

# DS 4400

## Machine Learning and Data Mining I Spring 2021

Alina Oprea

Associate Professor

Khoury College of Computer Science

Northeastern University

March 23 2021

# Announcements

- HW 4 is due on Friday, March 26
- Project milestone due on March 31
  - Template in Gradescope
- Final exam on Tuesday, April 6
  - Review on Thursday, April 1
- Last homework on ethics
  - After ethics class (April 8)
  - Group assignment

# Outline

- Feed Forward Neural Networks
  - Forward Propagation
  - Hyper-parameters
  - Activations
- Multi-class classification
  - The softmax classifier
- Examples
- Keras tutorial

# References

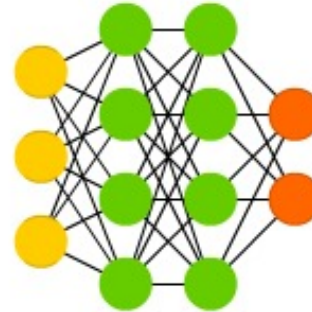
- Deep Learning books
  - <https://d2l.ai/> (D2L)
  - <https://www.deeplearningbook.org/> (advanced)
- Stanford notes on deep learning
  - [http://cs229.stanford.edu/summer2020/cs229-notes-deep\\_learning.pdf](http://cs229.stanford.edu/summer2020/cs229-notes-deep_learning.pdf)

# Neural Network Architectures

## Feed-Forward Networks

- Neurons from each layer connect to neurons from next layer

Deep Feed Forward (DFF)

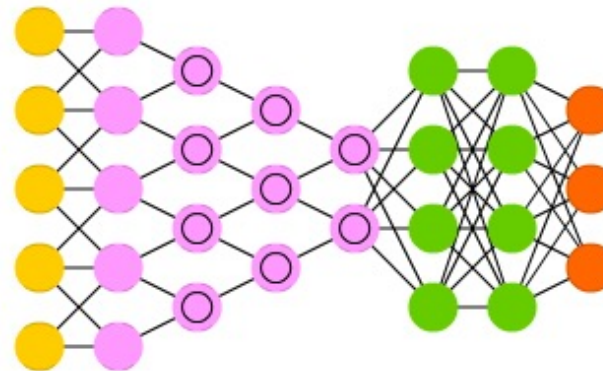


MLP

## Convolutional Networks

- Includes convolution layer for feature reduction
- Learns hierarchical representations

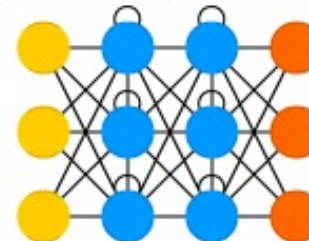
Deep Convolutional Network (DCN)



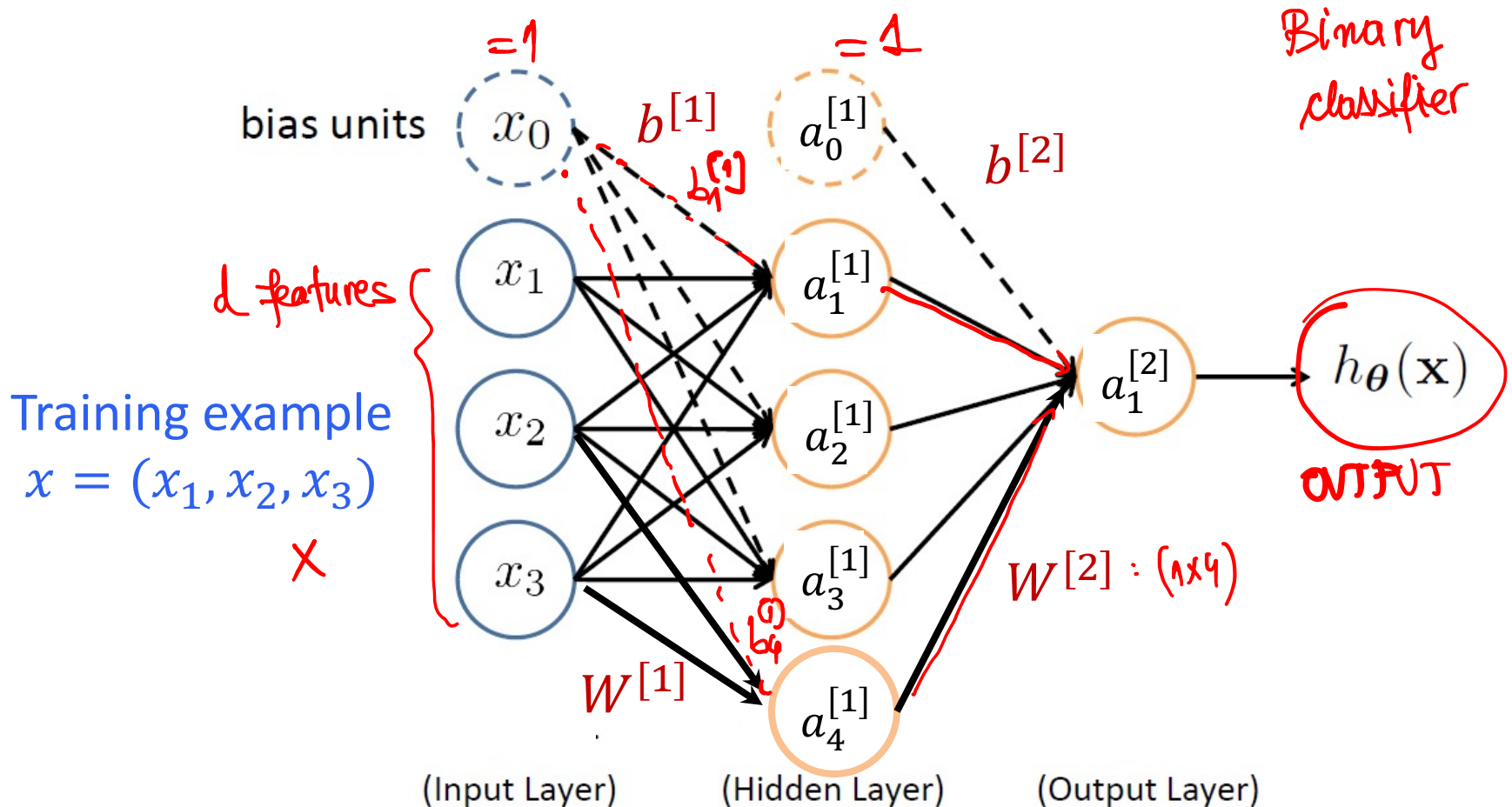
## Recurrent Networks

- Keep hidden state
- Have cycles in computational graph

Recurrent Neural Network (RNN)



# Feed-Forward Neural Network



No cycles

$$\theta = (b^{[1]}, W^{[1]}, b^{[2]}, W^{[2]})$$

vector      matrix

# Hyperparameters

- Architecture
  - # layers (hidden)
  - Activation function
  - # neurons per layer
- Hyper-param for GD
  - learning rate; stopping condition
- Regularization

# Vectorization: First Layer

$$z^{[1]} = W^{[1]} \cdot x + b^{[1]}$$

LINEAR

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_4^{[1]} \end{bmatrix}$$

Size: (4x1)

$$W^{[1]} = \begin{bmatrix} w_1^{[1]} \\ w_2^{[1]} \\ \vdots \\ w_4^{[1]} \end{bmatrix}$$

Size: (4x3)  
↓   ↓  
# neurons   Layer 1

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

(3x1)

$$b^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix}$$

(4x1)

$$a^{[1]} = g(z^{[1]}) = \begin{bmatrix} g(z_1^{[1]}) \\ \vdots \\ g(z_4^{[1]}) \end{bmatrix}$$

Size (4x1)

NON-LINEAR  
ACTIVATION



# Vectorization: Second Layer

$$z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]} \quad \text{LINEAR}$$

$\downarrow \quad \quad \downarrow \quad \quad \downarrow$   
 $(1 \times 4) \quad (4 \times 1) \quad (1 \times 1)$   
 $\vdots$

$$z^{[j]} = W^{[j]} \cdot a^{[j-1]} + b^{[j]}$$

$$a^{[2]} = g(z^{[2]}) \quad \text{ACTIVATION}$$

$\downarrow$   
OUTPUT :  $h_{\theta}(x) = a^{[2]}$

# Terminology

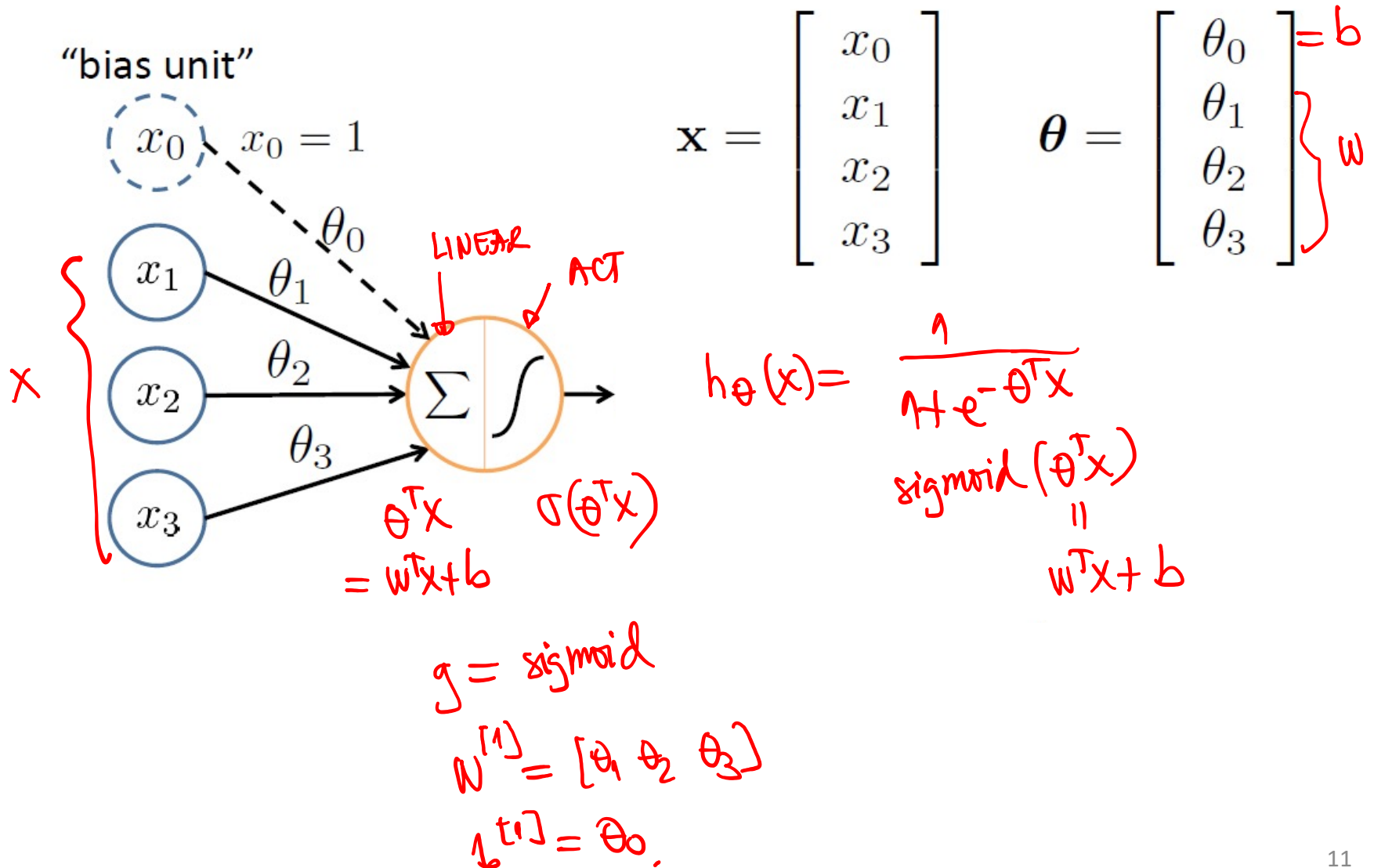
$g$  = act. function

$a_i^{[j]}$  = act. of unit  $i$  at layer  $j$

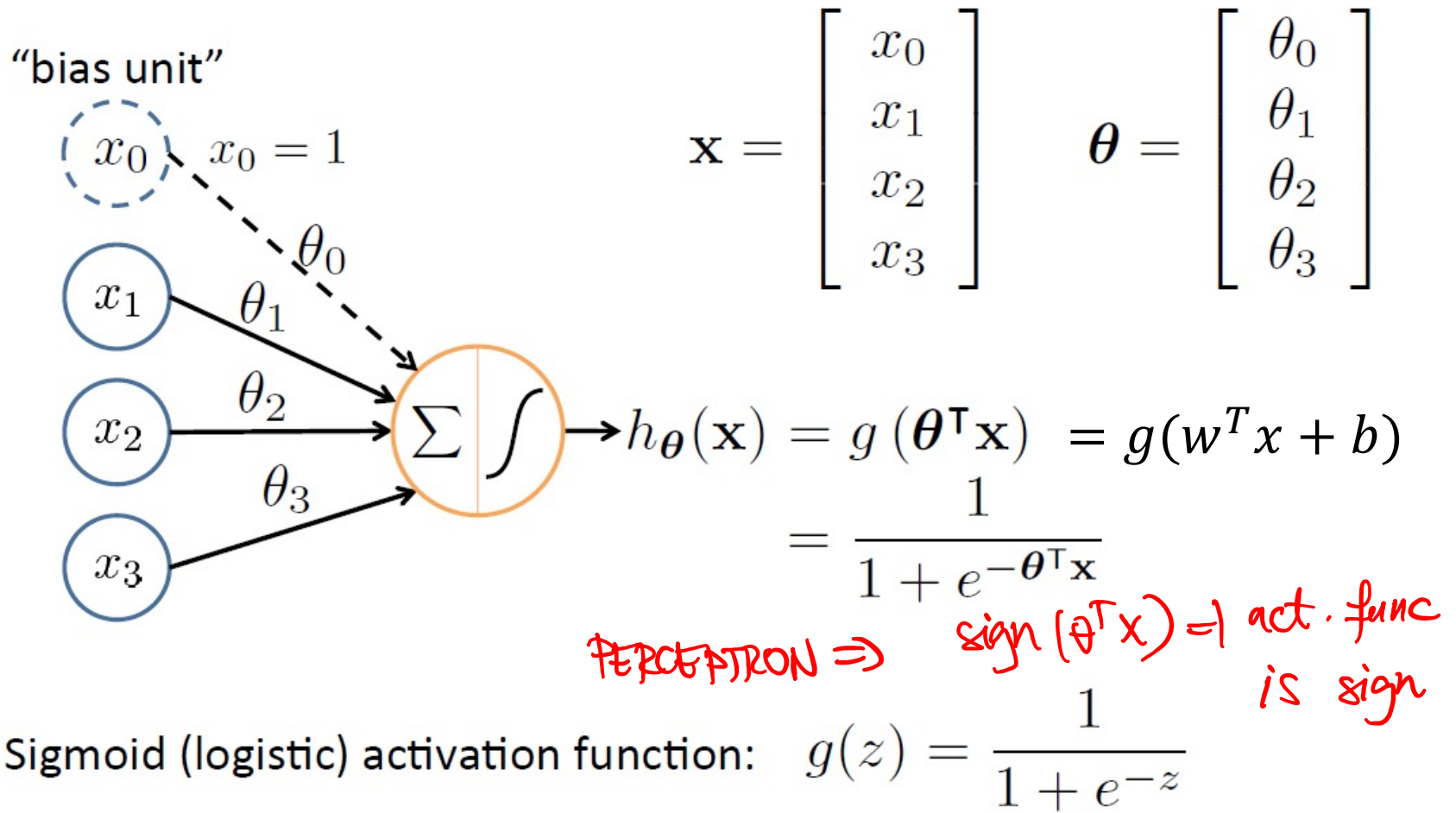
$W^{[j]}$  = weight matrix layer  $j$

$b^{[j]}$  = bias vector layer  $j$

# Logistic Unit: A simple NN



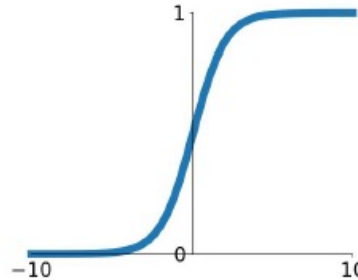
# Logistic Unit: A simple NN



No hidden layers

# Activation Functions

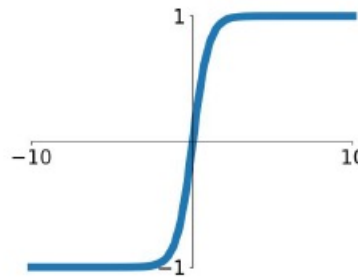
- **Sigmoid**  
 $\sigma(x) = \frac{1}{1+e^{-x}}$



*last layer*

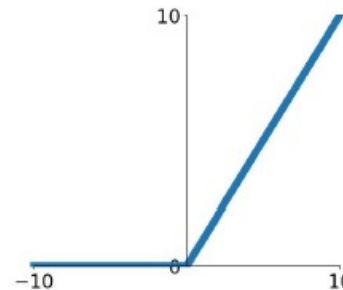
Binary  
Classification

**tanh**  
 $\tanh(x)$



Regression

- **ReLU**  
 $\max(0, x)$



Intermediary  
layers

- *softmax*

*Multi-class  
classification*

# Why Non-Linear Activations?

- Assume  $g$  is linear:  $g(z) = Uz$  ;  $U$  matrix.

Layer 1:  $z^{[1]} = W^{[1]}x + b^{[1]}$

$$a^{[1]} = g(z^{[1]}) = Uz^{[1]} = UW^{[1]}x + Ub^{[1]} \quad \text{LINEAR}$$

Layer 2:  $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$

$$a^{[2]} = g(z^{[2]}) = Uz^{[2]} = UW^{[2]}a^{[1]} + Ub^{[2]}$$

$$= \underbrace{UW^{[2]}UW^{[1]}}_A x + \underbrace{UW^{[2]}Ub^{[1]} + Ub^{[2]}}_b$$

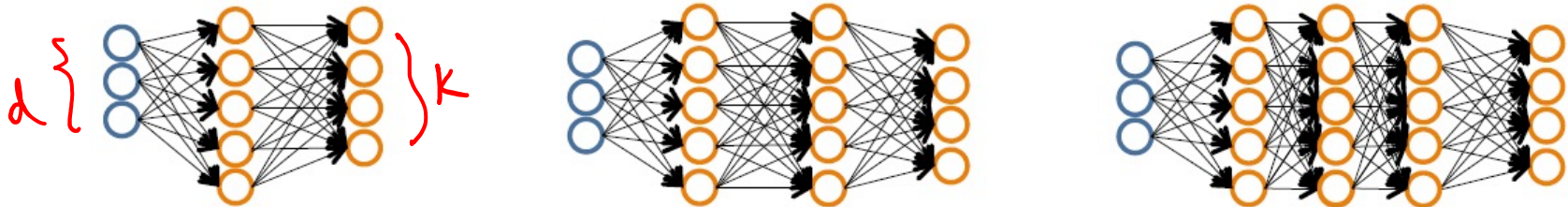
$$= Ax + b$$

LINEAR FUNCTION of INPUT  $x$

Similar to  $x \begin{Bmatrix} 0 & \dots & b \\ 0 & & \\ 0 & & \\ 0 & & \end{Bmatrix} \xrightarrow{A} \begin{Bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{Bmatrix}$  ; No hidden layers

# How to pick architecture?

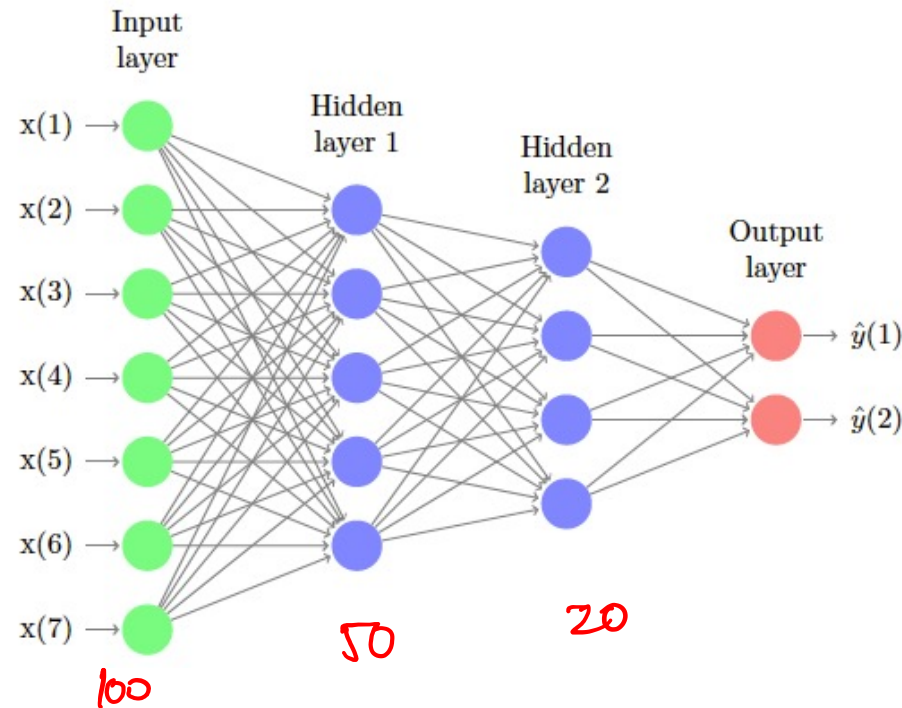
Pick a network architecture (connectivity pattern between nodes)



- # input units = # of features in dataset
- # output units = # classes  $k$

**Reasonable default:** 1 hidden layer

# FFNN Architectures



- Input and Output Layers are completely specified by the problem domain
- In the Hidden Layers, number of neurons in Layer  $i+1$  is usually smaller or equal to the number of neurons in Layer  $i$



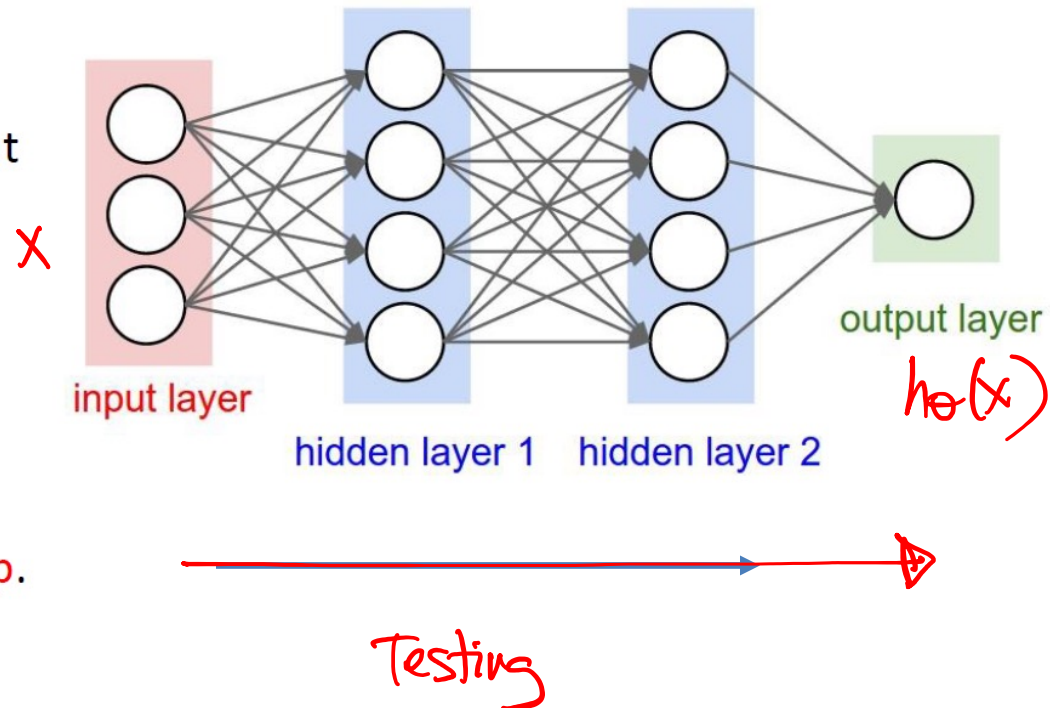
# Training Neural Networks

- Input training dataset  $D$ 
  - Number of features:  $d$
  - Labels from  $K$  classes
- First layer has  $d+1$  units (one per feature and bias)
- Output layer has  $K$  units
- Training procedure determines parameters that optimize loss function
  - Backpropagation *: training*
  - Learn optimal  $W^{[i]}, b^{[i]}$  at layer  $i$
- Evaluation of a point done by forward propagation *testing*

# Forward Propagation

Fix  $w^{[j]}$ ,  $b^{[j]}$

- The input neurons first receive the data features of the object. After processing the data, they send their output to the first hidden layer.
- The hidden layer processes this output and sends the results to the next hidden layer.
- This continues until the data reaches the final output layer, where the output value determines the object's classification.
- This entire process is known as **Forward Propagation**, or **Forward prop.**



# Multi-Class Classification



Pedestrian



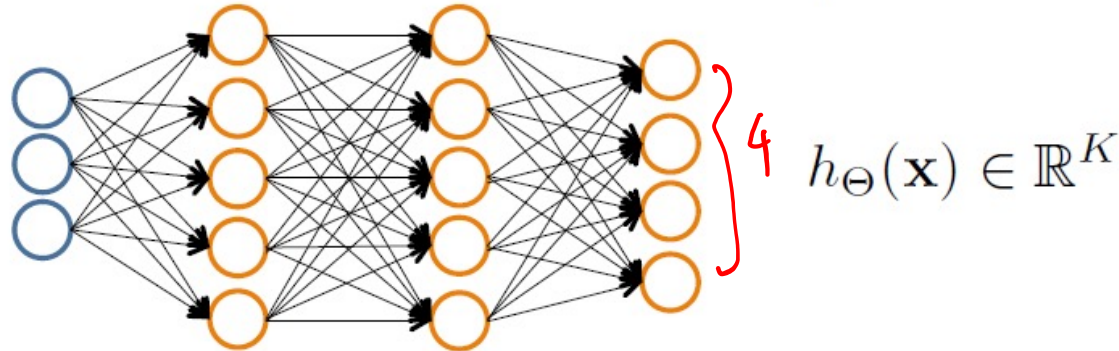
Car



Motorcycle



Truck



We want:

*y<sub>i</sub> labels.*

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} \textcircled{1} \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

when pedestrian

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

when car

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

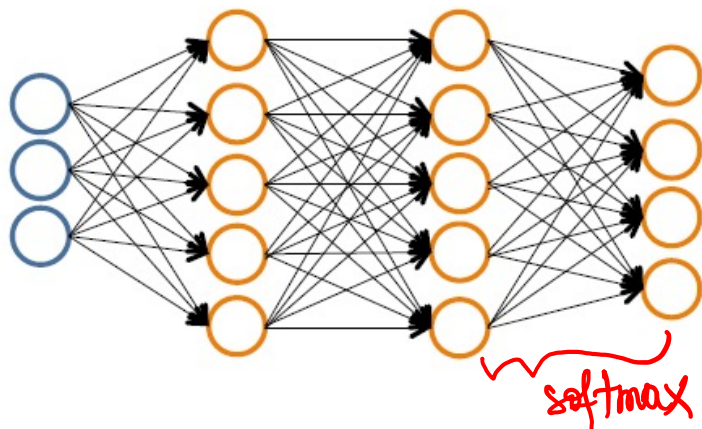
when motorcycle

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

when truck

*one-hot encoding for label*

# Neural Network Classification



**Given:**

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

$\mathbf{s} \in \mathbb{N}^{+L}$  contains # nodes at each layer

–  $s_0 = d$  (# features)

## Binary classification

$y = 0$  or  $1$

1 output unit ( $s_{L-1} = 1$ )

classes 0 / 1

Sigmoid

Prob of class 1

## Multi-class classification ( $K$ classes)

$$\mathbf{y} \in \mathbb{R}^K \quad \text{e.g.} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

pedestrian   car   motorcycle   truck

$K$  output units ( $s_{L-1} = K$ )

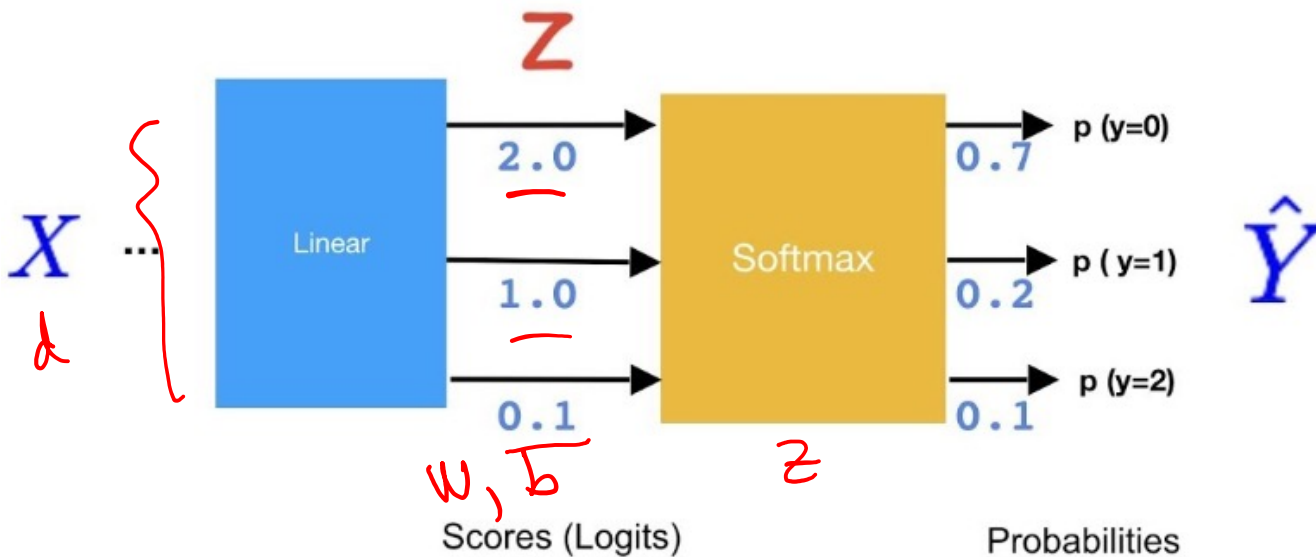
$$\begin{pmatrix} 0.1 \\ 0 \\ 0.8 \\ 0.1 \end{pmatrix}$$

Softmax

vector of prob. for each class  
Output class of max prob.

# Softmax classifier

LINEAR.



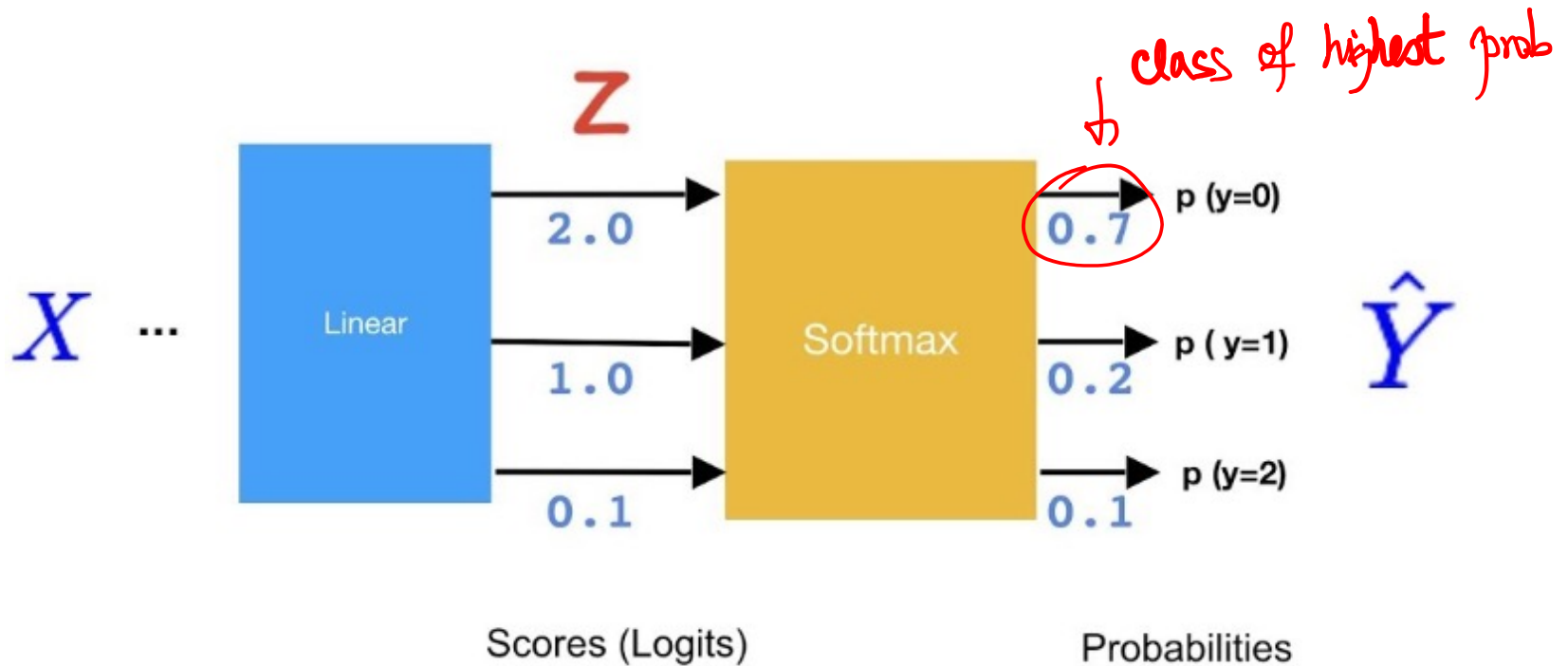
$$z = WX + b \quad \text{logits}$$

$$\begin{bmatrix} z_1 \\ \vdots \\ z_k \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix}$$

$$y_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

softmax is  
an act-function

# Softmax classifier

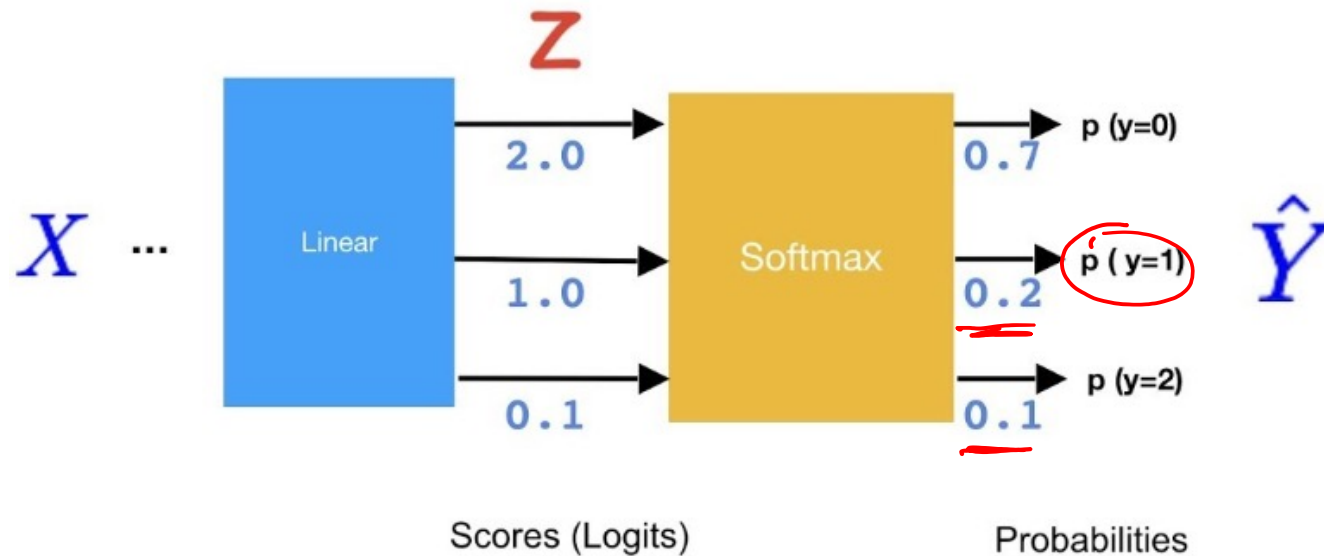


$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

✓ Softmax function  
Generalizes the sigmoid



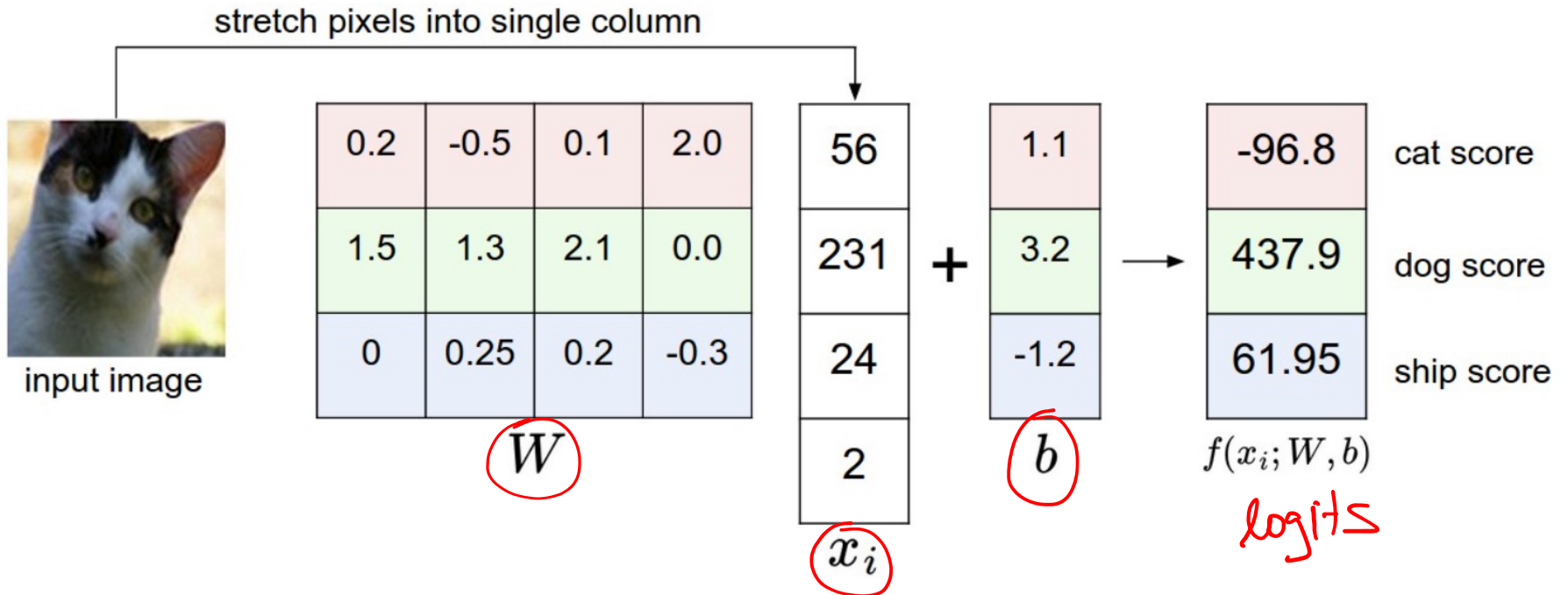
# Cross-entropy loss



$$L_j = -\log\left(\frac{y_{ij}}{\sum_{k=1}^K y_{ik}}\right)$$

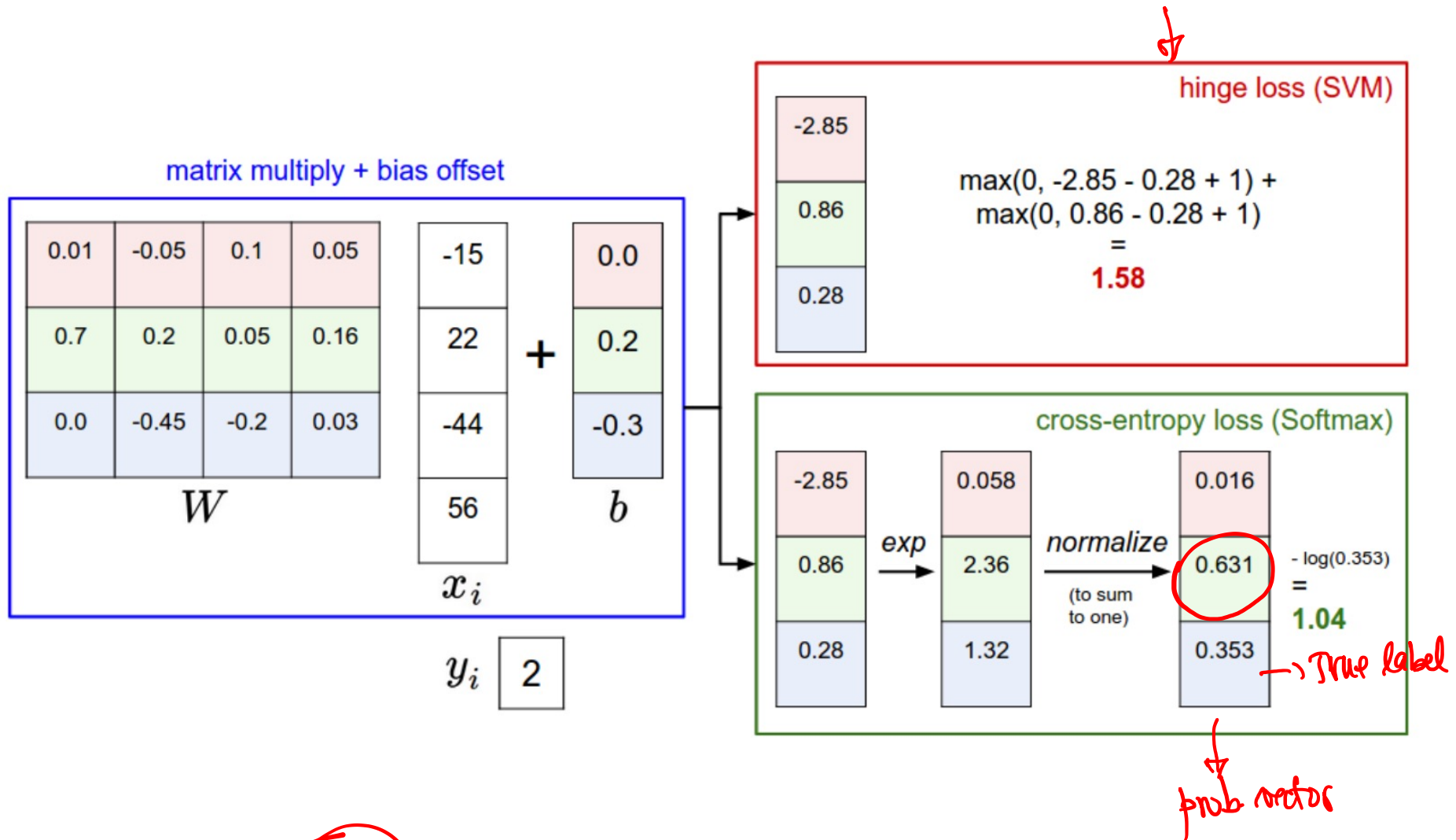
$y_{ij}$  = true label of training example

# Softmax Example



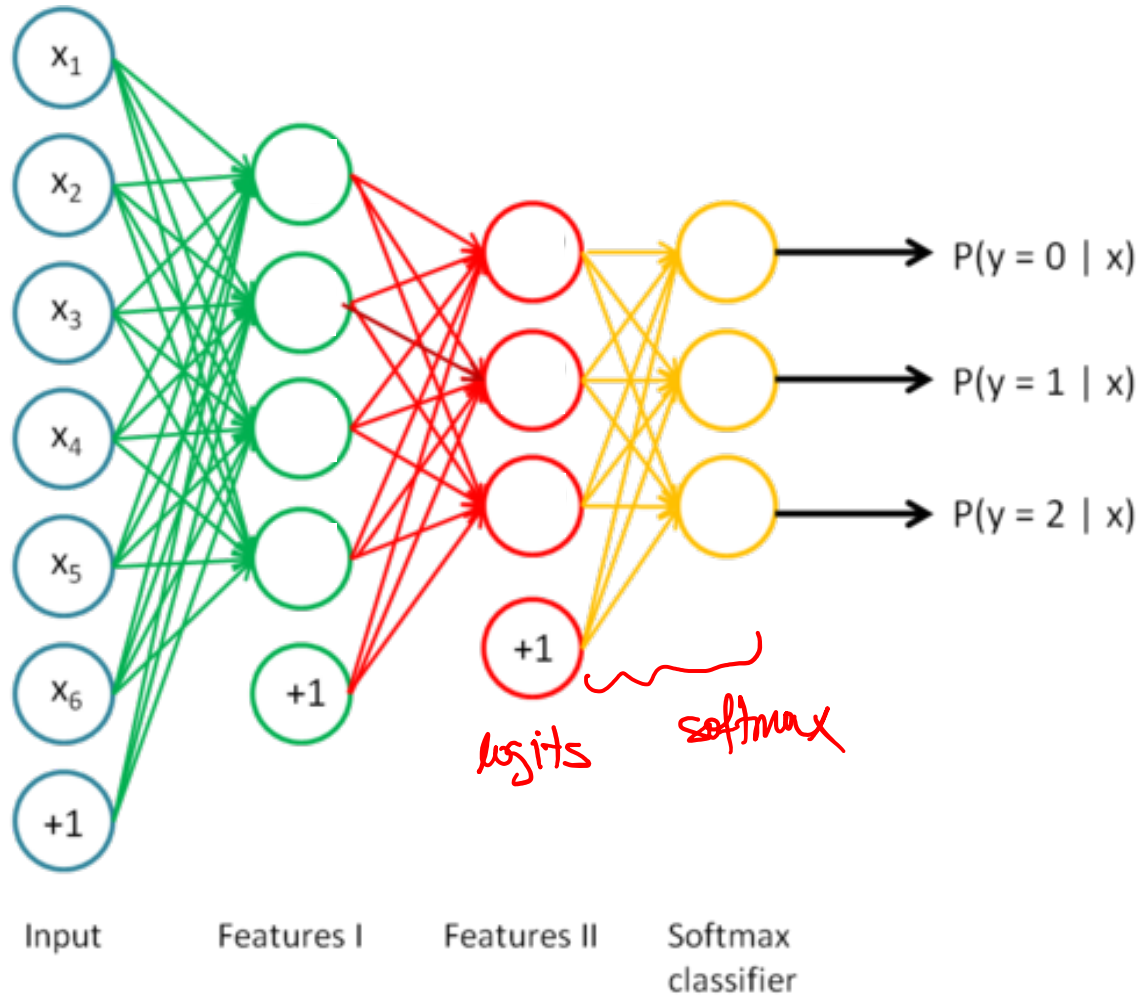


# Softmax Example

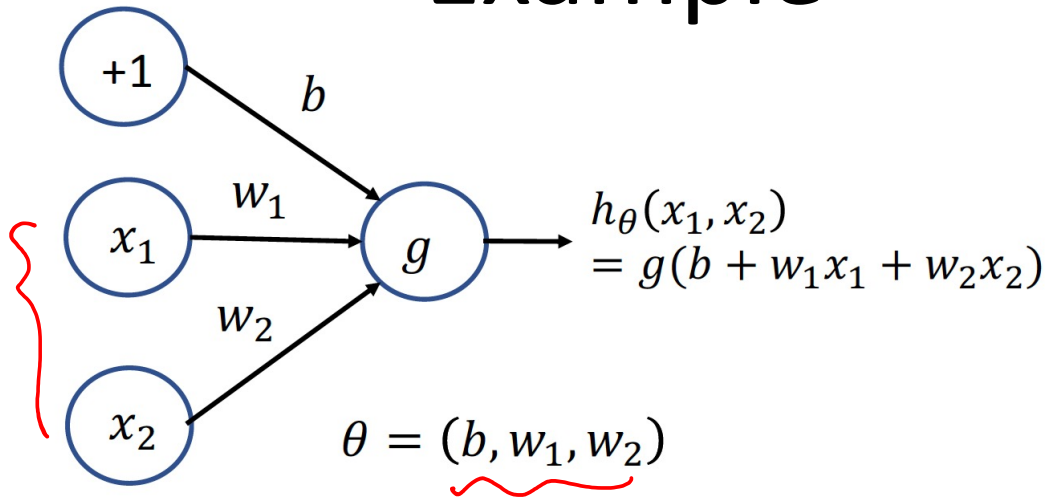


From: <https://cs231n.github.io/linear-classify/>

# Multi-class classification



# Example



1. Given  $b = -10, w_1 = 12, w_2 = 5$

Activation  $g(z) = \text{sign}(z)$

Compute the output:

2. Find out the weights  $b, w_1, w_2$  and activation function to get the following output:

$x_1$	$x_2$	$h(x_1, x_2)$
0	0	-1
0	1	-1
1	0	1
1	1	1

$\text{sign}(-10)$   
 $\text{sign}(-10+5)$   
 $\text{sign}(-10+12)=1$   
 $\text{sign}(-10+12+5)=1$

$x_1$	$x_2$	$h(x_1, x_2)$
0	0	1
0	1	1
1	0	1
1	1	0

X AND

①

$$b = +12$$

$$w_1 = -8$$

$$w_2 = -4$$

$$g(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

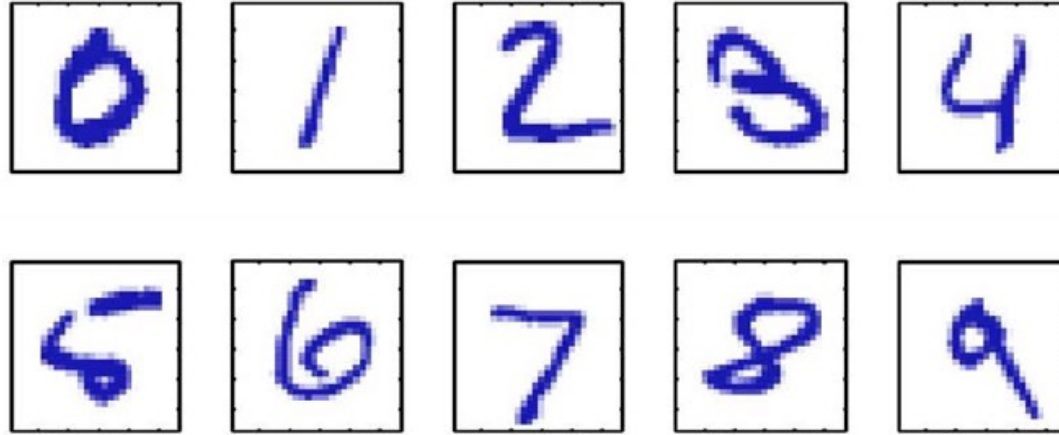
②

$$b = 1$$

$$w_1 = -0.5$$

$$w_2 = -0.5$$

# MNIST: Handwritten digit recognition



Images are 28 x 28 pixels

Represent input image as a vector  $\mathbf{x} \in \mathbb{R}^{784}$

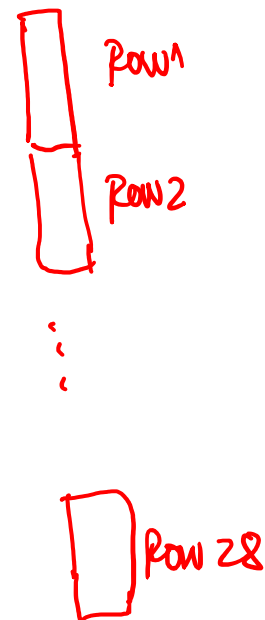
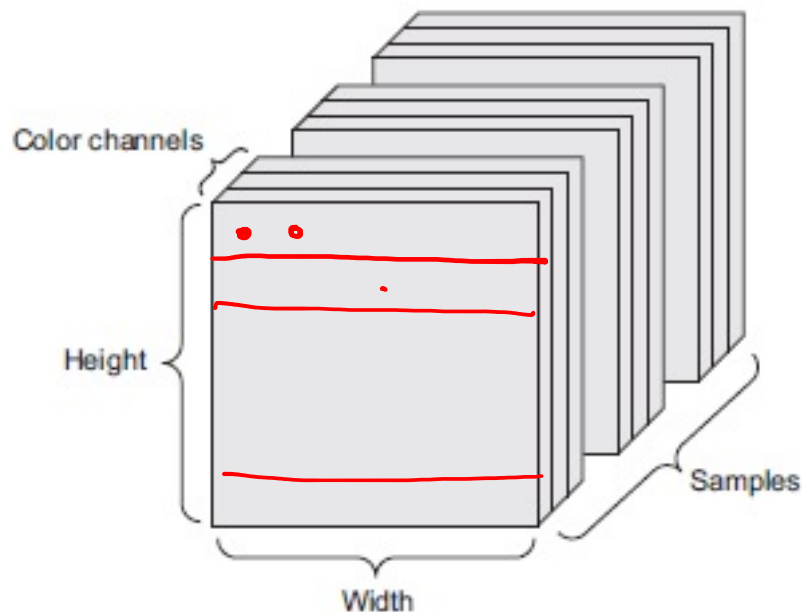
Learn a classifier  $f(\mathbf{x})$  such that,

$$f : \mathbf{x} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

Predict the digit  
Multi-class classifier

# Image Representation

- Image is 3D “tensor”: height, width, color channel (RGB)
- Black-and-white images are 2D matrices:  
height, width  
– Each value is pixel intensity  
– Pixel  $\in [0, 255]$



$$28 \times 28 = 784$$

# Lab – Feed Forward NN

```
import time
import numpy as np
#!/pip install tensorflow
#!/pip install keras

from keras.utils import np_utils
import keras.callbacks as cb
from keras.models import Sequential
from keras.layers.core import Dense, Dropout, Activation
from keras.optimizers import RMSprop
from keras.datasets import mnist

import matplotlib
import matplotlib.pyplot as plt
```

```
def load_data():
    print("Loading data")
    (X_train, y_train), (X_test, y_test) = mnist.load_data()

    X_train = X_train.astype('float32')
    X_test = X_test.astype('float32')

    # Normalize
    { X_train /= 255
      X_test /= 255 } ∈ [0,1]

    { y_train = np_utils.to_categorical(y_train, 10)
      y_test = np_utils.to_categorical(y_test, 10) }

    { X_train = np.reshape(X_train, (60000, 784))
      X_test = np.reshape(X_test, (10000, 784)) } → VECTOR

    print("Data loaded")
    return [X_train, X_test, y_train, y_test]
```

# Neural Network Architecture

```
def init_model1():  
    start_time = time.time()  
  
    print("Compiling Model")  
    model = Sequential()  
    model.add(Dense(10, input_dim=784))  
    model.add(Activation('relu'))  
  
    model.add(Dense(10))  
    model.add(Activation('softmax'))  
  
    rms = RMSprop()  
    model.compile(loss='categorical_crossentropy', optimizer=rms, metrics=['accuracy'])  
  
    print("Model finished "+format(time.time() - start_time))  
    return model
```

*Handwritten annotations:*

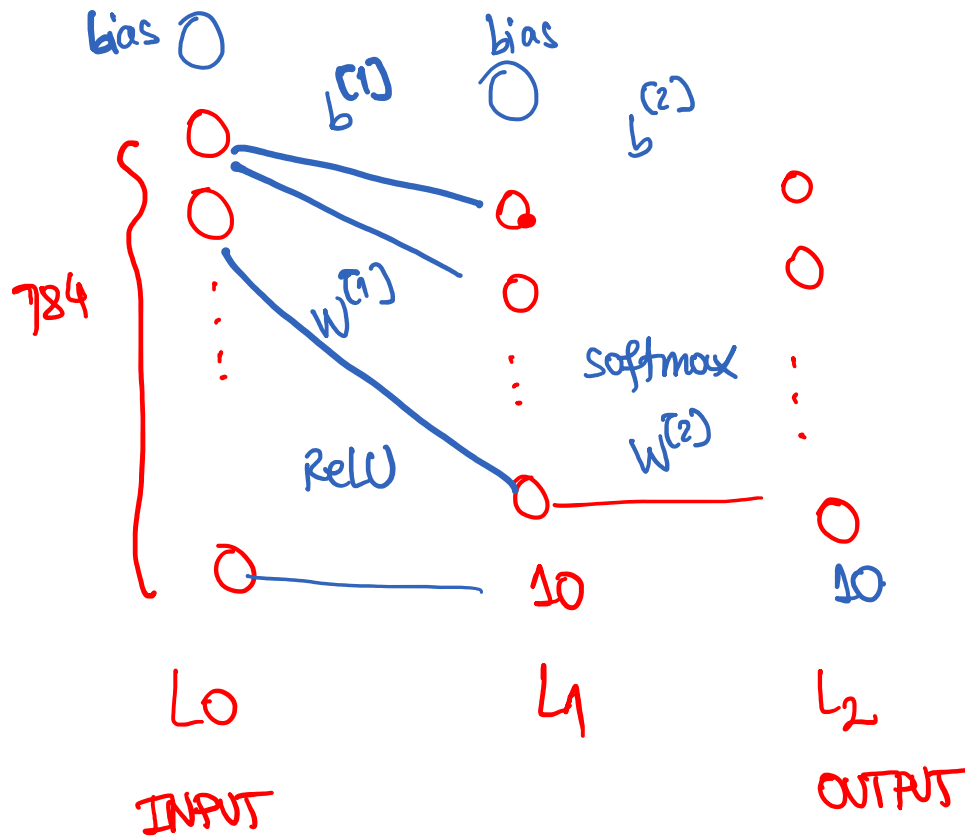
- L1** (next to the first Dense layer)
- OUTPUT L2** (next to the second Dense layer)
- FF-NN INPUT** (above the first Dense layer)
- relu** (underlined)
- rms** (circled)

## Feed-Forward Neural Network Architecture

- 1 Hidden Layer ("Dense" or Fully Connected)
- 10 neurons
- Output layer uses softmax activation



# Number of Parameters



$L_0 \rightarrow L_1$

$$\begin{aligned} W^{(1)}: 10 \times 784 &= 7840 \\ b^{(1)}: 10 \end{aligned} \} 7850$$

$L_1 \rightarrow L_2$

$$\begin{aligned} W^{(2)}: 10 \times 10 &= 100 \\ b^{(2)}: 10 \end{aligned} \} 110$$

# Number of Parameters

```
model1.summary()
```

Model: "sequential\_6"

Layer (type)	Output Shape	Param #
dense_16 (Dense)	(None, 10)	7850
activation_16 (Activation)	(None, 10)	0
dense_17 (Dense)	(None, 10)	110
activation_17 (Activation)	(None, 10)	0
Total params: 7,960		
Trainable params: 7,960		
Non-trainable params: 0		

# Train and evaluate

```
def run_network(data=None, model=None, epochs=20, batch=256):
    try:
        start_time = time.time()
        if data is None:
            X_train, X_test, y_train, y_test = load_data()
        else:
            X_train, X_test, y_train, y_test = data

        print("Training model")
        history = model.fit(X_train, y_train, epochs=epochs, batch_size=batch,
                           validation_data=(X_test, y_test), verbose=2)

        print("Training duration:" + format(time.time() - start_time))
        score = model.evaluate(X_test, y_test, batch_size=16)

        print("\nNetwork's test loss and accuracy:" + format(score))
        return history
    except KeyboardInterrupt:
        print("KeyboardInterrupt")
        return history
```

Handwritten annotations: A blue bracket on the left groups the data loading logic. A blue arrow points from the text "CV." to the `validation_data` parameter in the `model.fit` call. A blue wavy line underlines the `score` variable.

# Training/testing results

```
Compiling Model
Model finished 0.04014420509338379
Loading data
Data loaded
Training model
Epoch 1/10
235/235 - 1s - loss: 0.9142 - accuracy: 0.7501 - val_loss: 0.4398 - val_accuracy: 0.8833
Epoch 2/10
235/235 - 0s - loss: 0.3856 - accuracy: 0.8959 - val_loss: 0.3392 - val_accuracy: 0.9050
Epoch 3/10
235/235 - 0s - loss: 0.3245 - accuracy: 0.9093 - val_loss: 0.3043 - val_accuracy: 0.9141
Epoch 4/10
235/235 - 0s - loss: 0.2992 - accuracy: 0.9165 - val_loss: 0.2890 - val_accuracy: 0.9178
Epoch 5/10
235/235 - 0s - loss: 0.2853 - accuracy: 0.9202 - val_loss: 0.2797 - val_accuracy: 0.9214
Epoch 6/10
235/235 - 0s - loss: 0.2755 - accuracy: 0.9234 - val_loss: 0.2735 - val_accuracy: 0.9217
Epoch 7/10
235/235 - 0s - loss: 0.2690 - accuracy: 0.9251 - val_loss: 0.2689 - val_accuracy: 0.9252
Epoch 8/10
235/235 - 0s - loss: 0.2634 - accuracy: 0.9263 - val_loss: 0.2658 - val_accuracy: 0.9271
Epoch 9/10
235/235 - 0s - loss: 0.2590 - accuracy: 0.9276 - val_loss: 0.2666 - val_accuracy: 0.9257
Epoch 10/10
235/235 - 0s - loss: 0.2554 - accuracy: 0.9284 - val_loss: 0.2616 - val_accuracy: 0.9284
Training duration: 3.1347107887268066 sec
625/625 [=====] - 1s 707us/step - loss: 0.2616 - accuracy: 0.9284
```

→ Network's test loss and accuracy: [0.2615792751312256, 0.9283999800682068]

loss                      Acc

# Training/testing results

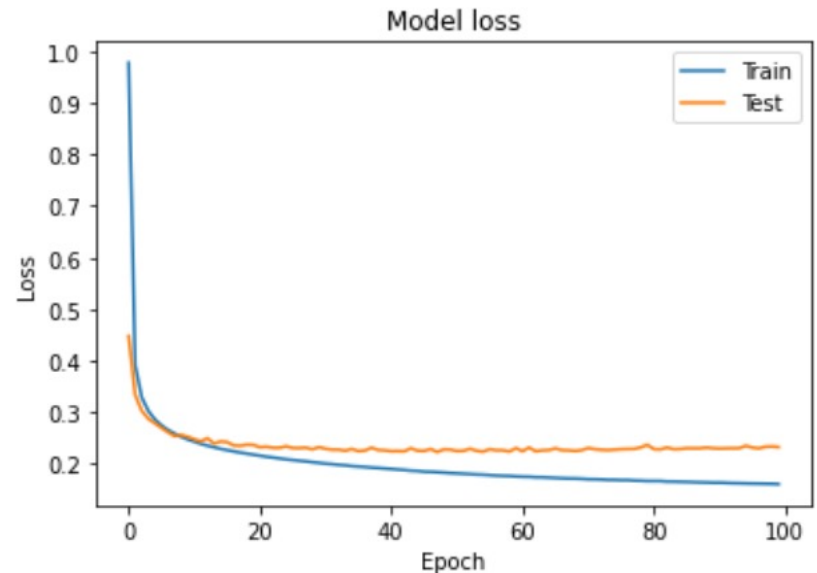
```
Epoch 98/100
235/235 - 0s - loss: 0.1611 - accuracy: 0.9552 - val_loss: 0.2329 - val_accuracy: 0.9411
Epoch 99/100
235/235 - 0s - loss: 0.1609 - accuracy: 0.9550 - val_loss: 0.2334 - val_accuracy: 0.9392
Epoch 100/100
235/235 - 0s - loss: 0.1603 - accuracy: 0.9550 - val_loss: 0.2323 - val_accuracy: 0.9401
Training duration: 22.163272857666016 sec.
625/625 [=====] - 1s 711us/step - loss: 0.2323 - accuracy: 0.9401

Network's test loss and accuracy:[0.23233847320079803, 0.9401000142097473]
```

# Monitor Loss

```
def plot_losses(hist):  
    plt.plot(hist.history['loss']) → Train  
    plt.plot(hist.history['val_loss']) → Val.  
    plt.title('Model loss')  
    plt.ylabel('Loss')  
    plt.xlabel('Epoch')  
    plt.legend(['Train', 'Test'], loc='upper right')  
    plt.show()
```

```
modell = init_model1()  
history1 = run_network(model = modell, epochs=100)  
plot_losses(history1)
```



# Review

- Feed-Forward Neural Networks are the common neural networks architectures
  - Fully connected networks are called Multi-Layer Perceptron
- Input, output, and hidden layers
  - Linear matrix operations followed by non-linear activations at every layer
- Activations:
  - ReLU, tanh, etc., for hidden layers
  - Sigmoid (binary classification) and softmax (for multi-class classification) at last layer
- Forward propagation: process of evaluating input through the network

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
  - Andrew Moore
  - Yann LeCun
- Thanks!