

# Optimization for Linear Models

MSE for Linear Reg:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T \vec{x}^{(i)})^2 \rightarrow \nabla J(\theta) = \dots \rightarrow \text{GD}$$

MSE for Perceptron:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \text{sign}(\theta^T \vec{x}^{(i)}))^2 \rightarrow \nabla J(\theta) = \text{FAIL}$$

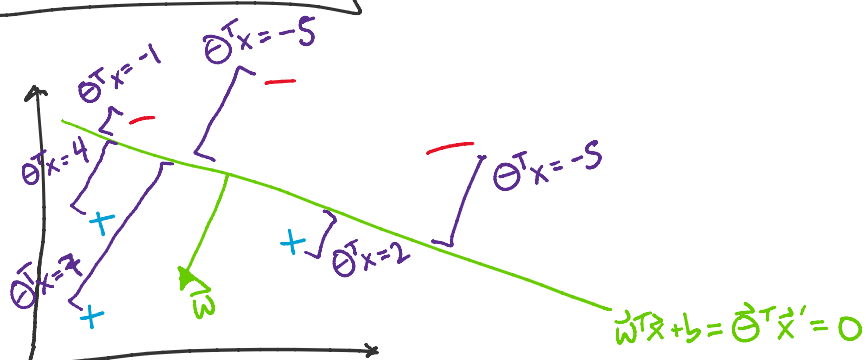
$y^{(i)}$  not differentiable

MSE for Logistic Regression

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\theta^T \vec{x}^{(i)}))^2 \rightarrow \nabla J(\theta) = \dots \rightarrow \text{GD}$$

$\sigma$  chosen to be differentiable

## What is $\theta^T \vec{x}$ ?



## Binary Logistic Regression

① Model:

$$y \sim \text{Bernoulli}(\phi)$$

$$\phi = \sigma(\theta^T \vec{x}) \text{ where } \sigma(u) = \frac{1}{1 + \exp(-u)}$$

$$p(y | \vec{x}, \theta) = \begin{cases} \sigma(\theta^T \vec{x}) & \text{if } y = 1 \\ 1 - \sigma(\theta^T \vec{x}) & \text{if } y = 0 \end{cases}$$

$1 - \phi$

② Objective:

$$\begin{aligned} \ell(\vec{\theta}) &= \log p(D | \vec{\theta}) = \log \prod_{i=1}^N p(y^{(i)} | \vec{x}^{(i)}, \vec{\theta}) \\ &= \sum_{i=1}^N \log p(y^{(i)} | \vec{x}^{(i)}, \theta) \end{aligned}$$

$$\log(a \cdot b) = \log a + \log b$$

\* negative average conditional log-likelihood,  $J(\theta)$  for Log. Reg. is convex!

$$= \sum_{i=1}^N \log p(y^{(i)} | \vec{x}^{(i)}, \theta)$$

$$J(\vec{\theta}) = -\frac{1}{N} \ell(\theta) = \frac{1}{N} \sum_{i=1}^N \underbrace{-\log p(y^{(i)} | \vec{x}^{(i)}, \theta)}_{J^{(i)}(\theta)}$$

likelihood,  $J(\theta)$   
for Log. Reg. is convex!

$$\theta_{MLE} = \operatorname{argmax} \ell(\theta)$$

$$= \operatorname{argmin} -\ell(\theta)$$

$$= \operatorname{argmin} -\frac{1}{N} \ell(\theta)$$

### ③ Derivatives

$$\frac{\partial J^{(i)}(\vec{\theta})}{\partial \theta_m} = \frac{\partial}{\partial \theta_m} -\log(p(y^{(i)} | \vec{x}^{(i)}, \vec{\theta}))$$

$$= \begin{cases} \partial/\partial \theta_m - \log(\sigma(\vec{\theta}^T \vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ \partial/\partial \theta_m - \log(1 - \sigma(\vec{\theta}^T \vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

$$= \dots$$

$$= \dots$$

$$= \dots$$

$$= -\left( \underbrace{y^{(i)}}_{\text{true value}} - \underbrace{\sigma(\vec{\theta}^T \vec{x}^{(i)})}_{\text{prob. of } y^{(i)}=1} \right) \underbrace{x_m^{(i)}}_{\text{with feature}}$$

$$\nabla J^{(i)}(\theta) = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = -(y^{(i)} - \sigma(\theta^T \vec{x}^{(i)})) \vec{x}^{(i)}$$

$$\nabla J(\vec{\theta}) = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\theta)$$

recitation  
or HW

④ Find  $\vec{\theta}$  by gradient descent or SGD

⑤ Predict the most probable class, for new  $\vec{x}$

$$\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} p(y | \vec{x}, \vec{\theta})$$

$$= \begin{cases} y=1 & \text{if } p(y=1 | \vec{x}, \vec{\theta}) \geq 0.5 \\ y=0 & \text{otherwise} \end{cases}$$

$$= \text{"sign"}(\vec{\theta}^T \vec{x}) \leftarrow \text{Exercise: Prove that this is true}$$

OH

$$w^T x + b = 0$$

$$x_2 = \left( \frac{-w_1}{w_2} \right) x_1 + \left( \frac{-b}{w_2} \right)$$

- repeat
- ① pick hyp.  $\alpha, \beta$
  - ② train w/  $\alpha$  and  $\beta$
  - ③ report val error



$$\left( y^{(i)} - (\theta_1 x_1 + \theta_2 x_2) \right)^2$$

$$\left( 7 - (\theta_1 \cdot 3 + \theta_2 \cdot 2) \right)^2$$

