

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	3
4	(15 points) Expanded Model	6
5	(15 points) State-Level Fixed Effects	9
6	(10 points) Consider a Random Effects Model	12
7	(10 points) Model Forecasts	14
8	(5 points) Evaluate Error	17
## Warning in system("timedatectl", intern = TRUE): running command 'timedatectl'		
## had status 1		

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

1. (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.)

```
# Create a new column 'speed_limit' and initialize it with 0
data$speed_limit <- 0

# Assign the corresponding speed limit value to 'speed_limit' based on the true condition
data$speed_limit[data$sl55 >= 0.5] <- 55
data$speed_limit[data$sl65 >= 0.5] <- 65
data$speed_limit[data$sl70 >= 0.5] <- 70
data$speed_limit[data$sl75 >= 0.5] <- 75
data$speed_limit[data$slnone >= 0.5] <- NA # DEL make this arbitrarily high?

# Drop the unnecessary speed limit columns
data <- subset(data, select = -c(sl55, sl65, sl70, sl75, slnone))

# Create a year_of_observation variable
data$year_of_observation <- factor(data$year)

# Drop the unnecessary year columns
data <- subset(data, select = -grep("^d\\d{2}$", names(data)))

# Factor state
data$state <- factor(data$state)

# Reencode one-hot variables
data$bac <- 0
data$bac[data$bac08>=0.5] <- .08
data$bac[data$bac10>=0.5] <- .1

data$zeroTolerance <- 0
data$zeroTolerance[data$zerotol>=0.5] <- 1
data$zeroTolerance[data$zerotol<0.5] <- 0

data$minAge <- 0
data$minAge[data$minage>=19.5] <- 21
data$minAge[data$minage<19.5] <- 18

data$perSe <- 0
data$perSe[data$perse>=0.5] = 1
data$perSe[data$perse<0.5] = 0

# Drop the unnecessary columns
data <- subset(data, select = -c(bac08, bac10, zerotol, minage, perse))

# Rename variables
```

```
data <- data %>%
  rename(total_fatality_rate = totfatrte,
         nighttime_fatality_rate = nghtfatrte,
         weekend_fatality_rate = wkndfatrte,
         total_fatalities = totfat,
         nighttime_fatalities = nghtfat,
         weekend_fatalities = wkndfat,
         total_fatalities_per_1mmiles = totfatpvm,
         nighttime_fatalities_per_1mmiles = nghtfatpvm,
         weekend_fatalities_per_1mmiles = wkndfatpvm)

# Log non-normal variables
data$log_fatality_rate = log(data$total_fatality_rate)
data$log_unem = log(data$unem)
data$log_vehicmilespc = log(data$vehicmilespc)
```

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.

```
# DELETE this chunk
```

```
# Missing 2, 9, 12 -- 48 states
table(data$state)
```

```
##
##  1  3  4  5  6  7  8 10 11 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29
## 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25
## 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51
## 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25
```

3 (15 points) Preliminary Model

Estimate a linear regression model of *totfatrte* 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?

```
# Fit linear model
```

```
lm_model <- lm(log_fatality_rate ~ year_of_observation, data = data)
```

```
stargazer(lm_model, title = "Preliminary Model Results", align = TRUE)
```

```
##
```

```
## % Table created by stargazer v.5.2.3 by Marek Hlavac, Social Policy Institute. E-mail: marek.hlavac@spu.cz
```

```
## % Date and time: Tue, Aug 08, 2023 - 02:43:29 AM
```

```
## % Requires LaTeX packages: dcolumn
```

```
## \begin{table}[!htbp] \centering
```

```
## \caption{Preliminary Model Results}
```

```
## \label{}
```

```
## \begin{tabular}{@{\extracolsep{5pt}}lD{.}{.}{-3} }
```

```
## \hline
```

```
## \hline
```

```
## & \multicolumn{1}{c}{\textit{Dependent variable:}} \\\
```

```
## \cline{2-2}
```

```
## \hline & \multicolumn{1}{c}{log\_fatality\_rate} \\\
```

```
## \hline
```

```
## year\_of\_observation1981 & -0.079 \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1982 & -0.200^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1983 & -0.235^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1984 & -0.226^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1985 & -0.243^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1986 & -0.197^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1987 & -0.199^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1988 & -0.189^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1989 & -0.248^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1990 & -0.268^{***} \\\
```

```
## & (0.066) \\\
```

```
## & \\\
```

```
## year\_of\_observation1991 & -0.344^{***} \\\
```

```

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1992 & -0.402^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1993 & -0.403^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1994 & -0.408^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1995 & -0.385^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1996 & -0.399^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1997 & -0.386^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1998 & -0.410^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation1999 & -0.414^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation2000 & -0.437^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation2001 & -0.435^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation2002 & -0.427^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation2003 & -0.440^{***} \\  

##      & (0.066) \\  

##      & \\  

## year\_of\_observation2004 & -0.449^{***} \\  

##      & (0.066) \\  

##      & \\  

## Constant & 3.196^{***} \\  

##      & (0.047) \\  

##      & \\  

## \hline \ll[-1.8ex]  

## Observations & \multicolumn{1}{c}{1,200} \\  

## R\^{2} & \multicolumn{1}{c}{0.126} \\  

## Adjusted R\^{2} & \multicolumn{1}{c}{0.108} \\  

## Residual Std. Error & \multicolumn{1}{c}{0.325 (df = 1175)} \\  

## F Statistic & \multicolumn{1}{c}{7.057^{***}} (df = 24; 1175) \\  

## \hline  

## \hline \ll[-1.8ex]  

## \textit{Note:} & \multicolumn{1}{r}{\^{*}}p<$0.1; \^{**}}p<$0.05; \^{***}}p<$0.01} \\  

## \end{tabular}

```

```
## \end{table}
```

```
# summary(lm_model)
```

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.

- 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?

```
# Fit expanded linear model
```

```
exp_lm_model <- lm(log_fatality_rate ~ year_of_observation + bac + perSe + sbprim + sbsecon +  
                    sl70plus + gdl + perc14_24 + log_unem + log_vehicmilespc, data = data)
```

```
stargazer(exp_lm_model, title = "Expanded Model Results", align = TRUE)
```

```
##
```

```
## % Table created by stargazer v.5.2.3 by Marek Hlavac, Social Policy Institute. E-mail: marek.hlavac@sp.i.cas.cz
```

```
## % Date and time: Tue, Aug 08, 2023 - 02:43:29 AM
```

```
## % Requires LaTeX packages: dcolumn
```

```
## \begin{table}[!htbp] \centering
```

```
## \caption{Expanded Model Results}
```

```
## \label{}
```

```
## \begin{tabular}{@{\extracolsep{5pt}}lD{.}{.}{-3} }
```

```
## \hline
```

```
## \hline \hline
```

```
## & \multicolumn{1}{c}{\textit{Dependent variable:}} \hline
```

```
## \cline{2-2}
```

```
## \hline & \multicolumn{1}{c}{log\_fatality\_rate} \hline
```

```
## \hline \hline
```

```
## year\_of\_observation1981 & -0.092^{**} \hline
```

```
## & (0.041) \hline
```

```
## & \hline
```

```
## year\_of\_observation1982 & -0.293^{***} \hline
```

```
## & (0.042) \hline
```

```
## & \hline
```

```
## year\_of\_observation1983 & -0.348^{***} \hline
```

```
## & (0.043) \hline
```

```
## & \hline
```

```

## year\_of\_observation1984 & -0.301^{***} \\
## & (0.044) \\
## & \\
## year\_of\_observation1985 & -0.340^{***} \\
## & (0.045) \\
## & \\
## year\_of\_observation1986 & -0.318^{***} \\
## & (0.046) \\
## & \\
## year\_of\_observation1987 & -0.355^{***} \\
## & (0.048) \\
## & \\
## year\_of\_observation1988 & -0.367^{***} \\
## & (0.051) \\
## & \\
## year\_of\_observation1989 & -0.454^{***} \\
## & (0.053) \\
## & \\
## year\_of\_observation1990 & -0.513^{***} \\
## & (0.054) \\
## & \\
## year\_of\_observation1991 & -0.628^{***} \\
## & (0.055) \\
## & \\
## year\_of\_observation1992 & -0.735^{***} \\
## & (0.056) \\
## & \\
## year\_of\_observation1993 & -0.727^{***} \\
## & (0.057) \\
## & \\
## year\_of\_observation1994 & -0.718^{***} \\
## & (0.058) \\
## & \\
## year\_of\_observation1995 & -0.699^{***} \\
## & (0.059) \\
## & \\
## year\_of\_observation1996 & -0.822^{***} \\
## & (0.061) \\
## & \\
## year\_of\_observation1997 & -0.843^{***} \\
## & (0.063) \\
## & \\
## year\_of\_observation1998 & -0.887^{***} \\
## & (0.064) \\
## & \\
## year\_of\_observation1999 & -0.889^{***} \\
## & (0.065) \\
## & \\
## year\_of\_observation2000 & -0.902^{***} \\
## & (0.066) \\
## & \\
## year\_of\_observation2001 & -0.958^{***} \\
## & (0.066) \\
## & \\

```

```

## year\_of\_observation2002 & -1.008^{***} \\
## & (0.066) \\
## & \\
## year\_of\_observation2003 & -1.035^{***} \\
## & (0.066) \\
## & \\
## year\_of\_observation2004 & -1.032^{***} \\
## & (0.067) \\
## & \\
## bac & -0.084 \\
## & (0.191) \\
## & \\
## perSe & -0.030^{**} \\
## & (0.014) \\
## & \\
## sbprim & 0.001 \\
## & (0.025) \\
## & \\
## sbsecon & 0.026 \\
## & (0.021) \\
## & \\
## sl70plus & 0.233^{***} \\
## & (0.022) \\
## & \\
## gdl & -0.026 \\
## & (0.026) \\
## & \\
## perc14\_24 & 0.017^{***} \\
## & (0.006) \\
## & \\
## log\_unem & 0.262^{***} \\
## & (0.024) \\
## & \\
## log\_vehicmilespc & 1.534^{***} \\
## & (0.044) \\
## & \\
## Constant & -11.215^{***} \\
## & (0.402) \\
## & \\
## \hline \\[-1.8ex]
## Observations & \multicolumn{1}{c}{1,200} \\
## R^{2}$ & \multicolumn{1}{c}{0.668} \\
## Adjusted R^{2}$ & \multicolumn{1}{c}{0.658} \\
## Residual Std. Error & \multicolumn{1}{c}{0.201 (df = 1166)} \\
## F Statistic & \multicolumn{1}{c}{70.951^{***}$ (df = 33; 1166)} \\
## \hline
## \hline \\[-1.8ex]
## \textit{Note:} & \multicolumn{1}{r}{$^{*}$p$<$0.1; $^{**}$p$<$0.05; $^{***}$p$<$0.01} \\
## \end{tabular}
## \end{table}

```

```
# summary(exp_lm_model)
```

```
# DEL confirmed there is nothing other than 0/1 indicators for sbprim and sbsecond
```



```

# table(data$sbprim)
# table(data$sbsecon)

# Calculate the effect of BAC
(1 - exp(coefficients(exp_lm_model)["bac"]/100)) * 100

##          bac
## 0.08405037

# Calculate the effect of per se
(1 - exp(coefficients(exp_lm_model)["perSe"])) * 100

##      perSe
## 2.930884

```

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?

```

# Convert data to a plm dataframe
# DEL test if get same result using plm.data()
data_plm <- pdata.frame(data, index = c("state", "year_of_observation"))

# Fit fixed effects model
fixed_model <- plm(log_fatality_rate ~ year_of_observation + bac + perSe + sbprim + sbsecon +
                    sl70plus + gdl + perc14_24 + log_unem + log_vehicmilespc,
                    data = data_plm,
                    model = "within",
                    effect = "individual")

stargazer(fixed_model, title = "Fixed Effects Model Results", align = TRUE)

##
## % Table created by stargazer v.5.2.3 by Marek Hlavac, Social Policy Institute. E-mail: marek.hlavac@sp.i.cas.cz
## % Date and time: Tue, Aug 08, 2023 - 02:43:29 AM
## % Requires LaTeX packages: dcolumn
## \begin{table}[!htbp] \centering
##   \caption{Fixed Effects Model Results}
##   \label{}
##   \begin{tabular}{@{\extracolsep{5pt}}lD{.}{.}{-3} }
##     \hline
##     \hline
##     & \multicolumn{1}{c}{\textit{Dependent variable:}} & \\
##     \cline{2-2}

```

```

## \[-1.8ex] & \multicolumn{1}{c}{\log\_fatality\_rate} \\
## \hline \[-1.8ex]
## year\_of\_observation1981 & -0.063^{***} \\
## & (0.018) \\
## & \\
## year\_of\_observation1982 & -0.135^{***} \\
## & (0.019) \\
## & \\
## year\_of\_observation1983 & -0.168^{***} \\
## & (0.020) \\
## & \\
## year\_of\_observation1984 & -0.208^{***} \\
## & (0.021) \\
## & \\
## year\_of\_observation1985 & -0.233^{***} \\
## & (0.021) \\
## & \\
## year\_of\_observation1986 & -0.197^{***} \\
## & (0.023) \\
## & \\
## year\_of\_observation1987 & -0.243^{***} \\
## & (0.025) \\
## & \\
## year\_of\_observation1988 & -0.274^{***} \\
## & (0.027) \\
## & \\
## year\_of\_observation1989 & -0.349^{***} \\
## & (0.029) \\
## & \\
## year\_of\_observation1990 & -0.358^{***} \\
## & (0.030) \\
## & \\
## year\_of\_observation1991 & -0.395^{***} \\
## & (0.031) \\
## & \\
## year\_of\_observation1992 & -0.456^{***} \\
## & (0.032) \\
## & \\
## year\_of\_observation1993 & -0.474^{***} \\
## & (0.033) \\
## & \\
## year\_of\_observation1994 & -0.507^{***} \\
## & (0.034) \\
## & \\
## year\_of\_observation1995 & -0.508^{***} \\
## & (0.035) \\
## & \\
## year\_of\_observation1996 & -0.559^{***} \\
## & (0.037) \\
## & \\
## year\_of\_observation1997 & -0.586^{***} \\
## & (0.038) \\
## & \\
## year\_of\_observation1998 & -0.639^{***} \\

```

```

## & (0.038) \\  

## & \\  

## year\_of\_observation1999 & -0.657^{***} \\  

## & (0.039) \\  

## & \\  

## year\_of\_observation2000 & -0.690^{***} \\  

## & (0.040) \\  

## & \\  

## year\_of\_observation2001 & -0.659^{***} \\  

## & (0.040) \\  

## & \\  

## year\_of\_observation2002 & -0.622^{***} \\  

## & (0.040) \\  

## & \\  

## year\_of\_observation2003 & -0.625^{***} \\  

## & (0.040) \\  

## & \\  

## year\_of\_observation2004 & -0.664^{***} \\  

## & (0.040) \\  

## & \\  

## bac & -0.099 \\  

## & (0.108) \\  

## & \\  

## perSe & -0.057^{***} \\  

## & (0.009) \\  

## & \\  

## sbprim & -0.040^{***} \\  

## & (0.015) \\  

## & \\  

## sbsecon & 0.006 \\  

## & (0.011) \\  

## & \\  

## sl70plus & 0.077^{***} \\  

## & (0.012) \\  

## & \\  

## gdl & -0.021^{*} \\  

## & (0.013) \\  

## & \\  

## perc14\_24 & 0.019^{***} \\  

## & (0.004) \\  

## & \\  

## log\_unem & -0.193^{***} \\  

## & (0.017) \\  

## & \\  

## log\_vehicmilespc & 0.678^{***} \\  

## & (0.051) \\  

## & \\  

## \hline \)[-1.8ex]  

## Observations & \multicolumn{1}{c}{1,200} \\  

## R^{2}$ & \multicolumn{1}{c}{0.729} \\  

## Adjusted R^{2}$ & \multicolumn{1}{c}{0.709} \\  

## F Statistic & \multicolumn{1}{c}{91.037^{***}}$ (df = 33; 1119)} \\  

## \hline  

## \hline \)[-1.8ex]

```

```
## \textit{Note:} & \multicolumn{1}{r}{\mathrel{\sim}^*p<0.1; \mathrel{\sim}^{**}p<0.05; \mathrel{\sim}^{***}p<0.01} \\
## \end{tabular}
## \end{table}

# summary(fixed_model)

# ?pdata.frame

# Calculate the effect of per se
(1 - exp(coefficients(fixed_model)["perSe"])) * 100

##      perSe
## 5.539437

# Calculate the effect of sbprim
(1 - exp(coefficients(fixed_model)["sbprim"])) * 100

##      sbprim
## 3.927938
```

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?

```
# DEL don't print out these results in the final report because don't actually need to estimate

# Fit random effects model
random_model <- plm(log_fatality_rate ~ year_of_observation + bac + perSe + sbprim + sbsecon +
                    sl70plus + gdl + perc14_24 + log_unem + log_vehicmilespc,
                    data = data_plm,
                    model = "random")

summary(random_model)

## Oneway (individual) effect Random Effect Model
##      (Swamy-Arora's transformation)
##
## Call:
## plm(formula = log_fatality_rate ~ year_of_observation + bac +
##      perSe + sbprim + sbsecon + sl70plus + gdl + perc14_24 + log_unem +
##      log_vehicmilespc, data = data_plm, model = "random")
##
## Balanced Panel: n = 48, T = 25, N = 1200
##
## Effects:
##              var  std.dev share
## idiosyncratic 0.007743 0.087992 0.241
```

```

## individual      0.024430 0.156301 0.759
## theta: 0.8881
##
## Residuals:
##      Min.      1st Qu.      Median      3rd Qu.      Max.
## -0.4132201 -0.0530907  0.0028669  0.0569146  0.2905780
##
## Coefficients:
##              Estimate Std. Error z-value Pr(>|z|)
## (Intercept)    -3.5957404   0.4536937  -7.9255 2.273e-15 ***
## year_of_observation1981 -0.0643840   0.0185681  -3.4674 0.0005254 ***
## year_of_observation1982 -0.1419541   0.0194893  -7.2837 3.248e-13 ***
## year_of_observation1983 -0.1767199   0.0202670  -8.7196 < 2.2e-16 ***
## year_of_observation1984 -0.2136545   0.0210815 -10.1347 < 2.2e-16 ***
## year_of_observation1985 -0.2404522   0.0220349 -10.9123 < 2.2e-16 ***
## year_of_observation1986 -0.2056173   0.0235850  -8.7181 < 2.2e-16 ***
## year_of_observation1987 -0.2535212   0.0255676  -9.9157 < 2.2e-16 ***
## year_of_observation1988 -0.2852297   0.0279528 -10.2040 < 2.2e-16 ***
## year_of_observation1989 -0.3614884   0.0297716 -12.1421 < 2.2e-16 ***
## year_of_observation1990 -0.3742853   0.0309269 -12.1022 < 2.2e-16 ***
## year_of_observation1991 -0.4142437   0.0316410 -13.0920 < 2.2e-16 ***
## year_of_observation1992 -0.4777209   0.0326760 -14.6199 < 2.2e-16 ***
## year_of_observation1993 -0.4950929   0.0332864 -14.8737 < 2.2e-16 ***
## year_of_observation1994 -0.5267332   0.0341531 -15.4227 < 2.2e-16 ***
## year_of_observation1995 -0.5279114   0.0351931 -15.0004 < 2.2e-16 ***
## year_of_observation1996 -0.5822055   0.0371701 -15.6633 < 2.2e-16 ***
## year_of_observation1997 -0.6093535   0.0383909 -15.8723 < 2.2e-16 ***
## year_of_observation1998 -0.6626534   0.0391542 -16.9242 < 2.2e-16 ***
## year_of_observation1999 -0.6812420   0.0396590 -17.1775 < 2.2e-16 ***
## year_of_observation2000 -0.7132495   0.0402404 -17.7247 < 2.2e-16 ***
## year_of_observation2001 -0.6871045   0.0403998 -17.0076 < 2.2e-16 ***
## year_of_observation2002 -0.6536951   0.0404178 -16.1735 < 2.2e-16 ***
## year_of_observation2003 -0.6579505   0.0404430 -16.2686 < 2.2e-16 ***
## year_of_observation2004 -0.6962559   0.0408052 -17.0629 < 2.2e-16 ***
## bac              -0.1139066   0.1106851  -1.0291 0.3034303
## perSe            -0.0551111   0.0096492  -5.7115 1.120e-08 ***
## sbprim           -0.0380542   0.0152926  -2.4884 0.0128315 *
## sbsecon           0.0070569   0.0112487   0.6274 0.5304255
## sl70plus          0.0808570   0.0120242   6.7245 1.762e-11 ***
## gdl              -0.0207167   0.0130676  -1.5854 0.1128867
## perc14_24         0.0198780   0.0042437   4.6841 2.812e-06 ***
## log_unem          -0.1743113   0.0174342  -9.9982 < 2.2e-16 ***
## log_vehicmilespc  0.7619057   0.0501253  15.2000 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Total Sum of Squares:    33.307
## Residual Sum of Squares: 9.5676
## R-Squared:              0.71274
## Adj. R-Squared:         0.70461
## Chisq: 2893.09 on 33 DF, p-value: < 2.22e-16
# Test for random effects
phtest(fixed_model, random_model)

```

```
##
## Hausman Test
##
## data: log_fatality_rate ~ year_of_observation + bac + perSe + sbprim + ...
## chisq = 83.713, df = 33, p-value = 2.731e-06
## alternative hypothesis: one model is inconsistent
```

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.

```
monthly_miles <- read_csv("../data/monthly_miles_driven.csv")
```

```
## Rows: 65 Columns: 2
## -- Column specification -----
## Delimiter: ","
## chr (1): date
## dbl (1): millions_of_miles
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
# Show first few rows
head(monthly_miles)
```

```
## # A tibble: 6 x 2
##   date      millions_of_miles
##   <chr>          <dbl>
## 1 1/1/2018      244736
## 2 2/1/2018      227759
## 3 3/1/2018      270705
## 4 4/1/2018      275127
## 5 5/1/2018      283713
## 6 6/1/2018      282648
```

```
monthly_population <- read_csv("../data/POPTHM.csv")
```

```
## Rows: 66 Columns: 2
## -- Column specification -----
## Delimiter: ","
```

```

## dbl (1): POPTHM
## date (1): DATE
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
# Show first few rows
head(monthly_population)

## # A tibble: 6 x 2
##   DATE      POPTHM
##   <date>    <dbl>
## 1 2018-01-01 327969
## 2 2018-02-01 328085
## 3 2018-03-01 328219
## 4 2018-04-01 328364
## 5 2018-05-01 328521
## 6 2018-06-01 328692

# Reformat data for comparisons
monthly_miles$date <- lubridate::mdy(monthly_miles$date)
monthly_miles$year <- lubridate::year(monthly_miles$date)
monthly_miles$month <- lubridate::month(monthly_miles$date)

monthly_population$year <- lubridate::year(monthly_population$DATE)
monthly_population$month <- lubridate::month(monthly_population$DATE)

monthly_miles <- monthly_miles %>%
  left_join(monthly_population, by=c("year", "month")) %>%
  mutate(miles_per_capita = millions_of_miles * 1000 / POPTHM)

monthly_miles_2018 <- monthly_miles %>%
  filter(date >= "2018-01-01" & date <= "2018-12-31") %>%
  select(date, year, month, millions_of_miles_nonpand = miles_per_capita)

monthly_miles_pand <- monthly_miles %>%
  filter(date >= "2020-03-01") %>%
  select(date, year, month, millions_of_miles_pand = miles_per_capita)

comparison_data <- monthly_miles_pand %>%
  left_join(monthly_miles_2018, by = "month") %>%
  select(date = date.x, year = year.x, month, millions_of_miles_pand, millions_of_miles_nonpand)

comparison_data$total_change <- (comparison_data$millions_of_miles_pand - comparison_data$millions_of_miles_nonpand)
comparison_data$total_log_change <- (log(comparison_data$millions_of_miles_pand) - log(comparison_data$millions_of_miles_nonpand))
comparison_data$perc_change <- (comparison_data$millions_of_miles_pand - comparison_data$millions_of_miles_nonpand) / comparison_data$millions_of_miles_nonpand

# What month demonstrated the largest decrease in driving? How much, in percentage terms, lower was this than the month before?
covid_bust_date <- comparison_data$date[which.min(comparison_data$perc_change)]
print(covid_bust_date)

## [1] "2020-04-01"

```

```

covid_bust_perc <- min(comparison_data$perc_change)
print(covid_bust_perc)

## [1] -39.68631

covid_bust_tot <- min(comparison_data$total_log_change)
print(covid_bust_tot)

## [1] -0.5056112

#What month demonstrated the largest increase in driving? How much, in percentage terms, higher was this than the bust?
covid_boom_date <- comparison_data$date[which.max(comparison_data$perc_change)]
print(covid_boom_date)

## [1] "2022-09-01"

covid_boom_perc <- max(comparison_data$perc_change)
print(covid_boom_perc)

## [1] 0.6897441

covid_boom_tot <- max(comparison_data$total_log_change)
print(covid_boom_tot)

## [1] 0.006873763

log_miles_pc_coef <- fixed_model$coefficients[33] %>% unname()

log_miles_pc_sd <- sqrt(fixed_model$vcov[33,33])
z <- qnorm(.975)

lower_bust <- (log_miles_pc_coef - z*log_miles_pc_sd) * covid_bust_tot
est_bust <- log_miles_pc_coef * covid_bust_tot
upper_bust <- (log_miles_pc_coef + z*log_miles_pc_sd) * covid_bust_tot

lower_boom <- (log_miles_pc_coef - z*log_miles_pc_sd) * covid_boom_tot
est_boom <- log_miles_pc_coef * covid_boom_tot
upper_boom <- (log_miles_pc_coef + z*log_miles_pc_sd) * covid_boom_tot

round_digits <-4

forecast_ci <- data.frame(
  scenario= c("Boom", "Bust"),
  month = c(covid_boom_date, covid_bust_date),
  lower = c(round(lower_boom, round_digits), round(lower_bust, round_digits)),
  estimate = c(round(est_boom, round_digits), round(est_bust, round_digits)),
  upper = c(round(upper_boom, round_digits), round(upper_bust, round_digits))
)

forecast_ci

##   scenario      month  lower estimate  upper
## 1      Boom 2022-09-01  0.0040   0.0047  0.0053
## 2      Bust 2020-04-01 -0.2923  -0.3426 -0.3929

```


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?

```
pcdtest(fixed_model, test = "lm")
```

```
##
```

```
## Breusch-Pagan LM test for cross-sectional dependence in panels
```

```
##
```

```
## data: log_fatality_rate ~ year_of_observation + bac + perSe + sbprim + sbsecon + sl70plus + gdl
```

```
## chisq = 2748.1, df = 1128, p-value < 2.2e-16
```

```
## alternative hypothesis: cross-sectional dependence
```

```
pbgtest(fixed_model)
```

```
##
```

```
## Breusch-Godfrey/Wooldridge test for serial correlation in panel models
```

```
##
```

```
## data: log_fatality_rate ~ year_of_observation + bac + perSe + sbprim + ...
```

```
## chisq = 243.21, df = 25, p-value < 2.2e-16
```

```
## alternative hypothesis: serial correlation in idiosyncratic errors
```

```
pdwtest(fixed_model)
```

```
##
```

```
## Durbin-Watson test for serial correlation in panel models
```

```
##
```

```
## data: log_fatality_rate ~ year_of_observation + bac + perSe + sbprim + ...
```

```
## DW = 1.2138, p-value < 2.2e-16
```

```
## alternative hypothesis: serial correlation in idiosyncratic errors
```