

ECOM90025 Advanced Data Analysis

Tutorial 2

Zheng Fan

The University of Melbourne

Introduction

Zheng Fan

- Ph.D. student in Economics
- Email: fan.z@unimelb.edu.au

Seek help?

- Ed discussion board
- Consultations: refer to Canvas for details

Shooting Stars

Consider a simple linear model

$$y_i = x_i\beta + \varepsilon_i,$$

where x 's dimension is $p = 100$ (hence β). The true value of β is a vector of 0's. Such a setting assumes all variables of x are noises. The error term has a standard normal distribution, $\varepsilon_i \sim N(0, 1)$.

Carry out a simulation study to learn how likely some variables are falsely classified to have a statistically significant impact on the output variable y .

- Significance level $\alpha = 0.05$. Assume a sample size of $n = 300$.
- Let each element x_{ik} for the i th observations and k th variable to be independent and randomly drawn from a standard normal distribution $N \sim N(0, 1)$.
- Discuss simulation results and their implication for empirical works.

Simulation as a Tool

Simulation is a cost-effective approach if you want a quick understanding of a model or a random variable. For example, if we do not know what a chi-square distribution looks like, we can simulate many variables from standard normal and square them to have some visual impression.

- 1 Draw $B = 10000$ standard normal random variables.
- 2 Square them.
- 3 Make a histogram and a kernel density plot for the simulated and transformed data.
- 4 What is the above distribution?

Order statistics

Definition (recap):

- Sample X_1, X_2, \dots, X_n
- Arrange them in increasing order:

$X_{(1)} =$ Smallest of the X_i

$X_{(2)} =$ 2nd smallest of the X_i

\vdots

$X_{(n)} =$ Largest of the X_i

- These are called the order statistics

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$

- $X_{(k)}$ is called the k th order statistic of the sample

Order statistics from the Uniform distribution

There are $p = 10$ uniformly random variables denoted by p_1, \dots, p_{10} on the interval $[0, 1]$. Use a simulation method to investigate the order statistics $p^{(1)}, \dots, p^{(10)}$ (ascendingly ordered).

- 1 What are the means of these statistics?
- 2 Draw these means via a scatter plot with a proper reference line.
- 3 Discuss this exercise's implication for significance tests.

Order statistics from the Uniform distribution

In theory:

- The pdf of k th order statistic is

$$g_k(x) = k \binom{n}{k} x^{k-1} (1-x)^{n-k}$$

There is a formula for this. Google it if you are interested

- This is a beta distribution,

$$F(X_{(k)}) \sim \text{Beta}(k, n - k + 1)$$

- So the theoretical mean is

$$\begin{aligned}\mathbb{E}(X_{(k)}) &= \frac{k}{n+1} \\ \text{mode}(X_{(k)}) &= \frac{k-1}{n-1}\end{aligned}$$

The end

Thanks for your attention!

