

Warming Up - Logistic Regression

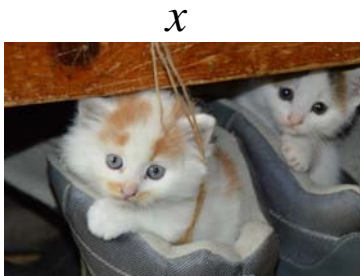
Most of this material is from Prof. Andrew Ng's slides.

Binary Classification



—————→ 1 (cat) vs 0 (non cat)

Binary Classification


$$\xrightarrow{y} 1 \text{ (cat) vs } 0 \text{ (non cat)}$$
[illegible]

input feature vector:

$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \end{bmatrix}$$

if $64 \times 64 \times 3 = 12288$
dimension of $x = n_x = 12288$

Notation

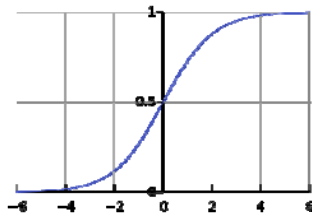
- single training example
 $(x, y) \quad x \in \mathbb{R}^{n_x}, \quad y \in \{0, 1\}$
- m training examples
 $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- more compact notation using matrix

$$\mathbf{X} = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \quad \begin{matrix} \updownarrow \\ n_x \end{matrix} \quad \mathbf{X} \in \mathbb{R}^{n_x \times m}$$

$$Y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}] \quad Y \in \mathbb{R}^{1 \times m}$$

Logistic Regression

- Given $x \in \mathbb{R}^{n_x}$, we want to estimate $\hat{y} = P(y = 1|x)$
 - We want \hat{y} to be probability: $0 \leq \hat{y} \leq 1$
- (Model) Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$
- (Model) Output: $\hat{y} = \sigma(w^T x + b)$
 - Sigmoid function $\sigma(z)$



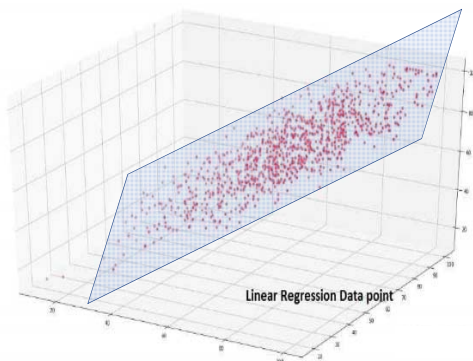
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

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

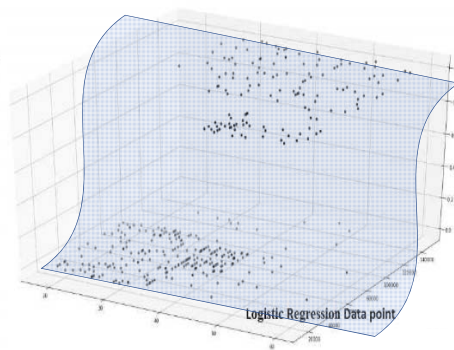
If z is large negative number: $\sigma(z) \approx \frac{1}{1 + \infty} \approx 0$

Linear Regression vs Logistic Regression

$$y = w_1 x_1 + w_2 x_2 + b$$



$$y = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + b)}}$$



Logistic Regression

- Given $x \in \mathbb{R}^{n_x}$, we want to estimate $\hat{y} = P(y|x)$
 - We want \hat{y} to be probability: $0 \leq \hat{y} \leq 1$
- (Model) Parameters: $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$
- (Model) Output: $\hat{y} = \sigma(w^T x + b)$
- The goal of logistic regression
 - try to learn the parameters w and b so that \hat{y} becomes a good estimate of the probability of y

Logistic Regression

- Given $x \in \mathbb{R}^{n_x}$, we want to estimate $\hat{y} = P(y|x)$
 - We want \hat{y} to be probability: $0 \leq \hat{y} \leq 1$
- (Model) Parameters: $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$
- (Model) Output: $\hat{y} = \sigma(w^T x + b)$
- The goal of logistic regression
- Notational convention
 - $x_0 = 1, x \in \mathbb{R}^{n_x+1}$

$$\begin{array}{lcl}
 - \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} & \begin{array}{l} \xrightarrow{\text{blue}} b \\ \text{blue bracket} \end{array} & \\
 & \text{blue bracket} & w
 \end{array}
 \qquad \hat{y} = \sigma(\theta^T x)$$

Logistic Regression Cost Function

- Model: $\hat{y} = \sigma(w^T x + b)$, where $\sigma(z) = \frac{1}{1+e^{-z}}$
- Given training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, we want $\hat{y}^{(i)} \approx y^{(i)}$

Logistic Regression Cost Function

- Model:
 $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$, where $\sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$ and $z^{(i)} = w^T x^{(i)} + b$
- Given training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, we want $\hat{y}^{(i)} \approx y^{(i)}$

Logistic Regression Cost Function

- Model:

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

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

- Loss (error) function:

$$- \mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2 \rightarrow \text{non-convex and multiple local optima}$$

Logistic Regression Cost Function

- Model:

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

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

- Loss (error) function:

$$- \cancel{\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2} \rightarrow \text{non-convex and multiple local optima}$$

$$- \mathcal{L}(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y})) : \text{cross-entropy loss}$$

$$\text{If } y \rightarrow 1 : \mathcal{L}(\hat{y}, y) = -y \log \hat{y} \rightarrow \text{want } \log \hat{y} \text{ large} \rightarrow \text{want } \hat{y} \text{ large} \rightarrow \text{want } \hat{y} \approx 1$$

$$\text{If } y \rightarrow 0 : \mathcal{L}(\hat{y}, y) = -\log(1 - \hat{y}) \rightarrow \text{want } \log(1 - \hat{y}) \text{ large} \rightarrow \text{want } \hat{y} \text{ small} \rightarrow \text{want } \hat{y} \approx 0$$

Logistic Regression Cost Function

- Model:

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

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

- Loss (error) function:

– ~~$\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$~~ → non-convex and multiple local optima

– $\mathcal{L}(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$: cross-entropy loss

If $y = 1$: $\mathcal{L}(\hat{y}, y) = -y \log \hat{y} \rightarrow \text{want } \log \hat{y} \text{ large} \rightarrow \text{want } \hat{y} \text{ large} \rightarrow \text{want } \hat{y} \approx 1$

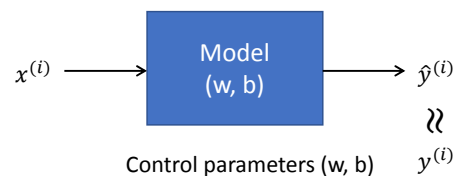
If $y = 0$: $\mathcal{L}(\hat{y}, y) = -\log(1 - \hat{y}) \rightarrow \text{want } \log(1 - \hat{y}) \text{ large} \rightarrow \text{want } \hat{y} \text{ small} \rightarrow \text{want } \hat{y} \approx 0$

- Cost function:

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

Recap

- Logistic regression model
- Loss function (cross-entropy loss)
 - measures how well your parameters w and b are doing on a single training example
- Cost function
 - measures how well your parameters w and b are doing on your entire training set
- Now let's investigate how to use the gradient descent algorithm to train the parameters w and b on your training set



Gradient Descent

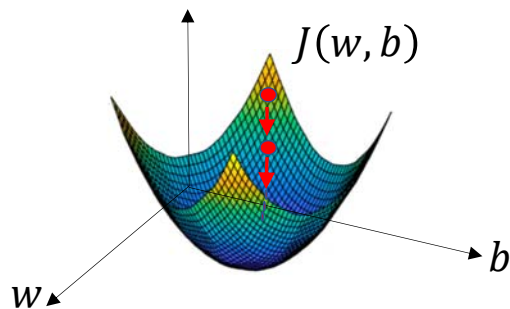
- Model: $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$, where $\sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$ and $z^{(i)} = w^T x^{(i)} + b$
- Cost function:
$$-J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

function of parameters (w, b)
- We want to find w, b that minimize J(w,b)

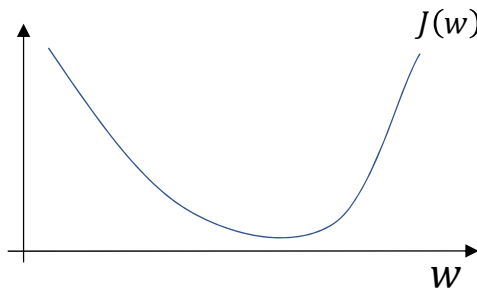
Gradient Descent

- Model: $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$, where $\sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$ and $z^{(i)} = w^T x^{(i)} + b$
- Cost function:
$$-J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

function of parameters (w, b)
- We want to find w, b that minimize J(w,b)



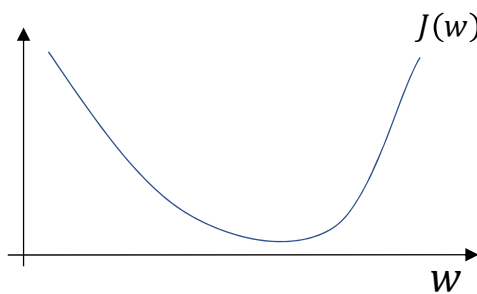
Gradient Descent



Repeat {
 $w \leftarrow w - \alpha \underbrace{\frac{dJ(w)}{dw}}_{\text{derivative}}$
}

learning rate

Gradient Descent



Repeat {
 $w \leftarrow w - \alpha \underbrace{\frac{dJ(w)}{dw}}_{\text{derivative}}$
}

learning rate

Repeat {
 $w \leftarrow w - \alpha \frac{dJ(w,b)}{dw}$
 $b \leftarrow b - \alpha \frac{dJ(w,b)}{db}$
}

Computation Graph

- Computation of neural network
 - Forward propagation: we compute the output of the neural network
 - Backward propagation: we compute gradients (derivatives)
- Computation graph explains why it is organized in this way

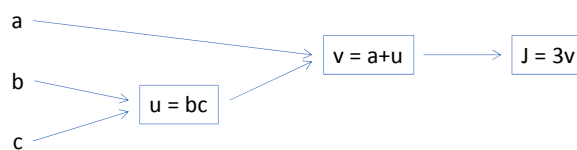
Computation Graph

- We're trying to compute a function J of three variables
- $J(a,b,c) = 3(a+bc)$

– Computing this function has three distinct steps

$$u = bc \rightarrow v = a+u \rightarrow J = 3v$$

– Computation graph:

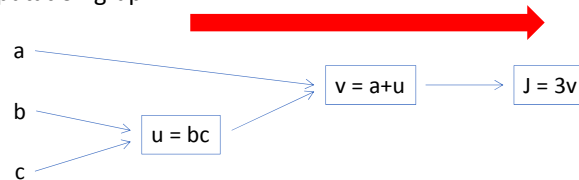


Computation Graph

- We're trying to compute a function J of three variables
- $J(a,b,c) = 3(a+bc)$

– Computing this function has three distinct steps
 $u = bc \rightarrow v = a+u \rightarrow J = 3v$

– Computation graph: Forward pass to compute cost function

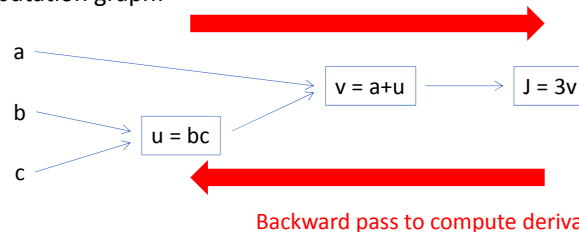


Computation Graph

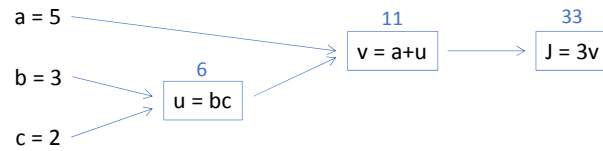
- We're trying to compute a function J of three variables
- $J(a,b,c) = 3(a+bc)$

– Computing this function has three distinct steps
 $u = bc \rightarrow v = a+u \rightarrow J = 3v$

– Computation graph: Forward pass to compute cost function

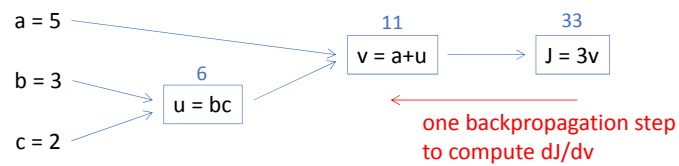


Derivatives with a Computation Graph



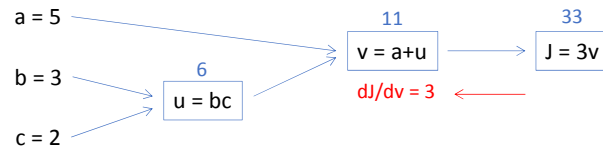
- $dJ/dv = ?$

Derivatives with a Computation Graph



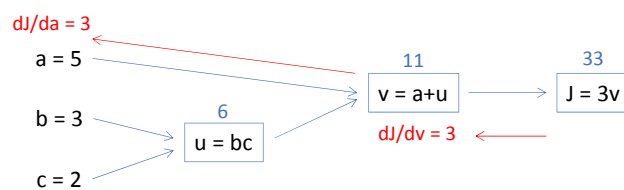
- $dJ/dv = 3$

Derivatives with a Computation Graph



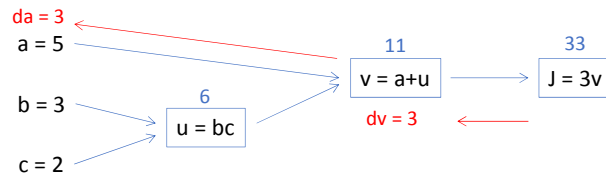
- $dJ/dv = 3$
- $dJ/da = ?$

Derivatives with a Computation Graph



- $dJ/dv = 3$
- $dJ/da = dJ/dv \cdot dv/da = 3 \cdot 1 = 3$ ("chain rule" in calculus)

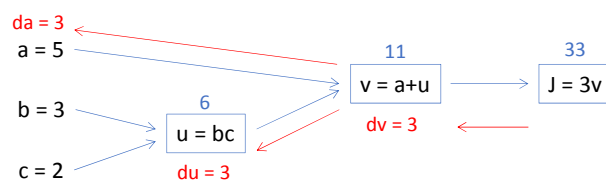
Derivatives with a Computation Graph



- $dJ/dv = 3$
- $dJ/da = dJ/dv \cdot dv/da = 3 \cdot 1 = 3$ ("chain rule" in calculus)

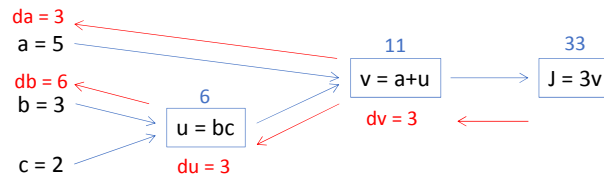
NOTE: $dJ/d(\text{var})$ is usually simplified into $d(\text{var})$ in implementing code

Derivatives with a Computation Graph



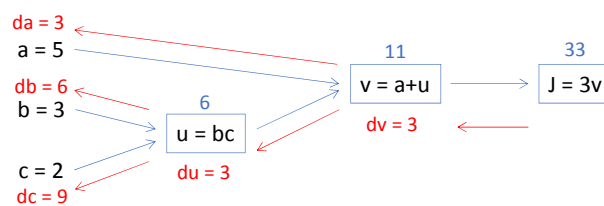
- $dJ/dv = 3$
- $dJ/da = dJ/dv \cdot dv/da = 3 \cdot 1 = 3$
- $dJ/du = dJ/dv \cdot dv/du = 3 \cdot 1 = 3$

Derivatives with a Computation Graph



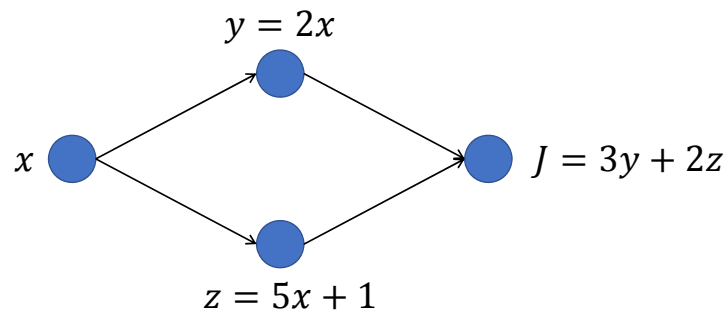
- $dJ/dv = 3$
- $dJ/da = dJ/dv \cdot dv/da = 3 \cdot 1 = 3$
- $dJ/du = dJ/dv \cdot dv/du = 3 \cdot 1 = 3$
- $dJ/db = dJ/dv \cdot dv/du \cdot du/db = 3 \cdot 1 \cdot c = 3 \cdot 1 \cdot 2 = 6$

Derivatives with a Computation Graph



- $dJ/dv = 3$
- $dJ/da = dJ/dv \cdot dv/da = 3 \cdot 1 = 3$
- $dJ/du = dJ/dv \cdot dv/du = 3 \cdot 1 = 3$
- $dJ/db = dJ/dv \cdot dv/du \cdot du/db = 3 \cdot 1 \cdot c = 3 \cdot 1 \cdot 2 = 6$
- $dJ/dc = dJ/dv \cdot dv/du \cdot du/dc = 3 \cdot 1 \cdot b = 3 \cdot 1 \cdot 3 = 9$

Exercise



$$\frac{\partial J}{\partial x} ?$$

Recap

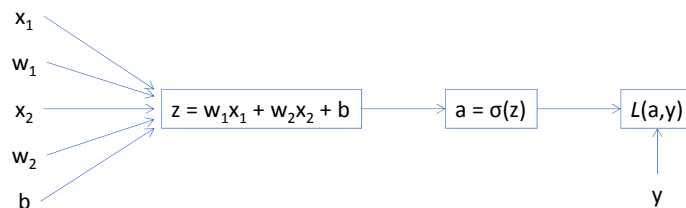
- Logistic regression model
- Loss function (cross-entropy loss)
 - measures how well your parameters w and b are doing on a single training example
- Cost function
 - measures how well your parameters w and b are doing on your entire training set
- Gradient descent
 - how to minimize the cost function by using gradient (derivatives)
 - gradient: how does a variable affect the value of the cost?
- Computing gradient
 - backpropagation with respect to computation graph

Logistic Regression Derivatives

- We will investigate how to compute derivatives for you to implement gradient descent for logistic regression
- Logistic regression:
 - $z = w^T x + b \rightarrow$ linear function
 - $\hat{y} = a = \sigma(z) \rightarrow$ nonlinear function
 - $\mathcal{L}(a, y) = -(y \log a + (1 - y) \log(1 - a)) \rightarrow$ loss for one example

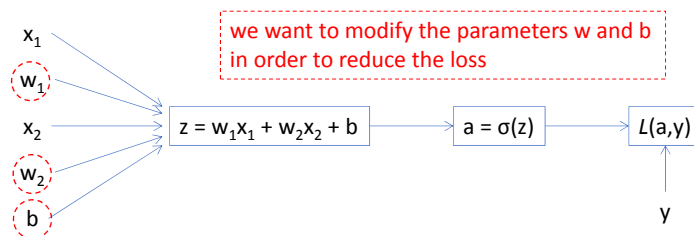
Logistic Regression Derivatives

- We will investigate how to compute derivatives for you to implement gradient descent for logistic regression
- Logistic regression:
 - $z = w^T x + b \rightarrow$ linear function
 - $\hat{y} = a = \sigma(z) \rightarrow$ nonlinear function
 - $\mathcal{L}(a, y) = -(y \log a + (1 - y) \log(1 - a)) \rightarrow$ loss for one example

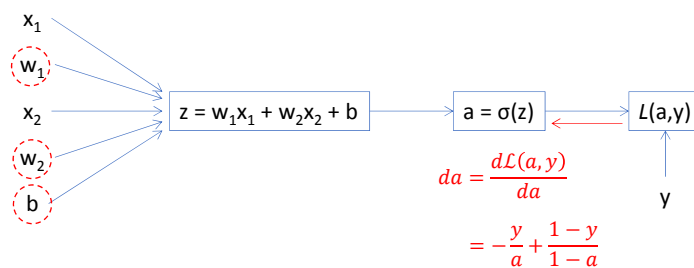


Logistic Regression Derivatives

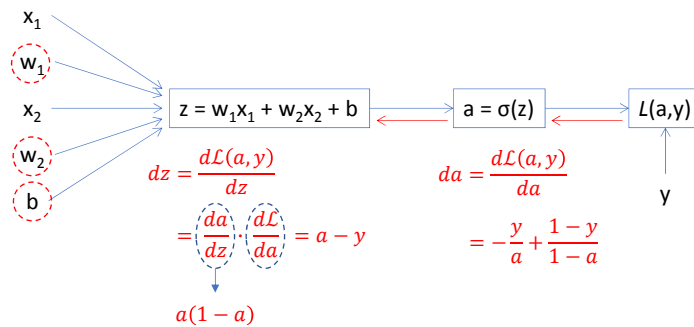
- We will investigate how to compute derivatives for you to implement gradient descent for logistic regression
- Logistic regression:
 - $z = w^T x + b \rightarrow$ linear function
 - $\hat{y} = a = \sigma(z) \rightarrow$ nonlinear function
 - $\mathcal{L}(a, y) = -(y \log a + (1 - y) \log(1 - a)) \rightarrow$ loss for one example



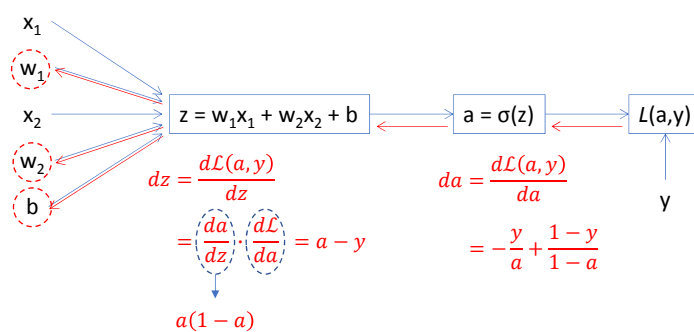
Logistic Regression Derivatives



Logistic Regression Derivatives



Logistic Regression Derivatives

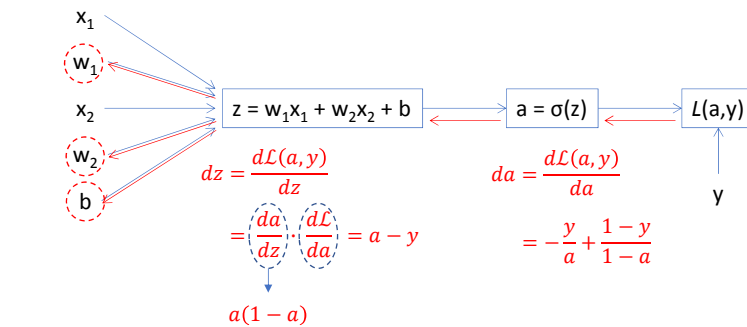


$$dw_1 = \frac{dL(a, y)}{dw_1} = \frac{dz}{dw_1} \cdot \frac{dL(a, y)}{dz} = x_1 \cdot dz = x_1 \cdot (a - y)$$

$$dw_2 = \frac{dL(a, y)}{dw_2} = \frac{dz}{dw_2} \cdot \frac{dL(a, y)}{dz} = x_2 \cdot dz = x_2 \cdot (a - y)$$

$$db = \frac{dL(a, y)}{db} = \frac{dz}{db} \cdot \frac{dL(a, y)}{dz} = dz = (a - y)$$

Logistic Regression Derivatives



$$dw_1 = \frac{dL(a, y)}{dw_1} = \frac{dz}{dw_1} \cdot \frac{dL(a, y)}{dz} = x_1 \cdot dz = x_1 \cdot (a - y)$$

$$dw_2 = \frac{dL(a, y)}{dw_2} = \frac{dz}{dw_2} \cdot \frac{dL(a, y)}{dz} = x_2 \cdot dz = x_2 \cdot (a - y)$$

$$db = \frac{dL(a, y)}{db} = \frac{dz}{db} \cdot \frac{dL(a, y)}{dz} = dz = (a - y)$$

Gradient Descent Algorithm

$$w_1 \leftarrow w_1 - \alpha \cdot dw_1$$

$$w_2 \leftarrow w_2 - \alpha \cdot dw_2$$

$$b \leftarrow b - \alpha \cdot db$$

Logistic Regression on m Examples

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

- Derivatives:

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

$$\vdots$$

$dw_1^{(i)}$ for $(x^{(i)}, y^{(i)})$

Logistic Regression on m Examples

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

For $i = 1$ to m

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

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

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

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

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

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

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

$J \neq m$

$dw_1 \neq m; dw_2 \neq m; db \neq m;$

Logistic Regression on m Examples

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

For $i = 1$ to m

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

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

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

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

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

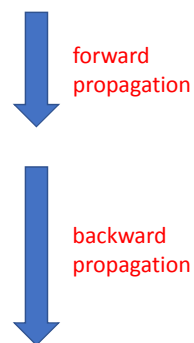
$$dw_2 += x_2^{(i)} dz^{(i)}$$

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

$J \neq m$

$dw_1 \neq m; dw_2 \neq m; db \neq m;$

gradient
computation



Logistic Regression on m Examples

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

For $i = 1$ to m

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

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

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

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

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

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

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

$J \neq m$

$dw_1 \neq m; dw_2 \neq m; db \neq m;$

gradient
computation

$$w_1 \leftarrow w_1 - \alpha \cdot dw_1$$

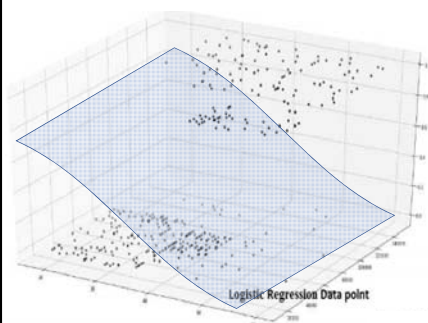
$$w_2 \leftarrow w_2 - \alpha \cdot dw_2$$

$$b \leftarrow b - \alpha \cdot db$$

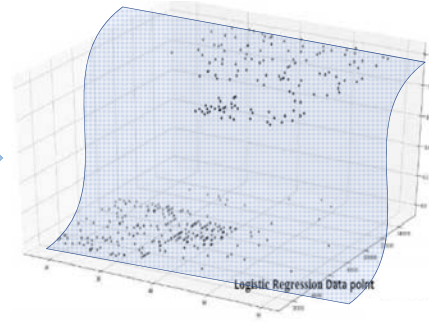
parameter
update

Logistic Regression by Gradient Descent Algorithm

$$y = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + b)}}$$



when w_1 , w_2 , and b are (randomly) initialized



when w_1 , w_2 , and b are converged by their sufficient updates