# GDALR: An efficient model duplication attack on black-box Machine Learning models

BY

REWANTH COOL (SECURITY CONSULTANT) , NIKHIL JOSHI (SECURITY RESEARCHER)

# About us

**REWANTH COOL**

**NIKHIL JOSHI**

Security Consultant

Security Researcher

ML Enthusiast

Build and Break Deep Learning systems

Full stack developer

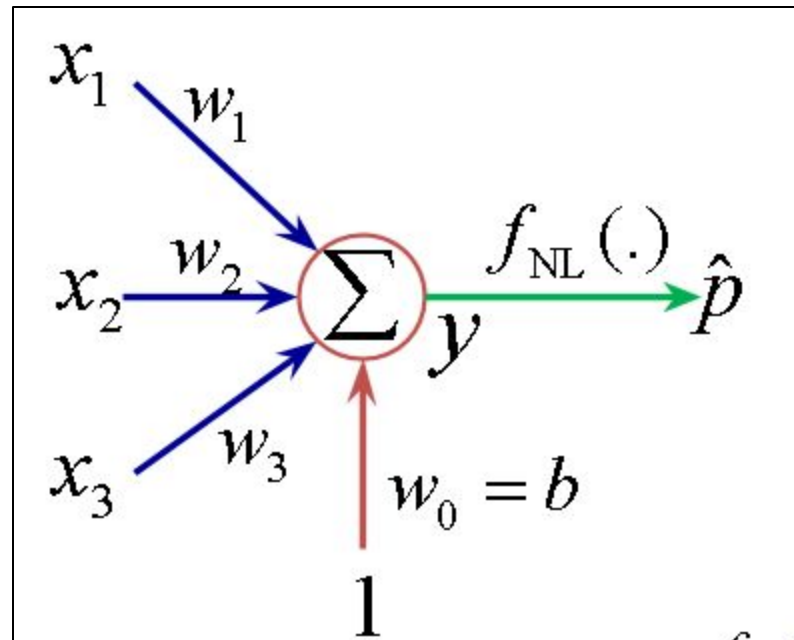Speaker at HITB, CRESTCON, BSIDES

- About Payatu
  - A boutique security testing company specializing in IoT, Mobile, Cloud – https://payatu.com
  - In-house Fuzz testing Infrastructure
  - Mobile/Windows kernel/IoT exploitation training – Blackhat, Brucon, Hack In Paris, HITB and Corporate trainings
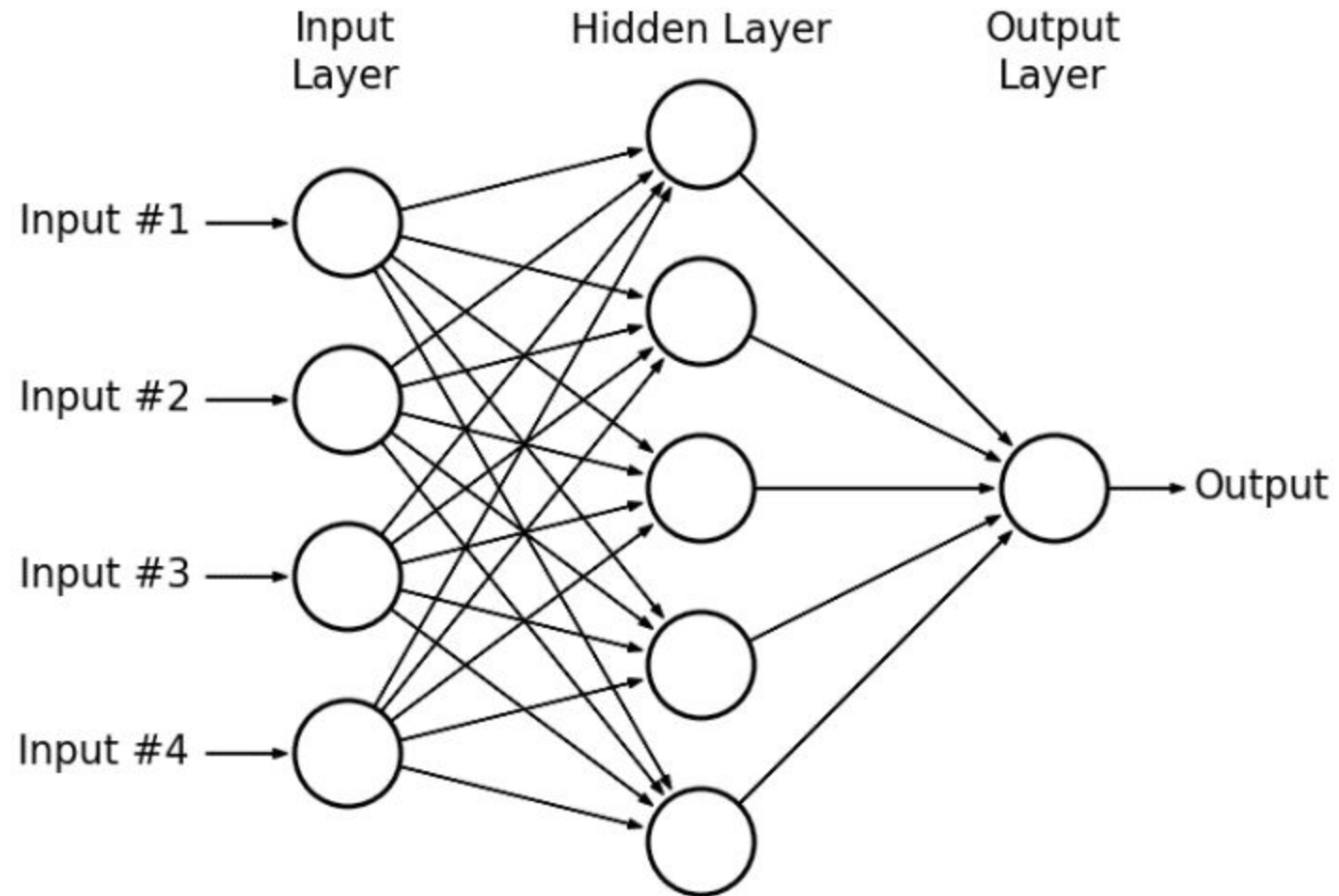
# Agenda

- End-to-end Machine Learning pipeline

- Model stealing/duplication techniques

- Abusing APIs to steal models deployed on cloud

- Present attack methodology

- Inefficiencies with present attack methodology

- Scope for Attack optimization

- Proposed approach (GDALR)

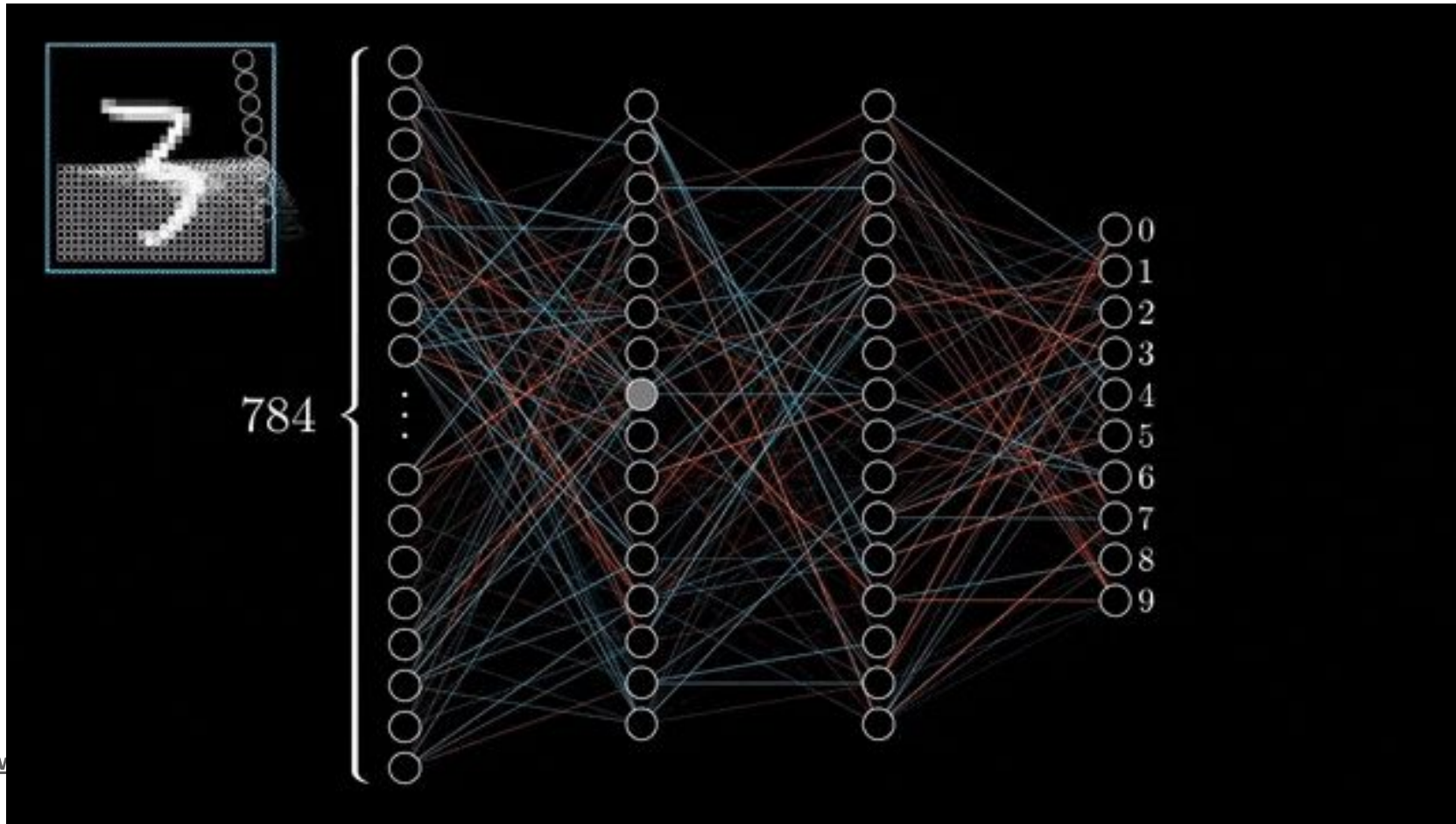- Results and conclusion

-

# PERCEPTRON

# MULTI LAYER PERCEPTRON (MLP)



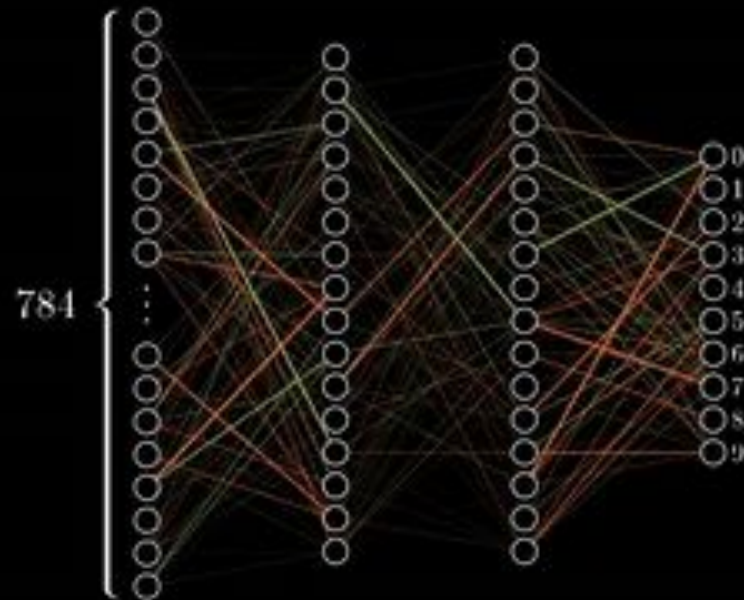Ref: https://www.researchgate.net/figure/A-hypothetical-example-of-Multilayer-Perceptron-Network_fig4_303875065
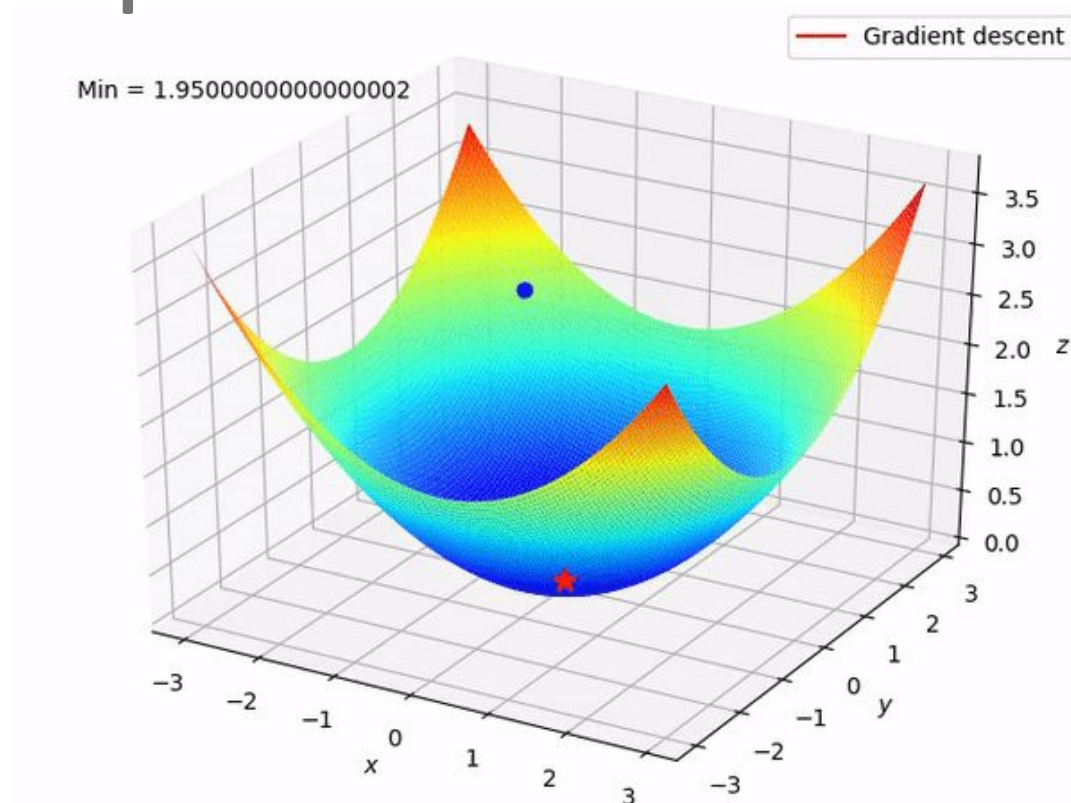
# MLP



Ref: https://www

# MLP



Ref: http

# Optimization

Gradient descent

Min = 1.9500000000000002

$$W_{i+1} = W_i + \alpha \frac{\partial}{\partial w} C$$

Ref: http://www.xpertup.com/2018/05/11/loss-functions-and-optimization-algorithms/

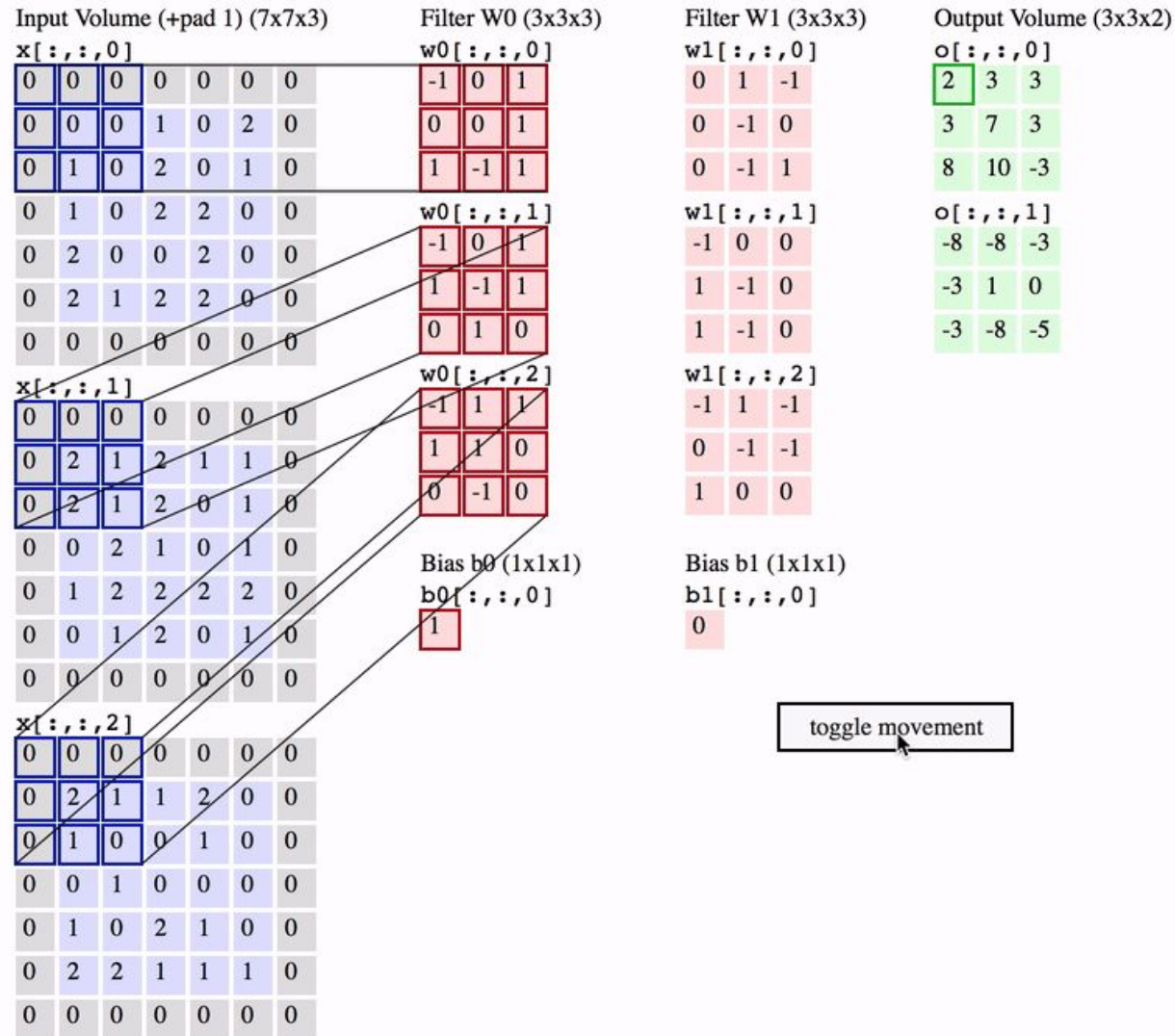# MLP


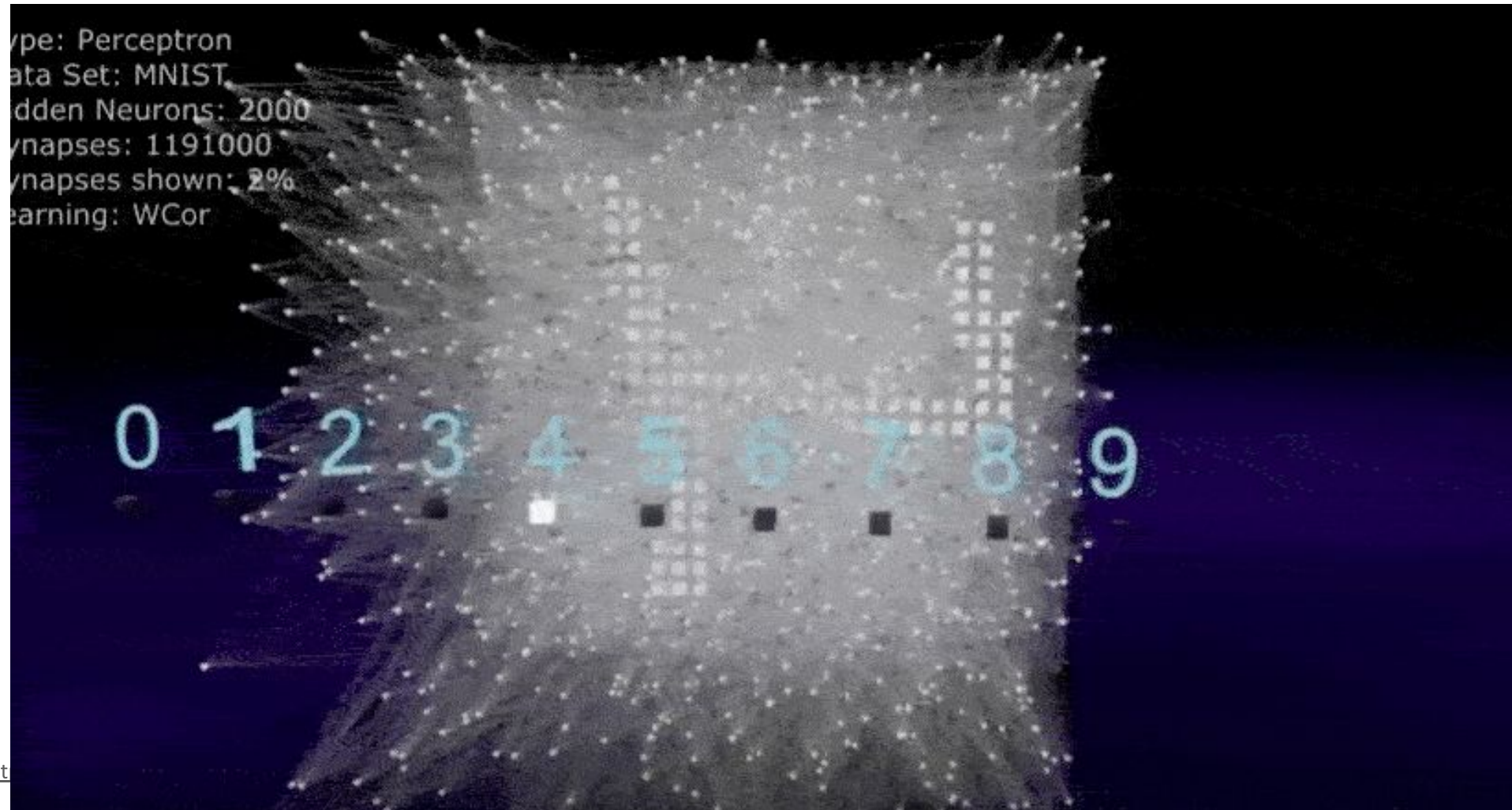
Ref: h...

# Convolutional Neural Networks

Ref: http://cs231n.github.io/convolutional-networks/

# Convolutional Neural Networks



Ref: htt

# Model stealing/duplication techniques

- Offline attacks

- Online attacks

# Offline attacks

Steps -

1. Reverse engineer the executable to find hidden gems
2. Locate the trained model stored on device
3. Analyse the serialized model
4. Own the model

# Offline attacks

```java
public Model loadModel(String modelFolder) {
    List<String> categories = loadCategories(modelFolder + "/categories.txt");
    if (categories == null) {
        Log.e(TAG, "Failed to load categories: " + modelFolder + "/categories.txt");
        return null;
    }
    ByteBuffer enginePtr = loadModelFromAssets(modelFolder + "/model.net", modelFolder + "/stat.t7");
    if (enginePtr != null) {
        return new Model(enginePtr, categories, 224);
    }
    Log.e(TAG, "Failed to load model");
    return null;
}
```

# Offline attacks

# Offline attacks

```
# Loading model
from torch.utils.serialization import load_lua
model = load_lua(model_path)
stat = load_lua(model_path[:-9]+'stat.t7')
model_op = predict(IMAGE_PATH)
```
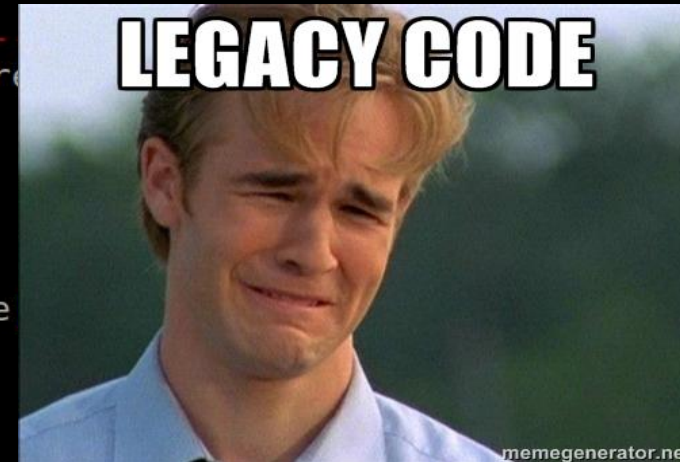
# Offline attacks

# Offline attacks

# Offline attacks



```
ithout bias
                |       (4): nn.SpatialBatchNormaliz
                |       (5): nn.SpatialDropout
                |     }
                |`-> (1): nn.Identity
                 +. -> output
           }
       (1): nn.CAddTable
       (2): nn.ReLU
     }
   }
 (8): nn.Identity
 (9): nn.SpatialAveragePooling(14x14, 1, 1)
 (10): nn.View(128)
 (11): nn.Linear(128 -> 696)
 (12): nn.SoftMax
}
```

```
            |`-> (1): nn.Identity
             +. -> output
         }
      (1): nn.CAddTable
      (2): nn.ReLU
     }
   }
 (8): nn.Identity
 (9): nn.SpatialAveragePooling(14x14, 1, 1)
 (10): nn.View(1, 128)
 (11): nn.Linear(128 -> 696)
 (12): nn.SoftMax
}
```

```
torch.legacy.nn.View(1,128)
```

# Offline attacks

# MLaaS Service Providers

# Azure ML business model



## What is Azure Machine Learning

**Data**
- Blobs and Tables
- Hadoop (HDInsight)
- Relational DB (Azure SQL DB)

ML STUDIO

Integrated development environment for Machine Learning

API

Model is now a web service that is callable

Monetize the API through our marketplace

**Clients**

Ref: **https:/**

# Online attacks

# Online attacks



**Input:** [f1, f2, f3, ..., fn]

**Internal F output:** P(class1), P(class2), P(class3), …, P(classN)

**Cloud API output:** max(P(class1), P(class2), P(class3), …, P(classN))

# Present attack methodology



Ref: https://www.usenix.org/conference/usenixsecurity16/technical-sessions/presentation/tramer
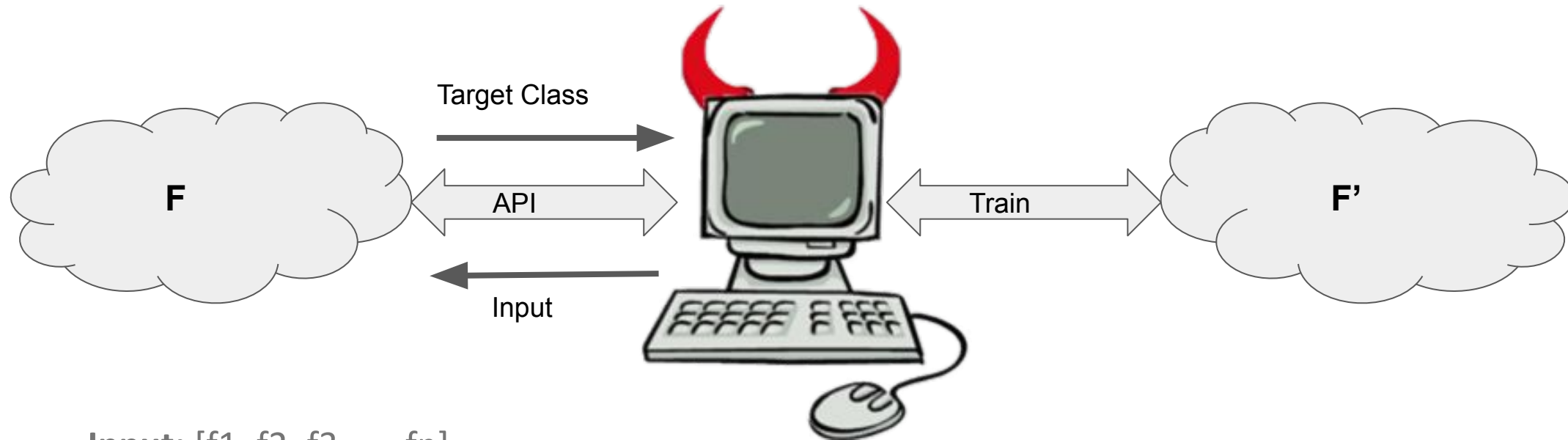
# Present attack methodology

In traditional approach, attackers train their local models based on Cloud API output

| Input | Cloud API output | Class (A/B/C) |
|---|---|---|
| x11, x12, x13, x14 | 0, 0, 1 | C |
| x21, x22, x23, x24 | 0, 1, 0 | B |
| x31, x32, x33, x34 | 0, 0, 1 | C |
| x41, x42, x43, x44 | 1, 0, 0 | A |
| x51, x52, x53, x54 | 1, 0, 0 | A |

# Inefficiencies with present attack methodology

Assumptions made by traditional/present attack methodology

**Input -**
[1, 2, 3, 4]

**Actual Output -**
[0.3, 0.2, 0.5]

**Output by Cloud API -**
[0, 0, 1]

**Assumption -**
[0, 0, 1]      [0.3, 0.2, 0.5]
            $\approx$

# Inefficiencies with present attack methodology

| Input | Cloud API output | Actual Output | Unconventional probability loss |
|---|---|---|---|
| x11, x12, x13, x14 | 0, 0, 1 | 0.2, 0.3, 0.5 | 0.2+0.3 |
| x21, x22, x23, x24 | 0, 1, 0 | 0.01, 0.9, 0.09 | 0.01+0.09 |
| x31, x32, x33, x34 | 0, 0, 1 | 0.1, 0.4, 0.5 | 0.1+0.4 |
| x41, x42, x43, x44 | 1, 0, 0 | 0.38, 0.32, 0.3 | 0.32+0.3 |
| x51, x52, x53, x54 | 1, 0, 0 | 0.45, 0.3, 0.25 | 0.3+0.25 |

# Scope for Attack optimization

1. Reconsider the way to analyze labels

      Having access to all the probability values will definitely help us to clone models in an efficient way

2. Learning parameters in hyperspace
    * To Duplicate the target model we need to learn the boundaries that the target model has learnt
    * Considering probability of predicted class as 1 and others to be 0 will cause unwanted loss and increase the gradient
    * Increased gradients cause the optimizer to change weights abruptly

$$W_{i+1} = W_i + \alpha \frac{\partial}{\partial w} C \uparrow$$

# Proposed approach (GDALR)
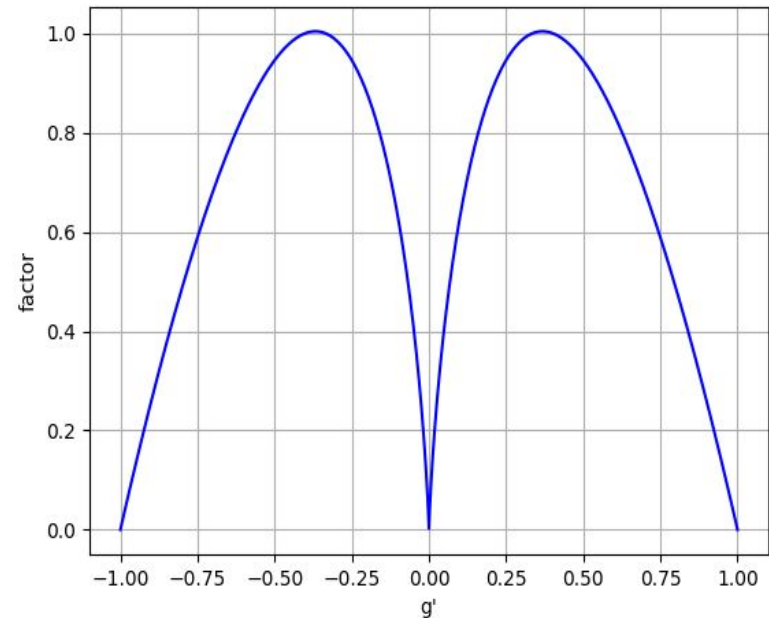
GDALR: Gradient Driven Adaptive Learning Rate

$$W_{i+1} = W_i + \alpha \frac{\partial}{\partial w} C$$

# Mathematical modification to current attack methodology

$$g_i' = \tanh(g_i) \qquad (7)$$

$$fact_i = \text{abs}\left(g_i' 2\pi \log_{10}\left(\text{abs}\left(g_i'\right)\right)\right) \qquad (8)$$
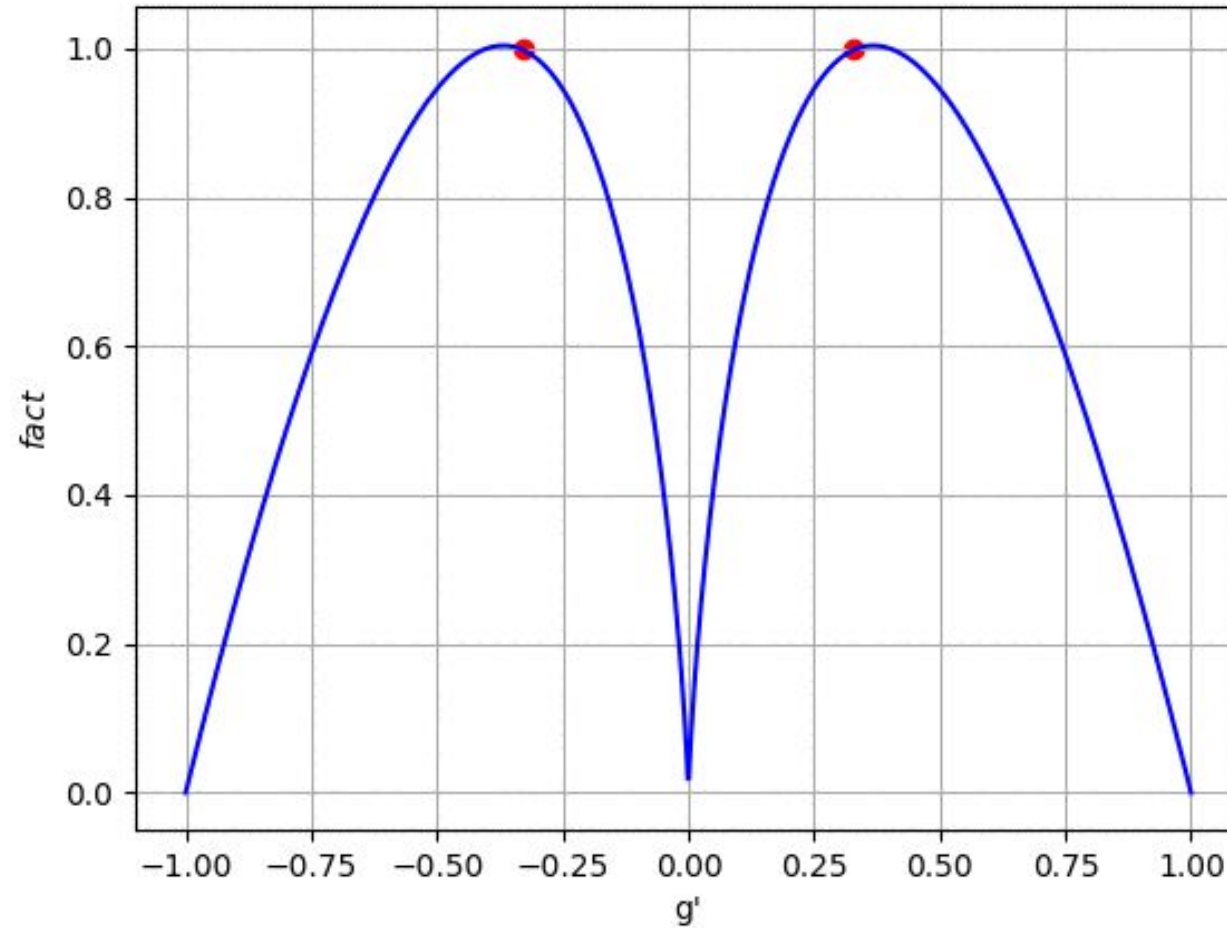
$$l_i' = l_i \cdot fact_i \qquad (9)$$

# GDALR in Action


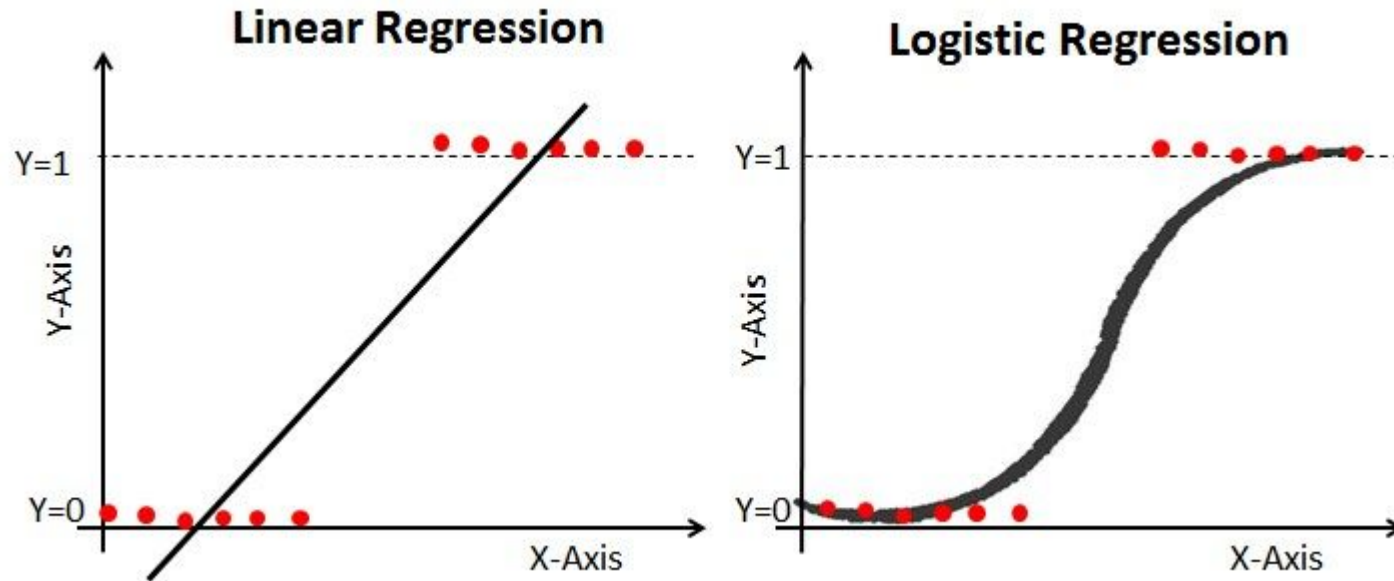
$$g' = 0.33 \qquad fact_i = 1.00$$

# Experimental setup

GDALR has been tested on multiple classifiers -

- LOGISTIC REGRESSION

- MULTI LAYER PERCEPTRON
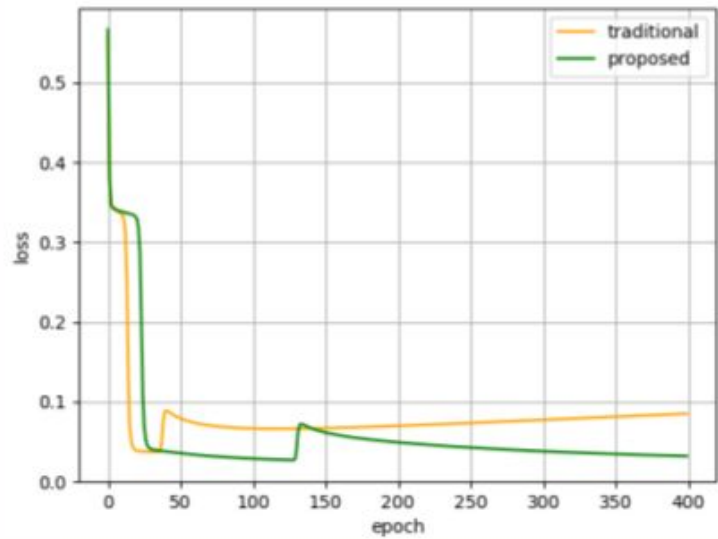
- CONVOLUTIONAL NEURAL NETs

# LOGISTIC REGRESSION



Ref: http://www.datacamp.com/

$$y_{linear} = wx + b$$

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

# LOGISTIC REGRESSION



| $l = 0.01$ | $l = 0.05$ |
|---|---|
| $T_{Loss} = 0.0849$ | $T_{Loss} = 0.1233$ |
| $P_{Loss} = 0.0317$ | $P_{Loss} = 0.0342$ |

# MULTI LAYER PERCEPTRON

| $l = 10^{-3}$ | $l = 10^{-5}$ |
|---|---|
|  |  |
| $T_{Loss} = 0.0014$ | $T_{Loss} = 0.0219$ |
| $P_{Loss} = 5.444 \times 10^{-5}$ | $P_{Loss} = 0.0007$ |

# CNN



| $l = 10^{-4}$ | $l = 10^{-5}$ |
|---|---|
| $T_{Loss} = 0.0011$ | $T_{Loss} = 0.0073$ |
| $P_{Loss} = 3.993 \times 10^{-5}$ | $P_{Loss} = 4.184 \times 10^{-5}$ |

# Thanks!

- Q & A
- Reach us at

Email - [rewanth|nikhilj]@payatu.com

Twitter - @Rewanth_Cool | @nikhilj_73