# Politics Are Afoot!

## Da Qi Ren

## The Setup

There is *a lot* of money that is spent in politics in Presidential election years. So far, estimates have the number at about $11,000,000,000 (11 billion USD). For context, in 2019 Twitter's annual revenue was about $3,500,000,000 (3.5 billion USD).

## The work

Install the package, `fec16`.

```
## install.packages('fec16')
```

This package is a compendium of spending and results from the 2016 election cycle. In this dataset are 9 different datasets that cover:

- `candidates`: candidate attributes, like their name, a unique id of the candidate, the election year under consideration, the office they're running for, etc.
- `results_house`: race attributes, like the name of the candidates running in the election, a unique id of the candidate, the number of `general_votes` garnered by each candidate, and other information.
- `campaigns`: financial information for each house & senate campaign. This includes a unique candidate id, the total receipts (how much came in the doors), and total disbursements (the total spent by the campaign), the total contributed by party central committees, and other information.

## Your task

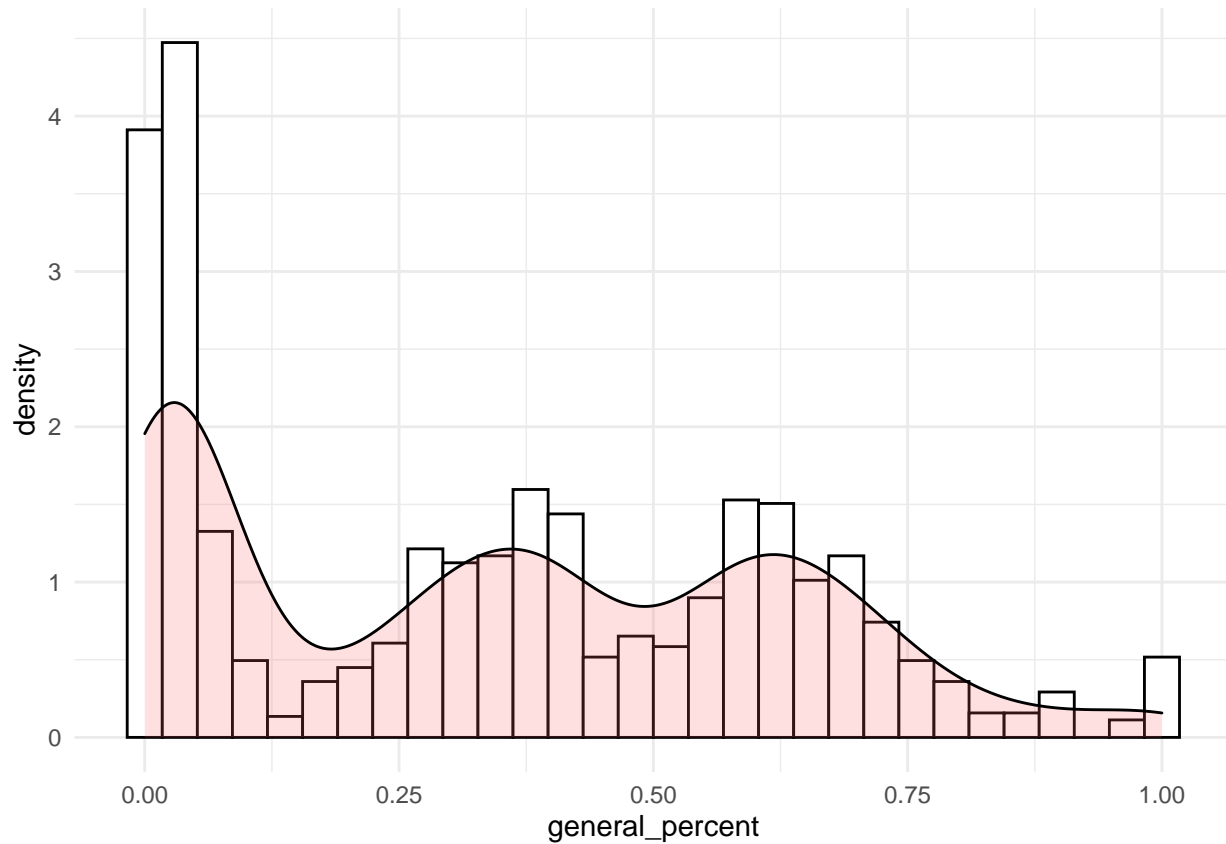Describe the relationship between spending on a candidate's behalf and the votes they receive.
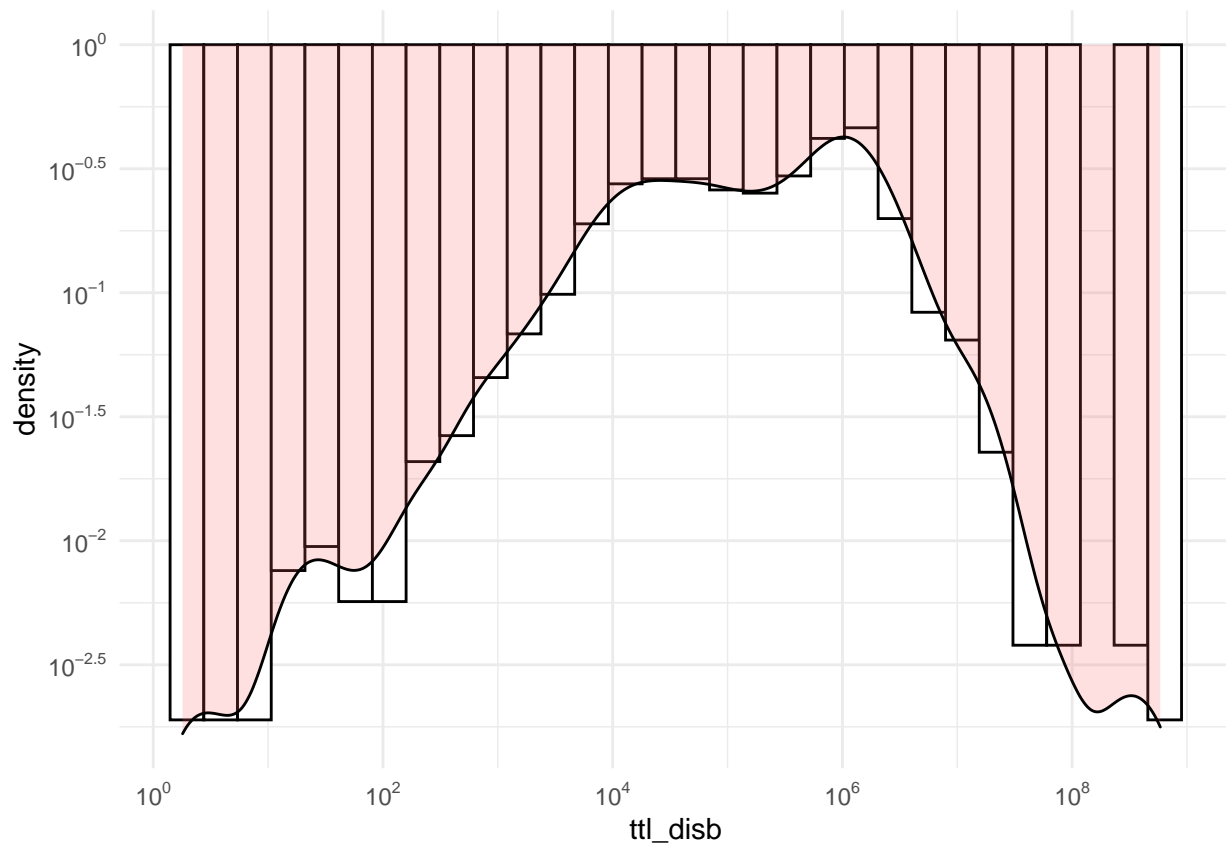
## Your work

- We want to keep this work *relatively* constrained, which is why we're providing you with data through the `fec16` package. It is possible to gather all the information from current FEC reports, but it would require you to make a series of API calls that would pull us away from the core modeling tasks that we want you to focus on instead.
- Throughout this assignment, limit yourself to functions that are within the `tidyverse` family of packages: `dplyr`, `ggplot`, `patchwork`, and `magrittr` for wrangling and exploration and `base`, `stats`, `sandwich` and `lmtest` for modeling and testing. You do not *have* to use these packages; but try to limit yourself to using only these.

```
candidates     <- fec16::candidates
results_house  <- fec16::results_house
campaigns      <- fec16::campaigns
```

# 1. What does the distribution of votes and of spending look like?

1. (3 points) In separate histograms, show both the distribution of votes (measured in
   `results_house$general_percent` for now) and spending (measured in `ttl_disb`). Use a log trans-
   form if appropriate for each visualization. How would you describe what you see in these two plots?

## 2. Exploring the relationship between spending and votes.

2. (3 points) Create a new dataframe by joining `results_house` and `campaigns` using the `inner_join` function from `dplyr`. (We use the format `package::function` – so `dplyr::inner_join`.)

```
nrow(results_house)
```

```
## [1] 2110
```

```
nrow(campaigns)
```

```
## [1] 1898
```

```
d1 <- inner_join(results_house, campaigns, by = NULL)
```

```
## Joining, by = "cand_id"
```

```
#d1 <- merge(results_house, campaigns, by = "cand_id")
#d2 <- merge(results_house, campaigns)
```

```
nrow(d1)
```
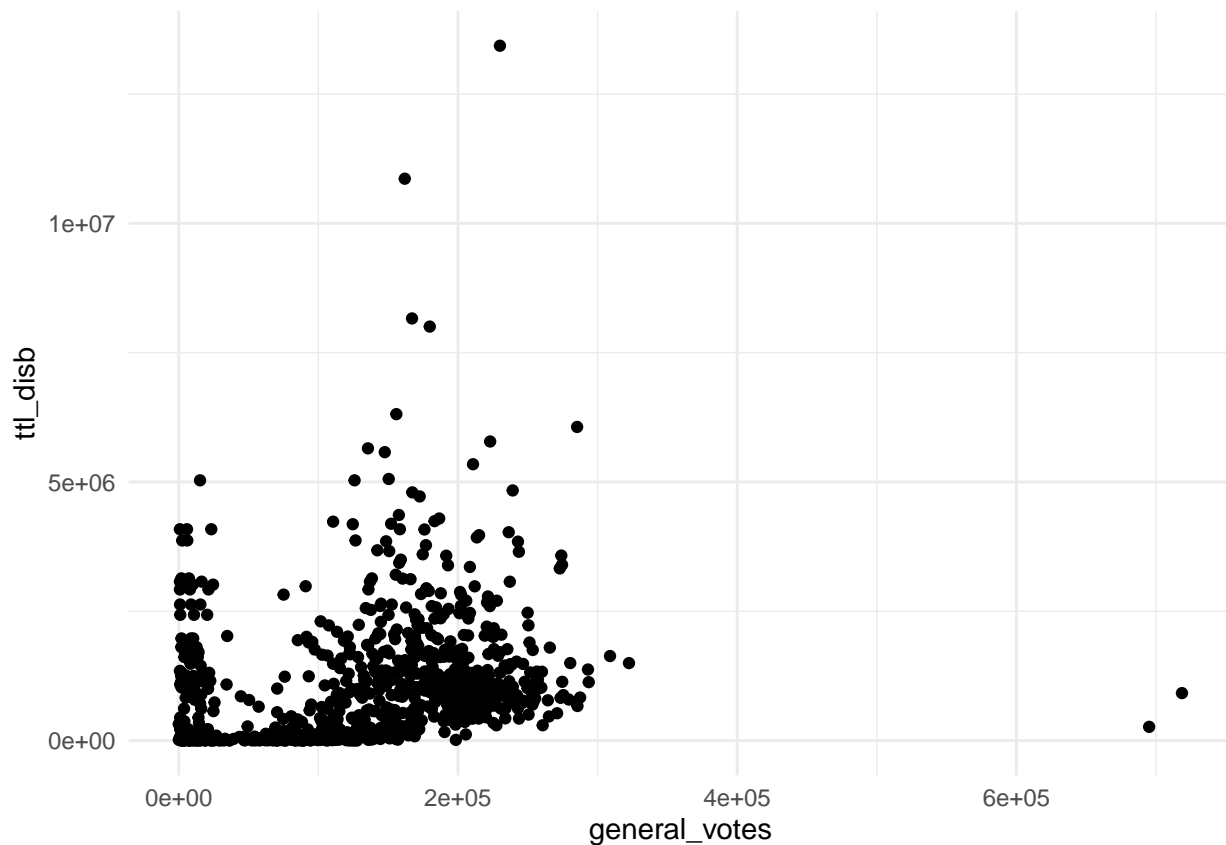
```
## [1] 1342
```

```
#nrow(d2)

#comparison <- compare(d1,d2,allowAll=TRUE)
#comparison

#summary(d1)
#summary(d2)
```

3. (3 points) Produce a scatter plot of `general_votes` on the y-axis and `ttl_disb` on the x-axis. What do you observe about the shape of the joint distribution?

```
ggplot(d1, aes(x=general_votes, y=ttl_disb)) + geom_point()
```

## Warning: Removed 462 rows containing missing values (geom_point).
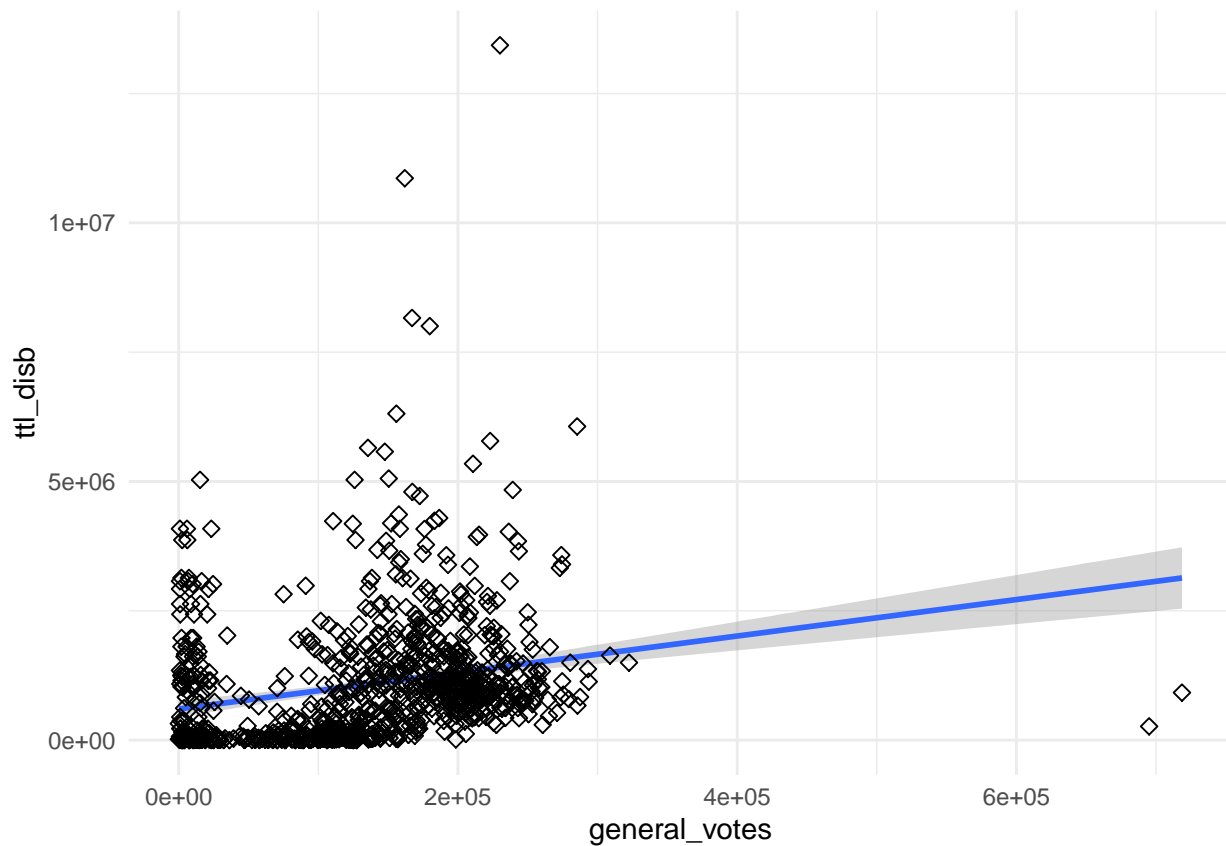


```
# Change the point size, and shape
sp <- ggplot(d1, aes(x=general_votes, y=ttl_disb )) +
  geom_smooth(method=lm)+
  geom_point(size=2, shape=23)

sp
```

## `geom_smooth()` using formula 'y ~ x'

4

```
## Warning: Removed 462 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 462 rows containing missing values (geom_point).
```



```
#sp + geom_density_2d()
```

4. (3 points) Create a new variable to indicate whether each individual is a "Democrat", "Republican" or "Other Party".

- Here's an example of how you might use `mutate` and `case_when` together to create a variable.

```
starwars %>%
  select(name:mass, gender, species) %>%
  mutate(
  type = case_when(
    height > 200 | mass > 200 ~ "large",
    species == "Droid"        ~ "robot",
    TRUE                      ~ "other"
    )
  )
```

Once you've produced the new variable, plot your scatter plot again, but this time adding an argument into the `aes()` function that colors the points by party membership. What do you observe about the distribution of all three variables?

```
d2<-d1 %>%
  select(cand_pty_affiliation, general_votes, ttl_disb, state) %>%
  na.omit() %>%
    mutate(
    can_party = case_when(
      cand_pty_affiliation=="REP" ~ "REP",
      cand_pty_affiliation=="DEM" ~ "DEM",
      TRUE ~ "Other"
    )
  )



write.csv(d2, "d2.csv")

head(d2)
```

```
## # A tibble: 6 x 5
##   cand_pty_affiliation general_votes ttl_disb state can_party
##   <chr>                        <dbl>    <dbl> <chr> <chr>
## 1 REP                         208083 1172750. AL    REP
## 2 REP                         134886 1850536. AL    REP
## 3 DEM                         112089   36844  AL    DEM
## 4 REP                         192164 1071289. AL    REP
## 5 DEM                          94549    7348  AL    DEM
## 6 REP                         235925 1394461. AL    REP
```
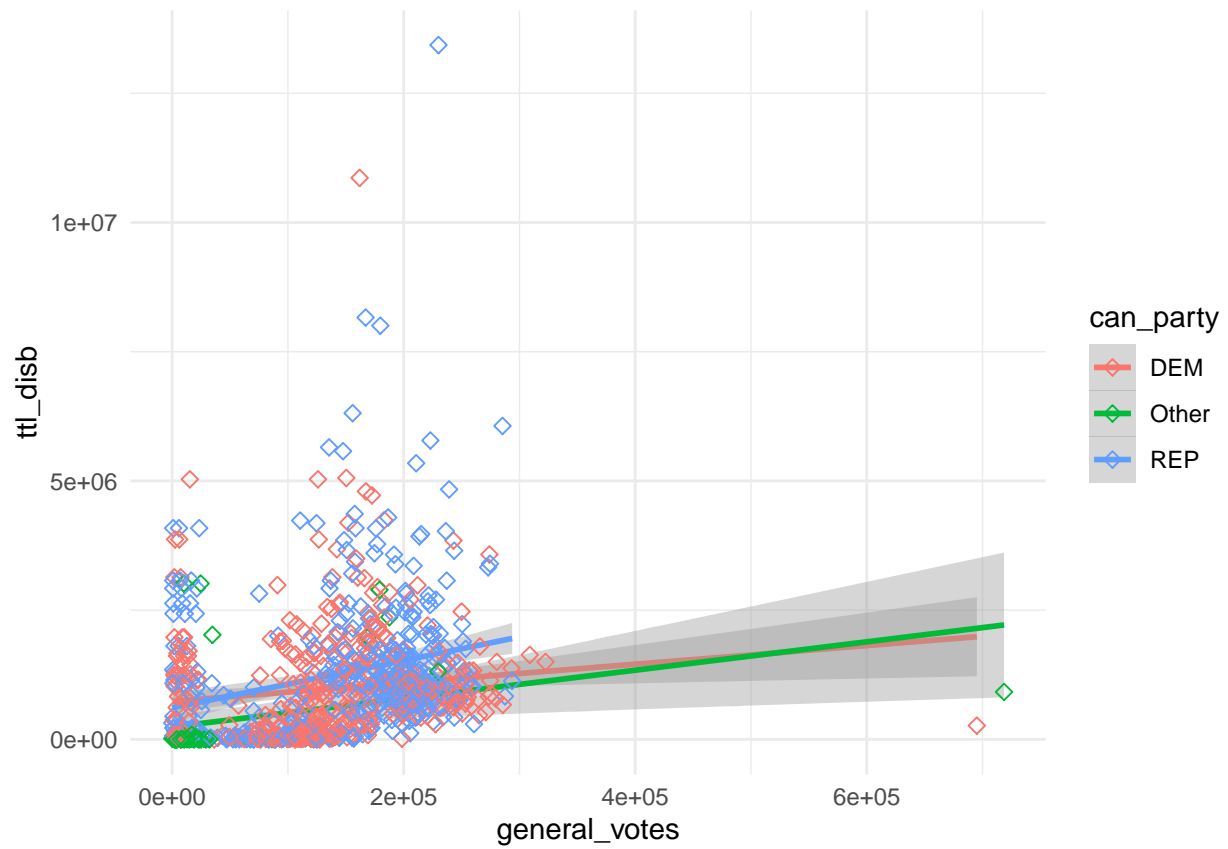
```
sp <- ggplot(d2, aes(x=general_votes, y=ttl_disb, color=can_party)) +
  geom_smooth(method=lm)+
  geom_point(size=2, shape=23)
sp
```
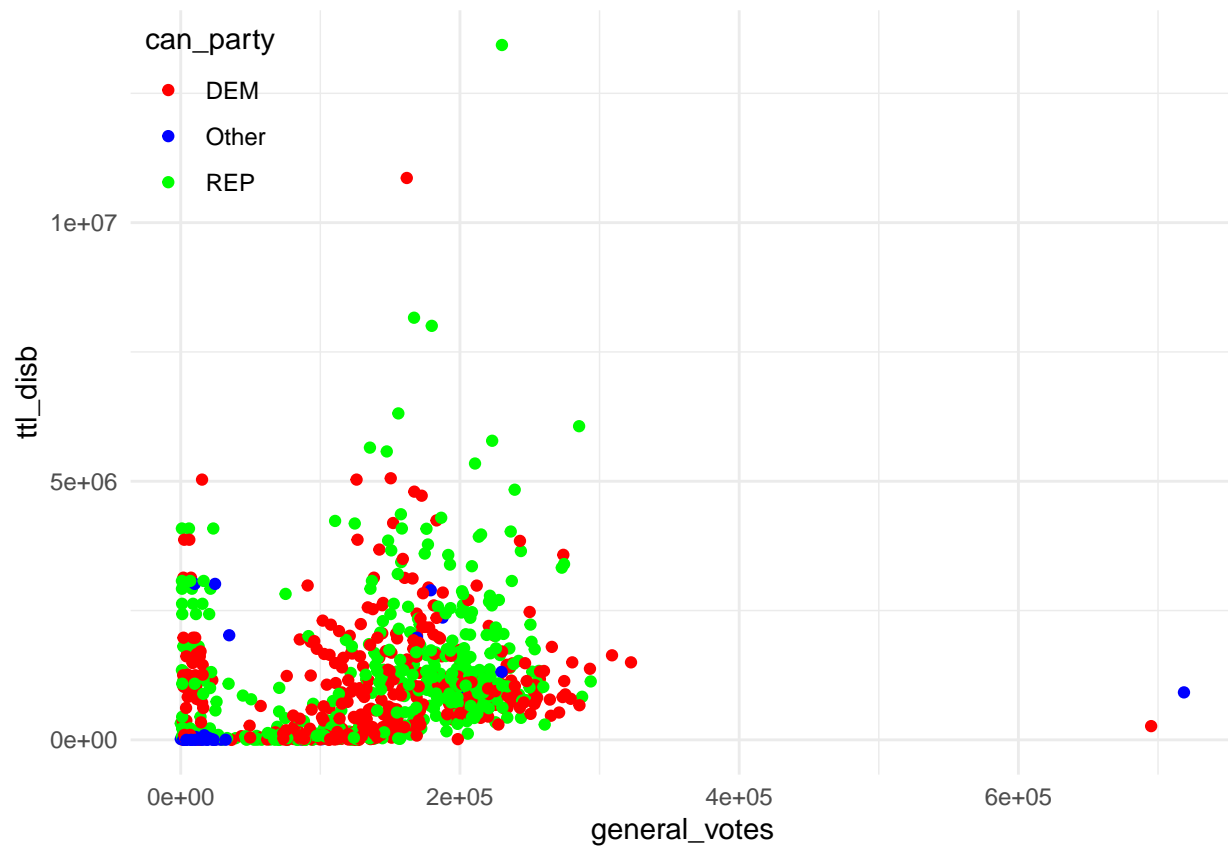
```
## `geom_smooth()` using formula 'y ~ x'
```
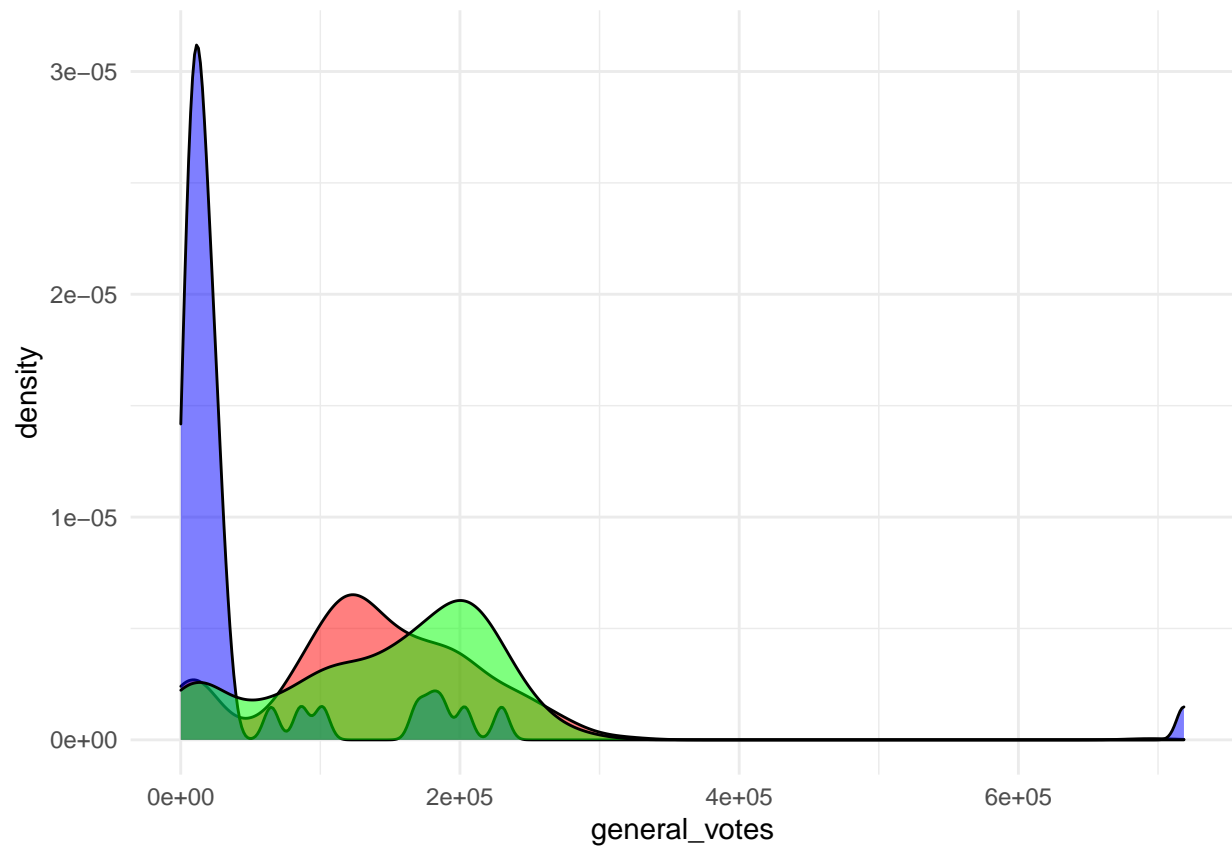
```
p1<-ggplot(d2, aes(x=general_votes, y=ttl_disb, color=can_party)) +
  geom_point() +
  scale_color_manual(values = c("red", "blue", "green")) +
  theme(legend.position=c(0,1), legend.justification=c(0,1))
p1
```
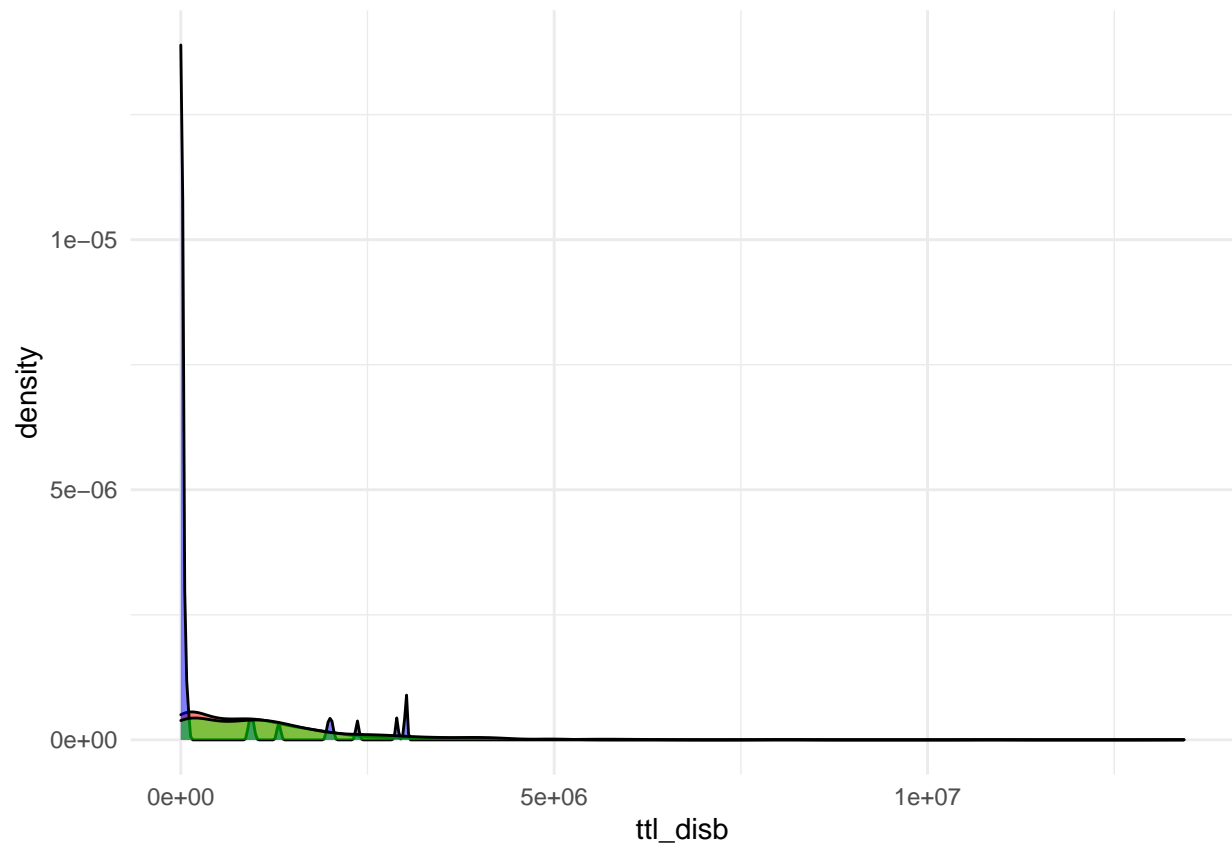
```
p2<-ggplot(d2, aes(x=general_votes, fill=can_party)) +
  geom_density(alpha=.5) +
  scale_fill_manual(values =  c("red", "blue", "green")) +
  theme(legend.position = "none")
p2
```

```
# Marginal density plot of y (right panel)
p3<-ggplot(d2, aes(x=ttl_disb, fill=can_party)) +
  geom_density(alpha=.5) +
  scale_fill_manual(values =  c("red", "blue", "green")) +
  theme(legend.position = "none")
p3
```
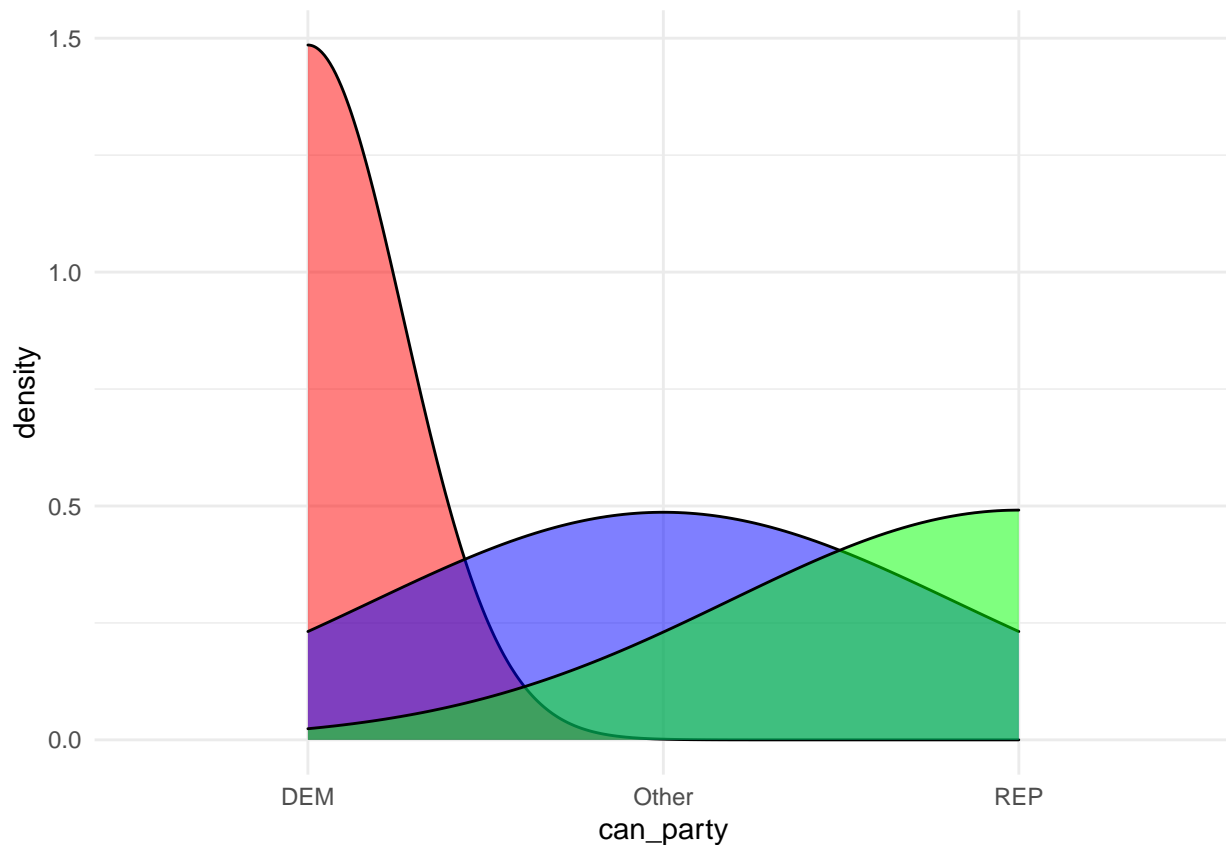
```
p3<-ggplot(d2, aes(x=can_party, fill=can_party)) +
  geom_density(alpha=.5) +
  scale_fill_manual(values =  c("red", "blue", "green")) +
  theme(legend.position = "none")
p3
```

```
#sp + geom_density_2d()
```

```
#summary(d1)
```

## Produce a Descriptive Model

5. (5 Points) Given your observations, produce a linear model that you think does a good job at describing the relationship between candidate spending and votes they receive. You should decide what transformation to apply to spending (if any), what transformation to apply to votes (if any) and also how to include the party affiliation.

```
summary(d2$state)
```

```
##    Length     Class      Mode
##       880 character character
```

```
d2$disb <- log(d2$ttl_disb)
d2$votes <- log(d2$general_votes)
```

```
write.csv(d2, "d2.csv")
```

```
#d2[which(!is.finite(d2))] <- 0
#d2 <- d2[is.finite(rowSums(d2)),]
```

```
d2[d2 == -Inf] <- 0

#data_new <- d2                                    # Duplicate data

#d2[is.na(d2$disb) | d2$disb == "Inf"] <- NA   # Replace NaN & Inf with NA

#d3 <- data_new

head(d2)
```

```
## # A tibble: 6 x 7
##    cand_pty_affiliation general_votes ttl_disb state can_party  disb votes
##    <chr>                        <dbl>    <dbl> <chr> <chr>      <dbl> <dbl>
## 1 REP                         208083 1172750. AL    REP         14.0  12.2
## 2 REP                         134886 1850536. AL    REP         14.4  11.8
## 3 DEM                         112089   36844  AL    DEM         10.5  11.6
## 4 REP                         192164 1071289. AL    REP         13.9  12.2
## 5 DEM                          94549    7348  AL    DEM          8.90 11.5
## 6 REP                         235925 1394461. AL    REP         14.1  12.4
```

```
head(d2$disb)
```

```
## [1] 13.974862 14.430986 10.514448 13.884374  8.902183 14.148019
```

```
#d3<-d3%>%na.omit()

fit <- lm(d2$general_votes ~ d2$disb + d2$state + d2$can_party)

summary(fit)
```

```
##
## Call:
## lm(formula = d2$general_votes ~ d2$disb + d2$state + d2$can_party)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -378949  -35379   -1422   30616  228002
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -78916.4    39004.1  -2.023   0.0434 *
## d2$disb        15114.2      880.1  17.174  < 2e-16 ***
## d2$stateAL     56260.9    40260.6   1.397   0.1627
## d2$stateAR     57802.7    43824.5   1.319   0.1876
## d2$stateAS    -72654.9    47899.7  -1.517   0.1297
## d2$stateAZ     29836.2    39421.4   0.757   0.4494
## d2$stateCA     17128.9    37399.7   0.458   0.6471
## d2$stateCO     64120.1    39580.5   1.620   0.1056
## d2$stateCT    -19862.7    39417.3  -0.504   0.6145
## d2$stateDC    146804.2    64142.8   2.289   0.0223 *
## d2$stateDE     81261.9    52351.8   1.552   0.1210
## d2$stateFL     47406.9    37680.6   1.258   0.2087
```

```
## d2$stateGA          76113.3     39053.8    1.949    0.0516 .
## d2$stateGU         -85174.9     52378.5   -1.626    0.1043
## d2$stateHI          24440.2     47823.8    0.511    0.6095
## d2$stateIA          56284.4     41384.9    1.360    0.1742
## d2$stateID          64117.5     45397.5    1.412    0.1582
## d2$stateIL          53099.1     38209.2    1.390    0.1650
## d2$stateIN          39081.9     39168.9    0.998    0.3187
## d2$stateKS          -5399.2     40934.3   -0.132    0.8951
## d2$stateKY          66803.2     40264.3    1.659    0.0975 .
## d2$stateLA         -36807.4     38941.3   -0.945    0.3448
## d2$stateMA          84732.1     39304.1    2.156    0.0314 *
## d2$stateMD          53793.9     39578.7    1.359    0.1745
## d2$stateME          59654.7     45336.0    1.316    0.1886
## d2$stateMI          50454.6     38311.6    1.317    0.1882
## d2$stateMN          67285.5     39467.7    1.705    0.0886 .
## d2$stateMO          74373.2     39785.7    1.869    0.0619 .
## d2$stateMP          -3586.2     64623.1   -0.055    0.9558
## d2$stateMS          70879.1     42832.0    1.655    0.0983 .
## d2$stateMT          94261.3     52358.5    1.800    0.0722 .
## d2$stateNC          72652.2     38502.5    1.887    0.0595 .
## d2$stateND          44297.4     47882.1    0.925    0.3552
## d2$stateNE          42828.1     45344.4    0.945    0.3452
## d2$stateNH         -13691.4     43832.0   -0.312    0.7548
## d2$stateNJ          38246.2     38684.9    0.989    0.3231
## d2$stateNM          12176.5     43801.2    0.278    0.7811
## d2$stateNV           1680.0     40934.3    0.041    0.9673
## d2$stateNY         -71680.4     37323.9   -1.920    0.0551 .
## d2$stateOH          59912.9     38261.6    1.566    0.1178
## d2$stateOK          52178.6     43815.3    1.191    0.2340
## d2$stateOR          91212.5     41974.8    2.173    0.0301 *
## d2$statePA          64973.1     38275.0    1.698    0.0900 .
## d2$statePR         434255.9     48096.7    9.029   < 2e-16 ***
## d2$stateRI           2223.8     45377.7    0.049    0.9609
## d2$stateSC           3311.5     39057.5    0.085    0.9325
## d2$stateSD          53147.5     52347.7    1.015    0.3103
## d2$stateTN          56307.3     39638.8    1.421    0.1558
## d2$stateTX          27554.8     37713.9    0.731    0.4652
## d2$stateUT           8442.1     41387.1    0.204    0.8384
## d2$stateVA          57819.4     38747.1    1.492    0.1360
## d2$stateVI         -98991.8     64146.4   -1.543    0.1232
## d2$stateVT         138293.3     64139.7    2.156    0.0314 *
## d2$stateWA          43843.3     39278.6    1.116    0.2647
## d2$stateWI          65675.6     39171.6    1.677    0.0940 .
## d2$stateWV             76.6     41999.6    0.002    0.9985
## d2$stateWY          16878.2     47890.3    0.352    0.7246
## d2$can_partyOther -72265.7      8767.5   -8.242 6.65e-16 ***
## d2$can_partyREP      683.6      3727.5    0.183    0.8545
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 52350 on 821 degrees of freedom
## Multiple R-squared:  0.6041, Adjusted R-squared:  0.5761
## F-statistic:  21.6 on 58 and 821 DF,  p-value: < 2.2e-16
```
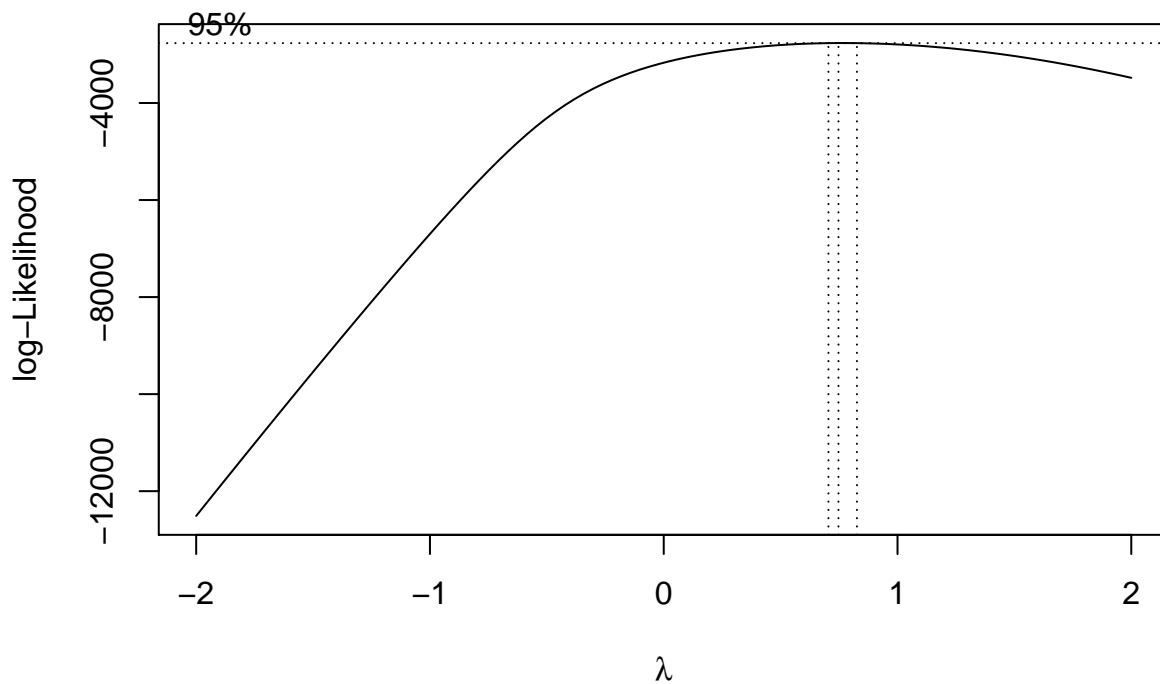
```
## boxcox test
library(MASS)
```

```
##
## Attaching package: 'MASS'

## The following object is masked from 'package:patchwork':
##
##     area

## The following object is masked from 'package:dplyr':
##
##     select
```
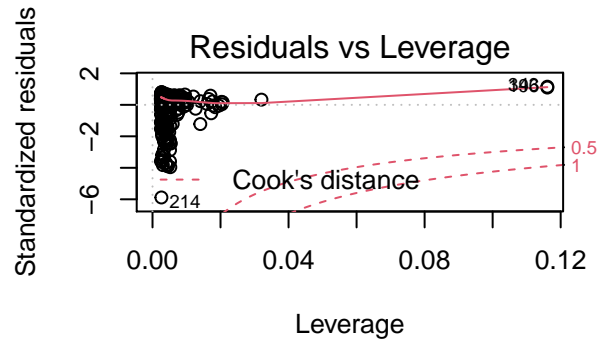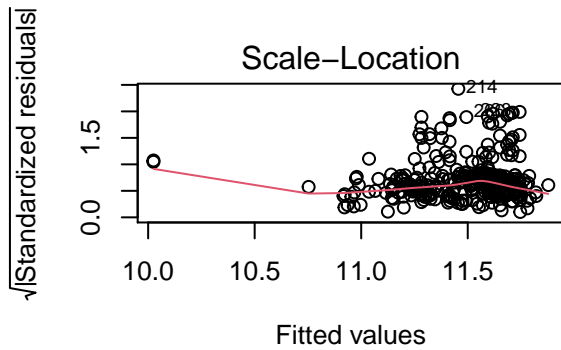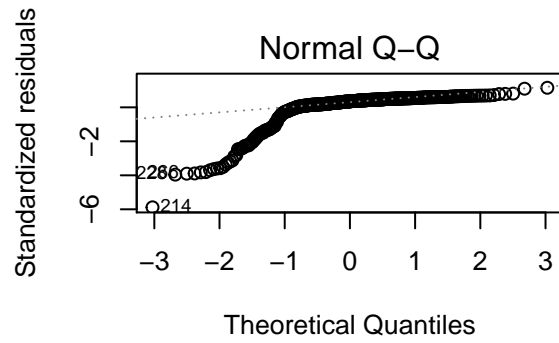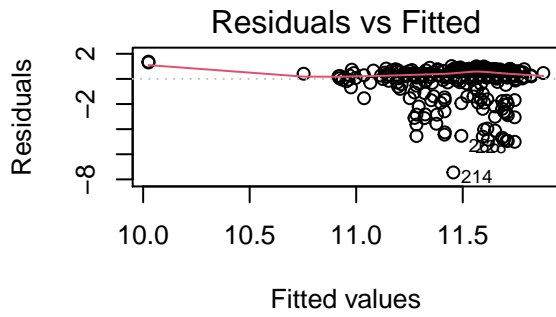
```
boxcox(general_votes~poly(disb,2),
       data = d2)
```
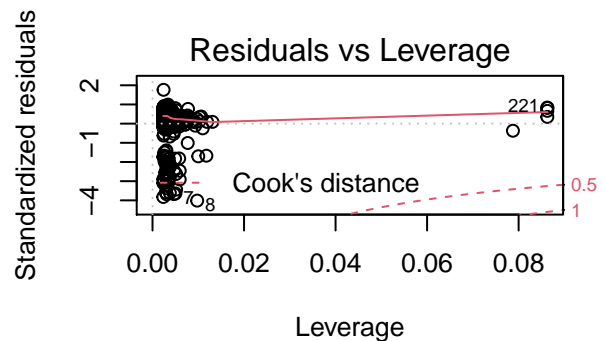


```
g1 <- filter(d2, can_party == "REP")
g2 <- filter(d2, can_party == "DEM")
g3 <- filter(d2, can_party == "Other")


fit <- lm(g1$votes ~ g1$disb)
par(mfrow=c(2,2))
plot (fit)
```
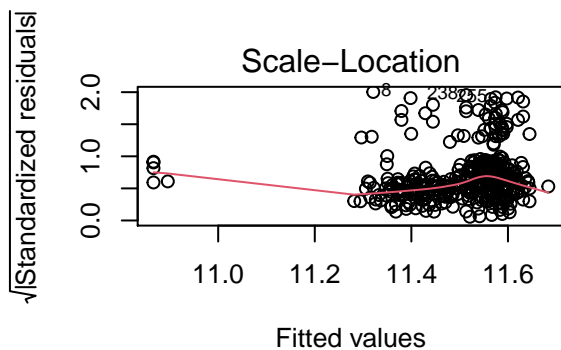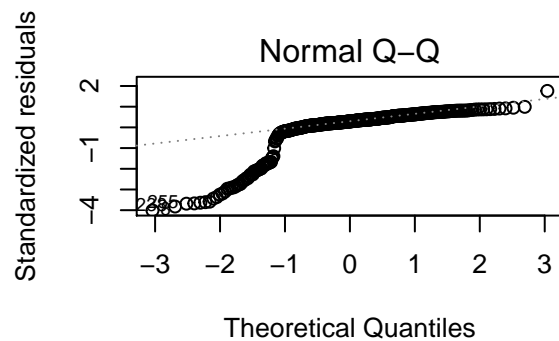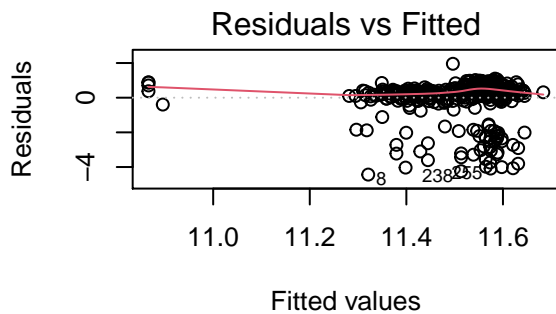
## Residuals vs Fitted

## Normal Q–Q

## Scale–Location

## Residuals vs Leverage

```
fit1 <- lm(g2$votes ~ g2$disb)
par(mfrow=c(2,2))
plot (fit1)
```

## Residuals vs Fitted

## Normal Q–Q

## Scale–Location

## Residuals vs Leverage

15

```
fit2 <- lm(g3$votes ~ g3$disb)
par(mfrow=c(2,2))
plot (fit2)
```



```
summary(fit)
```

```
##
## Call:
## lm(formula = g1$votes ~ g1$disb)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -7.4485  0.1288  0.4484  0.6419  1.3730
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 10.02502    0.43189  23.212  < 2e-16 ***
## g1$disb      0.11290    0.03245   3.479 0.000557 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.268 on 404 degrees of freedom
## Multiple R-squared:  0.0291, Adjusted R-squared:  0.02669
## F-statistic: 12.11 on 1 and 404 DF,  p-value: 0.0005571
```

```
summary(fit1)
```

```
##
## Call:
## lm(formula = g2$votes ~ g2$disb)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -4.4355  0.0615  0.3208  0.5981  1.9554
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 10.86598    0.32627  33.304   <2e-16 ***
## g2$disb      0.05047    0.02509   2.011   0.0449 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.111 on 421 degrees of freedom
## Multiple R-squared:  0.009519,   Adjusted R-squared:  0.007166
## F-statistic: 4.046 on 1 and 421 DF,  p-value: 0.04491
```

```
summary(fit2)
```

```
##
## Call:
## lm(formula = g3$votes ~ g3$disb)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -5.3535 -0.5197  0.1090  0.7988  2.6582
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.08213    0.63215  11.203 4.06e-15 ***
## g3$disb      0.27274    0.06218   4.387 6.10e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.257 on 49 degrees of freedom
## Multiple R-squared:  0.282,  Adjusted R-squared:  0.2673
## F-statistic: 19.24 on 1 and 49 DF,  p-value: 6.103e-05
```

```
d2[d2 == -Inf] <- 0
```

```
head(d2)
```

```
## # A tibble: 6 x 7
##   cand_pty_affiliation general_votes ttl_disb state can_party  disb votes
##   <chr>                        <dbl>    <dbl> <chr> <chr>     <dbl> <dbl>
## 1 REP                         208083 1172750. AL    REP        14.0  12.2
## 2 REP                         134886 1850536. AL    REP        14.4  11.8
## 3 DEM                         112089   36844  AL    DEM        10.5  11.6
## 4 REP                         192164 1071289. AL    REP        13.9  12.2
## 5 DEM                          94549    7348  AL    DEM         8.90  11.5
## 6 REP                         235925 1394461. AL    REP        14.1  12.4
```
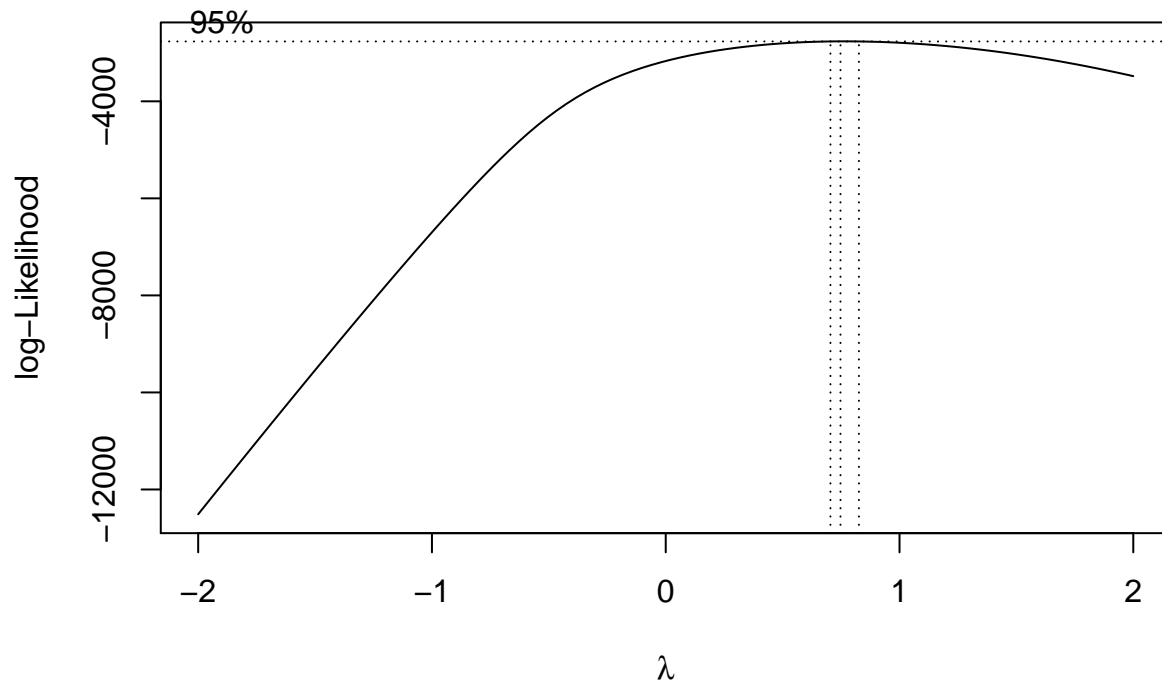
17

```
head(d2$disb)
```

```
## [1] 13.974862 14.430986 10.514448 13.884374  8.902183 14.148019
```

```
#d3<-d3%>%na.omit()
```

```
fit <- lm(d2$general_votes ~ d2$disb)
summary(fit)
```

```
##
## Call:
## lm(formula = d2$general_votes ~ d2$disb)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -170750  -34066    7653   45029  568412
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -46697      14420  -3.238  0.00125 **
## d2$disb        14339       1109  12.928  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 73740 on 878 degrees of freedom
## Multiple R-squared:  0.1599, Adjusted R-squared:  0.159
## F-statistic: 167.1 on 1 and 878 DF,  p-value: < 2.2e-16
```

```
## boxcox test
library(MASS)
boxcox(general_votes~poly(disb,2),
        data = d2)
```

```r
# g0 <- d2
# g0$votes <- log10(g0$general_votes)
# g0$disb <- log10(g0$ttl_disb)
# g0[g0 == -Inf] <- 0


g0 <- d2
g0$votes <- g0$general_votes
g0$disb <- g0$ttl_disb
g0[g0 == -Inf] <- 0

g1 <- filter(d2, can_party == "REP")
g1$votes <- g1$general_votes*g1$general_votes
g1$disb <- log(g1$ttl_disb)
g1[g1 == -Inf] <- 0



g2 <- filter(d2, can_party == "DEM")
g2$votes <- g2$general_votes*g2$general_votes
g2$disb <- log(g2$ttl_disb)
g2[g2 == -Inf] <- 0

g3 <- filter(d2, can_party == "Other")
g3$votes <- g3$general_votes
g3$disb <- log(g3$ttl_disb)
g3[g3 == -Inf] <- 0


write.csv(g1, "g1.csv")
write.csv(g2, "g2.csv")
write.csv(g3, "g3.csv")
```

```
fit0 <- lm(g0$votes ~ g0$disb + g0$state + g0$can_party )
par(mfrow=c(2,2))
plot (fit0)
```

```
## Warning: not plotting observations with leverage one:
##    168, 640, 815, 837
```



```
fit1 <- lm(g1$votes ~ g1$disb + g1$state )
par(mfrow=c(2,2))
plot (fit1)
```

```
## Warning: not plotting observations with leverage one:
##    7, 8, 75, 113, 205, 293, 338, 406
```

```
fit2 <- lm(g2$votes ~ g2$disb + g2$state )
par(mfrow=c(2,2))
plot (fit2)
```

```
## Warning: not plotting observations with leverage one:
##   6, 91, 92, 126, 130, 211, 212, 307, 322, 341, 356, 389, 401, 423
```

```
## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs produced
```

```
## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs produced
```

## Residuals vs Fitted

## Normal Q–Q

## Scale–Location

## Residuals vs Leverage

```r
fit3 <- lm(g3$votes ~ g3$disb + g3$state )
par(mfrow=c(2,2))
plot (fit3)
```

```
## Warning: not plotting observations with leverage one:
##   1, 7, 15, 16, 31, 32, 39, 40, 45, 46, 47, 51
```
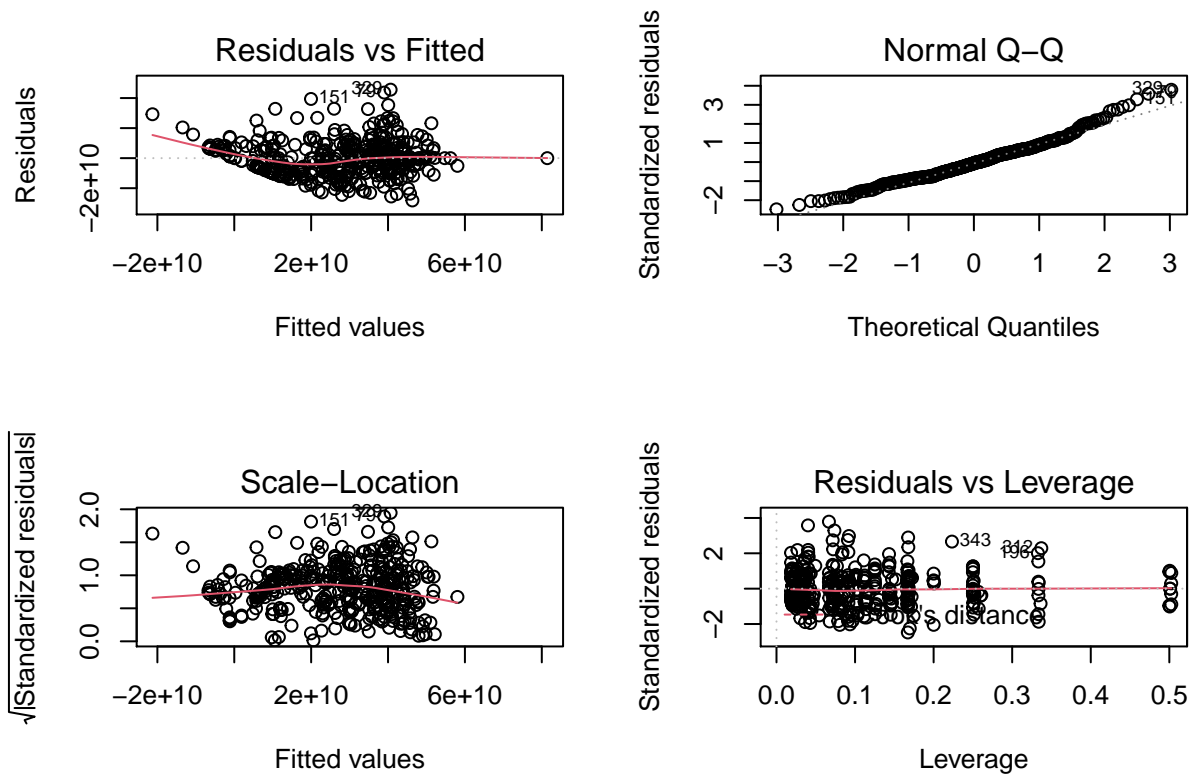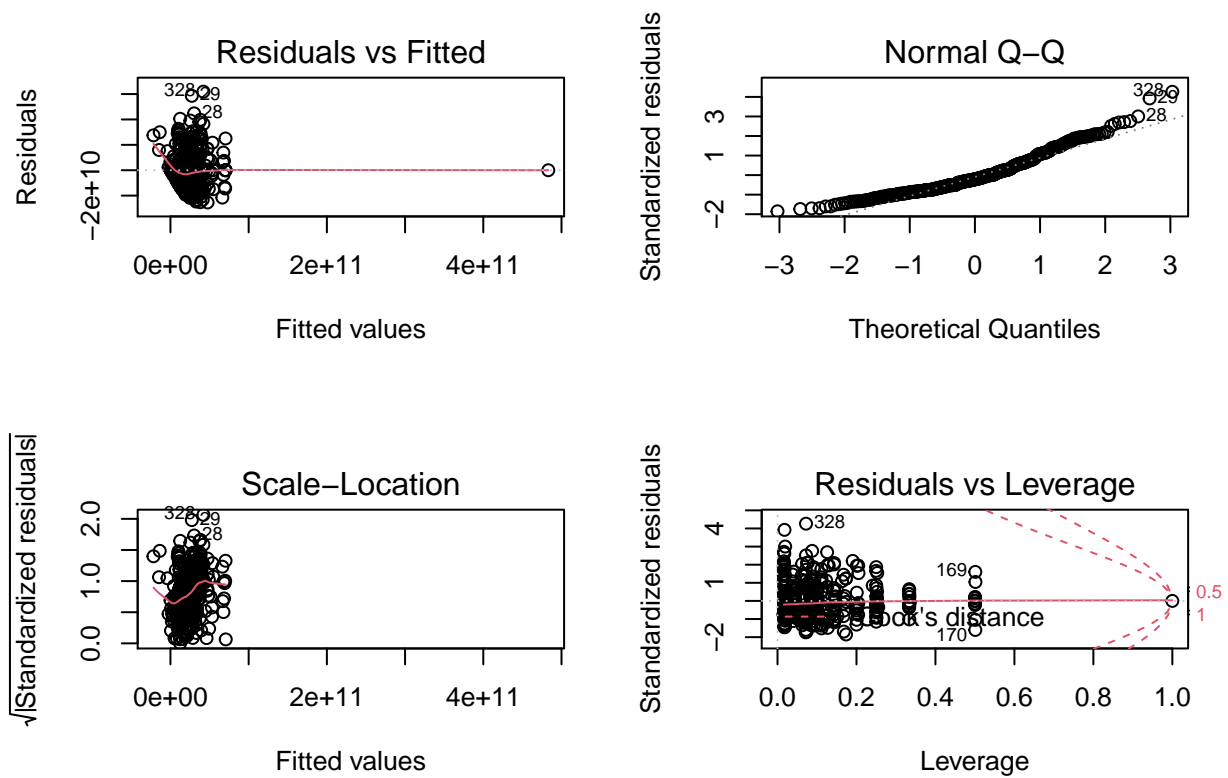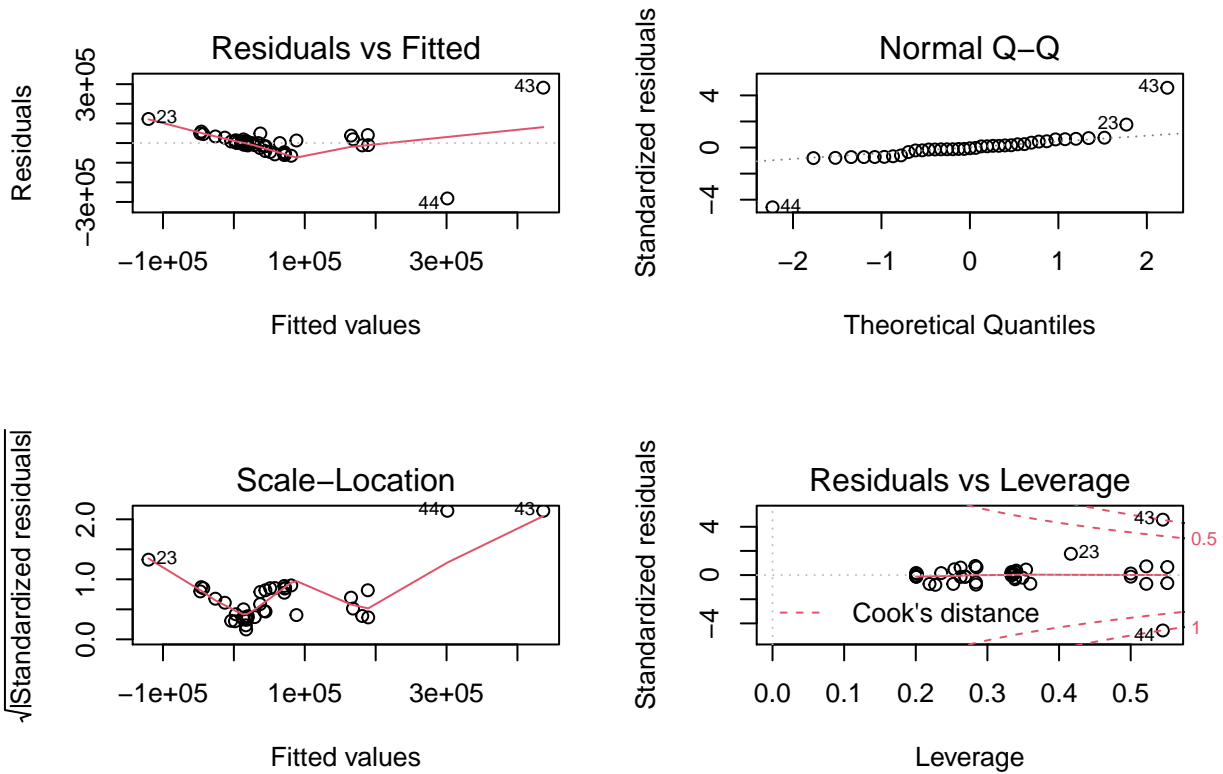
```
summary(fit0)
```

```
##
## Call:
## lm(formula = g0$votes ~ g0$disb + g0$state + g0$can_party)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -415756  -39794   -5242   36879  269903
##
## Coefficients:
##                 Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.139e+05  4.120e+04   2.763  0.00585 **
## g0$disb        1.518e-02  1.661e-03   9.138  < 2e-16 ***
## g0$stateAL     4.215e+04  4.471e+04   0.943  0.34607
## g0$stateAR     4.851e+04  4.867e+04   0.997  0.31926
## g0$stateAS    -1.112e+05  5.312e+04  -2.094  0.03660 *
## g0$stateAZ     1.136e+04  4.377e+04   0.260  0.79531
## g0$stateCA    -1.482e+02  4.153e+04  -0.004  0.99715
## g0$stateCO     5.012e+04  4.395e+04   1.140  0.25447
## g0$stateCT    -2.964e+04  4.378e+04  -0.677  0.49848
## g0$stateDC     1.442e+05  7.125e+04   2.024  0.04325 *
## g0$stateDE     7.554e+04  5.815e+04   1.299  0.19425
## g0$stateFL     2.945e+04  4.184e+04   0.704  0.48167
## g0$stateGA     5.646e+04  4.335e+04   1.302  0.19316
## g0$stateGU    -1.001e+05  5.817e+04  -1.721  0.08558 .
## g0$stateHI     2.308e+04  5.312e+04   0.434  0.66406
## g0$stateIA     4.883e+04  4.597e+04   1.062  0.28847
```

```
## g0$stateID         5.123e+04  5.041e+04   1.016  0.30988
## g0$stateIL         4.114e+04  4.244e+04   0.969  0.33260
## g0$stateIN         2.570e+04  4.350e+04   0.591  0.55475
## g0$stateKS        -9.865e+03  4.547e+04  -0.217  0.82828
## g0$stateKY         5.083e+04  4.471e+04   1.137  0.25585
## g0$stateLA        -4.898e+04  4.324e+04  -1.133  0.25764
## g0$stateMA         7.588e+04  4.365e+04   1.738  0.08254 .
## g0$stateMD         4.579e+04  4.396e+04   1.042  0.29784
## g0$stateME         4.369e+04  5.037e+04   0.867  0.38598
## g0$stateMI         3.268e+04  4.254e+04   0.768  0.44261
## g0$stateMN         5.922e+04  4.386e+04   1.350  0.17730
## g0$stateMO         6.812e+04  4.419e+04   1.541  0.12358
## g0$stateMP         6.315e+03  7.178e+04   0.088  0.92991
## g0$stateMS         3.385e+04  4.749e+04   0.713  0.47613
## g0$stateMT         6.441e+04  5.838e+04   1.103  0.27026
## g0$stateNC         5.234e+04  4.274e+04   1.225  0.22102
## g0$stateND         2.870e+04  5.317e+04   0.540  0.58947
## g0$stateNE         3.875e+04  5.037e+04   0.769  0.44191
## g0$stateNH        -9.872e+03  4.869e+04  -0.203  0.83938
## g0$stateNJ         1.872e+04  4.295e+04   0.436  0.66310
## g0$stateNM         5.723e+03  4.865e+04   0.118  0.90638
## g0$stateNV        -5.810e+03  4.546e+04  -0.128  0.89835
## g0$stateNY        -7.894e+04  4.146e+04  -1.904  0.05726 .
## g0$stateOH         4.231e+04  4.248e+04   0.996  0.31949
## g0$stateOK         4.615e+04  4.867e+04   0.948  0.34321
## g0$stateOR         8.356e+04  4.662e+04   1.792  0.07343 .
## g0$statePA         5.607e+04  4.251e+04   1.319  0.18754
## g0$statePR         4.313e+05  5.342e+04   8.074 2.42e-15 ***
## g0$stateRI        -2.106e+04  5.037e+04  -0.418  0.67599
## g0$stateSC        -1.508e+04  4.336e+04  -0.348  0.72814
## g0$stateSD         4.393e+04  5.815e+04   0.756  0.45016
## g0$stateTN         2.770e+04  4.397e+04   0.630  0.52899
## g0$stateTX         1.064e+04  4.187e+04   0.254  0.79956
## g0$stateUT        -3.352e+03  4.597e+04  -0.073  0.94188
## g0$stateVA         4.654e+04  4.303e+04   1.082  0.27975
## g0$stateVI        -1.045e+05  7.125e+04  -1.466  0.14300
## g0$stateVT         1.387e+05  7.124e+04   1.947  0.05184 .
## g0$stateWA         2.911e+04  4.361e+04   0.667  0.50465
## g0$stateWI         4.214e+04  4.350e+04   0.969  0.33301
## g0$stateWV        -8.468e+03  4.665e+04  -0.182  0.85599
## g0$stateWY        -8.747e+03  5.316e+04  -0.165  0.86935
## g0$can_partyOther -1.104e+05  9.287e+03 -11.891  < 2e-16 ***
## g0$can_partyREP    1.661e+03  4.157e+03   0.400  0.68957
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 58140 on 821 degrees of freedom
## Multiple R-squared:  0.5116, Adjusted R-squared:  0.477
## F-statistic: 14.82 on 58 and 821 DF,  p-value: < 2.2e-16
```

```
summary(fit1)
```

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

```
## lm(formula = g1$votes ~ g1$disb + g1$state)
##
## Residuals:
##         Min        1Q     Median        3Q        Max
## -2.815e+10 -8.813e+09 -3.549e+08  6.976e+09  4.558e+10
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.780e+10  1.339e+10  -2.824  0.00502 **
## g1$disb      4.389e+09  3.491e+08  12.570  < 2e-16 ***
## g1$stateAL   2.081e+10  1.345e+10   1.548  0.12261
## g1$stateAR   1.472e+10  1.392e+10   1.057  0.29121
## g1$stateAS  -1.141e+10  1.764e+10  -0.647  0.51798
## g1$stateAZ   7.702e+09  1.321e+10   0.583  0.56021
## g1$stateCA  -1.213e+09  1.262e+10  -0.096  0.92352
## g1$stateCO   1.745e+10  1.332e+10   1.310  0.19091
## g1$stateCT  -2.819e+09  1.348e+10  -0.209  0.83451
## g1$stateDE   1.347e+10  1.762e+10   0.765  0.44500
## g1$stateFL   1.699e+10  1.270e+10   1.337  0.18193
## g1$stateGA   1.724e+10  1.296e+10   1.330  0.18436
## g1$stateGU  -1.353e+10  1.763e+10  -0.767  0.44343
## g1$stateIA   1.872e+10  1.392e+10   1.345  0.17955
## g1$stateID   2.803e+10  1.525e+10   1.838  0.06692 .
## g1$stateIL   1.186e+10  1.301e+10   0.912  0.36251
## g1$stateIN   1.220e+10  1.321e+10   0.924  0.35610
## g1$stateKS  -1.339e+09  1.364e+10  -0.098  0.92184
## g1$stateKY   2.486e+10  1.331e+10   1.867  0.06273 .
## g1$stateLA  -8.544e+09  1.289e+10  -0.663  0.50788
## g1$stateMA   2.607e+09  1.396e+10   0.187  0.85204
## g1$stateMD   4.015e+09  1.346e+10   0.298  0.76563
## g1$stateME   1.164e+10  1.525e+10   0.763  0.44567
## g1$stateMI   1.207e+10  1.301e+10   0.928  0.35396
## g1$stateMN   1.309e+10  1.331e+10   0.983  0.32633
## g1$stateMO   2.717e+10  1.345e+10   2.020  0.04412 *
## g1$stateMS   2.436e+10  1.398e+10   1.742  0.08234 .
## g1$stateMT   5.069e+10  1.762e+10   2.878  0.00425 **
## g1$stateNC   1.711e+10  1.292e+10   1.323  0.18653
## g1$stateND   3.177e+10  1.761e+10   1.804  0.07201 .
## g1$stateNE   1.409e+10  1.438e+10   0.980  0.32792
## g1$stateNH  -8.923e+09  1.525e+10  -0.585  0.55887
## g1$stateNJ   5.136e+09  1.313e+10   0.391  0.69583
## g1$stateNM  -2.122e+09  1.526e+10  -0.139  0.88946
## g1$stateNV  -2.764e+09  1.392e+10  -0.199  0.84277
## g1$stateNY  -1.740e+10  1.257e+10  -1.384  0.16710
## g1$stateOH   2.203e+10  1.289e+10   1.709  0.08835 .
## g1$stateOK   1.505e+10  1.392e+10   1.081  0.28045
## g1$stateOR   2.319e+10  1.438e+10   1.612  0.10787
## g1$statePA   1.731e+10  1.286e+10   1.346  0.17919
## g1$stateRI   7.234e+07  1.533e+10   0.005  0.99624
## g1$stateSC   1.359e+10  1.345e+10   1.010  0.31296
## g1$stateSD   2.849e+10  1.761e+10   1.618  0.10655
## g1$stateTN   1.653e+10  1.315e+10   1.257  0.20954
## g1$stateTX   6.193e+09  1.266e+10   0.489  0.62515
## g1$stateUT   7.225e+09  1.392e+10   0.519  0.60410
```

```
## g1$stateVA    1.569e+10   1.306e+10     1.201   0.23042
## g1$stateWA    5.199e+09   1.348e+10     0.386   0.69987
## g1$stateWI    2.383e+10   1.345e+10     1.772   0.07729 .
## g1$stateWV   -3.932e+08   1.438e+10    -0.027   0.97820
## g1$stateWY   -1.792e+09   1.761e+10    -0.102   0.91901
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.245e+10 on 355 degrees of freedom
## Multiple R-squared:  0.6461, Adjusted R-squared:  0.5962
## F-statistic: 12.96 on 50 and 355 DF,  p-value: < 2.2e-16
```

```
summary(fit2)
```

```
##
## Call:
## lm(formula = g2$votes ~ g2$disb + g2$state)
##
## Residuals:
##        Min        1Q     Median        3Q       Max
## -2.538e+10 -1.073e+10 -2.444e+09  7.456e+09  6.188e+10
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.076e+10  1.600e+10  -1.922 0.055324 .
## g2$disb      3.097e+09  3.868e+08   8.008 1.56e-14 ***
## g2$stateAL   1.405e+10  1.654e+10   0.850 0.396151
## g2$stateAR   1.135e+10  2.136e+10   0.531 0.595464
## g2$stateAS   4.557e+08  1.852e+10   0.025 0.980386
## g2$stateAZ   1.033e+10  1.612e+10   0.641 0.522153
## g2$stateCA   1.432e+10  1.520e+10   0.942 0.346879
## g2$stateCO   2.555e+10  1.611e+10   1.586 0.113565
## g2$stateCT   6.029e+09  1.588e+10   0.380 0.704452
## g2$stateDC   6.065e+10  2.131e+10   2.846 0.004681 **
## g2$stateDE   4.136e+10  2.131e+10   1.941 0.053040 .
## g2$stateFL   1.665e+10  1.536e+10   1.084 0.278957
## g2$stateGA   3.523e+10  1.632e+10   2.158 0.031575 *
## g2$stateGU  -6.583e+09  2.132e+10  -0.309 0.757684
## g2$stateHI   1.150e+10  1.740e+10   0.661 0.508951
## g2$stateIA   1.720e+10  1.685e+10   1.021 0.307958
## g2$stateID   4.271e+09  2.131e+10   0.200 0.841305
## g2$stateIL   2.302e+10  1.553e+10   1.482 0.139199
## g2$stateIN   1.132e+10  1.628e+10   0.695 0.487225
## g2$stateKS   2.107e+09  1.741e+10   0.121 0.903706
## g2$stateKY   1.378e+10  1.687e+10   0.817 0.414583
## g2$stateLA   6.148e+09  1.654e+10   0.372 0.710317
## g2$stateMA   5.687e+10  1.588e+10   3.580 0.000389 ***
## g2$stateMD   3.498e+10  1.611e+10   2.172 0.030530 *
## g2$stateME   2.645e+10  1.845e+10   1.433 0.152633
## g2$stateMI   1.801e+10  1.564e+10   1.151 0.250297
## g2$stateMN   2.327e+10  1.685e+10   1.381 0.168017
## g2$stateMO   1.891e+10  1.652e+10   1.145 0.253021
## g2$stateMS   1.375e+10  1.846e+10   0.745 0.456887
## g2$stateMT   2.729e+10  2.131e+10   1.280 0.201256
```

```
## g2$stateNC    2.698e+10   1.572e+10    1.716 0.087011 .
## g2$stateND   -7.935e+08   2.132e+10   -0.037 0.970329
## g2$stateNE    4.028e+09   2.131e+10    0.189 0.850209
## g2$stateNH    1.443e+10   1.846e+10    0.782 0.434738
## g2$stateNJ    1.708e+10   1.575e+10    1.085 0.278757
## g2$stateNM    1.056e+10   1.740e+10    0.607 0.544415
## g2$stateNV    4.560e+09   1.685e+10    0.271 0.786798
## g2$stateNY   -1.247e+09   1.519e+10   -0.082 0.934636
## g2$stateOH    1.800e+10   1.563e+10    1.151 0.250283
## g2$stateOK    6.966e+09   2.134e+10    0.326 0.744237
## g2$stateOR    4.134e+10   1.685e+10    2.454 0.014594 *
## g2$statePA    2.885e+10   1.560e+10    1.849 0.065201 .
## g2$statePR    4.752e+11   2.132e+10   22.293  < 2e-16 ***
## g2$stateRI    5.501e+09   1.845e+10    0.298 0.765818
## g2$stateSC    3.419e+09   1.571e+10    0.218 0.827864
## g2$stateSD    8.423e+09   2.131e+10    0.395 0.692931
## g2$stateTN    1.212e+10   1.687e+10    0.718 0.473015
## g2$stateTX    9.048e+09   1.541e+10    0.587 0.557360
## g2$stateUT   -2.441e+08   1.685e+10   -0.014 0.988453
## g2$stateVA    2.226e+10   1.574e+10    1.414 0.158272
## g2$stateVI   -8.469e+09   2.131e+10   -0.397 0.691344
## g2$stateVT    5.865e+10   2.131e+10    2.752 0.006209 **
## g2$stateWA    2.067e+10   1.580e+10    1.308 0.191848
## g2$stateWI    2.545e+10   1.599e+10    1.592 0.112293
## g2$stateWV    1.758e+09   1.743e+10    0.101 0.919686
## g2$stateWY   -1.132e+09   2.132e+10   -0.053 0.957693
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.507e+10 on 367 degrees of freedom
## Multiple R-squared:  0.7772, Adjusted R-squared:  0.7438
## F-statistic: 23.27 on 55 and 367 DF,  p-value: < 2.2e-16
```

```
summary(fit3)
```

```
##
## Call:
## lm(formula = g3$votes ~ g3$disb + g3$state)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -282119  -11977       0   13659  282119
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -153401     110453  -1.389  0.17667
## g3$disb         22085       6306   3.502  0.00169 **
## g3$stateFL     -54719      99904  -0.548  0.58856
## g3$stateID      18427     130043   0.142  0.88841
## g3$stateIL      -2650     102013  -0.026  0.97947
## g3$stateIN     -26225     105545  -0.248  0.80572
## g3$stateKS     -78968     129323  -0.611  0.54675
## g3$stateMA     -45073     105309  -0.428  0.67217
## g3$stateMD     -44706     129001  -0.347  0.73171
```

```
## g3$stateMI      32926     101862    0.323   0.74909
## g3$stateMN      14323     103691    0.138   0.89120
## g3$stateMO      31203     111852    0.279   0.78248
## g3$stateMP     -77640     129149   -0.601   0.55294
## g3$stateND       3729     129615    0.029   0.97727
## g3$stateNH    -132536     132253   -1.002   0.32551
## g3$stateNJ     -27077     112109   -0.242   0.81104
## g3$stateNV     -56156     128984   -0.435   0.66689
## g3$stateNY    -104945     103055   -1.018   0.31790
## g3$stateOH     -48102     111696   -0.431   0.67027
## g3$statePR     286648     111806    2.564   0.01648 *
## g3$stateTN      30708     131593    0.233   0.81731
## g3$stateTX      -1954     129482   -0.015   0.98807
## g3$stateWI      -8531     105822   -0.081   0.93636
## g3$stateWV     -81526     129332   -0.630   0.53396
## g3$stateWY      23580     130886    0.180   0.85843
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 91200 on 26 degrees of freedom
## Multiple R-squared:  0.6474, Adjusted R-squared:  0.322
## F-statistic: 1.989 on 24 and 26 DF,  p-value: 0.04476
```
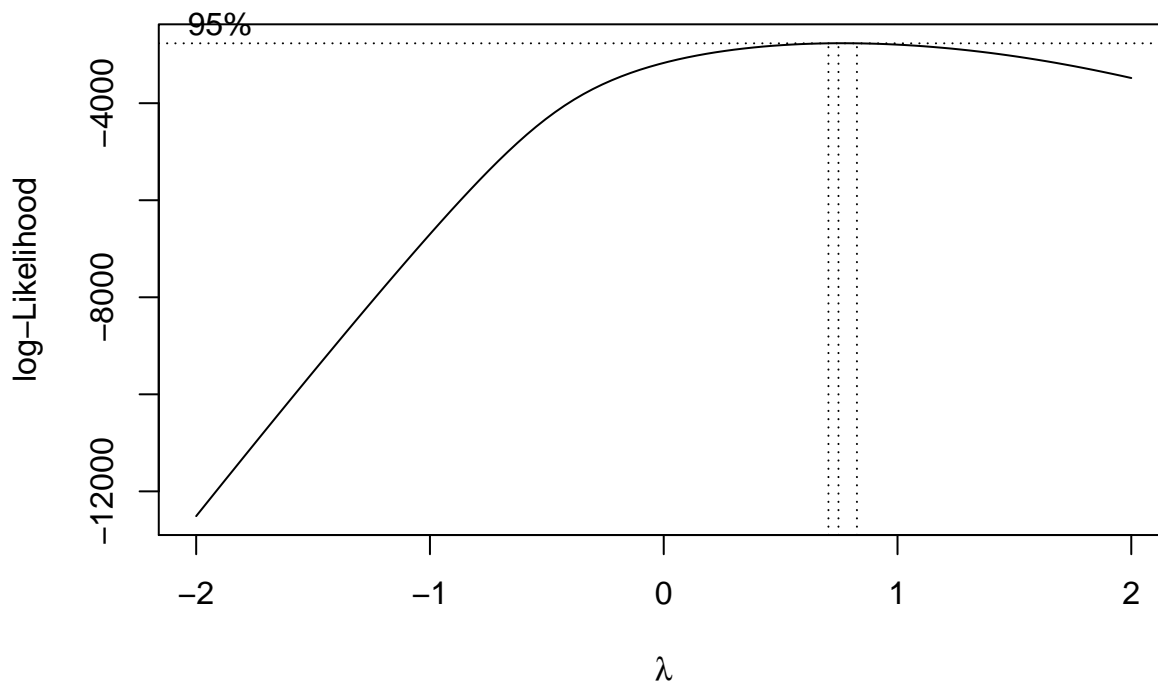
```r
#d2$disb <- log(d2$ttl_disb)
#d2$votes <- log(d2$general_votes)

write.csv(d2, "d2.csv")

#d2[which(!is.finite(d2))] <- 0
#d2 <- d2[is.finite(rowSums(d2)),]
d2[d2 == -Inf] <- 0

#data_new <- d2                              # Duplicate data

#d2[is.na(d2$disb) | d2$disb == "Inf"] <- NA  # Replace NaN & Inf with NA

#d3 <- data_new

head(d2)
```

```
## # A tibble: 6 x 7
##   cand_pty_affiliation general_votes ttl_disb state can_party  disb votes
##   <chr>                        <dbl>    <dbl> <chr> <chr>     <dbl> <dbl>
## 1 REP                         208083 1172750. AL    REP        14.0  12.2
## 2 REP                         134886 1850536. AL    REP        14.4  11.8
## 3 DEM                         112089   36844  AL    DEM        10.5  11.6
## 4 REP                         192164 1071289. AL    REP        13.9  12.2
## 5 DEM                          94549    7348  AL    DEM         8.90  11.5
## 6 REP                         235925 1394461. AL    REP        14.1  12.4
```

```r
head(d2$disb)
```

```
## [1] 13.974862 14.430986 10.514448 13.884374  8.902183 14.148019
```

```
#d3<-d3%>%na.omit()

fit <- lm(d2$general_votes ~ d2$disb)

summary(fit)
```

```
##
## Call:
## lm(formula = d2$general_votes ~ d2$disb)
##
## Residuals:
##      Min      1Q  Median      3Q     Max
## -170750  -34066    7653   45029  568412
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -46697      14420  -3.238  0.00125 **
## d2$disb         14339       1109  12.928  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 73740 on 878 degrees of freedom
## Multiple R-squared:  0.1599, Adjusted R-squared:  0.159
## F-statistic: 167.1 on 1 and 878 DF,  p-value: < 2.2e-16
```

```
## boxcox test
library(MASS)
boxcox(general_votes~poly(disb,2),
        data = d2)
```

```r
g0 <- d2
g0$votes <- log10(g0$general_votes)
g0$disb <- log10(g0$ttl_disb)
g0[g0 == -Inf] <- 0

g1 <- filter(d2, can_party == "REP")
g1$votes <- g1$general_votes*g1$general_votes
g1$disb <- log(g1$ttl_disb)
g1[g1 == -Inf] <- 0


g2 <- filter(d2, can_party == "DEM")
g2$votes <- g2$general_votes*g2$general_votes
g2$disb <- log(g2$ttl_disb)
g2[g2 == -Inf] <- 0

g3 <- filter(d2, can_party == "Other")
g3$votes <- g3$general_votes
g3$disb <- log(g3$ttl_disb)
g3[g3 == -Inf] <- 0


write.csv(g1, "g1.csv")
write.csv(g2, "g2.csv")
write.csv(g3, "g3.csv")

fit0 <- rlm(g0$votes ~ g0$disb)
par(mfrow=c(2,2))
plot (fit)

fit1 <- rlm(g1$votes ~ g1$disb)
par(mfrow=c(2,2))
plot (fit)
```
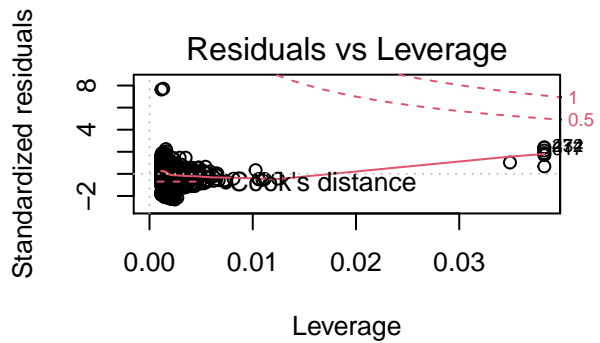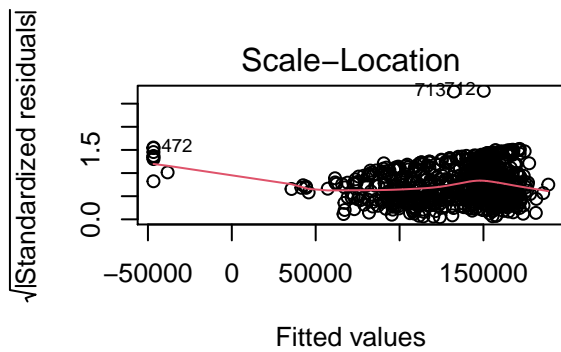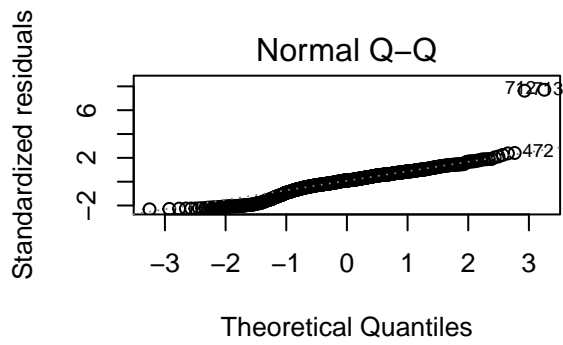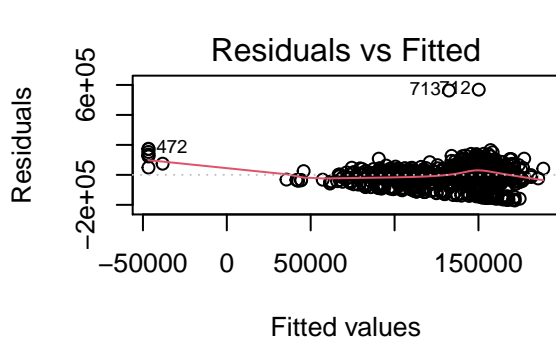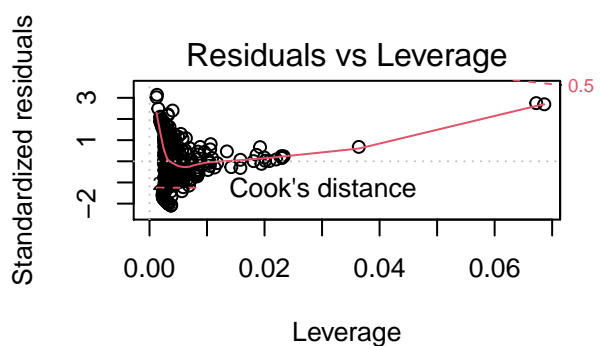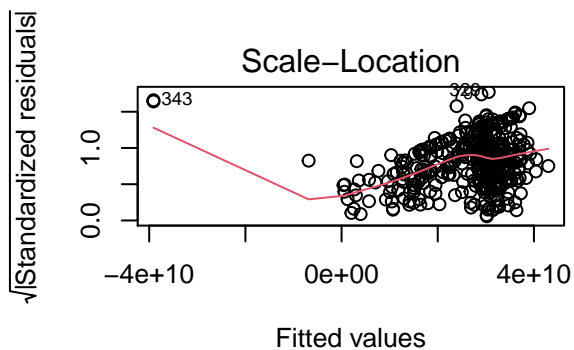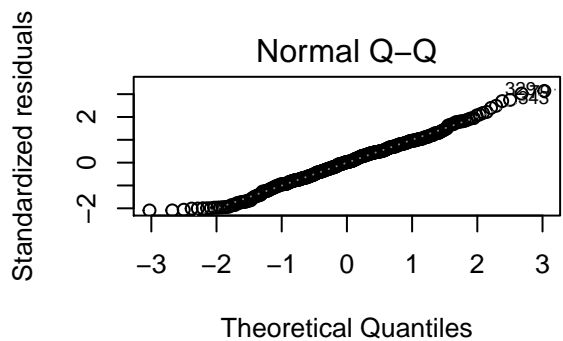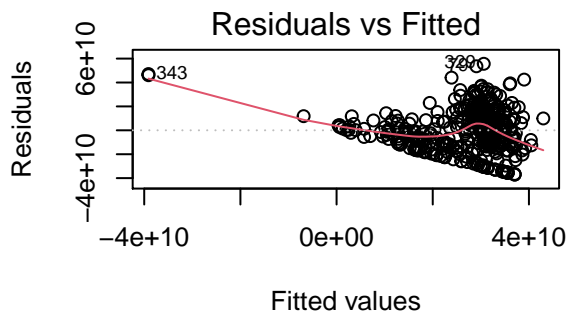
## Residuals vs Fitted
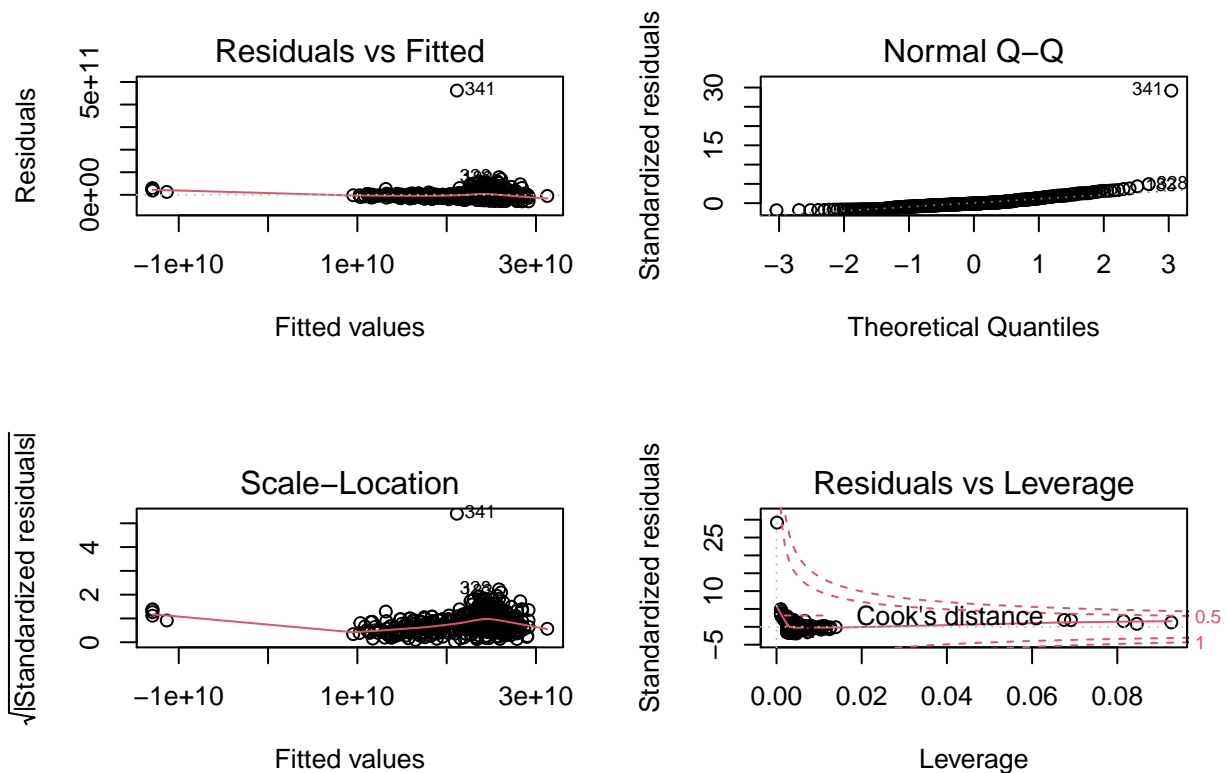
## Normal Q–Q

## Scale–Location

## Residuals vs Leverage

```
fit2 <- rlm(g2$votes ~ g2$disb)
par(mfrow=c(2,2))
plot (fit1)
```



## Residuals vs Fitted

## Normal Q–Q

## Scale–Location

## Residuals vs Leverage

31

```
fit3 <- rlm(g3$votes ~ g3$disb)
par(mfrow=c(2,2))
plot (fit2)
```



```
summary(fit0)
```

```
##
## Call: rlm(formula = g0$votes ~ g0$disb)
## Residuals:
##      Min       1Q    Median        3Q       Max
## -3.36130 -0.12353  0.03335   0.13016   0.88942
##
## Coefficients:
##              Value    Std. Error t value
## (Intercept)  4.2162   0.0414     101.8323
## g0$disb      0.1609   0.0073      21.9398
##
## Residual standard error: 0.1917 on 878 degrees of freedom
```

```
summary(fit1)
```

```
##
## Call: rlm(formula = g1$votes ~ g1$disb)
## Residuals:
##        Min         1Q      Median          3Q          Max
## -3.704e+10  -1.201e+10   4.285e+08   1.181e+10   5.564e+10
##
```

32

```
## Coefficients:
##              Value        Std. Error    t value
## (Intercept) -3.911640e+10  6.067250e+09 -6.447100e+00
## g1$disb      5.002444e+09  4.558440e+08  1.097400e+01
##
## Residual standard error: 1.772e+10 on 404 degrees of freedom
```

`summary(fit2)`

```
##
## Call: rlm(formula = g2$votes ~ g2$disb)
## Residuals:
##        Min         1Q      Median         3Q        Max
## -2.893e+10 -9.863e+09 -2.154e+09  1.217e+10  4.620e+11
##
## Coefficients:
##              Value        Std. Error    t value
## (Intercept) -1.294200e+10  5.318627e+09 -2.433300e+00
## g2$disb      2.728824e+09  4.089823e+08  6.672200e+00
##
## Residual standard error: 1.583e+10 on 421 degrees of freedom
```

`summary(fit3)`

```
##
## Call: rlm(formula = g3$votes ~ g3$disb)
## Residuals:
##     Min     1Q Median     3Q     Max
## -30841 -10191  -2653  11235 682215
##
## Coefficients:
##              Value      Std. Error  t value
## (Intercept) -15663.5931   8732.4093   -1.7937
## g3$disb       3790.2361    858.8929    4.4129
##
## Residual standard error: 15330 on 49 degrees of freedom
```

6. (3 points) Interpret the model coefficients you estimate.

- Tasks to keep in mind as you're writing about your model:
    - At the time that you're writing and interpreting your regression coefficients you'll be *deep* in the analysis. Nobody will know more about the data than you do, at that point. *So, although it will feel tedious, be descriptive and thorough in describing your observations.*
    - It can be hard to strike the balance between: on the one hand, writing enough of the technical underpinnings to know that your model meets the assumptions that it must; and, on the other hand, writing little enough about the model assumptions that the implications of the model can still be clear. We're starting this practice now, so that by the end of Lab 2 you will have had several chances to strike this balance.

`#lm(d2$general_votes ~ b1*d2$ttl_disb + b2)`