

# *Deep Learning - Theory and Practice*

---

Linear Regression, Least Squares  
Classification and Logistic Regression

20-02-2020

<http://leap.ee.iisc.ac.in/sriram/teaching/DL20/>

[deeplearning.cce2020@gmail.com](mailto:deeplearning.cce2020@gmail.com)





# Least Squares for Classification

- ❖ K-class classification problem

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}$$

- ❖ With 1-of-K hot encoding, and least squares regression

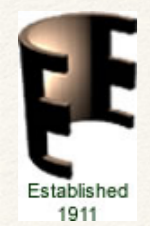
$$\mathbf{t} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\sum_{n=1}^N$$

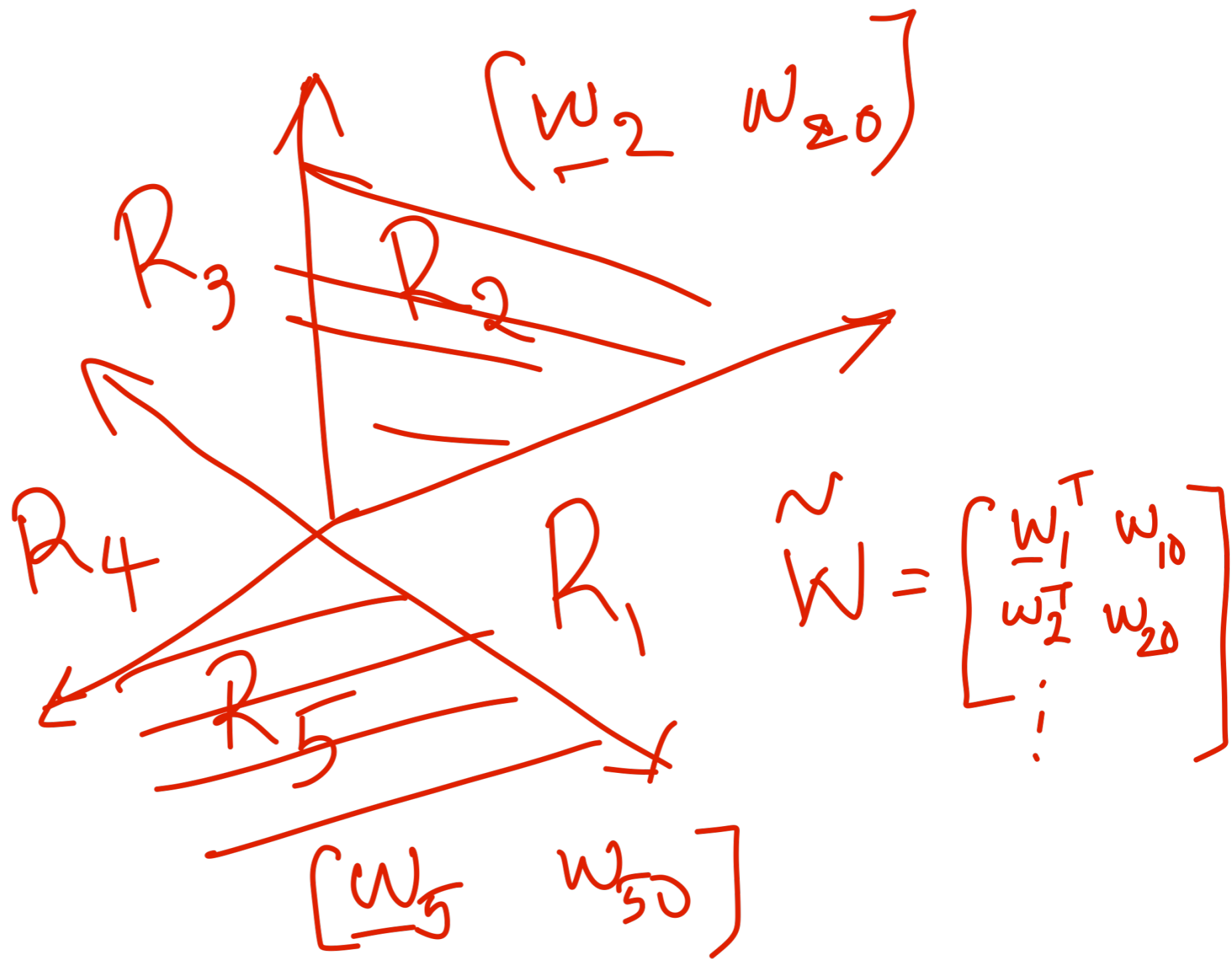
$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\}$$

$N \times N$

Bishop - PRML book (Chap 3)



$$T_{91} \left( \begin{array}{c} \|y(x_1) - t_1\|^2 \\ \|y(x_2) - t_2\|^2 \\ \vdots \\ \|y(x_N) - t_N\|^2 \end{array} \right)_{N \times N}$$





# Logistic Regression (classification)

- ❖ 2- class logistic regression

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

- ❖ Maximum likelihood solution

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

- ❖ K-class logistic regression

$$p(C_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

- ❖ Maximum likelihood solution

$$a_k = \mathbf{w}_k^T \phi.$$

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

# Issues

\* Outputs are real numbers

\* Model targets are binary.

\* Cost Function - least squares

Next Model. → make outputs  
close to prob.

$$P(C_1/\phi(x)) = \frac{\text{Binary Classification } (C_1, C_2)}{P(\phi(x)/C_1) P(C_1)}$$

$$\phi(x) \quad P(\phi(x)) = P(\phi(x), C_1) + P(\phi(x), C_2)$$

↑ fixed non-linear transformation

$$P(C_1/\phi(x)) = \frac{1}{1 + e^{-f(x)}} ;$$

$$1 + \frac{P(\phi(x), C_2)}{P(\phi(x), C_1)}$$

$$\underline{\underline{f(x) = -\ln \frac{P(\phi(x), C_2)}{P(\phi(x), C_1)}}}$$



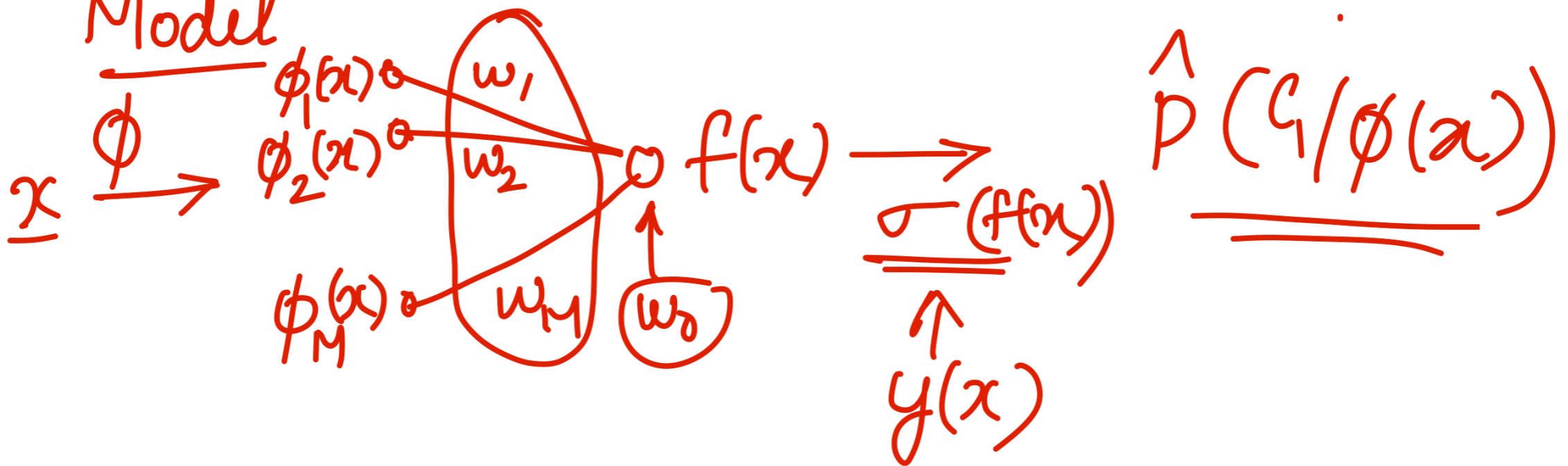
$$p(C_1 | \phi(x)) = \frac{1}{1 + e^{-f(x)}} = \sigma(f(x))$$

↑  
logistic  
function

Approx.

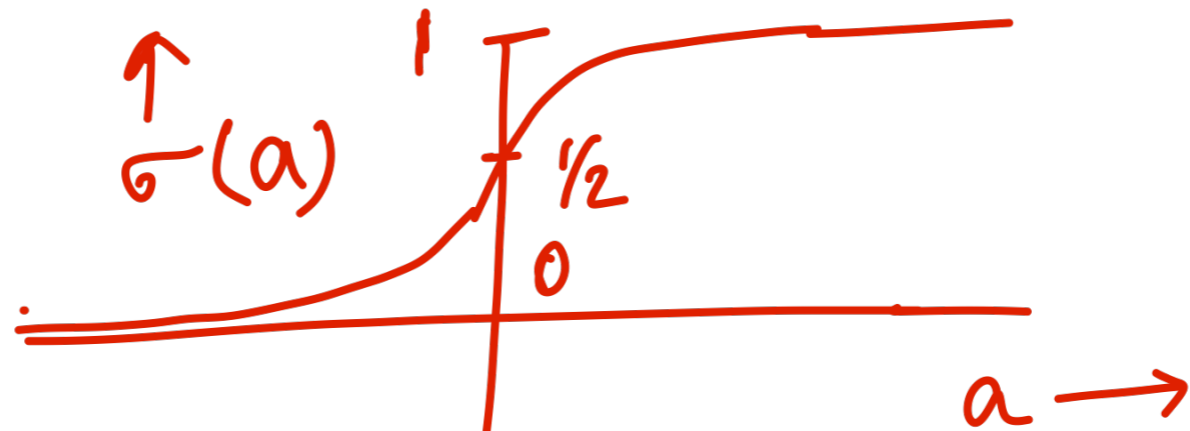
$$f(x) = \underline{\underline{w}}^T \phi(x) + w_0 = \underline{\underline{\tilde{w}}} = [\underline{w}^T \ w_0]^T$$

Model





$\sigma(a)$



$$\begin{aligned} p(C_2 | \phi(x)) &= 1 - \sigma(f(x)) \\ &= \sigma(-f(x)) \end{aligned}$$

$$t_n = \begin{cases} 0 & \text{if } x_n \in C_2 \\ 1 & \text{if } x_n \in C_1 \end{cases}$$

$$y(x) = \underline{\sigma(f(x))} = p(C_1 | \phi(x)) \quad x_1, x_2, \dots, x_N.$$

$$\begin{aligned} \text{Likelihood} &= \prod_{n=1}^N (y(x_n))^{t_n} (1 - y(x_n))^{1-t_n} \\ \uparrow L(x) & \quad \quad \quad \text{w.r.t } \underline{w}, w_0 \end{aligned}$$

$$\underline{\log L(x)} = \sum_{n=1}^N t_n \log y_n + (1-t_n) \log (1-y_n)$$

B.C.E (binary cross entropy) =  $-\log L(x)$   
 between  $\begin{pmatrix} y(x_n) \\ 1-y(x_n) \end{pmatrix}, \begin{pmatrix} t_n \\ 1-t_n \end{pmatrix}$   
 $y_n = y(x_n)$

$$y_n = y(x_n) = \frac{1}{1 + e^{-\{w^T \phi(x_n) + w_0\}}}$$

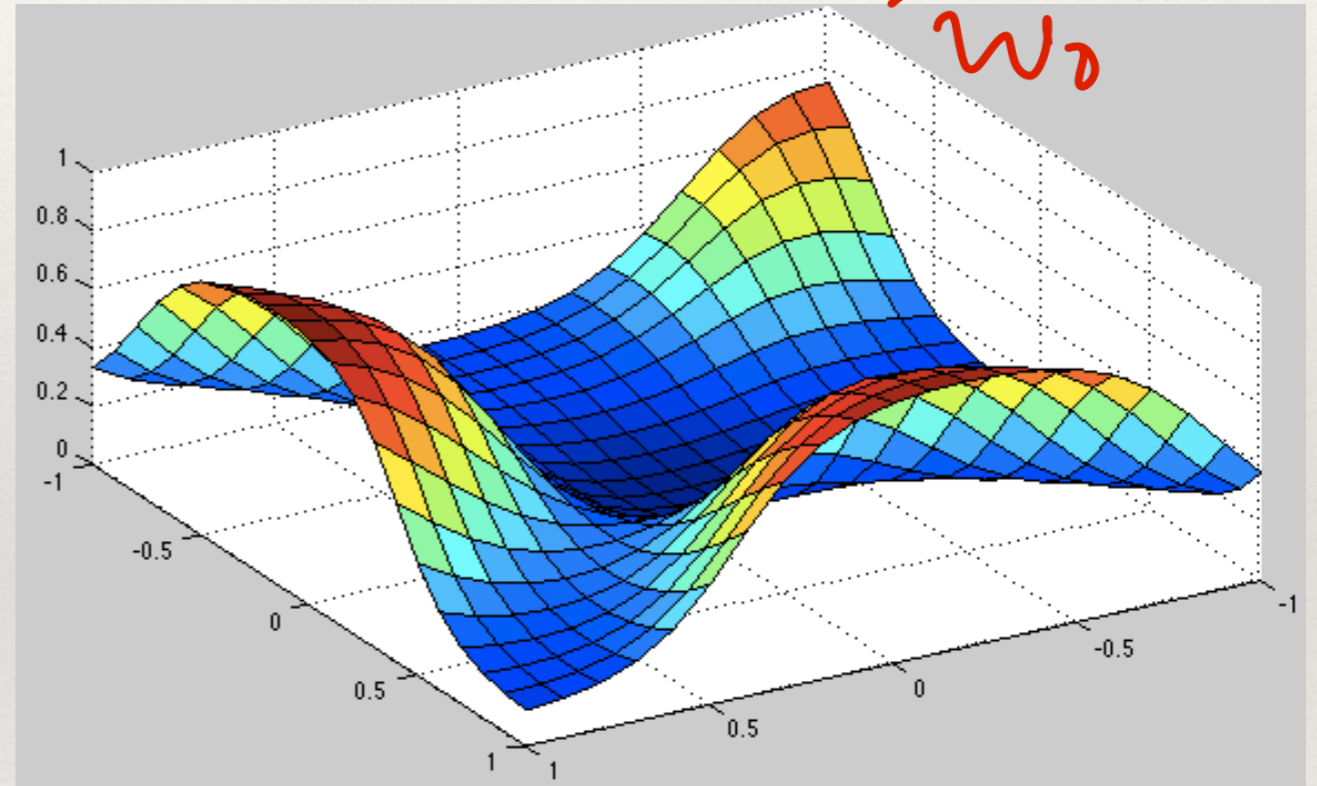
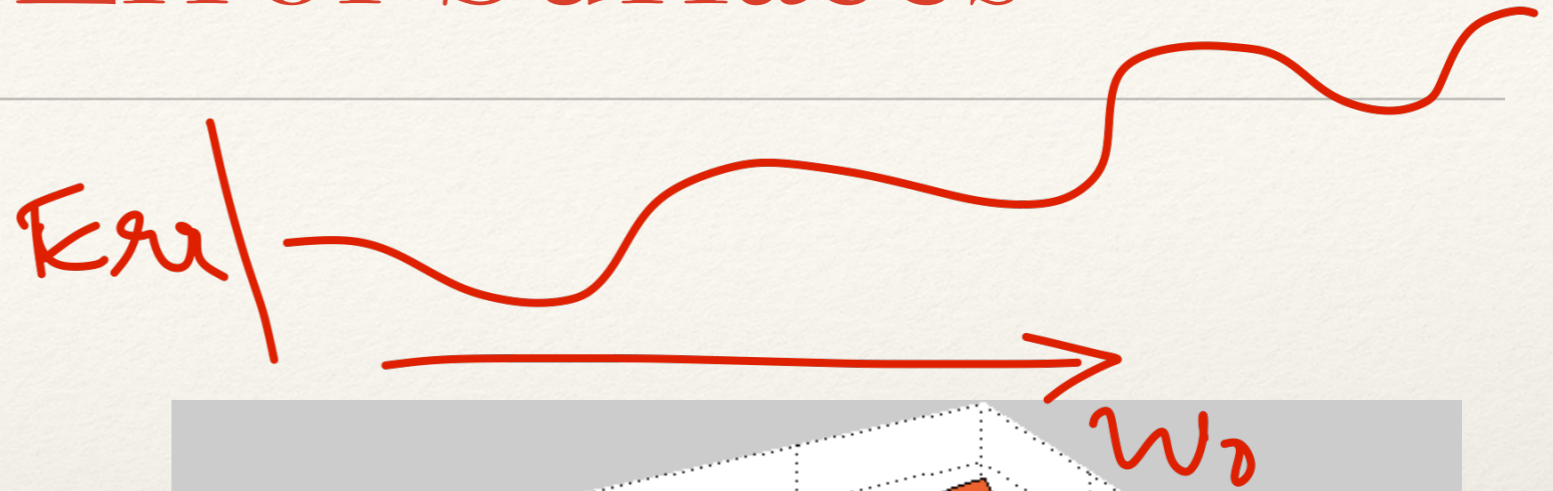
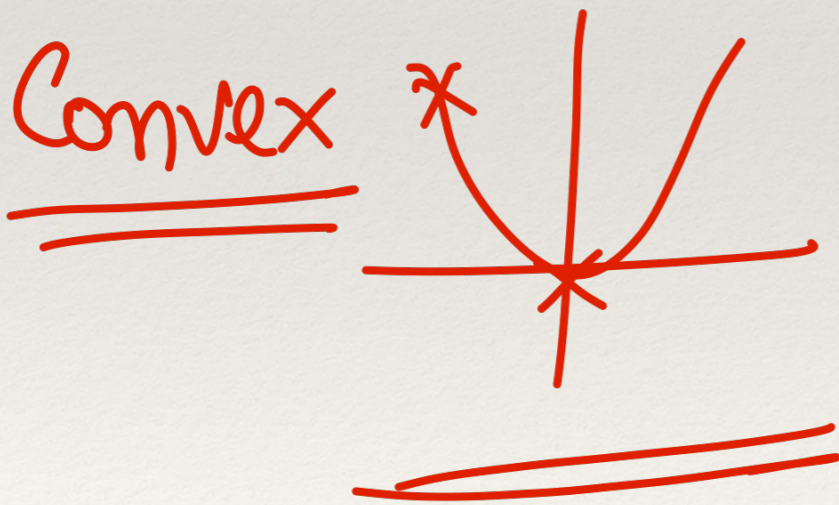
Minimize B.C.E  
 w.r.t  $w, w_0$

$$x = \{x_1 \dots x_N\}$$

# Typical Error Surfaces

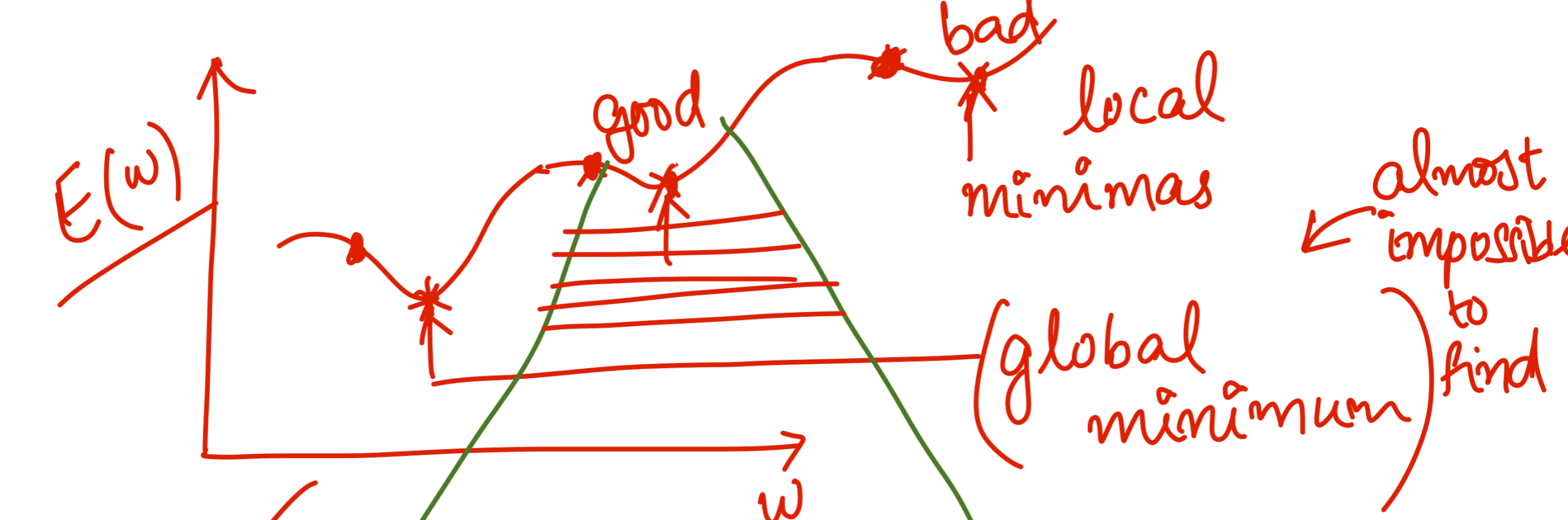
B.C.E  
↑

Typical Error Surface as a function of parameters (weights and biases)

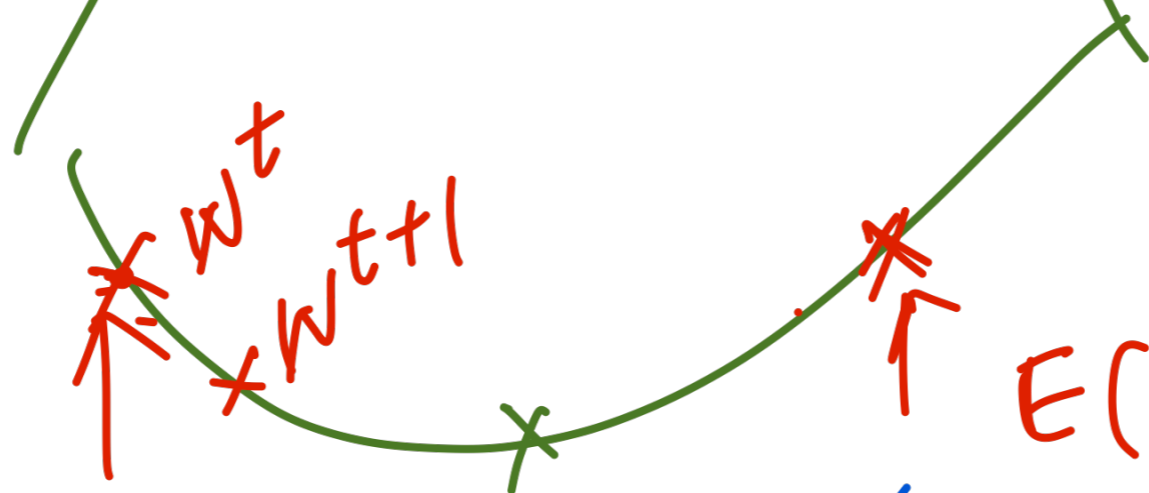


Non-convex  
Multiple max and min





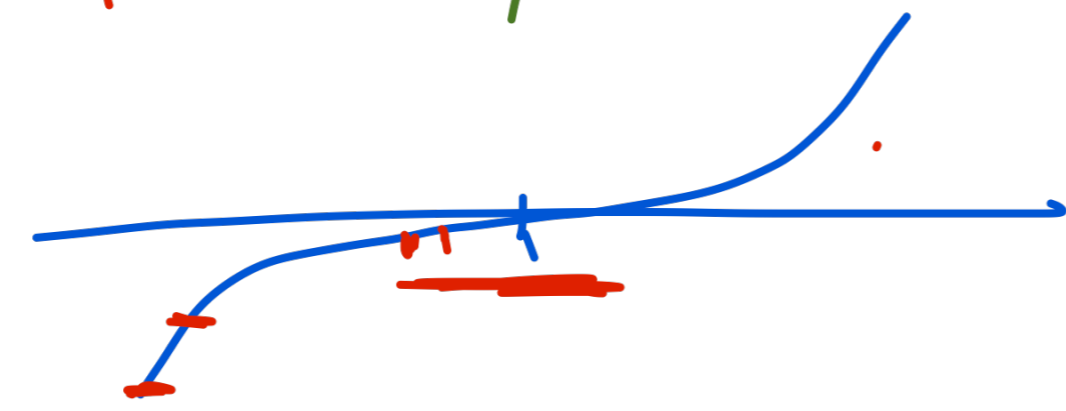
zero saddle point



Iterative

$$\underline{\underline{E(w^{t+1}) \leq E(w^t)}}$$

$\frac{\partial E}{\partial w}$



Direction (Error surfaces are smooth)

Derivative  $\frac{\partial E}{\partial w}$  (-ve direction of derivative)

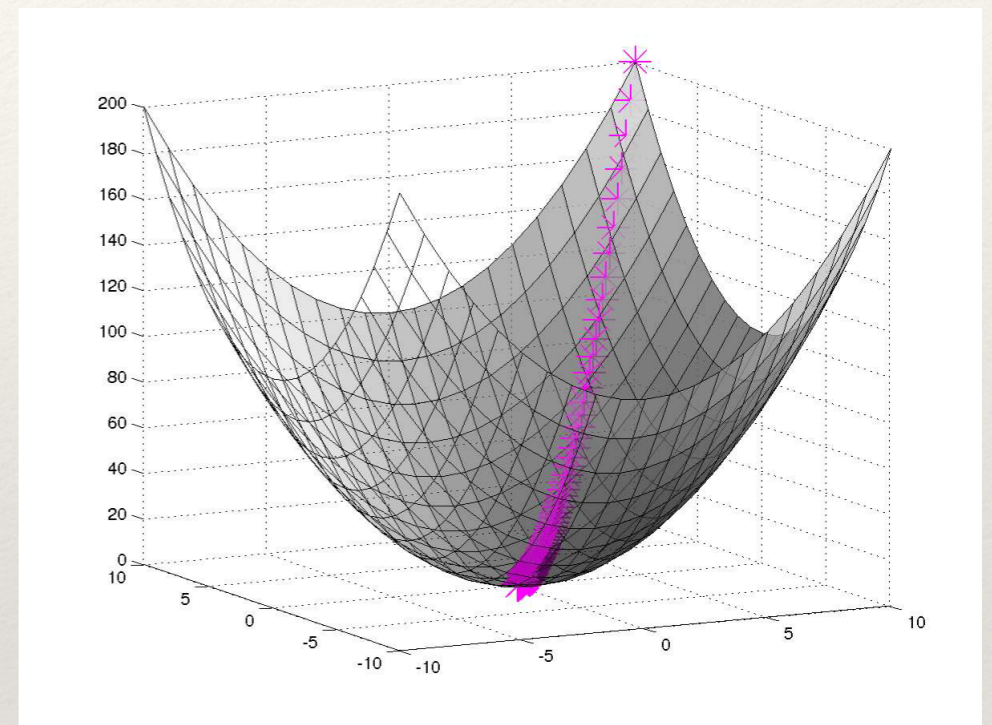
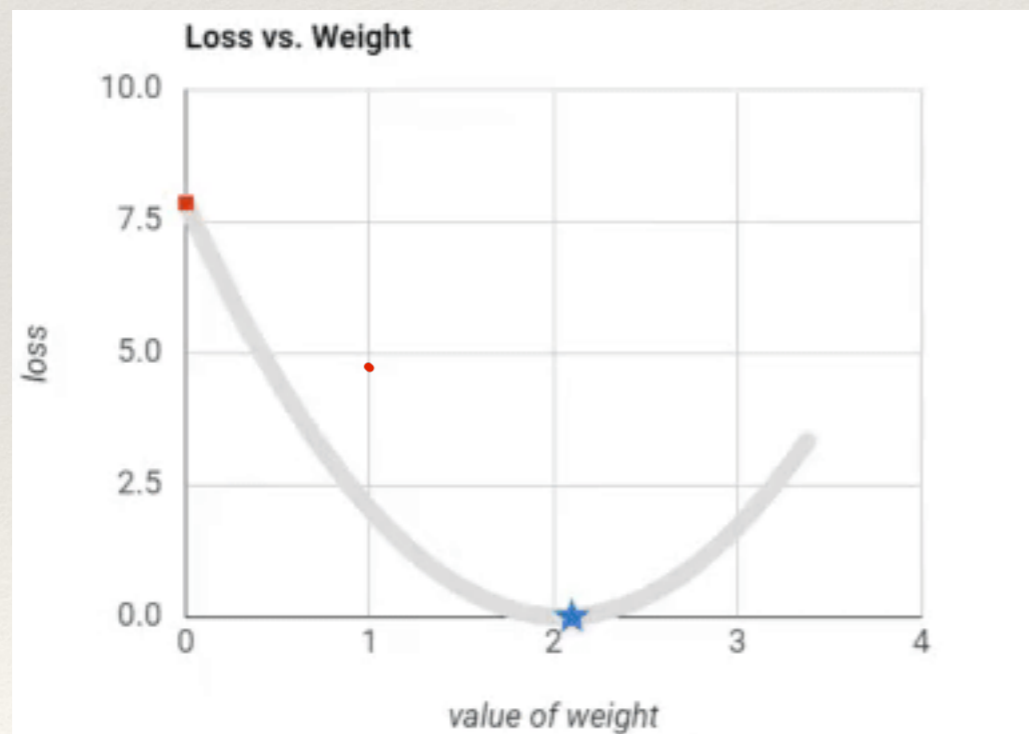
Magnitude -  $\eta$  (learning rate) +ve  
small

$$w^{t+1} = w^t - \eta \left. \frac{\partial E}{\partial w} \right|_{w=w^t}$$

Gradient Descent Algorithm

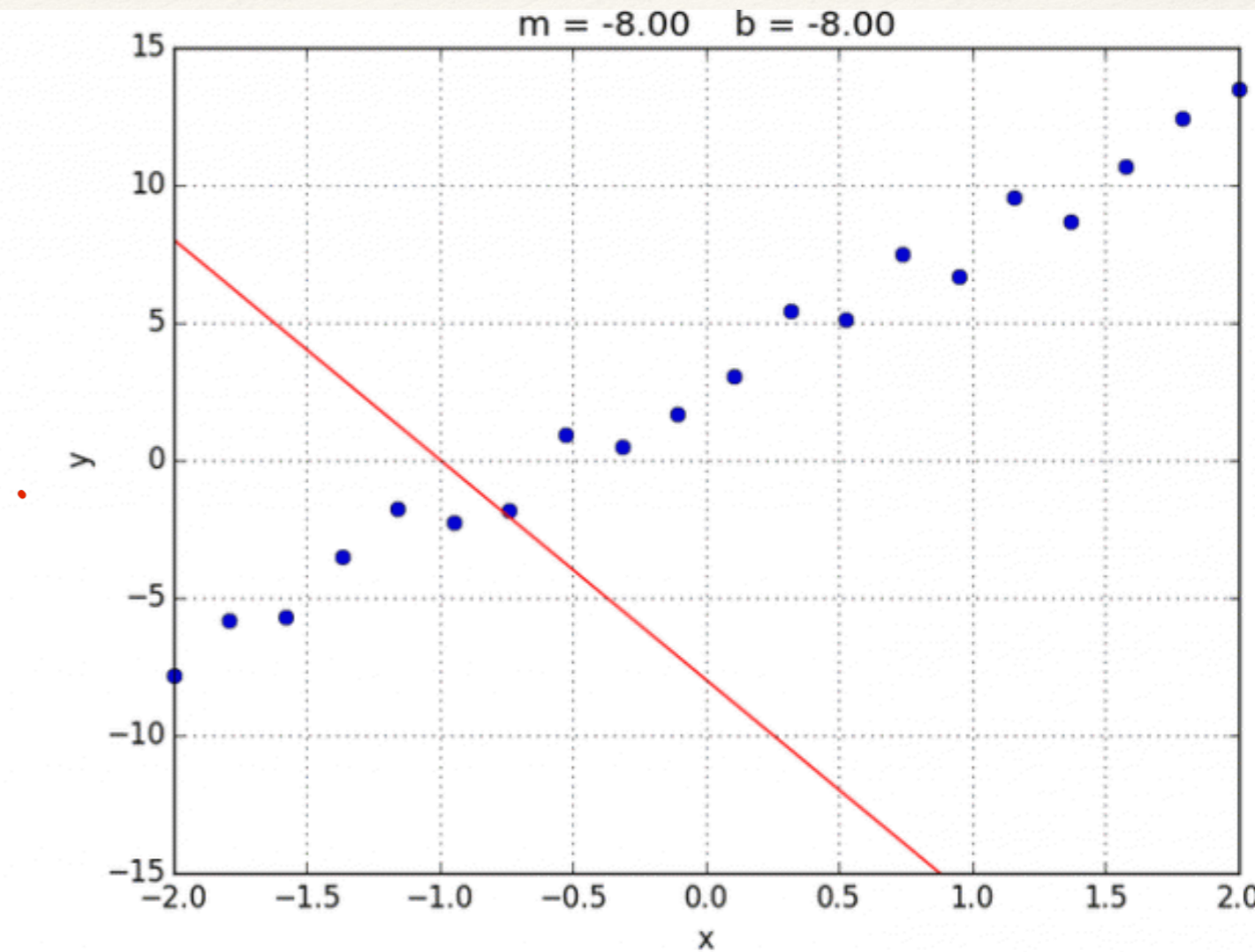
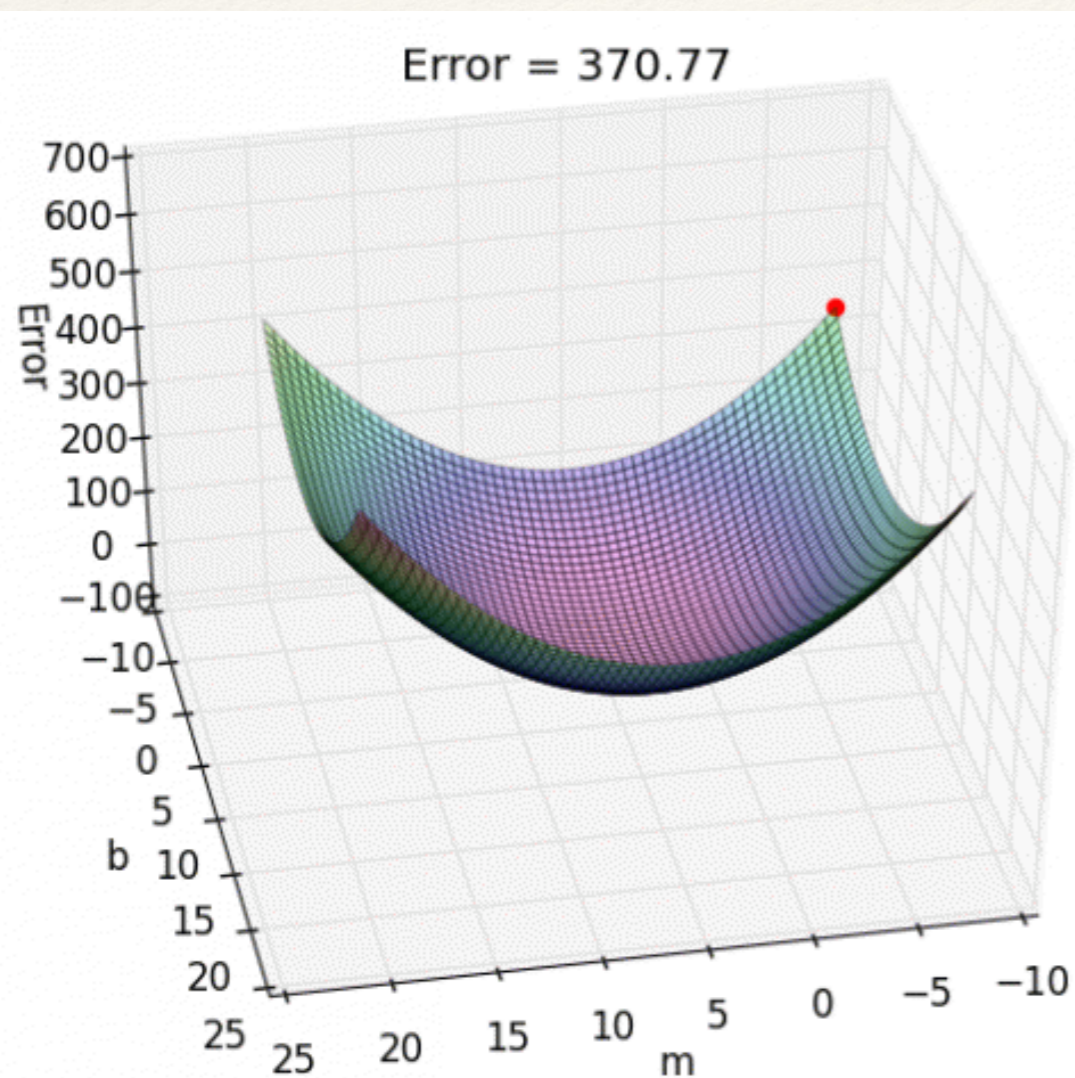
# Learning with Gradient Descent

Error surface close to a local



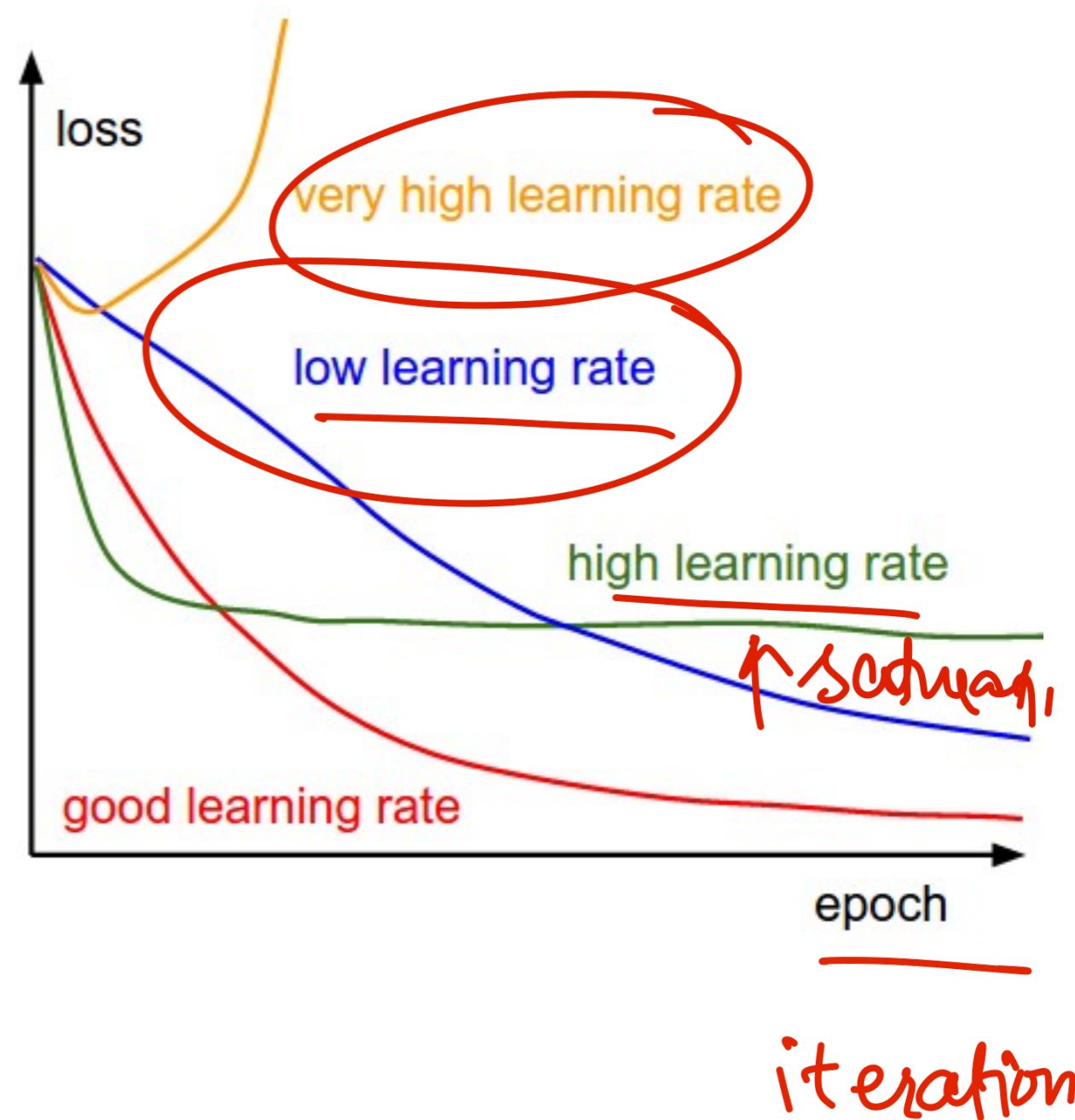


# Learning Using Gradient Descent



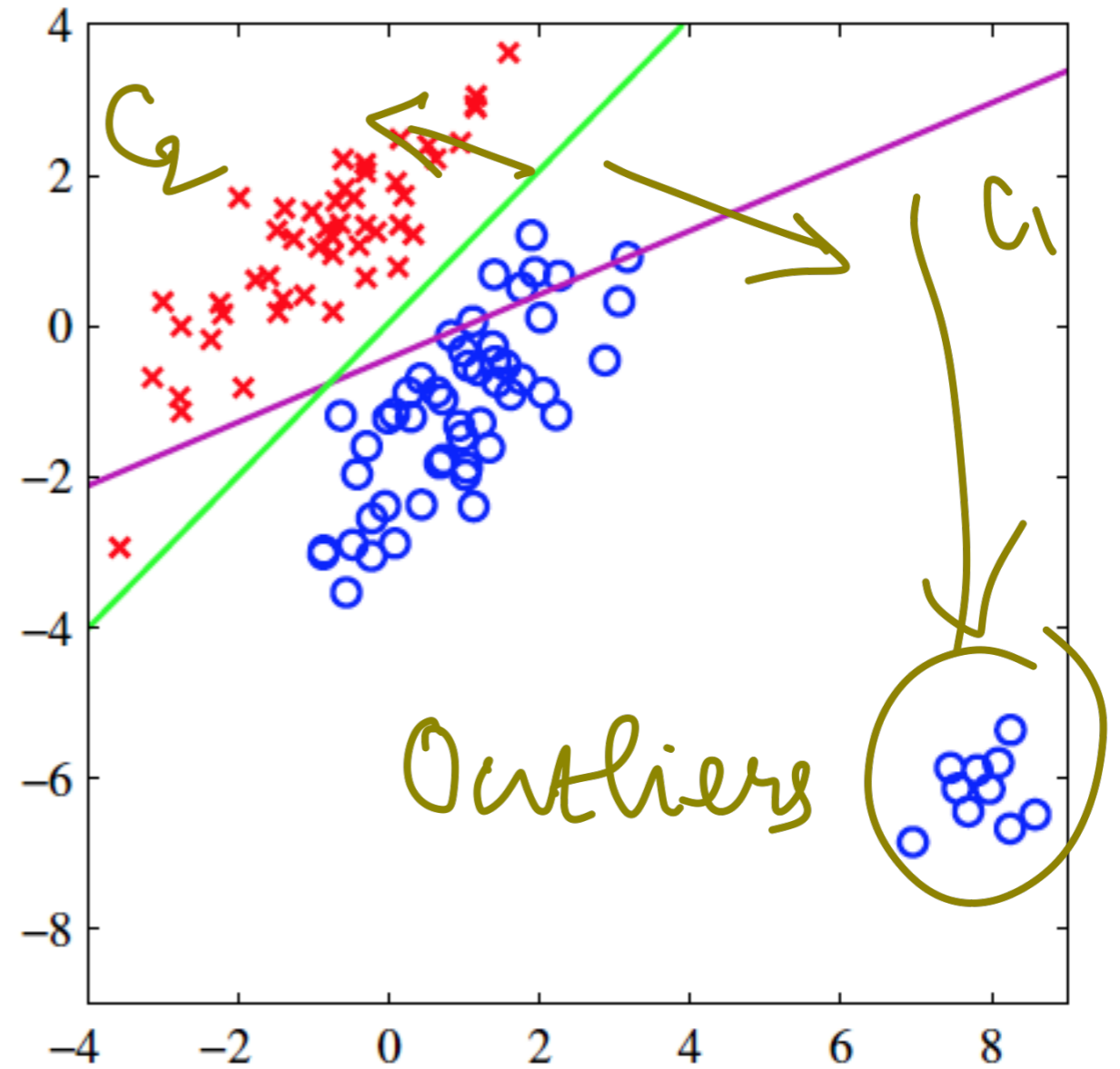
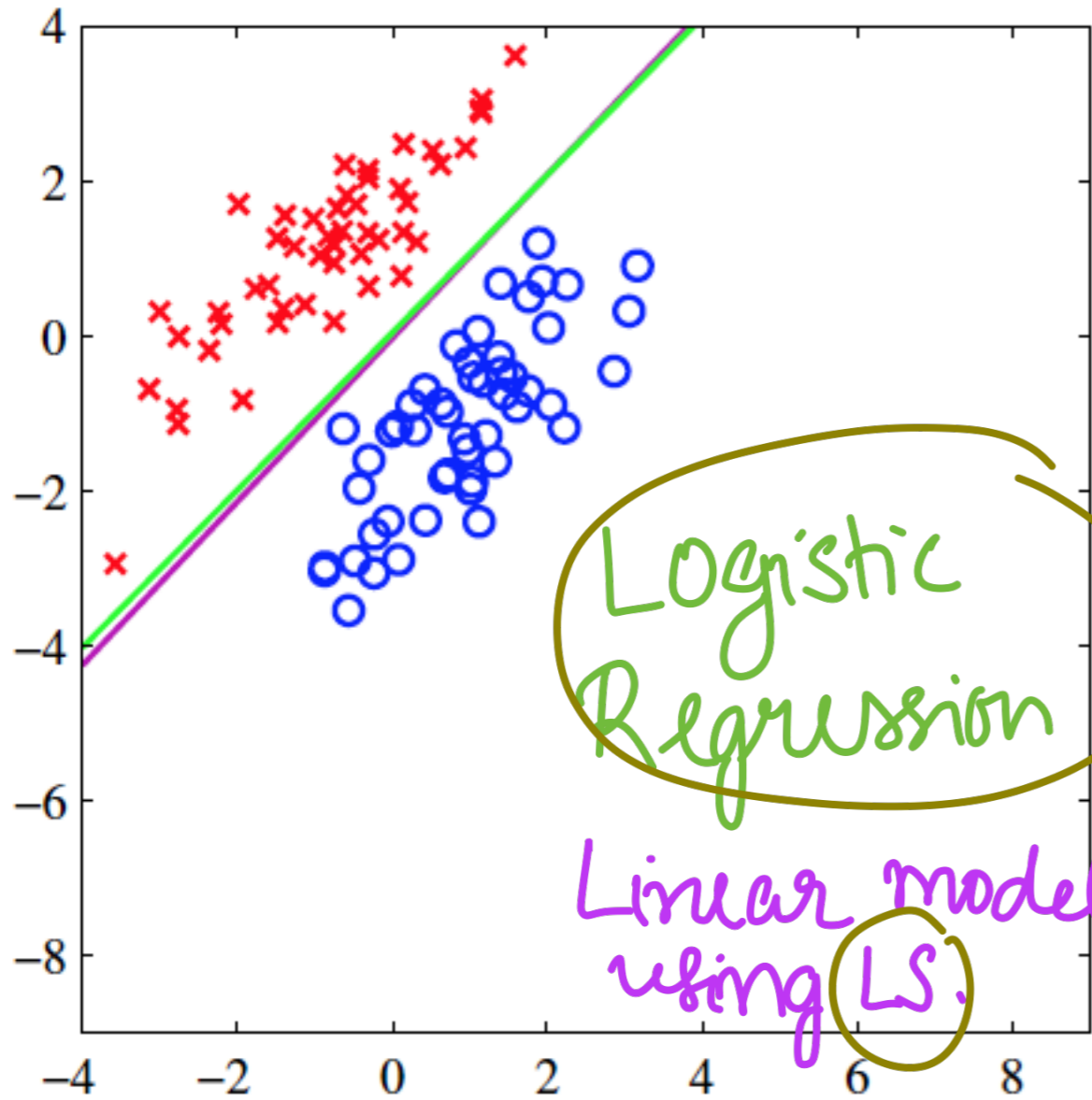
# Parameter Learning

- Solving a non-convex optimization.
- Iterative solution.
- Depends on the initialization.
- Convergence to a local optima.
- Judicious choice of learning rate



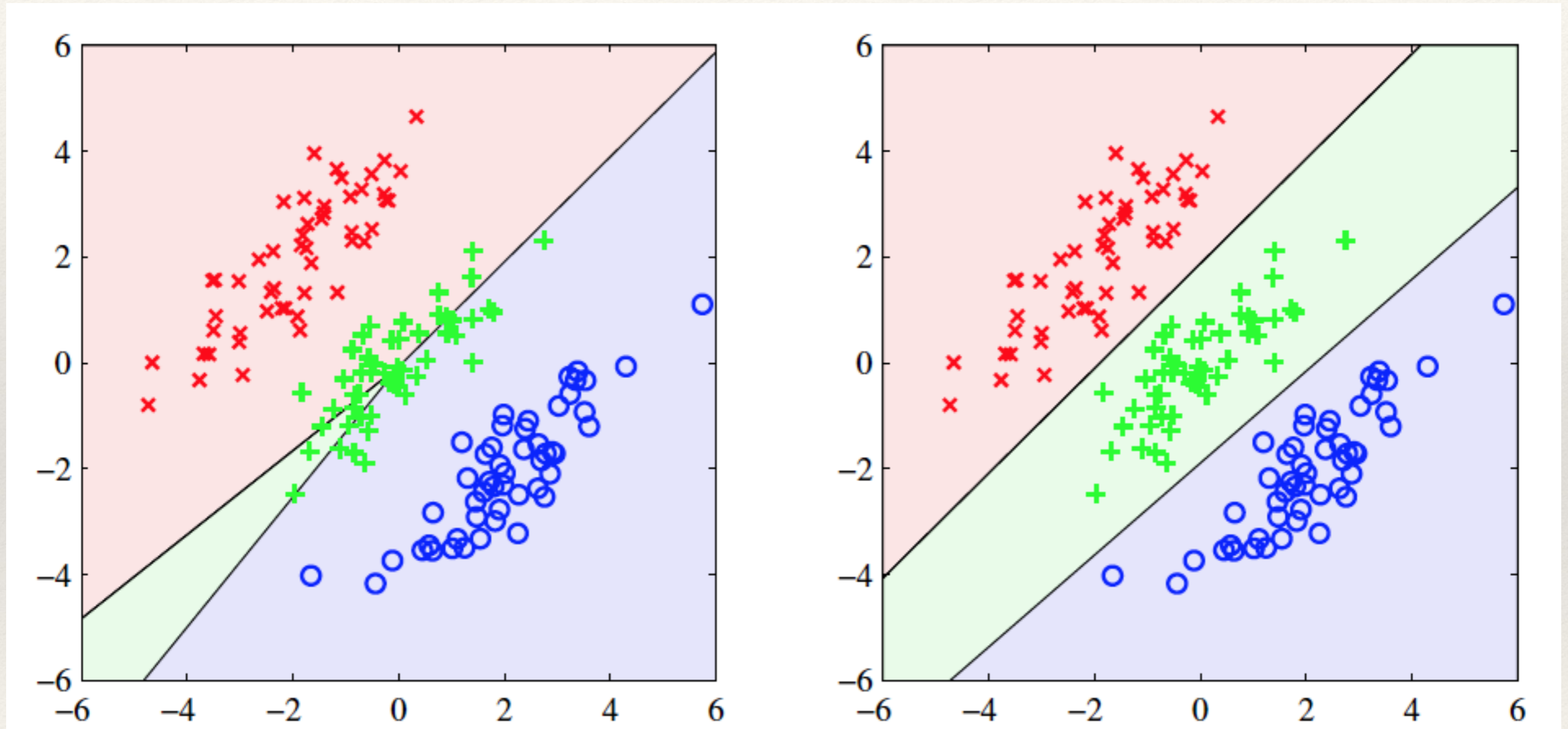


# Least Squares versus Logistic Regression





# Least Squares versus Logistic Regression





---

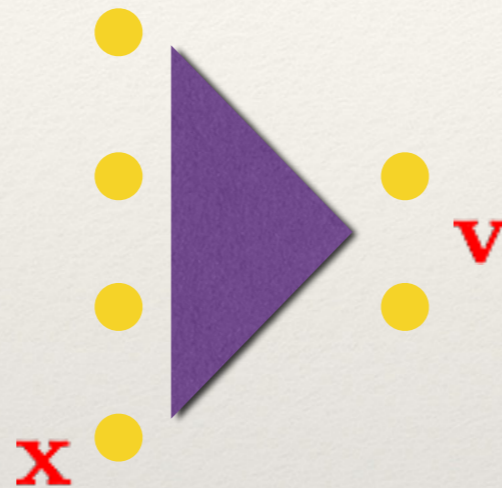
# Neural Networks

---



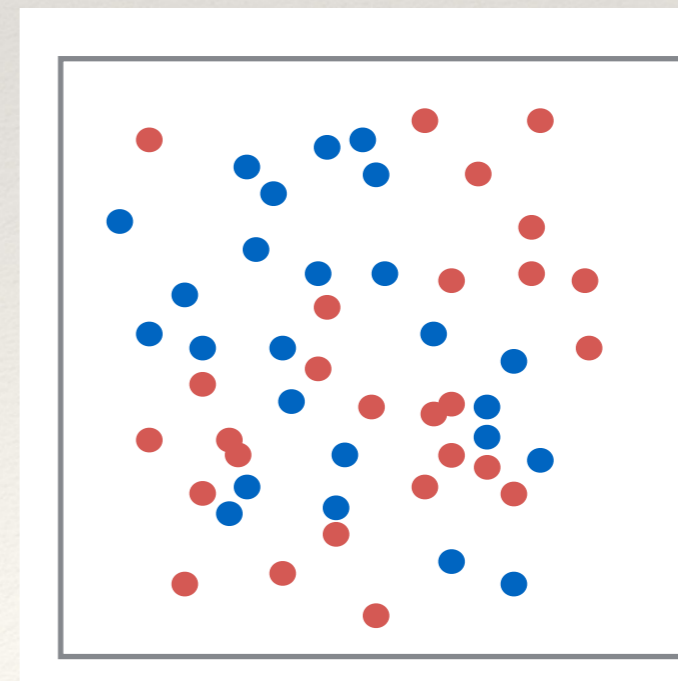
# Perceptron Algorithm

Perceptron Model [McCulloch, 1943, Rosenblatt, 1957]



Targets are binary classes  $[-1, 1]$

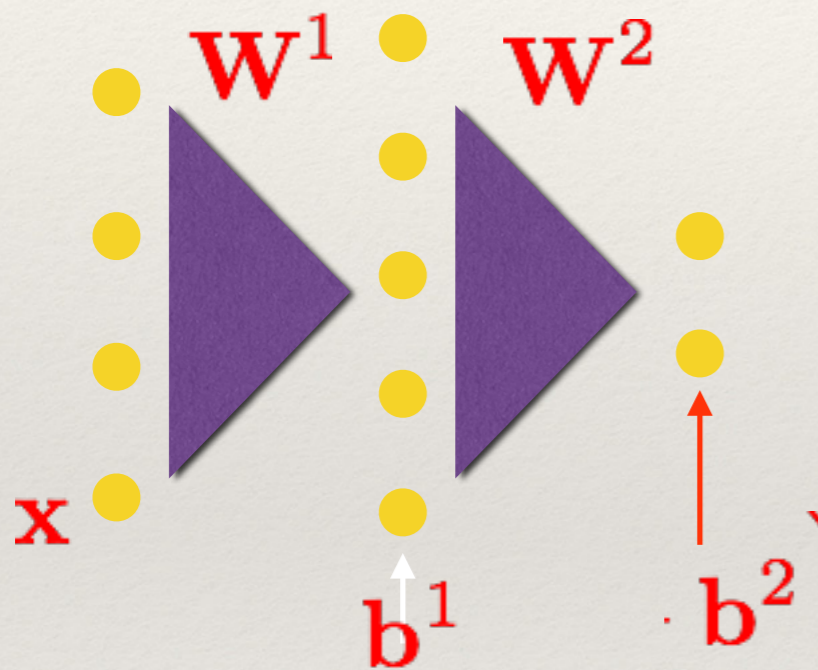
What if the data is not  
linearly separable





# Multi-layer Perceptron

Multi-layer Perceptron [Hopfield, 1982]



$$\mathbf{v}^2 = \psi \left( \mathbf{W}^2 \phi(\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2 \right)$$

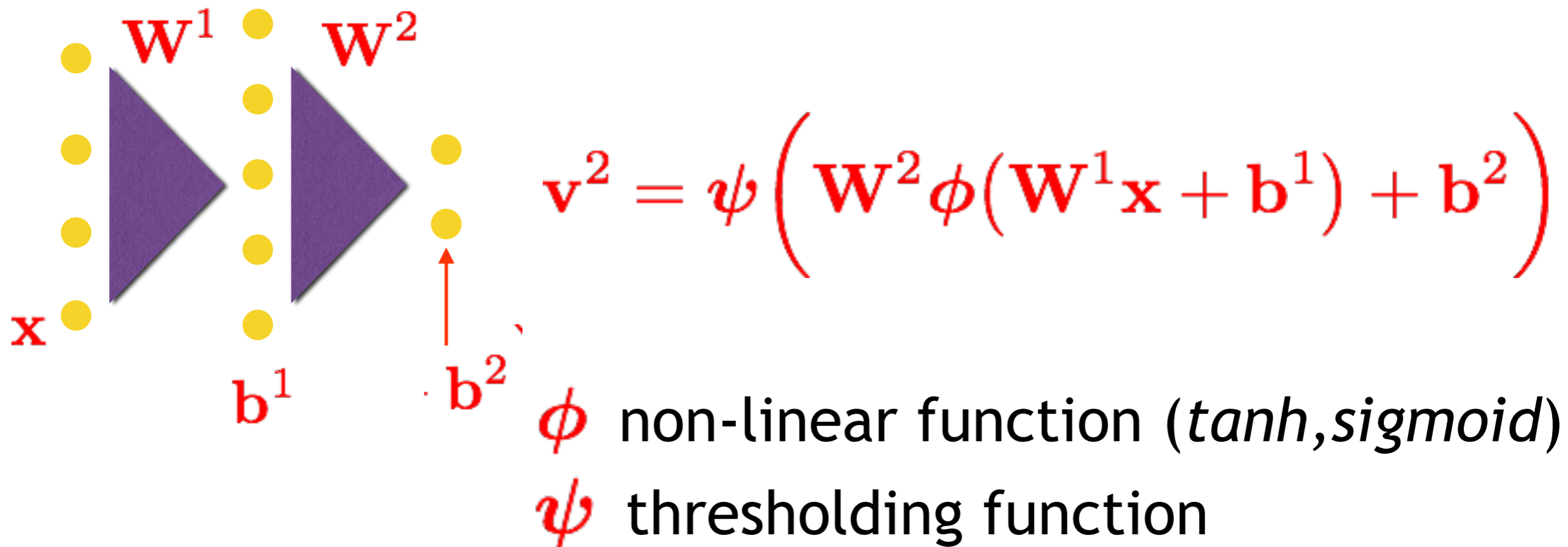
$\phi$  non-linear function (*tanh, sigmoid*)

$\psi$  thresholding function



# Neural Networks

## Multi-layer Perceptron [Hopfield, 1982]

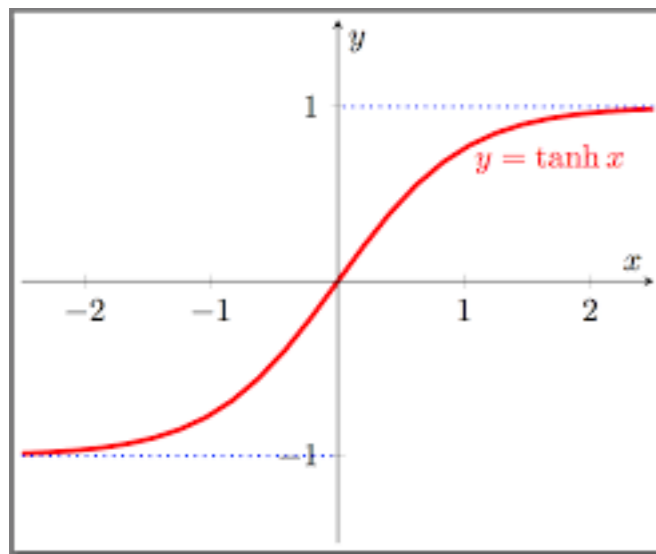


- Useful for classifying **non-linear data boundaries** - non-linear class separation can be realized given enough data.

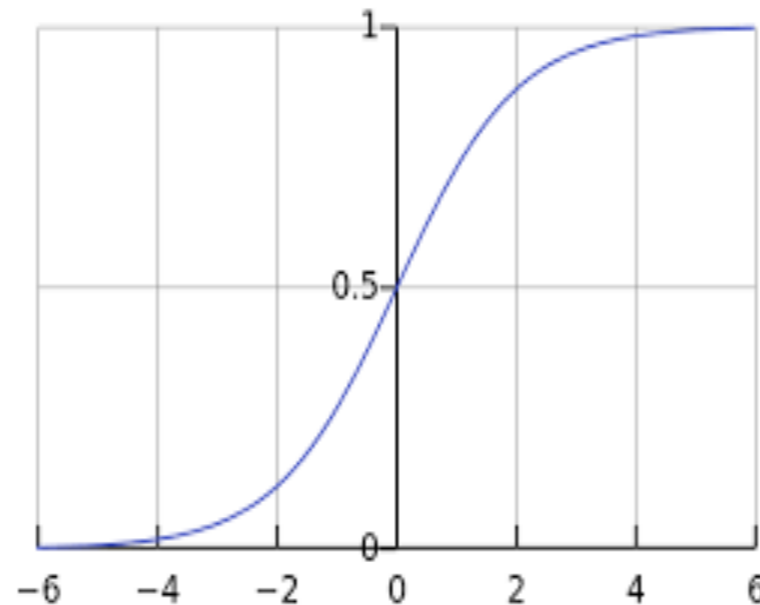
# Neural Networks

## Types of Non-linearities $\phi$

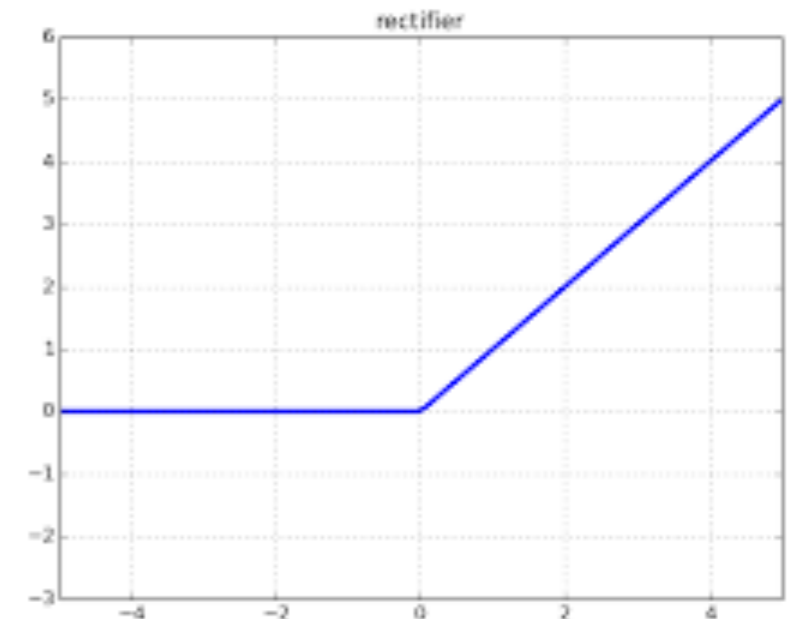
tanh



sigmoid



ReLu



## Cost-Function

Mean Square Error

$$J_{MSE} = \sum_{i=1}^M ||\mathbf{v}_i - \mathbf{y}_i||^2$$

Cross Entropy

$$J_{CE} = - \sum_{i=1}^M \mathbf{y}_i^T \log(\mathbf{v}_i)$$

$\mathbf{y}_i$  are the desired outputs



# Learning Posterior Probabilities with NNs

## Choice of target function $\psi$

- Softmax function for classification

$$\psi(v_i) = \frac{e^{v_i}}{\sum_i e^{v_i}}$$

- Softmax produces positive values that sum to 1
- Allows the interpretation of outputs as posterior probabilities