

# ST404 Assignment 2

Frank Or, Remos Gong, Sam Glanfield, Thomas Broadbent

2022-03-09

## Contents

<b>1 Findings</b>	<b>2</b>
1.1 Key Findings . . . . .	2
1.2 Summary of EDA . . . . .	3
1.3 Modelling Approach . . . . .	3
1.4 Limitations of Analysis . . . . .	4
1.5 Further Improvements to the Model . . . . .	5
1.6 Areas that do not conform the general pattern . . . . .	5
<b>2 Statistical Methodology</b>	<b>6</b>
2.1 Outcomes of EDA . . . . .	6
2.2 Modelling Approach and Variable Selection . . . . .	7
2.3 Statistical Interpretation and Validation . . . . .	12
<b>3 Author Contributions</b>	<b>14</b>
<b>4 References</b>	<b>14</b>
<b>5 Appendix</b>	<b>15</b>
5.1 Outcomes of EDA . . . . .	15
5.2 Modelling Approach and Variable Selection . . . . .	15
5.3 Statistical Interpretation and Validation . . . . .	18

# 1 Findings

## 1.1 Key Findings

Given our preferred model the most important predictors for cancer death rates for a given state were incidence rates, median income and private health care coverage. All predictors for the model are in the table below with the parameter estimate and a scaled model to show importance of predictors. This shows that `incidenceRate` and `log medianIncome` followed by `PctPrivateCoverage` have the largest influence over death rates. The table shows in the first row the model unscaled parameters and in the second the scaled model parameters, the larger the values the more the predictor affects death rates

Intercept	IncidRT	PctEmp	PctUnemp	PrivCov	EmpPrivCov	Edu	lPctBlack	lMedInc
470.47	0.226	-0.251	0.291	-0.315	0.212	-10.056	0.904	-32.388
0.00	0.426	-0.078	0.036	-0.130	0.088	-0.082	0.057	-0.284

Understandably cancer diagnoses are a key predictor for death rates, with an increase in the number of cancer diagnoses by 50 per 100,000 increasing cancer death rates by 10 per 100,000. This of course does not imply that there should be a policy to decrease cancer diagnoses. Increasing cancer diagnoses will lead to more people with cancer on death certificates even if cancer may not have been the primary cause of death. Wealthier north western states have similar mean incidence rates to poorer south western states but much lower death rates showing that diagnoses are a good baseline additional predictors for poverty need to be included.

Due to the strong collinearity of many variables and thus our choice of regression method, we can not give a strong interpretation for keeping all but one variable the same and giving the associated change in death rates. An increased Death Rate is more likely if a county has a low median income, low private coverage, these though will be dependent on low employment and high unemployment percentages which may be because of a low education rate. Due to many variables explaining the same information about death rates the LASSO model reduces the complexity, for example even though Percent Public Coverage is correlated well with death rates it can be left out of the model as it is strongly correlated with both percent private coverage and employee provided coverage.

## 1.2 Summary of EDA

### 1.2.1 Missing/Incorrect Values and Outliers

We identified 152 missing values in **Percent Employed 16 and Over**. As these values were Missing Completely at Random we imputed them back into our data set with a linear model. We also identified values in **Average Household Size** that were unreasonably small. We believed this is an error in data entry and we scaled them by 100. For outliers, We removed ‘Union County, Florida’ and ‘Williamsburg City, Virginia’ because of there high **Incidence Rates**.

### 1.2.2 Transformations

In order to fulfill the model assumptions, namely linearity, normality, and constant variance, we performed log transforms for two variables: **Median Income** and **Percent Black** which suffered the most among all variables. Although the transformations did not cure the problems, they result in improvements for both. For other variables, we deduced that they are sufficient for the model assumptions and hence did not perform any transformations for model simplicity.

## 1.3 Modelling Approach

We used Stepwise Regression, RIDGE Regression and LASSO Regression to build our model. We compared the outputs and the variables selections of these models. We also analyse the goodness of fits of these models using Leave-one-out cross validation,  $R^2$  statistics and residuals analysis. We did not include **Geography** as it is an id variable. **Binned Income** was not used in the model as more information was included in **Median Income**.

### 1.3.1 Stepwise Regression

We used **Bayesian Information Criteria(BIC)** in our analysis. This penalises additional parameters harder than **Akaike information criterion(AIC)** which agrees with parsimony in our modelling. We performed forward, backward and hybrid stepwise regression to account for local minimums. In all of the models that were generated we notice that we have groups of parameters that are similar and hence induce multicollinearity in the models. For example, we see both **Percent Married** and **Percent Married Households** in the models and as stepwise methods do not account for multicollinearity this causes increases in the variance of our coefficient estimates and makes the model sensitive to changes, thus reducing drastically the predictive power of the model.

In order to address the issue of multicollinearity we apply Ridge and Lasso techniques which account for this.

### **1.3.2 RIDGE and LASSO Regression**

The suggested model contains eight variables which is the simplest model among our approaches. Since the purpose of this report is to reveal patterns in the mortality rate, we believed it is best to choose a simple model for stronger explanatory power.

### **1.3.3 Final Model Choice and Diagnostics**

We used residuals plot and QQ plot to diagnose the stepwise regression model. Both plots are satisfactory, agreeing with model assumptions.

We also performed Leave-one-out cross-validation for this stepwise model, which gives  $R^2$  value of 0.4621. We did the same cross-validation for our Ridge and Lasso model. The Ridge model gives  $R^2$  value of 0.4436. The Lasso model gives  $R^2$  value of 0.4441.

The stepwise model fits the observed data the best among all our models. The Lasso and the Ridge model is similar in this sense. However, the Lasso model contains the least number of predictor variables and we believed this is the most important criterion. Thus from our analysis briefly outlined above and in the next section we recommend the LASSO model due to its mix of predictive and explanatory power as well as being a simple model.

## **1.4 Limitations of Analysis**

One limitation of using LASSO regression is that it does not completely handle the problem of multicollinearity, however it reduces the multicollinearity present in the model by a large amount. Often it selects one feature from a group of correlated features, which happens arbitrarily in nature. In our model, this happens on the three health care coverage variables, where Percent Public Coverage was dropped. However, in our case, this is not a severe problem as we focus more on the explanatory power than the predictive power. LASSO regression reduces our model to 8 predictor variables and this leads to a loss of predictive power. At the same time, the log transformed variables are difficult to interpret.

Another issue is that LASSO models are difficult to diagnose. Unlike stepwise regressions, it is difficult to produce hypothesis tests to check the statistical significance of the predictors.

## **1.5 Further Improvements to the Model**

In order to develop this model in the future and improve its capabilities we make the following suggestions. For a more complex model with a large number of predictors including the state as a variable will mean each state's unique qualities can be considered by this coefficient, but will greatly increase the complexity of the model. Another improvement is to add different health coverage costs for different states, as for example a low income state with cheap private insurance is likely to also have low death rates, but the model does not account for this.

## **1.6 Areas that do not conform the general pattern**

Similarly, we see from the average residual map that in particular the 3 states Utah (Average of -22), Idaho (Average of -19) and Colorado (Average of -18) all have high residual averages and as they are negative this implies that the model is overfitting for these states and thus these states don't conform to the general pattern. However, as the model is overfitting and we are estimating death rate this isn't a cause for concern as we know the average death rates for these states is highly likely to be lower than what the model predicts. It is also important to note that these all lie in a similar region and in general we see that the majority of states on the west coast have an average negative residual, so the model overfits for these states, whereas on the east coast and central we see positive average residuals, the model is underfitting these states. With the main states of concern for overfitting with an average residual of 14 each for Oklahoma and Arkansas and an average of 13 in the District of Columbia.

## 2 Statistical Methodology

### 2.1 Outcomes of EDA

#### 2.1.1 Missing Or Incorrect Values

We saw counties with missing values in Percentage Employed 16 and Over and we concluded that the data is Missing Completely at Random. In order to rectify this we impute this data by fitting a linear regression model of Percentage Employed 16 and Over on the remaining variables to estimate what these values would be. (See Appendix 2.1.1).

For counties with Average Household Size less than one we took the decision to scale these values by 100 and kept them in the data set. This fixed the normality of Average Household Size as shown in the histogram.

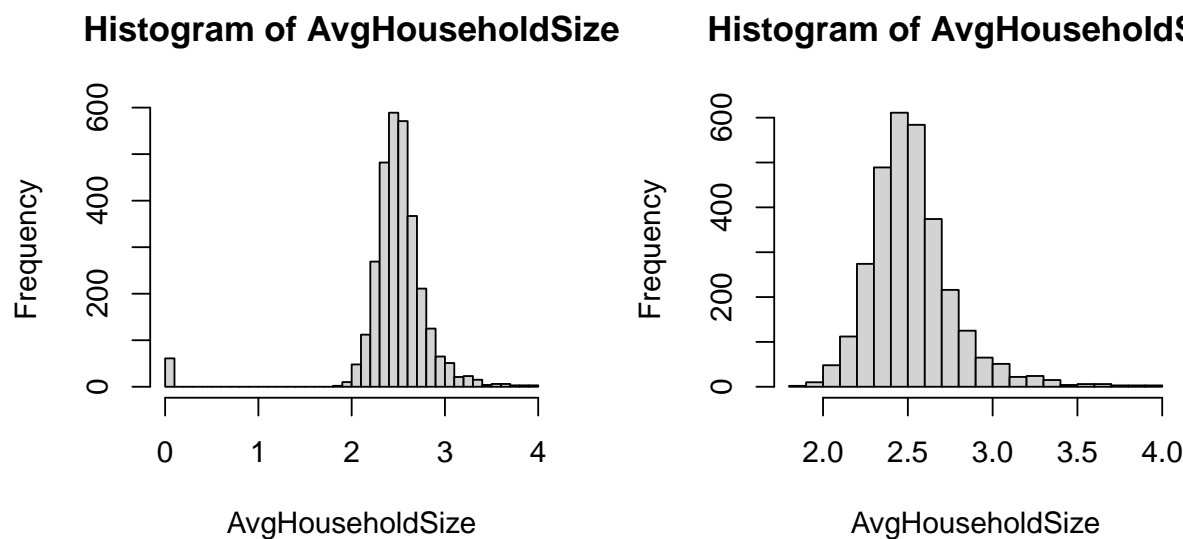


Figure 1: Histograms before(Left) and after(Right)

#### 2.1.2 Outliers

We found counties with high incidence rates, namely **Union County, Florida** and **Williamsburg City, Virginia**. We looked into the cook's distance plots (Fig.2) and noticed only 'Williamsburg City, Virginia' has large cook's distance and hence influential. The first cook's distance plot used a linear model with only incidenceRate as the predictor variable. The second used all the numerical variables. We concluded although **Union County, Florida** has high leverage, it is not influential

and hence should be kept in our data set. (See Appendix 2.1.2)

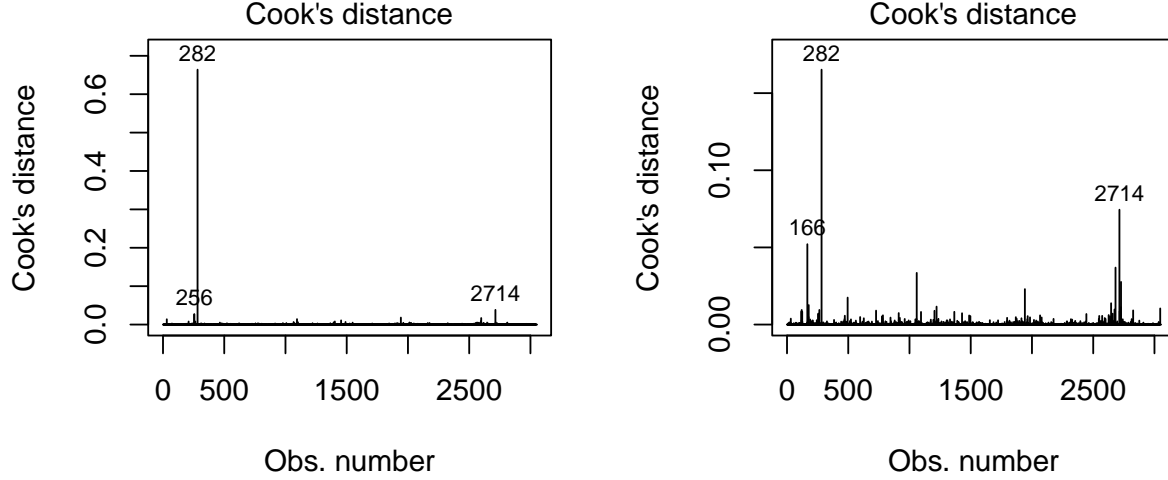


Figure 2: Cook's distance plots of models with only incidence Rate(Left) and all variables(Right)

### 2.1.3 Transformations

We transformed Percent Black by first shifting the values upwards by 0.05, to ensure we have no zero values, then take a log transform. We also transformed Median Income by again taking a log transformation. We performed these transformations to ensure the data was not heavily skewed and allowed for a more accurate model. (See Appendix 2.1.3)

The following residual plots (Fig.3) show the improvements in homoscedasticity in Percent Black and Median Income after log-transform respectively.

## 2.2 Modelling Approach and Variable Selection

### 2.2.1 AIC and BIC Forward and Backward Variable Selection

We performed forward, backward and hybrid stepwise regression according to BIC to ensure a more parsimonious model. We saw that the models generated for BIC are the same. It has 11 variables and the estimates are shown in Table 1.

For the model generated we can see from the summary output (See Appendix 5.2.1), with strong evidence, that all the coefficients are different from zero. From the output we would expect

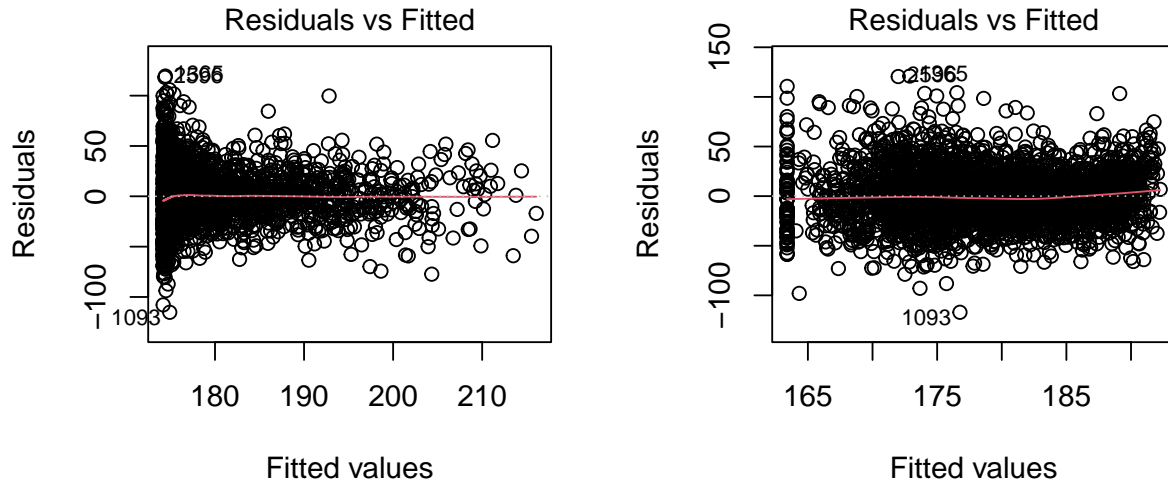


Figure 3: Residuals plot Percent Black(Left) and Log Percent Black(Right)

multicollinearity in these models due to variables that measure similar or opposite quantities, for example Percent Married and Percent Married Households.

We computed the Variance Inflation Factors (VIF) for this model and found that there are numerous values that are at least 5 indicating multicollinearity in the predictors. In particular, we saw a VIF of at least 10 in Percent Married and a VIF of 8 in Percent Married Households. As well as in  $\log(\text{Median Income})$ , Percent Employed 16 and Over, Percent Employers Private Coverage and Percent Private Coverage. Due to this frequent multicollinearity between the predictors in the stepwise model, the stepwise model is not a suitable approach as multicollinearity was not taken into account and thus we proceeded with Ridge and Lasso Regression in the following sections, hoping for a model with reduced multicollinearity.

To complete the section on stepwise regression we performed leave-one-out cross-validation (in order to compute the  $R^2$  of the stepwise models and also the Root Mean Squared Error (See Appendix 5.2.1)). This allowed us to compare these models with the Ridge and Lasso models we generate in the next sections. Performing this we see in the stepwise model that we have an  $R^2$  value of 0.4623.



### 2.2.2 RIDGE Regression

To address the multicollinearity present in the data we assessed the viability of a Ridge Regression Model. For this we fitted all continuous variables as predictors using the `glmnet` package.

The trace plot shows that the Ridge penalisation method reduced many variables close to 0 but did not remove any variables from the model itself suggesting that these predictors are insignificant when it comes to prediction. (See Appendix 5.2.2)

Fitting a ridge model and using the value of  $\lambda$  one standard error further away from the minimum does not remove any of the terms from the model and gives the coefficients as seen in Table 1. We also include the parameter estimates when the parameters are scaled to highlight how significant the parameter is.

Performing a leave-one-out cross-validation gives an  $R^2$  statistic of 0.4438. Which is similar to the previous stepwise regression method and is more difficult to diagnose. Moreover, the Ridge model kept all variables which made a complex model which is against our preference of parsimony. Therefore we concluded that the Ridge regression is not a suitable model.

### 2.2.3 LASSO Regression

We used cross-validation from the `glmnet` library which leaves out a 10th of the data every time. We produced the following plot (Fig. 4) of mean-squared error against  $\log(\lambda)$ . We decided to use the 1-standard-error- $\lambda$  because this likely to shrink some predictor variables to zero, performing variable selection. We prefer a simpler model. (See Appendix 5.2.3)

We produced a plot (Fig.5) that shows a trace of each parameter estimate for different values of  $\log(\lambda)$ . We scaled the model for a better visual interpretation.

The Lasso regression using 1-standard-error- $\lambda$  produced the following parameter estimates. Alongside we see the parameter estimates for Ridge regression and stepwise regression to allow for easier comparison of the models.

### 2.2.4 Interpretation

From the estimates output above, the Lasso regression eliminates a considerable number of predictor variables. The Lasso regression model has the following non-zero predictor estimates: **Incidence Rate, Percent Employed 16 and Over, Percent Private Coverage, Percent Employer**

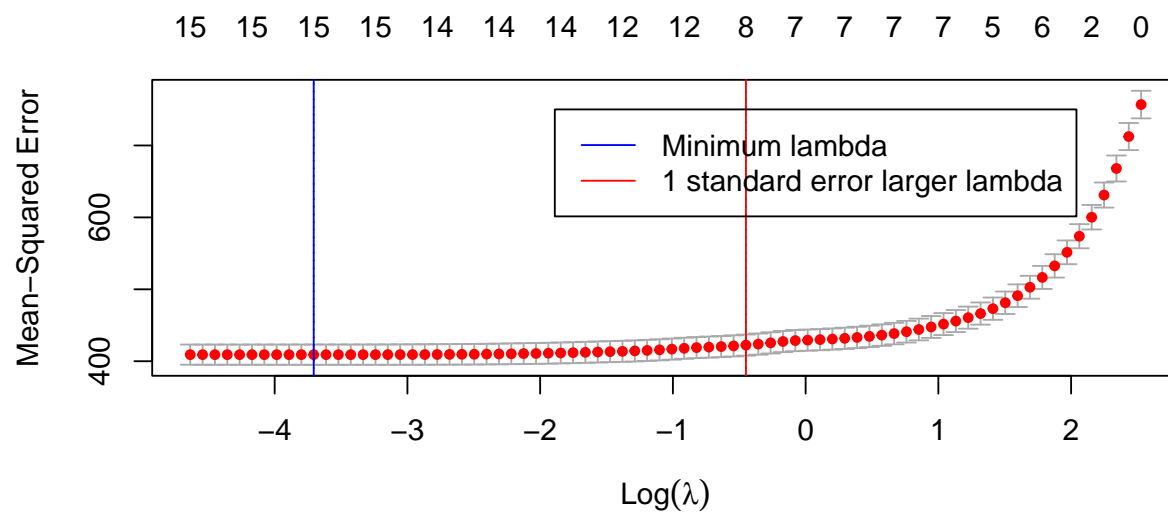


Figure 4: Mean-squared error against  $\log(\lambda)$

(#fig:Lasso cv)

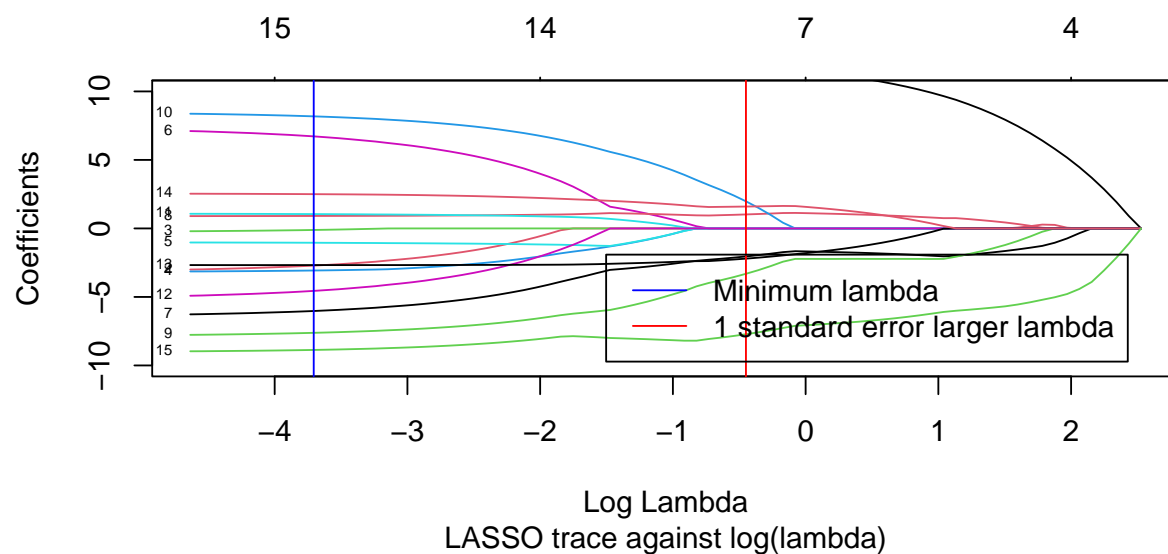


Figure 5: Trace of parameter estimates for different  $\lambda$

(#fig:Lasso trace)

Table 1: Parameter Estimates of the Models

	Ridge	LASSO	Stepwise
(Intercept)	353.941	470.470	507.558
incidenceRate	0.184	0.226	0.227
povertyPercent	0.228	0.000	0.000
MedianAgeMale	-0.125	0.000	0.000
MedianAgeFemale	-0.214	0.000	-0.421
AvgHouseholdSize	-5.098	0.000	0.000
PercentMarried	0.177	0.000	1.285
PctEmployed16_Over	-0.299	-0.251	-0.821
PctUnemployed16_Over	0.458	0.291	0.000
PctPrivateCoverage	-0.258	-0.315	-0.685
PctEmpPrivCoverage	0.258	0.212	0.899
PctPublicCoverage	0.282	0.000	0.000
PctMarriedHouseholds	-0.093	0.000	-0.973
Edu18_24	-11.649	-10.056	-11.839
logpctblack	1.060	0.904	1.405
logmedincome	-19.070	-32.388	-32.750

**Provided Private Coverage, Education Levels, Log of Percent Black, Log of Median Income.**

The two age variables and Average Household Size were removed from the model. This agreed with our EDA which showed that they have almost no correlation with death rate. Both employment variables and two out of three healthcare coverage variables are included in this model. We believe this is reasonable because although they showed evidence of collinearity/multicollinearity, we did not have strong arguments to remove any of them. Therefore, we agreed with the variable selection suggested by the Lasso Regression.

Most predictor estimates are smaller than one. However the estimates for Education Levels and Log Median Income are -10.05 and -32.39 respectively which are exceptionally large and negative. This means that they have much larger impact on the predicted death rates than other predictor variables. It makes sense that Median Income has the largest impact because patients in counties with higher median income tend to be able to afford better treatment which will reduce mortality rates.

We calculated the  $R^2$  statistics using leave-one-out cross-validation which is 0.4441.

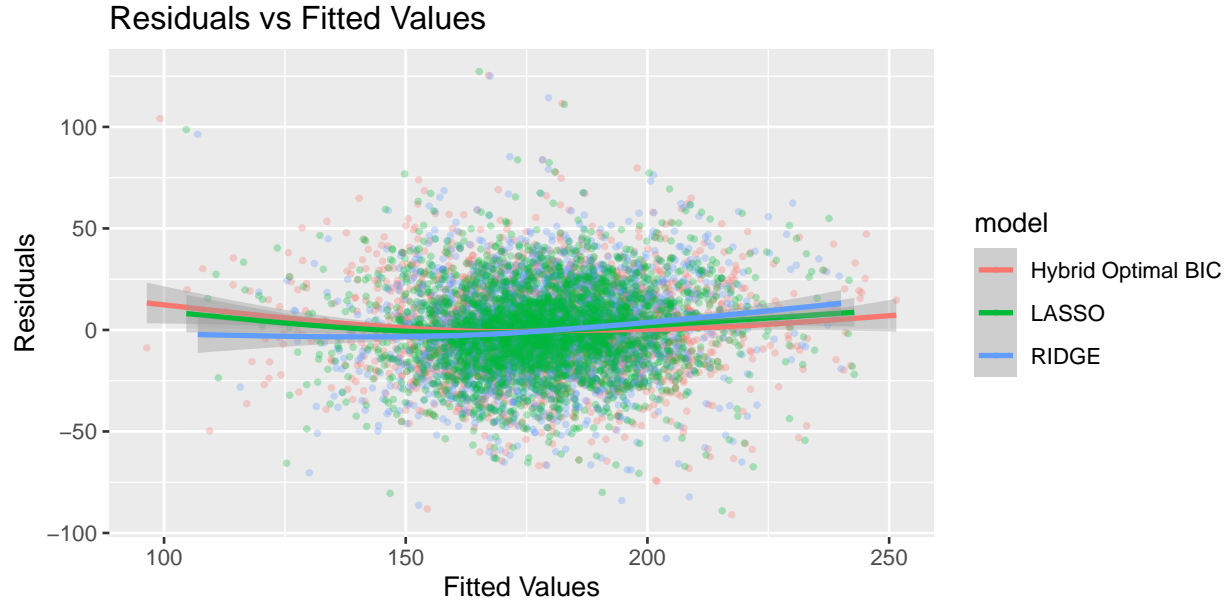
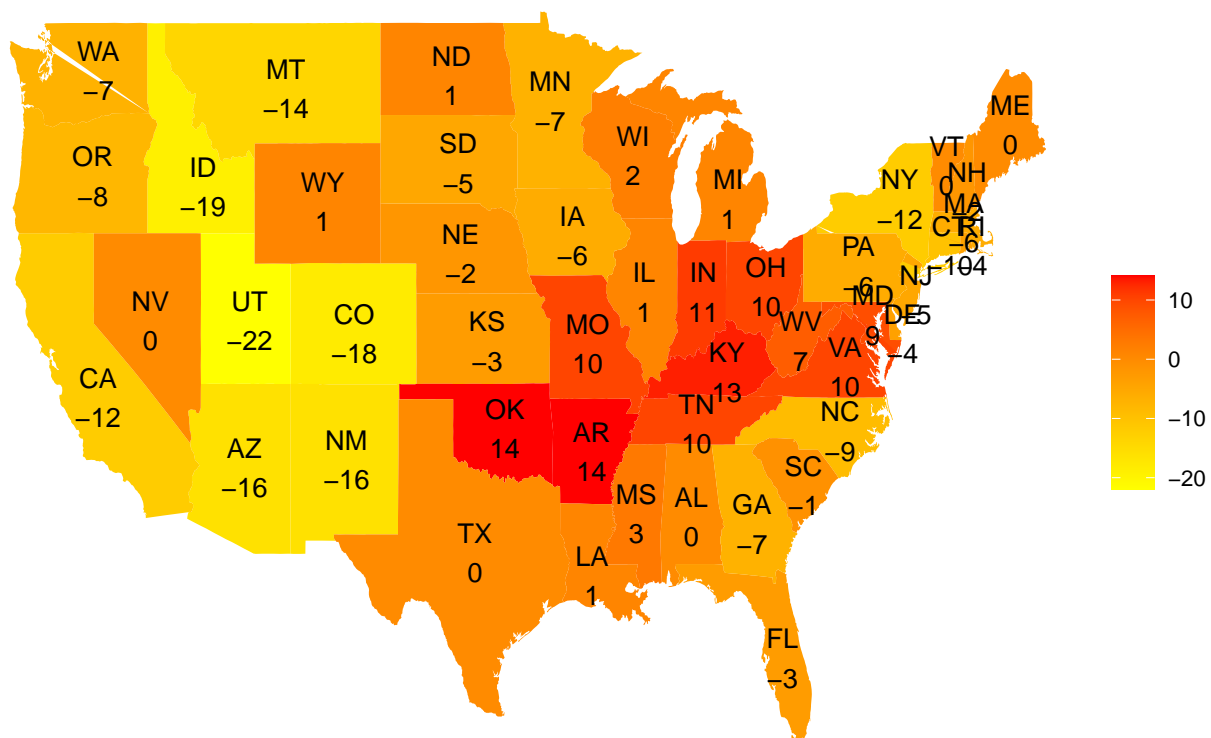


Figure 6: Residuals vs Fitted

### 2.3 Statistical Interpretation and Validation

We see from the plot above that all three models appear to exhibit constant variance with an even spread of points as well as having a mean of zero with the distributions symmetrical about zero. Thus satisfying our key linear regression assumptions of homoscedasticity and residuals having a mean of zero and thus providing further evidence that the LASSO model is a suitable model with respect to its residuals.

We also produced a heatmap of the US according to the average residual for each state to identify if there were any states/regions that did not conform to the general pattern of our model. As we can see and as highlighted in our findings we do see regions that do not conform and in general we see states further west are being overfitted by the model e.g. Utah, Idaho and Colorado all with high negative average residuals of -22, -19 and -18 respectively. We also see the central states being underfitted by the model for states such as Oklahoma and Arkansas with average residuals of 14 and 13 respectively. However, in general we see that the model fits the data relatively well with no large residuals other than those pointed out above and again gives us further evidence on the predictive and explanatory power of the LASSO model.



### **3 Author Contributions**

### **4 References**

## 5 Appendix

### 5.1 Outcomes of EDA

#### 5.1.1 Missing Or Incorrect Values

```
# Impute the missing data seen in the data set
mod1=lm(PctEmployed16_Over~incidenceRate+medIncome+binnedInc+povertyPercent+MedianAgeMale+MedianAgeFemale+AvgHouseholdSize
        +PercentMarried+PctUnemployed16_Over+PctPrivateCoverage+PctEmpPrivCoverage+PctPublicCoverage+PctBlack
        +PctMarriedHouseholds+Edu18_24,cancer)
missdf = cancer[which(is.na(cancer$PctEmployed16_Over)==TRUE),]
imputed = predict(mod1,missdf)
cancer[which(is.na(cancer$PctEmployed16_Over)==TRUE),"PctEmployed16_Over"] = imputed
# Scale average household sizes that are less than 1 by 100
cancer[which(cancer$AvgHouseholdSize < 1), ]$AvgHouseholdSize <-
  100*cancer[which(cancer$AvgHouseholdSize < 1), ]$AvgHouseholdSize
hist(cancer$AvgHouseholdSize, breaks=30, xlab="AvgHouseholdSize", main="Histogram of AvgHouseholdSize")
```

#### 5.1.2 Outliers

```
# Cook's distance Plot
par(mfrow=c(1,2))
plot(lm(deathRate ~ incidenceRate,data=cancer),4)
plot(lm(deathRate ~ .,data=cancer[-c(1,4)]),4)
# Removing outlier incidence rate 'Williamsburg City, Virginia'
cancer <- filter(cancer, incidenceRate <= 850)
```

#### 5.1.3 Transformations

```
# Log transforming the heavily skewed distributions of PctBlack and medIncome
cancer$logpctblack = log(cancer$PctBlack+0.05)
cancer$logmedincome = log(cancer$medIncome)
# Showing improvements in homoscedasticity in PctBlack
# Similar plots can be produced for medIncome
par(mfrow=c(1,2))
plot(lm(deathRate~PctBlack,data=cancer),1)
plot(lm(deathRate~logpctblack,data=cancer),1)
```

### 5.2 Modelling Approach and Variable Selection

#### 5.2.1 AIC and BIC Forward and Backward Variable Selection

```
# We perform stepwise regression for BIC
cancermodel <- cancer %>% select(
```

```

!c("Geography", "medIncome", "binnedInc", "PctBlack"))
c0=lm(deathRate~1,cancermodel)
cmax=lm(deathRate~.,cancermodel)
forwardoptimalBIC = step(c0,direction="forward",
scope=list("lower"=c0,"upper"=cmax),trace=0,k=log(3045))
backwardoptimalBIC = step(cmax,direction="backward",
scope=list("lower"=c0,"upper"=cmax),trace=0,k=log(3045))
hybridoptimalBIC = step(c0,direction="both",
scope=list("lower"=c0,"upper"=cmax),trace=0,k=log(3045))
# An example code for summary and coefficient of models
summary(hybridoptimalBIC)
coef(hybridoptimalBIC)
# Computing the VIF of the BIC stepwise regression model
vif(hybridoptimalBIC)
library(caret)
#specify the cross-validation method
#fit a regression model and use LOOCV to evaluate performance
loocv <- function(lm1, data=cancer) {
  ctrl <- trainControl(method = "LOOCV")
  xnam <- names(lm1$coefficients)[-1]
  fmla <- as.formula(paste("deathRate ~ ", paste(xnam, collapse= "+")))
  model <- train(fmla, data = data, method = "lm", trControl = ctrl)
  return(model)
}
hybridoptimalBIC.loocv <- loocv(hybridoptimalBIC,data=cancermodel)

```

## 5.2.2 RIDGE Regression

```

library(glmnet)
# Producing a 1-standard-error-lambda RIDGE regression model
# We also used similar codes for LASSO regression
cancer.dm <- cancer %>%
  select(!c("Geography", "medIncome", "binnedInc", "PctBlack", "deathRate")) %>%
  data.matrix()
lm.ridge <- glmnet(cancer.dm, cancer$deathRate, alpha=0)
lm.ridge.cv <- cv.glmnet(cancer.dm, cancer$deathRate, alpha=0)
# Create Model Plot Func
traceLogLambda <- function(lm1, lm1.cv, ylim=NULL, sub=NULL) {
  plot(lm1,"lambda",label = T, ylim=ylim)
  abline(v=log(lm1.cv$lambda.1se),col="red")
  abline(v=log(lm1.cv$lambda.min),col="blue")
  legend("bottomright",legend=c("Minimum lambda", "1 standard error larger lambda"),lty=c(1,1),col=c("blue","red"), ins=0.05)
  title(sub=sub)
}
traceLogLambda(lm.ridge, lm.ridge.cv, ylim=c(-20, 1), sub="RIDGE trace against log(lambda)")

```



```

lm.ridge.1se <- glmnet(cancer.dm, cancer$deathRate, alpha = 0,
                      lambda = lm.ridge.cv$lambda.1se)
lm.ridgefitted <- predict(lm.ridge.1se, newx=cancer.dm)
scatter.smooth(cancer$deathRate - lm.ridgefitted, x=lm.ridgefitted,
              xlab="Fitted", ylab="Residuals",
              sub="Residuals vs Fitted for RIDGE Model")

# Leave-one-out cross-validation for RIDGE Regression
# Calculate R-squared statistics
# We also performed the same procedures for LASSO Regression

library(parallel)
library(foreach)
library(doParallel)
numCores <- detectCores()
registerDoParallel(numCores)
n <- dim(cancer)[1]
dev.ratios <- rep(NA, n)
dev.ratios1 <- rep(NA, n)
lm.ridge.cv <- cv.glmnet(cancer.dm, cancer$deathRate, alpha=0, nfolds=n, parallel=TRUE)
lm.ridge.cv$lambda.1se
for (i in 1:n) {
  lm.ridge.1se <- glmnet(cancer.dm[-i,], cancer$deathRate[-i], alpha = 0,
                      lambda = lm.ridge.cv$lambda.1se)
  dev.ratios1[i] <- lm.ridge.1se$dev.ratio
}
mean(dev.ratios1)

##Produce column of table for coefficients table
library(knitr)

# This produces a table of parameter estimate for RIDGE model
# We also produced similar table for LASSO model and combined these with the stepwise coefficients into a 'kable'
a <- data.frame(sapply(round(coef(lm.ridge.1se), 3), FUN=identity))
colnames(a) <- "Parameter Estimates"
rownames(a) <- rownames(round(coef(lm.ridge.1se), 3))

```

## 5.2.3 LASSO Regression

```

# Produce a LASSO Regression model
# Produce a plot of mean-squared error against log(lambda)
set.seed((934))
lm.lasso <- glmnet(cancer.lasso.dm, cancer$deathRate, alpha=1)
lm.lasso.cv <- cv.glmnet(cancer.lasso.dm, cancer$deathRate, alpha=1)
plot(lm.lasso.cv)
abline(v=log(lm.lasso.cv$lambda.1se), col="red")
abline(v=log(lm.lasso.cv$lambda.min), col="blue")
legend("topright", legend=c("Minimum lambda", "1 standard error larger lambda"), lty=c(1,1), col=c("blue", "red"), lins=0.1)
lm.lasso.cv$lambda.1se

```

```
lm.lasso.1se <- glmnet(cancer.lasso.dm, cancer$deathRate, lambda = lm.lasso.cv$lambda.1se, alpha=1)
lm.lasso.1se.fitted <- predict(lm.lasso.1se, newx=cancer.lasso.dm)
lm.lasso.1se.residuals <- cancer$deathRate - lm.lasso.1se.fitted
# Produce a plot that shows the trace of each parameter estimate for different lambda
lm.lasso.scaled <- glmnet(scale(cancer.lasso.dm), cancer$deathRate, alpha=1)
traceLogLambda(lm.lasso.scaled, lm.lasso.cv, ylim=c(-10,10), sub="LASSO trace against log(lambda)")
```

## 5.3 Statistical Interpretation and Validation

```
# Create plot of Residual vs Fitted
library(ggplot2)
multResidualPlot <- function(residual.list, fitted.list, models) {
  a <- data.frame()
  for (i in 1:length(models)) {
    x <- fitted.list[[i]]
    y <- residual.list[[i]]
    model <- rep(models[i], times=length(x))
    df.temp <- data.frame(x, y, model)
    colnames(df.temp) <- c("x", "y", "model")
    a <- rbind(a, df.temp)
  }
  mrp <- ggplot(a, aes(x=x, y=y, colour=model)) +
    geom_point(alpha=0.3, size=0.75) + geom_smooth() +
    labs(x="Fitted Values", y="Residuals", title="Residuals vs Fitted Values", model="Models")
  return(mrp)
}
lm.ridge.residuals <- cancer$deathRate - lm.ridgefitted
residual.list <- list(hybridoptimalBIC$residuals, lm.ridge.residuals, lm.lasso.1se.residuals)
fitted.list <- list(hybridoptimalBIC$fitted.values, lm.ridgefitted, lm.lasso.1se.fitted)
mrp <- multResidualPlot(residual.list, fitted.list, c("Hybrid Optimal BIC", "RIDGE", "LASSO"))
```

### 5.3.1 Produce map of each states average residual value for the LASSO model.

```
#average for each state
averageresidual=round(tapply(lm.lasso.1se.residuals, cancermodel2$State, mean)[-c(2,12)],0)
states <- map_data("state")
averagedf <- data.frame(region=unique(states$region), averageresidual)
mergedf <- merge(states, averagedf, by="region")
statenames <- data.frame(region=tolower(state.name), clong=state.center$x, clat=state.center$y)
statenames <- merge(statenames, averagedf, by="region")
statenames$lab <- paste(c(state.abb, 'DC')[match(statenames$region, c(tolower(state.name), 'District of Columbia'))], '\n', state.name)
#Produce heatmap of average residual
qplot(long, lat, data=mergedf, geom="polygon", fill=averageresidual, group=region) +
  scale_fill_gradient(averageresidual, low="yellow", high="red") +
```

```
geom_text(data=statenames,aes(clong,clat,label=statenames$lab,inherit.aes = FALSE,label.size=0.001)) +  
  theme_void() + theme(legend.title = element_blank(),plot.title = element_text(hjust = 0.5))  
+ ggtitle("Heatmap of Average Residuals of US States")
```