

# R Applications

Tan Sein Jone

University of British Columbia

August 8, 2024

## Table of contents

1. Piping
2. Hypothesis Testing
3. Regression Analysis
4. Research Design

# Table of Contents

## 1. Piping

## 2. Hypothesis Testing

## 3. Regression Analysis

## 4. Research Design

# Piping

- Piping is a way to simplify your code by passing the output of one function directly to the input of another function.
- Piping is done using the pipe operator, which is represented by the `%>%` symbol.
- Piping makes your code more readable and easier to understand by breaking it down into smaller, more manageable steps.

# Piping

- Piping is especially useful when you have a series of operations that need to be performed on the same data.
- Piping allows you to chain together multiple functions in a single line of code.
- Piping is a powerful tool that can help you write more efficient and concise code.

# Piping

- Piping is done using the pipe operator, which is represented by the `%>%` symbol.
- The pipe operator takes the output of the function on the left and passes it as the first argument to the function on the right.
- The pipe operator can be used to chain together multiple functions in a single line of code.

## Example of Piping vs. Non-Piping

- Without Piping:

```
# Load necessary library  
library(dplyr)
```

```
# Generate data  
set.seed(123)  
data <- data.frame(  
  x = rnorm(100, mean = 5, sd = 2),  
  y = rnorm(100, mean = 6, sd = 2)  
)
```

```
# Non-piping approach  
data <- mutate(data, z = x + y)  
data <- filter(data, z > 10)  
summary(data)
```

## Example of Piping vs. Non-Piping

- Without Piping:

```
# Load necessary library  
library(dplyr)
```

```
# Generate data  
set.seed(123)  
data <- data.frame(  
  x = rnorm(100, mean = 5, sd = 2),  
  y = rnorm(100, mean = 6, sd = 2)  
)
```

```
# Piping approach  
data %>%  
  mutate(z = x + y) %>%  
  filter(z > 10) %>%  
  summary()
```



# Breakdown of Piping

- **Load Necessary Library:**

- `library(dplyr)`: Loads the `dplyr` package, which provides functions for data manipulation.

- **Generate Data:**

- `set.seed(123)`: Sets the seed for random number generation to ensure reproducibility of results.
- `data <- data.frame(x = rnorm(100, mean = 5, sd = 2), y = rnorm(100, mean = 6, sd = 2))`: Generates 100 random numbers from a normal distribution with specified means and standard deviations, and assigns them to variables `x` and `y`.

- **Non-Piping Approach:**

- `data <- mutate(data, z = x + y)`: Adds a new variable `z` to the data frame by summing variables `x` and `y`.
- `data <- filter(data, z > 10)`: Filters the data frame to include only rows where variable `z` is greater than 10.
- `summary(data)`: Displays a summary of the data frame.

# Tips for Piping

- Use piping to chain together multiple functions in a single line of code.
- Use piping to break down complex operations into smaller, more manageable steps.
- Use piping to make your code more readable and easier to understand.
- Use piping to improve the efficiency and conciseness of your code.



# Hypothesis Testing

- Hypothesis testing is a statistical method that is used in making statistical decisions using experimental data.
- A hypothesis test evaluates two mutually exclusive statements about a population to determine which statement is best supported by the sample data.
- These two statements are called the null hypothesis and the alternative hypothesis.

# Hypothesis Testing

- The null hypothesis is the statement being tested. Usually, the null hypothesis states that there is no effect or no difference.
- The alternative hypothesis is the statement that is accepted if the sample data provide enough evidence that the null hypothesis is false.
- The alternative hypothesis states that there is an effect or a difference.

# Hypothesis Testing

- The hypothesis test is conducted by comparing the value of the test statistic to a critical value.
- The critical value is a value that determines whether the null hypothesis can be rejected.
- If the test statistic is more extreme than the critical value, then the null hypothesis is rejected.

# Types of Hypothesis Tests

- One-Sample t-Test: Used to test whether the mean of a single sample is significantly different from a specified value.
- Two-Sample t-Test: Used to test whether the means of two independent samples are significantly different from each other.
- Paired t-Test: Used to test whether the means of two paired samples are significantly different from each other.

## Types of Hypothesis Tests

- One-Sample z-Test: Used to test whether the mean of a single sample is significantly different from a specified value when the sample size is large.
- Two-Sample z-Test: Used to test whether the means of two independent samples are significantly different from each other when the sample sizes are large.
- Chi-Square Test: Used to test whether the observed frequencies in a contingency table are significantly different from the expected frequencies.



## Types of Hypothesis Tests

- One-Way ANOVA: Used to test whether the means of three or more independent samples are significantly different from each other.
- Two-Way ANOVA: Used to test whether the means of two or more independent samples are significantly different from each other, taking into account two independent variables.
- Goodness of Fit Test: Used to test whether the observed frequencies in a sample are consistent with the expected frequencies.

# Types of Hypothesis Tests

- Test for Equal Variances: Used to test whether the variances of two or more samples are equal.
- Test for Normality: Used to test whether the data in a sample comes from a normal distribution.
- Test for Independence: Used to test whether the observations in a sample are independent of each other.

# Steps in Hypothesis Testing

- Step 1: State the null hypothesis and the alternative hypothesis.
- Step 2: Choose the significance level.
- Step 3: Collect the sample data and calculate the test statistic.
- Step 4: Determine the critical value or the p-value.
- Step 5: Make a decision to reject or fail to reject the null hypothesis.
- Step 6: Interpret the results of the hypothesis test.

# Null and Alternative Hypotheses

- The null hypothesis is denoted by  $H_0$  and states that there is no effect or no difference.
- The alternative hypothesis is denoted by  $H_1$  and states that there is an effect or a difference.
- The null hypothesis is the default assumption that is tested against the alternative hypothesis.
- The null hypothesis is rejected if the sample data provide enough evidence that the null hypothesis is false.

# Assumptions of Hypothesis Testing

- The sample data are independent and identically distributed.
- The sample data are drawn from a population that is normally distributed.
- The sample data are drawn from a population that has a constant variance.
- The sample data are drawn from a population that is randomly selected.
- The sample data are drawn from a population that is representative of the population of interest.

## Significance Level

- The significance level is the probability of rejecting the null hypothesis when it is true.
- The significance level is denoted by  $\alpha$  and is usually set to 0.05.
- The significance level is the threshold for determining whether the null hypothesis should be rejected.
- If the p-value is less than the significance level, then the null hypothesis is rejected.

# Test Statistic

- The test statistic is a numerical value that is used to determine whether the null hypothesis should be rejected.
- The test statistic is calculated from the sample data and is compared to a critical value or a p-value.
- The test statistic measures the strength of the evidence against the null hypothesis.
- The test statistic is used to make a decision to reject or fail to reject the null hypothesis.

# Critical Value

- The critical value is a value that determines whether the null hypothesis can be rejected.
- The critical value is determined by the significance level and the degrees of freedom.
- The critical value is compared to the test statistic to determine whether the null hypothesis should be rejected.
- If the test statistic is more extreme than the critical value, then the null hypothesis is rejected.



# P-Value

- The p-value is the probability of observing a test statistic as extreme as the one computed from the sample data, assuming that the null hypothesis is true.
- The p-value is a measure of the strength of the evidence against the null hypothesis.
- The p-value is compared to the significance level to determine whether the null hypothesis should be rejected.
- If the p-value is less than the significance level, then the null hypothesis is rejected.

## Type I and Type II Errors

- Type I Error: Occurs when the null hypothesis is rejected when it is true.
- Type II Error: Occurs when the null hypothesis is not rejected when it is false.
- The probability of a Type I Error is denoted by  $\alpha$  and is equal to the significance level.
- The probability of a Type II Error is denoted by  $\beta$  and is equal to 1 minus the power of the test.

## Power of the Test

- The power of the test is the probability of rejecting the null hypothesis when it is false.
- The power of the test is equal to 1 minus the probability of a Type II Error.
- The power of the test is a measure of the ability of the test to detect an effect or a difference when it exists.
- The power of the test is affected by the sample size, the effect size, and the significance level.

## Confidence Interval

- A confidence interval is a range of values that is likely to contain the true value of the population parameter.
- A confidence interval is calculated from the sample data and is used to estimate the population parameter.
- The confidence interval is calculated from the sample mean and the standard error of the mean.
- The confidence interval is used to make inferences about the population parameter with a specified level of confidence.

## Confidence Level

- The confidence level is the probability that the confidence interval contains the true value of the population parameter.
- The confidence level is denoted by  $(1 - \alpha) \times 100\%$  and is usually set to 95%.
- The confidence level is the proportion of confidence intervals that contain the true value of the population parameter.
- The confidence level is used to make inferences about the population parameter with a specified level of confidence.

# Hypothesis Testing in R

- In R, the `t.test()` function is used to perform hypothesis tests.
- The `t.test()` function takes in the sample data and the null hypothesis as arguments.
- The function returns the test statistic, the p-value, and the confidence interval.

# Hypothesis Testing in R

- The `t.test()` function can be used to perform one-sample t-tests, two-sample t-tests, and paired t-tests.
- The `t.test()` function can also be used to perform one-sample z-tests and two-sample z-tests.
- The `t.test()` function can be used to perform hypothesis tests for means, proportions, and variances.

# Hypothesis Testing in R

```
# Generate some data
set.seed(123)
x <- rnorm(100, mean = 5, sd = 2)
y <- rnorm(100, mean = 6, sd = 2)

# One-sample t-test
t.test(x, mu = 0)

# Two-sample t-test
t.test(x, y)

# Paired t-test
t.test(x, y, paired = TRUE)
```



# Breakdown of Hypothesis Testing in R

- **Set Seed:**

- `set.seed(123)`: Sets the seed for random number generation to ensure reproducibility of results.

- **Generate Data:**

- `x <- rnorm(100, mean = 5, sd = 2)`: Generates 100 random numbers from a normal distribution with a mean of 5 and a standard deviation of 2, and assigns them to vector `x`.
- `y <- rnorm(100, mean = 6, sd = 2)`: Generates 100 random numbers from a normal distribution with a mean of 6 and a standard deviation of 2, and assigns them to vector `y`.

- **One-Sample t-Test:**

- `t.test(x, mu = 0)`: Performs a one-sample t-test to check if the mean of vector `x` is significantly different from 0.

- **Two-Sample t-Test:**

- `t.test(x, y)`: Performs a two-sample t-test to check if the means of vectors `x` and `y` are significantly different from each other.



# Regression Analysis

- Regression analysis is a statistical method that is used to model the relationship between a dependent variable and one or more independent variables.
- The goal of regression analysis is to estimate the parameters of the regression model that best fit the data.
- The regression model is a mathematical equation that describes the relationship between the dependent variable and the independent variables.

# Types of Regression Models

- Linear Regression: Used to model the relationship between a continuous dependent variable and one or more continuous independent variables.
- Logistic Regression: Used to model the relationship between a binary dependent variable and one or more continuous independent variables.
- Polynomial Regression: Used to model the relationship between a continuous dependent variable and one or more continuous independent variables using polynomial functions.

## Types of Regression Models

- Ridge Regression: Used to model the relationship between a continuous dependent variable and one or more continuous independent variables with multicollinearity.
- Lasso Regression: Used to model the relationship between a continuous dependent variable and one or more continuous independent variables with feature selection.
- Elastic Net Regression: Used to model the relationship between a continuous dependent variable and one or more continuous independent variables with both ridge and lasso regularization.

# Linear Regression

- Linear regression is a regression model that assumes a linear relationship between the dependent variable and the independent variables.
- The linear regression model is a mathematical equation that describes the relationship between the dependent variable and the independent variables.
- The linear regression model is estimated using the method of least squares, which minimizes the sum of the squared differences between the observed values and the predicted values.

# Linear Regression

- The linear regression model is specified by the equation  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$ , where  $y$  is the dependent variable,  $x_1, x_2, \dots, x_n$  are the independent variables,  $\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are the coefficients of the independent variables, and  $\epsilon$  is the error term.
- The coefficients of the independent variables are estimated using the method of least squares, which minimizes the sum of the squared differences between the observed values and the predicted values.

# Ordinary Least Squares

- Ordinary least squares (OLS) is a method that is used to estimate the parameters of a linear regression model.
- OLS minimizes the sum of the squared differences between the observed values and the predicted values.
- OLS is used to estimate the coefficients of the independent variables that best fit the data.
- OLS is used to estimate the coefficients of the independent variables that minimize the sum of the squared differences between the observed values and the predicted values.



## Formula for OLS

- The formula for OLS is  $\hat{\beta} = (X^T X)^{-1} X^T y$ , where  $\hat{\beta}$  is the estimated coefficients of the independent variables,  $X$  is the matrix of the independent variables, and  $y$  is the vector of the dependent variable.
- The formula for OLS is derived by minimizing the sum of the squared differences between the observed values and the predicted values.
- The formula for OLS is used to estimate the coefficients of the independent variables that best fit the data.

# Polynomial Regression

- Polynomial regression is a regression model that is used when the relationship between the dependent variable and the independent variables is not linear.
- The polynomial regression model is a mathematical equation that describes the relationship between the dependent variable and the independent variables using polynomial functions.
- The polynomial regression model is estimated using the method of least squares, which minimizes the sum of the squared differences between the observed values and the predicted values.

# Polynomial Regression

- The polynomial regression model is specified by the equation  $y = \beta_0 + \beta_1x + \beta_2x^2 + \dots + \beta_nx^n + \epsilon$ , where  $y$  is the dependent variable,  $x$  is the independent variable,  $\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are the coefficients of the independent variables, and  $\epsilon$  is the error term.
- The coefficients of the independent variables are estimated using the method of least squares, which minimizes the sum of the squared differences between the observed values and the predicted values.
- Having a polynomial of degree  $n - 1$  would mean that the model will perfectly fit the data. So, be careful of overfitting.

# Regression Analysis

- The goodness of fit of the regression model is measured using the coefficient of determination, which is the proportion of the variance in the dependent variable that is explained by the independent variables.
- The coefficient of determination ranges from 0 to 1, with higher values indicating a better fit. (This is your  $R^2$  value)
  - A higher  $R^2$  does not mean your regression is better. It just means that your model explains more of the variance in the dependent variable.
  - It could also be an indication of overfitting.
- The significance of the regression model is tested using the F-test, which tests whether the regression model is a better fit than a model with no independent variables.

# Regression Analysis in R

- In R, the `lm()` function is used to fit linear regression models.
- The `lm()` function takes in the formula for the regression model and the data as arguments.
- The formula specifies the dependent variable and the independent variables in the regression model.

# Regression Analysis in R

- The `summary()` function is used to display the results of the regression analysis.
- The `summary()` function displays the estimated coefficients, the standard errors, the t-values, and the p-values of the regression model.
- The `summary()` function also displays the coefficient of determination and the results of the F-test.

# Regression Analysis in R

```
# Generate some data
set.seed(123)
data <- data.frame(
  y = rnorm(100, mean = 5, sd = 2),
  x1 = rnorm(100, mean = 6, sd = 2),
  x2 = rnorm(100, mean = 7, sd = 2)
)

# Fit linear regression model
model <- lm(y ~ x1 + x2, data = data)

# Display results
summary(model)
```

# Breakdown of Regression Analysis in R

- **Set Seed:**

- `set.seed(123)`: Sets the seed for random number generation to ensure reproducibility of results.

- **Generate Data:**

- `data <- data.frame(y = rnorm(100, mean = 5, sd = 2), x1 = rnorm(100, mean = 6, sd = 2), x2 = rnorm(100, mean = 7, sd = 2))`: Generates 100 random numbers from a normal distribution with specified means and standard deviations, and assigns them to variables `y`, `x1`, and `x2`.

- **Fit Linear Regression Model:**

- `model <- lm(y ~ x1 + x2, data = data)`: Fits a linear regression model to predict variable `y` using variables `x1` and `x2` as predictors.

- **Display Results:**

- `summary(model)`: Displays the results of the linear regression analysis, including the estimated coefficients, standard errors, t-values, p-values, coefficient of determination, and F-test results.



# Best Practices for Displaying Results

- Use tables to display the results of the regression analysis.
- Include the estimated coefficients, standard errors, t-values, and p-values in the table.
- Highlight the statistically significant coefficients in the table.
- Include the coefficient of determination and the results of the F-test in the table.
- Use visualizations, such as scatter plots and residual plots, to display the relationship between the dependent variable and the independent variables.

## Best Practices for Interpreting Results

- Interpret the estimated coefficients of the independent variables in the regression model.
- Interpret the coefficient of determination, which is the proportion of the variance in the dependent variable that is explained by the independent variables.
- Interpret the results of the F-test, which tests whether the regression model is a better fit than a model with no independent variables.
- Interpret the p-values of the estimated coefficients, which indicate the statistical significance of the coefficients.



## IV in a Nutshell

$$y = \beta_0 + \beta_1 x + \epsilon \quad (1)$$

- $x_1$  is not randomly assigned in general ( $\text{cov}(x_1, \epsilon) \neq 0$ ), but ...
- Some small part of the variation in  $x_1$  is idiosyncratic
- If you could identify and isolate that idiosyncratic variance, you could use it to estimate the causal effect of  $x_1$  on  $y$
- This is the basic idea behind instrumental variables

## Two Stage Least Squares

$$x_1 = \pi_1 z + \pi_2 X_2 + \nu \quad (2)$$

$$y = \beta_1 \hat{x}_1 + \beta_2 X_2 + \epsilon \quad (3)$$

- Regress  $x_1$  on the instrument  $z$  to get  $\hat{x}_1$ , be sure to add any additional controls in the regression.
- Use the predicted values of  $x_1$  in the second stage regression along with any other controls.

# What Makes a Good Instrument

- Exogeneity:  $\text{cov}(z, \epsilon) = 0$ 
  - The instrument is uncorrelated with the error term
  - You can test this using observables
  - But this is necessary and not sufficient, you would have to argue this. There may be unobservables that are correlated with the instrument and the error term.
- Excludability: The instrument is not correlated with the dependent variable ( The only way the instrument affects the dependent variable is through the independent variable. )
- Relevance: The instrument is correlated with the independent variable ( This is testable )

# Local Average Treatment Effect

- Homogenous Treatment Effects
  - This means that the treatment effect is the same for everyone
- If we assume that the treatment effect is homogenous for everyone, then the instrumental variable estimates the local average treatment effect (LATE)
- This means the average treatment effect for the subpopulation of individuals who are affected by the instrument

# Subgroups

- Always takers: People who would take the treatment regardless of the instrument
- Never takers: People who would never take the treatment regardless of the instrument
- Compliers: People who would take the treatment if the instrument is high or low enough



# LATE Assumptions

- Independence: The instrument is as good as randomly assigned
- Exclusion restriction: The instrument only affects the dependent variable through the independent variable
- Monotonicity: There are no defiers (non compliers)
  - Either  $\pi_{1i} \geq 0 \forall i$  or  $\pi_{1i} \leq 0 \forall i$

# Instrumental Variables in R

- The `ivreg()` function estimates the parameters of the regression model using the method of instrumental variables.
- The `ivreg()` function returns the estimated coefficients, the standard errors, the t-values, and the p-values of the regression model.
- The `ivreg()` function is used to estimate the causal effect of the independent variable on the dependent variable by removing the correlation between the independent variable and the error term.

# Instrumental Variables in R

```
library(AER)
# Generate some data
set.seed(123)
n <- 100
x <- rnorm(n)
z <- rnorm(n)
y <- 1 + 2 * x + 3 * z + rnorm(n)

# Fit instrumental variables regression model
model <- ivreg(y ~ x | z)
summary(model)
```

# Breakdown of Instrumental Variables in R

- **Load Package:**

- `library(AER)`: Loads the AER package, which contains the `ivreg()` function for estimating instrumental variables regression models.

- **Generate Data:**

- `x <- rnorm(n)`: Generates `n` random numbers from a normal distribution and assigns them to vector `x`.
- `z <- rnorm(n)`: Generates `n` random numbers from a normal distribution and assigns them to vector `z`.
- `y <- 1 + 2 * x + 3 * z + rnorm(n)`: Generates `n` random numbers from a normal distribution and assigns them to vector `y`.

- **Fit Instrumental Variables Regression Model:**

- `model <- ivreg(y ~ x | z)`: Fits an instrumental variables regression model to predict variable `y` using variable `x` as a predictor and variable `z` as an instrumental variable.

- **Display Results:**

- `summary(model)`: Displays the results of the instrumental variables regression analysis, including the estimated coefficients, standard errors, t-values, and p-values.