

Practice for Concepts in CASI Chapters 4 and 5

Try these questions with those around you. Important points from chapters 4 and 5 are covered.

1. Binomial Info

Find the Fisher information for a single observation drawn from a Binomial(n, p) distribution. Remember that the pmf takes the form: $P(X = x) = \binom{n}{x} (p)^x (1 - p)^{n-x}, x = 0, \dots, n, 0 < p < 1$.

Now, the computation above gave the theoretical (expected) information. The text suggests considering the observed Fisher information (evaluate based on the MLE for unknown parameters). What is the observed Fisher information here? How does this value relate to our typical confidence intervals for a proportion?

2. Randomization/Permutation Procedures

The Leukemia data set for gene 136 was used in the analysis below. The mean expression values for ALL and AML are being compared. Use the output to address the questions that follow.

```
leukemia_big <- read.csv("http://web.stanford.edu/~hastie/CASI_files/DATA/leukemia_big.csv")
gene136 <- t(leukemia_big[136, ]) #says pictures from row 136
type <- c(rep("ALL", 20), rep("AML", 14), rep("ALL", 27), rep("AML", 11))
leukemia <- data.frame(gene136, type)
leukemia <- rename(leukemia, gene136 = x136)
```

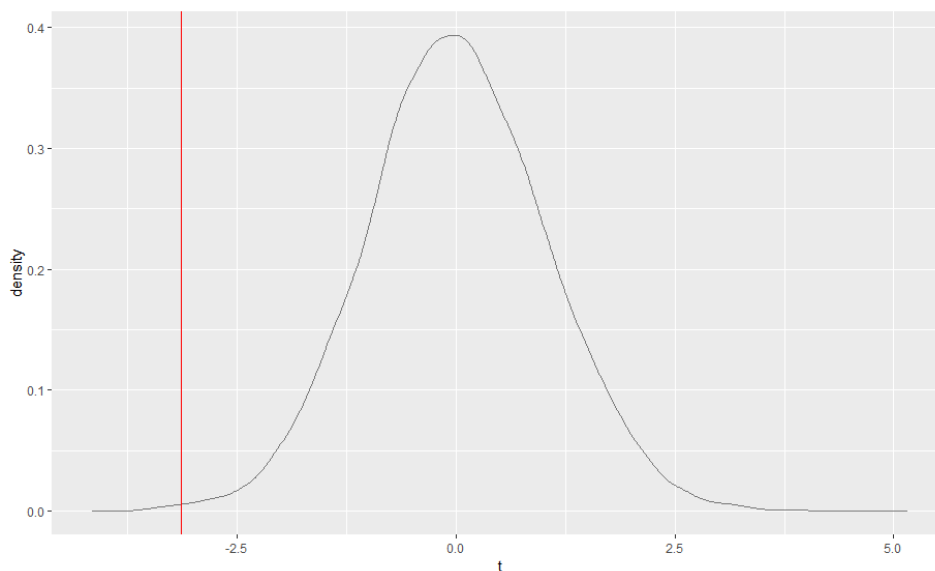
```
t.test(gene136 ~ type, data = leukemia)
```

welch Two Sample t-test

```
data: gene136 by type
t = -3.1323, df = 54.667, p-value = 0.002786
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.32386725 -0.07112012
sample estimates:
mean in group ALL mean in group AML
      0.7524794      0.9499731
```

```
set.seed(500) #make the results reproducible
ttest <- do(10000) * (t.test(gene136 ~ shuffle(type), data = leukemia)$statistic)
ttest <- as.data.frame(ttest)
```

```
gf_dens(~ t, data = ttest) %>% gf_vline(xintercept = -3.1323, color = "red")
```



```
pdata(~ t, -3.1323, data = ttest, lower.tail = TRUE)
[1] 0.0017
```

a. Explain what the t-test results on the original data suggest. (Assuming appropriate conditions hold).

b. There were concerns about the original t-test analysis. In your own words, explain what distribution is being created in the density plot, and how it can be used to perform an appropriate analysis.

c. What do you conclude from the randomization/permutation test results?

3. Chapter 5 Concepts

a. What explains the classical preference for parametric models?

b. Name two distributions that are part of an exponential family.

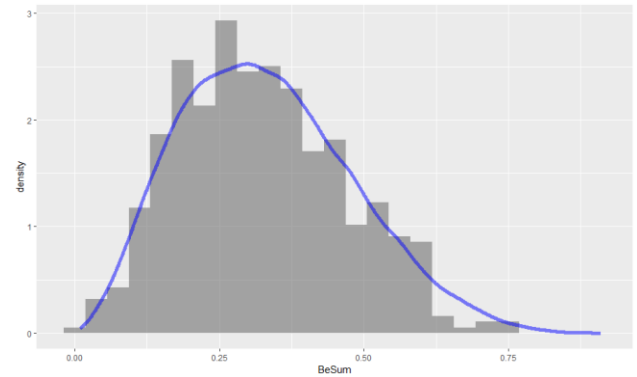
c. Name one of the two multivariate distributions examined in depth in the chapter.

4. Creating Betas - The textbook states that there are many relationships among the distributions listed in table 5.1. They give an example of creating a Beta from two independent Gamma RVs with the same second parameter. This is illustrated below:

```
set.seed(495)
Gam1 <- rgamma(500, 3, 5)
Gam2 <- rgamma(500, 6, 5)
BeSum <- Gam1/(Gam1+Gam2)
```

```
BeMatch <- rbeta(50000, 3, 6)
```

```
gf_dhistogram(~ BeSum) %>% gf_dens(~ BeMatch, color = "blue", size = 2)
```



Verify, using appropriate variable transformation techniques, that if X and Y are independent Gamma RVs with the same second parameter, then $W = X/(X+Y)$ follows a Beta distribution with parameters inherited from X and Y. The statement from the text is provided for your reference (CASI page 54). Hint: Jacobians.

<i>Gamma</i> $\text{Gam}(v, \sigma)$	$\frac{x^{v-1} e^{-x/\sigma}}{\sigma^v \Gamma(v)}$	$x \geq 0$	$v > 0$ $\sigma > 0$	$\frac{\sigma v}{\sigma^2 v}$
<i>Beta</i> $\text{Be}(v_1, v_2)$	$\frac{\Gamma(v_1 + v_2)}{\Gamma(v_1) \Gamma(v_2)} x^{v_1-1} (1-x)^{v_2-1}$	$0 \leq x \leq 1$	$v_1 > 0$ $v_2 > 0$	$\frac{v_1/(v_1 + v_2)}{(v_1 + v_2)^2 (v_1 + v_2 + 1)}$

Relationships abound among the table's families. For instance, independent gamma variables $\text{Gam}(v_1, \sigma)$ and $\text{Gam}(v_2, \sigma)$ yield a beta variate according to

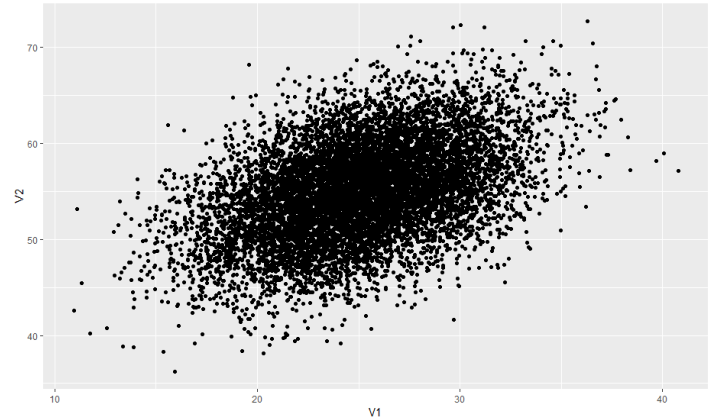
$$\text{Be}(v_1, v_2) \sim \frac{\text{Gam}(v_1, \sigma)}{\text{Gam}(v_1, \sigma) + \text{Gam}(v_2, \sigma)}. \quad (5.3)$$

Second page for #4. Hint: The sum $X+Y$ should have a Gamma distribution. This is verifiable by mgfs, but you can take it as fact. You should be able to use W and the sum in the Jacobian transformation process fairly easily, and then integrate out the sum from the resulting joint distribution. For the integration, the trick is to recognize it as a distribution you know, but with pieces missing. (Integrate without integrating.)

5. Multivariate Normals

Multivariate normals can be simulated in R using the `mvtnorm` package. You need a mean vector and the covariance matrix. An example in 2-D is provided.

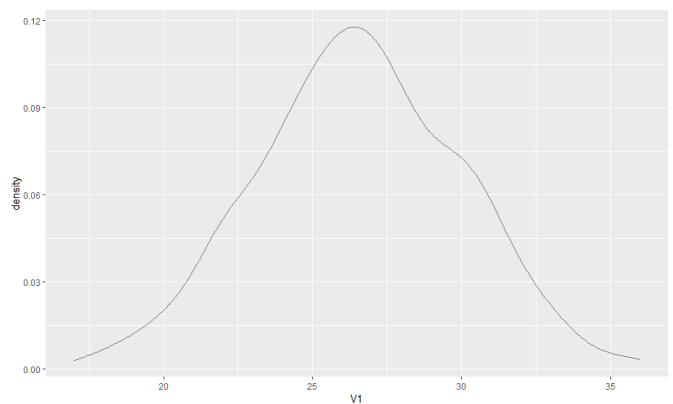
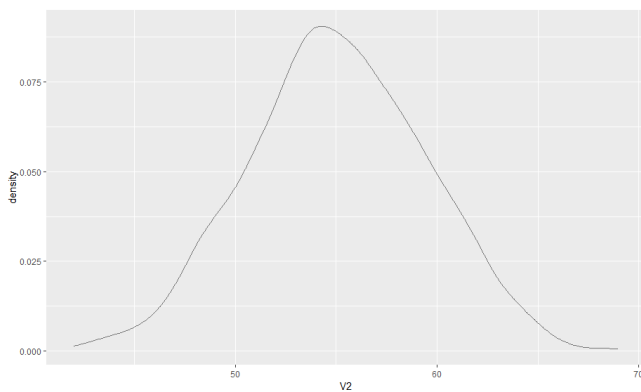
```
set.seed(52)
sigma <- matrix(c(16, 9, 9, 25), ncol=2)
NormData <- as.data.frame(rmvnorm(10000,
mean = c(25, 55), sigma = sigma))
gf_point(V2 ~ V1, data = NormData)
with(NormData, cor(V1, V2))
[1] 0.4474522
```



a. Find the theoretical correlation between V1 and V2 based on the provided covariance matrix. How well does it match the observed correlation? (Be sure you understand what each entry is in the covariance matrix, ask if unsure!)

b. The following code and plots illustrate a neat property of the Multivariate normal. What property is being shown?

```
rounded <- round(NormData, 0)
rounded %>% filter(V1 == 25) %>% gf_dens(~ V2)
rounded %>% filter(V2 == 60) %>% gf_dens(~ V1)
```



6. Multinomial Distributions

The second multivariate distribution covered in the chapter is the multinomial, a generalization of the binomial. There are multiple functions available in R to simulate from this distribution, depending on what you want. You can obtain a summary of the counts from a random draw of a certain number of objects from the distribution, or get the category numbers as data to use.

```
rmultinom(1, 40, c(0.1, 0.2, 0.3, 0.4)) #gives count summaries
      [,1]
[1,]     4
[2,]     9
[3,]     9
[4,]    18
```

A related function in the Hmisc package gives you a random draw of the category numbers as a vector.

```
m <- rMultinom(rbind(c(0.1, 0.2, 0.3, 0.4)), 40)
m #shows all forty entries are 1:4
t(apply(m, 1, table)/40)
      1      2      3      4
[1,] 0.05 0.225 0.35 0.375
```

a. Explain why X_1 = number of observations in category 1 and X_2 = number of observations in category 2 have a negative covariance/correlation in a multinomial setting.

b. What distribution has neat relationships to the multinomial as described in the text?

7. Exponential Families

The text presents exponential families via a formula relating any two densities in the family via a renormalized exponential tilt. There are other ways to recognize that a density is in an exponential family.

Suppose X is a random variable with density given by: $f(x | \theta) = a(\theta)b(x)\exp[c(\theta)d(x)]$, where $a()$ and $c()$ are functions only of the parameter θ , and $b()$ and $d()$ are functions of x (the data). If X 's density can be written in this form, then the density is in an exponential family.

For example, for the Bernoulli distribution, some re-writing enables us to see that:

$$f(x | p) = p^x(1-p)^{1-x} = (1-p)\left(\frac{p}{1-p}\right)^x = (1-p)\exp\left[x\log\left(\frac{p}{1-p}\right)\right].$$

Here, $a(p) = 1-p$, $b(x) = 1$, $c(p) = \log\left[\frac{p}{1-p}\right]$, $d(x) = x$. Hint: $a = \exp(\log(a))$.

If you have a random sample of observations from a distribution in this family, then the sample mean (or sum of the sample observations) is a sufficient statistic for the parameter(s). This is based on the form of $d()$, so it can differ from distribution to distribution.

Verify that the Poisson density can be written in this form.

Poisson
 $\text{Poi}(\mu)$

$$\frac{e^{-\mu}\mu^x}{x!}$$

Not all distributions we have seen exist in exponential families.

An example of a distribution in a *curved* exponential family is the _____ distribution.