

Math of Big Data, Summer 2018

Prof: Gu

Name: Forest Kobayashi

HW #: 1

Day: Mon. Tue. Wed. Thu. Fri.

Date: 05/15/2018

No.	Points	Acknowledgments
1		
2		
Total		

This Assignment is (check one):



On Time



Late, without deduction



Late, with deduction

Comments: Feel free to work with other students, but make sure you write up the homework and code on your own (no copying homework *or* code; no pair programming). Feel free to ask students or instructors for help debugging code or whatever else, though.

The starter code for problem 2 part c and d can be found under the Resource tab on course website.

Note: You need to create a Github account for submission of the coding part of the homework. Please create a repository on Github to hold all your code and include your Github account username as part of the answer to problem 2.

Problem 1. (Linear Transformation)

Let $\mathbf{y} = A\mathbf{x} + \mathbf{b}$ be a random vector. Show that expectation is linear:

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[A\mathbf{x} + \mathbf{b}] = A\mathbb{E}[\mathbf{x}] + \mathbf{b}.$$

Also show that

$$\text{cov}[\mathbf{y}] = \text{cov}[A\mathbf{x} + \mathbf{b}] = A \text{cov}[\mathbf{x}] A^\top = A\Sigma A^\top.$$

Solution:

(a) We have

$$\begin{aligned}\mathbb{E}[\mathbf{y}] &= \mathbb{E}[A\mathbf{x} + \mathbf{b}] \\ &= \mathbb{E}[A\mathbf{x}] + \mathbb{E}[\mathbf{b}] \\ &= \mathbb{E}[A\mathbf{x}] + \mathbf{b}\end{aligned}$$

(we got from the second line to the third by the fact that $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for random variables X and Y , and from the second to the third by the fact that the expected value of a constant is the constant itself). It remains to show that $\mathbb{E}[A\mathbf{x}] = A\mathbb{E}[\mathbf{x}]$. Suppose \mathbf{x} is an n -dimensional random vector in a space X :

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

where each of the $x_i \in \mathbb{R}$. By definition, the expectation value of \mathbf{x} is given by

$$\mathbb{E}[\mathbf{x}] = \begin{bmatrix} \mathbb{E}[x_0] \\ \mathbb{E}[x_1] \\ \vdots \\ \mathbb{E}[x_n] \end{bmatrix}$$

and for any continuous random variable, we have

$$\mathbb{E}[Y] = \int_{-\infty}^{\infty} y f(y) \, dy$$

where $f(y)$ is some probability density function. For each of the x_i , define $f_i(x_i)$ to be the corresponding probability density function. Then

$$\mathbb{E}[\mathbf{x}] = \begin{bmatrix} \int_{-\infty}^{\infty} x_0 f_0(x_0) \, dx_0 \\ \int_{-\infty}^{\infty} x_1 f_1(x_1) \, dx_1 \\ \vdots \\ \int_{-\infty}^{\infty} x_n f_n(x_n) \, dx_n \end{bmatrix}$$

Let $A \in M_{m \times n}(\mathbb{R})$, with rows $\mathbf{a}_0, \mathbf{a}_1, \dots, \mathbf{a}_m$.

$$\begin{aligned} A\mathbf{x} &= \begin{bmatrix} a_{0,0} & a_{0,1} & \cdots & a_{0,n} \\ a_{1,0} & a_{1,1} & \cdots & a_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,0} & a_{m,1} & \cdots & a_{m,n} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= \begin{bmatrix} \langle \mathbf{a}_0, \mathbf{x} \rangle \\ \langle \mathbf{a}_1, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{a}_m, \mathbf{x} \rangle \end{bmatrix} \end{aligned}$$

Using the standard Euclidean inner product, we have

$$= \begin{bmatrix} \sum_{i=0}^n a_{0,i}x_i \\ \sum_{i=0}^n a_{1,i}x_i \\ \vdots \\ \sum_{i=0}^n a_{m,i}x_i \end{bmatrix}$$

Problem 2.

Given the dataset $\mathcal{D} = \{(x, y)\} = \{(0, 1), (2, 3), (3, 6), (4, 8)\}$

- (a) Find the least squares estimate $y = \boldsymbol{\theta}^\top \mathbf{x}$ by hand using Cramer's Rule.
 - (b) Use the normal equations to find the same solution and verify it is the same as part (a).
 - (c) Plot the data and the optimal linear fit you found.
 - (d) Find randomly generate 100 points near the line with white Gaussian noise and then compute the least squares estimate (using a computer). Verify that this new line is close to the original and plot the new dataset, the old line, and the new line.
-

Solution: