

Day 2 Last time:

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$\text{Loss/cost function } J(\theta) = \sum_{i=1}^{2x_1} (y_i - h_\theta(x_i))^2$$

Goal: find  $\hat{\theta} = \arg \min \theta J(\theta)$

Normal regression:  $y_i \stackrel{\text{iid}}{\sim} N(h_\theta(x_i), \sigma^2)$  or equivalently  $y_i = h_\theta(x_i) + \epsilon_i$ ,  $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$   
 Question: if we assume Normal distribution, how does it change the loss function?

Estimation by Maximum Likelihood:

$$f(y_i | x_i, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - h_\theta(x_i))^2}$$

Reminder: probability distr =  $f(\text{data}/\text{params})$   
 Likelihood =  $f(\text{params}/\text{data})$

$$L(\theta | y_1, \dots, y_N, x_1, \dots, x_N) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - h_\theta(x_i))^2} = \frac{1}{(\sqrt{2\pi\sigma^2})^N} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - h_\theta(x_i))^2}$$

$$\arg \max_{\theta} L(\theta | y_1, \dots, y_N, x_1, \dots, x_N) = \arg \max_{\theta} \log L(\theta, y_1, \dots, y_N, x_1, \dots, x_N)$$

$$= \arg \max_{\theta} \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - h_\theta(x_i))^2 \right] = \arg \min_{\theta} \sum_{i=1}^N (y_i - h_\theta(x_i))^2 \Rightarrow \text{same criterion as ordinary least squares}$$

- Comments
- Assumption of Normal distribution is not required for OLS (only ML)
  - For ML estimation, the marginal distribution of  $y$  is not Normal (only conditional  $y|x$ )

Least squares optimization with one parameter

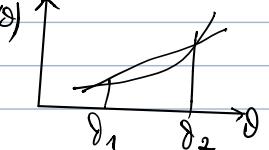
$$J(\theta) = \sum_{i=1}^N (y_i - \theta x_i)^2$$

- continuous wrt  $\theta$
- smooth (has non-zero 1st and 2nd derivatives)
- convex

A function  $J(\theta) : \mathbb{R}^P \rightarrow \mathbb{R}$  is convex if for any  $\theta_1, \theta_2 \in \mathbb{R}^P$  and  $t \in [0, 1]$

$$J(t\theta_1 + (1-t)\theta_2) \leq t J(\theta_1) + (1-t) J(\theta_2)$$

$J(\theta)$  has a single stationary point, i.e.  $\theta: \frac{dJ(\theta)}{d\theta} = 0$   
 $\Rightarrow$  the stationary point is the global minimum



Determining the stationary point:

$$\frac{dJ(\theta)}{d\theta} = \frac{d}{d\theta} \sum_{i=1}^N (y_i - \theta x_i)^2 = - \sum_{i=1}^N 2(y_i - \theta x_i) x_i \stackrel{\text{set to 0}}{=} 0$$

$$\Rightarrow - \sum_{i=1}^N y_i x_i + \theta \sum_{i=1}^N x_i^2 = 0 \quad \Rightarrow \hat{\theta}_1 = \frac{\sum y_i x_i}{\sum x_i^2} \quad \text{- closed-form solution}$$

For strictly convex functions, the stationary point is uniquely identifiable

Optimization with multiple weights:

$$\underset{w \in \mathbb{R}^p}{J(\theta)} = \sum_{i=1}^N (y_i - \theta_0 - \theta_1 x_i)^2$$

Gradient  $\nabla J(\theta) = [\frac{\partial J(\theta)}{\partial \theta_1}, \frac{\partial J(\theta)}{\partial \theta_2}]^T$ . Written as a system of equations:

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_0} &= -\sum_{i=1}^N 2(y_i - \theta_0 - \theta_1 x_i) = 0 \\ \frac{\partial J(\theta)}{\partial \theta_1} &= -\sum_{i=1}^N 2(y_i - \theta_0 - \theta_1 x_i)x_i = 0 \end{aligned}$$

2 equations  
2 unknowns  
 $\Rightarrow$  unique analytical solution under some conditions

### Matrix notation

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix}$$

$y_{N \times 1} \quad X_{N \times p} \quad \theta_{p \times 1} \quad e_{N \times 1}$

$$\begin{aligned} J(\theta) &= \sum_{i=1}^N (y_i - x_i^\top \theta)^2 = \|y - X\theta\|_2^2 \\ &= \underbrace{(y - X\theta)^\top}_{1 \times N} \underbrace{(y - X\theta)}_{N \times 1} \\ &= y^\top y - \theta^\top X^\top y - y^\top X\theta + \theta^\top X^\top X\theta \end{aligned}$$

Gradient in matrix notation (i.e., derivative in each dimension of  $\theta$ )

$$\nabla J(\theta) = -X^\top y - (X^\top X)^\top + 2X^\top X\theta = -2X^\top y + 2X^\top X\theta \stackrel{\text{set to } 0}{=} 0$$

$$X^\top X\theta - X^\top y = 0 \quad \Rightarrow \quad X^\top X\theta = X^\top y \quad \} \text{Normal equations}$$

$$\theta = (X^\top X)^{-1} X^\top y \quad \Rightarrow \quad \text{analytical solution}$$

Note:  $X^\top X$  cannot be singular. I.e.,  $N > p$

$X$  full rank (no linear dependencies among columns of  $X$ )

### Prediction

$$\hat{y} = X\hat{\theta} = \underbrace{X(X^\top X)^{-1}}_{H - \text{hat matrix}} X^\top y = Hy$$

-projection matrix (projects  $y$  on the subspace defined by columns of  $X$ )

Note:  $\hat{\theta}$  and  $\hat{y}$  are linear combinations of  $y$ . This is why it is called linear regression

Cost of the inverse:  $O(\underbrace{p^2 N}_{X^\top X \text{ inverse}} + \underbrace{p^3})$  - expensive if  $N$  or  $p$  are large

## Numeric optimization in high dimensions: gradient descent

Goal: locally compute direction of descent; iterate

Back to the one-parameter case:  $J(\delta) = \sum_{i=1}^N (y_i - \delta x_i)^2$

Based on Taylor series approximation:

$$J(\delta) = \sum_{t=0}^{\infty} \frac{1}{t!} \left. \frac{d^t J(\delta)}{d\delta^t} \right|_{\delta_0} (\delta - \delta_0)^t$$

Only keep the first two derivatives:

$$\bullet J(\delta) \approx J(\delta_0) + (\delta - \delta_0) \underbrace{\left. \frac{d J(\delta)}{d\delta} \right|_{\delta_0}}_{\substack{\text{second-order} \\ \text{approximation}}} + \underbrace{\frac{1}{2} (\delta - \delta_0)^2 \left. \frac{d^2 J(\delta)}{d\delta^2} \right|_{\delta_0}}_{\substack{\text{first} \\ \text{derivative}}} \text{ second derivative}$$

• Stationary point: derivative of the approximation set to 0

$$\frac{d J(\delta)}{d\delta} \approx \left. \frac{d J(\delta)}{d\delta} \right|_{\delta_0} + (\delta - \delta_0) \left. \frac{d^2 J(\delta)}{d\delta^2} \right|_{\delta_0} = 0$$

$$\Rightarrow \delta_1 = \delta_0 - \frac{\left. \frac{d J(\delta)}{d\delta} \right|_{\delta_0}}{\left. \frac{d^2 J(\delta)}{d\delta^2} \right|_{\delta_0}} \quad \text{or} \quad \delta_1 = \delta_0 - \frac{J'(\delta_0)}{J''(\delta_0)}$$

Since the second-order approximation was inaccurate in the first place, we need to iterate through iterations  $t=1, 2, \dots$  until converging

$$\delta^{(t+1)} = \delta^{(t)} - \frac{J'(\delta^{(t)})}{J''(\delta^{(t)})} \quad \begin{array}{l} \text{- Newton-Raphson method, or} \\ \text{second-order gradient descent} \end{array}$$

Simplify the algorithm to avoid the second derivative

$$J(\delta) \approx J(\delta_0) + (\delta - \delta_0) \left. \frac{d J(\delta)}{d\delta} \right|_{\delta_0} + \frac{1}{2} (\delta - \delta_0)^2 \cdot \frac{1}{2}$$

$$\Rightarrow \delta^{(t+1)} = \delta^{(t)} - \cancel{\lambda} \cdot J'(\delta^{(t)}) \quad \begin{array}{l} \text{arbitrary constant replacing} \\ \text{learning rate} \end{array} \quad \begin{array}{l} \text{the second derivative} \\ \text{replacing} \end{array}$$

- this is first-order gradient descent

Extending to multiple params of linear regression:

$$\underset{p \times 1}{\delta^{(t+1)}} = \underset{p \times 1}{\delta^{(t)}} - \cancel{\lambda} \nabla J(\delta) \Big|_{\delta^{(t)}} = \underset{p \times 1}{\delta^{(t)}} - \cancel{\lambda} \underset{p \times N}{X^T} \underset{N \times p}{(X \delta^{(t)} - y)} \underset{N \times 1}{}$$

Cost:  $O(N)$   
expensive for large  $N$

Algo: batch gradient descent for linear regression:

- Initialize  $\delta$
- Repeat until convergence {
  - For  $j = 1, \dots, p$  {
$$\delta_j \leftarrow \delta_j + \lambda \frac{d}{d\delta_j} \sum_{i=1}^N (y_i - x_i^T \delta)^2$$
} can replace with line search}

Algo: stochastic gradient descent:

- Initialize  $\delta$
- Repeat until convergence {
  - shuffle the order of observations  $i = 1, \dots, N$
  - For  $i = 1, \dots, N$  {
    - For  $j = 1, \dots, p$  {
$$\delta_j \leftarrow \delta_j + \lambda \frac{d}{d\delta_j} (y_i - x_i^T \delta)^2$$
} can replace with line search}}