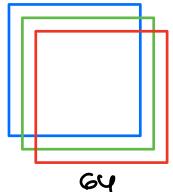


lec 1. Binary classification:

64  \rightarrow 1(cat) vs 0(non cat)

RGB ($64 \times 64 \times 3$) \rightarrow feature vector

 \rightarrow stack them $x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}_{64 \times 64 \times 3}$ $n_x = 64 \times 64 \times 3 = 12288$

$n_x = n$ = dimensions of input feature vector (x)

Binary classification

$X \rightarrow Y$
image (0,1)

Notation:

single training example $(x_i, y_i) \quad x_i \in \mathbb{R}^{n_x}, y_i \in \{0, 1\}$

m training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

training set $m = N_{\text{train}}$ n_x = dimension of feature vector

test set $m = N_{\text{test}}$

$$X = \left[\begin{array}{cccc} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | & | \end{array} \right] \quad m \quad n_x \quad X \in \mathbb{R}^{n_x \times m}$$

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

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} (1, m)$$

$$X.\text{shape} (n_x, m)$$

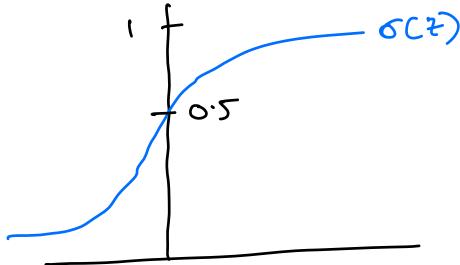
Lec 2: Logistic Regression :-

Given x , want $\hat{y} = P(y=1|x)$ ($0 \leq \hat{y} \leq 1$)

$$x \in \mathbb{R}^{n_x}$$

parameters of NN : $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

$$\text{output : } \hat{y} = \sigma(w^T x + b)$$



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

$$\text{if } z \text{ large } \sigma(z) = \frac{1}{1+0} = 1$$

$$\text{if } z \text{ large negative } \sigma(z) = \frac{1}{1+0} \approx 0$$

$$\text{if } z = 0 \quad \sigma(z) = \frac{1}{1+1} = 0.5$$

lec 3:- logistic Regression Cost function

$$\hat{y} = \sigma(\omega^T x + b) \quad \text{where } \sigma(z) = \frac{1}{1+e^{-z}}$$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$
 pred ground truth

$$z^{(i)} = \omega^T x^{(i)} + b \quad \text{i-th example}$$

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

loss(error) function: $\ell(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

GD doesn't work well
 with squared error

$$\ell(\hat{y}, y) = - \left[y \log \hat{y} + (1-y) \log (1-\hat{y}) \right] \quad \text{for i-th example}$$

if $y=1$ $\ell(\hat{y}, 1) = -\log \hat{y} \Rightarrow \log \hat{y}$ as big as possible
 $\Rightarrow \hat{y}$ has to be large ≈ 1

if $y=0$ $\ell(\hat{y}, 0) = -\log(1-\hat{y}) \Rightarrow \log(1-\hat{y})$ as big as possible
 $\Rightarrow 1-\hat{y}$ has to be large ≈ 1
 $\Rightarrow \hat{y}$ has to be min.

Cost Function:

$$\begin{aligned} J(\omega, b) &= \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right] \end{aligned}$$

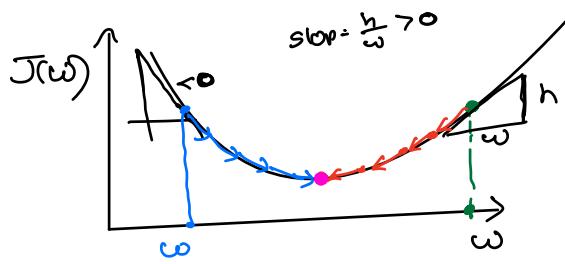
\Rightarrow average loss function for entire training set.

dec 4: Gradient Descent :-

$$\hat{y} = \sigma(\omega^T x + b) \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

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

want to find (ω, b) that minimizes $J(\omega, b)$



Repeat $\tilde{\omega}$

$$\omega := \omega - \alpha \frac{dJ(\omega)}{d\omega}$$

\tilde{g} α - learning rate

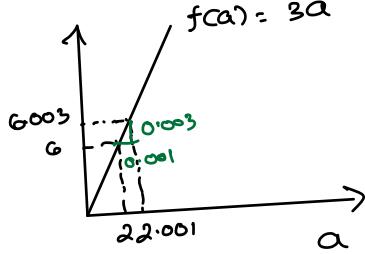
$$d\omega = \frac{dJ(\omega)}{d\omega} = \text{gradient}$$

$$J(\omega, b) \quad \omega := \omega - \alpha \frac{\partial J(\omega, b)}{\partial \omega}$$

$$b := b - \alpha \frac{\partial J(\omega, b)}{\partial b}$$

$J(\omega, b)$	$\omega := \omega - \alpha d\omega$
	$b := b - \alpha db$

dec 5: Derivatives:



$$\begin{aligned} a &= 2 & f(a) &= 6 \\ a &= 2.001 & f(a) &= 6.003 \end{aligned}$$

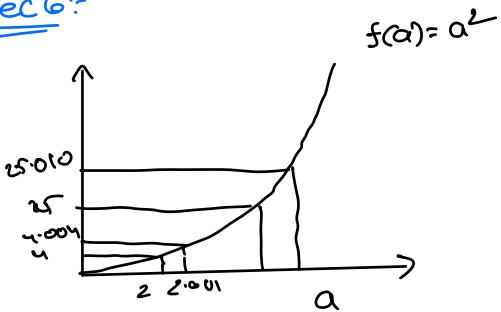
derivative of $f(a)$ at $a=2 = 3$

$$f'(2) = \frac{\Delta f(a)}{\Delta a} = \frac{0.003}{0.001} = 3$$

$$\begin{aligned} a &= 5 & f(a) &= 15 \\ a &= 5.001 & f(a) &= 15.003 \end{aligned}$$

$$f'(5) = \frac{\Delta f(a)}{\Delta a} = \frac{0.003}{0.001} = 3$$

dec 6:



$$f(a) = a^2$$

$$f'(a) = 2a$$

$$f'(2) = 4 \quad f'(5) = 10$$

$$\begin{aligned} a &= 2 & f(a) &= 4 \\ a &= 2.001 & f(a) &= 4.004 \end{aligned}$$

$$f'(2) = \frac{0.004}{0.001} = 4$$

$$\begin{aligned} a &= 5 & f(a) &= 25 \\ a &= 5.001 & f(a) &= 25.010 \end{aligned}$$

$$f'(5) = \frac{0.01}{0.001} = 10$$

Dec 7 :- Computation Graph:

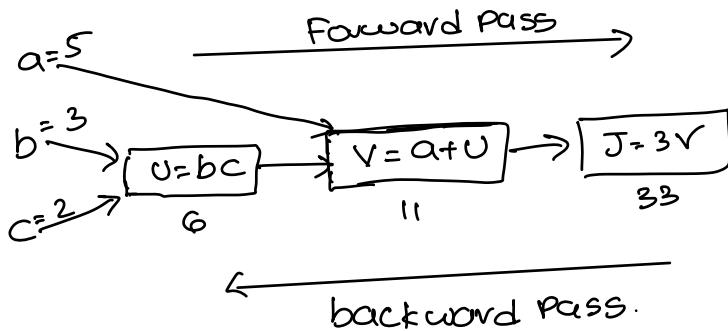
$$J(a, b, c) = 3(a + \frac{bc}{v})$$

$\boxed{\frac{v}{\sigma}}$

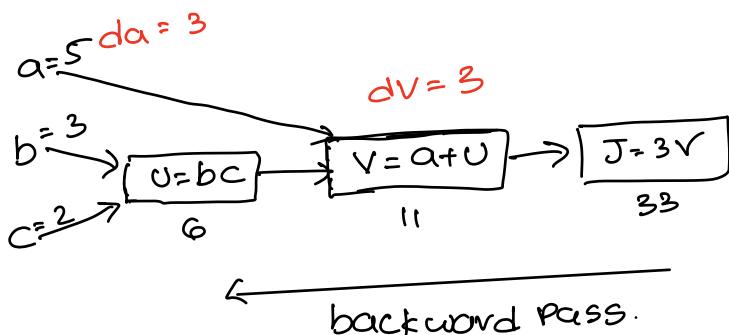
$$U = bc$$

$$V = a + U$$

$$J = 3V$$



Dec 8 :- Derivatives with a computation graph.



$$\frac{dJ}{dy} = ?$$

$$J = 3V$$

$$dV = \frac{dJ}{dV} = 3$$

$$\frac{dJ}{da} = ?$$

$$a = 5 \rightarrow 5.001$$

$$v = 11 \rightarrow 11.001$$

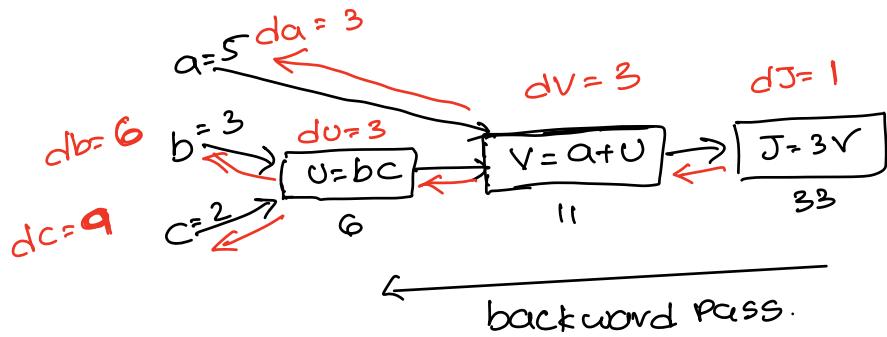
$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{da} = \frac{0.003}{0.001} = 3$$

$$\frac{dJ}{da} = \frac{dJ}{dV} \frac{dV}{da} = \frac{d(3V)}{dV} \frac{d(a+U)}{da} = (3)(1) = 3$$

$$\therefore \frac{dJ}{da} = 3$$

$$\frac{d \text{Final output variable}}{d \text{var}} = d \text{var}$$



$$du = \frac{dJ}{du} = \frac{dJ}{dv} \frac{dv}{du} = \frac{d(3v)}{dv} \frac{d(a+u)}{du} = (3)(1) = 3$$

$$db = \frac{dJ}{db} = \frac{dJ}{du} \frac{du}{db} = (3) \frac{d(bc)}{db} = (3)(c) = 3 \times 2 = 6$$

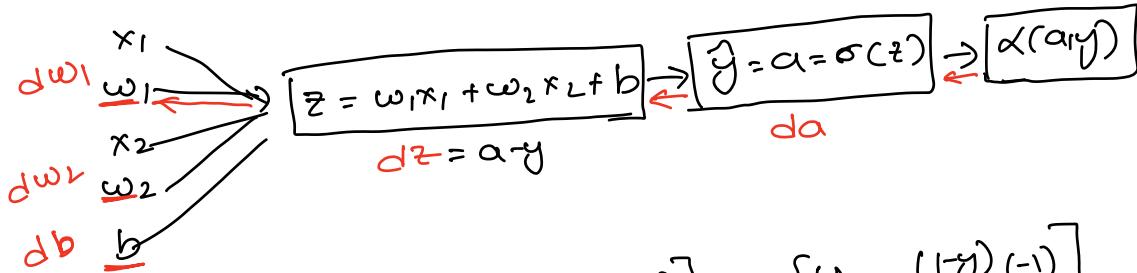
$$dc = \frac{dJ}{dc} = \frac{dJ}{du} \frac{du}{dc} = (3) \frac{d(bc)}{db} = (3)(b) = 3 \times 3 = 9$$

dec a :- logistic regression (GD with single training example).

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

$$\hat{y} = a = \sigma(z)$$

$$\alpha(a, y) = -[y \log a + (1-y) \log(1-a)]$$



$$\frac{d\alpha}{da} = \frac{dL}{da} = \frac{d - [y \log a + (1-y) \log(1-a)]}{da} = -\left[\frac{y}{a} + \frac{(1-y)}{(1-a)}(-1)\right]$$

$$\frac{d\alpha}{da} = -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

$$dz = \frac{dL}{dt} = \frac{dL}{da} \frac{da}{dt} = (da) \frac{d\sigma(t)}{dt} =$$

$$\sigma(t) = \frac{1}{1+e^{-t}}$$

$$\begin{aligned} \frac{d\sigma(t)}{dt} &= \frac{-(-e^{-t})}{(1+e^{-t})^2} = \frac{e^{-t}}{(1+e^{-t})} \left(\frac{1}{1+e^{-t}} \right) \\ &= \left(1 - \frac{1}{1+e^{-t}} \right) \left(\frac{1}{1+e^{-t}} \right) \\ &= (1 - \sigma(t)) (\sigma(t)) \end{aligned}$$

$$d\hat{z} = (da)(1-a)(a) = (1-a)(a)$$

$$\begin{aligned} &= \left(-\frac{y}{a} + \frac{(1-y)}{(1-a)} \right) (1-a)a = -y(1-a) + (1-y)a \\ &= -y + ya + a - ya \\ &= a - y \end{aligned}$$

$$d\hat{z} = a - y$$

$$d\omega_1 = \frac{dJ}{d\omega_1} = \frac{dJ}{dz} \frac{dz}{d\omega_1} = (dz) \frac{d(x_1\omega_1 + x_2\omega_2 + b)}{d\omega_1} = x_1 dz$$

$$d\omega_2 = \frac{dJ}{d\omega_2} = \frac{dJ}{dz} \frac{dz}{d\omega_2} = (dz) \frac{d(x_1\omega_1 + x_2\omega_2 + b)}{d\omega_2} = x_2 dz$$

$$db = \frac{dJ}{db} = \frac{dJ}{dt} \frac{dt}{db} = (dt) \frac{d(x_1\omega_1 + x_2\omega_2 + b)}{db} = dz$$

$$dt = \alpha - y$$

$$\omega_1 := \omega_1 - \alpha d\omega_1$$

$$d\omega_1 = x_1 dz$$

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

$$d\omega_2 = x_2 dz$$

$$b := b - \alpha db$$

$$db = dz$$

decior: logistic Regression with m examples:-

Cost function

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y^{(i)}) \quad (x^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^\top x^{(i)} + b) \quad d\omega_1^{(i)}, d\omega_2^{(i)}, db^{(i)}$$

$$\frac{\partial J(\omega, b)}{\partial \omega_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \ell(a^{(i)}, y^{(i)})}{\partial \omega_1^{(i)}}$$

Initialize $J=0, d\omega_1=0, d\omega_2=0, db=0$

for $i=1$ to m

$$z^{(i)} = \omega^\top x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$d\hat{z}^{(i)} = a^{(i)} - y^{(i)}$$

$$d\omega_1 += x_1^{(i)} d\hat{z}^{(i)}$$

$$d\omega_2 += x_2^{(i)} d\hat{z}^{(i)}$$

$$db += d\hat{z}^{(i)}$$

$$J / m$$

$$d\omega_1 / m$$

$$d\omega_2 / m$$

$$db / m$$

$$\omega_1 := \omega_1 - \alpha d\omega_1$$

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

$$b := b - \alpha db$$

Single
Step of GD

Vectorization: dec11

$$z = \omega^T x + b$$

$$\omega \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

Non vector:

$$z = 0$$

```
for i in range(len(nx)):
    z += omega[i] * x[i]
```

$$z += b$$

400ms

vector:

$$z = np.dot(\omega, x) + b$$

$$\omega^T x$$

1.5ms

dec12: More examples (avoid for loops, whenever possible)

$$U = A \cdot V$$

$$U = np.dot(A, V)$$

$$V = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \quad U = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$$U = np.exp(V)$$

Vectorized inside loop :-

$$J=0, d\omega = \text{np.zeros}((n_x, 1)), db = 0$$

for i in range(len(m)):

$$\hat{z}^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(\hat{z}^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$d\hat{z}^{(i)} = a^{(i)} - y^{(i)}$$

$$d\omega += x^{(i)} d\hat{z}^{(i)}$$

$$db += d\hat{z}^{(i)}$$

$$J /= m, d\omega /= m, db /= m$$

dec 13: Vectorizing logistic Regression :- (Predictions)

$$\hat{z}^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(\hat{z}^{(i)})$$

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \in \mathbb{R}^{(n_x, m)}$$

$$z = \begin{bmatrix} \hat{z}^{(1)} & \hat{z}^{(2)} & \dots & \hat{z}^{(m)} \end{bmatrix}_{(1, m)} = \underbrace{w^T x}_{(1, n_x) \times (n_x, m)} + \underbrace{\begin{bmatrix} b & b & \dots & b \end{bmatrix}}_{(1, m)}$$

$$z = \text{np.dot}(w.T, x) + b_{(1, 1)} \quad \text{broadcasting}$$

$$\text{Af } \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix} = \sigma(z)$$

dec 14: Vectorizing backward pass:

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)}$$

$$dz = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]_{(1,m)}$$

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}]$$

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

$$dz = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}]$$

Initialize $d\omega = 0$

$$\left. \begin{array}{l} d\omega += x^{(1)} dz^{(1)} \\ d\omega += x^{(2)} dz^{(2)} \\ \vdots \\ d\omega /= m \end{array} \right| \left. \begin{array}{l} db = 0 \\ db += dz^{(1)} \\ db += dz^{(2)} \\ \vdots \\ db /= m \end{array} \right| \text{looping through } m \text{ examples}$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} = \frac{1}{m} np\text{-sum}(dz) \quad (1,1)$$

$$d\omega = \frac{1}{m} X dz^T = \frac{1}{m} \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$d\omega = \frac{1}{m} \left\{ x_{(n_x,1)}^{(1)} dz^{(1)} + x_{(n_x,1)}^{(2)} dz^{(2)} + \dots + x_{(n_x,1)}^{(m)} dz^{(m)} \right\} \quad (n_x, 1)$$

Final Implementation of Logistic Regression : dec15

For i in range (iter):

$$z = \omega^T x + b$$

$$= npdot(\omega.T, X) + b$$

$$A = \sigma(z)$$

$$dz = A - Y$$

$$d\omega = \frac{1}{m} \times dz^T$$

$$db = \frac{1}{m} npsum(dz)$$

$$\omega := \omega - \alpha d\omega$$

$$b := b - \alpha db$$

dec 16: Broadcasting example cal in 100g of food

$$\begin{array}{c}
 \text{Apples} & \text{Beef} & \text{Eggs} & \text{Potatoes} \\
 \text{carbs} & 56.0 & 0.0 & 4.4 \\
 \text{Protein} & 1.2 & 104.0 & 52.0 \\
 \text{Fat} & 1.8 & 135.0 & 8.0
 \end{array}
 \quad
 \begin{array}{c}
 68.0 \\
 8.0 \\
 0.9
 \end{array}$$

'.' of cal from carbs, protein, Fat.

$$A = \begin{bmatrix} 56.0 & 0.0 & 4.4 & 68.0 \\ 1.2 & 104.0 & 52.0 & 8.0 \\ 1.8 & 135.0 & 9.0 & 0.9 \end{bmatrix} \quad \text{cal} = A \cdot \text{sum}(\text{axis}=0) = [59.0 \ 231.0 \ 155.4 \ 86.9]$$

$$\text{per} = [A/\text{cal}] \text{ (broadcasting)}$$

$$A(3,4) \quad \text{cal}(1,4)$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}_{(4,1)} + 100_{(1,1)} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}_{(4,1)} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}_{(4,1)} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}_{(4,1)}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(2,3)} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{(1,3)} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}_{(2,3)}$$

General principle:-

$$(m, n) \xrightarrow[\star]{+} (1, n) \rightarrow (m, n)$$

$$(m, n) \xrightarrow[\star]{+} (m, 1) \rightarrow (m, n)$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

$$\{1, 2, 3\} + 100 = \{101, 102, 103\}$$

dec17: Numpy Vectors:-

$\text{assert } (\mathbf{a}.shape == (s, 1))$
 don't use rank 1 array $(n,)$
 use row $[:, n]$ vector
 use column $[n, :]$ vector.
 use veshare operation.

dec18:- Cost Function:-

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

interpret $\hat{y} = P(y=1|x)$

$$\text{if } y=1: P(y|x) = \hat{y}$$

$$y=0: P(y|x) = 1 - \hat{y}$$

$$P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\text{if } y=1 P(y|x) = \hat{y}$$

$$\text{if } y=0 P(y|x) = (1-\hat{y})$$

$$\log P(y|x) = y \log \hat{y} + (1-y) \log (1-\hat{y}) = -d(y, \hat{y})$$

$\min \text{ loss} \Rightarrow \max P(y|x)$
 (iid) (identical independent dist)

cost on m examples:-

$$\log P(\text{ables in training set}) = \log \prod_{i=1}^m P(y^{(i)} | x^{(i)})$$

$$\begin{aligned} \log P_C &= \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) \\ &= - \sum_{i=1}^m d(\hat{y}^{(i)}, y^{(i)}) \end{aligned}$$

$$\text{cost } J(\mathbf{w}, b) = \underset{\text{minimize}}{\frac{1}{m} \sum_{i=1}^m d(\hat{y}^{(i)}, y^{(i)})}$$