

IE 330 Project 2
Due date: Oct 31st, 2020 at 11:59 pm EST

Group Name: Team 1

Group Members and PUIDs: Kyla Hardy(0030134630), Matthew Lanum (0031373892), Namita Mekala (003072819), Ryan Pfanstiel (0030444959), Zaid Qubain (0030730474)

Introduction

Our goal in this project was to create a multiple linear regression model to estimate the number of COVID-19 deaths in the summer by state. To generate this model we gathered the Proj2_data dataset which contains a row for each of the 46 states we could find data on, and columns for statistics on each state such as the number of hospital beds, the median salary, total deaths in the spring, etc. Links to the resources we used to compile this dataset can be found in the Resources section at the bottom of this report. The goal of our model was to predict the summer_tot_deaths column which was the total number of deaths in the summer due to COVID 19. Below is a list of the potential covariates used.

Covariates Used

Variable	Description
num_poverty	The number of people below the federal poverty threshold in each state
num_hosp_beds	The number of hospital beds available in each state.
num_lungs_death	The number of lung related deaths in 2019
Pop_abv_65	The number of people in each state older than 65
media_salary	The median_salary in each state
population	The population of each state
pop_density	The population density for each state
recovered	The total people that recovered by the beginning of summer
Spring_death	The total deaths during spring. This helps with adjusting the data since not all states got their infection at the same time.

Model Creation

In order to select the best possible model, we used the both the stepwise, and best subsets approaches.

Stepwise Regression

```
#model without variable selection
model1 <- lm(summer_tot_deaths~., data = specificregoin)

#stepwise variable selection
summary(step(model1, direction="both", scope=~.))

## Call:
## lm(formula = summer_tot_deaths ~ Spring_death + num_poverty +
##     num_hosp_beds + Pop_abv_65 + media_salary, data = specificregoin)
```

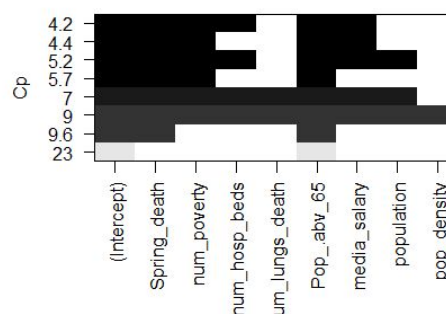
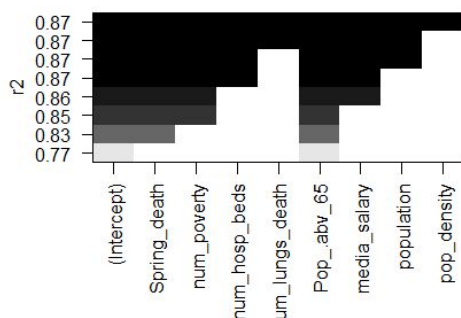
```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1238.9  -406.1  -179.4   297.1  2249.3
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.454e+03  8.140e+02  -1.786  0.082005 .
## Spring_death   1.627e-01  4.538e-02   3.586  0.000943 ***
## num_poverty    -3.209e-04  1.208e-04  -2.657  0.011478 *
## num_hosp_beds   1.277e-02  8.356e-03   1.528  0.134867
## Pop_.abv_65     1.401e+00  1.195e-01  11.730  3.39e-14 ***
## media_salary    1.650e-02  8.464e-03   1.950  0.058616 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 757.1 on 38 degrees of freedom
## (4 observations deleted due to missingness)
## Multiple R-squared:  0.8702, Adjusted R-squared:  0.8531
## F-statistic: 50.94 on 5 and 38 DF, p-value: 8.029e-16
```

Based on this method the selected variables will be Spring_death, num_poverty, num_hosp_beds, pop_.abv_65, and media_salary. Next we will try the best subsets approach.

Best Subsets Regression

```
#best subsets variable selection
best_subsets <- regsubsets(summer_tot_deaths~., data = specificregoin)
plot(best_subsets, scale = "r2")

plot(best_subsets, scale = "Cp")
```



Judging by the plots shown above, the model with 6-variables and 5-variable has the highest R-squared and model with 5- variables as well has the highest Cp.

Best Subsets Options

Type	R-Square	Mallows's Cp
6-variable	0.87	5.7
5-variable	0.87	4.2
4-variable	0.86	4.4

The 6-variable model has a high R-squared and high Cp. The 4-variable model has a lower Cp than 6-variable but also has lower R-squared. The 5-variable has the highest R-squared and the lowest Mallows's Cp and is therefore the best model between the options

Final Variable Selection

Regression subset and stepwise regression both favored the 5-variable model. Our final model includes the variable Spring_death, num_poverty, num_hosp_beds, Pop_.abv_65, and media_salary.

Final Model Generation and Summary

```
#Selected model
model1 <- lm(formula = summer_tot_deaths ~ Spring_death + num_poverty +
              num_hosp_beds + Pop_.abv_65 + media_salary, data = specificregoin)

#summary
summary(model1)

##
## Call:
## lm(formula = summer_tot_deaths ~ Spring_death + num_poverty +
##     num_hosp_beds + Pop_.abv_65 + media_salary, data = specificregoin)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1238.9   -406.1   -179.4    297.1   2249.3
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.454e+03  8.140e+02  -1.786  0.082005 .
## Spring_death   1.627e-01  4.538e-02   3.586  0.000943 ***
## num_poverty   -3.209e-04  1.208e-04  -2.657  0.011478 *
## num_hosp_beds  1.277e-02  8.356e-03   1.528  0.134867
## Pop_.abv_65    1.401e+00  1.195e-01  11.730  3.39e-14 ***
## media_salary   1.650e-02  8.464e-03   1.950  0.058616 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 757.1 on 38 degrees of freedom
## (4 observations deleted due to missingness)
## Multiple R-squared:  0.8702, Adjusted R-squared:  0.8531
## F-statistic: 50.94 on 5 and 38 DF, p-value: 8.029e-16
```

Question 1: Is the regression significant?

According to the F-test, the overall regression model is significant with a p-value of 8.029e-16. This is less than $\alpha = 0.05$ and leads us to reject the null hypothesis that the fit of the intercept only model and predictor model are equal. Spring deaths, the population in poverty, and the population above 65 were found significant at an $\alpha = 0.5$ using t-tests. The low p-values means we reject the null hypothesis that the coefficients of these predictor variables are equal to zero, and therefore changes in the predictor variables are related to the changes in the response variable (COVID-19 deaths in summer). The other variables had p-values over 0.05. A 95% confidence interval on the coefficients can be found in the table below.

95% Confidence Interval on Coefficients

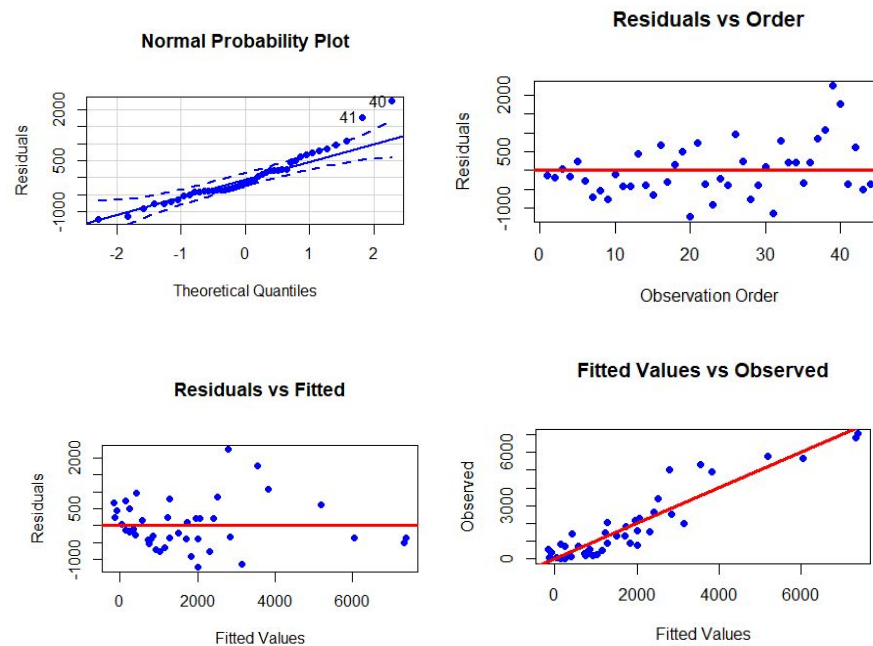
Coefficient	Lower Bound	Upper Bound
(Intercept)	-3.10E+03	1.94E+02

Spring_death	7.09E-02	2.55E-01
num_poverty	-5.65E-04	-7.64E-05
num_hosp_beds	-4.15E-03	2.97E-02
Pop_abv_65	1.16E+00	1.64E+00
media_salary	-6.32E-04	3.36E-02

Question 3: Is the model valid?

To check the validity of the model, we ran a series of tests on the residuals of the model in R. The output is shown below.

Model Diagnostic Plots



Analysis of Diagnostic Plots

- Normality:** The Normal Probability plot shows that most data points stay within the 95% CI. There are a few points outside of the CI interval, but not enough to invalidate the entire model. This upholds the normality assumption.
Fix: Take out the outliers.
- Independence of residuals:** There is a slight pattern across the x-axis of the Residuals vs Order plot, indicating the independence of residuals assumption is invalid.
Fix: Reduce the collinearity by removing variables that are strongly related to one another thus creating more independent residuals.
- Linearity and Homoscedasticity:** The Residuals vs. Fitted graph seems to have an N-shape pattern signalling a higher order relationship that has not been accounted for. This violates the

linearity assumption. In addition, a funnel shape is present on the Residuals vs. Fitted plot which indicates a violation of homoscedasticity.

Fix: Transform the data using a log, exponential, or other function.

Question 3: Does your model fit the data well?

The adjusted R-squared value of our model is 0.8531 which shows a good fit as it is above 0.85. In addition, the difference between Multiple R-squared (0.8702), and the adjusted R-squared is small thus showing minimal to no overfitting in our data.

Question 4: What is the takeaway message from your model?

Equation:

$$(0.1627 * \text{Spring_death}) - (0.0003209 * \text{num_poverty}) + (0.01277 * \text{num_hosp_beds}) + (1.401 * \text{Pop_abv_65}) + (\text{media_salary} * 0.0165) - 1454$$

Multicollinearity and Impact on inference:

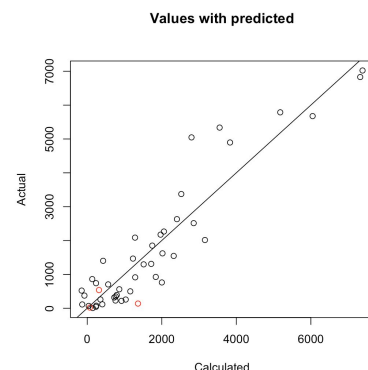
```
> vif(model1)
Spring_death  num_poverty num_hosp_beds  Pop_.abv_65  media_salary
    1.239035     1.234183     1.051372     1.454807     1.068912
```

All multicollinearity variables are below 10. This points to low inference which means the predictors can be trusted.

Testing the model

We left out 3 of the states when we created the model. We can now test our model on this set of data the model has not been trained on. The predicted (calculated) values are displayed as a line, and the actual values of the 3 new states are plotted in red.

We can see that our model's prediction is close to what the actual values are; this further validates the model.



We found our model was significant, and a good predictor of the COVID-19 related deaths in summer. This shows that the number of COVID-19 related deaths in a given state during the summer can be accurately predicted given the number of deaths in Spring in the state, the number of people in poverty in the state, the quantity of available hospital beds in the state, the population above 65 in the state, and the median salary within the state. Furthermore, due to the low VIF values, we can interpret the coefficients of the model to determine a negative or positive correlation with the COVID-19 related deaths in summer. However, because many of the coefficients have both positive and negative values for their coefficients within the 95% CI, this does not give much insight. Overall, our model is a good predictor of the number of COVID-19 related deaths for a state in summer.

References

- ahd. (2020, 5 10). *Hospital Statistics by State*. Retrieved from ahd:
https://www.ahd.com/state_statistics.html
- Amadeo, K. (2020, 9 17). *US Poverty Rate by State*. Retrieved from the balance:
<https://www.thebalance.com/us-poverty-rate-by-state-4585001>
- Atlantic, T. (2022, 10 30). *US Historical Data*. Retrieved from The Covid Tracking Project:
<https://covidtracking.com/data>
- CDC. (n.d.). *CDC*. Retrieved from CDC: <https://www.cdc.gov/nchs/fastats/copd.htm>
- PK. (2019, 12 12). *Average Income by State plus Median, Top 1%, and All Income Percentiles in 2020*. Retrieved from DQYDJ: <https://dqydj.com/average-income-by-state-median-top-percentiles/>
- PRB. (2019, 3 19). *Which U.S. States Have the Oldest Populations*. Retrieved from PRB:
<https://www.prb.org/which-us-states-are-the-oldest/>

File Description

Proj2_data.csv: Contains all the data used to calculate this model

Project2_FIN: Contains all the R code used with comments for clarification

zoom_5.mp4: is our extra credit video