

Final_Jeffrey

Jack Jeffrey

Final

```
# Set working directory and CRAN mirror
setwd("/Users/jackjeffrey/Documents/Poli502_Jeffrey/final502")
options(repos = c(CRAN = "https://cran.rstudio.com/"))

# Install and load required packages
required_packages <- c("dplyr", "ggplot2", "knitr", "kableExtra", "stargazer", "pROC", "tidy")

# Check and install missing packages
for (pkg in required_packages) {
  if (!requireNamespace(pkg, quietly = TRUE)) {
    install.packages(pkg)
  }
  library(pkg, character.only = TRUE)
}
```

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

Attaching package: 'kableExtra'

The following object is masked from 'package:dplyr':

group_rows

Please cite as:

Hlavac, Marek (2022). stargazer: Well-Formatted Regression and Summary Statistics Tables.

R package version 5.2.3. <https://CRAN.R-project.org/package=stargazer>

Type 'citation("pROC")' for a citation.

Attaching package: 'pROC'

The following objects are masked from 'package:stats':

cov, smooth, var

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v forcats   1.0.0      v stringr   1.5.1
v lubridate 1.9.3      v tibble    3.2.1
v purrr     1.0.2      v tidyr     1.3.1
v readr     2.1.5
```

```
-- Conflicts ----- tidyverse_conflicts() --
```

```
x dplyr::filter()      masks stats::filter()
x kableExtra::group_rows() masks dplyr::group_rows()
x dplyr::lag()          masks stats::lag()
```

i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become

#1.

```
# Load the Titanic dataset
td <- read.csv("/Users/jackjeffrey/Documents/Pol502_Jeffrey/Data/titanic2.csv")

# View the structure of the dataset
str(td)
```



```
# Check for missing values in the relevant columns
sum(is.na(td$female)) # Missing values in gender
```

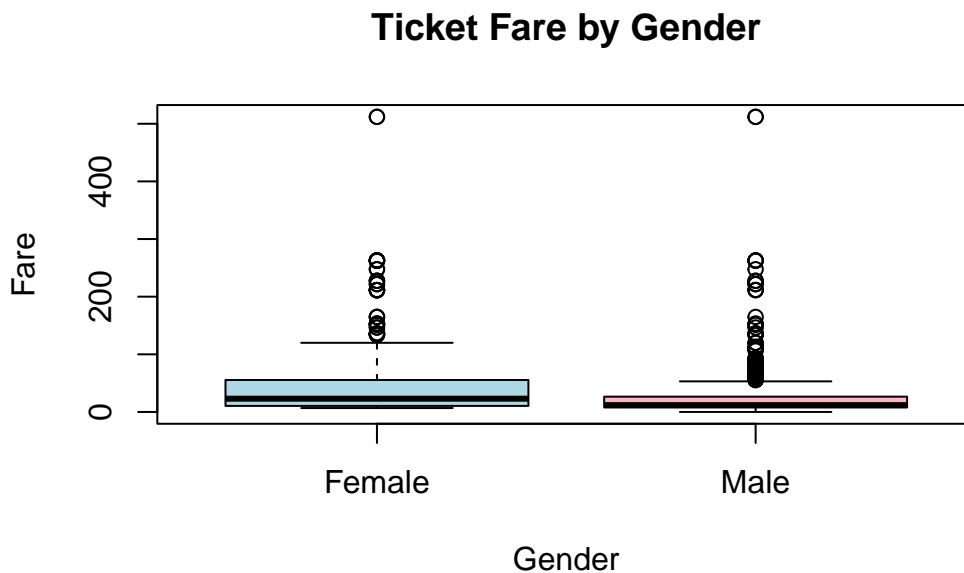
```
[1] 0
```

```
sum(is.na(td$fare)) # Missing values in fare
```

```
[1] 1
```

```
# Clean dataset by removing rows with NA values in the relevant columns
td_clean <- td[complete.cases(td[, c("survived", "fare", "female", "child")]), ]

# Boxplot to explore the relationship between gender and fare
boxplot(fare ~ female, data = td,
        main = "Ticket Fare by Gender",
        xlab = "Gender",
        ylab = "Fare",
        col = c("lightblue", "lightpink"))
```



```
# View summary statistics of fare by gender
tapply(td$fare, td$female, summary)
```

\$Female

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
6.75	10.50	23.00	46.20	55.33	512.00

\$Male

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0.000	7.876	11.887	26.154	26.550	512.000	1

2.

```
# Perform a two-sample t-test to compare the fares between female and male passengers
t_test_result <- t.test(fare ~ female, data = td)
t_test_result
```

Welch Two Sample t-test

data: fare by female

t = 6.1168, df = 701.7, p-value = 1.585e-09

alternative hypothesis: true difference in means between group Female and group Male is not 0

95 percent confidence interval:

13.60960 26.47613

sample estimates:

mean in group Female	mean in group Male
46.19668	26.15382

```
# Interpretation of results
```

```
# The p-value of 1.585e-09 is extremely small,
```

```
# the difference in means between female and male passengers
```

```
# is statistically significant so the null hypothesis can be rejected.
```

```
# The confidence interval for the difference in means
```

```
# does not contain 0, which confirms a significant difference
```

```
# between the fares for females and males.
```

#3a.

```
# Estimate the first logit model (original fare)
```

```
model1 <- glm(survived ~ fare + female + child, data = td_clean, family = "binomial")
model1
```

```
Call: glm(formula = survived ~ fare + female + child, family = "binomial",
          data = td_clean)
```

Coefficients:

(Intercept)	fare	femaleMale	childChild
0.640380	0.009279	-2.362196	0.675223

Degrees of Freedom: 1044 Total (i.e. Null); 1041 Residual

Null Deviance: 1414

Residual Deviance: 1058 AIC: 1066

```
# Clean dataset to remove rows with fare <= 0
td_clean <- td_clean[td_clean$fare > 0, ]

# Estimate the second logit model (log of fare)
model2 <- glm(survived ~ log(fare) + female + child, data = td_clean, family = "binomial")
model2
```

```
Call: glm(formula = survived ~ log(fare) + female + child, family = "binomial",
          data = td_clean)
```

Coefficients:

(Intercept)	log(fare)	femaleMale	childChild
-0.7479	0.5585	-2.3341	0.5597

Degrees of Freedom: 1036 Total (i.e. Null); 1033 Residual

Null Deviance: 1404

Residual Deviance: 1036 AIC: 1044

```
# Display the results using stargazer
stargazer(model1, model2, type = "text", title = "Logit Models of Titanic Passenger Survival")
```

Logit Models of Titanic Passenger Survival

```
=====
Dependent variable:
-----
survived
(1)          (2)
```

```

-----
fare                0.009***
                   (0.002)

log(fare)                                0.559***
                                         (0.082)

femaleMale          -2.362***           -2.334***
                   (0.156)             (0.157)

childChild           0.675***           0.560**
                   (0.235)             (0.234)

Constant             0.640***           -0.748***
                   (0.137)             (0.280)

-----
Observations         1,045              1,037
Log Likelihood        -528.894          -518.122
Akaike Inf. Crit.    1,065.788          1,044.245
=====
Note:                 *p<0.1; **p<0.05; ***p<0.01

```

3b.

```

# Based on the model fit statistics,
# Model 2 performs better. This is because it has a
# lower Akaike Information Criterion (AIC) of 976.974 compared
# to Model 1's AIC of 1000.716,
# indicating a better trade-off between model fit and complexity.

```

3c.

```

# Calculate the median values of the independent variables (female and child)
median_values <- data.frame(female = "Female", child = "Adult", fare = median(td$fare, na.rm=TRUE))

# Model 1: Using the original fare variable

```

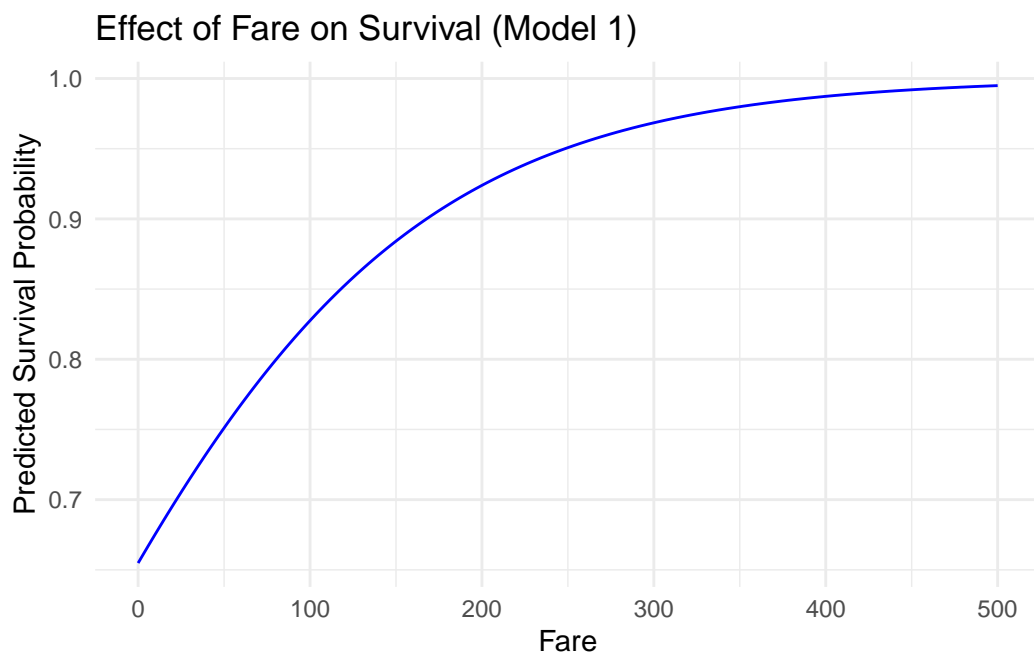
```

fare_values_1 <- seq(0, 500, by = 1) # Range for fare
predicted_1 <- predict(model1, newdata = expand.grid(fare = fare_values_1, female = "Female")

# Model 2: Using the log-transformed fare variable
fare_values_2 <- seq(0, 500, by = 1) # Range for fare
predicted_2 <- predict(model2, newdata = expand.grid(fare = fare_values_2, female = "Female")

# First plot (Model 1)
ggplot(data.frame(fare = fare_values_1, survival_prob = predicted_1), aes(x = fare, y = survival_prob)) +
  geom_line(color = "blue") +
  labs(x = "Fare", y = "Predicted Survival Probability", title = "Effect of Fare on Survival Probability") +
  theme_minimal()

```

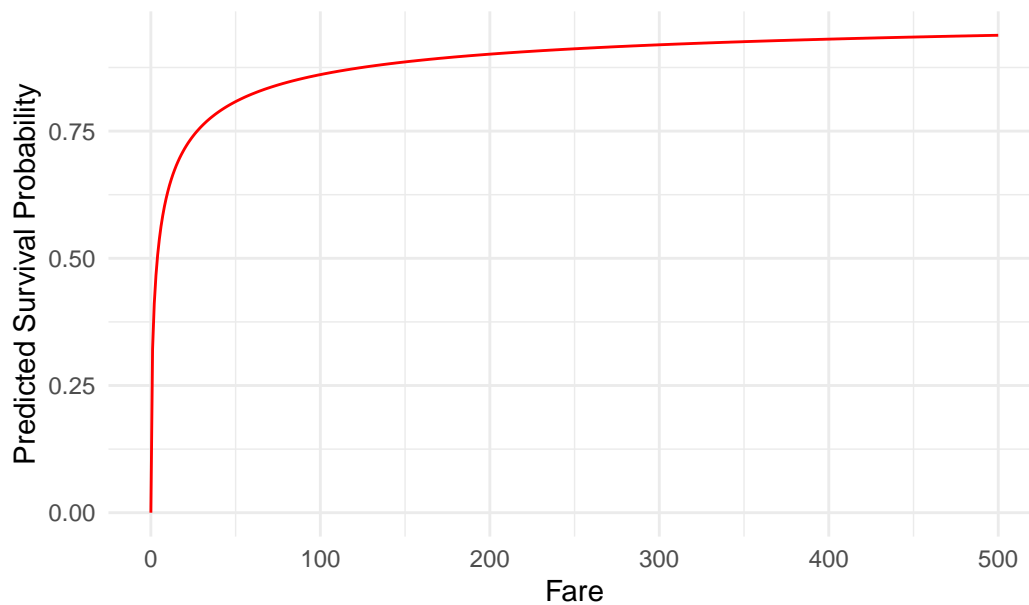


```

# Plot Effect of fare on passenger survival for Model 2 (log-transformed fare)
ggplot(data.frame(fare = fare_values_2, survival_prob = predicted_2), aes(x = fare, y = survival_prob)) +
  geom_line(color = "red") +
  labs(x = "Fare", y = "Predicted Survival Probability", title = "Effect of Fare on Survival Probability") +
  theme_minimal()

```


Effect of Fare on Survival (Model 2)



#3d.

```
# Set seed for reproducibility
set.seed(123)

# Remove rows with missing values for fare
td_clean_non_na <- td_clean[!is.na(td_clean$fare), ]

# Create a random 80-20 split
train_id <- sample(1:nrow(td_clean_non_na), nrow(td_clean_non_na) * 0.8)

# Split the data into training and test sets
train_data <- td_clean_non_na[train_id, ]
test_data <- td_clean_non_na[-train_id, ]

# Fit a logit model without the logged fare variable
model_no_log_fare <- glm(survived ~ fare + female + child,
                        data = train_data,
                        family = "binomial")

# Fit a logit model with the logged fare variable
model_with_log_fare <- glm(survived ~ log(fare) + female + child,
                          data = train_data,
                          family = "binomial")
```

```
# Predict survival probabilities using both models
pred_no_log_fare <- predict(model_no_log_fare, newdata = test_data, type = "response")
pred_with_log_fare <- predict(model_with_log_fare, newdata = test_data, type = "response")

# Generate ROC curve for both models
roc_no_log_fare <- roc(test_data$survived, pred_no_log_fare)
```

Setting levels: control = 0, case = 1

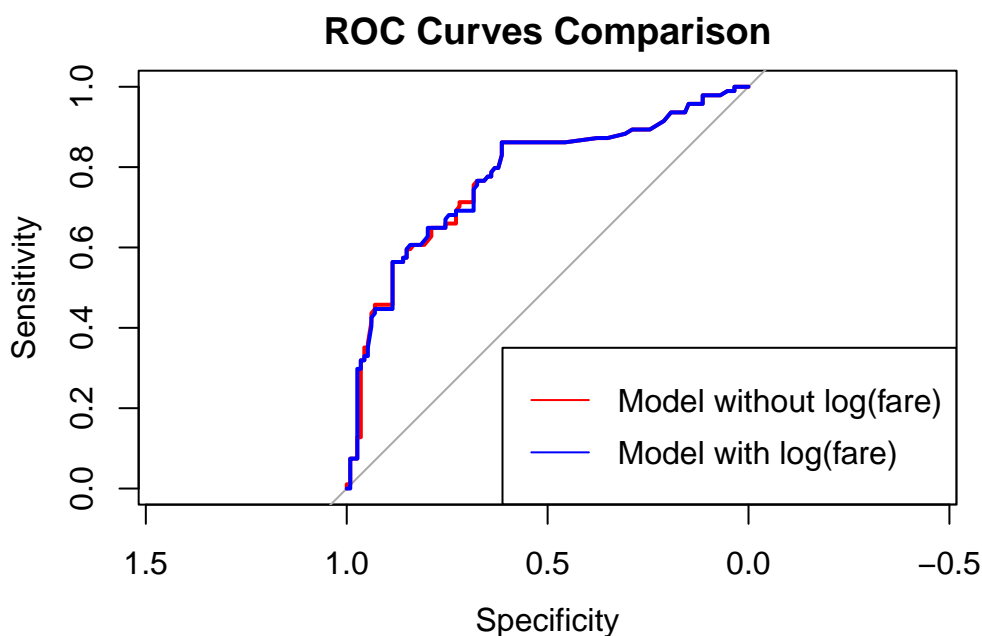
Setting direction: controls < cases

```
roc_with_log_fare <- roc(test_data$survived, pred_with_log_fare)
```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

```
# Plot both ROC curves
plot(roc_no_log_fare, col = "red", main = "ROC Curves Comparison")
plot(roc_with_log_fare, col = "blue", add = TRUE)
legend("bottomright", legend = c("Model without log(fare)", "Model with log(fare)"),
      col = c("red", "blue"), lty = 1)
```



```
# Print AUC values
auc_no_log_fare <- auc(roc_no_log_fare)
auc_with_log_fare <- auc(roc_with_log_fare)

cat("AUC for model without log(fare):", auc_no_log_fare, "\n")
```

AUC for model without log(fare): 0.7750093

```
cat("AUC for model with log(fare):", auc_with_log_fare, "\n")
```

AUC for model with log(fare): 0.7756626

3e.

```
# Fit the first logit model (female and child as predictors)
model_female_child <- glm(survived ~ female + child,
                          data = td_clean_non_na,
                          family = "binomial")

# Fit the second logit model (log(fare), female, and child as predictors)
model_log_fare_female_child <- glm(survived ~ log(fare) + female + child,
                                   data = td_clean_non_na,
                                   family = "binomial")

# Generate a stargazer table comparing both models
stargazer(model_female_child, model_log_fare_female_child,
          type = "text",
          title = "Comparison of Logit Models",
          column.labels = c("Model 1: Female + Child", "Model 2: Log(Fare) + Female + Child"),
          out = "logit_models_comparison.txt")
```

Comparison of Logit Models

```
=====
                        Dependent variable:
-----
                        survived
Model 1: Female + Child Model 2: Log(Fare) + Female + Child
```

	(1)	(2)
log(fare)		0.559*** (0.082)
femaleMale	-2.457*** (0.153)	-2.334*** (0.157)
childChild	0.660*** (0.237)	0.560** (0.234)
Constant	1.031*** (0.121)	-0.748*** (0.280)
Observations	1,037	1,037
Log Likelihood	-542.342	-518.122
Akaike Inf. Crit.	1,090.683	1,044.245
Note:	*p<0.1; **p<0.05; ***p<0.01	

3f.

```
# Fit the first logit model (female and child as predictors)
model_female_child <- glm(survived ~ female + child,
                          data = td_clean_non_na,
                          family = "binomial")

# Fit the second logit model (log(fare), female, and child as predictors)
model_log_fare_female_child <- glm(survived ~ log(fare) + female + child,
                                   data = td_clean_non_na,
                                   family = "binomial")

# Fit the second logit model (log(fare), female, and child as predictors)
model_log_fare_female_child <- glm(survived ~ log(fare) + female + child,
                                   data = td_clean_non_na,
                                   family = "binomial")

# Generate a stargazer table comparing both models
stargazer(model_female_child, model_log_fare_female_child,
```

```

type = "text",
title = "Comparison of Logit Models",
column.labels = c("Model 1: Female + Child", "Model 2: Log(Fare) + Female + Child"),
out = "logit_models_comparison.txt")

```

Comparison of Logit Models

Dependent variable:		
survived		
	Model 1: Female + Child (1)	Model 2: Log(Fare) + Female + Child (2)
log(fare)		0.559*** (0.082)
femaleMale	-2.457*** (0.153)	-2.334*** (0.157)
childChild	0.660*** (0.237)	0.560** (0.234)
Constant	1.031*** (0.121)	-0.748*** (0.280)
Observations	1,037	1,037
Log Likelihood	-542.342	-518.122
Akaike Inf. Crit.	1,090.683	1,044.245
Note:	*p<0.1; **p<0.05; ***p<0.01	

The second model seems to perform better bases on the lower AIC.

3g.

```

# Get predicted probabilities for both models
pred_female_child <- predict(model_female_child, type = "response")

```

```

pred_log_fare_female_child <- predict(model_log_fare_female_child, type = "response")

# Generate ROC curves for both models
roc_female_child <- roc(td_clean_non_na$survived, pred_female_child)

```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

```

roc_log_fare_female_child <- roc(td_clean_non_na$survived, pred_log_fare_female_child)

```

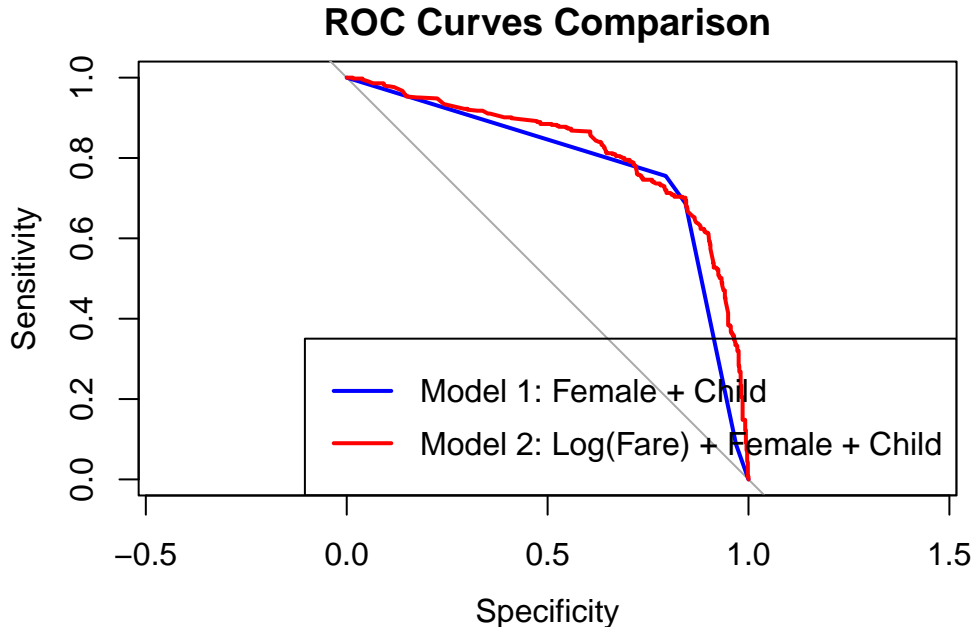
Setting levels: control = 0, case = 1

Setting direction: controls < cases

```

# Plot ROC curves for both models in one graph
plot(roc_female_child, col = "blue", main = "ROC Curves Comparison",
     lwd = 2, xlim = c(0, 1), ylim = c(0, 1))
lines(roc_log_fare_female_child, col = "red", lwd = 2)
legend("bottomright", legend = c("Model 1: Female + Child", "Model 2: Log(Fare) + Female + Child"),
      col = c("blue", "red"), lwd = 2)

```



3h.

```
# Generate ROC curves for both models (if not already done)
roc_female_child <- roc(td_clean_non_na$survived, pred_female_child)
```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

```
roc_log_fare_female_child <- roc(td_clean_non_na$survived, pred_log_fare_female_child)
```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

```
# Calculate AUC for both models
auc_female_child <- auc(roc_female_child)
auc_log_fare_female_child <- auc(roc_log_fare_female_child)

# Report AUC scores
auc_female_child
```

Area under the curve: 0.7824

```
auc_log_fare_female_child
```

Area under the curve: 0.8194

```
# The second model controlling for logfare has a higher
# predictive power than the first model.
```

3i.

```
# Based on the ROC curves and AUC scores, the model
# that includes the logged fare variable performs better.
# This is evident from its higher AUC score, which indicates
# that it is more effective in distinguishing between survival
# and non-survival outcomes compared to the model that only
# includes female and child variables. Therefore, the
# inclusion of the logged fare variable enhances
# the model's predictive ability.
```

4.

```
setwd("/Users/jackjeffrey/Documents/Poli502_Jeffrey/Data")
# Load dataset
putnam <- read.csv("putnam.csv")
# Preview data
head(putnam)
```

	Region	InstPerform	CivicCommunity	NorthSouth	EconModern
1	Ab	7.5	8	South	7.0
2	Ba	7.5	4	South	3.0
3	Cl	1.5	1	South	3.0
4	Cm	2.5	2	South	6.5
5	Em	16.0	18	North	13.0
6	Fr	12.0	17	North	14.5

```
# Explore the dataset
str(putnam) # Get the structure of the data
```

```
'data.frame': 20 obs. of 5 variables:
 $ Region      : chr  "Ab" "Ba" "Cl" "Cm" ...
 $ InstPerform : num  7.5 7.5 1.5 2.5 16 12 10 11 11 9 ...
 $ CivicCommunity: num  8 4 1 2 18 17 13 16 17 15.5 ...
 $ NorthSouth  : chr  "South" "South" "South" "South" ...
 $ EconModern   : num  7 3 3 6.5 13 14.5 12.5 15.5 19 10.5 ...
```

```
summary(putnam) # Summary statistics of the dataset
```


Region	InstPerform	CivicCommunity	NorthSouth
Length:20	Min. : 1.50	Min. : 1.000	Length:20
Class :character	1st Qu.: 6.25	1st Qu.: 3.875	Class :character
Mode :character	Median :10.00	Median :15.000	Mode :character
	Mean : 9.15	Mean :11.350	
	3rd Qu.:11.25	3rd Qu.:16.250	
	Max. :16.00	Max. :18.000	

EconModern

Min. : 2.50

1st Qu.: 6.25

Median :11.75

Mean :10.43

3rd Qu.:14.50

Max. :19.00

```
# Check column names
colnames(putnam)
```

```
[1] "Region"      "InstPerform"  "CivicCommunity" "NorthSouth"
[5] "EconModern"
```

```
# Create a dummy variable for North
putnam$North <- ifelse(putnam$NorthSouth == "North", 1, 0)

# (a) Simple linear regression model
model_a <- lm(InstPerform ~ CivicCommunity, data = putnam)
model_a
```

Call:

```
lm(formula = InstPerform ~ CivicCommunity, data = putnam)
```

Coefficients:

```
(Intercept)  CivicCommunity
      2.7112         0.5673
```

```
# (b) Fit additive model using dummy variable North
model_b <- lm(InstPerform ~ CivicCommunity + North, data = putnam)
model_b
```

Call:

```
lm(formula = InstPerform ~ CivicCommunity + North, data = putnam)
```

Coefficients:

(Intercept)	CivicCommunity	North
2.69850	0.57094	-0.04781

```
# the regression equation would be as follows -
```

```
# InstPerform=2.69850+0.57094×CivicCommunity-0.04781×North
```

```
# Fit the interactive model
```

```
model_c <- lm(InstPerform ~ CivicCommunity * North, data = putnam)
```

```
model_c
```

Call:

```
lm(formula = InstPerform ~ CivicCommunity * North, data = putnam)
```

Coefficients:

(Intercept)	CivicCommunity	North
2.82828	0.54040	-1.19414

CivicCommunity:North

0.09374

```
# Now we need to graph these models to understand the results
```