# Letter Recognition

# 1.
# Introduction

# Members

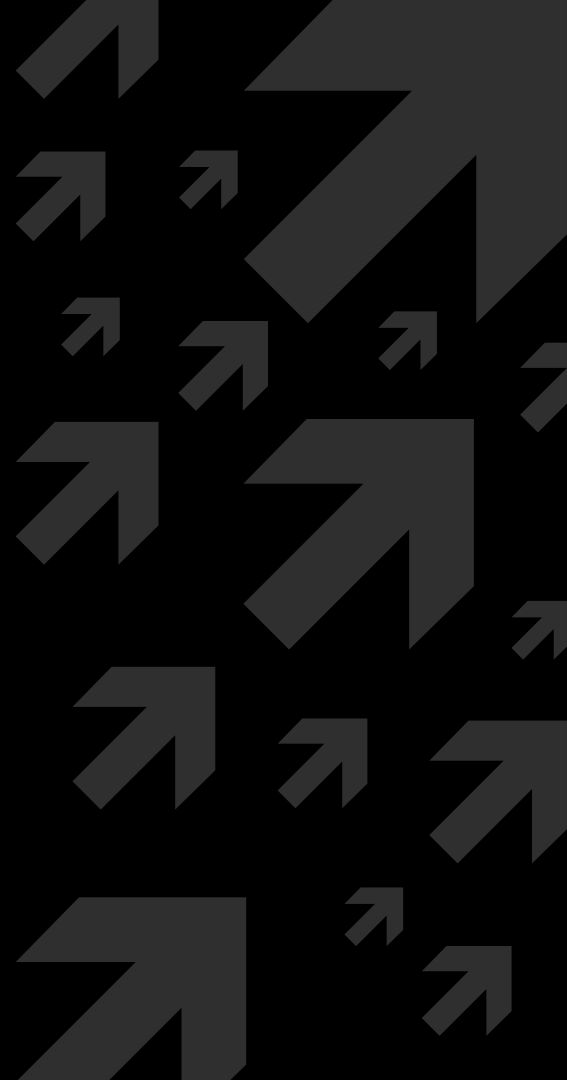| MSSV | Name | Work |
|------|------|------|
| HE150402 | Tô Văn Đức | Model + Application |
| HE150303 | Trần Đức Tuấn | Dataset + Model |
| HE150258 | Nguyễn Trung Nghĩa | Dataset + Application |

# Problem

➜ Letter Recognition:
  ⇥ Recognize letters from 28x28 images
  ⇥ Each image contains only one letter

# 2.
# Dataset

Data preparation and exploration

# **Dataset**

➔ Source:

https://www.nist.gov/itl/products-and-services/emnist-dataset

## Dataset Summary

There are six different splits provided in this dataset. A short summary of the dataset is provided below:

- EMNIST ByClass: 814,255 characters. 62 unbalanced classes.
- EMNIST ByMerge: 814,255 characters. 47 unbalanced classes.
- EMNIST Balanced: 131,600 characters. 47 balanced classes.
- EMNIST Letters: 145,600 characters. 26 balanced classes.
- EMNIST Digits: 280,000 characters. 10 balanced classes.
- EMNIST MNIST: 70,000 characters. 10 balanced classes.

The full complement of the NIST Special Database 19 is available in the ByClass and ByMerge splits. The EMNIST Balanced dataset contains a set of characters with an equal number of samples per class. The EMNIST Letters dataset merges a balanced set of the uppercase and lowercase letters into a single 26-class task. The EMNIST Digits and EMNIST MNIST dataset provide balanced handwritten digit datasets directly compatible with the original MNIST dataset.

Please refer to the EMNIST paper [PDF , BIB] for further details of the dataset structure.

# Dataset

➔ Sample:

# Dataset

➜ Overview:
- → Total: **145,600** images
- → Number of classes: **26**
- → Training Set:
  - ■ **124,800** images
  - ■ **4,800** images / class
- → Test Set:
  - ■ **20,800** images
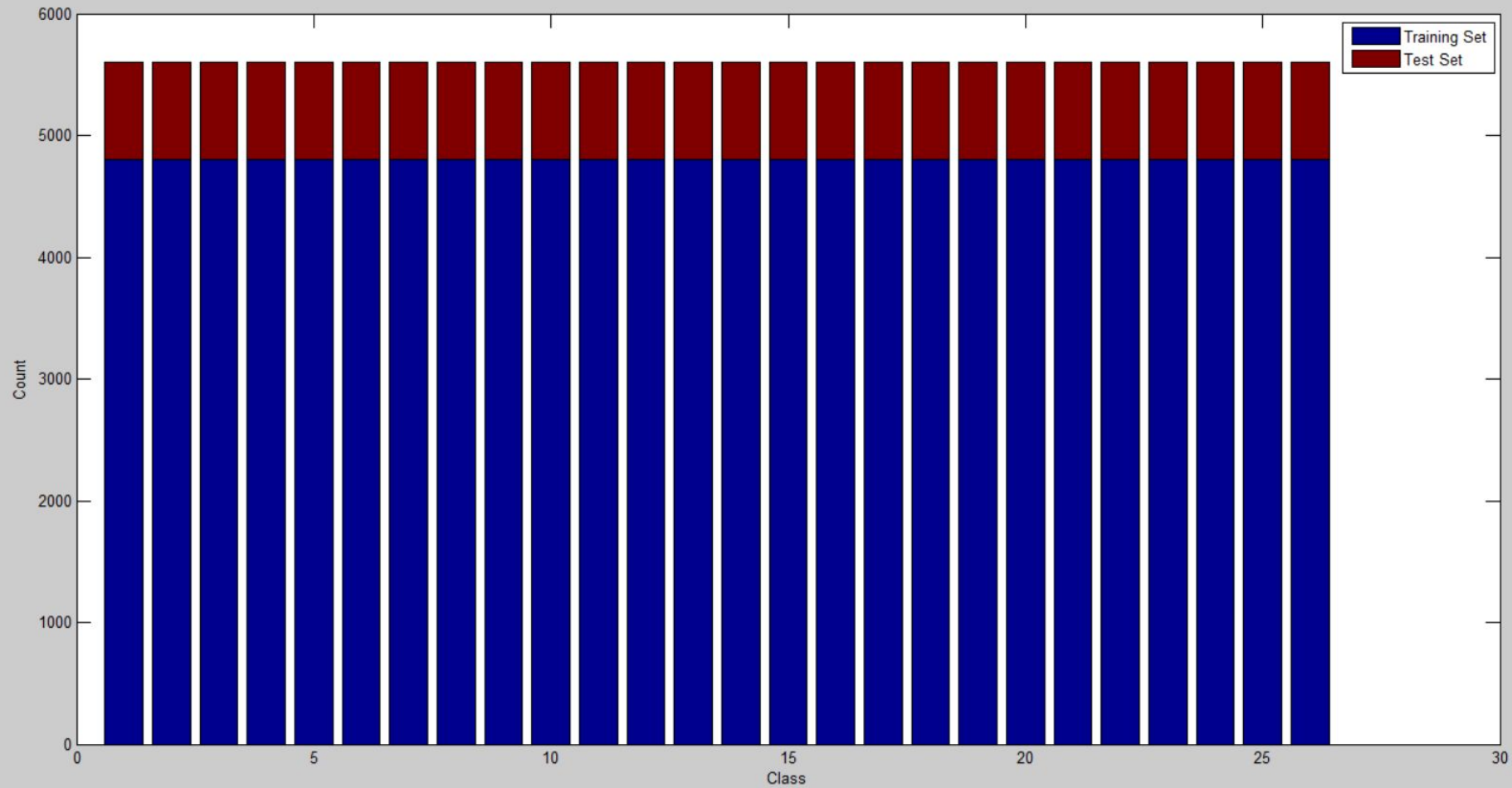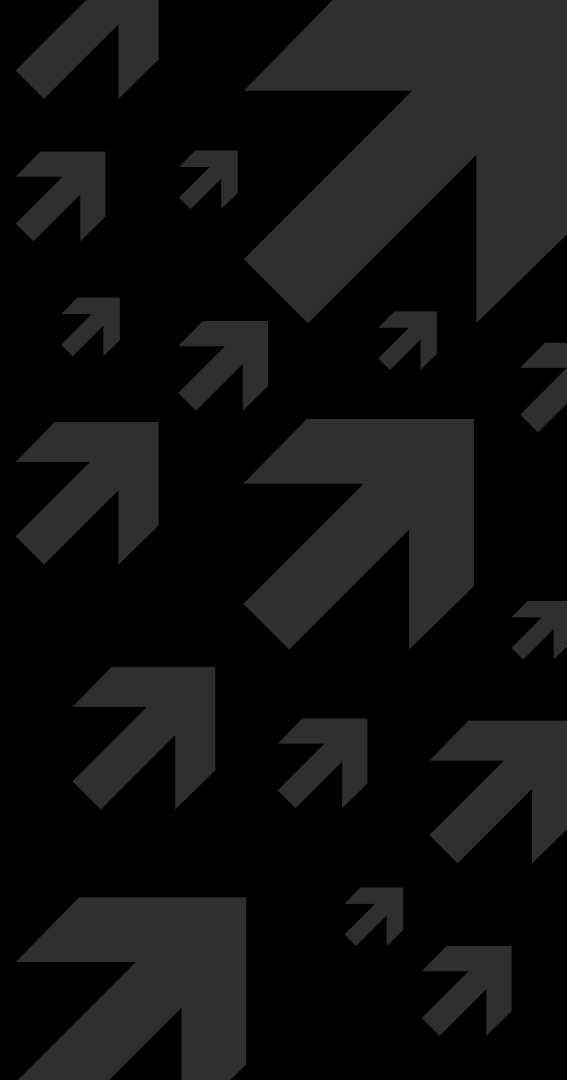  - ■ **800** images / class

## Dataset Summary

There are six different splits provided in this dataset. A short summary of the dataset is provided below:

- EMNIST ByClass: 814,255 characters. 62 unbalanced classes.
- EMNIST ByMerge: 814,255 characters. 47 unbalanced classes.
- EMNIST Balanced: 131,600 characters. 47 balanced classes.
- EMNIST Letters: 145,600 characters. 26 balanced classes.
- EMNIST Digits: 280,000 characters. 10 balanced classes.
- EMNIST MNIST: 70,000 characters. 10 balanced classes.

The full complement of the NIST Special Database 19 is available in the ByClass and ByMerge splits. The EMNIST Balanced dataset contains a set of characters with an equal number of samples per class. The EMNIST Letters dataset merges a balanced set of the uppercase and lowercase letters into a single 26-class task. The EMNIST Digits and EMNIST MNIST dataset provide balanced handwritten digit datasets directly compatible with the original MNIST dataset.

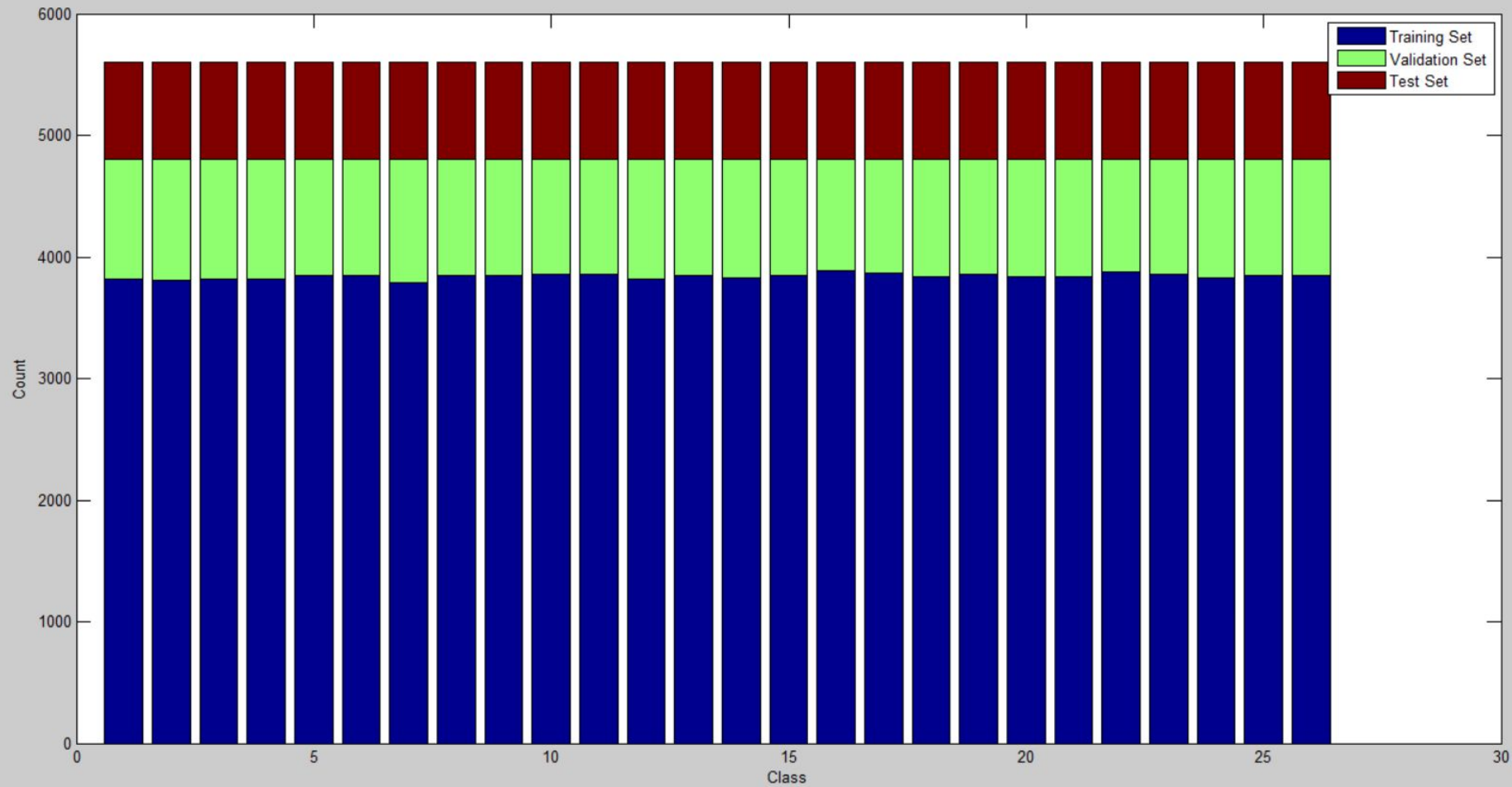Please refer to the EMNIST paper [PDF⧉ , BIB]for further details of the dataset structure.

# 3.
# Preparation

# Data preparation

➔ Split data into training set and validation set:

→ Training set:  80%

→ Validation set:  20%

➔ Normalize data:

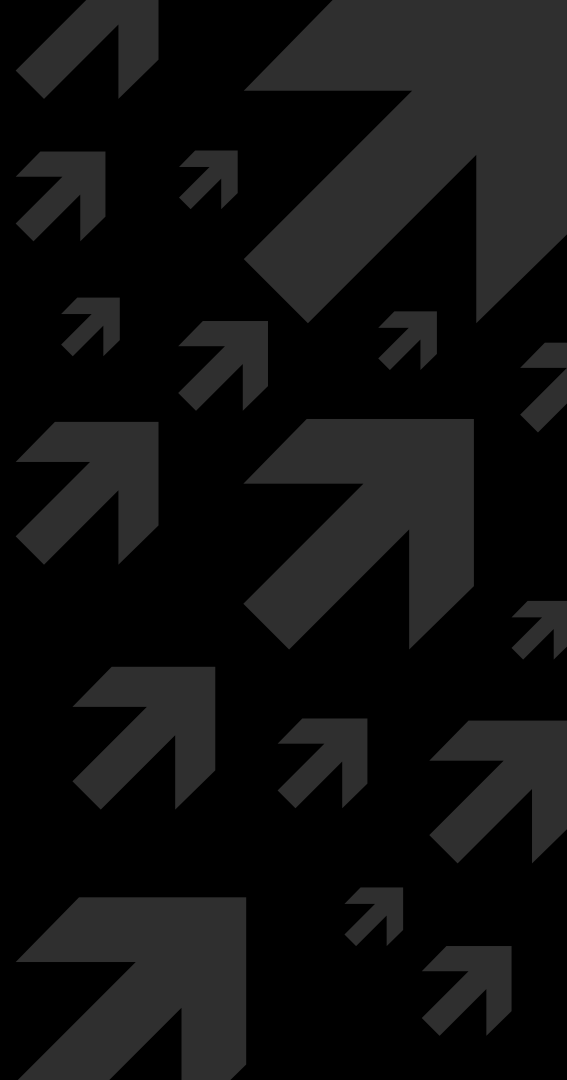→ Scale each value to the range of [0, 1]

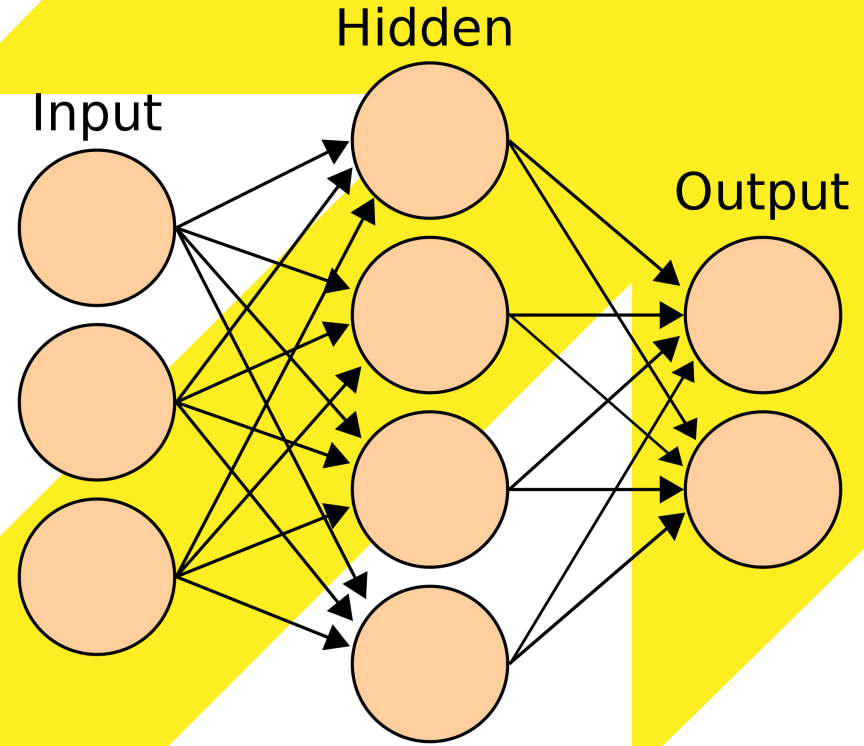# Label preparation

➜ One-hot encoding the labels

| Label |
|-------|
| 1 (A) |
| 2 (B) |
| ... |
| 26 (Z) |

| 1 (A) | 2 (B) | ... | 26 (Z) |
|-------|-------|-----|--------|
| 1 | 0 | ... | 0 |
| 0 | 1 | ... | 0 |
| ... | ... | ... | ... |
| 0 | 0 | ... | 1 |

# 4.
# Progress

# Our Approach

➔ Model:
  → Neural Network

➔ Evaluation Metrics:
  → Training Set Accuracy
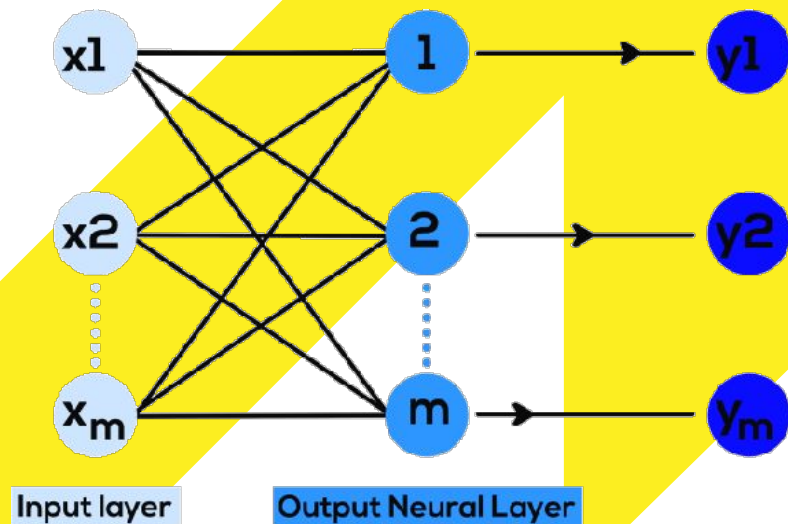  → Validation Set Accuracy
  → Test Set Accuracy

Input

Hidden

Output

# Model parameters

➔ We try to find the best combination of:

   → Activation function

   → Number of nodes

   → Number of hidden layers

# Single Layer Perceptron

➔ Architecture:

→ 1 input layer (784 nodes)

→ 0 hidden layer

→ 1 output layer (26 nodes)



Input layer    Output Neural Layer

**17**

# Weight initialization

➔ **Random initialization:**

→ Initialize each weight to the same value will cause all nodes in neural network to learn the same feature while training => worsen model's performance

→ **Symmetry breaking:** initialize each weight randomly so each node will learn different features

# Weight initialization

➔ **Xavier initialization:** Initialize each weight to a random value between [-ϵ, ϵ]:

$$\epsilon = \frac{\sqrt{6}}{\sqrt{L_{in} + L_{out}}}$$

$$\Theta^{(l)} = 2\epsilon * rand(L_{out}, L_{in} + 1) - \epsilon$$

# Single Layer Perceptron

➔ Output layer:
  → Sigmoid + BCE Loss
  → Softmax + CCE Loss



Input layer     Output Neural Layer

# Sigmoid + BCE Loss

➔ Formula:

$$\hat{y}_i = \frac{1}{1 + e^{-z}}$$

➔ Binary Cross Entropy Loss:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{K} [-y_k^{(i)} \log \hat{y}_k^{(i)} - (1 - y_k^{(i)}) \log(1 - \hat{y}_k^{(i)})]$$

# Sigmoid + BCE Loss

➔ Derivative:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$
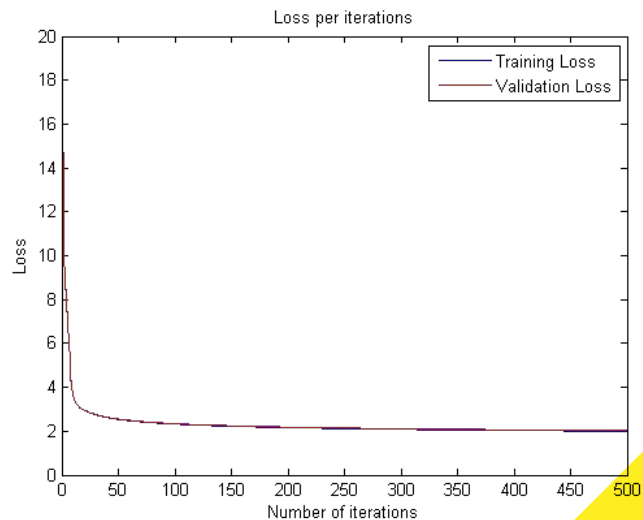
# Softmax + CCE Loss

→ Formula:

$$\hat{y}_i = \frac{e^{z_i}}{\sum_{k=1}^{K} e^{z_k}}$$

→ Categorical Cross Entropy Loss:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log \hat{y}_k^{(i)}$$

# Softmax + CCE Loss

➔ Derivative:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

# Let's train

➔ Learning rate: **0.5**

➔ Number of iterations: **500**

➔ Save the weights that minimize validation loss

➔ **Early stopping technique**: stop training if validation loss does not decrease



Error

Under-fitting

Over-fitting

Validation set

"sweet spot"

Training set

Number of iterations

# Single Layer Perceptron

| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| Perceptron_sigmoid | 20,410 | 68.99 % | 68.15 % | 68.66 % |
| Perceptron_softmax | 20,410 | 70.15 % | 69.76 % | 69.88 % |

# Single Layer Perceptron



**Perceptron_sigmoid**



**Perceptron_softmax**

# Perceptron_softmax

Confusion Matrix

# Neural Network

➔ Architecture:
  - 1 input layer (784 nodes)
  - 1 hidden layer
  - 1 output layer (26 nodes)

Input

Hidden

Output

# Neural Network

➔ Forward propagation:

⇢ $a^{(1)} = x$

⇢ $z^{(l+1)} = a^{(l)} * \Theta^{(l)}$ $\qquad (l : 1 \rightarrow L)$

⇢ $a^{(l+1)} = g(z^{(l+1)})$ $\qquad (l : 1 \rightarrow L)$

⇢ $J(\theta) = -\dfrac{1}{m}\displaystyle\sum_{i=1}^{m}\sum_{k=1}^{K} y_k^{(i)} \log \hat{y}_k^{(i)}$ $\quad (\hat{y} = a^{(L+1)})$

# Neural Network

➔ Backward propagation:

→ $\delta^{(L+1)} = \dfrac{\partial J}{\partial z^{(L+1)}} = a^{(L+1)} - y$

→ $\dfrac{\partial J}{\partial a^{(l)}} = \dfrac{\partial J}{\partial z^{(l+1)}} * \dfrac{\partial z^{(l+1)}}{\partial a^{(l)}} = \delta^{(l+1)} * \Theta^{(l)}$ $\qquad (l : L \rightarrow 1)$
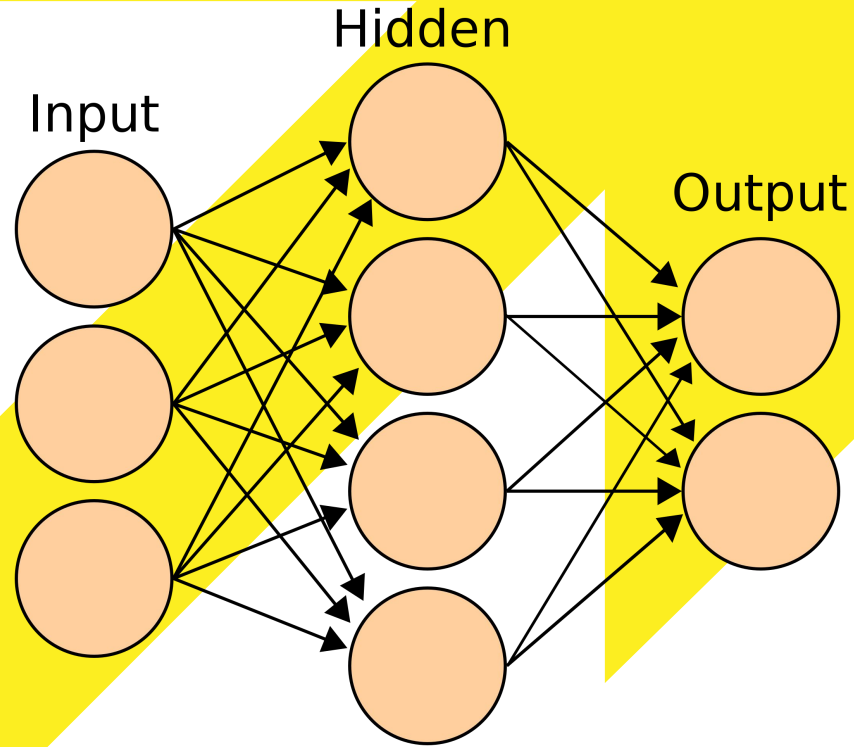
→ $\delta^{(l)} = \dfrac{\partial J}{\partial z^{(l)}} = \dfrac{\partial J}{\partial a^{(l)}} * \dfrac{\partial a^{(l)}}{\partial z^{(l)}} = (\delta^{(l+1)} * \Theta^{(l)}) \;.* \; g'(z^{(l)})$ $\qquad (l : L \rightarrow 1)$

→ $\dfrac{\partial J}{\partial \Theta^{(l)}} = \dfrac{\partial J}{\partial z^{(l+1)}} * \dfrac{\partial z^{(l+1)}}{\partial \Theta^{(l)}} = \delta^{(l+1)} * a^{(l)}$ $\qquad (l : L \rightarrow 1)$

→ $\Delta^{(l)} = \dfrac{1}{m} * \dfrac{\partial J}{\partial \Theta^{(l)}}$

# **Neural Network**

➔ Hidden layer:

- → Sigmoid
- → Tanh
- → ReLU
- → Leaky ReLU
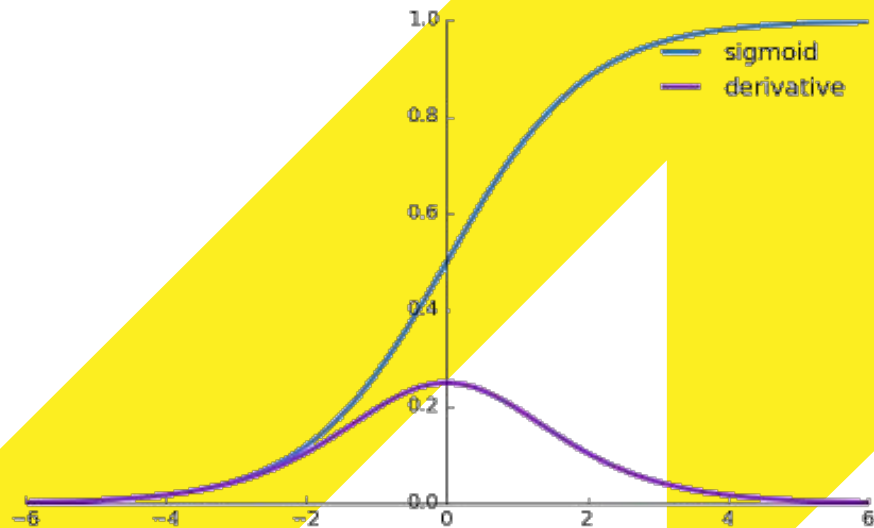- → Swish
- → Mish

Hidden

Input

Output

# Sigmoid activation

→ Formula:

$$g(z) = \frac{1}{1 + e^{-z}}$$

→ Derivative:
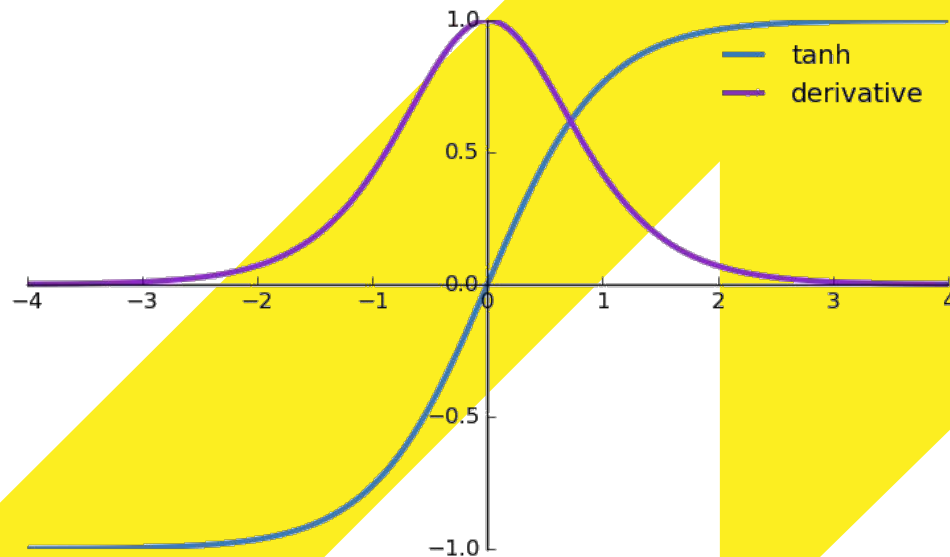
$$g'(z) = g(z)(1 - g(z))$$

# Tanh activation

➔ Formula:

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

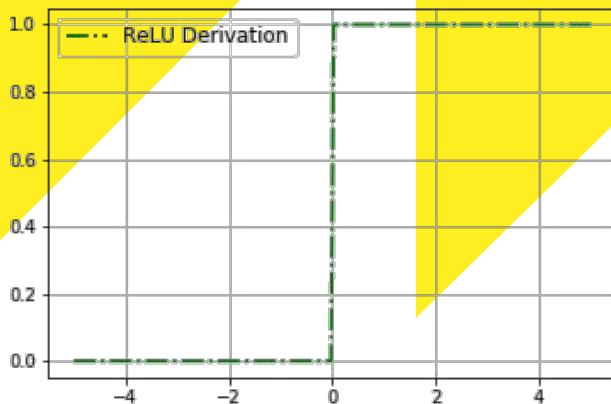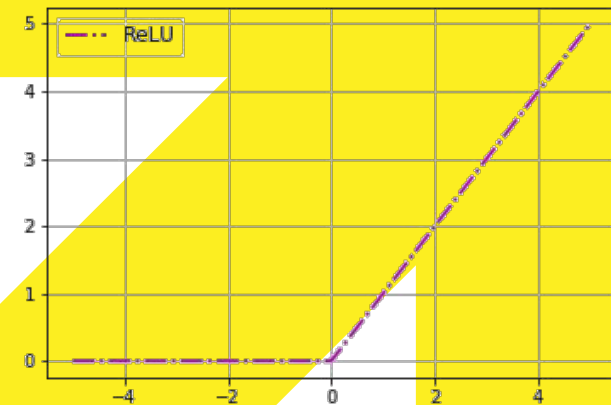➔ Derivative:

$$g'(z) = 1 - g(z)^2$$

# ReLU activation

→ Formula:

$$g(z) = \begin{cases} 0 & \text{for} \quad z < 0 \\ z & \text{for} \quad z \geq 0 \end{cases}$$

→ Derivative:

$$g'(z) = \begin{cases} 0 & \text{for} \quad z < 0 \\ 1 & \text{for} \quad z \geq 0 \end{cases}$$
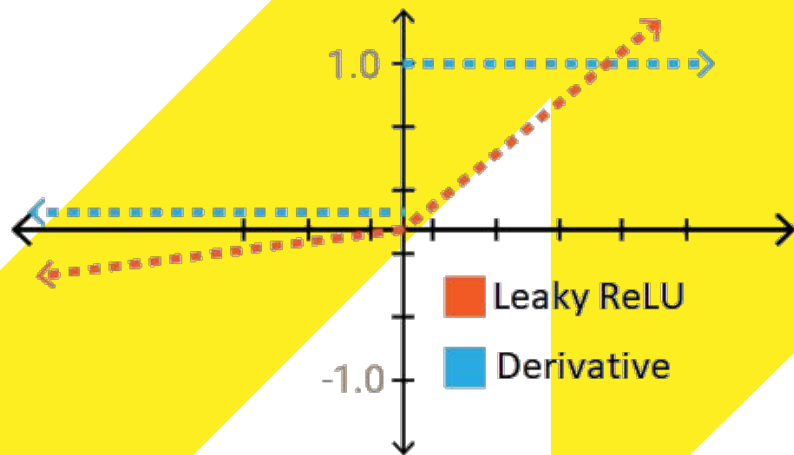
# Leaky ReLU activation

→ Formula:

$$g(z) = \begin{cases} \alpha z & \text{for} & z < 0 \\ z & \text{for} & z \geq 0 \end{cases}$$

→ Derivative:

$$g'(z) = \begin{cases} \alpha & \text{for} & z < 0 \\ 1 & \text{for} & z \geq 0 \end{cases}$$
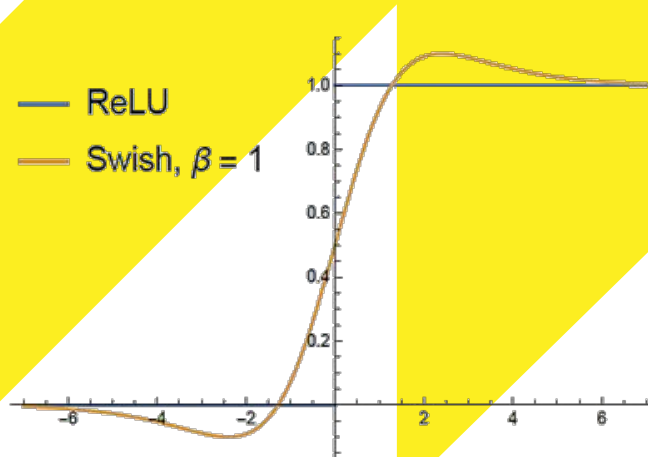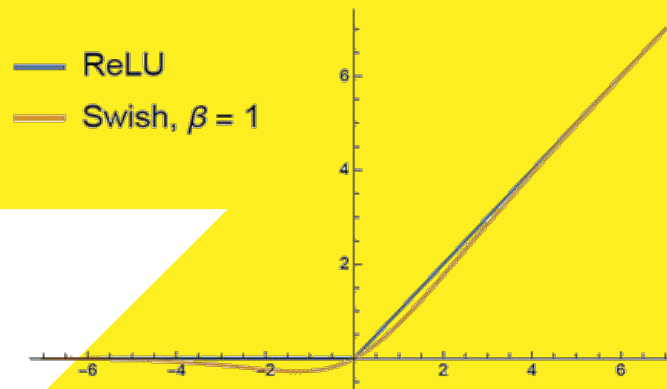
Leaky ReLU

Derivative

# Swish activation

→ Formula:

$$g(z) = z * sigmoid(z)$$

→ Derivative:

$$g'(z) = g(z) + sigmoid(z) * (1 - g(z))$$



ReLU
Swish, $\beta = 1$

ReLU
Swish, $\beta = 1$

# Mish activation
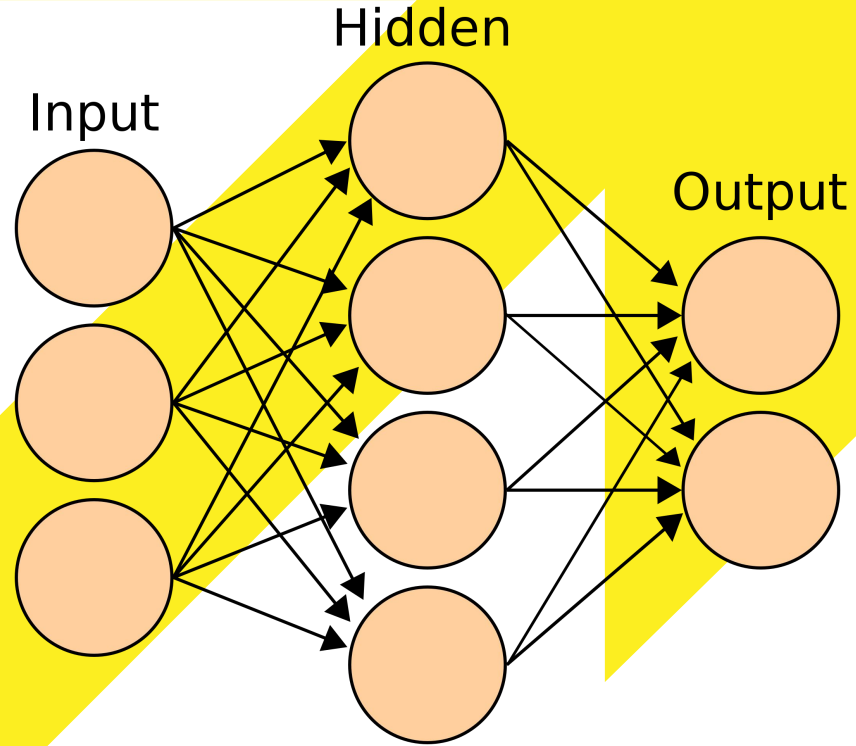
➜ Formula:

$$g(z) = z * \tanh(\ln(1 + e^z))$$

➜ Derivative:

$$g'(z) = \frac{e^z \omega}{\delta^2} = \frac{e^z(4e^{2z} + e^{3z} + 4(1 + z) + e^z(6 + 4z))}{(2 + 2e^z + e^{2z})^2}$$
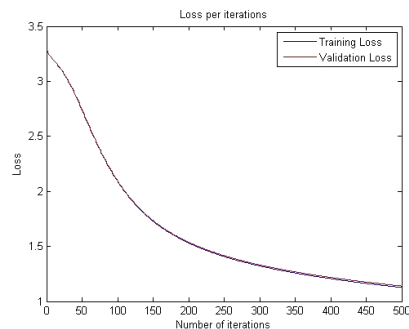
# Neural Network

➔ Hidden layer:
  → 64 nodes
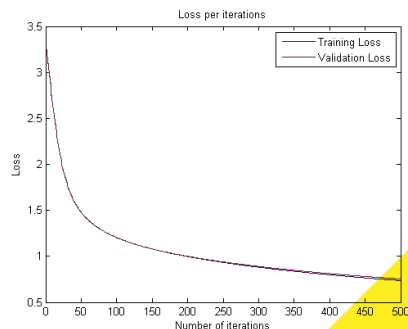  → 128 nodes
  → 256 nodes
  → 512 nodes



Input

Hidden

Output

# Neural Network

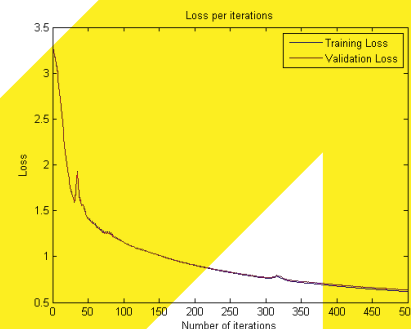| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| NN_sigmoid_softmax_64 | 51,930 | 67.95 % | 68.27 % | 67.81 % |
| NN_tanh_softmax_64 | 51,930 | 79.13 % | 78.97 % | 78.53 % |
| NN_relu_softmax_64 | 51,930 | 82.12 % | 81.98 % | 81.74 % |
| NN_leakyrelu_softmax_64 | 51,930 | 81.64 % | 81.91 % | 80.96 % |
| NN_swish_softmax_64 | 51,930 | 81.77 % | 80.86 % | 80.92 % |
| NN_mish_softmax_64 | 51,930 | 82.31 % | 81.85 % | 81.75 % |

# Neural Network



**NN_sigmoid_softmax_64**
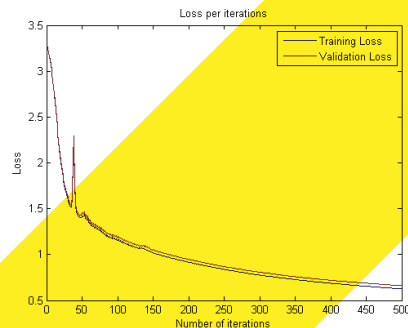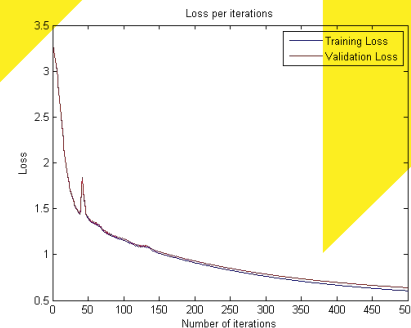
**NN_tanh_softmax_64**

**NN_relu_softmax_64**

**NN_leakyrelu_softmax_64**

**NN_swish_softmax_64**

**NN_mish_softmax_64**

# Neural Network

| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| NN_relu_softmax_128 | 103,834 | 84.06 % | 83.37 % | 83.19 % |
| NN_leakyrelu_softmax_128 | 103,834 | 84.14 % | 82.67 % | 82.95 % |
| NN_swish_softmax_128 | 103,834 | 83.71 % | 82.84 % | 83.02 % |
| NN_mish_softmax_128 | 103,834 | 83.78 % | 83.05 % | 83.14 % |

# Neural Network

| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| NN_relu_softmax_256 | **207,642** | **85.57 %** | **85.18 %** | **84.60 %** |
| NN_leakyrelu_softmax_256 | **207,642** | **85.37 %** | **84.70 %** | **84.15 %** |
| NN_swish_softmax_256 | **207,642** | **84.77 %** | **83.85 %** | **83.55 %** |
| NN_mish_softmax_256 | **207,642** | **85.16 %** | **84.43 %** | **84.30 %** |

# Neural Network

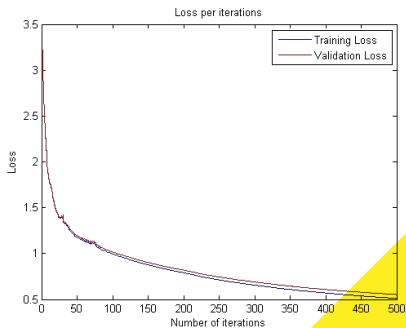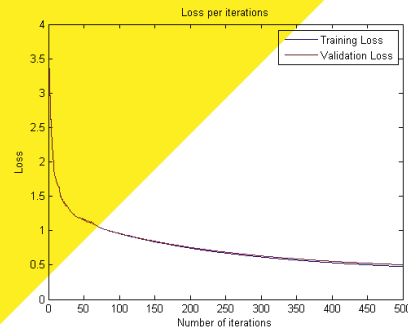| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| NN_relu_softmax_512 | 415,258 | 87.03 % | 86.02 % | 85.51 % |
| NN_leakyrelu_softmax_512 | 415,258 | 86.79 % | 85.00 % | 85.30 % |
| NN_swish_softmax_512 | 415,258 | 85.20 % | 83.89 % | 84.13 % |
| NN_mish_softmax_512 | 415,258 | 86.21 % | 85.65 % | 85.01 % |

# Neural Network



NN_relu_softmax_512



NN_leakyrelu_softmax_512



NN_swish_softmax_512



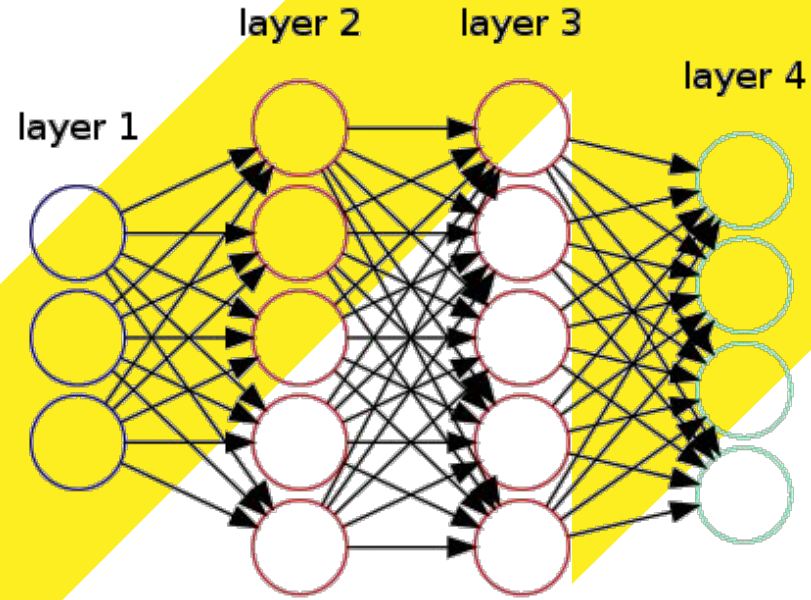NN_mish_softmax_512

# NN_relu_softmax_512

## Confusion Matrix

# Deeper Neural Network

➔ Architecture:

  → 1 input layer (784 nodes)

  → 2 hidden layers

  → 1 output layer (26 nodes)

layer 1

layer 2

layer 3

layer 4

# Deeper Neural Network

➔ Hidden layer 1:

  → 512 nodes

➔ Hidden layer 2:

  → 64 nodes

  → 128 nodes

  → 256 nodes

# Deeper Neural Network

| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| NN_relu_softmax_512_64 | 436,442 | 88.46 % | 86.85 % | 86.71 % |
| NN_leakyrelu_softmax_512_64 | 436,442 | 88.84 % | 87.31 % | 87.28 % |
| NN_swish_softmax_512_64 | 436,442 | 87.58 % | 86.68 % | 86.42 % |
| NN_mish_softmax_512_64 | 436,442 | 87.94 % | 87.17 % | 87.01 % |

# Deeper Neural Network



NN_relu_softmax_512_64



NN_leakyrelu_softmax_512_64



NN_swish_softmax_512_64



NN_mish_softmax_512_64

# Deeper Neural Network

| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| NN_relu_softmax_512_128 | 470,938 | 89.40 % | 88.05 % | 87.56 % |
| NN_leakyrelu_softmax_512_128 | 470,938 | 89.25 % | 87.36 % | 87.65 % |
| NN_swish_softmax_512_128 | 470,938 | 87.86 % | 87.02 % | 86.81 % |
| NN_mish_softmax_512_128 | 470,938 | 88.34 % | 87.37 % | 87.24 % |

# Deeper Neural Network

| Model | Total Parameters | Training Set | Validation Set | Test Set |
|---|---|---|---|---|
| NN_relu_softmax_512_256 | 539,930 | 89.98 % | 88.10 % | 88.29 % |
| NN_leakyrelu_softmax_512_256 | 539,930 | 89.99 % | 88.34 % | 88.31 % |
| NN_swish_softmax_512_256 | 539,930 | 87.56 % | 86.71 % | 86.44 % |
| NN_mish_softmax_512_256 | 539,930 | 88.74 % | 87.69 % | 87.50 % |

# Deeper Neural Network



NN_relu_softmax_512_256

NN_leakyrelu_softmax_512_256

NN_swish_softmax_512_256

NN_mish_softmax_512_256

# NN_leakyrelu_softmax_512_256

## Confusion Matrix



Confusion Matrix

# Even deeper NN, but wait…

OOP-style Python realization.

Naive implementation → exploding gradient, NaN loss, network learning nothing.

Solution:

➔ Stable softmax (max subtraction)
➔ Good initialization (Xavier/He)

# Optimization technique

Gradient descent with momentum → faster convergence, avoid bad local minima.

**Crucial**: mini-batch GD. Much faster convergence.

For the rest of training process, we use lr=0.1, momentum=0.9, batch size 256.

# New results

Every run easily reached 88-89% accuracy on test set and almost converges within **1-2 minutes** of training time.



After 6 epochs with training time of 1m47s, in this run a 512-256-128 ReLU model achieved 0.299 training loss and 88.33% accuracy on test set.

# About regularization

Simple regularization prevents overfitting, but:

➔ It is essentially limiting the model's complexity.
➔ Cherry-picking lambda can be hard.



Underfit     Optimal     Overfit

# Data shuffling

Serving the data in a meaningful order or the same order every epoch might introduce some bias.

Thus the training data is shuffled every epoch.

# New results

Every run reached ~90% accuracy on test set.



Loss graph of a run with 90.92% accuracy on test set (512-256-128 Swish). The run was automatically stopped after 26 epochs.

# On data augmentation

Low image resolution + simple model (NN only) →
Applying augmentation is **generally bad**.

However, small rotations might be useful.



20 degree right rotation of an 'a' letter. In training we applied on-the-fly 50% chance 10-15 degree rotation for each image.

# Tanh Exponential activation

➔ Formula:

$$g(z) = z * tanh(e^z)$$

➔ Derivative:

$$g'(z) = \tanh(e^z) - ze^z(\tanh^2(e^z) - 1)$$



62

# The result

We reached ~91% test accuracy on the most runs.



Loss graph of a run with 91.25% accuracy on test set (512-256-128 TanhExp). The run was automatically stopped after 38 epochs.

# We tried, didn't work

→ Deeper/wider network.
→ Multi-angle prediction.
→ Learning rate decay.

# 5.
# Evaluation

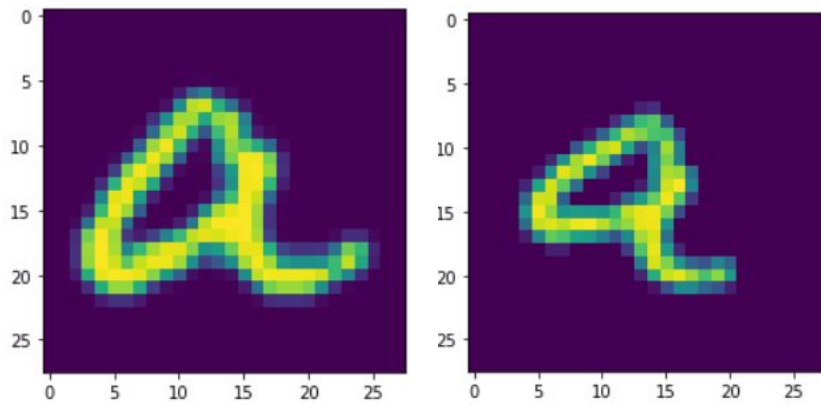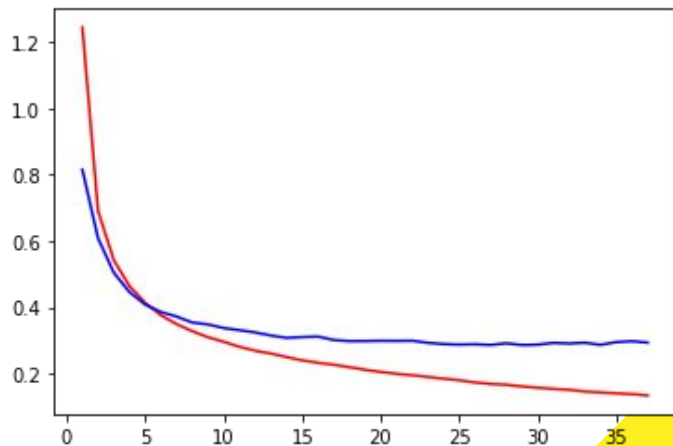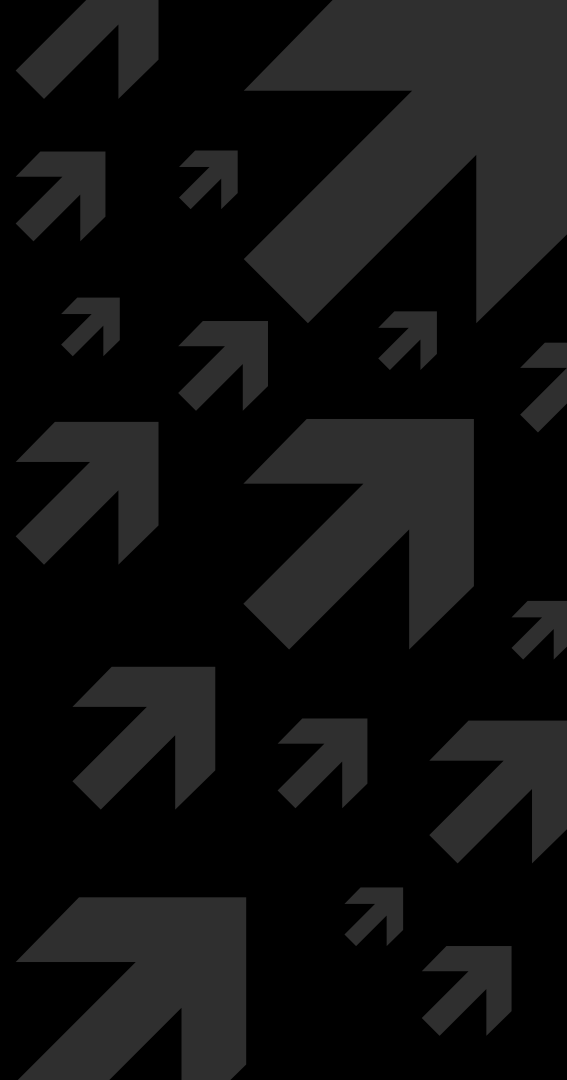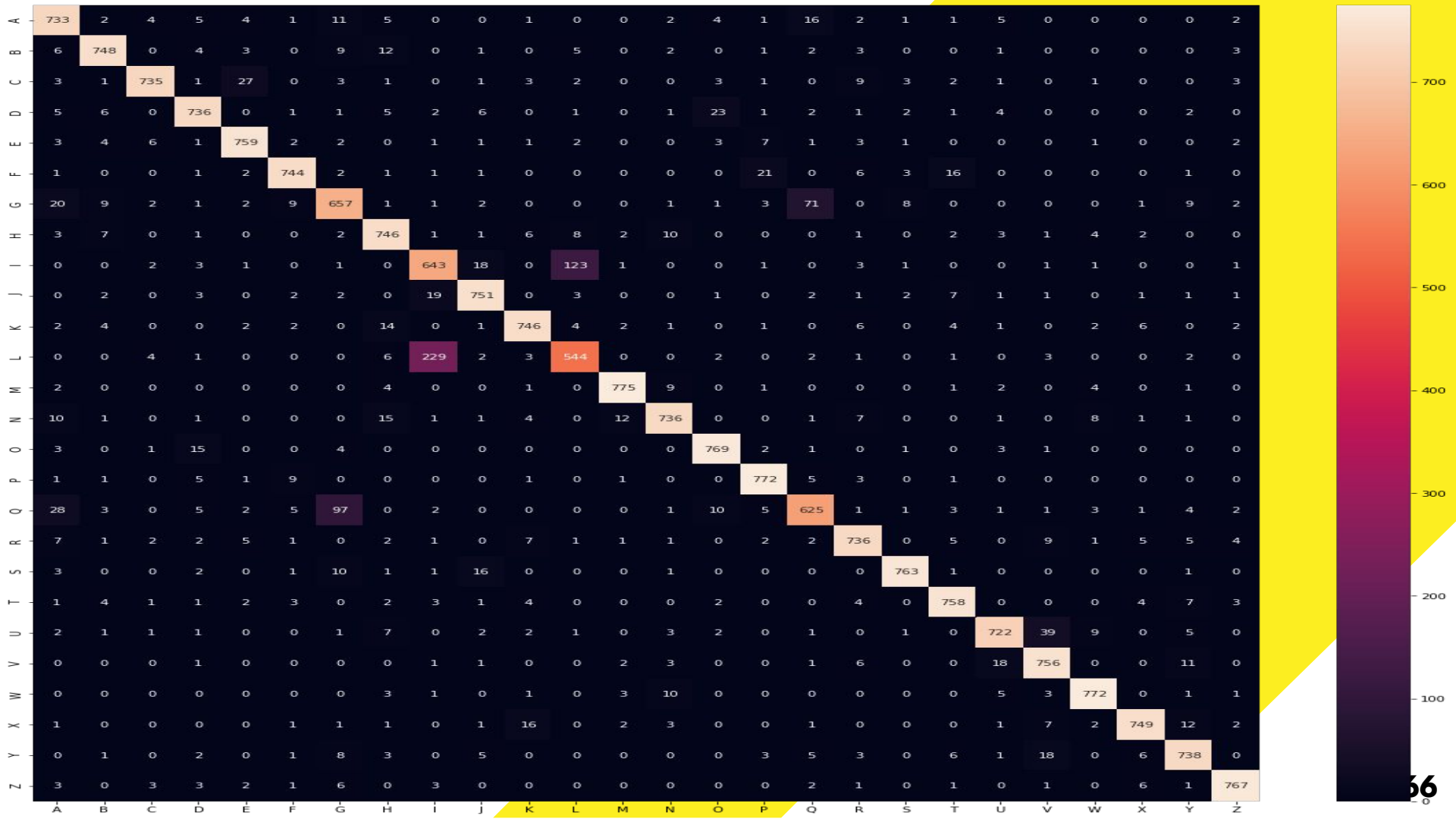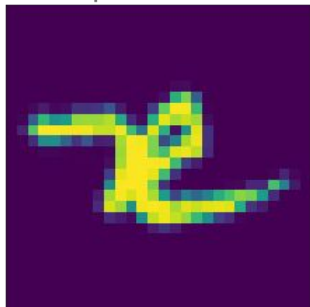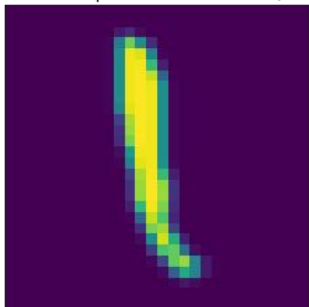|   | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 733 | 2 | 4 | 5 | 4 | 1 | 11 | 5 | 0 | 0 | 1 | 0 | 0 | 2 | 4 | 1 | 16 | 2 | 1 | 1 | 5 | 0 | 0 | 0 | 0 | 2 |
| B | 6 | 748 | 0 | 4 | 3 | 0 | 9 | 12 | 0 | 1 | 0 | 5 | 0 | 2 | 0 | 1 | 2 | 3 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 3 |
| C | 3 | 1 | 735 | 1 | 27 | 0 | 3 | 1 | 0 | 1 | 3 | 2 | 0 | 0 | 3 | 1 | 0 | 9 | 3 | 2 | 1 | 0 | 1 | 0 | 0 | 3 |
| D | 5 | 6 | 0 | 736 | 0 | 1 | 1 | 5 | 2 | 6 | 0 | 1 | 0 | 1 | 23 | 1 | 2 | 1 | 2 | 1 | 4 | 0 | 0 | 0 | 2 | 0 |
| E | 3 | 4 | 6 | 1 | 759 | 2 | 2 | 0 | 1 | 1 | 1 | 2 | 0 | 0 | 3 | 7 | 1 | 3 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 2 |
| F | 1 | 0 | 0 | 1 | 2 | 744 | 2 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 21 | 0 | 6 | 3 | 16 | 0 | 0 | 0 | 0 | 1 | 0 |
| G | 20 | 9 | 2 | 1 | 2 | 9 | 657 | 1 | 1 | 2 | 0 | 0 | 0 | 1 | 1 | 3 | 71 | 0 | 8 | 0 | 0 | 0 | 0 | 1 | 9 | 2 |
| H | 3 | 7 | 0 | 1 | 0 | 0 | 2 | 746 | 1 | 1 | 6 | 8 | 2 | 10 | 0 | 0 | 0 | 1 | 0 | 2 | 3 | 1 | 4 | 2 | 0 | 0 |
| I | 0 | 0 | 2 | 3 | 1 | 0 | 1 | 0 | 643 | 18 | 0 | 123 | 1 | 0 | 0 | 1 | 0 | 3 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 |
| J | 0 | 2 | 0 | 3 | 0 | 2 | 2 | 0 | 19 | 751 | 0 | 3 | 0 | 0 | 1 | 0 | 2 | 1 | 2 | 7 | 1 | 1 | 0 | 1 | 1 | 1 |
| K | 2 | 4 | 0 | 0 | 2 | 2 | 0 | 14 | 0 | 1 | 746 | 4 | 2 | 1 | 0 | 1 | 0 | 6 | 0 | 4 | 1 | 0 | 2 | 6 | 0 | 2 |
| L | 0 | 0 | 4 | 1 | 0 | 0 | 0 | 6 | 229 | 2 | 3 | 544 | 0 | 0 | 2 | 0 | 2 | 1 | 0 | 1 | 0 | 3 | 0 | 0 | 2 | 0 |
| M | 2 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 1 | 0 | 0 | 775 | 9 | 0 | 1 | 0 | 0 | 1 | 2 | 0 | 4 | 0 | 1 | 0 | 0 |
| N | 10 | 1 | 0 | 1 | 0 | 0 | 0 | 15 | 1 | 1 | 4 | 0 | 12 | 736 | 0 | 0 | 1 | 7 | 0 | 0 | 1 | 0 | 8 | 1 | 1 | 0 |
| O | 3 | 0 | 1 | 15 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 769 | 2 | 1 | 0 | 1 | 3 | 1 | 0 | 0 | 0 | 0 | 0 |
| P | 1 | 1 | 0 | 5 | 1 | 9 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 772 | 5 | 3 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Q | 28 | 3 | 0 | 5 | 2 | 5 | 97 | 0 | 2 | 0 | 0 | 0 | 0 | 1 | 10 | 5 | 625 | 1 | 1 | 3 | 1 | 1 | 3 | 1 | 4 | 2 |
| R | 7 | 1 | 2 | 2 | 5 | 1 | 0 | 2 | 1 | 0 | 7 | 1 | 1 | 1 | 0 | 2 | 2 | 736 | 0 | 5 | 0 | 9 | 1 | 5 | 5 | 4 |
| S | 3 | 0 | 0 | 2 | 0 | 1 | 10 | 1 | 1 | 16 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 763 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| T | 1 | 4 | 1 | 1 | 2 | 3 | 0 | 2 | 3 | 1 | 4 | 0 | 0 | 2 | 0 | 0 | 4 | 0 | 0 | 758 | 0 | 0 | 0 | 4 | 7 | 0 |
| U | 2 | 1 | 1 | 1 | 0 | 0 | 1 | 7 | 0 | 2 | 2 | 1 | 0 | 3 | 2 | 0 | 1 | 0 | 1 | 0 | 722 | 39 | 9 | 0 | 5 | 0 |
| V | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 2 | 3 | 0 | 1 | 6 | 0 | 0 | 0 | 18 | 756 | 0 | 0 | 11 | 0 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 1 | 0 | 1 | 0 | 3 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 3 | 772 | 0 | 1 | 1 |
| X | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 16 | 0 | 2 | 3 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 7 | 2 | 749 | 12 | 2 |
| Y | 0 | 1 | 0 | 2 | 0 | 1 | 8 | 3 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 3 | 5 | 3 | 0 | 6 | 1 | 18 | 0 | 6 | 738 | 0 |
| Z | 3 | 0 | 3 | 3 | 2 | 1 | 6 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 1 | 0 | 1 | 0 | 6 | 1 | 767 |

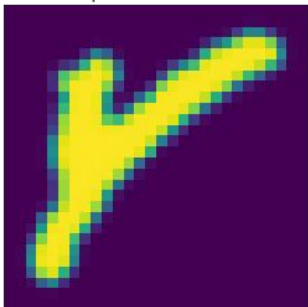# Some misclassified examples



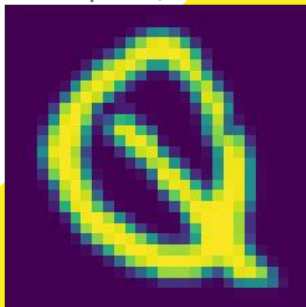Truth: R, pred: ['Z', 'X', 'R', 'K', 'E']

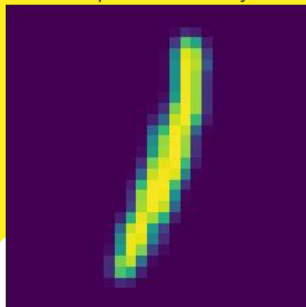Truth: L, pred: ['I', 'L', 'T', 'Y', 'Q']
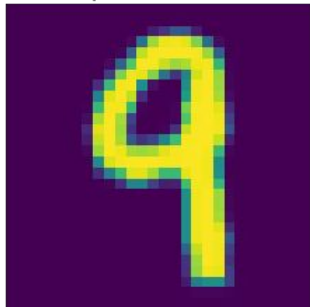
Truth: Y, pred: ['R', 'Y', 'P', 'V', 'F']

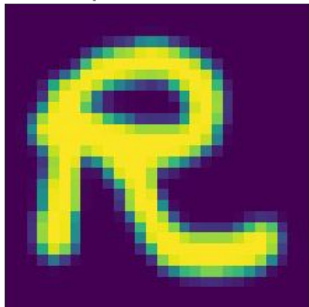Truth: G, pred: ['Q', 'G', 'A', 'O', 'S']
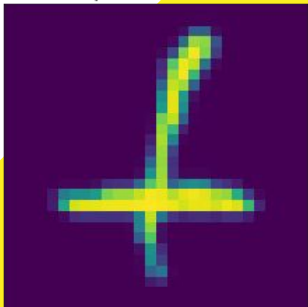
Truth: I, pred: ['L', 'I', 'T', 'J', 'V']
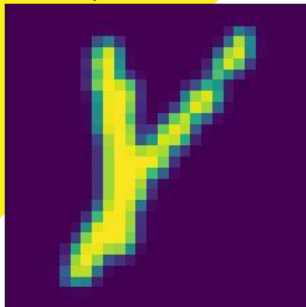
Truth: G, pred: ['Q', 'G', 'Y', 'R', 'A']

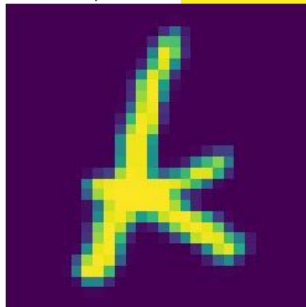Truth: R, pred: ['E', 'R', 'Q', 'K', 'P']

Truth: T, pred: ['F', 'T', 'S', 'I', 'G']

Truth: X, pred: ['Y', 'X', 'V', 'K', 'R']

Truth: K, pred: ['H', 'K', 'L', 'B', 'T']

# Top-K accuracy
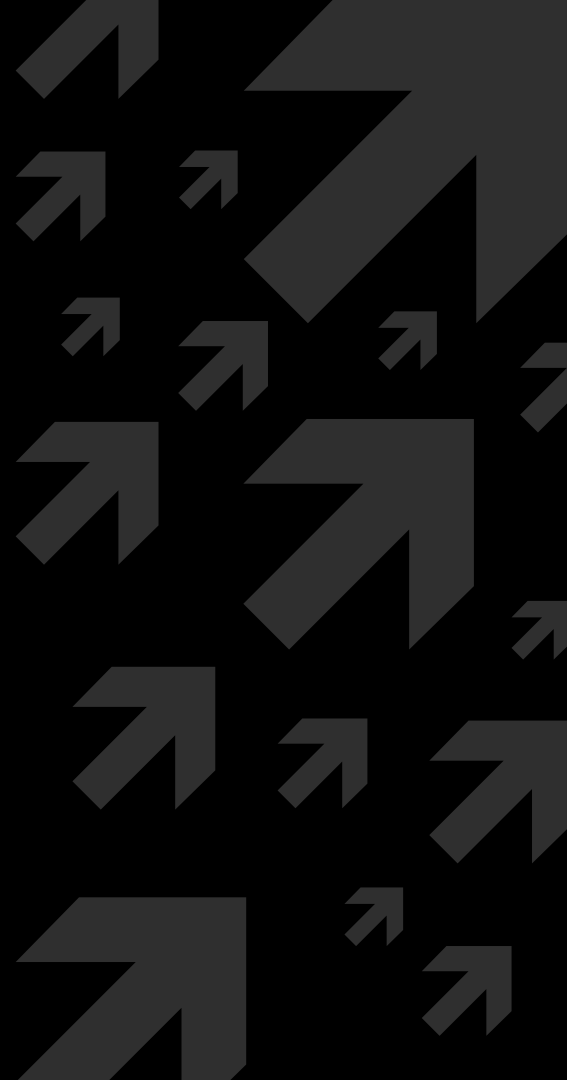
A correct prediction have the ground truth value in top K most confident classes of the prediction.

Our best model:

➔ Top-1: 91.25%
➔ Top-2: 97.63%
➔ Top-5: 99.25%

# 6.
# Application

# Predicting digital images of letters

Problems:

➔ The model only takes grayscale inputs.
➔ The image size must be 28x28.
➔ Train/test images have a common standard: black background and letter foreground.

# Predicting digital images of letters

Straight conversation to grayscale and resizing just doesn't work!



The 28x28 'S' letter image above was predicted to be [Q, G, A, M, W].

# Our workaround

1. Resize the image to a fixed size (300x300).
2. Apply K-mean segmentation on the image with K=2.
3. Find the 'background' and 'foreground' classes, assuming the background has more pixels.
4. Set all the pixels in background class to 0 and foreground class to 1.
5. Rescale to 28x28.
6. Apply Gaussian blur.

# Our workaround

In most cases, the example becomes closer to what had been seen in the training set, thus allowing more precise prediction.



The 'S' letter is now correctly predicted.
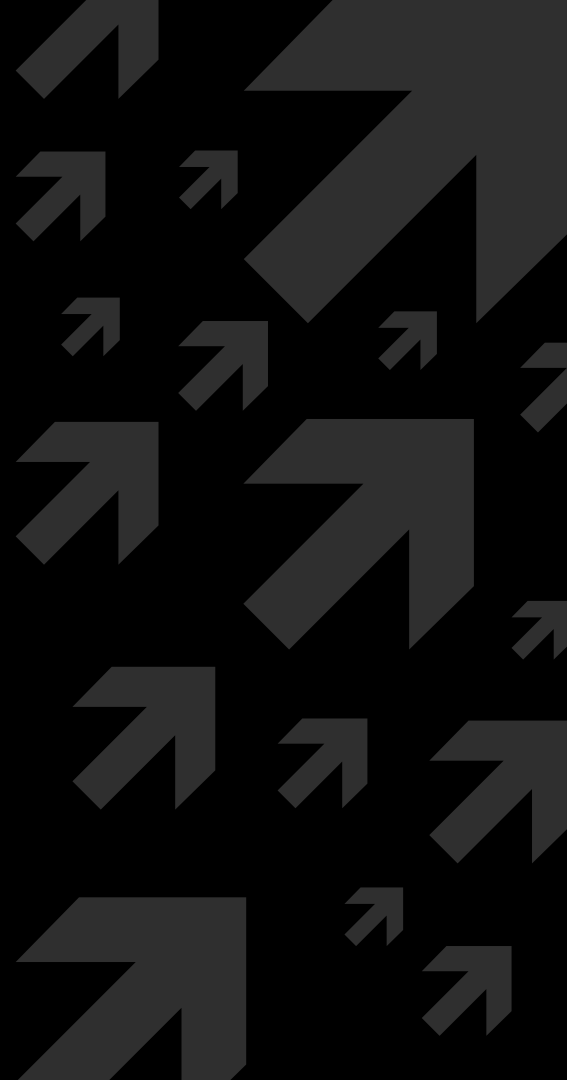
# Some correctly predicted samples

# Better use with letter segmentation

In the example below, all 4 letters once segmented were correctly predicted.

# 7.
# Final words

# Summary

We implemented and trained NNs without any high-end libraries on the EMNIST letters dataset. Top result achieved 91.25% accuracy on the test set.

We showed that the simple model can be used for some real-world images application.

# Takeaway

We gained a better understanding of neural networks and related concepts.

We now have more experience in actually debugging and improving a model.

We learned to put our ML model to potential real-world applications.

# Thanks!

Any questions?