in granting you access.

[Hint]

```
# You can make use of dumped machine learning model to conduct your targated attack. Are you smart enough to fool the "intelligent" mlAuth?
```

# Next Steps

