

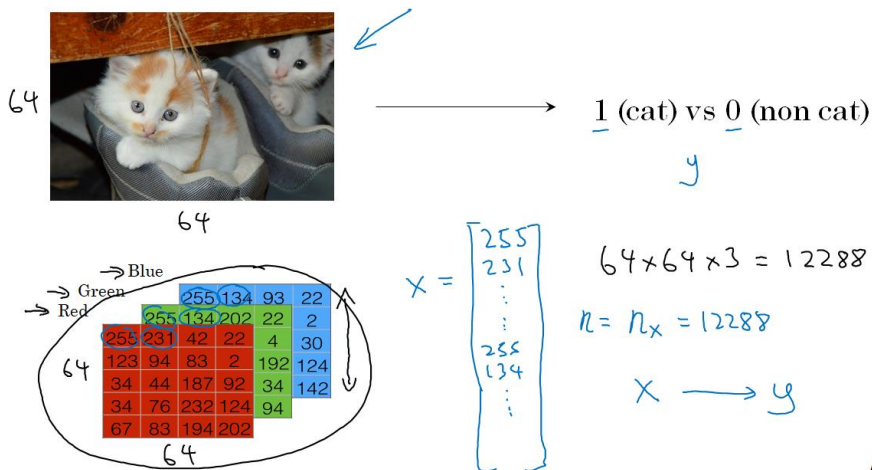


deeplearning.ai

# Basics of Neural Network Programming

## Binary Classification

### Binary Classification



Andrew Ng

### Notation

$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

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

$M = M_{\text{train}} \quad M_{\text{test}} = \# \text{test examples.}$

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \quad \begin{matrix} \uparrow \\ n_x \\ \downarrow \end{matrix}$$

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

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

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

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

Andrew Ng

## Binary Classification

In a binary classification problem, the result is a discrete value output.

For example    - account hacked (1) or compromised (0)  
                  - a tumor malign (1) or benign (0)

Example: Cat vs Non-Cat

The goal is to train a classifier that the input is an image represented by a feature vector,  $x$ , and predicts whether the corresponding label  $y$  is 1 or 0. In this case, whether this is a cat image (1) or a non-cat image (0).



An image is store in the computer in three separate matrices corresponding to the Red, Green, and Blue color channels of the image. The three matrices have the same size as the image, for example, the resolution of the cat image is 64 pixels X 64 pixels, the three matrices (RGB) are 64 X 64 each.

The value in a cell represents the pixel intensity which will be used to create a feature vector of  $n$ -dimension. In pattern recognition and machine learning, a feature vector represents an object, in this case, a cat or no cat.

To create a feature vector,  $x$ , the pixel intensity values will be “unroll” or “reshape” for each color. The dimension of the input feature vector  $x$  is  $n_x = 64 \times 64 \times 3 = 12\,288$ .

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix} \begin{array}{l} \text{red} \\ \text{green} \\ \text{blue} \end{array}$$



deeplearning.ai

# Basics of Neural Network Programming

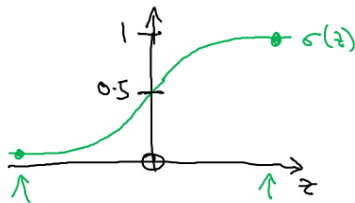
## Logistic Regression

### Logistic Regression

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

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

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



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

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

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \begin{matrix} \} b \leftarrow \\ \} w \leftarrow \end{matrix}$$

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

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

If  $z$  large negative number

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

Andrew Ng



deeplearning.ai

# Basics of Neural Network Programming

## Logistic Regression cost function

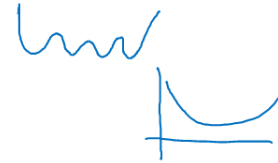
# Logistic Regression cost function

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

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

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

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



**Cost** function:  $\mathcal{L}(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$  ←

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

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

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

Andrew Ng

## Logistic Regression: Cost Function

To train the parameters  $w$  and  $b$ , we need to define a cost function.

Recap:

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

$x^{(i)}$  the  $i$ -th training example

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

Loss (error) function:

The loss function measures the discrepancy between the prediction ( $\hat{y}^{(i)}$ ) and the desired output ( $y^{(i)}$ ). In other words, the loss function computes the error for a single training example.

$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2}(\hat{y}^{(i)} - y^{(i)})^2$

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

- If  $y^{(i)} = 1$ :  $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$  where  $\log(\hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 1
- If  $y^{(i)} = 0$ :  $L(\hat{y}^{(i)}, y^{(i)}) = -\log(1 - \hat{y}^{(i)})$  where  $\log(1 - \hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 0

Cost function

The cost function is the average of the loss function of the entire training set. We are going to find the parameters  $w$  and  $b$  that minimize the overall cost function.

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



deeplearning.ai

## Basics of Neural Network Programming

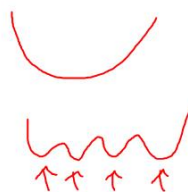
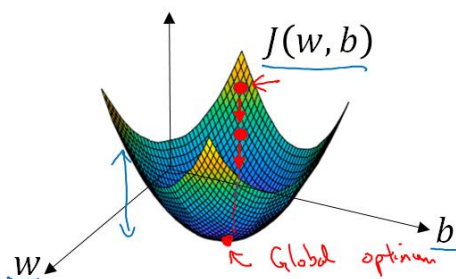
### Gradient Descent

## Gradient Descent

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

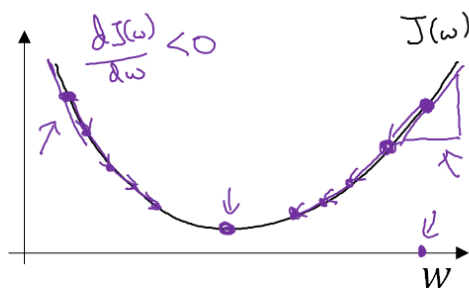
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

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



Andrew Ng

## Gradient Descent



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

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

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

"partial derivative"  
 $\frac{\partial}{\partial w}$   
 $\frac{\partial}{\partial b}$   
 $dw$   
 $db$

Andrew Ng

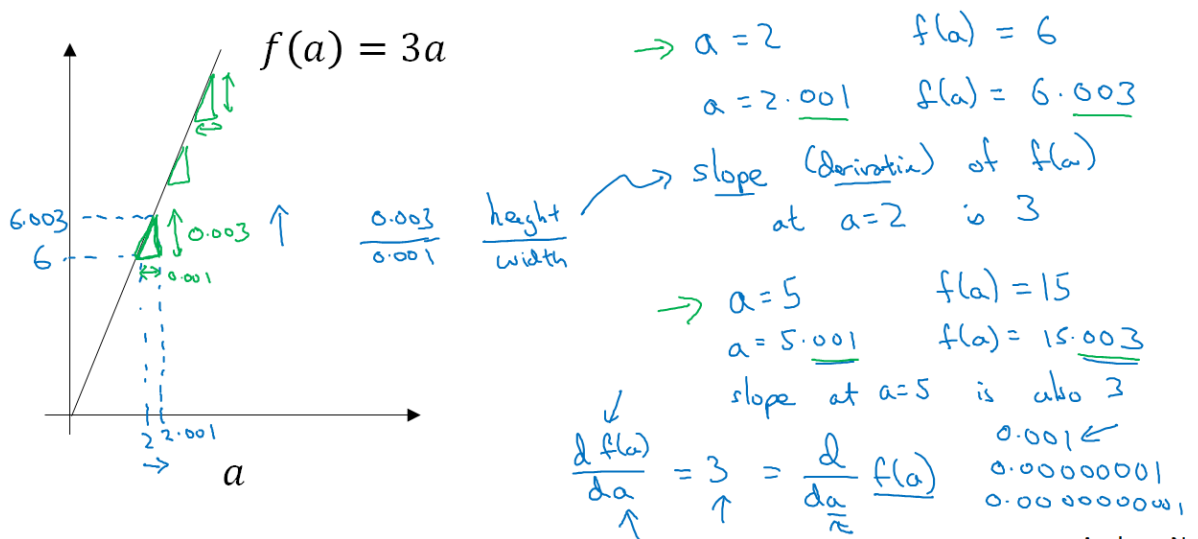


deeplearning.ai

# Basics of Neural Network Programming

## Derivatives

### Intuition about derivatives



Andrew Ng

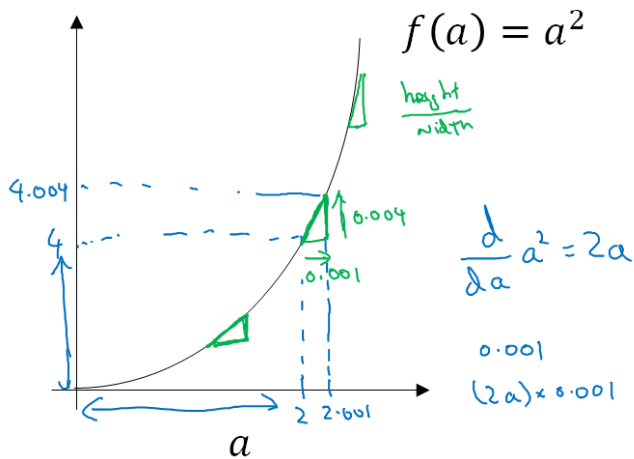


deeplearning.ai

# Basics of Neural Network Programming

## More derivatives examples

# Intuition about derivatives



$a = 2$        $f(a) = 4$   
 $a = 2.001$        $f(a) \approx 4.004$   
 $(4.004 - 4) / (2.001 - 2) = 0.004 / 0.001 = 4$   
 slope (derivative) of  $f(a)$  at  $a = 2$  is  $4$ .  
 $\frac{d}{da} f(a) = 4$  when  $a = 2$ .  
 $a = 5$        $f(a) = 25$   
 $a = 5.001$        $f(a) \approx 25.010$   
 $\frac{d}{da} f(a) = 10$  when  $a = 5$ .  
 $\frac{d}{da} f(a) = \frac{d}{da} a^2 = 2a$

Andrew Ng

## More derivative examples

$f(a) = a^2$        $\frac{d}{da} f(a) = \frac{2a}{4}$

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

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

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

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

$a = 2$        $f(a) \approx 0.69315$   
 $a = 2.001$        $f(a) \approx 0.69365$   
 $(0.69365 - 0.69315) / (2.001 - 2) = 0.0005 / 0.001 = 0.5$

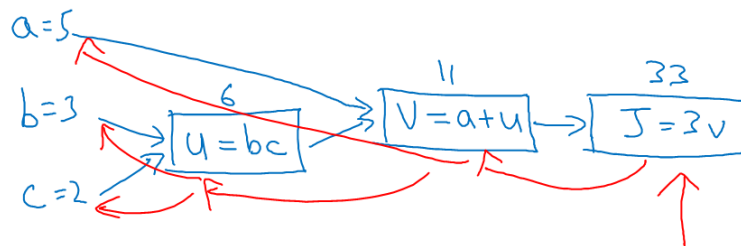
Andrew Ng



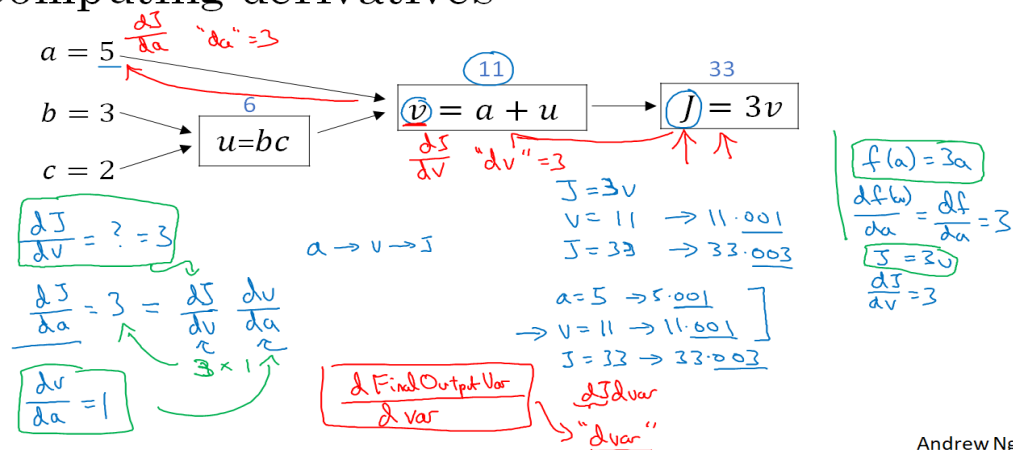
# Computation Graph

$$J(a, b, c) = 3(\underbrace{a + \underbrace{b, c}_u}_v) = 3(5 + 3 \cdot 2) = 33$$

$$\begin{aligned} u &= bc \\ v &= at + u \\ j &= 3v \end{aligned}$$

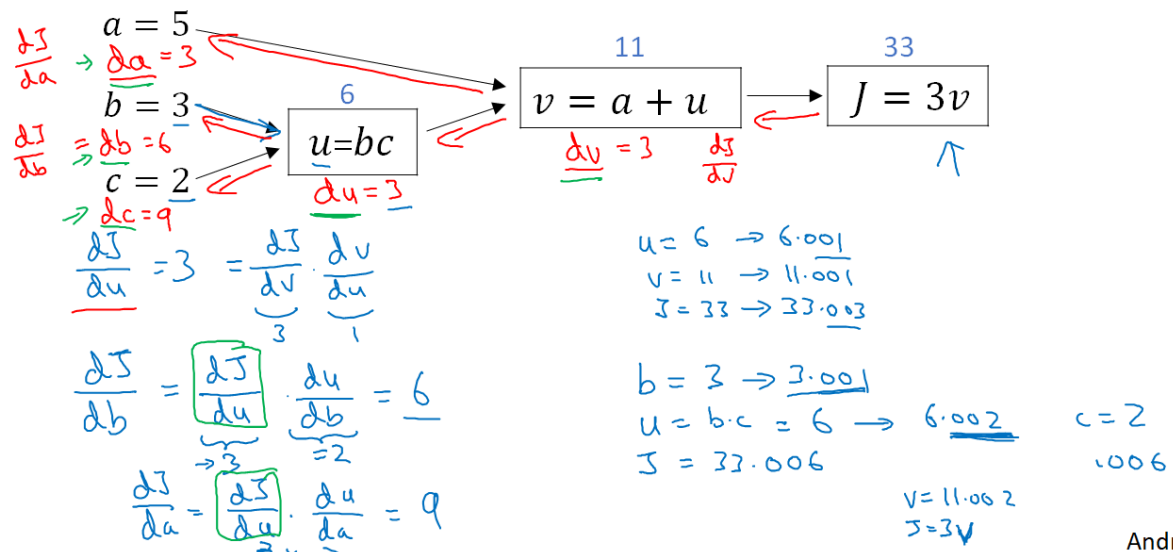


## Computing derivatives





# Computing derivatives



Andrew Ng



Basics of Neural  
Network Programming

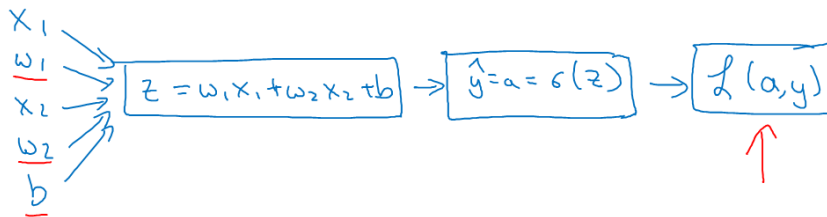
Logistic Regression  
Gradient descent

## Logistic regression recap

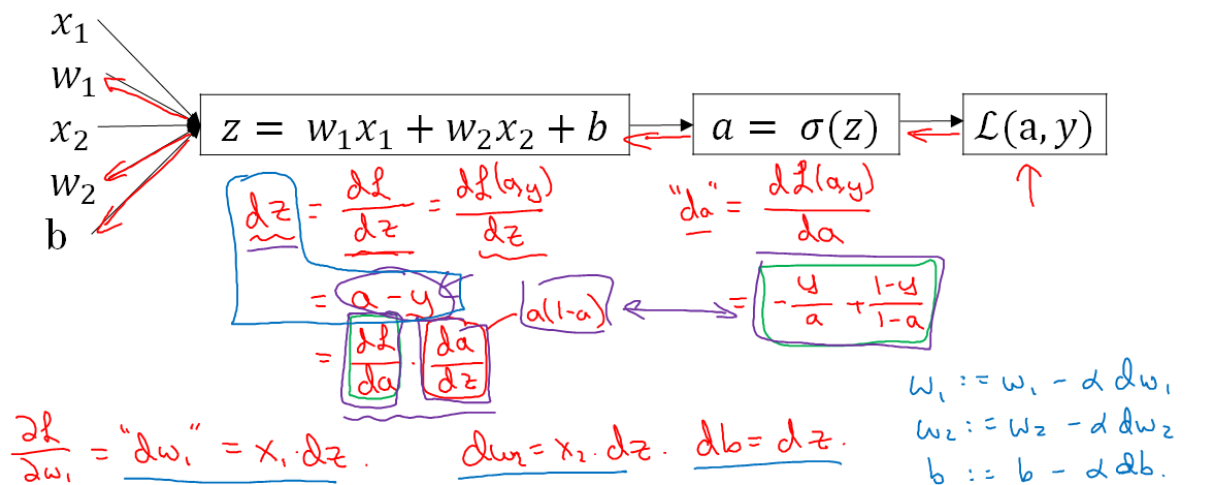
$$\rightarrow z = w^T x + b$$

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

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



# Logistic regression derivatives



Andrew Ng



deeplearning.ai

## Basics of Neural Network Programming

### Gradient descent on $m$ examples

## Logistic regression on $m$ examples

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

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

Derivatives:  $\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial b}$

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

Derivatives:  $\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial b}$

# Logistic regression on $m$ examples

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

→ For  $i=1$  to  $m$

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

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

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

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

$$\begin{array}{l} \uparrow \\ dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \\ \downarrow \end{array} \quad \begin{array}{l} \uparrow \\ n=2 \\ \downarrow \end{array}$$

$J /= m \leftarrow$

$$\begin{array}{ccc} dw_1 /= m & ; & dw_2 /= m; db /= m. \leftarrow \\ \uparrow & & \uparrow \quad \uparrow \end{array}$$

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

$$w_1 := w_1 - \alpha \underline{dw_1}$$

$$w_2 := w_2 - \alpha \underline{dw_2}$$

$$b := b - \alpha \underline{db}$$

Vectorization