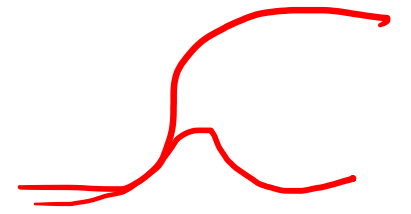


# Vanishing gradient problem

Yesterday → Caused bcz of Sigmoid → 0 to 1



$$\boxed{\eta = 1}$$

$$\frac{\partial \phi(z)}{\partial z} = \underline{0 \text{ to } 0.25}$$

$$w_{ij}^{\text{new}} = w_{ij}^{\text{old}} - \eta \frac{\partial L}{\partial w_{ij}^{\text{old}}} \Rightarrow \frac{\partial L}{\partial w_{ij}^{\text{old}}} = \frac{\partial}{\partial} \cdot \frac{\partial}{\partial} \cdot \frac{\partial}{\partial}$$

$\frac{\partial L}{\partial w_{ij}^{\text{old}}} \approx 0.0004$

$0.2 \quad 0.1 \quad 0.02 \quad = 0.0004$

v.v. small.

$$\boxed{w_{ij}^{\text{new}} \approx w_{ij}^{\text{old}}}$$

# Exploding Gradient Descent

(RELU is not invented yet)  
we only as researchers know  
SIGMOID.


→ Calprit: WEIGHTS.

How? 
 small whole | small whole | big whole  
 $2 \times 4 = 8$  |  $3 \times 12 = 36$ 
 } Exploding Gradient

Vanishing Gradient

small point | small point | more smaller no.  
 $0.2 \times 0.4 = 0.08$   
 $0.3 \times 0.12 = 0.036$

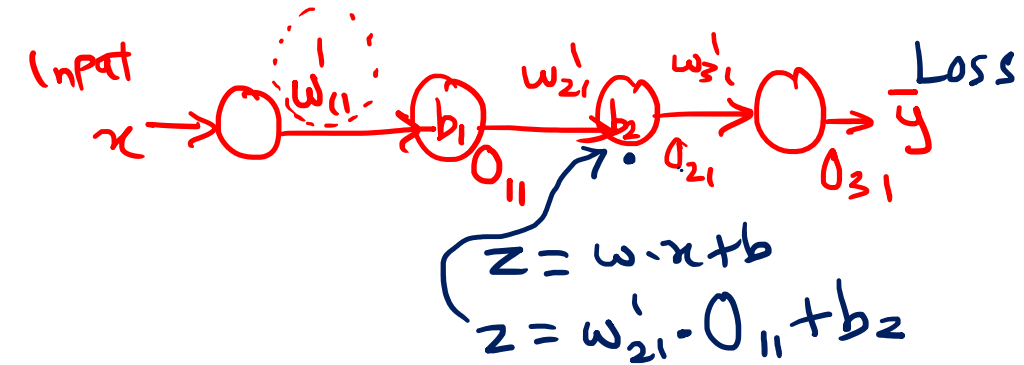
\* In EX. Gr. Pb, the larger assigned weight when mult. by derivative of Sigmoid returns a v. big No.

assume:  $450$   
 $w_{new} = w_{old} - \eta \cdot \frac{\partial L}{\partial w_{old}}$  → chain rule  


# Exploding Gradient Descent

$$0 \leq \phi(z) \leq 0.25$$

$$\eta = 1$$



$$w'_{11, \text{new}} = w'_{11, \text{old}} - \eta \frac{\partial L}{\partial w'_{11, \text{old}}}$$

$w'_{11, \text{new}} \approx w'_{11, \text{old}} - \eta \cdot \text{diff}$

large

v.v. hugeno

$$\frac{\partial L}{\partial w'_{11, \text{old}}} = \frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w'_{11, \text{old}}}$$

150  $\times$  100  $\times$  75

$$\frac{\partial o_{21}}{\partial o_{11}} = \frac{\partial \phi(z)}{\partial z} \cdot \frac{\partial z}{\partial o_{11}}$$

0 to 0.25  $\times$   $\frac{\partial}{\partial o_{11}} (w'_{21} \cdot o_{11} + b_2)$

$$= (0 \text{ to } 0.25) \times w'_{21}$$

$$= (0 \text{ to } 0.25) \times \text{large no}$$

$$= 0.25 \times 400 = 100$$

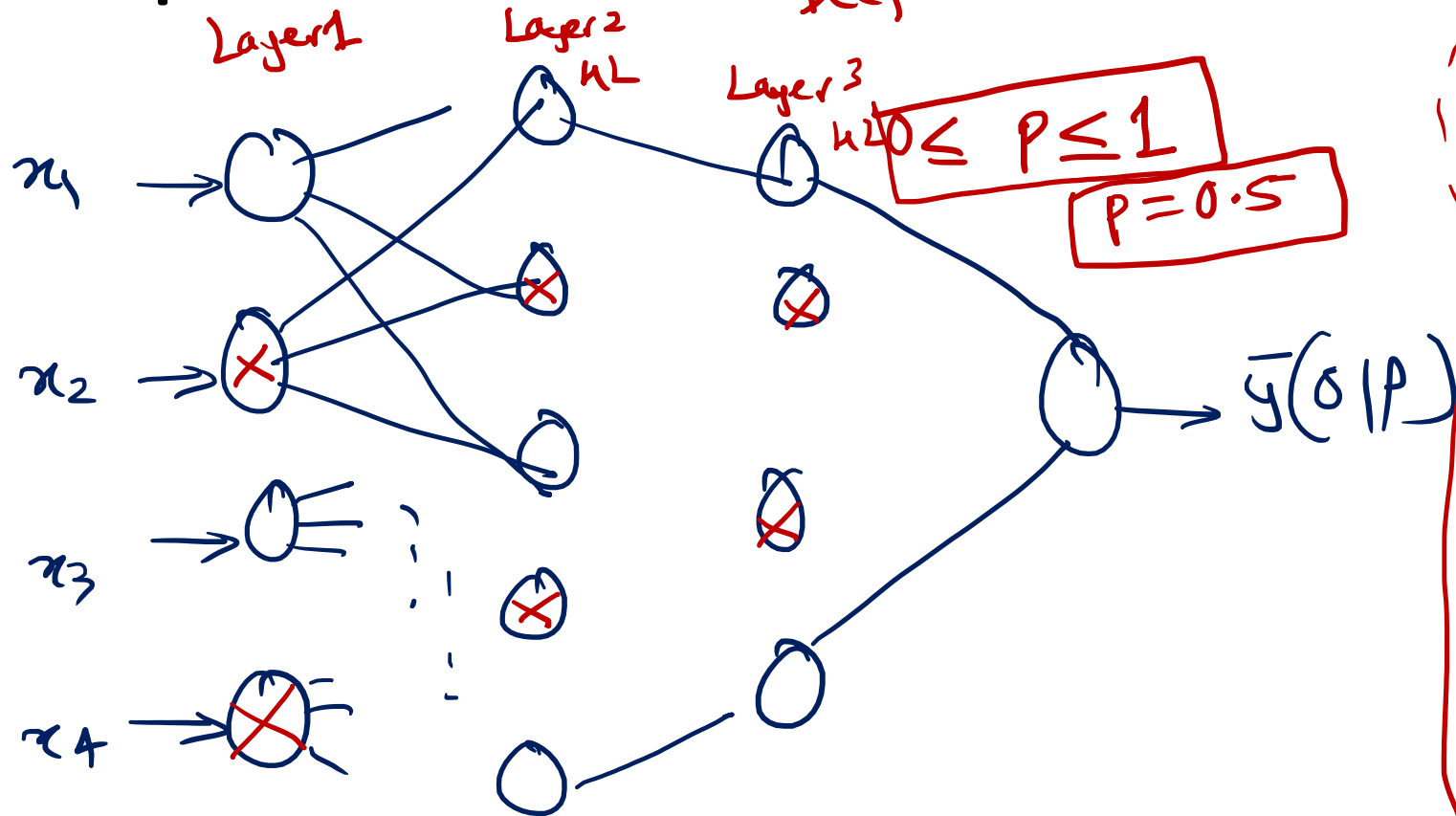


Weights  
 $\Downarrow$   
 Exploding Gradient Problem

# Exploding Gradient Descent

# Dropout in MNN

M.L.  $\Rightarrow$  Random Forest  
is used for feature selection.

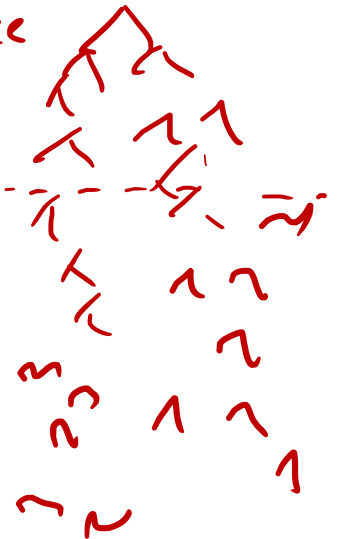


OVERFITTING

Decision Tree  
 $\hookrightarrow$  Overfitting

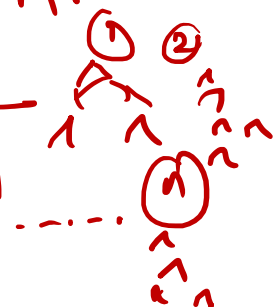
H. Variance  
Overcome!

$\text{max\_depth} = 3$   
Hyperparameter



RF  $\rightarrow$  Many D.T.

x	1	2	3	4	5	6	7	8	y



Solve! 2 ways:

- Regularization:  $L_1, L_2$
- Dropout: Outcome of a PhD student (Syris)

Nitish Shrivastava & Jeffery Hinton (2014)

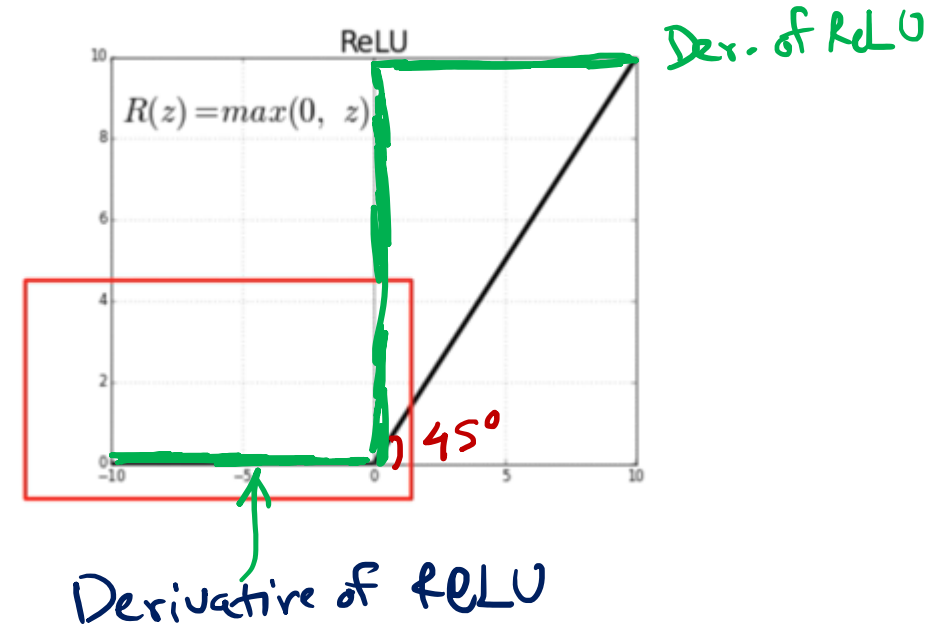
# Dropout in MNN

How to select p-value? → Hyperparameter Optimization  
(K-Fold C.V.)

# Dropout in MNN

# ReLU

$$= \max(0, z)$$



Simply Saying: 1 if  $z > 0$   
0 if  $z < 0$

Calculus: Der. of 0 is not differentiable

$$w_{\text{new}} = w_{\text{old}} - \eta \left( \frac{\partial L}{\partial w} \right)$$

$$\frac{\partial}{\partial 1} \cdot \frac{\partial}{\partial 1} \cdot \frac{\partial}{\partial 1} = 1$$

Vanishing Gradient  
Prob. is thus  
solved by ReLU



# ReLU

$$w_{new} = w_{old} - \eta \cdot \frac{\partial L}{\partial w_{old}} \Rightarrow \frac{\partial}{\partial} \cdot \frac{\partial}{\partial} \cdot \frac{\partial}{\partial}$$

0 : 1 : 1

DEAD NEURON/deactivated

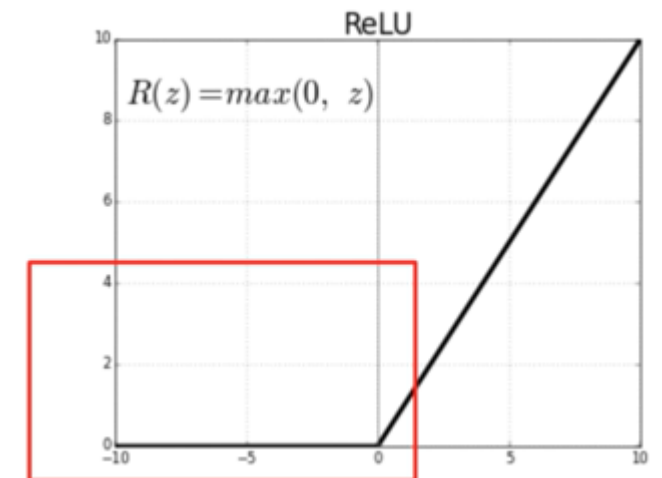
$w_{new} = w_{old}$

- Now if we apply this in the chain rule, works perfectly.
- But if any derivative gives value as 0, chain rule will straightaway make the calculation value as 0 thereby making it a dead neuron or a dead activation function.
- Therefore we go for Leaky ReLU
- With ReLU, always use weight initializer as he normal or he uniform.
- With Sigmoid, use weight initializer as glorot\_uniform

Der. of ReLU

$$1 \text{ if } z > 0$$

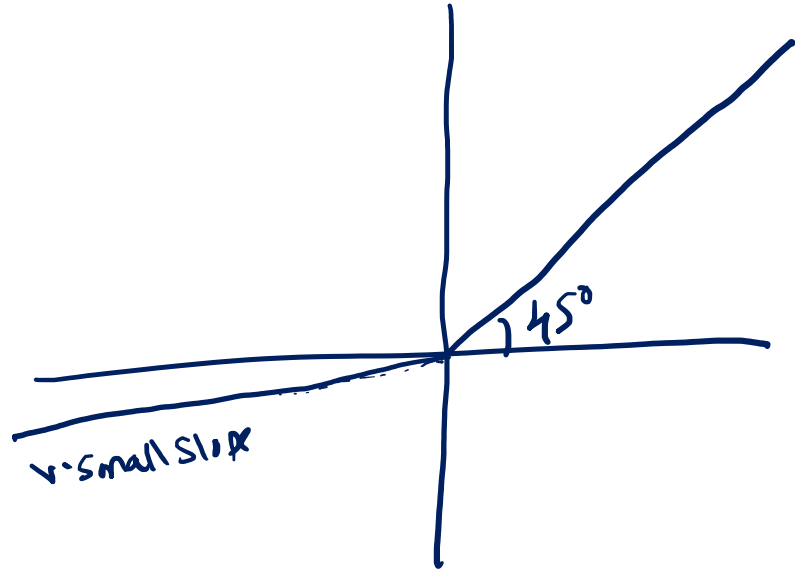
$$0 \text{ if } z < 0$$



# ReLU

# ReLU

# Leaky ReLU



$$\frac{\partial}{\partial z} \cdot \frac{\partial}{\partial y} \cdot \frac{\partial}{\partial x}$$
$$0.01 \cdot 1 \cdot 1$$

Der. of ~~ReLU~~ Leaky

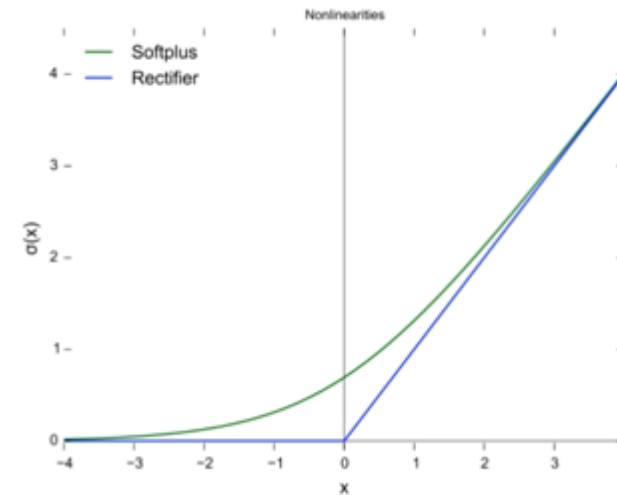
$$\text{ReLU} : \begin{cases} 1 & \text{if } z \geq 0 \\ 0.01 \cdot (z) & \text{if } z < 0 \end{cases}$$

add some  
small value for -ve  
derivate

# Softplus Activation Function

- Another option which is very similar is the soft plus activation which because you're taking the log of the exponent looks very linear when the input is reasonably large.

$$f(x) = \log(1 + e^x)$$



- But for both of these previous activation functions, there is the vanishing gradient on the left side but we've established that it's not so much of a problem since we know that the real you already works and it has gradients that are equal to zero also.

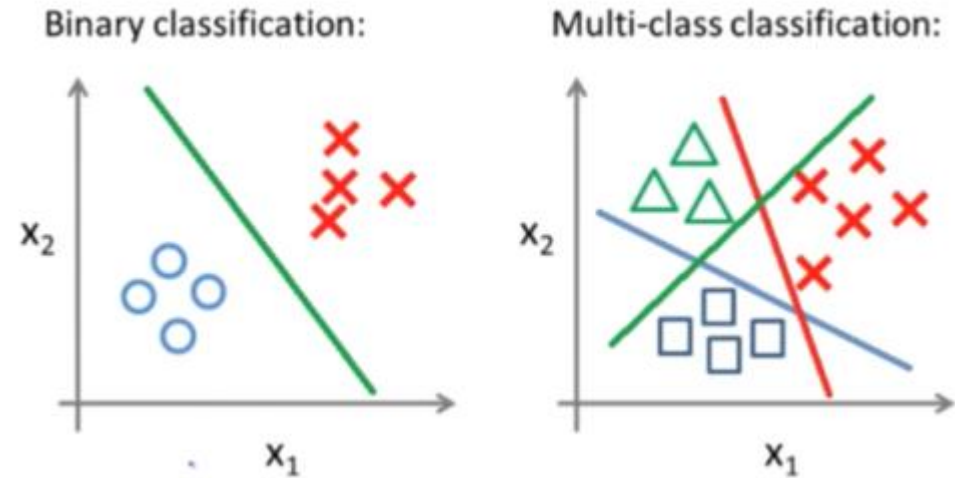
# Softplus Activation Function

- Although we initially stated that we would like the inputs at each layer to be centered around zero we can see that the ReLU and Softplus do not accomplish this.
- For soft plus and the ReLU the minimum value is zero while the maximum value is infinity.
- This definitely means they won't be centered around zero. So is the ~~ReLU~~ <sup>ReLU</sup> ~~you~~ not a good choice in the end?
- Now despite all this work to find alternatives to the value activation these days most people still use the ~~ReLU~~ <sup>ReLU</sup> as a reasonable default choice.
- It works well and sometimes you'll find that using other alternatives such as the leaky ReLU or the ELU offer no benefit. But Sometimes they do, which is why you always have to experiment for yourself.

# Softplus Activation Function

- My motto which a lot of my students are tired of hearing by this point is that machine learning is experimentation and not philosophy.
- Never use your mind to try and predict the outcome of a computer program.
- If you have a computer that is always the suboptimal course of action why not simply run the computer program with a computer.
- Your mind is not suitable for running computer programs but computers are therefore follow the rule.
- Don't use philosophy use experimentation.

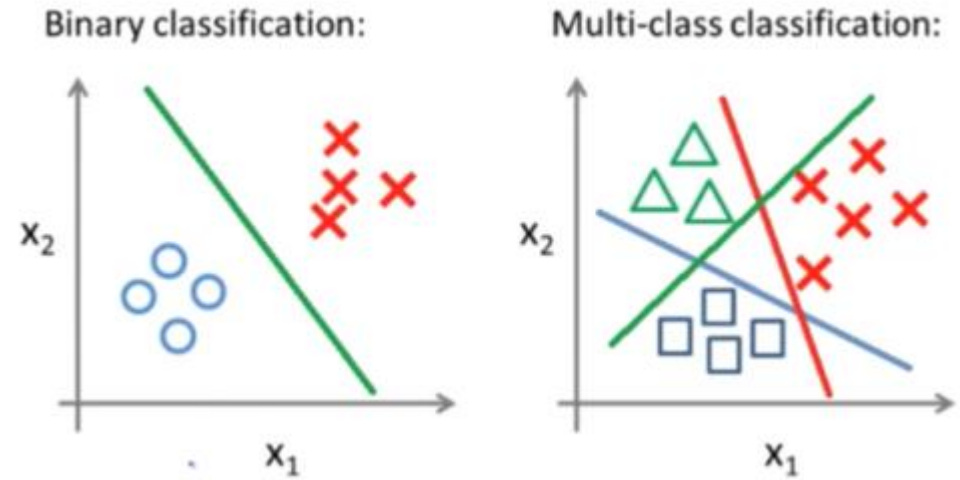
# Multiclass Classification



- We saw that for a binary classification, we use sigmoid at the output.
- We replaced sigmoid with ReLUs in the hidden layer
- However, for output, sigmoid is still the right choice for binary classification
- Applications of Binary Classification:
  - Disease vs No Disease
  - Fraud vs No Fraud
  - Click vs No Click
  - Accept Friend Request



# Multiclass Classification



- But there are situations that binary classification cannot handle such as when we might have multiple categorical outcomes.

# Multiclass Classification

Example :

① Character Recognition/Handwriting Recognition ⇒

② Speech Recognition: Each word counts as a separate category.

③ Image Classification: ImageNet (million of images of thousands of categories like bagels, vending m/c, soccerball, etc)



# Multiclass Classification



# Using Softmax in Tensorflow

- Well just like most of the other functions we've applied so far it's very simple.
- We just pass in the string 'softmax'.
- In other words the only requirement is that you spell it correctly.
- Of course you can always implement the softmax yourself but in Tensorflow, there's no need.
- As a side note, the softmax is considered an activation function but unlike the ReLU, sigmoid and tanh, it is not meant for hidden layer activations.
- If you want to try using the softmax as a hidden layer activation, you are most welcome to but you'll generally find that it doesn't work that well.
- So the softmax is technically an activation function but we normally only use it when we're trying to get an output probability from a vector of activation values such as at the end of a neural network.

```
tf.keras.layers.Dense(K, activation='softmax')
```

Trainer: Dr. Darshan Ingle.

# Task Summary

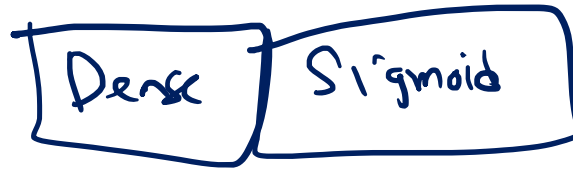
Task	A.F
Binary	Sigmoid
Multi-class	Softmax
Regression	None/Linear

# The Model Type doesn't matter

Linear Reg.



Binary Logistic Reg.



Multi-class Logistic Reg.



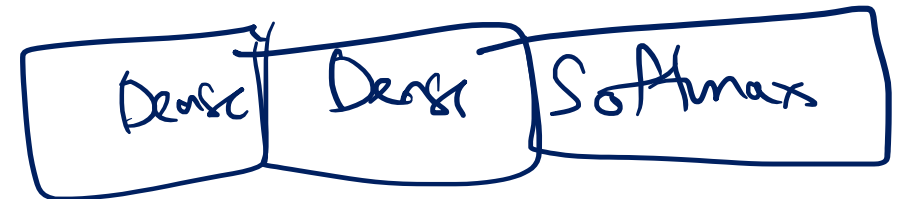
ANN - Reg.



ANN. Binary Class.



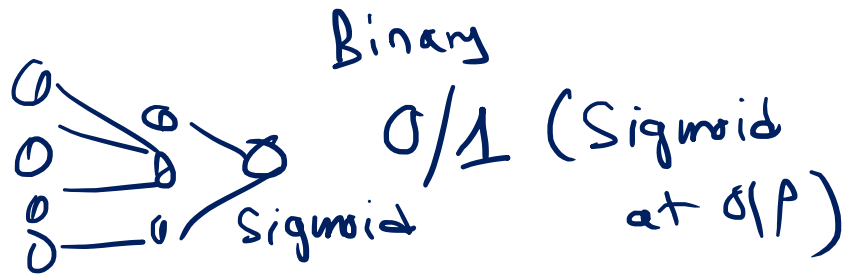
ANN Multi-class Classification



# Softmax is more general

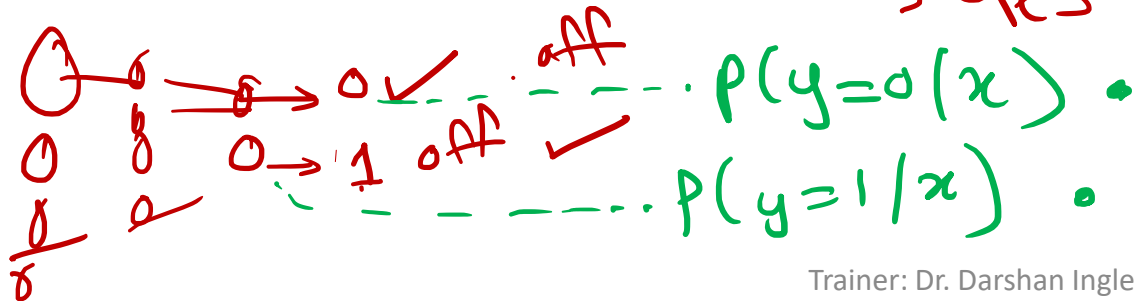
Softmax — Multiclass

Can it be also used for Binary Classification?



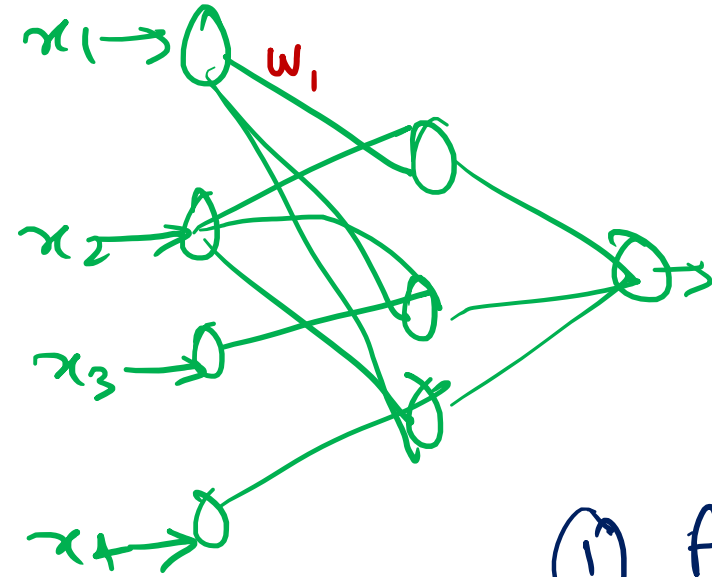
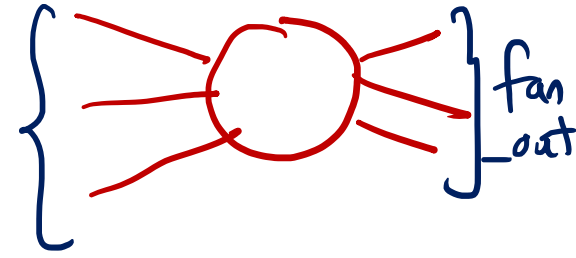
Can I use Softmax here? Softmax = Multiclass (greater than 1)

Yes (Class  $k=2$ )



# Weight Initialization

Research Paper Convention/  
Notation!



Important Points:

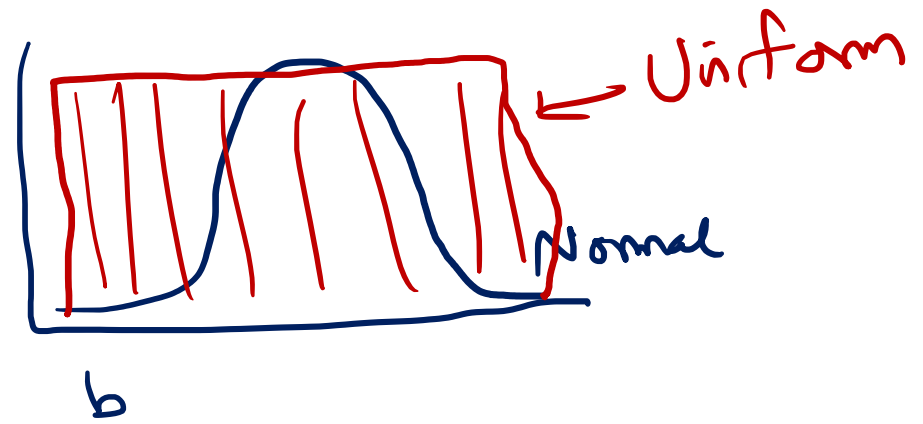
- ① Always initialize weights to a small value.
- ② Don't keep all weights same.
- ③ All weights should have good amount of variance.



# Weight Initialization

① Uniform Distribution:

$$w_{ij} \sim \text{Uniform} \left[ \frac{-1}{\sqrt{\text{fan-in}}}, \frac{1}{\sqrt{\text{fan-in}}} \right]$$



# Weight Initialization

② Xavier/Glorot : Use Sigmoid

Two Types:

① Xavier/Glorot Normal

$$w_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{(fan\_in + fan\_out)}}$$

② Xavier/Glorot Uniform

$$w_{ij} \sim U(0, \sigma)$$

$$\sigma = \left[ \frac{\sqrt{6}}{fan\_in + fan\_out}, \frac{\sqrt{6}}{\sqrt{fan\_in + fan\_out}} \right]$$

# Weight Initialization

③ He init '(ReLU)':

Two Types:

① He Uniform

$$w_{ij} \approx U \left[ -\sqrt{\frac{6}{\text{fan-in}}}, \sqrt{\frac{6}{\text{fan-out}}} \right]$$

② He Normal

$$w_{ij} \approx N(0, \sigma)$$
$$\sigma = \sqrt{\frac{2}{\text{fan-in}}}$$