**Gisma University of Applied Sciences**

**Assessment Submission Form**

| | |
|---|---|
| **Student Number** (If this is group work, please include the student numbers of all group participants) | GH1021715 |
| **Assessment Title** | Billionaires' Net Worth Analysis |
| **Module Code** | B105 |
| **Module Title** | Applied Statistical Modelling |
| **Module Tutor** | William Baker Morrison |
| **Date Submitted** | 26/09/2024 |

# Billionaires' Net Worth Analysis

Stanley Osei-Wusu (GH1021715)

September 26, 2024

# Contents

# 1 Introduction

This report presents a comprehensive statistical analysis of the Billionaires' dataset, which contains information about 2,640 billionaires, including their net worth, age, gender, country of residence, and industry. The goal of this analysis is to answer key business questions, test specific hypotheses, and interpret statistical findings to provide actionable insights for stakeholders.

## 1.1 Key Business Questions

1. What are the key factors that determine the net worth of billionaires?

2. How do industry and age affect a billionaire's net worth?

3. Are there significant differences in net worth based on gender?

4. Can we predict whether a billionaire belongs to the technology industry based on their age, gender, and net worth?

# 2 Hypotheses

1. **H0 (Null Hypothesis)**: There is no significant relationship between age, industry, or gender and the net worth of billionaires.

2. **H1 (Alternative Hypothesis)**: There is a significant relationship between age, industry, or gender and the net worth of billionaires.

A second hypothesis for the logistic regression is:

- **H0**: A billionaire's membership in the technology industry is not significantly predicted by their age, gender, or net worth.

- **H1**: A billionaire's membership in the technology industry can be significantly predicted by their age, gender, and net worth.

# 3 Data Preparation

The data was sourced from the *Billionaires Statistics Dataset* on Kaggle, containing 2,640 records. Several steps were performed to ensure the data was clean and ready for analysis.

## 3.1 Data Cleaning

- Missing values in the `age` column were imputed using the median age.

- The `finalWorth` column was cleaned by removing non-numeric characters.

- Categorical variables such as `gender`, `industry`, and `country` were converted to factors.

- Rows with missing or zero values in critical variables such as `finalWorth` were removed.

## 3.2 Sampling

To ensure a balanced analysis, the full dataset was used without sampling, as the size of the dataset (2,640 records) was manageable for the analysis.

# 4 Exploratory Data Analysis

## 4.1 Descriptive Statistics

The dataset contains a wide range of billionaire net worths, with the minimum being $1 billion and the maximum reaching $211 billion. The mean net worth is approximately $4.6 billion, and the distribution is heavily skewed toward lower values.

```
Summary statistics:
Min. Net Worth: 1 billion
Max. Net Worth: 211 billion
Median Net Worth: 2.3 billion
Mean Net Worth: 4.6 billion
Age Range: 18 - 101 years
```

## 4.2 Distribution of Net Worth

The net worth distribution is right-skewed, with the majority of billionaires having a net worth below $10 billion, while a small number hold extremely high values.

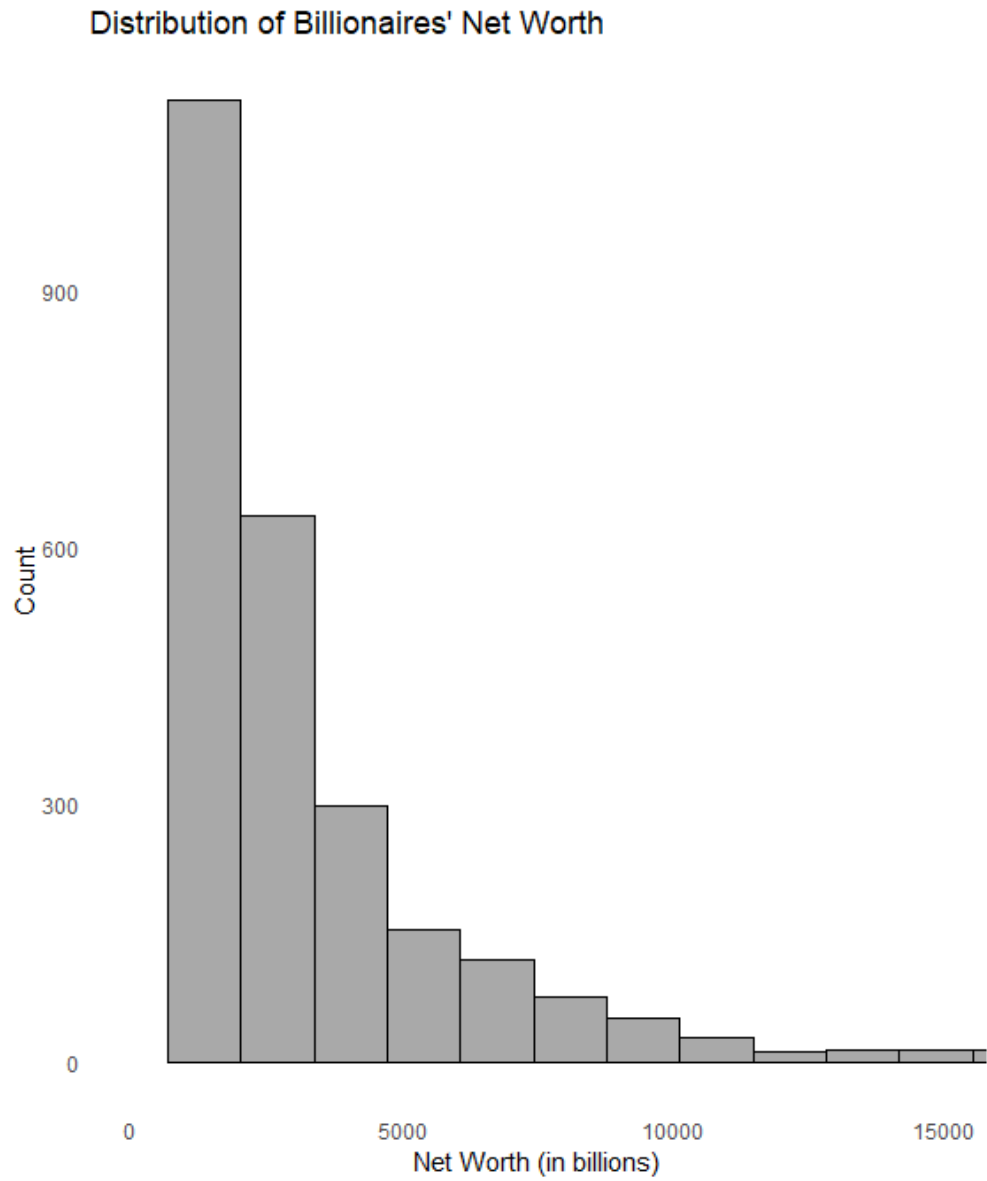Figure 1: Distribution of Billionaires' Net Worth

## 4.3   Age vs Net Worth

A scatter plot of age vs. net worth reveals a weak positive correlation, with older individuals generally having more wealth.



Figure 2: Age vs Net Worth

## 4.4   Net Worth by Industry

The boxplot below shows the distribution of net worth across different industries, highlighting that sectors like Energy and Finance tend to have a broader range of wealth.



Figure 3: Net Worth by Key Industries (Filtered)

## 4.5    Average Net Worth by Gender

The bar chart indicates that male billionaires have a slightly higher average net worth compared to their female counterparts, but the difference is not statistically significant.



Figure 4: Average Net Worth by Gender

# 5 Inferential Statistics
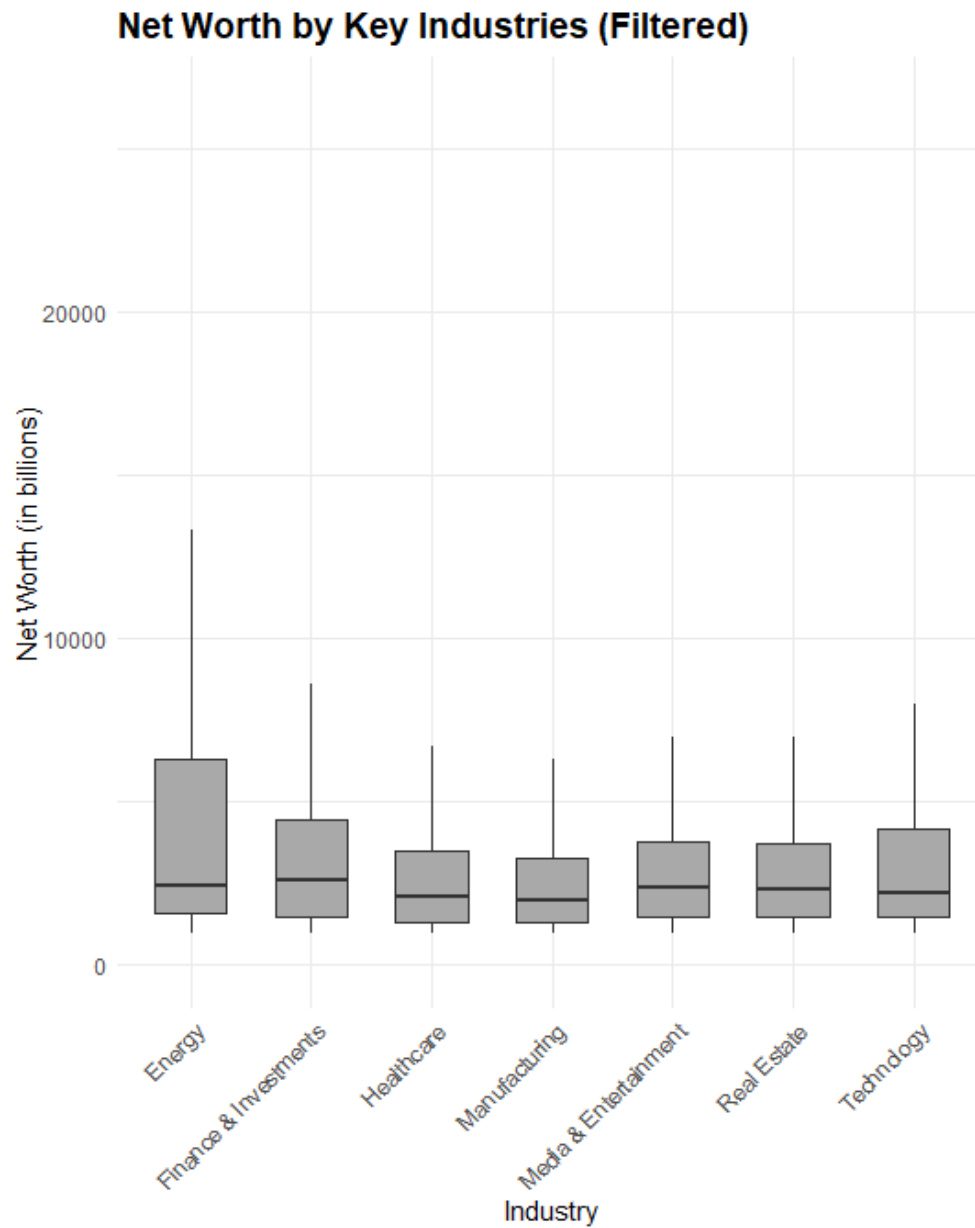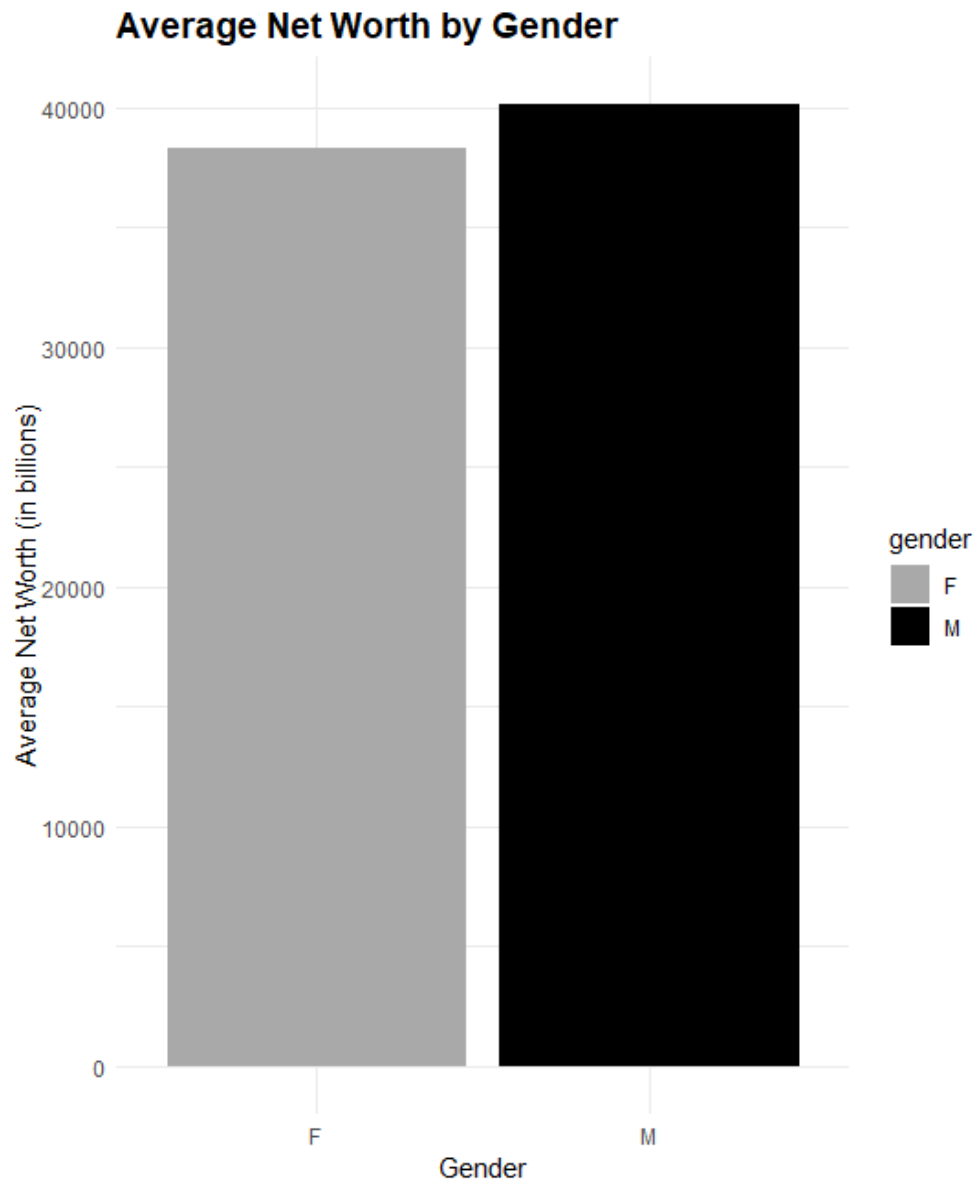
## 5.1 Hypothesis Testing: T-test for Gender and Net Worth

A Welch Two Sample t-test was conducted to compare the mean net worth between male and female billionaires.

```
Welch Two Sample t-test
t = 0.77808, df = 422.13, p-value = 0.437
Mean in group F = 4023.952, Mean in group M = 3803.423
```

Since the p-value is greater than 0.05, we fail to reject the null hypothesis, indicating no significant difference in net worth based on gender.

## 5.2 ANOVA for Industry

An ANOVA test was performed to determine if there are significant differences in net worth across industries.

```
ANOVA results:
F-value: 3.316, p-value: 4.87e-06
```

The p-value is significant, suggesting that industry has a statistically significant effect on net worth.

## 5.3 Logistic Regression

A logistic regression model was used to predict whether a billionaire belongs to the technology industry based on age, gender, and net worth.

```
Logistic Regression Results:
Age: p < 0.001, Gender: p = 0.027, Net Worth: p = 0.366 (not significant)
```

The results indicate that age ($p < 0.001$) and gender ($p = 0.027$) are statistically significant predictors of technology industry membership, while net worth ($p = 0.366$) is not a significant predictor.

5.4 ROC Curve

## 5.4 ROC Curve

The ROC curve for the logistic regression model shows an AUC of approximately 0.7, indicating a moderately good fit for predicting membership in the technology industry.
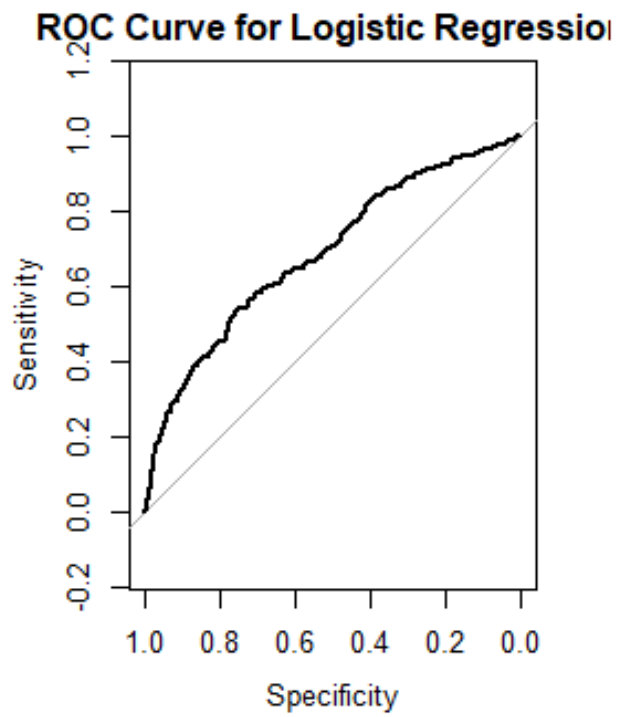
Figure 5: ROC Curve for Logistic Regression

# 6 Conclusion

The analysis provided several key insights:

- Age and industry significantly influence a billionaire's net worth.

- There is no significant difference in net worth between male and female billionaires.

- Logistic regression models indicate that age and gender are predictive of whether a billionaire belongs to the technology industry, but net worth is not.

## 6.1 Limitations

The analysis was limited by the lack of more granular data on billionaire behaviors or financial strategies. Additionally, there may be unobserved variables that impact wealth accumulation not included in this dataset.

## 6.2 Future Work

Future analyses could explore non-linear models and include additional variables such as investment behavior, philanthropy, and inheritance patterns, which may provide deeper insights into wealth accumulation among billionaires.

# 7 GitHub Repository

The full R code and scripts for this analysis are available on GitHub at the following link: https://github.com/stanleymay20/Bi

# 8 Screenshots of R coding

```r
 1  # Load required libraries
 2  library(tidyverse)
 3  library(dplyr)
 4  library(ggplot2)
 5
 6  # Load the dataset
 7  data <- read.csv("Billionaires Statistics Dataset.csv")
 8
 9  # Check for missing values
10  summary(data)
11
12  # Handle missing values (example: replace missing age with median)
13  data$age[is.na(data$age)] <- median(data$age, na.rm = TRUE)
14
15  # Convert categorical variables to factors
16  data$gender <- as.factor(data$gender)
17  data$industry <- as.factor(data$industries)
18  data$country <- as.factor(data$country)
19
20  # Remove any outliers from the finalWorth column for cleaner visualization
21  qnt <- quantile(data$finalWorth, probs = c(0.25, 0.75), na.rm = TRUE)
22  caps <- quantile(data$finalWorth, probs = c(0.01, 0.99), na.rm = TRUE)
23  IQR <- qnt[2] - qnt[1]
24
25  data <- data %>%
26    filter(finalWorth >= caps[1] & finalWorth <= caps[2])
27
28  # Creating a histogram of billionaire net worth
29  ggplot(data, aes(x = finalWorth)) +
30    geom_histogram(bins = 30, fill = "darkgrey", color = "black") +
31    labs(title = "Distribution of Billionaires' Net Worth",
32         x = "Net Worth (in billions)",
33         y = "Count") +
34    theme_minimal() +
35    coord_cartesian(xlim = c(0, 15000)) +  # Adjust the limit based on actual data to focus on common values
36    theme(panel.grid.major = element_blank(),  # Remove major grid lines
37          panel.grid.minor = element_blank())   # Remove minor grid lines
38
```

```r
39  # Scatter plot of Age vs Net Worth
40  ggplot(data, aes(x = age, y = finalWorth)) +
41    geom_point(alpha = 0.5, size = 2, color = "darkgrey") +  # Adjust transparency (alpha) and point size
42    labs(title = "Age vs Net Worth",
43         x = "Age",
44         y = "Net Worth (in billions)") +
45    theme_minimal() +
46    coord_cartesian(ylim = c(0, 15000)) +  # Set y-axis limits to focus on main distribution
47    theme(panel.grid.major = element_blank(),  # Remove major grid lines
48          panel.grid.minor = element_blank())   # Remove minor grid lines
49
50  # Remove outliers beyond the 99th percentile
51  q99 <- quantile(data$finalWorth, 0.99)
52  # Select seven key industries to focus on
53  key_industries <- c("Technology", "Finance & Investments", "Healthcare", "Manufacturing",
54                      "Media & Entertainment", "Real Estate", "Energy")
55
56  data_filtered_key <- data_filtered %>% filter(industry %in% key_industries)
57
58  # Boxplot: Net Worth by Key Industries with cleaner visualization
59  ggplot(data_filtered_key, aes(x = industry, y = finalWorth)) +
60    geom_boxplot(fill = "darkgrey", outlier.shape = NA, width = 0.6) +  # Wider boxplot for better visualization
61    coord_cartesian(ylim = c(0, q99)) +  # Limit y-axis to remove extreme outliers
62    theme_minimal() +
63    labs(title = "Net Worth by Key Industries (Filtered)",
64         x = "Industry", y = "Net Worth (in billions)") +
65    theme(axis.text.x = element_text(angle = 45, hjust = 1),  # Rotate x-axis labels
66          plot.title = element_text(size = 14, face = "bold"))  # Bold title for emphasis
67
68
69  # Bar plot: Average Net Worth by Gender with black and grey colors
70  ggplot(data, aes(x = gender, y = finalWorth, fill = gender)) +
71    geom_bar(stat = "identity", position = "dodge") +
72    scale_fill_manual(values = c("F" = "darkgrey", "M" = "black")) +
73    theme_minimal() +
74    labs(title = "Average Net Worth by Gender",
75         x = "Gender", y = "Average Net Worth (in billions)") +
76    theme(plot.title = element_text(size = 14, face = "bold"))
```

File   Edit   Code   View   Plots   Session   Build   Debug   Profile   Tools   Help

Go to file/function      Addins

Billionaires_Net_Worth_Analysis.Rmd ×   df ×   data ×   Billionaires_Net_Worth_Analysis.Rmd* ×   Untitled3* ×   Untitled2* ×

Source on Save                                                                              Run      Source

```r
 73     theme_minimal() +
 74     labs(title = "Average Net Worth by Gender",
 75          x = "Gender", y = "Average Net Worth (in billions)") +
 76     theme(plot.title = element_text(size = 14, face = "bold"))
 77
 78   # Hypothesis testing: T-test for gender and net worth
 79   t.test(finalWorth ~ gender, data = data)
 80
 81   # Hypothesis testing: ANOVA for industry and net worth
 82   anova_model <- aov(finalWorth ~ industry, data = data)
 83   summary(anova_model)
 84
 85   # Example: Combine industries with fewer than a certain number of observations
 86   data$industry_combined <- ifelse(table(data$industry)[data$industry] < 10, "Other", as.character(data$industry))
 87
 88   # Update table
 89   table_gender_industry_combined <- table(data$gender, data$industry_combined)
 90
 91   # Perform Chi-Square test on the updated table
 92   chisq.test(table_gender_industry_combined)
 93   chisq.test(table_gender_industry, simulate.p.value = TRUE, B = 10000)
 94   # Linear regression model
 95   lm_model <- lm(finalWorth ~ age + gender + industry + country, data = data)
 96   summary(lm_model)
 97
 98   # Diagnostic Plots
 99   par(mfrow = c(2, 2))
100   plot(lm_model)
101   # Logistic regression model to predict if the billionaire is in the technology industry
102   data$tech_industry <- ifelse(data$industry == "Technology", 1, 0)
103   logit_model <- glm(tech_industry ~ age + gender + finalWorth, data = data, family = binomial)
104   summary(logit_model)
105
106   # Plot ROC curve
107   install.packages("pROC")
108   library(pROC)
109   roc_curve <- roc(data$tech_industry, fitted(logit_model))
110   plot(roc_curve, main = "ROC Curve for Logistic Regression")
111
```

28:48   (Top Level)                                                                          R Script

13