

## PRACTICAL 3 - LINEAR & LOGISTIC REGRESSION

### QUESTION 1

$$\begin{aligned} x^{(1)} &= [-2, 0] & y^{(1)} &= 2.0 \\ x^{(2)} &= [-1, 0] & y^{(2)} &= 3.0 \\ x^{(3)} &= [0, 0] & y^{(3)} &= 1.0 \\ x^{(4)} &= [2, 0] & y^{(4)} &= -1.0 \end{aligned}$$

more computation  
more expensive  
alternative to a closed form  
solution in linear regression  
an iterative is  
such as gradient descent

$$\hat{y} = w^T x + b = w^T x' \uparrow \\ w/bias$$

1. Find the closed form solution for a Linear regression that minimizes the sum of squared errors on the training data.

design matrix  $n \times (d+1)$   
 4x2 num of inputs  
 4x1+1 num of features  
 $x = \begin{bmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 2 \end{bmatrix}$   $y = \begin{bmatrix} 2 \\ 3 \\ 1 \\ -1 \end{bmatrix}$

Now find the  $w = [w_0, w_1]^T$  that minimizes sum of squared errors

$y = w^T x \Leftrightarrow$   
 $\Leftrightarrow y^T = x^T w$   
 $\Leftrightarrow x^T w = y^T$   
 $\Leftrightarrow w = (x^T)^{-1} y^T$

if  $x$  not invertible  
 we use the pseudo inverse

$w = x^+ y^T$   
 $= (x^T x)^{-1} x^T y$  ? what opt?  
 $= (x^T x)^{-1} x^T y$   
 pseudo inverse  $x^+$

 $w = (x^T x)^{-1} x^T y$   
 $= \left[ \begin{array}{cccc} 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 2 \end{array} \right] \left[ \begin{array}{c} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 2 \end{array} \right]^{-1} \times \left[ \begin{array}{cccc} 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 2 \end{array} \right] \left[ \begin{array}{c} 2 \\ 3 \\ 1 \\ -1 \end{array} \right]$   
 $= \left[ \begin{array}{cc} 4 & -1 \\ -1 & 9 \end{array} \right]^{-1} \left[ \begin{array}{cccc} 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 2 \end{array} \right] \left[ \begin{array}{c} 2 \\ 3 \\ 1 \\ -1 \end{array} \right]$   
 $= \frac{1}{\det(A)} \left[ \begin{array}{cc} d & -b \\ -c & a \end{array} \right] (\dots)$   
 $= \frac{1}{36-1} (\dots)$   
 $= \frac{1}{35} \left[ \begin{array}{cc} 9 & 1 \\ 1 & 4 \end{array} \right] (\dots)$   
 $= \left[ \begin{array}{c} 1.0286 \\ -0.8857 \end{array} \right]$

$$\begin{aligned} L(w) &= \frac{1}{2} \|y - xw\|^2 & \nabla_b \quad a^T b = a \\ & \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \\ & \text{labels} \quad \text{inputs} \quad \boxed{b^T b = A b + A^T b} \\ & = \frac{1}{2} (y - xw)^T (y - xw) \\ & = \frac{1}{2} (y^T y - 2y^T xw + w^T x^T w) \end{aligned}$$

$$\begin{aligned} \nabla_w L(w) &= \nabla_w \left[ \frac{1}{2} (y^T y - 2y^T xw + w^T x^T x w) \right] \\ &= \frac{1}{2} \underbrace{\nabla_w y^T y}_0 - \underbrace{\nabla_w y^T x w}_{x^T y} + \frac{1}{2} \cancel{(x^T x)} \\ &\stackrel{\cancel{(x^T x)}}{=} \frac{1}{2} \nabla_w (w^T x^T x w) \\ &= x^T x w - x^T y \\ &\downarrow \\ \nabla_w L(w) &= 0 \Leftrightarrow \end{aligned}$$

$$\begin{aligned} &\Leftrightarrow x^T x w - x^T y = 0 \\ &\Leftrightarrow w = (x^T x)^{-1} x^T y \end{aligned}$$

2. Predict the target value for  $x_{query} = [1]$

$\hat{y} = w^T x$   
 $\Rightarrow \text{dot product}$

$$= \left[ \begin{array}{c} 1.0286 \\ -0.8857 \end{array} \right]^T \left[ \begin{array}{c} 1 \\ 1 \end{array} \right] = 0.1429$$

3. Sketch the hyperplane along which the linear regression predicts points will fall.

hyperplane equation is given for a general input:

$$x = \begin{bmatrix} 1 \\ x_1 \end{bmatrix} \text{ by setting the } \hat{y} = w \cdot x \text{ to zero } \hat{y} = w \cdot x = 0 \Rightarrow \begin{bmatrix} 1.0286 \\ -0.8857 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ x_1 \end{bmatrix} = 0 \Rightarrow 1.0286 - 0.8857x_1 = 0$$

4. Compute the mean squared error produced by the linear regression

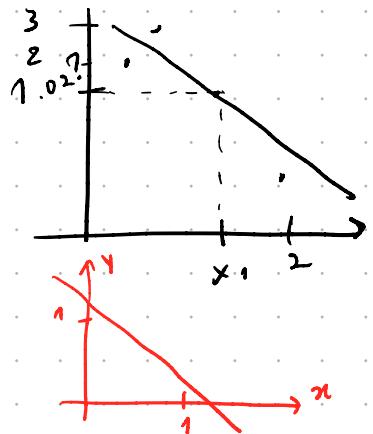
$$\begin{aligned} (y^{(1)} - \hat{y}^{(1)})^2 &= (y^{(1)} - w \cdot x^{(1)})^2 = (2.0 - 2.8)^2 = 0.64 \\ (y^{(2)} - \hat{y}^{(2)})^2 &= (y^{(2)} - w \cdot x^{(2)})^2 = 1.1783 \\ (y^{(3)} - \hat{y}^{(3)})^2 &= 0.0008 \\ &= 0.0661 \end{aligned}$$

$$MSE = 0.4714$$

$$\text{number of training examples } N$$

$$MSE = \frac{1}{N} \sum_{k=1}^N (\hat{y}^{(k)} - y^{(k)})^2$$

$$L(w) = \frac{1}{2} \|y - xw\|^2$$



$$\text{binary logistic regression } P_w(y=1|x) = \sigma(w \cdot x) = \frac{1}{1 + \exp(-w \cdot x)}$$

$$\text{cross entropy loss function } L(w) = - \sum_{i=1}^N \log(P_w(y^{(i)} | x^{(i)})) = - \sum_{i=1}^N (y^{(i)} \log \sigma(w \cdot x^{(i)}) + (1-y^{(i)}) \log(1-\sigma(w \cdot x^{(i)})))$$

gradient descent w/ update rule of size  $\eta$  towards the opposite direction of the gradient of the error function w/ respect to the weights

$$w \leftarrow w - \eta \frac{\partial L(w)}{\partial w}$$

$$L_i(w) = - [y^{(i)} \log(\sigma(w \cdot x^{(i)})) + (1-y^{(i)}) \log(1-\sigma(w \cdot x^{(i)}))]$$

$$\frac{\partial \log(\sigma(w \cdot x^{(i)}))}{\partial w} = \underset{\text{chain rule}}{\frac{1}{\sigma(w \cdot x^{(i)})}} \frac{\partial \sigma(w \cdot x^{(i)})}{\partial w}$$

$$\begin{aligned} \sigma(z) &= \frac{1}{1+e^{-z}} & \frac{\partial \sigma(z)}{\partial z} &= \sigma(z)(1-\sigma(z)) = \underset{\text{chain rule}}{\frac{\partial \sigma(w \cdot x^{(i)})}{\partial w}} = \\ && &= \sigma(w \cdot x^{(i)}) (1-\sigma(w \cdot x^{(i)})) \frac{\frac{\partial w \cdot x^{(i)}}{\partial w} x^{(i)}}{\sigma(w \cdot x^{(i)})} = \end{aligned}$$

$$\frac{\partial \log(\sigma(w \cdot x^{(i)}))}{\partial w} = (1-\sigma(w \cdot x^{(i)})) x^{(i)} \quad (1)$$

*probability model*

$$\frac{\partial}{\partial w} \log(1 - \sigma(w \cdot x^{(i)})) = \frac{1}{1 - \sigma(w \cdot x^{(i)})} \cdot \frac{\partial(-\sigma(w \cdot x^{(i)}))}{\partial w} = -\sigma(w \cdot x^{(i)})x^{(i)} \quad (2)$$

so by (1) & (2):

$$\begin{aligned} \frac{\partial L_i(w)}{\partial w} &= x^{(i)} [-y^{(i)}(1 - \sigma(w \cdot x^{(i)})) + (1 - y^{(i)})\sigma(w \cdot x^{(i)})] \\ &= x^{(i)} [\sigma(w \cdot x^{(i)}) - y^{(i)}] \end{aligned}$$

$$\frac{\partial L(w)}{\partial w} = \sum_{i=1}^N \frac{\partial L_i(w)}{\partial w} = + \sum_{i=1}^N x^{(i)} [y^{(i)} - \sigma(w \cdot x^{(i)})]$$

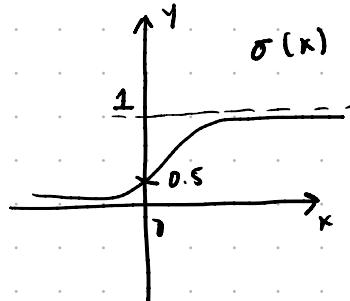
$$w = w - \eta \frac{\partial L(w)}{\partial w} = w + \sum_{i=1}^N x^{(i)} [y^{(i)} - \sigma(w \cdot x^{(i)})]$$

(better solution to this below)

2) compute the stochastic gradient descent update assuming initialization of all zeros. Assume learning rate = 1.

In SGD we make one update for each training sample instead of summing across all data points.

$$\begin{aligned} w &= w + n \cdot x^{(i)} (y^{(i)} - \sigma(w \cdot x^{(i)})) = \\ &= \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} (0 - \sigma(0)) = \\ &= \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -0.5 \\ 0.5 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0.5 \\ 0 \end{bmatrix} \end{aligned}$$



*(shorter  
form)*

- 2) we compute the MSE and the one-hot error is the better the least error is the better
- 3) lr-batch-gd (i.e. basically the sigmoid and then we use the function) so basically the update term
- 3 (a) stochastic (what's the difference??)
- 3 (?) softmax - converting values into probabilities

NOTES  
from code  
Python

## QUESTION 2 REMAKE

### Logistic Regression (Binary)

↓ the gradient of the sum is the sum of the gradients

$$w^T x = w_1 x_1 + \dots + w_N x_N$$

$$w^T \phi(x) = w_1 x_1 + \dots + w_{N+1} x_1$$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_N \\ x_i \end{bmatrix}$$

get higher order, it's where we transition from regular ML to deeper learning

$$\frac{\partial L(w)}{\partial w} = \frac{\partial}{\partial w} \left[ - \sum_{i=1}^N y^{(i)} \log(\sigma(w^T x^{(i)})) + (1-y^{(i)}) \log(1-\sigma(w^T x^{(i)})) \right]$$

Simpler deriving each entry and then try to write everything in matrix notation

$$\frac{\partial}{\partial w_i} \left[ y \log(\sigma(w^T x)) + (1-y) \log(1-\sigma(w^T x)) \right]$$

$$\frac{\partial \log(f(x))}{\partial x} = \frac{1}{f(x)} \cdot \frac{\partial f(x)}{\partial x}$$

$$= y \frac{\partial}{\partial w_i} \log(\sigma(w^T x)) + (1-y) \frac{\partial}{\partial w_i} \log(1-\sigma(w^T x))$$

$$= y \frac{1}{\sigma(w^T x)} \frac{\partial}{\partial w_i} \sigma(w^T x) + (1-y) \frac{1}{1-\sigma(w^T x)} \frac{\partial}{\partial w_i} (1-\sigma(w^T x))$$

$$\hat{\sigma}(z) = \left( \frac{1}{1+e^{-z}} \right)' = \underbrace{\frac{0 - (-e^{-z})}{(1+e^{-z})^2}}_{\sigma(z)} = \underbrace{\frac{1}{1+e^{-z}}}_{\sigma(z)} \cdot \frac{-(-e^{-z})}{1+e^{-z}} = \sigma(z) \cdot \frac{1-e^{-z}-1}{1+e^{-z}}$$

$$= \sigma(z) (1 - \sigma(z))$$

$$\sigma'(z) = \sigma(z) (1 - \sigma(z))$$

$$y \frac{1}{\sigma(w^T x)} \frac{\partial}{\partial w_i} \sigma(w^T x) + (1-y) \frac{1}{1-\sigma(w^T x)} \frac{\partial}{\partial w_i} (1-\sigma(w^T x))$$

$$= y \frac{1}{\sigma(w^T x)} (\sigma(w^T x) (1-\sigma(w^T x))) \frac{\partial}{\partial w_i} w^T x + (1-y) \frac{1}{1-\sigma(w^T x)} (-1) \cdot \sigma(w^T x) (1-\sigma(w^T x))$$

$$\textcircled{*} \frac{\partial}{\partial w_i} \underbrace{w^T x}_{x_i} g(f(x)) = g'(f(x)) f'(x)$$

$$= y (1-\sigma(w^T x)) x_i + (1-y) (-1) \sigma(w^T x) x_i$$

$$= y x_i - y \sigma(w^T x) x_i - \sigma(w^T x) x_i + y \sigma(w^T x) x_i$$

$$= (y - \sigma(w^T x)) x_i$$

$$\frac{\partial L}{\partial w_i} = - \sum_{k=1}^N (y^{(k)} - \sigma(w^T x^{(k)})) x_i^{(k)}$$

$$\nabla_w L(w) = \sum_{k=1}^N (\sigma(w^T x^{(k)}) - y^{(k)}) x^{(k)}$$

$$\nabla_w L(w) = \begin{bmatrix} \frac{\partial}{\partial w_1} L(w) \\ \vdots \\ \frac{\partial}{\partial w_D} L(w) \end{bmatrix}$$

$$w \leftarrow w - \eta \nabla_w L(w) = w - \eta \sum_{k=1}^N [\sigma(w^T x^{(k)}) - y^{(k)}] x^{(k)}$$

↑  
learning rate

Stochastic GD

$$w \leftarrow w - \eta [\sigma(w^T x^{(i)}) - y^{(i)}] x^{(i)}$$

$$w = [-0.5 \ 0.5 \ 0]^T$$

adapted to stochastic  
GD