

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	1
3	(15 points) Preliminary Model	4
4	(15 points) Expanded Model	4
5	(15 points) State-Level Fixed Effects	5
6	(10 points) Consider a Random Effects Model	5
7	(10 points) Model Forecasts	5
8	(5 points) Evaluate Error	6

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="./data/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 `s155`, `s165`, `s170`, `s175`, 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.)
2. (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?
3. (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.

```
# For the fractions, we are taking the majority as a speed limit
# We skipped year_of_observation since there a year column which aligns with dx
df <- data %>%
  mutate(speed_limit = ifelse(sl55 >= 0.5, 55,
                              ifelse(sl65 >= 0.5, 65,
                              ifelse(sl70 >= 0.5, 70,
                              ifelse(sl75 >= 0.5, 75, 0)))),
          blood_alcohol_limit_10 = ifelse(bac10 >= 0.5, 1, 0),
          blood_alcohol_limit_08 = ifelse(bac08 >= 0.5, 1, 0),
          perse = ifelse(perse >= 0.5, 1, 0)) %>%
  select(!c((sl55:slnone), (d80:d04), bac10, bac08)) %>% # Excluding
  rename(minimum_drinking_age = minage, zero_tolerance_law = zerotol,
          graduated_drivers_license_law = gdl, per_se_law = perse,
          total_fatalities = totfat, nighttime_fatalities = nghtfat,
          weekend_fatalities = wkndfat, total_fatalities_per_100M_miles = totfatpvm,
          nighttime_fatalities_per_100M_miles = nghtfatpvm,
          weekend_fatalities_per_100M_miles = wkndfatpvm,
          state_population = statepop, total_fatalities_rate = totfatrte,
          nighttime_fatalities_rate = nghtfatrte,
          weekend_fatalities_rate = wkndfatrte,
          vehicle_miles_traveled = vehicmiles, unemployment_rate = unem,
          population_aged_14_to_24_rate = perc14_24,
          speed_limit_70_plus = sl70plus, seat_belt = seatbelt,
          primary_seatbelt_law = sbprim, secondary_seatbelt_law = sbsecon,
          miles_driven_per_capita = vehicmilespc)

# Adding states to the dataframe
state_df <- data.frame("index" = 1:51,
                       "state_name" = sort(c(state.name, "District of Columbia")))
```

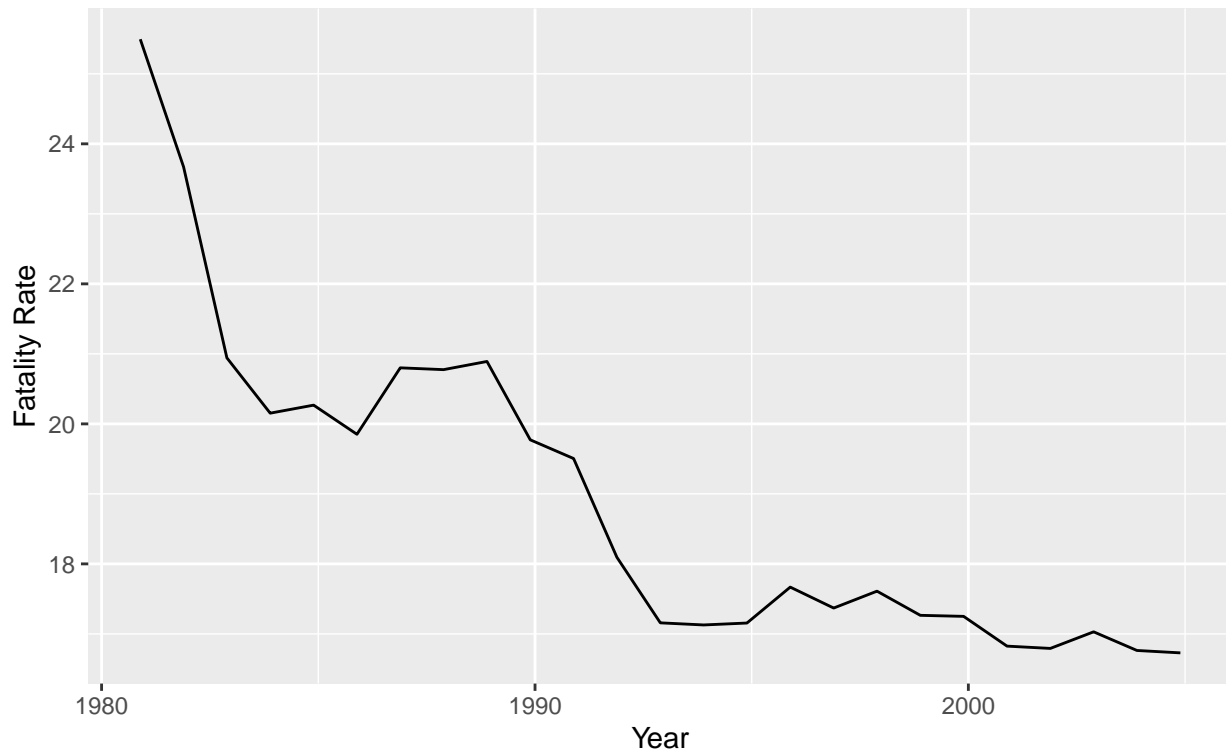
```
df <- merge(df, state_df, by.x = 'state', by.y = 'index')

# Converting data frame to pdata.frame
pdata <- pdata.frame(df, index=c("state", "year"))
#head(pdata)

# Average mean for total_fatalities_rate over years
pdata %>%
  group_by(year) %>%
  summarise(mean = mean(total_fatalities_rate), n = n()) %>%
  ggplot(aes(x = as.Date(year,"%Y"), y = mean)) +
  geom_line() +
  labs(title = "Average mean fatality rate across US",
        subtitle = "Fatality rate is going down",
        x = "Year", y = "Fatality Rate") +
  theme(legend.position = "none")
```

Average mean fatality rate across US

Fatality rate is going down



'In the past few decades, the average fatality rate has gone down in the US. During the late 80s, there was an increase in the fatality rate, which was followed by a decline in a few years.'

```
cut_point <- c(-1, 12, 24, 36, 48)
plots <- vector('list', 4)

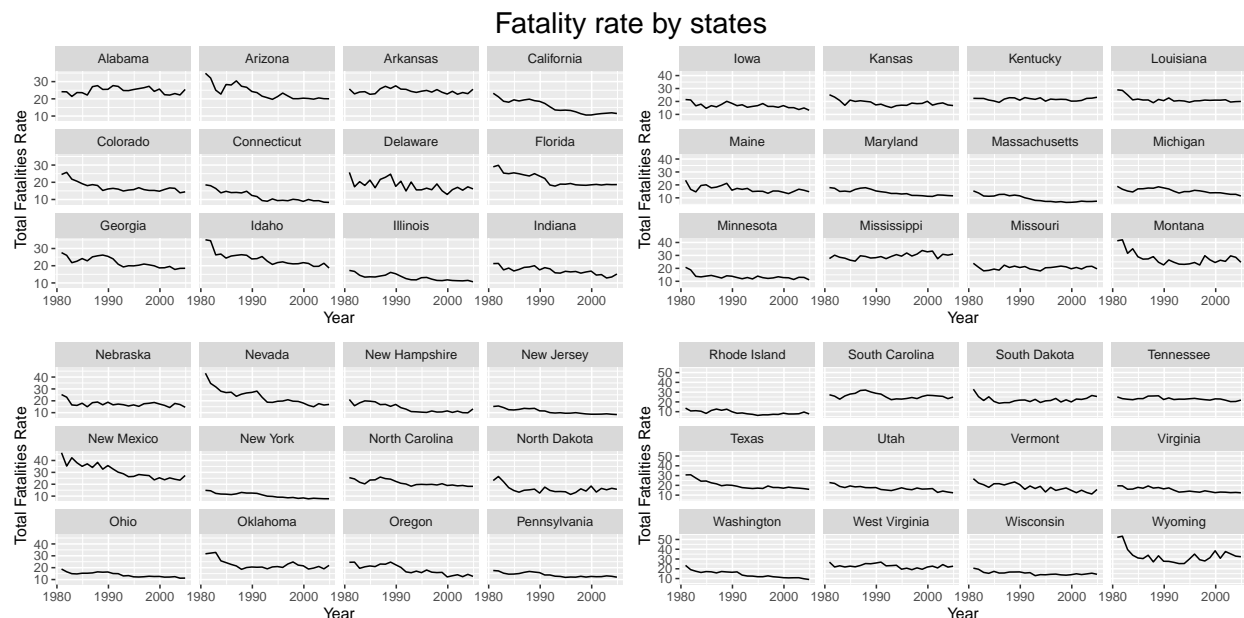
for (i in 2:5) {
  plots[[i-1]] <- (pdata %>%
    filter(as.integer(state) > cut_point[i-1] & as.integer(state) <= cut_point[i]) %>%
    ggplot(aes(x = as.Date(year,"%Y"), y = total_fatalities_rate)) +
```

```

geom_line() +
facet_wrap(~ state_name, nrow = 3, ncol=4) +
labs(x = "Year", y = "Total Fatalities Rate") +
theme(legend.position = "none"))
}

grid.arrange(plots[[1]], plots[[2]], plots[[3]], plots[[4]], nrow = 2, ncol = 2,
top = textGrob("Fatality rate by states", gp=gpar(fontsize=20)))

```



> ‘For most states, fatality rates go down over the years, but some states like Alabama and Arkansas do not show many changes. Surprisingly, Mississippi has an increase in the fatality rate.’

```

# traffic laws that we are exploring are seat_belt, minimum_drinking_age,
# zero_tolerance_law, graduated_drivers_license_law, per_se_law, speed_limit,
# speed_limit_70_plus, primary_seatbelt_law, secondary_seatbelt_law,
# blood_alcohol_limit_10, blood_alcohol_limit_08

```

3 (15 points) Preliminary Model

Estimate a linear regression model of *totfatrate* 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?
- What does this model explain, and what do you find in this model?
- Did driving become safer over this period? Please provide a detailed explanation.
- 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?

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 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.

- How are the blood alcohol variables defined? Interpret the coefficients that you estimate for this concept.
- Do *per se laws* have a negative effect on the fatality rate?
- Does having a primary seat belt law?

5 (15 points) State-Level Fixed Effects

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

- 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?

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?

6 (10 points) Consider a Random Effects Model

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

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

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:

- 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?

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.

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?