

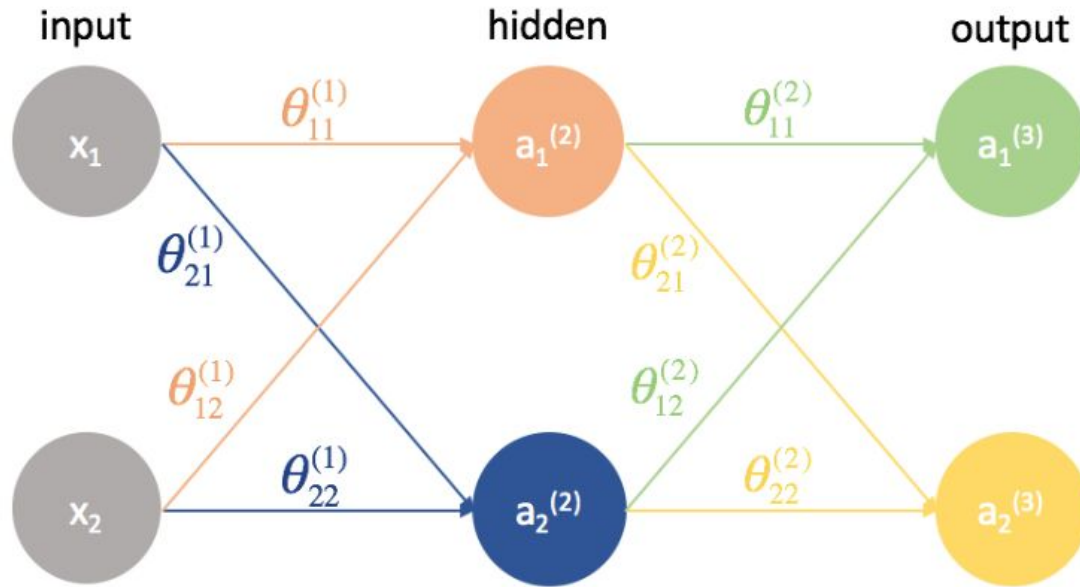
THIS IS CS4045!

**GCR:dxuxugo**

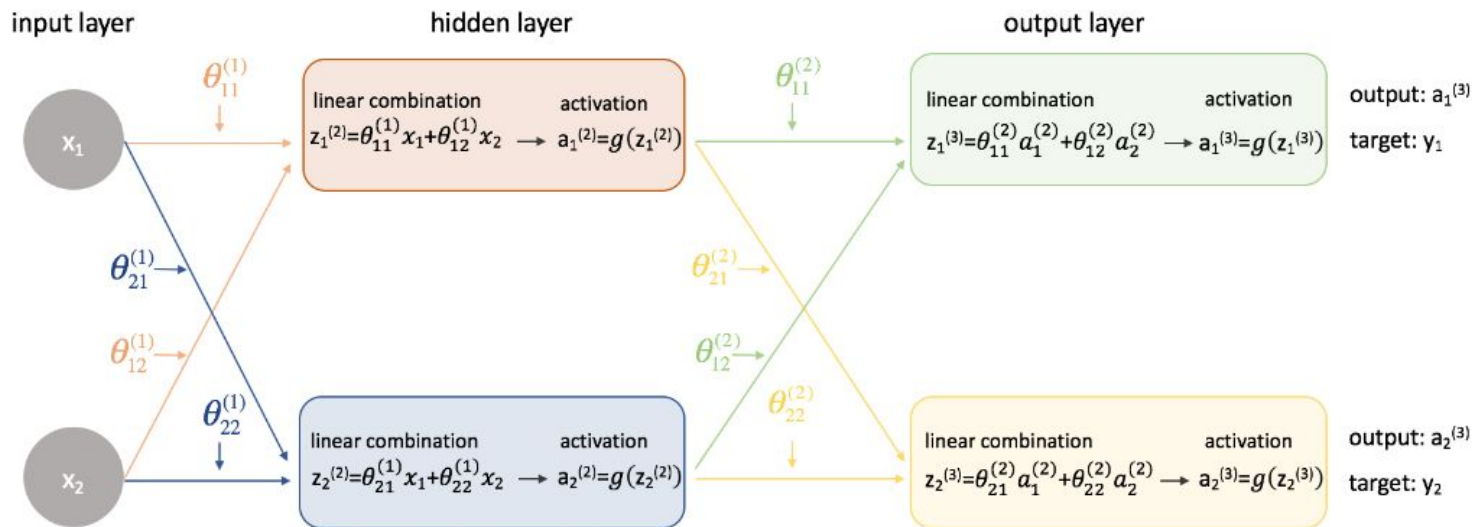
P.S. THESE SLIDES ARE USELESS IF YOU DO  
NOT ATTEND CLASSES

# NEURAL NETWORKS

# NEURAL NETWORK WITH MULTI OUTPUT



# NEURAL NETWORK WITH MULTI OUTPUT



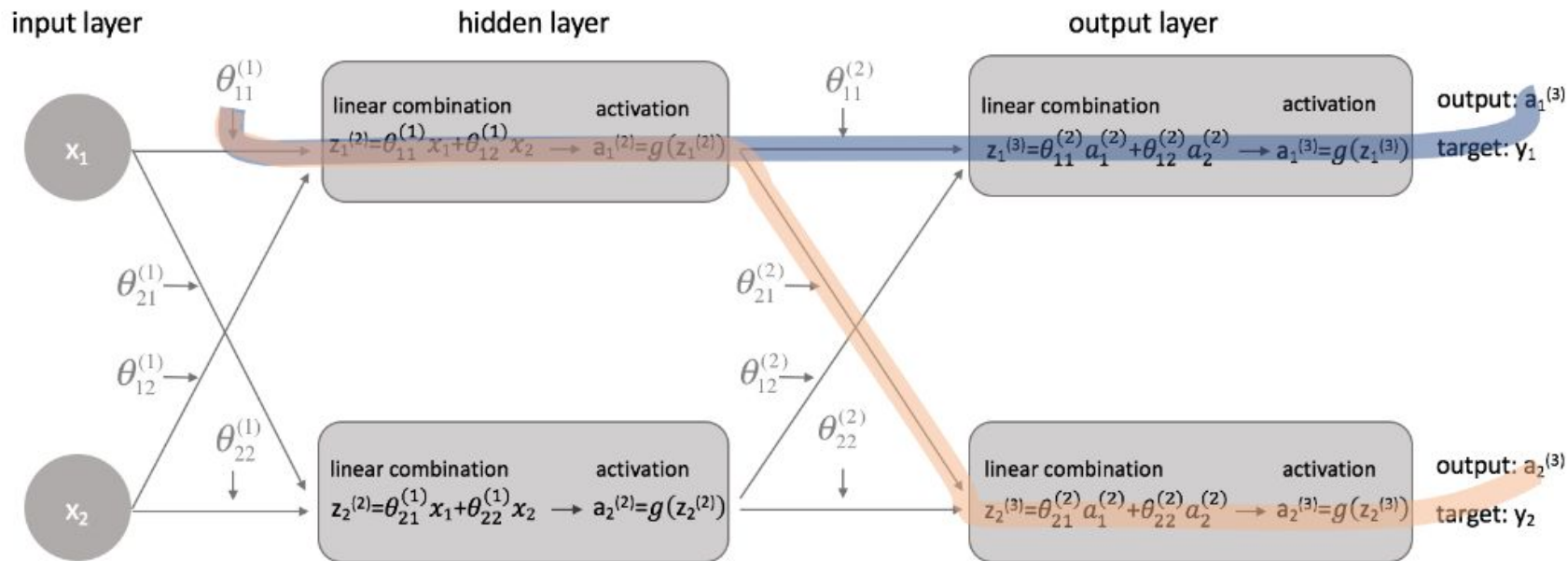
# NEURAL NETWORK WITH MULTI OUTPUT

Loss calculation

m=number of samples

$$J(\theta) = \frac{1}{2m} \sum \left( y_i - a_i^{(2)} \right)^2$$

# NEURAL NETWORK WITH MULTI OUTPUT



The derivative chain for the blue path is:

$$\left( \frac{\partial J(\theta)}{\partial \mathbf{a}_1^{(3)}} \right) \left( \frac{\partial \mathbf{a}_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial \mathbf{a}_1^{(2)}} \right) \left( \frac{\partial \mathbf{a}_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right)$$

The derivative chain for the orange path is:

$$\left( \frac{\partial J(\theta)}{\partial \mathbf{a}_2^{(3)}} \right) \left( \frac{\partial \mathbf{a}_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial \mathbf{a}_1^{(2)}} \right) \left( \frac{\partial \mathbf{a}_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right)$$

Combining these, we get the total expression for  $\frac{\partial J(\theta)}{\partial \theta_{11}^{(1)}}$ .

$$\frac{\partial J(\theta)}{\partial \theta_{11}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial \mathbf{a}_1^{(3)}} \right) \left( \frac{\partial \mathbf{a}_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial \mathbf{a}_1^{(2)}} \right) \left( \frac{\partial \mathbf{a}_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial \mathbf{a}_2^{(3)}} \right) \left( \frac{\partial \mathbf{a}_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial \mathbf{a}_1^{(2)}} \right) \left( \frac{\partial \mathbf{a}_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right)$$



## Layer 2 Parameters

$$\frac{\partial J(\theta)}{\partial \theta_{11}^{(2)}} = \left( \frac{\partial J(\theta)}{\partial \mathbf{a}_1^{(3)}} \right) \left( \frac{\partial \mathbf{a}_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial \theta_{11}^{(2)}} \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_{12}^{(2)}} = \left( \frac{\partial J(\theta)}{\partial \mathbf{a}_1^{(3)}} \right) \left( \frac{\partial \mathbf{a}_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial \theta_{12}^{(2)}} \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_{21}^{(2)}} = \left( \frac{\partial J(\theta)}{\partial \mathbf{a}_2^{(3)}} \right) \left( \frac{\partial \mathbf{a}_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial \theta_{21}^{(2)}} \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_{22}^{(2)}} = \left( \frac{\partial J(\theta)}{\partial \mathbf{a}_2^{(3)}} \right) \left( \frac{\partial \mathbf{a}_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial \theta_{22}^{(2)}} \right)$$



# DISCUSSION ON ASSIGNMENT NETWORK

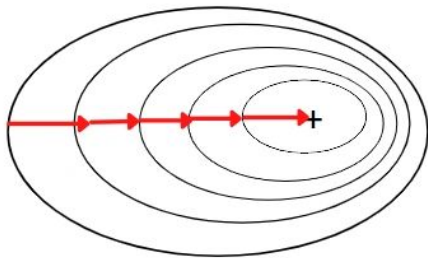
1. [20 Points] Draw 3-layer Neural Networks (Input layer is not counted here) with 2 nodes in each hidden layer and 1 output node. Input size is 3. Activation function is relu for 1st hidden layer and sigmoid for 2nd hidden and last layer. Write mathematical equations of all hidden and output layer nodes for a forward pass. Write down mathematical equations of the partial derivative of Loss w.r.t. all weights and bias for a backward pass. Simplify all these equations. Use (i) Squared Loss (ii) Cross Entropy Loss

# DISCUSSION ON ASSIGNMENT NETWORK

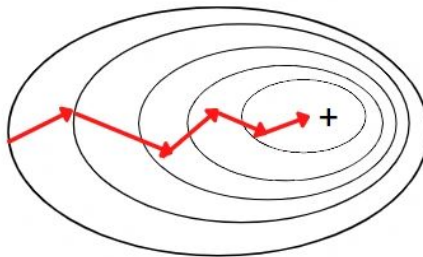
2. [10 Points] Draw 2-layer Neural Networks with 1 node in the hidden layer and 2 output nodes. Input size is 2. The activation function in the hidden layer is tanh while in the output layer, sigmoid is used. Write mathematical equations of all hidden and output layer nodes for a forward pass. Write down mathematical equations of the partial derivative of Loss w.r.t. all weights and bias for a backward pass. Simplify all these equations. Use Cross Entropy Loss.

# GRADIENT DESCENT VARIATIONS - HOW TO UPDATE PARAMETERS?

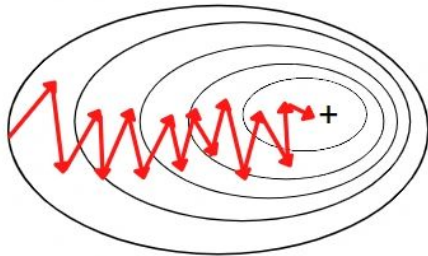
**Batch Gradient Descent**



**Mini-Batch Gradient Descent**



**Stochastic Gradient Descent**



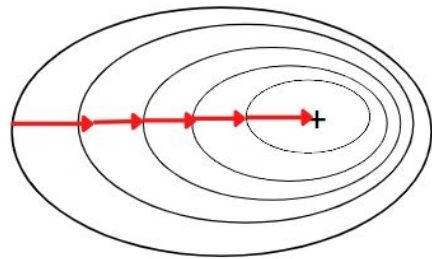
# GRADIENT DESCENT VARIATIONS - HOW TO UPDATE PARAMETERS?

Batch gradient descent, also known as vanilla gradient descent

Calculates the error for each example within the training dataset.

Still, the model is not changed until every training sample has been assessed. The entire procedure is referred to as a cycle and a training epoch.

Batch Gradient Descent

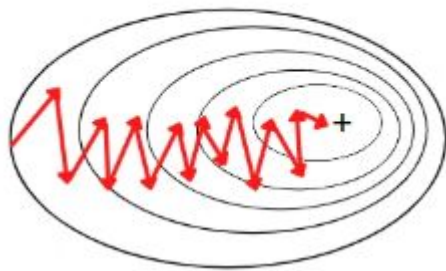


# GRADIENT DESCENT VARIATIONS - HOW TO UPDATE PARAMETERS?

SGD changes the parameters for each training sample one at a time for each training example in the dataset.

It makes SGD faster than batch gradient descent.

Stochastic Gradient Descent



One benefit is that the regular updates give us a fairly accurate idea of the rate of improvement.

The frequency of such updates can also produce noisy gradients, which could cause the error rate to fluctuate rather than gradually go down.

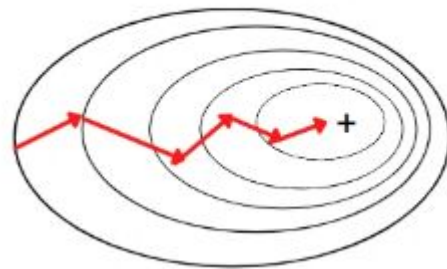
# GRADIENT DESCENT VARIATIONS - HOW TO UPDATE PARAMETERS?

Mini-batch gradient descent combines ideas of batch gradient descent with SGD, it is the preferred technique.

It divides the training dataset into manageable groups and updates each separately.

This strikes a balance between batch gradient descent's effectiveness and stochastic gradient descent's durability.

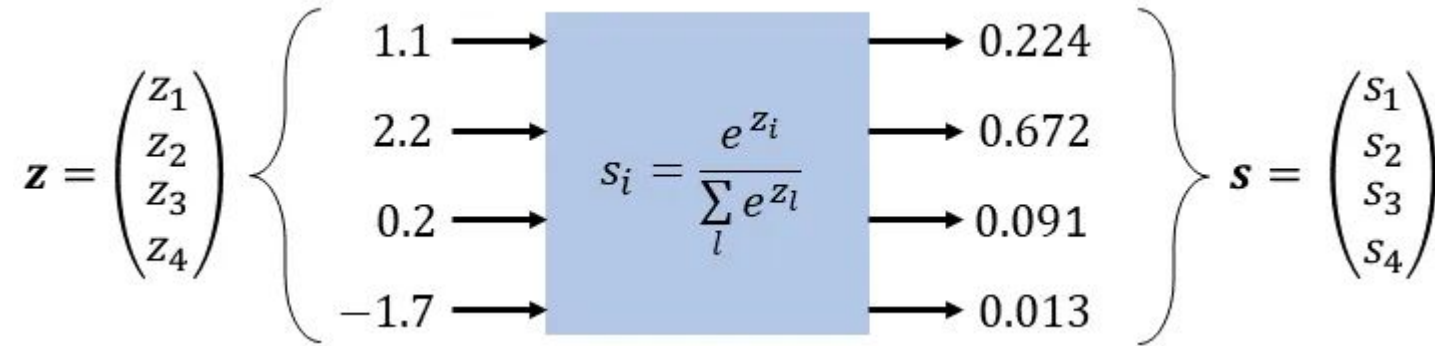
Mini-Batch Gradient Descent





# SOFTMAX ACTIVATION FUNCTION FOR MULTICLASS

$$\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$



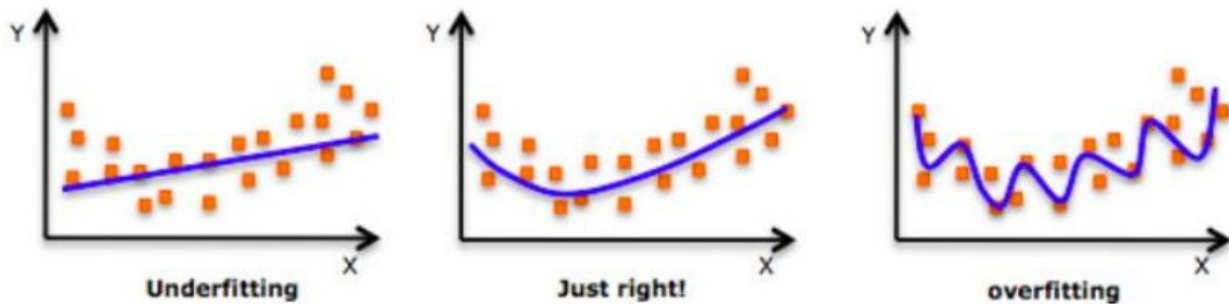
# SOFTMAX ACTIVATION FUNCTION FOR MULTICLASS - DERIVATIVE

$$J_{softmax} = \begin{pmatrix} \frac{\partial s_1}{\partial z_1} & \frac{\partial s_1}{\partial z_2} & \dots & \frac{\partial s_1}{\partial z_n} \\ \frac{\partial s_2}{\partial z_1} & \frac{\partial s_2}{\partial z_2} & \dots & \frac{\partial s_2}{\partial z_n} \\ \frac{\partial s_3}{\partial z_1} & \frac{\partial s_3}{\partial z_2} & \dots & \frac{\partial s_3}{\partial z_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial s_n}{\partial z_1} & \frac{\partial s_n}{\partial z_2} & \dots & \frac{\partial s_n}{\partial z_n} \end{pmatrix}$$

$$J_{softmax} = \begin{pmatrix} s_1 \cdot (1 - s_1) & -s_1 \cdot s_2 & -s_1 \cdot s_3 & -s_1 \cdot s_4 \\ -s_2 \cdot s_1 & s_2 \cdot (1 - s_2) & -s_2 \cdot s_3 & -s_2 \cdot s_4 \\ -s_3 \cdot s_1 & -s_3 \cdot s_2 & s_3 \cdot (1 - s_3) & -s_3 \cdot s_4 \\ -s_4 \cdot s_1 & -s_4 \cdot s_2 & -s_4 \cdot s_3 & s_4 \cdot (1 - s_4) \end{pmatrix}$$

# REGULARIZATION

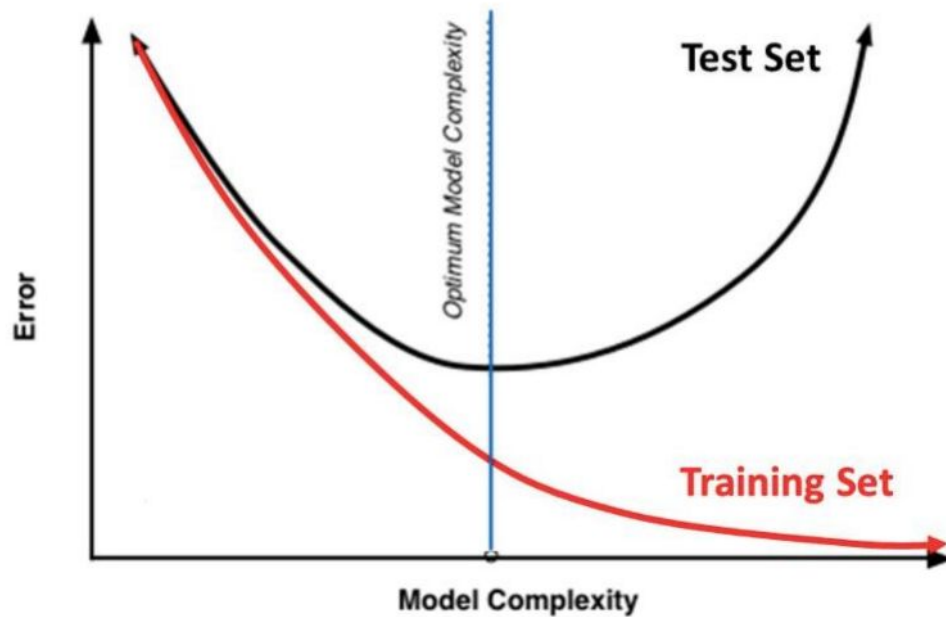
# REGULARIZATION



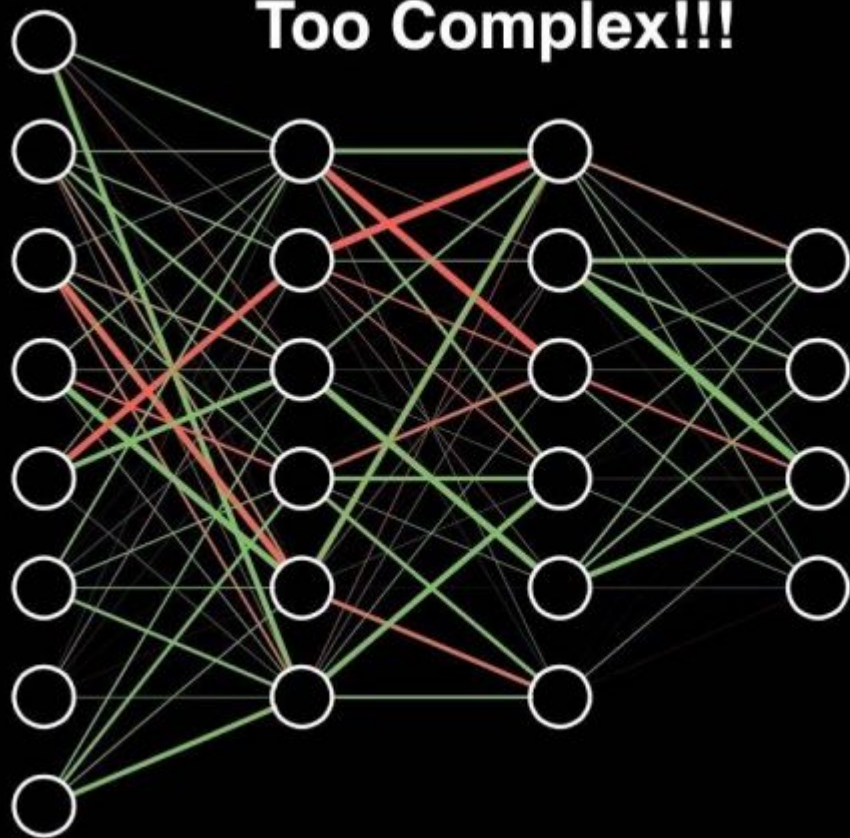
As move towards right, poor performance on unseen data

# REGULARIZATION

## Training Vs. Test Set Error



**Too Complex!!!**

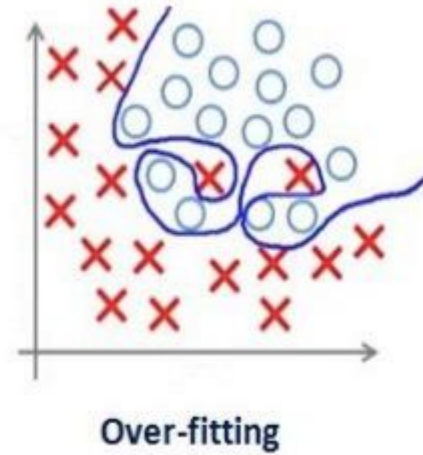
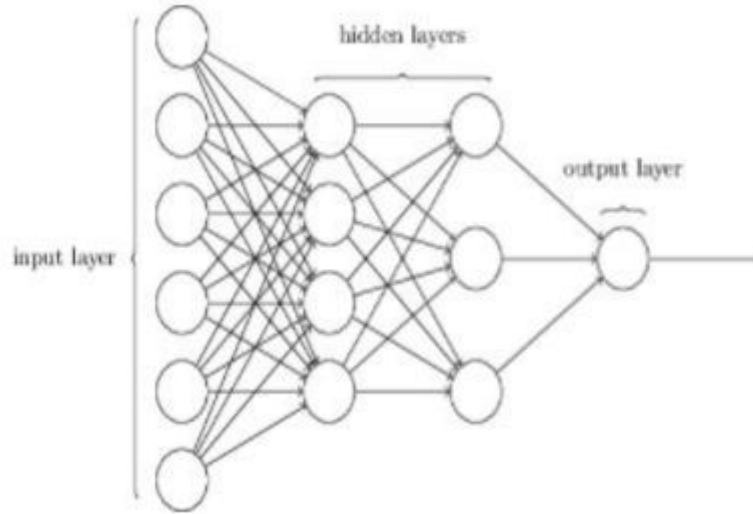


# WHAT IS REGULARIZATION ?

Regularization is a technique which makes slight modifications to the learning algorithm such that the model generalizes better

This in turn improves the model's performance on the unseen data as well

# REGULARIZATION





# REGULARIZATION

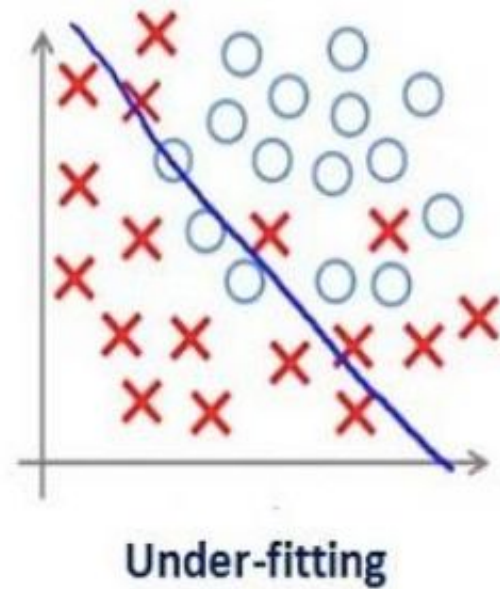
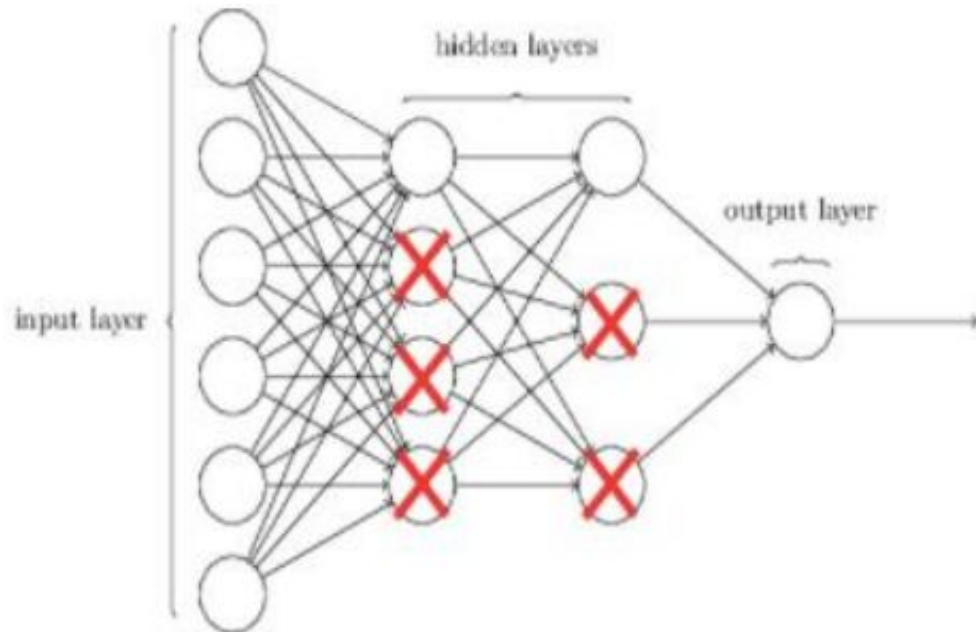
In machine learning, regularization penalizes the coefficients

In deep learning, it actually penalizes the weight matrices of the nodes

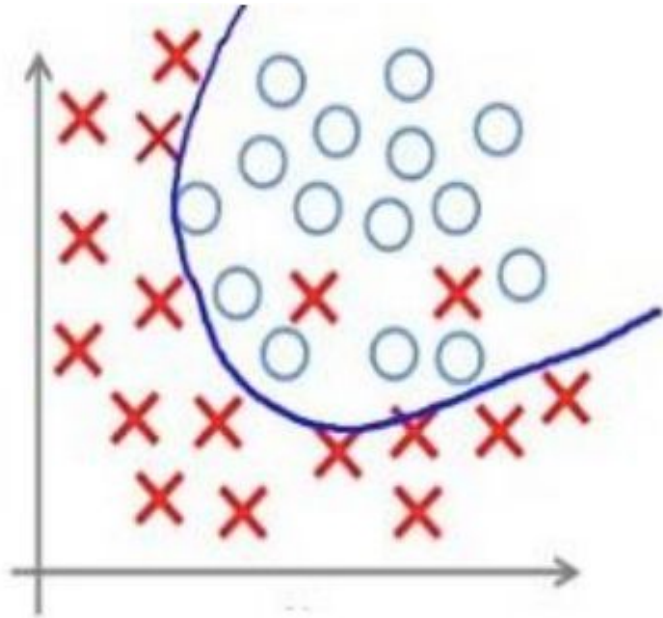
Assume that our regularization coefficient is so high that some of the weight matrices are nearly equal to zero

This will result in a much simpler linear network and slight underfitting of the training data.

# REGULARIZATION



# REGULARIZATION



**Appropriate-fitting**

Such a large value of the regularization coefficient is not that useful

We need to optimize the value of regularization coefficient in order to obtain a well-fitted model as shown in the image below

# REGULARIZATION

- L1 and L2 regularization
- DropOut
- Data Augmentation
- Early Stopping

# READING ASSIGNMENT

Read Neural Network L2 Regularization Using Python -- Visual Studio Magazine.pdf

Read Neural Network L1 Regularization Using Python -- Visual Studio Magazine.pdf

# L1 L2 REGULARIZATION

## L1 Regularization

$$\begin{array}{l} \text{Modified loss} \\ \text{function} \end{array} = \text{Loss function} + \lambda \sum_{i=1}^n |W_i|$$

## L2 Regularization

$$\begin{array}{l} \text{Modified loss} \\ \text{function} \end{array} = \text{Loss function} + \lambda \sum_{i=1}^n W_i^2$$

11

$$E = \underbrace{\frac{1}{2} * \sum (t_k - o_k)^2}_{\text{squared error}} + \underbrace{\lambda * \sum |w_i|}_{\text{L1 weight penalty}}$$

$$\Delta w_{jk} = -1 * \underbrace{\eta}_{\text{learning rate}} * \underbrace{\left[ x_j * (o_k - t_k) * o_k * (1 - o_k) \right] \pm \lambda}_{\text{signal} \quad \frac{\partial E}{\partial w_{jk}} \text{ gradient}}$$

$$w_{jk} = w_{jk} + \Delta w_{jk}$$

12

$$E = \frac{1}{2} * \sum (t_k - o_k)^2 + \frac{\lambda}{2} * \sum w_i^2$$

plain error

weight penalty

elegant math



simple math



$$\frac{\partial E}{\partial w_{jk}} \text{ gradient}$$

$$\Delta w_{jk} = \eta * [x_j * (o_k - t_k) * o_k * (1 - o_k)] + [\lambda * w_{jk}]$$

learning  
rate

signal



# DATA AUGMENTATION



4

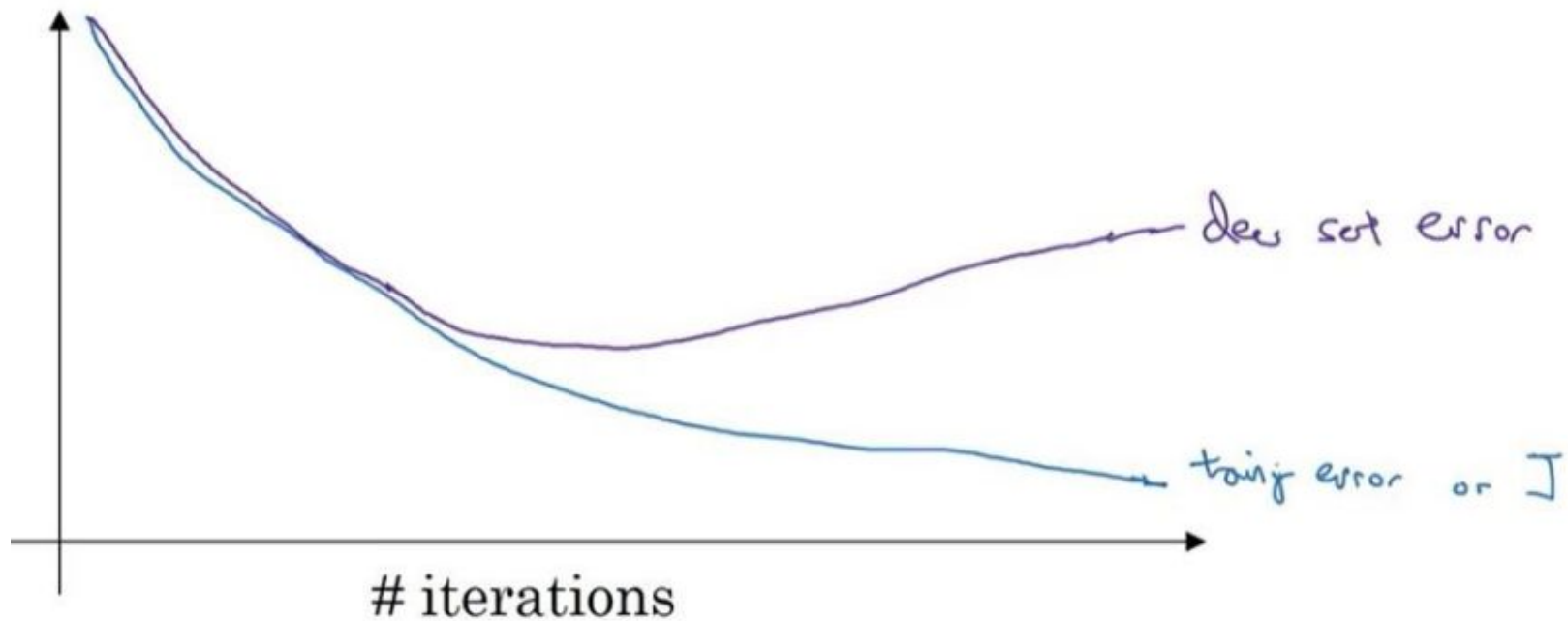


4

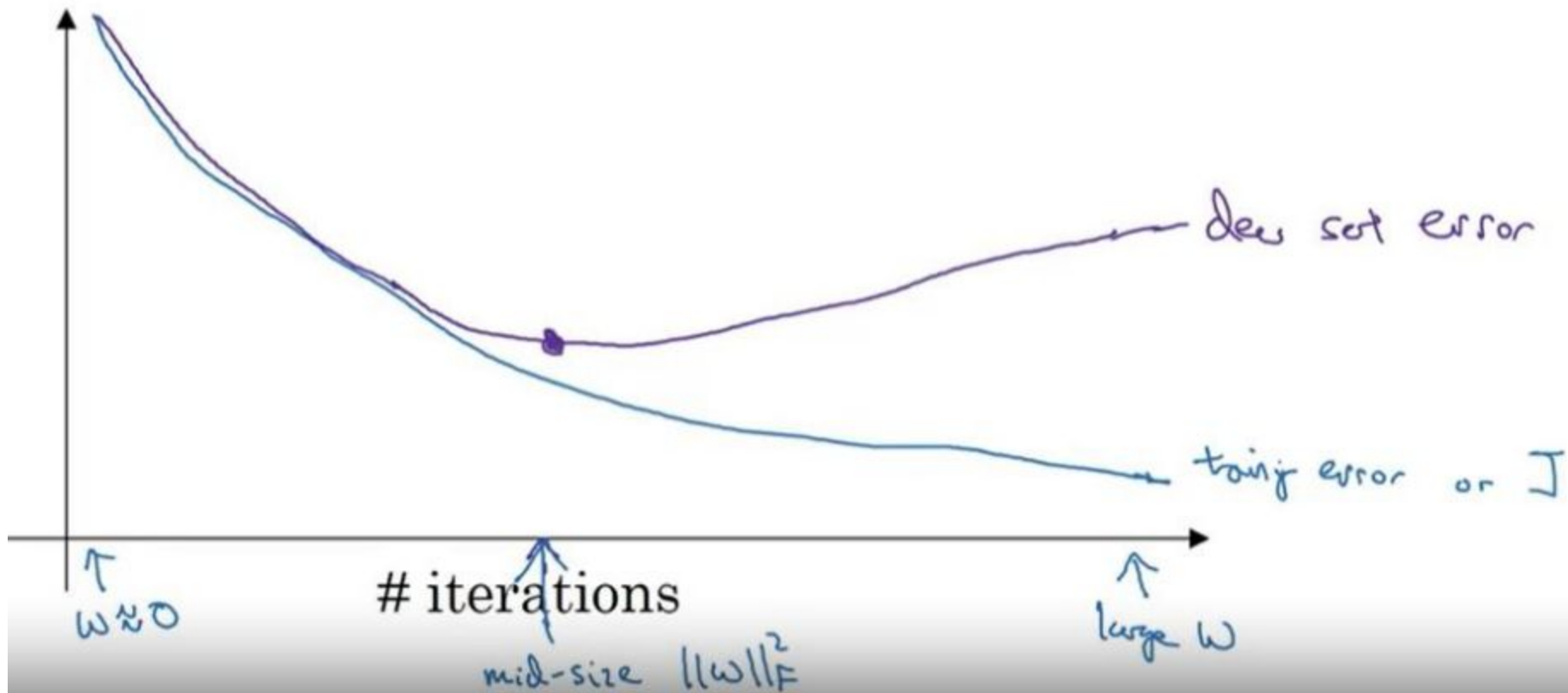
4

4

# EARLY STOPPING

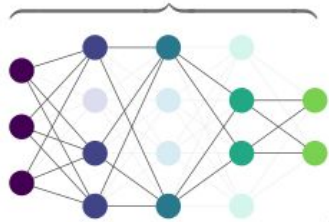


# BARLEY STOPPING



# DROP OUT

The network with **dropout** during a single forward pass



Nodes set to zero during forward passes



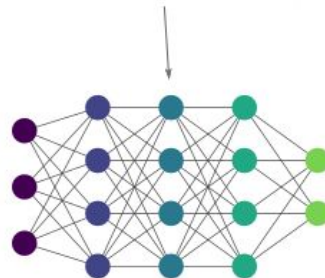
**Dropout** is the equivalent of training several independent, smaller networks on the same task. The final model is like an ensemble of smaller networks, reducing variance and providing more robust predictions.

## Dropout

Network regularization

For each forward pass during training, set the output of each node to zero with probability  $P$ .

For testing and inference use the entire network



# REFERENCES

<https://www.jeremyjordan.me/neural-networks-training/>

[https://www.reddit.com/r/learnmachinelearning/comments/x89qsi/dropout in neural networks what it is and how it/?rdt=55837](https://www.reddit.com/r/learnmachinelearning/comments/x89qsi/dropout_in_neural_networks_what_it_is_and_how_it/?rdt=55837)

<https://www.analyticsvidhya.com/blog/2018/04/fundamentals-deep-learning-regularization-techniques/>

<https://medium.com/towards-data-science/derivative-of-the-softmax-function-and-the-categorical-cross-entropy-loss-ffceefc081d1>