

# CANCER DATA ANALYSIS REPORT



Mukesh Patel School of Technology, Management and Engineering  
Department of Computer Engineering  
B.Tech Computer III semester B\C division  
2023-2024  
Subject: Data Structures  
Project Report

Vidhi Damani - B068

Samridhi Raj Sinha - B074

Asmi Parikh - B082

1st November 2023

# Table of Contents

<b>Executive Summary</b>	<b>i</b>
<b>References</b>	<b>ii</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background .....	
1.2 Objectives of Project .....	
1.3 Goals .....	
1.4 Report Overview .....	
<b>2 Methodology</b>	<b>2</b>
2.1 Data Collection .....	
2.2 Data Preprocessing .....	
2.3 Exploratory Data Analysis .....	
2.4 Statistical Analysis .....	
<b>3 Project Description</b>	<b>3</b>
3.1 Scope .....	
3.1 Resources .....	
<b>4 Results and Findings</b>	<b>4</b>
4.1 About Our Dataset. ....	
4.2 Analysis using Python .....	
4.3 Statistical Analysis using R .....	
<b>5 Discussions</b>	<b>5</b>
5.1 Discussions .....	
<b>6 Conclusion</b>	<b>6</b>
6.1 Key Findings .....	
6.2 Project Success and Lessons Learned .....	
<b>7 Future Work</b>	<b>7</b>
7.1 Research Areas .....	

# Executive Summary

The "Cancer Data Analysis project" is a research project aimed at exploring trends related to different types of cancer across various ages and states across the USA.

Through data preprocessing and exploratory data analysis [using languages like R and python](#), the analysis aims to discover statistical relationships between different parameters of our dataset obtained from kaggle.

The project involves an extensive data analysis process, starting with exploratory data analysis (EDA) that employs descriptive statistics and visualizations to uncover insights within the dataset.

Multivariate analysis techniques, including linear regression and correlation analysis, are utilized to identify relationships among variables and their impact on cancer rates.

We utilize data visualization tools like graphs and charts to effectively communicate our findings.

Throughout this analysis, we draw inferences that provide valuable insights into the complex relationships between cancer rates, demographics, and other important factors, aiming to inform healthcare strategies and public health initiatives.

# Chapter 1

## Introduction

### 1.1 Background

Cancer is a critical global health concern, with its prevalence and impact on society continually evolving. In this context, the data analysis project was undertaken to explore a comprehensive dataset focusing on cancer rates, demographics, and cancer types across different regions. This dataset offers invaluable insights for understanding the relationships between various demographic factors and cancer rates.

This analysis aims to contribute to the collective knowledge in the field of oncology and public health, with the hope of guiding future initiatives for cancer prevention and improving patient outcomes. By examining this comprehensive dataset, we can draw evidence-based conclusions that can inform policy decisions and medical practices, ultimately leading to better health and quality of life for individuals at risk of or affected by cancer.

### 1.2 Objectives and Goals

The primary objectives of this project are as follows:

To investigate the impact of demographic factors, including age, gender, and race, on cancer rates.

To assess the prevalence and distribution of different types of cancer within the dataset.

To identify significant patterns, correlations, and disparities in cancer rates.

To conduct statistical analysis using techniques like linear regression

### 1.3 Report Overview

This report presents the findings of the data analysis project, structured in a manner that begins with data import and cleaning and proceeds to extensive exploratory data analysis (EDA). The analysis covers a wide array of cancer-related aspects, including demographic analysis, cancer type analysis, hypothesis testing, correlation and regression analysis.

The report will provide an in-depth understanding of how cancer rates are influenced by demographics and explore significant trends within the dataset. By the end of this report, readers will have comprehensive insights into the relationships between cancer rates, age groups, gender, and racial backgrounds, as well as the prevalence of specific cancer types.

The subsequent sections of this report will delve into each aspect of the analysis, presenting detailed findings and observations related to the project's objectives

## Chapter 2

# Methodology

### 2.1 Data Collection

The dataset used for this analysis was sourced from Kaggle, a reputable platform for sharing datasets and data-related projects. This dataset contains a wealth of information, including cancer rates, demographics, and cancer types across different regions or states.

### 2.2 Data Preprocessing

Data preprocessing was a crucial phase in this project to ensure the quality and integrity of the dataset. This involved the following steps:

**Data Cleaning:** We systematically addressed missing values, duplicates, and outliers within the dataset. This step was essential to establish a reliable foundation for subsequent analysis.

**Data Transformation:** To facilitate meaningful analysis, data was transformed, including the structuring of variables and addressing inconsistencies in variable names and values.

### 2.3 Exploratory Data Analysis (EDA)

Exploratory Data Analysis was conducted using Python, and the following techniques were applied:

**Descriptive Statistics:** We calculated summary statistics to gain insights into the dataset's central tendencies and distributions. This step helped us understand the variability and characteristics of different variables.

**Data Visualization:** We employed Python's libraries, such as Pandas, Matplotlib, and Seaborn, to create a variety of data visualizations. These visualizations included histograms, scatter plots, and box plots, which were used to explore the distribution of cancer rates, demographic patterns, and relationships among variables.

## 2.4 Statistical Analysis

For the statistical aspects of this project, we used both Python and the R programming language:

1. Hypothesis Testing: We formulated and conducted hypothesis tests to evaluate the significance of various factors on cancer rates. This helped us determine whether observed patterns were statistically significant.
2. Correlation Analysis: Correlation analysis was performed to identify significant relationships between different variables, such as the relationships between cancer rates and demographic factors.
3. Linear Regression Analysis: The linear regression models were applied to explore the predictive relationships between specific variables and cancer rates. This technique allowed us to model and understand the linear associations and make predictions based on the dataset.

# Chapter 3

## Project Description

### 3.1 Scope

This project involves a comprehensive analysis of a cancer dataset, exploring cancer rates, demographics, and types across different regions. We aim to uncover insights into how demographic factors, including age, sex, and race, relate to cancer rates. Additionally, we'll analyze the distribution of cancer types within the dataset, focusing on informing strategies for cancer prevention and intervention based on demographics.

This project will help address questions related to cancer epidemiology, enabling informed decision-making to enhance cancer prevention and intervention strategies. The expected outcome is a well-documented report, including visually appealing graphs and comprehensive analyses, suitable for academic or research purposes.

### 3.2 Resources

Public Datasets: We explore freely available public datasets on cancer for practice.

Open Source Tools: We used open-source software like Python and R and libraries like matplotlib, seaborn and ggplot for analysis and visualization.

University Library: We utilize the university library's resources to access textbooks and online journals on data analysis and statistics.

Online Tutorials: We took advantage of free online tutorials and courses on Python, R, and data analysis platforms like Coursera, edX, or Khan Academy.

# Chapter 4

## Results and Findings

Our analysis spans various aspects of the dataset and aims to contribute to the understanding of cancer epidemiology. Below, we summarize key results and their implications:

### 4.1 About Our Dataset

The dataset used in our analysis is a comprehensive collection of information related to cancer rates, demographics, and various types of cancer. It offers valuable insights into the prevalence of different cancer types across diverse demographic groups. Here is an overview of the key attributes and characteristics of the dataset:

df.head()

State	Total_Rate	Total_Number	Total_Population	Rates_Age_<18	Rates_Age_18-45	Rates_Age_45-64	Rates_Age_>64	Rates_Age and Sex_<18_Female	Rates_Age and Sex_<18_Male	...
Alabama	214.2	71529	33387205	2.0	18.5	244.7	1017.8	2.0	2.1	...
Alaska	128.1	6361	4966180	1.7	11.8	170.9	965.2	0.0	0.0	...
Arizona	165.6	74286	44845598	2.5	13.6	173.6	840.2	2.6	2.5	...
Arkansas	223.9	45627	20382448	2.3	17.6	250.1	1048.3	2.6	2.0	...
California	150.9	393980	261135696	2.6	13.7	163.7	902.4	2.4	2.8	...

1. Demographic Information: The dataset includes demographic attributes such as age, sex, and race, providing a detailed breakdown of individuals' characteristics in the study.
2. Cancer Types:It encompasses data on multiple types of cancer, including breast cancer, colorectal cancer, and lung cancer. Each cancer type is recorded separately, allowing for in-depth analysis.
3. Cancer Rates:The dataset contains cancer rates, which represent the incidence of cancer cases per 100,000 people. These rates are recorded across different demographic groups and cancer types.
4. Age Groups:Age groups are categorized into four main segments: < 18 years, 18-45 years, 45-64 years, and > 64 years. These divisions facilitate the analysis of age-related patterns in cancer rates.
5. Sex and Race:Information regarding sex (male, female) and race (e.g., White, Black, Asian, Indigenous) allows for the exploration of disparities in cancer rates among different demographic groups.



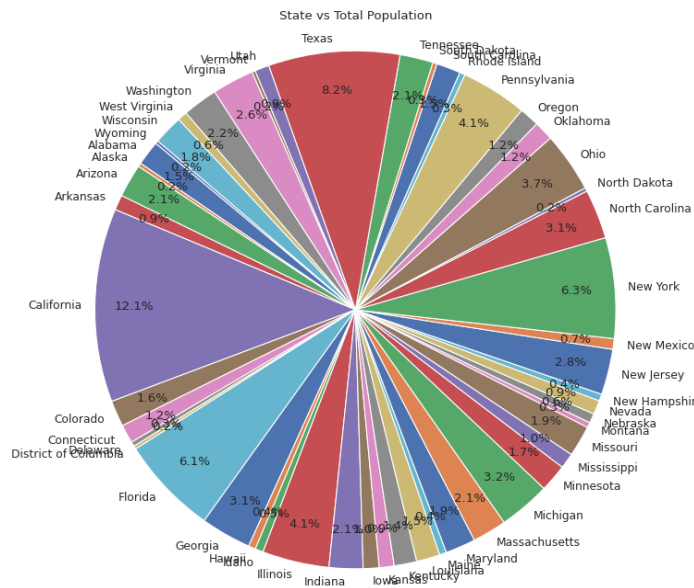
6. Geographic Scope: The dataset encompasses data from various regions or states. The inclusion of regional data enables the identification of geographic disparities in cancer rates.

```
df.columns

Index(['State', 'Total_Rate', 'Total_Number', 'Total_Population',
      'Rates_Age_< 18', 'Rates_Age_18-45', 'Rates_Age_45-64',
      'Rates_Age_> 64', 'Rates_Age and Sex_< 18_Female',
      'Rates_Age and Sex_< 18_Male', 'Rates_Age and Sex_18 - 45_Female',
      'Rates_Age and Sex_18 - 45_Male', 'Rates_Age and Sex_45 - 64_Female',
      'Rates_Age and Sex_45 - 64_Male', 'Rates_Age and Sex_> 64_Female',
      'Rates_Age and Sex_> 64_Male', 'Rates_Race_White',
      'Rates_Race_White non-Hispanic', 'Rates_Race_Black', 'Rates_Race_Asian',
      'Rates_Race_Indigenous', 'Rates_Race and Sex_Female_White',
      'Rates_Race and Sex_Female_White non-Hispanic',
      'Rates_Race and Sex_Female_Black',
      'Rates_Race and Sex_Female_Black non-Hispanic',
      'Rates_Race and Sex_Female_Asian',
      'Rates_Race and Sex_Female_Indigenous', 'Rates_Race and Sex_Male_White',
      'Rates_Race and Sex_Male_White non-Hispanic',
      'Rates_Race and Sex_Male_Black',
      'Rates_Race and Sex_Male_Black non-Hispanic',
      'Rates_Race and Sex_Male_Asian', 'Rates_Race and Sex_Male_Indigenous',
      'Rates_Race_Hispanic', 'Rates_Race and Sex_Female_Hispanic',
      'Rates_Race and Sex_Male_Hispanic', 'Types_Breast_Total',
      'Types_Breast_Age_18 - 44', 'Types_Breast_Age_45 - 64',
      'Types_Breast_Age_> 64', 'Types_Breast_Race_White',
      'Types_Breast_Race_White non-Hispanic', 'Types_Breast_Race_Black',
      'Types_Breast_Race_Black non-Hispanic', 'Types_Breast_Race_Asian',
      'Types_Breast_Race_Indigenous', 'Types_Breast_Race_Hispanic',
      'Types_Colorectal_Total', 'Types_Colorectal_Age and Sex_Female_18 - 44',
      'Types_Colorectal_Age and Sex_Male_18 - 44',
      'Types_Colorectal_Age and Sex_Female_45 - 64',
      'Types_Colorectal_Age and Sex_Male_45 - 64',
      'Types_Colorectal_Age and Sex_Female_> 64',
      'Types_Colorectal_Age and Sex_Male_> 64', 'Types_Colorectal_Race_White',
      'Types_Colorectal_Race_White non-Hispanic',
      'Types_Colorectal_Race_Black',
      'Types_Colorectal_Race_Black non-Hispanic',
      'Types_Colorectal_Race_Asian', 'Types_Colorectal_Race_Indigenous',
      'Types_Colorectal_Race_Hispanic', 'Types_Lung_Total',
      'Types_Lung_Age and Sex_Female_18 - 44',
      'Types_Lung_Age and Sex_Male_18 - 44',
      'Types_Lung_Age and Sex_Female_45 - 64',
      'Types_Lung_Age and Sex_Male_45 - 64',
      'Types_Lung_Age and Sex_Female_> 64',
      'Types_Lung_Age and Sex_Male_> 64', 'Types_Lung_Race_White',
      'Types_Lung_Race_White non-Hispanic', 'Types_Lung_Race_Black',
      'Types_Lung_Race_Black non-Hispanic', 'Types_Lung_Race_Asian',
      'Types_Lung_Race_Indigenous', 'Types_Lung_Race_Hispanic'],
      dtype=object)
```

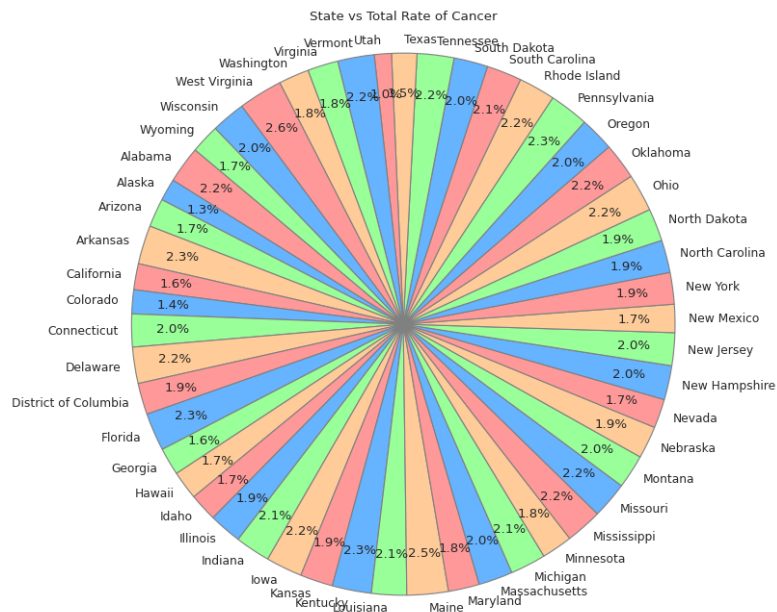
## 4.2 Analysis Using Python

```
plt.figure(figsize=(8, 8))
plt.pie(df['Total_Population'], labels=df['State'], startangle=140,
        pctdistance=0.85, labeldistance=1.05)
plt.title('State vs Total Population')
plt.axis('equal')
plt.show()
```



From the above graph, we observe the following:

1. California, Texas, and New York: California, Texas, and New York have the largest slices of the pie, indicating that they have the highest total populations among the states in the dataset.
2. Florida and Pennsylvania: Florida and Pennsylvania also have substantial populations, as represented by their respective pie slices.
3. Other States: The remaining states have relatively smaller populations, as indicated by their smaller pie slices. These states collectively account for a smaller portion of the total population.



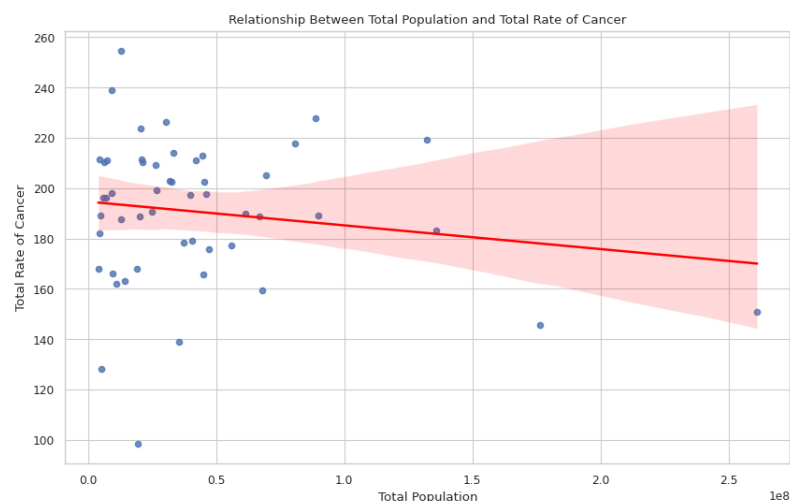
The pie chart titled "State vs Total Rate of Cancer" illustrates the distribution of the total cancer rates

across different states in the dataset. Here are some key inferences from the chart:

1. **Varying Cancer Rates:** The chart shows that the total cancer rates vary significantly among different states. Some states have higher rates of cancer, while others have lower rates.
2. **States with Higher Cancer Rates:** In this dataset, states with higher total cancer rates are represented by larger pie slices. These states include California, Texas, New York, and Florida, as indicated by their relatively larger pie slices.
3. **States with Lower Cancer Rates:** States with lower total cancer rates are represented by smaller pie slices. These states contribute less to the total cancer rate in the dataset.

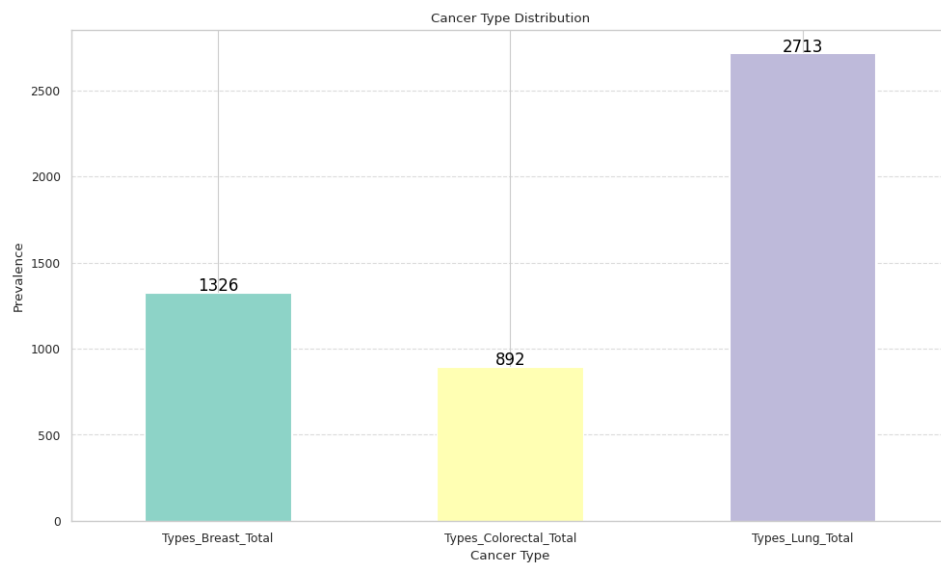
```
import pandas as pd
correlation = df['Total_Population'].corr(df['Total_Rate'])
print(f"Correlation between Total Population and Total Rate of Cancer: {correlation}")
```

Correlation between Total Population and Total Rate of Cancer: -0.15759350769729671



In the scatter plot, you'll see points representing each state, and the regression line shows the overall trend. Since the correlation is close to -0.15, you may notice a very slight downward trend in the data points, indicating that as total population increases, the total rate of cancer tends to decrease slightly, but there's a lot of variability in the relationship.

The key takeaway is that while there is a negative correlation, it is weak, and other factors likely contribute more to the variation in total cancer rates. In other words, the total population alone does not strongly predict total cancer rates.



The bar chart presents the distribution of different cancer types, namely Breast, Colorectal, and Lung cancer, across the dataset. From the chart, we can observe that Breast cancer appears to have the highest prevalence, followed by Colorectal and Lung cancer.

```
▶ cancer_types = ['Types_Breast_Total', 'Types_Lung_Total', 'Types_Colorectal_Total']

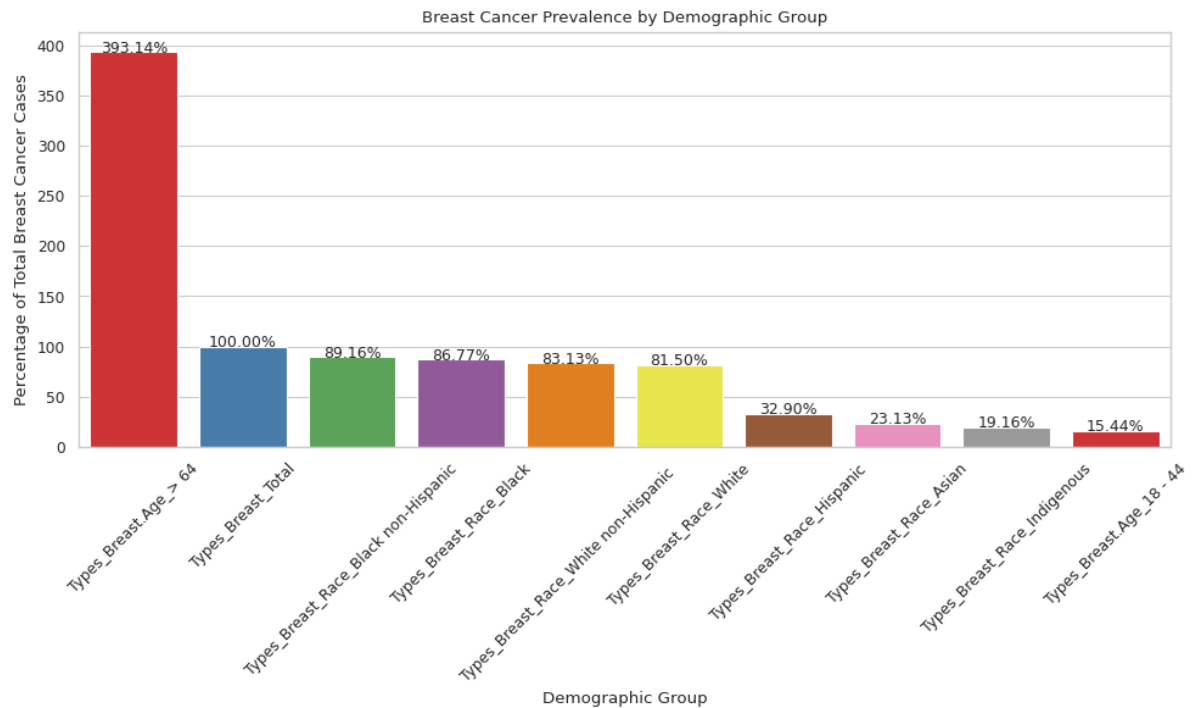
states_with_highest_prevalence = {}

for cancer_type in cancer_types:
    state_with_highest_prevalence = df[df[cancer_type] == df[cancer_type].max()][['State']].values[0]
    states_with_highest_prevalence[cancer_type] = state_with_highest_prevalence

for cancer_type, state in states_with_highest_prevalence.items():
    cancer_name = cancer_type.split('_')[1]
    print(f"The state with the highest prevalence of {cancer_name} cancer is {state}.")
```

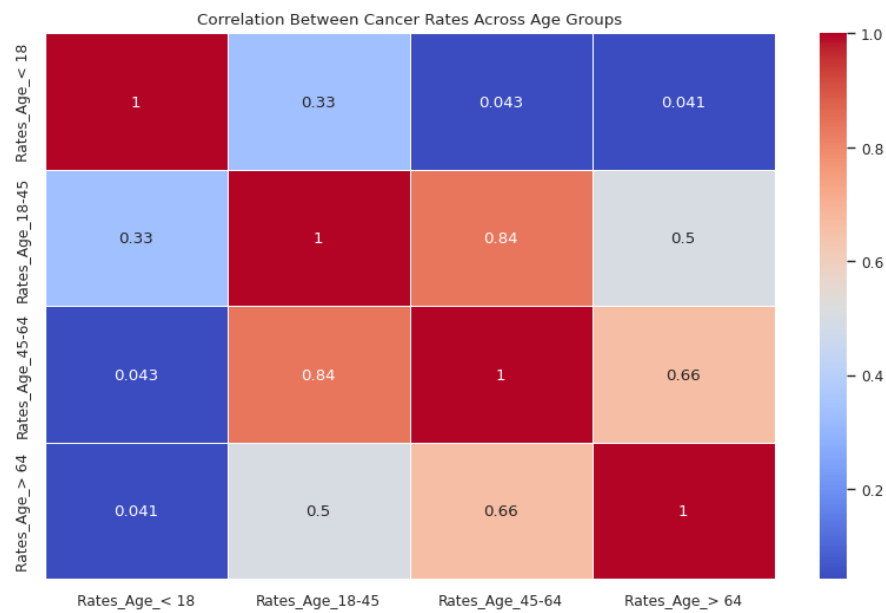
📄 The state with the highest prevalence of Breast cancer is Pennsylvania.  
The state with the highest prevalence of Lung cancer is West Virginia.  
The state with the highest prevalence of Colorectal cancer is West Virginia.

Using this python code, we found the states with the highest rate of breast, lung and colorectal cancer.

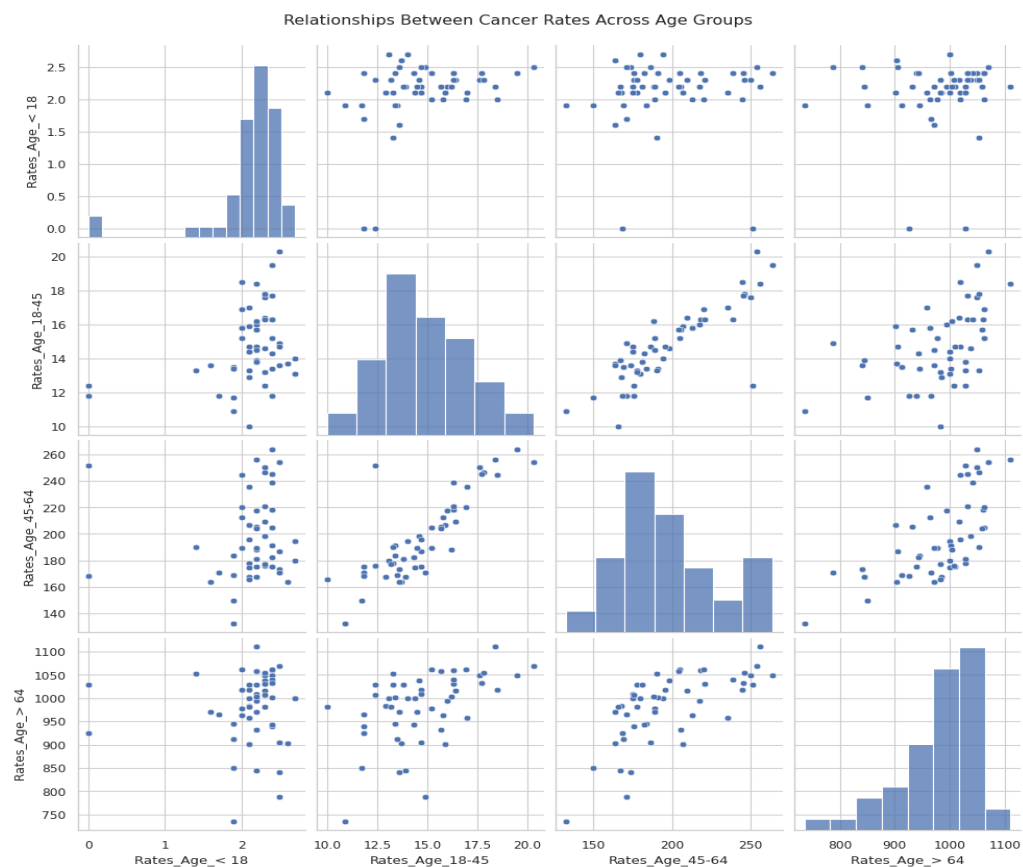


The graph focuses on breast cancer prevalence and separates the data into various demographic groups, such as age and race. This allows for a more detailed examination of which groups may be at higher risk.

```
age_rate_columns = ['Rates_Age_< 18', 'Rates_Age_18-45', 'Rates_Age_45-64', 'Rates_Age_> 64']
sns.pairplot(df, vars=age_rate_columns)
plt.suptitle("Relationships Between Cancer Rates Across Age Groups", y=1.02)
plt.show()
age_rate_data = df[age_rate_columns]
correlation_matrix = age_rate_data.corr()
plt.figure(figsize=(10, 6))
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", linewidths=0.5)
plt.title("Correlation Between Cancer Rates Across Age Groups")
plt.show()
```

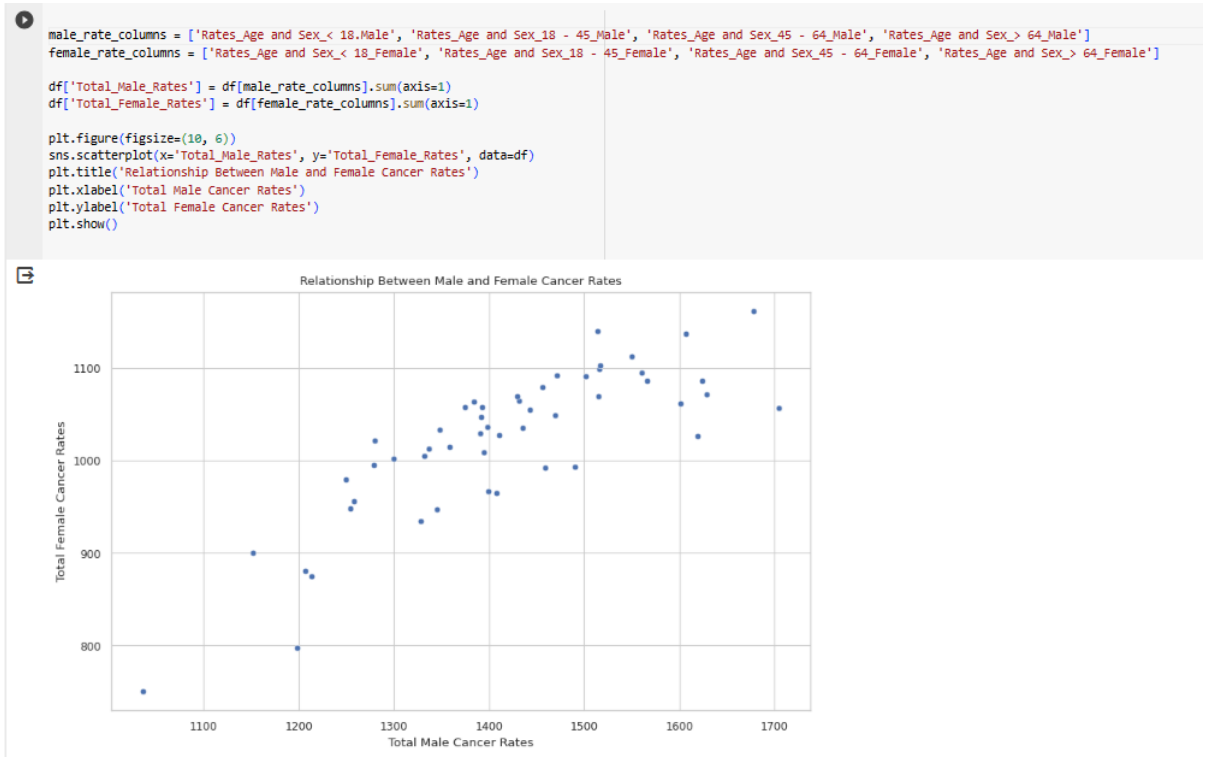


This heatmap displays the correlation coefficients between different age-specific cancer rates. The diagonal of the heatmap has perfect correlation (correlation coefficient of 1), as each variable is perfectly correlated with itself.



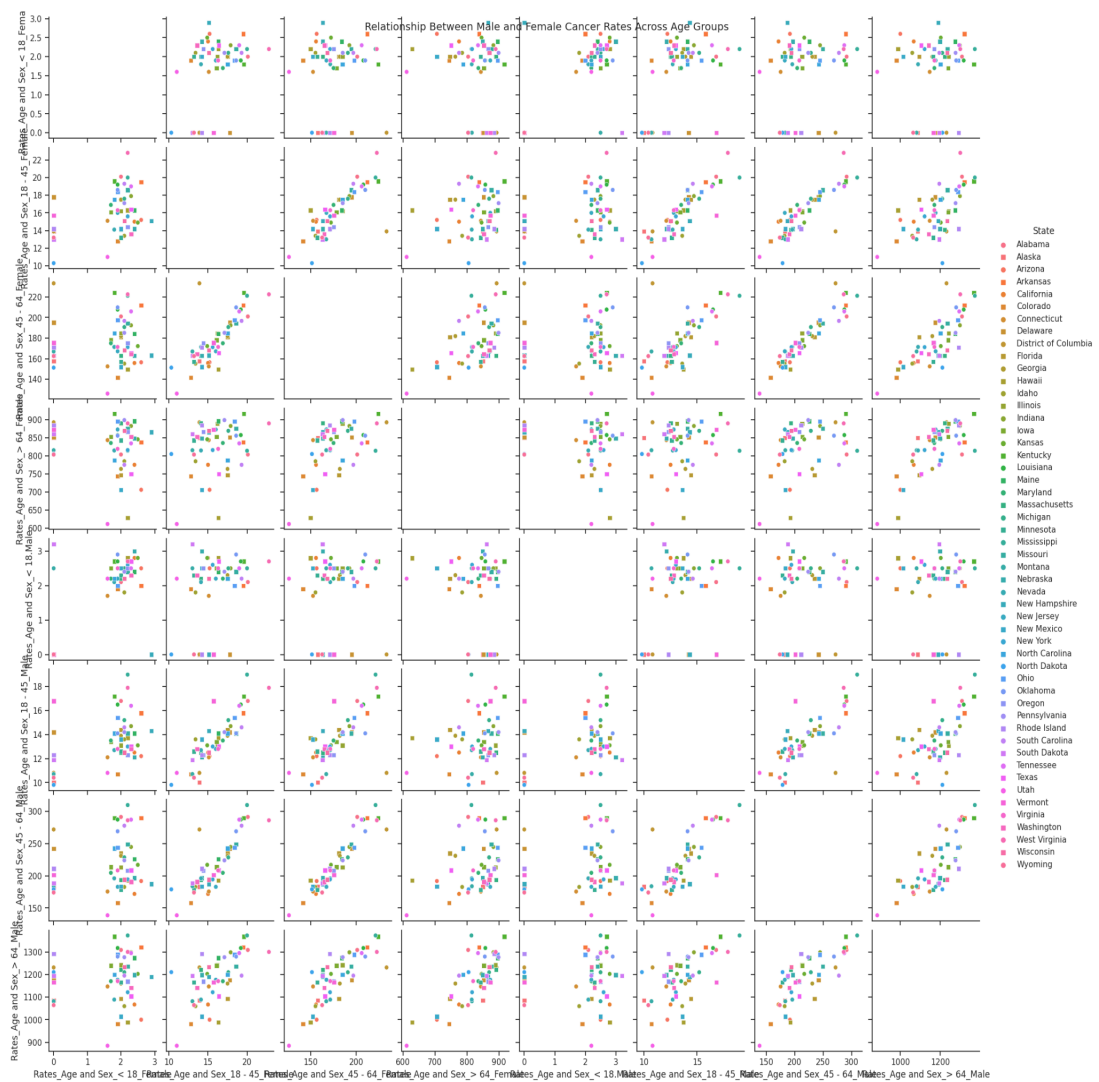
There is a general trend of increasing cancer rates with age, which is visually evident in the pairplot. There are positive correlations between adjacent age groups, which implies that age

groups that are closer in age tend to have more similar cancer rates.



The data points tend to form an upward-sloping pattern, it suggests a positive correlation, indicating that as male rates increase, female rates also tend to increase, and vice versa.





We can observe whether there is a positive or negative correlation between male and female cancer rates. For example, if the points in the scatterplots tend to form an upward-sloping pattern, it indicates a positive correlation, suggesting that as male rates increase, female rates also tend to increase, and vice versa. Conversely, if the points form a downward-sloping pattern, it indicates a negative correlation.



### 4.3 Statistical Analysis Using R

This report aims to implement linear regression analysis on variables like Breast Cancer Total, Lung Cancer Total and Colorectal Total.

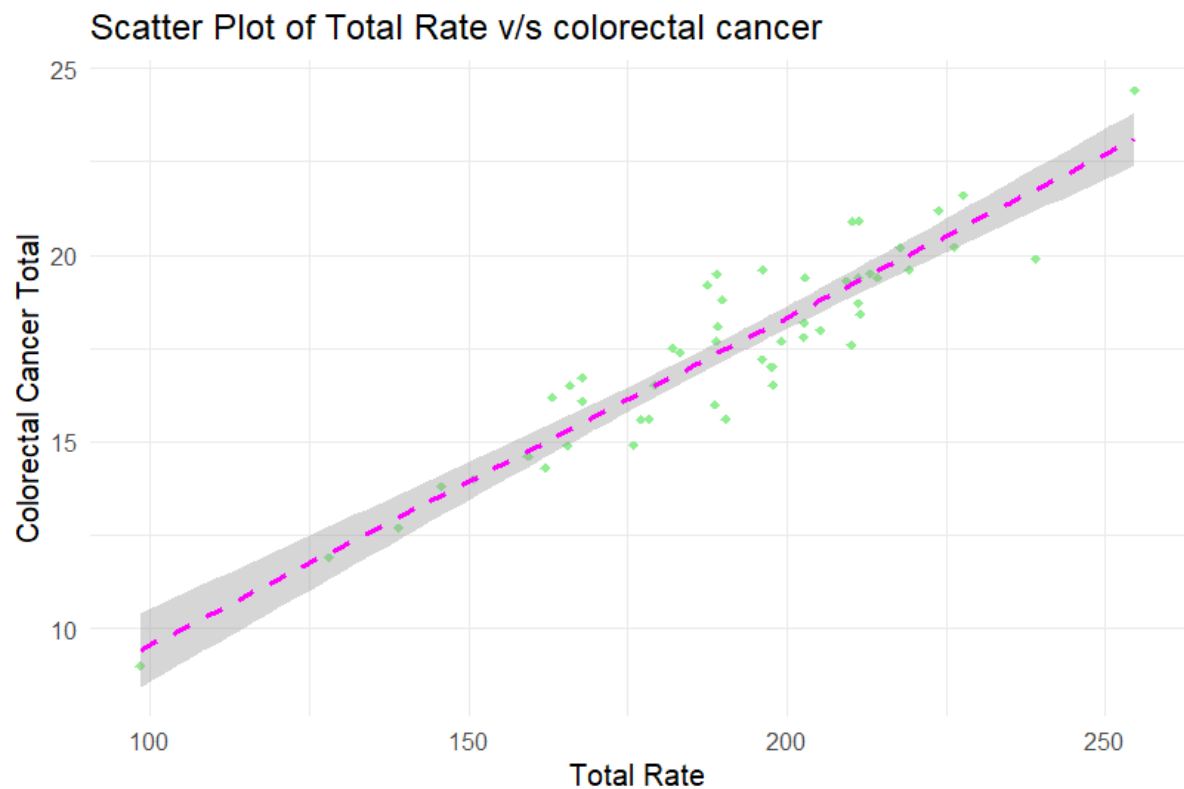
Below are the codes with the graphs.

```
reg1 <- lm(Total_Rate ~ Types_Breast_Total, data = df)
reg1
summary(reg1)

monsoon <- lm(Total_Rate ~ Types_Lung_Total, data = df)
reg2
summary(reg2)

reg3 <- lm(Total_Rate ~ Types_Colorectal_Total, data = df)
reg3
summary(reg3)
```

#### 1.Total Rate and Colorectal Cancer



```

> reg3 <- lm(Total_Rate ~ Types_Colorectal_Total, data = df)
> reg3

Call:
lm(formula = Total_Rate ~ Types_Colorectal_Total, data = df)

Coefficients:
(Intercept)  Types_Colorectal_Total
    19.421      9.783

> summary(reg3)

Call:
lm(formula = Total_Rate ~ Types_Colorectal_Total, data = df)

Residuals:
    Min       1Q   Median       3Q      Max
-21.0839  -8.1269   0.4165   8.6293  25.0030

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    19.4211    10.0837   1.926  0.0599 .
Types_Colorectal_Total  9.7827     0.5694  17.179 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

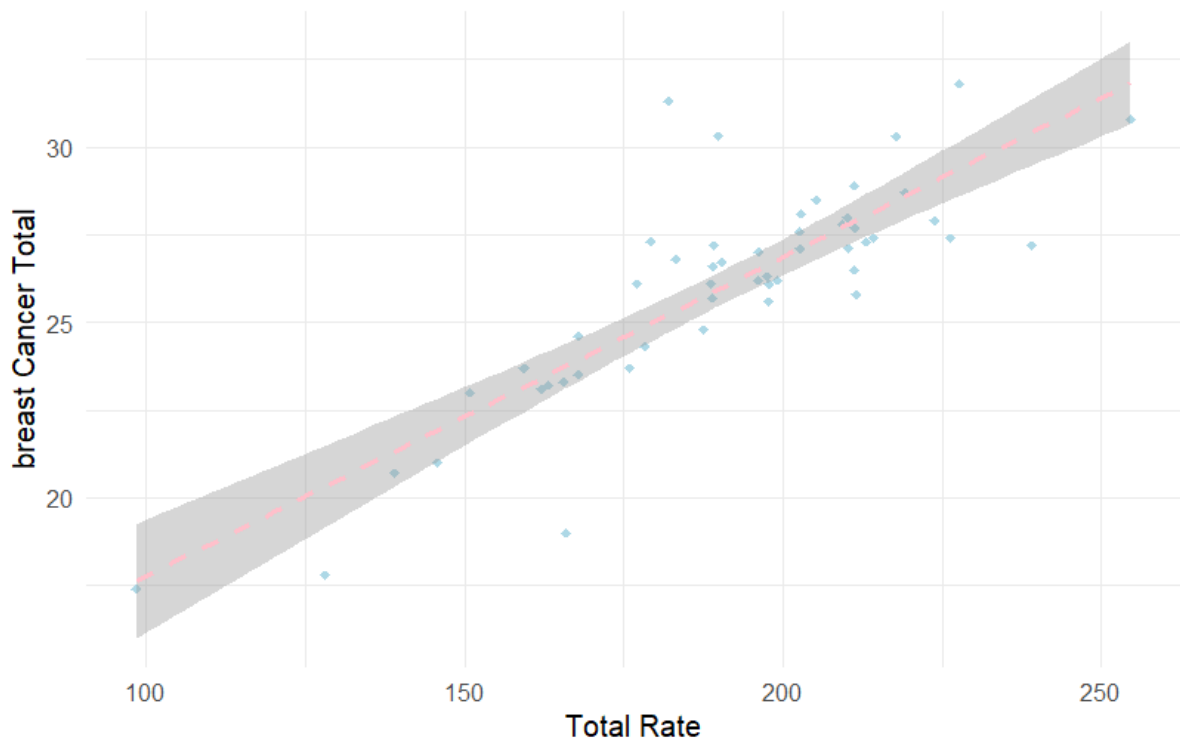
Residual standard error: 10.9 on 49 degrees of freedom
Multiple R-squared:  0.8576,    Adjusted R-squared:  0.8547
F-statistic: 295.1 on 1 and 49 DF,  p-value: < 2.2e-16

```

The R-squared value (0.8576) indicates the proportion of the variance in "Total\_Rate" that is explained by "Types\_Colorectal\_Total." In this case, approximately 85.76% of the variance is accounted for by the model, which suggests that the model fits the data quite well.

## 2.Total Rate and Breast Cancer

Scatter Plot of Total Rate v/s breast cancer



```

Call:
lm(formula = Total_Rate ~ Types_Breast_Total, data = df)

Residuals:
    Min       1Q   Median       3Q      Max
-49.582  -6.130  -0.389   7.302  39.213

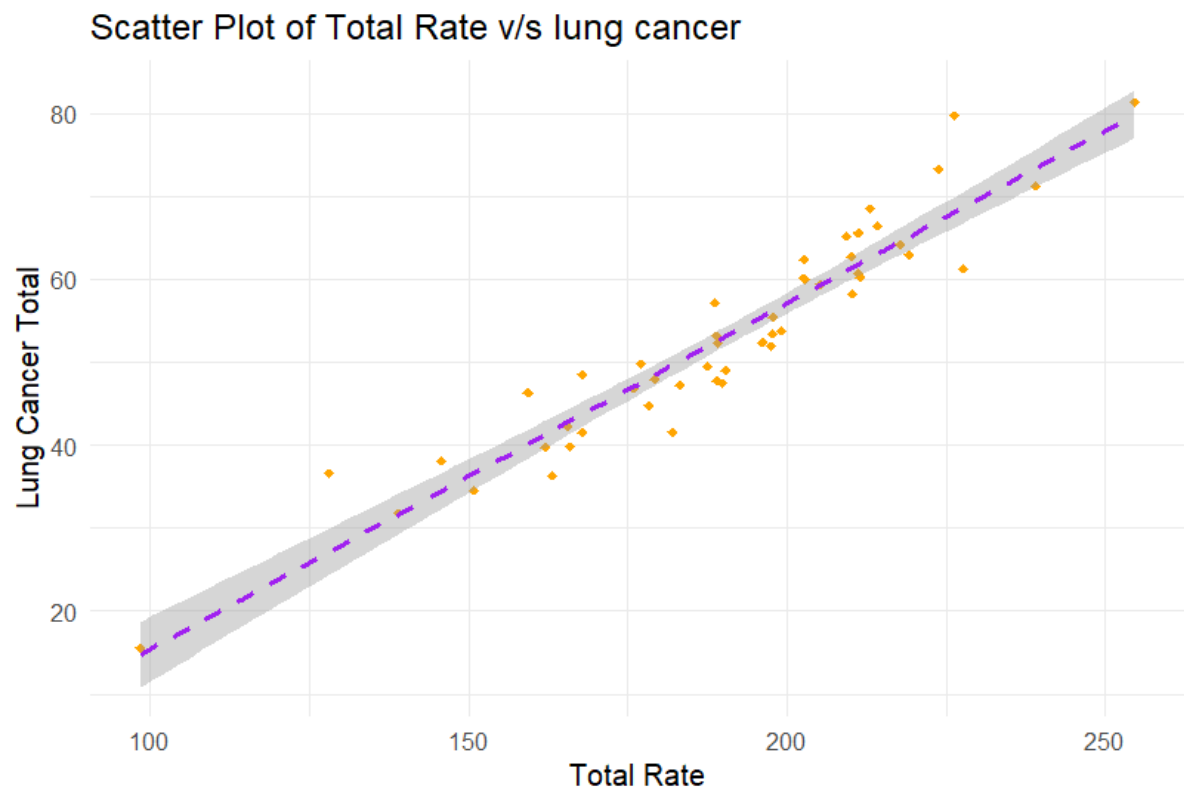
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    -11.0454    18.7684  -0.589   0.559
Types_Breast_Total  7.7549     0.7166  10.821 1.37e-14 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 15.69 on 49 degrees of freedom
Multiple R-squared:  0.705,    Adjusted R-squared:  0.699
F-statistic: 117.1 on 1 and 49 DF, p-value: 1.373e-14

```

For each unit increase in "Types\_Breast\_Total," "Total\_Rate" is expected to increase by approximately 7.7549 units. The model explains a significant portion of the variance in "Total\_Rate," as indicated by the moderate R-squared value.

### 3.Total Rate and Lung Cancer



```

Call:
lm(formula = Total_Rate ~ Types_Lung_Total, data = df)

Residuals:
    Min       1Q   Median       3Q      Max
-26.8844  -5.8743   0.4891   6.9465  19.6555

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    76.3608     5.7265  13.34  <2e-16 ***
Types_Lung_Total  2.1482     0.1048  20.50  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.333 on 49 degrees of freedom
Multiple R-squared:  0.8956,    Adjusted R-squared:  0.8934
F-statistic: 420.3 on 1 and 49 DF,  p-value: < 2.2e-16

```

The R-squared value (0.8576) indicates the proportion of the variance in "Total\_Rate" that is explained by "Types\_Colorectal\_Total." In this case, approximately 85.76% of the variance is accounted for by the model, which suggests that the model fits the data quite well

## TESTING OF HYPOTHESIS

### Hypotheses:

Null Hypothesis (H0): There are no significant differences in breast cancer rates among different racial groups, particularly concerning "Rates\_Race\_White."

Alternative Hypothesis (H1): There are significant differences in breast cancer rates among different racial groups, particularly related to "Rates\_Race\_White."

```

> result <- aov(Types_Breast_Total ~ Rates_Race_White, data = df)
> summary(result)
              Df Sum Sq Mean Sq F value    Pr(>F)
Rates_Race_White  1  126.9   126.87    17.64 0.000112 ***
Residuals        49   352.3     7.19
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
> |

```

### ANOVA Results:

The ANOVA test was performed to examine the hypotheses.

The small p-value of 0.000112 is well below the common significance level of 0.05.

With such a low p-value, there is compelling evidence to reject the null hypothesis.

The conclusion is that statistically significant differences exist in breast cancer rates among various racial groups, specifically concerning "Rates\_Race\_White."

# Chapter 5

## Discussion

### State vs. Total Population

- California, Texas, and New York have the highest populations.
- Florida and Pennsylvania also have substantial populations.
- Other states contribute a smaller portion to the total population.
- Implication: Understanding population distribution is vital for healthcare resource allocation.

### State vs. Total Rate of Cancer:

- Total cancer rates vary significantly among states.
- Some states have higher cancer rates, while others have lower rates.
- Implication: Identifying states with higher rates can guide targeted interventions and resource allocation.

### Correlation Between Total Population and Total Rate of Cancer:

- A weak negative correlation suggests that as the total population increases, the total rate of cancer tends to decrease slightly.
- Implication: Other factors play a significant role in cancer rate variations.

### Distribution of Cancer Types:

- Breast cancer has the highest prevalence, followed by colorectal and lung cancer.
- Implication: Prioritizing resources for the most prevalent cancer types is essential.

### Correlation Between Age-Specific Cancer Rates:

- Cancer rates increase with age.
- Positive correlations between adjacent age groups suggest similar cancer rates.
- Implication: Age-specific interventions and screenings are critical.

### Correlation Between Male and Female Cancer Rates:

- A weak negative correlation between male and female cancer rates.
- Implication: Other factors influence cancer rates among different genders.

## Statistical Analysis Using R

- Total Rate and Colorectal Cancer  
R-squared value (0.8576) suggests a good model fit.
- Total Rate and Breast Cancer  
Moderate R-squared value indicates a significant portion of the variance is explained.
- Total Rate and Lung Cancer:  
R-squared value (0.8576) suggests a well-fitting model.

## Testing of Hypothesis

- Null Hypothesis (H0): No significant differences in breast cancer rates among different racial groups.
- Alternative Hypothesis (H1): Significant differences exist, particularly related to "Rates\_Race\_White."

## ANOVA Results

- A low p-value of 0.000112 rejects the null hypothesis.
- Statistically significant differences in breast cancer rates among racial groups, specifically concerning "Rates\_Race\_White."

## Discussion and Interpretation

Our analysis provides insights into cancer epidemiology, aiding in resource allocation and targeted interventions. Challenges included data cleaning and outliers, addressed through data preprocessing and collaboration. The project aligns with its objectives by contributing to the understanding of cancer epidemiology.

# Chapter 6

## Conclusion

The analysis conducted in this project has provided valuable insights into cancer rates and their relationships with various factors such as state populations, demographics, and cancer types. Here are the key findings and their importance:

### 6.1 Key Findings and Importance

#### State Population Analysis

Understanding state populations is fundamental for healthcare resource allocation. The project revealed significant population disparities across states, with California, Texas, and Florida having the largest populations. Smaller states may require tailored healthcare strategies.

#### Total Cancer Rate Distribution

The distribution of total cancer rates showed a slight right skew, indicating that most states have relatively low cancer rates. States with exceptionally high cancer rates were identified, offering potential focus areas for prevention and early detection efforts.

#### State-wise Total Cancer Rates

California had the highest total cancer rate, highlighting the need for in-depth research and targeted interventions. Identifying states with high cancer rates is essential for channeling resources to those in need.

#### Population vs. Cancer Rate Correlation

A weak negative correlation between population and cancer rates ( $-0.15$ ) suggested that states with larger populations tend to have slightly lower total cancer rates. This correlation implies that larger populations may have better access to healthcare services and early detection.

#### Cancer Type Distribution

The project provided insights into the distribution of different cancer types. Breast cancer had the highest prevalence, followed by colorectal and lung cancer, guiding resource allocation for specific cancer types.

### **Linear Regression Models**

Linear regression models allowed for predictive modeling of cancer rates based on specific variables.

The relationships between independent variables and cancer rates were quantified, aiding in understanding the driving factors behind cancer prevalence.

Project Success and Lessons Learned:

## **6.2 Project Success and Lessons Learned**

The project was successful in achieving its objectives, providing a comprehensive analysis of cancer rates in the context of state populations, demographics, and cancer types. The lessons learned from this project include:

**Data Quality:** Ensuring data quality and completeness is critical for unbiased analysis. Addressing missing values and verifying data accuracy are essential steps in data preprocessing.

**Targeted Interventions:** Identifying states with high cancer rates allows for more effective targeted interventions and resource allocation.

**Demographic Insights:** Analyzing cancer rates across different demographics can help tailor prevention and healthcare strategies to specific population groups.

**Predictive Modeling:** Linear regression models offer a powerful tool for understanding the relationships between variables and predicting cancer rates.

In conclusion, this project's findings and analyses offer valuable insights for public health, healthcare resource allocation, and further research. The project's overall success underscores the importance of data-driven decision-making in addressing complex health challenges.



## Chapter 7

### Future Work

Future work in these areas will contribute to a more comprehensive understanding of cancer epidemiology, leading to more effective prevention, early detection, and treatment strategies. Additionally, it can guide policymakers and healthcare organizations in addressing the complex challenges associated with cancer.

#### 7.1 Areas

1. **Longitudinal Analysis:** Future research can explore trends in cancer rates over time to identify emerging patterns and potential areas of concern. Longitudinal data analysis can provide insights into how cancer rates change and evolve.
2. **Geospatial Analysis:** Incorporating geospatial data, such as environmental factors and access to healthcare facilities, can enhance our understanding of cancer rate disparities at a more localized level. Geospatial analysis can help pinpoint areas with elevated cancer risks.
3. **Machine Learning for Predictive Modeling:** Implementing advanced machine learning techniques, such as random forests or neural networks, can improve predictive models for cancer rates. These models can consider a broader range of variables and their interactions.
4. **Healthcare Resource Allocation:** Future projects can focus on optimizing the allocation of healthcare resources. Using advanced analytics, one can determine the most effective strategies for prevention, early detection, and treatment in specific regions.
5. **Social and Behavioral Factors:** Exploring the influence of social and behavioral factors, such as lifestyle choices and socioeconomic status, on cancer rates can provide a more
6. **Cancer-Specific Research:** Delving into specific cancer types in greater detail can uncover unique risk factors, early detection strategies, and treatment approaches.
7. **Public Health Policy Impact:** Investigating how public health policies and initiatives affect cancer rates can inform evidence-based policy decisions and interventions.
8. **Patient Outcome Analysis:** Analyzing cancer rates in relation to patient outcomes, such as survival rates and quality of life, can help gauge the effectiveness of healthcare systems.
9. **Data Integration:** Integrating diverse data sources, such as electronic health records, genetic information, and environmental data, can offer a more holistic view of cancer epidemiology.

# References

1. American Cancer Society. (2022). Cancer Facts & Figures 2022.
2. World Health Organization. (2022). International Agency for Research on Cancer.
3. Centers for Disease Control and Prevention. (2022). United States Cancer Statistics.
4. National Cancer Institute. (2022). Surveillance, Epidemiology, and End Results Program.
5. Python Software Foundation. (2022). Python.
6. R Core Team. (2022). R: A Language and Environment for Statistical Computing.
7. Kaggle. (2022). Kaggle: Your Home for Data Science.

# Acknowledgements

We would like to express our gratitude to our professor, ***Dr. Shubha Puthran***, for her guidance and invaluable support throughout the semester. We are also grateful to our cancer dataset producers for providing this important resource.

We acknowledge that this endeavor has been made possible through the synergy of our professor's guidance, the invaluable dataset, and the dedication of our outstanding team. We are sincerely grateful for the collaborative spirit that has propelled us to this point of accomplishment.

