

27.09.2024

# STATISTICAL METHODS IN AI (CS7.403)

## Lecture-14: Neural Networks-2

Ravi Kiran ([ravi.kiran@iiit.ac.in](mailto:ravi.kiran@iiit.ac.in))

<https://ravika.github.io>

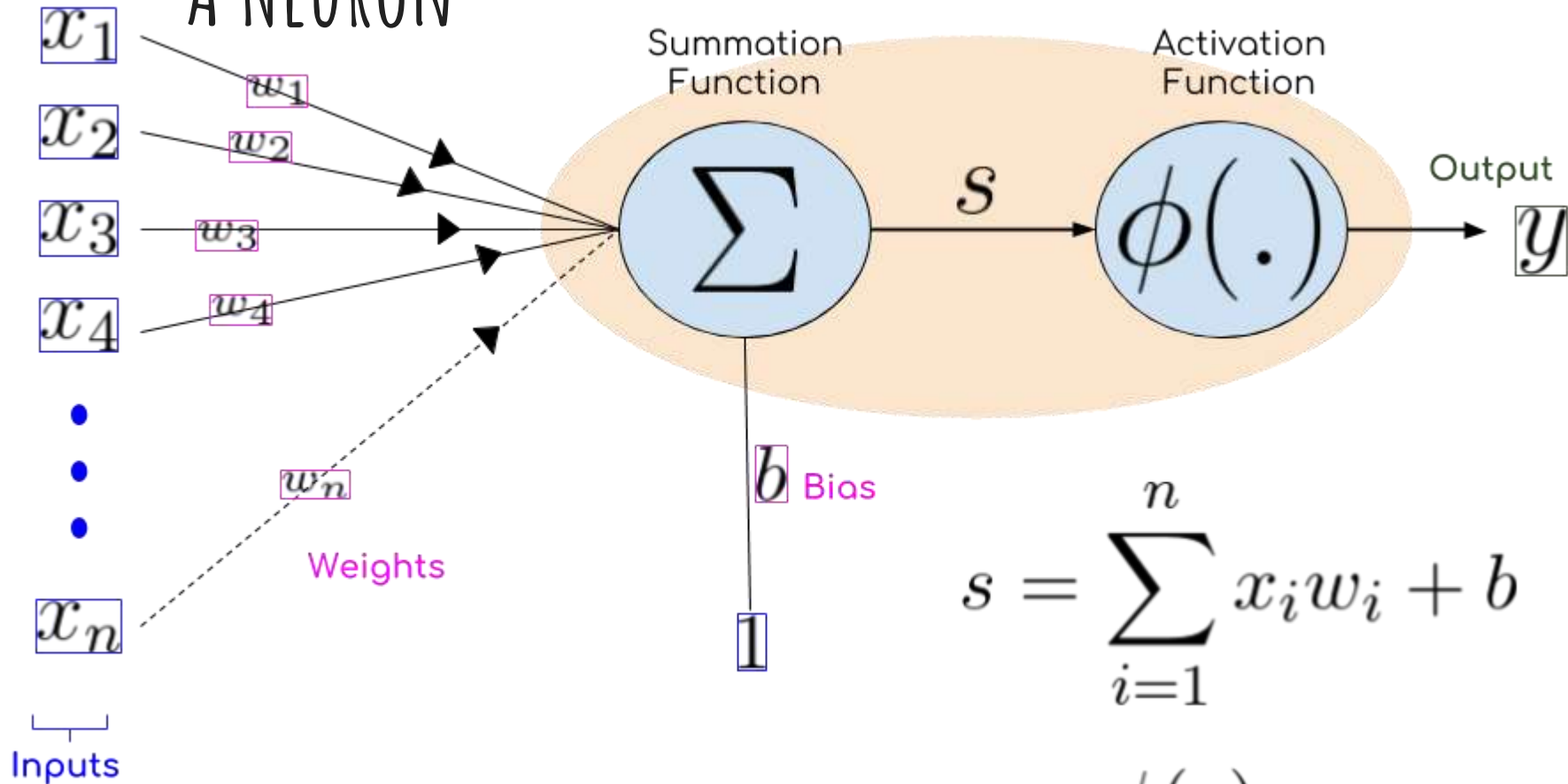


@vikataravi



Center for Visual Information Technology (CVIT)  
IIIT Hyderabad

# A NEURON



$$s = \sum_{i=1}^n x_i w_i + b$$

$$y = \phi(s)$$

# The Perceptron Cost Function

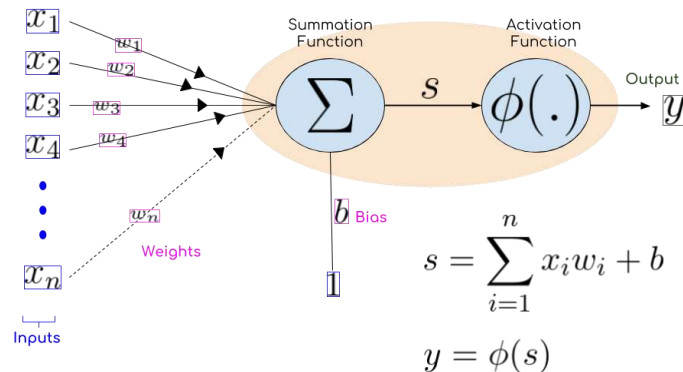
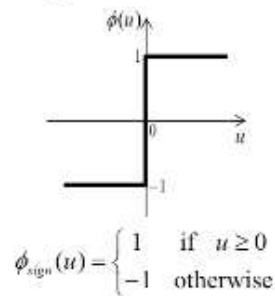
- Prediction is correct if  $y^{(i)} \mathbf{x}^{(i)} \boldsymbol{\theta} > 0$

- 0/1 loss

$$J_{0/1}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(\text{sign}(\mathbf{x}^{(i)} \boldsymbol{\theta}), y^{(i)})$$

where  $\ell()$  is 0 if the prediction is correct, 1 otherwise

sign function



# The Perceptron Cost Function

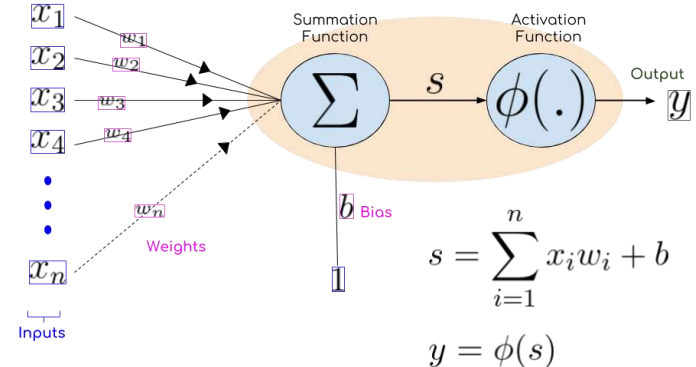
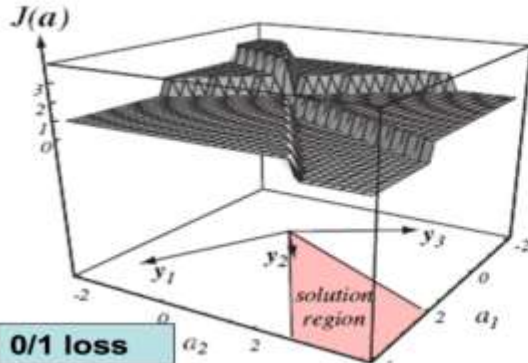
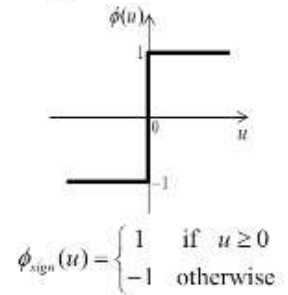
- Prediction is correct if  $y^{(i)} \mathbf{x}^{(i)} \boldsymbol{\theta} > 0$

- 0/1 loss

$$J_{0/1}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(\text{sign}(\mathbf{x}^{(i)} \boldsymbol{\theta}), y^{(i)})$$

where  $\ell()$  is 0 if the prediction is correct, 1 otherwise

sign function



# The Perceptron Cost Function

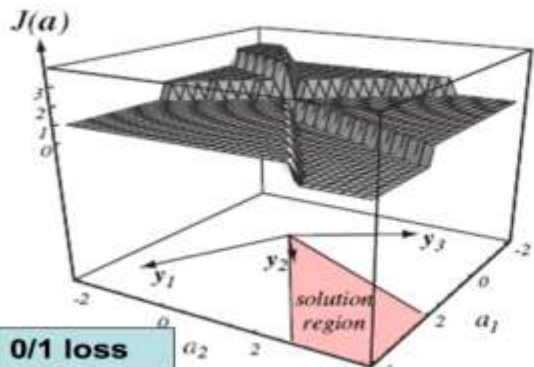
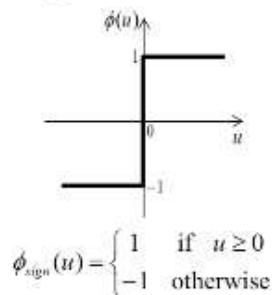
- Prediction is correct if  $y^{(i)} \mathbf{x}^{(i)} \boldsymbol{\theta} > 0$

- 0/1 loss

$$J_{0/1}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(\text{sign}(\mathbf{x}^{(i)} \boldsymbol{\theta}), y^{(i)})$$

where  $\ell()$  is 0 if the prediction is correct, 1 otherwise

sign function



Doesn't produce a useful gradient



No gradient  $\rightarrow$  We cannot use gradient-based optimization methods. Does NOT mean 'unsolvable'

# Improving the perceptron

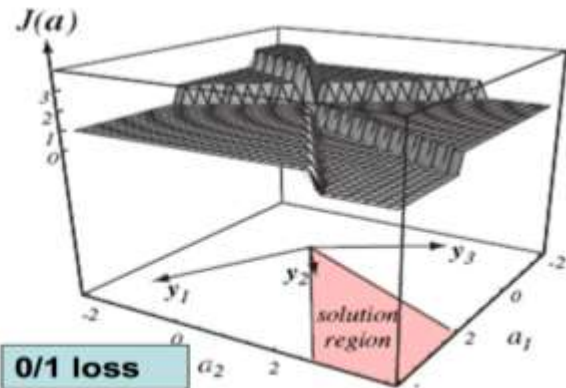
Problem #0: Non-differentiable loss function



- Could have used 0/1 loss

$$J_{0/1}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(\text{sign}(x^{(i)} \boldsymbol{\theta}), y^{(i)})$$

where  $\ell()$  is 0 if the prediction is correct, 1 otherwise



0/1 loss

Doesn't produce a useful gradient

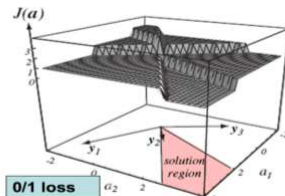
# Improving the perceptron

Problem #0: Non-differentiable loss function

- Could have used 0/1 loss

$$J_{0/1}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\text{sign}(x^{(i)}\theta), y^{(i)})$$

where  $\ell()$  is 0 if the prediction is correct, 1 otherwise



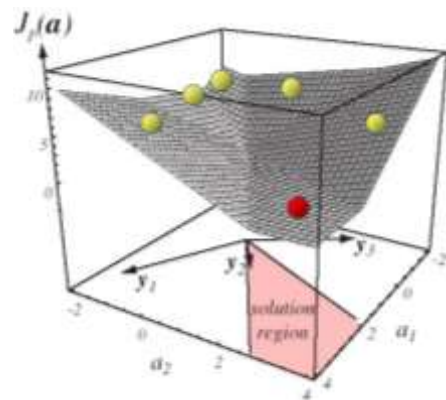
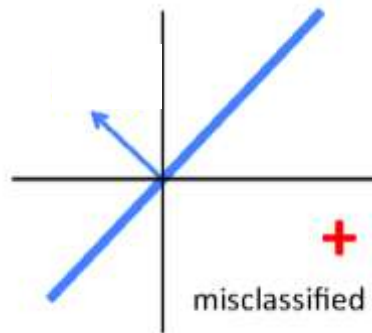
Doesn't produce a useful gradient



Solution: Minimize the 'misclassification distance'

$$J_p(\mathbf{a}) = \sum_{\mathbf{x} \in \mathbb{M}} (-\mathbf{a}^T \mathbf{x})$$

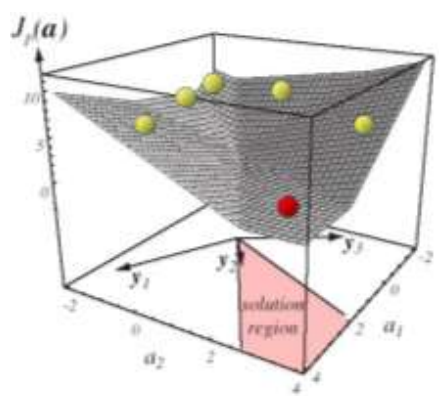
$\mathbb{M} \rightarrow$  Set of misclassified samples



Solution: Minimize the ‘misclassification distance’

$$J_p(\mathbf{a}) = \sum_{\mathbf{x} \in \mathbb{M}} (-\mathbf{a}^T \mathbf{x})$$

$\mathbb{M} \rightarrow$  Set of misclassified samples



$$\nabla J_p(\mathbf{a}) = \sum_{\mathbf{x} \in \mathbb{M}} (-\mathbf{x})$$

$$\mathbf{a}^{(k+1)} = \mathbf{a}^k - \eta^k \nabla J_p(\mathbf{a})$$

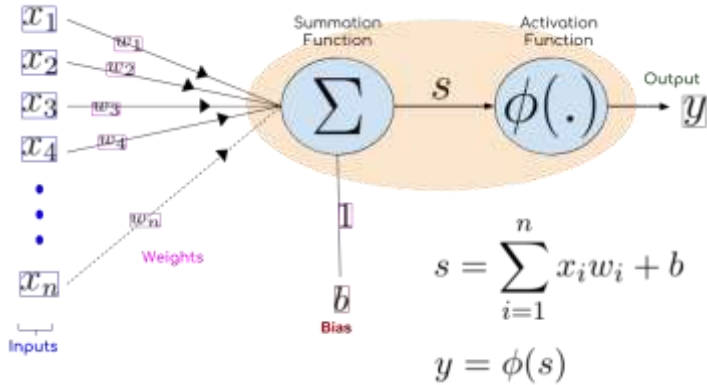
$$\mathbf{a}^{(k+1)} = \mathbf{a}^k + \eta^k \sum_{\mathbf{x} \in \mathbb{M}} \mathbf{x}$$

c.f. logistic regression

$$\nabla_a J(\hat{y}^{(i)}, y^{(i)}) = (\hat{y}^{(i)} - y^{(i)}) \hat{x}^{(i)}$$

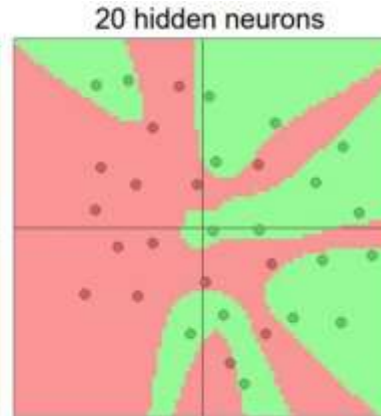
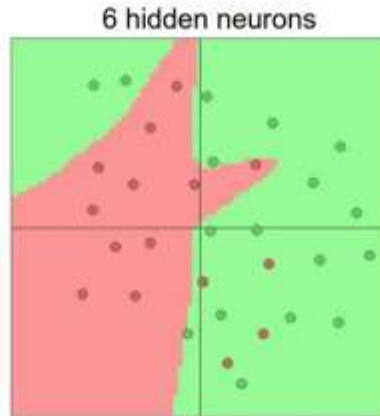
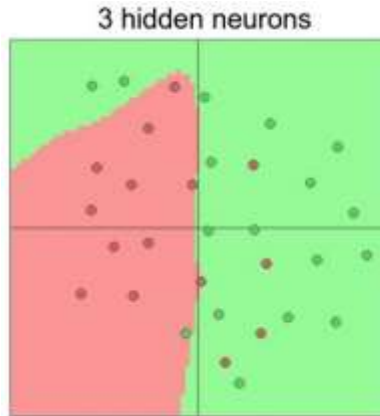
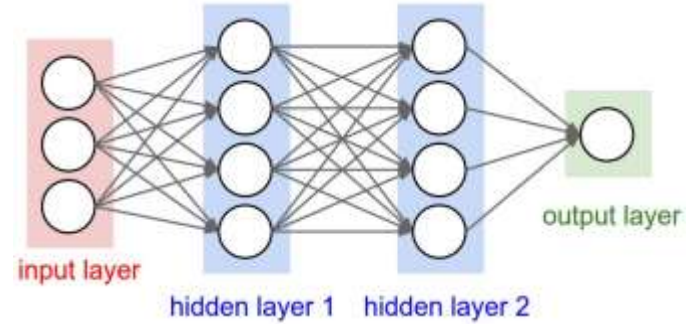
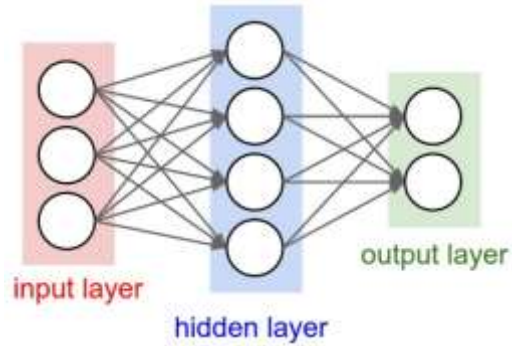


# ACTIVATION FUNCTIONS

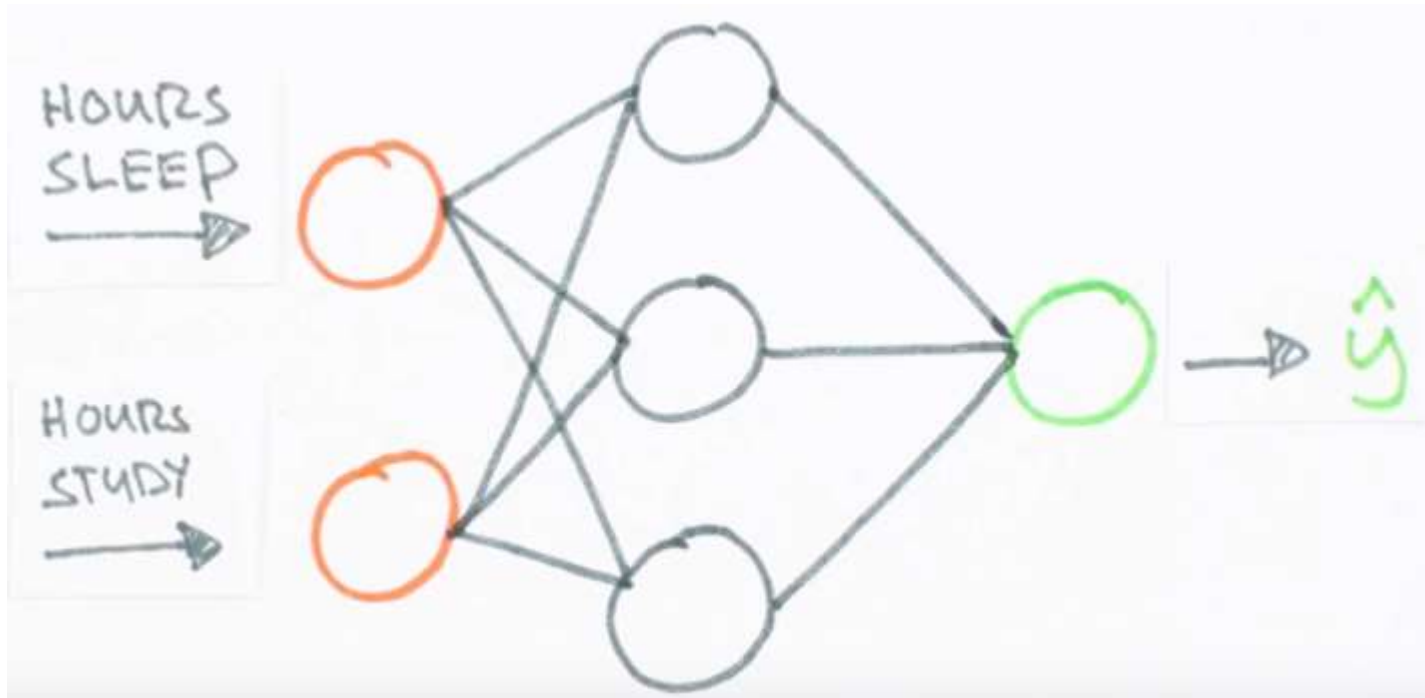


Activation Function	Equation	Example	1D Graph
Linear	$\phi(z) = z$	Adaline, linear regression	
Unit Step (Heaviside Function)	$\phi(z) = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$	Perceptron variant	
Sign (signum)	$\phi(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ 1 & z > 0 \end{cases}$	Perceptron variant	
Piece-wise Linear	$\phi(z) = \begin{cases} 0 & z \leq -1/2 \\ z + 1/2 & -1/2 \leq z \leq 1/2 \\ 1 & z \geq 1/2 \end{cases}$	Support vector machine	
Logistic (sigmoid)	$\phi(z) = \frac{1}{1 + e^{-z}}$	Logistic regression,	
Hyperbolic Tangent (tanh)	$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$		
ReLU	$\phi(z) = \begin{cases} 0 & z < 0 \\ z & z > 0 \end{cases}$		

# WHY USE ONLY ONE NEURON ?

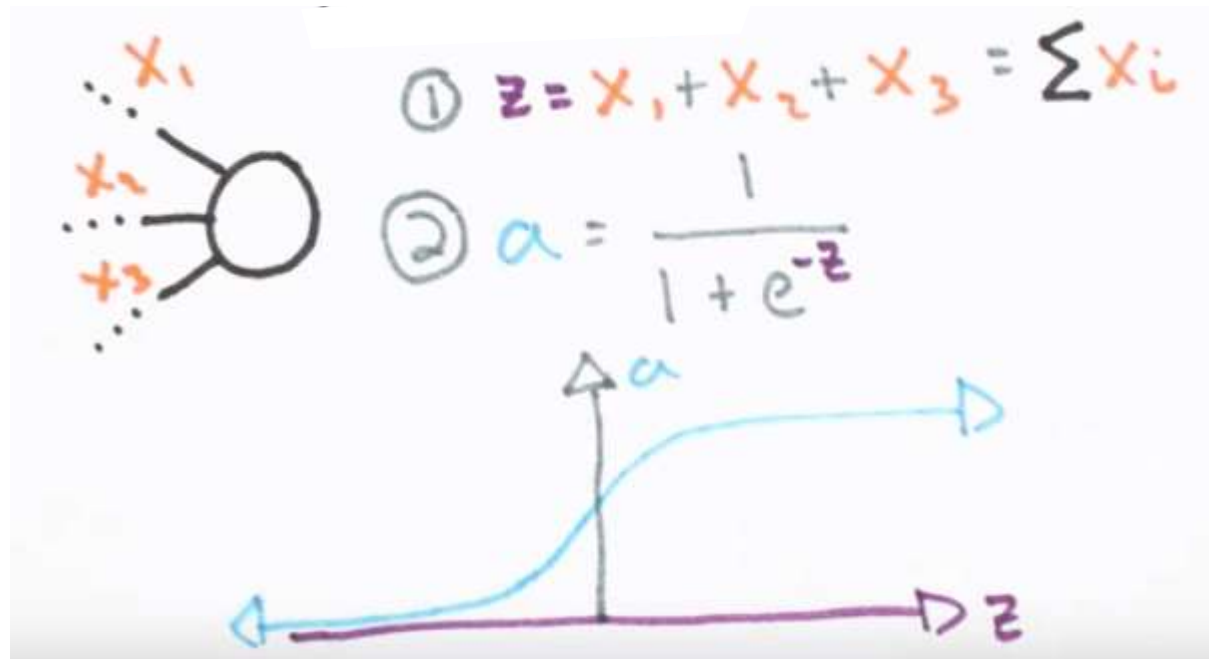


# MULTI-NEURON NETWORKS :: ARCHITECTURE



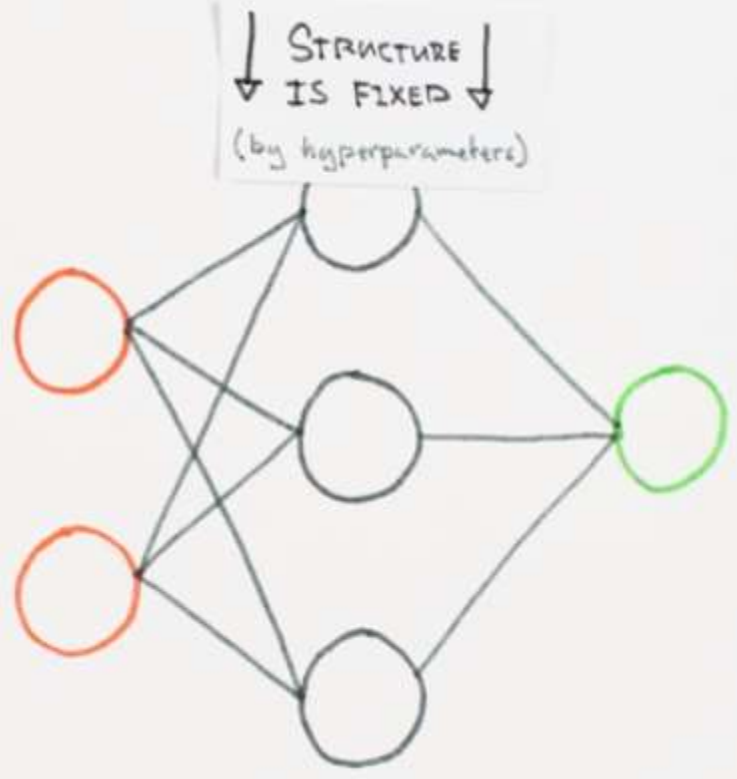
# MULTI-NEURON NETWORKS :: ARCHITECTURE

NEURON



# MULTI-NEURON NETWORKS :: ARCHITECTURE

```
class Neural_Network(object):  
    def __init__(self):  
        #Define Hyperparameters  
        self.inputLayerSize = 2  
        self.outputLayerSize = 1  
        self.hiddenLayerSize = 3  
  
        #Weights (parameters)  
        self.W1 = np.random.randn(self.inputLayerSize, self.hiddenLayerSize)  
        self.W2 = np.random.randn(self.hiddenLayerSize, self.outputLayerSize)
```



# MULTI-NEURON NETWORKS :: TRAINING

INITIALIZE NETWORK WITH RANDOM WEIGHTS

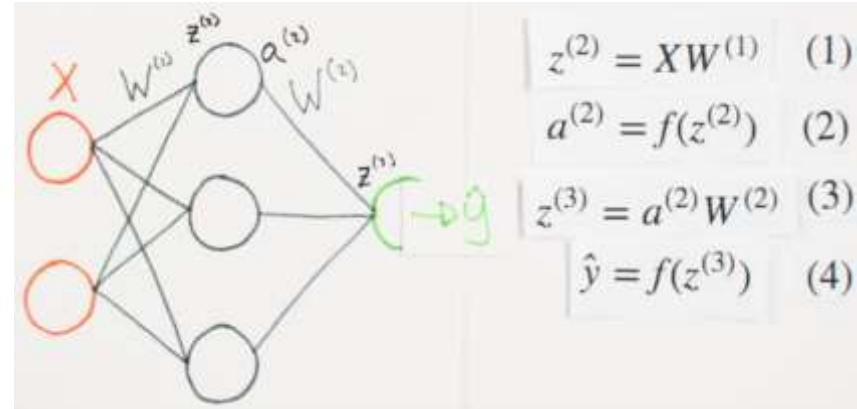
WHILE [NOT CONVERGED]

DO FORWARD PROP

DO BACKPROP AND DETERMINE CHANGE IN WEIGHTS

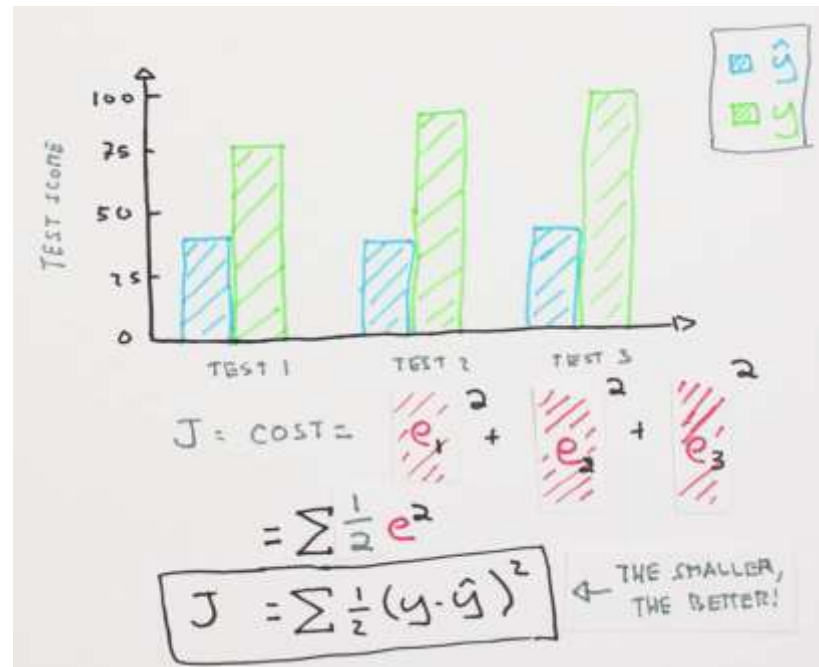
UPDATE ALL WEIGHTS IN ALL LAYERS

# MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



```
def forward(self, X):  
    #Propagate inputs though network  
    self.z2 = np.dot(X, self.W1) # z2 = X * W1  
    self.a2 = self.sigmoid(self.z2) # a2 = sigmoid(z2)  
    self.z3 = np.dot(self.a2, self.W2) # z3 = a2 * W2  
    yHat = self.sigmoid(self.z3) # yHat = sigmoid(z3)  
    return yHat  
  
def sigmoid(self, z):  
    #Apply sigmoid activation function to scalar, vector, or matrix  
    return 1/(1+np.exp(-z))
```

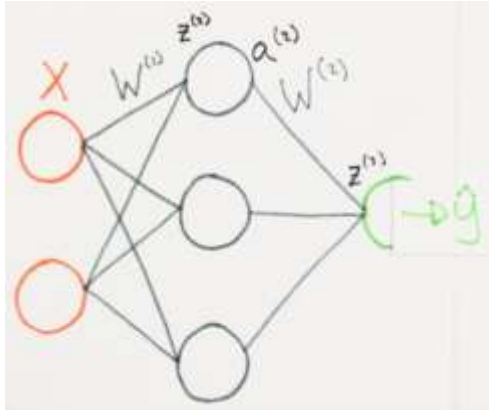
# MULTI-NEURON NETWORKS :: GRADIENT DESCENT



```
def costFunction(self, X, y):  
    #Compute cost for given X,y, use weights already stored in class.  
    self.yHat = self.forward(X)  
    J = 0.5*sum((y-self.yHat)**2)  
    return J
```



# MULTI-NEURON NETWORKS :: BACKPROPAGATION



$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

↑ HOW DOES THIS CHANGE IF I CHANGE THESE? ↑

$$\frac{\partial J}{\partial W}$$

$$W^{(1)} = \begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{bmatrix}$$

$$\frac{\partial J}{\partial W^{(1)}} = \begin{bmatrix} \frac{\partial J}{\partial W_{11}^{(1)}} & \frac{\partial J}{\partial W_{12}^{(1)}} & \frac{\partial J}{\partial W_{13}^{(1)}} \\ \frac{\partial J}{\partial W_{21}^{(1)}} & \frac{\partial J}{\partial W_{22}^{(1)}} & \frac{\partial J}{\partial W_{23}^{(1)}} \end{bmatrix}$$

$$\frac{\partial J}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial J}{\partial W_{11}^{(2)}} \\ \frac{\partial J}{\partial W_{21}^{(2)}} \\ \frac{\partial J}{\partial W_{31}^{(2)}} \end{bmatrix}$$

# MULTI-NEURON NETWORKS :: BACKPROPAGATION

$$J = \sum \frac{1}{2} (y - \hat{y})^2$$

$$\frac{\partial J}{\partial W^{(2)}} = (a^{(2)})^T \delta^{(3)}$$

$$\delta^{(3)} = -(y - \hat{y})f'(z^{(3)})$$

$$\frac{\partial J}{\partial W^{(1)}} = X^T \delta^{(2)}$$

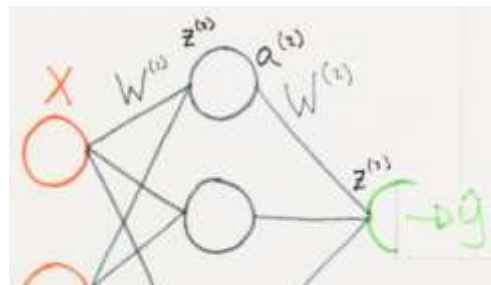
$$\delta^{(2)} = \delta^{(3)} (W^{(2)})^T f'(z^{(2)})$$

$$\begin{aligned} z^{(2)} &= XW^{(1)} & (1) \\ a^{(2)} &= f(z^{(2)}) & (2) \\ z^{(3)} &= a^{(2)}W^{(2)} & (3) \\ \hat{y} &= f(z^{(3)}) & (4) \end{aligned}$$

$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

How does this change if I change these?

$$\frac{\partial J}{\partial W}$$



```
# backpropagation
def costFunctionPrime(self, X, y):
    #Compute derivative with respect to W1 and W2 for a given X and y:
    self.yHat = self.forward(X)

    delta3 = np.multiply(-(y-self.yHat), self.sigmoidPrime(self.z3))
    dJdW2 = np.dot(self.a2.T, delta3)

    delta2 = np.dot(delta3, self.W2.T)*self.sigmoidPrime(self.z2)
    dJdW1 = np.dot(X.T, delta2)

    return dJdW1, dJdW2
```

$$\left[ \begin{matrix} \frac{\partial J}{\partial W_1} \\ \frac{\partial J}{\partial W_2} \end{matrix} \right]$$

# MULTI-NEURON NETWORKS :: TRAINING

INITIALIZE NETWORK WITH RANDOM WEIGHTS

WHILE [NOT CONVERGED]

DO FORWARD PROP

DO BACKPROP AND DETERMINE CHANGE IN WEIGHTS

UPDATE ALL WEIGHTS IN ALL LAYERS

```
NN = Neural_Network()  
cost1 = NN.costFunction(X, y)  
print('cost1=',cost1)  
dJdw1, dJdw2 = NN.costFunctionPrime(X, y)  
print('dJ/dw1=',dJdw1)  
print('dJ/dw2=',dJdw2)  
eta = 0.01  
NN.W1 = NN.W1 - eta * dJdw1  
NN.W2 = NN.W2 - eta * dJdw2
```

One  
Iteration

# MULTI-NEURON NETWORKS :: BACKPROPAGATION

$$J = \sum \frac{1}{2} (y - \hat{y})^2$$

$$\frac{\partial J}{\partial W^{(2)}} = (a^{(2)})^T \delta^{(3)}$$

$$\delta^{(3)} = -(y - \hat{y})f'(z^{(3)})$$

$$\frac{\partial J}{\partial W^{(1)}} = X^T \delta^{(2)}$$

$$\delta^{(2)} = \delta^{(3)}(W^{(2)})^T f'(z^{(2)})$$

```

NN = Neural_Network()
cost1 = NN.costFunction(X, y)
print('cost1=', cost1)
dJdW1, dJdW2 = NN.costFunctionPrime(X, y)
print('dJ/dW1=', dJdW1)
print('dJ/dW2=', dJdW2)
eta = 0.01
NN.W1 = NN.W1 - eta * dJdW1
NN.W2 = NN.W2 - eta * dJdW2
cost2 = NN.costFunction(X, y)
print('cost2=', cost2)

cost1= [0.44735371]
dJ/dW1= [[-0.08913117 -0.04750461 -0.00562623]
          [-0.05862425 -0.03130539 -0.00351033]]
dJ/dW2= [[-0.4110688 ]
          [-0.37530217]
          [-0.4590466 ]]
cost2= [0.44202336]
    
```

$$z^{(2)} = XW^{(1)} \quad (1)$$

$$a^{(2)} = f(z^{(2)}) \quad (2)$$

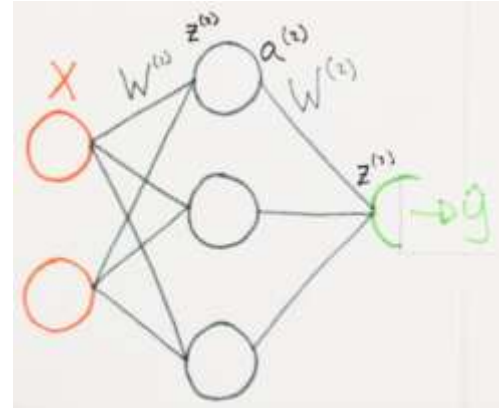
$$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$$

$$\hat{y} = f(z^{(3)}) \quad (4)$$

$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

How does this change if I change these?

$$\frac{\partial J}{\partial W}$$



$$W^{(1)} = \begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \end{bmatrix} \quad \frac{\partial J}{\partial W^{(1)}} = \begin{bmatrix} \frac{\partial J}{\partial W_{11}^{(1)}} & \frac{\partial J}{\partial W_{12}^{(1)}} & \frac{\partial J}{\partial W_{13}^{(1)}} \\ \frac{\partial J}{\partial W_{21}^{(1)}} & \frac{\partial J}{\partial W_{22}^{(1)}} & \frac{\partial J}{\partial W_{23}^{(1)}} \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{bmatrix} \quad \frac{\partial J}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial J}{\partial W_{11}^{(2)}} \\ \frac{\partial J}{\partial W_{21}^{(2)}} \\ \frac{\partial J}{\partial W_{31}^{(2)}} \end{bmatrix}$$

# MULTI-NEURON NETWORK – NUMERICAL GRADIENT CHECKING

# MULTI-NEURON NETWORK – NUMERICAL GRADIENT CHECKING

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x+\Delta x) - f(x)}{\Delta x}$$

$$f(x) = x^2$$

# MULTI-NEURON NETWORK – NUMERICAL GRADIENT CHECKING

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x+\Delta x) - f(x)}{\Delta x}$$

$$f(x) = x^2$$

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{(x+\Delta x)^2 - x^2}{\Delta x}$$

$$= \lim_{\Delta x \rightarrow 0} \frac{x^2 + 2x\Delta x + \Delta x^2 - x^2}{\Delta x}$$

$$= \lim_{\Delta x \rightarrow 0} \frac{\Delta x^2 + 2x\Delta x}{\Delta x}$$

$$= \lim_{\Delta x \rightarrow 0} \Delta x + 2x = \boxed{2x}$$

# MULTI-NEURON NETWORK – NUMERICAL GRADIENT CHECKING

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x+\Delta x) - f(x)}{\Delta x}$$

$$f(x) = x^2$$

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{(x+\Delta x)^2 - x^2}{\Delta x}$$

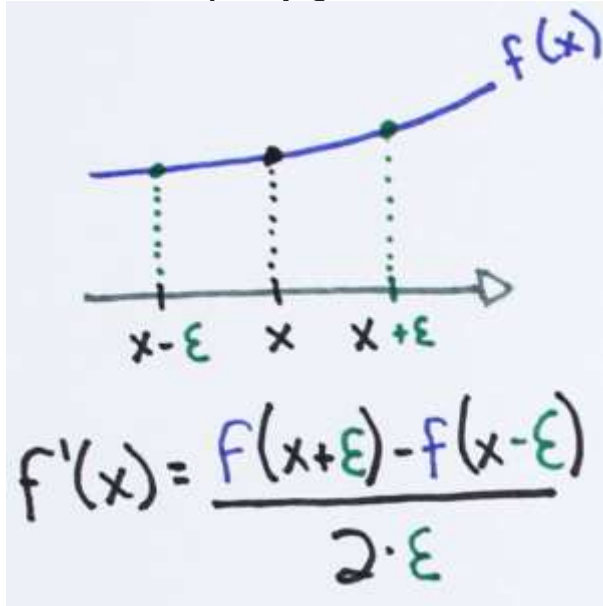
$$= \lim_{\Delta x \rightarrow 0} \frac{x^2 + 2x\Delta x + \Delta x^2 - x^2}{\Delta x}$$

$$= \lim_{\Delta x \rightarrow 0} \frac{\Delta x^2 + 2x\Delta x}{\Delta x}$$

$$= \lim_{\Delta x \rightarrow 0} \Delta x + 2x = \boxed{2x}$$



# MULTI-NEURON NETWORK – NUMERICAL GRADIENT



# MULTI-NEURON NETWORK – NUMERICAL GRADIENT



$$f'(x) = \frac{f(x+\epsilon) - f(x-\epsilon)}{2 \cdot \epsilon}$$

```
In [4]: def f(x):  
        return x**2
```

```
In [5]: epsilon = 1e-4  
        x = 1.5
```

```
In [6]: numericGradient = (f(x+epsilon) - f(x-epsilon)) / (2*epsilon)
```

```
In [7]: numericGradient, 2*x
```

```
Out[7]: (2.9999999999996696, 3.0)
```

```
In [ ]: |
```

# MULTI-NEURON NETWORK – NUMERICAL GRADIENT CHECKING

**Parameter vector  $\theta$**

→  $\theta \in \mathbb{R}^n$  (E.g.  $\theta$  is “unrolled” version of  $\Theta^{(1)}$ ,  $\Theta^{(2)}$ ,  $\Theta^{(3)}$ )

→  $\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$

$$\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$$

⋮

# MULTI-NEURON NETWORK – NUMERICAL GRADIENT CHECKING



$$f'(x) = \frac{f(x+\epsilon) - f(x-\epsilon)}{2 \cdot \epsilon}$$

```
In [3]: numgrad = computeNumericalGradient(NN, X, y)
```

```
In [4]: grad = NN.computeGradients(X, y)
```

```
In [5]: numgrad
```

```
Out[5]: array([ -7.10752568e-03,  -6.30194392e-03,  -4.96392693e-03,
                -6.55946987e-04,   7.57595597e-05,  -1.11297012e-03,
                -8.81243102e-03,  -4.45550176e-03,  -1.93471143e-02])
```

```
In [6]: grad
```

```
Out[6]: array([ -7.10752569e-03,  -6.30194393e-03,  -4.96392693e-03,
                -6.55946985e-04,   7.57595604e-05,  -1.11297012e-03,
                -8.81243100e-03,  -4.45550176e-03,  -1.93471142e-02])
```

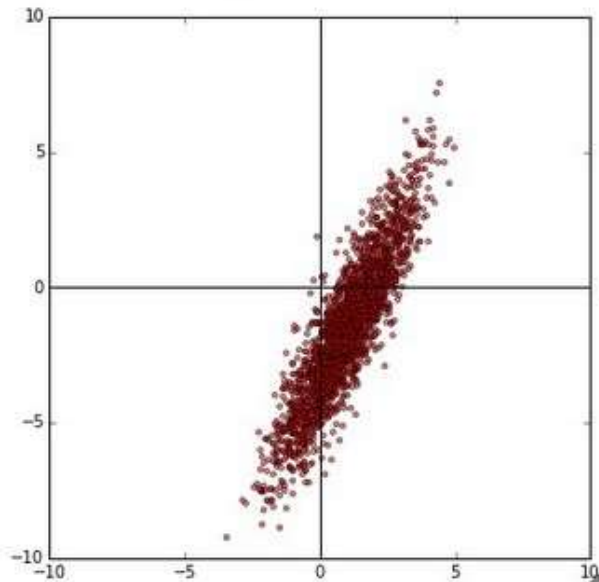
```
In [7]: norm(grad-numgrad)/norm(grad+numgrad)
```

```
Out[7]: 1.9824969610227768e-09
```

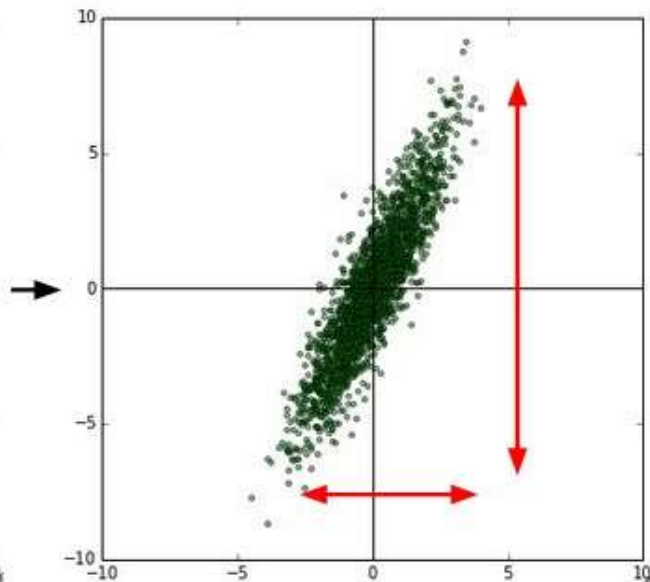
# DATA SETUP

- Preprocessing:

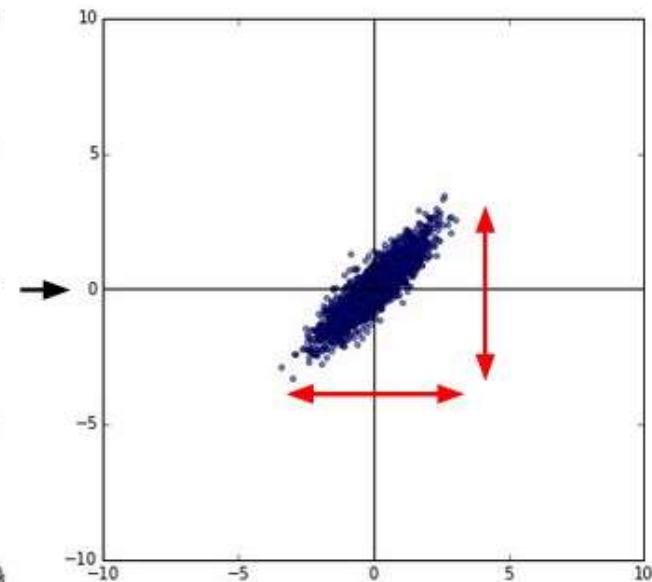
original data



zero-centered data



normalized data

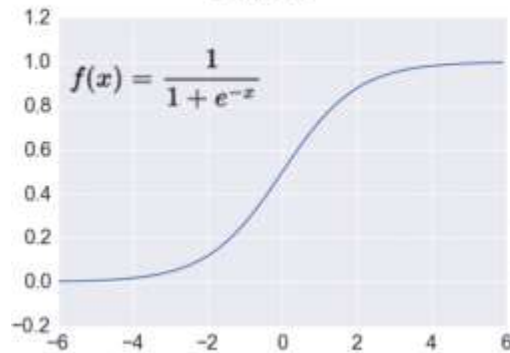


# WEIGHT INITIALIZATION

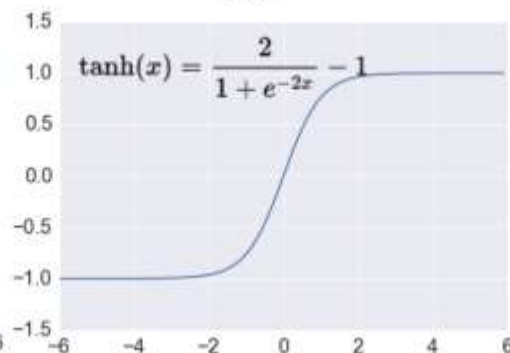
- ALL ZEROS
- RANDOM  $[0,1]$
- RANDOM  $[-1,1]$
- $w = \text{np.random.randn}(n) * \text{sqrt}(2.0/n)$ ,  $n = \#$  of inputs to neuron

# ACTIVATION FUNCTIONS

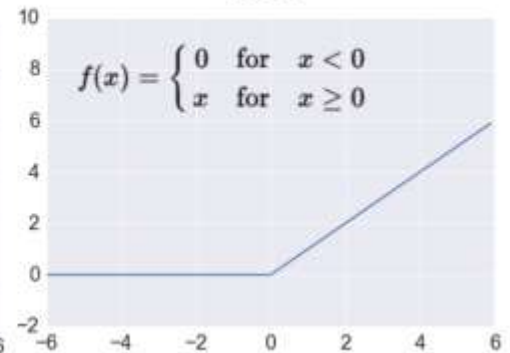
Sigmoid



TanH

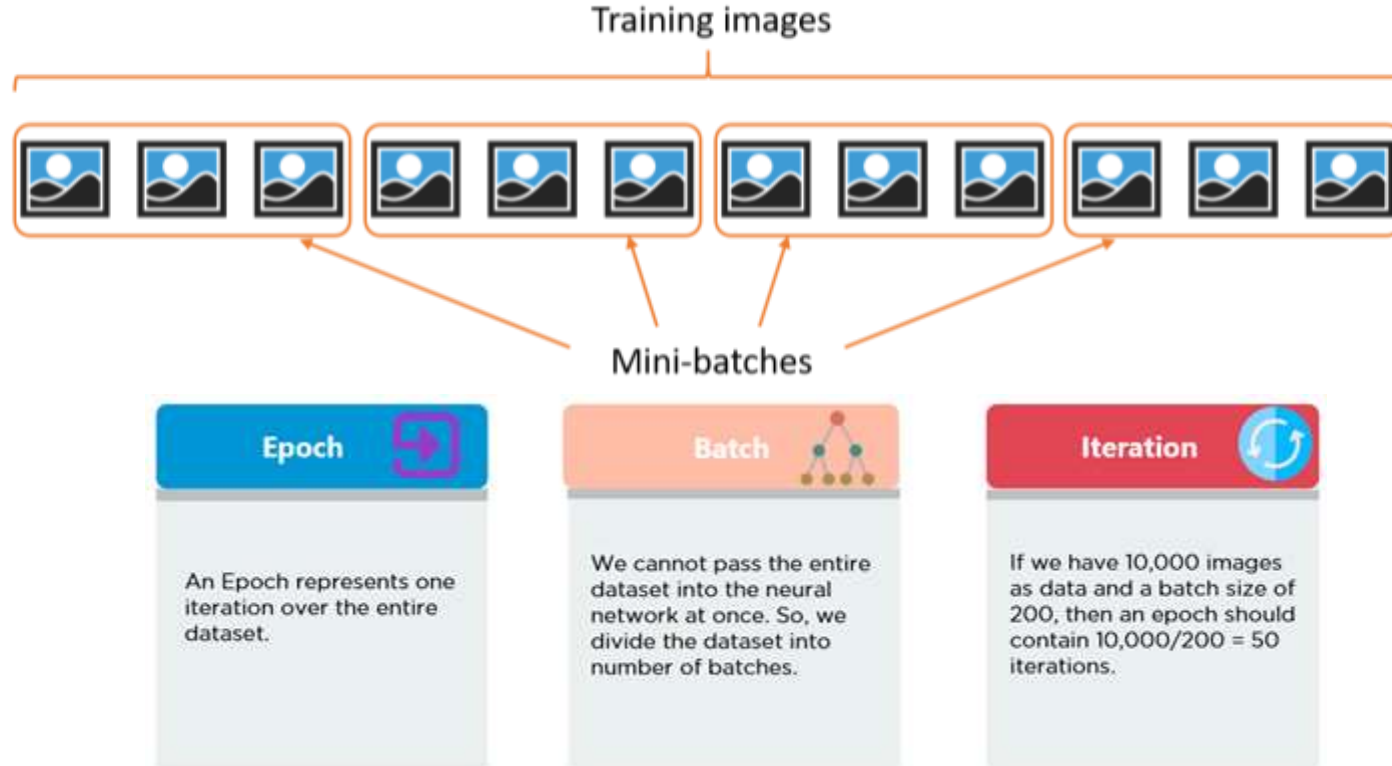


ReLU



# MINIBATCH VS SINGLE

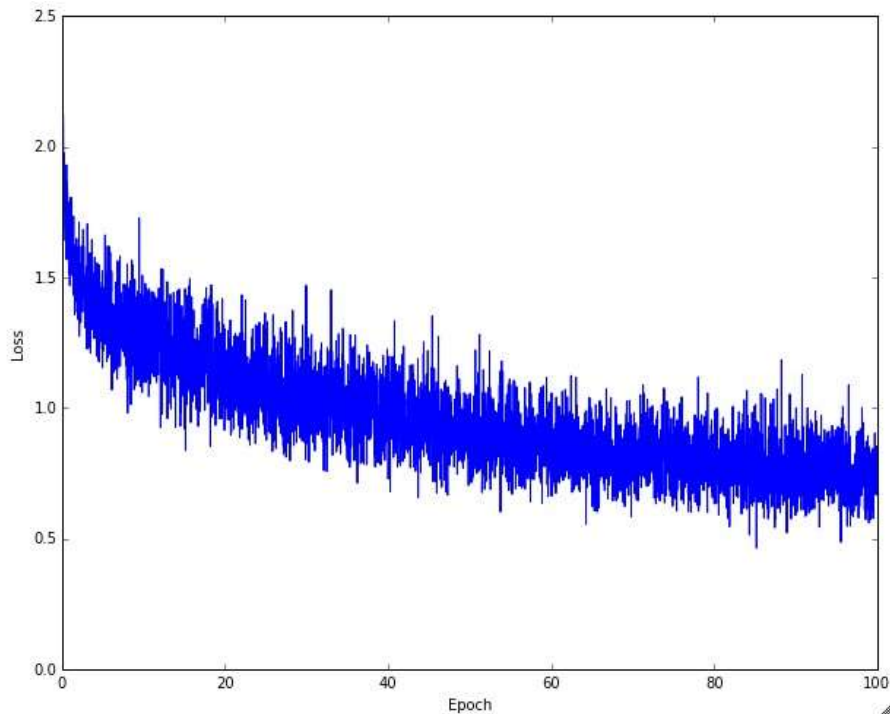
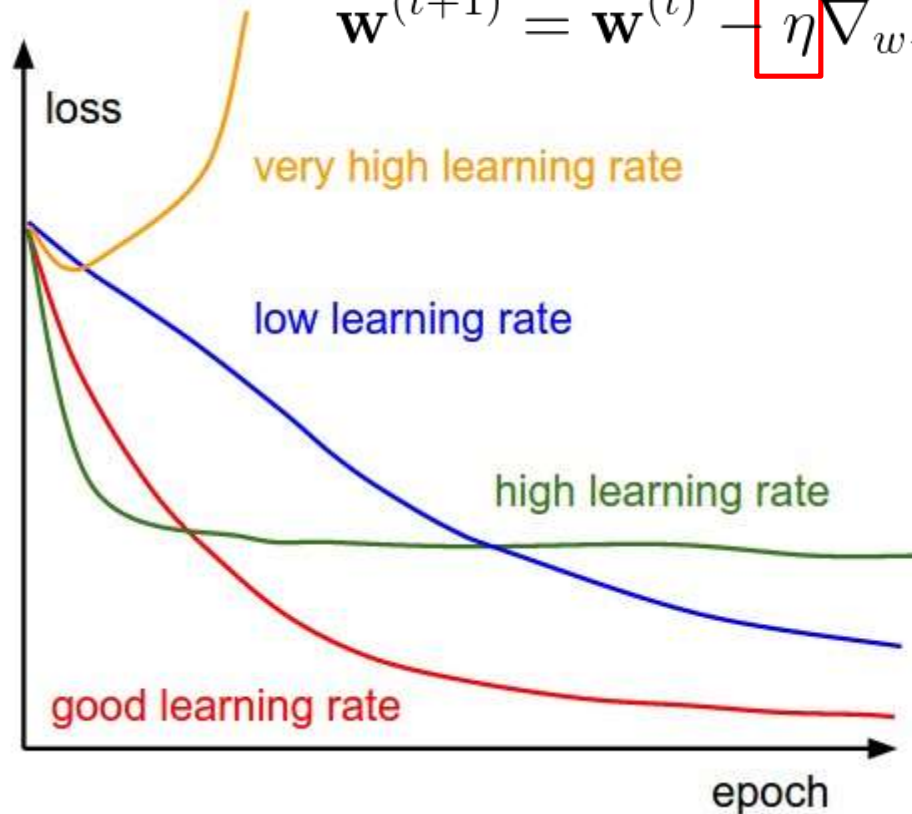
- Average error, gradients





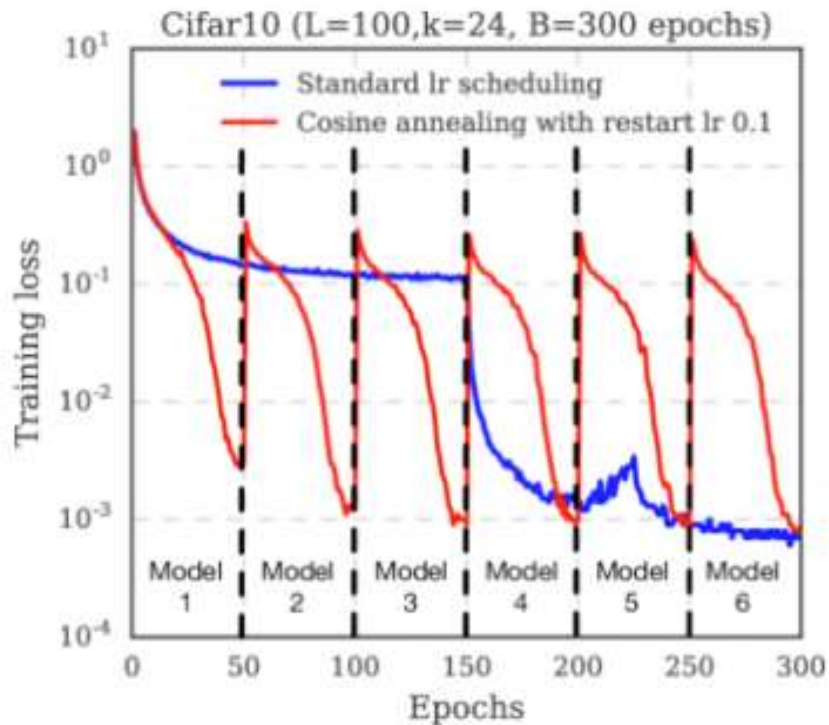
# TRAINING - SETTING LEARNING RATE

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$



# TRAINING - SETTING LEARNING RATE

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta^{(t)} \nabla_{\mathbf{w}} \mathbf{J}(\mathbf{w})$$



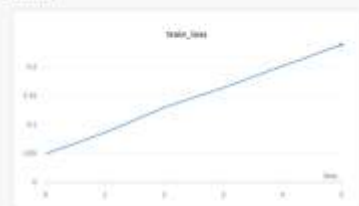
# lossfunctions

They are a window to your model's heart

Compute loss functions to @kharzhy. It doesn't matter if your loss functions are flat, converge, diverge, stop or oscillate (or any combination of the above). All loss functions are computed beautiful in their own way.

[POSTS](#) [ARCHIVE](#)

Chart: 4



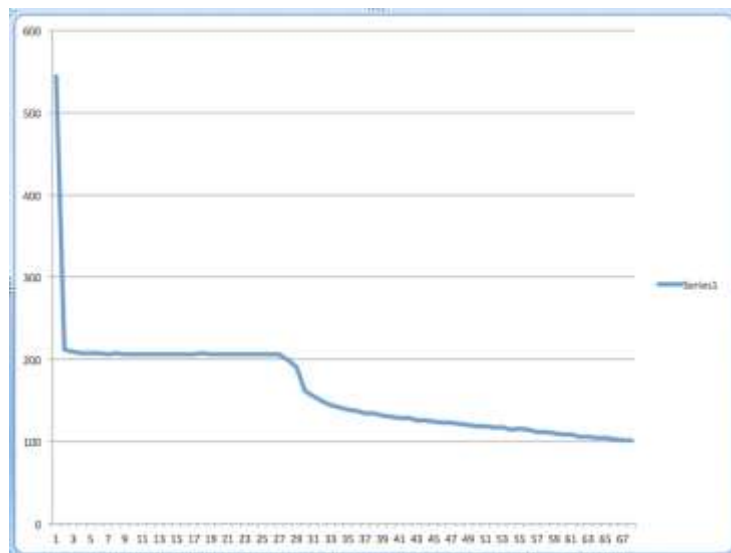
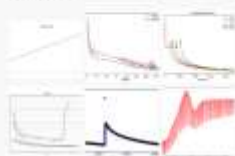
Trained loss minimizers. Wrote loss maximizers.

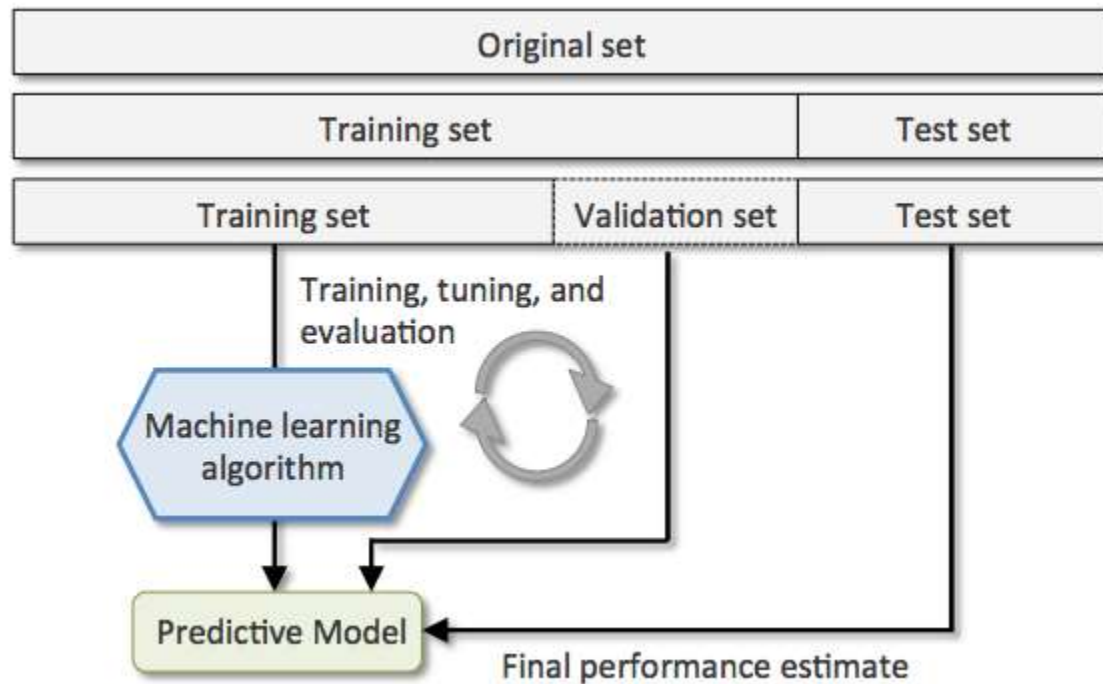
By @charisharveem

5 replies

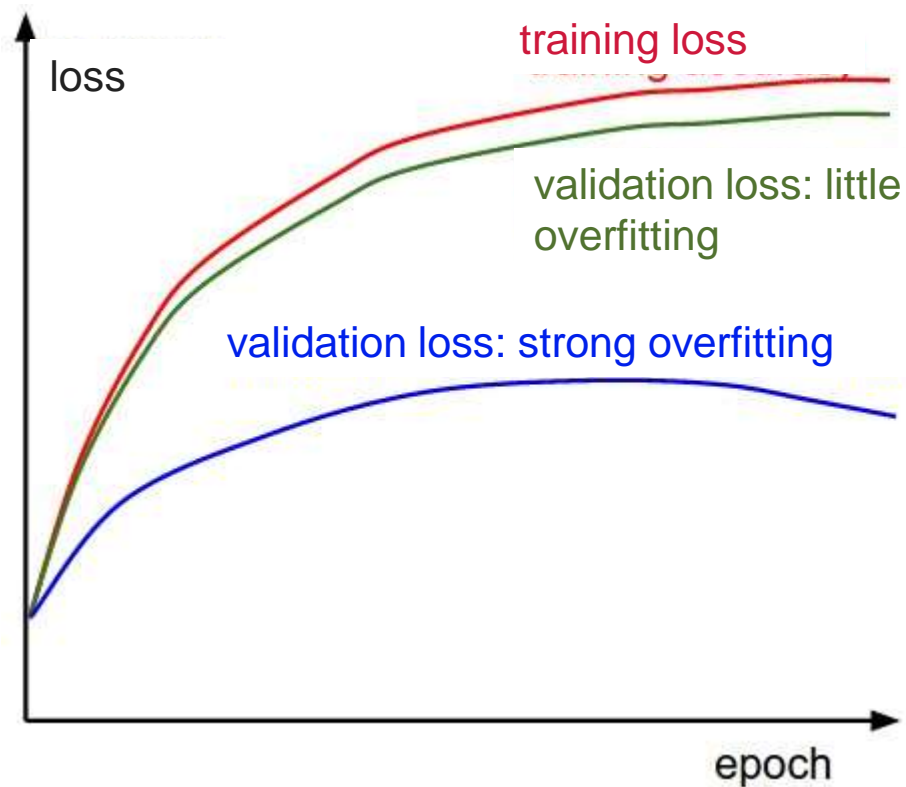


TOP PHOTOS



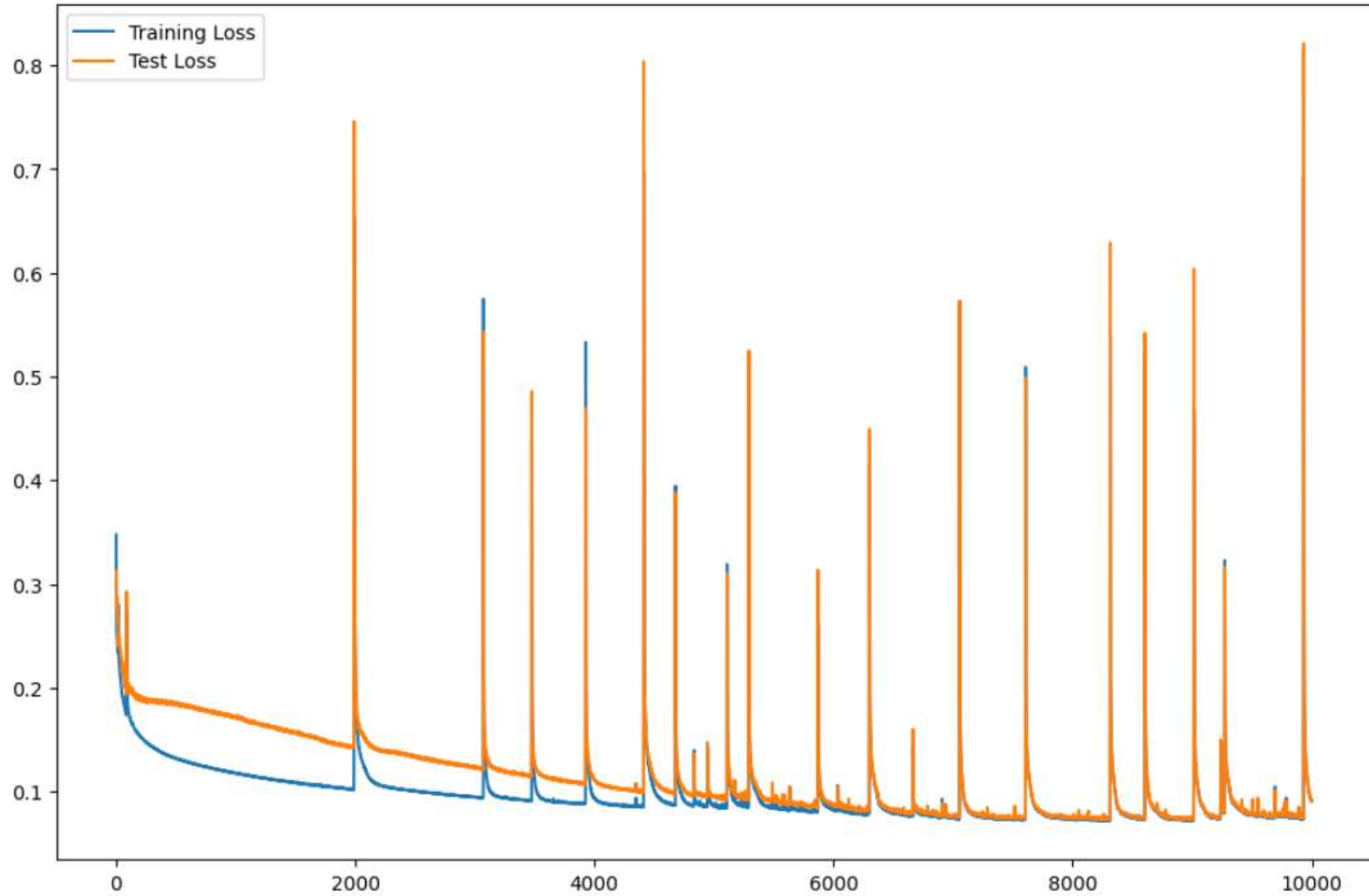


# WHEN TO STOP TRAINING



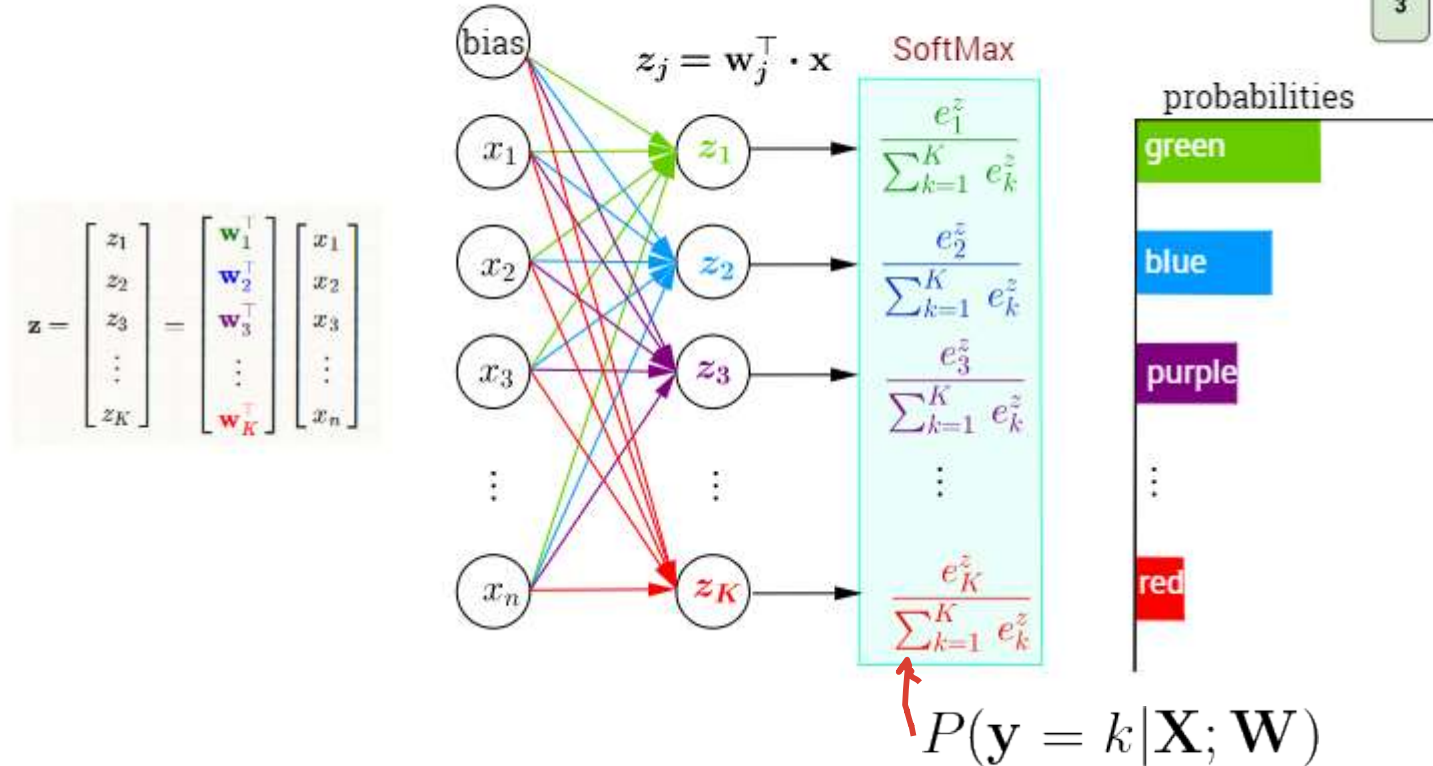
# WHEN TO STOP TRAINING





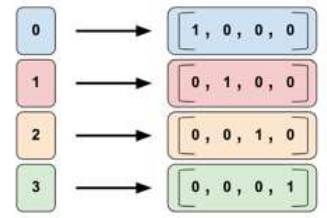
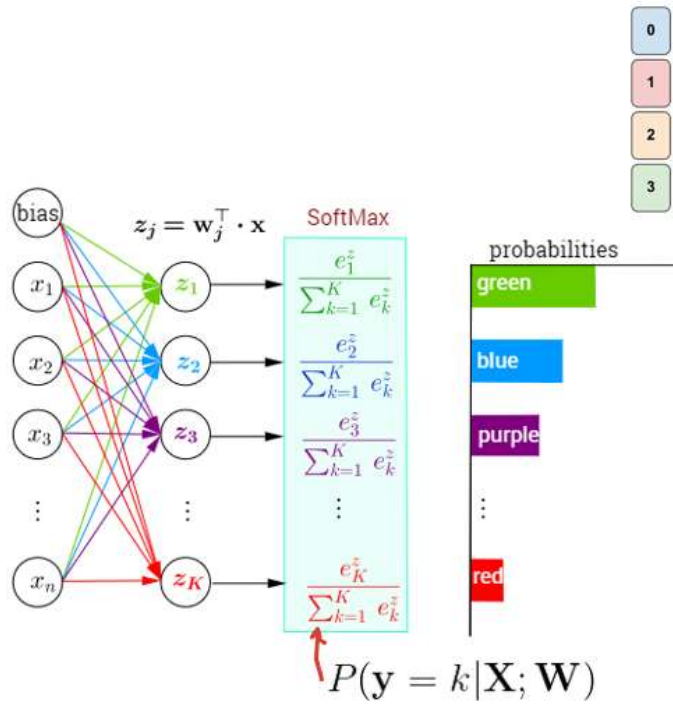
# CLASSIFICATION LOSS

## Multi-Class Classification with NN and SoftMax Function



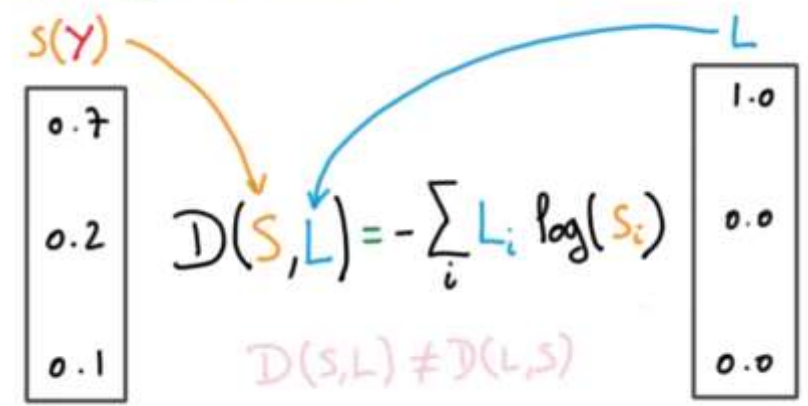


# CLASSIFICATION LOSS



$$\begin{aligned}
 D_{\text{KL}}(p|q) &= \sum_i p_i \log \frac{p_i}{q_i} \\
 &= \sum_i (-p_i \log q_i + p_i \log p_i) \\
 &= -\sum_i p_i \log q_i + \sum_i p_i \log p_i \\
 &= -\sum_i p_i \log q_i - \sum_i p_i \log \frac{1}{p_i} \\
 &= -\sum_i p_i \log q_i - H(p) \\
 &= \sum_i p_i \log \frac{1}{q_i} - H(p)
 \end{aligned}$$

## CROSS-ENTROPY



# CLASSIFICATION LOSS



computed				targets				correct?
-----								
0.3	0.3	0.4		0	0	1		yes
0.3	0.4	0.3		0	1	0		yes
0.1	0.2	0.7		1	0	0		no

Average Classification Error ?

# CLASSIFICATION LOSS

A

computed				targets				correct?
-----								
0.3	0.3	0.4		0	0	1		yes
0.3	0.4	0.3		0	1	0		yes
0.1	0.2	0.7		1	0	0		no

Average Classification Error ?

B

computed				targets				correct?
-----								
0.1	0.2	0.7		0	0	1		yes
0.1	0.7	0.2		0	1	0		yes
0.3	0.4	0.3		1	0	0		no

Average Classification Error ?

# CLASSIFICATION LOSS

A

computed			targets			correct?
-----			-----			
0.3	0.3	0.4	0	0	1	yes
0.3	0.4	0.3	0	1	0	yes
0.1	0.2	0.7	1	0	0	no

Which classifier is better ?

B

computed			targets			correct?
-----			-----			
0.1	0.2	0.7	0	0	1	yes
0.1	0.7	0.2	0	1	0	yes
0.3	0.4	0.3	1	0	0	no

# CLASSIFICATION LOSS

A

computed	targets	correct?
0.3 0.3 0.4	0 0 1	yes
0.3 0.4 0.3	0 1 0	yes
0.1 0.2 0.7	1 0 0	no

Which classifier is better ?

B

computed	targets	correct?
0.1 0.2 0.7	0 0 1	yes
0.1 0.7 0.2	0 1 0	yes
0.3 0.4 0.3	1 0 0	no



Classification accuracy is a crude way to measure how well NN has been trained!

# CLASSIFICATION LOSS

## A

computed	targets	correct?
0.3 0.3 0.4	0 0 1	yes
0.3 0.4 0.3	0 1 0	yes
0.1 0.2 0.7	1 0 0	no

CROSS-ENTROPY

$$D(S, L) = - \sum_i L_i \log(S_i)$$

$D(S, L) \neq D(L, S)$

Cross-entropy error ?

## B

computed	targets	correct?
0.1 0.2 0.7	0 0 1	yes
0.1 0.7 0.2	0 1 0	yes
0.3 0.4 0.3	1 0 0	no



Classification accuracy is a crude way to measure how well NN has been trained!

# CLASSIFICATION LOSS

A

MSE ?

computed	targets	correct?
0.3 0.3 0.4	0 0 1	yes
0.3 0.4 0.3	0 1 0	yes
0.1 0.2 0.7	1 0 0	no

B

computed	targets	correct?
0.1 0.2 0.7	0 0 1	yes
0.1 0.7 0.2	0 1 0	yes
0.3 0.4 0.3	1 0 0	no



$\ln()$  function in cross-entropy takes into account the closeness of a prediction and is a more granular way to compute error.

# RESOURCES

- Videos

- Example from lecture: <https://www.youtube.com/watch?v=bxe2T-V8XR&list=PLiaHhY2iBX9hdHaRr6b7XevZtgZR1PoU>
- 3Blue1Brown: [https://www.youtube.com/watch?v=aircAruvnKk&list=PLZHQ0b0WTQDNU6R1\\_67000Dx\\_ZCJB-3pi](https://www.youtube.com/watch?v=aircAruvnKk&list=PLZHQ0b0WTQDNU6R1_67000Dx_ZCJB-3pi)
- StatQuest: <https://www.youtube.com/watch?v=zxagGtF9MeU&list=PLblh5JK0oLUIxGDQs4LFFD--41Vzf-ME1>
- NN zero to hero: <https://www.youtube.com/watch?v=VMj-3S1tku0&list=PLAqhIrjkxbuWI23v9cThsA9GvCAUhRvKZ>
- BP in terms of computational graph (cs221n): <https://www.youtube.com/watch?v=i940vYb6noo>
- <https://theaisummer.com/weights-and-biases-tutorial/>