# Intro to Big Data Science: Assignment 1

Due Date: Mar 3, 2022

## Exercise 1 (Self-learning)

Log into "cookdata.cn", and enroll the course "数据科学导引". The online homework is given there. The system will judge your answers.

## Exercise 2 (Written Problem)

Given the ordered data  $\{x_{(i)}\}_{i=1}^{2n-1}$  with increasing order. Show that the median of the data set is equal to the minimizer of the following  $L^1$  minimization problem:

$$x_{(n)} = \arg\min_{c} \sum_{i=1}^{2n-1} |x_{(i)} - c|.$$

#### Exercise 3 (Written Problem)

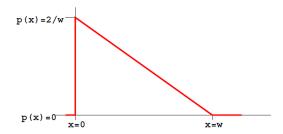
Consider the probability density function (PDF) shown in the following figure and equations:

$$p(x) = \begin{cases} 0, & \text{if } x < 0, \\ \frac{2}{w} - \frac{2x}{w^2}, & \text{if } 0 \le x \le w, \\ 0, & \text{if } w < x. \end{cases}$$

1. Which of the following expression is true? (Only one truth.)

(A) 
$$E[X] = \int_{-\infty}^{\infty} (\frac{2}{w} - \frac{2x}{w^2}) dx;$$

(B) 
$$E[X] = \int_{-\infty}^{\infty} x(\frac{2}{w} - \frac{2x}{w^2}) dx;$$



(C) 
$$E[X] = \int_{-\infty}^{\infty} w(\frac{2}{w} - \frac{2x}{w^2}) dx;$$

(D) 
$$E[X] = \int_0^w (\frac{2}{w} - \frac{2x}{w^2}) dx;$$

(E) 
$$E[X] = \int_0^w x(\frac{2}{w} - \frac{2x}{w^2}) dx;$$

(F) 
$$E[X] = \int_0^w w(\frac{2}{w} - \frac{2x}{w^2}) dx;$$

- 2. What is  $\mathbb{P}(x = 1 | w = 2)$ ?
- 3. When w = 2, what is p(1)?

### Exercise 4 (Written Problem)

Let X and Y be two continuous random variables. The conditional expectation of Y on X = x is defined as the expectation of Y with respect to the conditional probability density p(Y|X):

$$E(Y|X=x) = \int_{\mathscr{Y}} y p(y|X=x) dy = \frac{\int_{\mathscr{Y}} y p(x,y) dy}{p_x(x)},$$

where  $p_x(x)$  is the marginal probability density of Y. Show the following properties of the conditional expectation:

1.  $E_{p_y}Y = E_{p_x}[E(Y|X)]$ , where  $E_{p_y}$  means taking the expectation with respect to the marginal probability density  $p_y$ .

Remark: This formula is sometimes called the tower rule.

2. If *X* and *Y* are independent, then E(Y|X=x)=E(Y).

#### Exercise 5 (Written Problem)

The Jaccard distance between two sets A and B is defined as  $J_{\delta}(A,B)=1-\frac{|A\cap B|}{|A\cup B|}=\frac{|A\triangle B|}{|A\cup B|}$ , where |S| stands for the number of elements in the set S. Show that the Jaccard distance  $J_{\delta}$  is actually a distance, i.e., it satisfies the three properties:

- 1. Positivity:  $J_{\delta}(A, B) \ge 0$ , and "=" if and only if A = B;
- 2. Symmetry:  $J_{\delta}(A, B) = J_{\delta}(B, A)$ ;
- 3. Triangle inequality:  $J_{\delta}(A, B) \leq J_{\delta}(A, C) + J_{\delta}(B, C)$ .