

MAST90104: Introduction to Statistical Learning

Week 3 Workshop and Lab

1. Suppose that X is a random variable with density function, f , given by

$$f(x) = \sum_{i=0}^{\infty} p(i)g(x; i)$$

where $p(0), p(1), \dots$ is a discrete probability mass function on $\{0, 1, \dots\}$ and each $g(x; i)$ is a probability density function. Suppose that $\mu(i), \sigma^2(i), M(t; i)$ are the mean, variance and moment generating function for the density $g(x; i)$. Let $M(t)$ be the moment generating function of X . Suppose also that N is a random variable with probability mass function $p(i), i = 0, 1, \dots$. Show that

- (a) $E(X) = E(\mu(N))$
- (b) $\text{var}(X) = E(\sigma^2(N)) + \text{var}(\mu(N))$
- (c) $M(t) = E(M(t; N))$.

(Hint: You may assume that interchange of infinite sums and integrals is justified.)

2. Section 2.3 states that the density function, f , of a non-central chi-square random variable, X , with degrees of freedom k and non-centrality parameter λ has the form of $f(x)$ in question 1. Further $g(x; i)$ is the central chi-square density with $k + 2i$ degrees of freedom and $p(0), p(1), \dots$ is the pmf of the Poisson distribution with parameter λ . Use the results in question 1 to show that:

- (a) $\text{var}(X) = 2k + 8\lambda$
- (b) $M(t) = (1 - 2t)^{-k/2} e^{2\lambda t / (1 - 2t)}$.

3. Use question 2 to prove Theorem 3.4.

4. Let $\mathbf{y} = \begin{bmatrix} y_1 & y_2 \end{bmatrix}^T$ be a normal random vector with mean and variance

$$\boldsymbol{\mu} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \quad V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Let

$$A = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad B = \frac{1}{2} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

- (a) Find the distributions of $\mathbf{y}^T A \mathbf{y}$ and $\mathbf{y}^T B \mathbf{y}$.
- (b) Are $\mathbf{y}^T A \mathbf{y}$ and $\mathbf{y}^T B \mathbf{y}$ independent?
- (c) What is the distribution of $\mathbf{y}^T A \mathbf{y} + \mathbf{y}^T B \mathbf{y}$?

5. Let y_1, \dots, y_n be an i.i.d. normal sample. Show that

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \text{ and } s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

are independent. (Hint: Express them as a random “vector” and quadratic form respectively.)

6. We model an individual’s income at age 30 against the number of years of formal education with a linear model. The following data is collected:

Years of formal education (x)	Income (\$k) (y)
8	8
12	15
14	16
16	20
16	25
20	40

Where possible, solve the following questions in two ways: using matrix calculations as detailed in the lectures, and using the `lm` command in R.

- Plot the data; is a linear model appropriate?
- Write down the linear model in matrix form.
- Find the normal equations for this model.
- Solve the normal equations to obtain the least squares estimates of the parameters. Add the fitted regression line to your plot (using `curve` for example).
- This model is a simple linear regression model. Use the standard linear regression formulae,

$$b_1 = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}, \quad b_0 = \bar{y} - b_1\bar{x},$$

to estimate the parameters again (where the bar indicates the mean). Check that you have the same answers as above.

- Calculate the sample variance s^2 .
- Estimate the average income of a person who has had 18 years of formal education.
- Calculate the standardised residuals, leverage, and Cook's distance for the first observation. You may need the R functions `rstandard`, `influence`, and `cooks.distance`. Check your numbers against the diagnostic plots produced by R.
- We know that the least squares estimator \mathbf{b} is an unbiased estimator for $\boldsymbol{\beta}$. Show that $\mathbf{t}^T \mathbf{b}$ is an unbiased estimator for $\mathbf{t}^T \boldsymbol{\beta}$, where \mathbf{t} is a vector of constants.

R exercises

Read Sections 5.1–5.3 of `spuRs` (which is available electronically in the library), then attempt the exercises below.

- The (Euclidean) length of a vector $v = (a_0, \dots, a_k)$ is the square root of the sum of squares of its coordinates, that is $\sqrt{a_0^2 + \dots + a_k^2}$. Write a function that returns the length of a vector.
- Last week you wrote a program to calculate $h(x, n)$, the sum of a finite geometric series. Turn this program into a *function* that takes two arguments, x and n , and returns $h(x, n)$. Make sure you deal with the case $x = 1$.
- In this question we simulate the rolling of a die. To do this we use the function `runif(1)`, which returns a 'random' number in the range (0,1). To get a random integer in the range $\{1, 2, 3, 4, 5, 6\}$, we use `ceiling(6*runif(1))`, or if you prefer, `sample(1:6, size=1)` will do the same job.
 - Suppose that you are playing the gambling game of the Chevalier de Méré. That is, you are betting that you get at least one six in four throws of a die. Write a program that simulates one round of this game and prints out whether you win or lose. Check that your program can produce a different result each time you run it.
 - Turn the program that you wrote in part (a) into a function `sixes`, which returns `TRUE` if you obtain at least one six in n rolls of a fair die, and returns `FALSE` otherwise. That is, the argument is the number of rolls n , and the value returned is `TRUE` if you get at least one six and `FALSE` otherwise. How would you give n the default value of 4?

- (c) Now write a program that uses your function `sixes` from part (b), to simulate N plays of the game (each time you bet that you get at least one six in n rolls of a fair die). Your program should then determine the proportion of times you win the bet. This proportion is an estimate of the *probability* of getting at least one six in n rolls of a fair die.

Run the program for $n = 4$ and $N = 100, 1000$, and 10000 , conducting several runs for each N value. How does the *variability* of your results depend on N ?

The probability of getting no 6's in n rolls of a fair die is $(5/6)^n$, so the probability of getting at least one is $1 - (5/6)^n$. Modify your program so that it calculates the theoretical probability as well as the simulation estimate and prints the difference between them. How does the *accuracy* of your results depend on N ?

- (d) In part (c), instead of processing the simulated runs as we go, suppose we first store the results of every game in a file, then later postprocess the results. You should read spuRs Chapter 4 to see how to read and write text files.

Write a program to write the result of all N runs to a textfile `sixes_sim.txt`, with the result of each run on a separate line. For example, the first few lines of the textfile could look like

```
TRUE
FALSE
FALSE
TRUE
FALSE
.
.
```

Now write another program to read the textfile `sixes_sim.txt` and again determine the proportion of bets won.

This method of saving simulation results to a file is particularly important when each simulation takes a very long time (hours or days), in which case it is good to have a record of your results in case of a system crash.