

MA 354: Data Analysis – Fall 2021 – Due 10/8 at 5p
Homework 2:

Complete the following opportunities to use what we've talked about in class. These questions will be graded for correctness, communication and succinctness. Ensure you show your work and explain your logic in a legible and refined submission.

The starting jobs will be applied in alphabetical order (last name) for question two.

1. **Solver:** provide a solution, if possible, and reasoning for the solution. **Due to group 10/5 or earlier.**
2. **Code Checker:** provides a first check of the solver's worked solutions and ensures they are correct with a solid interpretation.
3. **Checker** checks the solution for completeness, proposes and implements changes if agreed upon by the group. Provides a short paragraph summarizing the discussion of proposals and their reason for acceptance or non-acceptance.
4. **Double Checker** checks the solution for completeness, communication and polish. The Double Checker ensures that the solution is correct and highly polished for submission.

For subsequent questions student roles will move down one position. The rolls change as follows.

1. **Solver \Rightarrow Code Checker**
2. **Code Checker \Rightarrow Checker**
3. **Checker \Rightarrow Double Checker**
4. **Double Checker \Rightarrow Solver**

While students have assigned jobs for each question I encourage students to help each other complete the homework in collaboration.

1. Select a continuous distribution (Not the uniform or exponential). It does not have to be one that we cover in the notes! To explore the PDF of your distribution, specify two sets of parameter(s) for your distribution.

- (a) **History** Discuss what types of random variables are modeled with your distribution. Be sure to include a discussion about the support and ensure to provide the density function, and CDF. This requires some internet research – what’s the history of the distribution, why was it created and named? What are some exciting applications of this distribution?

Cite all of your sources in LaTeX by adding a BibTeX citation to the .bib file. To help, I’ve cited R (R Core Team, 2021) in parentheses here. R Core Team (2021) provides helpful tools for the rest of the questions below. BibTeX citations are available through Google Scholar by clicking the cite button below the article of interest and selecting the BibTeX option.

Solution:

Our group has chosen to use a normal distribution for the Q1 and Q2!

To start with, I will define what a normal distribution is. A continuous variable X has a normal distribution — with mean μ and σ^2 — $X \sim N(\mu, \sigma^2)$, if it has the following properties:

$$\mu \in \mathbb{R}; \sigma \in \mathbb{R}^+ \quad (\text{Parameters})$$

$$X = \{x : x \in \mathbb{R}\} \quad (\text{Support})$$

$$f_X(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (\text{PDF})$$

$$F_X(x|\mu, \sigma) = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (\text{CDF})$$

Description: Normal distribution is one of the most important distributions in the world of statistics, and it’s used in many fields. The graph of this probability function is a symmetric, bell-shaped curved. Symmetry of this function implies that its mean is going to be equal to its median value. We are going to prove it later on.

History: One of the first major great discoveries in the field of statistics happened in 1713, as Jacob Bernoulli published his work on proving the Weak Law of Large Numbers. According to Patel and Read (1996), initially normal distribution appeared in 1733 as an approximation to the probability for sums of binomially distributed quantities to lie between two values.

Today: According to Ahsanullah et al. (2014), normal distribution plays an important role in many applied problems in biology, economics, engineering, genetics, hydrology, mechanics, medicine, number theory, statistics, physics, psychology and so on.

- (b) Show that you have a valid PDF. You will find the `integrate()` function in R helpful.

Solution:

Here’s the equation of PDF for normal distribution!

$$f_X(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (\text{PDF})$$

Now, let’s prove that it’s valid!

Since normal (or gaussian — whatever you prefer!) distribution is a continuous probability distribution, it implies that the area under the curve is equal to 1 (or 100%!). This statement takes its roots from the second Kolmogorov axiom that states that the entirety of sample space is equal to one.

Therefore, let’s use CDF to prove that our PDF formula is going to return one! Since the support for normal distribution contain all rational numbers, our integral is going to be from negative infinity to positive infinity.

$$F_X(x|\mu, \sigma) = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (\text{CDF})$$

I am going to use `integrate()` function in order to compute this equation.

```
mean<-1 #mu
sd<-3 #sigma
pdf <- function(x){
```

```

    (1/(sd*sqrt(2*pi)))*exp(-(x-mean)^2)/(2*sd^2)) #our PDF
  }
  integrate(pdf, -Inf, Inf) #CDF for total area under the curve
## 1 with absolute error < 2.1e-07

```

- (c) Find the median for your two sets of parameter(s). Conduct some research to find the median based on our PDF to confirm that your numerical approach is correct.

Solution:

As we have established in part (a) of this problem, one of the key features of the normal distribution is the fact that it's symmetrical. Let's set out to prove it through the direct proof!

Let m be median!

$$P(X \leq m) = P(X \geq m) = \frac{1}{2} \quad (\text{Definition of the median})$$

Let normal distribution be symmetric. If it's symmetric, then the statement $\mu = m$ holds true.

Therefore, the following equation should also hold true:

$$\int_{-\infty}^{\mu} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \int_{\mu}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{2} \quad (\text{Assumption})$$

Let's check it through R!

```

mean<-1 #mu
sd<-3 #sigma
func <- function(x){
  (1/(sd*sqrt(2*pi)))*exp(-(x-mean)^2)/(2*sd^2))
}
firstPart <- integrate(func, -Inf, mean)
secondPart <- integrate(func, mean, Inf)

(firstPart$value)

## [1] 0.5

(firstPart$value==secondPart$value)

## [1] TRUE

```

It would appear that $\mu = m$! Therefore, mean is equal to median within the normal distribution. Therefore, since $\mu_1 = 1$ and $\mu_2 = -2$, the medians are 1 and -2 respectively!

- (d) Graph the PDF for several values of the parameter(s) including the two sets you specified. What does changing the parameter(s) do to the shape of the PDF?

Solution:

Let's take various values and plot different PDFs! I am going to use the ggplot2 (Wickham, 2016) library for it!

```
library(ggplot2)
plot.df <- data.frame(
  x=seq(-15, 15, 0.001),
  f1=dnorm(x=seq(-15, 15, 0.001), mean=1, sd=1),
  f2=dnorm(x=seq(-15, 15, 0.001), mean=-2, sd=4),
  f3=dnorm(x=seq(-15, 15, 0.001), mean=0, sd=5),
  f4=dnorm(x=seq(-15, 15, 0.001), mean=3, sd=3)
)

ggplot(plot.df, aes(x=x))+
  geom_line(aes(y=f1, color="m=1, sd=1"))+
  geom_line(aes(y=f2, color="m=-2, sd=4"))+
  geom_line(aes(y=f3, color="m=0, sd=5"))+
  geom_line(aes(y=f4, color="m=3, sd=3"))+
  theme_bw()
```

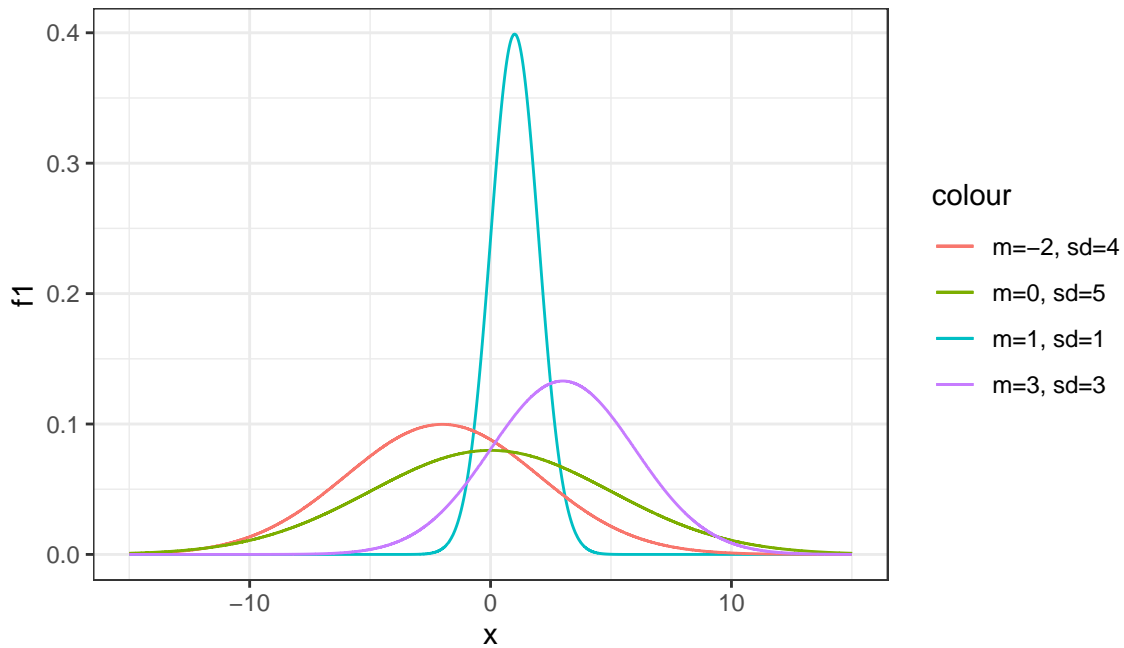


Figure 1: Gaussian PDF with various sets of parameters

As we can see on Figure 1, changing our parameters does indeed change the graph we have. Changing the μ transforms graph to the right (if the mean is positive) and to the left (if the mean is negative). Changing σ changes the kurtosis of the graph. Thus, distributions with lower value of σ show higher peaks.

- (e) Graph the CDF for the same values of the parameter(s) as you did in Question 1d. What does changing the parameter(s) do to the shape of the CDF? Comment on the aspects of the CDFs that show that the CDF is valid.

Solution:

```
plot.df <- data.frame(
  x=seq(-15, 15, 0.001),
  f1=pnorm(q=seq(-15, 15, 0.001), mean=1, sd=1),
  f2=pnorm(q=seq(-15, 15, 0.001), mean=-2, sd=4),
  f3=pnorm(q=seq(-15, 15, 0.001), mean=0, sd=5),
  f4=pnorm(q=seq(-15, 15, 0.001), mean=3, sd=3)
)

ggplot(plot.df, aes(x=x))+
  geom_line(aes(y=f1, color="m=1, sd=1"))+
  geom_line(aes(y=f2, color="m=-2, sd=4"))+
  geom_line(aes(y=f3, color="m=0, sd=5"))+
  geom_line(aes(y=f4, color="m=3, sd=3"))+
  theme_bw()
```

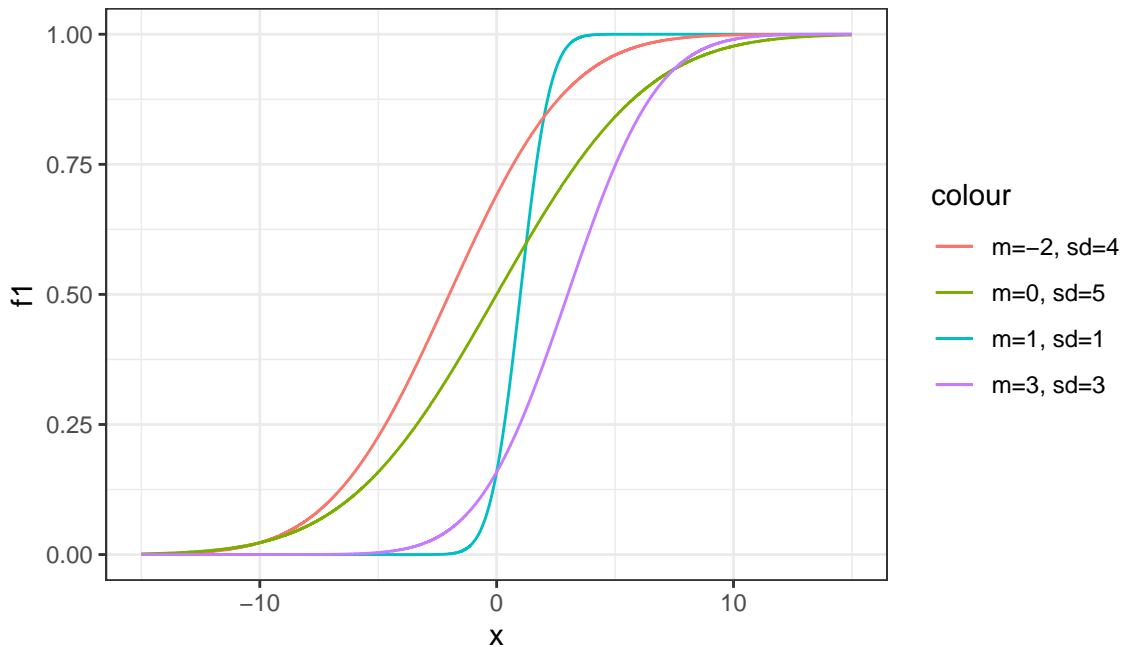


Figure 2: Gaussian CDF with various sets of parameters

We see that Figure 2 is a graph of CDF function with different sets of parameters. This function shows the probability that a probability that a random variable X will take a value less than or equal to x . We can see that our PDF is valid, for when we went over all values of X , the probability is 1.00 or 100%. Changing the value of μ changes the position of a graph (moves it to the right or to the left), while the change in σ increases the slope of the function: the higher the value of σ , the faster the function will go over all the values and reach 100%.

- (f) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). In a 4×2 grid, plot a histogram of each set of data and superimpose the true density function at the specified parameter values. Interpret the results.

Solution:

```
library(patchwork)
library(tidyverse)

sample.df <- list(x1=rnorm(10, mean=1, sd=1),
                  x2=rnorm(25, mean=1, sd=1),
                  x3=rnorm(100, mean=1, sd=1),
                  x4=rnorm(1000, mean=1, sd=1),
                  y1=rnorm(10, mean=-2, sd=4),
```

```

y2=rnorm(25, mean=-2, sd=4),
y3=rnorm(100, mean=-2, sd=4),
y4=rnorm(1000, mean=-2, sd=4))

buildingPlot <- function(source, mu, sigma){
  df <- data.frame(value=source) #turning values from the list into a df
  colnames(df) <- "value" #changing the name of the column

  answer<-ggplot(df, aes(value))+
    geom_histogram(aes(y=..density..), bins=10,
                   color="black")+ #building a histogram
    geom_function(fun=dnorm, args = list(mean = mu, sd = sigma),
                 color="red")+
    theme_bw() +#superimposing the function
    labs(x="Value", y="Density")
  answer
}

x1<-buildingPlot(sample.df[1], 1, 1)+labs(title="Sample=10",
                                           subtitle = "mean=1, sd=1")
x2<-buildingPlot(sample.df[2], 1, 1)+labs(title="Sample=25",
                                           subtitle = "mean=1, sd=1")
x3<-buildingPlot(sample.df[3], 1, 1)+labs(title="Sample=100",
                                           subtitle = "mean=1, sd=1")
x4<-buildingPlot(sample.df[4], 1, 1)+labs(title="Sample=1000",
                                           subtitle = "mean=1, sd=1")

y1<-buildingPlot(sample.df[5], -2, 4)+labs(title="Sample=10",
                                           subtitle = "mean=-2, sd=4")
y2<-buildingPlot(sample.df[6], -2, 4)+labs(title="Sample=25",
                                           subtitle = "mean=-2, sd=4")
y3<-buildingPlot(sample.df[7], -2, 4)+labs(title="Sample=100",
                                           subtitle = "mean=-2, sd=4")
y4<-buildingPlot(sample.df[8], -2, 4)+labs(title="Sample=1000",
                                           subtitle = "mean=-2, sd=4")

(x1+y1)/(x2+y2)/(x3+y3)/(x4+y4)

```

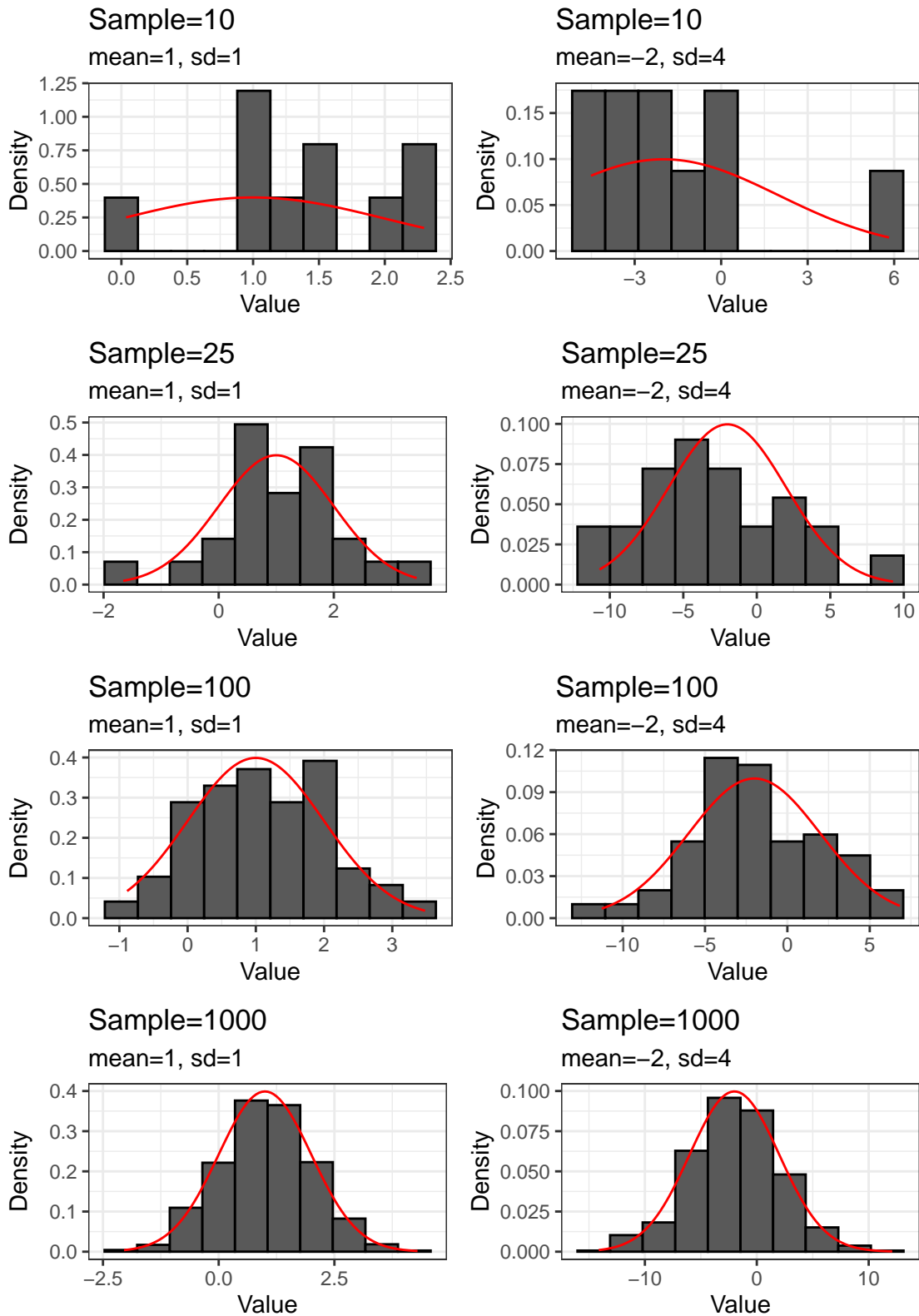


Figure 3: Histograms of each set of data with superimposed true density function

As we can see, the bigger sample we have, the closer the resulting histogram is to the real density function that we built on our set of parameters. We can notice that it's particularly true for $n > 30$.

- Our PDF is such a good fit for the data because we are using `rnorm()` function that randomly generates numbers for the normal distribution.

- ii. A good fit can also be explained via Central Limit Theorem: as our sample enlarges, the distribution of our random variable follows a normal distribution.

Checker Response (Ishraque): After our group meeting on Thursday I analyzed Chris' work and gave a few suggestions. I suggested some edits to the formatting of equations and descriptions in question 1. Next, I suggested some edits to the aesthetics of the plots of PDF and CDF and the histograms for the sections d, e, and f. These proposals were accepted promptly since they would improve the overall readability of the code and the solutions. I also suggested that for problem 1(f), he switches his method of plotting to a function based method instead. In general the code and the initial responses were already very well thought out and formatted and further edits only helped make it more concise and clear.

2. Continue with the continuous distribution you selected for Question 1.

- (a) Provide the mean, standard deviation, skewness, and kurtosis of the PDF. Ensure to interpret each.

$E(x) = \mu$	[Mean]
$var(X) = \sigma^2$	[Variance]
$skew(X) = 0$	[Skewness]
$kurt(X) = 0$	[Kurtosis]

The population skewness of the PDF is 0, which indicates that the Normal distribution is symmetric (with line of symmetry at the mean). The population excess kurtosis of the PDF is 0, which indicates that the Normal distribution is mesokurtic.

- (b) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). Calculate the sample mean, standard deviation, skewness, and kurtosis. Interpret the results.

```
library(e1071)
library(tidyverse)
library(patchwork)
obs <- c(10, 25, 100, 1000)

s1 <- data.frame(x = rnorm(n = obs[1], mean = 0, sd = 1))
s1_stats <- s1 %>%
  summarize(Mean = mean(x), SD = sd(x),
            Skewness = skewness(x),
            "Excess Kurtosis" = kurtosis(x))
s1_stats

##           Mean           SD Skewness Excess Kurtosis
## 1 0.2245916 0.8970171 0.189067      -0.9200314
```

Interpretation:

Center: The sample mean of a normal distribution is the balancing point of the distribution. In this case, this will be equal to around -

```
s1_stats$Mean
## [1] 0.2245916
```

Spread: The average distance from the sample mean is -

```
s1_stats$SD
## [1] 0.8970171
```

Skewness: The skewness of a dataset describes the symmetry of a distribution.

```
s1_stats$Skewness
## [1] 0.189067
```

Since skewness $\neq 0$, there are more observations for low values and fewer with high values.

Kurtosis describes how the peaked data is in relation to the Gaussian distribution (the bell-shaped curve).


```
s1_stats$'Excess Kurtosis'
## [1] -0.9200314
```

Since the excess kurtosis $\neq 0$, it indicates that the data is platykurtic.

```
s2 <- data.frame(x = rnorm(n = obs[2], mean = 0, sd = 1))
s2_stats <- s2 %>% summarize(Mean = mean(x),
                             SD = sd(x), Skewness = skewness(x),
                             "Excess Kurtosis" = kurtosis(x))

s2_stats

##           Mean          SD   Skewness Excess Kurtosis
## 1 -0.1281396  1.010546  0.04255577        -0.9013308
```

Interpretation:

Center: The sample mean of a normal distribution is the balancing point of the distribution. In this case, this will be equal to around -

```
s2_stats$Mean
## [1] -0.1281396
```

Spread: The average distance from the sample mean is -

```
s2_stats$SD
## [1] 1.010546
```

Skewness: The skewness of a dataset describes the symmetry of a distribution.

```
s2_stats$Skewness
## [1] 0.04255577
```

Since skewness is slightly greater than 0, there are more observations for low values and fewer with high values. However, note that the skewness for this sample is closer to 0 than s1.

Kurtosis describes how the peaked data is in relation to the Gaussian distribution (the bell-shaped curve).

```
s2_stats$'Excess Kurtosis'
## [1] -0.9013308
```

Since the excess kurtosis $\neq 0$, it indicates that the data is platykurtic. However, note that the excess kurtosis for this sample is closer to 0 than s1.

```
s3 <- data.frame(x = rnorm(n = obs[3], mean = 0, sd = 1))
s3_stats <- s3 %>% summarize(Mean = mean(x), SD = sd(x),
                             Skewness = skewness(x),
                             "Excess Kurtosis" = kurtosis(x))

s3_stats

##           Mean          SD   Skewness Excess Kurtosis
## 1  0.2391397  1.060967  0.3076019        0.03810334
```

Interpretation:

Center: The sample mean of a normal distribution is the balancing point of the distribution. In this case, this will be equal to around -

```
s3_stats$Mean
## [1] 0.2391397
```

Spread: The average distance from the sample mean is -

```
s3_stats$SD
## [1] 1.060967
```

Skewness: The skewness of a dataset describes the symmetry of a distribution.

```
s3_stats$Skewness
## [1] 0.3076019
```

Since skewness is slightly greater than 0, there are more observations for low values and fewer with high values. However, note that the skewness for this sample is closer to 0 than s2. Kurtosis describes how the peaked data is in relation to the Gaussian distribution (the bell-shaped curve).

```
s3_stats$`Excess Kurtosis`
## [1] 0.03810334
```

Since the excess kurtosis is > 0 , it indicates that the data is platykurtic. However, note that the excess kurtosis for this sample is closer to 0 than s2.

```
s4 <- data.frame(x = rnorm(n = obs[4], mean = 0, sd = 1))
s4_stats <- s4 %>% summarize(Mean = mean(x),
                             SD = sd(x), Skewness = skewness(x),
                             "Excess Kurtosis" = kurtosis(x))
```

Interpretation:

Center: The sample mean of a normal distribution is the balancing point of the distribution. In this case, this will be equal to around -

```
s4_stats$Mean
## [1] -0.004868836
```

Spread: The average distance from the sample mean is -

```
s4_stats$SD
## [1] 0.9724235
```

Skewness: The skewness of a dataset describes the symmetry of a distribution.

```
s4_stats$Skewness
## [1] 0.01798884
```

Since skewness is almost 0 (slightly greater), the distribution is almost symmetric. Kurtosis describes how the peaked data is in relation to the Gaussian distribution (the bell-shaped curve).

```
s4_stats$`Excess Kurtosis`
## [1] -0.3858467
```

Since the excess kurtosis is almost 0 (slightly lesser), it indicates that the data is almost mesokurtic.

As the sample size increases, the distribution moves towards becoming more Normal - the mean tends to equal 0, the SD tends to equal 1, and both the skewness and excess kurtosis tend to equal 0.

- (c) Generate a random sample of size $n = 10$ for your two sets of parameter(s).

```
dat1 <- data.frame(x = rnorm(n = obs[1], mean = 0, sd = 1))
```

Calculate the method of moments estimator(s) and maximum likelihood estimator(s).

```
library(nleqslv)
#####
# MOM estimator
```

```
#####
norm.mom<-function(par, data){
  mu <- par[1]
  sigma <- par[2]

  EX1 <- mu          # Expected value of a normal distribution
  EX2 <- sigma       # Variance of a normal distribution

  xbar1 <- mean(data)
  xbar2 <- mean(data^2)

  c(EX1-xbar1, EX2-xbar2)
}

# Entering the starting guess, the function(s) we want to solve for c(0, 0),
# and the dataframe a arguments of the non-linear equation solver
mom1 <- nleqslv(x = c(0,1),fn = norm.mom, data = dat1$x)

#####
# MLE
#####
norm.ll<-function(par, data, neg=T){
  mu <- par[1]
  sigma <- par[2]
  ll <- sum(dnorm(x=data, mean=mu, sd=sigma, log = T))
  ifelse(neg, -ll, ll)
  # Since the optim() function minimizes, we use neg because on multiplying by negative, it will
}

mle1 <- optim(par = c(0,1), fn = norm.ll, data=dat1$x)
```

In a 1×2 grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

```
mom_dat1 <- data.frame(x = rnorm(n = obs[1], mean = mom1$x[1], sd = mom1$x[2]))
mom1_p <- ggplot(data=mom_dat1, aes(x=x))+
  geom_histogram(color="black", fill = "blue",
    alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+ geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Method Of Moments (n = 10)") +
  geom_histogram(aes(x = dat1$x), color="black",
    fill = "green", alpha = 0.2) +
  geom_vline(xintercept = mean(dat1$x), alpha=0.35, color="red",size=1)

mle_dat1 <- data.frame(x = rnorm(n = obs[1],
    mean = mle1$par[1], sd = mle1$par[2]))
mle1_p <- ggplot(data=mle_dat1, aes(x=x))+
  geom_histogram(color="black", fill = "blue",
    alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Maximum Likelihood Estimator (n = 10)") +
  geom_histogram(aes(x = dat1$x, y=..density..),
    color="black", fill = "green", alpha = 0.2) +
  geom_vline(xintercept = mean(dat1$x), alpha=0.35, color="red",size=1)

mom1_p + mle1_p
```

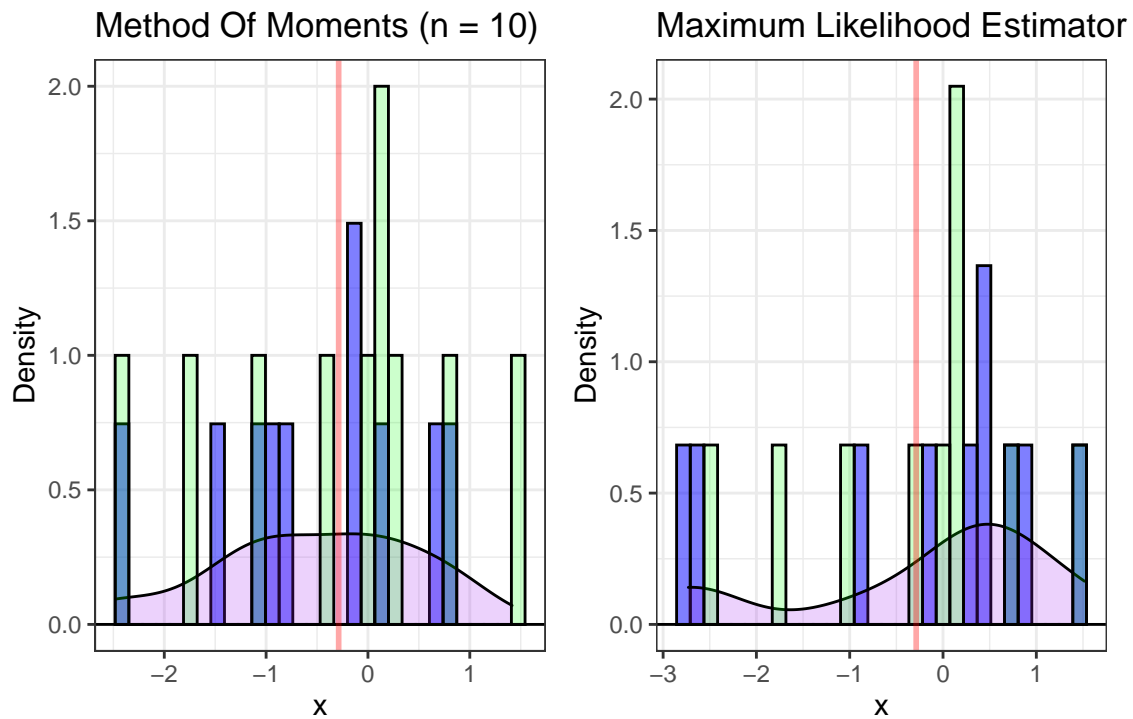


Figure 4: Cool caption!

```
dat1_h <- ggplot(data=dat1, aes(x=x))+
  geom_histogram(color="black", fill = "lightblue",
    alpha = 0.5, binwidth = 0.2, aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Normal Distrubution (n = 10)") +
  geom_vline(xintercept = mean(dat1$x), alpha=0.35,
    color="red",size=1)

dat1<-dat1 %>% mutate(n = seq(1,obs[1],1))

dat1_p <- ggplot(data=dat1, aes(x=n, y=x))+
  geom_line()+
  geom_hline(yintercept=0)+
  geom_hline(yintercept = mean(dat1$x), alpha=0.35, color="red",size=1) +
  theme_bw()+ xlab("n")+ ylab("x")+
  ggtitle("Normal Distrubution (n = 10)")

dat1_h + dat1_p
```

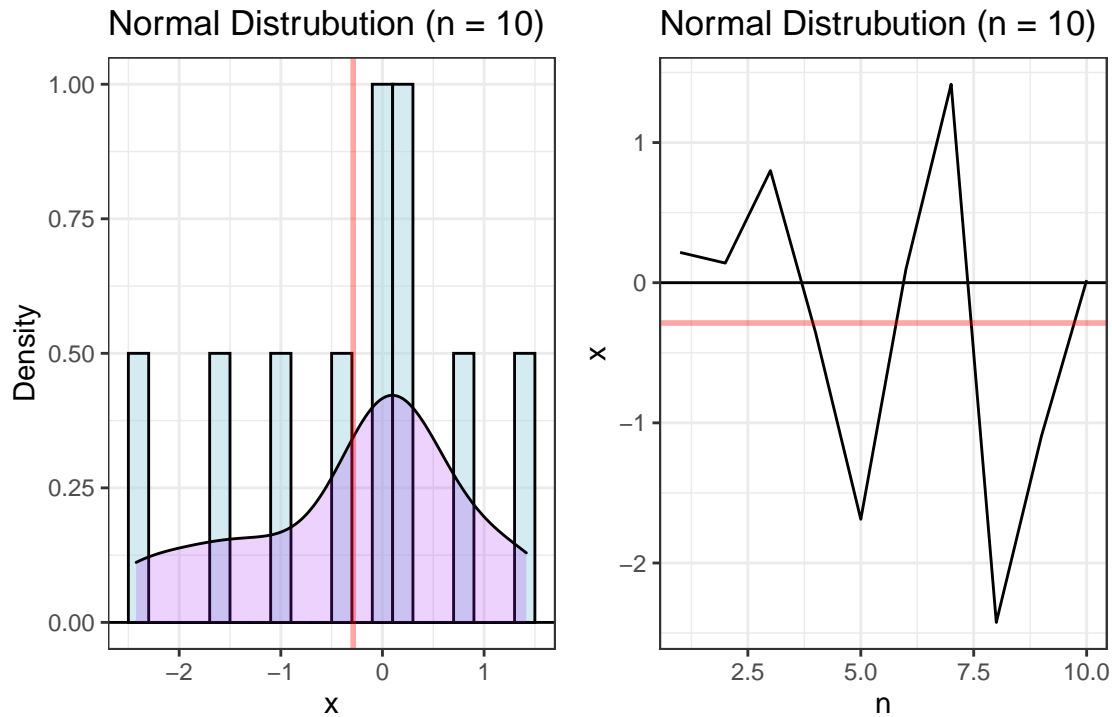


Figure 5: Cool caption!

- (d) Generate a random sample of size $n = 25$ for your two sets of parameter(s).

```
dat2 <- data.frame(x = rnorm(n = obs[2], mean = 0, sd = 1))
```

Calculate the method of moments estimator(s) and maximum likelihood estimator(s).

```
mom2 <- nleqslv(x = c(0,1), fn = norm.mom, data = dat2$x)
mle2 <- optim(par = c(0,1), fn = norm.ll, data=dat2$x)
```

In a 1×2 grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

```
mom_dat2 <- data.frame(x = rnorm(n = obs[2],
                                mean = mom2$x[1], sd = mom2$x[2]))
mom2_p <- ggplot(data=mom_dat2, aes(x=x))+
  geom_histogram(color="black", fill = "blue",
                 alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+ geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Method Of Moments (n = 25)") +
  geom_histogram(aes(x = dat2$x, y=..density..),
                 color="black", fill = "green", alpha = 0.2) +
  geom_vline(xintercept = mean(dat2$x),
             alpha=0.35, color="red",size=1)

mle_dat2 <- data.frame(x = rnorm(n = obs[2],
                                mean = mle2$par[1], sd = mle2$par[2]))
mle2_p <- ggplot(data=mle_dat2, aes(x=x))+
  geom_histogram(color="black", fill = "blue",
                 alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
```

```
ggtitle("Maximum Likelihood Estimator (n = 25)") +
geom_histogram(aes(x = dat2$x, y=..density..),
               color="black", fill = "green", alpha = 0.2) +
geom_vline(xintercept = mean(dat2$x), alpha=0.35, color="red",size=1)
```

```
mom2_p + mle2_p
```

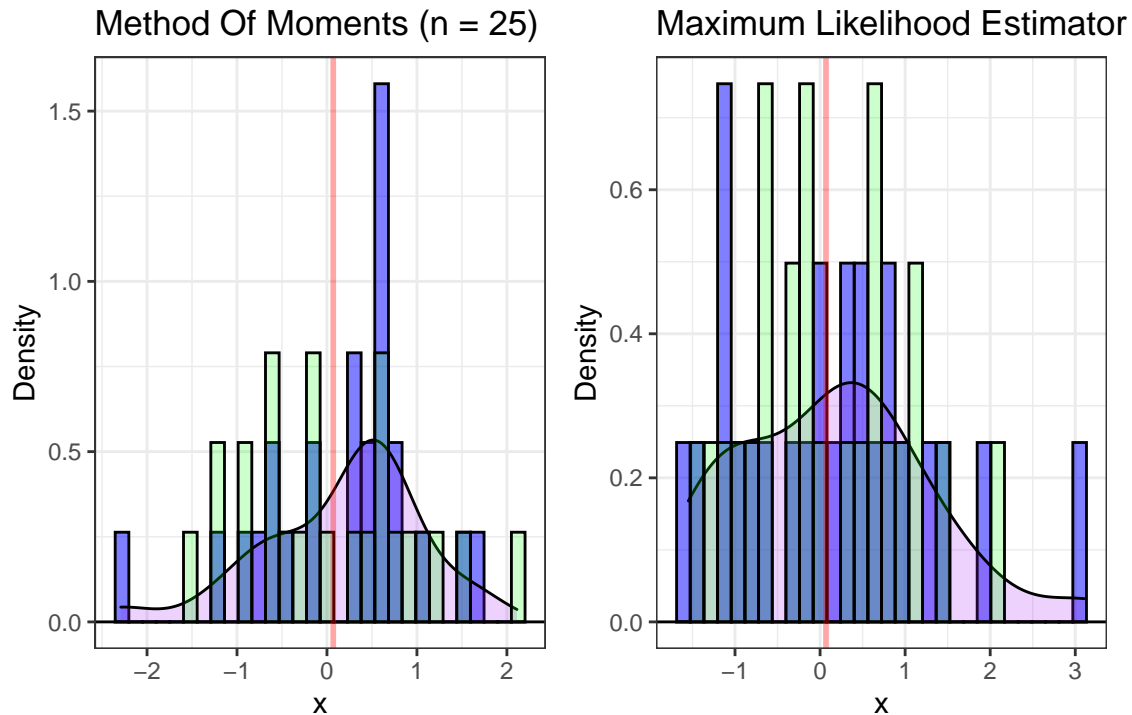


Figure 6: Cool caption!

```
dat2_h <- ggplot(data=dat2, aes(x=x))+
  geom_histogram(color="black", fill = "lightblue",
                alpha = 0.5, binwidth = 0.2,
                aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+
  xlab("x")+
  ylab("Density")+
  ggtitle("Normal Distrubution (n = 25)") +
  geom_vline(xintercept = mean(dat2$x), alpha=0.35, color="red",size=1)

dat2 <- dat2 %>% mutate(n = seq(1,obs[2],1))
dat2_p <- ggplot(data=dat2, aes(x=n, y=x))+
  geom_line()+
  geom_hline(yintercept=0)+
  geom_hline(yintercept = mean(dat2$x), alpha=0.35, color="red",size=1) +
  theme_bw()+
  xlab("n")+
  ylab("x")+
  ggtitle("Normal Distrubution (n = 25)")

dat2_h + dat2_p
```

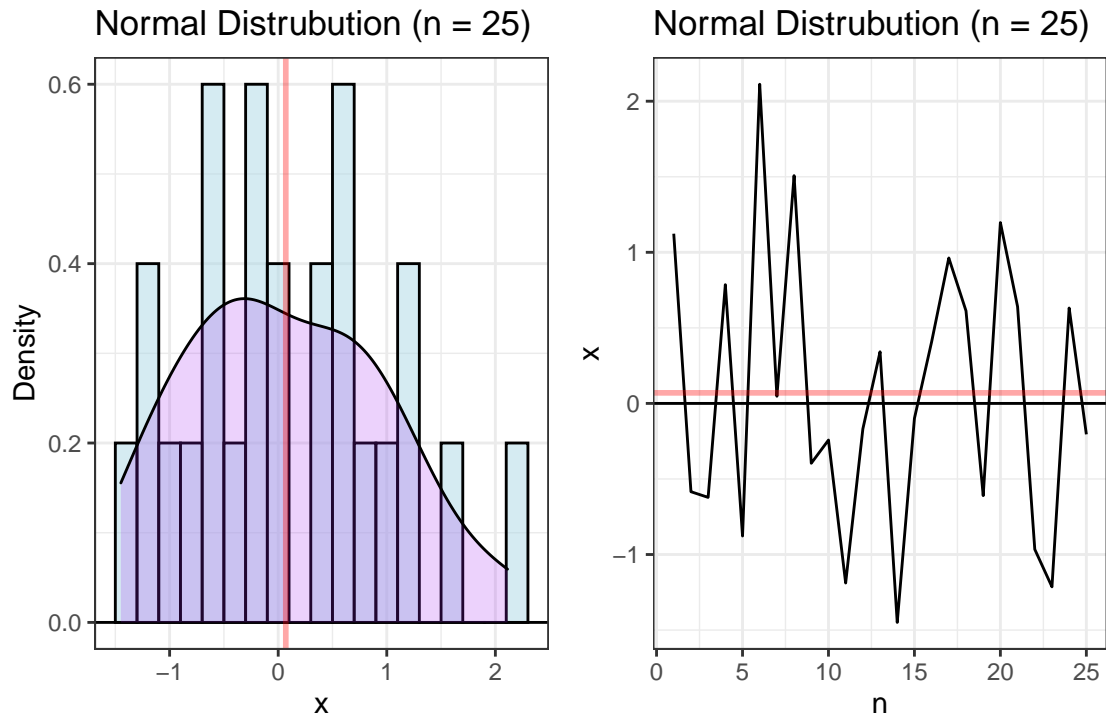


Figure 7: Cool caption!

- (e) Generate a random sample of size $n = 100$ for your two sets of parameter(s).

```
dat3 <- data.frame(x = rnorm(n = obs[3], mean = 0, sd = 1))
```

Calculate the method of moments estimator(s) and maximum likelihood estimator(s).

```
mom3 <- nleqslv(x = c(0,1),fn = norm.mom, data = dat3$x)
mle3 <- optim(par = c(0,1), fn = norm.ll, data=dat3$x)
```

In a 1×2 grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

```
mom_dat3 <- data.frame(x = rnorm(n = obs[3], mean = mom3$x[1], sd = mom3$x[2]))
mom3_p <- ggplot(data=mom_dat3, aes(x=x))+
  geom_histogram(color="black",
                 fill = "blue", alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+
  ylab("Density")+
  ggtitle("Method of Moments (n = 100)") +
  geom_histogram(aes(x = dat3$x, y=..density..),
                 color="black", fill = "green", alpha = 0.2) +
  geom_vline(xintercept = mean(dat3$x), alpha=0.35, color="red",size=1)

mle_dat3 <- data.frame(x = rnorm(n = obs[3],
                                mean = mle3$par[1], sd = mle3$par[2]))
mle3_p <- ggplot(data=mle_dat3, aes(x=x))+
  geom_histogram(color="black", fill = "blue",
                 alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
```

```
ggtitle("Maximum Likelihood Estimator (n = 100)") +
geom_histogram(aes(x = dat3$x, y=..density..),
               color="black", fill = "green", alpha = 0.2) +
geom_vline(xintercept = mean(dat3$x), alpha=0.35, color="red",size=1)
```

```
mom3_p + mle3_p
```

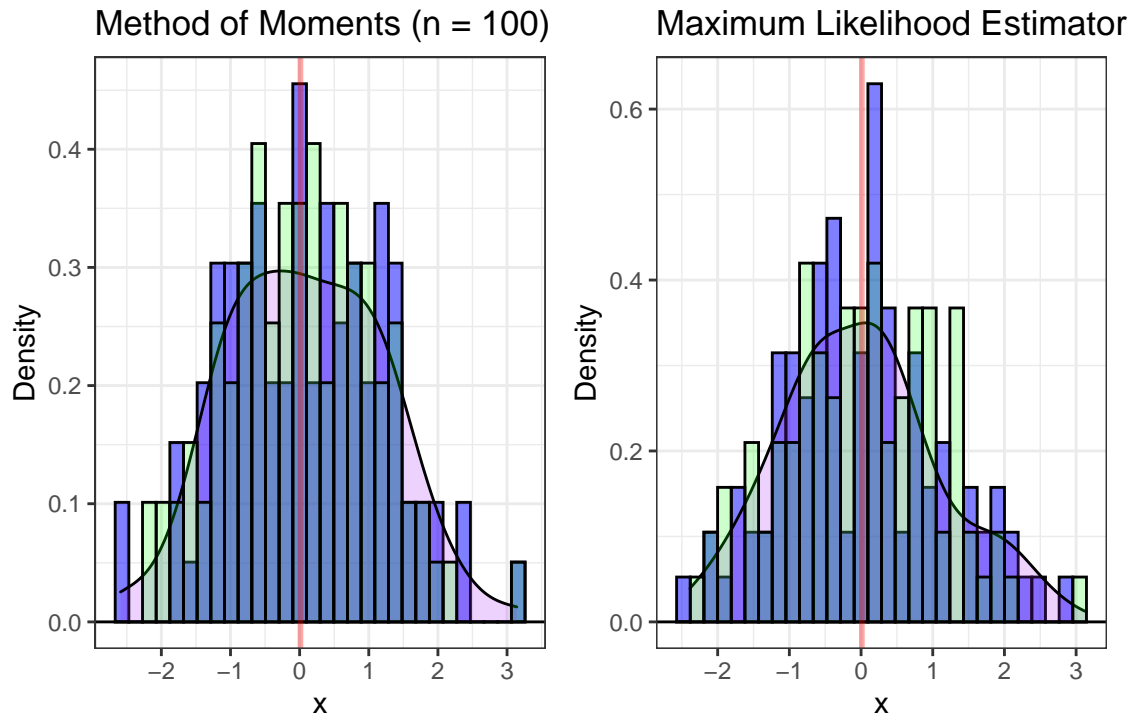


Figure 8: Cool caption!

```
dat3_h <- ggplot(data=dat3, aes(x=x))+
  geom_histogram(color="black", fill = "lightblue",
                alpha = 0.5, binwidth = 0.2, aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Normal Distrubution (n = 100)")

dat3 <- dat3 %>% mutate(n = seq(1,obs[3],1))
dat3_p <- ggplot(data=dat3, aes(x=n, y=x))+
  geom_line()+
  geom_hline(yintercept=0)+
  geom_hline(yintercept = mean(dat3$x), alpha=0.35, color="red",size=1) +
  theme_bw()+ xlab("n")+ ylab("x")+
  ggtitle("Normal Distrubution (n = 100)")

dat3_h + dat3_p
```

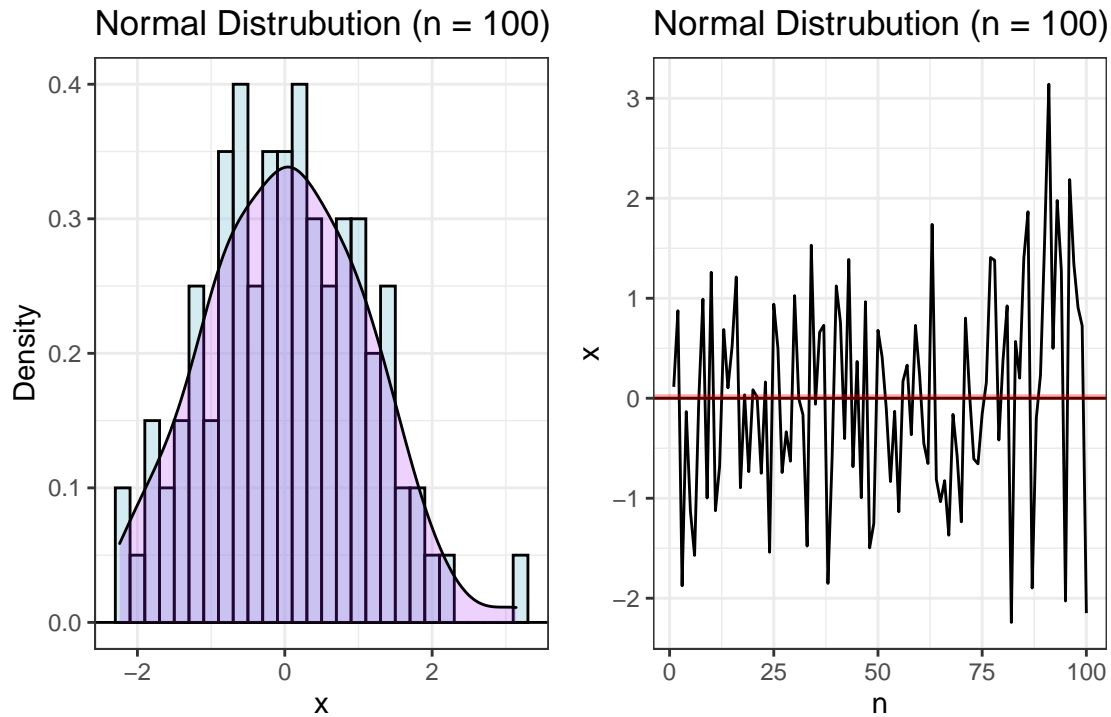



Figure 9: Cool caption!

- (f) Generate a random sample of size $n = 100$ for your two sets of parameter(s).

```
dat4 <- data.frame(x = rnorm(n = obs[4], mean = 0, sd = 1))
```

Calculate the method of moments estimator(s) and maximum likelihood estimator(s).

```
mom4 <- nleqslv(x = c(0,1),fn = norm.mom, data = dat4$x)
mle4 <- optim(par = c(0,1), fn = norm.ll, data=dat4$x)
```

In a 1×2 grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

```
mom_dat4 <- data.frame(x = rnorm(n = obs[4], mean = mom4$x[1], sd = mom4$x[2]))
mom4_p <- ggplot(data=mom_dat4, aes(x=x))+
  geom_histogram(color="black", fill = "blue",
    alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+
  geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Method of Moments (n = 1000)") +
  geom_histogram(aes(x = dat4$x, y=..density..),
    color="black", fill = "green", alpha = 0.2) +
  geom_vline(xintercept = mean(dat4$x), alpha=0.35, color="red",size=1)

mle_dat4 <- data.frame(x = rnorm(n = obs[4],
  mean = mle4$par[1], sd = mle4$par[2]))
mle4_p <- ggplot(data=mle_dat4, aes(x=x))+
  geom_histogram(color="black", fill = "blue",
    alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+ geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Maximum Likelihood Estimator (n = 1000)") +
  geom_vline(xintercept = mean(dat4$x), alpha=0.35,
```

```

    color="red",size=1) +
  geom_histogram(aes(x = dat4$x, y=..density..),
    color="black", fill = "green", alpha = 0.2)

```

```
mom4_p + mle4_p
```

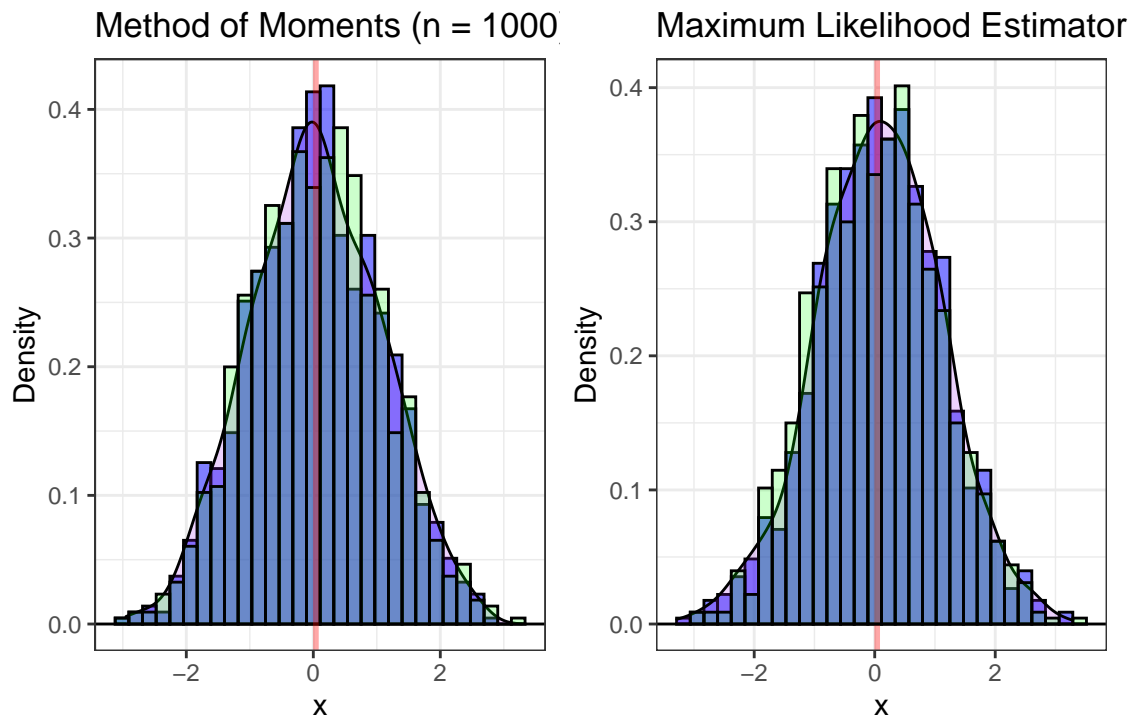


Figure 10: Cool caption!

```

dat4_h <- ggplot(data=dat4, aes(x=x))+
  geom_histogram(color="black", fill = "lightblue",
    alpha = 0.5, aes(y=..density..))+
  geom_hline(yintercept=0)+ geom_density(alpha=.2, fill="purple") +
  theme_bw()+ xlab("x")+ ylab("Density")+
  ggtitle("Normal Distrubution (n = 1000)") +
  geom_vline(xintercept = mean(dat4$x), alpha=0.35, color="red",size=1)

dat4 <- dat4 %>% mutate(n = seq(1,obs[4],1))
dat4_p <- ggplot(data=dat4, aes(x=n, y=x))+
  geom_line()+ geom_hline(yintercept=0)+
  geom_hline(yintercept = mean(dat3$x), alpha=0.35, color="red",size=1) +
  theme_bw()+ xlab("n")+ ylab("x")+ ggtitle("Normal Distrubution (n = 1000)")

dat4_h + dat4_p

```

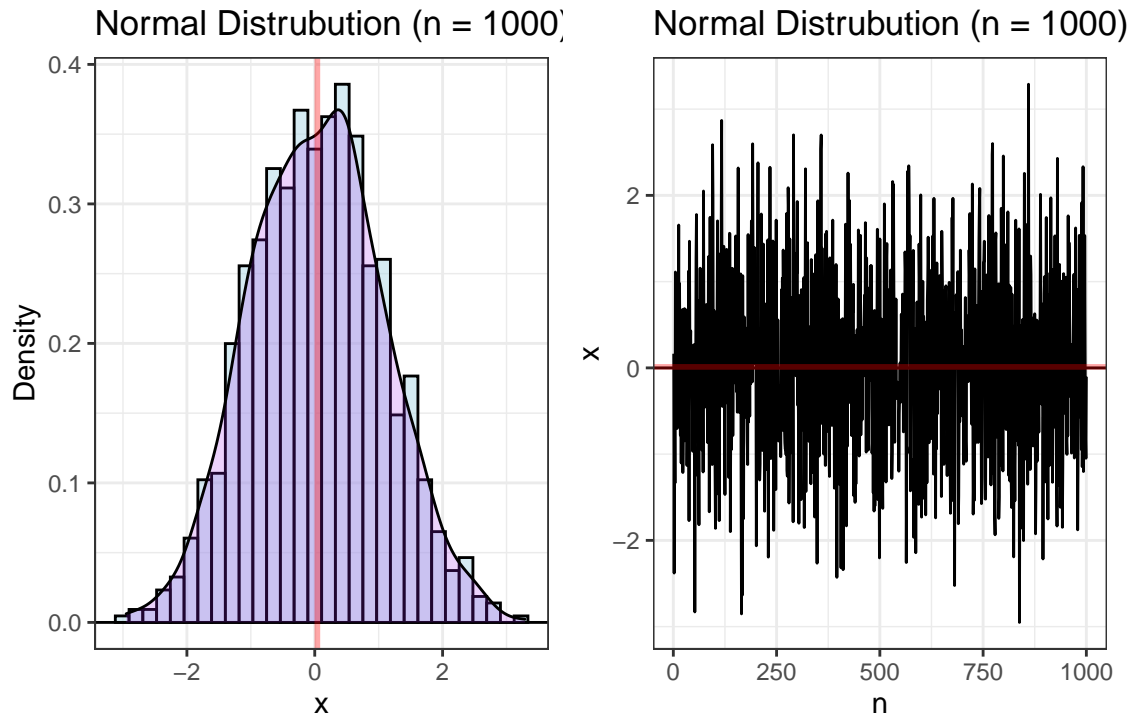


Figure 11: Cool caption!

(g) Comment on the results of parts (c)-(f).

Parts (c)-(f) illustrate the weak law of large numbers. As the sample size n increases, we see that both our computed estimators (MOM and MLE) tend to overlap more with their corresponding probability distributions - there is very less overlap when $n = 10$; however, at $n = 1000$, both the histograms are essentially superimposed on top of one another. Thus, as the sample size increases, our estimators do a better job at estimating the population statistics of the Normal distribution. We also notice that as the sample size increases, the sample mean gets closer to 0.

3. Select a discrete distribution (not the Poisson). It does not have to be one that we cover in the notes! To explore the PMF of your distribution, specify two sets of parameter(s) for your distribution.
 - (a) **History** Discuss what types of random variables are modeled with your distribution. Be sure to include a discussion about the support and ensure to provide the mass function, and CDF. This requires some internet research – what’s the history of the distribution, why was it created and named? What are some exciting applications of this distribution? Cite all of your sources.

Solution:

The discrete distribution we have chosen is the Bernoulli distribution. The Bernoulli distribution is a discrete probability distribution for a Bernoulli trial - a probabilistic experiment that can have one of two outcomes, success ($x = 1$) and failure ($x = 0$), and in which the probability of success is p . Often p is called the Bernoulli probability parameter (Forbes et al., 2011). In Bernoulli distribution, the random variable X can have only one of two values: 0 or 1. This means that the support of our discrete random variable is the set $\{0, 1\}$. This distribution can be summarized as follows:

$p \in (0, 1)$	[Parameter]
$\mathcal{X} = \{x : x \in \{0, 1\}\}$	[Support]
$f_X(x p) = p^x(1 - p)^{1-x}I(x \in \{0, 1\})$	[PMF]
$F_X(x p) = P(X \leq \lfloor x \rfloor)$ $= [(1 - p)I(\lfloor x \rfloor = 0)] + I(\lfloor x \rfloor \geq 1)$	[CDF]

Simply, the Bernoulli distribution can be thought of as a model for the set of possible outcomes of any single experiment that asks a yes-no question (Vukadin et al., 2021). This distribution is named after the 17th century Swiss mathematician Jacob Bernoulli, because he was the one who explicitly defined the concept of Bernoulli trial (in his book *Ars Conjectandi*) described above. Bernoulli distribution serves as a building block for discrete distributions which model Bernoulli trials, such as Binomial distribution and geometric distribution (Cthae, 2020). Logistic regression, a widely used classification model that models a binary outcome, also takes advantage of the Bernoulli distribution (Kleinbaum and Klein, 2010).

Since the Bernoulli distribution is not cataloged by R, we have to introduce it into our calculations with the following functions:

```
# Bernoulli PMF
dbern<-function(x,prob){
  if(prob<0 | prob>1){
    errormsg <- "This function is only valid for success probabilities between 0 and 1."
    stop(errormsg)
  }
  indicator <- rep(0, length(x))
  indicator[x==0] <- 1 # indicator should be one if x=0
  indicator[x==1] <- 1 # indicator should be one if x=1
  fx <- (prob^x * (1-prob)^(1-x)) * indicator # PMF formula
  return(fx)
}

# Bernoulli CDF
pbern<-function(q, prob){
  if(prob<0 | prob>1){
    errormsg<-"This function is only valid for success probabilities between 0 and 1."
    stop(errormsg)
  }
  indicator1 <- rep(1, length(q))
  indicator1[q != 0] <- 0 #indicator should be zero if x!=0
  indicator2 <- rep(1, length(q))
  indicator2[q < 1] <- 0 #indicator should be zero if x<1
  Fx <- (1-prob) * indicator1 + indicator2
  return(Fx)
}
```

The R packages Tidyverse (Wickham et al., 2019) and Patchwork (Pedersen, 2020) will be used in all of the following plots in question 3.

- (b) Show that you have a valid PMF. You can show this approximately by calculating the series in a repeat loop until probability mass evaluations are infinitesimally small.

Solution: For a PMF to be valid, it has to fulfill the following statements:

1. $0 \leq f_X(x) \leq 1$ for all $x \in \mathbb{R}$
2. $\sum_{-\infty}^{\infty} f_X = \sum_{\mathcal{X}} f_X = 1$

By definition, the Bernoulli distribution satisfies statement 1, since the support is 0,1 and the probability parameter $p \in (0, 1)$, so the term for PMF (see above) cannot have a value smaller than 0 or greater than 1.

We can also show that the statement 2 is true for the Bernoulli distribution:

$$\begin{aligned} \sum_{\mathcal{X}} f_X &= \sum_{x=0}^1 f_X(x) = \sum_{x=0}^1 p^x (1-p)^{1-x} \\ &= p^0 (1-p)^{1-0} + p^1 (1-p)^{1-1} \\ &= (1)(1-p) + p(1) \\ &= 1 - p + p \\ &= 1 \end{aligned}$$

Therefore, we have a valid PMF for our distribution.

- (c) Find the median for your two sets of parameter(s). Conduct some research to find the median based on our PMF to confirm that your numerical approach is correct.

Solution: In the Bernoulli distribution, where p is the probability parameter and $q = 1 - p$, when $p > q$, there are more "successes" (that is when X is 1) than failure (when X is 0), so the median must also be a success. Similar, when $q > p$, there are more "failures" than "successes", hence resulting in the median being 0. When $p = q$, there is an equal probability of X being 0 or 1, thus the median is ambiguous. Therefore, the median for when $p = 0.4$ is zero and the median for when $p = 0.6$ is 1. Further inquiry (Ber, 2021) suggests that this really is the case for the median of the Bernoulli distribution:

$$Median = \begin{cases} 0 & \text{if } p < 1/2 \\ [0, 1] & \text{if } p = 1/2 \\ 1 & \text{if } p > 1/2 \end{cases}$$

- (d) Graph the PMF for several values of the parameter(s) including the two sets you specified. What does changing the parameter(s) do to the shape of the PMF?

Solution: We can plot the PMF of the Bernoulli distribution for different parameters. Note that the initially chosen parameters were $p = 0.4$ and $p = 0.6$. Let us also look at the PMF when the parameter values are 0.5, and 0.8.

```
library(tidyverse)
library(patchwork)
plotbernPMF <- function(prob){ # Pass in the success probability
  ggdat <- data.frame(x = (-1:2),
                     f = dbern(x = (-1:2), prob = prob),
                     F = pbern(q = (-1:2), prob = prob))

  ## Plot PMF
  PMF <- ggplot(data = ggdat, aes(x = x)) +
    geom_linerange(aes(ymin = 0), ymax = f) +
    geom_hline(yintercept = 0) +
    theme_bw() +
    ylim(0, 1) +
    xlab("X") +
    ylab(bquote(f[x](x))) +
    ggtitle("Bernoulli Distribution", subtitle = paste("p =", prob))

  return(PMF)
}

plotbernPMF(0.4) + plotbernPMF(0.5) + plotbernPMF(0.6) + plotbernPMF(0.8)
```

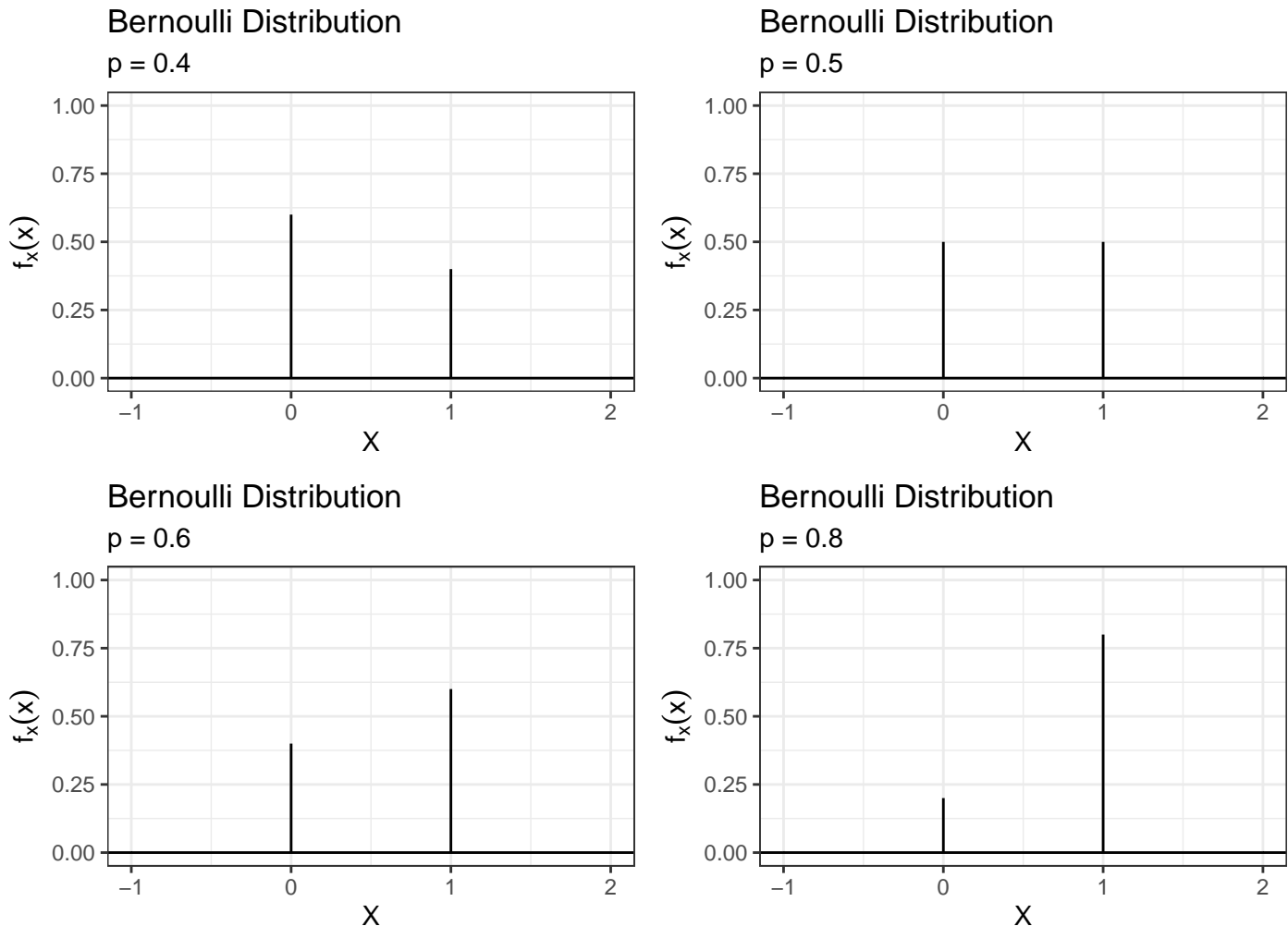


Figure 12: The CDF of the Bernoulli distribution for different probability parameters p

We observe the same characteristic of the Bernoulli distribution in Figure 12 as we discussed about its median above. If the $p < 0.5$ then there are more failures than successes, indicated by the taller PMF at 0. On the other hand, if $p > 0.5$ then there are more successes than failures, which is reflected in the PMF plot being taller at 1. The plot of the PMF when $p = 0.5$ suggests that there is an equal chance for success as there is for failure.

- (e) Graph the CDF for the same values of the parameter(s) as you did in Question 3d. What does changing the parameter(s) do to the shape of the CDF? Comment on the aspects of the CDFs that show that the CDF is valid.

Solution: We can plot the CDF of the Bernoulli distribution for the different parameters 0.4, 0.5, 0.6, and 0.8 below:

```
plotbernCDF <- function(prob) { # Pass in the success probability
  ggdat <- data.frame(x = (-1:2),
    f = dbern(x = (-1:2), prob = prob),
    F = pbern(q = (-1:2), prob = prob))
  ggdat.openpoints <- data.frame(x = ggdat$x,
    y = pbern(ggdat$x-1, prob = prob))
  ggdat.closedpoints <- data.frame(x = ggdat$x,
    y = pbern(ggdat$x, prob = prob))
  CDF <- ggplot(data = ggdat, aes(x = x, y = F)) +
    geom_step() +
    geom_point(data = ggdat.openpoints, aes(x = x, y = y), shape = 1) +
    geom_point(data = ggdat.closedpoints, aes(x = x, y = y)) +
    geom_hline(yintercept = 0.5, linetype = "dotted", color = "red") +
```

```

theme_bw()+
xlab("X")+
ylab(bquote(F[x](x)))+
ggtitle("Bernoulli CDF", subtitle=(paste("p =", prob)))
return(CDF)
}
plotbernCDF(0.4) + plotbernCDF(0.5) + plotbernCDF(0.6) + plotbernCDF(0.8)

```

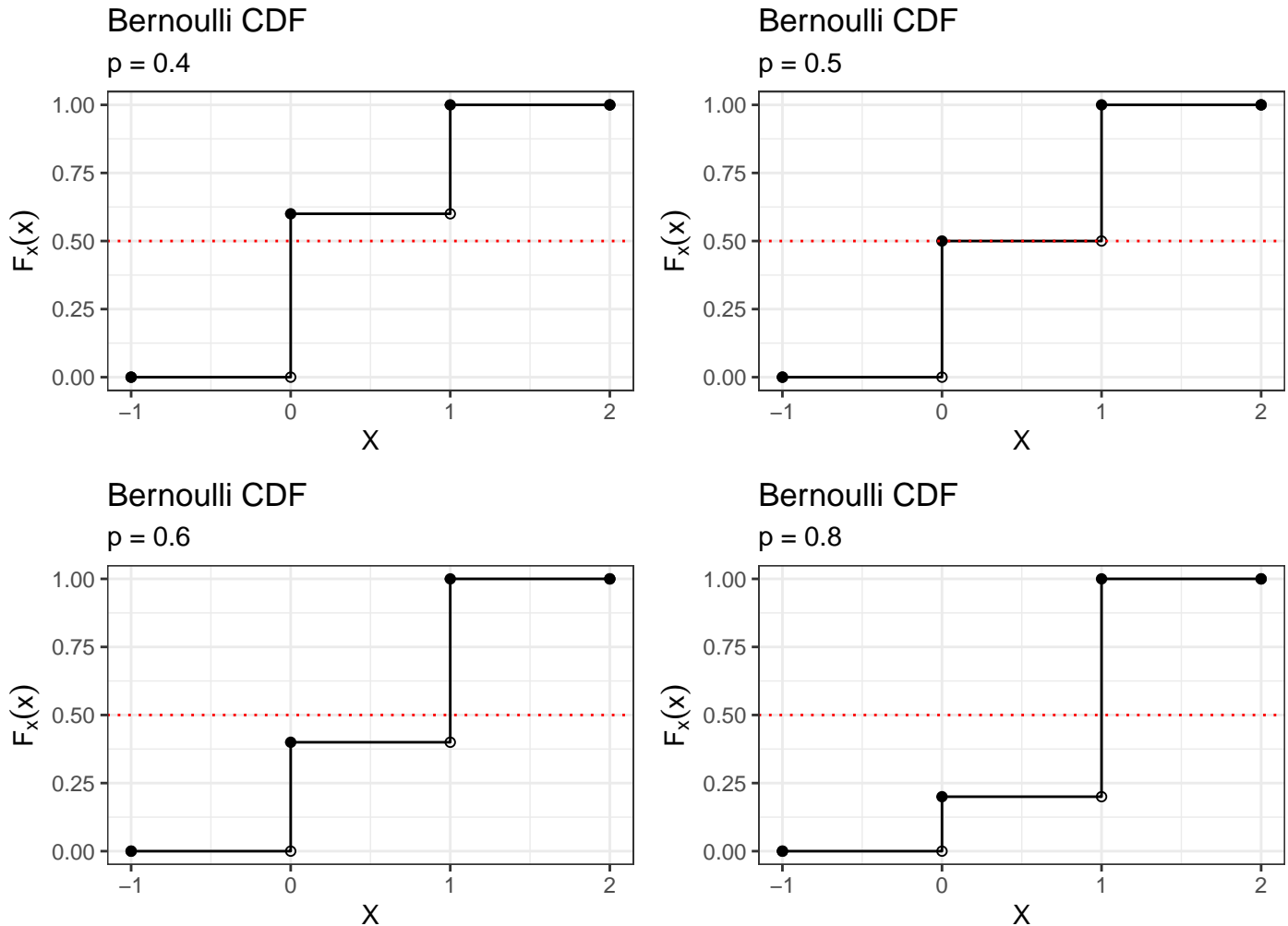


Figure 13: The CDF of the Bernoulli distribution for different probability parameters p

In the Figure 13, we can see that with increasing p , the most of the "area" under the CDF is moving towards 1, which correctly indicates that the number of successes increase as the value of p increases. Furthermore, we can see that in all the CDF plots, the functions always add up to 1. The CDF are also monotonically increasing towards the right. All of these above characteristics together indicate that our CDF is valid. Notice that in all of the plots of Figure 13, we have drawn a horizontal line at $1/2$. This is another way we can determine the median of the Bernoulli distribution using the CDF. Where the horizontal line intersects the CDF indicates where the median is for the given value of p . These results agree with our results from problem 3(c).

- (f) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). In a 4×2 grid, plot a histogram (with bin size 1) of each set of data and superimpose the true mass function at the specified parameter values. Interpret the results.

Solution: The two parameters we are working with are $p = 0.4$ and $p = 0.6$

```

# Using sampling with replacement to generate arrays of 0 and 1

library(tidyverse)
library(patchwork)

plot_df <- list( x1 = bern_sample(10, 0.4),
                 x2 = bern_sample(25, 0.4),
                 x3 = bern_sample(100, 0.4),
                 x4 = bern_sample(1000, 0.4),
                 y1 = bern_sample(10, 0.6),
                 y2 = bern_sample(25, 0.6),
                 y3 = bern_sample(100, 0.6),
                 y4 = bern_sample(1000, 0.6))

buildingPlot2 <- function(source, prob, sample_size){
  df <- data.frame(value=source) #turning values from the list into a df
  colnames(df) <- "value" #changing the name of the column
  df_PMF <- data.frame(x = (-1:2),
                      f = dbern(x = (-1:2), prob = prob))
  answer<-ggplot(df, aes(value))+
    geom_histogram(data = df, aes(y=..density..), binwidth=1,
                  color="black")+
    geom_linerange(data=df_PMF, aes(x=x, ymax = f), ymin = 0, size=2, color="red")+
    theme_bw() +
    ggtitle("Bernoulli Distribution", subtitle = paste("Sample Size =", sample_size, ", Prob =", prob))
  answer
}

x1 <- buildingPlot2(plot_df[1], 0.4, 10)
x2 <- buildingPlot2(plot_df[2], 0.4, 25)
x3 <- buildingPlot2(plot_df[3], 0.4, 100)
x4 <- buildingPlot2(plot_df[4], 0.4, 1000)

y1 <- buildingPlot2(plot_df[5], 0.6, 10)
y2 <- buildingPlot2(plot_df[6], 0.6, 25)
y3 <- buildingPlot2(plot_df[7], 0.6, 100)
y4 <- buildingPlot2(plot_df[8], 0.6, 1000)

(x1+y1)/(x2+y2)/(x3+y3)/(x4+y4)

```

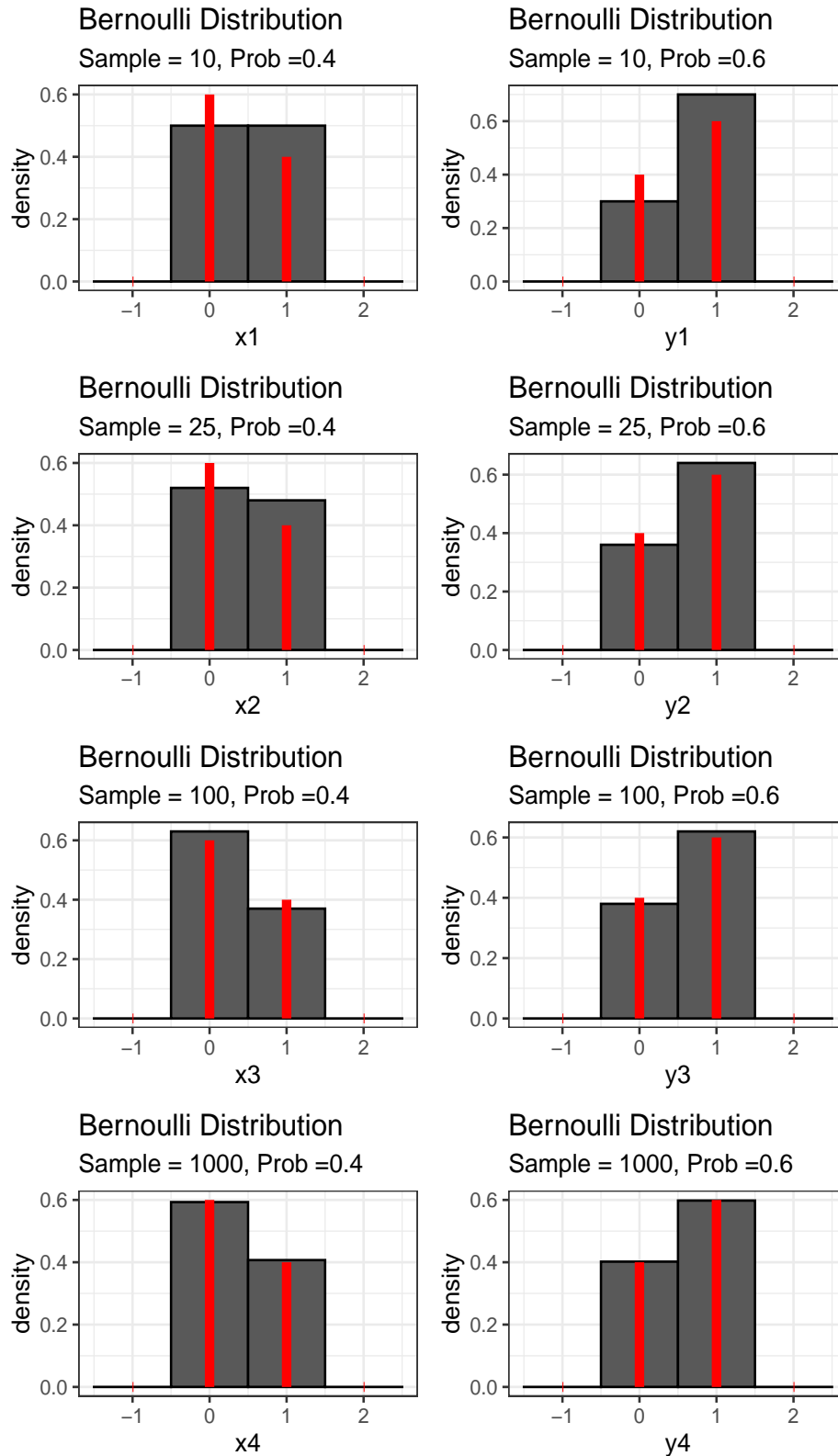



Figure 14: (Left Column) Samples of 10, 25, 100, 1000 for the Bernoulli distribution with $p = 0.4$. (Right Column) Samples of 10, 25, 100, 1000 for the Bernoulli distribution with $p = 0.6$. In each plot the red lines indicate the true PMF function of the corresponding Bernoulli distribution. Notice that the y-axis indicates the proportional frequency in the histograms

For each of the cases of p (see the first and the second row of plots in Figure 14), the more we sample, the better our histograms agree with the true PMF of the Bernoulli distribution. These histograms above represent only one variation of the random sampling done on the Bernoulli distribution. While the histograms for $n=1000$ does not stray

much from the PMF as we take newer samples, the smaller samples (10, 25) vary significantly from the PMF. From sample to sample, the results of the histograms vary (in terms of agreement with the true PMF). The larger the size of the sample, the more consistent is its agreement with the true PMF.

Checker Response (Chris): After having a group meeting on Thursday, I analyzed Ishraque's work and gave a couple of suggestions on the best way to improve his part of the assignment. First off, I suggested putting his separate data frames into one list instead: it would clean up his code and make it more readable. Also, instead of making new plots for each variable he's creating, I suggested writing a function that would take in variables and give him a plot as an output. This would improve the overall readability of the code. The team agreed with this decision, and we decided to implement them!

4. Continue with the discrete distribution you selected for Question 3.
 - (a) Provide the mean, standard deviation, skewness, and kurtosis of the PMF. Ensure to interpret each.
 - (b) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). Calculate the sample mean, standard deviation, skewness, and kurtosis. Interpret the results.
 - (c) Generate a random sample of size $n = 10$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a 1×2 grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.
 - (d) Generate a random sample of size $n = 25$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a 1×2 grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.
 - (e) Generate a random sample of size $n = 100$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a 1×2 grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.
 - (f) Generate a random sample of size $n = 100$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a 1×2 grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.
 - (g) Comment on the results of parts (c)-(f).

References

(2021).

Ahsanullah, M., Kibria, B. G., and Shakil, M. (2014). Normal distribution. In *Normal and Student st Distributions and Their Applications*, pages 7–50. Springer.

Cthaeh, T. (2020). The bernoulli distribution: Intuitive understanding.

Forbes, C., Evans, M., Hastings, N., and Peacock, B. (2011). *Statistical Distributions*. John Wiley & Sons, Inc., Hoboken, New Jersey.

Kleinbaum, D. G. and Klein, M. (2010). Introduction to Logistic Regression. In Kleinbaum, D. G. and Klein, M., editors, *Logistic Regression: A Self-Learning Text*, Statistics for Biology and Health, pages 1–39. Springer, New York, NY.

Patel, J. K. and Read, C. B. (1996). *Handbook of the normal distribution*, volume 150. CRC Press.

Pedersen, T. L. (2020). *patchwork: The Composer of Plots*. R package version 1.1.1.

R Core Team (2021). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.

Vukadin, O., Lee, S., Williams, C., and Khim, J. (2021). Bernoulli distribution | Brilliant.

Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.

Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Golemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., Takahashi, K., Vaughan, D., Wilke, C., Woo, K., and Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43):1686.