

NDD 1.  $\rightarrow$  rough work.  
single feature  $\rightarrow h_0(u) = \theta_0 + \theta_1 u_1$ .

$$h_0(u) = \theta_0 + \theta_1 u_1 + \theta_2 u_2 + \dots + \theta_n u_n. \quad x = \begin{bmatrix} u_0 \\ u_1 \\ \vdots \\ u_n \end{bmatrix} \in \mathbb{R}^n$$

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_0 u = \theta_0 x_0 + \theta_1 u_1 + \dots + \theta_n u_n = \Theta^T \begin{bmatrix} 1 & x_{1:n} \end{bmatrix}$$

Cost function:

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_0(u^{(i)}) - y^{(i)})^2$$

Concept of Gradient Descent. (One feature)

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1) \quad \text{say } \theta_0 = 0, \theta_1 = 1$$

Change  $\theta_0 / \theta_1$  to minimize  $J(\theta_0, \theta_1)$

repeat until convergence.

$$\Rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0 \text{ or } j=1)$$

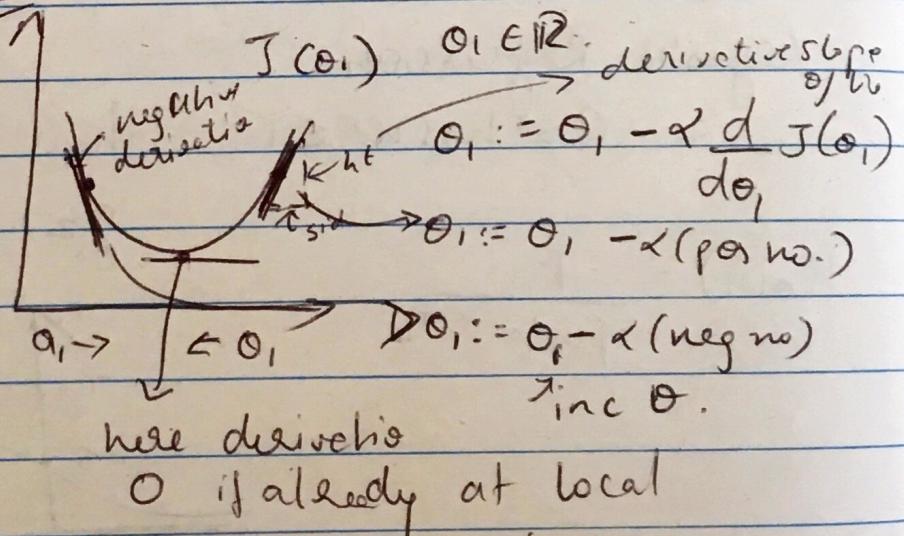
Simultaneously update  $\theta_0$  &  $\theta_1$

learning rate

$$\min_{\theta_1} J(\theta_1)$$

$\alpha$  small  $\rightarrow$  slow

$\alpha$  large  $\rightarrow$  overshoot  
(fall to convex)  
diverge



$\theta$  if already at local

Gradient descent for  $n$  variables.  $n \geq 1$ .

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(u^{(i)}) - y^{(i)}) u_j^{(i)}$$

simultaneously update  $\theta_j$  for all  $j$ .

Stochastic gradient descent.

$$\text{cost}(\theta, (u^{(i)}, y^{(i)})) = \frac{1}{2} (h_\theta(u^{(i)}) - y^{(i)})^2$$

$$J_{\text{stoch}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (u^{(i)}, y^{(i)}))$$

1) Randomly shuffle dataset.

2) Repeat {  
for  $i = 1, \dots, m$  {  
 $\theta_j = \theta_j - \alpha \frac{\partial \text{cost}(\theta, (u^{(i)}, y^{(i)}))}{\partial \theta_j}$   
1-10 times of } for  $j = 0, \dots, n$ .  
depends on size of dataset. } }  $\downarrow$   
 $\frac{\partial}{\partial \theta_j} \text{cost}(\theta, (u^{(i)}, y^{(i)}))$

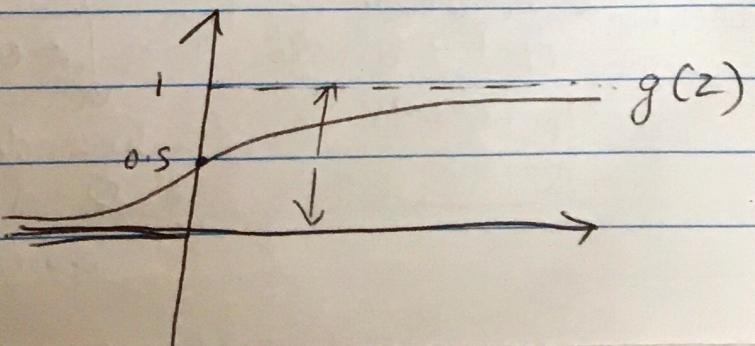
Checking for convergence. If divergence we smoothen.  
Can slowly decrease  $\alpha$ .

Logistic Regression.

$$0 \leq h_\theta(u) \leq 1 \quad h_\theta(u) = g(\theta^\top u)$$

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

$$h_\theta(u) = \frac{1}{1 + e^{-\theta^\top u}}$$



$$h_0(u) = P(y=1 | u; \theta)$$

Suppose "y=1" if  $h_0(u) \geq 0.5 \rightarrow \theta^T u \geq 0$   
 "y=0" if  $h_0(u) < 0.5 \rightarrow \theta^T u < 0$

$$h_0(u) \geq 0.5 \rightarrow \theta^T u \geq 0$$

Decision boundary:  $h_0(u) = g(\theta_0 + \theta_1 u_1 + \theta_2 u_2)$

$$\text{if } \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \text{ predict } y=1 \text{ if } -3 + u_1 + u_2 \geq 0$$

$$\theta^T u \downarrow u_1 + u_2 \geq 3$$

decision boundary where  $h_0(u) = 0.5$

$$\text{if } h_0(u) = g(\theta_0 + \theta_1 u_1 + \theta_2 u_2 + \theta_3 u_1^2 + \theta_4 u_2^2)$$

$$y=1 \text{ if } -1 + u_1^2 + u_2^2 \geq 0. \quad \theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

(cost func.  $\rightarrow$  training set  $\{(u^{(1)}, y^{(1)}), (u^{(2)}, y^{(2)}) \dots (u^{(m)}, y^{(m)})\}$ )  
 m.e.g.  $n \in \begin{bmatrix} \vdots \\ u_n \end{bmatrix} \subset \mathbb{R}^{n+1}$ ,  $u=1, y \in \{0, 1\}$ .

$$h_0(u) = \frac{1}{1 + e^{-\theta^T u}} \quad J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_0(u^{(i)}), y^{(i)})$$

$$\text{cost}(h_0(u), y) = \frac{1}{2} (h_0(u) - y)^2 \leftarrow \text{not convex}$$

$$\text{cost}(h_0(u), y) = \begin{cases} -\log(h_0(u)) & \text{if } y=1 \\ -\log(1-h_0(u)) & \text{if } y=0 \end{cases}$$

Given  $u$ , want  $\hat{y} = P(y=1 | u)$ .  $u \in \mathbb{R}^{n_u}$ .

$$\text{OP } \hat{y} = g(w^T u + b) \quad \begin{array}{c} \text{if } w^T u + b \geq 0 \\ \frac{1}{1+e^{-z}} \end{array} \quad \begin{array}{c} \text{if } w^T u + b < 0 \\ G(z) \end{array}$$

Decision Trees.

$$H(S) = -P_{C+} \log_2 P_{C+} - P_{C-} \log_2 P_{C-} \quad \text{for } S$$

$P_{C+} + P_{C-} \rightarrow$  in one subset not whole dataset

Information gain  $\rightarrow$  avg of entropy of subset

$$= \sum_i p_i \log_2 (p_i)$$

$$\text{Gain}(S, A) = H(S) - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} H(S_v)$$

Gain ratio:

$$\text{Split Entropy}(S, A) = - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \log \frac{|S_v|}{|S|}$$

or continuous.

$$\text{Gain ratio} = \frac{\text{Gain}(S, A)}{\text{Split Entropy}(S, A)}$$

$$\hookrightarrow - \int_{y \in Y} p(y) \log(p(y)) dy$$

Multiclass & regression.

$$H(S) = - \sum_c p_{Cc} \log(p_{Cc})$$

Random Forest:

$$p(C|U) = \frac{1}{T} \sum_{t=1}^T p_t(C|U)$$

Randomization

$$T_j \in T$$

$$o_j^* = \arg \max_{o_j \in T_j} f_j$$