

Statistical Methods for Discrete Response, Time Series, and Panel Data (W271): Group Lab 3

Instructions (Please Read Carefully):

- **Due Date: Sunday 04/19/20 11:59pm**
- 20 page limit
- Do not modify fontsize, margin or line-spacing settings
- One student from each group should submit the lab to their student github repo by the deadline; submission and revisions made after the deadline will not be graded
- Answers should clearly explain your reasoning; do not simply ‘output dump’ the results of code without explanation
- Submit two files:
 1. A pdf file that details your answers. Include all R code used to produce the answers. Do not suppress the codes in your pdf file
 2. The R markdown (Rmd) file used to produce the pdf file

The assignment will not be graded unless **both** files are submitted

- Name your files to include all group members names. For example the students’ names are Stan Cartman and Kenny Kyle, name your files as follows:
 - StanCartman_KennyKyle_Lab3.Rmd
 - StanCartman_KennyKyle_Lab3.pdf
- Although it sounds obvious, please write your names on page 1 of your pdf and Rmd files
- For statistical methods that we cover in this course, use the R libraries and functions that are covered in this course. If you use libraries and functions for statistical modeling that we have not covered, you must provide an explanation of why such libraries and functions are used and reference the library documentation. For data wrangling and data visualization, you are free to use other libraries, such as `dplyr`, `ggplot2`, etc.
- Your report needs to include:
 - A thorough analysis of the given dataset, which include examination of anomalies, missing values, potential of top and/or bottom code, and other potential anomalies, in each of the variables.
 - A comprehensive Exploratory Data Analysis (EDA) analysis, which includes both graphical and tabular analysis, as taught in this course. Output-dump (that is, graphs and tables that don’t come with explanations) will result in a very low, if not zero, score. Be selective when choosing visuals and tables to illustrate your key points and concise with your explanations (please do not ramble).
 - A proper narrative for each question answered. Make sure that your audience can easily follow the logic of your analysis and the rationale of decisions made in your modeling,

supported by empirical evidence. Use the insights generated from your EDA step to guide your modeling approach.

- Clear explanations of all steps used to arrive at a final model, with conclusions that summarize results with respect to the question(s) being asked and key takeaways from the analysis.
- For mathematical formulae, type them in your R markdown file. Do not e.g. write them on a piece of paper, snap a photo, and use the image file.
- Incorrectly following submission instructions results in deduction of grades
- Students are expected to act with regard to UC Berkeley Academic Integrity

U.S. traffic fatalities: 1980-2004

In this lab, you are asked to answer the question “**Do changes in traffic laws affect traffic fatalities?**” To do so, you will conduct the tasks specified below using the data set *driving.Rdata*, which 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 come with the dataste.

Exercises:

1. (40%) Load the data. Provide a description of the basic structure of the dataset, as we have done throughout the semester. Conduct a very thorough EDA, which should include both graphical and tabular techniques, on the dataset, including both the dependent variable *totfatrte* and the potential explanatory variables. You need to write a detailed narrative of your observations of your EDA. *Reminder: giving an “output dump” (i.e. providing a bunch of graphs and tables without description and hoping your audience will interpret them) will receive a zero in this exercise.*

```
load(file = "driving.RData")
driving <- data
sum(is.na(driving))
```

```
## [1] 0
```

```
nrow(driving)
```

```
## [1] 1200
```

```
colnames(driving)
```

```
## [1] "year"      "state"      "sl55"       "sl65"
## [5] "sl70"      "sl75"       "slnone"     "seatbelt"
## [9] "minage"    "zerotol"    "gdl"        "bac10"
## [13] "bac08"     "perse"      "totfat"     "nghtfat"
## [17] "wkndfat"   "totfatpvm"  "nghtfatpvm" "wkndfatpvm"
## [21] "statepop"  "totfatrte"  "nghtfatrte" "wkndfatrte"
## [25] "vehicmiles" "unem"       "perc14_24"  "sl70plus"
## [29] "sbprim"    "sbsecon"    "d80"        "d81"
## [33] "d82"       "d83"        "d84"        "d85"
## [37] "d86"       "d87"        "d88"        "d89"
## [41] "d90"       "d91"        "d92"        "d93"
## [45] "d94"       "d95"        "d96"        "d97"
## [49] "d98"       "d99"        "d00"        "d01"
## [53] "d02"       "d03"        "d04"        "vehicmilespc"
```

There are no missing values in the dataset and there are 1200 observations. Here the subjects are indicated by the state variable. For each state there are multiple observations, one for each year.

```
head(driving,10)
```

```
##      year state  sl55  sl65 sl70 sl75 slnone seatbelt minage zerotol gdl
## 1  1980      1 1.000 0.000    0    0      0      0      18      0    0
## 2  1981      1 1.000 0.000    0    0      0      0      18      0    0
## 3  1982      1 1.000 0.000    0    0      0      0      18      0    0
## 4  1983      1 1.000 0.000    0    0      0      0      18      0    0
## 5  1984      1 1.000 0.000    0    0      0      0      18      0    0
## 6  1985      1 1.000 0.000    0    0      0      0      20      0    0
## 7  1986      1 1.000 0.000    0    0      0      0      21      0    0
## 8  1987      1 0.542 0.458    0    0      0      0      21      0    0
## 9  1988      1 0.000 1.000    0    0      0      0      21      0    0
## 10 1989      1 0.000 1.000    0    0      0      0      21      0    0
##      bac10 bac08 perse totfat nghtfat wkndfat totfatpvm nghtfatpvm
## 1      1      0      0    940    422    236      3.200      1.437
## 2      1      0      0    933    434    248      3.350      1.558
## 3      1      0      0    839    376    224      2.810      1.259
## 4      1      0      0    930    397    223      3.000      1.281
## 5      1      0      0    932    421    237      2.830      1.278
## 6      1      0      0    882    358    224      2.510      1.019
## 7      1      0      0   1080    500    279      3.177      1.471
## 8      1      0      0   1111    499    300      2.970      1.334
## 9      1      0      0   1024    423    226      2.580      1.066
## 10     1      0      0   1029    418    247      2.520      1.024
##      wkndfatpvm statepop totfatrte nghtfatrte wkndfatrte vehicmiles unem
## 1      0.803  3893888      24.14      10.84      6.06  29.37500  8.8
## 2      0.890  3918520      24.07      11.08      6.33  27.85200 10.7
## 3      0.750  3925218      21.37      9.58      5.71  29.85765 14.4
## 4      0.719  3934109      23.64      10.09      5.67  31.00000 13.7
## 5      0.720  3951834      23.58      10.65      6.00  32.93286 11.1
## 6      0.637  3972527      22.20      9.01      5.64  35.13944  8.9
## 7      0.821  3991569      27.08      12.53      6.99  33.99371  9.8
## 8      0.802  4015261      27.67      12.43      7.47  37.40741  7.8
## 9      0.569  4023858      25.45      10.51      5.62  39.68992  7.2
## 10     0.605  4030229      25.53      10.37      6.13  40.83333  7.0
##      perc14_24 sl70plus sbprim sbsecon d80 d81 d82 d83 d84 d85 d86 d87 d88
## 1      18.9      0      0      0    1    0    0    0    0    0    0    0    0
## 2      18.7      0      0      0    0    1    0    0    0    0    0    0    0
## 3      18.4      0      0      0    0    0    1    0    0    0    0    0    0
## 4      18.0      0      0      0    0    0    0    1    0    0    0    0    0
## 5      17.6      0      0      0    0    0    0    0    1    0    0    0    0
## 6      17.3      0      0      0    0    0    0    0    0    1    0    0    0
## 7      17.0      0      0      0    0    0    0    0    0    0    1    0    0
## 8      16.6      0      0      0    0    0    0    0    0    0    0    1    0
## 9      16.2      0      0      0    0    0    0    0    0    0    0    0    1
## 10     15.8      0      0      0    0    0    0    0    0    0    0    0    0
##      d89 d90 d91 d92 d93 d94 d95 d96 d97 d98 d99 d00 d01 d02 d03 d04
```

```
## 1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 2  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 3  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 4  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 5  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 6  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 7  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 8  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 9  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
## 10 1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
##      vehicmilespc
## 1      7543.874
## 2      7107.785
## 3      7606.622
## 4      7879.802
## 5      8333.562
## 6      8845.614
## 7      8516.377
## 8      9316.308
## 9      9863.649
## 10     10131.764
```

```
table(driving$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
## 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25 25
## 29 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 25 25
```

```
table(driving$year)
```

```
##
## 1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994
##   48   48   48   48   48   48   48   48   48   48   48   48   48   48   48
## 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004
##   48   48   48   48   48   48   48   48   48   48
```

For each state there are 25 observations one for each year. For each year there are 48 observations, one for each state. *This panel is balanced.* There are a total of 1200 observations and since $1200 = 25 \times 48$ we have a balanced panel. d80 through d04 are the indicator variables for each of the time periods.

We cannot use OLS here because we suspect violation of the independence assumptions. Let us confirm that using the Durbin-Watson test.

```
library(lmtest)
```

```
## Loading required package: zoo
```

```
##
```

```
## Attaching package: 'zoo'

## The following object is masked from 'package:tsibble':
##
##      index

## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric

dwtest(totfatrte ~ seatbelt + zerotol + slnone, data=driving)
```

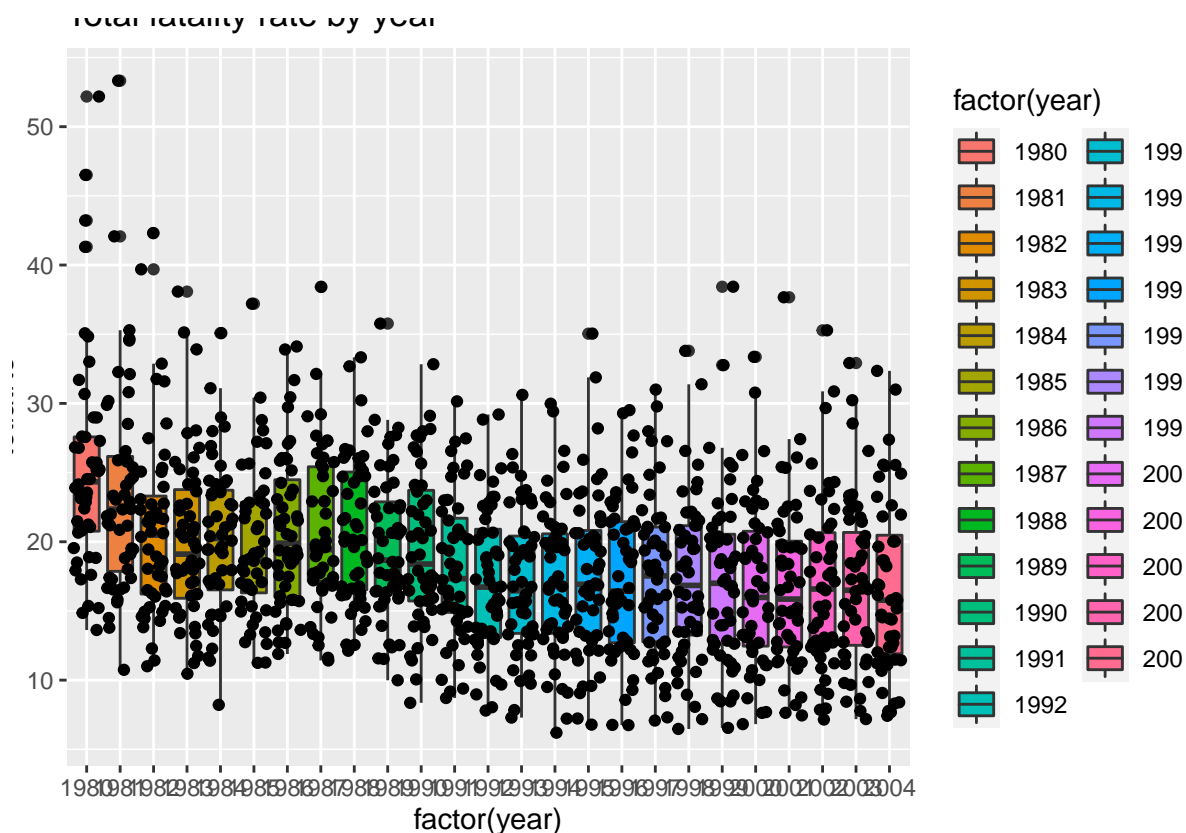
```
##
## Durbin-Watson test
##
## data: totfatrte ~ seatbelt + zerotol + slnone
## DW = 0.24975, p-value < 2.2e-16
## alternative hypothesis: true autocorrelation is greater than 0
```

Null hypothesis is rejected, this confirms the violation of the independence assumption.

So we will have to use panel methods here. There are 25 panels in this dataset, one for each year of observations.

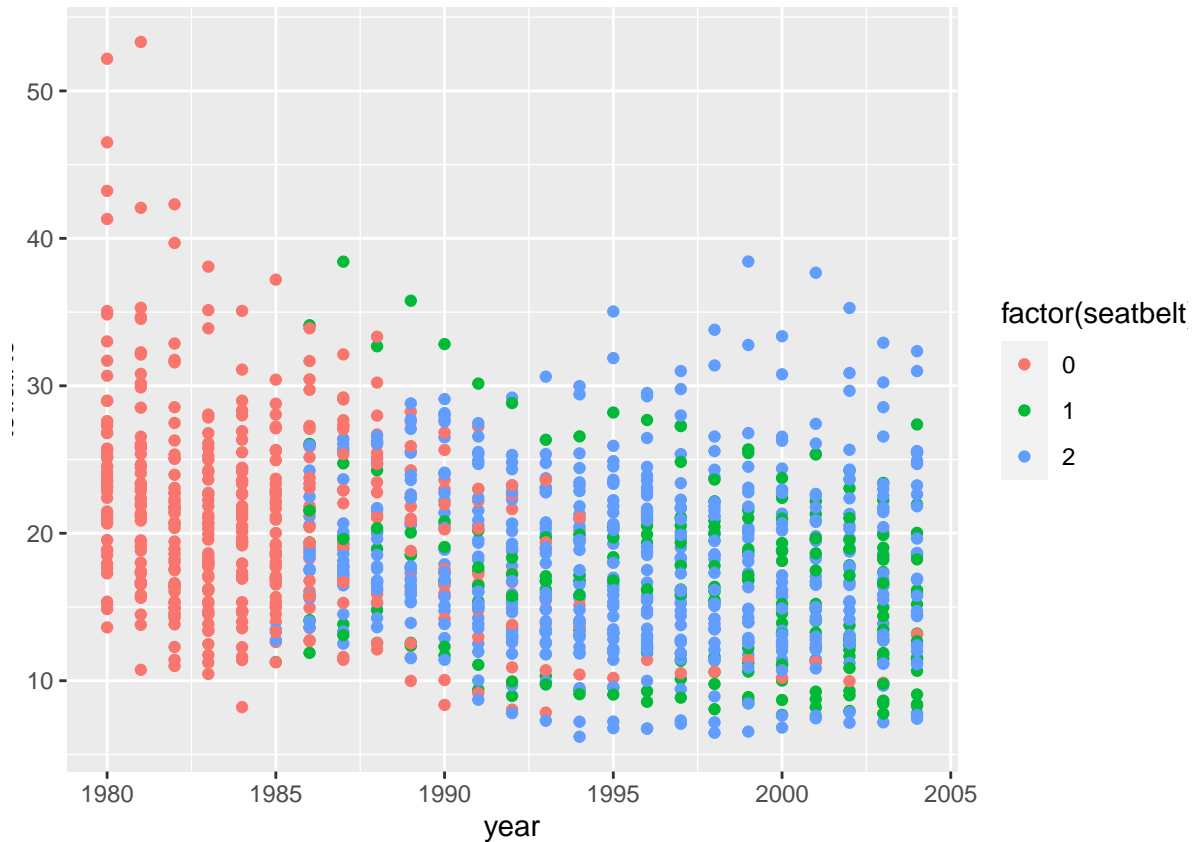
Let's see the response variable (total fatality rate) distribution broken down by panels

```
ggplot(driving, aes(factor(year), totfatrte)) + geom_boxplot(aes(fill = factor(year))) + geom_jitter
```



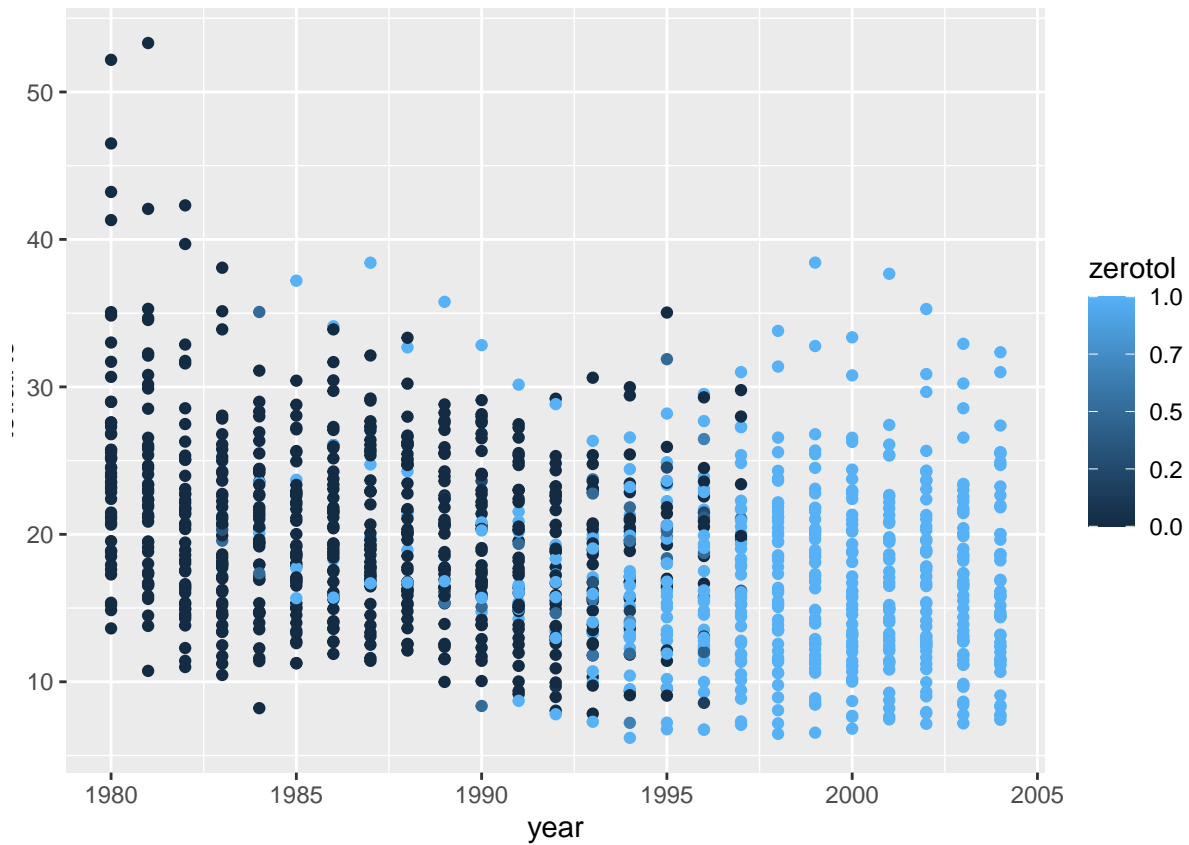
Overall, the total fatality rate has dropped over the years.

```
ggplot(driving, aes(year, totfatrte)) + geom_point(aes(color = factor(seatbelt)))
```



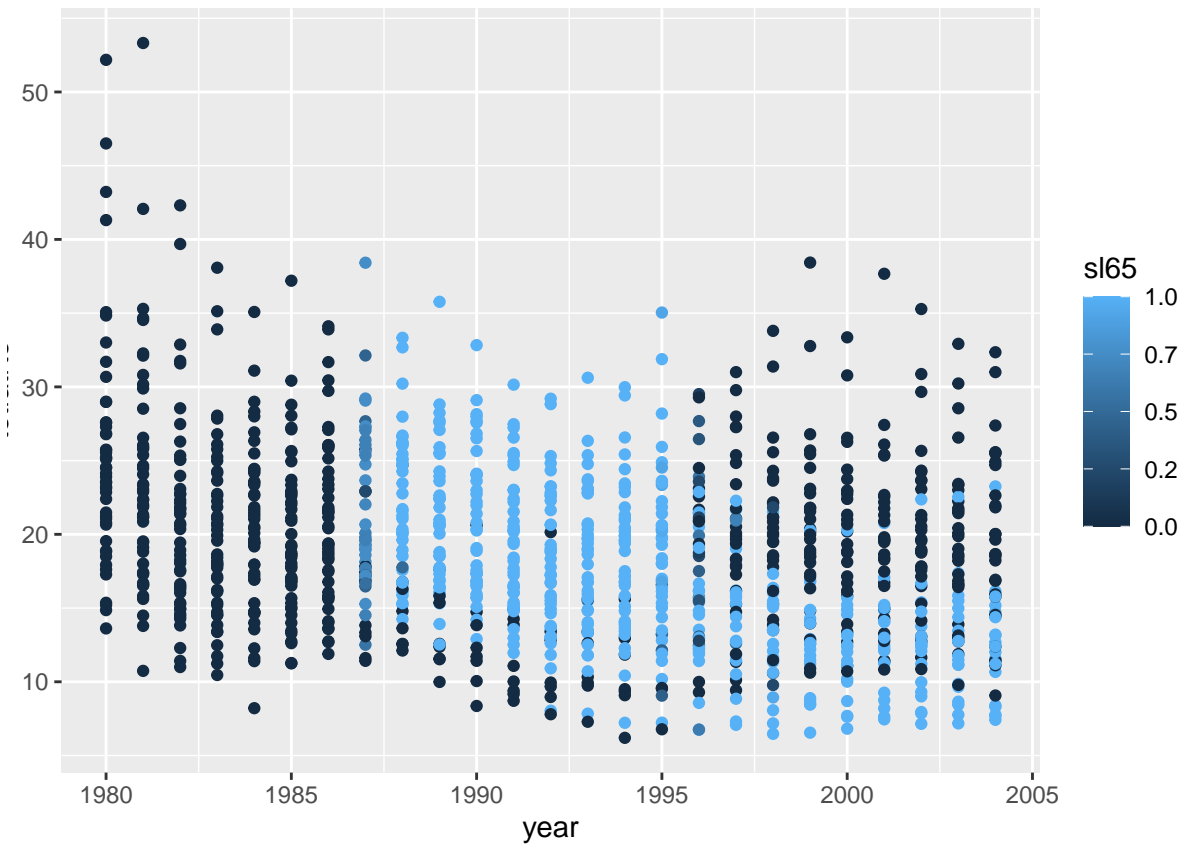
Early years, pre-1985 there seems to have been no seatbelt law. The fatality rate dropped over the years as 1, 2 seatbelt laws were introduced. States started with introduction of seatbelt rule 2 and then majority of them moved to seatbelt rule 1.

```
ggplot(driving, aes(year, totfatrte)) + geom_point(aes(color = zerotol))
```



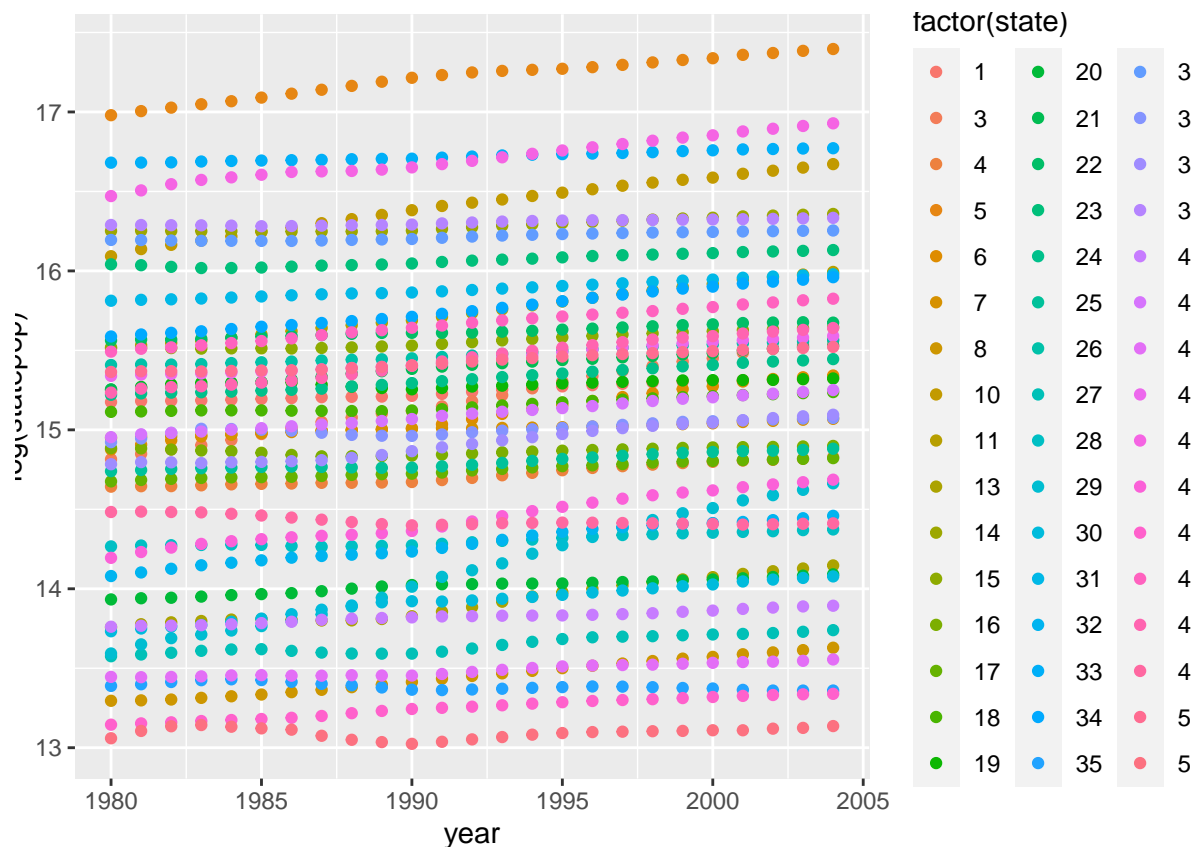
Most of the states moved to a zero tolerance policy over the years. In fact all the states seems to have introduced theh zero tolerance policy.

```
ggplot(driving, aes(year, totfatrte)) + geom_point(aes(color = sl65))
```

Most of the states seemed to have moved to speed limit 65 policy which correlates strongly with the drop in fatality rate. However several states seemed to have moved away to a possibly lower limit (55 mph) post 1996.

```
ggplot(driving, aes(year, log(statepop))) + geom_point(aes(color = factor(state)))
```



TODO: find the variables that are time invariant and the ones which are time varying. Perform specific EDA on those variables.

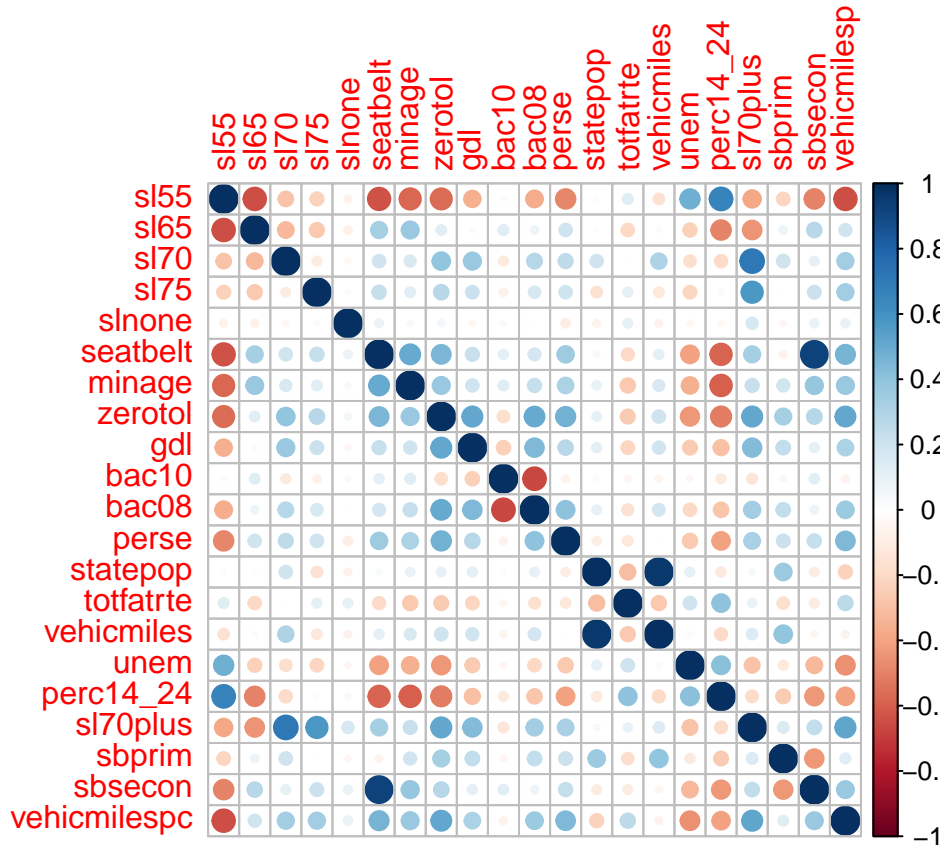
We could take a look at the panel data after inserting a structure into the dataset. The indices here are the “state” and the “year”.

We will have to select the appropriate explanatory variables to look at the correlation between those variables and the fatality rate. We would drop some of variables which are likely to have a higher correlation such as “fatality rate” with “weekend fatality rate” and “night fatality rate” etc. from the dataset used for the correlation matrix.

```
library(plm)
library(corrplot)

## corrplot 0.84 loaded

driving.panel <- pdata.frame(driving, c("state","year"), drop.index = TRUE)
M <- cor(driving.panel[,c(1:12,19:20,23:28,54)])
corrplot(M, method='circle')
```



We see that the “perc_14_24” which is the percent population between 14 and 24 has a very strong correlation with the “fatality rate”. The ‘minimum age’ has a negative correlation with the fatality rate, so is the ‘zero tolerance’ even though it is not very strong.

2. (15%) How is the our dependent variable of interest *totfatrte* defined? What is the average of this variable in each of the years in the time period covered in this dataset? Estimate a linear regression model of *totfatrte* on a set of dummy variables for the years 1981 through 2004. What does this model explain? Describe what you find in this model. Did driving become safer over this period? Please provide a detailed explanation.

Running OLS only on data for one panel only, the equation can be written as below

Fixed year effect

$$totfatrte_{it} = \beta_0 + \delta_{81}d81 + \delta_{82}d82 + \dots + \delta_{04}d04 + u_{it}$$

δ_t the change that is common to every city in year t . It estimates the common change in the fatality rate in year t relative to the base / reference year 1980. We assume that this change is common across all cities for a given year t .

```
driving.ols <- lm(totfatrte ~ d81 + d82 + d83 + d84 + d85 + d86 + d87 + d88 + d89 + d90 + d91 +
summary(driving.ols)
```

```
##
## Call:
```

```
## lm(formula = totfatrte ~ d81 + d82 + d83 + d84 + d85 + d86 +
##      d87 + d88 + d89 + d90 + d91 + d92 + d93 + d94 + d95 + d96 +
##      d97 + d98 + d99 + d00 + d01 + d02 + d03 + d04, data = driving)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.9302  -4.3468  -0.7305   3.7488  29.6498
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  25.4946     0.8671  29.401  < 2e-16 ***
## d81          -1.8244     1.2263  -1.488  0.137094
## d82          -4.5521     1.2263  -3.712  0.000215 ***
## d83          -5.3417     1.2263  -4.356  1.44e-05 ***
## d84          -5.2271     1.2263  -4.263  2.18e-05 ***
## d85          -5.6431     1.2263  -4.602  4.64e-06 ***
## d86          -4.6942     1.2263  -3.828  0.000136 ***
## d87          -4.7198     1.2263  -3.849  0.000125 ***
## d88          -4.6029     1.2263  -3.754  0.000183 ***
## d89          -5.7223     1.2263  -4.666  3.42e-06 ***
## d90          -5.9894     1.2263  -4.884  1.18e-06 ***
## d91          -7.3998     1.2263  -6.034  2.14e-09 ***
## d92          -8.3367     1.2263  -6.798  1.68e-11 ***
## d93          -8.3669     1.2263  -6.823  1.43e-11 ***
## d94          -8.3394     1.2263  -6.800  1.66e-11 ***
## d95          -7.8260     1.2263  -6.382  2.51e-10 ***
## d96          -8.1252     1.2263  -6.626  5.25e-11 ***
## d97          -7.8840     1.2263  -6.429  1.86e-10 ***
## d98          -8.2292     1.2263  -6.711  3.01e-11 ***
## d99          -8.2442     1.2263  -6.723  2.77e-11 ***
## d00          -8.6690     1.2263  -7.069  2.67e-12 ***
## d01          -8.7019     1.2263  -7.096  2.21e-12 ***
## d02          -8.4650     1.2263  -6.903  8.32e-12 ***
## d03          -8.7310     1.2263  -7.120  1.88e-12 ***
## d04          -8.7656     1.2263  -7.148  1.54e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.008 on 1175 degrees of freedom
## Multiple R-squared:  0.1276, Adjusted R-squared:  0.1098
## F-statistic: 7.164 on 24 and 1175 DF,  p-value: < 2.2e-16
```

From the above coefficients, we can say that for the year 2000, the fatality rate has dropped by 8.66 units as compared to the year 1980. Except for the year 1981, the coefficients for all the year dummies are significant. From this we can also say that the fatality rate has been dropping over the years. Driving did seem to become safer over the years.

3. (15%) Expand your model in *Exercise 2* by adding variables *bac08*, *bac10*, *perse*, *sbprim*, *sbsecon*, *sl70plus*, *gdl*, *perc14_24*, *unem*, *vehicmilespc*, and perhaps *transformations of some*

or all of these variables. Please 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 variables *bac8* and *bac10* defined? Interpret the coefficients on *bac8* and *bac10*. Do *per se laws* have a negative effect on the fatality rate? What about having a primary seat belt law? (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.)

4. (15%) Reestimate the model from *Exercise 3* using a fixed effects (at the state level) model. How do the coefficients on *bac08*, *bac10*, *perse*, and *sbprim* compare with the pooled OLS estimates? Which set of estimates do you think is more reliable? What assumptions are needed in each of these models? Are these assumptions reasonable in the current context?
5. (5%) Would you prefer to use a random effects model instead of the fixed effects model you built in *Exercise 4*? Please explain.
6. (5%) Suppose that *vehicmilespc*, the number of miles driven per capita, increases by 1,000. Using the FE estimates, what is the estimated effect on *totfatrte*? Please interpret the estimate.
7. (5%) If there is serial correlation or heteroskedasticity in the idiosyncratic errors of the model, what would be the consequences on the estimators and their standard errors?