

# Statistical Inference

## Project Phase-2

### University of Tehran

ECE Department

Fall 2022

Teaching Assistant	Email
Sarmad Zandi	<a href="mailto:sarmadzandi@ut.ac.ir">sarmadzandi@ut.ac.ir</a>
Sheyda Eshaghi	<a href="mailto:sheydaes@gmail.com">sheydaes@gmail.com</a>
Zahra Mohaghegh Rad	<a href="mailto:zah96rad@gmail.com">zah96rad@gmail.com</a>
Kimia Afrazande	<a href="mailto:kimiaa96@gmail.com">kimiaa96@gmail.com</a>



# Table of Contents

<b>Introduction</b>	<b>1</b>
<b>Datasets Description</b>	<b>2</b>
<b>Important Notices</b>	<b>5</b>
<b>Question 1</b>	<b>6</b>
<b>Question 2</b>	<b>7</b>
<b>Question 3</b>	<b>8</b>
<b>Question 4</b>	<b>9</b>
<b>Question 5</b>	<b>11</b>
<b>Question 6</b>	<b>12</b>
<b>Question 7</b>	<b>13</b>
<b>References</b>	<b>14</b>



## Introduction

In this project, we intend to study and analyze a series of real datasets with what you learned in this course. The first step to begin analyzing a dataset is to get familiar with it. In the first step, this acquaintance can be made by observing the dataset features and distribution of the values and visualizing the data to make initial guesses about it. In the next step, by performing statistical tests, we make sure our guesses are correct and make our claims with certainty.

To answer each question, you have to fully explain the meaning of your analysis and interpret the generated plot and what you observe, even when it is not explicitly stated in the question. The more reasonable your analysis is, the more positive effect it has on the acquired grade of the corresponding question.

Note that there is no one way to solve each question correctly. Furthermore, whenever you need to do hypothesis testing, you must check all of the pre-requisite conditions (such as sample size, skewness, etc.). Finally, the validity of your results should be discussed.

## Datasets Description

One of the following datasets is assigned to you. for more information about your dataset, please refer to the mentioned references.

Dataset Name	Description
Breast Cancer	This dataset includes information about breast cancer features as well as the diagnosis of the type of cancer. (This dataset contains two features of race and marital status. For more information, go to <a href="#">[3]</a> )
Heart Disease	This dataset includes information about some heart disease features as well as the "target" field, which refers to the presence of heart disease in the patient. <a href="#">[4]</a>
Non-Voters	This directory contains the data behind the story Why Many Americans Don't Vote. The poll was conducted among a sample of U.S. citizens that oversampled young, Black and Hispanic respondents, with 8,327 respondents, and was weighted according to general population benchmarks for U.S. citizens. A voter company matched what respondents said and what they actually did. The data included here is the final sample we used: 5,239 respondents who matched to the voter file and whose verified vote history we have, and 597 respondents who did not match to the voter file and described themselves as voting "rarely" or "never," all of whom have been eligible for at least 4 elections. <a href="#">[5]</a>
Nutrition Studies	This directory contains data and code behind the story You Can't Trust What You Read About Nutrition. Many diet and nutrition studies include multiple variables with vast amounts of data, making it easy to p-hack your way to false results. We learned this firsthand when we invited readers to take a survey about their eating habits known as the food frequency questionnaire and answer a few other questions about themselves, so we ended up with 54 complete responses. <a href="#">[6]</a>

Dataset Name	Description
Airbnb Open Data	Airbnb, Inc is an American company that operates an online marketplace for lodging, primarily homestays for vacation rentals and tourism activities. Based in San Francisco, California, the platform is accessible via website and mobile app. Airbnb does not own any of the listed properties; instead, it profits by receiving commission from each booking. The company was founded in 2008. Airbnb is a shortened version of its original name, AirBedandBreakfast.com. Since 2008, guests and hosts have used Airbnb to travel in a more unique, personalized way. As part of the Airbnb Inside initiative, this dataset describes the listing activity of homestays in New York City. <a href="#">[7]</a>
Car Insurance Claim Prediction	This is the training dataset that contains all the independent and target features. It contains information on policyholders having the attributes like policy tenure, age of the car, age of the car owner, the population density of the city, make and model of the car, power, engine type, etc., and the target variable indicating whether the policyholder files a claim in the next 6 months or not. <a href="#">[8]</a>
Housing Price in Beijing	This dataset describes features like Longitude, Latitude, DOM (days on market), Price, SquareMeter, Living Room, Kitchen, Bathroom, Building Type, Construction time, etc., for houses in Beijing. These features can help predict the price of a house with specific characteristics. <a href="#">[9]</a>
Maximum Credit	This dataset is about the Maximum Credit that is available for different people with different backgrounds and characteristics, like their current loan amount, their annual income, how many years they have used credit, whether they have been bankrupt or not, etc. <a href="#">[10]</a>



## Important Notices

- Use the R language in answering questions. Submit your codes in a separate file next to your report. Reports without R codes are pointless.
- In some datasets, you need to clean the data and convert the format and data type to more appropriate formats. So do this before answering the questions and explain the steps at the beginning of your report.
- If you need more categorical variables, you can add a new one to the dataset using some of your numerical variables. In this case, you need to describe how you created the categorical variable from the numerical variable.
- In most questions, you should use the ggplot2 library to visualize and produce the desired charts. You can find more about this elegant visualization package in references [\[1\]](#), [\[2\]](#).
- For each question, you need to fully explain your answer. An important part of the score will be attributed to your description. Drawing charts and performing calculations without sufficient explanations will result in losing the score. These descriptions show how much you understand the dataset. If you see interesting things in the diagrams, don't forget to mention them.
- When performing statistical tests, be sure to check the requirements for that test and write it down in your answer.



## Question 1

Consider two categorical variables in your dataset for which at least one of them has more than two levels. Using these, follow these steps:

- A. Derive a 95% confidence interval for the difference of these two variables and interpret it.
- B. By hypothesis testing, determine if the two variables are independent or not.



## Question 2

Choose a binary categorical variable and randomly select a small sample of your data (small sample size, e.g.,  $n \leq 15$ ). Then, perform a hypothesis test for the variable's success rate by means of the Simulation method.



### Question 3


Answer the following questions: (Note: To answer the following questions, first implement them by yourself in R and then use R functions to solve them.)

- A. Choose a categorical variable that has more than two levels, calculate its probability distribution. Then choose two samples of size 100 from your dataset. One of the samples should be randomly selected and the other should be biased on purpose. Compare each sample with the real distribution using  $\chi^2$  (goodness of fit) and interpret your results. (be sure to check the requirements for the test and write them down in your answer.)
- B. Pick up another categorical variable and compare it to the one you chose in part (a). Using the  $\chi^2$  test, check if the two variables are independent or not.

## Question 4

From your dataset choose a numerical variable that predicts its future value is meaningful within the context of your dataset. next, choose two explanatory variables which you believe are the best predictors for your response variable:

- A. Without building a model yet, which explanatory variable do you guess is the more significant predictor and why? (use your knowledge from phase 1)
- B. for each explanatory variable:
  - a. Check the Linearity, Nearly Normal Residuals, and Constant Variability conditions in R.
  - b. Compute the least squares regression.
  - c. Write the predictive equation for the response variable and interpret its parameters.
  - d. Draw a scatter plot of the relation between these two variables overlaid with this least-squares fit as a dashed line.
- C. By using the previous part results, try to explain which explanatory variable is the more significant predictor.
- D. Now, Compare your models, once using adjusted  $R^2$  and another time by ANOVA table. Explain results.
- E. According to the results that you found in the previous parts, list the features of a good predictor.

- 
- F. Choose a random sample of 100 data points from the dataset.
- a. By 90 percent of data, Build two Linear Regression models and design hypothesis tests to see if these explanatory variables are a significant predictor of the response variable or not.
  - b. Calculate the 95% confidence interval for the slope of the relationship between response variable and explanatory variables. Interpret these CIs.
  - c. Use your models to predict the values of the response variable for the remaining percent of samples.
  - d. Compare the predicted values with actuals. Report success rate.

## Question 5

Consider the response variable you selected in the previous question. You can use as many explanatory variables as you deem necessary:

- A. Plot a correlogram for explanatory variables and discuss the correlation between them. Could you find which explanatory variable plays a more significant role in prediction
- B. Develop a multiple linear regression model for the response variable using explanatory variables you found in part A.
- C. How well do you think your model fits the data?
- D. Develop the “best” possible multiple linear regression model for the response variable using different approaches and metrics.
- E. Use 5-fold cross-validation and compare the model’s RMSE (part B and C). How do you interpret these values?
- F. Check diagnostics for your model in part C (Three conditions: 1. Linearity, 2. Nearly normal residuals, and 3. Constant variability) and explain if this is a reliable model or not.
- G. What percent of the variation in the response variable is explained by the model (part B and C)?

## Question 6

Choose a binary categorical variable from your dataset as a response variable and choose several categorical and numerical variables which you think can best explain the response variable.

- A. Construct a logistic regression model and interpret the intercept and the slopes in terms of log odds and log odds ratio.
- B. Draw the ROC curve for the model. What does this diagram signify? Discuss the goodness of the model based on the AUC.
- C. Choose a categorical variable in your model among the explanatory variables and plot the odds ratio curve for that variable. Interpret the plot.
- D. Calculate a 95% confidence interval for the odds ratio

## Question 7

Please answer the following questions in respect to the model in the previous question.

- A. Which explanatory variable in the model plays the most significant role in the prediction? Why?
- B. Select another categorical variable except the response variable from the model, draw the OR (odd ratio) curve for this categorical variable and interpret the plot.
- C. Select explanatory variables with the most meaningful roles in the model prediction, and construct the new Logistic Regression model, and then interpret the result.
- D. Draw the utility curve for the model you've created in part C (define the utility of different outcomes yourself). What is the best threshold for this model?



## References

- [1] Intro to Data Visualization with R & ggplot2 [\(Link\)](#)
- [2] Data visualization with R and ggplot2: the R Graph Gallery [\(Link\)](#)
- [3] Breast Cancer DataSet [\(Link\)](#)
- [4] Heart Disease DataSet [\(Link\)](#)
- [5] Non-Voters DataSet [\(Link\)](#)
- [6] Nutrition Studies DataSet [\(Link\)](#)
- [7] Airbnb Open Data DataSet [\(Link\)](#)
- [8] Car Insurance Claim Prediction DataSet [\(Link\)](#)
- [9] Housing price in Beijing DataSet [\(Link\)](#)
- [10] Maximum Credit DataSet [\(Link\)](#)