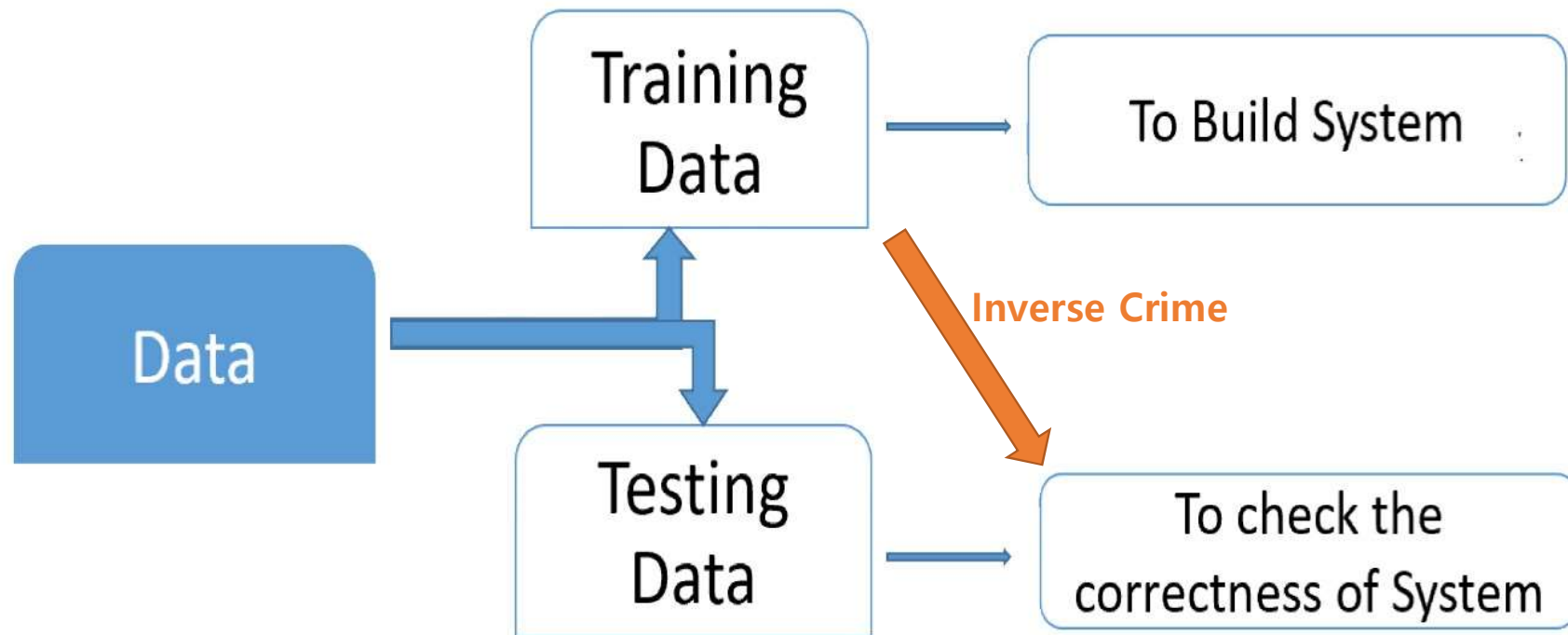
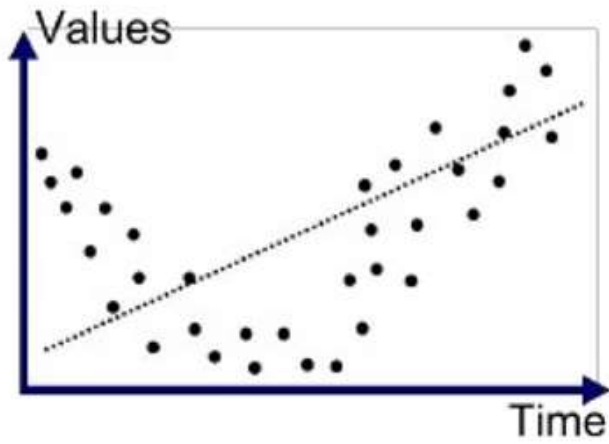


MDL Ch. 1 Handout

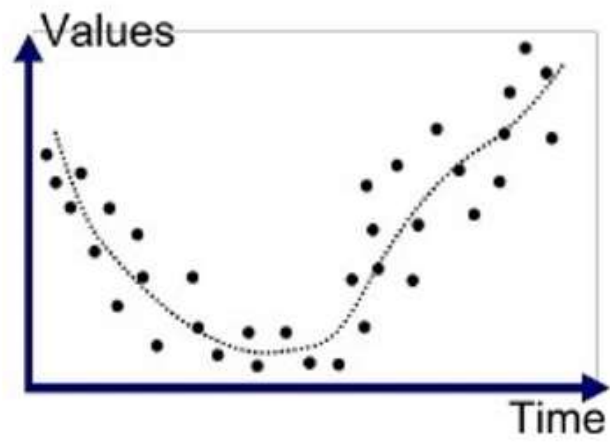
Inverse Crime



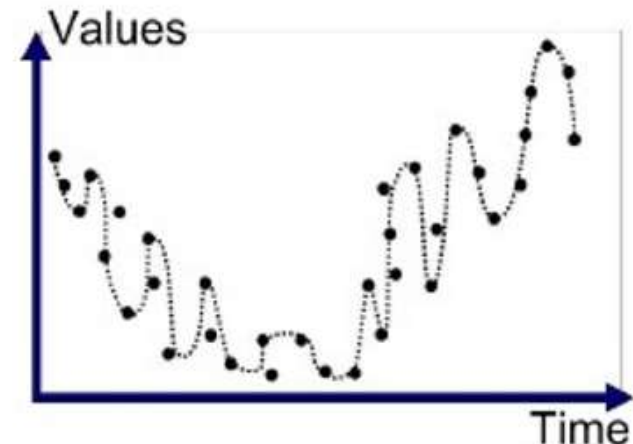
Overfitting (1D)



Underfitted



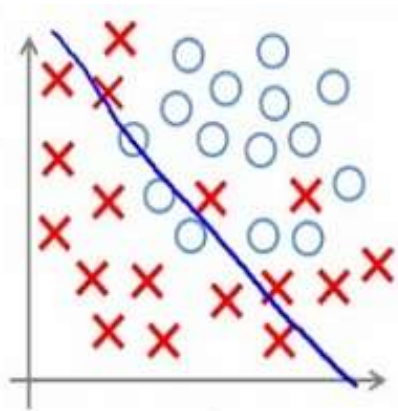
Good Fit/Robust



Overfitted

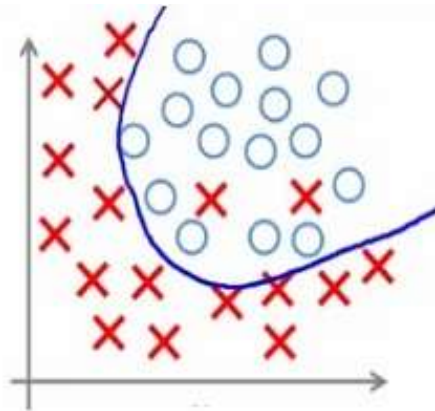
Overfitting (2D)

Feature Space

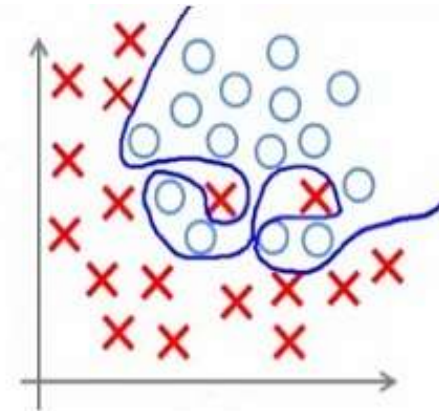


Under-fitting

(too simple to explain the variance)

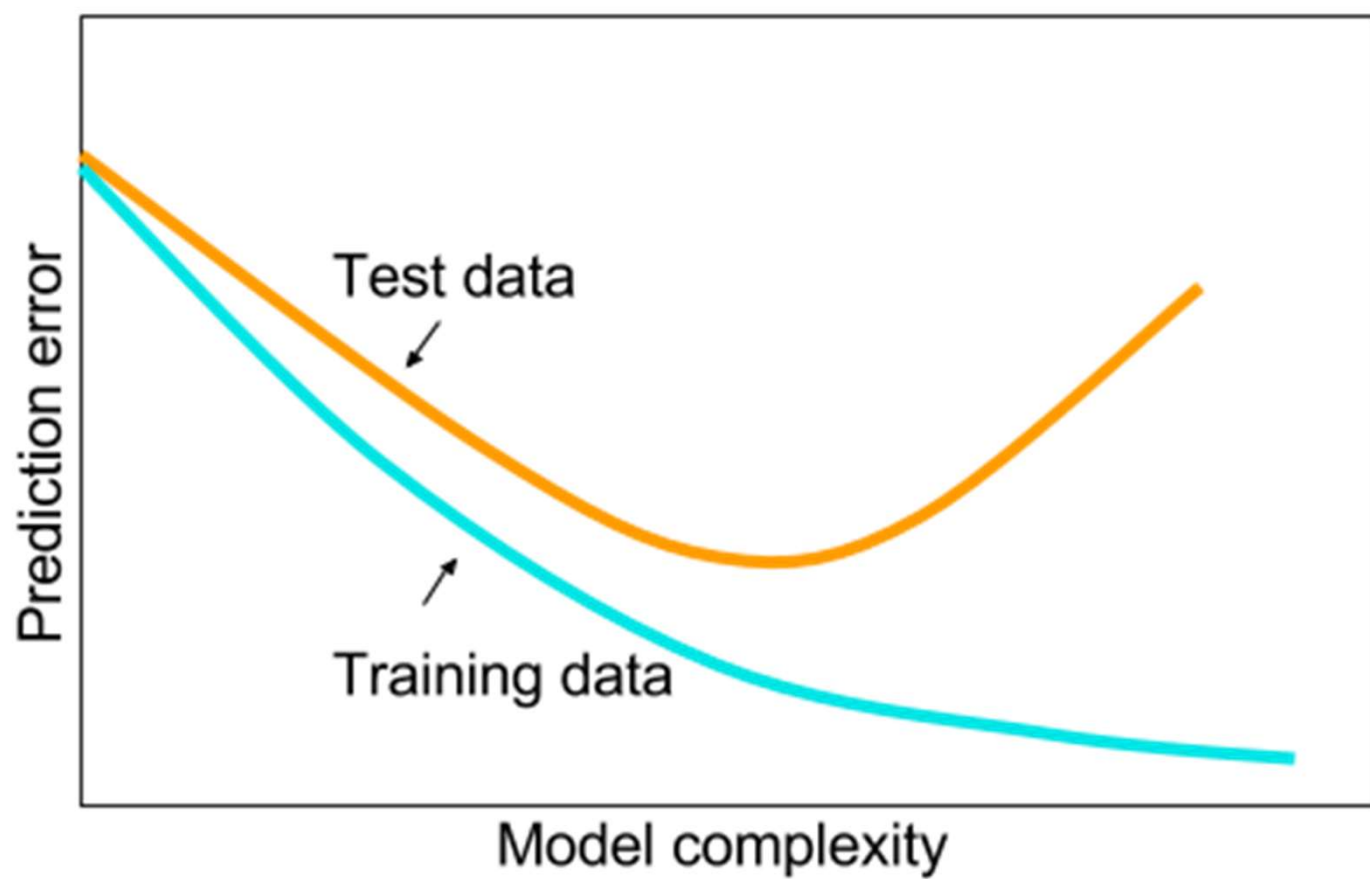


Appropriate-fitting



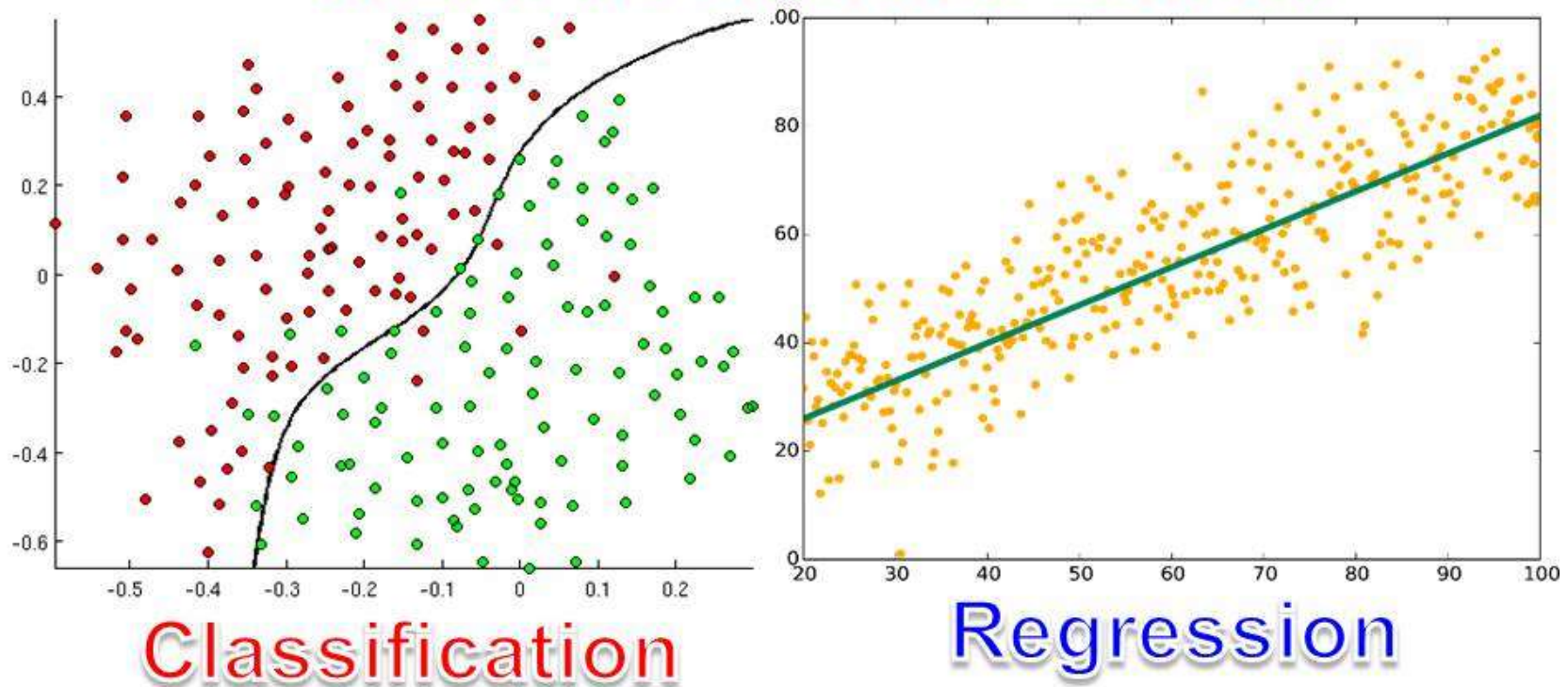
Over-fitting

(forcefitting -- too good to be true)



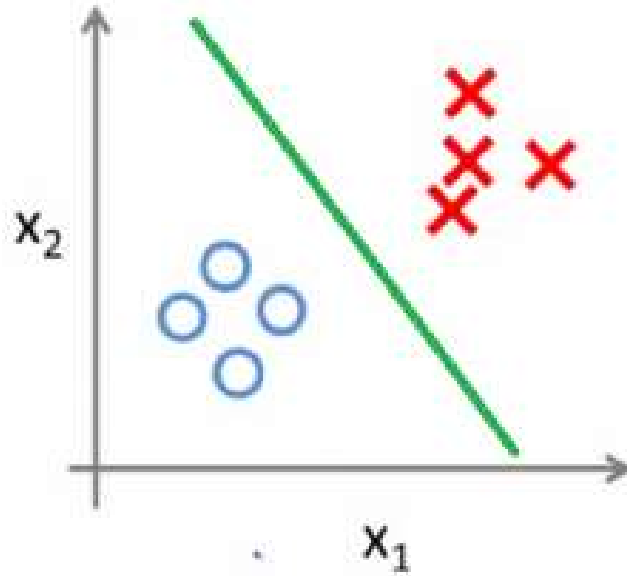
Classification vs. Regression

What is the Difference Between

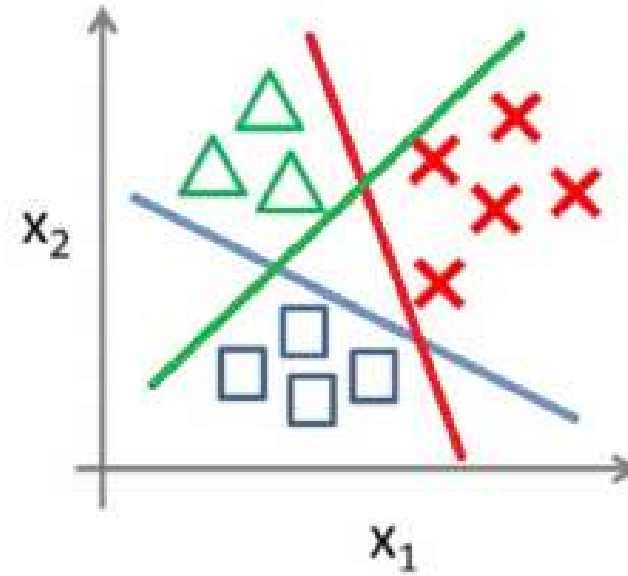


What is Classification?



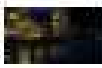


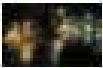






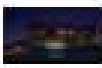

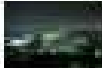
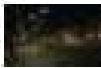




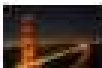

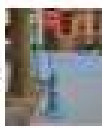
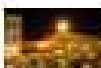
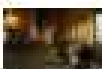



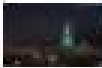

Binary classification:



Multi-class classification:



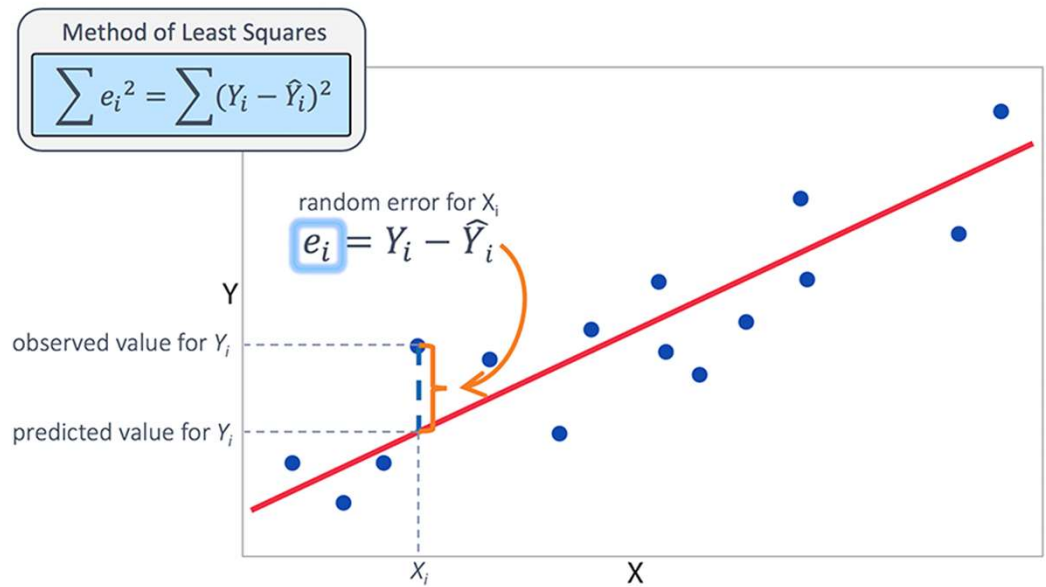
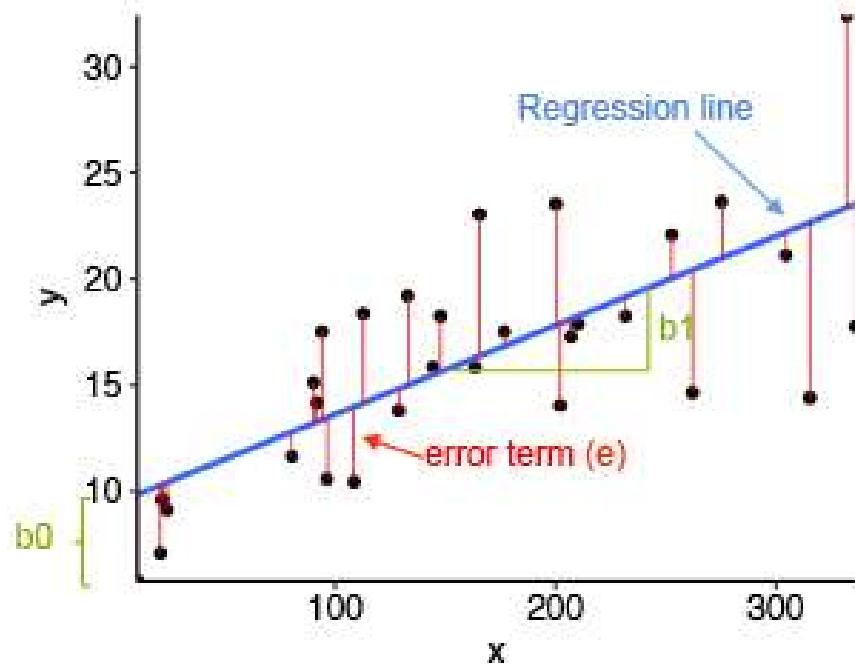
Classification (Ex.)

```
daynight = Classify[{  
   → "Night",  → "Day",  → "Night",  → "Night",  
   → "Day",  → "Night",  → "Day",  → "Day",  
   → "Night",  → "Night",  → "Day",  → "Night",  
   → "Night",  → "Day",  → "Night",  → "Night",  
  
   → "Day",  → "Day",  → "Day",  → "Day",  
  
   → "Night",  → "Night",  → "Day",  → "Night",  
   → "Night",  → "Day",  → "Day",  → "Day",  
   → "Night",  → "Day"}]
```

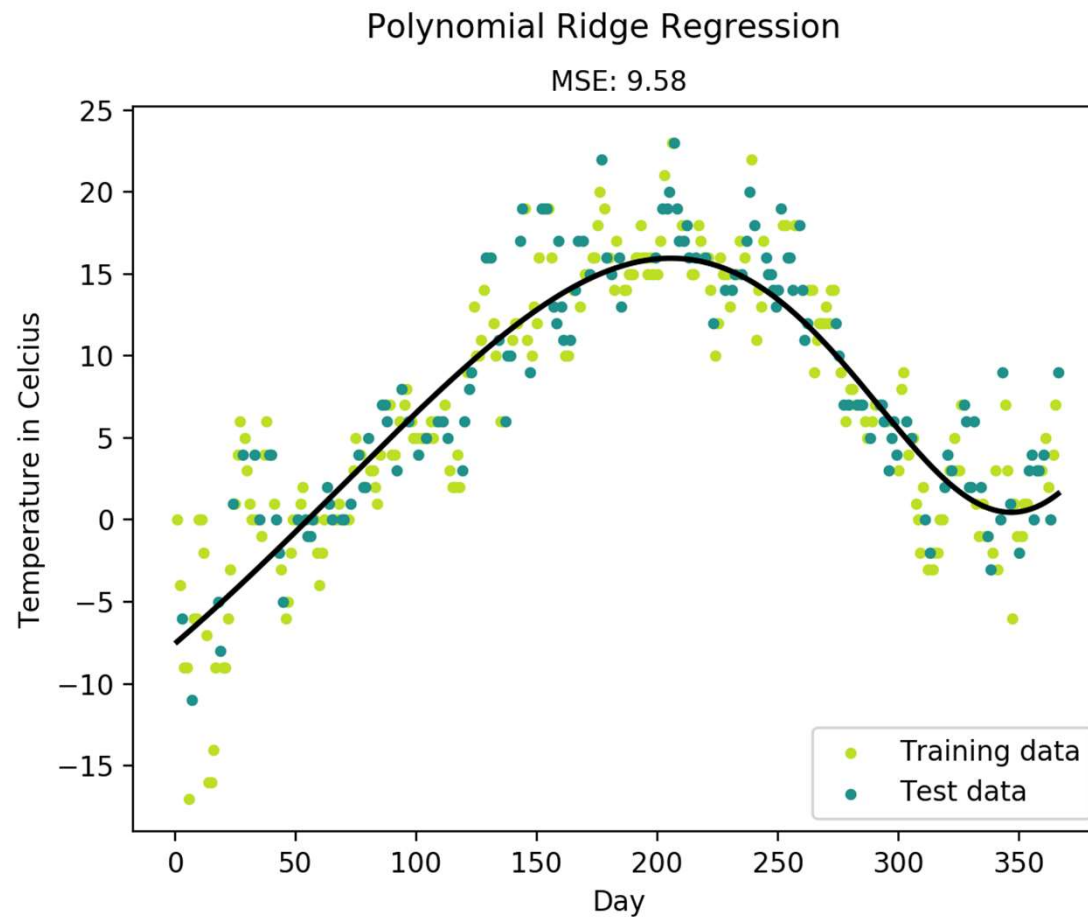

Classification (Ex.)



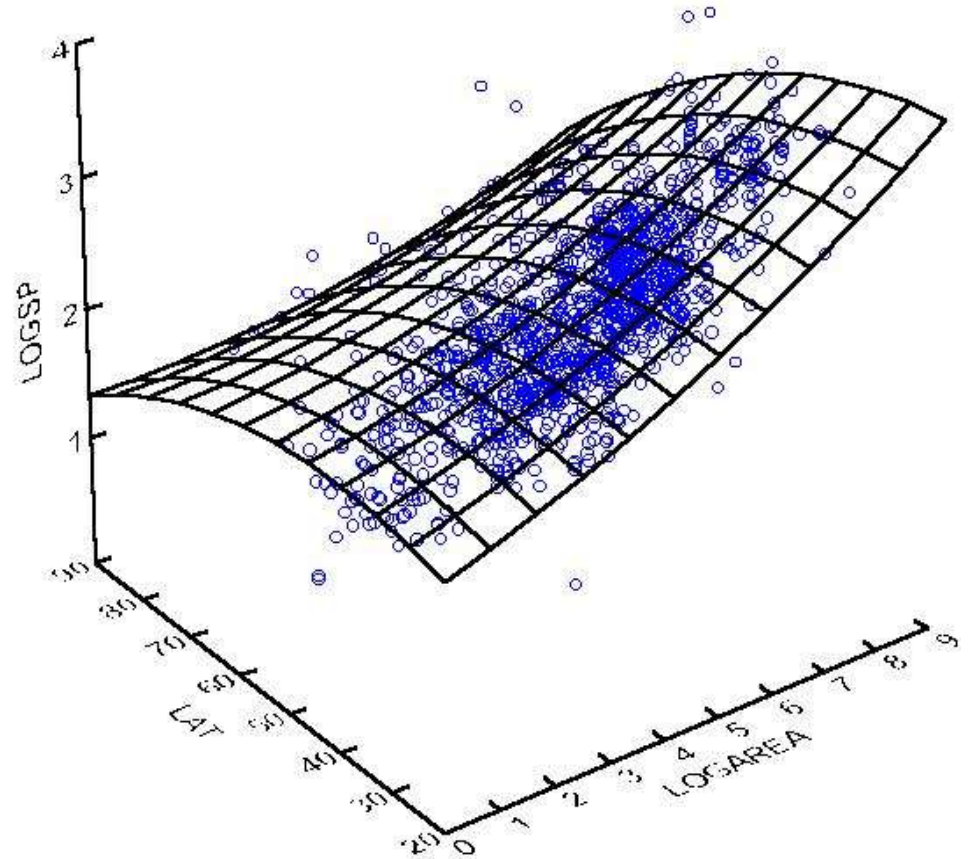
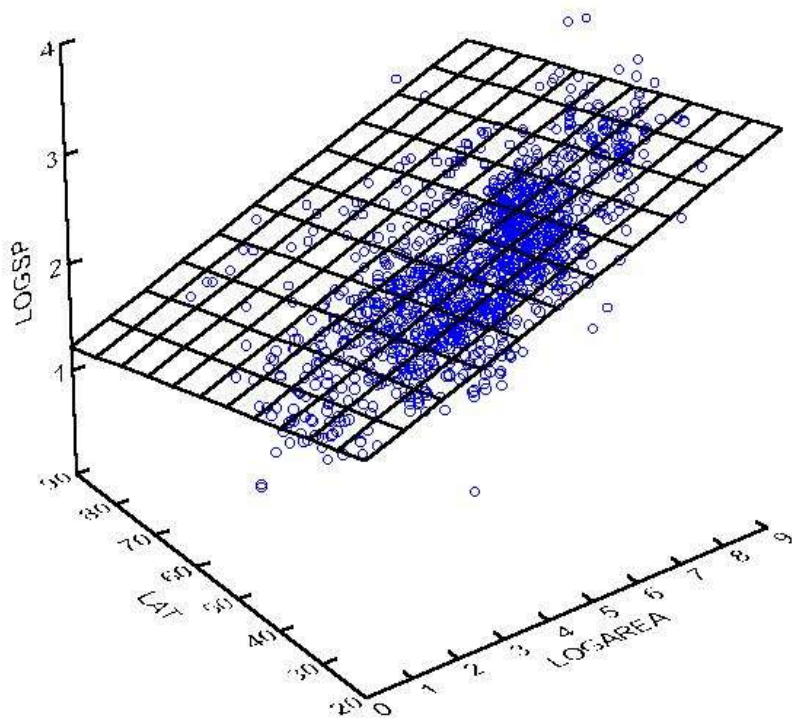
What is Regression?



Regression



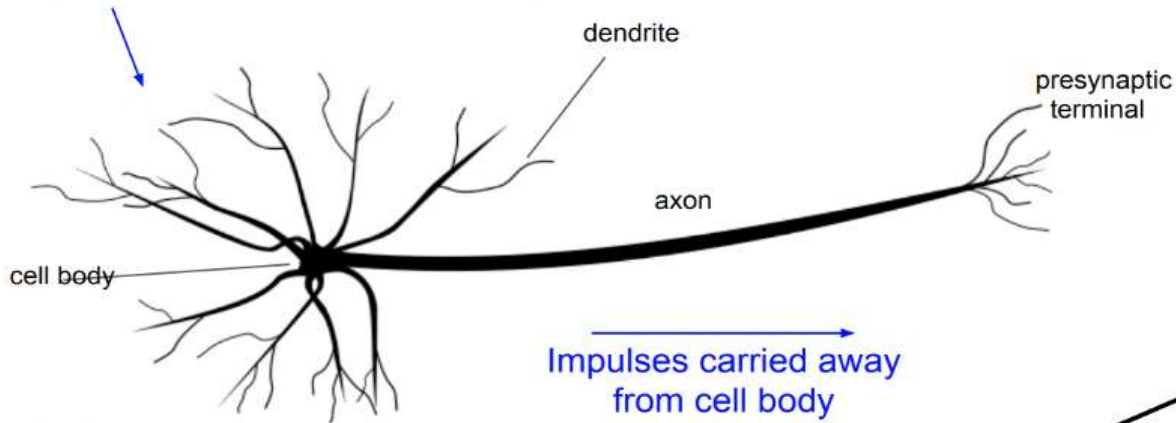
Multi-dimensional Regression



MDL Ch. 2 Handout

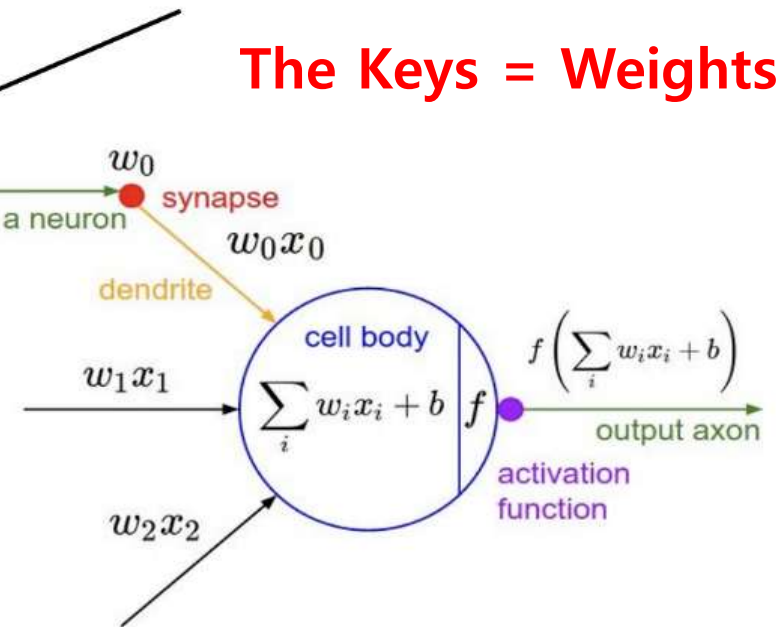
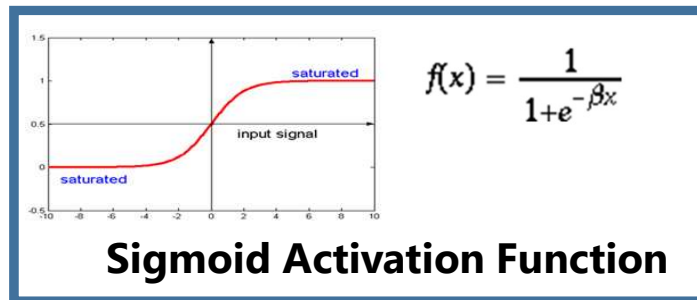
Biological Neuron vs. Artificial Neuron

Impulses carried toward cell body



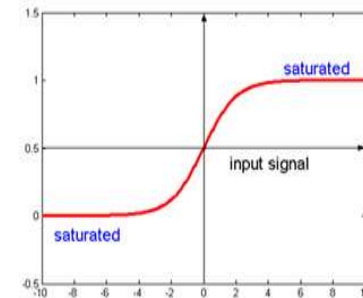
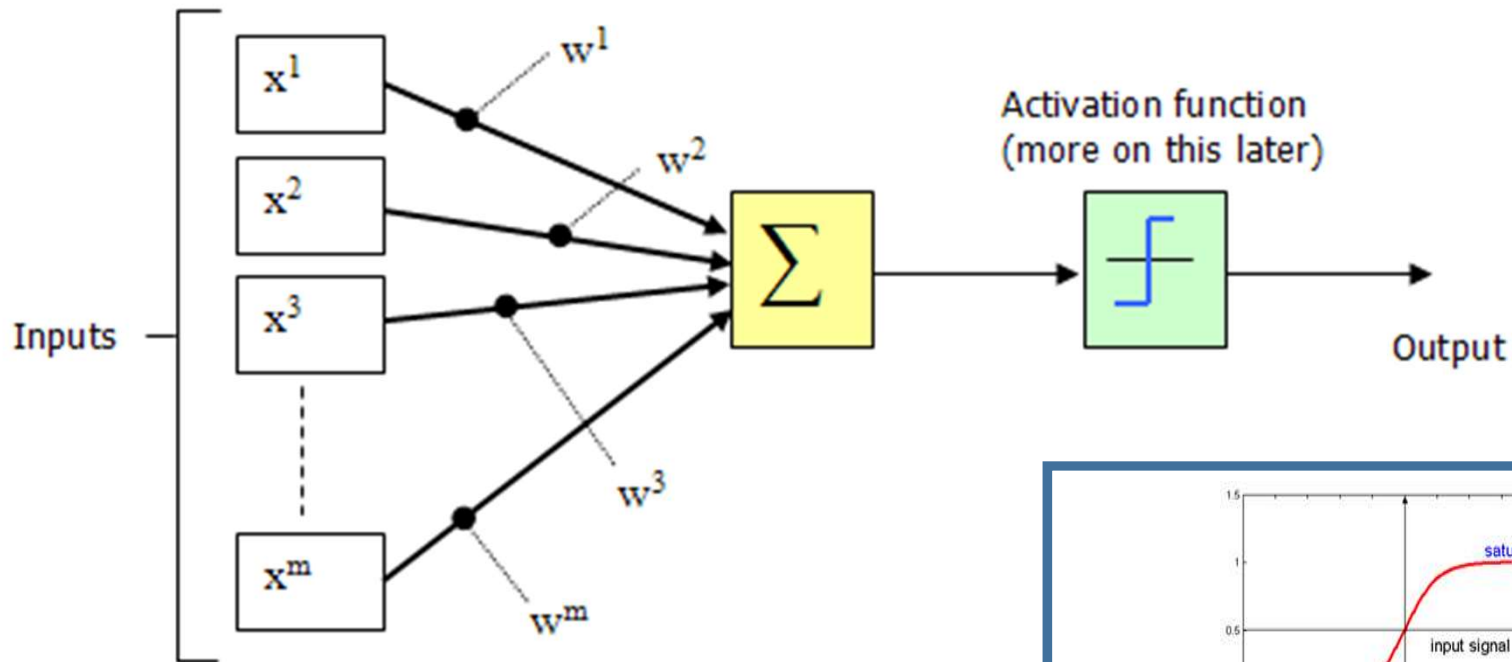
This image by Felipe Perucho is licensed under [CC-BY 3.0](https://creativecommons.org/licenses/by/3.0/)

The Keys = Weights !!



Artificial Neural Network (ANN)

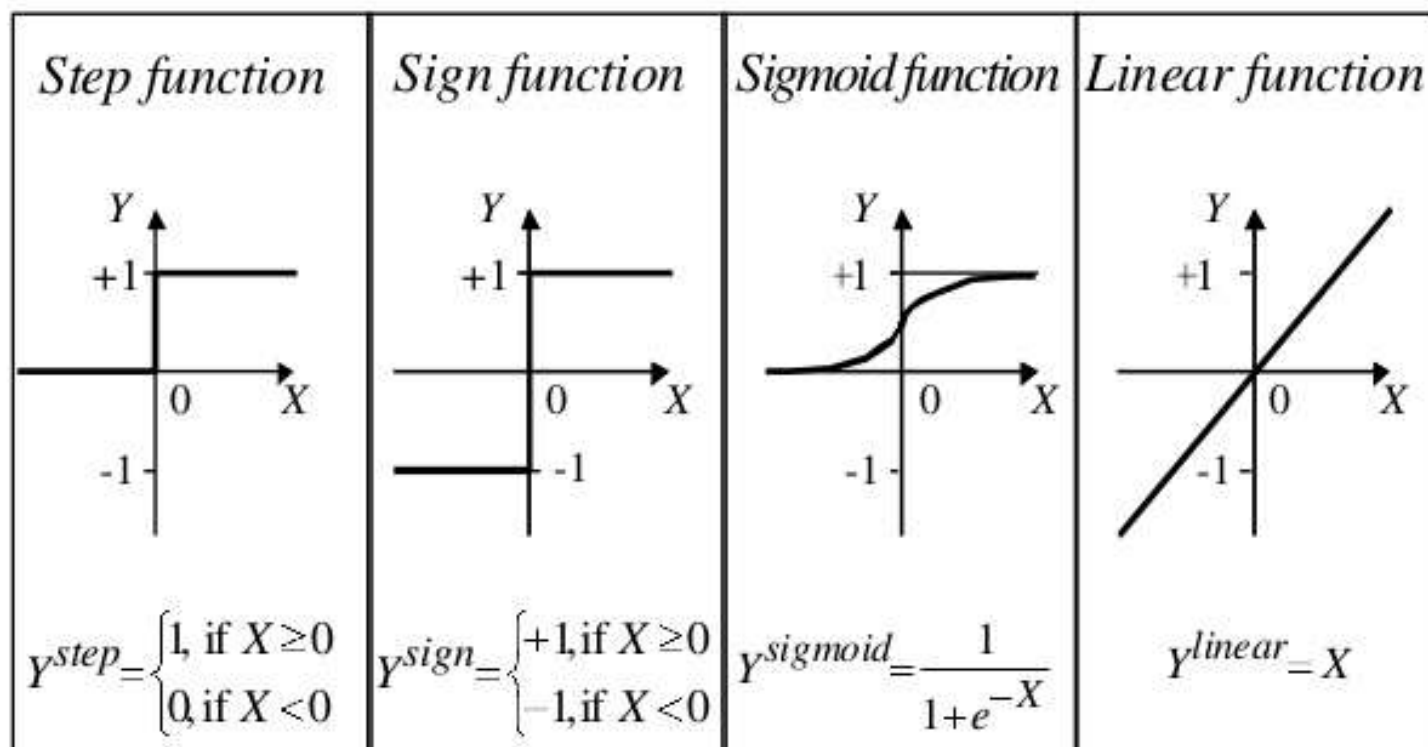
The Keys = Weights !!!



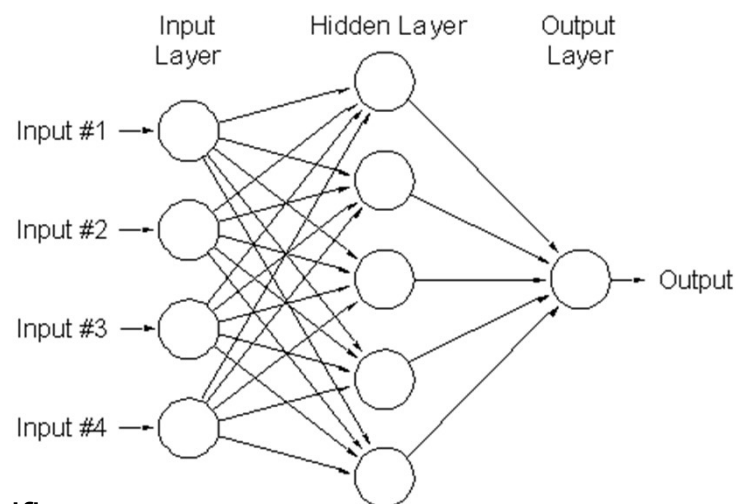
$$f(x) = \frac{1}{1+e^{-\beta x}}$$

Sigmoid Activation Function

Some Activation functions of a neuron



Multi-Layer Neural Networks



- Nonlinear classifier
- **Training: find network weights w to minimize the error** between true training labels y_i and estimated labels $f_w(\mathbf{x}_i)$:

$$E(\mathbf{w}) = \sum_{i=1}^N (y_i - f_w(\mathbf{x}_i))^2$$

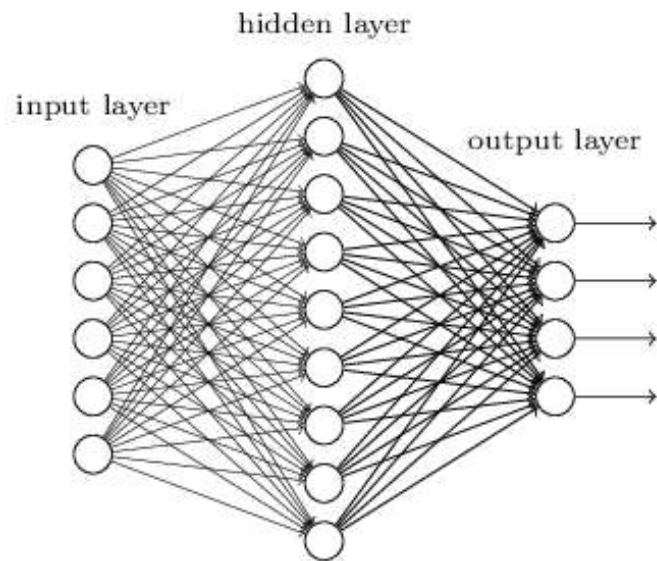
- Minimization can be done by gradient descent provided f is differentiable
 - This training method is called **back-propagation**

Shallow Network vs. Deep Network

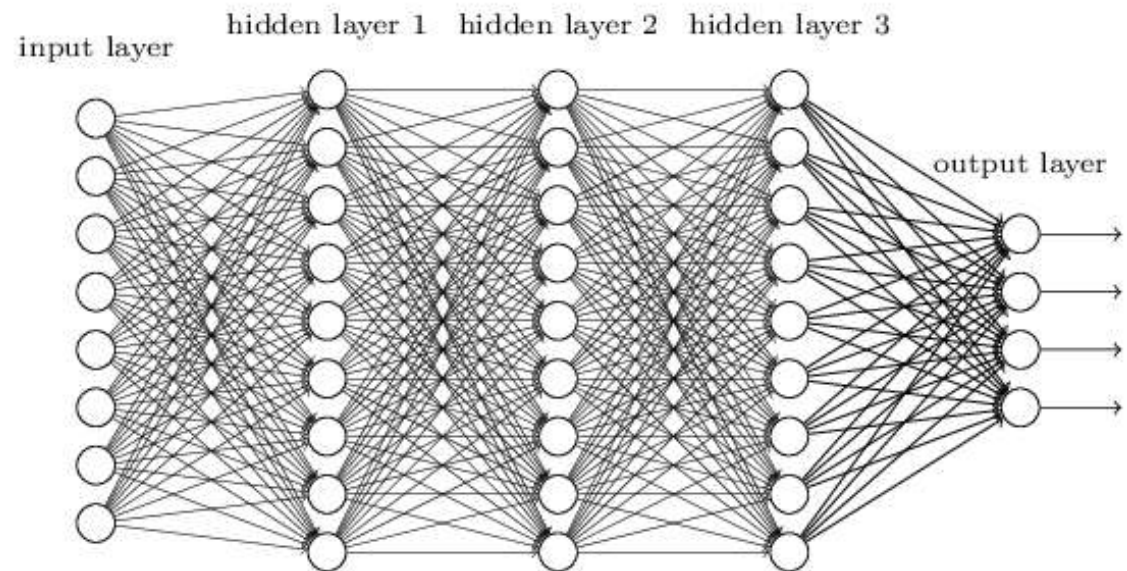


"Non-deep" feedforward neural network

Deep neural network



of Hidden Layer ≤ 1 (i.e., shallow network)

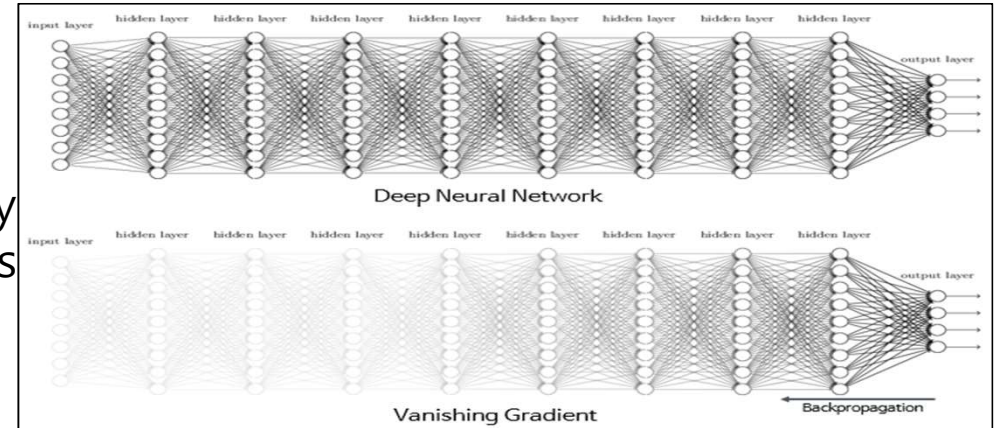


of Hidden Layer ≥ 2 (i.e., deep network)

Difficulties of Training Deep Neural Network (or Multi Layer NN)

- **Vanishing Gradient in Backpropagation**

- Problem with nonlinear activation function
- Gradient (error signal) decreases exponentially with the number of layers and the front layers train very slowly.



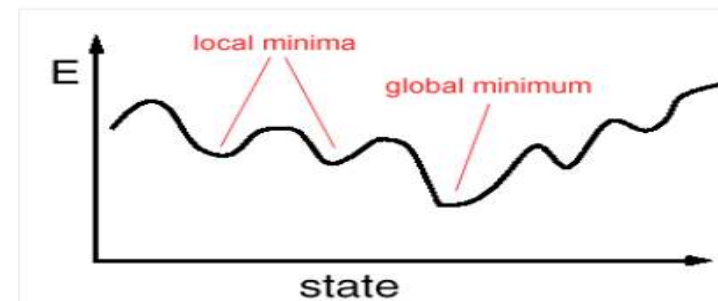
- **Over-fitting (No Generalization)**

- Given limited amounts of labeled data, training via backpropagation does not work well



- **Local Minima**

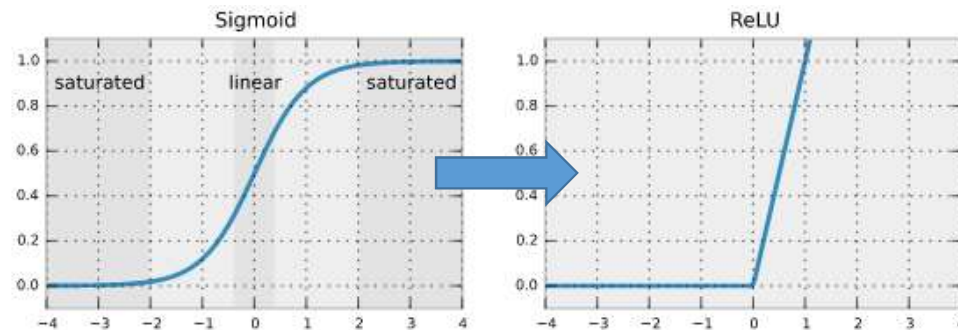
- Difficulty in optimization



New Solutions for Deep Neural Network

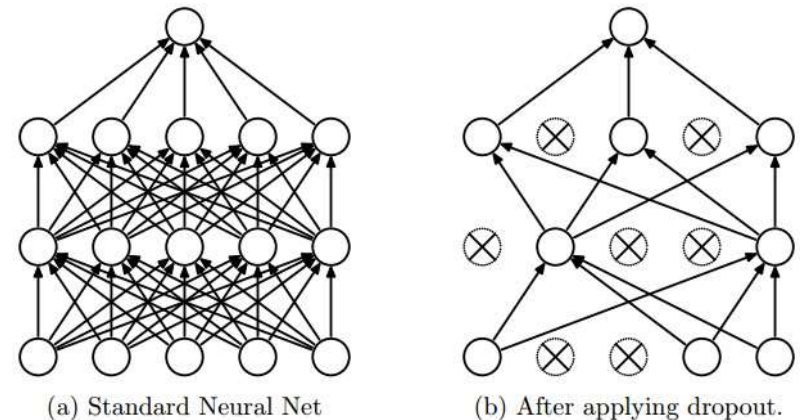
- **Vanishing Gradient**

- Solved by a new non-linear activation function: **Rectified Linear Unit (ReLU)** in 2010 & 2011



- **Over-fitting**

- Solved by new regularization methods: **dropout** (Hinton et al., 2012) etc.



- **Local minima**

- Solved by **high-dimensional non-convex optimization**: local minima are all similar
- Local minima are good and close to global minima

Delta Rule : Perceptron Rule

Perceptron Training Rule

- Weights modified for each example
- Update Rule:

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta(t - o)x_i$$

learning rate target value perceptron output input value

In textbook

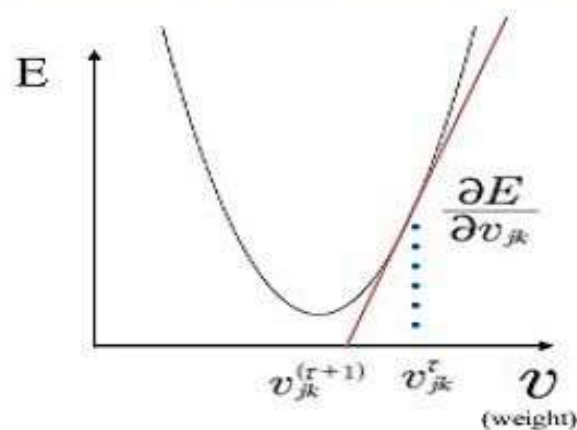
$$w_{ij} \leftarrow w_{ij} + \alpha e_i x_j$$

(Equation 2.2)

Gradient Descent in Delta Rule

How to Learn Perceptron?

Delta Learning Rule



$$v_{jk}^{(\tau+1)} = v_{jk}^{\tau} + \Delta v_{jk}$$

$$\Delta v_{jk} = -\eta \frac{\partial E}{\partial v_{jk}}$$

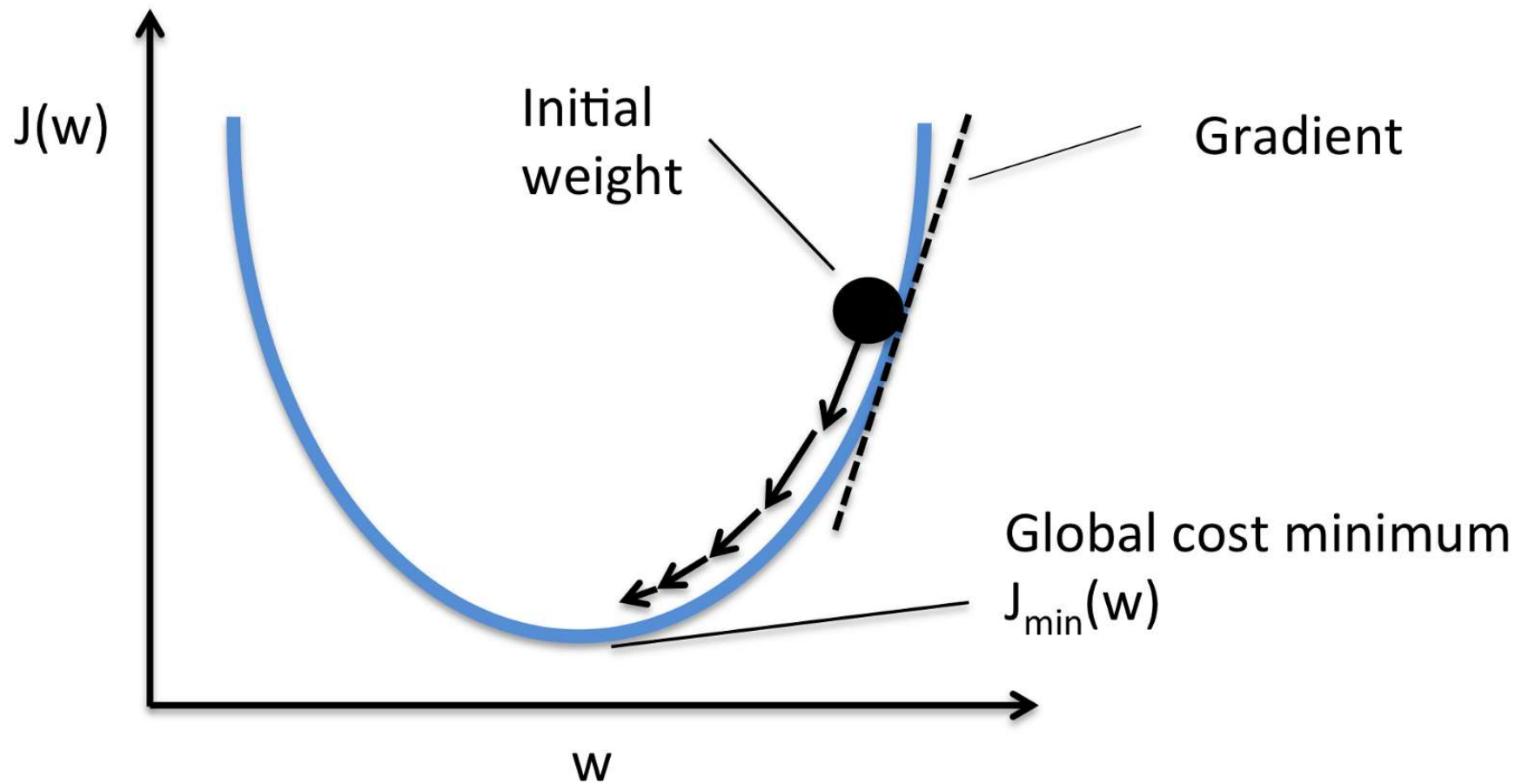
$v_{jk}^{(\tau+1)}$ new weight
 v_{jk}^{τ} current weight
 η learning rate
 E Error Function

$$v = \omega$$

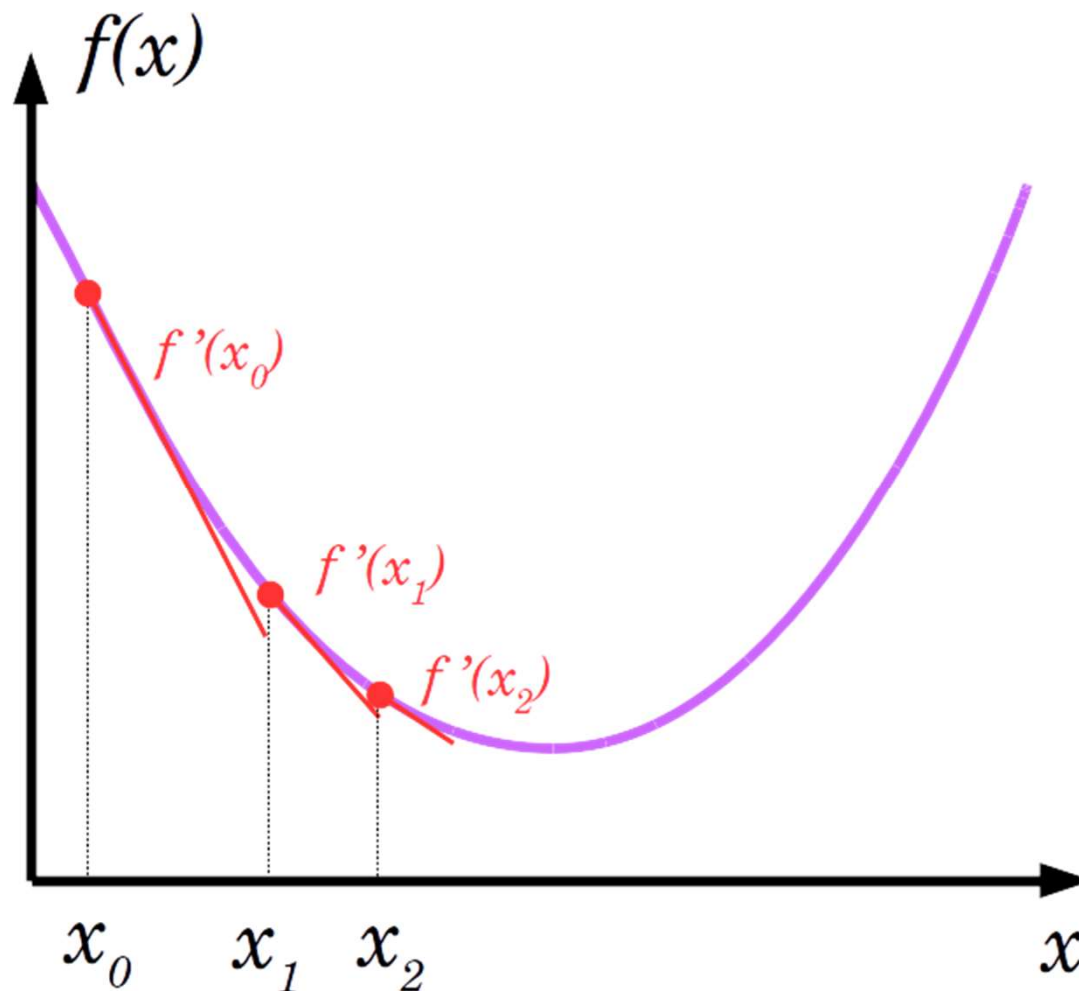
$$E_p = \frac{1}{2} \sum_n (t_{j_n} - a_{j_n})^2$$

t = target
a = net output = f(v)

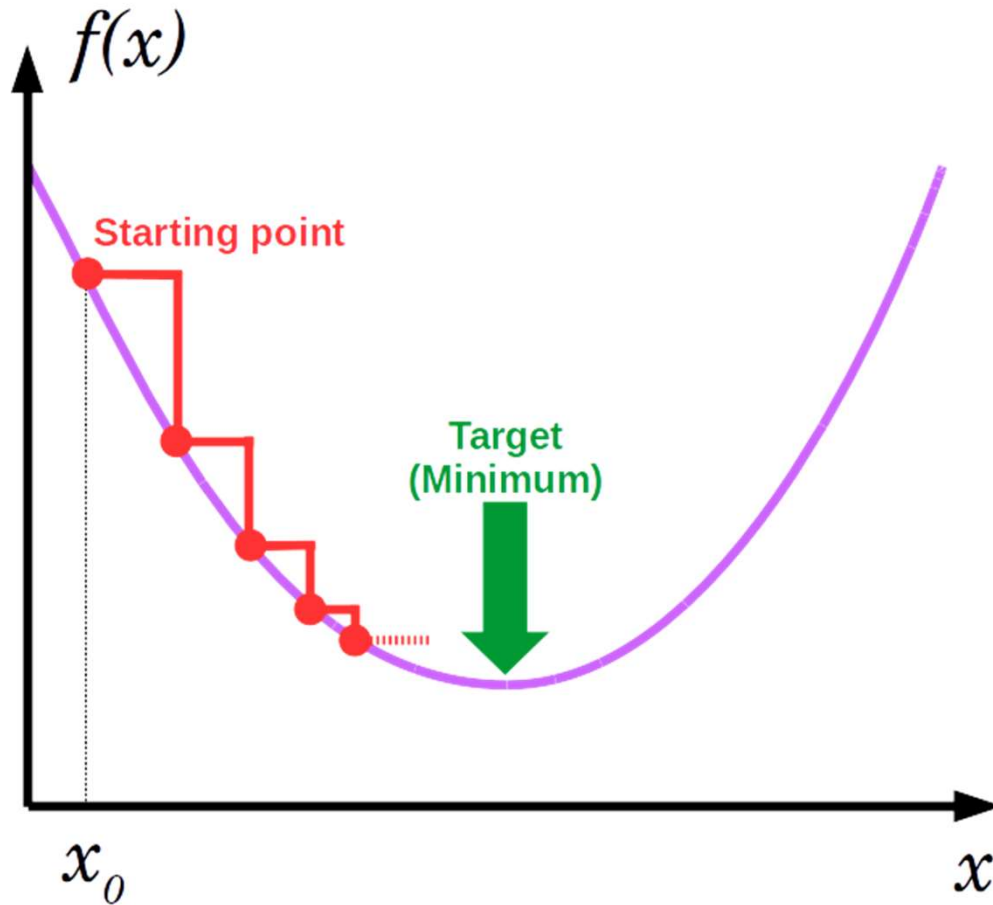
Gradient Descent (1)



Gradient Descent (2)



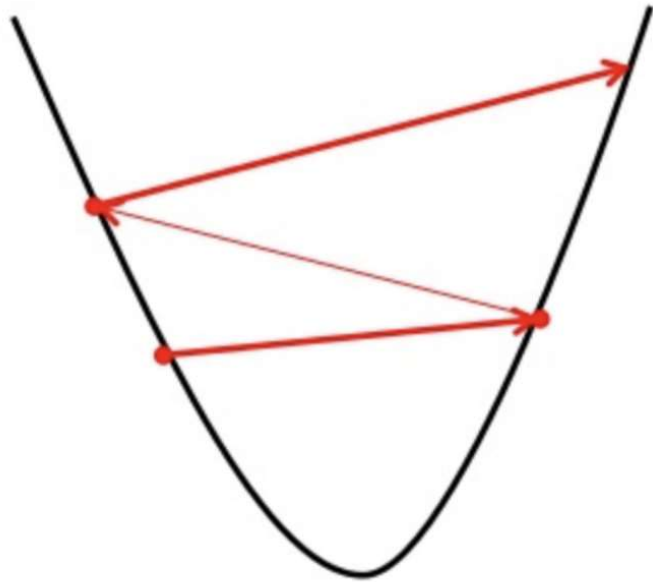
Stepping (Learning Rate) in Gradient Descent



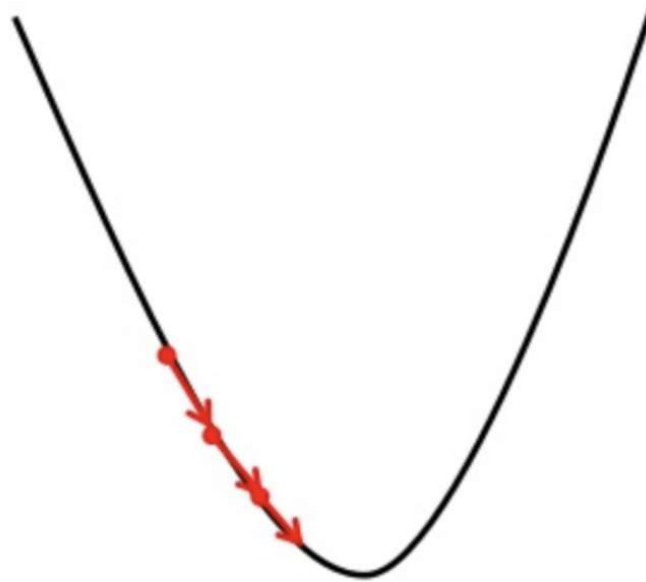
<https://www.neural-networks.io/en/single-layer/gradient-descent.php>

Effect of Learning Rate or Step Size

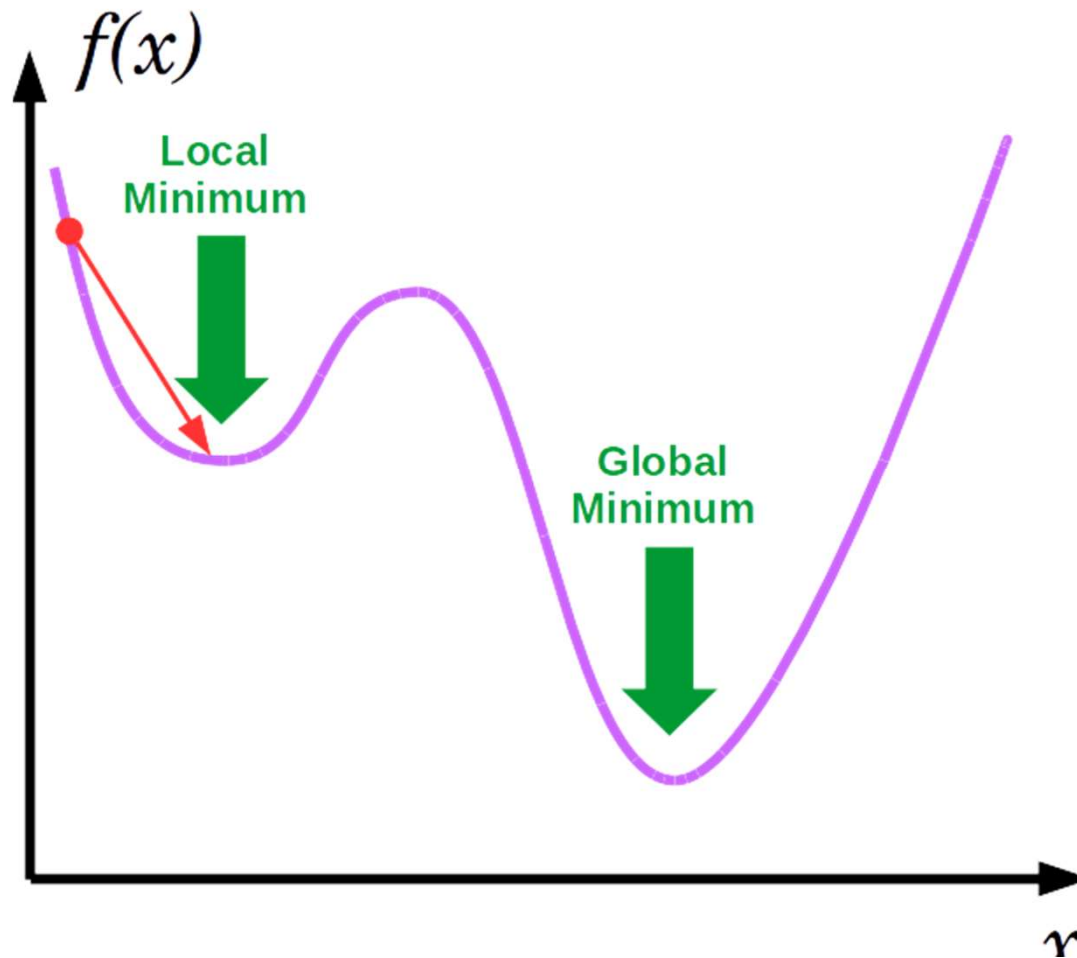
Big learning rate



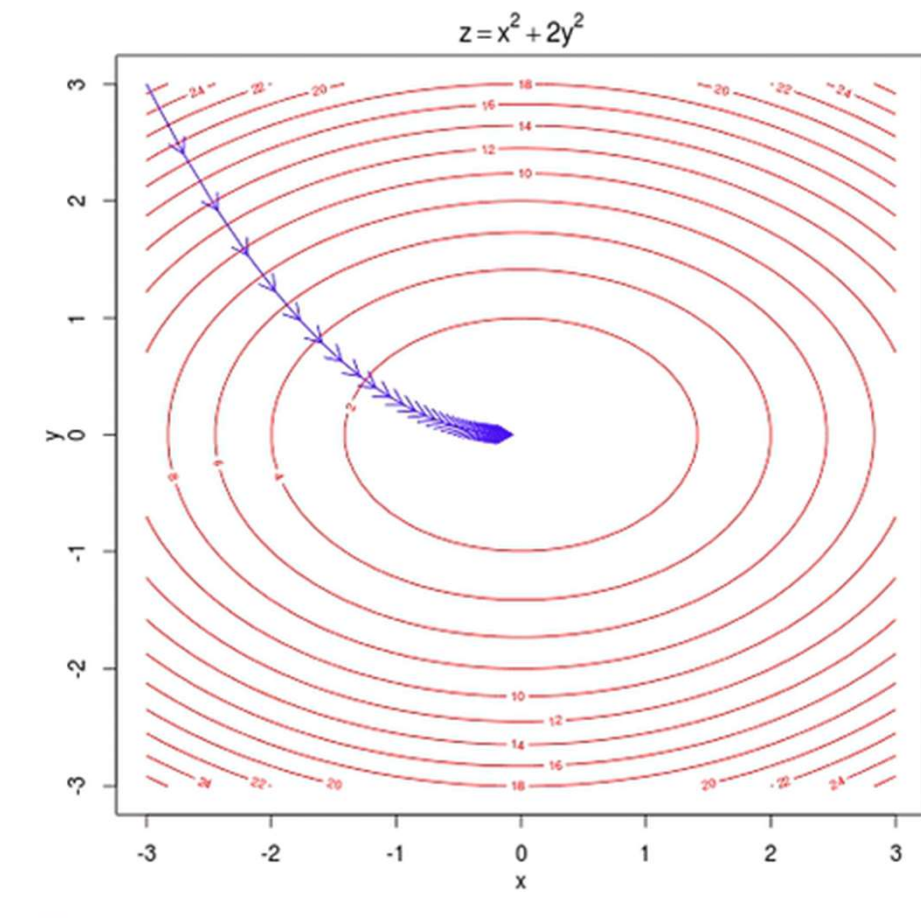
Small learning rate



Effect of Local Minima



Multidimensional Optimization



Generalized Delta Rule

Why activation function then?

- If neural networks had no activation functions, they would fail to learn the complex non-linear patterns that exist in real-world data. Activation functions enable neural networks to learn these non-linear relationships by introducing non-linear behaviors through activation functions.
- Checkout <https://www.datacamp.com/tutorial/introduction-to-activation-functions-in-neural-networks>

$$w_{ij} \leftarrow w_{ij} + \alpha \delta_i x_j \quad (\text{Equation 2.3})$$

$$\delta_i = \varphi'(v_i) e_i \quad (\text{Equation 2.4})$$

where

e_i = The error of the output node i

v_i = The weighted sum of the output node i

φ' = The derivative of the activation function φ of the output node i

e from the cost function.

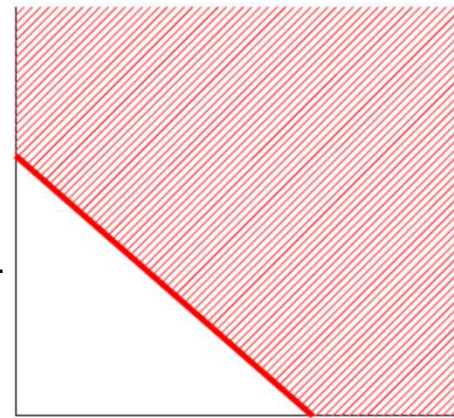
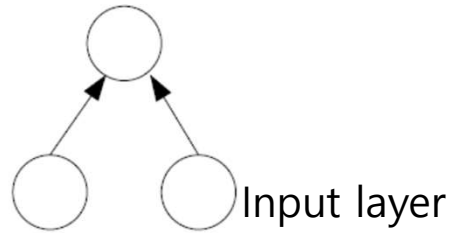
If activation fun. is x, then its derivative is 1 (check Equation 2.2)

Single vs. Multi Layer Network: XOR Problem

Multi Layer Network: XOR Problem

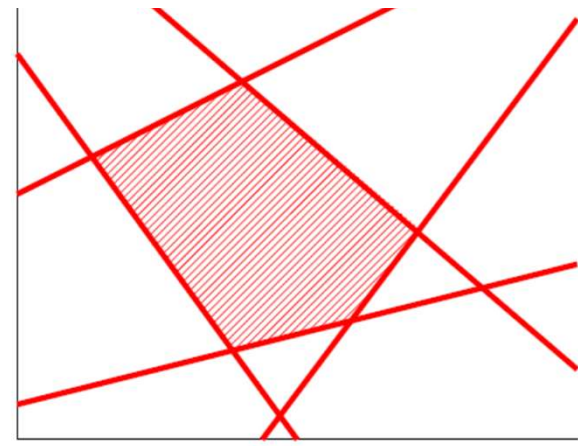
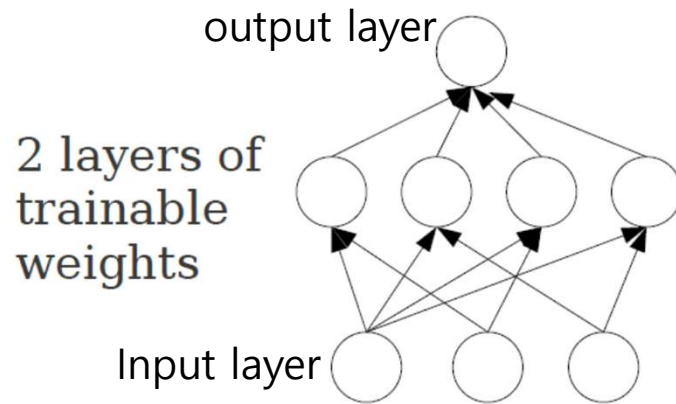
Single Layer NN and Decision Boundary

1 layer of trainable weights



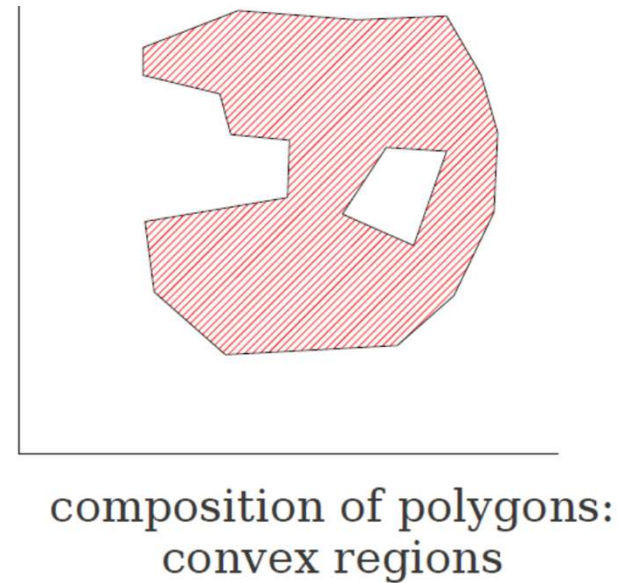
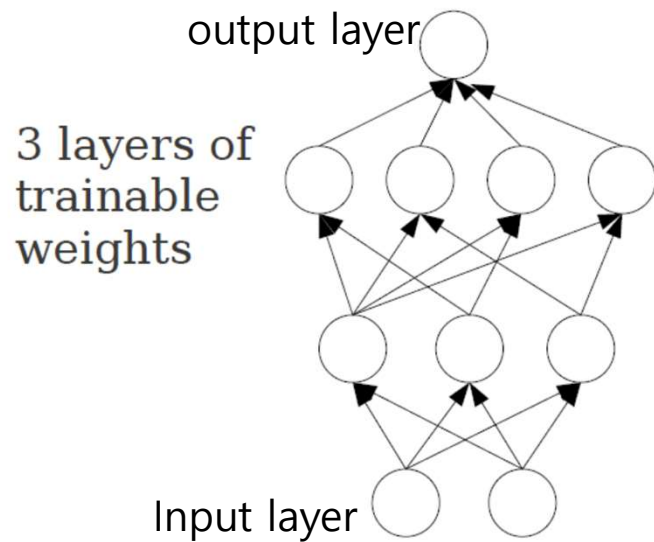
separating hyperplane

2 Layer NN and Decision Boundary



convex polygon region

3 Layer NN and Decision Boundary



Solving XOR with 2 Layer NN

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

Two solutions:

$$x_1 \bar{x}_2 \vee \bar{x}_1 x_2$$
$$(x_1 \vee x_2) \wedge \overline{x_1 \wedge x_2}$$

