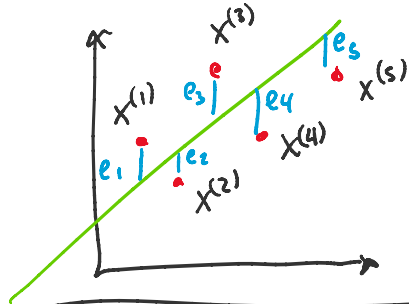


$$D = \{(\vec{x}^{(i)}, y^{(i)})\}_{i=1}^N \quad \hat{y}^{(i)} = h(\vec{x}^{(i)})$$

Residuals

Def: a residual is the vertical distance from observed value $y^{(i)}$ to predicted output $\hat{y}^{(i)}$



$$e_i = |y^{(i)} - h(\vec{x}^{(i)})|$$

$$= |y^{(i)} - (\vec{w}^T \vec{x}^{(i)} + b)|$$

Key Idea of Lin. Reg.

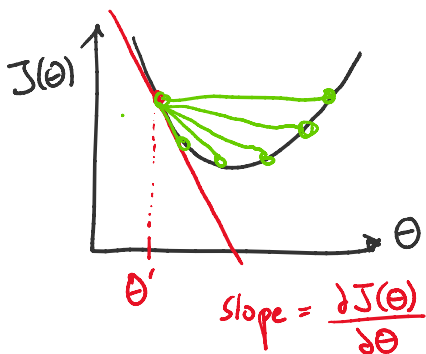
Find the linear function h (w/parameters w and b) that minimizes the squares of the residuals for a training set. *one option*

Def: mean squared error (MSE)

$$J_D(\vec{w}, b) = \frac{1}{N} \sum_{i=1}^N (e_i)^2 = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (\vec{w}^T \vec{x}^{(i)} + b))^2$$

Derivatives

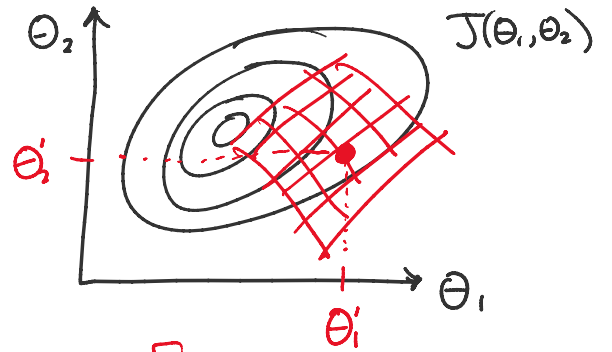
① Deriv. as Slope of Tangent



② Deriv. as a Limit

$$\frac{\partial J(\theta)}{\partial \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta)}{\epsilon}$$

③ Partial Deriv. as Tangent Planes



$$\left[\begin{array}{l} \partial J(\theta_1, \theta_2) / \partial \theta_1 \\ \partial J(\theta_1, \theta_2) / \partial \theta_2 \end{array} \right]$$

$$\frac{\partial J(\theta)}{\partial \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta)}{\epsilon}$$

limit of secants is the tangent

$$\left[\frac{\partial J(\theta_1, \theta_2)}{\partial \theta_2} \right]$$

Gradient

Def: the gradient of $J: \mathbb{R}^M \rightarrow \mathbb{R}$

$$\nabla J(\vec{\theta}) = \begin{bmatrix} \frac{\partial J(\vec{\theta})}{\partial \theta_1} \\ \frac{\partial J(\vec{\theta})}{\partial \theta_2} \\ \vdots \\ \frac{\partial J(\vec{\theta})}{\partial \theta_M} \end{bmatrix}$$

first order partial derivative

Gradient Descent

Algorithm:

- ① Choose an initial point $\vec{\theta}$
- ② Repeat $t=1, 2, 3, \dots$
 - a) Compute gradient $\vec{g} = \nabla J(\vec{\theta})$
 - b) Choose a step size γ_t
 - c) Update $\vec{\theta} \leftarrow \vec{\theta} - \gamma_t \vec{g}$
- ③ Return $\vec{\theta}$ when stopping criterion is reached

Remarks:

Initial Point:

- randomly
- $\vec{\theta} = \text{all zeros}$

Step Size:

- fixed value $\gamma = 0.1$
- schedule $\gamma_t = \frac{\gamma_0}{(t-1)\gamma_0 + 1}$
- line search

Stopping Criterion:

- when $\vec{g} \approx \vec{0}$
- $\|\nabla J(\vec{\theta})\|_2 < \epsilon$
for $\epsilon = 10^{-8} (f(\theta))^2$

Gradient for Lin. Reg.

$$MSE: \tau(\vec{x}) = \frac{1}{N} \sum J^{(i)}(\vec{\theta})$$

$$\text{where } J^{(i)}(\vec{\theta}) = \frac{1}{2} \left(y^{(i)} - \underbrace{\vec{\theta}^T \vec{x}^{(i)}}_{f(\theta)} \right)^2$$

Gradient for LM. ...

$$\text{MSE: } J(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\vec{\theta})$$

$$\text{where } J^{(i)}(\vec{\theta}) = \frac{1}{2} (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})^2$$

↑ does not change the argmin

Partial Derivatives:

$$\begin{aligned} \frac{\partial J(\vec{\theta})}{\partial \theta_j} &= \frac{1}{2} \cdot 2 (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) \frac{\partial}{\partial \theta_j} (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) \\ &= (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) \frac{\partial}{\partial \theta_j} \left(y^{(i)} - \sum_{m=1}^M \theta_m x_m^{(i)} \right) \\ &= - (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) x_j^{(i)} \end{aligned}$$

Gradient:

$$\nabla J^{(i)}(\vec{\theta}) = \begin{bmatrix} \partial J^{(i)}(\vec{\theta}) / \partial \theta_1 \\ \partial J^{(i)}(\vec{\theta}) / \partial \theta_2 \\ \vdots \\ \partial J^{(i)}(\vec{\theta}) / \partial \theta_M \end{bmatrix} = - \underbrace{(y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})}_{\text{scalar}} \underbrace{\vec{x}^{(i)}}_{\text{vector}}$$

$$\nabla J(\vec{\theta}) = \nabla \left(\frac{1}{N} \sum_{i=1}^N J^{(i)}(\vec{\theta}) \right) = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\vec{\theta})$$

$$= \frac{1}{N} \sum_{i=1}^N - (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) \vec{x}^{(i)}$$

known b/c in training
known b/c passed it in

OH



