

Core Notations (Binary Classification)

each training example is (x, y) where:

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

$$y \in \{0, 1\}$$

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

$$m = m_{\text{train}}$$

Matrix Representation

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

rows = num features

cols = num examples

⇒ OUR COLUMNS ARE training example.

$$X \in \mathbb{R}^{n_x \times m}$$

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

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

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

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

This is a different notation from 10-31S.

Logistic Regression

when output is (0 or 1). This is a classification algorithm.

given x , want $\hat{y} = P(y=1|x)$

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

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

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

$$\equiv$$

if z is large, $\sigma(z) \approx 1$

if z is small, $\sigma(z) \approx 0$

(large -ve num)

⇒ sigmoid function :

continuous function $\mathbb{R} \rightarrow (0, 1)$

Alternative Convention

⇒ feed bias term in

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

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix}$$

inferior, so

in NN, we don't use this notation.

Logistic Regression Cost Function

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b) \quad \text{we want } \hat{y}^{(i)} \approx y^{(i)}$$

Loss (error) function: (for one example)

Square error:

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 \Rightarrow \begin{aligned} &\text{non-convex optimization problem} \\ &\text{multiple local minima} \\ &\Rightarrow \text{gradient descent doesn't work well.} \end{aligned}$$

Cross-entropy loss: (aka logarithmic loss)

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

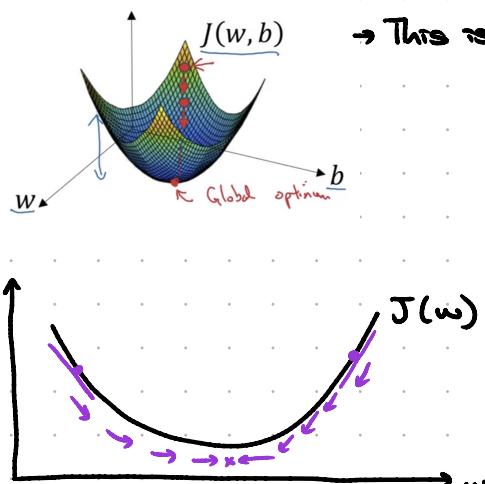
if $y=1$, then we want $\log \hat{y}$ to be large, so \hat{y} as close to 1 as possible

if $y=0$, then we want $\log(1-\hat{y})$ to be large, so \hat{y} as close to 0 as possible

Cost Function: (the cost of parameters - entire training set)

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

Gradient Descent



→ This is a good way to illustrate gradient descent.

REPEAT {

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

}

Slope of function
⇒ as (-) in front,
we move down
↓
move in opposite way

When generalized: $J(w, b)$

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

denoted
 $\frac{\partial}{\partial w}$

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

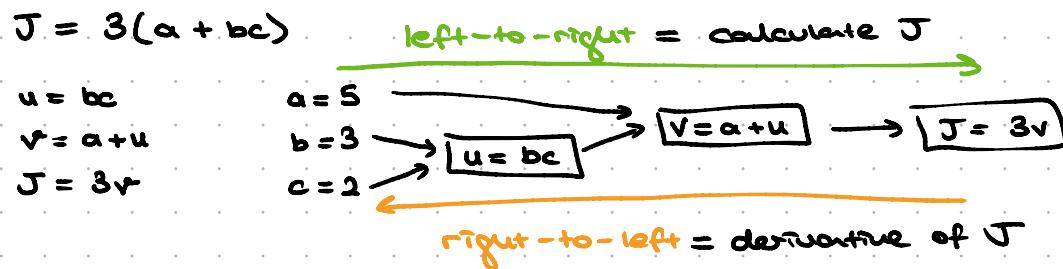
denoted
 $\frac{\partial}{\partial b}$

∂ : partial derivative symbol.

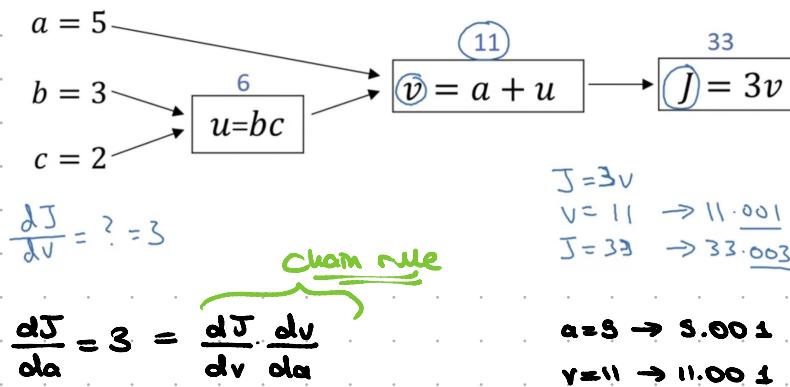
∂ : when 2 or more variables

$\frac{\partial}{\partial}$: when only 1 variable
in function we are
taking derivative of

Computation Graph



Computation Graph: Computing Derivatives



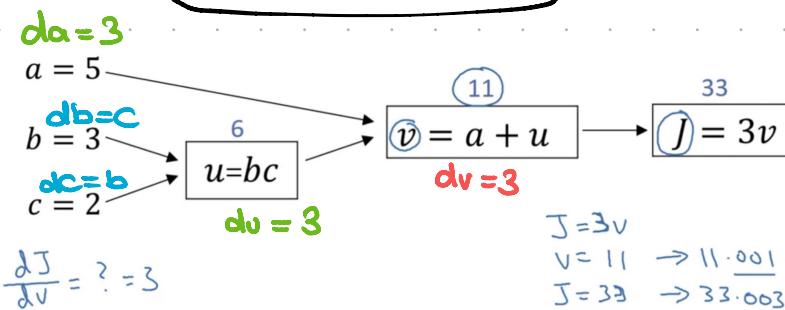
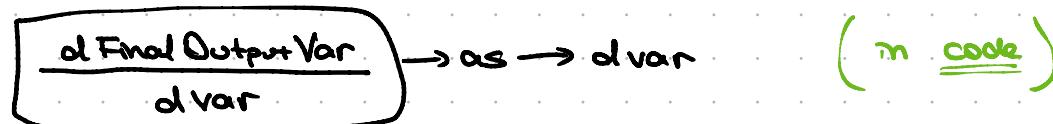
if we want to calculate derivative of final output

w.r.t. $v \Rightarrow 1$ step back prop
w.r.t. $a \Rightarrow 2$ step back prop

$$\begin{aligned} J &= 3v \\ v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned}$$

$$\begin{aligned} a &= 5 \rightarrow 5.001 \\ v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned}$$

Note: because we are always taking derivative of the final output we can name



$$\begin{aligned} J &= 3v \\ v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned}$$

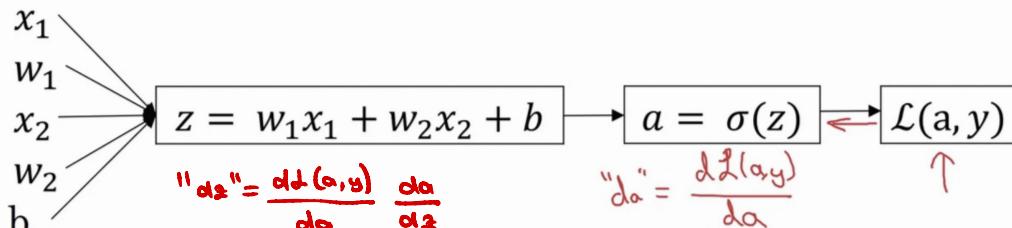
Logistic Regression Gradient Descent

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

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

$$\mathcal{L}(a, y) = - (y \log(a) + (1-y) \log(1-a))$$

Consider that we have two features.



$$\begin{aligned}
 dw_1 &= x_1 \cdot dz &= \hat{y} - y &= a - y \\
 dw_2 &= x_2 \cdot dz & &= -\frac{y}{a} + \frac{1-y}{1-a} \\
 db &= dz
 \end{aligned}$$

Extra: Derivation of dz

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

Step 1: $\frac{d\mathcal{L}}{da} \Rightarrow \mathcal{L} = - (y \log a + (1-y) \log(1-a))$

$$\frac{d\mathcal{L}}{da} = -\frac{y}{a} - \left(\frac{1-y}{1-a} (-1) \right) = -\frac{y}{a} + \frac{(1-y)}{1-a}$$

Let's give both terms in same denominator

$$\begin{aligned}
 \frac{d\mathcal{L}}{da} &= \frac{-y(1-a)}{a(1-a)} + \frac{(1-y)a}{a(1-a)} \\
 &= \frac{-y + ay + a - ay}{a(1-a)} \\
 &= \frac{a - y}{a(1-a)}
 \end{aligned}$$

Step 2: $\frac{da}{dz}$

$$\begin{aligned}
 \frac{da}{dz} &= \frac{d}{dz} \sigma(z) = \sigma(z) \times (1 - \sigma(z)) \\
 &= a(1-a) \quad [\text{as } a = \sigma(z)]
 \end{aligned}$$

Step 3: $\frac{d\mathcal{L}}{dz}$

$$\frac{d\mathcal{L}}{dz} = \underbrace{\frac{d\mathcal{L}}{da} \times \frac{da}{dz}}_{\text{you can just multiply outcomes}} = \frac{a-y}{a(1-a)} \cdot a \cdot (1-a) = a - y$$

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

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

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

Gradient Descent for M-Examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) \quad a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^\top x^{(i)} + b)$$

$$\frac{\partial}{\partial w_j} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_j} L(a^{(i)}, y^{(i)})}_{\alpha(w^{(i)})}$$

initialize $J=0$, $dw_1=0$, $dw_2=0$, $db=0$

For $i=1$ to m : → For loop over training examples

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

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

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

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

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

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

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

$$J := m$$

$$dw_1 / m, dw_2 / m, db / m$$

For loop
over features

} assuming we have
2 features
 $k=2$

The left side is ONE step of gradient descent

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

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

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

To get rid of For loops
↳ Vectorization