

Mathematical Statistics - Section 1 - NYU Spring 2019 - Final exam

NAME:

- In order to get full credit, you must show your work / justify your answers. This usually means that you should use words alongside your computations.
- No justification is needed for the True/False questions, but pay close attention to the statement of the questions - some are tricky.

Some (approximate) values you can use, where X is a standard normal random variable

x	$\mathbb{P}(X \geq x)$ where $X \sim \mathcal{N}(0, 1)$
0	0.5
0.3	0.39
0.5	0.21
0.9	0.19
1.3	0.1
1.6	0.06
1.7	0.045
2.4	0.01
3.9	5/100 000

Some inequalities you can use:

- Markov's inequality: if X is a random variable with finite first moment, for any $t > 0$ we have

$$\mathbb{P}(|X| \geq t) \leq \frac{1}{t} \mathbb{E}[|X|].$$

- Bienaymé-Tchebychev's inequality: if X is a random variable with finite second moment, for any $t > 0$ we have

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq t) \leq \frac{1}{t^2} \text{Var}[X].$$

- Hölder's inequality (weak form): if X_1, \dots, X_m are i.i.d. Bernoulli random variables with parameter τ , we have

$$\mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m X_i - \tau\right| \geq \varepsilon\right) \leq 2 \exp(-2\varepsilon^2 m).$$

MLE and Fisher information (from HW5) In this exercise, we will denote by M the constant

$$M := \int_{-\infty}^{+\infty} e^{-x^4} dx,$$

which we will not try to compute.

We consider the family of probability distribution functions $x \mapsto f(x; \gamma)$ on $(-\infty, \infty)$ defined by

$$f(x; \gamma) := \frac{1}{M} \gamma e^{-\gamma^4 x^4},$$

where γ is some real number.

1. Let X_1, \dots, X_n be an observation, with true parameter γ_* . **Compute** the maximum likelihood estimator of γ_*

2. **Show** that the Fisher information of this statistical model can be written as

$$I(\gamma) = \frac{A}{\gamma^2},$$

for some constant A that we will not try to compute.

Testing the parameter of an exponential distribution (adapted from HW7)

Let X_1, \dots, X_n be identically distributed, independent random variables.

We assume that they are distributed according to the exponential distribution on $(0, +\infty)$ given by

$$x \mapsto f(x; \theta) := \theta e^{-\theta x},$$

for a certain value of the parameter $\theta > 0$.

We believe that $\theta = 1$, this is our null hypothesis H_0 .

1. **Build** a test statistic T , and an appropriate rejection region R such that

- The probability of making a type I error is bounded by 10%
- If $\theta \neq 1$, the probability that T is in R tends to 1 as $n \rightarrow \infty$.

2. If we remove either one of these two properties, why is the question much simpler?

CLT and applications

1. Let X_1, \dots, X_n be i.i.d. random variables such that $\mathbb{E}[X_1^2]$ is finite. **Compute**

$$\text{Var} \left(\frac{1}{n} \sum_{i=1}^n X_i \right).$$

2. **State** the central limit theorem.

3. Application: Two candidates are running for election (the one who gets the most votes is elected) in a country with a billion voters. After counting a million votes, the results are as follows:

$$A : 532\,000, B : 468\,000.$$

Who do you think is going to win? Justify your answer with a quantitative argument, e.g. a confidence interval, a p -value for a certain test...

Statistics Below is an excerpt of a data table from the “Demographic and Health Surveys”, concerning the mean height (in centimeters) for adult women in various countries. SD stands for (empirical) “standard deviation”. “Percent Urban“ is the percentage of the women in the sample who live in a city.

Country	Sample size	Empirical mean height	SD height	Percent Urban
Azerbaijan	5,412	158.4	5.9	52.9
Benin	11,015	159.3	6.5	40.3
Colombia	22,947	155.0	6.2	76.4

1. Do you think that the average height among all women in Azerbaijan is greater than the average height among all women in Colombia? **Justify** your answer with a quantitative argument, e.g. a confidence interval, a p -value for a certain test...

2. If we gather the data from these three samples, will the (empirical) standard deviation of heights be $\frac{5.9+6.5+6.2}{3} = 6.2$?

3. Can we say that living in cities causes people to grow taller?

True or false? Note: a statement that is *meaningless* should be treated as false.

TF1	TF2	TF3	TF4	TF5	TF6	TF7	TF8	TF9	TF10

- Let g be a quantity and \hat{g} be an estimator of g . If \hat{g} is consistent, then it is unbiased.
- Let $(X_i)_{i \geq 1}$ be random variables, such that $\mathbb{E}[X_i] = 10$ for all i . Then the following convergence holds

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n X_i \right] = 10.$$

- Let $(X_i)_{i \geq 1}$ be random variables, such that $\mathbb{E}[X_i] = 10$ for all i . Then the following convergence holds *in probability*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = 10.$$

- Let X, Y be two random variables, with real values. We assume that for every real number t , we have

$$\mathbb{P}[X \geq t, Y \geq 0] = \mathbb{P}[X \geq t] \mathbb{P}[Y \geq 0].$$

Then X and Y are independent.

- Let X_0 be some random variable, let N_1, N_2 be two random variables, and let

$$X_1 := X_0 + N_1, \quad X_2 := X_0 + N_2.$$

If N_1 and N_2 are independent, then X_1 and X_2 are independent.

- Let N_1, N_2 be two random variables, and let

$$X_1 := \frac{1}{2} + N_1, \quad X_2 := \frac{1}{4} + N_2.$$

If N_1 and N_2 are independent, then X_1 and X_2 are independent.

- The larger the p -value, the more confidence we have when rejecting the null hypothesis.
- A smaller variance for an unbiased estimator means a narrower confidence interval for the estimated quantity.
- The Fisher information is the best estimator we have for the empirical mean of a parametric predictor.
- Let X, Y be two standard normal random variables. Then

$$\text{Var}(10 + X) \geq \text{Var}(-5 + 2Y).$$

Linear regression through the origin Here is a data set of four data points:

$$(X_1, Y_1) = (1, 3), \quad (X_2, Y_2) = (2, 5), \quad (X_3, Y_3) = (3, 6), \quad (X_4, Y_4) = (5, 12).$$

Find the linear regression through the origin that minimizes the residual sum of squares. In other words, find the coefficient α such that

$$RSS := \sum_{i=1}^4 (Y_i - \alpha X_i)^2$$

is as small as possible.

About the empirical cdf (inspired by HW3) Let $x \mapsto f(x)$ be a continuous pdf on \mathbb{R} . Let X_1, \dots, X_n be i.i.d. random variables with common distribution f .

For any real numbers t, s with $t < s$, we let $G(t, s)$ be the quantity

$$G(t, s) := \int_t^s f(x) dx.$$

and we define the statistic $\hat{G}_n(t, s)$ as

$$\hat{G}_n(t, s) := \frac{\text{number of data points in } [t, s)}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{t \leq X_i < s}.$$

We recall that $\mathbf{1}_{t \leq X < s}$ is equal to 1 if $X \in [t, s)$ and to 0 if $X \notin [t, s)$. The variables $\mathbf{1}_{t \leq X_i < s}$ ($i = 1 \dots n$) are thus i.i.d. Bernoulli random variables.

1. For any $t < s$ fixed, **compute** the expectation and the variance of $\hat{G}_n(t, s)$.

2. **Prove** that $\hat{G}_n(t, s)$ is an unbiased, consistent estimator of $G(t, s)$. *Hint: you can rely on results that we have proved in class, or you can re-do the proof.*

Bonus question: pick the question you prefer (and answer it):

1. Continuation of the last exercise: **compute** the covariance of $\hat{G}_n(0, 1)$ and $\hat{G}_n(0, 2)$.
2. Three candidates A, B, C are running for election (the one who gets the most votes is elected) in a country with a million voters. After counting 10 000 votes, the results are as follows:

$$A : 3245, B : 4765, C : 1990.$$

Who is going to win?

3. Let f be some probability distribution function. We **believe** that it is an exponential distribution of parameter 1. **Explain** how you would design a test for this hypothesis.

