

## Binary Classification

### Logistic Regression

Given  $X$ , want  $\hat{y} = P(y=1|X)$   
 output:  $\hat{y} = \sigma(w^T x + b)$   
 $\sigma(z) = \frac{1}{1+e^{-z}}$

### Cost function

Given  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

want  $\hat{y}^{(i)} \approx y^{(i)}$

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

### Gradient Descent

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

### Derivatives

### Computation Graph

example:

$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$

$$\begin{array}{l}
 a = 5 \\
 b = 3 \\
 c = 2
 \end{array}
 \rightarrow
 \begin{array}{l}
 V = a + u \\
 u = bc
 \end{array}
 \rightarrow
 J = 3V$$

Derivatives with a computation Graph.

$$\begin{array}{l}
 a = 5 \\
 b = 3 \\
 c = 2
 \end{array}
 \rightarrow
 \begin{array}{l}
 V = a + u \\
 u = bc
 \end{array}
 \rightarrow
 J = 3V$$

chain rule:

$$\begin{aligned}
 &\text{change } a \rightarrow \text{change } V \rightarrow \text{change } J \\
 \frac{\partial J}{\partial a} &= \frac{\partial J}{\partial V} \frac{\partial V}{\partial a} = 3
 \end{aligned}$$

$$\begin{array}{l}
 a = 5 \\
 b = 3 \\
 c = 2
 \end{array}
 \rightarrow
 \begin{array}{l}
 V = a + u \\
 u = bc
 \end{array}
 \rightarrow
 J = 3V$$

$\frac{\partial J}{\partial V}$        $\frac{\partial V}{\partial a}$

Logistic regression gradient descent

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

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

$$J = -\frac{1}{m} \sum (y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

$$\frac{\partial J}{\partial a} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial a} = \left( -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \cdot \frac{-1}{e^{-z} \cdot (-1)}$$

$$\begin{aligned}
 \partial z / \partial w &= -\frac{y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \\
 &= -y(1-\hat{y}) + (1-y)\hat{y} \\
 &= \hat{y} - y
 \end{aligned}$$

$$dz = \frac{\partial z}{\partial w} dw + \frac{\partial z}{\partial b} db = (\hat{y} - y) \times dw + db$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w} = (\hat{y} - y) \times$$

$$\begin{aligned}
 w &:= w - \alpha \frac{\partial L}{\partial w} \\
 &= w - \alpha (\hat{y} - y) \times \\
 b &:= b - \alpha \frac{\partial L}{\partial b}
 \end{aligned}$$

## Gradient descent examples

m examples

for  $i = 1 \dots m$ :

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

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

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

$$"dz^{(i)}" = \hat{y}^{(i)} - y^{(i)}$$

$$"dw_1" + = x_1^{(i)} \cdot "dz^{(i)}" = x_1^{(i)} (\hat{y}^{(i)} - y^{(i)})$$

$$"dw_2" + = x_2^{(i)} \cdot "dz^{(i)}" = x_2^{(i)} (\hat{y}^{(i)} - y)$$

$$"db" + = 1 \cdot "dz^{(i)}" = \hat{y}^{(i)} - y$$

$$J = J/m \quad "dw_1" = "dw_1" / m \quad "dw_2" = "dw_2" / m \quad "db" = \frac{"db"}{m}$$

$$w_1 := w_1 - \alpha "dw_1" = w_1 - \frac{\alpha}{m} \sum_i X_1^{(i)} (\hat{y}^{(i)} - y^{(i)})$$

$$w_2 := w_2 - \alpha "dw_2" = w_2 - \frac{\alpha}{m} \sum_i X_2^{(i)} (\hat{y}^{(i)} - y^{(i)})$$

$$b := b - \alpha "db" = b - \frac{\alpha}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})$$

## Vectorization

$\left. \begin{matrix} \text{GPU} \\ \text{CPU} \end{matrix} \right\}$  SIMD  $\Rightarrow$  "Single instruction multiple data"

## Neural network programming guidelines

Whenever possible, avoid explicit for loops

$$u = A v \Rightarrow u = \text{np.dot}(A, v)$$

$$u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix} \Rightarrow u = \text{np.exp}(v)$$

## Vectorize logistic regression

$$X = \begin{bmatrix} 1 & 1 & 1 \\ X_1^{(1)} & X_1^{(2)} & \dots & X_1^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ X_n^{(1)} & X_n^{(2)} & \dots & X_n^{(m)} \end{bmatrix}_{n \times m} \quad \begin{matrix} m: \text{trading examples} \\ n: \text{features} \end{matrix}$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}_{n \times 1} \quad b = [b, b, \dots, b]$$

$$[z^{(1)}, z^{(2)}, \dots, z^{(m)}] = w^T X + b = Z$$

$$Z = \text{np.(w.T, X)} + b \quad \text{broadcasting}$$

logistic regression prediction:

$$\hat{Y} = \begin{bmatrix} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(m)} \end{bmatrix} = \frac{1}{1 + e^{-w^T x + b}}$$

Vectorize Gradient Descent

$$dz = \hat{Y} - Y$$

$$db = \frac{1}{m} \text{np.sum}(dz)$$

$$dw = \frac{1}{m} X(dz)^T$$

$$\textcircled{1} z = w^T x + b = \text{np.dot}(w.T, x) + b$$

$$\textcircled{2} \hat{Y} = \sigma(z)$$

$$\textcircled{3} dz = \hat{Y} - Y$$

$$\textcircled{4} dw = \frac{1}{m} X(dz)^T$$

$$\textcircled{5} db = \frac{1}{m} \text{np.sum}(dz)$$

$$w := w - \alpha dw$$

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

Broadcasting In Python

example

	Apple	Beef	Eggs	Potatoes
Carb.				
Protein				
Fat				

---

question: calculate % of carbons from carb, Protein, Fat.  
can you do that w/o for loop?

Broadcasting Example

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}_{1 \times n} \xrightarrow{(m \times n)} = \begin{bmatrix} 101 & 102 & 103 \\ 104 & 105 & 106 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}_{m \times 1} \xrightarrow{m \times n} = \begin{bmatrix} 101 & 102 & 103 \\ 201 & 202 & 203 \end{bmatrix}$$

General Principle

$$(m, n) \xrightarrow{\quad} (1, n) \rightarrow (m, n)$$
$$\xrightarrow{\quad} (n, 1) \rightarrow (m, n)$$

Python vector s denotes

$a = np.random.rand(5) \rightarrow \text{shape } (5, ) \rightarrow \text{rank 1 array}$

$a = \dots \dots \dots (5, 1) \rightarrow \text{shape } (5, 1) \rightarrow \text{col. vector}$

$a = \dots \dots \dots (1, 5) \rightarrow \text{shape } (1, 5) \rightarrow \text{row vector}$

Explanation cost function of logistic regression

$$\hat{y} = \frac{1}{1+e^{-v^T x}}$$

if  $y=1 \Rightarrow P(y|X) = \hat{y}$   $\Rightarrow \text{prob.}$

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

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

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

then minimize the cost function  $\mathcal{L}$  actually maximize  
the  $\log(P(y|x))$  func

Logistic regression is maximum likelihood estimation

$$\begin{aligned} P(\text{labels in target set}) &= \log \prod_{i=1}^m P(y^{(i)}|x^{(i)}) \\ \log P(\dots) &= \sum_{i=1}^m \log (P(y^{(i)}|x^{(i)})) = -\sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \\ \text{cost } J(w, b) &= \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \end{aligned}$$

minimize  $J \rightarrow \max P$   
 $\rightarrow$  maximum likelihood optimization