

ESE 402/542 :Statistics for Data Science
Instructor: Hamed Hassani
Fall 2020

Midterm Examination

NAME	
------	--

Note: Each Multiple Choice (MC) question has 10 points, Problem 1 has 30 points, and Problem 2 has 40 points. For the multiple choice questions, you are only required to write the item that you've chosen (e.g. Question 1: (d) or Question 5: (e)).

Additional Information: If X is distributed according to the Gaussian distribution with mean μ and variance σ^2 , then $\mathbb{E}[X^2] = \mu^2 + \sigma^2$.

Also, for $\beta > 0$ and $n > 1$ we have $\int_0^\infty \frac{1}{(x+\beta)^n} dx = \frac{1}{(n-1)\beta^{n-1}}$.

	Grade (y/n)	Score	Max. Score
Multiple Choice			30
Problem 1			30
Problem 2			40
TOTAL			100

MC-Question 1: Assume we use the sample mean estimator to estimate the true mean of data that is generated i.i.d. from some distribution. Consider the following statements about the sample mean estimator:

- (1) It is always an unbiased estimator of the mean.
- (2) It is always a minimum-variance-unbiased-estimator of the mean.
- (3) It always achieves the Cramer-Rao bound.

Which of the above is wrong?

- (a) 1, 3
- (b) 2
- (c) 3
- d 2,3
- (e) 1,2

MC-Question 2: The formula that we derived in class for the confidence interval of the mean (using the sample mean estimator) has the variance of the data in it. However, we oftentimes do not know the true variance of the data. Hence, to fix the formula we need to:

- (a) Estimate the variance from data and further incorporate the error due to the estimation of the variance into the formula.
- (b) Just estimate the variance from data and use it instead of the true variance.
- (c) Use half of the data to estimate the variance and the other half for the confidence interval.
- (d) We need to know the true variance; Otherwise, the formula is highly inaccurate.

MC-Question 3: Which of the following is correct?

- 1. The maximum-likelihood estimator is unbiased.
- 2. The method-of-moments estimator is unbiased.
- 3. The maximum-likelihood estimate will become close to the true parameter as the number of samples grows.
- 4. The maximum-likelihood estimator always achieves the Cramer-Rao bound.

Problem 1. [30 pts] We have access to a data set X_1, X_2, \dots, X_n where X_i 's are generated i.i.d. according to a distribution with the following pdf:

$$f(x|a) = \frac{2a^2}{(x+a)^3} \mathbf{1}\{x \geq 0\},$$

where the parameter a is known to be positive.

(a) Use the method of moments to estimate the parameter a from data.

(b) Use the method of maximum likelihood to estimate a from data.

(Note that you can express the outcome of the estimator as the solution of some specific equation.)

Problem 2. [40 pts] We know that μ is a positive number but we do not know its value. We have access to a data set X_1, X_2, \dots, X_n where X_i 's are generated i.i.d. according to a distribution with the following pdf:

$$f(x|\mu) = \frac{1}{2} \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}} + \frac{1}{\sqrt{2\pi}} e^{-\frac{(x+\mu)^2}{2}} \right).$$

- (a) How does the pdf look like? Draw $f(x|\mu)$ as a function of x .
- (b) Find the mean and variance of the distribution in terms of μ . Note that you don't really need to do any integration here. To compute the variance, pay attention to the fact that the pdf is the average of two Gaussian pdfs: one with mean μ and the other with mean $-\mu$.
- (c) Design an unbiased estimator for the parameter $\theta = \mu^2$ using the sample data X_1, \dots, X_n .

(d) Compute the variance of your estimator.

(e) Provide a 95% confidence interval for the parameter θ using your estimator and the data sample (assume that the number of data points, n , is large).

(f) Can you comment on how good your estimator is? i.e. do you think you might be able to find another estimator with smaller variance?