

Drugs and Jobs: The effect of unemployment on drug overdose deaths in America

Evan Arnold and Caleb Ren

12/06/2019

Introduction and Motivation

Overdose deaths in the US have increased dramatically since 1999. Opioid use deaths in particular have reached epidemic levels, with a 200% increase in overdose deaths between 2000 and 2014 (Ghertner and Groves 2018). The Centers for Disease Control and Prevention estimates that over 60,000 drug overdose deaths occurred in 2016, three times the rate of drug deaths as 1999 (Hedegaard and Warner 2018)

Our analysis has two goals. First, we will determine the significance of unemployment when predicting overdose death rates, and validate this significance in the presence of other predictors. Second, we will determine which model from a set of possible, well-tuned models has the most predictive ability.

Exploratory Data Analysis Methods

Data Summary

As a next step in our project, we collected data from the CDC in the form of the Vital Statistics Rapid Release dataset (VSRR). The VSRR data contains provisional counts of drug overdose deaths in the US as reported by agencies from all 50 states and the District of Columbia. The data is collected on a monthly basis.

The data of import to this project is the number of deaths in each state as a result of drug overdose. Drug overdoses are counted by state agencies in accordance to World Health Organization standards, which lay out the basic guides for reporting agencies to code and classify causes of death. Drug categories that are represented in this dataset include the major drivers of the opioid epidemic like heroin (coded by T40.1), natural opioid analgesics (morphine and codeine), synthetic opioids (oxycodone, hydrocodone, oxycodone; T40.2), methadone (T40.3), other synthetics (fentanyl, tramadol; T40.4) and other drugs like cocaine, methamphetamine, etc.

There were over 26052 data points from the VSRR dataset. Of those data points, many are individual observation of different coded deaths from different drugs; after reshaping and data cleaning, there are now 2652 individual observations. The data ranges from January 1, 2015 to April 1, 2019, with each state reporting 52 observations (once per month). Overdose deaths range from 55 deaths in the month of May 2015 in South Dakota to a high of 5697 in Pennsylvania in September of 2017.

Unemployment data was sourced from the Bureau of Labor Statistics. Unemployment data is published in monthly increments from the Bureau of Labor Statistics by state. Data is published beginning in 1976 and is published on the first of each month describing the previous month's unemployment rate.

There is a very specific definition of who in the labor force is considered *unemployed*. According to the BLS, those who are currently unemployed are those who are "jobless, looking for a job, and available for work." People who are incarcerated, in a nursing home, or in a mental health care facility are not considered unemployed as they are not fit for work.

Using this definition, data was scraped from the BLS website and aggregated by each state and the District of Columbia. The unemployment rate in percent is given by the `unemployment` column. The lowest

unemployment rate in a given state and month is Vermont in 2019 with a 2.1% unemployment rate. The highest rate is DC in 2015 with a 7.4% unemployment rate. The data itself is roughly Normally distributed with a mean of 4.2% and a median of 4.31%.

Here we consider a series of variables from the Federal Reserve Bank of St. Louis (Federal Reserve Economic Data or FRED for short). All such datasets are originally from Bureau of Labor Statistics or the U.S. Census Bureau, however, FRED has their data in a more convenient format.

First, we consider the number of new private housing units authorized by building permit. This timeseries is reported monthly. In our analysis, this variable serves as a proxy to housing development as well as the health of the housing market; more permits implies a healthy housing market and ample housing options. We thus expect a negative association between overdose deaths and new housing permits.

Next, we consider the money spent on imports of manufactured and non-manufactured commodities in millions of dollars. This timeseries is reported monthly. This feature will offer our analysis an insight into the health of the state-local manufacturing sector. In the past few years, the states with some of the worst opioid overdose numbers are those in the Rust Belt (<https://www.drugabuse.gov/drugs-abuse/opioids/opioid-summaries-by-state>). One potential narrative is that the increased use of overseas labor strips entire communities of their jobs, and many of those effected turn to opioids. We can reasonably quantify this effect by state with the number of dollars spent on imports. We expect an increase in spending on imports to imply a decrease in availability of manufacturing jobs and the health of the local manufacturing sector, and thus we expect a positive association between import spending and overdose deaths.

Now, we consider the value of manufactured and non-manufactured exports in millions of dollars. This timeseries is reported monthly. We hope to use this variable as another perspective on the health of manufacturing and the availability of jobs in the manufacturing sector. We expect the value of exports to have a negative association with overdose deaths.

Finally, we consider per capita personal income in dollars. This timeseries is reported annually. It is reasonable to assume that the average personal income does not change dramatically intrayear. Per capita income affords us a glimpse into the personal financial freedom of state residents. We chose to include per capita income instead of median household income in order to capture the influence of the wealthy. We expect a negative association between income and overdose deaths.

Below, we include annual population estimates from the Census Bureau. We use annual data for two reasons: the availability of state population estimates is limited and we can reasonably assume that state population does not change dramatically intrayear. Furthermore, we use the population estimate from the previous year as a given year's population variable. We do so because population estimates are generated late into the year, and thus for any given year, the previous year's Census Bureau estimate is likely more accurate. We will use population to normalize the other predictors.

R inherently has several statistics about states. Most of these are static estimates from the early 70's (population estimates, income per capita, illiteracy,...). We can, however, use the datasets which give information about the region a state is in. This is a categorical variable with four categories: northeast, northcentral, south, and west. Though the nomenclature is dated, these categories reasonably divide states by potentially important factors. For example, we expect the Rust Belt (which R includes in the region northcentral) to have a positive association with overdose deaths (as noted above). Here we note that we have to encode Washington DC ourselves as the data in R does not contain the District of Columbia (in accordance with R's encoding, Washington DC is in the south).

R also provides us with the area of each state in square miles. Again, we must encode Washington DC manually (source: <https://www.britannica.com/place/Washington-DC>).

Before we begin any modeling or analysis, we must normalize overdose deaths, permits, imports, and exports. By doing so, we convert each to a rate which can be directly compared across states. In order to avoid working with very small numbers, we convert each rate to: overdose deaths per 100000 people, permits per 100000 people, spending on imports in millions of dollars per 100000 people, and value of exports in millions of dollars per 100000 people.

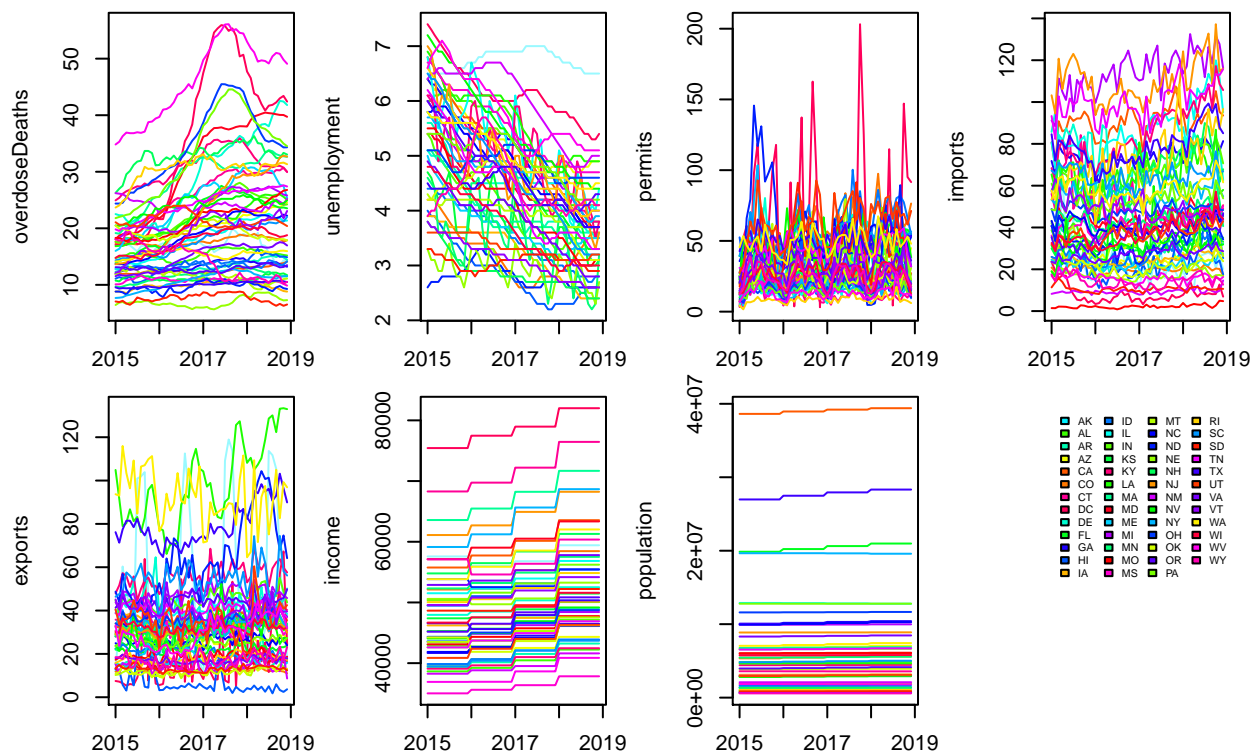
EDA

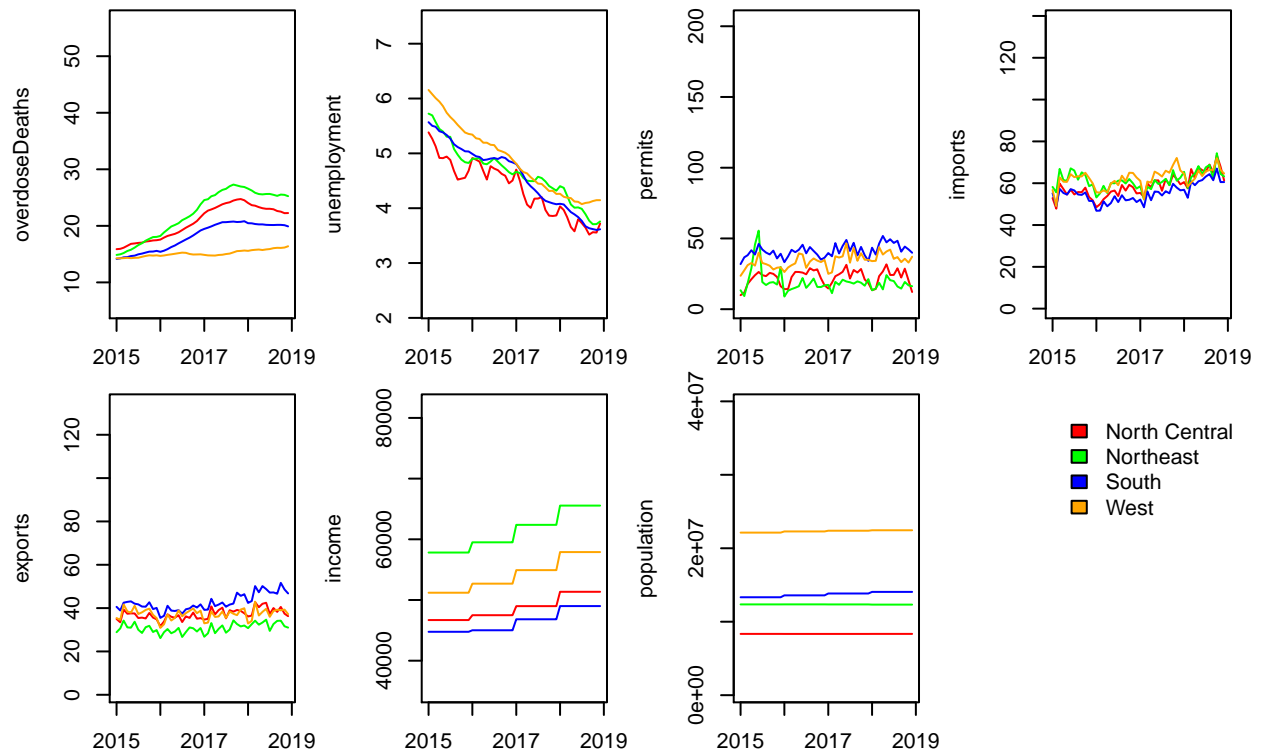
Below we check to see if any cells of the dataframe are NA, NULL, NaN, or infinite. No such values exist. We have a relatively clean dataset.

```
## [1] FALSE
```

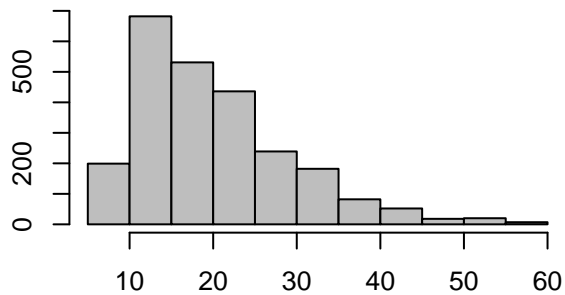
We now consider the possibility of multicollinearity. As shown below, some of the quantitative features of our dataset are highly correlated (population and imports, imports and exports,...). In order to handle this multicollinearity, we will consider ensemble models such as random forests.

```
##      unemployment    permits    imports    exports    income
## unemployment    1.00000000 -0.19668769  0.09159131  0.14169833 -0.08702747
## permits         -0.19668769  1.00000000 -0.08778809  0.10420801 -0.05360477
## imports          0.09159131 -0.08778809  1.00000000  0.44711152  0.07188101
## exports          0.14169833  0.10420801  0.44711152  1.00000000 -0.02178304
## income          -0.08702747 -0.05360477  0.07188101 -0.02178304  1.00000000
## population       0.15069205 -0.02456117  0.43294469  0.19201412  0.12074508
## area            0.25436018  0.02988105 -0.13252477  0.19145599 -0.04051548
##
##      population      area
## unemployment    0.15069205  0.25436018
## permits         -0.02456117  0.02988105
## imports          0.43294469 -0.13252477
## exports          0.19201412  0.19145599
## income          0.12074508 -0.04051548
## population       1.00000000  0.14993987
## area            0.14993987  1.00000000
```

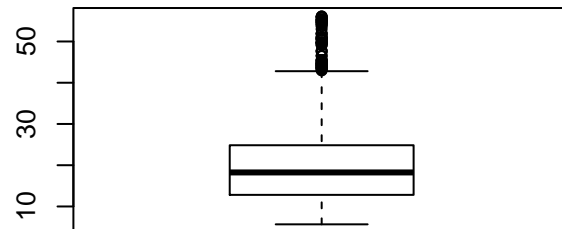




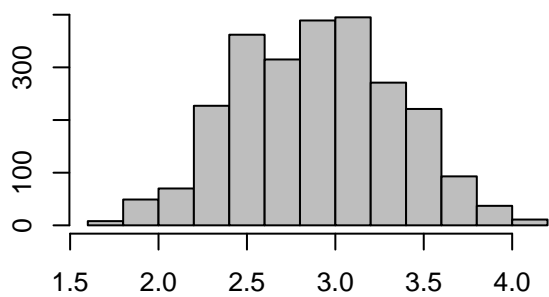
Histogram of Overdose Deaths



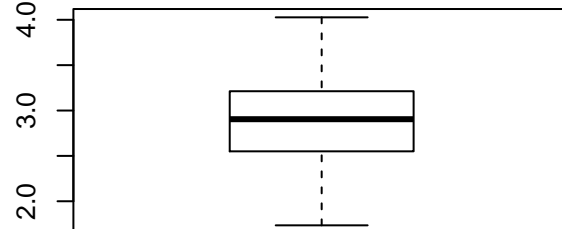
Boxplot of Overdose Deaths



Histogram of Log-Overdose Deaths

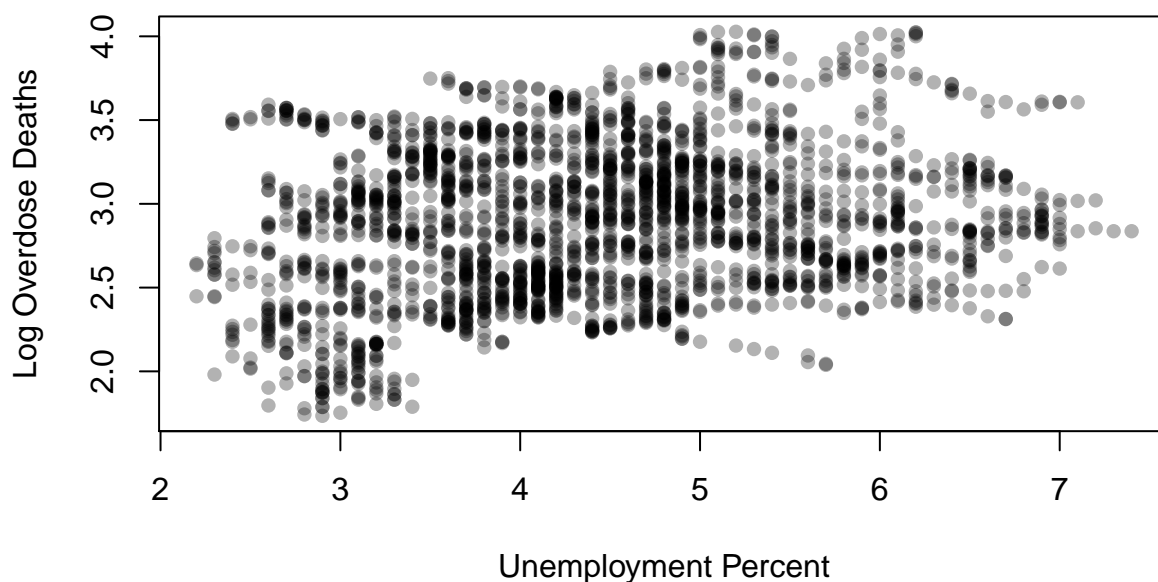


Boxplot of Log-Overdose Deaths



We see that the data is much closer to a Normal distribution if we apply a log transformation. We shall predict the log-rate of overdose deaths in our linear models.

Log Overdose Deaths vs. Unemployment

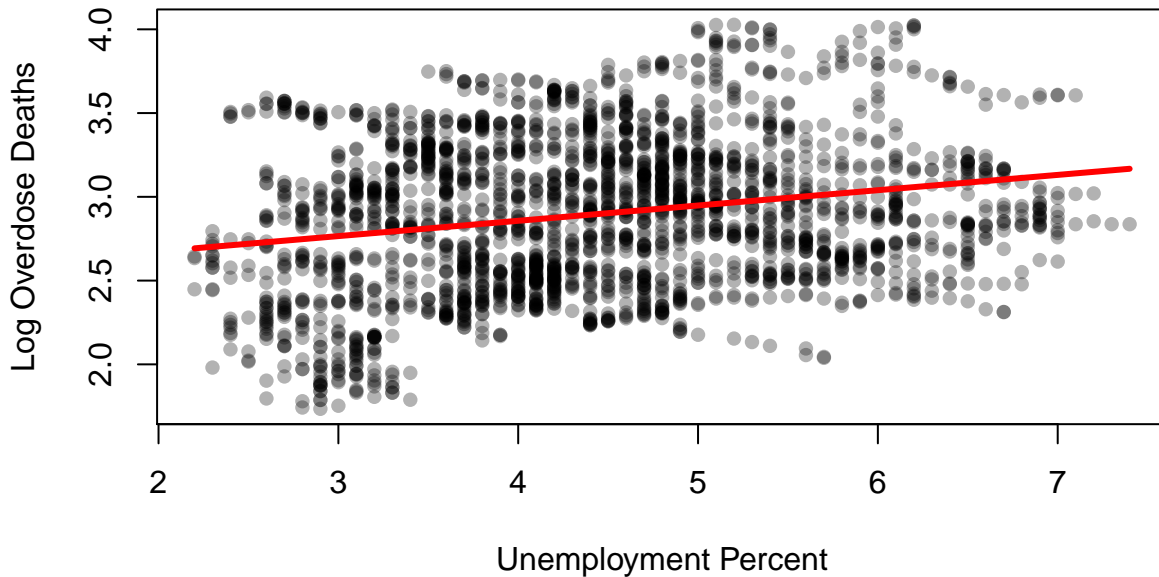


Baseline Model

```
##
## Call:
## lm(formula = log(overdoseDeaths) ~ unemployment, data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.02191 -0.36262 -0.00615  0.31632  1.06045
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.491641   0.041337  60.28  <2e-16 ***
## unemployment  0.091326   0.009067  10.07  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4377 on 1936 degrees of freedom
## Multiple R-squared:  0.04979,    Adjusted R-squared:  0.0493
## F-statistic: 101.4 on 1 and 1936 DF,  p-value: < 2.2e-16
```

The simple regression model has a positive coefficient for unemployment (0.09179). With a t -statistic of 10.11 (p -value ≈ 0), this coefficient is very significant. The model has a positive association between unemployment and overdose deaths.

Log-Overdose Deaths vs. Unemployment



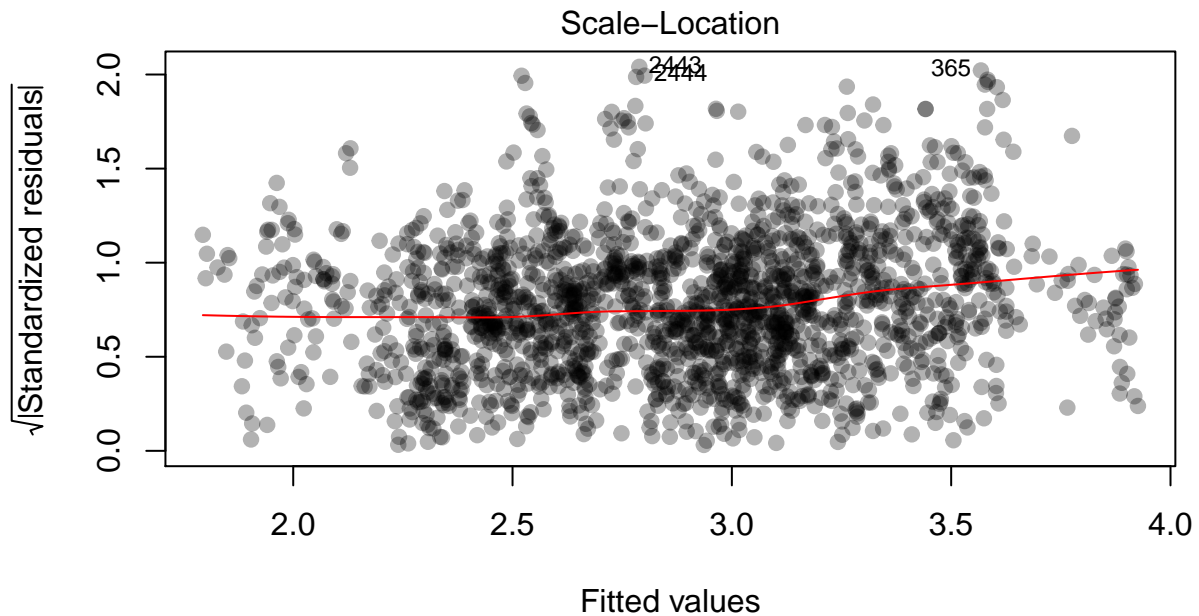
In order to better evaluate the importance of unemployment, we apply an ESS F-test. First we fit two linear models. The first model include all main effects except unemployment, while the second includes all main effects. Then, we fit two quadratic models. The first model contains all main effects and their respective quadratic variants except unemployment. The second model contains the same predictors, except that it includes unemployment and its quadratic effect. We again verify the assumptions of a linear model via a plot of residuals vs. fitten values.

```
## Analysis of Variance Table
##
## Model 1: log(overdoseDeaths) ~ state + month + year + permits + imports +
##   exports + income + population + region + area
## Model 2: log(overdoseDeaths) ~ state + month + year + permits + imports +
##   exports + income + population + region + area + unemployment
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1    1868 22.807
## 2    1867 22.634   1   0.17353 14.314 0.0001596 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

With an F-statistic of 13.387 (with 1887, 1 degrees of freedom) and a corresponding p-value < 0.001, unemployment provides significant predictive ability. We now consider the same ESS F-test with quadratic models.

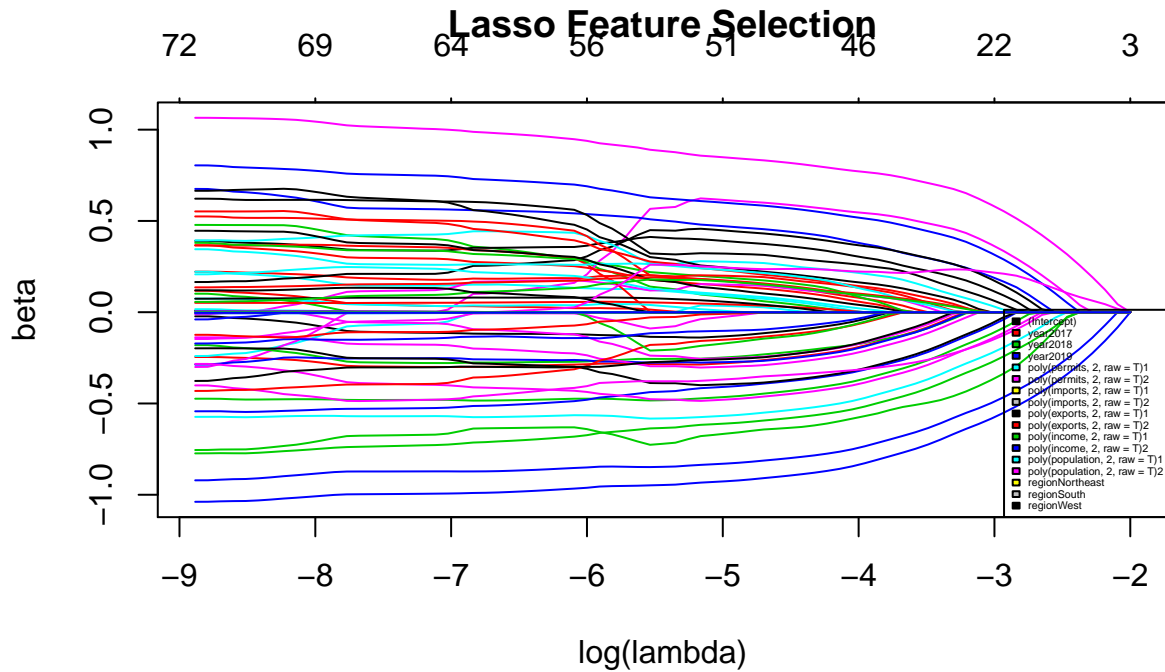
```
## Analysis of Variance Table
##
## Model 1: log(overdoseDeaths) ~ state + month + year + poly(permits, 2,
##   raw = T) + poly(imports, 2, raw = T) + poly(exports, 2, raw = T) +
##   poly(income, 2, raw = T) + poly(population, 2, raw = T) +
##   region + area
## Model 2: log(overdoseDeaths) ~ state + month + year + poly(permits, 2,
##   raw = T) + poly(imports, 2, raw = T) + poly(exports, 2, raw = T) +
##   poly(income, 2, raw = T) + poly(population, 2, raw = T) +
##   region + area + poly(unemployment, 2, raw = T)
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
```

```
## 1 1865 25.280
## 2 1863 24.347 2 0.93256 35.678 6.229e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



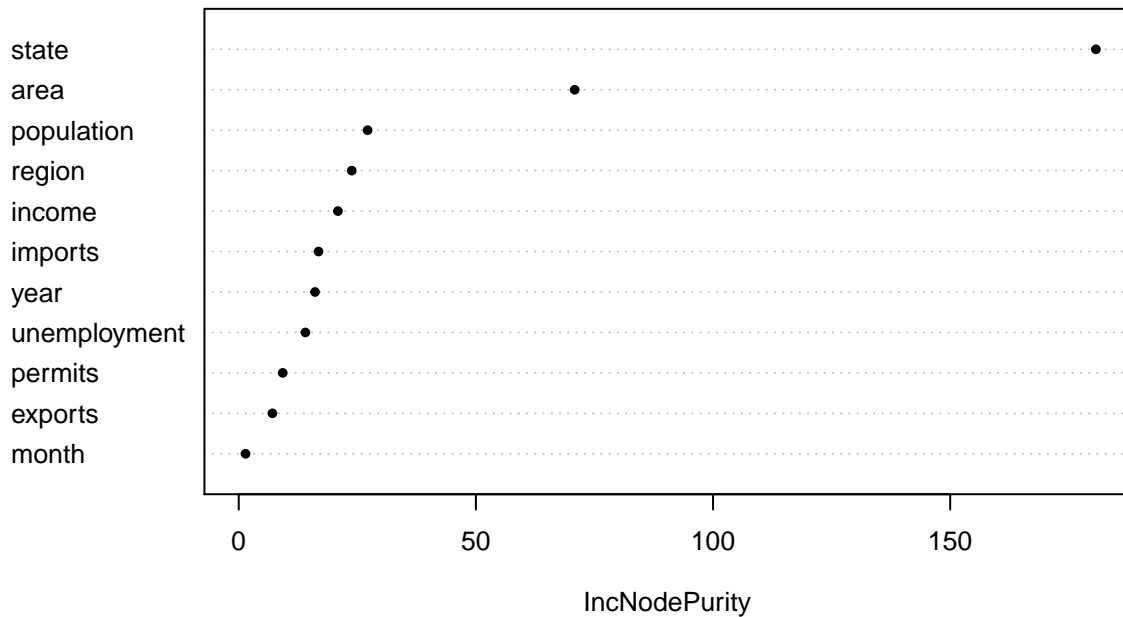
With an F-statistic of 17.359 (with 1881, 2 degrees of freedom) and a corresponding p-value < 0.001 , unemployment provides significant predictive ability. Additionally, there seems to be no underlying structure to the residuals, the residuals appear to have reasonably constant variance, and are reasonably normally distributed: the assumptions of our models (linearity, homoscedasticity of residuals, and normality) appear to be satisfied.

In order to be thorough, we examine the importance of unemployment in two additional ways: lasso regression, and random forest regression. We begin with lasso regression. Due to l1 penalty of lasso regression, less important variables quickly converge to zero as the regularization parameter increases. Below, we fit a lasso model (with the same design matrix as that of the full quadratic kitchen sink model) with a series of possible regularization parameters. We then consider the order in which the variable coefficients shrink to zero.



Lasso feature selection agrees with the results of the ESS F-test. Unemployment (the quadratic effect in this case) is one of the last coefficients to shrink to zero along with population. We further test this conclusion with a random forest model. Below, we fit a random forest model with all main effects and then examine its relative feature importance.

Random Forest Feature Importance



Unemployment is relatively not important in the random forest model. It appears as though state and area (which is correlated with state) are the most important predictors. In order to handle the grouped nature of time, state, area, and so forth, we fit a mixed effects model.

Below, we fit a three-layered mixed-effects model. The first two layers are state and year respectively. Theoretically, we would prefer to fit a four-layered model which includes month below year. This is impossible

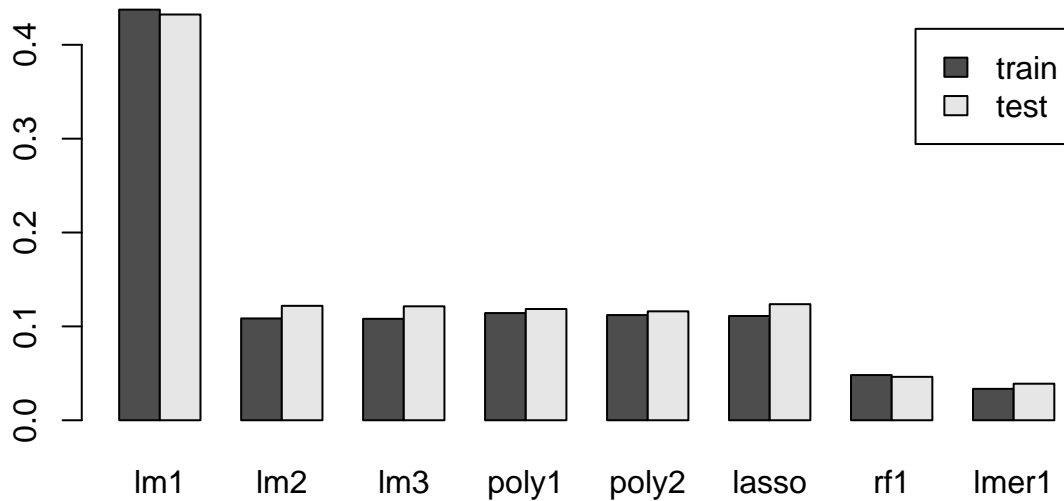
with our dataset ($n = 2448$). With 51 states, 5 years, and 12 months, we would need a minimum of $51 \times 5 \times 12 = 3060$ samples in order to fit such a model. We include a random intercept for states and years. The only random effect we include in the model is unemployment. These are purposeful design choices considering the limited size of our dataset. More complex models would not be reasonably fit.

```
## Analysis of Variance Table
##           Df    Sum Sq   Mean Sq F value
## month      11  0.305188  0.0277444  20.0454
## permits     1  0.003387  0.0033866   2.4469
## imports     1  0.002240  0.0022399   1.6183
## exports     1  0.000181  0.0001809   0.1307
## income      1  0.012506  0.0125057   9.0354
## population  1  0.002036  0.0020364   1.4713
## region      3  0.009944  0.0033146   2.3948
## area        1  0.003675  0.0036749   2.6551
## unemployment 1  0.000444  0.0004445   0.3211
```

We are not able to calculate p-values for the anova table (the statistics from a mixed model are not F-distributed), however, we can reasonably say that month is the most significant predictor followed by per capita income and region. Unemployment is relatively significant, however, we can show its significance using a p-value.

Results

```
##           train      test
## lm1    0.43747067 0.43228048
## lm2    0.10848281 0.12190404
## lm3    0.10806934 0.12138662
## poly1  0.11421185 0.11852647
## poly2  0.11208546 0.11606524
## lasso  0.11109847 0.12365852
## rf1    0.04815521 0.04625506
## lmer1  0.03348137 0.03896440
```



Conclusions and Decisions

The Significance of Unemployment:

Unemployment is a significant predictor of overdose death rates. Our simple linear model agrees with this, however, our goal is to show that this relationship holds when other predictors are considered. We thus fit two kitchen sink linear models: one which includes unemployment, and one which does not. An ESS F-test showed that unemployment provides significant predictive power. We corroborated this result with the same analysis on two quadratic kitchen sink models (again, one included unemployment while the other did not). The resulting ESS F-test agreed with that of the linear models.

We further explored this result with two additional models: a lasso regression model and a random forest regression model. The lasso model showed that unemployment is one of the later predictors to shrink to zero. This corroborates its significance. The random forest model, however, gave much more importance to state, and other non-monthly variables. We can control for these variables with a mixed effects model.

Our mixed effects model controlled for state and year grouping (with random intercepts) as well as included a random slope for unemployment. The anova table of the resulting model (though limited due to data constraints) showed that unemployment is indeed likely significant.

Prediction:

The least predictive model is the simple linear model. The rest of the linear models (including the quadratic and lasso-regularized models) performed very similarly on the test and train sets. This indicates that overfitting is not a serious concern for any of our models. The random forest model and the mixed effects model performed the best reducing RMSE by more than half that of the multiples linear models.

Direction of Future Research

Sample Size:

Month is the most important fixed effect in the mixed model. The predictive ability of the model may increase if month is included as a grouping variable. In our analysis, this was not possible due to limited sample size. In the future, it will be possible to fit such a model.

Predictors:

The insignificance of many of our predictors in the mixed model indicates that we are likely missing important variables for predicting overdose death rates. Some potential predictors are: average temperature, crime rate, high school graduation rate, and literacy rate among others. As of now, these predictors are either not readily available online, not current, or not available in a reasonably frequent timeserie. As more data comes out, variables such as these may add significant predictive power.

References

- Ghertner, R., and L. Groves. 2018. "The Opioid Crisis and Economic Opportunity: Geographic and Economic Trends." *Office of the Assistant Secretary for Planning and Execution* 24 (September).
- Hedegaard, A. M., H. M.D. Miniño, and M. Ph.D. Warner. 2018. "Drug Overdose Deaths in the United States, 1999 - 2017." *Centers for Disease Control and Prevention* 329 (November).