# Binary Classification



64

64

1 (cat) vs 0 (non cat)

$y$

Blue
Green
Red

| 255 | 134 | 93 | 22 |
| 255 | 134 | 202 | 22 | 2 |
| 255 | 231 | 42 | 22 | 4 | 30 |
| 123 | 94 | 83 | 2 | 192 | 124 |
| 34 | 44 | 187 | 92 | 34 | 142 |
| 34 | 76 | 232 | 124 | 94 |
| 67 | 83 | 194 | 202 |

64

64

$$X = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}$$

$64 \times 64 \times 3 = 12288$

$n = n_x = 12288$

$X \longrightarrow y$

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix} \updownarrow n_x$$

$\longleftarrow m \longrightarrow$

$X \in \mathbb{R}^{n_x \times m}$     $X.shape = (n_x, m)$

$Y = [\, y^{(1)} \quad y^{(2)} \quad \cdots \quad , y^{(m)} \,]$

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

$Y.shape = (1, m)$

7:54 / 8:23

# Logistic Regression

Given $x$, want $\hat{y} = P(y=1 \mid x)$

$x \in \mathbb{R}^{n_x}$      $0 \leq \hat{y} \leq 1$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$.

Output $\hat{y} = \sigma(\underbrace{w^T x + b}_{z})$

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

If $z$ large   $\sigma(z) \approx \dfrac{1}{1+0} = 1$

If $z$ large negative number

$\sigma(z) = \dfrac{1}{1 + e^{-z}} \approx \dfrac{1}{1 + \text{Bignum}} \approx 0$

# Logistic Regression cost function

$\rightarrow \hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$, where $\sigma(z^{(i)}) = \dfrac{1}{1 + e^{-z^{(i)}}}$    $z^{(i)} = w^T x^{(i)} + b$

$x^{(i)}$
$y^{(i)}$   $i$-th example.
$z^{(i)}$

Given $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$.

Loss (error) function: $\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

$\mathcal{L}(\hat{y}, y) = -\left(\boxed{y \log \hat{y}} + (1-y)\log(1-\hat{y})\right) \leftarrow$

If $y=1$: $\mathcal{L}(\hat{y}, y) = -\log \hat{y}$   $\leftarrow$ Want $\log \hat{y}$ large, want $\hat{y}$ large.

If $y=0$: $\mathcal{L}(\hat{y}, y) = -\log(1-\hat{y})$   $\leftarrow$ Want $\log 1-\hat{y}$ large .... want $\hat{y}$ small

Cost function: $J(w,b) = \dfrac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$
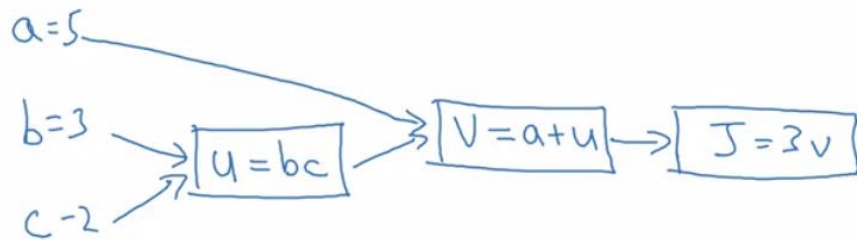
Derivative = Slope 😊

# Computation Graph

$$J(a,b,c) = 3(a + \underbrace{\underbrace{\overbrace{bc}^{u}}_{v}}_{J})$$
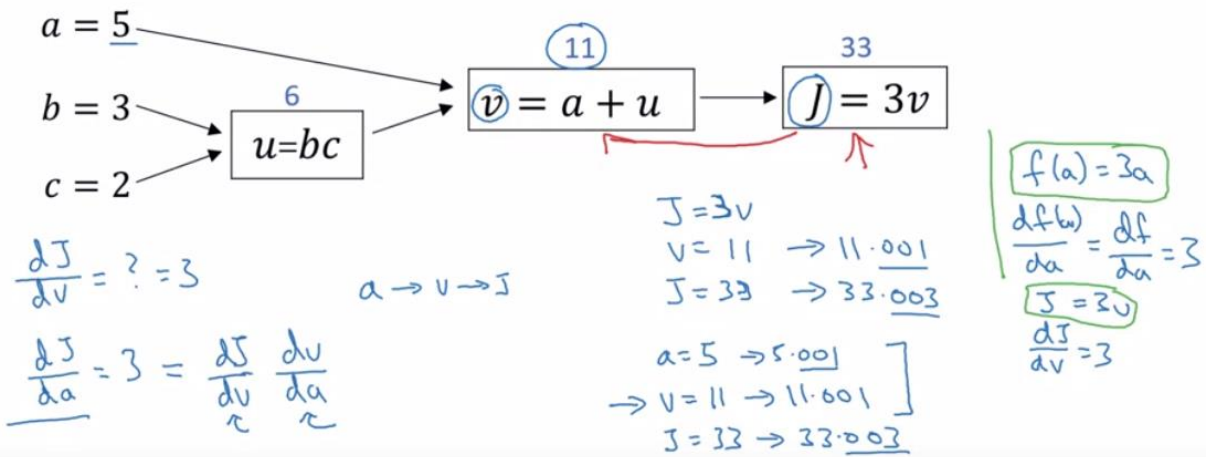
$u = bc$
$V = a + u$
$J = 3v$

$a = 5$

$b = 3$ → $u = bc$ → $V = a + u$ → $J = 3v$

$c - 2$

# Computing derivatives

$a = 5$

$b = 3$

$c = 2$

6

$u = bc$

(11)

$v = a + u$

33

$J = 3v$

$\frac{dJ}{dv} = ? = 3$

$a \to v \to J$

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

$J = 3v$
$v = 11 \to 11.001$
$J = 33 \to 33.003$

$a = 5 \to 5.001$
$\to v = 11 \to 11.001$
$J = 33 \to 33.003$

$f(a) = 3a$

$\frac{df(a)}{da} = \frac{df}{da} = 3$

$J = 3v$
$\frac{dJ}{dv} = 3$

Chain Rule of Calculus!

# Computing derivatives

$\frac{dJ}{da}$

$a = 5$
$da = 3$

$b = 3$

$c = 2$

$u=bc$
$du=3$

$6$

$11$
$v = a + u$
$dv = 3 \qquad \frac{dJ}{dJ}$

$33$
$J = 3v$

$\frac{dJ}{du} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{du}$

$\qquad 3 \qquad 1$

$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db}$

$\qquad = 2$

$u = 6 \rightarrow 6.001$
$v = 11 \rightarrow 11.001$
$J = 33 \rightarrow 33.003$

$b = 3 \rightarrow 3.001$
$u = b \cdot c = 6 \rightarrow 6.002 \qquad c = 2$

# Logistic regression derivatives

$x_1$
$w_1$
$x_2$
$w_2$
$b$

$z = w_1 x_1 + w_2 x_2 + b \quad \rightarrow \quad a = \sigma(z) \quad \rightarrow \quad \mathcal{L}(a,y)$

$dz = \frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}(a,y)}{dz}$

$"da" = \frac{d\mathcal{L}(a,y)}{da}$

$= a - y$

$a(1-a)$

$= -\frac{y}{a} + \frac{1-y}{1-a}$

$= \frac{d\mathcal{L}}{da} \cdot \frac{da}{dz}$

$\frac{d\mathcal{L}}{dw_1} = "dw_1" = x_1 \cdot dz . \qquad dw_2 = x_2 \cdot dz. \quad db = dz.$

$w_1 := w_1 - \alpha \, dw_1$
$w_2 := w_2 - \alpha \, dw_2$
$b .$

# Logistic regression on $m$ examples

$J = 0$ ; $dw_1 = 0$ ; $dw_2 = 0$ ; $db = 0$

For $i = 1$ to $m$

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

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

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

$\qquad dz^{(i)} = a^{(i)} - y^{(i)}$

$\qquad dw_1 += x_1^{(i)} dz^{(i)}$ $\quad \left.\right\} n = 2$

$\qquad dw_2 += x_2^{(i)} dz^{(i)}$

$\qquad db += dz^{(i)}$

$J /= m$

$dw_1 /= m$ ; $\quad dw_2 /= m$ ; $db /= m$.

```
250286.989866
Vectorized version:1.5027523040771484ms
250286.989866
For loop:474.29513931274414ms
```

# Logistic regression derivatives

$J = 0,$ $\boxed{dw1 = 0, \quad dw2 = 0,}$ $db = 0$ $\qquad dw = np.zeros((n\text{-}x, 1))$

→ for i = 1 to m:

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

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

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

$\qquad dz^{(i)} = a^{(i)}(1 - a^{(i)})$

for $j = 1 \ldots n_x$
$dw_j$

$\qquad \boxed{\begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \end{array}}$ $\quad n_x = 2$ $\qquad dw += x^{(i)} dz^{(i)}$

$\qquad db += dz^{(i)}$

$J = J/m,$ $\boxed{dw_1 = dw_1/m, \quad dw_2 = dw_2/m,}$ $db = db/m$

$\qquad\qquad\qquad dw \;/= m.$

---

# Vectorizing Logistic Regression

→ $z^{(1)} = \boxed{w^T x^{(1)} + b}$ $\qquad z^{(2)} = \boxed{w^T x^{(2)} + b}$ $\qquad z^{(3)} = w^T x^{(3)} + b$

→ $\boxed{a^{(1)}} = \sigma(z^{(1)})$ $\qquad \boxed{a^{(2)}} = \sigma(z^{(2)})$ $\qquad \boxed{a^{(3)}} = \sigma(z^{(3)})$

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix} \qquad \frac{(n_x, m)}{\mathbb{R}^{n_x \times m}} \qquad w^T \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$Z = \boxed{[z^{(1)} \;\; z^{(2)} \cdots z^{(m)}]} = w^T X + \underbrace{[b \;\; b \cdots b]}_{1 \times m} = [\;\boxed{w^T x^{(1)} + b} \;\; \boxed{w^T x^{(1)} + b} \cdots w^T x^{(m)} + b\;]$$

$$\quad 1 \times m$$

→ $Z = np.dot(w.T, X) + b \quad (1,1) \quad \mathbb{R}$ $\qquad$ "Broadcasting"

$A = [a^{(1)} \;\; a^{(2)} \cdots a^{(m)}] = \sigma(Z)$

$$Z = \omega^T X + b$$
$$= np.dot(\omega.T, X) + b$$
$$A = \sigma(Z)$$
$$dZ = A - Y$$
$$d\omega = \frac{1}{m} X dZ^T$$
$$db = \frac{1}{m} np.sum(dZ)$$

$$\omega := \omega - \alpha \, d\omega$$
$$b := b - \alpha \, db$$

$(m, n)$     $\frac{+}{*}$     $(1, n)$   $\rightsquigarrow$   $(m, n)$

matrix

      $(m, 1)$   $\rightsquigarrow$   $(m, n)$

$(m, 1)$     $+$     $\mathbb{R}$

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

$[1 \quad 2 \quad 3]$    $+$    $100$   $= [101 \quad 102 \quad 103]$

Matlab/Octave: $\underline{bsxfun}$

```
a = np.random.randn(5,1)
print(a)
```

```
[[-0.0967311 ]
 [-2.38617377]
 [-0.3243588 ]
 [-0.96216349]
 [ 0.54410384]]
```

```
print(a.T)
```

```
[[-0.0967311  -2.38617377 -0.3243588  -0.96216349  0.54410384]]
```

# Python/numpy vectors

```
a = np.random.randn(5)
```
$a.shape = (5,)$
"rank 1 array"
} Don't use

```
a = np.random.randn(5,1)
```
$\rightarrow a.shape = (5,1)$   column vector ✓

```
a = np.random.randn(1,5)
```
$\rightarrow a.shape = (1,5)$   row vector ✓

```
assert(a.shape == (5,1))
```

# Logistic regression cost function

$\rightarrow$ If $y = 1$: $\qquad p(y|x) = \hat{y}$

$\rightarrow$ If $y = 0$: $\qquad p(y|x) = 1 - \hat{y}$ $\Big\} \quad p(y|x)$

$$p(y|x) = \hat{y}^{y} (1-\hat{y})^{(1-y)} \quad \leftarrow$$

If $y = 1$: $\quad p(y|x) = \hat{y} \quad (1-\hat{y})^{0}$
$$\underset{=1}{\underbrace{\phantom{(1-\hat{y})}}}$$

If $y = 0$: $\quad p(y|x) = \hat{y}^{0} \quad (1-\hat{y})^{(1-y)} = 1 \times (1-\hat{y}) = 1 - \hat{y}$

$$\log p(y|x) = \log \hat{y}^{y} (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log(1-\hat{y})$$

$$= - \mathcal{L}(\hat{y}, y)$$

# Cost on $m$ examples

$$\log p(\text{labels in trainy set}) = \log \prod_{i=1}^{m} p(y^{(i)}|x^{(i)})$$

$$\log p(----) = \sum_{i=1}^{m} \log p(y^{(i)}|x^{(i)})$$
$$\underset{-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})}{\underbrace{\phantom{\log p(y^{(i)}|x^{(i)})}}}$$

Maximum likelihood estimation

$$= -\sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

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