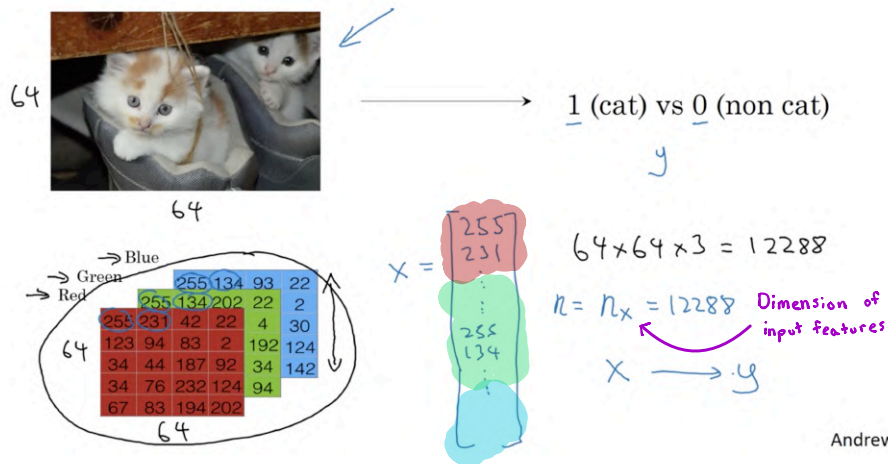


Binary Classification



Andrew Ng

Notation

(x, y) $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$
 m training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
 $M = M_{\text{train}}$ $M_{\text{test}} = \# \text{ test examples}$
 $X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$ n_x
 $X \in \mathbb{R}^{n_x \times m}$ $X.\text{shape} = (n_x, m)$
 $Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$
 $Y \in \mathbb{R}^{1 \times m}$ $Y.\text{shape} = (1, m)$

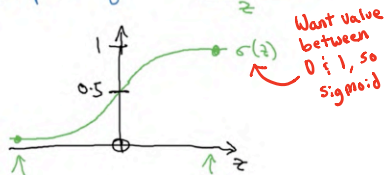
Andrew Ng

Logistic Regression

Given x , want $\hat{y} = \frac{P(y=1|x)}{0 \leq \hat{y} \leq 1}$
 $x \in \mathbb{R}^{n_x}$

Parameters: $\boxed{w} \in \mathbb{R}^{n \times 1}$, $\boxed{b} \in \mathbb{R}$.

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



$$x_0 = 1, \quad x \in \mathbb{R}^{n_x+1}$$

$$\hat{y} = \sigma(\theta^T x)$$

$$\hat{y} = \sigma(\theta^T x)$$

$$\Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \\ \vdots \\ \Theta_{N_x} \end{bmatrix} \begin{matrix} \} b \leftarrow \\ \\ \\ \} w \leftarrow \end{matrix}$$

Don't use this notation in this course

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

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

If z large negative number

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{Big num}} \approx 0$$

Andrew Ng

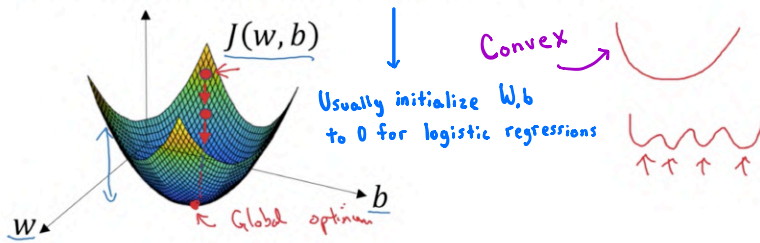
Gradient Descent

Recap: $\hat{y} = \sigma(w^T x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$ ←

Cost
Function
↓

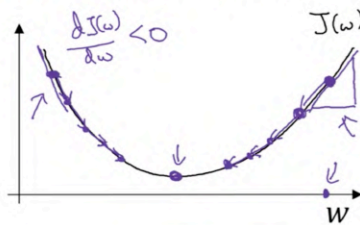
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Want to find w, b that minimize $J(w, b)$



Andrew Ng

Gradient Descent



Repeat {
 $w := w - \alpha \frac{dJ(w)}{dw}$
 }
 $w := w - \alpha dw$

learning rate
 "dw"

$\frac{dJ(w)}{dw} = ?$

$J(w, b)$

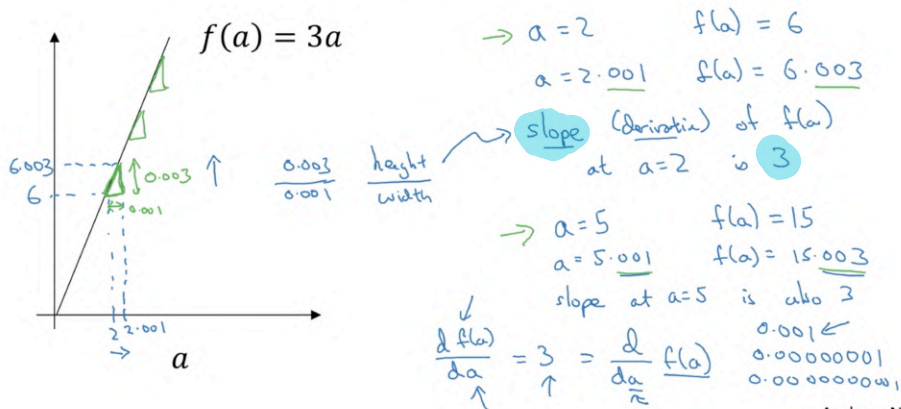
$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

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

partial derivative
 $\frac{\partial J(w, b)}{\partial w}$
 $\frac{\partial J(w, b)}{\partial b}$
 dw
 db

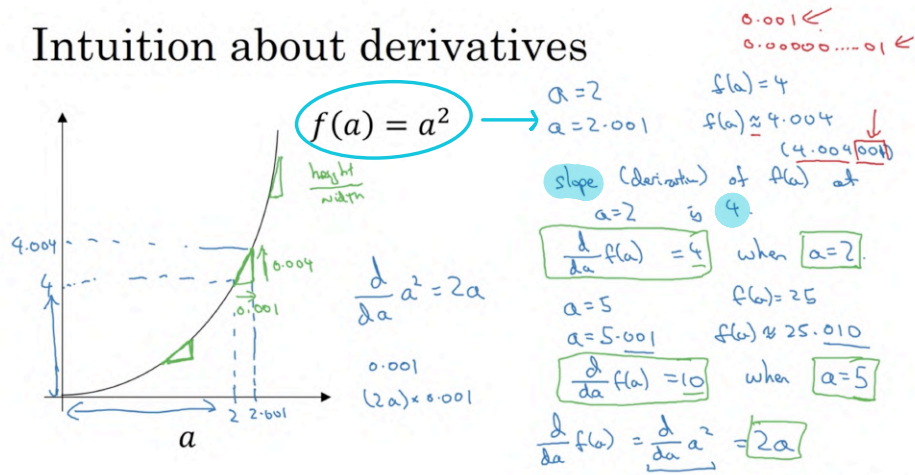
Andrew Ng

Intuition about derivatives a.k.a. slope



Andrew Ng

Intuition about derivatives



Andrew Ng

More derivative examples

$f(a) = a^2$
 $\frac{d}{da} f(a) = 2a$
 Plug in 2
 $\frac{d}{da} f(a) = \frac{2 \cdot 2}{4} = 1$

$f(a) = a^3$
 $\frac{d}{da} f(a) = 3a^2$
 $3 \cdot 2^2 = 12$

$f(a) = \log_e(a)$
 $\ln(a)$
 $\frac{d}{da} f(a) = \frac{1}{a}$
 $\frac{d}{da} f(a) = \frac{1}{2}$

$a=2$
 $a=2.001$
 $f(a)=4$
 $f(a) \approx 4.004$

$a=2$
 $a=2.001$
 $f(a)=8$
 $f(a) \approx 8.012$

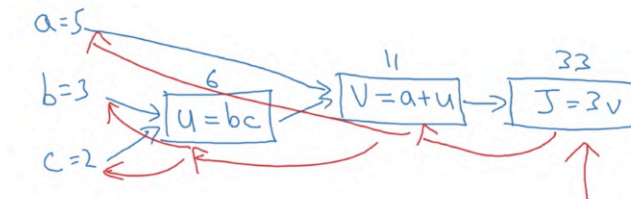
$a=2$
 $a=2.001$
 $f(a) \approx 0.69315$
 $f(a) \approx 0.69365$
 0.0005
 0.0005

Andrew Ng

Computation Graph

$$J(a, b, c) = 3(a + bc) = 3(5 + 3 \cdot 2) = 33$$

$u = bc$
 $V = a + u$
 $J = 3V$



Andrew Ng

Handwritten notes on a whiteboard explaining the chain rule for differentials.

Variables and their values:

- $a = 5$
- $b = 3$
- $c = 2$
- $u = bc$
- $v = a + u$
- $J = 3v$

Calculated differentials:

- $\frac{dJ}{dv} = 3$
- $\frac{dv}{da} = 1$
- $\frac{dJ}{da} = 3$

Summary of the chain rule:

$$\frac{d(\text{Final Output Var})}{d \text{ var}} = \frac{dJ}{d \text{ var}} \cdot \frac{d \text{ var}}{da}$$

Note: Bumping A bumps v which bumps J

$a = 5$
 $\frac{dJ}{da} \rightarrow \underline{\frac{da}{da}} = 3$
 $b = 3$
 $\frac{dJ}{db} \Rightarrow \underline{db} = 6$
 $c = 2$
 $\frac{dJ}{dc} \Rightarrow \underline{dc} = 9$

$u = bc$
 $\underline{\frac{du}{dc}} = 3$

$v = a + u$
 $\underline{\frac{dv}{du}} = 3$

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

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

$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} = 6$

$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 9$

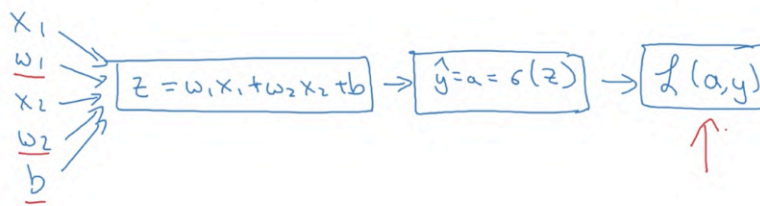
$u = 6 \rightarrow 6.001$
 $v = 11 \rightarrow 11.001$
 $J = 33 \rightarrow 33.003$

$b = 3 \rightarrow 3.001$
 $u = b \cdot c = 6 \rightarrow 6.002$
 $J = 33.006$

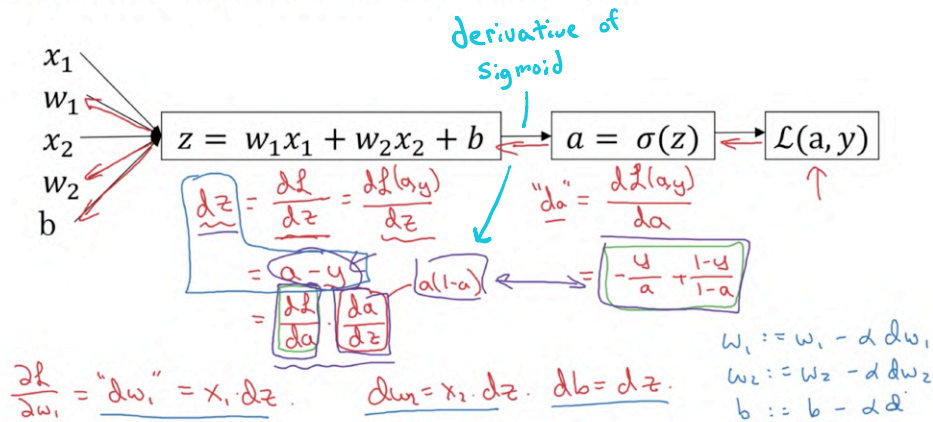
$v = 11.002$
 $J = 33v$

Andrew Ng

- $z = w^T x + b$
- $\hat{y} = a = \sigma(\underline{z})$
- $\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$



Logistic regression derivatives



Andrew Ng

Logistic regression on m examples

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

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

$(x^{(i)}, y^{(i)})$

$\underline{dw_1^{(i)}}, \underline{dw_2^{(i)}}, \underline{db^{(i)}}$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \mathcal{L}(a^{(i)}, y^{(i)})$$

$\underline{dw_1^{(i)}} - (x^{(i)}, y^{(i)})$

Andrew Ng

Logistic regression on m examples

$$J=0; \underline{dw_1}=0; \underline{dw_2}=0; \underline{db}=0$$

For $i=1$ to m

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

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

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

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

$$dw_1 += x_1^{(i)} dz^{(i)}$$

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

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

$\frac{dw_1}{dw_2}$
 $\frac{dw_1}{dw_2}$

$J/=m \leftarrow$

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

$$dw_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

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

Vectorization

Andrew Ng

Vectorizing Logistic Regression

$$\begin{aligned}
 dz^{(1)} &= a^{(1)} - y^{(1)} & dz^{(2)} &= a^{(2)} - y^{(2)} & \dots \\
 \underline{dz} &= [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}] \quad 1 \times m \\
 A &= [a^{(1)} \ \dots \ a^{(m)}] & Y &= [y^{(1)} \ \dots \ y^{(m)}] \\
 \rightarrow dz &= A - Y = [\underline{a^{(1)} - y^{(1)}} \ \underline{a^{(2)} - y^{(2)}} \ \dots]
 \end{aligned}$$

$$\begin{aligned}
 \left[\begin{array}{l} \rightarrow dw = 0 \\ dw += \frac{1}{m} dz^{(1)} \\ dw += \frac{1}{m} dz^{(2)} \\ \vdots \\ dw /= m \end{array} \right] & \quad \left[\begin{array}{l} db = 0 \\ db += dz^{(1)} \\ db += dz^{(2)} \\ \vdots \\ db += dz^{(m)} \\ db /= m \end{array} \right]
 \end{aligned}$$

$$\begin{aligned}
 db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\
 &= \frac{1}{m} \text{np.sum}(dz) \\
 dw &= \frac{1}{m} X dz^T \\
 &= \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \\ \vdots & & \vdots \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} \\
 &= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}] \quad n \times 1
 \end{aligned}$$

Andrew Ng

Implementing Logistic Regression

$$\begin{aligned}
 J &= 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0 \\
 \text{for } i &= 1 \text{ to } m: \\
 z^{(i)} &= w^T x^{(i)} + b \\
 a^{(i)} &= \sigma(z^{(i)}) \\
 J &+= -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})] \\
 dz^{(i)} &= a^{(i)} - y^{(i)} \\
 \left[\begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{array} \right] & \quad dw += x^{(i)} * dz^{(i)} \\
 J &= J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m \\
 db &= db/m
 \end{aligned}$$

$$\begin{aligned}
 \text{for } \text{iter in range}(1000): \\
 Z &= w^T X + b \\
 &= \text{np.dot}(w, X) + b \\
 A &= \sigma(Z) \\
 dz &= A - Y \\
 dw &= \frac{1}{m} X dz^T \\
 db &= \frac{1}{m} \text{np.sum}(dz) \\
 w &:= w - \alpha dw \\
 b &:= b - \alpha db
 \end{aligned}$$

Andrew Ng

Broadcasting example

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

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}_{(m,1) \rightarrow (m,n)} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

General Principle

$$\begin{array}{ccc} (m, n) & + & (1, n) \rightarrow (m, n) \\ \text{matrix} & * & \\ & / & (m, 1) \rightarrow (m, n) \end{array}$$

$$\begin{array}{ccc} (m, 1) & + & \mathbb{R} \\ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} & + & 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix} \\ [1 \ 2 \ 3] & + & 100 = [101 \ 102 \ 103] \end{array}$$

Matlab/Octave: bsxfun

Python/numpy vectors

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

`a = np.random.randn(5, 1)` → `a.shape = (5, 1)` Column vector ✓

`a = np.random.randn(1, 5)` → `a.shape = (1, 5)` Row vector ✓

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

Andrew Ng

Logistic regression cost function

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

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

If $y = 1$: $p(y|x) = \hat{y} \cdot \underbrace{(1-\hat{y})^0}_{=1}$

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

$$\begin{aligned} \uparrow \log p(y|x) &= \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -\frac{1}{\epsilon} f(\hat{y}, y) \downarrow \end{aligned}$$

Andrew Ng

Cost on m examples

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

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

maximum likelihood
estimator \nwarrow

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

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