

Lecture 8: 2/13/17

Optimization for Linear Regression

#1) Closed Form Solution:

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_M^{(1)} \\ \vdots & \vdots & & \vdots \\ 1 & x_1^{(N)} & \dots & x_M^{(N)} \end{bmatrix}$$

← $\vec{x}^{(1)}$
← $\vec{x}^{(N)}$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - (\theta^T x^{(i)}))^2$$

← "Design Matrix"
← "Ordinary Least Squares"

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y})$$

$J(\theta)$

$$= (X^T X)^{-1} (X^T \vec{y})$$

the "Normal Equations"

$$\nabla J(\theta) = X^T X \theta - X^T \vec{y} = 0$$

Computational Complexity of Closed Form

Background: Matrix Multiplication

Matrices A and B

If A is $q \times r$ and B is $r \times s$: AB takes $O(qrs)$

If A and B are $q \times q$: AB takes $O(q^{2.373})$

If A is $q \times q$: A^{-1} takes $O(q^{2.373})$

OLS

$$\underbrace{\underbrace{(X^T X)^{-1}}_{M \times M}}_{M \times M} \underbrace{\underbrace{(X^T y)}_{M \times 1}}_{M \times 1}$$

$X^T X$	$O(M^2 N)$
$()^{-1}$	$O(M^{2.373})$
$X^T y$	$O(MN)$
$()()$	$O(M^2)$
total	$O(M^2 N + M^{2.373})$

⇒ Linear in # examples, N
Polynomial in # features, M

Stability

Can we invert $X^T X$?

Case #1: When $N \gg M$; Yes so long as X has full column rank

Case #2: When $N \ll M$: X does not have full rank
 \Rightarrow problem is non-identifiable

#2) SGD For Linear Regression

Called "Least Mean Squares (LMS)" algorithm ← doesn't affect argmin

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

$$\nabla J_i(\theta) = \begin{bmatrix} \frac{dJ_i(\theta)}{d\theta_1} \\ \vdots \\ \frac{dJ_i(\theta)}{d\theta_M} \end{bmatrix}$$

$$\begin{aligned} \frac{dJ_i(\theta)}{d\theta_j} &= (y^{(i)} - \theta^T \vec{x}^{(i)}) \frac{d}{d\theta_j} \left(y^{(i)} - \sum_{n=1}^M \theta_n x_n^{(i)} \right) \\ &= - \underbrace{(y^{(i)} - \theta^T \vec{x}^{(i)})}_{\text{scalar}} \underbrace{x_j^{(i)}}_{\text{vector}} \end{aligned}$$

$$\star = - \underbrace{(y^{(i)} - \theta^T \vec{x}^{(i)})}_{\text{scalar}} \underbrace{\vec{x}^{(i)}}_{\text{vector}}$$

① Choose initial $\vec{\theta}$ (zeros, random)

② For $t = 1 \dots T$:

a) Sample $i \sim \text{Uniform}(\{1, \dots, N\})$

b) compute $\nabla J_i(\theta)$

c) Choose learning rate $\gamma^{(t)}$

d) Update $\vec{\theta} = \vec{\theta} - \gamma^{(t)} \nabla J_i(\theta)$

for m in $\{1, \dots, M\}$:

$$\theta_m = \theta_m + \gamma (\quad)$$

for fast code it's better to use the vector version \star

#3 Gradient Descent for Lin. Reg.

- Straightforward application of GD to solve: $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$

- How to compute $\nabla J(\theta)$?

$$\Rightarrow \nabla J(\theta) = \sum_{i=1}^N \nabla J_i(\theta) = \sum_{i=1}^N -(y^{(i)} - \theta^T \vec{x}^{(i)}) \vec{x}^{(i)}$$

Convex

Prob. Interp. of Linear Regression

Def: Generative probabilistic model is $p(x, y | \theta)$

Def: Discriminative probabilistic model is $p(y | x, \theta)$

Generative vs. Discriminative Views

	Gen	Disc.
① Assume unk. data dist.	$(x^{(i)}, y^{(i)}) \sim p^*(x, y \theta^*)$	$(x^{(i)}, y^{(i)}) \sim p^*(x, y \theta^*)$
② Choose model family	$p(x, y \theta)$	$p(y x, \theta)$
③ Learn	Maximize Likelihood $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log \prod_{i=1}^N p(x^{(i)}, y^{(i)} \theta)$	Maximize Cond. Likelihood. $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log \prod_{i=1}^N p(y^{(i)} x^{(i)}, \theta)$
④ Predict.	$\hat{y} = h(x) = \underset{y}{\operatorname{argmax}} p(y x, \theta)$ $= \underset{y}{\operatorname{argmax}} p(x, y \theta)$	$\hat{y} = h(x) = \underset{y}{\operatorname{argmax}} p(y x, \theta)$

★ Replace $p(x, y)$ and $p(y | x)$ with probability density functions if y continuous

★ Takeaway: $p(y | x)$ is used for predict, so just learn it directly.

Recall: Gaussian Dist.

Let X be a Gaussian r.v.: $X \sim \text{Gaussian}(\mu, \sigma^2)$
 $\mu \in \mathbb{R}$ $\sigma^2 > 0$

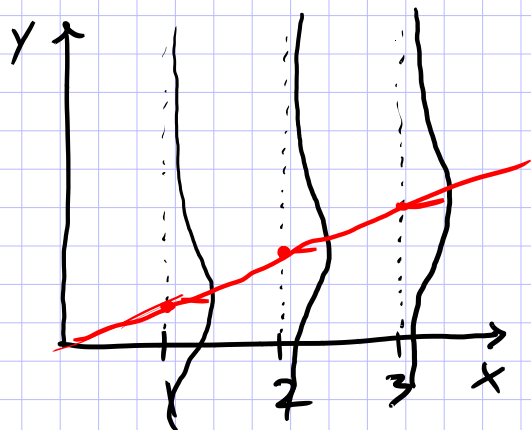
Def: probability density fn.

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\sigma^2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

Case #1: 1D Lin. Reg.

Assume: two r.v.s X, Y

Goal: learn a function $f: X \rightarrow P(Y|X)$



Case #2: Multiple Linear Regression

Assume: r.v.s $X_1^{(i)}, \dots, X_M^{(i)}, Y^{(i)}$

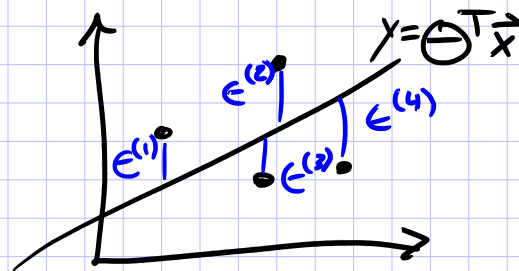
data $D = \{(\vec{x}^{(i)}, y^{(i)})\}_{i=1}^n$

Story: (generate $y^{(i)}$ given $x^{(i)}$)

$\epsilon^{(i)} \sim \text{Gaussian}(\mu=0, \sigma^2)$ ← 1D Gaussian

$$y^{(i)} = \vec{\theta}^T \vec{x}^{(i)} + \epsilon^{(i)}$$

← Is $Y^{(i)}$ still a r.v.?



Equivalent Story:

$$y^{(i)} \sim \text{Gaussian}(\mu = \vec{\theta}^T \vec{x}^{(i)}, \sigma^2)$$

Q: What is the pdf of $y^{(i)}$?

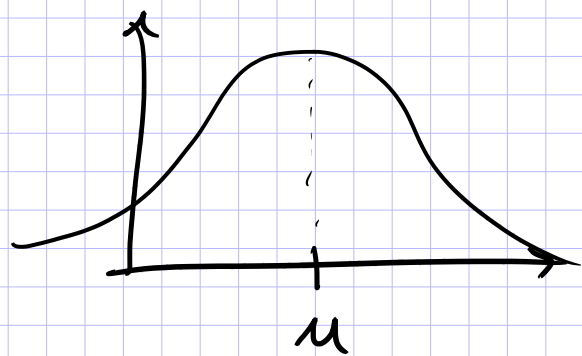
$$g(y^{(i)} | x^{(i)}, \vec{\theta}, \sigma^2) = f_{\text{Gaussian}}(y^{(i)} | \mu = \vec{\theta}^T \vec{x}^{(i)}, \sigma^2)$$

$$= \frac{1}{\sqrt{2\sigma^2\pi}} \exp\left\{-\frac{(y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})^2}{2\sigma^2}\right\}$$

Predictions

$$\hat{y} = h(\vec{x}) = \underset{y \in \mathbb{R}}{\operatorname{argmax}} g(y | \vec{x}, \vec{\Theta}, \sigma^2)$$

The mode of a Gaussian is its mean.
(most probable value)



$$\hat{y} = h(\vec{x}) \triangleq \vec{\Theta}^T \vec{x}$$

Equivalence #1: Given the prob. model $p(y|x, \theta)$ above the predictions are just as for our "fn. approx." linear reg. approach.

Learning

Cond. Log. Like

$$\begin{aligned} l(\theta) &= \log \prod_{i=1}^N g(y^{(i)} | x^{(i)}, \theta, \sigma^2) \\ &= \sum_{i=1}^N \log g(y^{(i)} | x^{(i)}, \theta, \sigma^2) \\ &= \sum_{i=1}^N \left[-\log(\sqrt{2\sigma^2\pi}) - \frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2 \right] \end{aligned}$$

MLE

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmax}} l(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \left[-\log(\sqrt{2\sigma^2\pi}) - \frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2 \right] \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N -\frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2 \\ &= \underset{\theta}{\operatorname{argmin}} \underbrace{\sum_{i=1}^N \frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2}_{\text{MSE!}} \end{aligned}$$

Equivalence #2:

MLE for above prob. $p(y|x, \theta) = \theta^T x + \epsilon$ is equivalent to minimizing MSE for $h(\vec{x}) = \theta^T x$