

Lab 3: Panel Models

US Traffic Fatalities: 1980 - 2004

Contents

1	U.S. traffic fatalities: 1980-2004	1
2	(30 points, total) Build and Describe the Data	2
2.1	Part 1: loading data and creating features	2
2.2	Part 2: Dataset Description	3
2.3	Part 3: EDA	3
3	(15 points) Preliminary Model	16
4	(15 points) Expanded Model	17
5	(15 points) State-Level Fixed Effects	19
6	(10 points) Consider a Random Effects Model	22
7	(10 points) Model Forecasts	26
8	(5 points) Evaluate Error	28

1 U.S. traffic fatalities: 1980-2004

In this lab, we are asking you to answer the following **causal** question:

“Do changes in traffic laws affect traffic fatalities?”

To answer this question, please complete the tasks specified below using the data provided in `data/driving.Rdata`. This data includes 25 years of data that cover changes in various state drunk driving, seat belt, and speed limit laws.

Specifically, this data set contains data for the 48 continental U.S. states from 1980 through 2004. Various driving laws are indicated in the data set, such as the alcohol level at which drivers are considered legally intoxicated. There are also indicators for “per se” laws—where licenses can be revoked without a trial—and seat belt laws. A few economics and demographic variables are also included. The description of the each of the variables in the dataset is also provided in the dataset.

```
load(file="driving.RData")

## please comment these calls in your work
# glimpse(data)
# desc
```

2 (30 points, total) Build and Describe the Data

- (5 points) Load the data and produce useful features. Specifically:
 - Produce a new variable, called `speed_limit` that re-encodes the data that is in `sl55`, `sl65`, `sl70`, `sl75`, and `slnone`;
 - Produce a new variable, called `year_of_observation` that re-encodes the data that is in `d80`, `d81`, ... , `d04`.
 - Produce a new variable for each of the other variables that are one-hot encoded (i.e. `bac*` variable series).
 - Rename these variables to sensible names that are legible to a reader of your analysis. For example, the dependent variable as provided is called, `totfatrte`. Pick something more sensible, like, `total_fatalities_rate`. There are few enough of these variables to change, that you should change them for all the variables in the data. (You will thank yourself later.)
- (5 points) Provide a description of the basic structure of the dataset. What is this data? How, where, and when is it collected? Is the data generated through a survey or some other method? Is the data that is presented a sample from the population, or is it a *census* that represents the entire population? Minimally, this should include:
 - How is the our dependent variable of interest `total_fatalities_rate` defined?
- (20 points) Conduct a very thorough EDA, which should include both graphical and tabular techniques, on the dataset, including both the dependent variable `total_fatalities_rate` and the potential explanatory variables. Minimally, this should include:
 - How is the our dependent variable of interest `total_fatalities_rate` defined?
 - What is the average of `total_fatalities_rate` in each of the years in the time period covered in this dataset?

As with every EDA this semester, the goal of this EDA is not to document your own process of discovery – save that for an exploration notebook – but instead it is to bring a reader that is new to the data to a full understanding of the important features of your data as quickly as possible. In order to do this, your EDA should include a detailed, orderly narrative description of what you want your reader to know. Do not include any output – tables, plots, or statistics – that you do not intend to write about.

2.1 Part 1: loading data and creating features

```
# creating columns
# speed limit column
data$speed_limit <- ifelse(data$sl55 > 0, '55', NA)
data$speed_limit <- ifelse(data$sl65 > 0, '65', data$speed_limit)
data$speed_limit <- ifelse(data$sl70 > 0, '70', data$speed_limit)
data$speed_limit <- ifelse(data$sl75 > 0, '75', data$speed_limit)
data$speed_limit <- ifelse(data$slnone > 0, 'None', data$speed_limit)
data$perse_binary <- ifelse(data$perse > 0.5, 1, 0)
```

```

data$gdl_binary <- ifelse(data$gdl > 0.5, 1, 0)
data$sbprim_binary <- ifelse(data$sbprim > 0.5, 1, 0)
data$sbsec_binary <- ifelse(data$sbsecon > 0.5, 1, 0)
data$sl70plus <- ifelse(data$sl70plus > 0.5, 1, 0)

col_num <- 31:55
years <- 1980:2004

year_of_observation <- rep(NA, 1200)

data$year_test <- year_of_observation

for (i in 1:25){
  data$year_test <- ifelse(data[, col_num[i]]== 1, years[i], data$year_test)
}

data$bac <- ifelse(data$bac08 >0.5, "0.08", "None")
data$bac <- ifelse(data$bac10 > 0.5, "0.1", data$bac)

data$total_fatalities_rate <- data$totfatrte
data$night_fatalities_rate <- data$nghtfatrte
data$weekend_fatalities_rate <- data$wkndfatrte
data <- data %>%
  select(-totfatrte, -nghtfatrte, -wkndfatrte)
#head(data)

```

2.2 Part 2: Dataset Description

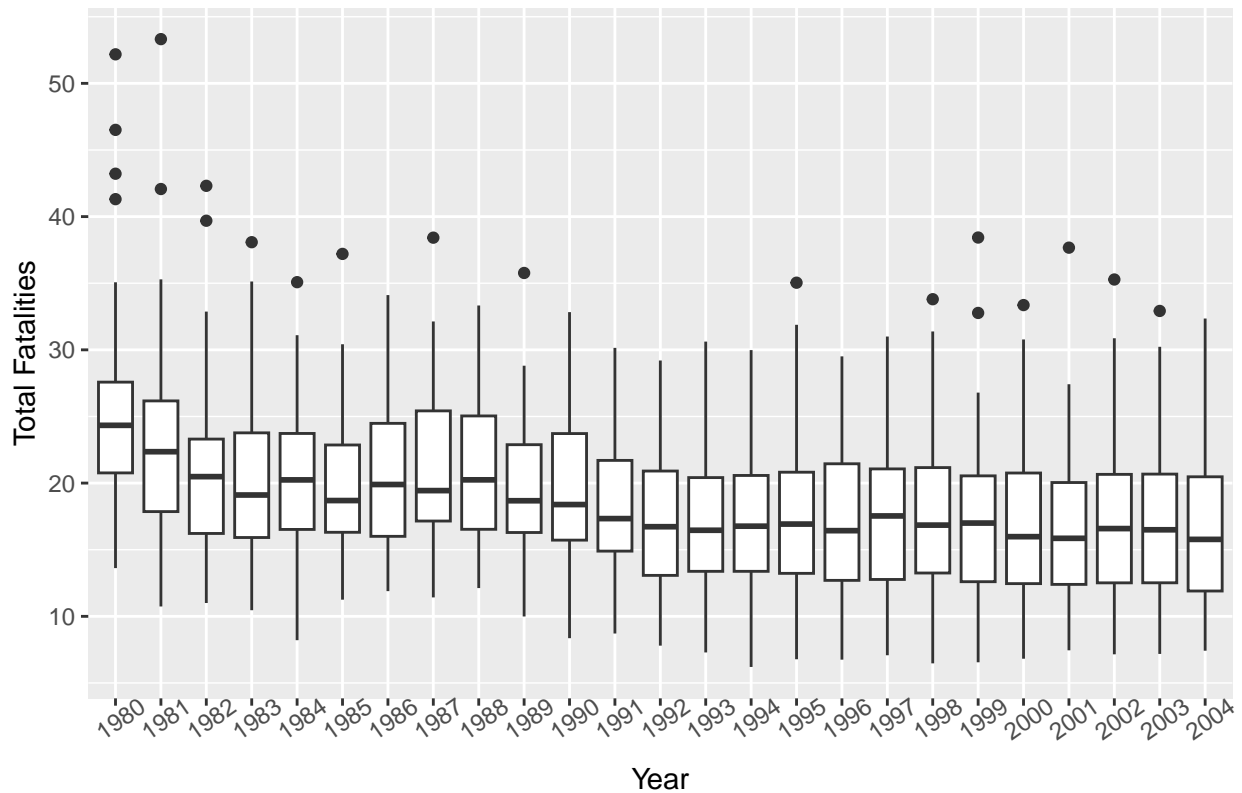
Answer The data comprises 1200 observations, each observation represents yearly data for a single contiguous US state from 1980 to 2004, so for each of the 48 states, there are 25 observations in the dataset and is therefore in panel data structure. There are 56 variables in total, relating to driving fatality statistics, driving laws, and some demographic factors. The fatality data is from the Federal Analysis Reporting System (FARS) by the National Highway and Traffic Safety Administration which records data on all crashes that are fatal, the reporting is standard across states. Every fatal car crash is recorded so this is census data that represents all traffic fatalities in the contiguous US. The variable of interest, 'total_fatalities_rate', is defined as the number of fatalities due to car crashes for a given state and year per 100,000 people, based on the state population for that year.

2.3 Part 3: EDA

2.3.1 Total Fatalities Rate

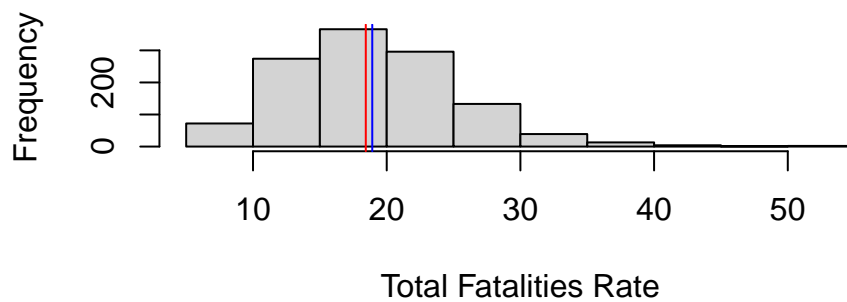
What is the average of `total_fatalities_rate` in each of the years in the time period covered in this dataset? - Based on the boxplot of total fatalities rate by year below, the mean total fatalities rate starts at around 24 per 100,000 and decreases slightly throughout the time period, ending at a mean of approximately 16 per 100,000.

Total Fatalities Rate by Year



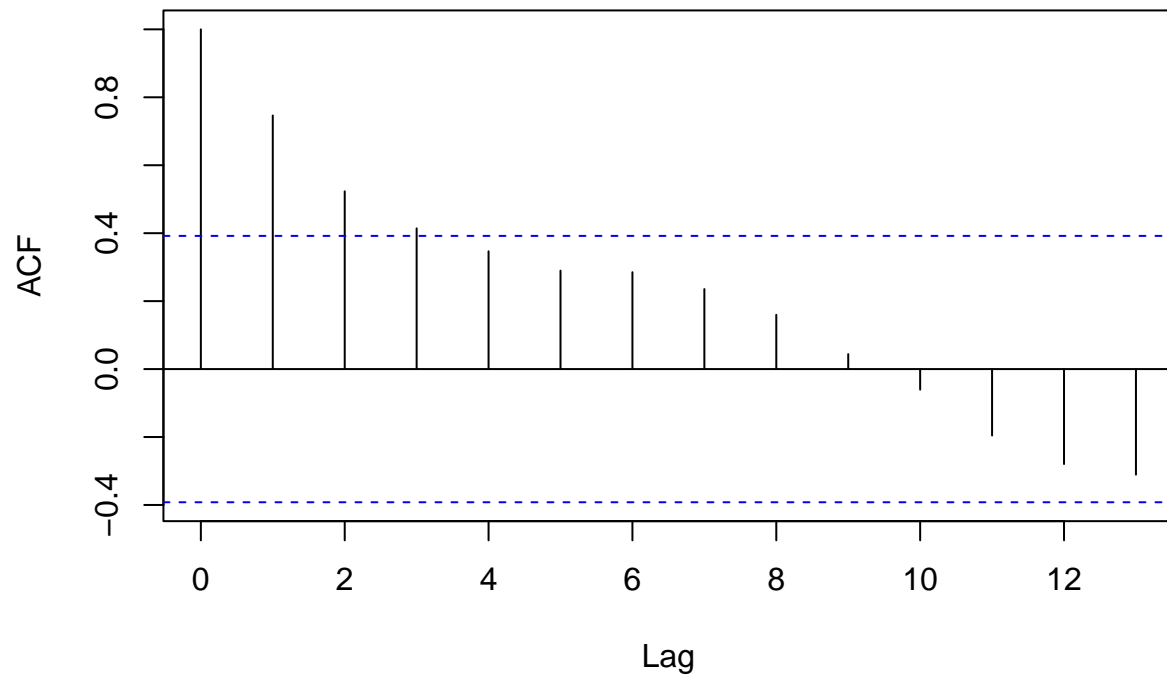
The main variable of interest is the total fatalities rate (TFR). The mean TFR across all time and observations is 18.92. This is slightly higher than the median TFR, 18.435, illustrated as the red line in the histogram. The fatalities distribution is skewed to the right. Outliers would lie either about the 95% quartile of a TFR greater or equal to 29.89 or below the 5th percentile ($TFR \leq 9.58$). In the following ACF diagram of annual average TFR across states also illustrates the autocorrelation across lagged TFR values.

Histogram of Total Fatalities Rate



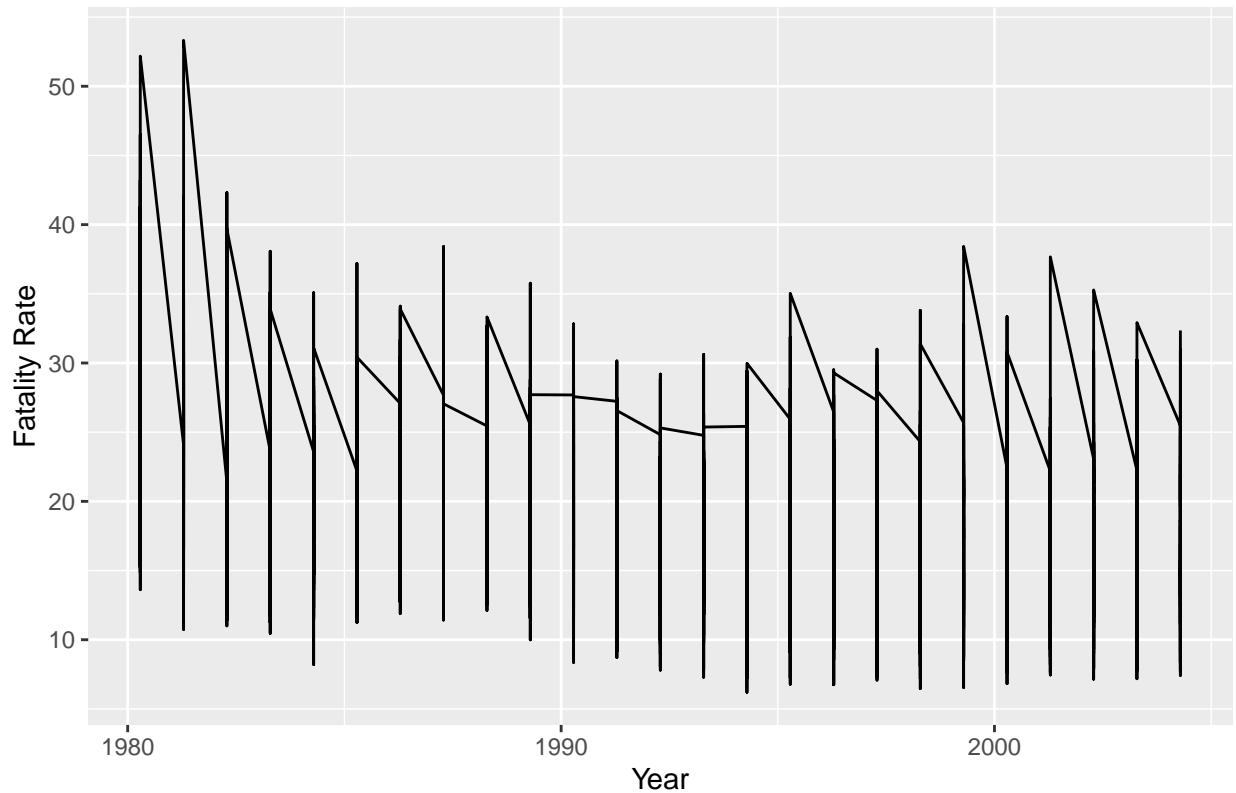
```
acf(annual_mean_total_fatalities_rate$avg_total_fatalities_rate, main = "ACF of total fatalities rate")
```

ACF of total fatalities rate



The total fatality rate is dynamic over time, as illustrated in the line plot of TFR across time.

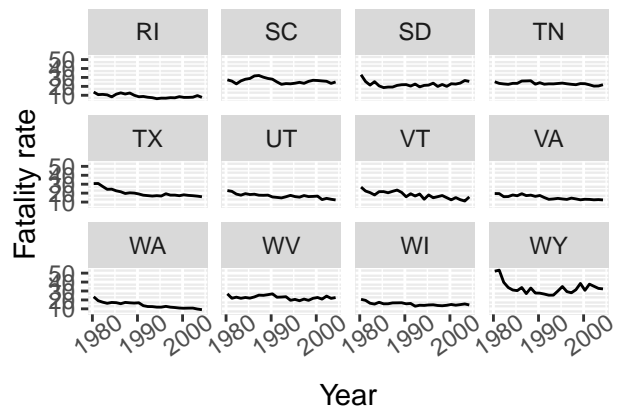
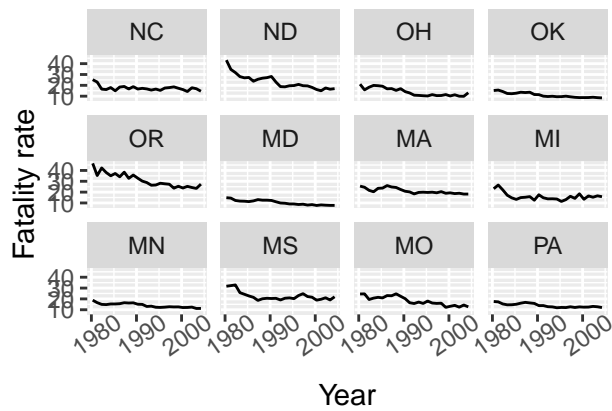
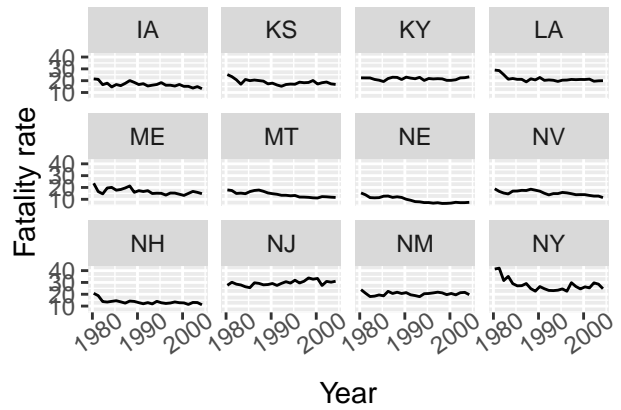
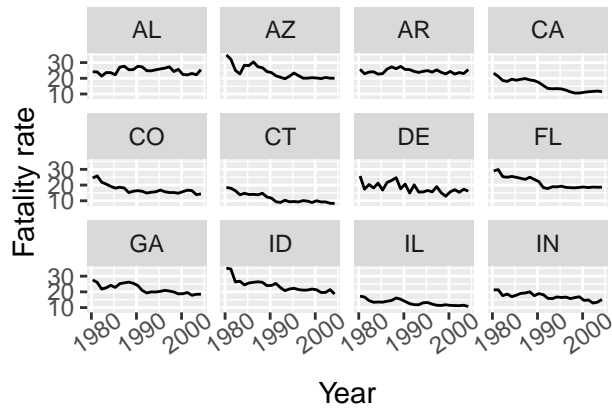
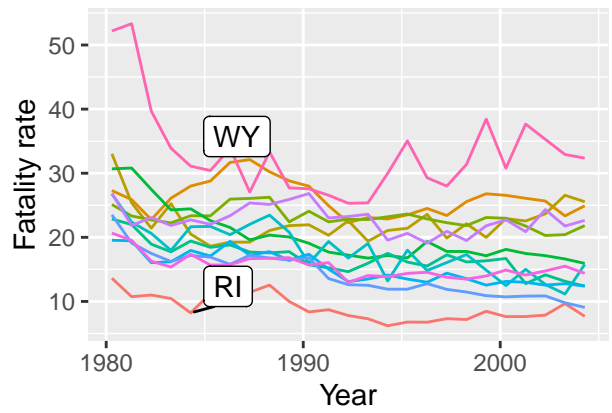
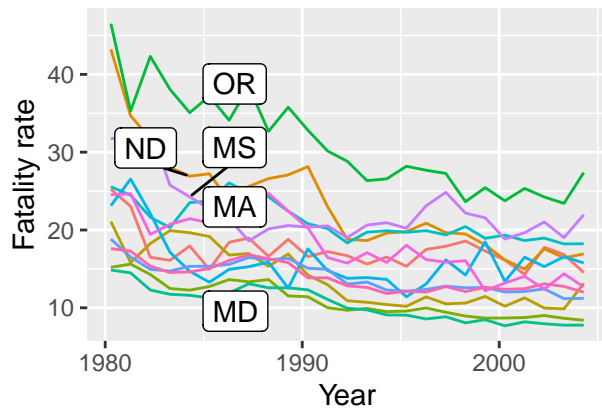
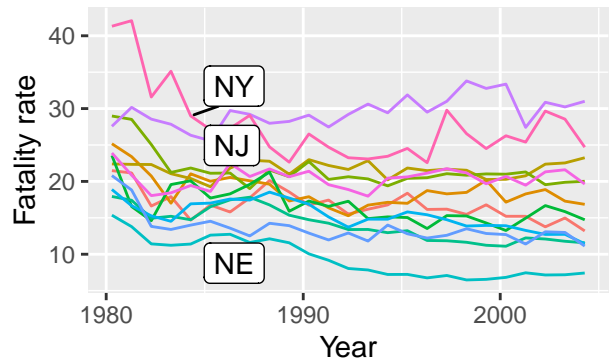
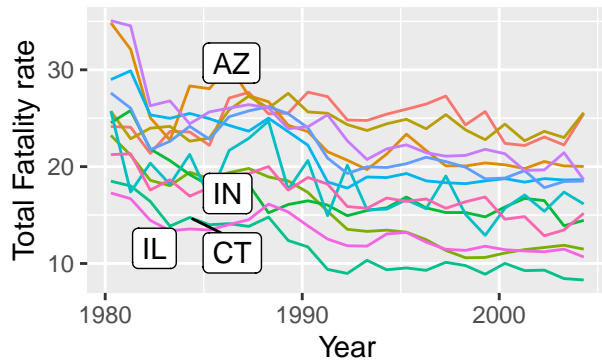
line plot of total fatalities rate over time



Further visualization of TFR is presented by state in the following diagrams. The four panels include different groupings of the states, by alphabetical order. In the upper left panel, the state with the highest initial TFR is Arizona and the lowest is Illinois. As time progresses, TFR drops across all states, especially Delaware and Colorado. For the states represented in the upper right-hand panel, New York has a TFR about 40, one of the highest of all time. Only Nebraska appears to maintain a low TFR among members of this group. The next group of states, represented in the lower left-hand panel, appear to be more likely to have TFRs below 20 per 100,000 persons. The next group, in the lower right-hand panel, has the highest TFR per 100,000 persons across all groups. Wyoming starts out with well over 40 TFR per 100,000 persons in the 1980s. Rhode Island, also in the group, is among the states with the lowest TFR per 100,000.

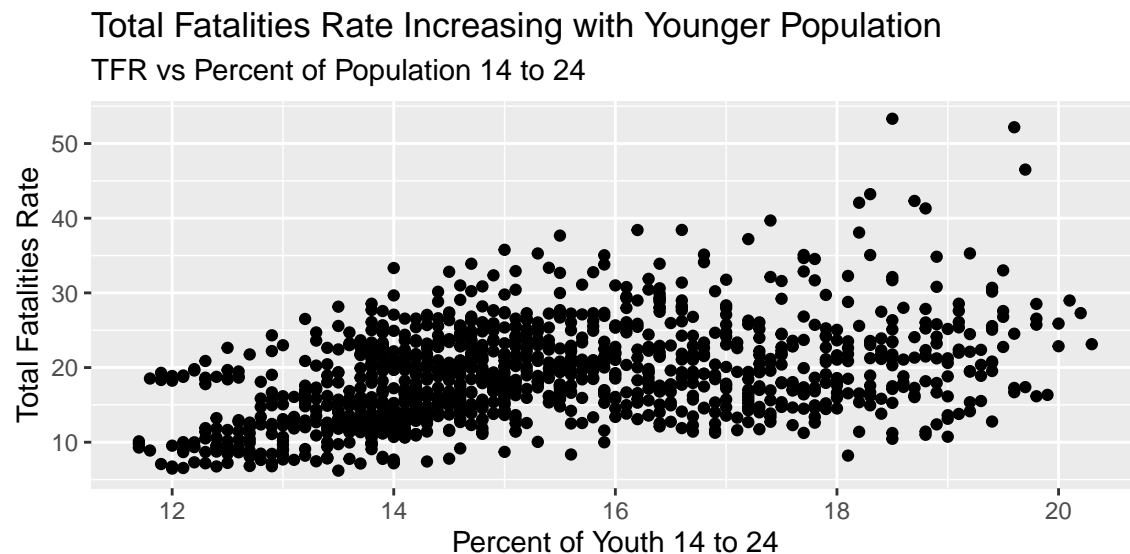
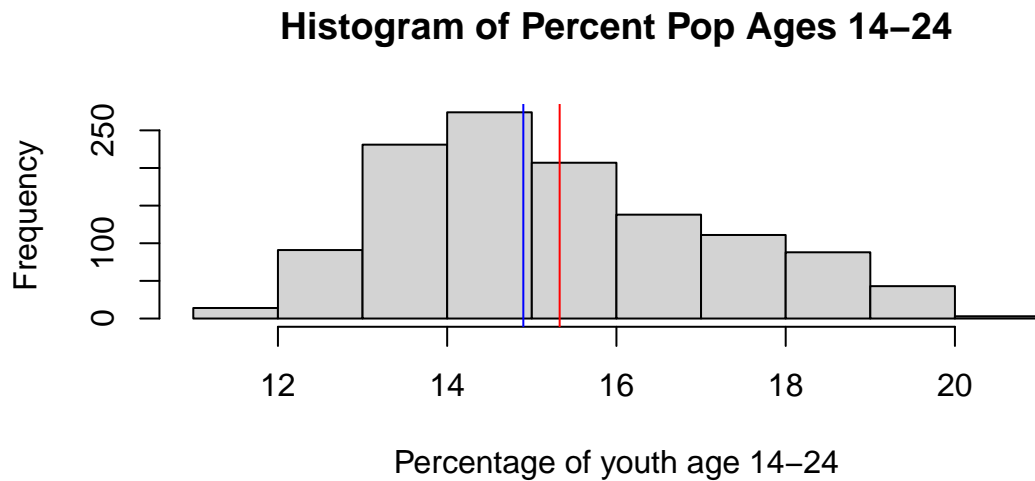
The next set of diagrams separates all TFRs across all states. This additional view of TFRs allows us to see the differences in TFR variation with more clarity. While states such as Delaware had relatively low TFR by 2004, the diagram indicates higher rates of variation than states such as Kansas, Kentucky, Nebraska, Minnesota, Pennsylvania, and Washington.

Total Fatality Rate by State

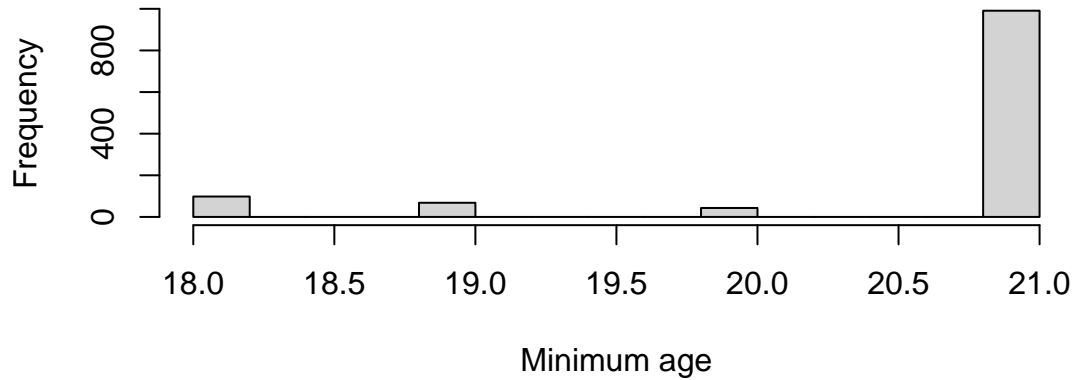


We now turn our attention to other variables of interest with relationship to TFR. These variables include measures of young people in the population, unemployment, miles driven per capita, seatbelt laws, and drunk driving regulations. Across the nation, the percent of the population between ages 14 and 24 years is 15.33 percent. This is almost half a percent higher than the median, 14.29 percent. The distribution is skewed to the right. The legal age to drink alcohol is primary 21 or in 991 of the cases minimum age is 21. It is 18 in 98 cases, 19 in 68 observations, and 20 in 42 observations. We may want to consider a dummy variable for minimum age is at or above 21 is one and zero otherwise.

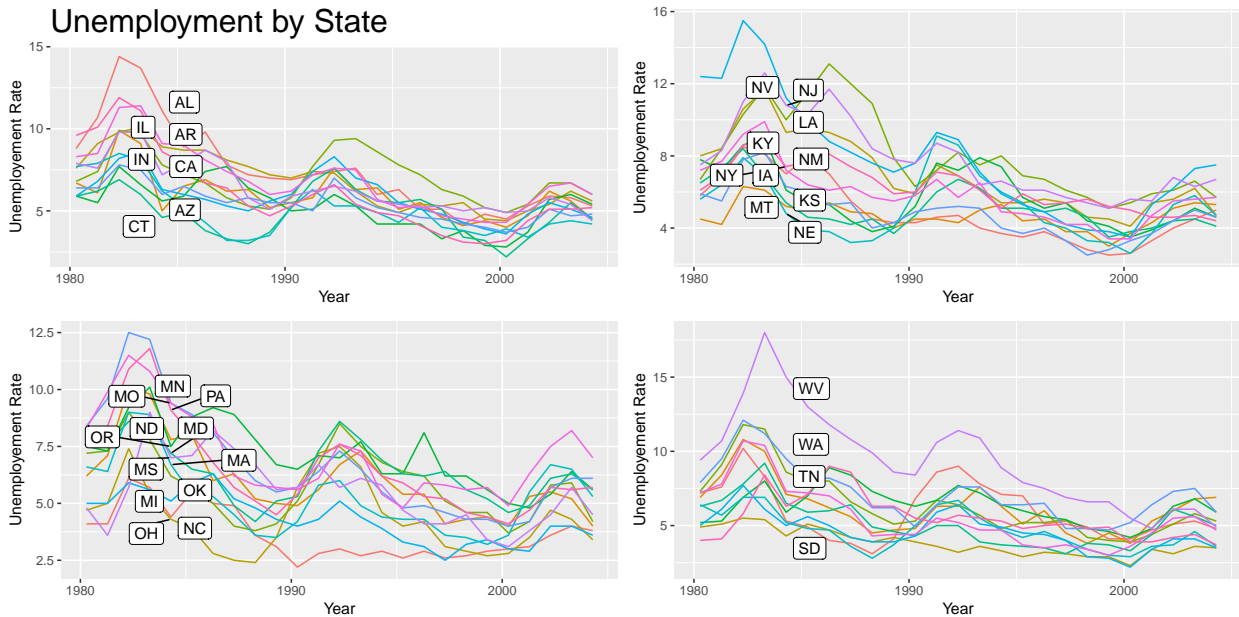
2.3.2 Youth



Histogram of minimum age

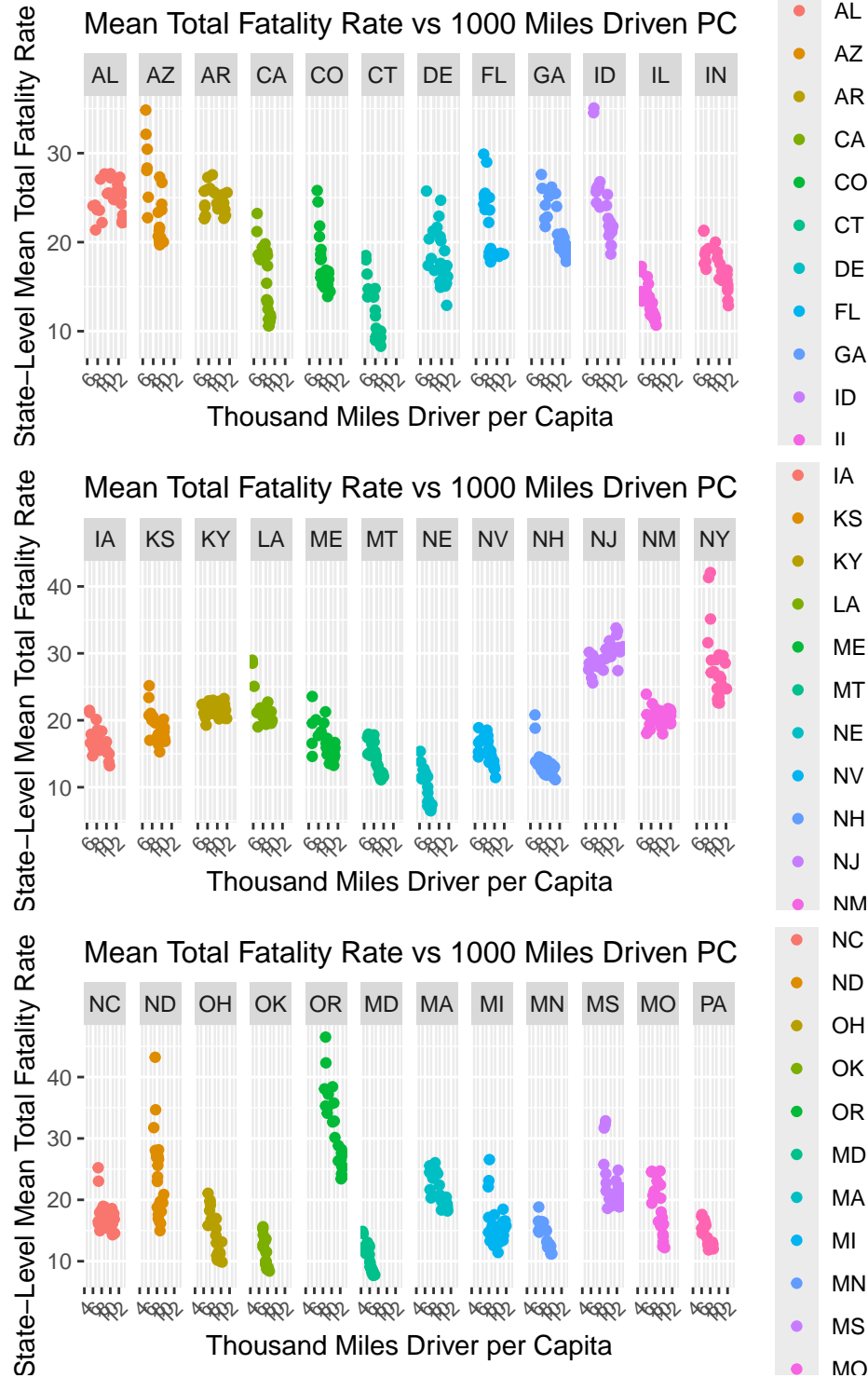


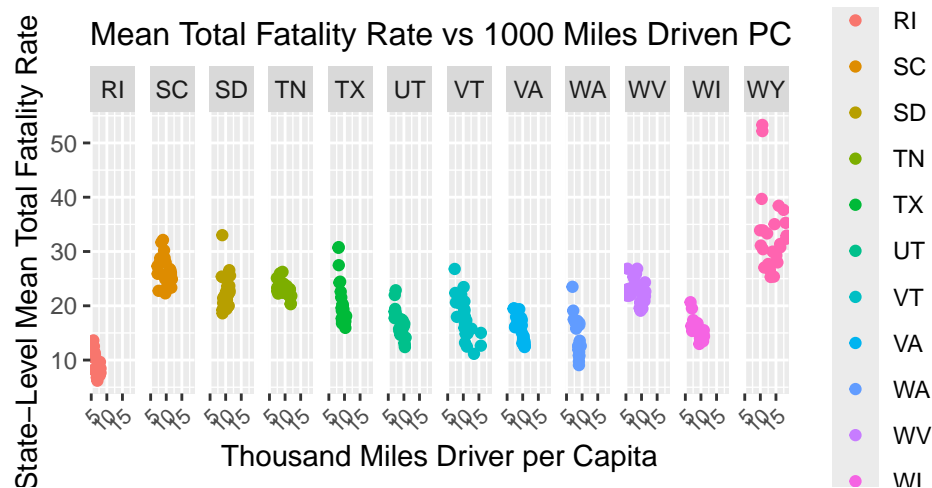
2.3.3 Unemployment



2.3.4 Miles Per Capita

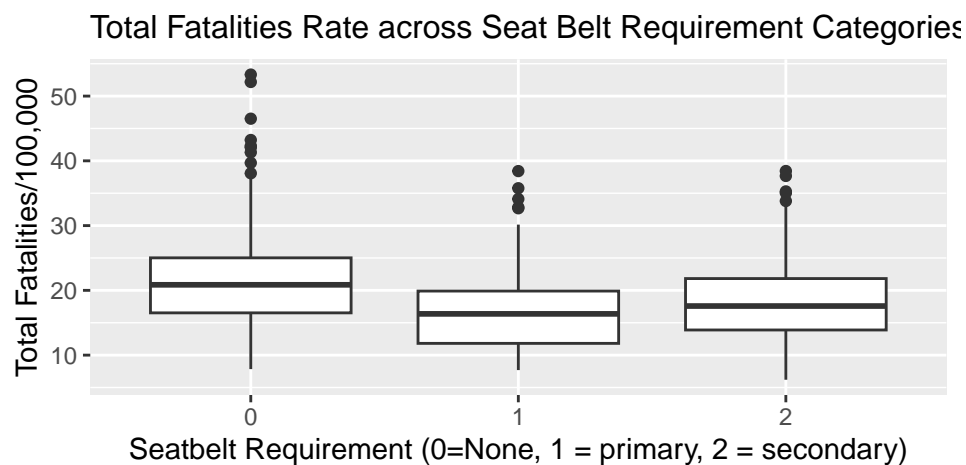
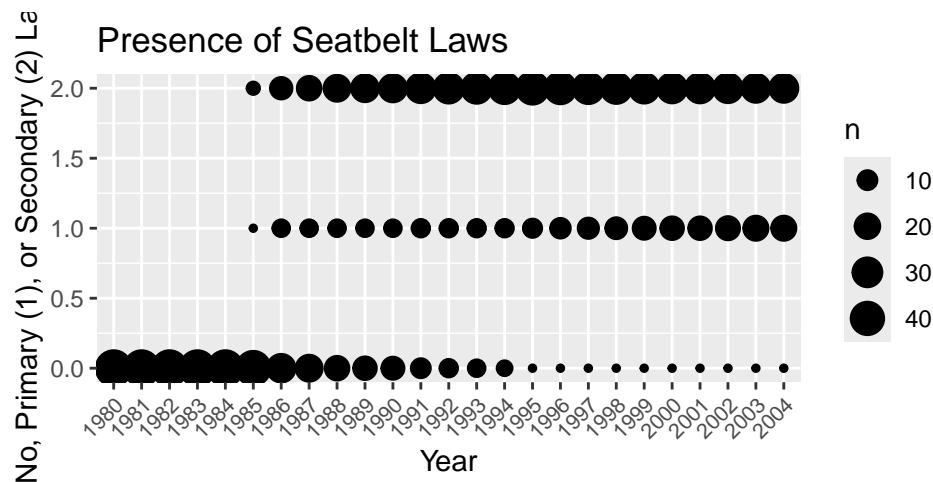
We converted the miles driven per state variable from billion miles per state to thousands of miles per capita for visualization. The four panel visualization illustrates a mixed relationship between miles driven per capita and total fatalities rate. There is an indication of a positive relationship between miles driven per capita and total fatalities rate when looking at extreme examples of high fatality states (i.e., Wyoming) and low fatality states (i.e., Delaware and Rhode Island). The correlation between the TFR and total vehicle miles driven is -0.26.





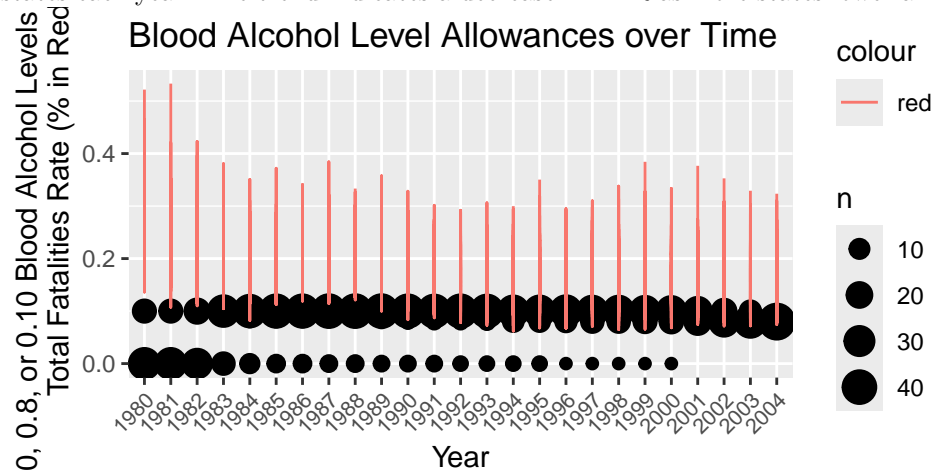
2.3.5 Seatbelt Regulation

There are a number of NAs in the Seatbelt regulation variable ($n > 900$). Still, over 99 percent of all panel observations included either primary or secondary seatbelt use. There is little variation in this variable. Essentially, it appears seatbelts were mandatory across locations by 1985.

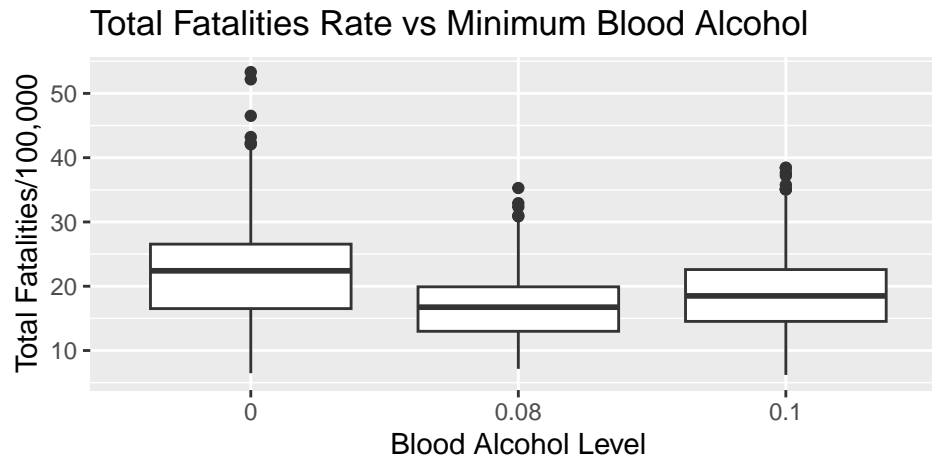


2.3.6 Drunk Driving

The blood alcohol variable indicates the level of constraint on drivers regarding the amount of alcohol they may consume and not be considered driving drunk. As the following figure illustrates, most states started to cap blood alcohol at 0.10 in prior to 1980. Then, starting in 1984, states increasingly capped the blood alcohol level at 0.85 for drunk driving. The red vertical bars represent the spread of TFR rates across states each year. The trend indicates a decrease in TFR as more states lower allowable blood alcohol levels.



The following box plot shows the relationship between blood alcohol restrictions and fatalities. Mean TFR is lowest when the blood alcohol limit is 0.85 in a state.

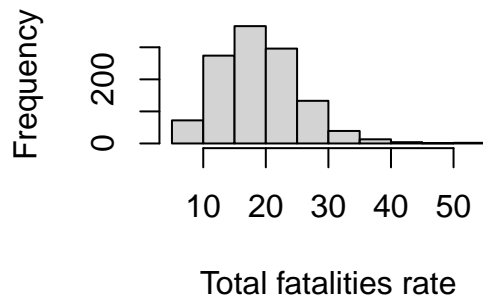


Below are histograms of variables that we found to be skewed and the effect of log transformation.

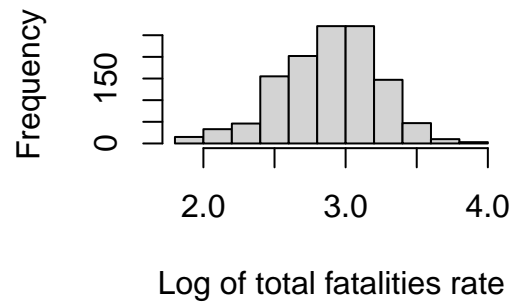
We found that a log transformation was helpful in normalizing the distributions of total fatalities rate, vehicle miles driven per capita, and unemployment rate.

A log transformation was not useful for percent of population ages 14-24.

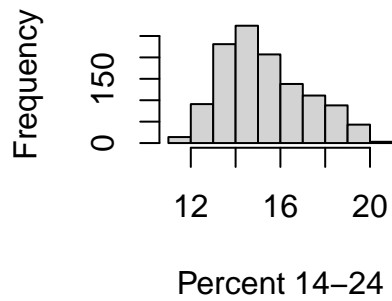
Histogram of total fatalities rate



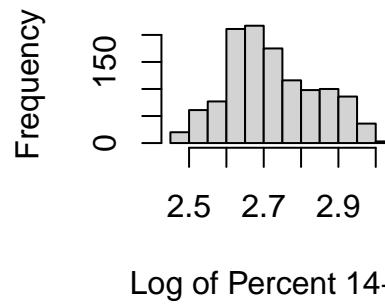
Histogram of log of total fatalities rate



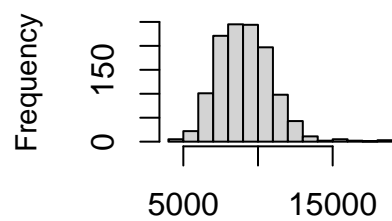
Histogram of Percent 14–24



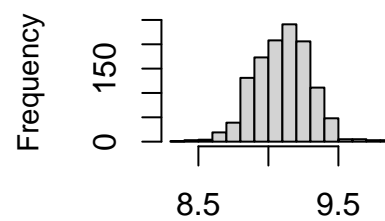
Histogram of log of Percent 14–24

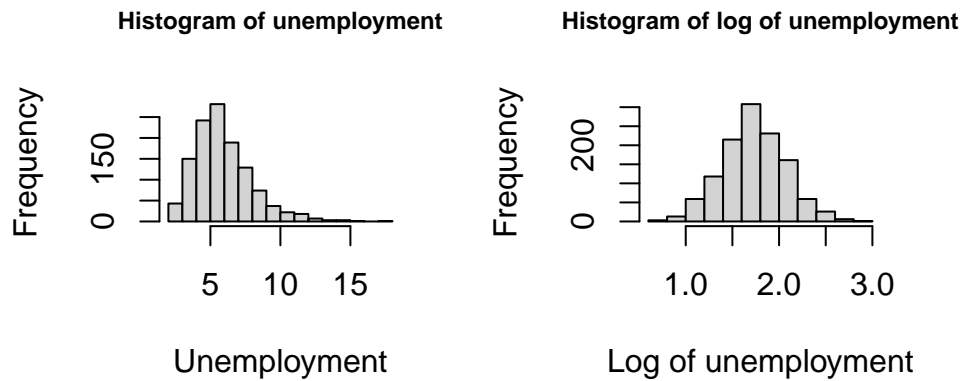


Histogram of vehicle miles per capita



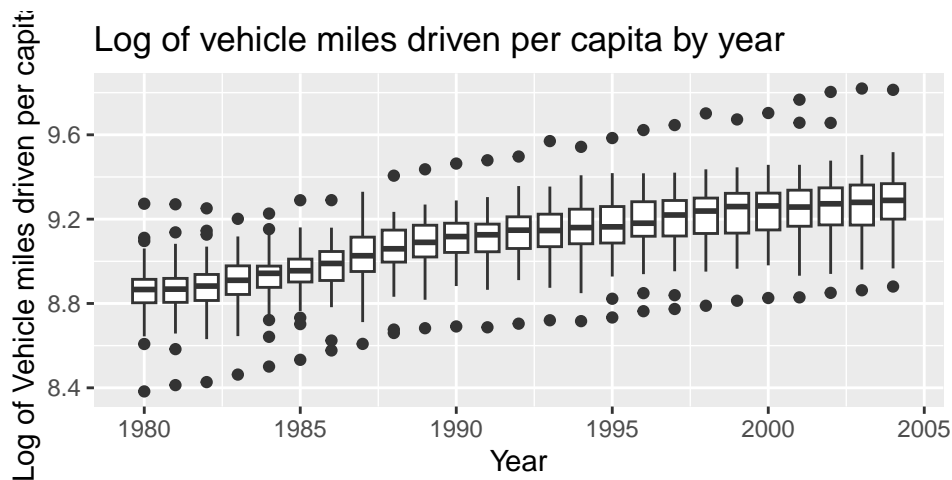
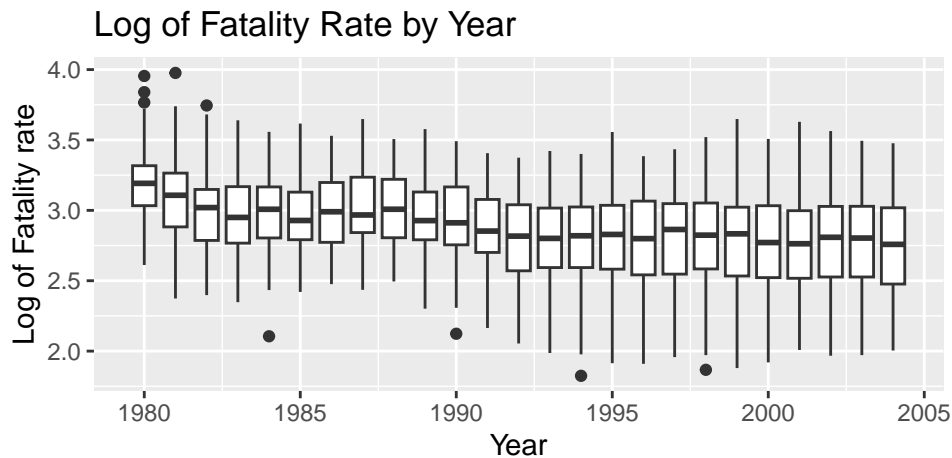
Histogram of log of vehicle miles per capita



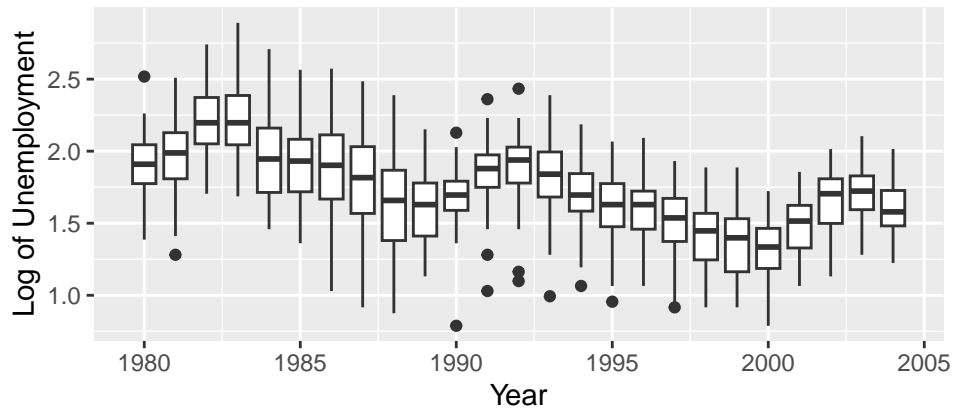


Below are boxplots showing the variables that we identified a log transformation to be useful. The log transformation was effective for linearizing the decrease in mean total fatalities rate over time, as well as linearizing the trend of vehicles miles driven over time. Unemployment is still variable over time after transformation.

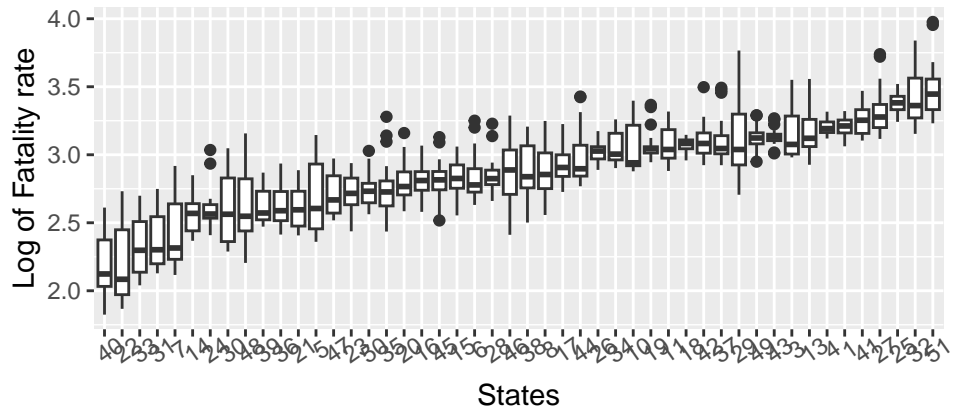
The transformed variables are also shown by state



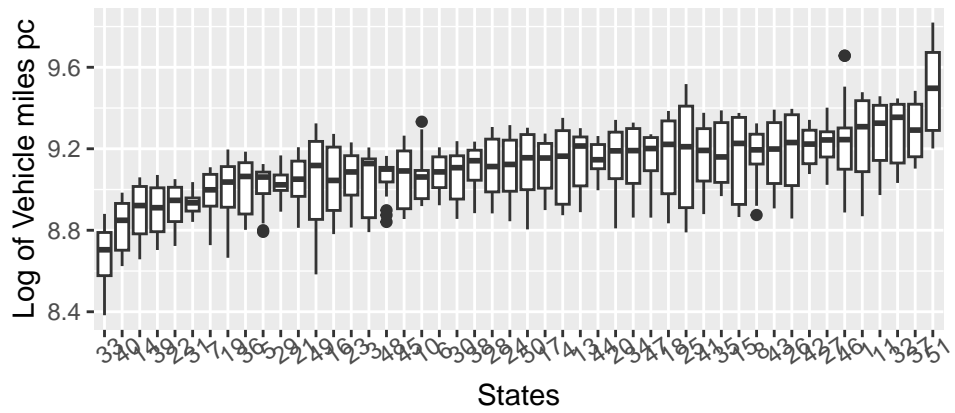
Log of unemployment by year

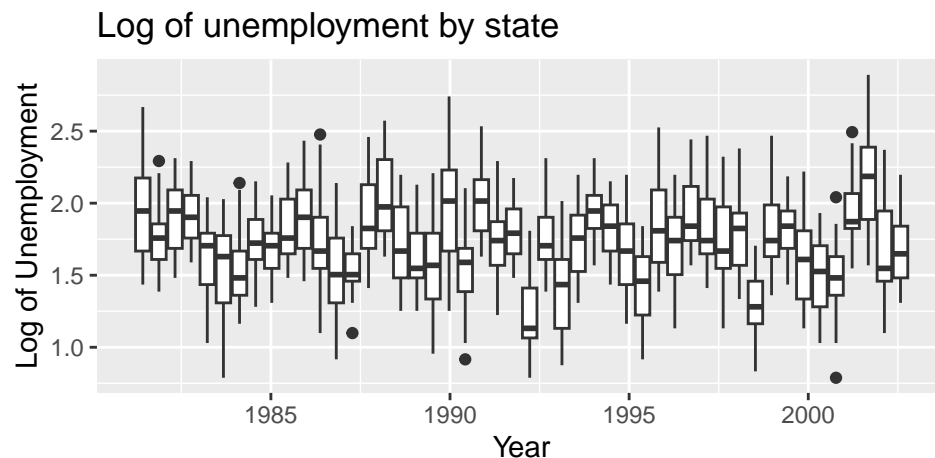


Log of fatality rate by state



Log of vehicle miles driven per capita by state





3 (15 points) Preliminary Model

Estimate a linear regression model of *totfatrtte* on a set of dummy variables for the years 1981 through 2004 and interpret what you observe. In this section, you should address the following tasks:

Why is fitting a linear model a sensible starting place?

A linear model makes sense as a starting place because it is an easy model to implement and interpret.

What does this model explain, and what do you find in this model?

The preliminary model is structured such that every year is represented by a dummy variable with 1980 as the baseline. A negative coefficient in this case represents a reduction in the mean rate of traffic fatalities for that year compared to the mean traffic fatalities rate in 1980. All of the coefficients are negative, meaning that there were less fatalities per 100,000 on average for that year compared to 1980.

In 1981, the reduction in traffic fatalities rate was 1.8 per 100,000 on average. From 1982 to 1990, the mean difference per year stayed between approximately 4.5 and 6, decreasing by more in the late 80s. In 1991 the mean difference in traffic fatalities rate since 1980 was -7.4 per 100,000. The difference stays mostly level, dropping a little throughout the rest of the data until 2004, when the mean difference was -8.8 per 100,000.

```
data$year_test_fac <- factor(data$year_test)
linear_mod <- lm(log(total_fatalities_rate)~year_test_fac, data = data)
summary(linear_mod)
```

```
##
## Call:
## lm(formula = log(total_fatalities_rate) ~ year_test_fac, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.96324 -0.22134  0.01005  0.23221  0.86830
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      3.19577    0.04697   68.035 < 2e-16 ***
## year_test_fac1981 -0.07878    0.06643   -1.186  0.235904
## year_test_fac1982 -0.19957    0.06643   -3.004  0.002719 **
```



```

## year_test_fac1983 -0.23523      0.06643   -3.541  0.000414 ***
## year_test_fac1984 -0.22585      0.06643   -3.400  0.000697 ***
## year_test_fac1985 -0.24301      0.06643   -3.658  0.000265 ***
## year_test_fac1986 -0.19681      0.06643   -2.963  0.003111 **
## year_test_fac1987 -0.19871      0.06643   -2.991  0.002836 **
## year_test_fac1988 -0.18885      0.06643   -2.843  0.004547 **
## year_test_fac1989 -0.24815      0.06643   -3.735  0.000196 ***
## year_test_fac1990 -0.26785      0.06643   -4.032  5.89e-05 ***
## year_test_fac1991 -0.34372      0.06643   -5.174  2.69e-07 ***
## year_test_fac1992 -0.40229      0.06643   -6.056  1.88e-09 ***
## year_test_fac1993 -0.40257      0.06643   -6.060  1.83e-09 ***
## year_test_fac1994 -0.40798      0.06643   -6.142  1.12e-09 ***
## year_test_fac1995 -0.38492      0.06643   -5.794  8.79e-09 ***
## year_test_fac1996 -0.39949      0.06643   -6.014  2.42e-09 ***
## year_test_fac1997 -0.38596      0.06643   -5.810  8.03e-09 ***
## year_test_fac1998 -0.40954      0.06643   -6.165  9.67e-10 ***
## year_test_fac1999 -0.41450      0.06643   -6.240  6.11e-10 ***
## year_test_fac2000 -0.43694      0.06643   -6.578  7.18e-11 ***
## year_test_fac2001 -0.43521      0.06643   -6.552  8.50e-11 ***
## year_test_fac2002 -0.42672      0.06643   -6.424  1.93e-10 ***
## year_test_fac2003 -0.43978      0.06643   -6.620  5.44e-11 ***
## year_test_fac2004 -0.44853      0.06643   -6.752  2.29e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3254 on 1175 degrees of freedom
## Multiple R-squared:  0.126, Adjusted R-squared:  0.1081
## F-statistic: 7.057 on 24 and 1175 DF, p-value: < 2.2e-16

```

Did driving become safer over this period? Please provide a detailed explanation.

The increasingly negative coefficients for each of the dummy variables shows that over time the trend of the total traffic fatalities is decreasing. This indicates that driving is becoming safer over the period

What, if any, are the limitation of this model. In answering this, please consider **at least**: - Are the parameter estimates reliable, unbiased estimates of the truth? Or, are they biased due to the way that the data is structured? - Are the uncertainty estimate reliable, unbiased estimates of sampling based variability? Or, are they biased due to the way that the data is structured?

The linear model makes a couple of assumptions of the data. The first assumption is that all of the observations are independent from one another. This cannot be the case for this data there are repeated measures of individual states, which creates correlation. OLS regression does not account for this and therefore the parameter estimates will not be reliable. Another common feature of panel data is heteroskedasticity in the errors, which violates the second OLS assumption that the errors are independent and normally distributed.

4 (15 points) Expanded Model

Expand the **Preliminary Model** by adding variables related to the following concepts:

- Blood alcohol levels
- Per se laws
- Primary seat belt laws (Note that if a law was enacted sometime within a year the fraction of the year is recorded in place of the zero-one indicator.)
- Secondary seat belt laws

- Speed limits faster than 70
- Graduated drivers licenses
- Percent of the population between 14 and 24 years old
- Unemployment rate
- Vehicle miles driven per capita.

If it is appropriate, include transformations of these variables. Please carefully explain carefully your rationale, which should be based on your EDA, behind any transformation you made. If no transformation is made, explain why transformation is not needed.

We performed log transformations on the fatalities rate, unemployment, and vehicle miles driven per capita variables to normalize them after seeing that the distributions of the variables were skewed.

```
expanded_mod <- lm(log(total_fatalities_rate)~year_test_fac + factor(bac) +
  factor(sl70plus) + factor(perse_binary) +
  factor(sbprim_binary) + factor(sbsec_binary) +
  factor(gdl_binary) + perc14_24 + log(unem) +
  log(vehicmilespc), data = data)
summary(expanded_mod)
```

```
##
## Call:
## lm(formula = log(total_fatalities_rate) ~ year_test_fac + factor(bac) +
##      factor(sl70plus) + factor(perse_binary) + factor(sbprim_binary) +
##      factor(sbsec_binary) + factor(gdl_binary) + perc14_24 + log(unem) +
##      log(vehicmilespc), data = data)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-0.58513	-0.12787	-0.00044	0.14047	0.62367

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-1.135e+01	4.017e-01	-28.267	< 2e-16 ***
year_test_fac1981	-9.192e-02	4.121e-02	-2.231	0.02589 *
year_test_fac1982	-2.956e-01	4.206e-02	-7.028	3.57e-12 ***
year_test_fac1983	-3.520e-01	4.269e-02	-8.247	4.35e-16 ***
year_test_fac1984	-2.997e-01	4.346e-02	-6.896	8.77e-12 ***
year_test_fac1985	-3.371e-01	4.443e-02	-7.586	6.75e-14 ***
year_test_fac1986	-3.133e-01	4.640e-02	-6.752	2.30e-11 ***
year_test_fac1987	-3.492e-01	4.839e-02	-7.216	9.60e-13 ***
year_test_fac1988	-3.602e-01	5.094e-02	-7.071	2.64e-12 ***
year_test_fac1989	-4.451e-01	5.291e-02	-8.412	< 2e-16 ***
year_test_fac1990	-5.053e-01	5.407e-02	-9.345	< 2e-16 ***
year_test_fac1991	-6.207e-01	5.517e-02	-11.252	< 2e-16 ***
year_test_fac1992	-7.266e-01	5.627e-02	-12.914	< 2e-16 ***
year_test_fac1993	-7.185e-01	5.694e-02	-12.617	< 2e-16 ***
year_test_fac1994	-7.055e-01	5.802e-02	-12.161	< 2e-16 ***
year_test_fac1995	-6.818e-01	5.950e-02	-11.459	< 2e-16 ***
year_test_fac1996	-8.080e-01	6.161e-02	-13.114	< 2e-16 ***
year_test_fac1997	-8.170e-01	6.266e-02	-13.038	< 2e-16 ***
year_test_fac1998	-8.654e-01	6.381e-02	-13.562	< 2e-16 ***
year_test_fac1999	-8.671e-01	6.459e-02	-13.424	< 2e-16 ***
year_test_fac2000	-8.794e-01	6.568e-02	-13.388	< 2e-16 ***

```
## year_test_fac2001      -9.349e-01  6.610e-02 -14.144 < 2e-16 ***
## year_test_fac2002      -9.794e-01  6.653e-02 -14.721 < 2e-16 ***
## year_test_fac2003      -1.003e+00  6.679e-02 -15.012 < 2e-16 ***
## year_test_fac2004      -9.839e-01  6.853e-02 -14.357 < 2e-16 ***
## factor(bac)0.1          4.494e-02  1.846e-02   2.435  0.01504 *
## factor(bac)None          6.190e-02  2.433e-02   2.545  0.01107 *
## factor(sl70plus)1        2.219e-01  2.162e-02  10.262 < 2e-16 ***
## factor(perse_binary)1    -1.882e-02  1.464e-02  -1.286  0.19869
## factor(sbprim_binary)1    9.419e-04  2.456e-02   0.038  0.96942
## factor(sbsec_binary)1     2.043e-02  2.144e-02   0.953  0.34084
## factor(gdl_binary)1      -2.129e-02  2.529e-02  -0.842  0.40000
## perc14_24                1.779e-02  6.111e-03   2.911  0.00367 **
## log(unem)                 2.673e-01  2.414e-02  11.071 < 2e-16 ***
## log(vehicmiles)          1.541e+00  4.432e-02  34.765 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2015 on 1165 degrees of freedom
## Multiple R-squared:  0.6678, Adjusted R-squared:  0.6581
## F-statistic: 68.87 on 34 and 1165 DF,  p-value: < 2.2e-16
```

How are the blood alcohol variables defined? Interpret the coefficients that you estimate for this concept.

The blood alcohol variable is a categorical variable defined with 3 levels. The most restrictive BAC law is set at 0.08, the next level describes states who prohibit driving with a BAC of 0.10, the next level describes states with no BAC restrictions.

The coefficient of $BAC_{0.10}$ is 0.04, this means that mean total fatality rate increases by 0.04 when a state has a law prohibiting driving with a BAC level of 0.10 compared to 0.08.

The coefficient of $BAC_{0.08}$ is 0.06, which means that the total fatality rate increases by 0.06 per 100,000 on average, when a state has no BAC restrictions compared to a restriction of 0.08.

Do *per se* laws have a negative effect on the fatality rate?

The coefficient representing if a state has *per se* laws has a coefficient of -0.02. The model indicates that having *per se* laws does have a negative effect on the fatality rate, however, this coefficient is not significant, so we can not conclude that the effect on mean fatality rate related to having *per se* laws is significantly different than 0.

Does having a primary seat belt law?

The dummy variable representing a state having a primary seatbelt law has a coefficient of 0.0009. The model indicates that having *per se* laws has a positive effect on the fatality rate, however, this coefficient is not significant, so we can not conclude that the effect on mean fatality rate related to having *per se* laws is significantly different than 0 from this model.

5 (15 points) State-Level Fixed Effects

Re-estimate the Expanded Model using fixed effects at the state level.

We ran within model and first difference models for state level fixed effects.

```
within_model <- plm(log(total_fatalities_rate) ~ year_test_fac + factor(bac) +
                    factor(sl70plus) + factor(perse_binary) +
```

```

        factor(sbprim_binary) + factor(sbsec_binary) +
        factor(gdl_binary) + perc14_24 + log(unem) +
        log(vehicmilespc),
    data = data,
    index = c("state", "year_test_fac"),
    effect = "individual", model = "within")

# summary(within_model)

```

```

fd_model<- plm(log(total_fatalities_rate)~year_test_fac + factor(bac) +
        factor(sl70plus) + factor(perse_binary) +
        factor(sbprim_binary) + factor(sbsec_binary) +
        factor(gdl_binary) + perc14_24 + log(unem) +
        log(vehicmilespc),
    data = data,
    index = c("state", "year_test_fac"),
    effect = "individual", model = "fd")

# summary(fd_model)

```

We fail to reject both of the null hypothesis for the Wooldridge first difference statistical test. When both null hypothesis are rejected it means that the residuals are serially correlated for both the within model and the first difference model. To correct for serial correlation, we need to estimate the model parameters with autocorrelation robust standard errors. We did so using cluster standard errors.

```

# first null hypothesis wooldridge test
pwfdtest(fd_model, data = data, h0="fe")

```

```

##
## Wooldridge's first-difference test for serial correlation in panels
##
## data: fd_model
## F = 6.4253, df1 = 1, df2 = 1102, p-value = 0.01139
## alternative hypothesis: serial correlation in original errors

```

```

# second null hypothesis wooldridge test
pwfdtest(fd_model, data = data, h0="fd")

```

```

##
## Wooldridge's first-difference test for serial correlation in panels
##
## data: fd_model
## F = 133.32, df1 = 1, df2 = 1102, p-value < 2.2e-16
## alternative hypothesis: serial correlation in differenced errors

```

```

# within model estimate with autocorrelation robust standard errors
within_model_robust <- coeftest(within_model, vcov=vcovHC(within_model,type="HC0",cluster="group"))
within_model_robust

```

```

##
## t test of coefficients:

```

```

##
##               Estimate Std. Error t value Pr(>|t|)
## year_test_fac1981    -0.0629920  0.0170299 -3.6989 0.0002270 ***
## year_test_fac1982    -0.1368192  0.0180319 -7.5876 6.840e-14 ***
## year_test_fac1983    -0.1721909  0.0220070 -7.8244 1.176e-14 ***
## year_test_fac1984    -0.2071493  0.0224198 -9.2396 < 2.2e-16 ***
## year_test_fac1985    -0.2307212  0.0274958 -8.3911 < 2.2e-16 ***
## year_test_fac1986    -0.1923439  0.0354120 -5.4316 6.850e-08 ***
## year_test_fac1987    -0.2380879  0.0399884 -5.9539 3.501e-09 ***
## year_test_fac1988    -0.2696855  0.0494726 -5.4512 6.153e-08 ***
## year_test_fac1989    -0.3426135  0.0547402 -6.2589 5.515e-10 ***
## year_test_fac1990    -0.3535455  0.0603720 -5.8561 6.223e-09 ***
## year_test_fac1991    -0.3895487  0.0645289 -6.0368 2.136e-09 ***
## year_test_fac1992    -0.4505470  0.0684003 -6.5869 6.898e-11 ***
## year_test_fac1993    -0.4677952  0.0693084 -6.7495 2.378e-11 ***
## year_test_fac1994    -0.5033180  0.0688585 -7.3095 5.096e-13 ***
## year_test_fac1995    -0.5000318  0.0743208 -6.7280 2.740e-11 ***
## year_test_fac1996    -0.5531910  0.0769218 -7.1916 1.170e-12 ***
## year_test_fac1997    -0.5754564  0.0788891 -7.2945 5.667e-13 ***
## year_test_fac1998    -0.6300575  0.0792083 -7.9544 4.383e-15 ***
## year_test_fac1999    -0.6497906  0.0814573 -7.9771 3.686e-15 ***
## year_test_fac2000    -0.6827377  0.0810812 -8.4204 < 2.2e-16 ***
## year_test_fac2001    -0.6538599  0.0851197 -7.6817 3.418e-14 ***
## year_test_fac2002    -0.6158438  0.0830093 -7.4190 2.329e-13 ***
## year_test_fac2003    -0.6188054  0.0855231 -7.2355 8.595e-13 ***
## year_test_fac2004    -0.6553754  0.0893742 -7.3329 4.312e-13 ***
## factor(bac)0.1        0.0042323  0.0176069  0.2404 0.8100822
## factor(bac)None       0.0204314  0.0236711  0.8631 0.3882473
## factor(sl70plus)1     0.0720100  0.0219278  3.2840 0.0010553 **
## factor(perse_binary)1 -0.0534343  0.0153619 -3.4784 0.0005239 ***
## factor(sbprim_binary)1 -0.0404917  0.0246488 -1.6427 0.1007174
## factor(sbsec_binary)1  0.0056980  0.0162028  0.3517 0.7251515
## factor(gdl_binary)1   -0.0154092  0.0202008 -0.7628 0.4457431
## perc14_24             0.0202582  0.0109021  1.8582 0.0634052 .
## log(unem)             -0.1926117  0.0230680 -8.3497 < 2.2e-16 ***
## log(vehicmiles)pc     0.6784007  0.1355082  5.0063 6.441e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

- What do you estimate for coefficients on the blood alcohol variables? How do the coefficients on the blood alcohol variables change, if at all?
- What do you estimate for coefficients on per se laws? How do the coefficients on per se laws change, if at all?
- What do you estimate for coefficients on primary seat-belt laws? How do the coefficients on primary seatbelt laws change, if at all?

The coefficient for bac = 0.1 variable in this fixed effects model is now 0.0042323, which is a smaller coefficient compared to the expanded model which has a positive coefficient for bac laws (4.494e-02). What is also notable is that this bac law was a significant explanatory variable in the expanded model, and in the fixed effects model it is not significant to explaining fatality rates.

In the fixed effects model, the coefficient for perse laws is -0.0534343, compared to -1.882e-02 in the expanded model. The direction is still negative, implying that the addition of per se laws decrease fatality rates. In the expanded model this coefficient wasn't significant, but in the fixed effects model it is highly significant.

The coefficient for the primary seatbelt law in this model is -0.0404917, and for the secondary seat belt law it is 0.0056980. The coefficient for the primary seatbelt law is small and positive for the expanded model, but negative in the fixed effect model. The coefficients for the secondary seat belt law are both small and positive in the fixed effects model and the expanded model. These coefficients are not significant in the expanded model, but in the fixed effects model the primary seatbelt law coefficient is significant. This implies that the primary seatbelt law has a negative effect on fatality rates.

Which set of estimates do you think is more reliable? Why do you think this?

- What assumptions are needed in each of these models?
- Are these assumptions reasonable in the current context?

The fixed effects model is likely more reliable since it is able to avoid omitted variable bias in the model. We can see from the difference in some coefficients that the expanded model was underestimating the effects of some of the variables. One clear example is how only the fixed effects model recognized the magnitude and significance of the per se laws coefficient.

The assumptions for the expanded model would be the OLS regression assumptions in order to get statistically correct estimates. Because of unobserved effects, repeated observations, and time-constant explanatory variables, this model is likely statistically inaccurate.

The assumptions for the fixed effect model are linearity (model is linear in parameters), strict exogeneity (zero conditional means), observations are IID across entities, and no perfect collinearity. The second assumption is necessary to avoid omitted variable bias. Using the within model for the fixed effects means we are removing the unobserved individual heterogeneity by using the time-demeaned model. This allows the equation to meet the assumptions for the fixed effect model since we are averaging by state, making the observations independent and removing the fixed effects coefficients.

We also ran a pFtest that rejected the null hypothesis (p-value = 2.2e-16) which proves that there are individual effects that are better handled with the fixed effects model. We further tested and made sure to use the within model instead of first differencing to estimate our fixed effects model.

```
pFtest(within_model, expanded_mod)
```

```
##
## F test for individual effects
##
## data: log(total_fatalities_rate) ~ year_test_fac + factor(bac) + factor(sl70plus) + ...
## F = 105.56, df1 = 47, df2 = 1118, p-value < 2.2e-16
## alternative hypothesis: significant effects
```

6 (10 points) Consider a Random Effects Model

Instead of estimating a fixed effects model, should you have estimated a random effects model?

```
re.model <- plm(log(total_fatalities_rate)~year_test_fac + factor(bac) +
               factor(sl70plus) + factor(perse_binary) +
               factor(sbprim_binary) + factor(sbsec_binary) +
               factor(gdl_binary) + perc14_24 +
               log(unem) + log(vehicmilespc),
               data = data,
               index = c("state", "year_test_fac"), model = "random", random.method = "walhus")
# summary(re.model)
```

Please state the assumptions of a random effects model, and evaluate whether these assumptions are met in the data.

The random effects model requires a strong assumption of independence between random effects and other predictors. We use this model when we think the unobserved effect is uncorrelated with all the explanatory variables. For these assumptions to be met, there can't be a correlation between the state effects and the other explanatory variables. However, it is apparent that each state and their own fatality rates are correlated with the passing of laws related to driving safety. Therefore, the assumptions for a random effects model aren't met.

If the assumptions are, in fact, met in the data, then estimate a random effects model and interpret the coefficients of this model. Comment on how, if at all, the estimates from this model have changed compared to the fixed effects model.

- If the assumptions are **not** met, then do not estimate the data. But, also comment on what the consequences would be if you were to *inappropriately* estimate a random effects model. Would your coefficient estimates be biased or not? Would your standard error estimates be biased or not? Or, would there be some other problem that might arise?

We ran the Hausman test for random versus fixed effects. The p-value was very small and much less than 0.05. We reject the null hypothesis that the random effects model is appropriate. The assumptions of the model are not met.

```
phtest(within_model, re.model, include = FALSE)
```

```
##
## Hausman Test
##
## data: log(total_fatalities_rate) ~ year_test_fac + factor(bac) + factor(sl70plus) + ...
## chisq = 392.57, df = 34, p-value < 2.2e-16
## alternative hypothesis: one model is inconsistent
```

If we inappropriately estimate a random effects model, we risk biased coefficient estimates due to omitted variable bias. There would also be correlation in the errors. The random effects model is a more efficient estimator, meaning that if the assumptions were held true, that the standard errors of the betas would be less than the fixed effect models. However, the assumption is not true, so the standard error estimates would be biased.

```
stargazer(linear_mod, expanded_mod, fd_model, within_model_robust,
style="qje", type="text", omit.stat=c("adj.rsq", "f"),
column.labels = c("Prelim Model", "Expanded Model", "FD", "Within Robust SE"))
```

```
##
## =====
##               log(total_fatalities_rate)
##               OLS               panel
##               OLS               linear
##               Prelim Model       FD       coefficient
##               (1)               (2)       test
##               (1)               (2)       Within Robust SE
##               (1)               (2)       (3)       (4)
## -----
```

## year_test_fac1981	-0.079	-0.092**	-0.047***	-0.063***
##	(0.066)	(0.041)	(0.013)	(0.017)
##				

## year_test_fac1982	-0.200***	-0.296***	-0.113***	-0.137***
##	(0.066)	(0.042)	(0.020)	(0.018)
##				
## year_test_fac1983	-0.235***	-0.352***	-0.112***	-0.172***
##	(0.066)	(0.043)	(0.025)	(0.022)
##				
## year_test_fac1984	-0.226***	-0.300***	-0.085***	-0.207***
##	(0.066)	(0.043)	(0.028)	(0.022)
##				
## year_test_fac1985	-0.243***	-0.337***	-0.074**	-0.231***
##	(0.066)	(0.044)	(0.032)	(0.027)
##				
## year_test_fac1986	-0.197***	-0.313***	0.002	-0.192***
##	(0.066)	(0.046)	(0.035)	(0.035)
##				
## year_test_fac1987	-0.199***	-0.349***	0.017	-0.238***
##	(0.066)	(0.048)	(0.039)	(0.040)
##				
## year_test_fac1988	-0.189***	-0.360***	0.045	-0.270***
##	(0.066)	(0.051)	(0.043)	(0.049)
##				
## year_test_fac1989	-0.248***	-0.445***	0.012	-0.343***
##	(0.066)	(0.053)	(0.046)	(0.055)
##				
## year_test_fac1990	-0.268***	-0.505***	0.023	-0.354***
##	(0.066)	(0.054)	(0.048)	(0.060)
##				
## year_test_fac1991	-0.344***	-0.621***	-0.013	-0.390***
##	(0.066)	(0.055)	(0.048)	(0.065)
##				
## year_test_fac1992	-0.402***	-0.727***	-0.049	-0.451***
##	(0.066)	(0.056)	(0.048)	(0.068)
##				
## year_test_fac1993	-0.403***	-0.718***	-0.040	-0.468***
##	(0.066)	(0.057)	(0.047)	(0.069)
##				
## year_test_fac1994	-0.408***	-0.706***	-0.042	-0.503***
##	(0.066)	(0.058)	(0.046)	(0.069)
##				
## year_test_fac1995	-0.385***	-0.682***	-0.010	-0.500***
##	(0.066)	(0.059)	(0.044)	(0.074)
##				
## year_test_fac1996	-0.399***	-0.808***	-0.017	-0.553***
##	(0.066)	(0.062)	(0.043)	(0.077)
##				
## year_test_fac1997	-0.386***	-0.817***	-0.007	-0.575***
##	(0.066)	(0.063)	(0.040)	(0.079)
##				
## year_test_fac1998	-0.410***	-0.865***	-0.031	-0.630***
##	(0.066)	(0.064)	(0.037)	(0.079)
##				
## year_test_fac1999	-0.414***	-0.867***	-0.035	-0.650***
##	(0.066)	(0.065)	(0.033)	(0.081)
##				

## year_test_fac2000	-0.437***	-0.879***	-0.051*	-0.683***
##	(0.066)	(0.066)	(0.030)	(0.081)
##				
## year_test_fac2001	-0.435***	-0.935***	-0.026	-0.654***
##	(0.066)	(0.066)	(0.025)	(0.085)
##				
## year_test_fac2002	-0.427***	-0.979***	0.006	-0.616***
##	(0.066)	(0.067)	(0.019)	(0.083)
##				
## year_test_fac2003	-0.440***	-1.003***	0.008	-0.619***
##	(0.066)	(0.067)	(0.014)	(0.086)
##				
## year_test_fac2004	-0.449***	-0.984***		-0.655***
##	(0.066)	(0.069)		(0.089)
##				
## factor(bac)0.1		0.045**	-0.008	0.004
##		(0.018)	(0.015)	(0.018)
##				
## factor(bac)None		0.062**	0.034**	0.020
##		(0.024)	(0.016)	(0.024)
##				
## factor(sl70plus)1		0.222***	0.020	0.072***
##		(0.022)	(0.020)	(0.022)
##				
## factor(perse_binary)1		-0.019	-0.012	-0.053***
##		(0.015)	(0.015)	(0.015)
##				
## factor(sbprim_binary)1		0.001	-0.013	-0.040
##		(0.025)	(0.023)	(0.025)
##				
## factor(sbsec_binary)1		0.020	-0.012	0.006
##		(0.021)	(0.014)	(0.016)
##				
## factor(gdl_binary)1		-0.021	0.013	-0.015
##		(0.025)	(0.015)	(0.020)
##				
## perc14_24		0.018***	0.040***	0.020*
##		(0.006)	(0.015)	(0.011)
##				
## log(unem)		0.267***	-0.091***	-0.193***
##		(0.024)	(0.023)	(0.023)
##				
## log(vehicmilespc)		1.541***	0.114	0.678***
##		(0.044)	(0.094)	(0.136)
##				
## Constant	3.196***	-11.355***	-0.013***	
##	(0.047)	(0.402)	(0.004)	
##				
## N	1,200	1,200	1,152	
## R2	0.126	0.668	0.177	
## Residual Std. Error	0.325 (df = 1175)	0.202 (df = 1165)		
## =====				
## Notes:		***Significant at the 1 percent level.		
##		**Significant at the 5 percent level.		

```
##
```

```
*Significant at the 10 percent level.
```

7 (10 points) Model Forecasts

The COVID-19 pandemic dramatically changed patterns of driving. Find data (and include this data in your analysis, here) that includes some measure of vehicle miles driven in the US. Your data should at least cover the period from January 2018 to as current as possible. With this data, produce the following statements:

```
miles_data <- read.csv('miles_driven.csv')
```

```
miles_data$miles <- miles_data$M12MTVUSM227NFWA  
miles_data <- na.omit(miles_data)  
miles_data$DATE <- as.Date(miles_data$DATE)
```

Comparing monthly miles driven in 2018 to the same months during the pandemic: - What month demonstrated the largest decrease in driving? How much, in percentage terms, lower was this driving?
- What month demonstrated the largest increase in driving? How much, in percentage terms, higher was this driving?

Month with largest decrease in driving during the pandemic: February 2021 Percentage decrease: -1.39 %

Month with largest increase in driving during the pandemic: February 2020 Percentage increase: 0.23 %

```
# Extracting data for 2018 and pandemic period  
  
# get population to convert to per capita  
# miles are measured in millions, multiple by 1 million  
# take log to reflect the "within" model vehicle miles per capita input  
  
USA_pop <- 333000000  
  
miles_data$miles_pc_log <- log((miles_data$miles*1000000)/USA_pop)  
  
miles_2018 <- subset(miles_data, year(DATE) == 2018)  
  
miles_pandemic <- subset(miles_data, year(DATE) >= 2020)  
  
# Finding corresponding months in 2018 for each month in the pandemic period  
corresponding_months_2018 <- as.Date(paste("2018",  
                                           format(miles_pandemic$DATE, "%m"), "01", sep = "-"))  
  
# Extracting miles driven for corresponding months in 2018  
miles_2018_corresponding <-  
  miles_2018[miles_2018$DATE %in% corresponding_months_2018, ]  
  
# Merging monthly miles data for 2018 and pandemic period  
# Extract month from DATE column  
miles_2018_corresponding$month <- format(miles_2018_corresponding$DATE, "%m")  
miles_pandemic$month <- format(miles_pandemic$DATE, "%m")  
  
# Merge by month  
merged_miles <- merge(miles_2018_corresponding, miles_pandemic, by = "month",  
                      suffixes = c("_2018", "_pandemic"))
```

```

# Calculating percentage change in monthly miles driven
merged_miles$percentage_change <- ((merged_miles$miles_pc_log_pandemic -
    merged_miles$miles_pc_log_2018) /
    merged_miles$miles_pc_log_2018) * 100

# Month with largest decrease in driving during the pandemic
max_decrease_month_pandemic <-
    merged_miles[which.min(merged_miles$percentage_change), "DATE_pandemic"]

# Month with largest increase in driving during the pandemic
max_increase_month_pandemic <-
    merged_miles[which.max(merged_miles$percentage_change), "DATE_pandemic"]

# Output results

# cat("decrease:", format(max_decrease_month_pandemic, "%B %Y"), "\n")
# cat("Percentage decrease:",
# round(min(merged_miles$percentage_change), 2), "%\n\n")
#
# cat("increase:", format(max_increase_month_pandemic, "%B %Y"), "\n")
# cat("Percentage increase:",
# round(max(merged_miles$percentage_change), 2), "%\n\n")

```

Now, use these changes in driving to make forecasts from your models.

- Suppose that the number of miles driven per capita, increased by as much as the COVID boom. Using the FE estimates, what would the consequences be on the number of traffic fatalities? Please interpret the estimate.
- Suppose that the number of miles driven per capita, decreased by as much as the COVID bust. Using the FE estimates, what would the consequences be on the number of traffic fatalities? Please interpret the estimate.

We multiplied the FE estimate by the percent increase in driving during the covid boom to find the impact on traffic fatalities. Estimated increase in traffic fatalities during the COVID boom: 0.1547862 %

We multiplied the FE estimate by the percent increase in driving during the covid bust to find the impact on traffic fatalities. " Estimated decrease in traffic fatalities during the COVID bust: -0.9365583 %

```

# Print the summary of the fixed effects model
# summary(within_model)

# Get the Fixed Effects estimate
fe_estimate <- coef(within_model)["log(vehicmiles_pc)"]

# Maximum increase and decrease in miles driven during the pandemic
max_increase_percentage <- max(merged_miles$percentage_change)

max_decrease_percentage <- min(merged_miles$percentage_change)

# Calculate the estimated increase in traffic fatalities during the COVID boom
fatalities_increase_boom <- fe_estimate * (max_increase_percentage / 100)

# Interpret the estimate for the COVID boom

```

```
# cat("fatalities, COVID boom:", (exp(fatalities_increase_boom) - 1)*100,
#      "%", "\n")

# Calculate the estimated decrease in traffic fatalities during the COVID bust
# fatalities_decrease_bust <- fe_estimate * (max_decrease_percentage/100)

# Interpret the estimate for the COVID bust

# cat("fatalities, COVID bust:", (exp(fatalities_decrease_bust) - 1)*100,
#      "%", "\n")
```

8 (5 points) Evaluate Error

If there were serial correlation or heteroskedasticity in the idiosyncratic errors of the model, what would be the consequences on the estimators and their standard errors? Is there any serial correlation or heteroskedasticity?

We have evidence to suggest that there is heteroscedasticity present in the errors of our fixed effects model after conducting the Breusch-Pagan test on the Fixed Effects model and getting a p-value of 1.795e-06. This means that the variance of the errors in our model is not constant across all levels of the independent variables, which may affect the reliability of our analysis. Heteroskedasticity occurs when the variance of the errors in predicting traffic fatality rates varies across different conditions, violating the assumption of homoscedasticity. This means that the coefficient estimates may not be as precise as they could be if the errors were homoskedastic. Additionally, the standard errors of the coefficient estimates can be biased. Underestimated standard errors can lead to an increased likelihood of Type I errors, while overestimated standard errors can result in a higher likelihood of Type II errors.

```
# Breusch-Pagan Test for Heteroskedasticity
bptest(within_model)
```

```
##
## studentized Breusch-Pagan test
##
## data: within_model
## BP = 86.572, df = 34, p-value = 1.795e-06
```

Answer We have evidence to suggest that there is serial correlation in idiosyncratic errors after conducting the Breusch-Godfrey/Wooldridge test for serial correlation in panel models and getting a p-value of 2.2e-16. Serial correlation means that the errors in predicting traffic fatality rate are not independent across observations. Positive serial correlation can cause coefficient estimates to be less variable than expected, potentially leading to an underestimation of the true effects of variables such as laws and driving behavior on traffic fatalities. Conversely, negative serial correlation can result in coefficient estimates that are more variable, potentially leading to overestimation of these effects. When serial correlation is present, standard errors tend to be underestimated. This means that the confidence intervals around our coefficient estimates may be too narrow, leading to an increased risk of Type I errors, where we mistakenly conclude that there is a significant effect when there is not. We also estimated the model with autocorrelation robust standard errors, which will try and correct for this.

```
pbgttest(within_model)
```

```
##
```

```
## Breusch-Godfrey/Wooldridge test for serial correlation in panel models
##
## data: log(total_fatalities_rate) ~ year_test_fac + factor(bac) + factor(sl70plus) + ...
## chisq = 243.23, df = 25, p-value < 2.2e-16
## alternative hypothesis: serial correlation in idiosyncratic errors
```