**Damilare Kolawole**
**X21235571**

# 1.0 INTRODUCTION

This study aims to investigate the socio-economic factors influencing the Cancer mortality rate in the US using multiple regression models to help discover socio-economic factor that has an higher influence.

It is important to understand multiple regression and the assumptions associated with it before this research is done.

Simply put, multiple linear regression is a statistical method used to examine the relationship between dependent and several independent variables. To predict the value of the dependent variable, a linear equation is derived using the values of the independent ones.

The different socio-economic factors are the independent variables while the mortality rate is the dependent variable in the US cancer mortality rate study. The following multiple linear regression assumptions should be met.[1]

- ❖ Linearity refers to the connection between the dependent variable and the independent variables.
- ❖ Independence: The findings should be distinct from one another.
- ❖ Homoscedasticity: The residual variance (the gap between predicted and actual values) should be constant across all levels of the independent variables.
- ❖ Normality: The residuals should be distributed normally.
- ❖ There should be no multicollinearity, which means that the independent variables should not be strongly correlated with one another.
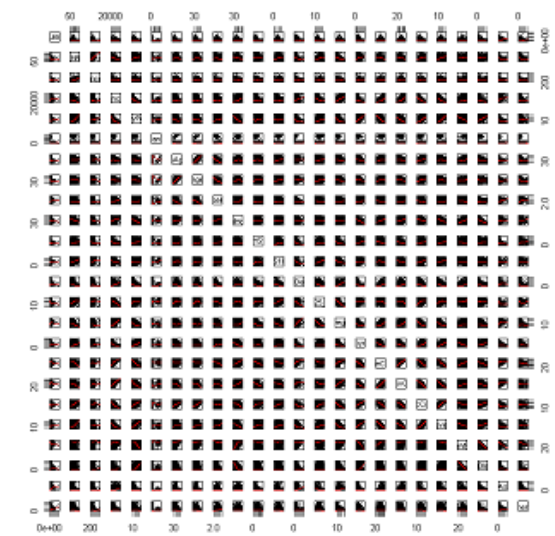
When these assumptions are satisfied, the multiple linear regression model can be constructed by calculating the coefficients of the independent factors in the linear equation. After determining the values of the independent variables, the model can be utilized in forecasting the mortality rate. Various socioeconomic factors associated with cancer can be used as independent variables in the instance of cancer mortality rates in the United States.[2] [3]

Finally, multiple linear regression is an effective statistical method for examining the relationship between a dependent variable and numerous independent variables. In the case of cancer mortality rates in the United States, different socio-economic factors can be used as independent variables to predict mortality rates. Before conducting multiple linear regression, the assumptions of linearity, and independence, homoscedasticity, normality, and no multicollinearity should be met.

## 2.0 DESCRIPTION OF DATASET AND ITS VARIABLES

The R programming language will be used to carry out this experiment. The dataset got consists of 3047 observations and 25 columns. Out of these twenty-five columns (25), 23 independent variables and 1 dependent variable. in this case, the use of multiple linear regression becomes more appropriate to analyze the relationship between 2 or more predictor variables and 1 continuous variable. Below are the socio-economic factors.

- ❖ Income related: the median income, the percentage of unemployed, and the poverty rate
- ❖ Age related: Median Age across the population, and for male and female separately
- ❖ Household related: Average Household Size and percentage of Married Households
- ❖ Education related: Percentage of the highest educational level attained (No High School / High School / Bachelor Degree) in the age groups 18-24 and over 25.
- ❖ Health Insurance related: Percentage of Private Insurance, Private Insurance paid by
- ❖ Employer, Public Insurance and Public Insurance Only, and
- ❖ Race related: percentage of White/Black/Asian/Other.

*Figure 3checking if the data is continuous.*

### A. VISUALIZATION OF THE DATA

In this section, we visualize the data using a boxplot and histogram. Each variable histogram was plotted to see if the data is normalized, and the boxplot function was used to check for outliers in the dataset.
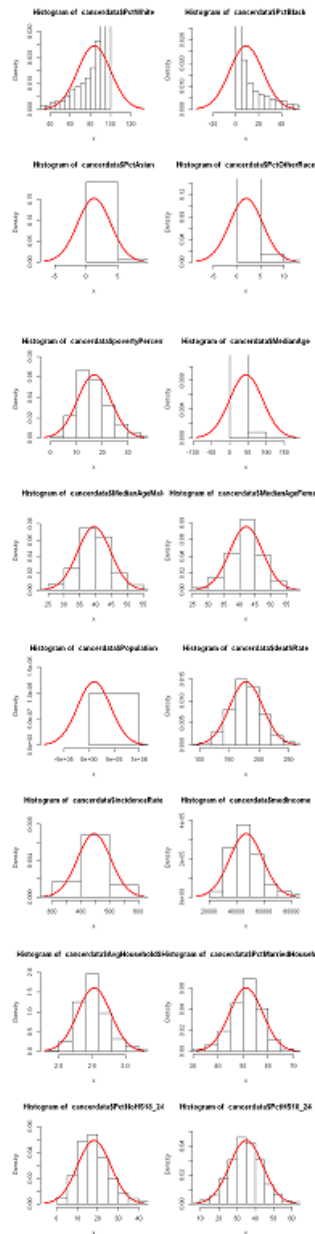
| County | Population | deathRate | incidenceRate | medIncome | povertyPercent | MedianAge | MedianAgeMale | MedianAgeFemale |
|--------|-----------|-----------|---------------|-----------|----------------|-----------|---------------|-----------------|
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |

*Figure 4Histogram of the various variable*

Now to check clearly for the outliers, having done the summary some specific variables were carefully selected. The figure below shows a boxplot of these variables.
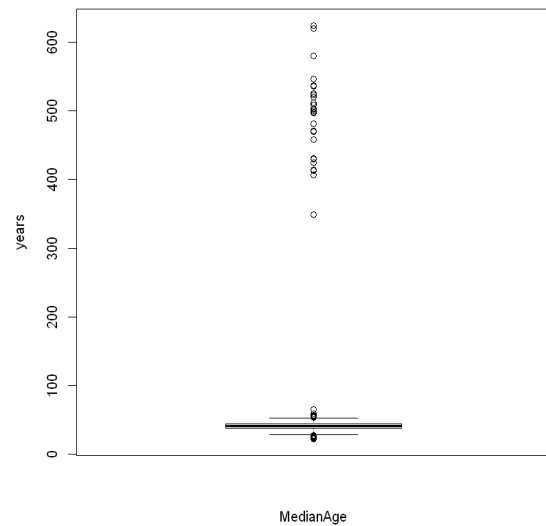


*Figure 5 Boxplot for Median Age*

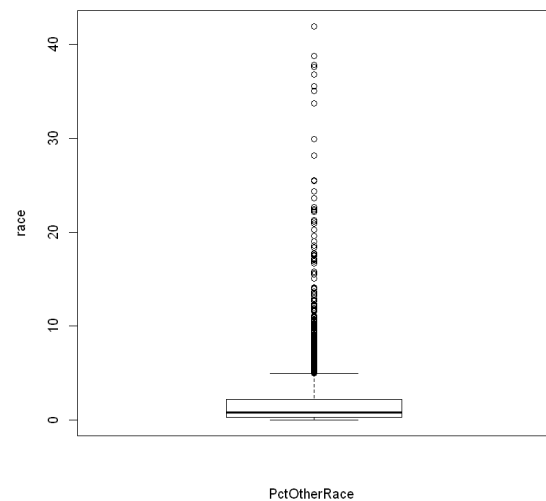Figure 5 shows the boxplot for the median age, we can see so many outliers in the age.



*Figure 6 Boxplot for other race*

The outliers in other-race variable is also much and must be taken care of and lastly, a boxplot for the percentage of degree bachelors between 18 and 24 is shown in the figure below:
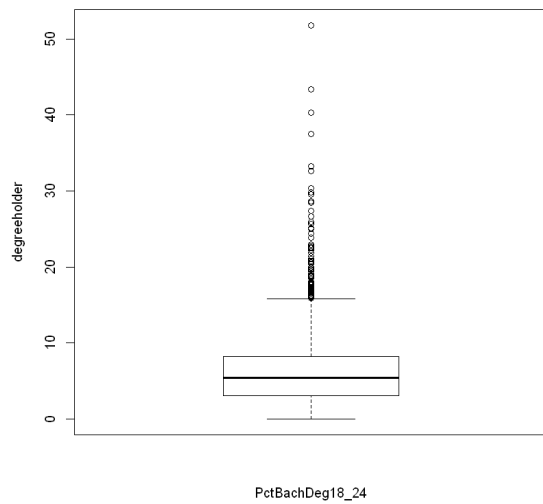
*Figure 7 Boxplot for PCTBct18-24*

## 3.0 MODEL

This section discusses the model used and the one which comes out best. After we have processed our data, by checking the missing values of which none was found, also checked if the data is continuous and made sure all assumptions have been met.

Our first model was built by comparing all variables as shown in the figure below:



*Figure 8: Model 1*

The model above shows that there is a relationship between the predictors and the response. Notice that some predictors do not have significant statistical effects. The r squared (R^2) value shows that 53% of the predictors can explain the changes in this multiple regression model.

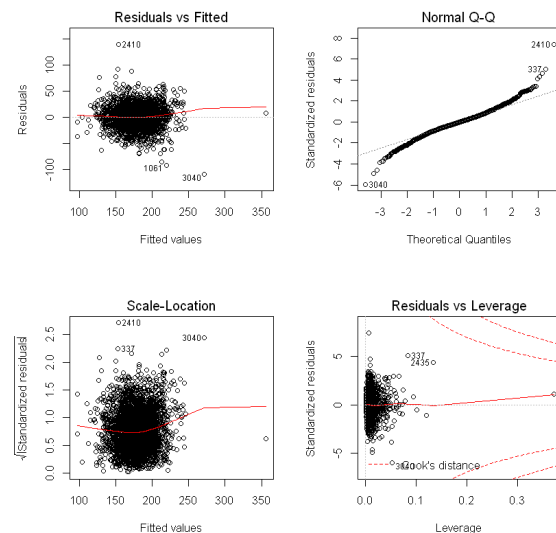The next figure below visualizes the model, checking the residuals, normal q-q



*Figure 9 Visualization of Model 1*

The figure above shows us visualized results for model 1. We have some outliers that need to be removed before plotting the next model.



*Figure 10 Model 2*

The model above shows that there is a relationship between the predictors and the response. Notice that some predictors have less significant statistical effects. The r squared (R^2) value shows that 53% of the predictors can explain the changes with an increased F statistics in this multiple regression model.

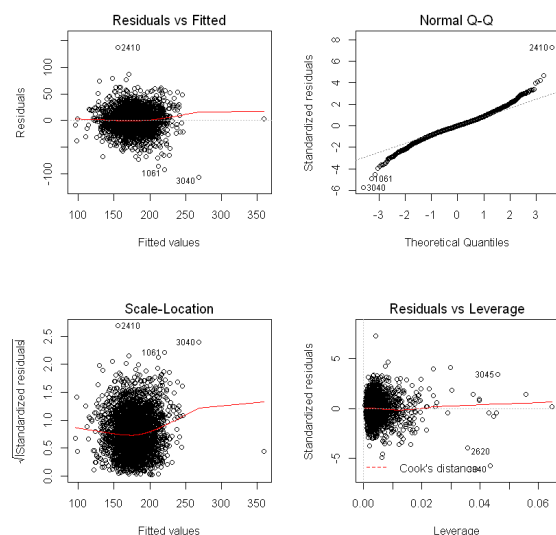The next figure below visualizes the model, checking the residuals, normal q-q

predictors have significant statistical effects. The r squared (R^2) value shows that 53% of the predictors and also a better F statistics, the residual error was also improved on. This improvement explain the changes in this multiple regression model.

The next figure below visualizes the model, checking the residuals, normal q-q



*Figure 11 the visualization of model 2*

```
Call:
lm(formula = deathRate ~ incidenceRate + povertyPercent + AvgHouseholdSize +
    PctHS18_24 + PctBachDeg25_Over + PctPrivateCoverage + PctWhite +
    PctEmpPrivCoverage + PctOtherRace, data = cancerdata)

Residuals:
     Min       1Q   Median       3Q      Max
-108.387  -10.353    0.125   10.505  135.633

Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)        144.320794   9.904224  14.572  < 2e-16 ***
incidenceRate        0.211355   0.006328  33.400  < 2e-16 ***
povertyPercent       0.567009   0.109019   5.201 2.11e-07 ***
AvgHouseholdSize    -9.305004   1.769947  -5.257 1.56e-07 ***
PctHS18_24           0.255031   0.043112   5.916 3.68e-09 ***
PctBachDeg25_Over   -1.459241   0.090516 -16.121  < 2e-16 ***
PctPrivateCoverage  -0.735072   0.090252  -8.145 5.50e-16 ***
PctWhite            -0.119089   0.027460  -4.337 1.49e-05 ***
PctEmpPrivCoverage   0.575576   0.074163   7.761 1.14e-14 ***
PctOtherRace        -0.942433   0.109439  -8.611  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.96 on 3037 degrees of freedom
Multiple R-squared:  0.5347,	Adjusted R-squared:  0.5333
F-statistic: 387.8 on 9 and 3037 DF,  p-value: < 2.2e-16
```
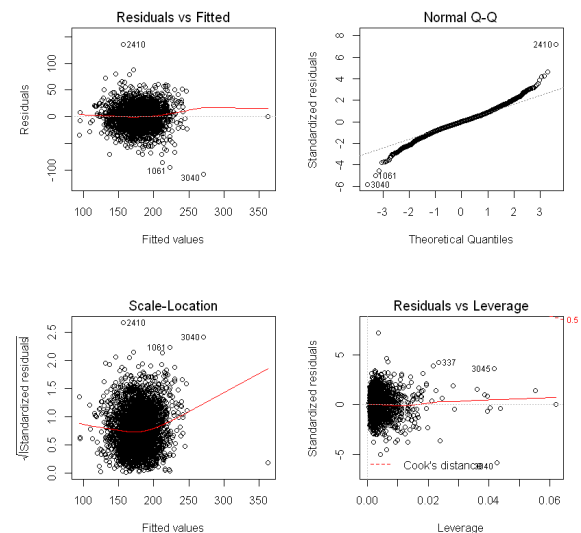
*Figure 12 Model 3*

The model above also shows that there is a relationship between the predictors and the response. In this model, notice that all



To see the model more clearly a visualization was done in and it is observed that we have some more outliers that needs to be handled.

The Outliers were all handled, and a new model was plotted. The next figure shows the performance of the new model.

```
Call:
lm(formula = deathRate ~ incidenceRate + povertyPercent + AvgHouseholdSize +
    PctHS18_24 + PctWhite + PctBachDeg25_Over + PctEmpPrivCoverage +
    PctOtherRace, data = cancerdata)

Residuals:
     Min       1Q   Median       3Q      Max
-115.739  -10.526    0.256   10.848  135.181

Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)        89.902583   7.389240  12.167  < 2e-16 ***
incidenceRate       0.212555   0.006394  33.243  < 2e-16 ***
povertyPercent      1.045378   0.092827  11.262  < 2e-16 ***
AvgHouseholdSize   -3.074254   1.613149  -1.906  0.05678 .
PctHS18_24          0.299996   0.043214   6.942 4.71e-12 ***
PctWhite           -0.124795   0.027744  -4.498 7.12e-06 ***
PctBachDeg25_Over  -1.609662   0.089559 -17.973  < 2e-16 ***
PctEmpPrivCoverage  0.171935   0.055763   3.083  0.00207 **
PctOtherRace       -0.805362   0.109294  -7.369 2.21e-13 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 19.16 on 3038 degrees of freedom
Multiple R-squared:  0.5245,	Adjusted R-squared:  0.5233
F-statistic: 418.9 on 8 and 3038 DF,  p-value: < 2.2e-16
```

*Figure 13 Model 4*

The model above shows that there is a relationship between the predictors and the response. But in this model notice that predictors have a very less significant statistical effects. The r squared ($R^2$) value shows that 52% of the predictors can explain the changes with a better F statistic in this multiple regression model.
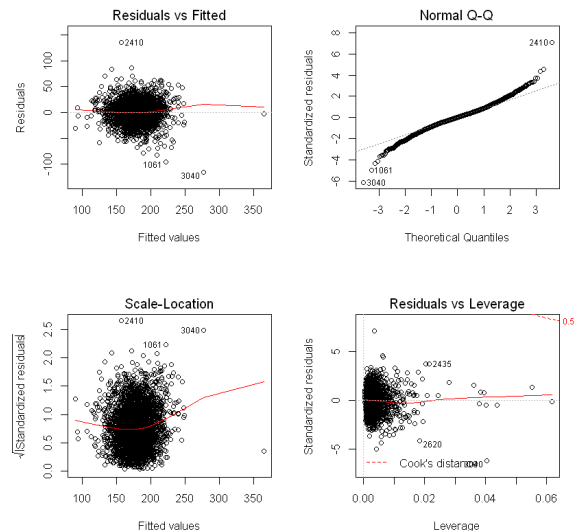


*Figure 14 Visualization of Model 4*

## 4.0 FINAL MODEL.

Model 5 was built by removing the less significant predictor.

```
Call:
lm(formula = deathRate ~ incidenceRate + povertyPercent + PctHS18_24 +
    PctWhite + PctBachDeg25_Over + PctEmpPrivCoverage + PctOtherRace,
    data = cancerdata)

Residuals:
    Min      1Q  Median      3Q     Max
-116.25  -10.46    0.25   10.85  135.04

Coefficients:
                     Estimate Std. Error t value Pr(>|t|)
(Intercept)         81.544163   5.949399  13.706  < 2e-16 ***
incidenceRate        0.214184   0.006339  33.786  < 2e-16 ***
povertyPercent       1.032691   0.092628  11.149  < 2e-16 ***
PctHS18_24           0.292278   0.043043   6.790 1.34e-11 ***
PctWhite            -0.110599   0.026737  -4.137 3.62e-05 ***
PctBachDeg25_Over   -1.594905   0.089262 -17.868  < 2e-16 ***
PctEmpPrivCoverage   0.149609   0.054542   2.743  0.00612 **
PctOtherRace        -0.867645   0.104338  -8.316  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 19.17 on 3039 degrees of freedom
Multiple R-squared:  0.524,	Adjusted R-squared:  0.5229
F-statistic: 477.9 on 7 and 3039 DF,  p-value: < 2.2e-16
```

Having removed the less significant variable, the model performed better. A better residual standard error was gotten a better F statistics
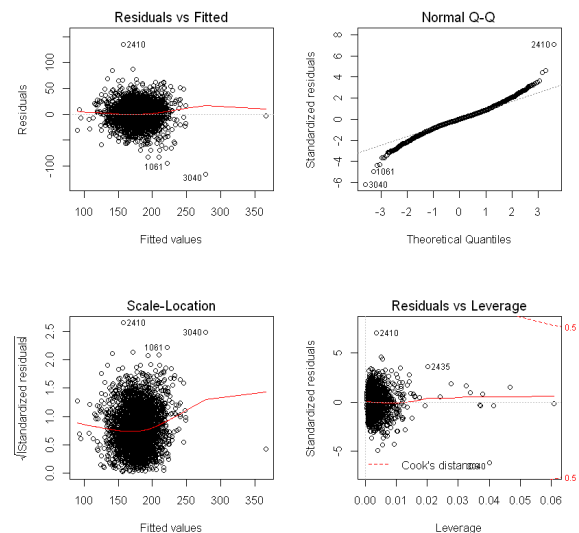
that show how fitted it is and a R-squared value of 52%.



*Figure 15 Visualization of the model*

In other to avoid overfitting, the outliers seen in the above visualized model would not be considered as all variables are highly significant predictors. For better clarity ANOVA was used to compare all five (5) models.



**n [84]:** anova(model1, model2, model3, model4, model5)

| Res.Df | RSS | Df | Sum of Sq | F | Pr(>F) |
|---|---|---|---|---|---|
| 3023 | 1075265 | NA | NA | NA | NA |
| 3035 | 1087694 | -12 | -12429.884 | 2.912116 | 5.034768e-04 |
| 3037 | 1091528 | -2 | -3833.722 | 5.389065 | 4.610215e-03 |
| 3038 | 1115370 | -1 | -23841.417 | 67.027790 | 3.903041e-16 |
| 3039 | 1116703 | -1 | -1333.403 | 3.748730 | 5.294075e-02 |

*Figure 16 ANOVA Result*

## 5.CONCLUSION

Having ran the test in 5 models, it is observed that the fifth model performed well with a R-squared value of 52% and F-statistics of 477 and a degree of freedom of -1 from the ANOVA comparison.

# REFERENCES

[1] C. N. Fru, T. Andrew, F. N. Cho, T. Tassang, and P. N. Fru, "Socio-economic Determinants Influencing Cervical Cancer Screening in Buea: A Cross-Sectional Study," *IJTDH*, pp. 14–22, Aug. 2020, doi: 10.9734/ijtdh/2020/v41i1130331.

[2] T. Akinyemiju, Q. Meng, and N. Vin-Raviv, "Race/ethnicity and socio-economic differences in colorectal cancer surgery outcomes: analysis of the nationwide inpatient sample," *BMC Cancer*, vol. 16, no. 1, p. 715, Dec. 2016, doi: 10.1186/s12885-016-2738-7.

[3] M. Mohebbi, R. Wolfe, D. Jolley, A. B. Forbes, M. Mahmoodi, and R. C. Burton, "The spatial distribution of esophageal and gastric cancer in Caspian region of Iran: An ecological analysis of diet and socio-economic influences," *Int J Health Geogr*, vol. 10, no. 1, p. 13, 2011, doi: 10.1186/1476-072X-10-13.