

Regression Model Showcase

1. Installations

1.1 install the libraries

```
library(ggplot2)
library(car)

# Loading required package: carData
```

2.1 load the dataset

```
EV = read.csv("Data/Amilias/Desktop/2022 Fall/DSO SIS/PROJECT/Final/All_Data_Final-Overview.csv") # File Name:
All_Data_Final - Overview.csv
```

2. Data Processing

2.1 change the column name

```
colnames(EV)[1] = 'EV_Registration_Num'
colnames(EV)[2] = 'Gas_Price'
colnames(EV)[4] = 'Population'
colnames(EV)[5] = 'Household_Income'
colnames(EV)[6] = 'Charging_Stations'
colnames(EV)[7] = 'Political_Party'
```

2.2 change the datatype

```
EV$EV_Registration_Num = as.numeric(gsub(',', '', EV$EV_Registration_Num))
EV$Household_Income = gsub(',', '', EV$Household_Income)
EV$Household_Income = as.numeric(gsub(',', '', EV$Household_Income))
EV$Charging_Stations = as.integer(gsub(',', '', EV$Charging_Stations))
EV$Political_Party = as.factor(EV$Political_Party)
EV$Population = as.numeric(gsub(',', '', EV$Population))
```

3. Visualization and Descriptive Analysis

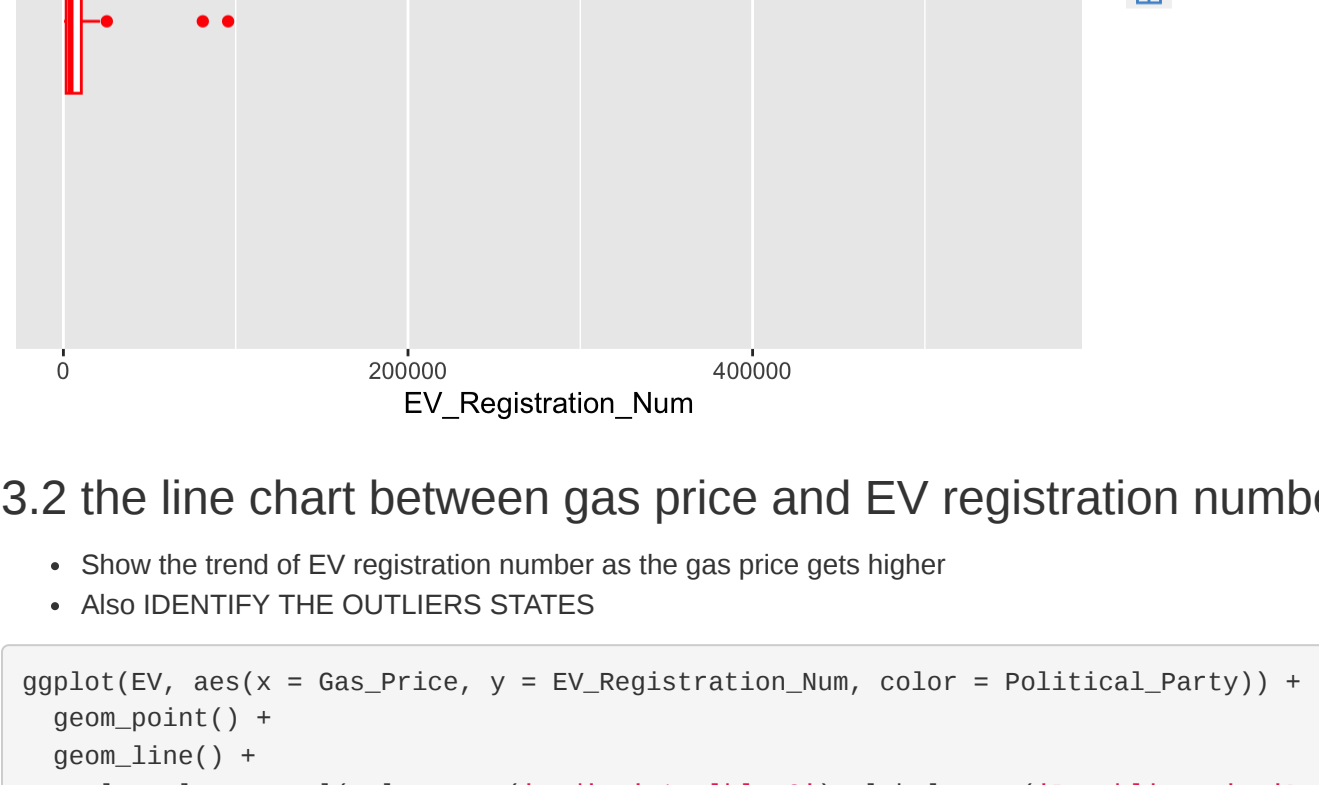
• Original Graph

```
options(scipen = 999)
```

3.1 the box plot of the EV registration number by Democratic/Republican states

```
ggplot(EV, aes(EV_Registration_Num, color = Political_Party)) +
  geom_boxplot() +
  scale_y_discrete(breaks = "NULL") +
  scale_color_manual(values = c('red', 'steelblue3'), labels = c('Republicans', 'Democrats')) +
  labs(title = 'Boxplot of EV Registration Number by Democratic/Republican States')
```

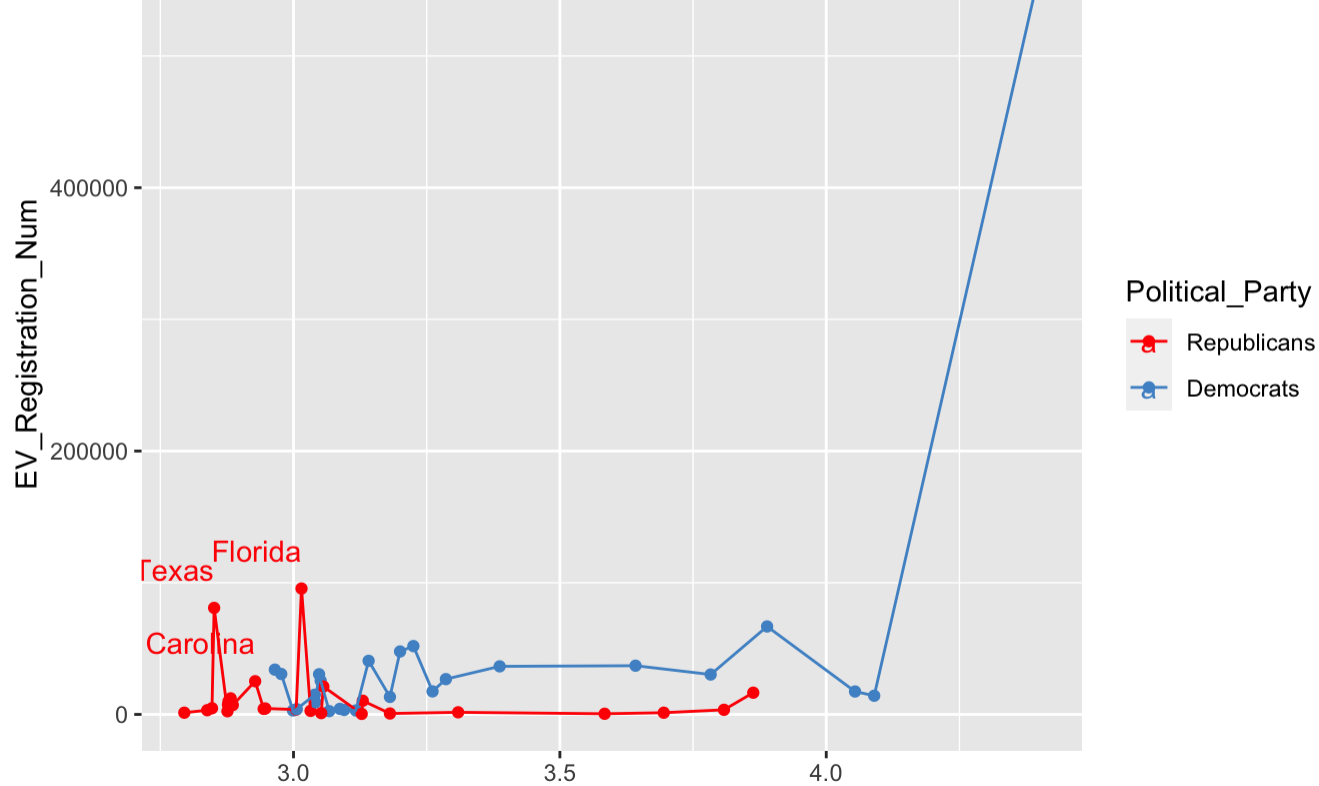
Boxplot of EV Registration Number by Democratic/Republican States



3.2 the line chart between gas price and EV registration number

• Show the trend of EV registration number as the gas price gets higher
• Also IDENTIFY THE OUTLIERS STATES

```
ggplot(EV, aes(x = Gas_Price, y = EV_Registration_Num, color = Political_Party)) +
  geom_point() +
  geom_smooth() +
  scale_color_manual(values = c('red', 'steelblue3'), labels = c('Republicans', 'Democrats')) +
  facet_warp(Political_Party) +
  labs(title = 'Line Chart between Gas Price and EV Registration Num')
```



3.3 exclude the outliers

```
EV = EV[EV$States != 'California',] # Exclude the Democratic States Outlier
EVNEW = EV[(EV$States %in% c('Florida', 'Texas', 'North Carolina'))], # Exclude the Republicans States Outliers
```

3.4 draw the linear trend line for both political states after excluding outliers

• Separate

```
ggplot(EVNEW, aes(x = Gas_Price, y = EV_Registration_Num, color = Political_Party)) +
  geom_point() +
  geom_smooth(method = 'lm', fill = NA) +
  scale_color_manual(values = c('red', 'steelblue3'), labels = c('Republicans', 'Democrats')) +
  facet_warp(Political_Party) +
  labs(title = 'Trendline after Excluding the Outliers')
```

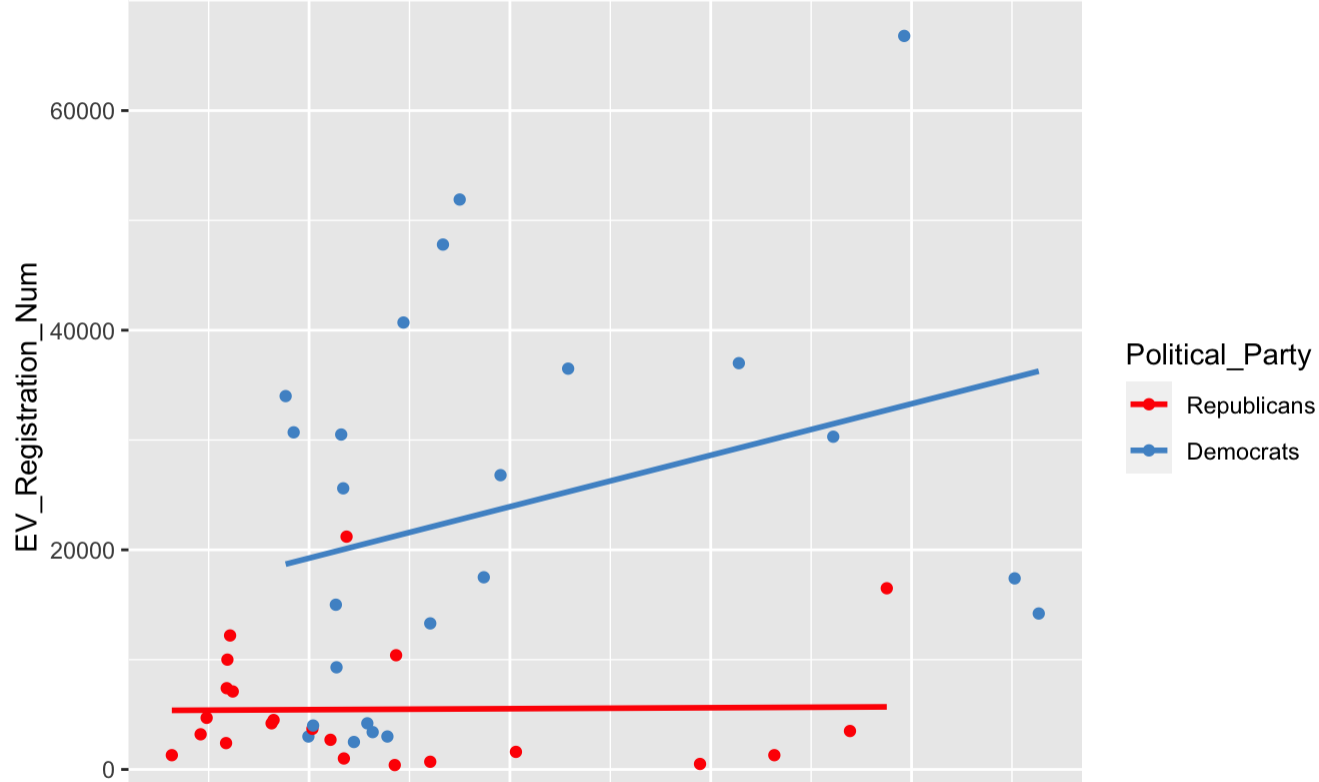
geom_smooth() using formula = 'y ~ x'



• Combined

```
ggplot(EVNEW, aes(x = Gas_Price, y = EV_Registration_Num, color = Political_Party)) +
  geom_point() +
  geom_smooth(method = 'lm', fill = NA) +
  scale_color_manual(values = c('red', 'steelblue3'), labels = c('Republicans', 'Democrats')) +
  labs(title = 'Trendline after Excluding the Outliers')
```

geom_smooth() using formula = 'y ~ x'



3.5 correlations between multiple potential numeric factors

```
cor(EVNEW[, c(2,3,4,5,6)])
```

	EV_Registration_Num	Gas_Price	Population
EV_Registration_Num	1.0000000	0.3894536	0.7691243
Gas_Price	0.3894536	1.0000000	-0.0882427
Population	0.7691243	-0.0882427	1.0000000
Household_Income	0.473384	-0.0259553	0.0753320
Charging_Stations	0.8448172	0.1392778	0.7957875
EV_Registration_Num	0.4473384	0.1392778	0.7957875
Gas_Price	0.425955	1.0000000	0.1392778
Population	0.0753320	0.0753320	1.0000000
Household_Income	0.0882427	0.3784218	0.0882427
Charging_Stations	0.7957875	0.1392778	1.0000000

• Conclusion:
= 1. Charging Station and Population are strongly and positively correlated with EV Registration No (v-r=0)
= 2. Charging Station and Population are correlated (v-r=0.4) - **need to check the VIF No. in the regression model

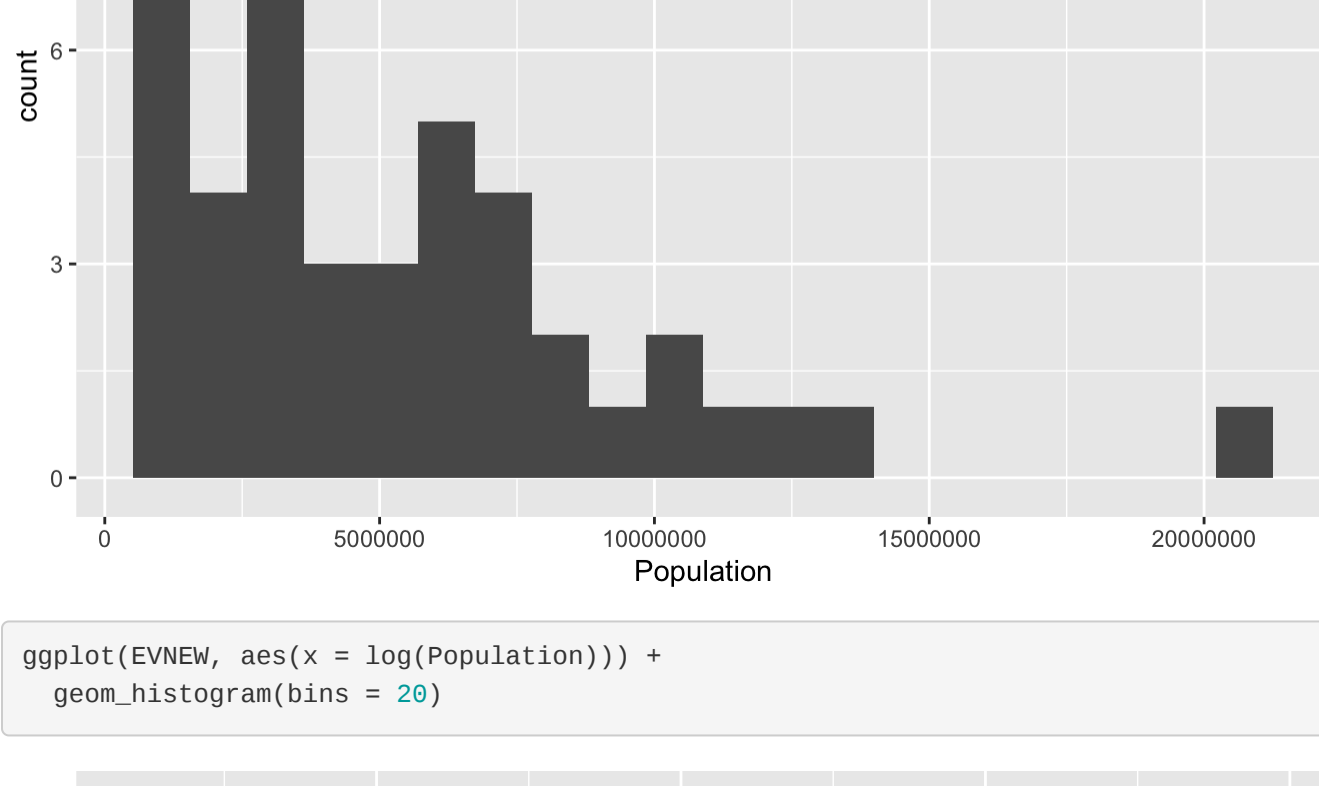
4. Regression Model

4.1 LOG functions:

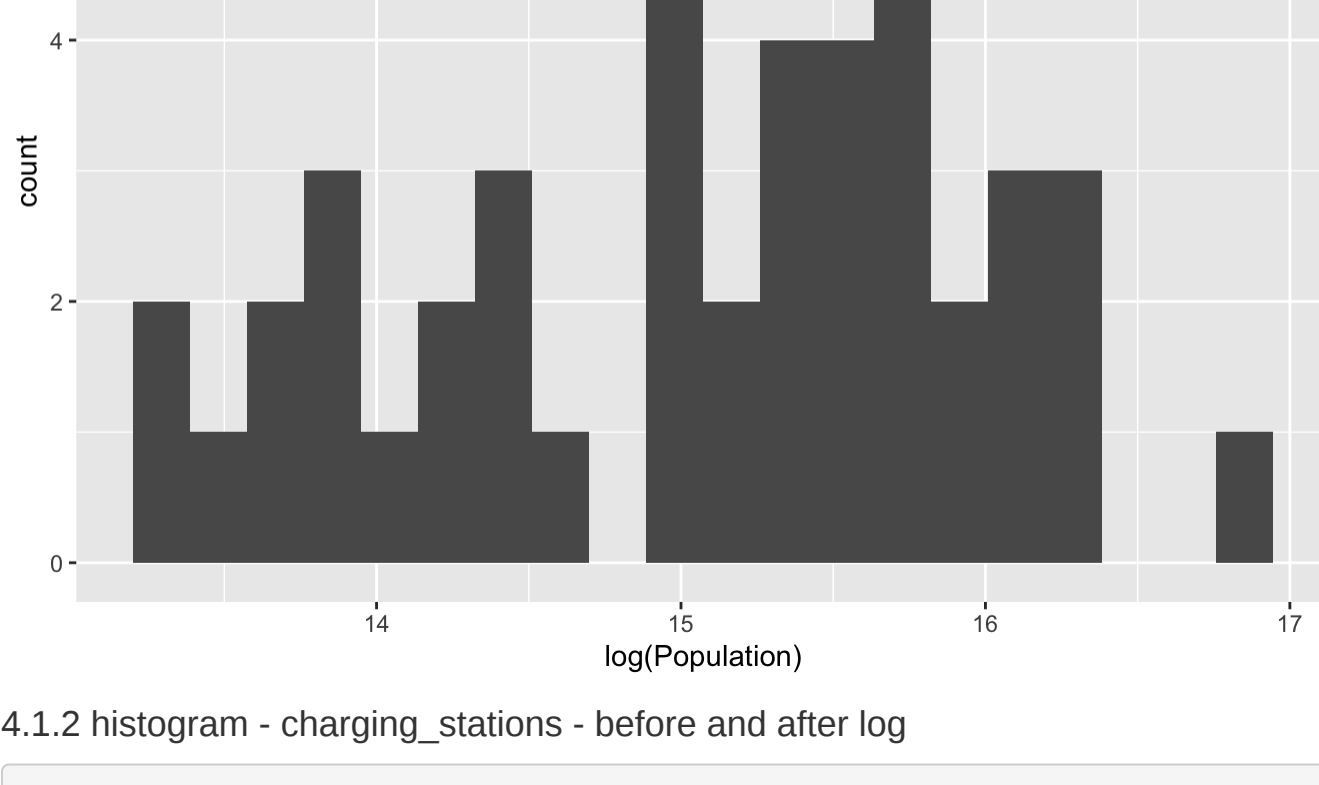
Because we check the histograms for each numeric independent variables and find out all variables except Household_Income are in skewed distribution. For better modeling, we will take the log for all numeric variables except Household_Income.

4.1.1 histogram - population - before and after log

```
ggplot(EVNEW, aes(x = Population)) +
  geom_histogram(bins = 20)
```

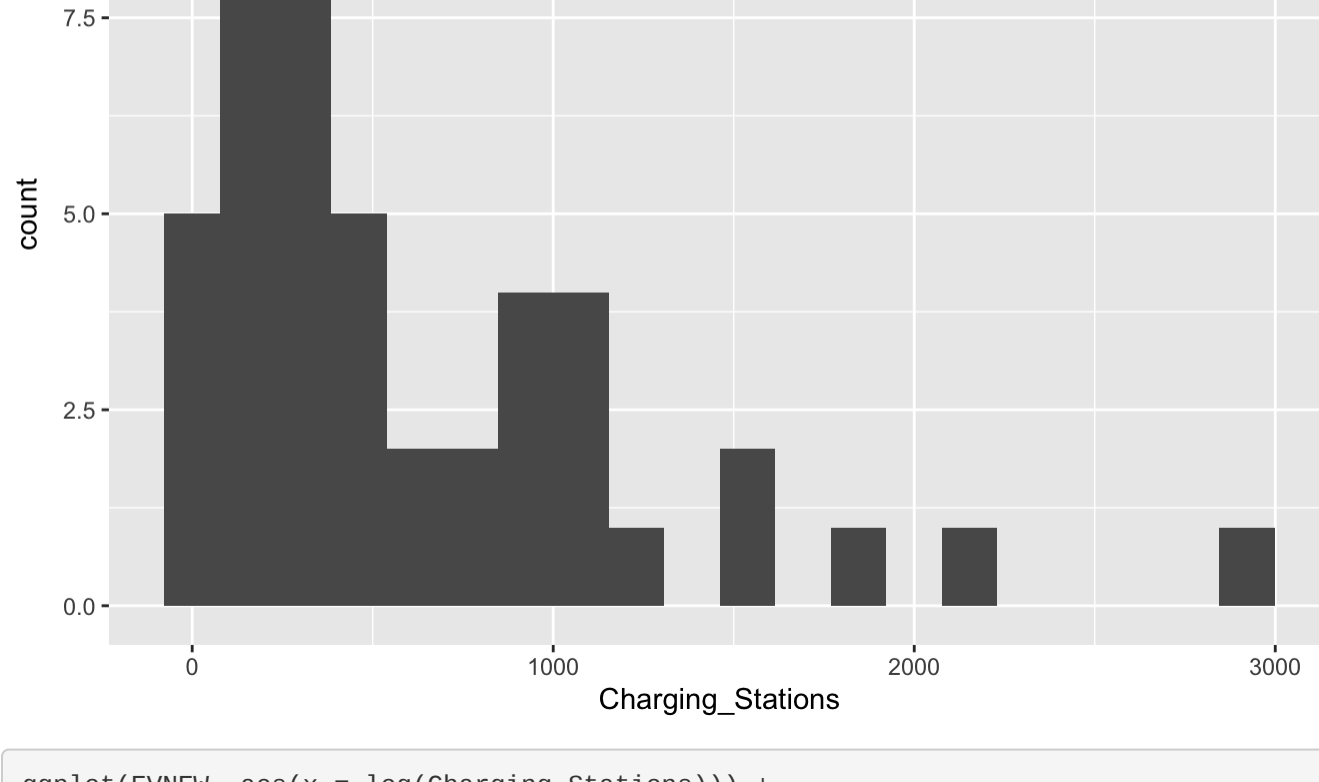


```
ggplot(EVNEW, aes(x = log(Population))) +
  geom_histogram(bins = 20)
```

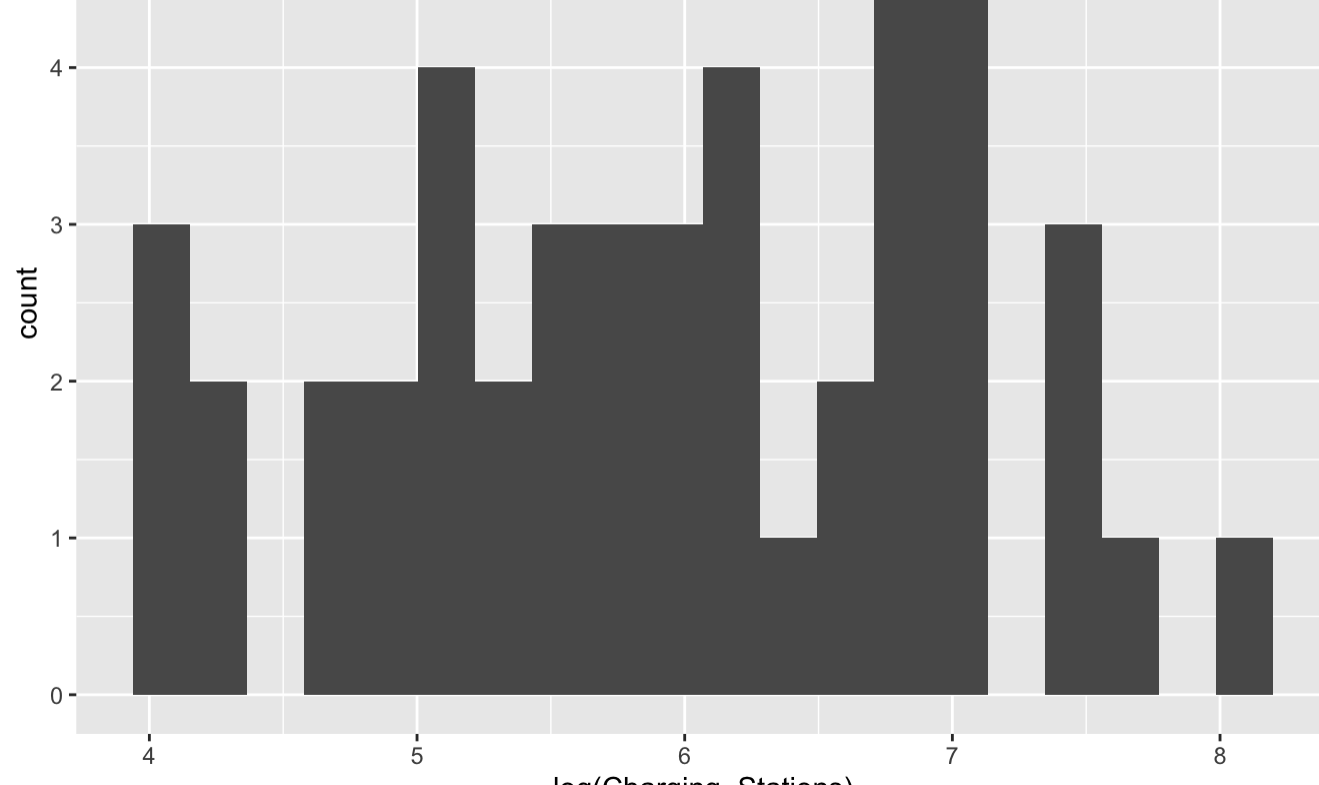


4.1.2 histogram - charging_stations - before and after log

```
ggplot(EVNEW, aes(x = Charging_Stations)) +
  geom_histogram(bins = 20)
```

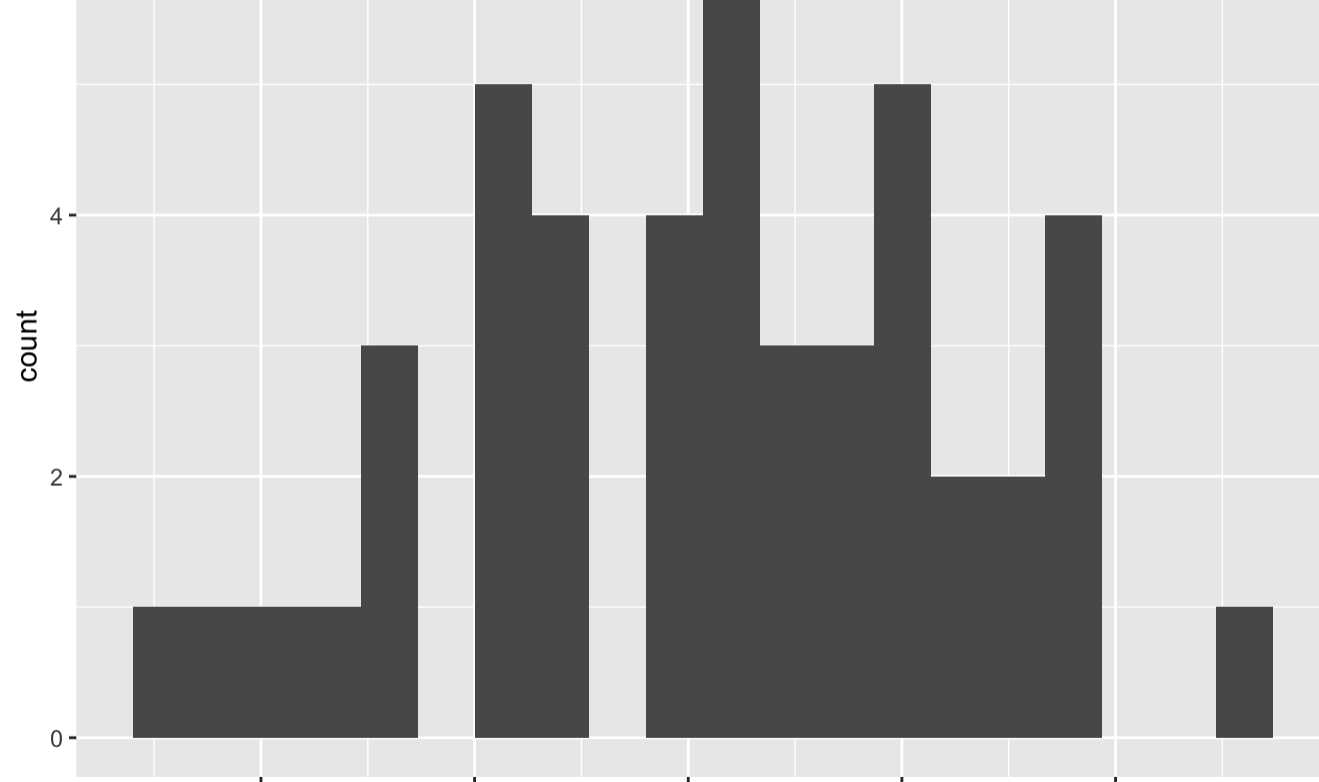


```
ggplot(EVNEW, aes(x = log(Charging_Stations))) +
  geom_histogram(bins = 20)
```



4.1.3 histogram - household_income - DONT need to log because its histogram is in normal distribution

```
ggplot(EVNEW, aes(x = Household_Income)) +
  geom_histogram(bins = 20)
```



5. Regression Model - Simple and Multiple

5.1 regression model 1 - single variable + dummy variable

• LOG(EV_Registration_Num) ~ 10 + b1Democrat + b2LOG(Gas Price) + b3Democrat*LOG(Gas Price)
= Democrat States: (b0 + b1) + (b2 + b3)LOG(Gas Price) - Hypothesis: b3=0
= Republican States: b0 + b2*LOG(Gas Price)

```
l_dummy_n_1 = lm(log(EV_Registration_Num)~Political_Party*log(Gas_Price), EVNEW)
summary(l_dummy_n_1)
```

```
# Call:
lm(formula = log(EV_Registration_Num) ~ Political_Party * log(Gas_Price),
    data = EVNEW)

# Residuals:
#      Min       1Q   median       3Q      Max
# -2.8853 -0.8759  0.1278  0.7782  2.6408

# Coefficients:
# (Intercept)              Estimate Std. Error t value Pr(>|t|)
# Political_Party1         -4.951      3.717   -1.325  0.193599
# log(Gas_Price)           -1.931      2.320   -0.832  0.418098
# Political_Party1:log(Gas_Price)  5.571      3.203   1.739  0.089373

# <--
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# Residual standard error: 1.858 on 42 degrees of freedom
# Multiple R-squared:  0.488, Adjusted R-squared:  0.3957
# F-statistic: 9.647 on 3 and 42 DF, p-value: 0.0000749
```

• Conclusion:
= 1. D: P increase, EV increase, R-P increase, EV decrease - Different than our expectation
= 2. R-squared: 0.48, IV P Value: not low enough to be statistically significant

5.2 regression model 2 - multiple variables + dummy variable

• LOG(EV_Registration_Num) ~ 10 + b1Democrat + b2LOG(Gas Price) + b3LOG(Population) + b4Income + b5LOG(Charging Stations) + b6Democrat*LOG(Gas Price)
= Democrat States: (b0 + b1) + (b2 + b3)LOG(Gas Price) + (b3)LOG(Population) + b4Income + b5LOG(Charging Stations) - controlling variables
= Republican States: b0 + b2*LOG(Gas Price) + (b3)LOG(Population) + b4Income + b5LOG(Charging Stations) - controlling variables

```
l_dummy_n_2 = lm(log(EV_Registration_Num)~Political_Party*log(Gas_Price)+log(Population)+Household_Income*log(Charging_Stations), EVNEW)
summary(l_dummy_n_2)
```

```
# Call:
lm(formula = log(EV_Registration_Num) ~ Political_Party * log(Gas_Price) +
    log(Population) + Household_Income * log(Charging_Stations),
    data = EVNEW)

# Residuals:
#      Min       1Q   median       3Q      Max
# -0.6844 -0.2378  0.0489  0.2157  0.6295

# Coefficients:
# (Intercept)              Estimate Std. Error t value Pr(>|t|)
# Political_Party1         -7.0290366   1.7803384   -3.950  0.000197 ***
# log(Gas_Price)           -0.8497819   1.2463389   -0.682  0.498341
# log(Population)           1.4157633   0.8597823   1.647  0.107632
# Household_Income          0.0703887   0.1146024    0.613  0.5400977 ***
# log(Charging_Stations)   0.08061261  0.0000055    2.155  0.037355 *
# Political_Party1:log(Gas_Price) 1.1218870   1.8756286    0.598  0.5509897 ***
# <--
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# Residual standard error: 0.3154 on 38 degrees of freedom
# Multiple R-squared:  0.9511, Adjusted R-squared:  0.9436
# F-statistic: 126.5 on 6 and 38 DF, p-value: < 0.0000000000000022
```

• Conclusion:
= 1. both D and R: P increase, EV increase
= 2. R-squared: 0.95, b6 P Value: not low enough to be statistically significant - accept the null hypothesis
= 3. Population, Charging Station, and Income are statistically significant. But the business implication for income is not significant
b5: increase in income, 0.0022% increase in EV

5.3 check the multicollinearity between IV

• Conclusion: $VIF \leq 10$, they are not highly correlated and should not be excluded from the model

```
vif(l_dummy_n_3)
```

```
# there are higher-order terms (interactions) in this model
# consider setting type = 'predictor'; see 'vif'
```

	Political_Party	log(Gas_Price)
Political_Party1	179.225239	3.453988
log(Population)	1.4157633	1.647
Household_Income	4.962884	2.155594
log(Charging_Stations)	0.0806126	1.850988

5.4 the relationship between each IV and DV, controlling other IVs

```
addPlot(l_dummy_n_3)
```

