

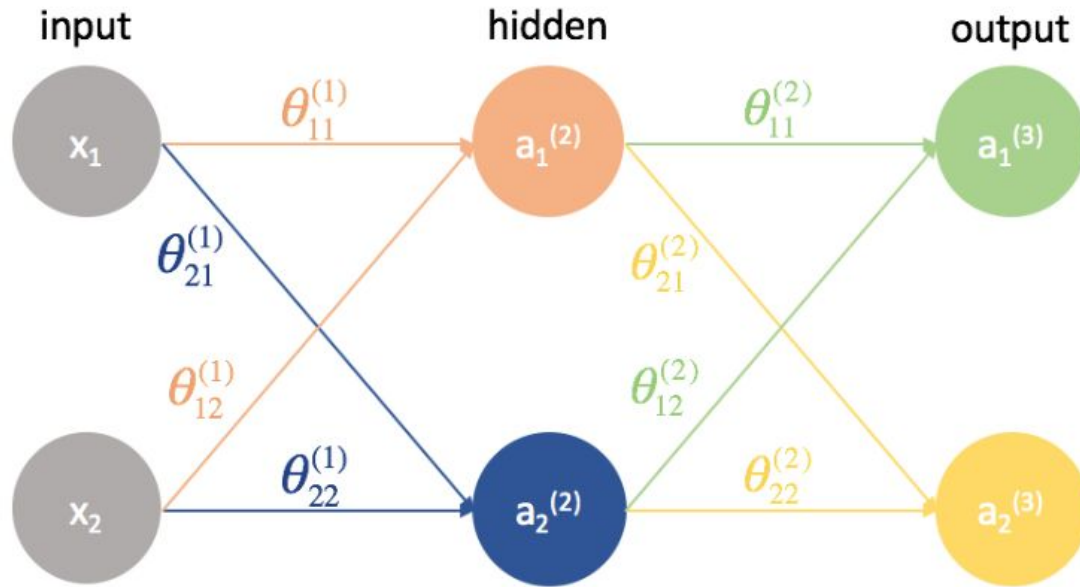
THIS IS CS5045!

GCR:ioc7cdl

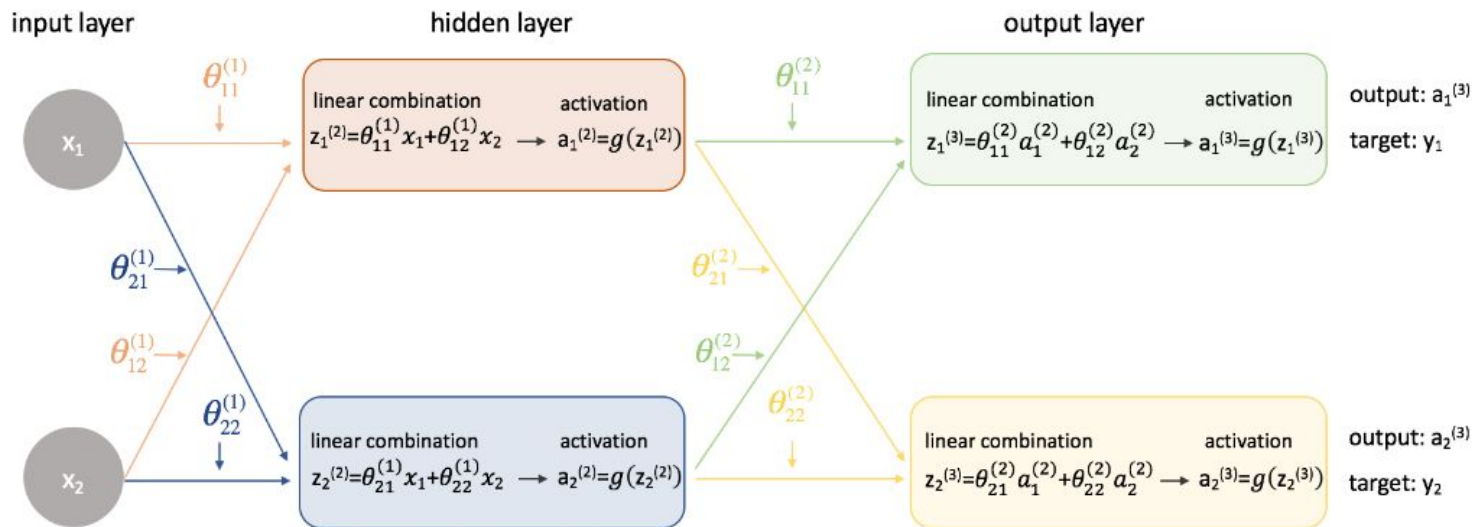
P.S. THESE SLIDES ARE USELESS IF YOU DO
NOT ATTEND CLASSES

NEURAL NETWORKS

NEURAL NETWORK WITH MULTI OUTPUT



NEURAL NETWORK WITH MULTI OUTPUT

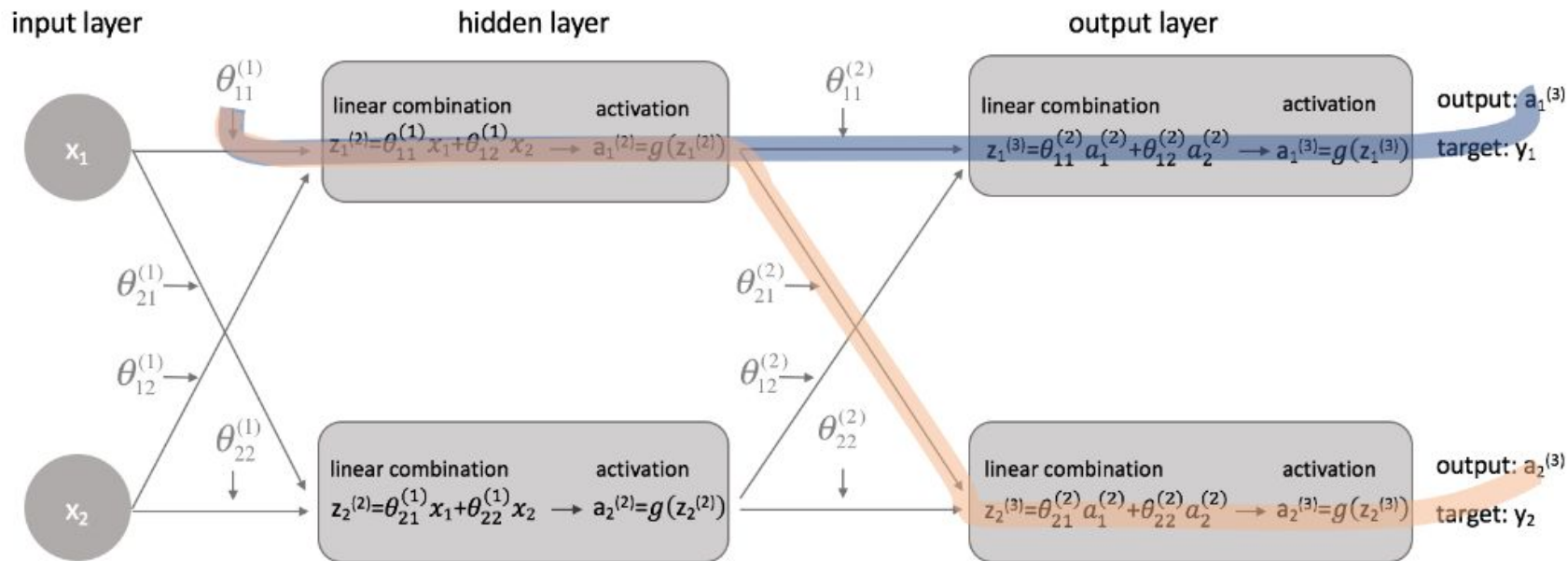


NEURAL NETWORK WITH MULTI OUTPUT

Loss calculation

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

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

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

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

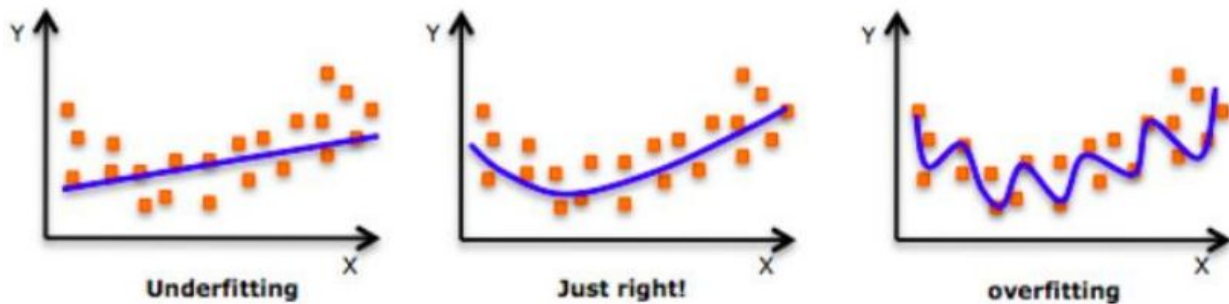
$$\frac{\partial J(\theta)}{\partial \theta_{21}^{(2)}} = \left(\frac{\partial J(\theta)}{\partial \mathbf{a}_2^{(3)}} \right) \left(\frac{\partial \mathbf{a}_2^{(3)}}{\partial z_2^{(3)}} \right) \left(\frac{\partial z_2^{(3)}}{\partial \theta_{21}^{(2)}} \right)$$

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

DISCUSSION ON ASSIGNMENT NETWORK

REGULARIZATION

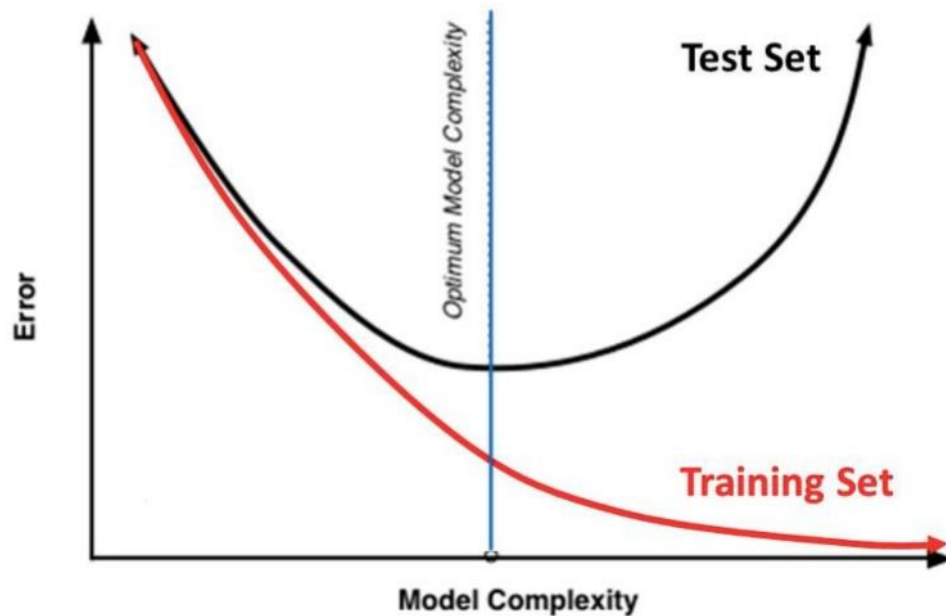
REGULARIZATION



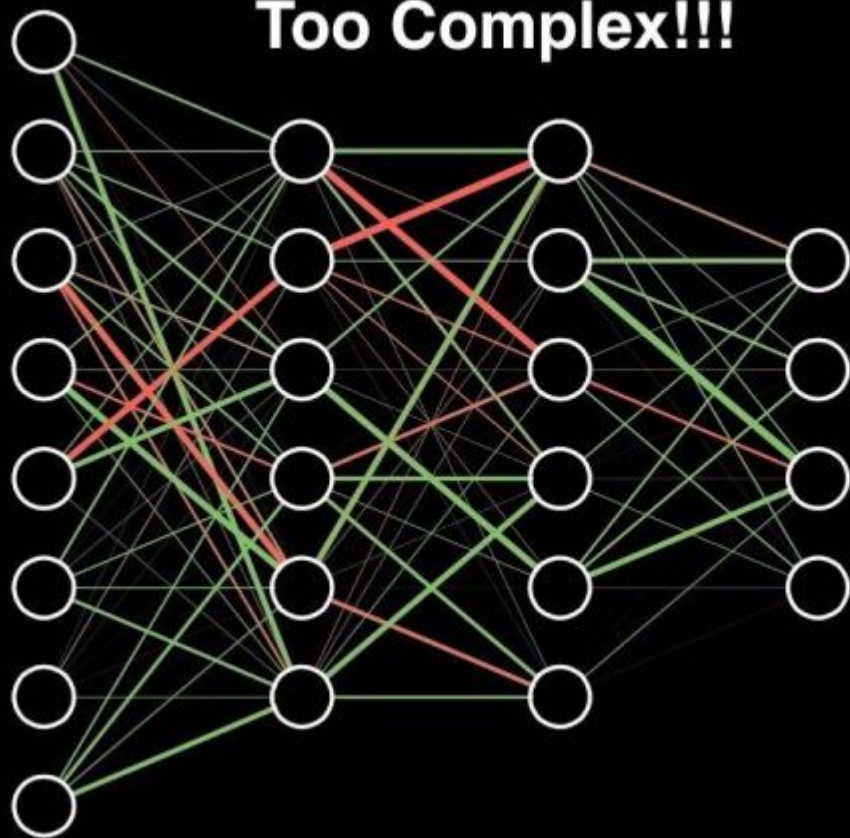
As move towards right, poor performance on unseen data

REGULARIZATION

Training Vs. Test Set Error



Too Complex!!!

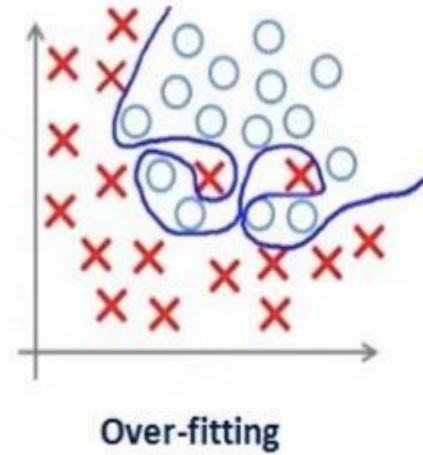
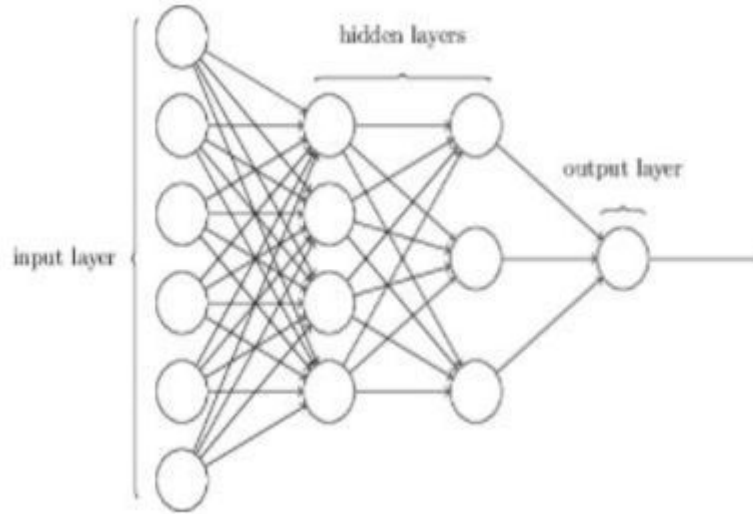


WHAT IS REGULARIZATION ?

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

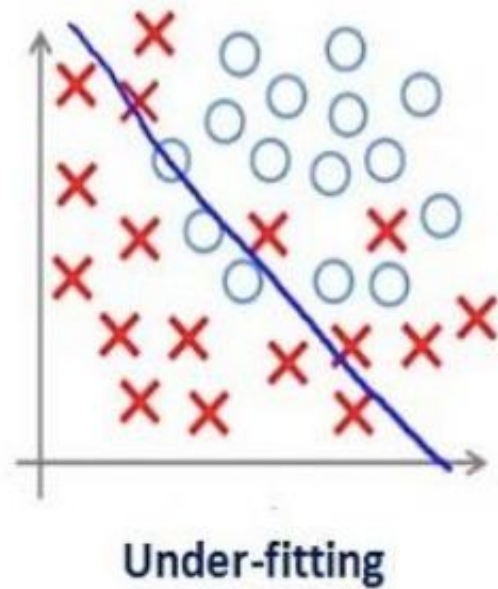
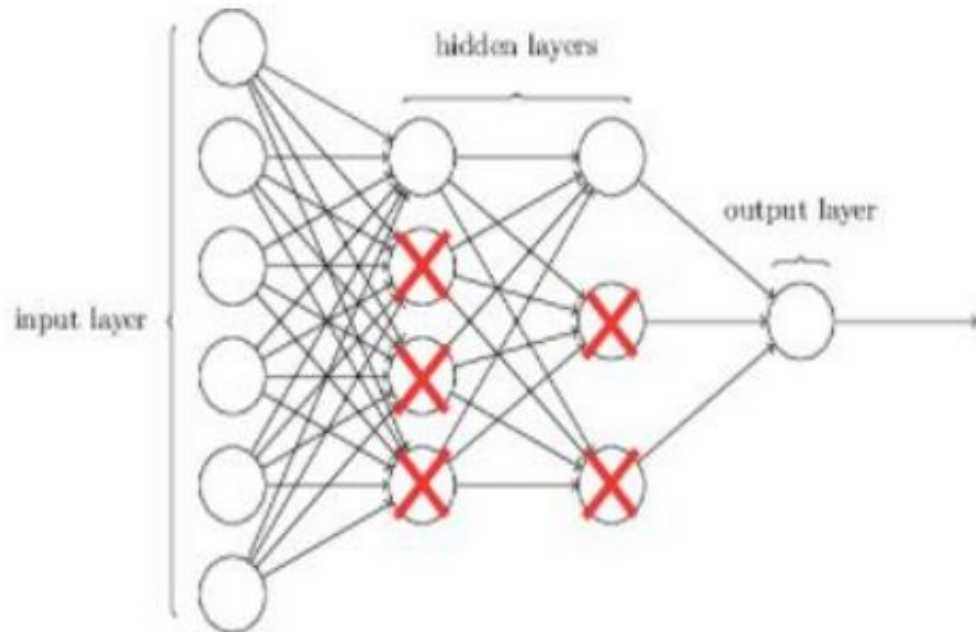
In machine learning, regularization penalizes the coefficients

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

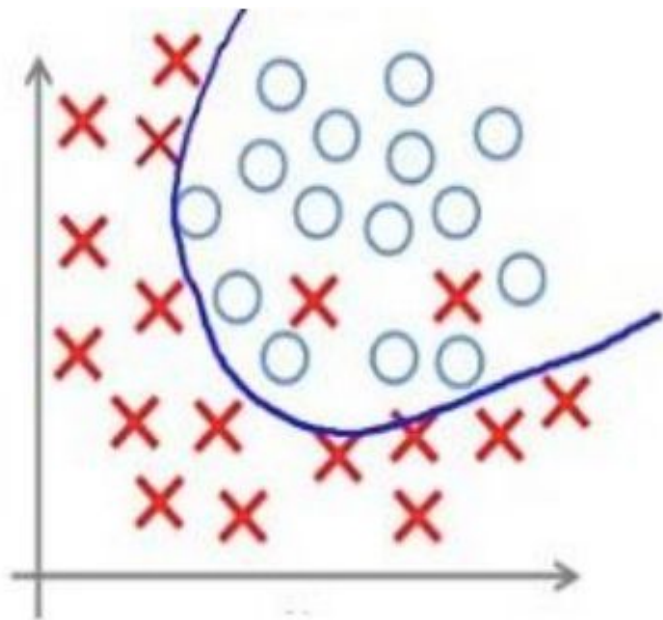
Assume that our regularization coefficient is so high that some of the weight matrices are nearly equal to zero

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



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

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

learning
rate

signal

DATA AUGMENTATION



4

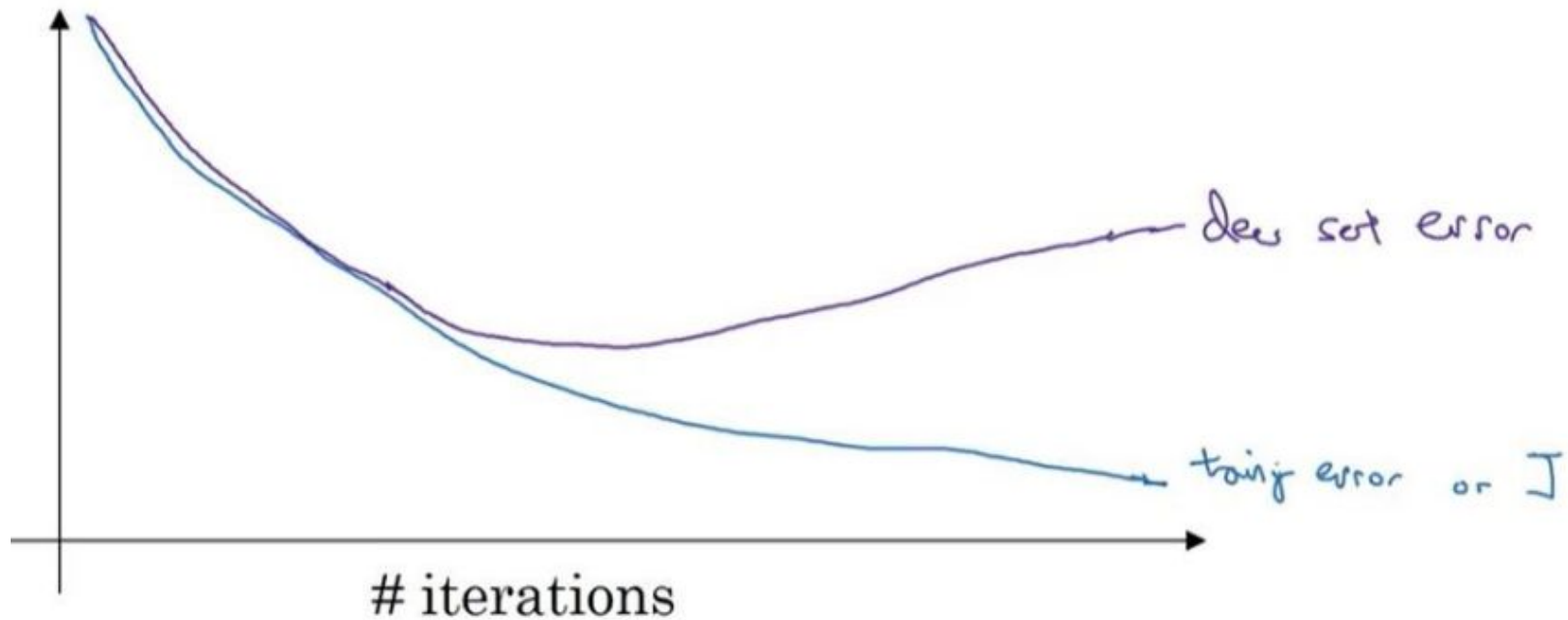


4

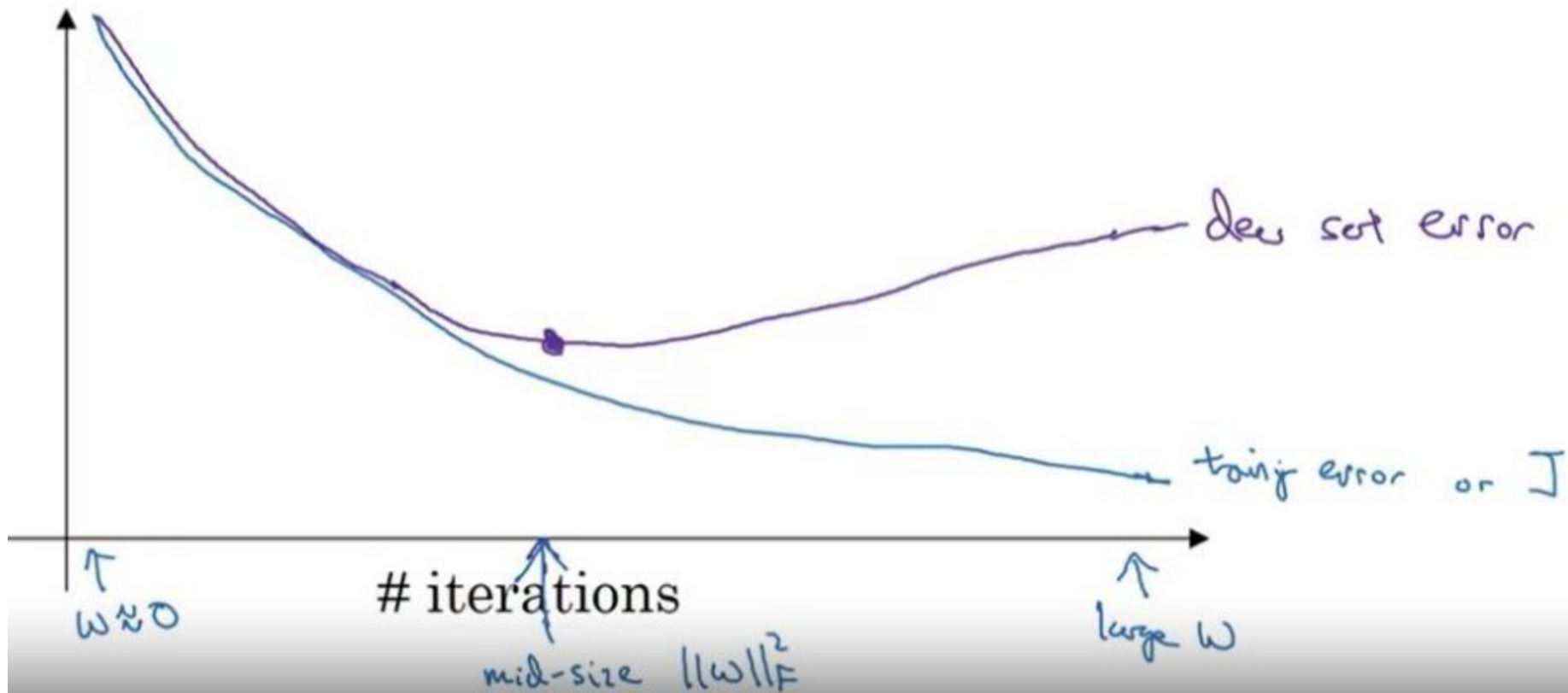
4

4

EARLY STOPPING

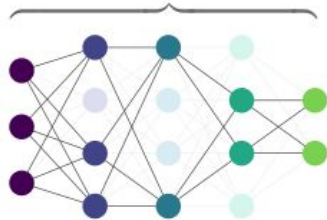


BARON STOPPING



DROP OUT

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



Nodes set to zero during forward passes



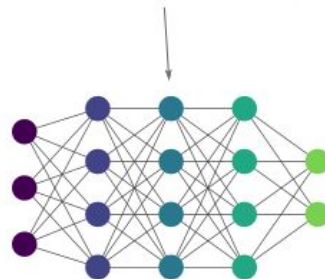
Dropout is the equivalent of training several independent, smaller networks on the same task. The final model is like an ensemble of smaller networks, reducing variance and providing more robust predictions.

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/>