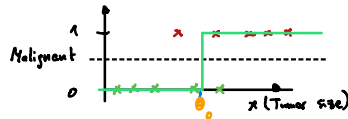


* Classification



MSE \rightarrow Rep-

$$\min_{\theta} \frac{1}{2} (y - g(h_{\theta}(x)))^2$$

$$g(h_{\theta}(x)) = \begin{cases} 0 & h_{\theta}(x) < 0 \\ 1 & h_{\theta}(x) \geq 0 \end{cases}$$

Thresholding

(classification) \rightarrow linearly

Logistic regression: Force to have $0 < g(h_{\theta}(x)) < 1$

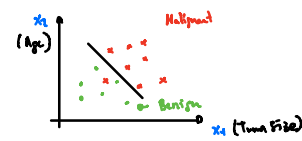
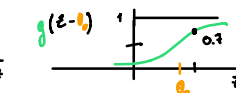
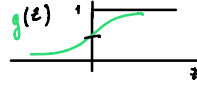
g \rightarrow Sigmoid function

$h_{\theta}(x) \rightarrow$ hypothesis function

$$w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n \leftarrow \text{Representation!}$$

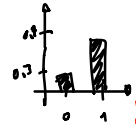
* Sigmoid function: (Logistic function)

$$g(z) = \frac{1}{1 + e^{-z}}$$



* Interpretation of $g(h_{\theta}(x))$: Estimated probability that $y=1$ on input x

For some x_i $g(h_{\theta}(x_i)) = 0.7 \Rightarrow 70\%$ chance of tumor being malignant



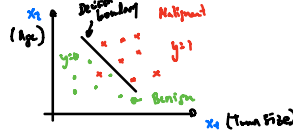
$$g(h_{\theta}(x)) = p(y=1 | x; \theta)$$

$$1 - g(h_{\theta}(x)) = 1 - p(y=1 | x; \theta) = p(y=0 | x; \theta) = 0.3$$

* Decision boundary

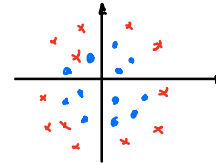
$$g(\theta^T x) = \begin{cases} 1 & \theta^T x \geq 0 \\ 0 & \theta^T x < 0 \end{cases}$$

Linear



* Non-linear decision boundary

$$g(\theta_1^T x + \theta_2^T x^2) = g([\theta_1 \ \theta_2] \begin{bmatrix} x_1 \\ x_1^2 \end{bmatrix}) = \begin{cases} 1 & \theta_1^T x + \theta_2^T x^2 \geq 0 \\ 0 & \theta_1^T x + \theta_2^T x^2 < 0 \end{cases}$$



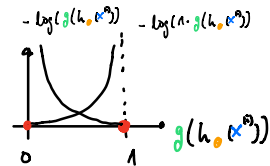
* Logistic regression model

Given $y \in \{0, 1\}$, Bernoulli: $\Pr(y | x; \theta) = g(h_{\theta}(x))^y (1 - g(h_{\theta}(x)))^{1-y}$

Likelihood (iid): $L(\theta | y; x) = \Pr(y | x; \theta) = \prod_i \Pr(y_i | x_i; \theta)$

$$= \prod_i g(h_{\theta}(x_i))^{y_i} (1 - g(h_{\theta}(x_i)))^{1-y_i}$$

Taking max log likelihood $\Rightarrow \min_{\theta} J(\theta)$ MLE!



$\begin{cases} \text{if } y^i = 0 \Rightarrow \text{penalty } 1 = g(h_{\theta}(x^i)) \\ \text{if } y^i = 1 \Rightarrow \text{penalty } 0 = g(h_{\theta}(x^i)) \end{cases}$

$$* \text{Cost: } J(\theta) = \frac{1}{m} \sum_{i=1}^m (1 - y^i) [-\log(1 - g(h_{\theta}(x^i)))] + y^i [-\log(g(h_{\theta}(x^i)))]$$

* Gradient

$y_i = 0 \Rightarrow g = 0$ $y_i = 1 \Rightarrow g = 1 \Rightarrow -\log(1-g) \rightarrow \infty$

$\frac{\partial h_{\theta}}{\partial \theta} = h = \theta^T x$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (1 - y^i) \frac{-1}{(1 - g(h_{\theta}(x^i)))} - (g(h_{\theta}(x^i)) (1 - g(h_{\theta}(x^i))) x^{(i)j})$$

$$+ (y^i) \frac{-1}{(g(h_{\theta}(x^i)))} - (g(h_{\theta}(x^i)) (1 - g(h_{\theta}(x^i))) x^{(i)j})$$

$$= \frac{1}{m} \sum_{i=1}^m [g(h_{\theta}(x^i)) - y^i] x^{(i)j}$$

Same as linear regression!
Easy!

* Playing with $\left\{ \begin{array}{l} \text{cost function definition} \\ + \\ \text{composition of non-linear functions} \end{array} \right\}$ \rightarrow best guess \rightarrow Linear sigmoid kernel

* Minimization

$$\begin{aligned} \min J(\theta) \Rightarrow & \begin{cases} \nabla J(\theta) = 0 \rightarrow \text{Eq. no linear in general} \\ \theta_{k+1} = \theta_k - \alpha \nabla J(\theta_k) \Rightarrow J(\theta_{k+1}) < J(\theta_k) \end{cases} \\ \text{s.t.} & \quad \downarrow \text{step size} \quad \uparrow \text{learning rate} \\ & \quad \quad \quad \text{step size} \rightarrow \text{line search} \end{aligned}$$

* Ill-posedness

Scenario 1: $X = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ $y = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

Scenario 2: $X = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ $y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Regularization: $\min_{\theta} J(\theta) = \|X\theta - y\|^2 + \lambda \|\theta\|^2$

Scenario 1: $\min_{\theta} J(\theta) = 0 \Rightarrow X^T X \theta = X^T y$

Scenario 2: $\min_{\theta} J(\theta) = 0 \Rightarrow X^T X \theta = X^T y$

Rank of $X^T X$ is 1

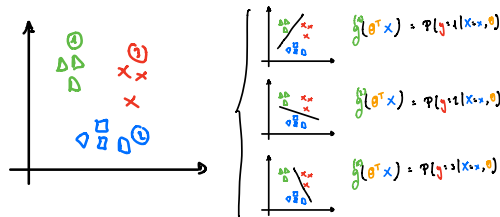
Rank of $X^T X$ is 2

lim. repr. \rightarrow One V_1 all \rightarrow NN

* Multi-classification

* One-vs-all:

$$y \in \{1, 2, 3\}$$



How is the cost?

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K (1 - y_k^i) (-\log(1 - g(h_{\theta}(x^i))) + y_k^i (-\log(g(h_{\theta}(x^i)))) \right] + \lambda \|\theta\|^2$$

$$y_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}; y_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}; \dots$$

* Practice: Extend logistic regression (classification) to multi-classification

* Neural model: Logistic unit

