

Part I. (140 points) Do all calculations in $\LaTeX + R + knitr$. Insert computer text output and graphics to support what you are saying. For this assignment, all R code should well commented and be visible (`echo=TRUE`) in the document where you have written it.

(70pts) **1. Multinomial sampling**

Suppose $\underline{X} = \{X_1, X_2, \dots, X_k\}$ is a discrete random variable. It is easy to see that the joint distribution of $\underline{X} = \{X_1, X_2, \dots, X_k\}$ can be obtained through a series of conditional distributions:

$$\begin{aligned} \Pr[X_1 = x_1, X_2 = x_2, \dots, X_k = x_k] &= \Pr[X_k = x_k | X_1 = x_1, X_2 = x_2, \dots, X_{k-1} = x_{k-1}] \\ &\quad \times \Pr[X_1 = x_1, X_2 = x_2, \dots, X_{k-1} = x_{k-1}] \\ \Pr[X_1 = x_1, X_2 = x_2, \dots, X_{k-1} = x_{k-1}] &= \Pr[X_{k-1} = x_{k-1} | X_1 = x_1, X_2 = x_2, \dots, X_{k-2} = x_{k-2}] \\ &\quad \times \Pr[X_1 = x_1, X_2 = x_2, \dots, X_{k-2} = x_{k-2}] \\ &\quad \vdots \end{aligned}$$

That is,

$$\begin{aligned} \Pr[X_1 = x_1, X_2 = x_2, \dots, X_k = x_k] &= \Pr[X_1 = x_1] \\ &\quad \times \prod_{i=2}^k \Pr[X_i = x_i | X_1 = x_1, X_2 = x_2, \dots, X_{i-1} = x_{i-1}]. \end{aligned}$$

When $\underline{X} = \{X_1, X_2, \dots, X_k\}$ has a Multinomial($m, \underline{\theta}$) distribution, with $\underline{\theta} = (\theta_1, \dots, \theta_k)$, it can be shown that

$$\begin{aligned} X_1 &\sim \text{Binomial}(m, \theta_1) \\ X_2 | X_1 = x_1 &\sim \text{Binomial}(m - x_1, \theta_2^*), \\ &\quad \text{where } \theta_2^* = \theta_2 / (1 - \theta_1) \\ &\quad \vdots \\ X_i = x_i | X_1 = x_1, X_2 = x_2, \dots, X_{i-1} = x_{i-1} &\sim \text{Binomial}(m - (x_1 + \dots + x_{i-1}), \theta_i^*), \\ &\quad \text{where } \theta_i^* = \theta_i / (1 - \theta_1 - \dots - \theta_{i-1}) \\ &\quad \vdots \end{aligned}$$

To be precise here, you need to recognize that:

- If $x_1 + \dots + x_{i-1} = m$ (the total sample size), then

$$X_i = x_i | X_1 = x_1, X_2 = x_2, \dots, X_{i-1} = x_{i-1} \sim \text{Binomial}(m - m, \theta_i^*)$$

means $X_i = 0$ with probability 1. Similarly, $X_{i+1} = \dots = X_k = 0$ conditional on previous X_i s.

- For the last cell X_k , we are constrained so that $x_1 + \dots + x_k = m$ with probability 1 and the conditional distribution

$$X_k = x_k | X_1 = x_1, X_2 = x_2, \dots, X_{k-1} = x_{k-1} \sim \text{Binomial}(m - (x_1 + \dots + x_{k-1}), \theta_k^*),$$

where $\theta_k^* = \theta_k / (1 - \theta_1 - \dots - \theta_{k-1})$
 $= \theta_k / \theta_k = 1.$

This implies that given X_1, X_2, \dots, X_{k-1} , that $x_k = m - (x_1 + \dots + x_{k-1})$ with probability 1, that is, all “successes” in this Binomial distribution!

This characterization of the multinomial distribution is commonly used to generate random samples from a multinomial distribution, given a routine to generate binomial random variables.

- (a) (10 pts) For this problem, I want you to write a function that will generate a *single* multinomial sample given an input sample size m and a probability vector $\underline{\theta} = (\theta_1, \dots, \theta_k)$. Some thoughts:

- You may assume the input arguments are “permissible”.
- The algorithm is naturally programmed using a `for()` loop or a `while()` loop.
- You need to be concerned if, for example, $x_1 + \dots + x_{k-1} = m$ (it can’t go above it). If this condition holds, you need to break out of the loop and return, assuming all x_i s are initialized to zero. With this in mind, I believe the `while()` loop may be more transparent, because you can loop *while* this condition is not satisfied, given you first generate x_1 .
- As an aside, a function `f()` that returns the k -by-1 vector x can have the following structure:

```
## This function demonstrates a particular structure.
## It is not intended to run as it is.
```

```
f <- function( inputs ) {
  x <- rep(0, k)
  ...
  x[1] <- s
  ...
  for (i in 2:(k-1)) {
    x[i] <- something
    if ( condition ) {
      some commands
      return(x)
    }
  }
  x[k] <- something
  return(x)
}
```

The important point here is that whenever you have a `return()` statement in a function, the execution of the function ceases, and the current value of x is returned. If you choose to use a `for()` loop to generate a multinomial random variable x then your code will likely mimic this structure.

- You can generate a Binomial(m, p) random variable with `rbinom()`. Or, if you prefer, you can write your own Binomial generator.
- (b) (10 pts) Present some visual evidence that the function above “works correctly”. What you might choose to do is generate 100 or so samples $\underline{X} = \{X_1, X_2, \dots, X_k\}$ for a given m and $\underline{\theta} = (\theta_1, \dots, \theta_k)$. Noting that $X_i \sim \text{Binomial}(m, \theta_i)$, you might make a histogram of the generated counts X_i for the i th cell and compare the shape to that of the binomial pdf (using `pbinom()`) for $i = 1, \dots, k$. You might think of how to turn counts for the histogram into proportions and plot those with the binomial probabilities (together).

- (c) (20 pts) For each combination of

n = number of Multinomial samples = 500, 1000

m = Multinomial sample size = 25, 100, 500

k = number of cells = 3, 6, 9, 12, 15

compute n samples from a multinomial distribution with given m and $\theta_1 = \dots = \theta_k = k^{-1}$ (equal cell probabilities). Keeping a record of the total clock time needed to complete this task for a given n , m , and k . For each value of n , make a plot of the $\text{time}(m, k)$ values. For example, with time on the vertical axis and k on the horizontal axis, plot the times where each m is grouped and connected a line of different colors and line types (solid, dashed, etc.).

Discussed the results. You might also consider plotting $\text{time}(m, k)/n$ in the same manner (average time per sample).

Note that there are at least three ways to measure the clock time for a process:

- see the example at the bottom of help page for `?proc.time`
 - see package `rbenchmark`
 - see package `microbenchmark`
- (d) (10 pts) Consider the following pseudo-code, where $\underline{\theta} = (\theta_1, \dots, \theta_k)$ is a vector of probabilities with $\theta_i > 0$ and $\sum_i^k \theta_i = 1$. Suppose I do the following, given $\underline{\theta}$ is defined: First, define $\underline{\theta}$ as a 1-by- k vector and define scalar m , then define the following function:

```
f.dmultigen <- function(m, theta) {
  theta.c <- c(0, cumsum(theta))
  u <- runif(m)
  x <- hist(u, breaks = theta.c, plot = FALSE)$counts
  return(x)
}

# or, (probably) faster, all n samples at once
f.dmultigen2 <- function(theta, m, n = 1) {
  theta.c <- c(0, cumsum(theta))
  x <- t(apply(matrix(.bincode(runif(m * n), breaks = theta.c), nrow = n), 1, tabulate))
  return(x)
}
```

Explain what this code is doing step-by-step and argue why this algorithm generates a single Multinomial($m, \underline{\theta}$) sample.

- (e) (10 pts) Using the new code (or a modification of it) in part (d), repeat the analysis done in part (c) and compare the results. Which method, if either, is faster for these combinations of n , m , and k ?
- (f) (10 pts) Devise a good way of modifying the strategy in part (a) to generate n Multinomial samples per call to the function. That is, pass n as input to the function and return n samples.

Do the same for the strategy in part (d).

Does this modification impact $\text{time}(m, k)$, the time needed to generate n samples from $\text{Multinomial}(m, \underline{\theta})$? If so, how (you don't have to redo part (c), just consider a few cases)?

(40pts) **2. Multinomial hypothesis testing**

Suppose $\underline{X} = \{X_1, X_2, \dots, X_k\}$ has a $\text{Multinomial}(m, \underline{\theta})$ distribution, where $\underline{\theta} = (\theta_1, \dots, \theta_k)$ such that $\theta_i > 0$ and $\sum_i \theta_i = 1$. We wish to test $H_0 : \theta_1 = \theta_{01}, \dots, \theta_k = \theta_{0k}$, versus $H_1 : \text{not } H_0$. Two standard test statistics are the Pearson statistic

$$P = \sum_{i=1}^k \frac{(x_i - m\theta_{0i})^2}{m\theta_{0i}}$$

and the likelihood ratio statistic

$$G^2 = 2 \sum_{i=1}^k x_i \log_e \left(\frac{x_i}{m\theta_{0i}} \right),$$

where $0 \log_e(0) \equiv 0$. For large n , if H_0 is true, then both P and $G^2 \sim \chi_{k-1}^2$. A standard large-sample test is to reject H_0 based on the p-value

$$\begin{aligned} \text{p-value}(P) &= \Pr(\chi_{k-1}^2 > P_{\text{obs}}) && \text{or} \\ \text{p-value}(G^2) &= \Pr(\chi_{k-1}^2 > G_{\text{obs}}^2), \end{aligned}$$

that is, the area under the χ_{k-1}^2 curve to the right of observed values of P and G^2 .

- (a) (10 pts) Write separate functions to compute P and G^2 given vector inputs $\underline{X} = \{X_1, X_2, \dots, X_k\}$ and $\underline{\theta}_0 = (\theta_{01}, \dots, \theta_{0k})$.

- assuming input vectors are “permissible”, that is, $X_i \geq 0$ and $\theta_{0i} > 0$.
- calculations should be based on vector or matrix calculations, not `for()` loops
- for G^2 , since $x_i \log_e(x_i) \equiv 0$ when $x_i = 0$, you might think of identifying the elements of $\{X_1, X_2, \dots, X_k\}$ that are positive and basing the statistic solely on these elements. Alternatively, you might think of slightly changing the definition to have components $x_i \log_e((x_i + \varepsilon)/m\theta_{0i})$ where ε is a very small number, say $\varepsilon = 10^{-10}$.

- (b) (10 pts) Write a script where you can input the cell counts $\underline{X} = \{X_1, X_2, \dots, X_k\}$ and the null probabilities $\underline{\theta}_0 = (\theta_{01}, \dots, \theta_{0k})$. Have your script check whether the input values are “permissible”, that is, you need to input k non-negative counts $X_i \geq 0$ and k probabilities $\theta_{0i} > 0$ that sum to 1. Given permissible input, the

script should call your functions to compute P and G^2 , then print out the observed values of P and G^2 with their p-values. The output should include labels with the printing, that is, return summaries and a message to the screen/command window. You can use the `pchisq()` function to compute the χ_{k-1}^2 cdf.

- (c) (10 pts) The Dean of Arts and Sciences at a certain university established grading guidelines of 10% As and Fs, 20% Bs and Ds, and 40% Cs for his faculty. In a statistics class consisting of 117 students, the number of individuals receiving the five letter grades were as follows:

grade	A	B	C	D	F
number	16	50	31	11	9

Does it appear that the professor is following the Dean's recommendation, that is, does it appear that the grades are a random sample from a population with the recommended grade distribution? Use the code designed in the earlier part of the problem to answer this question.

- (d) (10 pts) Generalize your P and G^2 functions so that they can compute each statistic for multiple samples, given by rows of a matrix. Illustrate their use on the data in this table.

grade	A	B	C	D	F
Prof 1	16	50	31	11	9
Prof 2	10	23	22	20	7
Prof 3	21	10	42	3	1
Prof 4	3	12	31	16	0

(30pts) **3. Goodness-of-fit and upper-tail probability approximation**

For the goodness-of-fit test, a standard approach is to reject H_0 if

$$P \geq \chi_{k-1,1-\alpha}^2$$

where

$$\Pr[\chi_{k-1}^2 \geq \chi_{k-1,1-\alpha}^2] = \alpha,$$

that is, $\chi_{k-1,1-\alpha}^2$ is the upper α percentile of the χ_{k-1}^2 distribution, and α is the desired size of the test. A similar rule is used for G^2 . Because the χ_{k-1}^2 is only an approximation, neither

$$\Pr(P \geq \chi_{k-1}^2) \quad \text{nor} \quad \Pr(G^2 \geq \chi_{k-1}^2)$$

may not be exactly equal to α in small samples. This problem seeks to answer “how good is the approximation?” for a few settings. In practice, $\alpha = 0.01, 0.05, \text{ or } 0.10$, so let's restrict attention to these three cases. The `qchisq()` function can be used to compute the quantile

$$t = \chi_{k-1,1-\alpha}^2$$

for any choices of k and α .

I would like you to write a script that uses the crude MC estimate (using the same stream of random numbers) to calculate

$$\Pr(P \geq \chi_{k-1}^2) \quad \text{and} \quad \Pr(G^2 \geq \chi_{k-1}^2)$$

for the specified values of t (based on α and k), using the same k , m , and ℓ combinations considered in Problem 2, and the number n of simulated samples so that the margin-of-error (MOE) in the estimated probabilities does not exceed (approximately) 0.01.

Summarize the results in tabular form, giving values for α (nominal level), m , k , ℓ , plus the two crude MC estimates plus their estimated standard error (that is, the square root of the estimated variance). Also give the estimate for $(\Pr(P \geq t) - \Pr(G^2 \geq t))$ and its standard error. For these parameter combinations, does it appear that the upper-tail approximation to either P or G^2 is accurate? Discuss the results.