

$$X = \begin{bmatrix} | & | & | & \dots & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \\ | & | & | & \dots & | \end{bmatrix} \quad \begin{array}{l} n \text{ features, } m \text{ examples.} \\ \underline{n \times m} \rightarrow \text{diff. from linear reg.} \end{array}$$

$$\underline{Y} = [ \underline{y}^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(m)} ] \quad \underline{1 \times m}$$

Logistic Regression:

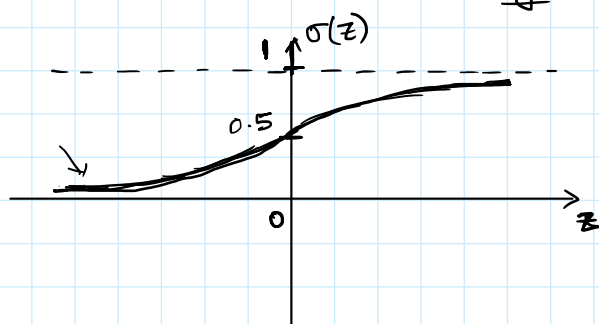
example  $x$ , we want  $\hat{y} = P(y=1 | x)$ .  $\hat{y}$  prediction.  $\square \rightarrow \begin{array}{l} \text{CAT } (y=1) \\ \text{NOT CAT } (y=0) \end{array}$

in linear reg.,  $h(x) = \theta_0 + \theta_1 x$  Here, define 2 parameters  $w, b$   $w \in \mathbb{R}^n$   
 $b \in \mathbb{R}$ .

$$\hat{y} = \boxed{w^T x + b} \quad \begin{array}{l} \leftarrow \text{APPLY FN.} \\ \downarrow \text{one example} \\ (1 \times n) \quad (n \times 1) \end{array}$$

we need  $0 \leq \hat{y} \leq 1$ . use SIGMOID

$$\hat{y} = \sigma \left( \underbrace{w^T x + b}_z \right). \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$



$z$  very large,  $z \rightarrow \infty$ ,  $\sigma(z) \rightarrow 1$   
 $z$  very small  $z \rightarrow -\infty$ ,  $\sigma(z) \rightarrow 0$

$\theta_0 \rightarrow$  we add a new feature  $x_0 = 1$   
 $\theta_1$   
 $\theta_2$   
 $\vdots$   
 $\underline{\theta^T x}$   $\times$  keep  $b$  separate

$$\hat{y} = \sigma \left( \frac{w^T x + b}{z} \right) \quad \text{COST function. } L(\hat{y}, y) = - \left( y \log \hat{y} + (1-y) \log(1-\hat{y}) \right)$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

function of parameters.

CROSSENTROPY LOSS. how diff.  $\hat{y}$  is from  $y$ .

Now we have: • hypothesis  $\rightarrow \hat{y} = \sigma(w^T x + b)$   
• cost  $\rightarrow J(w, b)$ .

GOAL: minimize  $J$ .  $\rightarrow$  gradient descent

Repeat until convergence:  $w = w - \alpha \cdot \frac{\partial J}{\partial w}$

$$dw = \frac{\partial J}{\partial w} = \frac{\partial L}{\partial w}, \quad db = \frac{\partial L}{\partial b}$$

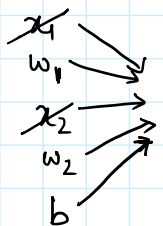
$$b = b - \alpha \frac{\partial J}{\partial b} \quad \rightarrow \text{LEARNING RATE}$$

Set  $n=2$ .

$$z = w^T x + b$$

$$= w_1 x_1 + w_2 x_2 + b$$

$$\hat{y} = \sigma(z) \rightarrow \text{predictions.}$$



$$z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = \sigma(z) \rightarrow L(\hat{y}, y).$$

computation graph.

$$dz = \frac{\partial L}{\partial z} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z}$$

$$d\hat{y} = \frac{\partial L}{\partial \hat{y}}$$

$$L = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

$$\hat{y} = \sigma(z) = \frac{1}{1+e^{-z}} = -\left(\frac{y}{\hat{y}} - \frac{(1-y)}{1-\hat{y}}\right) = -\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})}$$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z)$$

$$= \hat{y}(1-\hat{y})$$

$$\frac{e^{-z}}{(1+e^{-z})^2} \left(\frac{1}{1+e^{-z}}\right) \cdot \frac{1+e^{-z}-1}{(1+e^{-z})} = \frac{1}{(1+e^{-z})} \left(1 - \frac{1}{1+e^{-z}}\right) = \hat{y}(1-\hat{y})$$

$$dz = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z}$$

$$= \left(-\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})}\right) \cdot \hat{y}(1-\hat{y})$$

$$= -y(1-\hat{y}) + (1-y)\hat{y} = -y + y\hat{y} + \hat{y} - y\hat{y} = \hat{y} - y$$

$$\frac{\partial L}{\partial z} = dz = \hat{y} - y$$

$$\frac{\partial z}{\partial w_1} = x_1, \quad \frac{\partial z}{\partial w_2} = x_2$$

$$\frac{\partial z}{\partial b} = 1$$

$$\frac{\partial L}{\partial w_1} = dw_1 = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_1} = x_1 dz = x_1(\hat{y} - y)$$

$$dw_2 = x_2 dz = x_2(\hat{y} - y)$$

$$db = dz = \hat{y} - y$$

FOR ONE STEP OF GRADIENT DESCENT:  $w_1 = w_1 - \alpha dw_1$ ,  $w_2 = w_2 - \alpha dw_2$   
 $b = b - \alpha db$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{bmatrix}_{n \times 1}$$

$$z^{(1)} = w^T x^{(1)} + b$$

$$z^{(2)} = w^T x^{(2)} + b$$

$$z^{(3)} = w^T x^{(3)} + b$$

$$\hat{y}^{(1)} = \sigma(z^{(1)})$$

$$\hat{y}^{(2)} = \sigma(z^{(2)})$$

$$\hat{y}^{(3)} = \sigma(z^{(3)})$$

$$Z = w^T X + b$$

$$\begin{matrix} \downarrow & \downarrow \\ (1,n) & (n,m) \\ & \searrow \\ & (1,m) \end{matrix}$$

$$[z^{(1)} \quad z^{(2)} \quad z^{(3)} \quad \dots]$$

$$\begin{matrix} \uparrow & \uparrow & \uparrow \\ \text{add same } b. & & \end{matrix} \quad dz = [dz^{(1)} \quad dz^{(2)} \quad \dots]_{(1,m)}$$

$$dz = \hat{y} - y$$

$$dw_1 = x_1 dz$$

$$dw_2 = x_2 dz$$

$$\hat{Y} = \sigma(Z)$$

$$dZ = \hat{Y} - Y$$

$$dw = X \cdot dZ^T \cdot \frac{1}{m}$$

ALL EXAMPLES, using J.

$$Y = \sigma(Z).$$

$$dZ = Y - 1.$$

$$\underline{dw} = \underset{(n,m)}{X} \cdot \underset{(m,1)}{dZ}^T \cdot \underset{(n,1)}{\frac{1}{m}}.$$

ALL EXAMPLES, using J.

$$J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}).$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} = \frac{1}{m} \text{np.sum}(dZ), \quad \frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m dw^{(i)} = \frac{1}{m} \cdot \underline{dw}$$

$$\frac{1}{m} \sum db = \frac{1}{m} \sum dz^{(i)}$$

$$dZ \checkmark, dw, db.$$

$$w = \underset{(n,1)}{w} - \alpha \underset{(n,1)}{dw}, \quad b = b - \alpha db. \rightarrow \text{one step of GD.}$$

$$\underline{J} = - \frac{1}{m} \left( \underline{Y} \cdot (\text{np.log}(\hat{Y}))^T + (1 - \underline{Y}) (\text{np.log}(1 - \hat{Y}))^T \right).$$