

R Practice 3 (Answers)

Ryan Safner

11/13/2018

Download the `speeding_tickets.csv` dataset from Blackboard (under Data). This data comes from a paper by Makowsky and Strattman (2009) that we will examine later. Even though state law sets a formula for tickets based on how fast a person was driving, police officers in practice often deviate from that formula. This dataset includes information on all traffic stops. An amount for the fine is given only for observations in which the police officer decided to assess a fine.

- `Amount`: Amount of fine assessed for speeding
- `Age`: Age of speeder in years
- `MPHover`: Miles per hour over speed limit
- `Black`: = 1 if driver was black, = 0 if not
- `Hispanice`: = 1 if driver was Hispanic, = 0 if not
- `Female`: = 1 if driver was female, = 0 if not
- `OutTown`: = 1 if driver was not from local town, = 0 if not
- `OutState`: = 1 if driver was not from local state, = 0 if not
- `StatePol`: = 1 if driver was stopped by State Police, = 0 if stopped by other (local)

We want to explore who gets fines, and how much.

1. Load the data and inspect it briefly with `str()` and `head()`. We will have to do a little bit of cleaning to get the data in a more usable form.

```
library("ggplot2")

# load data
speed<-read.csv("../Data/speeding_tickets.csv")

str(speed)

## 'data.frame': 68357 obs. of  9 variables:
## $ Black    : int  0 0 0 0 0 0 0 0 ...
## $ Hispanic : int  0 0 0 0 0 0 0 0 ...
## $ Female   : int  1 1 1 0 0 0 1 0 1 0 ...
## $ Amount   : int  NA NA NA NA NA NA NA NA NA ...
## $ MPHover  : int  14 15 15 13 12 17 15 15 15 15 ...
## $ Age      : int  22 43 32 24 54 30 18 53 51 33 ...
## $ OutTown  : int  1 1 0 1 1 1 0 0 1 1 ...
## $ OutState : int  0 0 0 0 0 0 0 0 0 ...
## $ StatePol : int  0 0 0 0 0 0 0 0 0 ...

head(speed)

##   Black Hispanic Female Amount MPHover Age OutTown OutState StatePol
## 1     0        0     1     NA      14    22       1       0      0
## 2     0        0     1     NA      15    43       1       0      0
## 3     0        0     1     NA      15    32       0       0      0
## 4     0        0     0     NA      13    24       1       0      0
## 5     0        0     0     NA      12    54       1       0      0
## 6     0        0     0     NA      17    30       1       0      0
```

a. What class of variable are Black, Hispanic, Female, OutTown, and OutState?

```
class(speed$Black)
## [1] "integer"
class(speed$Hispanic)
## [1] "integer"
class(speed$Female)
## [1] "integer"
class(speed$OutState)
## [1] "integer"
class(speed$OutTown)
## [1] "integer"
```

b. Notice that when importing the data from the .csv file, R interpreted these variables as **integer**, but we want them to be **factor** variables, to ensure R recognizes that there are two groups (categories), 0 and 1. Convert each of these variables into factors by reassigning it according to the format: `df$var.name<-as.factor(df$var.name)`, where

- df is the name of your data frame
- var.name is the name of the variable

```
speed$Black<-as.factor(speed$Black)
speed$Hispanic<-as.factor(speed$Hispanic)
speed$Female<-as.factor(speed$Female)
speed$OutTown<-as.factor(speed$OutTown)
speed$OutState<-as.factor(speed$OutState)
```

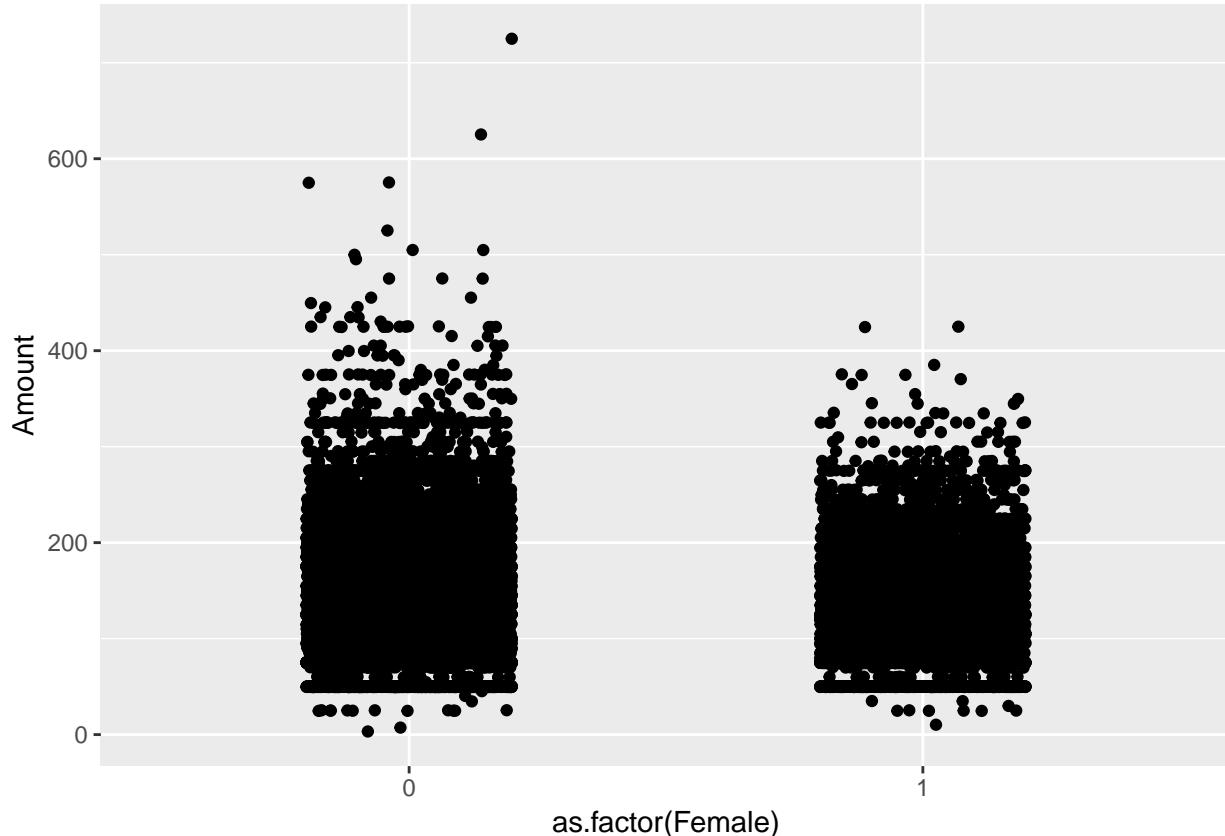
c. Confirm they are each now factors by checking their class again.

```
class(speed$Black)
## [1] "factor"
class(speed$Hispanic)
## [1] "factor"
class(speed$Female)
## [1] "factor"
class(speed$OutState)
## [1] "factor"
class(speed$OutTown)
## [1] "factor"
```

2. Create a scatterplot between Amount and Female. Use `geom_jitter()` instead of `geom_point()` to plot the points, and play around with width settings inside `geom_jitter()`.

```
ggplot(data = speed, aes(x = as.factor(Female), y = Amount)) +  
  geom_jitter(width=0.2)
```

Warning: Removed 36683 rows containing missing values (geom_point).



3. Check the distribution of Amount with `summary()`.

```
summary(speed$Amount)  
  
##      Min. 1st Qu. Median     Mean 3rd Qu.    Max.    NA's  
##        3       75     115     122     155     725   36683
```

a. If you notice, Amount has a lot of missing values (for people that did not get fined). Let's keep only data for which Amount is a positive number. Use the `subset()` command and overwrite your data (or make a new object) with `df1<-subset(df, condition)`

```
speed<-subset(speed, speed$Amount>0)
```

b. Double check this worked by checking the `summary()` of Amount again.

```
summary(speed$Amount)
```

```
##      Min. 1st Qu. Median     Mean 3rd Qu.    Max.
##        3       75     115     122     155     725
```

4. Find the mean and standard deviation of Amount, for male drivers and for female drivers.

```
# get mean for males
mean(speed$Amount[speed$Female==0])
```

```
## [1] 124.6654
```

```
# get sd for males
sd(speed$Amount[speed$Female==0])
```

```
## [1] 58.28387
```

```
# get mean for females
mean(speed$Amount[speed$Female==1])
```

```
## [1] 116.726
```

```
# get sd for females
sd(speed$Amount[speed$Female==1])
```

```
## [1] 51.48665
```

a. What is the difference between the average speed for Males and Females? (Calculate this manually)

```
mean(speed$Amount[speed$Female==0]) - mean(speed$Amount[speed$Female==1])
```

```
## [1] 7.939397
```

b. Use `t.test` to check if this is a statistically significant difference. The syntax is similar for regression: `t.test(y~d, data=df)` where `y` is the variable we are testing (`Amount`) and `d` is the dummy variable (`Female`)

```
t.test(Amount~Female, data=speed)
```

```
##
## Welch Two Sample t-test
##
## data: Amount by Female
## t = 12.356, df = 23400, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  6.679941 9.198853
## sample estimates:
## mean in group 0 mean in group 1
##           124.6654          116.7260
```

5. Now run the following regression to ensure we get the same result

$$\text{Amount}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{Female}_i$$

```

reg<-lm(Amount~Female, data=speed)
summary(reg)

##
## Call:
## lm(formula = Amount ~ Female, data = speed)
##
## Residuals:
##    Min     1Q Median     3Q    Max
## -121.67 -49.67 -6.73 30.33 600.33
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 124.6654    0.3857 323.23 <2e-16 ***
## Female1      -7.9394    0.6698 -11.85 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 56.12 on 31672 degrees of freedom
## Multiple R-squared:  0.004416,   Adjusted R-squared:  0.004385
## F-statistic: 140.5 on 1 and 31672 DF, p-value: < 2.2e-16

```

- a. Write out the estimated regression equation.

$$\widehat{\text{Amount}}_i = 124.67 - 7.94\text{Female}_i$$

- b. Use the regression coefficients to find (i) the average `Amount` for men, (ii) the average `Amount` for women, and (iii) the difference in average `Amount` between men and women

- Males get fined \$124.67 ($\hat{\beta}_0$)
- Females get fined $124.67 - 7.94 = 116.73$ ($\hat{\beta}_0 + \hat{\beta}_1$)
- The difference is $-\$7.94$ ($\hat{\beta}_3$)

6. Let's recode the sex variable. Make a new variable called `Male` and use the `ifelse()` function to define it as 1 when `df$Female==0` and 0 otherwise.

```

speed$Male<-ifelse(speed$Female==0,1,0)

```

- a. Run the same regression as in question 5, but use `Male` instead of `Female`.

```

regm<-lm(Amount~Male, data=speed)
summary(regm)

##
## Call:
## lm(formula = Amount ~ Male, data = speed)
##
## Residuals:
##    Min     1Q Median     3Q    Max
## -121.67 -49.67 -6.73 30.33 600.33
## 
```

```

## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 116.7260    0.5477 213.13 <2e-16 ***
## Male        7.9394    0.6698   11.85 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 56.12 on 31672 degrees of freedom
## Multiple R-squared:  0.004416, Adjusted R-squared:  0.004385
## F-statistic: 140.5 on 1 and 31672 DF, p-value: < 2.2e-16

```

b. Write out the estimated regression equation.

$$\widehat{\text{Amount}} = 116.73 + 7.94\text{Male}$$

c. Use the regression coefficients to find (i) the average Amount for men, (ii) the average Amount for women, and (iii) the difference in average Amount between men and women

- Females get fined \$116.73 ($\hat{\beta}_0$)
- Males get fined $116.73 + 7.94 = \$124.67$ over ($\hat{\beta}_0 + \hat{\beta}_1$)
- The difference is 7.94 ($\hat{\beta}_3$)

7. Run a regression of Amount on Male and Female. What happens, and why?

Male and Female are perfectly multicollinear, as for every person i , $\text{Male}_i + \text{Female}_i = 1$. We can confirm this by seeing the correlation between Male and Female is exactly -1. To run a regression, we must exclude one of the dummies, and as we've seen, it makes no difference which one we exclude.

8. Age probably has a lot to do with differences in fines, perhaps also age affects fines differences between males and females. Run a regression of Amount on Age and Female. How does the coefficient on Female change?

```

reg2<-lm(Amount~Female+Age, data=speed)
summary(reg2)

##
## Call:
## lm(formula = Amount ~ Female + Age, data = speed)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -125.77  -45.63   -5.91   33.67  597.66
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 134.19965    0.90969 147.52 <2e-16 ***
## Female1     -7.85172    0.66849 -11.74 <2e-16 ***
## Age         -0.28598    0.02472 -11.57 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##

```

```

## Residual standard error: 56 on 31671 degrees of freedom
## Multiple R-squared:  0.008604,   Adjusted R-squared:  0.008542
## F-statistic: 137.4 on 2 and 31671 DF,  p-value: < 2.2e-16

```

- a. Now let's see if the difference in fine between men and women are different depending on the driver's age. Run the regression again, but add an interaction term between Female and Age interaction term.

```

reg3<-lm(Amount~Female+Age+Female*Age, data=speed)
summary(reg3)

```

```

##
## Call:
## lm(formula = Amount ~ Female + Age + Female * Age, data = speed)
##
## Residuals:
##    Min      1Q  Median      3Q     Max
## -126.35  -44.68   -5.49   33.49  597.29
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 135.54592   1.07428 126.173 < 2e-16 ***
## Female1     -12.02331   1.89296 -6.352 2.16e-10 ***
## Age         -0.32636   0.03008 -10.848 < 2e-16 ***
## Female1:Age  0.12435   0.05279   2.355  0.0185 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 56 on 31670 degrees of freedom
## Multiple R-squared:  0.008778,   Adjusted R-squared:  0.008684
## F-statistic: 93.49 on 3 and 31670 DF,  p-value: < 2.2e-16

```

- b. Write out your estimated regression equation.

$$\widehat{\text{Amount}}_i = 135.55 - 12.02\text{Female}_i - 0.33\text{Age}_i + 0.12\text{Female}_i * \text{Age}_i$$

- c. Interpret the interaction effect. Is it statistically significant?

The coefficient on the interaction term, $\hat{\beta}_3$ is 0.12. For every additional year of age, females can expect their fine to increase by \$0.12 *more* than males gain for every additional year of age.

$\hat{\beta}_3$ has a standard error of 0.52, which means it has a *t*-statistic of 2.355 and *p*-value of 0.0185, so it is statistically significant.

- d. Plugging in 0 or 1 as necessary, rewrite (on your paper) this regression as *two separate* equations, one for Males and one for Females.

For Males (Female=0):

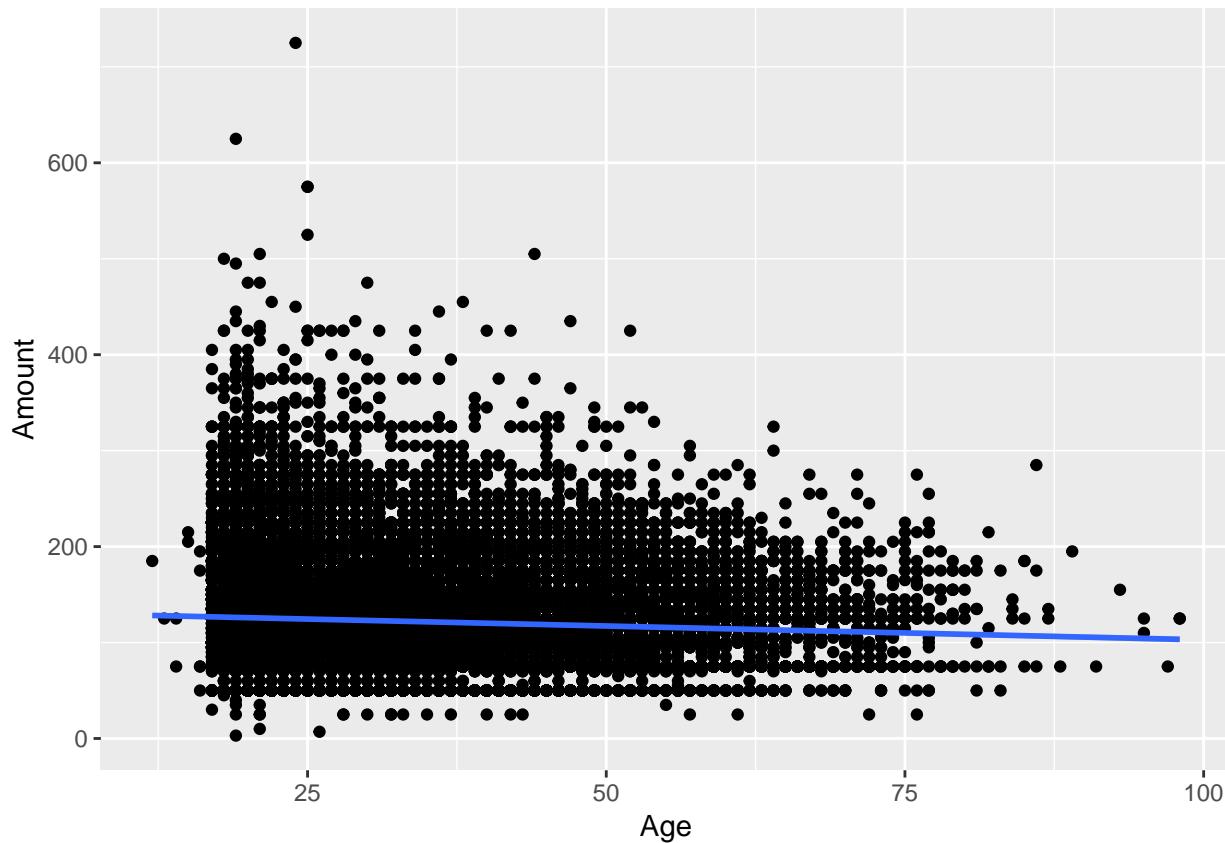
$$\begin{aligned}
\widehat{\text{Amount}} &= 135.55 - 0.33\text{Age} - 12.02\text{Female} + 0.12\text{Female} \times \text{Age} \\
&= 135.55 - 0.33\text{Age} - 12.02(0) + 0.12(0)\text{Age} \\
&= 135.55 - 0.33\text{Age}
\end{aligned}$$

For Females (Female=1):

$$\begin{aligned}
\widehat{\text{Amount}} &= 135.55 - 0.33\text{Age} - 12.02\text{Female} + 0.12\text{Female} \times \text{Age} \\
&= 135.55 - 0.33\text{Age} - 12.02(1) + 0.12(1)\text{Age} \\
&= (135.55 - 12.02) + (-0.33 + 0.12)\text{Age} \\
&= 123.53 - 0.21\text{Age}
\end{aligned}$$

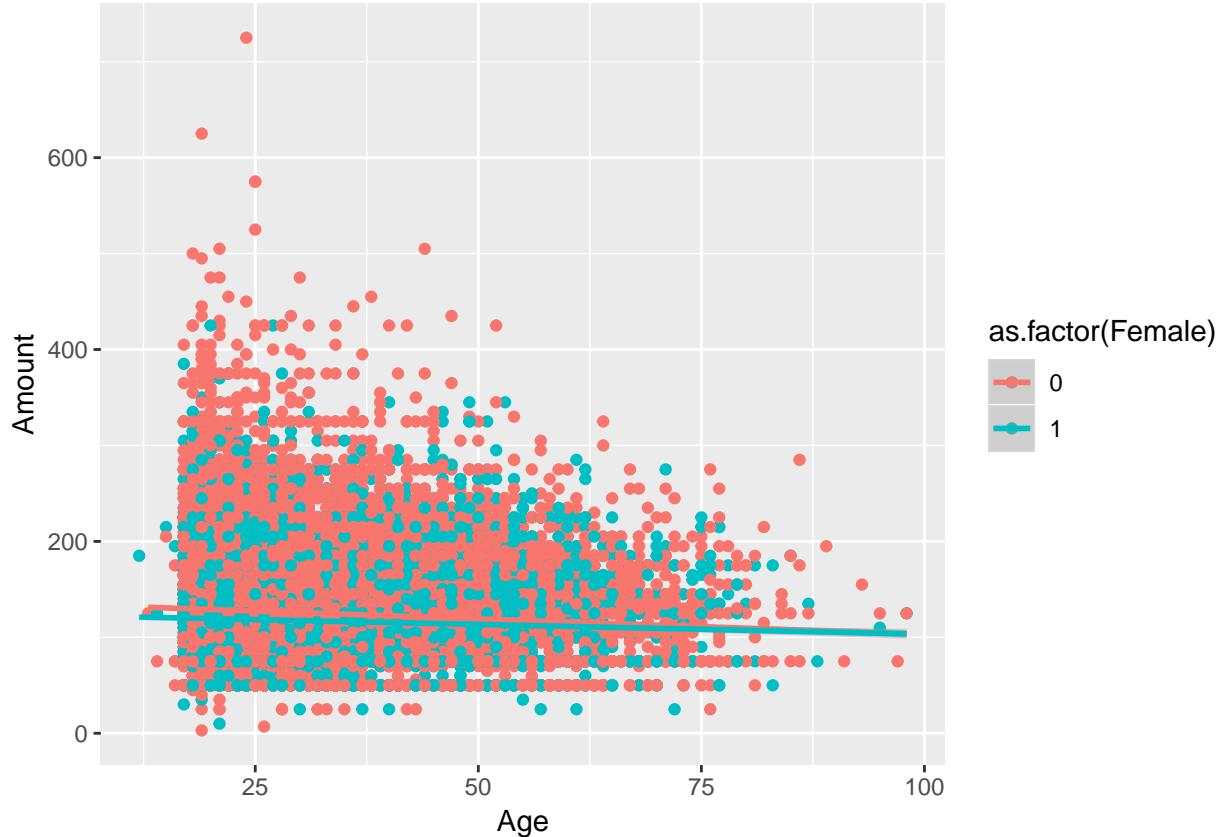
f. Let's try to visualize this. Make a scatterplot of `Age` (*X*) and `Amount` (*Y*) and include a regression line.

```
p1<-ggplot(data = speed, aes(x = Age, y = Amount))+
  geom_point()+
  geom_smooth(method="lm")
p1
```



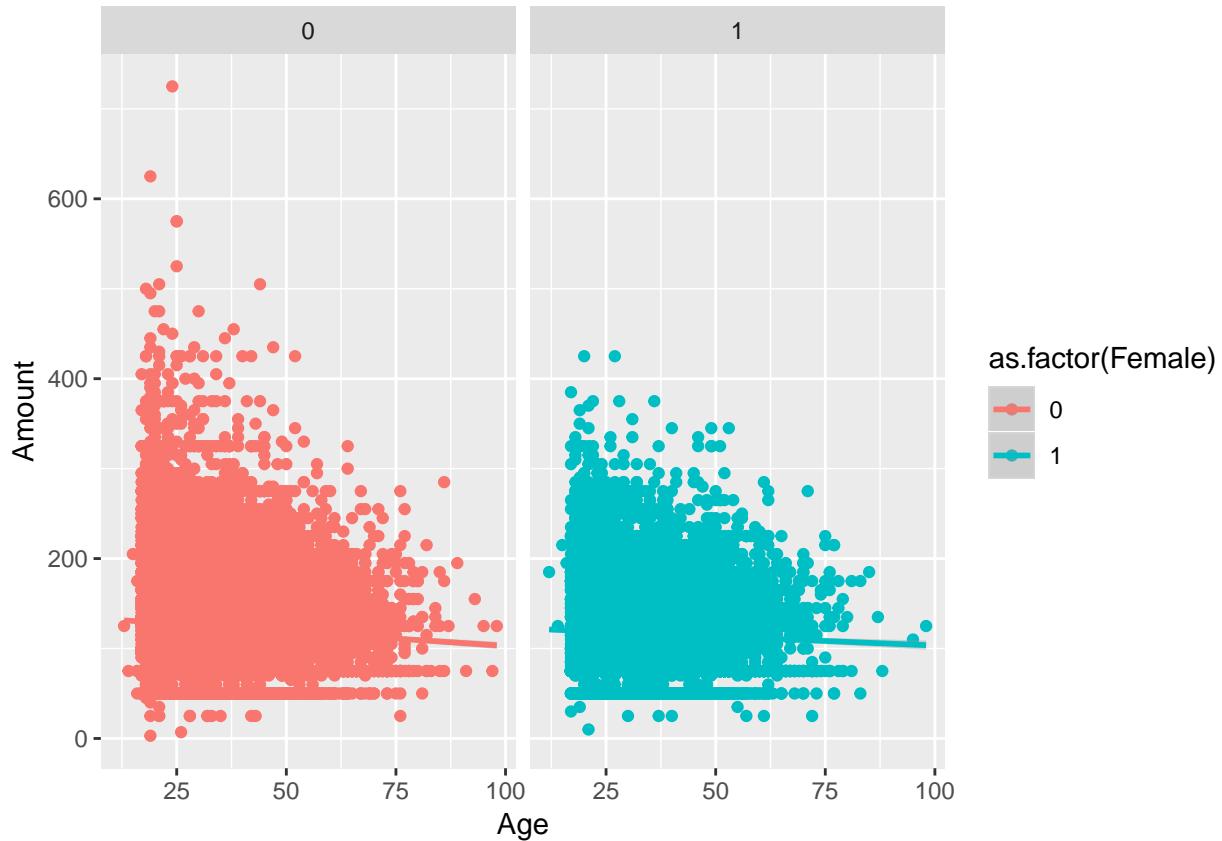
g. Try adding to your base layer `aes()`, set `color=Female`. This will produce two lines and color the points by Female. Sometimes we may also need to remind R that Female is a factor with `as.factor(Female)`.

```
p2<-ggplot(data = speed, aes(x = Age, y = Amount, color=as.factor(Female)))+  
  geom_point() +  
  geom_smooth(method="lm")  
p2
```



h. Add a facet layer to make two different scatterplots with an additional layer
+facet_grid(cols=vars(Female))

```
p2+facet_grid(cols=vars(Female))
```



9. Now let's look at the possible interaction between Sex (Male or Female) and whether a driver is from In-State or Out-of-State (OutState).

a. Use R to examine the data and find the mean for (i) Males In-State, (ii) Males Out-of-State, (iii) Females In-State, and (iv) Females Out-of-State. Hint: use & to join multiple conditions!

```
# get mean for In-State Males
mean(speed$Amount[speed$Female==0 & speed$OutState==0])
## [1] 123.6775

# get mean for Out-State Males
mean(speed$Amount[speed$Female==0 & speed$OutState==1])
## [1] 127.969

# get mean for In-State Females
mean(speed$Amount[speed$Female==1 & speed$OutState==0])
## [1] 114.7985

# get mean for Out-State Females
mean(speed$Amount[speed$Female==1 & speed$OutState==1])
## [1] 124.2657
```

b. Now run a regression of the following model:

```

Amounti =  $\hat{\beta}_0 + \hat{\beta}_1 Female_i + \hat{\beta}_2 OutState_i + \hat{\beta}_3 Female_i * OutState_i$ 

reg4<-lm(Amount~Female+Female*OutState, data=speed)
summary(reg4)

##
## Call:
## lm(formula = Amount ~ Female + Female * OutState, data = speed)
##
## Residuals:
##    Min     1Q Median     3Q    Max
## -120.68 -48.68 -8.68  31.32 601.32
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 123.6775   0.4391 281.650 < 2e-16 ***
## Female1      -8.8790   0.7541 -11.775 < 2e-16 ***
## OutState1       4.2915   0.9152   4.689 2.76e-06 ***
## Female1:OutState1  5.1757   1.6381   3.160  0.00158 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 56.06 on 31670 degrees of freedom
## Multiple R-squared:  0.006629, Adjusted R-squared:  0.006535
## F-statistic: 70.44 on 3 and 31670 DF, p-value: < 2.2e-16

```

c. Write out the estimated regression equation.

$$\widehat{Amount} = 123.68 - 8.88Female + 4.29OutState + 5.17Female \times OutState$$

d. What does each coefficient mean?

- $\hat{\beta}_0 = \$123.68$; mean for in-state males
- $\hat{\beta}_1 = -\$8.88$: difference between in-state males and females
- $\hat{\beta}_2 = \$4.29$: difference between males in-state vs. out-of-state
- $\hat{\beta}_3 = \$5.17$: difference between effect of being in-state vs. out-of-state between males vs. females (or, equivalently, difference between effect of being male vs. female between in-state vs. out-of-state)

e. Using the regression equation, what are the means for (i) Males In-State, (ii) Males Out-of-State, (iii) Females In-State, and (iv) Females Out-of-State? Compare to your answers in part a.

$$\widehat{Amount} = 123.68 - 8.88Female + 4.29OutState + 5.17Female * OutState$$

- Males In-State: $\hat{\beta}_0 = \$123.68$
- Males Out-of-State: $\hat{\beta}_0 + \hat{\beta}_2 = 123.68 + 4.29 = \127.97
- Females In-State: $\hat{\beta}_0 + \hat{\beta}_1 = 123.68 - 8.88 = \114.80
- Females Out-of-State: $\hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_2 + \hat{\beta}_3 = 123.68 - 8.88 + 4.29 + 5.17 = \124.26

10. Collect your regressions from questions 5, 6a, 8, 8a, and 9b and output them in a regression table with stargazer.

```
library("stargazer")

##
## Please cite as:

## Hlavac, Marek (2018). stargazer: Well-Formatted Regression and Summary Statistics Tables.
## R package version 5.2.2. https://CRAN.R-project.org/package=stargazer
stargazer(reg, regm, reg2, reg3, reg4, header=FALSE, float=FALSE, type="latex", font.size="tiny")
```

| | Dependent variable: Amount | | | | |
|-------------------------|-------------------------------|----------------------------|----------------------------|---------------------------|---------------------------|
| | (1) | (2) | (3) | (4) | (5) |
| Female1 | -7.939*** (0.670) | | -7.852*** (0.668) | -12.023*** (1.893) | -8.879*** (0.754) |
| Male | | 7.939*** (0.670) | | | |
| Age | | | -0.286*** (0.025) | -0.326*** (0.030) | |
| Female1:Age | | | | 0.124** (0.053) | |
| OutState1 | | | | | 4.291*** (0.915) |
| Female1:OutState1 | | | | | 5.176*** (1.638) |
| Constant | 124.665*** (0.386) | 116.726*** (0.548) | 134.200*** (0.910) | 135.546*** (1.074) | 123.678*** (0.439) |
| Observations | 31,674 | 31,674 | 31,674 | 31,674 | 31,674 |
| R ² | 0.004 | 0.004 | 0.009 | 0.009 | 0.007 |
| Adjusted R ² | 0.004 | 0.004 | 0.009 | 0.009 | 0.007 |
| Residual Std. Error | 56.122 (df = 31672) | 56.122 (df = 31672) | 56.004 (df = 31671) | 56.000 (df = 31670) | 56.061 (df = 31670) |
| F Statistic | 140.483*** (df = 1; 31672) | 140.483*** (df = 1; 31672) | 137.434*** (df = 2; 31671) | 93.485*** (df = 3; 31670) | 70.444*** (df = 3; 31670) |

Note:

* p<0.1; ** p<0.05; *** p<0.01