

COMP 4220: Machine Learning, Fall 2018

Exam 1

Date: October 15, 2018

75
100

1. Let X and Y be two random variables, β a constant, and $\epsilon \sim \mathcal{N}(0, \sigma^2)$ be a Gaussian random variable with zero mean and variance σ^2 . We assume that $Y = \beta X + \epsilon$, and that ϵ is independent of X .

- (a) Show that given $X = x$, the distribution of Y is $\mathcal{N}(\beta x, \sigma^2)$.
 (b) Let $\{(X^{(i)}, Y^{(i)}), i = 1, \dots, n\}$ be n independent samples from the model above. Show that the maximum likelihood estimation of β has the following:

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (Y^{(i)} - \beta X^{(i)})^2$$

- (c) Show that solution of the above problem is:

$$\hat{\beta} = \frac{\sum_{i=1}^n Y^{(i)} X^{(i)}}{\sum_{i=1}^n (X^{(i)})^2}$$

a) $P(Y|X=x, \beta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(Y-\beta x)^2}$ ← This might be wrong. I forgot the Gaussian Distribution formula in Cheat sheet.

b) The above will yield something in the form of $P(Y|X=x, \beta) \sim \mathcal{N}(\beta x, \sigma^2)$, I didn't have the Gaussian distribution formula. My closest guess is the following:

$$\beta = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(Y^{(i)} - \beta X^{(i)})^2}$$

$$\log \beta = \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(Y^{(i)} - \beta X^{(i)})^2} \right) = \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \sum_{i=1}^n \log \left(e^{-\frac{1}{2\sigma^2}(Y^{(i)} - \beta X^{(i)})^2} \right)$$

Drop this because is constant.

$$\hat{\beta} = \arg \max_{\beta} -\frac{1}{2} \sum_{i=1}^n (Y^{(i)} - \beta X^{(i)})^2$$

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \sum_{i=1}^n (Y^{(i)} - \beta X^{(i)})^2$$

$$\hat{\beta} = \frac{P(Y|X, \beta) \cdot P(X)}{P(Y)}$$

- 10/30 2. Let X and Y be two random variables, and $Y = \beta X + \epsilon$, β is a constant, and $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Given n independent sample points $(X^{(1)}, Y^{(1)}), \dots, (X^{(n)}, Y^{(n)})$, instead of ordinary least squares, here we estimate β with "ridge regression", by solving the following problem:

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \left(\sum_{i=1}^n (Y^{(i)} - \beta X^{(i)})^2 + \lambda \beta^2 \right)$$

where $\lambda \geq 0$ is a tuning parameter.

(a) Give a solution in explicit form for $\hat{\beta}$.

(b) When λ goes to infinity, how does $\hat{\beta}$ change? Give a brief explanation.

0/10

$$\nabla_{\beta} \frac{1}{2} \left(\sum_{i=1}^n (y^{(i)} - \beta x^{(i)})^2 + \lambda \beta^2 \right)$$

$$\nabla_{\beta} \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \beta x^{(i)})^2 + \nabla_{\beta} \lambda \beta^2$$

$$\nabla_{\beta} \frac{1}{2} (\beta x^{(i)} - y^{(i)}) + \nabla_{\beta} (\lambda \beta^2)$$

10/20

$$0 = -\beta^T y + \beta^T \beta x + 2\lambda \beta$$

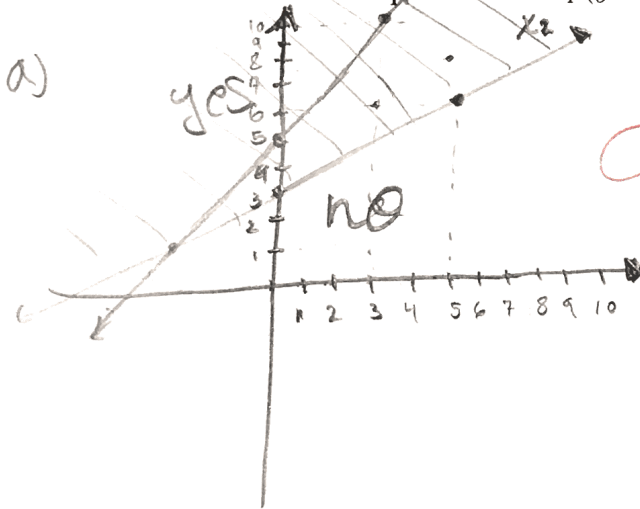
$$x = (\beta^T \beta + 2\lambda \mathbf{I})^{-1} \beta^T y$$

3. Suppose you are given the following classification task: predict the target $y \in \{0, 1\}$ given two real valued features x_1 and x_2 . After some training, you learn the following decision rule:

$$"y = 1 \text{ if } \underbrace{w_1 x_1 + w_2 x_2 + w_0}_{\substack{3 \quad 5 \quad -15}} \geq 0 \text{ and } y = 0 \text{ otherwise}"$$

where $w_1 = 3, w_2 = 5, w_0 = -15$.

- (a) Plot the decision boundary and label the region where we would predict $y = 1$ and $y = 0$.
 (b) Suppose that we learned the above weights using logistic regression. Using this model, what would be our prediction for $p(y = 1 | x_1, x_2)$?



$$3x_1 + 5x_2 - 15 \geq 0$$

$$3x_1 + 5x_2 \geq 15$$

$$x_1 \geq \frac{5}{3}x_2 + 5$$

$$x_2 \geq \frac{3}{5}x_1 + 3$$

$$\frac{10}{20}$$

b)

$$\sigma(w_0 + w_1 x_1 + w_2 x_2) = \frac{1}{1 - \exp(-(w_0 + w_1 x_1 + w_2 x_2))}$$

$$\frac{5}{10}$$

- 30
30
4. A set of data points is generated by the following process: $Y = w_0 + w_1X + w_2X^2 + w_3X^3 + w_4X^4 + \epsilon$, where X is a real-valued random variable and ϵ is a Gaussian noise variable. You use two models to fit the data: $E[\hat{Y}] - Y$

Model 1: $Y = a_0 + a_1X + \epsilon$

Model 2: $Y = a_0 + a_1X + \dots + a_9X^9 + \epsilon$

- (a) Model 1, when compared to Model 2 using a fixed number of training examples, has a *bias* which is:

- ☐ Lower
- ☒ Higher
- ☐ The Same

- (b) Model 1, when compared to Model 2 using a fixed number of training examples, has a *variance* which is:

- ☒ Lower
- ☐ Higher
- ☐ The Same

- (c) Given 12 training examples, which model is more likely to *overfit* the data?

- ☐ Model 1
- ☒ Model 2