

On The Association Between Dietary Factors and Mortality Due to Stroke and Diabetes:

A Retrospective Examination of Western and Asian High-Income Countries.

By

Faaiza Castellanos, B.A. (Economics), University of California at Berkeley, 2000 and
Jill DeWitt, M.S (Secondary Math Education), Grand Valley State University, 2006

Advisors:

Dr. Jossy Uvah

Dr. Achraf Cohen

A Graduate Capstone Project
In Partial Fulfillment of the Degree of
Master of Science in Mathematical Sciences
University of West Florida

Contents

List of Figures	iv
List of Tables	iv
Abstract	1
1 Introduction	2
1.1 Statement of Problem	2
1.2 Relevance of Problem	2
1.3 Literature Review	2
1.4 Limitations	4
2 Dietary Factors and Mortality	5
2.1 The Data Set	5
2.2 Principle Components Analysis	6
2.3 Methodology: Justification, Assumptions and Development of the Model	11
2.4 Testing and Analysis	13
2.4.1 Model Formulation for Stroke Rate	13
2.4.2 Model Formulation for Diabetes Rate	18
3 Conclusions	22
3.1 Interpretation and Summary	22
3.2 Suggestions for Further Study	22
References	24
Appendixes	25
Appendix A: Full Data Set	25
Appendix B: R Code	28

List of Figures

1	PCA Biplot for Year 1990.	8
2	PCA Biplot for Year 2000.	9
3	PCA Biplot for Year 2010.	10
4	Histogram for Stroke rate	14
5	95% Confidence Interval for Lamda, Stroke Rate.	14
6	Diagnostic Plot for Log Model, Response Variable = Stroke.	15
7	Histogram for Stroke Rate after <i>log</i> transformation.	15
8	Histogram for diabetes rate.	18
9	95% Confidence Interval for Lamda, Diabetes Rate.	19
10	Diagnostic Plot for Box-Cox Model, Response Variable = Diabetes.	19
11	Histogram for Diabetes Rate after <i>Box-Cox</i> transformation.	20

List of Tables

1	Summary Statistics for Stroke Rate	13
2	Linear Regression results for Stroke Rate with Log transformation	16
3	Comparison of full and final models for Stroke Rate	16
4	Final linear regression model results for Stroke Rate, with Log Transformation	17
5	Summary Statistics for diabetes Rate	18
6	Linear Regression results for Diabetes Rate with Box-Cox transformation	20
7	Comparison of full and final models for Diabetes Rate	20
8	Final linear regression model results for Diabetes Rate, with Box-Cox Transformation	21

Abstract

A Retrospective Look at the Association Between Dietary Factors and Mortality due to Stroke and Diabetes in Western and Asian High-Income Countries

The purpose of this research paper was to explore the relationship between diet, specifically the consumption of fruit, non-starchy vegetables, bean and legumes, nuts and seeds, unprocessed red meats, sugar-sweetened beverages, fruit juices, milk, protein, calcium, and protein, and rate of stroke and diabetes in Western, high-income countries. Our study took a retrospective approach by looking at data from three years-1990, 2000, and 2010. We used Principal Component Analysis and Multiple Linear Regression to determine which foods or nutrients had a significant effect on stroke or diabetes rates. Gender and year were also considered as factors in the study. The study found that in general, only a few nutrients had a positive or negative impact on disease rates.

1 Introduction

1.1 Statement of Problem

Diet and how it pertains to the development and progression of chronic disease is a growing topic of interest. The increase in prevalence of deadly diseases, such as diabetes and stroke, has led many to examine their diet in a more profound manner, searching for clues that can prevent, halt the progression of, or even reverse such conditions. In response to the escalating interest in diet and disease, many doctors and self-help gurus are promoting a variety of diets, such as keto, paleo, low-fat, low-carb, and many more. Many ancient civilizations have also promoted a variety of foods based on culture and environment. It seems that nutritional intake and its impact on diseases has been a major concern for most people for thousands of years.

With the variety of recommendations and guidelines circulating, the question remains-What foods or nutrients are linked with some of the main killers in our society, namely diabetes and stroke? In order to study diet, it is important to examine a variety of foods and food groups, rather than just focusing on one or a few. In this paper, we investigate a variety of nutrients intake-fruit, non-starchy vegetables, bean and legumes, nuts and seeds, unprocessed red meats, sugar-sweetened beverages, fruit juices, milk, protein, calcium, and protein, and their link to stroke and diabetes/kidney disease. In addition, we restrict our study to Western and Asian high-income countries in an effort to control factors such as high poverty and low healthcare that contribute to lack of resources and education to buy and choose healthy foods.

1.2 Relevance of Problem

The study of food and its link to diabetes and stroke is of paramount importance to the World Health Organization (WHO) [2]. Every year societies spend a large amount of money on healthcare in order to treat diabetes and stroke. Unfortunately, many individuals also die from these diseases, leading to familial emotional turmoil as well as loss of economic contribution from the deceased individual. If such diseases can be prevented, reversed, and even treated with diet, individuals could decrease healthcare cost, benefit the economy, and be socially and emotionally healthier. Diet is a daily part of everyone's lives, so modifications to diet are available to almost all at a much lower price than medication or surgery.

1.3 Literature Review

The importance of diet in preventing death due to diabetes is a major concern of the World Health Organization (WHO) [2], with valid concern. According to the WHO, the number of people with diabetes rose from 108 million in 1980 to 422 million in 2014 and between 2000 and 2016, there was a 5% increase in premature mortality from diabetes. Diabetes is a disease that then can lead to a myriad of other deadly diseases, such as kidney failure, heart attacks, and stroke. Overall, the WHO estimates that diabetes was the seventh leading cause of death in 2016. However, they do stress the importance of a healthy diet in order to prevent and even reverse diabetes, recommending avoiding sugar and saturated fats. The American

Diabetes Association [3] recommends to fill half your plate with non-starchy vegetables and to also include fruits, lean meats and plant-based sources of protein, less added sugar, and less processed foods.

Strokes are another major cause of deaths worldwide, and the WHO [8] estimates that 15 million people suffer from a stroke worldwide annually, of which 5 million die and another 5 million remain permanently disabled. The major cause of stroke is high blood pressure, which is often associated with diet. The Centers for Disease Control and Prevention (CDC) [7] states that up to 80% of strokes can be prevented through healthy lifestyle changes, and diet is a key component of those changes. According to the CDC, eating plenty of fresh fruits and vegetables, foods low in saturated fats, trans fat, and cholesterol, and foods high in fiber can prevent high cholesterol, which in turn lowers your chances of suffering a stroke. Also, limiting sodium can lower blood pressure, which as indicated before, increases the likelihood of a stroke.

The WHO in Europe [1] provides even further research on the connection between diet and disease and offers recommendations for prevention. The importance of preventing disease with diet is a major issue of study and increasing prevalence of these diseases places a huge stress on the healthcare system, impacting national economies and health service budgets negatively. With correct knowledge and programs to education and support populations in implementing healthy dietary guidelines, many countries can increase longevity, improve their economies, de-stress their healthcare systems, and improve mental health. [15]

To study our data in terms of descriptive statistics, we used Principal Component Analysis, or PCA. PCA [9] is a technique for feature extraction, meaning that it combines our variables in a way that we can omit the least influential variables while retaining their most valuable parts. There are three key assumptions that are behind PCA [16] and can indicate when PCA would not provide a strong analysis. The assumptions are linearity, large variances have important structure, and the principal components are orthogonal. The linearity assumption organizes the analysis as a change of basis. As for variances, principal components with larger variances are indicative of structure, while those with smaller variances just represent noise, which is irrelevant to the strength of the signal being analyzed. The assumption that principal components are orthogonal is suggestive that PCA can be solved using linear algebra decomposition techniques.

We then used Multiple Linear Regression with a Box-Cox transformation to model the data. Halinski and Feldt [11] provide a framework for choosing the best procedure to pursue while keeping two goals in mind. First, the best model should produce an equation that yields the best predictions for the population. Second, the best model should contain an optimal number of explanatory variables. We began our model-building using a general linear model using a Gamma, then Poisson distribution, but ultimately settled on a normal multiple linear regression model with a transformation of the response variable. In general, the multiple regression models in this study strive to balance accuracy and parsimony, that is, the models should accurately describe both the systematic and random components and be as simple as possible. There are four basic assumptions that must be met to use multiple linear regression analysis [12]:

1. There must be a linear relationship between the response variable and the explanatory variables.

2. The model error term (referred to as the model “residuals”) must be normally distributed. Residuals can be thought of as the information not explained by the model. A residual plot and QQ plot can be used to determine if this assumption is met.
3. There should not be any multicollinearity. This means that explanatory variables are not correlated with each other. This assumption can be tested by examining the Variance Inflation Factor.
4. Homoscedasticity—This means that the variance of the error terms are consistent across the explanatory variables. A Scale-Location plot can be used to determine if this assumption is met.

The regression analysis in this study addresses each of these assumptions for each proposed multiple linear regression model. Nested models were compared using the Akaike’s Information Criterion (AIC) and Bayesian Information Criterion (BIC). Our goal was to consider these two measures together and to choose models that were favored by both criterion, as argued by Kuha [13].

1.4 Limitations

There were many limitations of the various methods of data collection and analysis. We collected data from the Global Dietary Database, or GDD, and they state that they collect data on dietary habits via surveys, which then are based on the volunteer’s bias. Some volunteers might not remember their exact diet, others might overstate or understate their food consumption, and yet others might just not wish to disclose honest information. Of course, the GDD uses a complex food description and classification system to address the issue of variation of description of diet, but this does not address the issue of volunteer dishonesty or lack of memory recall.

In addition, the data collected was based on observational studies, where the volunteers were not affected in any manner and were just questioned on their diet. The data that was studied could easily indicate correlation, but not necessarily causation. A more promising study would be an experimental study on diet, where volunteers would be divided into an experimental and control group. Each nutrient could be then studied separately to see how it affects development of diabetes and stroke. Of course, such a study would be very time-consuming, as human lifespan can last upwards of 100 years. In addition, to study a specific nutrient in such a manner would require that the experimental group consume that nutrient consistently for years on end. Certain foods being studied, such as sugar-sweetened beverages, are commonly considered to be unhealthy, and to require a group of volunteers to consume the beverages consistently for years and years would be unethical.

The data collected was international, indicating a large variety of genetics, spices, culinary techniques, exercise, and social behavior that play a part in disease. Genetics are a factor that can affect someone’s chance of developing diabetes or a stroke, as well as influence how sensitive they are to certain food items. The use of spices may also affect how foods react in the body and lead to the development of disease. Culinary techniques, such as deep frying, grilling,

boiling, and sauteing, can also alter foods and contribute to diabetes and strokes. Other than dietary factors, exercise is a main agent in health, affecting diet metabolism and also disease. Lastly, social behavior can lead to change in emotion and eating habits which can greatly affect disease. There is much research on emotions, neurotransmitters, and disease, and even more to study.

2 Dietary Factors and Mortality

2.1 The Data Set

We obtained the data regarding nutrition from The Global Dietary Database (GDD) [4], which is a project of the Gerald J. and Dorothy R. Friedman School of Nutrition Science and Policy at Tufts University. The program is also supported by the Bill and Melinda Gates Foundation, which is actively engaged in programs to support public policy. The goal of the program is to understand and improve diet through data collection, analysis, and recommendation, thus leading to public policies that aim to prevent disease and improve healthcare. GDD collects its data from [6] items, there is a large variation on their descriptions. In order to address this issue, the GDD applied FoodEx2, which is a complex food description and classification system developed by the European Food Safety Authority. This system allows GDD to standardize the global dietary intake, thus leading to data which is more reliable.

We obtained the data regarding disease from the Global Health Data Exchange (GHDx) [5]. The GHDx is a data catalog created and supported by an independent global health research center at the University of Washington. The population figures used to calculate the disease rates are estimated based on World Population Prospects: 2015 Revision, from the United Nations Population Division and disease mortality figures are obtained from the WHO Human Mortality Database.

We are most interested in discovering which nutrients had a positive or negative impact on disease rates and whether gender or year was a significant factor. The data set consisted of the following variables:

- **Country**-Each row contained data from one of twenty-eight countries. All countries were Western or Asian high-income countries. A full list of the countries can be found in Appendix A.
- **Gender** was a coded continuous factor where Female = (1) and Male = (2)
- **Nutrient Intake**- Eleven nutrients were included in the data set- fruit, non-starchy vegetables, bean and legumes, nuts and seeds, unprocessed red meats, sugar-sweetened beverages, fruit juices, milk, protein, calcium, and protein. All measurements were in average grams per day, except calcium and potassium, which were measured in average milligrams per day.
- **Disease Rates**- Rates for stroke and diabetes for each country were measured as counts per 100,000 people.

- **Year-** Data from each country was included for three different years, 1990, 2000, and 2010. The year was a coded continuous factor where 1990= (1), 2000= (2), and 2010= (3).

In summary, each of the 28 countries has 6 rows of data, a row for each gender and for each of the three years. We gathered the data from two online sources. The Global Dietary Database [4] provided the nutrient data, and the Global Health Data Exchange [5] provided the disease rate data. We then used the statistical computing and programming language R to read the data, and for all statistical graphics and regression analysis. The R code can be viewed in Appendix B.

2.2 Principle Components Analysis

In an attempt to factor a variety of aspects of diet into our study, we obtained data for 11 nutrients. This is a substantial amount of data, so we utilized Principal Component Analysis (PCA) [16] to identify which variables impact disease the most. We utilized PCA as a purely descriptive statistics method, thus we used it to identify key foods and nutrients for the years 1990, 2000, and 2010, rather than to omit variables in our analysis. Biplots were created, using R, to help visualize the principle components and to see if any clusters of countries are revealed.

We first examine the PCA biplot for 1990. A PCA biplot [14] shows the PCA score plot and the loading plot, where the PCA score plot displays the PCA scores and the loading plot portrays how strongly each of the variables impacts a principal component. The PCA Biplots for all three years are shown at the end of this section (Figures 1-3).

As seen in Figure 1, the PCA score plots are the individual countries where the data was gathered from and the PCA loading plot vectors are the food variables. Examining the PCA scores, clusters of countries represent countries that impact the two PCAs in a similar manner, meaning they exhibit similar characteristics. The country represented by 43 and 44 is the Republic of Korea (male and female), so this country impacts PCA 1 strongly but PCA 2 negligibly. The countries represented by 15, 16, 23, and 24 are Finland and Iceland. These countries impact both PCA 1 and PCA 2. Both of these countries are high-income European nations, and they have similarities in their food intake. Examining the loading plot, vectors that are close to each other and have small angle between them are positively correlated. Non-starchy vegetables and beans and legumes are positively correlated, and so are sugar-sweetened beverages and fruit juice. If vectors meet at a 90° angle, they are not correlated, such as nuts and seeds and milk, as well as fruits juice and non-starchy vegetables. Lastly, if two vectors diverge at a large angle, such as 180°, they are negatively correlated. Non-starchy vegetables, beans and legumes, and fruit are negatively correlated with milk.

The PCA biplot for 2000 (Figure 2) shows us similar information. Again, the Republic of Korea was a cluster by itself, along with a cluster shown for Finland and Iceland. Non-starchy vegetables and beans and legumes are again positively correlated, and so are sugar-sweetened beverages and fruit juice. Nuts and seeds and milk, as well as fruits and non-starchy vegetables are not correlated. Again, non-starchy vegetables, beans and legumes, and fruit are negatively correlated with milk.

Lastly, the PCA biplot for 2010 (Figure 3) shows similarities to the plots for 1990 and 2000 but does vary. The country represented by 35 and 36, the Netherlands, is a cluster by itself. The countries 43, 44, and 45 are also clustered, referring to the Republic of Korea and Singapore. As for the loading plot, non-starchy vegetables and beans, as well as sugar-sweetened beverages, fruit juice, and nuts and seeds are positively correlated. Milk, protein, and potassium are also positively correlated. Protein and fruit juice exhibit no correlation, as well as fruit juice and non-starchy vegetables and beans. Protein and milk are negatively correlated with beans and non-starchy vegetables. The PCA biplots provide us invaluable information on the relationships between the variables being studied, specifically the countries and the food intake.

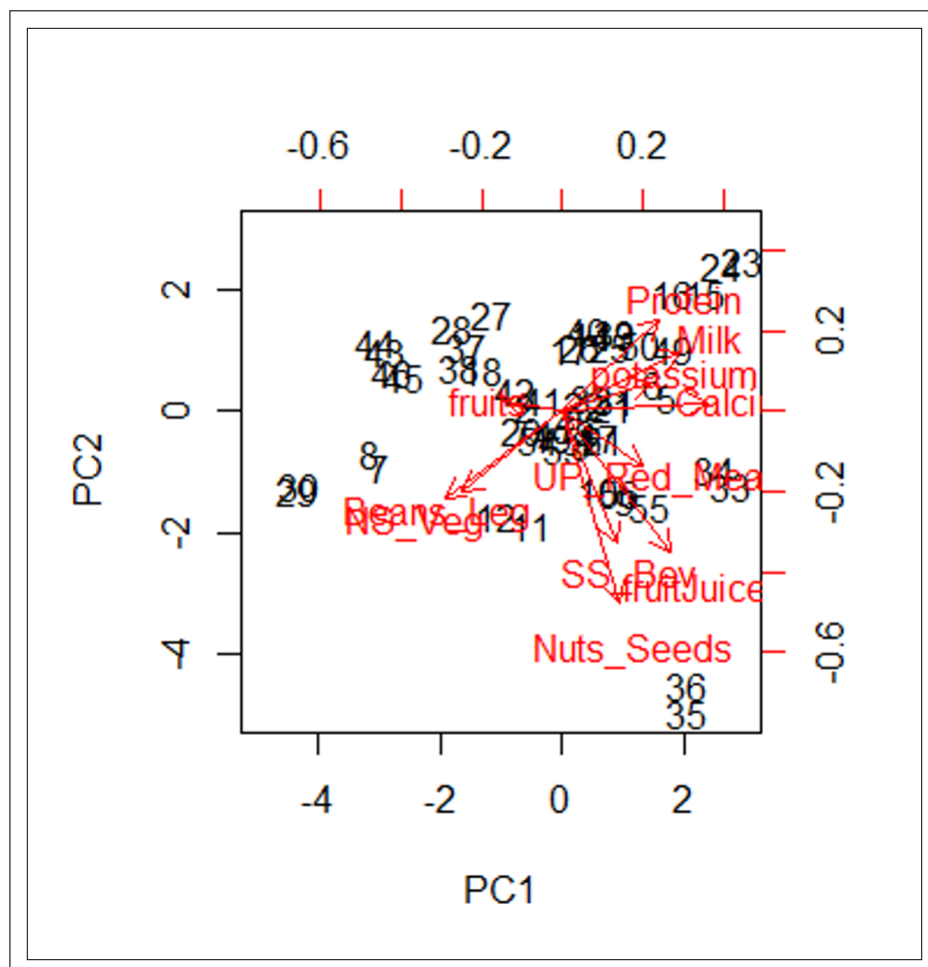


Figure 1: PCA Biplot for Year 1990.

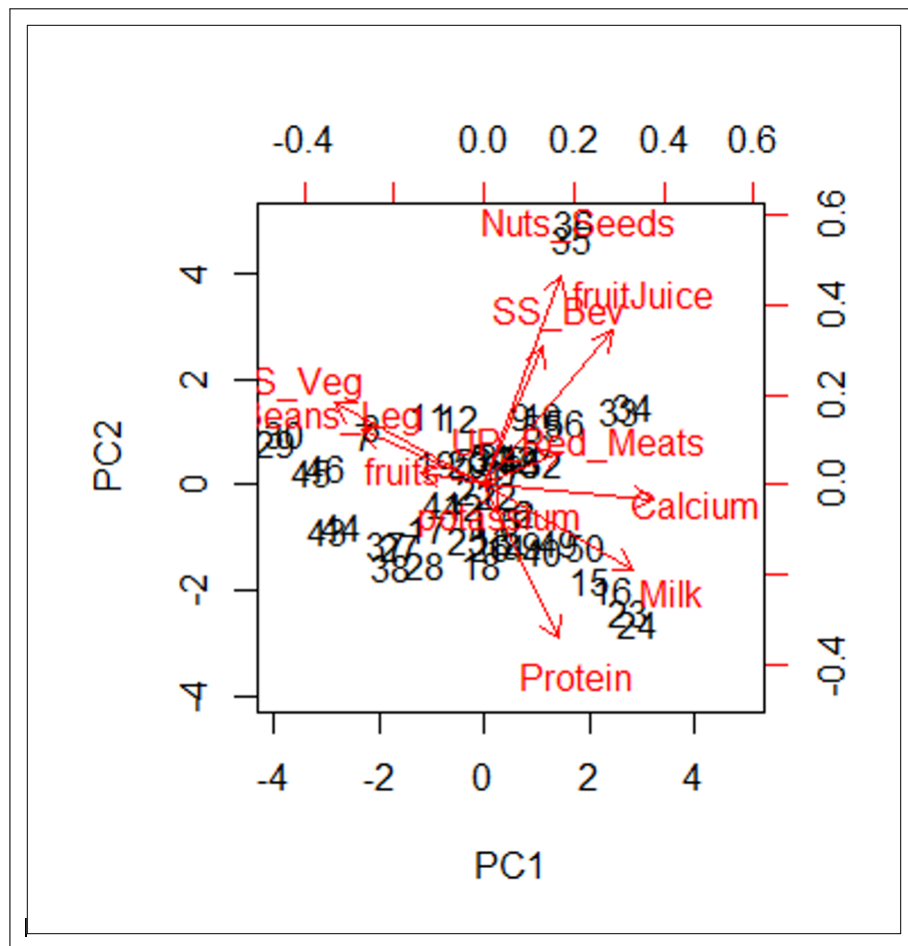


Figure 2: PCA Biplot for Year 2000.

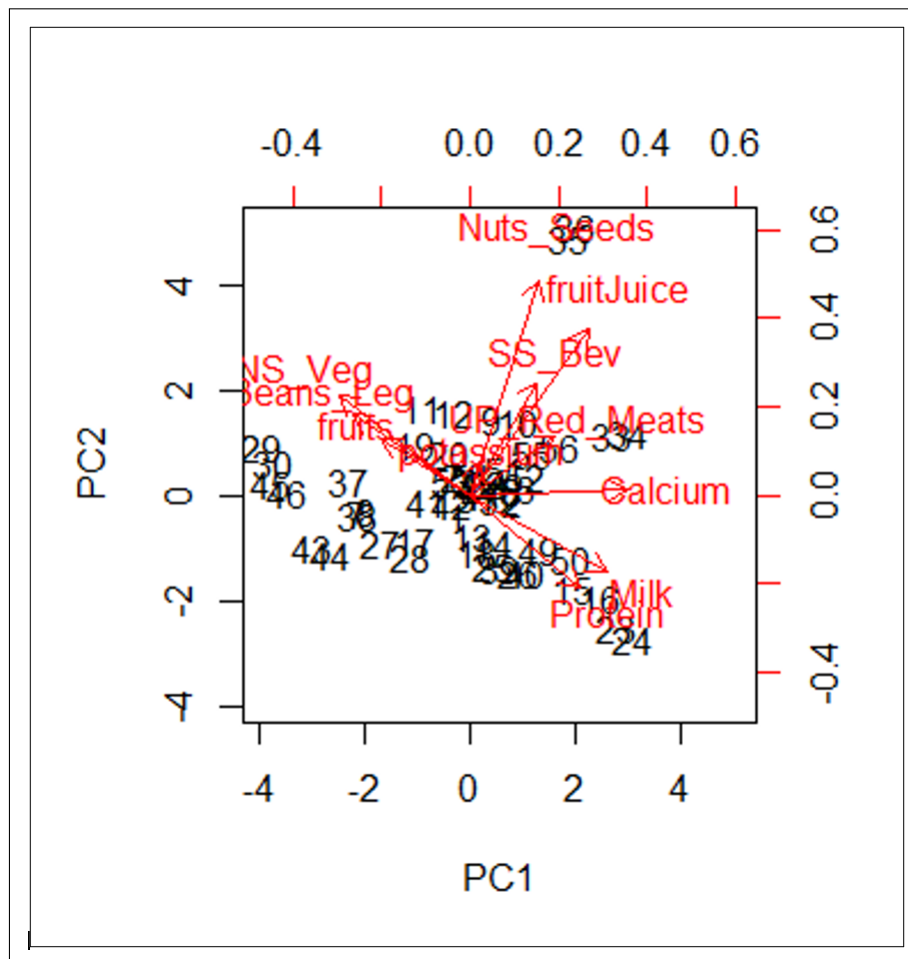


Figure 3: PCA Biplot for Year 2010.

2.3 Methodology: Justification, Assumptions and Development of the Model

The regression analysis focused on two response variables—Stroke rate and diabetes rate. Each response variable is recorded as a count per 100,000. Each model included thirteen potential explanatory variables—eleven nutrients, gender, and year. Several linearized and general linearized models were investigated, analyzed and compared to determine which model had the best fit. Since the response data was count data, two general linearized models were considered. A general linearized model is linear in its parameters. There are three components of a general linearized model. The systematic component is of the form

$$\eta_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_3 x_{pi} \quad (1)$$

The random component follows one of the distributions in the family of Exponential Dispersion Models (EDMs). The final component is the link function component, g , which links the mean, μ , to the linear predictor (1), such that

$$g(\mu_i) = \eta_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_3 x_{pi}$$

The general linearized model assumes that the response comes from an EDM. Initially, for this study two EDM's were considered—Poisson and Gamma. The Poisson distribution [10] was considered because the data was recorded as counts (per 100,000), which could be considered a rate, but since population data was not available to be used as a meaningful offset it was determined that the Poisson distribution was not appropriate. Next, the Gamma distribution was considered. The Gamma distribution can be considered when the response data is positive and continuous which typically means that the data is skewed to the right. Although, the responses in the data under study was positive and skewed right, it was determined that a transformation of the response did a fairly good job of normalizing the response variable. Ultimately, a general linearized model using the Normal distribution as the EDM was used. This can be justified when the response comes from a normal distribution or a transformation of the response makes it approximately normal. A multiple linear regression model is a special case of the general linearized model. Multiple linear regression models with a response variable y , and p explanatory variables, x_1, x_2, \dots, x_p , consist of two components—a systematic component and a random component. The systematic component is the expected value of the response variable which is linearly related to the explanatory variables x_j such as:

$$\mu_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ji} \quad (2)$$

In equation (2), $\mu_i = E[y_i]$ is the expected value of the response variable and the intercept is β_0 . The regression parameters, β_j, x_{ji} represent each of the p explanatory variables. The response variable, y_i , is the random component in the model. it is assumed to have a constant

variance, σ^2 , and to be normally distributed, that is, $y_i \sim N(\mu_i, \sigma^2)$. It should be noted that linear models are a special case of general linearized models. The goal of the statistical model is to mathematically represent the most important systematic and random components of the data. With a good model, one can understand how variables are related to each other and how the explanatory variables significantly affect the response variable. All model assumptions were checked and analyzed using R. The “Residuals vs. Fitted” plot is used to determine constant variance. Ideally, the plot should show a random pattern around the horizontal line at zero. The “Normal Q-Q plot” can be used to show whether the random component is normally distributed. The “Residuals vs. Leverage” plot can be used to determine if there are outliers or influential observations in the data. The R code for all statistical analysis is located in Appendix B.

Model appropriateness was determined using an ANOVA F-test. This is important because it should be determined whether the explanatory variables are useful predictors of the response variable. This can be determined by testing whether the regression sum of squares (SSReg) is larger than the residual sum of squares (RSS). The ANOVA F-test produces the following results:

- F Statistic:

$$F = \frac{SSReg/(p)}{RSS/(n-p-1)} = \frac{MSReg}{MS_E} \sim F_{(p, n-p-1)}$$

- The Coefficient of Determination, R^2 . This is the proportion of the total variation explained by regression:

$$R^2 = 1 - \frac{RSS}{SS_T}$$

<https://www.overleaf.com/project/60541973b5a9aa428084ba62>

- Adjusted R^2 . This is the proportion of the total variation explained by the regression adjusted for the number of explanatory variables:

$$R^2_{adjusted} = 1 - \frac{RSS/(n-p-1)}{SS_T/(n-1)}$$

In the spirit of parsimony, several nested models were considered. Also, a transformations of the response variables was considered. The goal was to get the simplest models with the best predictive and interpretive value. To compare nested models, Akaike’s Information Criterion (AIC) was calculated:

$$AIC = n\log(RSS/n) + 2(p+1) \tag{3}$$

Smaller AIC values (closer to $-\infty$) represent better models. In equation (3), the term, $2(p+1)$, is called the penalty. The Bayesian Information Criterion (BIC) was also calculated:

$$BIC = n\log(RSS/n) + \log(n)(p+1) \tag{4}$$

The BIC is inclined to select more parsimonious models than AIC. Smaller BIC values (closer to $-\infty$) represent better models. In equation (4), the term, $\log(n)(p+1)$, is called the penalty.

If the assumption of normality of the random variable y_i is not satisfied, a transformation of y_i can be considered. Common transformations include a logarithmic or square root transformation. Another useful method is the Box-Cox transformation. The Box-Cox transformation is used to determine the best lambda, λ , that should be used to transform the response variable, y , where the transformed model becomes $y^\lambda = \beta x + \epsilon$. Often, a Box-Cox transformation can improve linearity and homoscedasticity so it can be a very useful tool. The Box-Cox transformation uses a maximum likelihood estimator for λ . If λ is equal to 1, then no transformation is needed. If λ is 0, then a \log transformation should be used, that is $y^\lambda = \log(y)$. If $\lambda \neq 0$, then

$$y^\lambda = \frac{y^\lambda - 1}{\lambda}.$$

Lambda should be relatively small, usually between -3 and 3. R was used to determine the best lambda for the Box-Cox transformation of the data.

2.4 Testing and Analysis

The data was regressed using the thirteen explanatory variables. All model assumptions were checked and addressed if necessary. After an acceptable model was produced, the significant explanatory variables were kept in the model and model adequacy and diagnostics were reevaluated. Ultimately, two models were developed for the response variables stroke rate and diabetes rate. All of the details for model development follow, starting with the model for stroke rate.

2.4.1 Model Formulation for Stroke Rate

Let's begin by examining some descriptive statistics for the response variable Stroke rate. From Table 1, we can see that the mean is greater than the median, so it appears that this data is skewed to the right. The histogram of the variable stroke rate also shows that the data is right skewed (see Figure 4). In order to satisfy the linearity assumption, a transformation of

Minimum	1st Quartile	Median	Mean	3rd Quartile	Maximum
16.88	84.32	123.44	163.55	249.60	451.71

Table 1: Summary Statistics for Stroke Rate

the response variable stroke rate was necessary. The Box-Cox method was used to determine the best transformation. R was used to determine the best lambda, λ , for the Box-Cox transformation. Figure 5 shows the 95% confidence interval for the maximum likelihood estimator, λ .

Using the *boxcox* function in R, the best lambda was calculated to be $\lambda = \frac{10}{99}$. Note that zero is within the 95% confidence interval for Lambda in Figure 5. Since zero is within the 95%

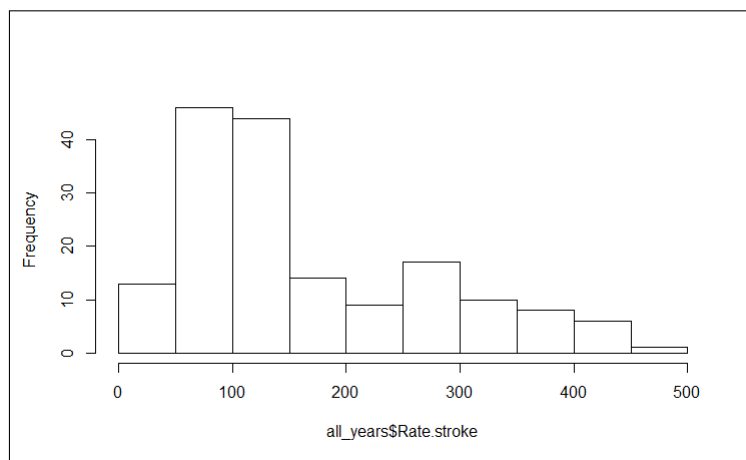


Figure 4: Histogram for Stroke rate

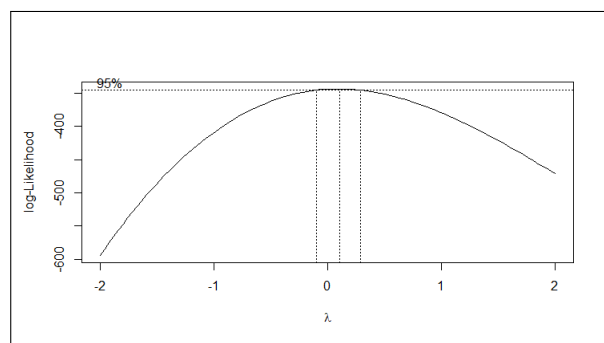


Figure 5: 95% Confidence Interval for Lamda, Stroke Rate.

confidence interval, A *log* transformation was used. The model diagnostics for the *log* model are shown on the next page in Figure 6. A visual analysis of the diagnostic plot in Figure 6 demonstrates that there are no major violations of the model assumptions for the transformed model. The *residuals vs fitted plot* shows a generally random pattern, the *normal QQ plot* shows some slight deviations from normality near the tails of the data, the *scale-location plot* shows that the data has homoscedasticity, and the *residuals vs. leverage plot* show that there are no outliers or influential data points.

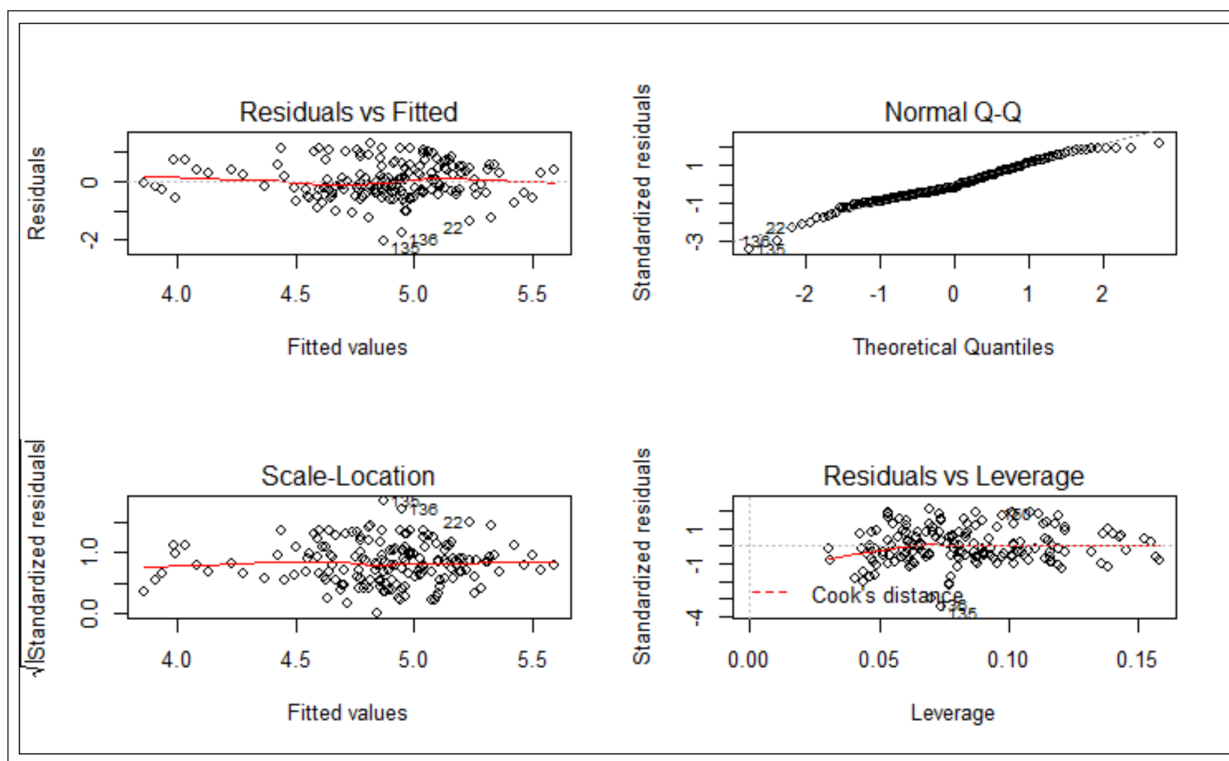


Figure 6: Diagnostic Plot for Log Model, Response Variable = Stroke.

To visualize the effect that a *log* transformation has on the response variable, we created a histogram of the transformed data. Figure 7 is a histogram of the stroke rate after a log transformation. We can see that the transformed data appears to have more normal distribution.

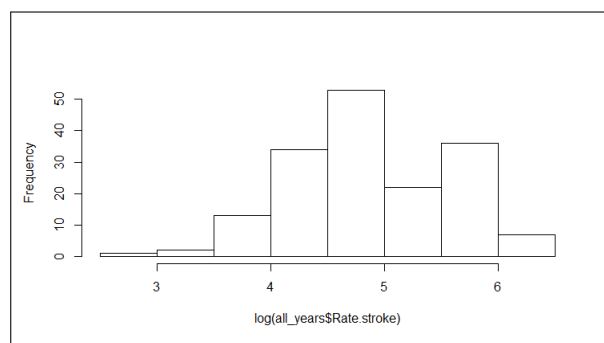


Figure 7: Histogram for Stroke Rate after *log* transformation.

Table 2 includes a full summary of the multiple linear regression model with a *log* transformation and all thirteen explanatory variables included in the model. The F-statistic for the model was 3.261, with a corresponding p-value of 0.0002194. The value of R^2 was 0.2155. Using a level of significance of 0.05, the regression model shows that the significant explanatory

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.4537	0.1146	12.69	0.0000
Gender	0.0023	0.0193	0.12	0.9072
year	-0.0130	0.0122	-1.07	0.2860
fruits	-0.0004	0.0003	-1.29	0.2006
NS_Veg	0.0004	0.0003	1.26	0.2101
Beans_Leg	0.0011	0.0006	1.79	0.0752
Nuts_Seeds	-0.0109	0.0043	-2.55	0.0116
UP_Red_Meats	0.0008	0.0004	2.15	0.0334
SS_Bev	-0.0000	0.0002	-0.10	0.9210
fruitJuice	-0.0001	0.0004	-0.21	0.8369
Protein	-0.0010	0.0011	-0.96	0.3395
Calcium	0.0001	0.0001	1.72	0.0879
potassium	0.0000	0.0000	1.86	0.0641
Milk	0.0002	0.0002	1.03	0.3044

Table 2: Linear Regression results for Stroke Rate with Log transformation

variables were nuts & seeds, and unprocessed red meat. Beans & legumes, calcium, and potassium were mildly significant predictors with p-values of 0.0752, 0.0879, and 0.0641 respectively. Next, the model was recalculated using these five explanatory variables. Potassium was no longer significant and was dropped from the model. The final model had an AIC of -757.22 and a BIC of -259.7149, which is evidence that this model is better than the full model containing all thirteen explanatory variables. Comparison of the full and final model are summarized in Table 3.

Model	AIC	BIC	R^2	F-Statistic	p-value
Full Model	-144.9	380.7262	0.2155	3.261	0.0002194
Final Model	-752.22	-259.7149	0.1631	7.939	0.000007153

Table 3: Comparison of full and final models for Stroke Rate

Table 4 summarizes the new multiple regression model with parameter estimates, with a *log* transformation of the response variable, stroke rate, with the following explanatory variables: nuts seeds, unprocessed red meat, beans legumes, calcium. The sign of the coefficient on the explanatory variable indicates a positive or negative effect on Stroke rate. A positive effect would indicate stroke rate decreased so the sign on the explanatory variable would be negative, likewise, a negative effect would indicate that stroke rate increased so the sign on the explanatory variable would be positive. Nuts & seeds have a positive effect on stroke rate, but unprocessed red meat, beans & legumes, and calcium have a negative effect on stroke rate. According to the CDC [7], foods low in saturated fats can prevent stroke, so nuts and seeds can have high saturated fats thus contributing to a higher risk of stroke, which our study supports. In addition, the CDC states that foods high in in fiber can prevent high cholesterol which can lead to a stroke, and beans and legumes are high in fiber, which our study shows has a negative effect on stroke. Unprocessed red meats can actually have high levels of saturated fat

and cholesterol, which would imply that they would increase the risk of stroke according to the CDC, but we did not find that in our results.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.4458	0.0493	29.32	0.0000
Nuts_Seeds	-0.0115	0.0029	-3.88	0.0002
Calcium	0.0002	0.0001	3.18	0.0018
Beans_Leg	0.0015	0.0005	3.07	0.0025
UP_Red_Meats	0.0009	0.0003	2.64	0.0090

Table 4: Final linear regression model results for Stroke Rate, with Log Transformation

2.4.2 Model Formulation for Diabetes Rate

Next, we investigated the analysis of the response variable, diabetes rate. The analysis followed the same process as the analysis for stroke rate, starting with a look at some descriptive statistics for the response variable, diabetes rate (Table 5).

Minimum	1st Quartile	Median	Mean	3rd Quartile	Maximum
17.99	46.75	72.05	116.67	149.87	473.51

Table 5: Summary Statistics for diabetes Rate

From Figure 8, we can see that, like stroke rate, the variable diabetes rate is skewed to the right. Similar to stroke rate, a transformation of the response variable, diabetes rate, was necessary.

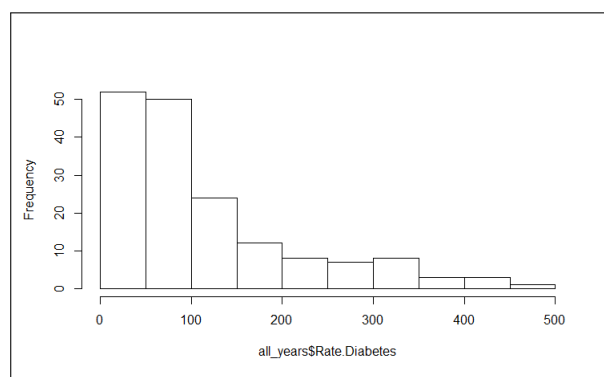


Figure 8: Histogram for diabetes rate.

We used R to determine the best lambda (λ) for the Box-Cox transformation. Figure 9 shows the 95% confidence interval for the maximum likelihood estimator, λ . Using the Box-Cox function in R, the best lambda was calculated to be $\lambda = -\frac{2}{9}$. The 95% confidence interval did not contain zero or another convenient estimate for Lambda, so $\lambda = -\frac{2}{9}$ was used for the transformation.

The model diagnostics for the Box-Cox model are shown in Figure 10. A visual analysis of the diagnostic plot shows that there are no major violations of the model assumptions for either model. The *residuals vs fitted plot* shows a generally random pattern, the *normal QQ plot* shows some slight deviations from normality near the tails of the data, the *scale-location plot* shows that the data has homoscedasticity, and the *residuals vs. leverage plot* shows that there are no outliers or influential data points.

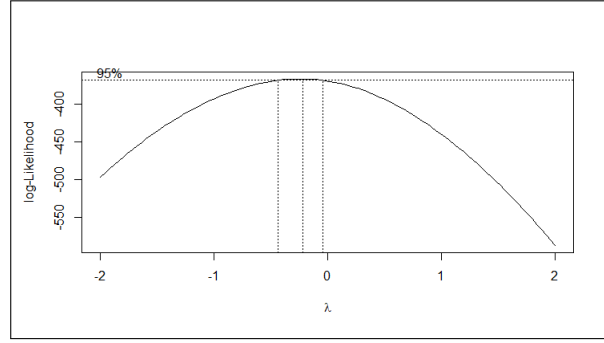


Figure 9: 95% Confidence Interval for Lamda, Diabetes Rate.

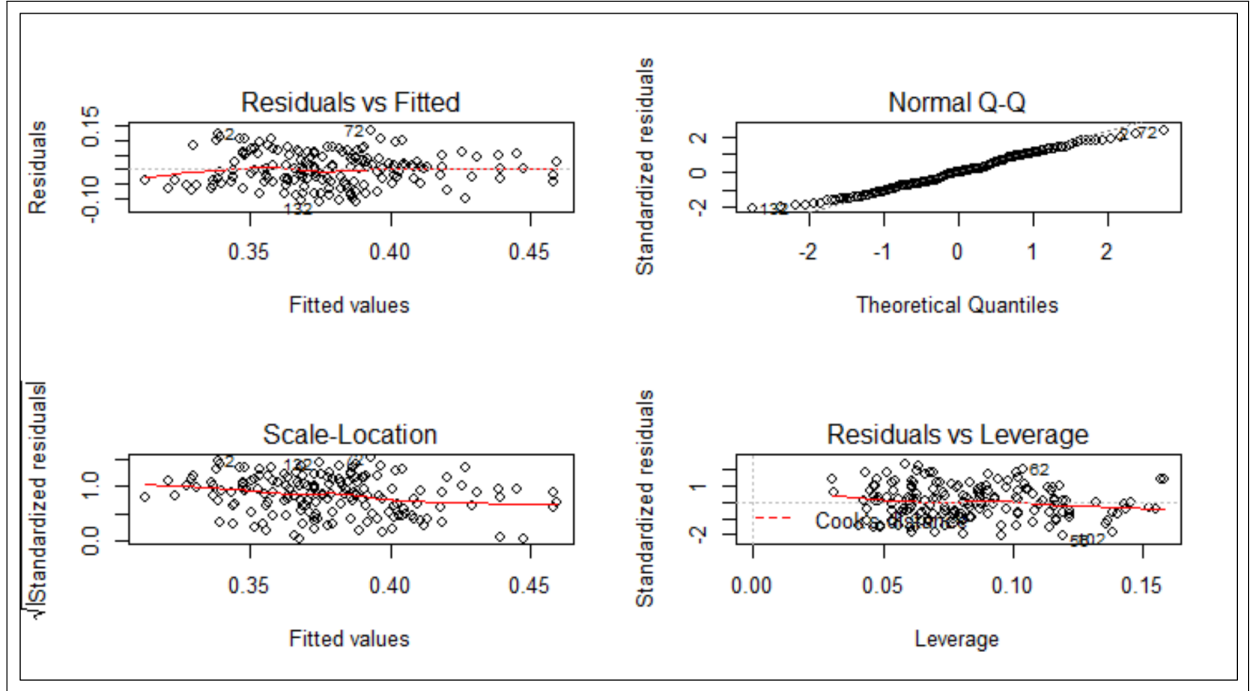


Figure 10: Diagnostic Plot for Box-Cox Model, Response Variable = Diabetes.

To visualize the effect that the Box-Cox transformation has on the response variable, we created a histogram of the transformed data. Figure 11 is histogram of the diabetes rate after a Box-Cox transformation (with $\lambda = -\frac{2}{9}$). We can see that the Box-Cox transformation does a good job “normalizing” the response variable.

Table 6 includes a full summary of the multiple linear regression model with a Box-Cox transformation and all thirteen explanatory variables included in the model. Using a level of significance of 0.05, the regression model shows that the significant explanatory variables were unprocessed red meat, protein, and milk. Year, beans & legumes, and nuts & seeds were mildly significant predictors with p-values of 0.07011, 0.09395, and 0.08111 respectively.

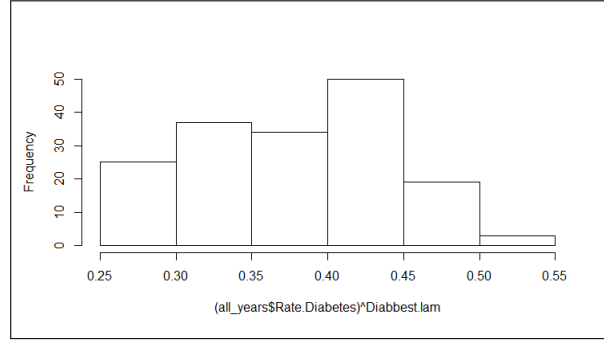


Figure 11: Histogram for Diabetes Rate after *Box-Cox* transformation.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.4098	0.0657	6.23	0.0000
Gender	-0.0175	0.0111	-1.58	0.1161
year	0.0127	0.0070	1.82	0.0701
fruits	0.0003	0.0002	1.48	0.1404
NS_Veg	-0.0003	0.0002	-1.30	0.1951
Beans_Leg	0.0006	0.0004	1.69	0.0940
Nuts_Seeds	0.0043	0.0024	1.76	0.0811
UP_Red_Meats	-0.0007	0.0002	-2.91	0.0042
SS_Bev	0.0001	0.0001	0.52	0.6006
fruitJuice	0.0000	0.0002	0.16	0.8762
Protein	-0.0014	0.0006	-2.32	0.0217
Calcium	0.0000	0.0000	0.38	0.7060
potassium	0.0000	0.0000	0.57	0.5676
Milk	0.0003	0.0001	2.10	0.0372

Table 6: Linear Regression results for Diabetes Rate with Box-Cox transformation

Next, the model was recalculated using these six explanatory variables. Unprocessed red meat and protein were no longer significant and were dropped from the final model. Table 7 summarizes other model statistics that can be used to determine which model is the best. A comparison of the AIC and BIC for both transformations shows that the model using a Box-Cox transformation is best. Based on the model diagnostics, histogram, AIC/BIC, and Coefficient of Determination R^2 , the linearized model using a Box-Cox transformation was determined to be the best and was used to develop the final parsimonious model.

Model	AIC	BIC	R^2	F-Statistic	p-value
Full Model	-457.7402	-410.8807	0.2019	2.996	0.0006074
Final Model	-937.5627	-440.0555	0.1172	5.41	0.0004082

Table 7: Comparison of full and final models for Diabetes Rate

The final model had an AIC of -937.56 and a BIC of -440.0555, which is evidence that the

final model is better than the full model containing all thirteen explanatory variables. Table 8 summarizes the new multiple regression model, including parameter estimates, with a Box-Cox transformation with the following explanatory variables: year, nuts & seeds, milk, beans & legumes. Nuts & seeds, milk, and beans & legumes all had a negative effect on the response variable, diabetes rate. The explanatory variable “year” was a coded numerical factor where year 1990 was coded as “1”, 2000 was coded as “2”, and 2010 was coded as “3”. In this model that means that for every decade, from 1990 to 2010, the rate of diabetes increased by about 0.28. According to the WHO [2], saturated fats increase the risk of diabetes and non-starchy vegetables, plant-based proteins, and lean meats decrease the risk of diabetes. Our study supports these recommendations as beans and legumes had a negative effect on diabetes and they are plant-based protein sources that the WHO recommends. Nuts and seeds are also plant based protein sources, yet they can also have high levels of saturated fats. Our study shows that nuts and seeds have a negative relationship. In addition, the WHO estimates that the number of diabetes cases rose from 108 million in 1980 to 422 million in 2014, and our study agrees with the increase in diabetes rate from 1990 to 2010.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.2873	0.0219	13.14	0.0000
year	0.0185	0.0061	3.03	0.0029
Nuts_Seeds	0.0046	0.0016	2.86	0.0047
Milk	0.0002	0.0001	2.30	0.0227
Beans_Leg	0.0005	0.0003	1.89	0.0602

Table 8: Final linear regression model results for Diabetes Rate, with Box-Cox Transformation

3 Conclusions

3.1 Interpretation and Summary

The aim of this project was to study if there is a relationship between diet and death by strokes and diabetes. We started by utilizing Principal Component Analysis to describe the data. We noticed that there were clusters of countries laying geographically close to each other who had similar food intake, such as Finland and Iceland. We also noticed that non-starchy vegetables and legumes were similar in characteristics, and so were sugar-sweetened beverages and fruit juice. After the descriptive statistics, we progressed to create a model for our variables by utilizing regression analysis. Multiple linear regression results provided parameter estimates for the multiple linear regression equations. Following are the final regression equations for the response variable stroke rate and diabetes rate, respectively

$$\begin{aligned} \log(y_{Stroke}) = & 1.4458 - 0.0115(NutsSeeds) + 0.0002(calcium) \\ & + 0.0015(BeansLegumes) + 0.0009(UPRedMeats) \end{aligned}$$

$$\begin{aligned} (y_{Diabetes})^\lambda = & 0.2873 - 0.0185(Year) + 0.0046(NutsSeeds) + 0.0002(milk) \\ & + 0.0005(BeansLegumes) \end{aligned}$$

Note that in each case, the response variable has been transformed. For the response variable, stroke rate a \log transformation was used, and for the response variable, diabetes rate, the Box-Cox transformation, with $\lambda = -\frac{2}{9}$, was used. In this context, interpretation of the intercept does not make sense since that would imply that consuming none of the eleven nutrients is within the parameters of the model. Both models show that consumption of beans and legumes is associated with slightly higher rates of both stroke and diabetes. Nuts and seeds have a positive effect on the rate of diabetes, but a negative effect on the rate of stroke. Consumption of unprocessed red meat has a slightly negative effect on the rate of stroke. Finally, the variable year is a significant factor in the model for diabetes. Because this effect is negative, we can conclude that over time, specifically each decade from 1990 to 2010, the rate of diabetes has decreased. Overall, these models shows some interesting associations between certain nutrients and disease rates for stroke and diabetes. These correlations do not imply causation and further research and experimentation is needed to state definitively that nutrition and consumption of certain foods causes higher or lower rates of stroke or diabetes.

3.2 Suggestions for Further Study

The conclusions made by the model presented demonstrate correlation, but in order to study if causation does in fact exist, an experimental study would need to be conducted. In

addition, nuts and seeds, unprocessed red meat, and beans and legumes were determined to influence diabetes and stroke, but the method in how these foods are prepared and seasoned can play a huge role on their effect on disease. It would be useful to study how such methods could affect the correlation, as well as how the addition of other food items, such as spices and condiments, not mentioned in our study could further influence the incidence of disease or stroke. Also, the food categories studied were extremely general, and studying more specific foods could uncover more accurate information on the link of food and disease. For example, when studying unprocessed red meat, specifying how the animals were raised, such as factory-farmed versus grass-fed, and the type of red meat, such as lamb, cattle, or bison, would be extremely relevant to study.

References

- [1] Cancer: Overview, prevention, and management. <https://www.who.int/health-topics/cancertab=tab2>, March 2021.
- [2] Diabetes: Key facts. <https://www.who.int/news-room/fact-sheets/detail/diabetes>, March 2021.
- [3] Eat good to feel good. <https://www.diabetes.org/nutrition/healthy-food-choices-made-easy>, March 2021.
- [4] The gdd 2015 beta-version. <https://www.globaldietarydatabase.org/gdd-2015-beta-version>, March 2021.
- [5] Global health data exchange. <http://ghdx.healthdata.org/>, March 2021.
- [6] Scope of current data collection. <https://www.globaldietarydatabase.org/methods/scope-current-data-collection>, March 2021.
- [7] Stroke. <https://www.cdc.gov/stroke/>, March 2021.
- [8] Stroke, cerebrovascular accident. <https://www.emro.who.int/health-topics/stroke-cerebrovascular-accident/index.html>, March 2021.
- [9] M. Brems. A one-stop shop or principle components analysis. *Towards Data Science*, 2017.
- [10] E. L. Frome. The analysis of rates using poisson regression models. *Biomettics*, Volume 39, No. 3, 1983.
- [11] R. Halinski. The selection of variables in multiple regression analysis. *Journal of Educational Measurement*, Volume 7:151-157, 1970.
- [12] Osborne Jason. Four assumptions of multiple linear regression that researchers should always test. *Practical Assessment, Research, and Evaluation*, Volume 8, 2003.
- [13] J. Kuha. Aic and bic: Comparisons of assumptions and performance. *Sociological Methods and Research*, Volume 33(2):188-229, 2004.
- [14] Linh Ngo. Principal component analysis explained simply. <https://blog.bioturing.com/2018/06/14/principal-component-analysis-explained-simply/>, June 2021.
- [15] A. Robertson. Diet and disease. *WHO Regional Publications, European Series*, 2021.
- [16] J. Shlens. A tutorial on principal component analysis. *Google Research*, 2014.

Year	High income countries including HIC in Asia	Country label	Country Name	Gender (female 1, Male 0)	Fruits, grains per day	Non-starchy vegetables, grams per day	Beans and legumes, grams per day	Nuts and seeds, grams per day	Unprocessed red meat, grams per day	Sugar-sweetened beverages, grams per day	Fruit juices, Total protein, grams per day	Calcium, milligrams (mg) per day	Potassium, milligrams (mg) per day	Total milk, grams per day	Cause of death, Stroke, per 100k	Cause of death, Diabetes, per 100k	Cause of death, Cancer, per 100k
1900 HIC	AUS	Australia	Australia	1	159.31151	150.48599	15.89856	3.04485	65.84836	106.93427	49.19272	85.91626	74.81263	27.479321	110.546853	35.14046125	223.894147
1900 HIC	AUS	Australia	Australia	0	140.82078	139.63774	133.46335	9.1065178	18.12635	9.1065178	4.9943973	106.93427	93.000708	93.000708	59.886736	67.451645	73.94214
1900 HIC	AUT	Austria	Austria	1	142.14287	133.46335	9.1065178	18.12635	9.1065178	4.9943973	106.93427	93.000708	93.000708	59.886736	67.451645	73.94214	28.101514
1900 HIC	AUT	Austria	Austria	0	116.27825	148.2383	4.978188	5.3503666	3.3793744	82.3862	192.41958	202.14958	64.338494	100.58065	64.338494	100.58065	64.338494
1900 HIC	BEL	Belgium	Belgium	1	154.42469	148.2383	4.978188	5.3503666	3.3793744	82.3862	192.41958	202.14958	64.338494	100.58065	64.338494	100.58065	64.338494
1900 HIC	BEL	Belgium	Belgium	0	128.57446	131.63817	13.698174	5.3503666	4.3079109	7.8884026	4.30479109	18.594023	192.9137	68.25137	68.25137	68.25137	68.25137
1900 Asia	BRN	Brunei Darussalam	Brunei Darussalam	1	92.637444	166.19374	166.19374	7.8884026	4.30479109	18.594023	192.9137	68.25137	68.25137	68.25137	68.25137	68.25137	68.25137
1900 Asia	BRN	Brunei Darussalam	Brunei Darussalam	0	159.10403	160.00567	159.10403	7.1396027	4.30479109	18.594023	192.9137	68.25137	68.25137	68.25137	68.25137	68.25137	68.25137
1900 HIC	CAN	Canada	Canada	1	149.75789	173.51277	14.888444	5.814889	57.04067	125.0845	113.6834	74.11222	67.83925	233.8076	153.00419	65.7792204	36.4222868
1900 HIC	CAN	Canada	Canada	0	125.68667	152.5697	15.665134	6.7143349	70.741796	104.00441	60.163132	81.217795	280.3403	186.93143	161.6571284	61.6571284	161.6571284
1900 HIC	CYP	Cyprus	Cyprus	1	138.48896	219.93274	69.82217	6.827087	14.64433	19.42856	52.92474	54.41647	72.82031	72.82031	72.82031	72.82031	72.82031
1900 HIC	CYP	Cyprus	Cyprus	0	182.84715	194.86051	182.84715	6.827087	14.64433	19.42856	52.92474	54.41647	72.82031	72.82031	72.82031	72.82031	72.82031
1900 HIC	DNK	Denmark	Denmark	1	132.31978	129.68385	132.31978	3.1080527	2.4302894	63.019428	83.198914	41.24245	58.933197	85.024388	117.4898	254.1.951	319.65489
1900 HIC	DNK	Denmark	Denmark	0	114.92109	125.75286	125.75286	7.8368676	2.4079608	55.91925	58.933197	85.024388	117.4898	254.1.951	319.65489	250.75526	280.75526
1900 HIC	FIN	Finland	Finland	1	108.38974	108.46956	108.46956	8.3397212	2.7361574	72.64850	53.135303	75.40768	61.423358	104.4.9868	255.4.465	319.65489	250.75526
1900 HIC	FIN	Finland	Finland	0	114.95436	200.68903	18.420555	2.6559777	44.47233	12.491199	60.66746	80.66746	80.66746	80.66746	80.66746	80.66746	80.66746
1900 HIC	FRA	France, Monaco	France, Monaco	1	104.55181	188.86267	18.420555	2.6559777	44.47233	12.491199	60.66746	80.66746	80.66746	80.66746	80.66746	80.66746	80.66746
1900 HIC	FRA	France, Monaco	France, Monaco	0	205.88083	221.08168	11.432217	2.4357324	52.780968	89.0818	78.33245	75.944861	1021.333	3431.9038	3431.9038	3431.9038	3431.9038
1900 HIC	DEU	Germany	Germany	1	167.06831	188.73515	12.60461	2.8753314	67.213448	132.15805	69.289834	78.160304	94.923946	330.30543	90.139378	131.875964	104.14078
1900 HIC	DEU	Germany	Germany	0	161.95056	165.04317	25.325388	5.7385657	101.78545	71.487383	22.922808	80.34323	95.156122	307.5.975	307.5.975	307.5.975	307.5.975
1900 HIC	GRC	Greece	Greece	1	138.70767	147.27157	26.824435	6.2817768	107.95437	93.037469	23.088633	86.355868	68.535868	68.535868	68.535868	68.535868	68.535868
1900 HIC	GRC	Greece	Greece	0	105.65095	90.401009	6.192303	0.4954561	63.2145	102.15916	67.972347	67.972347	67.972347	67.972347	67.972347	67.972347	67.972347
1900 HIC	ISL	Iceland	Iceland	1	85.100311	79.846334	5.5842638	0.55479509	63.2145	102.15916	67.972347	67.972347	67.972347	67.972347	67.972347	67.972347	67.972347
1900 HIC	ISL	Iceland	Iceland	0	111.6491	174.50885	39.341621	1.5248671	43.003651	108.9362	49.471588	85.864393	66.91021	66.91021	66.91021	66.91021	66.91021
1900 HIC	IRL	Ireland	Ireland	1	98.877159	157.08861	4.016762	0.8143813	50.391045	138.79421	44.570477	62.258036	101.46411	307.8.1543	307.8.1543	307.8.1543	307.8.1543
1900 HIC	IRL	Ireland	Ireland	0	222.10255	166.62082	22.882709	0.8143813	62.514748	19.297182	44.570477	62.258036	101.46411	307.8.1543	307.8.1543	307.8.1543	307.8.1543
1900 HIC	ITA	Italy	Italy	1	199.17184	146.46224	23.90389	0.9061005	69.745476	61.049808	17.33408	66.443428	2937.4844	2774.4497	2774.4497	2774.4497	2774.4497
1900 Asia	JPN	Japan	Japan	1	146.55228	284.8472	2.2024666	5.515905	66.524361	92.148926	66.524361	66.524361	66.524361	66.524361	66.524361	66.524361	66.524361
1900 Asia	JPN	Japan	Japan	0	108.78735	242.59707	71.306387	1.924291	56.28198	117.54391	72.058584	65.001036	47.068033	2334.9571	97.743011	132.830534	316.833637
1900 HIC	LUX	Luxembourg	Luxembourg	1	119.13246	141.62083	12.21556	3.570286	95.58528	103.4871	70.36907	66.91031	87.011877	32.672073	108.61927	108.61927	108.61927
1900 HIC	LUX	Luxembourg	Luxembourg	0	100.51255	128.28574	12.05235	4.802504	100.8773	136.9431	60.644188	81.424684	32.881532	106.94019	106.94019	106.94019	106.94019
1900 HIC	MLT	Malta	Malta	1	148.21938	126.03325	25.852381	8.426155	102.4653	132.6171	77.82355	106.8392	22.53142	263.13715	135.91374	135.91374	135.91374
1900 HIC	MLT	Malta	Malta	0	126.20165	112.37483	27.974649	9.200297	85.49326	168.26156	92.43926	66.207176	96.681909	22.042288	250.07888	119.215989	67.18882304
1900 HIC	NLD	Netherlands	Netherlands	1	168.65682	172.63239	46.96714	12.43968	15.45049	18.31442	181.658743	27.764285	183.142	131.658743	27.764285	27.764285	27.764285
1900 HIC	NLD	Netherlands	Netherlands	0	145.16576	160.22549	53.699447	14.155262	93.261192	157.24187	175.86302	67.804375	88.330524	37.676482	36.789131	37.676482	36.789131
1900 HIC	NZL	New Zealand	New Zealand	1	261.75635	170.40376	83.28253	1.586189	70.79304	54.3835	44.510426	94.773536	71.568371	287.223	193.0836	124.430719	31.28607023
1900 HIC	NZL	New Zealand	New Zealand	0	204.96317	146.7221	81.669334	1.7226901	59.07643	63.03654	44.22032	101.63755	2505.2008	160.8012	160.8012	160.8012	160.8012
1900 HIC	NOR	Norway, Svalbard and Jan Mayen	Norway, Svalbard and Jan Mayen	1	115.48391	115.7555	3.458766	2.356442	97.54007	71.19194	74.11975	25.681969	25.681969	25.681969	25.681969	25.681969	25.681969
1900 HIC	NOR	Norway, Svalbard and Jan Mayen	Norway, Svalbard and Jan Mayen	0	96.211073	101.62535	3.723541	2.7910262	105.40411	98.584694	36.161125	85.101128	363.9308139	179.80372	156.6131725	220.2659459	363.9308139
1900 HIC	PRT	Portugal	Portugal	1	148.22816	158.43887	44.818266	2.3619786	107.1513	98.472931	9.531082	83.724684	362.10405	145.51485	386.602086	145.51485	386.602086
1900 HIC	PRT	Portugal	Portugal	0	127.17575	141.14912	46.82361	8.598573	79.040718	76.901778	35.681192	144.94517	54.1523413	178.650401	154.523413	178.650401	154.523413
1900 Asia	KOR	Republic of Korea	Republic of Korea	1	146.37234	164.55528	0.5750272	0.5750272	19.487236	28.420128	19.487236	74.834068	465.02017	27.51354	27.51354	176.854001	10.3170824
1900 Asia	KOR	Republic of Korea	Republic of Korea	0	107.97582	160.81822	64.69567	0.551151	88.04391	41.35413	19.48883	74.754173	44.535971	255.134	99.735306	153.5351276	321.2441562
1900 Asia	SGP	Singapore	Singapore	1	169.23194	168.42163	48.586716	3.1226339	37.36786	97.072838	17.79151	68.753411	2619.7888	51.169756	51.169756	51.169756	51.169756
1900 Asia	SGP	Singapore	Singapore	0	142.83807	158.47678	43.004539	3.033291	38.971302	7.4823904	42.412317	2550.8516	46.877908	146.877908	146.877908	146.877908	146.877908
1900 HIC	ESP	Spain	Spain	1	130.89755	129.3286	35.975891	7.960426	73.190536	99							

2000	HIC	DNK	Denmark	137.13255	120.50188	23.569337	2.277905	65.01864	60.34412	41.848892	75.57239	97.222675	3108.6472	136.37762	153.890974	44.07682789	387.835604
2000	HIC	DNK	Denmark	114.87169	107.38815	24.782774	2.543200	73.19708	83.519646	81.432627	81.432656	3127.2278	134.5615	117.617915	51.7176656	407.637188	
2000	HIC	FIN	Finland	114.93607	119.36099	6.2678731	2.962007	48.00675	50.244759	54.931602	84.820368	1153.1167	282.43368	164.7281965	21.953732	259.489033	
2000	HIC	FIN	Finland	108.65964	104.61882	67.462226	2.672232	64.083461	73.68939	49.3176942	91.411972	1078.4974	2614.2607	261.04664	104.9933223	17.9993366	289.204723
2000	HIC	FRA	France, Monaco	118.25168	190.27296	14.073345	1.7579942	39.275307	39.275307	22.335496	79.237938	83.58952	139.41356	270.3482058	137.41356	116.684122	44.7665527
2000	HIC	FRA	France, Monaco	107.97547	117.52371	15.162725	3.181762	32.02327	98.419628	19.991577	98.34361	3319.1694	137.14151	54.5274308	66.5732403	41.15958537	
2000	HIC	DEU	Germany	202.32656	212.42178	8.8726176	2.4223804	9.996278	8.8726176	73.984361	73.984361	10551.739	3540.0554	80.627731	35.3101282	98.1769633	62.3205921
2000	HIC	DEU	Germany	166.96608	184.18611	9.9884539	2.870861	65.059708	168.44485	73.786871	73.786871	10651.739	3540.0554	80.627731	35.3101282	98.1769633	42.07160657
2000	HIC	GRC	Greece	156.94447	153.26363	21.689384	2.970891	79.61327	74.805106	84.88991	76.156031	980.20508	3454.1055	80.578456	46.93440255	280.2242339	
2000	HIC	GRC	Greece	133.93078	137.40349	22.941869	6.005348	64.45372	92.269821	22.94742	88.264338	9119.1309	4394.3184	90.207253	219.2955072	42.7911888	418.4056735
2000	HIC	ISL	Iceland	109.7618	87.140566	4.9691563	0.4888268	51.900678	101.8592	58.08796	96.85978	1067.0776	2630.8057	286.1246	270.344852	100.3444182	117.1062431
2000	HIC	ISL	Iceland	88.68891	77.162382	5.3503151	1.9456598	58.505409	72.87632	51.946588	98.871734	282.0591	266.46494	282.0591	74.7637079	19.26004482	
2000	HIC	IRL	Ireland	100.2801	149.22032	25.93013	5.429513	34.537098	116.7834	81.88639	98.738458	867.65271	3670.6889	68.6181251	26.1807187	32.24001579	
2000	HIC	IRL	Ireland	86.283165	134.4777	27.765566	1.9067928	39.877517	142.4305	37.773233	100.13203	97.45425	3590.5797	163.08932	136.6031226	66.9801457	37.3256452
2000	HIC	ITA	Italy	224.95158	157.35288	17.542099	0.8071988	55.485383	45.071633	127.77632	76.409225	187.87678	3231.8543	111.54497	295.4156917	66.98741839	166.5190228
2000	HIC	ITA	Italy	202.22233	140.53589	18.414808	0.9002974	61.984358	60.816288	16.03463	83.995177	94.643868	3097.3887	102.46165	441.499321	53.41087719	126.0305243
2000	HIC	JPN	Japan	144.44579	248.51274	71.121567	2.4862759	45.46479	68.80667	22.249532	66.88049	596.6275	2800.3015	90.746361	120.313869	34.28659075	242.6123162
2000	HIC	JPN	Japan	108.77936	227.93599	62.133442	2.062207	47.689767	61.95164	97.186794	65.807991	480.07666	65.807991	90.746361	120.313869	34.28659075	76.36464149
2000	HIC	LUX	Luxembourg	129.33229	126.97251	17.085117	2.1237295	77.282934	110.2277	58.639342	73.81719	909.73814	3294.7524	90.780373	167.7080468	46.72951462	280.0531312
2000	HIC	LUX	Luxembourg	110.1596	113.03977	18.276754	2.4395239	84.541412	144.0769	62.738739	79.45452	850.7381	3323.005	89.167709	100.3828988	32.0052166	38.8895662
2000	HIC	MLT	Malta	146.15157	112.23044	17.526529	8.2423746	9.2423347	140.01265	86.557259	85.55573	1026.4799	3116.4209	123.1382167	123.1382167	15.1901238	219.1830212
2000	HIC	MLT	Malta	124.51641	98.58006	18.673586	9.2423347	101.62215	178.52522	77.943002	90.797775	1020.1349	3145.1926	196.89159	196.89159	61.4539371	287.573781
2000	HIC	NLD	Netherlands	159.44373	199.88632	34.633797	14.912676	72.677235	172.6396	172.6396	74.75547	98.916928	3165.4744	137.87008	144.2128509	58.6108602	285.275291
2000	HIC	NLD	Netherlands	138.13588	149.12082	39.29763	1.5501611	79.505771	201.2481	156.66229	76.03735	914.78168	3105.1768	133.83802	92.3387948	37.4068457	
2000	HIC	NZL	New Zealand	272.64822	161.52966	58.157494	1.8501611	66.763241	55.417229	45.151615	83.190398	124.62012	3092.6438	141.80777	253.181823	113.5187618	38.3246781
2000	HIC	NZL	New Zealand	214.37637	141.36971	57.22464	1.885606	63.865618	40.474728	68.483795	54.838914	3129.2424	132.28464	300.5934808	30.6504808	45.288890909	
2000	HIC	NOR	Norway, Svalbard and Jan Mayen	111.09414	101.99855	2.3652024	2.2705204	73.690475	83.405202	34.817741	76.86258	830.44098	2723.6248	41.87833	35.4707371	308.5866178	172.1672709
2000	HIC	NOR	Norway, Svalbard and Jan Mayen	92.112432	89.607136	2.533752	2.8912515	70.74408	102.30785	31.053767	82.457157	786.13049	2754.2212	138.7773	37.9802342	124.133271	
2000	HIC	PRT	Portugal	146.01425	178.9445	38.844679	3.5667768	80.372053	98.96397	81.943088	75.123353	707.0451	5375.3435	130.5726	29.12265061	81.9427745	246.6317014
2000	HIC	PRT	Portugal	126.97682	129.96989	38.768881	4.0271187	89.300606	127.88326	7.927033	79.82837	3401.4138	129.834	268.2727489	73.8877775	363.69033074	
2000	Asia	KOR	Republic of Korea	154.00068	141.60341	50.431366	0.6274838	4.925982	27.0531	16.589808	73.043137	45.524023	3067.917	95.019936	14.616328	44.1897571	150.8311281
2000	Asia	KOR	Republic of Korea	113.48779	138.4635	45.165766	0.8242216	52.480964	42.220142	16.54518	72.97057	472.26046	2886.0369	80.429169	11.01989891	45.7171856	253.770251
2000	Asia	SGP	Singapore	153.89455	36.860383	36.860383	2.8004155	19.43201	89.782166	7.6196732	61.061516	67.196732	575.1062	61.061516	57.515424	54.2171184	11.126525
2000	Asia	SGP	Singapore	175.7309	145.01447	32.782055	2.884057	20.244188	111.94642	7.8639447	61.289586	45.415424	364.67913	50.028236	54.74202184	153.4882551	18.8309886
2000	HIC	ESP	Spain	124.41893	119.46782	26.234034	7.742384	58.541115	99.919197	23.242201	76.33480	87.679346	32.12552	162.62575	14.65717882	68.54707828	213.133585
2000	HIC	ESP	Spain	109.95933	107.3229	27.77431	8.1897603	65.571016	129.7416	20.951403	82.012474	82.012474	3237.252	160.3012	101.6731394	34.9892840	
2000	HIC	SWE	Sweden	137.94211	132.68839	9.9513531	2.3461229	70.74667	70.982667	65.023178	77.73556	88.843308	2853.3821	286.5842	32.13963728	179.3199129	40.07806564
2000	HIC	SWE	Sweden	106.96545	114.81633	11.19322	2.948363	88.75936	93.24637	58.46515	82.493095	86.539732	2807.4971	291.82789	38.438362	132.220567	41.78903861
2000	HIC	CHE	Switzerland, Liechtenstein	144.81605	134.86286	4.8277945	4.5149992	59.565728	94.91428	98.749378	66.665718	86.665718	2744.9191	76.188283	76.41386306	39.2686973	334.771534
2000	HIC	CHE	Switzerland, Liechtenstein	119.48265	119.8285	5.1066332	5.2020954	65.887093	124.78244	72.765697	62.164869	95.61009	2769.0156	124.74404	114.870712	11.5270091	239.421551
2000	HIC	GBR	United Kingdom	120.88939	139.30351	34.843378	4.5475909	36.446335	153.44857	36.51432	68.516889	89.756682	3534.6196	174.5383	17.02727409	28.04674814	345.3795347
2000	HIC	GBR	United Kingdom	100.15926	122.80782	36.819392	4.9860034	40.03592	165.61911	32.688461	71.611867	89.756682	3534.6196	174.5383	17.02727409	28.04674814	345.3795347
2000	HIC	USA	USA, Puerto Rico and US Virgin Islands	90.816201	132.38176	18.141495	3.1196719	38.52104	89.280665	81.013123	81.013123	82.53347	2956.9187	149.45209	10.5940485	66.5449642	27.70763239
2000	HIC	USA	USA, Puerto Rico and US Virgin Islands	81.8899	122.18653	121.89526	4.318814	47.59406	313.0304	80.74865	86.205268	94.130196	3034.658	189.74077	70.29134253	59.4474234	34.18874986
2000	HIC	AUS	Australia	138.82275	120.8683	22.1348	4.4053764	87.143446	169.5478	35.506111	85.237022	80.44629	2919.1914	18.28766	91.8986403	47.9029632	231.6818119
2000	HIC	AUS	Australia	156.75948	130.10384	19.378057	3.7807921	51.936079	120.16269	39.769291	68.82337	77.2052	78.62781	121.86072	57.96501779	47.81671951	32.1463732
2000	HIC	AUT	Austria	112.29446	103.88571	6.610344	4.0921454	109.68809	124.8773	45.17857	68.844789	77.41223	3117.5896	81.531965	27.84562779	91.65789803	78.68912684
2000	HIC	AUT	Austria	134.07485	117.62232	6.2303546	3.5690949	87.939385	95.449593	50.37204	63.52792	77.179779	3084.0129	82.62853	34.33653966	67.1539719	62.67347708
2000	HIC	BEL	Belgium	109.80388	116.84699	3.6786447	3.6388621	75.930275	204.3196	46.330525	98.541628	743.1217	4321.1167	72.912041	12.93161467	56.01081919	307.8456234
2000	HIC	BEL	Belgium	129.45607	130.76671	3.409324	3.0592546	70.593951	160.1274	51.887321	84.88381	43.86562	76.689102	86.839889	41.82156745	42.7791604	
2000	Asia	BRN	Brunei Darussalam	87.397346	149.27086	7.0795889	4.6099592	14.934725	216.2895	6.090847	69.130196	48.83209	1844.098	24.225843	43.89328013	88.662275	140.7731612
2000	Asia	BRN	Brunei Darussalam	99.911453	156.21239	7.7686339	4.7720242	14.52898	175.5044	62.996529	68.82337	50.1892	1897.7913	24.483322	41.2751941	77.5067524	135.6793725
2000	HIC	CAN	Canada	129.30704	139.78021	7.0166655	7.3259076	58.407188	156.87065	89.629923	77.2052	84.810862	3270.7629	143.9976	67.9141577	46.75687394	289.4340051
2000	HIC	CAN	Canada	153.05698	159.34947	6.6501212	6.3131089	47.611103	120.29539	99.921781	71.647514	930.80149	3246.2351	148.80614	48.3633087	46.46263909	317.8812184
2000																	

2010	Asia	JPN	Japan	0	108,221,889	2,16,14,207	56,869,757	2,16,017,1	43,360,737	112,183,79	20,324,27	64,751,068	499,417,45	2,196,31,94	83,622,955	287,339,042	42,309,020	120,186,883
2010	Asia	JPN	Japan	1	143,80,027	226,20,352	65,309,776	2,686,536	41,110,466	85,40,7616	21,319,603	68,538,383	611,173,15	257,91,563	97,62,421	451,709,969	44,584,4253	118,256,822
2010	HIC	LUX	Luxembourg	0	101,718,02	105,20,293	11,524,368	2,49,248	74,674,673	146,05,989	48,70,518	80,890,305	875,37,512	35,26,5974	75,43,9641	112,151,1778	277,97,23169	47,24,559132
2010	HIC	LUX	Luxembourg	1	118,397,26	116,85,345	10,72,389	2,16,93,908	68,696,35	113,100,76	54,23,204	73,0,67036	593,63,392	34,95,3975	78,57,9056	74,64,77,413	328,87,54419	34,53,61,2412
2010	HIC	MLT	Malta	0	115,671,61	92,907,547	16,62,9344	8,74,0593	87,193,488	162,7,658	70,83,983	90,154,648	1040,60,05	31,56,52,12	167,7,9219	213,254,034	96,200,4283	56,69,92,1152
2010	HIC	MLT	Malta	1	135,754,19	105,42,561	15,414,862	7,701,1871	81,192,978	143,241,55	78,83,824	84,887,866	1149,21,72	31,24,4,683	170,751,62	81,184,27,3236	58,33,94,7441	61,184,27,3236
2010	HIC	NLD	Netherlands	0	136,018,02	142,08,874	37,890,409	15,99,6328	77,12,9761	200,03,16	146,56,447	81,43,0672	94,2,92,761	324,6,446	120,34,307	102,62,18,649	52,304,61698	338,60,03,2818
2010	HIC	NLD	Netherlands	1	156,006,77	151,85,458	32,883,099	14,01,3341	70,872,177	173,45,276	162,7,1309	80,071,655	1019,72,16	32,96,1702	123,86,318	65,24,08,678	40,462,4599	412,60,05,99
2010	HIC	NZL	New Zealand	0	203,082,35	131,3,988	42,71,2078	2,05,421,9	49,662,956	67,40,6983	36,40,6117	95,94,6273	96,2,57,033	31,99,4,007	114,3,9547	255,02,82,469	93,37,26,5516	38,22,90,3639
2010	HIC	NZL	New Zealand	1	258,201,61	150,47,093	43,42,1249	1,494,3386	59,617,67	58,52,277	40,62,941	89,208,193	742,53,801	317,3,8241	121,846,73	307,104,626	62,7701,9719	44,3470,4244
2010	HIC	NOR	Norway, Svalbard and Jan Mayen	0	92,75,251	83,008,415	0,94,05,385	2,652,2877	68,03,942	105,281,39	28,291,484	83,930,605	809,51,969	297,1,2676	117,79,86	39,69,462,718	12,15,62,768	304,47,90117
2010	HIC	NOR	Norway, Svalbard and Jan Mayen	1	111,070,14	94,12,8227	0,87,65,136	2,235,2827	62,921,833	76,77,5162	31,72,06,85	78,052,193	654,83,472	29,35,47,24	120,30,702	39,60,920,066	81,19,862,9365	348,62,31167
2010	HIC	PRT	Portugal	0	123,607,11	123,20,934	33,584,518	3,47,35,98	78,956,13	130,75,31	71,63,9064	79,95,911	815,50,995	3407,3,677	112,09,917	254,154,9746	200,47,31,894	95,58,113
2010	HIC	PRT	Portugal	1	143,56,954	137,7,6947	31,838,121	3,048,9988	70,898,42	100,87,45	7,97,9941	74,500,344	888,02,987	3377,0,74	42,11,94,73	42,11,94,73	169,72,8755	84,90,40,2736
2010	Asia	KOR	Republic of Korea	0	108,887,98	124,996	36,058,375	0,6707,7637	40,552,27	40,23,9349	14,34,4538	70,16,065	493,67,964	2859,11,6	61,86,715	84,76,98,878	44,74,37,658	157,162,063
2010	Asia	KOR	Republic of Korea	1	147,613,92	127,79,382	40,311,671	0,685,9629	37,77,653	25,651,65	14,35,333	70,208,51	57,814,21	3037,8,237	73,162,62	77,74,54,684	46,531,0849	271,50,04,48
2010	Asia	SGP	Singapore	0	151,55,37	131,3,6411	22,617,411	3,863,2862	23,39,972	114,15,536	6,86,62,229	86,25,061	47,51,6536	2647,4,302	33,10,1646	39,25,94,677	123,58,7374	20,50,03,4701
2010	Asia	SGP	Singapore	1	179,780,38	139,37,749	25,481,451	3,881,2694	22,15,99,94	91,42,1913	6,87,24,475	85,890,518	53,15,0629	2721,3,353	36,16,9983	30,69,054,172	134,71,66318	15,70,82,56
2010	HIC	ESP	Spain	0	101,637,09	102,00,794	22,848,114	8,580,1945	63,37,0623	126,88,862	19,33,0283	816,19,867	19,33,0283	3113,8,892	138,64,375	103,3,62,343	64,58,85,2597	218,37,45179
2010	HIC	ESP	Spain	1	114,59,21	113,2,1651	21,594,639	7,489,4307	56,47,9423	98,24,2876	21,484,127	76,47,956	90,3,4,6273	30,83,9256	141,0,9637	76,36,073,073	368,621,0923	373,987,593
2010	HIC	SWE	Sweden	0	108,98,652	107,61,472	8,130,7516	3,030,9141	75,007,376	96,40,3137	54,22,554	84,40,8211	874,81,073	3127,84,67	256,15,945	144,82,44,49	46,93,55,6259	311,32,0987
2010	HIC	SWE	Sweden	1	140,036,65	124,12,985	7,255,6391	2,681,0267	59,314,97	72,53,8139	60,30,545	79,60,893	90,81,0583	3222,2,259	261,702,64	101,1,650,923	49,18,60,6207	343,613,3018
2010	HIC	CHE	Switzerland, Liechtenstein	0	121,26,073	113,8,601	3,424,9625	5,035,2001	56,54,6917	117,401,24	69,72,0734	70,829,619	994,01,95	27,61,8,665	112,692,45	67,65,268,391	41,4,64,6586	317,21,21276
2010	HIC	CHE	Switzerland, Liechtenstein	1	146,252,37	127,8,6882	3,234,9955	4,414,5555	51,507,417	93,04,2747	78,23,5161	66,861,388	1056,32,42	2726,6,53	114,4,0075	67,99,679,554	53,81,30,9671	244,68,3638
2010	HIC	GBR	United Kingdom	0	106,790,31	118,46,258	39,104,225	5,282,1541	35,330,527	172,657,77	32,57,944	70,890,526	862,84,344	3770,0,779	158,39,336	121,50,22,718	25,40,31,1796	325,61,60,369
2010	HIC	GBR	United Kingdom	1	128,058,17	134,04,881	37,272,766	4,581,4557	32,164,341	143,841,38	36,34,651	67,86,953	920,23,492	35,53,7256	165,54,646	80,91,801,163	23,16,72,7928	373,987,593
2010	HIC	USA	USA, Puerto Rico and US Virgin Islands	0	77,611,382	117,36,357	18,014,244	4,475,8754	44,56,9722	308,321,69	76,65,2855	86,65,6631	829,84,277	2890,50,49	145,44,727	265,230,389	83,047,9107	67,16,09,7482
2010	HIC	USA	USA, Puerto Rico and US Virgin Islands	1	85,610,626	126,75,591	14,936,326	3,852,3872	35,930,305	282,28,59	84,83,5745	81,422,737	84,51,4433	2941,55,96	135,788,13	301,8,941,794	67,207,7577	68,04,83,3031

Appendix B: RCode

```
setwd("H:/CPC STA6950")
install.packages("xtable")
install.packages("ggplot")
library(ggplot2)
library(xtable)
library(MASS)
library(plyr)
```

#combine three years into one large data set

```
y1990<-read.csv("H:/CPC STA6950/STA6950 Project Data Year 1990.csv",
               header = T)
y1990$year<-c(1)
y2000<-read.csv("H:/CPC STA6950/STA6950 Project Data Year 2000.csv",
               header = T)
y2000$year<-c(2)
y2010<-read.csv("H:/CPC STA6950/STA6950 Project Data Year 2010.csv",
               header = T)
y2010$year<-c(3)
```

#rename columns Year 1990

```
colnames(y1990)[colnames(y1990) %in% c("v01_wt_median","v02_wt_median")]<-c("fruits","NS_Veg")
colnames(y1990)[colnames(y1990) %in% c("v05_wt_median","v06_wt_median")]<-c("Beans_Leg","Nuts_Seeds")
colnames(y1990)[colnames(y1990) %in% c("v10_wt_median","v15_wt_median")]<-c("UP_Red_Meats","SS_Bev")
colnames(y1990)[colnames(y1990) %in% c("v16_wt_median","v23_wt_median")]<-c("fruitJuice","Protein")
colnames(y1990)[colnames(y1990) %in% c("v36_wt_median","v41_wt_median")]<-c("Calcium","potassium")
colnames(y1990)[colnames(y1990) %in% c("v57_wt_median")]<-c("Milk")
colnames(y1990)[colnames(y1990) %in% c("female.....1.")]<-c("Gender")
colnames(y1990)[colnames(y1990) %in% c("country.name")]<-c("countryname")
```

#rename columns Year 2000

```
colnames(y2000)[colnames(y2000) %in% c("v01_wt_median","v02_wt_median")]<-c("fruits","NS_Veg")
colnames(y2000)[colnames(y2000) %in% c("v05_wt_median","v06_wt_median")]<-c("Beans_Leg","Nuts_Seeds")
```

```

colnames(y2000)[colnames(y2000) %in% c("v10_wt_median","v15_wt_median")]<-c("UP_Red_Meats","SS_Bev")
colnames(y2000)[colnames(y2000) %in% c("v16_wt_median","v23_wt_median")]<-c("fruitJuice","Protein")
colnames(y2000)[colnames(y2000) %in% c("v36_wt_median","v41_wt_median")]<-c("Calcium","potassium")
colnames(y2000)[colnames(y2000) %in% c("v57_wt_median")]<-c("Milk")
colnames(y2000)[colnames(y2000) %in% c("female.....1.")]<-c("Gender")

```

```

#rename columns Year 2010

```

```

colnames(y2010)[colnames(y2010) %in% c("v01_wt_median","v02_wt_median")]<-c("fruits","NS_Veg")
colnames(y2010)[colnames(y2010) %in% c("v05_wt_median","v06_wt_median")]<-c("Beans_Leg","Nuts_Seeds")
colnames(y2010)[colnames(y2010) %in% c("v10_wt_median","v15_wt_median")]<-c("UP_Red_Meats","SS_Bev")
colnames(y2010)[colnames(y2010) %in% c("v16_wt_median","v23_wt_median")]<-c("fruitJuice","Protein")
colnames(y2010)[colnames(y2010) %in% c("v36_wt_median","v41_wt_median")]<-c("Calcium","potassium")
colnames(y2010)[colnames(y2010) %in% c("v57_wt_median")]<-c("Milk")
colnames(y2010)[colnames(y2010) %in% c("female.....1.")]<-c("Gender")

```

```

all_years<-rbind(y1990,y2000,y2010)

```

```

#histogram for response variables

```

```

#hist(all_years$Rate.Cancer)

```

```

hist(all_years$Rate.stroke, main=" ")

```

```

hist(all_years$Rate.Diabetes, main=" ")

```

```

##Summary stats for stroke and diabetes rate

```

```

s=summary(all_years$Rate.stroke)

```

```

d=summary(all_years$Rate.Diabetes)

```

```

#####response var treated as a continuous random var.

```

```

#####STROKE

```

```

##LM model for stroke

```

```

modelStr<-lm(formula = Rate.stroke ~ Gender+year+fruits+NS_Veg+Beans_Leg+

```

```

    Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice

```

```

    +Protein+Calcium+potassium+Milk,

```

```

    data = all_years)

```

```

summary(modelStr)

```



```
par(mfrow=c(2,2))  
plot(modelStr)  
extractAIC(modelStr)
```

```
#BC transformation
```

```
bcStr=boxcox(modelStr,lambda = seq(-2,2))  
Strbest.lam=bcStr$x[which(bcStr$y==max(bcStr$y))]
```

```
##LM model with BC trans
```

```
bcmodelstr<-lm((Rate.stroke)^Strbest.lam ~ Gender+year+fruits+NS_Veg+Beans_Leg+  
  Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice  
  +Protein+Calcium+potassium+Milk,  
  data = all_years)  
summary(bcmodelstr)  
par(mfrow=c(2,2))  
plot(bcmodelstr)  
extractAIC(bcmodelstr)  
hist((all_years$Rate.stroke)^Strbest.lam,main = " ")  
BIC(bcmodelstr)
```

```
#LM log model
```

```
log_modelStr<-lm(formula = log(Rate.stroke) ~ Gender+year+fruits+NS_Veg+Beans_Leg+  
  Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice  
  +Protein+Calcium+potassium+Milk,  
  data = all_years)  
summary(log_modelStr)  
par(mfrow=c(2,2))  
plot(log_modelStr)  
extractAIC(log_modelStr)  
hist(log(all_years$Rate.stroke),main = " ")  
BIC(log_modelStr)
```

```
#gamma model link inverse, log trans
```

```

gamma_modelstr=glm(log(Rate.stroke) ~ Gender+year+fruits+NS_Veg+Beans_Leg+
  Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice
  +Protein+Calcium+potassium+Milk,
  data = all_years, family = Gamma(link="inverse"))
summary(gamma_modelstr)
par(mfrow=c(2,2))
plot(gamma_modelstr)
extractAIC(gamma_modelstr)
#####End Stroke

```

#####Diabetes

#LM model no trans

```

modelDiab<-lm(formula = Rate.Diabetes ~ Gender+year+fruits+NS_Veg+Beans_Leg+
  Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice
  +Protein+Calcium+potassium+Milk,
  data = all_years)
summary(modelDiab)
par(mfrow=c(2,2))
plot(modelDiab)
extractAIC(modelDiab)

```

#BC TRans

```

bcDiab=boxcox(modelDiab,lambda = seq(-2,2))
Diabbest.lam=bcDiab$x[which(bcDiab$y==max(bcDiab$y))]

```

#LM model with BC trans

```

bcmodeldiab<-lm((Rate.Diabetes)^Diabbest.lam ~ Gender+year+fruits+NS_Veg+Beans_Leg+
  Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice
  +Protein+Calcium+potassium+Milk,
  data = all_years)
summary(bcmodeldiab)
par(mfrow=c(1,1))
plot(bcmodeldiab)

```

```
extractAIC(bcmodeldiab)
hist((all_years$Rate.Diabetes)^Diabbest.lam)
BIC(bcmodeldiab)
```

```
#LM Log Model
```

```
log_modeldiab<-lm(formula = log(Rate.Diabetes) ~ Gender+year+fruits+NS_Veg+Beans_Leg+
  Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice
  +Protein+Calcium+potassium+Milk,
  data = all_years)
summary(log_modeldiab)
par(mfrow=c(2,2))
plot(log_modeldiab)
extractAIC(log_modeldiab)
BIC(log_modeldiab)
```

```
#Gamma model link inverse, log trans
```

```
gamma_modeldiab=glm(log(Rate.Diabetes) ~ Gender+year+fruits+NS_Veg+Beans_Leg+
  Nuts_Seeds+UP_Red_Meats+SS_Bev+fruitJuice
  +Protein+Calcium+potassium+Milk,
  data = all_years, family = Gamma(link="inverse"))
summary(gamma_modeldiab)
par(mfrow=c(2,2))
plot(gamma_modeldiab)
extractAIC(gamma_modeldiab)
```

```
#####END Diabetes
```

```
##best models
```

```
summary(bcmodelstr)
summary(bcmodeldiab)
summary(log_modelCan)
```

```
#####new Models with significant Explanatory Variables only
```

####STROKE####

```
beststroke<-lm((Rate.stroke)^Strbest.lam ~ Nuts_Seeds+Calcium
               +Beans_Leg+UP_Red_Meats,
               data = all_years)
summary(beststroke)
par(mfrow=c(2,2))
plot(beststroke)
extractAIC(beststroke)
BIC(beststroke)
```

####Diab####

```
bestdiab<-lm((Rate.Diabetes)^Diabbest.lam ~ year+Nuts_Seeds
              +Milk+Beans_Leg,
              data = all_years)
summary(bestdiab)
par(mfrow=c(2,2))
plot(bestdiab)
extractAIC(bestdiab)
BIC(bestdiab)
```

###export model summaries to Latex###

```
xtable(beststroke)
xtable(bestdiab)
```

#####PCA Analysis #####

#PCA for 1990

```
pr1990<-prcomp(y1990[,5:15],scale=TRUE)
pr1990
summary(pr1990)
par(mfrow=c(1,1))
plot(pr1990,type="l") #Scree plot
biplot(pr1990,scale=0)
```

```
#PCA for 2000  
pr2000<-prcomp(y2000[,5:15],scale=TRUE)  
pr2000  
summary(pr2000)  
par(mfrow=c(1,1))  
plot(pr2000,type="l") #Scree plot  
biplot(pr2000,scale=0)
```

```
#PCA for 2010  
pr2010<-prcomp(y2010[,5:15],scale=TRUE)  
pr2010  
summary(pr2010)  
par(mfrow=c(1,1))  
plot(pr2010,type="l") #Scree plot  
biplot(pr2010,scale=0)
```

```
plot(modelCan)#PCA for 2010  
pr2010<-prcomp(y2010[,5:15],scale=TRUE)  
pr2010  
summary(pr2010)  
par(mfrow=c(1,1))  
plot(pr2010,type="l") #Scree plot  
biplot(pr2010,scale=0)
```

```
##All Years  
prAll<-prcomp(all_years[,5:15],scale=TRUE)  
prAll  
summary(prAll)  
par(mfrow=c(1,1))  
plot(prAll,type="l") #Scree plot  
biplot(prAll,scale=0)
```

```
#####End
```