

ECS 171 Machine Learning

Lecture4: Logistic Regression, Newton's method, Perceptron Model

Instructor: Dr. Setareh Rafatirad

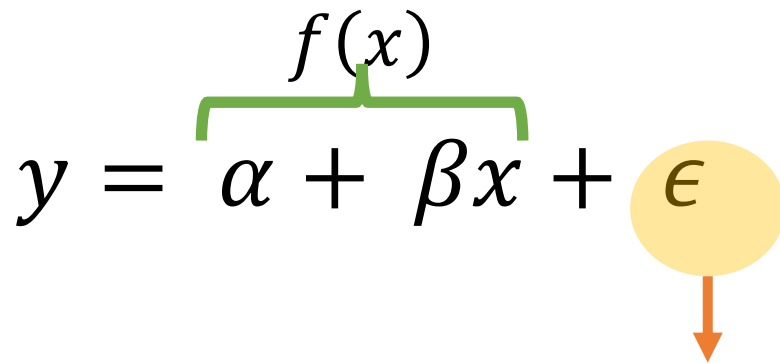
MLE for LR Recap

Sample : X_1, X_2, \dots, X_n

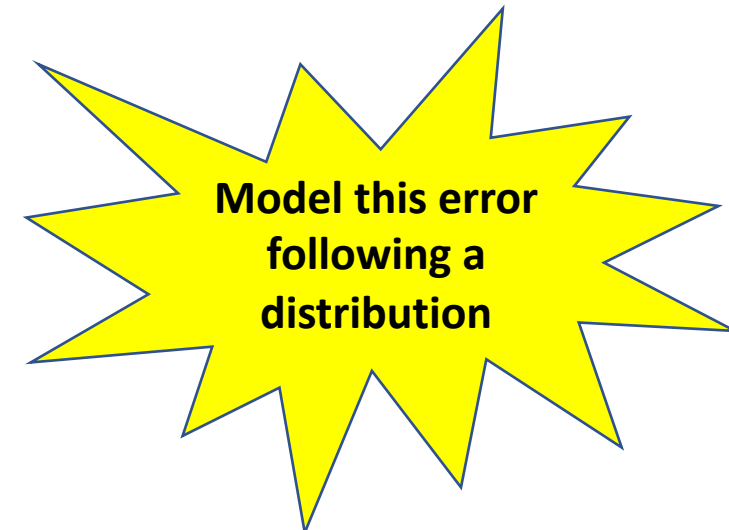
y : output (response variable)

$$y = \alpha + \beta x$$

when residual error $\epsilon = 0$.

$$y = \overbrace{\alpha + \beta x}^{f(x)} + \epsilon$$


Random noise



$$\epsilon \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

MLE for LR Recap

$$y = f(x) + \epsilon$$

$$\epsilon \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\downarrow$$
$$\alpha + \beta x$$

\downarrow
variance

Sample : X_1, X_2, \dots, X_n

For a fixed sample : X_i , the distribution of Y_i is equal to:

$$\mathcal{N}(f(X_i), \sigma^2)$$

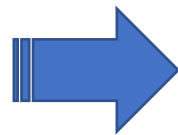
$$f(X) = w^T X \rightarrow \mathcal{N}(w^T X, \sigma^2)$$

\downarrow
mean

\downarrow
variance

p.d.f

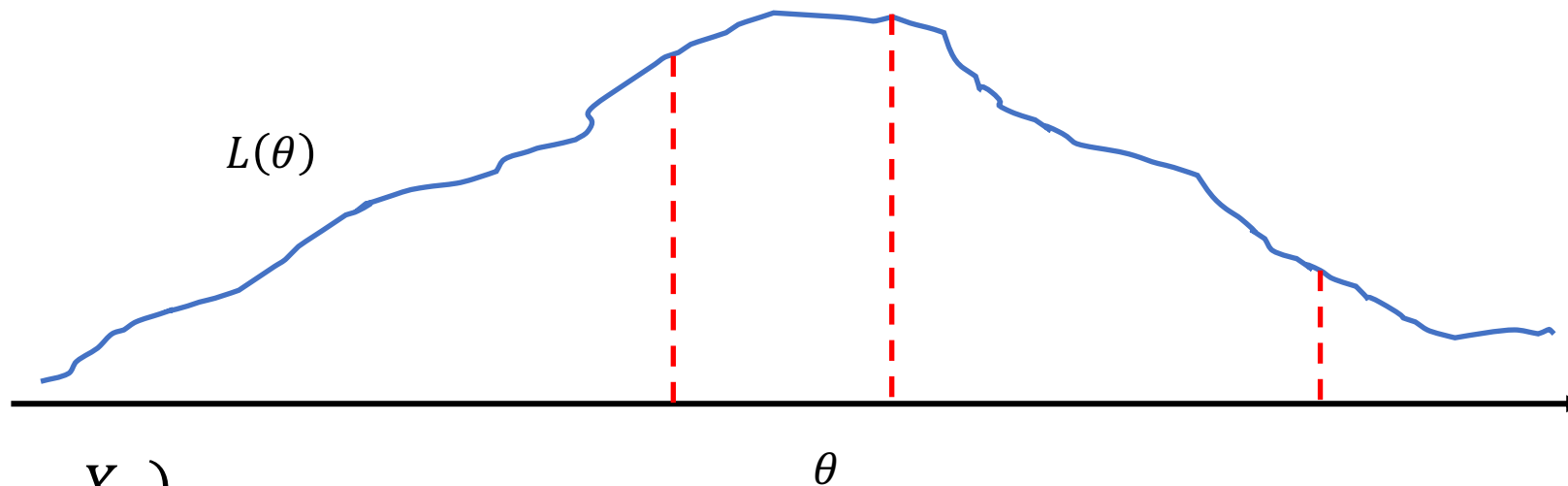
$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$



$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-f(x_i))^2}{2\sigma^2}}$$

Likelihood Function Recap

$L(\theta)$, where θ is the unknown parameter



$$\begin{aligned} L(\theta|X_1, \dots, X_n) &= f(X_1, \dots, X_n | \theta) \\ &= f(X_i | \theta) \\ &= f(X_1 | \theta) \dots f(X_n | \theta) \\ &= \prod_{i=1}^n f(X_i | \theta) \end{aligned}$$

$$\begin{aligned} \log(a \cdot b) &= \log a + \log b \\ \downarrow \\ L(\theta|X_1, \dots, X_n) &\cong -\log L(\theta|X_1, \dots, X_n) \end{aligned}$$

Negative Log Likelihood (NLL)

MLE Recap

$$L(\theta) \triangleq -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (Y_i - w^T X_i)^2$$

$$RSS(w) \triangleq \sum_{i=1}^N (Y_i - w^T X_i)^2 = ||\epsilon||^2$$

sum of squared errors (SSE) squared norm of residual errors.

$$NLL(\theta) \triangleq -\log f(X_i | \theta) = -\sum_{i=1}^N \log p(Y_i | X_i, \theta)$$

Logistic Regression

Logistic Regression or Linear Regression?

- Classification problem
- Example
 - Spam detection (1 or 0)
 - Tumor detection (1 or 0)

if the hypothesis is a linear regression model:

$$Z = w^T x$$

$$y^{(i)} = \begin{cases} 0 & ; \text{predicted value} < \text{threshold} \\ 1 & ; \text{predicted value} \geq \text{threshold} \end{cases}$$

Drawbacks of using Linear Regression for classification (predicting labels):

- **Sensitivity to outliers**
- **Sensitivity to the selected threshold**



Inputs: X_1, X_2, X_3 || Weights: $\theta_1, \theta_2, \theta_3$ || Outputs: **malignant or benign**

@dataaspirant.com

Figure 1: Logistic Regression Model

(Source: <http://dataaspirant.com/2017/03/02/how-logistic-regression-model-works/>)

Logistic Regression is a better way to perform classification within the regression framework.

Logistic Regression

- When the goal is to classify the data points (samples) into categories (or labels), you can use Logistic Regression.

- Output: 0 or 1 , or a probability estimate

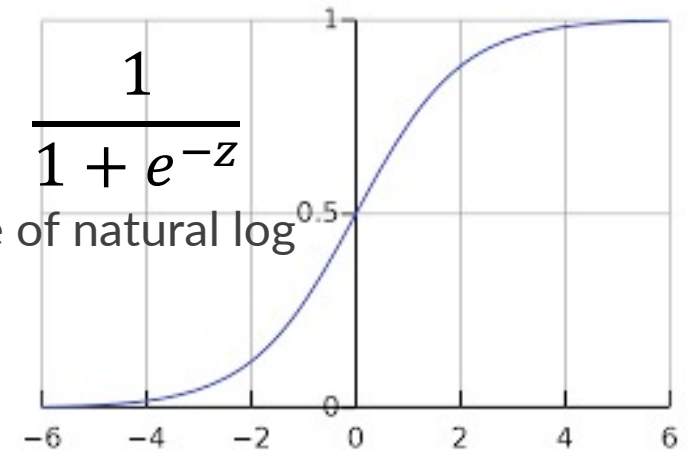
$$z = f(x; w) = w^T x$$

sigmoid function: $\text{sigm}(z) = \frac{1}{1 + e^{-z}}$

```
def sigmoid(z):
```

```
    return 1.0 / (1 + np.exp(-z))
```

Returns a probabilistic estimate p



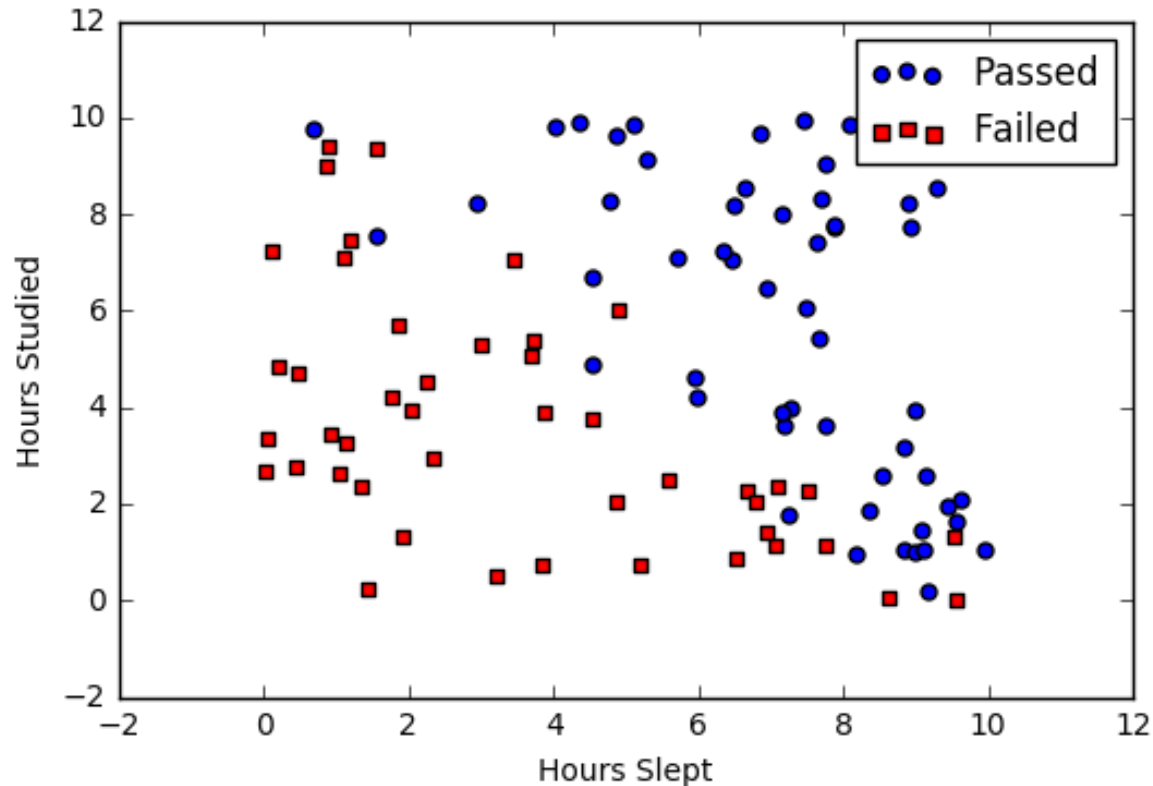
$$y^{(i)} = \begin{cases} 0 & \text{if } p < \text{threshold} \\ 1 & \text{if } p \geq \text{threshold} \end{cases}$$

$$g(x; w) = \text{sigm}(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

$$y^{(i)} = \begin{cases} 0 & \text{if } g(w^T x^{(i)}) < \text{threshold} \\ 1 & \text{if } g(w^T x^{(i)}) \geq \text{threshold} \end{cases}$$

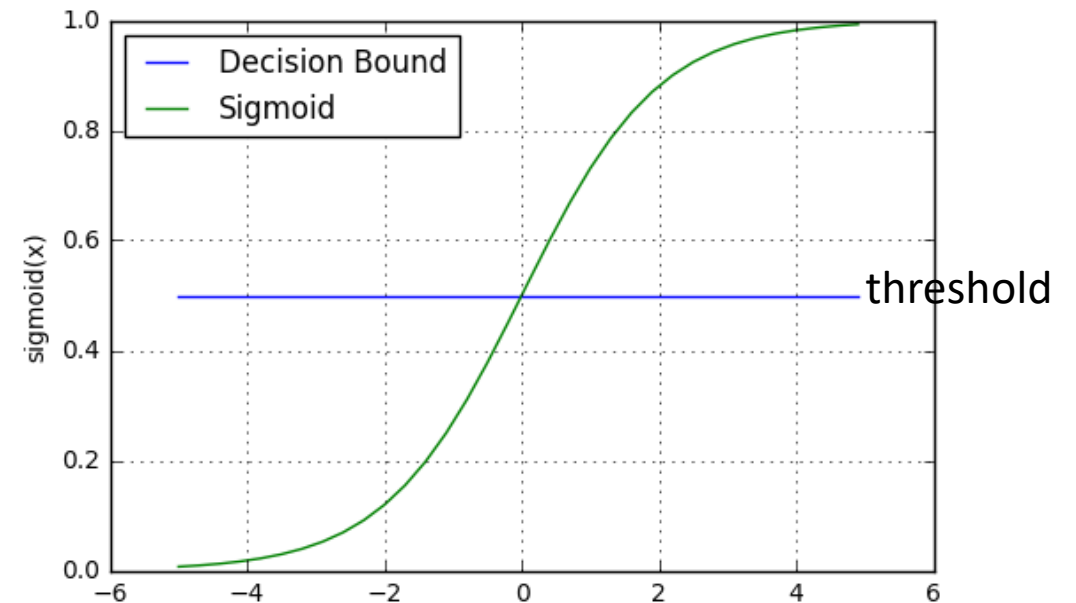
given the S curve which can capture more complex relationships compared to linear regression and its less sensitivity to outliers

Logistic Regression cont.



© Copyright 2017 Revision 43e12019. [Read the Docs.](#)

Studied	Slept	Passed
4.85	9.63	1
8.62	3.23	0
5.43	8.23	1
9.21	6.34	0



$$y^{(i)} = \begin{cases} 0 & \text{if } p < \text{threshold} \\ 1 & \text{if } p \geq \text{threshold} \end{cases}$$

MLE for Logistic Regression

1. Formulate Logistic regression using the sigmoid function
2. Formulate logistic regression as a MLE problem
3. Use Gradient Descent to find the optimal parameters of the model

$$\begin{aligned} p(y^{(i)} = 1 | x^{(i)}; w) &= g(x^{(i)}; w) \\ p(y^{(i)} = 0 | x^{(i)}; w) &= 1 - g(x^{(i)}; w) \end{aligned}$$

OLS does not do a good job in non-linear complex problems. Where you don't have a close form solution. In that case, Gradient Descent is a commonly used to tune the model parameters.

$$p(y^{(i)} | x^{(i)}; w) = g(x^{(i)}; w)^{y^{(i)}} (1 - g(x^{(i)}; w))^{1-y^{(i)}}$$

MLE for Logistic Regression cont.

1. Formulate Logistic regression using the sigmoid function
2. Formulate logistic regression as a MLE problem
3. Use Gradient Descent to find the optimal parameters of the model

$$p(y^{(i)} | x^{(i)}; w) = g(x^{(i)}; w)^{y^{(i)}} (1 - g(x^{(i)}; w))^{1-y^{(i)}}$$

$$l(w) \triangleq \log p(D|w) = \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}; w)$$

Assumption: Samples are iid

$$= \sum_{i=1}^N \log(g(x^{(i)}; w)^{y^{(i)}} (1 - g(x^{(i)}; w))^{1-y^{(i)}})$$

$$= \sum_{i=1}^N (y^{(i)} \log g(x^{(i)}; w) + (1 - y^{(i)}) \log(1 - g(x^{(i)}; w)))$$

MLE for Logistic Regression cont.

1. Formulate Logistic regression using the sigmoid function
2. Formulate logistic regression as a MLE problem
3. Use Gradient Descent to find the optimal parameters of the model

$$l(w) = \sum_{i=1}^N (y^{(i)} \log g(x^{(i)}; w) + (1 - y^{(i)}) \log(1 - g(x^{(i)}; w)))$$

$$\frac{\partial l(w)}{\partial w_j} = \frac{\partial (\sum_{i=1}^N (y^{(i)} \log g(x^{(i)}; w) + (1 - y^{(i)}) \log(1 - g(x^{(i)}; w))))}{\partial w_j}$$

$$= (y^{(i)} - g(x^{(i)}; w)) x_j^{(i)}$$

Now take Gradient Ascent!

Gradient Descent (Ascent)

$$\frac{\partial l(w)}{\partial w_j} = \left(y^{(i)} - g(x^{(i)}; w) \right) x_j^{(i)}$$

$$w_j = w_j + a \frac{\partial l(w)}{\partial w_j}$$

Step size or learning rate

Updated weights

$$w_j = w_j + a \left(y^{(i)} - g(x^{(i)}; w) \right) x_j^{(i)}$$

Is there another way to tune the parameters ?

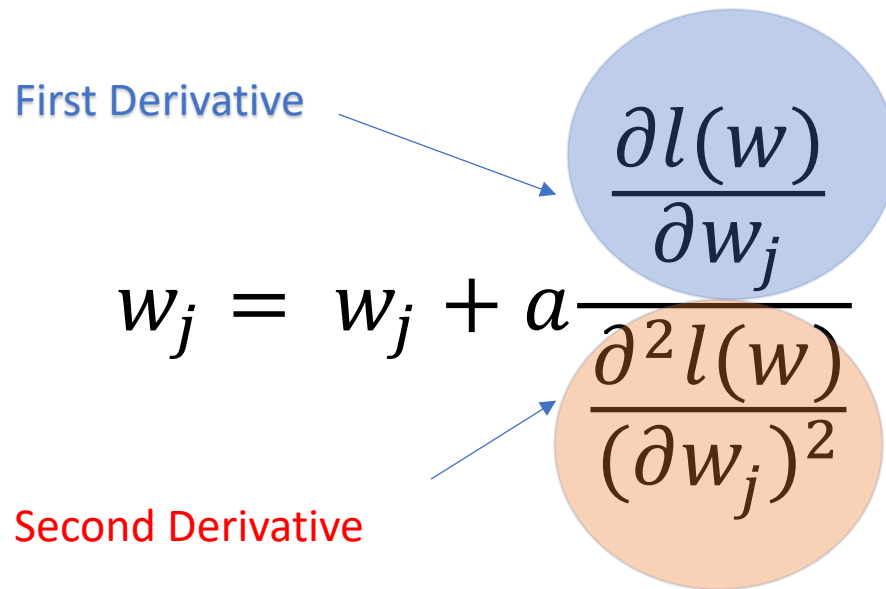
Newton's Method

- Another method we could use instead of GD to update the weights

First Derivative

$$w_j = w_j + a \frac{\frac{\partial l(w)}{\partial w_j}}{\frac{\partial^2 l(w)}{(\partial w_j)^2}}$$

Second Derivative

The diagram shows the Newton's Method update formula for weight w_j. The formula is w_j = w_j + a * (first derivative) / (second derivative). The first derivative, ∂l(w)/∂w_j, is enclosed in a light blue circle and labeled 'First Derivative' with a blue arrow pointing to it. The second derivative, ∂²l(w)/(∂w_j)², is enclosed in a light orange circle and labeled 'Second Derivative' with a red arrow pointing to it.

First derivative shows the slope of the tangent line.

Second derivative measures the instantaneous rate of change of the first derivative. The sign of the second derivative tells whether the slope of the tangent line is increasing or decreasing.

Newton's Method cont.

- Newton's is also called the **Newton-Raphson** method
- It is a root-finding algorithm that uses the first few terms of the **Taylor series** of a function $f(x)$ about a point x

$$w_{n+1} = w_n - \eta_k \frac{\frac{\partial l(w)}{\partial w_j}}{\frac{\partial^2 l(w)}{(\partial w_j)^2}}$$

Step size

$$w_{n+1} = w_n - \eta_k \frac{g_k}{H_k} = w_n + a d_k$$

g_k

H_k

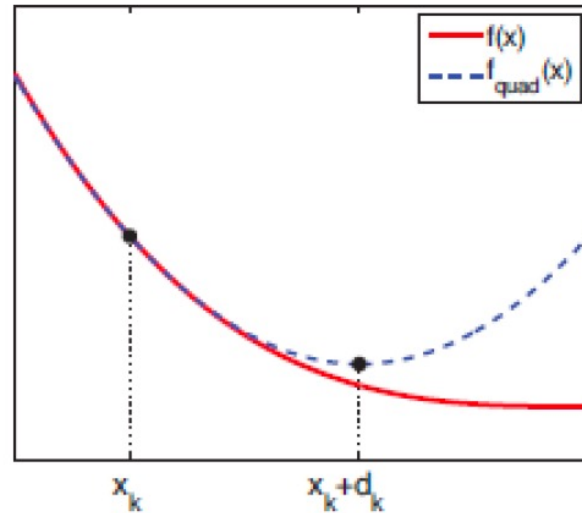
$-H_k^{-1} g_k = d_k$

Algorithm 8.1: Newton's method for minimizing a strictly convex function

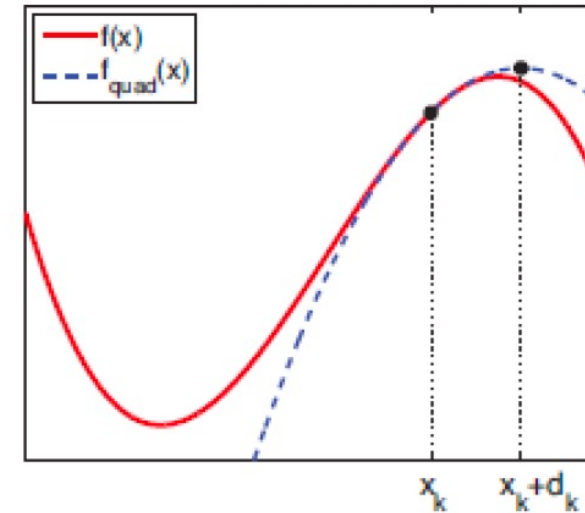
```
1 Initialize  $\theta_0$ ;  
2 for  $k = 1, 2, \dots$  until convergence do  
3   Evaluate  $\mathbf{g}_k = \nabla f(\theta_k)$ ;  
4   Evaluate  $\mathbf{H}_k = \nabla^2 f(\theta_k)$ ;  
5   Solve  $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$  for  $\mathbf{d}_k$ ;  
6   Use line search to find stepsize  $\eta_k$  along  $\mathbf{d}_k$ ;  
7    $\theta_{k+1} = \theta_k + \eta_k \mathbf{d}_k$ ;
```

$$w_j := w_j + \alpha \frac{\frac{\partial l(w)}{\partial w_j}}{\frac{\partial^2 l(w)}{(\partial w_j)^2}}$$

$$\theta_{k+1} = \theta_k - \eta_k \mathbf{H}_k^{-1} \mathbf{g}_k$$



(a)



(b)

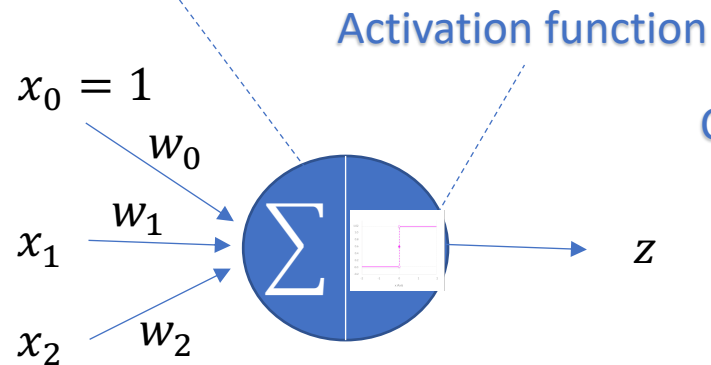
@IliasTagkopoulos

Newton's Method cont.

- Newton method is good for minimizing a strictly convex function.
- Newton's method make a quadratic approximation to the original function.
- A good read:
 - To get more information on how Taylor Expansion is used to formulate Newton's method. Read section "**2 Newton's Method for Numerical Optimization**" from the following article: <http://www.stat.cmu.edu/~cshalizi/350/lectures/26/lecture-26.pdf>

Perceptron Learning Algorithm

- The perceptron model takes an input, aggregates it to calculate the weighted sum, and then returns 1 if the weighted sum is more than a threshold, or else returns 0.



$$w_j = w_j + a \left(y^{(i)} - g(x^{(i)}; w) \right) x_j^{(i)}$$

z = output

Combination function: $z = w^T x = w_0 + w_1 x_1 + w_2 x_2$

Threshold can be any scalar value (such as 0)

$$g(z) = \begin{cases} 0 & \text{if } z < \text{threshold} \\ 1 & \text{if } z \geq \text{threshold} \end{cases}$$

- In general, Stochastic Gradient Descent on logistic regression is faster and has a better performance.