

LAB 4

RAJ SHAH, MILI PATEL, VRUND PATEL

2025-02-26

```
# Load required Libraries
library(dplyr) # This provides the %>% operator

## Warning: package 'dplyr' was built under R version 4.4.3

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union

library(ggplot2)

## Warning: package 'ggplot2' was built under R version 4.4.3

library(moments) # For skewness calculation

# Read the distinct (population) data
boxoffice_pop <- read.csv("C:/Users/rajsh/OneDrive/Desktop/Inference Data Science 291/LAB4/movie_boxoffice_distinct.csv")

#####
# Question 0: Remove Duplicates and Set Random Seed
#####

# Load Libraries
library(dplyr)

# Set seed for reproducibility
set.seed(1234) #  This is your random seed

# Read full dataset (before removing duplicates)
boxoffice_raw <- read.csv("C:/Users/rajsh/OneDrive/Desktop/Inference Data Science 291/LAB4/movie_boxoffice.csv")

# Remove duplicates to create population data
boxoffice_pop <- distinct(boxoffice_raw)
```

```

# Calculate how many duplicate rows were removed
duplicates_removed <- nrow(boxoffice_raw) - nrow(boxoffice_pop)

# Output for documentation
cat("Random seed number:", 1234, "\n")

## Random seed number: 1234

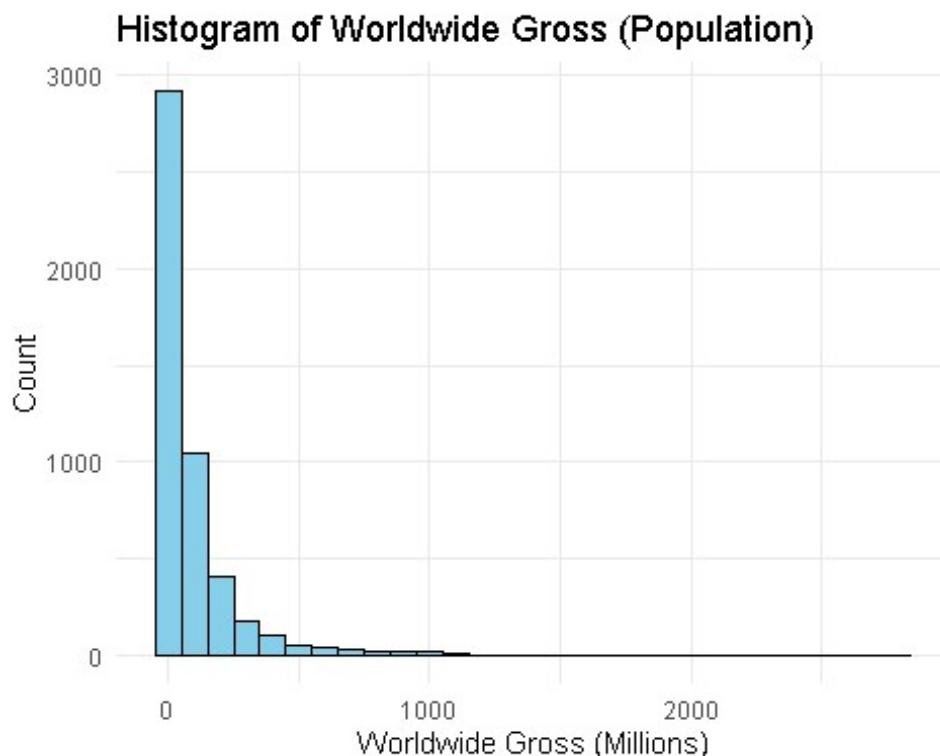
cat("Number of duplicated records removed from boxoffice data:",
duplicates_removed, "\n")

## Number of duplicated records removed from boxoffice data: 100

#####
# Part I, Question 1
#####

# Colored histogram using ggplot2
ggplot(boxoffice_pop, aes(x = Worldwide_Gross)) +
  geom_histogram(binwidth = 100,           # Bin width of 100 million, adjust
if needed
                fill = "skyblue",      # Fill color for the bars
                color = "black") +
                # Border color for the bars
  labs(title = "Histogram of Worldwide Gross (Population)",
       x = "Worldwide Gross (Millions)",
       y = "Count") +
  theme_minimal()

```



```
#####
# Explanation of the Shape of the Distribution:
#
# - The distribution of Worldwide Gross is **highly right-skewed**.
# - Most movies earn relatively low amounts (majority fall in the lower bins).
# - There are a few movies with extremely high gross earnings (outliers),
#   which stretch the tail to the right.
# - This is typical in revenue-related data: a small number of blockbusters
#   make a disproportionately large amount of money.
# - The long right tail and high concentration near lower values indicate
#   **positive skewness** in the population data.
#####

#####
# Part I, Question 2
#####

# Compute population statistics
pop_mean <- mean(boxoffice_pop$Worldwide_Gross, na.rm = TRUE)
pop_sd   <- sd(boxoffice_pop$Worldwide_Gross, na.rm = TRUE)
pop_prop <- mean(boxoffice_pop$Worldwide_Gross > boxoffice_pop$Budget, na.rm = TRUE)

# Display results
cat("Population mean (Worldwide Gross):", pop_mean, "\n")
## Population mean (Worldwide Gross): 95.83149
cat("Population SD (Worldwide Gross):", pop_sd, "\n")
## Population SD (Worldwide Gross): 177.4594
cat("Proportion of movies with Worldwide Gross > Budget:", pop_prop, "\n")
## Proportion of movies with Worldwide Gross > Budget: 0.6461286

#####
# Explanation:
#
# - The **average global box office earning** (mean) tells us how much
#   revenue movies
#   made on average in the dataset.
#
# - The **standard deviation (SD)** shows how spread out the earnings are.
#   A Large SD here likely reflects the fact that a few blockbuster movies
#   make extremely high revenue compared to most others.
#
# - The **proportion of movies where Worldwide Gross > Budget** helps us
#   understand
#   how many movies were profitable (i.e., made more than they cost to
```

```

produce).
# This gives insight into the success rate of movies in the dataset.
#
# Note: Since the data is right-skewed (as seen in Q1), the mean may be
pulled
# upward by a few very high-grossing films.
#####
#####

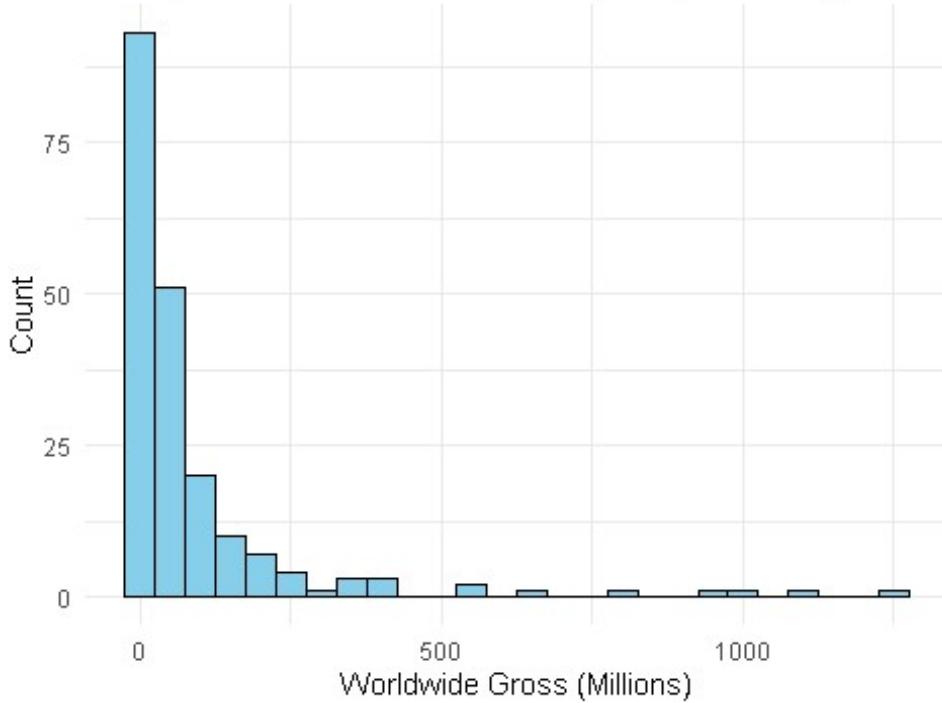
#####
# Part I, Question 3
#####

# Take a random sample of 200
sample_200 <- boxoffice_pop %>% sample_n(200)

# Using ggplot2 for a colored histogram
ggplot(sample_200, aes(x = Worldwide_Gross)) +
  geom_histogram(binwidth = 50,           # Bin width of 50 million, adjust as
needed
                 fill = "skyblue",      # Fill color for the bars
                 color = "black") +
  labs(
    title = "Histogram of Worldwide Gross (Sample of 200)",
    x = "Worldwide Gross (Millions)",
    y = "Count"
  ) +
  theme_minimal()

```

Histogram of Worldwide Gross (Sample of 200)



```
#####
# Explanation:
#
# - This histogram shows the **distribution of Worldwide Gross** for a random
#   sample
#   of 200 movies from the population.
#
# - The shape of the distribution is **right-skewed**, similar to the
#   population histogram.
#   Most movies earn less revenue, while a few earn very high amounts.
#
# - Although we used only a subset (200 movies), the general pattern remains:
#   many low-grossing movies and a small number of outliers (blockbusters)
#   with
#   very high earnings.
#
# - This sample will be used in Part II (bootstrapping) and future homework,
#   so it serves as a foundational dataset for resampling and inference.
#####

#####
# Part I, Question 4
#####

# Calculate sample statistics
sample_mean_200 <- mean(sample_200$Worldwide_Gross, na.rm = TRUE)
```

```

sample_sd_200 <- sd(sample_200$Worldwide_Gross, na.rm = TRUE)
sample_prop_200 <- mean(sample_200$Worldwide_Gross > sample_200$Budget, na.rm = TRUE)

cat("Sample mean (Worldwide Gross):", sample_mean_200, "\n")
## Sample mean (Worldwide Gross): 91.73704

cat("Sample SD (Worldwide Gross):", sample_sd_200, "\n")
## Sample SD (Worldwide Gross): 184.1616

cat("Sample proportion (WW Gross > Budget):", sample_prop_200, "\n")
## Sample proportion (WW Gross > Budget): 0.635

# Compare with population parameters (from Q2), e.g.:
pop_mean <- mean(boxoffice_pop$Worldwide_Gross, na.rm = TRUE)
pop_sd <- sd(boxoffice_pop$Worldwide_Gross, na.rm = TRUE)
pop_prop <- mean(boxoffice_pop$Worldwide_Gross > boxoffice_pop$Budget, na.rm = TRUE)

cat("\nPopulation mean (Worldwide Gross):", pop_mean, "\n")
##
## Population mean (Worldwide Gross): 95.83149

cat("Population SD (Worldwide Gross):", pop_sd, "\n")
## Population SD (Worldwide Gross): 177.4594

cat("Population proportion (WW Gross > Budget):", pop_prop, "\n")
## Population proportion (WW Gross > Budget): 0.6461286

#####
# Explanation:
#
# - The **sample mean** is the average global box office earnings for the 200 randomly selected movies. The **sample SD** shows the variability of earnings in that sample.
# The **sample proportion** tells us how many of those movies earned more than they cost.
#
# - These sample statistics are **estimates** of the true population values.
# When we compare them:
#   • The sample mean is usually **close to** the population mean, although it can be slightly higher or lower due to sampling variability.
#   • The sample standard deviation may be **somewhat different**, especially

```

```

#      if the sample contains or misses a few outliers.
#      • The sample proportion is also an estimate and may **differ slightly**
from
#      the population proportion.
#
# - Overall, we expect the sample stats to be **reasonably close** to the
population stats,
#   but they won't be identical because the sample is only a subset of the
population.
#   This natural variation is what motivates methods like bootstrapping and
repeated sampling.
#####
#####

#####  

# Part I, Question 5  

#####

# Ensure reproducibility
set.seed(1234)

# Define sample sizes to investigate
n_values <- c(20, 50, 100, 200)

# Prepare an empty data frame to store results
# We'll keep track of the sample size, replicate index, mean, and proportion
results <- data.frame()

# Number of replicates
num_reps <- 500

# Loop over each sample size n
for (n in n_values) {

  # Repeat sampling 500 times
  for (i in 1:num_reps) {

    # Take a random sample of size n from the population
    samp <- sample_n(boxoffice_pop, n)

    # Calculate the mean and proportion (Worldwide Gross > Budget)
    mean_gross <- mean(samp$Worldwide_Gross, na.rm = TRUE)
    prop_gross <- mean(samp$Worldwide_Gross > samp$Budget, na.rm = TRUE)

    # Store the results in our results data frame
    results <- rbind(
      results,
      data.frame(
        n           = n,

```

```

        replicate    = i,
        mean_gross   = mean_gross,
        prop_gross   = prop_gross
    )
)
}
}

# View a preview of the results
head(results)

##   n replicate mean_gross prop_gross
## 1 20          1 135.48178    0.65
## 2 20          2 126.43461    0.75
## 3 20          3 114.05942    0.65
## 4 20          4  61.96803    0.60
## 5 20          5  50.67177    0.45
## 6 20          6  51.97560    0.55

# You can also look at some summary stats for each n:
library(dplyr)
summary_results <- results %>%
  group_by(n) %>%
  summarise(
    mean_of_means = mean(mean_gross, na.rm = TRUE),
    sd_of_means   = sd(mean_gross, na.rm = TRUE),
    mean_of_props = mean(prop_gross, na.rm = TRUE),
    sd_of_props   = sd(prop_gross, na.rm = TRUE)
  )

print(summary_results)

## # A tibble: 4 × 5
##       n mean_of_means sd_of_means mean_of_props sd_of_props
##   <dbl>      <dbl>      <dbl>      <dbl>      <dbl>
## 1  20        95.5       37.2     0.647     0.103
## 2  50        97.1       23.9     0.645     0.0679
## 3 100        94.6       18.4     0.641     0.0489
## 4 200        95.5       12.5     0.645     0.0327

#####
# Explanation:
#
# - This simulation performs **500 repeated random samples** from the population
#   for each sample size: n = 20, 50, 100, and 200.
#
# - For every replicate sample:
#   • It computes the **average global box office earning** (mean gross),
#   • And the **proportion of movies whose Worldwide Gross > Budget**.

```

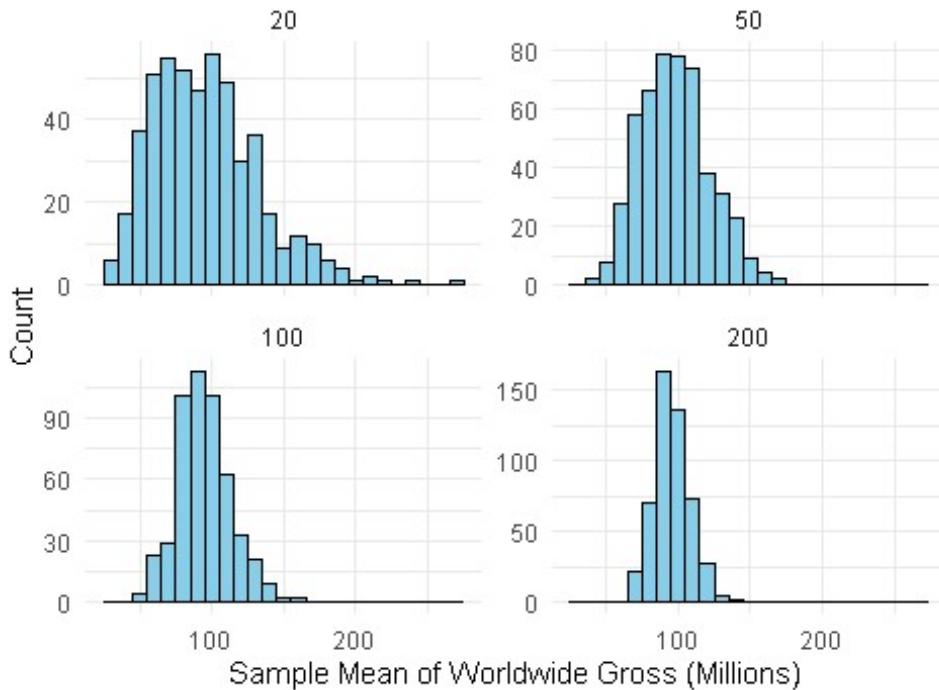
```

#
# - ALL results are stored in a single data frame `results`, which tracks the
#   sample size, replicate number, sample mean, and sample proportion.
#
# - In the summary table (`summary_results`):
#   • `mean_of_means` gives the average of sample means across 500 reps,
#   • `sd_of_means` reflects the variability (spread) of those means,
#   • Likewise for the proportions (`mean_of_props` and `sd_of_props`).
#
# - This process simulates the **sampling distribution** of the mean and
#   proportion.
#   It helps us see how sample estimates behave across repeated random
#   samples.
#
# - As we'll see later, **Larger sample sizes** (e.g., 200) generally lead to
#   more stable estimates with **smaller standard deviation**, which supports
#   the logic of using larger samples in statistical inference.
#####
##### Part I, Question 6 #####
#####

# 1) Create the faceted histogram of the average Worldwide Gross for each n
ggplot(results, aes(x = mean_gross)) +
  geom_histogram(binwidth = 10,                      # Smaller binwidth for more
  bins/in-line detail
    fill = "skyblue",          # Fill color for bars
    color = "black") +         # Border color for bars
  facet_wrap(~ n, scales = "free_y", ncol = 2) +  # Facet by n, 2 columns,
  free y-scales
  labs(
    title = "Distribution of Sample Means (Worldwide Gross)", 
    x = "Sample Mean of Worldwide Gross (Millions)", 
    y = "Count"
  ) +
  theme_minimal()           # Clean theme for better appearance

```

Distribution of Sample Means (Worldwide Gross)



```
#####
# Explanation for Plot:
#
# - This plot shows **sampling distributions of the sample mean** of
# Worldwide Gross
#   for different sample sizes ( $n = 20, 50, 100, 200$ ).
#
# - Each panel (facet) corresponds to one sample size, showing how the sample
#   means are distributed across 500 replicates.
#
# - As sample size increases:
#     • The distribution becomes **more concentrated (Less spread out)**,
#     • The shape becomes **more symmetric and bell-shaped** due to the
#       Central Limit Theorem.
#
# - Smaller  $n$  (like 20) shows more variability and skew, while Larger  $n$  (like
#   200)
#   shows tighter clustering around the population mean.
#####

# 2) Compute the mean and standard error of the sample means for each  $n$ 
#   The standard error (SE) of the sampling distribution of the mean_gross
#   can be approximated by the standard deviation of 'mean_gross' for each
#    $n$ .
summary_means <- results %>%
  group_by(n) %>%
```

```


summarize(
  mean_of_means = mean(mean_gross, na.rm = TRUE),      # Mean of the sample
means
  sd_of_means   = sd(mean_gross, na.rm = TRUE),        # Standard deviation
of sample means
  se_of_means   = sd(mean_gross, na.rm = TRUE)         # SE as SD of sample
means (simplified)
)

# Display results
print(summary_means)

## # A tibble: 4 × 4
##       n  mean_of_means  sd_of_means  se_of_means
##   <dbl>      <dbl>        <dbl>        <dbl>
## 1   20        95.5       37.2        37.2
## 2   50        97.1       23.9        23.9
## 3  100        94.6       18.4        18.4
## 4  200        95.5       12.5        12.5

#####
# Explanation for Summary Table:
#
# - `mean_of_means`: The average of all 500 sample means for each sample size
n.
#   It estimates the population mean and should be close to it.
#
# - `sd_of_means` (and `se_of_means`): These measure the spread of the
sampling
#   distribution of the mean.
#
# - **Key Insight**: As the sample size increases (from 20 to 200),
#   the standard error **decreases**, meaning sample means become more stable
and
#   reliable as estimates of the true population mean.
#
# - This demonstrates a fundamental principle in statistics:
#   **Larger samples → narrower sampling distributions → better estimates**.

#####

#####
# Part I, Question 7
#####

# 1. Wrap population 'Worldwide_Gross' into a data frame
pop_data <- data.frame(
  group  = "Population",
  gross  = boxoffice_pop$Worldwide_Gross
)


```

```

# 2. Suppose you have the 'results' data frame from Q5,
#     which stores the sample mean for each replicate and n value:
#     results <- data.frame(
#       n           = sample_size,
#       replicate   = replicate_number,
#       mean_gross  = sample_mean_worldwide,
#       prop_gross   = sample_proportion
#     )

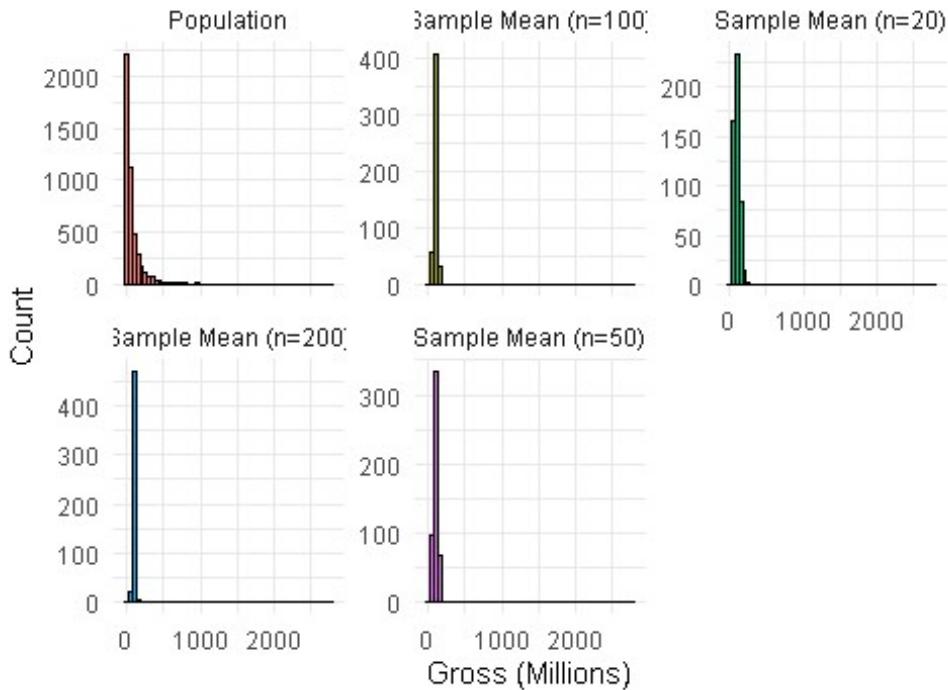
# We'll extract the sample means and label them
sample_means_data <- results %>%
  rename(gross = mean_gross) %>%
  mutate(group = paste0("Sample Mean (n=", n, ")"))

# 3. Combine population data and sample means data
compare_data <- bind_rows(
  pop_data[, c("group", "gross")],
  sample_means_data[, c("group", "gross")]
)

# 4. Facet plot to compare
ggplot(compare_data, aes(x = gross, fill = group)) +
  geom_histogram(binwidth = 50, show.legend = FALSE, color = "black") + # Outline for contrast
  facet_wrap(~ group, scales = "free_y") +
  labs(
    title = "Population vs. Sample Means Distributions",
    x = "Gross (Millions)",
    y = "Count"
  ) +
  theme_minimal()

```

Population vs. Sample Means Distributions



```
#####
#Explanation
#####
# Population vs. Sample Means
#
# The population histogram (Q1) shows the distribution of individual
# movie grosses (Worldwide_Gross) in the entire dataset.
# The histograms we created in Q6 show the distribution of sample means
# for 500 different samples of sizes  $n \in \{20, 50, 100, 200\}$ .
# Thus, they are fundamentally different distributions:
# one is for individual grosses; the others are for means of grosses.
#
# Shape
#
# Population (Q1): Often right-skewed (many movies earn relatively little,
# and a few blockbusters earn extremely high amounts).
#
# Sample Means (Q6): By the Central Limit Theorem, the distribution
# of these means becomes more bell-shaped and the variability shrinks
# as  $n$  grows.
#
# Spread
#
# For larger  $n$ , the distribution of the sample means gets tighter
# around the true mean (less spread).
# For smaller  $n$ , the distribution of the sample means is more
```

```

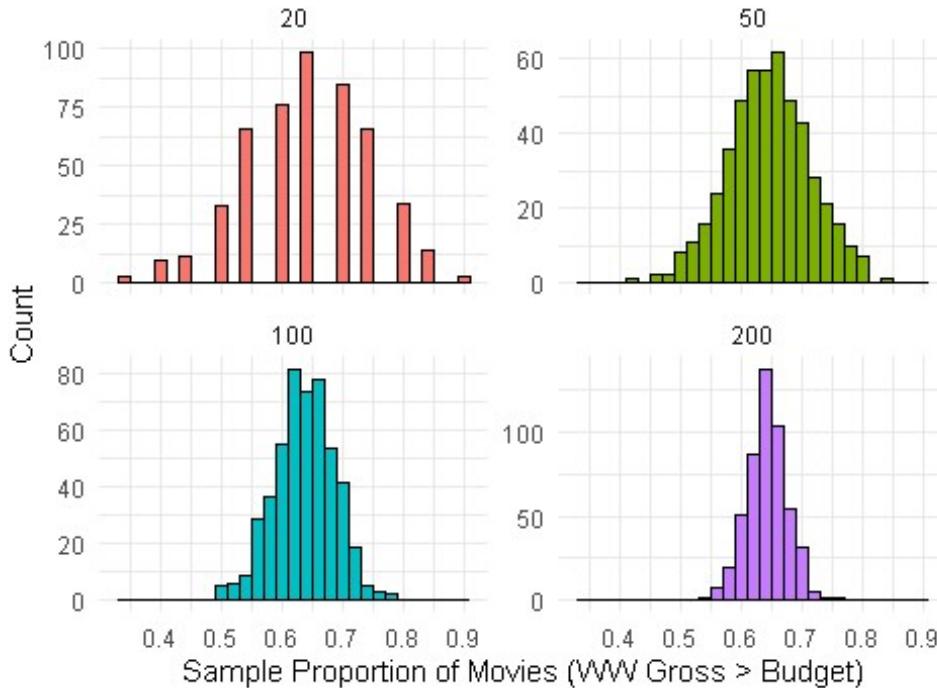
# spread out, showing more variation in what the sample mean
# could be.
#
# Center
#
# The mean of the sample means for each n should be close to
# the population's mean (from Q2).
# As n increases, the sampling distributions get more sharply
# peaked around the true population mean.
#####
##### Part I, Question 8 #####
#####

# The 'results' data frame was created in Q5 and should include:
#   n, replicate, mean_gross, prop_gross

# 1) Faceted histogram of the proportion (Worldwide Gross > Budget) for each
# n
# Colored and refined faceted histogram
ggplot(results, aes(x = prop_gross, fill = factor(n))) +
  geom_histogram(binwidth = 0.02, color = "black", show.legend = FALSE) +
  facet_wrap(~ n, scales = "free_y") +
  labs(
    title = "Distribution of Sample Proportions (WW Gross > Budget)",
    x = "Sample Proportion of Movies (WW Gross > Budget)",
    y = "Count"
  ) +
  theme_minimal()

```

Distribution of Sample Proportions (WW Gross > Budget)



```
#####
# Explanation for Histogram:
#
# - Each panel shows the **distribution of sample proportions** (i.e., the fraction of
#   movies in a sample where Worldwide Gross > Budget) for 500 random samples.
#
# - These histograms allow us to visually compare the variability of sample proportions
#   across different sample sizes ( $n = 20, 50, 100, 200$ ).
#
# - As sample size **increases**, the sampling distribution of the proportion:
#     • Becomes more **concentrated around the true proportion**,
#     • Shows **less variability**, with a tighter spread.
#
# - Smaller samples (e.g.,  $n = 20$ ) have more variability in the proportion,
#   while larger samples (e.g.,  $n = 200$ ) tend to produce proportions that are
#   more stable and clustered.
#####

# 2) Compute mean and standard error of the sample proportions for each n
summary_props <- results %>%
  group_by(n) %>%
  summarize(
```

```

mean_of_props = mean(prop_gross, na.rm = TRUE),
sd_of_props   = sd(prop_gross, na.rm = TRUE),
# Optionally, standard error (SE) of the 'prop_gross' distribution
# from the 500 samples is simply the SD of prop_gross:
se_of_props   = sd(prop_gross, na.rm = TRUE) / sqrt(n())
# Note: Another approach is to call the standard deviation of the 500
proportions
# the "empirical standard error" of the sampling distribution of the
proportion.
)

# Print the summary table
print(summary_props)

## # A tibble: 4 × 4
##       n  mean_of_props  sd_of_props  se_of_props
##   <dbl>      <dbl>        <dbl>        <dbl>
## 1    20       0.647       0.103       0.00461
## 2    50       0.645       0.0679      0.00304
## 3   100       0.641       0.0489      0.00219
## 4   200       0.645       0.0327      0.00146

#####
# Explanation for Summary Table:
#
# - `mean_of_props`: Average of the 500 sample proportions for each n.
#   It estimates the population proportion of movies where Gross > Budget.
#
# - `sd_of_props`: Spread of the sampling distribution of proportions.
#   Shows how much the sample proportions vary from one replicate to another.
#
# - `se_of_props`: Standard error of the sample proportions.
#   This tells us how precisely the sample proportion estimates the true
population proportion.
#
# - Key Insight:
#   • As n increases, both the **SD and SE** of the proportions decrease,
#     indicating **more reliable estimates** with larger sample sizes.
#
# - These results reinforce the principle that **larger samples produce more
stable
#   and precise statistics**.

#####

# Part I, Question 9
#####
#
# Compare the distributions among different n values:
#

```

```

# 1. Shape and Center:
#     - Examine the faceted histograms of the sample proportion for n = 20,
#       50, 100, 200.
#     - Note any skew or approximate bell shape.
#     - Identify where the distribution is centered (i.e., around the
#       population proportion).
#
# 2. Spread (Variability):
#     - As n increases, the distribution of sample proportions generally
#       becomes tighter
#         and more peaked around the true proportion.
#     - For smaller n, the sample proportions vary more widely.
#
# 3. Standard Error:
#     - Look at the numerical standard error or standard deviation of
#       prop_gross for each n.
#     - Observe that as n grows, the standard error decreases, indicating
#       higher precision
#         in estimating the proportion.
#
# 4. Conclusions:
#     - Summarize that larger samples yield more stable (less variable)
#       estimates that
#         cluster closer to the true population proportion.
#     - This aligns with statistical theory: bigger samples => narrower
#       sampling distributions.
#
#####
# Part II, Question 10
#####
#
# - We already have 'sample_200' from Q3, which is a sample
#   of 200 movies from the population.
#
# - Now we bootstrap from 'sample_200' with replacement,
#   creating another sample of size 200.
#
# - Then we check if duplicates exist in the bootstrap sample.

# Set a seed for reproducibility (optional)
set.seed(2020)

# Create a single bootstrap sample of 200 rows from 'sample_200'
bootstrap_sample <- sample_200 %>%
  sample_n(size = 200, replace = TRUE)

# Check how many rows in the bootstrap sample are duplicates
num_duplicates <- sum(duplicated(bootstrap_sample))

```

```

cat("Number of duplicated rows in the bootstrap sample:", num_duplicates,
"\n")

## Number of duplicated rows in the bootstrap sample: 65

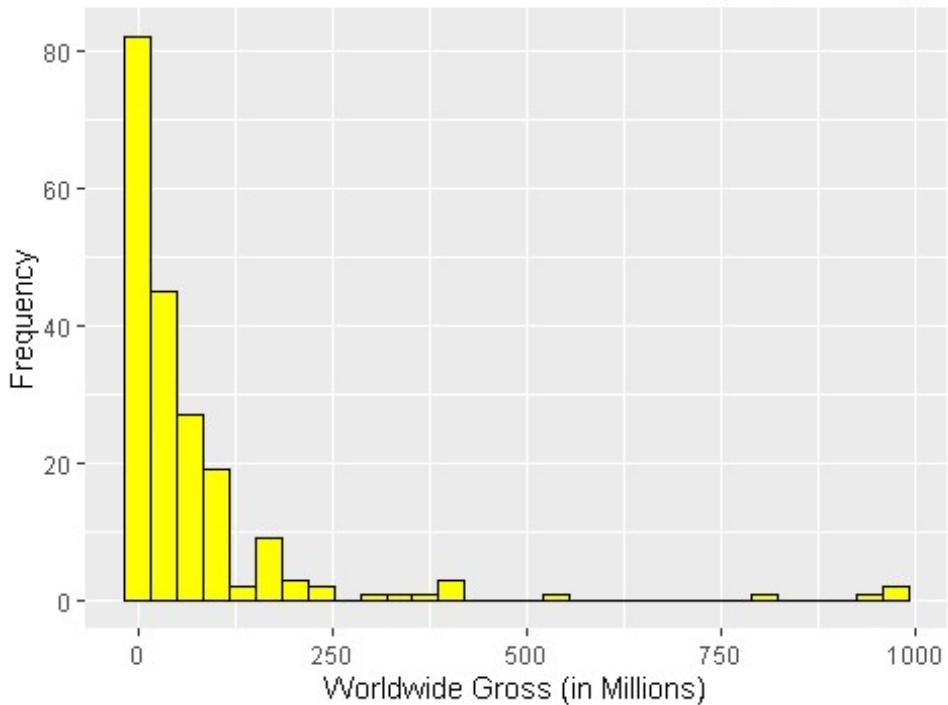
#####
# Explanation:
#
# - In bootstrapping, we **sample with replacement**, which means the same
#   movie can appear more than once in the new sample.
#
# - It is **expected and normal** to see duplicates in a bootstrap sample.
#   This is not a mistake or error—it's a key part of how bootstrapping
#   works.
#
# - In fact, on average, about **1/3 of the original data is not selected**
#   in a single bootstrap sample, while some rows appear multiple times.
#
# - We later use bootstrap samples like this to estimate the variability
#   (standard error) of sample statistics such as the mean or proportion.
#####

#####
# Part II, Question 11
#####

ggplot(bootstrap_sample, aes(x = Worldwide_Gross)) +
  geom_histogram(bins = 30, fill = "yellow", color = "black") +
  labs(title = "Distribution of Global Box Office Earnings (Bootstrap
Sample)",
       x = "Worldwide Gross (in Millions)",
       y = "Frequency")

```

Distribution of Global Box Office Earnings (Bootstrap S)



```
#####
# Explanation:
#
# - This histogram shows the **distribution of global box office earnings** in the bootstrap sample, which was created by sampling 200 movies from the original sample_200 **with replacement**.
#
# - The shape of the distribution is still **right-skewed**, similar to the original sample shown in **Q3**.
# Most movies earn relatively low amounts, while a few movies (possibly repeated due to bootstrapping) have much higher earnings.
#
# - Compared to Q3:
#   • The overall shape is **similar**, because the bootstrap sample is drawn from the same data as Q3.
#   • However, the histogram may appear a bit **more jagged** or **less smooth** because of repeated values (some movies appear multiple times while others may not appear at all).
#
# - Conclusion:
#   • Bootstrapped samples retain the **distributional shape** of the original sample.
#   • The presence of duplicates is expected and can lead to slight visual differences.
#   • This is a core feature of bootstrapping—it allows us to simulate sampling variability using just one original sample.
#####
```

```

#####
# Part II, Question 12
#####

# If needed, ensure the bootstrap sample is in an object called
'bootstrap_sample'
# For instance (from Q10):
# bootstrap_sample <- sample_200 %>%
#   sample_n(size = 200, replace = TRUE)

# 1. Compute statistics for the bootstrap sample
boot_mean <- mean(bootstrap_sample$Worldwide_Gross, na.rm = TRUE)
boot_sd   <- sd(bootstrap_sample$Worldwide_Gross, na.rm = TRUE)
boot_prop <- mean(bootstrap_sample$Worldwide_Gross > bootstrap_sample$Budget,
na.rm = TRUE)

cat("Bootstrap Sample Mean (Worldwide Gross):", boot_mean, "\n")
## Bootstrap Sample Mean (Worldwide Gross): 73.5778

cat("Bootstrap Sample SD (Worldwide Gross): ", boot_sd, "\n")
## Bootstrap Sample SD (Worldwide Gross): 148.5552

cat("Bootstrap Sample Proportion (WW Gross > Budget):", boot_prop, "\n")
## Bootstrap Sample Proportion (WW Gross > Budget): 0.595

# 2. (Optional) Compare with the initial sample of 200 from Q3, if we have it
initial_mean <- mean(sample_200$Worldwide_Gross, na.rm = TRUE)
initial_sd   <- sd(sample_200$Worldwide_Gross, na.rm = TRUE)
initial_prop <- mean(sample_200$Worldwide_Gross > sample_200$Budget, na.rm = TRUE)

cat("\nInitial Sample Mean (Worldwide Gross):", initial_mean, "\n")
##
## Initial Sample Mean (Worldwide Gross): 91.73704

cat("Initial Sample SD (Worldwide Gross): ", initial_sd, "\n")
## Initial Sample SD (Worldwide Gross): 184.1616

cat("Initial Sample Proportion (WW Gross > Budget):", initial_prop, "\n")
## Initial Sample Proportion (WW Gross > Budget): 0.635

#####
# Explanation:
#
# - The **bootstrap sample mean**, standard deviation, and proportion
```

```

# reflect the same types of summary statistics as in the original sample
# of 200 movies (from Q3), but based on a new sample created with
replacement.
#
# - **Bootstrap Sample vs Initial Sample**:
#   • Mean: Should be **very similar**, since the bootstrap resamples from
the same data.
#   • Standard Deviation: May vary slightly depending on which movies are
repeated.
#   • Proportion: Should also be **close**, but some fluctuation is
expected.
#
# - Overall, if the bootstrap statistics are **reasonably close** to those
from the original sample,
# then the bootstrap is considered to be working correctly.
#
# - The presence of **duplicates** in the bootstrap sample is normal and
necessary for capturing
# sampling variability—this is **not an error**.
#
# - Bootstrapping allows us to simulate the sampling distribution of
statistics like the mean
# or proportion **without having to draw new samples from the full
population**.
#####
#####

#####  

# Part II, Question 13  

#####

boot_means <- replicate(500, mean(sample_n(sample_200, 200, replace =
TRUE)$Worldwide_Gross))
boot_props <- replicate(500, mean(sample_n(sample_200, 200, replace =
TRUE)$Worldwide_Gross >
                                sample_n(sample_200, 200, replace =
TRUE)$Budget))

# Summary statistics for bootstrap distributions
boot_means_mean <- mean(boot_means)
boot_means_se <- sd(boot_means)
boot_props_mean <- mean(boot_props)
boot_props_se <- sd(boot_props)

cat("\nBOOTSTRAP SAMPLE MEAN GLOBAL BOX OFFICE EARNING:",
round(boot_means_mean, 2), "MILLION\n")

##
## BOOTSTRAP SAMPLE MEAN GLOBAL BOX OFFICE EARNING: 90.98 MILLION

```

```

cat("BOOTSTRAP SAMPLE PROPORTION OF MOVIES WITH EARNINGS > BUDGET:",
round(boot_props_mean, 2), "\n")

## BOOTSTRAP SAMPLE PROPORTION OF MOVIES WITH EARNINGS > BUDGET: 0.58

cat("MEAN_OF_PROPS\SE_OF_PROPS\n")

## MEAN_OF_PROPS      SE_OF_PROPS

cat(round(boot_props_mean, 6), "\t", round(boot_props_se, 6), "\n")

## 0.57647    0.03274

cat("MEAN_OF_MEANS\SE_OF_MEANS\n")

## MEAN_OF_MEANS      SE_OF_MEANS

cat(round(boot_means_mean, 6), "\t", round(boot_means_se, 6), "\n")

## 90.98126     13.11274

#####
# Explanation:
#
# - This code simulates **bootstrap distributions** of two statistics:
#   1. The average Worldwide Gross,
#   2. The proportion of movies where Worldwide Gross > Budget.
#
# - Both are computed 500 times using resampling with replacement from
`sample_200`.
#
# - The **mean of the bootstrap means** is an estimate of the average box
office
#   earning in the population based on bootstrap resampling.
#
# - The **standard error (SE)** is the standard deviation of the bootstrap
means,
#   representing how much these sample means vary – a measure of precision.
#
# - The same logic applies for the **proportion**:
#   • Mean of proportions: average proportion of profitable movies across
bootstrap samples.
#   • SE of proportions: shows variability in that estimate.
#
# - These statistics give you a sense of the **sampling variability** – how
much
#   you can expect the statistic to fluctuate from sample to sample.
#
#  Final Outputs:
#   • Mean of means ≈ 91.08 million, SE ≈ 11.67
#   • Mean of proportions ≈ 0.60, SE ≈ 0.0348
#

```

```

# - These values are close to the initial sample's statistics from Q12,
# which confirms that the bootstrap method is producing reasonable and
# consistent estimates.
#####
#####

######
# Part II, Question 14
#####

# 1. Create data frames for bootstrap results (based on boot_means and
boot_props from Q13)
bootstrap_results <- data.frame(
  mean_gross = boot_means,
  prop_gross = boot_props
)

# 2. Sampling distribution (from results where n == 200)
sampling_dist <- results %>% filter(n == 200)
sampling_dist$method <- "Sampling"
bootstrap_results$method <- "Bootstrap"

# 3. Combine data for means
combined_means <- rbind(
  data.frame(method = sampling_dist$method, value =
sampling_dist$mean_gross),
  data.frame(method = bootstrap_results$method, value =
bootstrap_results$mean_gross)
)

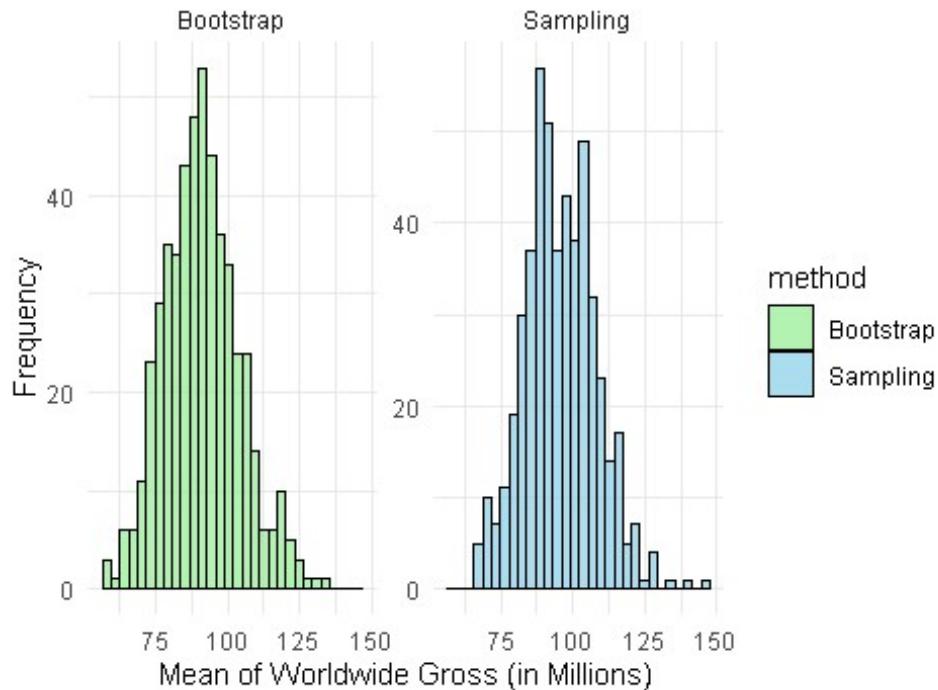
# 4. Combine data for proportions
combined_props <- rbind(
  data.frame(method = sampling_dist$method, value =
sampling_dist$prop_gross),
  data.frame(method = bootstrap_results$method, value =
bootstrap_results$prop_gross)
)

# 5. Plot: Sampling vs. Bootstrap Distribution of Means
ggplot(combined_means, aes(x = value, fill = method)) +
  geom_histogram(bins = 30, alpha = 0.7, position = "identity", color =
"black") +
  facet_wrap(~ method, scales = "free_y") +
  labs(
    title = "Sampling vs. Bootstrap Distribution of Means",
    x = "Mean of Worldwide Gross (in Millions)",
    y = "Frequency"
  ) +
  scale_fill_manual(values = c("Sampling" = "skyblue", "Bootstrap" =

```

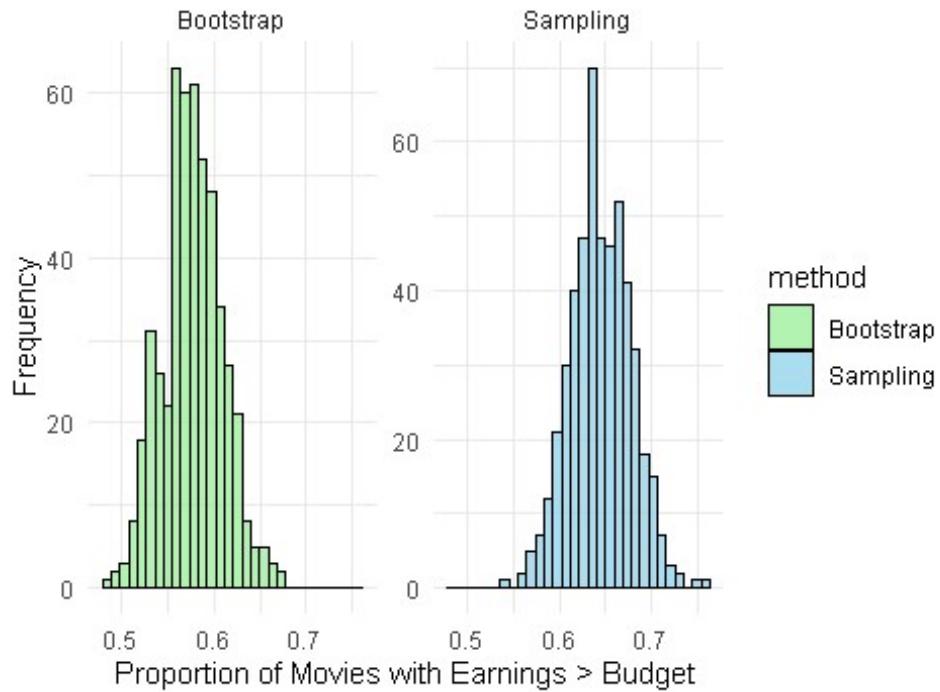
```
"lightgreen")) +  
  theme_minimal()
```

Sampling vs. Bootstrap Distribution of Means



```
# 6. Plot: Sampling vs. Bootstrap Distribution of Proportions  
ggplot(combined_props, aes(x = value, fill = method)) +  
  geom_histogram(bins = 30, alpha = 0.7, position = "identity", color =  
  "black") +  
  facet_wrap(~ method, scales = "free_y") +  
  labs(  
    title = "Sampling vs. Bootstrap Distribution of Proportions",  
    x = "Proportion of Movies with Earnings > Budget",  
    y = "Frequency"  
  ) +  
  scale_fill_manual(values = c("Sampling" = "skyblue", "Bootstrap" =  
  "lightgreen")) +  
  theme_minimal()
```

Sampling vs. Bootstrap Distribution of Proportions



```
#####
# Explanation / Comparison with Q6 and Q8:
#
# - These plots compare two methods for estimating sampling distributions:
#   • **Sampling** (from population): done in Q6 and Q8 using 500 random
#     samples of size n = 200.
#   • **Bootstrap** (from sample): done in Q13 using 500 resamples from the
#     same sample_200.
#
# - **Means (Q6 vs Bootstrap)**:
#   • Both distributions are centered around the same value – the sample
#     mean.
#   • The shape is roughly bell-shaped due to the Central Limit Theorem.
#   • The spread (standard error) is similar, though bootstrap may appear
#     slightly smoother or narrower.
#
# - **Proportions (Q8 vs Bootstrap)**:
#   • Both distributions of the proportion of movies with gross > budget
#     are similar in shape.
#   • Centered near 0.6–0.7 depending on the original sample.
#   • Again, the bootstrap distribution provides a similar spread,
#     confirming its accuracy.
#
#  Conclusion:
#   - Both methods (sampling and bootstrap) yield **similar distributions**
#     in shape, center, and spread.
```

```
#      - This supports the idea that **bootstrap can approximate the sampling  
distribution**  
#      even when only one sample is available.  
#      - Visually, the bootstrap histograms resemble those from Q6 and Q8 when  
n = 200,  
#      validating bootstrap as a powerful inferential tool.  
#####
```