

Statistics for Data Analytics CA1

Multiple Linear Regression Analysis

Harshitha Poolakanda Somanna

Student ID: x22150366

Statistics for Data Analytics

National College of Ireland, Dublin

x22150366@student.ncirl.ie

Abstract— Over the last few decades, there is a prevalent increase in cancer around the world. In this report, we estimate the best suited Multiple Linear Regression model that can be used to investigate the significant relationship of socio economic factors such as median income, poverty percentage, average household size, private insurance coverage and others influencing the death rate for Cancer patients in the US. The dataset sourced has data for over 3000 cancer patients in the US counties with about 25 attributes. SPSS has been used for descriptive statistics, visualizations, modeling and diagnostics.

Keywords—independent variables; dependent variable; linear regression; normality test; linearity test; multicollinearity; Durbin-Watson test; Gauss-Markov assumptions.

I. INTRODUCTION

The aim of this analysis is to perform Multiple Linear Regression on the cancer dataset and to investigate socio-economic factors influencing the cancer mortality in the US using the best suited regression model. Multiple regression is a statistical technique used to examine the association between a dependent variable and multiple independent variables simultaneously. The independent variables whose values are known are used to predict the dependent variable in this technique and the formula is as follows:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon \quad [1]$$

where:

- Y is the dependent variable that we aim to estimate
- β_0 is the intercept
- $\beta_1, \beta_2, \dots, \beta_p$ are the regression coefficients corresponding to the independent variables X_1, X_2, \dots, X_p
- ε denotes the amount of variation in the dependent variable that cannot be accounted for by the independent variables.

The general foundations of multiple linear regression include the following:

- Assumptions: Ensure that the relationship between the variables is linear, errors are independent and normally distributed, and there is no multicollinearity.
- Model fitting: Estimate the coefficients of the independent variables using a statistical method that minimizes the sum of squared errors.
- Interpretation of coefficients: Coefficients indicate the change in the dependent variable for a one-unit change in the independent variable, holding all other variables constant.

- Model evaluation: Use metrics like R-squared, adjusted R-squared to evaluate the model's fit and ability to generalize to new data.

This report includes descriptive statistics and suitable visualizations to facilitate comprehension of the variables in the dataset. Also, the steps taken to construct the regression model leading to the final selection will be discussed along with details on the reasons for choosing the independent variables, the data transformation, outlier treatment, etc. Also, the procedures used to confirm the Gauss Markov assumptions for multiple regression, such as assessing normality, identifying outliers, and detecting multicollinearity among variables that may impact the precision of the model, are discussed for each model.

A. Data Description

The cancer.csv file with 3047 records and 25 attributes has been used as the source dataset. To predict the mortality rate in the cancer dataset, we can use the incidence rate and death rate columns as key factors in our regression model. In this context, the variables have been categorized into two groups: independent variables and dependent variables. The variable that we aim to examine and understand is the death rate, which is considered as the dependent variable. To investigate the death rate, we need to consider several independent variables. The socio economic variables and medical variables used in the dataset is shown in Fig.1.

deathRate
Population
incidenceRate
medIncome
povertyPercent
MedianAge
MedianAgeMale
MedianAgeFemale
AvgHouseholdSize
PctMarriedHouseholds
PctBachDeg18_24
PctBachDeg25_Over
PctUnemployed16_Over
PctPrivateCoverage
PctEmpPrivCoverage
PctPublicCoverage
PctPublicCoverageAlone

Table.1. List of Variables in the Cancer dataset

II. EXPLORATIVE DATA ANALYSIS

Explorative Data Analysis (EDA) is an essential first step to understand and summarize the key characteristics of data. It helps us to identify any potential outliers or anomalies, detecting missing values and to develop hypothesis for further testing. It typically involves various visual and

statistical techniques, such as scatter plots, histograms, and correlation matrices.

A. Descriptive Statistics

Descriptive statistics is carried to summarize and describe the key features of a dataset such as the mean, median, mode, range, standard deviation, and variance. It also identifies outliers, patterns and trends in data and also helps us to explore the relationship between different variables in the data.

The initial stage of data preprocessing involves verifying the absence of any missing, empty, or null values, which were not identified in the dataset. When loading the dataset into SPSS, any preceding or succeeding spaces in string values were eliminated.

Table 2. shows the descriptive statistics of all possible variables in the cancer dataset where the characteristics of a dataset, such as its central tendency, variability, and distribution are described.

Descriptive Statistics										
	N	Range	Minimum	Maximum	Mean	Std. Deviation	Variance	Skewness	Kurtosis	
Population	2047	10169485	827	10170262	102837.37	229059.224	1.0805E+11	14.2460	.044	237.437 .089
deathRate	2047	303.1	59.7	362.8	178.654	27.7515	772.146	.275	.044	1.355 .089
incidenceRate	2047	1005.6	201.3	1208.9	445.654	57.4566	3301.259	.751	.044	13.794 .089
medianIncome	2047	102995	22640	125635	47063.28	12040.091	144963787.33	1.408	.044	3.713 .089
povertyPercent	2047	44.2	3.2	47.4	16.878	6.4091	41.076	.931	.044	1.276 .089
MedianAge	2047	607.7	22.3	624.0	45.572	45.3845	2052.496	9.940	.044	199.816 .089
MedianAgeMales	2047	42.3	22.4	64.7	38.571	5.2268	27.311	1.52	.044	.676 .089
MedianAgeFemales	2047	43.4	22.3	65.7	42.145	5.2928	28.014	-.208	.044	.577 .089
AvgHouseholdSize	2047	2.11	1.86	3.97	2.5287	.24845	.062	1.287	.044	3.674 .089
PctMarriedHouseholds	2047	55.08280084	22.89248889	78.07539683	51.243872141	6.5728137943	43.202	-.522	.044	1.414 .089
PctWhite18_24	2047	64.1	0	84.1	18.224	8.0691	65.488	.973	.044	2.211 .089
PctHis18_24	2047	72.5	0	72.5	35.002	9.0807	82.260	.179	.044	.534 .089
PctBlack18_24	2047	51.8	0	51.8	6.158	4.5291	20.512	1.956	.044	9.139 .089
PctHis25_Over	2047	47.3	7.5	54.8	34.805	7.0349	49.490	-.334	.044	1.119 .089
PctBlack25_Over	2047	39.7	2.5	42.2	13.282	5.3648	28.103	1.095	.044	1.737 .089
PctHisover18_Over	2047	79.0	4	79.0	7.952	3.4524	11.919	.891	.044	2.287 .089
PctPublicCoverage	2047	79.0	22.3	92.3	64.305	10.8471	113.360	-.384	.044	-.854 .089
PctEmpPvtCoverage	2047	57.2	13.5	70.7	41.196	9.4477	89.259	.089	.044	-.302 .089
PctPublicCoverage	2047	53.9	11.2	65.1	36.253	7.8417	61.493	-.005	.044	-.089 .089
PctPublicCoverageNone	2047	44.0	2.6	46.8	19.240	6.1130	37.369	.471	.044	.362 .089
PctWhite	2047	85.80084480	16.18415500	100.00000000	63.6452461205	16.380525228	268.395	-1.881	.044	2.691 .089
PctBlack	2047	85.847788500	.000000000	85.847788500	9.1070770146	14.5345379322	211.253	2.258	.044	5.939 .089
PctAsian	2047	42.619424540	.000000000	42.619424540	1.2539649642	2.8162763927	6.814	7.418	.044	78.397 .089
PctOtherRace	2047	41.830251420	.000000000	41.830251420	1.9835238038	3.517101375	12.374	4.452	.044	35.537 .089

Table 2. Descriptive Statistics

B. Visualisation and discussion of out of range values

1) Residuals

Residual statistics as shown in Table 3. summarizes the difference between actual and predicted values. They help evaluate the model's goodness-of-fit, including mean, standard deviation, minimum and maximum, skewness, kurtosis, Durbin-Watson, and Cook's Distance. These statistics assist in identifying outliers, influential cases, and assumptions of the model.

Residuals Statistics ^a					
	Minimum	Maximum	Mean	Std. Deviation	N
Predicted Value	110.506	354.491	178.747	19.2679	2827
Std. Predicted Value	-3.542	9.121	.000	1.000	2827
Standard Error of Predicted Value	.465	6.119	1.059	.513	2827
Adjusted Predicted Value	110.771	353.714	178.749	19.2738	2827
Residual	-51.4206	55.7293	.0000	15.6062	2827
Std. Residual	-3.286	3.561	.000	.997	2827
Stud. Residual	-3.292	3.566	.000	1.000	2827
Deleted Residual	-51.5896	55.8610	-.0021	15.6865	2827
Stud. Deleted Residual	-3.297	3.573	.000	1.000	2827
Mahal. Distance	1.495	431.091	14.995	25.483	2827
Cook's Distance	.000	.009	.000	.001	2827
Centered Leverage Value	.001	.153	.005	.009	2827

a. Dependent Variable: deathRate

Table 3. Residual Statistics

2) Scatter Plot

A scatter plot as shown in Fig1 for multiple linear regression helps visualize the relationship between variables and detect outliers and influential points that can impact the regression results. Clusters around a straight line indicate a strong linear relationship.

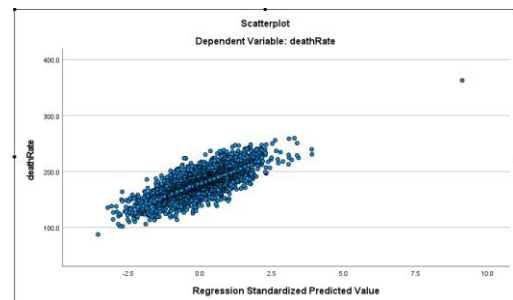


Fig 1. Scatter Plot

3) Histogram

Fig2. Shows a histogram plot that can identify distribution issues of independent variables in multiple linear regression by showing frequency distribution through grouped bins and count observations.

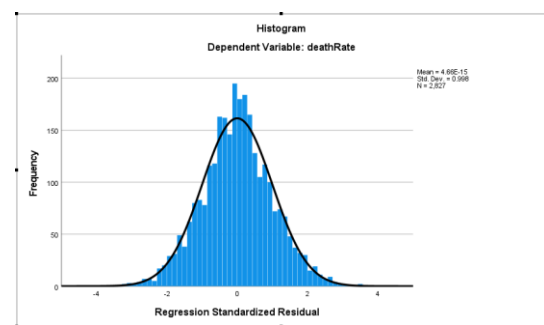


Fig 2. Histogram

4) Normal Probability plot

Normal plot as shown in Fig3. is useful for verifying the validity of model assumptions and the reliability of regression estimates in multiple linear regression by assessing the normality of the residuals[2]. A non-linear pattern in the normal plot of residuals indicates that the model assumptions have been violated and the regression estimates may be unreliable, necessitating a review of the model specification or alternative regression techniques.

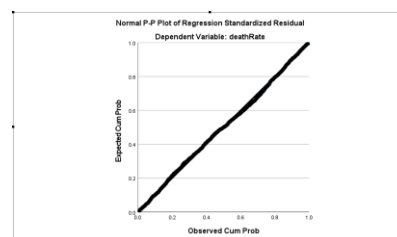
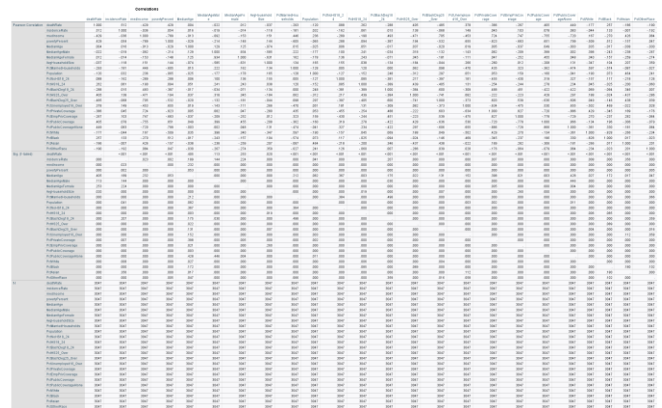


Fig 3. Normal P-P plot

5) Multicollinearity

Correlation matrix as displayed in Table4. depicts the correlation coefficients between variables in a dataset. In

Correlations										
	deathRate	Population	incidenceRate	medicines	povertyPercent	MedianAge	MedianAgeMale	MedianAgeFemale	AgeGroupSize15to64	AgeGroupSize65plus
Pearson Correlation	deathRate	-1.00	-120	512	-429	429	0.04	-0.02	0.01	-0.037
	Population	-120	1.000	0.02	238	-0.065	-0.025	-0.177	-0.178	-0.165
	incidenceRate	512	0.02	1.000	0.036	0.54	0.16	-0.019	-0.014	-0.118
	medicines	-429	238	-0.036	1.000	-0.789	0.13	-0.082	-0.153	-0.151
	povertyPercent	-429	-0.065	0.036	-0.789	1.000	-0.028	-0.144	-0.144	-0.144
	MedianAge	429	0.04	-0.025	0.13	-0.028	1.000	0.129	0.125	0.074
	MedianAgeMale	-0.022	-0.177	-0.019	-0.082	-0.144	0.129	1.000	0.934	-0.085
	MedianAgeFemale	-0.014	-0.178	0.014	-0.153	-0.144	0.125	0.934	1.000	-0.014
	AgeGroupSize15to64	-0.037	-0.165	-0.118	-0.151	-0.144	-0.074	-0.085	-0.031	1.000
	AgeGroupSize65plus	-0.033	-0.158	-0.181	-0.446	-0.055	0.15	0.222	0.62	1.000
PearsonR0181_24	deathRate	0.088	-0.127	-0.162	-0.289	0.088	0.01	0.100	0.136	-0.155
PearsonR0181_24	Population	-0.127	1.000	0.021	-0.180	-0.061	-0.051	-0.180	-0.243	-0.088
PearsonR0181_24	incidenceRate	-0.289	0.021	1.000	0.051	0.493	-0.387	-0.071	-0.034	-0.071
PearsonR0181_24	medicines	-0.446	-0.180	0.051	1.000	-0.789	0.01	-0.085	-0.153	-0.151
PearsonR0181_24	povertyPercent	-0.446	-0.180	0.051	-0.789	1.000	-0.028	-0.144	-0.144	-0.144
PearsonR0181_24	MedianAge	0.088	-0.127	-0.162	-0.289	0.088	0.01	0.100	0.136	-0.155
PearsonR0181_24	MedianAgeMale	-0.012	-0.178	-0.012	-0.082	-0.144	0.129	1.000	0.934	-0.085
PearsonR0181_24	MedianAgeFemale	-0.012	-0.178	0.012	-0.153	-0.144	0.125	0.934	1.000	-0.012
PearsonR0181_24	AgeGroupSize15to64	-0.037	-0.165	-0.118	-0.151	-0.144	-0.074	-0.085	-0.031	1.000
PearsonR0181_24	AgeGroupSize65plus	-0.033	-0.158	-0.181	-0.446	-0.055	0.15	0.222	0.62	1.000
PearsonR0181_24	deathRate	0.088	-0.127	-0.162	-0.289	0.088	0.01	0.100	0.136	-0.155
PearsonR0181_24	Population	-0.127	1.000	0.021	-0.180	-0.061	-0.051	-0.180	-0.243	-0.088
PearsonR0181_24	incidenceRate	-0.289	0.021	1.000	0.051	0.493	-0.387	-0.071	-0.034	-0.071
PearsonR0181_24	medicines	-0.446	-0.180	0.051	1.000	-0.789	0.01	-0.085	-0.153	-0.151
PearsonR0181_24	povertyPercent	-0.446	-0.180	0.051	-0.789	1.000	-0.028	-0.144	-0.144	-0.144
PearsonR0181_24	MedianAge	0.088	-0.127	-0.162	-0.289	0.088	0.01	0.100	0.136	-0.155
PearsonR0181_24	MedianAgeMale	-0.012	-0.178	-0.012	-0.082	-0.144	0.129	1.000	0.934	-0.085
PearsonR0181_24	MedianAgeFemale	-0.012	-0.178	0.012	-0.153	-0.144	0.125	0.934	1.000	-0.012
PearsonR0181_24	AgeGroupSize15to64	-0.037	-0.165	-0.118	-0.151	-0.144	-0.074	-0.085	-0.031	1.000
PearsonR0181_24	AgeGroupSize65plus	-0.033	-0.158	-0.181	-0.446	-0.055	0.15	0.222	0.62	1.000
PearsonR0181_24	deathRate	0.088	-0.127	-0.162	-0.289	0.088	0.01	0.100	0.136	-0.155
PearsonR0181_24	Population	-0.127	1.000	0.021	-0.180	-0.061	-0.051	-0.180	-0.243	-0.088
PearsonR0181_24	incidenceRate	-0.289	0.021	1.000	0.051	0.493	-0.387	-0.071	-0.034	-0.071
PearsonR0181_24	medicines	-0.446	-0.180	0.051	1.000	-0.789	0.01	-0.085	-0.153	-0.151
PearsonR0181_24	povertyPercent	-0.446	-0.180	0.051	-0.789					



From Fig 4., it is observed that the standardized residual value of -5.814 indicates that the corresponding observation has a much lower value than what is predicted by the model and may need to be further examined for any errors or anomalies. The histogram plot, scatter plot and Normal P-P plot reveals the distribution of the variables here and it can

be summarized that the dataset has outliers which must be further removed.

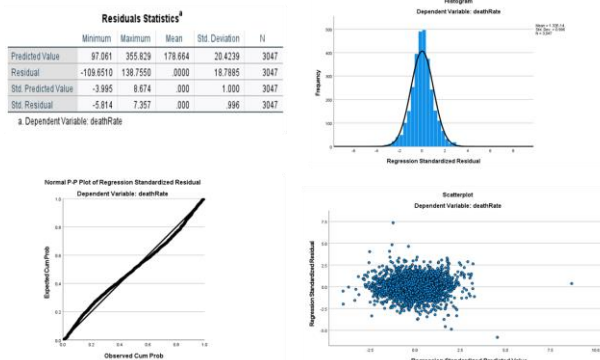


Fig 4. Residual and Normality graphs for Model 1

B. Model 2

For Model 2, it was necessary to eliminate variables that exhibited high correlation with each other, as well as those that had a Variance Inflation Factor (VIF) greater than 10 for all the 3047 rows of data.

1) Model Summary

The Table 9 shows that the model's R value of .729 indicates a strong positive correlation between the dependent variable and the independent variables. The R Square value of .532 suggests that the model can explain 53% of the variance in the outcome variable based on the selected independent variables. The adjusted R Square value of .529 indicates that 53% of the data points closely fit the linear regression curve. Additionally, the Durbin-Watson statistic of 1.961 suggests that the residuals exhibit good autocorrelation. Overall, considering the 19 socio-economic variables in the dataset, the model can be considered a moderately good fit for predicting deathrate.

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Change Statistics					Durbin-Watson
					R Square Change	F Change	df1	df2	Sig. F Change	
1	.729 ^a	.532	.529	19.0522	.532	180.824	19	3027	.000	1.961

a. Predictors: (Constant), PctEmpPrivCoverage, AvgHouseholdSize, MedianAge, incidenceRate, Population, PctHS18_24, PctBlack, PctOtherRace, PctNoHS18_24, PctAsian, PctBachDeg18_24, PctHS25_Over, PctUnemployed16_Over, MedianAgeMale, PctMarriedHouseholds, povertyPercent, PctBachDeg25_Over, medIncome, PctWhite

b. Dependent Variable: deathRate

Table 9. Model Summary for Model 2

2) ANOVA

The ANOVA analysis performed on Model 2, as shown in Table 10, evaluates how much the independent variables explain the variation in the dependent variable. The F-ratio of 180.8 and the corresponding p-value of 0 indicate that the independent variables have a statistically significant impact on the death rate outcome variable, with a much higher significance compared to the previous model.

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1247100.980	19	65636.894	180.824	.000 ^b
	Residual	1098764.895	3027	362.988		
	Total	2345865.875	3046			

Table 10. ANOVA for Model 2

3) Multicollinearity

Table 11 shows the high correlation coefficient of 0.7 between PctBachDeg25_Over, PctEmpPrivCoverage, and medIncome indicates a strong association among these variables, which could lead to redundancy. To avoid such redundancy, it is necessary to eliminate one of these variables. Additionally, the variables PctUnemployed16_Over and povertyPercent exhibit a correlation value of 0.66, which is indicative of a strong association. Furthermore, from Table 12, it is clear that PctBachDeg18_24 shows a correlation of 0.6 with PctBachDeg25_Over, which could suggest collinearity. Therefore, it is necessary to eliminate one of these variables to prevent collinearity issues.

Model		Unstandardized Coefficients		Standardized Coefficients		t	Sig.	Collinearity Statistics	
		B	Std. Error	Beta				Tolerance	VIF
1	(Constant)	122.790	12.474			9.844	<.001		
	PctEmpPrivCoverage	.041	.078	.014	.559	.576		.243	4.123
	AvgHouseholdSize	-7.417	2.423	-.066	-3.061	.002		.329	3.042
	MedianAge	-.004	.008	-.006	-.484	.628		.979	1.021
	incidenceRate	.205	.007	.425	31.219	<.001		.834	1.199
	Population	-7.955E-7	.000	-.009	-.642	.521		.716	1.396
	PctHS18_24	.292	.047	.095	6.168	<.001		.646	1.549
	PctBlack	-.116	.053	-.061	-2.177	.030		.197	5.085
	PctOtherRace	-.785	.121	-.099	-6.468	<.001		.655	1.528
	PctNoHS18_24	-.018	.053	-.005	-.337	.736		.644	1.553
	PctAsian	-.067	.084	-.006	-.364	.716		.514	1.946
	PctBachDeg18_24	.001	.106	.000	.009	.993		.521	1.918
	PctHS25_Over	.293	.094	.074	3.120	.002		.272	3.674
	PctUnemployed16_Over	.564	.151	.070	3.723	<.001		.436	2.295
	MedianAgeMale	-.422	.110	-.079	-3.849	<.001		.363	2.752
	PctMarriedHouseholds	-.161	.095	-.038	-1.703	.089		.307	3.256
	povertyPercent	.788	.134	.182	5.899	<.001		.162	6.158
	PctBachDeg25_Over	-1.414	.146	-.275	-9.697	<.001		.193	5.191
	medIncome	6.612E-5	.000	.029	.847	.397		.135	7.411
	PctWhite	-.154	.057	-.091	-2.681	.007		.135	7.428

a. Dependent Variable: deathRate

Table 11. Coefficient Matrix for Model 2

	AvgHouseholdSize	MedianAge	incidenceRate	Population	PctHS18_24	PctBlack	PctOtherRace	PctNoHS18_24	PctAsian	PctBachDeg18_24	PctHS25_Over	PctUnemployed16_Over	MedianAgeMale	PctMarriedHouseholds	povertyPercent	PctBachDeg25_Over	medIncome	PctWhite
PctEmpPrivCoverage	.158	.125	.012	.439	-.022	-.197	.257	.189	.168	.439	.084	.007	-.285	.089	.282	.405	.282	.405
AvgHouseholdSize	1.000	.002	.286	-.157	-.180	.075	.484	-.241	.149	.149	.022	.085	-.136	-.127	.152	.152	.152	.152
MedianAge	.002	1.000	-.036	-.019	-.044	.133	.087	-.162	.014	.014	-.119	-.101	-.191	-.142	.061	.106	.106	.106
incidenceRate	.439	.286	1.000	.002	-.047	-.272	.424	.084	.084	.084	.022	.085	-.136	-.127	.152	.152	.152	.152
Population	-.157	-.180	-.047	1.000	.006	-.243	.238	-.247	-.174	-.174	.139	.085	.222	.189	.241	.218	.218	.218
PctHS18_24	.149	.014	-.044	.006	1.000	-.024	.036	-.134	.104	.104	.002	.007	-.167	-.167	.040	.040	.040	.040
PctBlack	.087	.087	-.272	.424	-.024	1.000	.017	.023	.012	.012	.007	.004	.014	.117	.025	.024	.024	.024
PctOtherRace	.257	.189	.424	.084	.036	.017	1.000	.001	.001	.001	.007	.004	.007	.004	.001	.001	.001	.001
PctNoHS18_24	.168	.014	-.162	-.174	-.134	-.023	.031	1.000	.007	.007	.004	.007	.004	.007	.004	.007	.004	.007
PctAsian	.149	.014	-.044	.006	.104	.104	.002	.007	1.000	.002	.007	.004	.007	.004	.007	.004	.007	.004
PctBachDeg18_24	.007	.007	.004	.007	.002	.004	.007	.004	.007	1.000	.002	.007	.004	.007	.004	.007	.004	.007
PctHS25_Over	.084	.022	.085	.007	.002	.007	.004	.007	.004	.002	1.000	.004	.007	.004	.007	.004	.007	.004
PctUnemployed16_Over	.282	.152	.152	.218	.007	.007	.004	.007	.004	.007	.004	1.000	.004	.007	.004	.007	.004	.007
MedianAgeMale	-.285	-.136	-.127	.222	-.167	.117	.124	.124	.085	.085	.007	.004	1.000	.004	.007	.004	.007	.004
PctMarriedHouseholds	.089	.106	.061	.241	.040	.025	.020	.040	.040	.040	.007	.004	.007	1.000	.004	.007	.004	.007
povertyPercent	.405	.152	.152	.218	.007	.007	.004	.007	.004	.007	.004	.007	.004	.007	1.000	.004	.007	.004
PctBachDeg25_Over	-.285	-.136	-.127	.222	-.167	.117	.124	.124	.085	.085	.007	.004	.007	.004	.007	1.000	.004	.007
medIncome	.282	.152	.152	.218	.007	.007	.004	.007	.004	.007	.004	.007	.004	.007	.004	.007	1.000	.004
PctWhite	.405	.152	.152	.218	.007	.007	.004	.007	.004	.007	.004	.007	.004	.007	.004	.007	.004	1.000

Table 12. Pearson's Correlation for Model 2

4) Residuals

According to the results obtained in Fig 5, the standardized residual value of -6.069 suggests that the residual value deviates 6.069 standard deviations from the expected value, indicating the presence of an extreme outlier. Additionally, the analysis of the histogram, scatter plot, and Normality plot displayed in Figure 5 also indicates the existence of a significant outlier in the dataset, even after eliminating variables with a strong relationship and high VIF value.

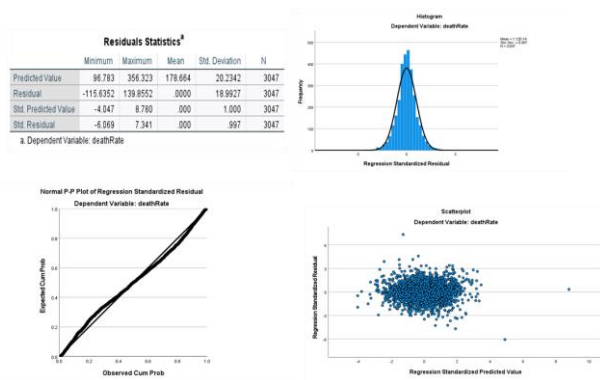


Fig 5. Residual and Normality graphs for Model 2

C. Model 3

In the current iteration, a variable from the previous model's set that showed strong correlation was removed, and outliers were eliminated using Cook's distance method.

1) Model Summary

The analysis utilizing Cook's method resulted in the identification of 220 rows of data as outliers, which were subsequently removed. The model was created on 2827 rows of data. Following this data cleansing, a new model was constructed using the remaining data.

Model Summary ^a										
Change Statistics										
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	R Square Change	F Change	df1	df2	Sig. F Change	Cohen's d
1	.777 ^a	.603	.562	15.6478	.004	235.656	1	3611	.000	1.885

a. Predictors: (Constant), PctMedicare, PctMedicaid, PctHIV_24, IncomeRate, PctHIV_18_24, Population, AgriHouseholdSize, PctOtherRace, medIncome, MedianAgeMale, PctAsian, PctH25_Over, PctMarriedHouseholds, povertyPercent, PctBlack

b. Dependent Variable: deathRate

Table 13. Model Summary

The newly constructed model as summarized in Table 13 shows an R value of .775 and an R Square value of .604, indicating that the model can accurately predict the outcome by 60.4% based on the selected independent variables. Additionally, the adjusted R Square value of .602 suggests that the data points conform well to the linear regression curve. The Durbin-Watson value of 1.885 indicates a strong autocorrelation between the residuals.

2) ANOVA

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1049160.520	15	69944.035	285.656	.000 ^b
	Residual	688283.853	2811	244.854		
	Total	1737444.373	2826			

a. Dependent Variable: deathRate

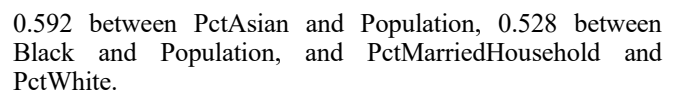
b. Predictors: (Constant), PctWhite, MedianAge, PctHS18_24, incidenceRate, PctHS18_24, Population, AvgHouseholdSize, PctOtherRace, medianIncome, MedianAgeMale, PctAsian, PctHS_Over35, PctMarriedHouseholds, povertyPercent, PctBlack

Table 14. ANOVA for Model 3

In this case, from Table 14, the higher F-ratio value of 285.6 is better at explaining the variation in the dependent variable and is better than the other two models.

3) Multicollinearity

The model exhibited a strong degree of colinearity, with the majority of variables demonstrating linearity below 5 as shown in Table 15. Despite this, there was still a notable correlation between certain variables, such as a correlation of



	Covariates										Explanatory Variables										Outcomes													
	Age	Gender	Marital	Religion	Ethnicity	Education	Income	Health	Employment	Family Size	Age	Gender	Marital	Religion	Ethnicity	Education	Income	Health	Employment	Family Size	Age	Gender	Marital	Religion	Ethnicity	Education	Income	Health	Employment	Family Size				
Physical Health	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2				
	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00				
	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1				
	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09				
	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00				
	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00				
	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00				
	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00				
	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00				
	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00				
Mental Health	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2				
	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00				
	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1				
	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09				
	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00				
	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00				
	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00				
	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00				
	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00				
	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00				
Family Size	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2	Age	Male	Married	Protestant	White	High School	Low	Good	Unemployed	2				
	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00	1980	0.84	0.17	0.22	0.36	0.29	0.17	0.17	0.47	0.00				
	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1	Gender	Female	Unmarried	Catholic	Black	College	High	Excellent	Employed	1				
	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09	Marital	0.17	0.13	0.09	0.13	0.09	0.09	0.09	0.09	0.09				
	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	Religion	0.22	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00				
	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	Ethnicity	0.36	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00				
	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	Education	0.29	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00				
	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00	Income	0.17	0.17	0.47	0.00	0.00	0.00	0.00	0.00	0.00				
	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Health	0.47	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00				
	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	Employment	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00				

Table 15. Pearson's Correlation for Model 3

Coefficients ^a								
Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	Collinearity Statistics	
		B	Std. Error	Beta			Tolerance	VIF
1	(Constant)	51.414	10.094		5.093	<.001		
	PctHS25_Ovr	.809	.062	.226	13.011	<.001	.466	2.144
	MedianAge	-.003	.009	-.005	-.392	.695	.981	1.020
	PctBlack	-.102	.058	-.058	-1.776	.076	.133	7.530
	IncidenceRate	.209	.006	.448	34.997	<.001	.861	1.162
	AvgHouseholdSize	5.435	2.239	.049	2.427	.015	.350	2.957
	PctNoHS18_24	.034	.045	.010	.749	.454	.736	1.359
	Population	-1.733E-6	.000	-.018	-1.150	.250	.603	1.659
	PctHS18_24	.373	.040	.131	9.224	<.001	.700	1.429
	PctOtherRace	-.968	.126	-.118	-7.666	<.001	.595	1.681
	povertyPercent	1.204	.118	.288	10.161	<.001	.176	5.692
	PctAsian	-.084	.206	-.007	-.406	.685	.424	2.357
	MedianAgeMale	-.262	.087	-.053	-3.019	.003	.454	2.203
	PctMarriedHouseholds	-.174	.085	-.044	-2.047	.041	.305	3.275
	medIncome	.000	.000	-.067	-2.314	.021	.170	5.895
	PctWhite	-.145	.061	-.089	-2.396	.017	.103	9.721

Table 16. Coefficient Matrix for Model 3

4) *Residual Statistics*

According to Figure 6, the regression model with a standardized residual of -3.286 performed better than the other two models. The histogram and scatter plot also indicate a significant reduction in the number of outliers in this model. Furthermore, the normal P-P plot shows that this model is almost aligned with the linear line, indicating a good fit to the normal distribution.

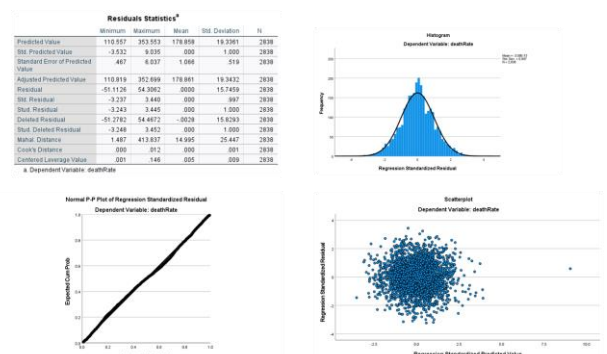


Fig 6. Residual and Normality graphs for Model 3

IV. MODEL SUMMARY

After analyzing Model 1, Model 2, and Model 3 in a regression analysis, it was determined that Model 3, which included multiple independent variables such as deathRate, Population, incidenceRate, medIncome, povertyPercent,

MedianAge, MedianAgeMale, PctWhite, PctBlack, PctAsian, PctOtherRace, AvgHouseholdSize, PctMarriedHouseholds, PctNoHS18_24, PctHS18_24, and PctHS25_over, outperformed the other two models in predicting the death rate.

To improve the accuracy of the analysis, outliers in these variables were removed based on a threshold for Cook's distance of 0.00131, which resulted in the elimination of 220 rows. Cook's distance is a measure of the influence of individual observations on a regression model, and a threshold value can be set to identify influential observations.

The removal of these outliers allowed for a more precise analysis of the relationship between the independent variables and the death rate. Based on these results, it can be concluded that Model 3 is the most effective in predicting the death rate, and the inclusion of multiple independent variables can provide a more comprehensive understanding of the factors that contribute to the outcome variable.

A) Gauss Markov Theorem

The Gauss-Markov theorem asserts that if the errors in the multiple linear regression model are normally distributed, have zero mean, constant variance (homoscedasticity), and are uncorrelated with each other, then the ordinary least squares (OLS) estimator is the best linear unbiased estimator (BLUE) [3].

1. Multicollinearity

The analysis indicates that there is no high degree of correlation between the independent variables in all three models. However, after transforming the model and conducting further analysis, Model 3 is found to be the best fitting model, as it has a correlation coefficient less than 6 and a Variance Inflation Factor (VIF) value less than 10. These findings suggest that Model 3 has the lowest level of multicollinearity among the three models and is therefore the most appropriate for predicting the outcome variable.

2. Independence of Error

Independence of errors is one of the key assumptions of the Gauss-Markov theorem in multiple regression. It requires that the errors in the model are independent of one another. Durbin-Watson test was done on all three models and all of them exhibited the results < 1.978 indicating acceptable levels of autocorrelation where Model 3 outperformed the other models by having result 1.885.

3. Homoscedasticity

Homoscedasticity is when the variance of residuals in a regression model is constant across all independent variable values. It's tested using a scatter plot to check for patterns that may indicate heteroscedasticity.

Model 1, 2 and 3 has a clustered pattern in the scatter plot which is more desirable than a funnel shape pattern, as it suggests that there are clear relationships between the independent variables and the dependent variable. However, in Model 3, the cluster is lying in a line with very few

outliers, it suggests that there is a strong linear relationship between the independent variable(s) and the dependent variable. This is a desirable pattern because it indicates that the model is able to explain a large proportion of the variation in the dependent variable based on the independent variable(s).

4. Normality

The normal probability plot (P-P plot) is a graphical tool that helps to assess the normality assumption of the Gauss-Markov theorem in multiple regression. It plots the standardized residuals against the expected values of a normal distribution, with a roughly straight line indicating that the residuals are normally distributed.

In this scenario, the analysis of all three models showed that the normal probability plot for Model 3 exhibited a roughly straight line with points evenly distributed along it. This suggests that the residuals for Model 3 are normally distributed and the normality assumption of the Gauss-Markov theorem is satisfied. Additionally, this finding indicates that Model 3 is the best fitting model among the three, as it satisfies the normality assumption.

5. Linearity

To satisfy linearity, the relationship between the dependent variable and the independent variables must be linear. In this scenario, it seems that all three models have shown a significant degree of linearity. However, Model 3 is the best fitting model in terms of linearity, which makes it a better choice for making predictions or drawing conclusions about the relationships between the variables.

V. CONCLUSION

The study employs a multiple regression analysis to examine how the incidence rate and a set of socio-economic variables impact the death rate of individuals diagnosed with cancer. My work has shown that the models have a respectable level of R square or adjusted R square of 53%. However, the adjusted R square value increases to 60% when outliers and a few correlated independent variables are removed from the dataset. The study identifies 11 variables that can be used to predict the death rate in the cancer dataset with relatively high accuracy. These variables include Population, incidenceRate, medIncome, povertyPercent, MedianAge, MedianAgeMale, PctWhite, PctOtherRace, AvgHouseholdSize, PctNoHS18_24, PctHS18_24, and PctHS25_over.

REFERENCES

- [1] M.H. Kutner, C.J. Nachtsheim, J. Neter, and W. Li, Applied Linear Statistical Models, McGraw-Hill, 2005.
- [2] G. Lind, W. Marchal, and S. Wathen, "Statistical Techniques in Business and Economics," 15th ed.
- [3] J. G. Dai and Z. Niu, "Gauss-Markov Theorem," in Encyclopedia of Biophysics, G. C. K. Roberts, Ed. Berlin, Germany: Springer, 2013, pp. 893-895. DOI: 10.1007/978-3-642-16712-6_740.