

# Homework 5

## Estimation

Jaiden Atterbury

---

### Instructions

Please answer the following questions in the order in which they are posed. Add a few empty lines below each and write your answers there. **Focus on answering in complete sentences and show work whether we ask for it or not.** You will also need scratch paper/pen to work out the answers before typing it.

For help with formatting documents in RMarkdown, please consult R Markdown: The Definitive Guide. Another option is to search using Google.

---

### Exercises

1. (MIAA basketball) The `MIAA05` dataset from the **fastR2** package contains statistics on each of the 134 players from the 2004-2005 season. In this problem we will consider modeling the players' free throw shooting percentage. The variable is called `FTPct`. It is the ratio of free throws made (FT) to free throws attempted (FTA). Please type `?MIAA05` in the Console for a description of the variables in the dataset.

Since `FTPct` is a proportion, we will consider the Beta distribution:

$$f(x) = \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0) \Gamma(\beta_0)} x^{\alpha_0-1} (1-x)^{\beta_0-1} \quad 0 < x < 1.$$

- a. Show that the method of moments estimators of  $\alpha_0$  and  $\beta_0$  are:

$$\hat{\alpha}_0^{mom} = \bar{x} \left[ \frac{\bar{x} - s}{s - \bar{x}^2} \right],$$
$$\hat{\beta}_0^{mom} = \hat{\alpha}_0^{mom} \frac{1 - \bar{x}}{\bar{x}}$$

where  $s = \frac{1}{n} \sum_{i=1}^n x_i^2$ .

*Hint:* You can use without proof from homework 1 key:

$$E[X] = \frac{\alpha_0}{\alpha_0 + \beta_0}$$
$$E[X^2] = \frac{\alpha_0(\alpha_0 + 1)}{(\alpha_0 + \beta_0)(\alpha_0 + \beta_0 + 1)}.$$

**Solving  $E[X]$  for  $\alpha_0$ :**

$$\begin{aligned}
E[X] &= \bar{x} \\
\frac{\alpha_0}{\alpha_0 + \beta_0} &= \bar{x} \\
\alpha_0 &= \bar{x}(\alpha_0 + \beta_0) \\
\alpha_0 &= \bar{x}\alpha_0 + \bar{x}\beta_0 \\
\alpha_0 - \bar{x}\alpha_0 &= \bar{x}\beta_0 \\
\alpha_0(1 - \bar{x}) &= \bar{x}\beta_0 \\
\alpha_0 &= \frac{\bar{x}\beta_0}{1 - \bar{x}}
\end{aligned}$$

**Solving  $E[X^2]$  for  $\beta_0$ :**

$$\begin{aligned}
E[X^2] &= s \\
\frac{\alpha_0(\alpha_0 + 1)}{(\alpha_0 + \beta_0)(\alpha_0 + \beta_0 + 1)} &= s \\
\frac{\bar{x}(\alpha_0 + 1)}{\alpha_0 + \beta_0 + 1} &= s \quad \text{Since } \frac{\alpha_0}{\alpha_0 + \beta_0} = \bar{x} \\
\bar{x}(\alpha_0 + 1) &= s(\alpha_0 + \beta_0 + 1) \\
\bar{x}\alpha_0 + \bar{x} &= s\alpha_0 + s\beta_0 + s \\
s\beta_0 - (\bar{x} - s)\alpha_0 &= \bar{x} - s \\
s\beta_0 - \frac{(\bar{x} - s)\bar{x}\beta_0}{1 - \bar{x}} &= \bar{x} - s \quad \text{Since } \alpha_0 = \frac{\bar{x}\beta_0}{1 - \bar{x}} \\
\beta_0\left(s - \frac{(\bar{x} - s)\bar{x}}{1 - \bar{x}}\right) &= \bar{x} - s \\
\beta_0 &= \frac{\bar{x} - s}{\left(s - \frac{(\bar{x} - s)\bar{x}}{1 - \bar{x}}\right)} \\
\beta_0 &= \frac{\bar{x} - s}{\frac{s - \bar{x}^2}{1 - \bar{x}}} \\
\beta_0 &= \frac{(\bar{x} - s)(1 - \bar{x})}{s - \bar{x}^2}
\end{aligned}$$

**Putting it all together:**

$$\begin{aligned}
\hat{\alpha}_0^{mom} &= \frac{\bar{x}\beta_0}{1 - \bar{x}} \\
&= \frac{\bar{x}(\bar{x} - s)(1 - \bar{x})}{(1 - \bar{x})(s - \bar{x}^2)} \quad \text{Substituting in } \beta_0 \\
&= \bar{x} \left[ \frac{\bar{x} - s}{s - \bar{x}^2} \right]
\end{aligned}$$

$$\begin{aligned}
\hat{\beta}_0^{mom} &= \frac{(\bar{x} - s)(1 - \bar{x})}{s - \bar{x}^2} \cdot \frac{\bar{x}}{\bar{x}} \\
&= \frac{\bar{x}(\bar{x} - s)}{s - \bar{x}^2} \cdot \frac{1 - \bar{x}}{\bar{x}} \\
&= \hat{\alpha}_0^{mom} \frac{1 - \bar{x}}{\bar{x}}
\end{aligned}$$

The remaining two parts involve coding. Be sure to show the code and output, however, suppress warnings and messages.

- b. Make a histogram and a QQplot to assess the goodness of fit.

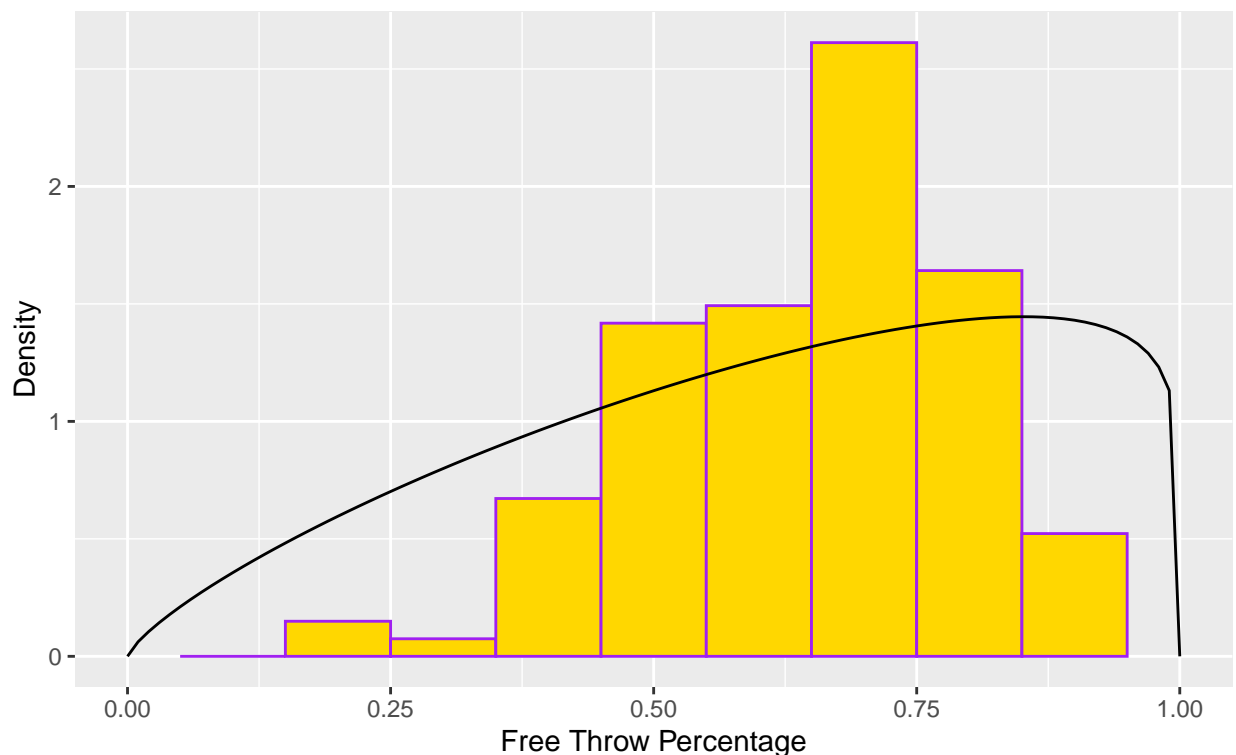
**Histogram:**

```
freethrow_df <- MIAA05
moments_df <- freethrow_df %>%
  summarise(xbar = mean(FTPct),
            sumsq = mean(FTPct^2),
            alpha_hat = xbar * (xbar - sumsq) / (sumsq - xbar^2),
            beta_hat = alpha_hat * (1 - xbar) / xbar)

ggplot(data = freethrow_df) +
  geom_histogram(mapping = aes(x = FTPct, y = ..density..),
                binwidth = 0.1,
                color = "purple",
                fill = "gold") +
  stat_function(fun = dbeta,
                args = list(shape1 = moments_df$alpha_hat,
                           shape2 = moments_df$beta_hat)) +
  xlim(c(0,1)) +
  labs(title = "Fitting a Beta Distribution",
       subtitle = "MIAA Free Throw Data",
       x = "Free Throw Percentage",
       y = "Density")
```

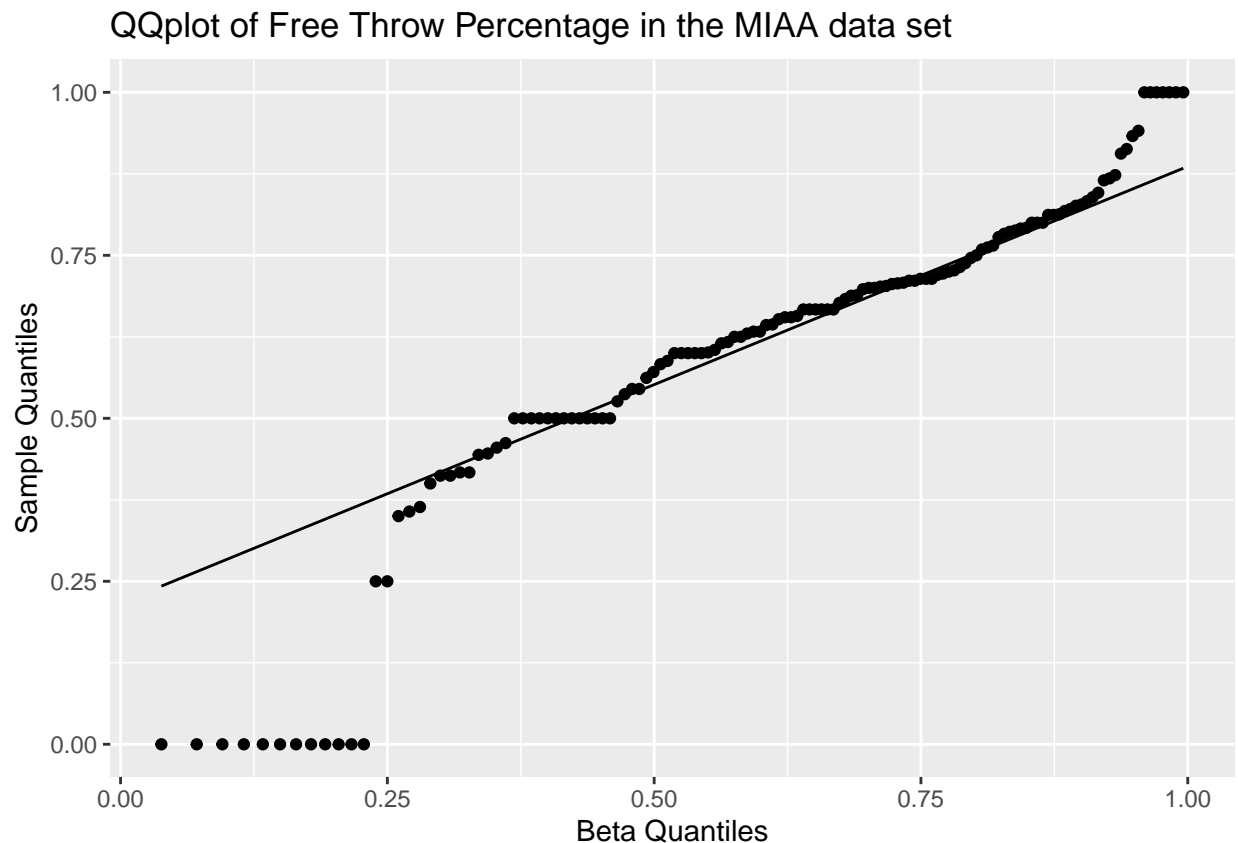
## Fitting a Beta Distribution

### MIAA Free Throw Data



### QQplot:

```
ggplot(data = freethrow_df,
       mapping = aes(sample = FTPct)) +
  stat_qq(distribution = qbeta,
         dparams = list(shape1 = moments_df$alpha_hat,
                        shape2 = moments_df$beta_hat)) +
  stat_qq_line(distribution = qbeta,
              dparams = list(shape1 = moments_df$alpha_hat,
                             shape2 = moments_df$beta_hat)) +
  labs(title = "QQplot of Free Throw Percentage in the MIAA data set",
       y = "Sample Quantiles",
       x = "Beta Quantiles",)
```



### Assessing the fit:

As we can see from the QQplot and the histogram of free throw percentages with the Beta distribution overlaid, this Beta distribution is not a good fit for the given data. In particular, we can see that the beta distribution underestimates  $x$  values between 0.4 and 0.8 while severely overestimating  $x$  values near the lower tail, while severely underestimating near the upper tail. For more confirmation that this isn't a good fit, we can look at the QQPlot and notice that there is no clear linear relationship between the sample and beta quantiles and the over/under estimating at the tails is also very apparent.

- c. Are there any players you should remove from the data before attempting the analysis? Decide on an elimination rule, apply it, and repeat the analysis. Do you like this fit better?

### Elimination rule:

```
summ_stats <- freethrow_df %>%
  summarise(mean = mean(FTA),
            median = median(FTA),
            min = min(FTA),
            max = max(FTA),
            q1 = quantile(FTA, .25),
            q3 = quantile(FTA, .75))

print(summ_stats)
```

```
##   mean median min max   q1   q3
## 1    30   20.5   0 191 7.25 44.25
```

As we can see from the above QQplot and histogram, there are many data points at the low and high end of the free throw percentage spectrum, which seems to be “ruining” our linear relationship in the QQplot. One reason that this is happening is because these players at the end of the spectrum have very few free throw attempts if any at all. Thus, their free throw percentages will be unreasonably high or low and will not match their theoretical “true” free throw percentage that we would observe if they were to shoot hundreds or thousands of free throws. After taking a look at the 5 number summary plus mean of the FTA variable (Free throws attempted) I want to choose a value of FTA that is high enough to start reaching a players true free throw percentage, while at the same time not getting rid of too many observations. With that in mind, we will take the conservative route and remove any player that hasn’t attempted 10 or more free throws during the season.

**New histogram:**

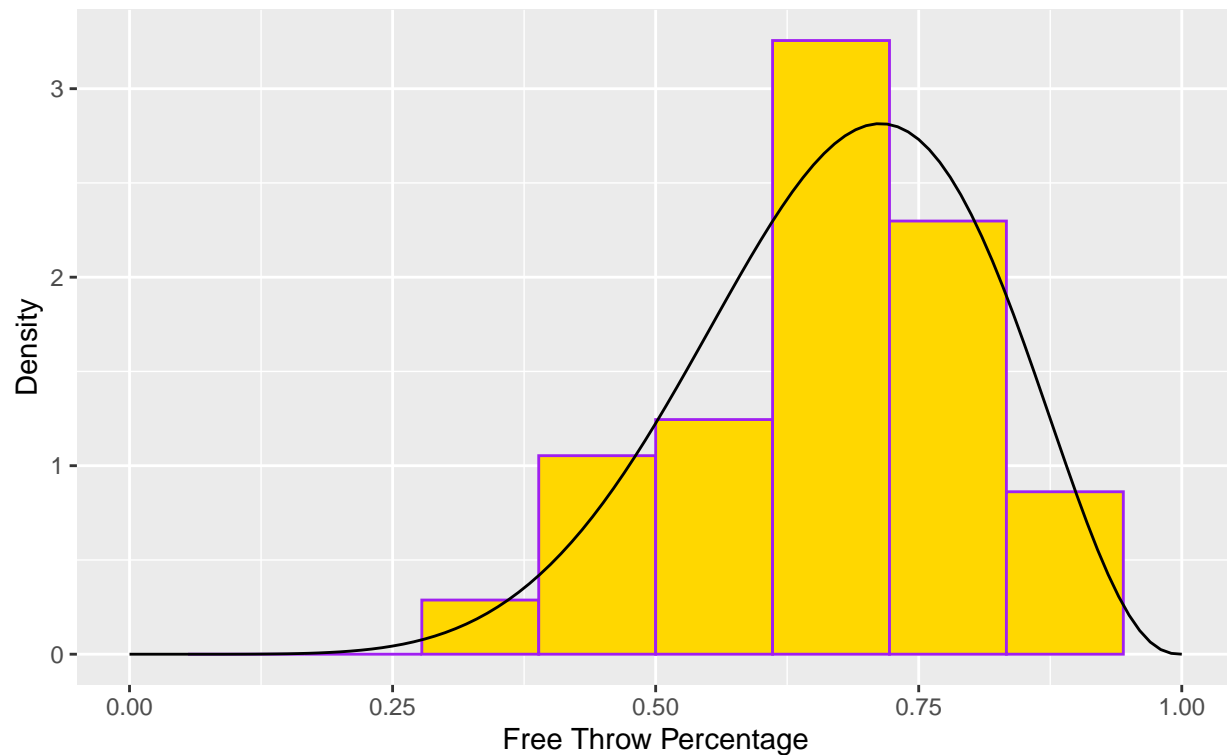
```
freethrow_updated <- freethrow_df %>%
  filter(FTA >= 10)

momentsup_df <- freethrow_updated %>%
  summarise(xbar_up = mean(FTPct),
            sumsq_up = mean(FTPct^2),
            alpha_hat_up = xbar_up * (xbar_up - sumsq_up) / (sumsq_up - xbar_up^2),
            beta_hat_up = alpha_hat_up * (1 - xbar_up) / xbar_up)

ggplot(data = freethrow_updated) +
  geom_histogram(mapping = aes(x = FTPct, y = ..density..),
                binwidth = 1/9,
                color = "purple",
                fill = "gold") +
  stat_function(fun = dbeta,
                args = list(shape1 = momentsup_df$alpha_hat_up,
                           shape2 = momentsup_df$beta_hat_up)) +
  xlim(c(0,1)) +
  labs(title = "Fitting a Beta Distribution",
       subtitle = "MIAA Free Throw Data (With FTA > 9)",
       x = "Free Throw Percentage",
       y = "Density")
```

## Fitting a Beta Distribution

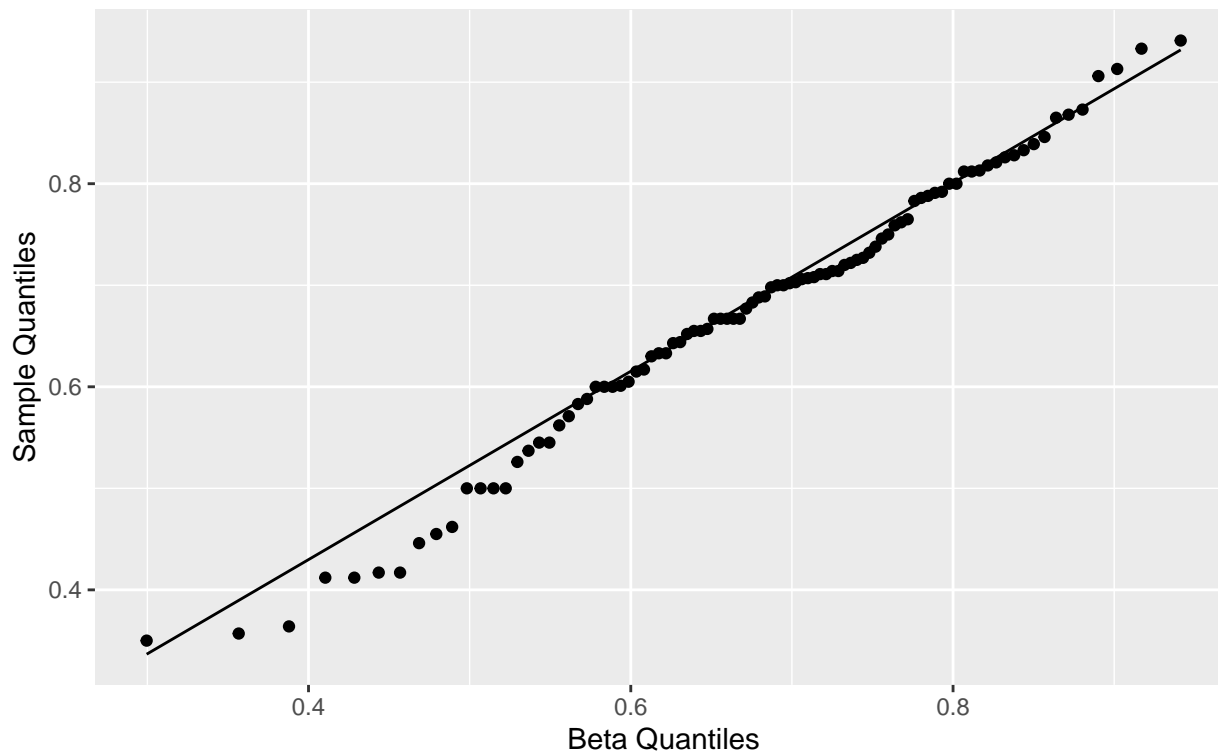
MIAA Free Throw Data (With FTA > 9)



New QQplot:

```
ggplot(data = freethrow_updated,
  mapping = aes(sample = FTPct)) +
  stat_qq(distribution = qbeta,
    dparams = list(shape1 = momentsup_df$alpha_hat,
      shape2 = momentsup_df$beta_hat)) +
  stat_qq_line(distribution = qbeta,
    dparams = list(shape1 = momentsup_df$alpha_hat_up,
      shape2 = momentsup_df$beta_hat_up)) +
  labs(title = "QQplot of Free Throw Percentage in the MIAA data set",
    subtitle = "Only Including Players With FTA > 9",
    y = "Sample Quantiles",
    x = "Beta Quantiles",)
```

QQplot of Free Throw Percentage in the MIAA data set  
Only Including Players With FTA > 9



#### Analysis:

With this new elimination rule in use we can see from the QQplot and the histogram with the beta distribution overlaid that this beta distribution fits the data much better now than it did before. We can see that more of the density under the beta distribution is better predicted and the QQplot shows a more linear relationship than the one before. Although this fit is far from perfect, the new beta distribution models the data as we'd expect. In particular, the model has little to no densities at the tails, while holding most of the density between 0.4 and 0.8. This matches what we'd expect from experienced basketball players as most players don't have free throw percentages below 40%, while at the same time it is very rare to see a player with a free throw percentage above 90%.

2. (Unbias your estimator) Let  $X \sim \text{Binom}(n, \pi_0)$ .

a. Show, with justification, that

$$E\left[\frac{X}{n} \left(1 - \frac{X}{n}\right)\right] = \frac{(n-1)\pi_0(1-\pi_0)}{n}$$

#### Setup:

Since we will need to know what  $E[X^2]$  is equal to in order to solve this problem, we will start by finding  $E[X^2]$  using the variance formula. If  $X \sim \text{Binom}(n, \pi_0)$ , then  $\text{Var}[X] = n\pi_0(1-\pi_0)$  and  $E[X] = n\pi_0$ . Thus it follows that:

$$\begin{aligned} \text{Var}[X] &= E[X^2] - (E[X])^2 \\ n\pi_0(1-\pi_0) &= E[X^2] - (n\pi_0)^2 \\ E[X^2] &= n\pi_0(1-\pi_0) + n^2\pi_0^2 \end{aligned}$$

**Solving  $E\left[\frac{X}{n}\left(1 - \frac{X}{n}\right)\right]$ :**

$$\begin{aligned}
 E\left[\frac{X}{n}\left(1 - \frac{X}{n}\right)\right] &= E\left[\frac{X}{n} - \frac{X^2}{n^2}\right] \\
 &= \frac{1}{n}E[X] - \frac{1}{n^2}E[X^2] \quad (\text{Linearity of expectation}) \\
 &= \frac{n\pi_0}{n} - \frac{n\pi_0(1 - \pi_0) + n^2\pi_0^2}{n^2} \\
 &= \pi_0 - \frac{\pi_0(1 - \pi_0) + n\pi_0^2}{n} \\
 &= \frac{n\pi_0 - \pi_0 + \pi_0^2 - n\pi_0^2}{n} \\
 &= \frac{\pi_0(n - 1 + \pi_0 - n\pi_0)}{n} \\
 &= \frac{(n - 1)\pi_0(1 - \pi_0)}{n}
 \end{aligned}$$

b. Suppose we want an unbiased estimator for  $\pi(1 - \pi)$ . Use your answer from (a) to construct such an estimator.

If we use an estimator of  $\frac{X}{n-1}\left(1 - \frac{X}{n}\right)$  we will get an unbiased estimate of  $\pi(1 - \pi)$ . To see this, we will show that if  $X \sim \text{Binom}(n, \pi_0)$ , then  $E\left[\frac{X}{n-1}\left(1 - \frac{X}{n}\right)\right] = \pi(1 - \pi)$ .

**Showing  $\frac{X}{n-1}\left(1 - \frac{X}{n}\right)$  is an unbiased estimator:**

$$\begin{aligned}
 E\left[\frac{X}{n-1}\left(1 - \frac{X}{n}\right)\right] &= E\left[\frac{n}{n-1}\left(\frac{X}{n}\left(1 - \frac{X}{n}\right)\right)\right] \\
 &= \frac{n}{n-1}E\left[\frac{X}{n}\left(1 - \frac{X}{n}\right)\right] \quad (\text{Linearity of expectation}) \\
 &= \frac{n}{n-1} \cdot \frac{(n-1)\pi_0(1 - \pi_0)}{n} \quad (\text{HW5 Problem 2a}) \\
 &= \pi_0(1 - \pi_0)
 \end{aligned}$$

Therefore,  $\frac{X}{n-1}\left(1 - \frac{X}{n}\right)$  is an unbiased estimator of  $\pi(1 - \pi)$ .

3. (Bayes estimator) Let  $X_1, X_2, \dots, X_n$  be independent Bernoulli random variables drawn from the PMF:

$$f(x) = \begin{cases} (1 - \pi_0) & x = 0 \\ \pi_0 & x = 1 \end{cases}$$

Consider the Bayesian estimator of  $\pi_0$ <sup>1</sup>:

$$\hat{\pi}_0^{\text{bayes}} = \frac{X + 1}{n + 2}$$

where  $X = X_1 + X_2 + \dots + X_n$  is  $\text{Binom}(n, \pi_0)$ .

a. Is  $\hat{\pi}_0^{\text{bayes}}$  an unbiased estimator of  $\pi_0$ ? If not, is it asymptotically unbiased?

**Setup:**

In order to find  $E\left[\frac{X+1}{n+2}\right]$ , we need to find  $E[X_i]$  which represents the expected value of an individual Bernoulli random variable. If  $X \sim \text{Bernoulli}(\pi_0)$ , then it follows that  $E[X] = \sum_{x=0}^1 x(1 - \pi_0)^{1-x}\pi_0^x = \pi_0$ .

---

<sup>1</sup>don't worry about how to derive this estimator



**Finding  $E[\hat{\pi}_0^{bayes}]$ :**

$$\begin{aligned}
E[\hat{\pi}_0^{bayes}] &= E\left[\frac{X+1}{n+2}\right] \\
&= \frac{1}{n+2}E[X] + \frac{1}{n+2} \quad (\text{Linearity of expectation}) \\
&= \frac{1}{n+2}E\left[\sum_{i=1}^n X_i\right] + \frac{1}{n+2} \\
&= \frac{1}{n+2} \sum_{i=1}^n E[X_i] + \frac{1}{n+2} \quad (\text{Linearity of expectation}) \\
&= \frac{n\pi_0}{n+2} + \frac{1}{n+2} \\
&= \frac{n\pi_0 + 1}{n+2}
\end{aligned}$$

Therefore, we can see that  $\hat{\pi}_0^{bayes}$  is not an unbiased estimator for  $\pi_0$ . However, if we remember infinite limit rules from calculus 1 we know that if  $f(x) = \frac{ax^m + \dots}{bx^n + \dots}$  and  $m = n$ , then  $\lim_{x \rightarrow \infty} f(x) = \frac{a}{b}$ . Thus we can see that  $\lim_{n \rightarrow \infty} \frac{n\pi_0 + 1}{n+2} = \frac{\pi_0}{1} = \pi_0$  and therefore Bayes estimator is asymptotically unbiased.

b. Is  $\hat{\pi}_0^{bayes}$  a consistent estimator?

**Setup:**

In order to find  $Var[\frac{X+1}{n+2}]$ , we need to find  $Var[X_i]$  which represents the variance of an individual Bernoulli random variable. If  $X \sim \text{Bernoulli}(\pi_0)$ , then it follows that  $Var[X] = E[X^2] - (E[X])^2 = \sum_{x=0}^1 x^2(1 - \pi_0)^{1-x}\pi_0^x - \pi_0^2 = \pi_0(1 - \pi_0)$ .

**Finding  $Var[\hat{\pi}_0^{bayes}]$ :**

$$\begin{aligned}
Var[\hat{\pi}_0^{bayes}] &= Var\left[\frac{X+1}{n+2}\right] \\
&= \frac{1}{(n+2)^2}Var[X] \quad (\text{Non-linearity of variance}) \\
&= \frac{1}{(n+2)^2}Var\left[\sum_{i=1}^n X_i\right] \\
&= \frac{1}{(n+2)^2} \sum_{i=1}^n Var[X_i] \quad (\text{Since the } X_i\text{'s are independent}) \\
&= \frac{n\pi_0(1 - \pi_0)}{n^2 + 4n + 4}
\end{aligned}$$

Again, if we remember infinite limit rules from calculus 1 we know that if  $f(x) = \frac{ax^m + \dots}{bx^n + \dots}$  and  $m < n$ , then  $\lim_{x \rightarrow \infty} f(x) = 0$ . Thus we can see that  $\lim_{n \rightarrow \infty} \frac{n\pi_0(1 - \pi_0)}{n^2 + 4n + 4} = 0$  and therefore Bayes estimator is consistent.

c. Based on the sample 1, 0, 1, 0, 1, calculate the value of  $\hat{\pi}_0^{bayes}$  for this sample. Also calculate its estimated standard error.

**Calculating  $\hat{\pi}_0^{bayes}$  from the sample:**

$$\begin{aligned}
\hat{\pi}_0^{bayes} &= \frac{x+1}{n+2} \\
&= \frac{(1+0+1+0+1)+1}{5+2} \\
&= \frac{4}{7}
\end{aligned}$$

Calculating the standard error of  $\hat{\pi}_0^{\text{bayes}}$ :

$$\begin{aligned}
\hat{SD}[\hat{\pi}_0^{\text{bayes}}] &= \hat{SD}\left[\frac{X+1}{n+2}\right] \\
&= \sqrt{\frac{n\hat{\pi}_0(1-\hat{\pi}_0)}{n^2+4n+4}} \quad (\text{Since } \text{Var}[\pi_0] = \frac{n\pi_0(1-\pi_0)}{n^2+4n+4}) \\
&= \sqrt{\frac{5(\frac{4}{7})(1-\frac{4}{7})}{5^2+4(5)+4}} \\
&= \sqrt{\frac{\frac{60}{49}}{49}} \\
&= \frac{\sqrt{60}}{49}
\end{aligned}$$

4. (Bias variance trade-off) In this problem we will continue working with the model described in problem 3. Our focus will be on comparing the mean square error (MSE) of  $\hat{\pi}_0^{\text{bayes}}$  with the MSE of

$$\hat{\pi}_0^{\text{mom}} = \frac{X}{n}.$$

- a. Give expressions for the MSE of each estimator. Show your work clearly.

To find  $MSE[\hat{\pi}_0^{\text{bayes}}]$ , we will need to follow the  $MSE$  formula which states  $MSE[\theta_0] = \text{bias}^2 + \text{Var}[\theta_0]$  where  $\text{bias} = E[\hat{\theta}_0] - \theta_0$ . Thus in order to find  $MSE[\hat{\pi}_0^{\text{bayes}}]$  we will need to find  $\text{Bias}[\hat{\pi}_0^{\text{bayes}}]$  and use the fact that  $E[\hat{\pi}_0^{\text{bayes}}] = \frac{n\pi_0+1}{n+2}$  and  $\text{Var}[\hat{\pi}_0^{\text{bayes}}] = \frac{n\pi_0(1-\pi_0)}{(n+2)^2}$  which we found in problem 3.

**Finding an expression for  $\text{Bias}[\hat{\pi}_0^{\text{bayes}}]$ :**

$$\begin{aligned}
\text{Bias}[\hat{\pi}_0^{\text{bayes}}] &= E[\hat{\pi}_0^{\text{bayes}}] - \pi_0 \\
&= \frac{n\pi_0+1}{n+2} - \pi_0 \\
&= \frac{n\pi_0+1-(n+2)\pi_0}{n+2} \\
&= \frac{n\pi_0+1-n\pi_0-2\pi_0}{n+2} \\
&= \frac{1-2\pi_0}{n+2}
\end{aligned}$$

**Finding an expression for  $MSE[\hat{\pi}_0^{\text{bayes}}]$ :**

$$\begin{aligned}
MSE[\hat{\pi}_0^{\text{bayes}}] &= \text{bias}^2 + \text{Var}[\hat{\pi}_0^{\text{bayes}}] \\
&= \left(\frac{1-2\pi_0}{n+2}\right)^2 + \frac{n\pi_0(1-\pi_0)}{(n+2)^2} \\
&= \frac{(1-2\pi_0)^2}{(n+2)^2} + \frac{n\pi_0(1-\pi_0)}{(n+2)^2} \\
&= \frac{(1-2\pi_0)^2 + n\pi_0(1-\pi_0)}{(n+2)^2} \\
&= \frac{1-4\pi_0+4\pi_0^2+n\pi_0-n\pi_0^2}{(n+2)^2} \\
&= \frac{(4-n)\pi_0^2+(n-4)\pi_0+1}{n^2+4n+4}
\end{aligned}$$

To find  $MSE[\hat{\pi}_0^{mom}]$ , where  $\hat{\pi}_0^{mom} = \frac{X}{n}$  we will need to follow the  $MSE$  formula, as we did for the previous estimator, which states  $MSE[\theta_0] = \text{bias}^2 + \text{Var}[\theta_0]$  where  $\text{bias} = E[\hat{\theta}_0] - \theta_0$ . Thus in order to find  $MSE[MSE[\hat{\pi}_0^{mom}]]$  we will need to find  $E[\hat{\pi}_0^{mom}]$ ,  $\text{Var}[\hat{\pi}_0^{mom}]$ , and  $\text{Bias}[\hat{\pi}_0^{mom}]$ .

**Finding an expression for  $E[\hat{\pi}_0^{mom}]$ :**

$$\begin{aligned}
 E[\hat{\pi}_0^{mom}] &= E\left[\frac{X}{n}\right] \\
 &= E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\
 &= \frac{1}{n} \sum_{i=1}^n E[X_i] \quad (\text{Linearity of expectation}) \\
 &= \frac{1}{n} \cdot n \cdot \pi_0 \quad (\text{Since } X_i \sim \text{Bernoulli}(\pi_0)) \\
 &= \pi_0 \\
 \therefore \hat{\pi}_0^{mom} &\text{ is an unbiased estimator}
 \end{aligned}$$

**Finding an expression for  $\text{Var}[\hat{\pi}_0^{mom}]$ :**

$$\begin{aligned}
 \text{Var}[\hat{\pi}_0^{mom}] &= \text{Var}\left[\frac{X}{n}\right] \\
 &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\
 &= \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n X_i\right] \quad (\text{Non-linearity of variance}) \\
 &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] \quad (\text{Since the } X_i \text{'s are independent}) \\
 &= \frac{1}{n^2} \cdot n \cdot \pi_0(1 - \pi_0) \quad (\text{Since } X_i \sim \text{Bernoulli}(\pi_0)) \\
 &= \frac{\pi_0(1 - \pi_0)}{n} \\
 \therefore \hat{\pi}_0^{mom} &\text{ is a consistent estimator} \quad (\text{Since } \lim_{n \rightarrow \infty} \frac{\pi_0(1 - \pi_0)}{n} = 0)
 \end{aligned}$$

**Finding an expression for  $\text{Bias}[\hat{\pi}_0^{mom}]$ :**

$$\begin{aligned}
 \text{Bias}[\hat{\pi}_0^{mom}] &= E[\hat{\pi}_0^{mom}] - \pi_0 \\
 &= \pi_0 - \pi_0 \\
 &= 0 \quad (\text{As we'd expect since } \hat{\pi}_0^{mom} \text{ is an unbiased estimator})
 \end{aligned}$$

**Finding an expression for  $MSE[\hat{\pi}_0^{mom}]$ :**

$$\begin{aligned}
 MSE[\hat{\pi}_0^{mom}] &= \text{bias}^2 + \text{Var}[\hat{\pi}_0^{mom}] \\
 &= 0^2 + \frac{\pi_0(1 - \pi_0)}{n} \\
 &= \frac{\pi_0(1 - \pi_0)}{n}
 \end{aligned}$$

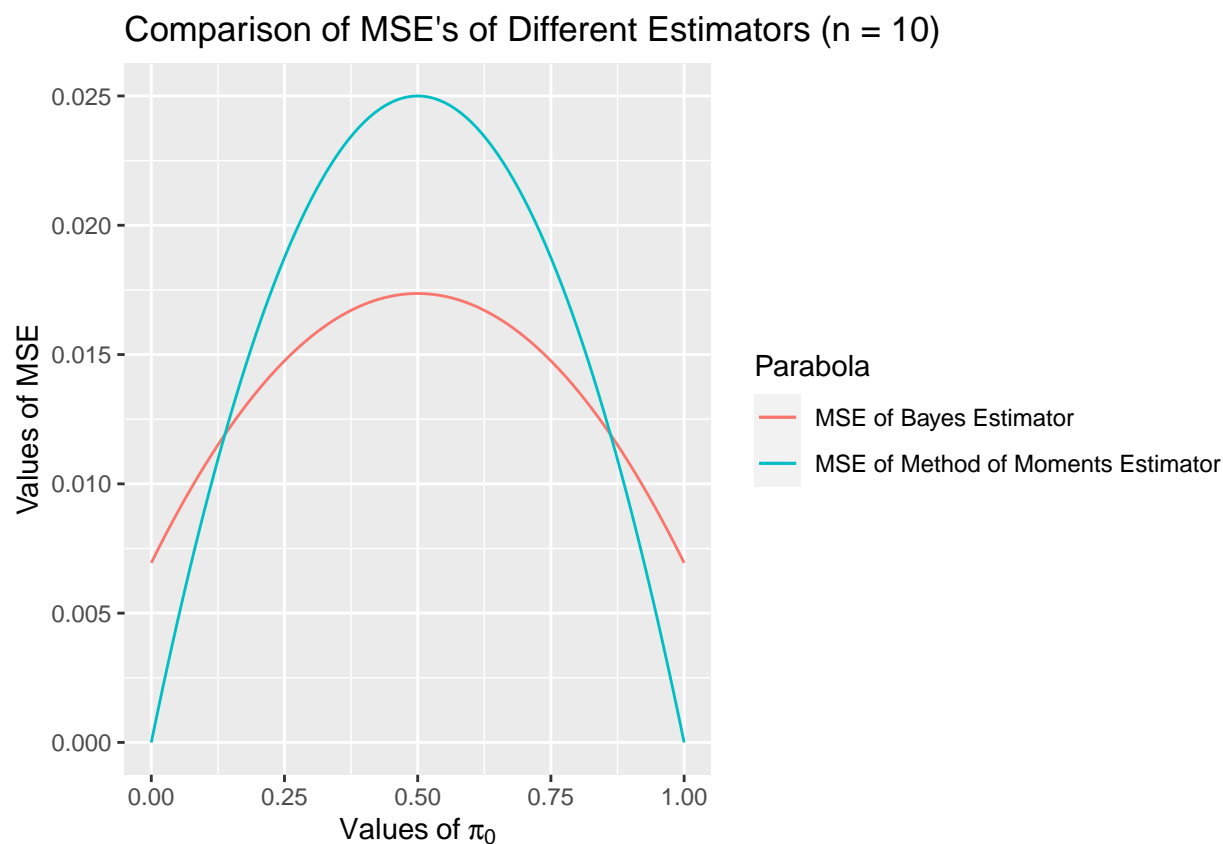
- b. For  $n = 10$ , plot the MSE of both estimators on the same graph as a function of  $\pi_0$ . Describe (just visually) the  $\pi_0$  values for which the MSE is smaller for the Bayes estimator. Repeat for  $n = 1,000$ .

Don't forget to label the plot, and make it easy for the reader to know which estimator is being represented by which curve.

*This example illustrates the bias variance tradeoff. Sometimes a biased estimator does better than an unbiased one, if it has a smaller variance!*

#### Plot of MSE's for n = 10:

```
ggplot() +
  geom_function(fun = function(x){(-6*x^2 + 6*x + 1) / 144},
               xlim = c(0, 1),
               mapping = aes(color = "MSE of Bayes Estimator")) +
  geom_function(fun = function(x){(x * (1 - x)) / 10},
               xlim = c(0, 1),
               mapping = aes(color = "MSE of Method of Moments Estimator")) +
  labs(title = "Comparison of MSE's of Different Estimators (n = 10)",
       x = expression(Values ~ of ~ pi[0]),
       y = "Values of MSE",
       color = "Parabola")
```



#### Analysis:

As can be seen from the graph, the  $MSE$  of  $\hat{\pi}_0^{bayes}$  is smaller than  $\hat{\pi}_0^{mom}$  for  $\pi_0$  values such that  $0.136 < \pi_0 < 0.862$ . Which means even though  $\hat{\pi}_0^{mom}$  is unbiased, the  $MSE$  for  $\hat{\pi}_0^{mom}$  is only smaller for “extreme” values of  $\pi_0$ .

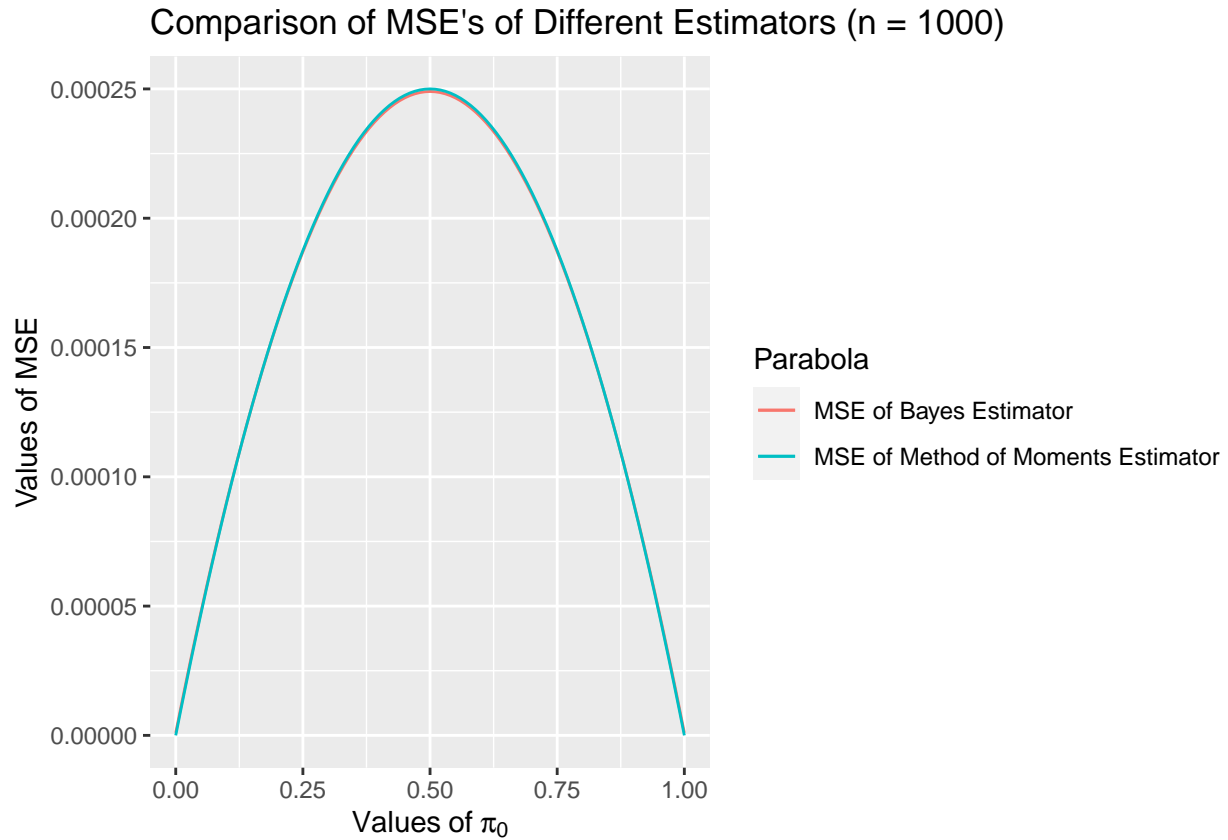
#### Plot of MSE's for n = 1000:

```
ggplot() +
  geom_function(fun = function(x){(-996*x^2 + 996*x + 1) / 1004004},
```

```

xlim = c(0, 1),
mapping = aes(color = "MSE of Bayes Estimator")) +
geom_function(fun = function(x){(x * (1 - x)) / 1000},
xlim = c(0, 1),
mapping = aes(color = "MSE of Method of Moments Estimator")) +
labs(title = "Comparison of MSE's of Different Estimators (n = 1000)",
x = expression(Values ~ of ~ pi[0]),
y = "Values of MSE",
color = "Parabola")

```



#### Analysis:

As can be seen from the graph above, there is barely any difference between the  $MSE$  curves of  $\hat{\pi}_0^{bayes}$  and  $\hat{\pi}_0^{mom}$ . This is because both estimators are consistent, and their variance approaches 0 as  $n$  approaches infinity, thus their  $MSE$  also approaches 0 as  $n$  approaches infinity. Since  $n = 1000$  is a large sample, both estimators have a nearly identical  $MSE$  curve and have  $MSE$  values of almost 0 for all values of  $\pi_0$ . For completeness, the values of  $\pi_0$  in which the  $MSE$  of  $\hat{\pi}_0^{bayes}$  is less than  $\hat{\pi}_0^{mom}$  are  $0.146 < \pi_0 < 0.854$ .