

# Data Exploration: Gender and World View

Kayla Huang

October 14, 2021

In this Data Exploration assignment, you will work with data that has been modified from the Barnhart et al. (2020) article. You will investigate whether certain types of countries are more likely to initiate conflicts with other countries. Note that you are only working with the data used to generate the *Monadic Findings* of the paper - that is, you will examine whether democracies initiate fewer conflicts than autocracies.

If you have a question about any part of this assignment, please ask! Note that the actionable part of each question is **bolded**.

## The Suffragist Peace

### Data Details:

- File Name: `suffrage_data.csv`
- Source: These data are from Barnhart et al. (2020).

Variable Name	Variable Description
<code>ccode1</code>	Unique country code
<code>country_name</code>	Country name
<code>year</code>	Year
<code>init</code>	The number of overall conflicts initiated by the country specified by <code>ccode1</code> during the year specified by <code>year</code>
<code>init_autoc</code>	The number of overall conflicts initiated by the country specified by <code>ccode1</code> <b>with autocracies</b> during the year specified by <code>year</code>
<code>init_democ</code>	The number of overall conflicts initiated by the country specified by <code>ccode1</code> <b>with democracies</b> during the year specified by <code>year</code>
<code>democracynosuff</code>	Indicator variable for a democracy without women's suffrage. 1 if the country is a democracy without women's suffrage, 0 otherwise.
<code>suffrage</code>	Indicator variable for a country with women's suffrage
<code>autocracy</code>	Indicator variable for a country with an autocratic government
<code>nuclear</code>	Indicator variable for whether the country is a nuclear power
<code>wcivillibs</code>	Measure of the degree of civil liberty women enjoy, ranging from 0-1, where higher values mean women have more civil liberties
<code>polity</code>	Polity score for the country specified by <code>ccode1</code> during the year specified by <code>year</code>

## Question 1

### Part a

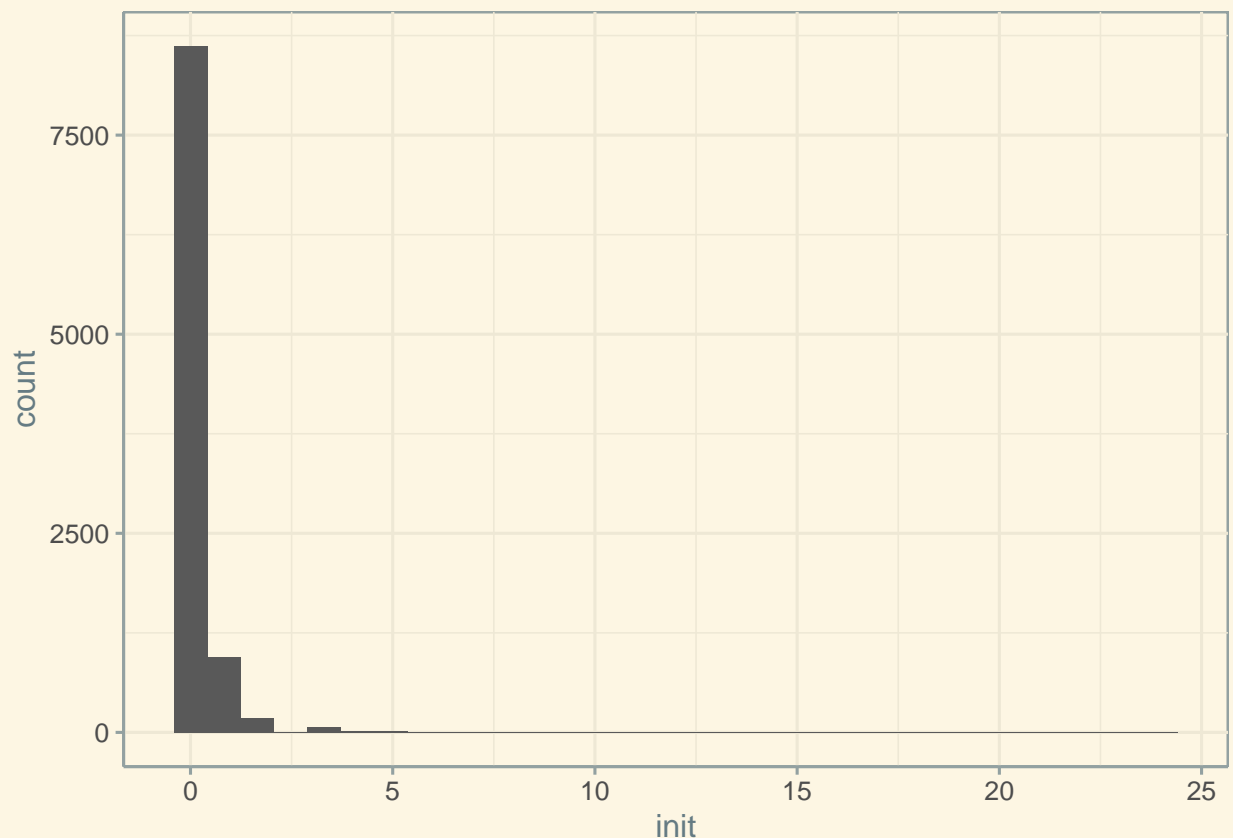
Before getting started, it is a good idea to take a look at the structure of the data. This data set is different from what we've seen so far. Until now, all the data we've looked at has had the individual as the unit of observation. This means that each row of the data corresponds to a single individual, and the columns correspond to some characteristics of that individual, like their responses to a survey. When working with data, it is important to understand the unit of observation, along with other characteristics of the data. The unit of observation is the object about which data is collected. That could be, say, an individual, a country, a football game, or an episode of TV. **Take a look at the data to determine the unit of observation. Note that the structure isn't exactly the same as the data used in Barnhart et al. (2020).**

### Part b

Is war rare or common? Make a histogram of the main dependent variable, `init`. Comment on what you see, being sure to keep the unit of observation and the definition of the `init` variable in mind. Is what you see surprising? What does it say about the frequency of initiating conflict?

```
library(ggthemes)
s_data %>%
  ggplot(aes(x=init)) + geom_histogram() + theme_solarized() + scale_fill_solarized(accent = "blue")

## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



## Question 2

How were the autocracy and suffrage variables defined? Can autocracies also have women's suffrage (at least in this coding scheme)? What reasons do the authors of the paper give for these coding decisions and how do you think it might affect their findings? (Hint: take a close look at pages 651 and 652 of the original article.)

These variables were defined as indicator variables. They are either 0 or 1 depending on whether or not the country has or does not have women's suffrage or an autocratic government, respectively. The way Barnard defines the existence of suffrage is as follows: "the variable is coded 1 if suffrage for national elections has been extended to women in a state and the state's Polity score is 1 or higher." Otherwise, the variable is coded as 0. For the autocracy variable: it "utilizes the standard coding of autocracy and is coded 1 if the state's Polity score is 5 or below" and is coded 0 otherwise.

Yes, autocracies can also have women's suffrage based on these two definitions (i.e. row 241, Cuba).

## Question 3

The democratic peace - i.e. the propensity for democracies to avoid conflict with each other, and to avoid conflict more generally - is an empirical regularity. The theory, as originally posed, is not gendered. Do the data support the democratic peace theory? **Ignoring suffrage status for now, do the data suggest that modern democracies initiate fewer conflicts than autocracies? Do democracies tend to initiate conflict more with autocracies or other democracies?**

```
dems <- filter(s_data, autocracy == 0)
auts <- filter(s_data, autocracy == 1)

print("Democracy initiated wars:")
```

```
## [1] "Democracy initiated wars:"
```

```
sum(dems$init)
```

```
## [1] 530
```

```
print("Autocracy initiated wars:")
```

```
## [1] "Autocracy initiated wars:"
```

```
sum(auts$init)
```

```
## [1] 1286
```

```
print("Democracy initiating with Autocracy")
```

```
## [1] "Democracy initiating with Autocracy"
```

```
mean(dems$init_autoc)
```

```
## [1] 0.07865169
```

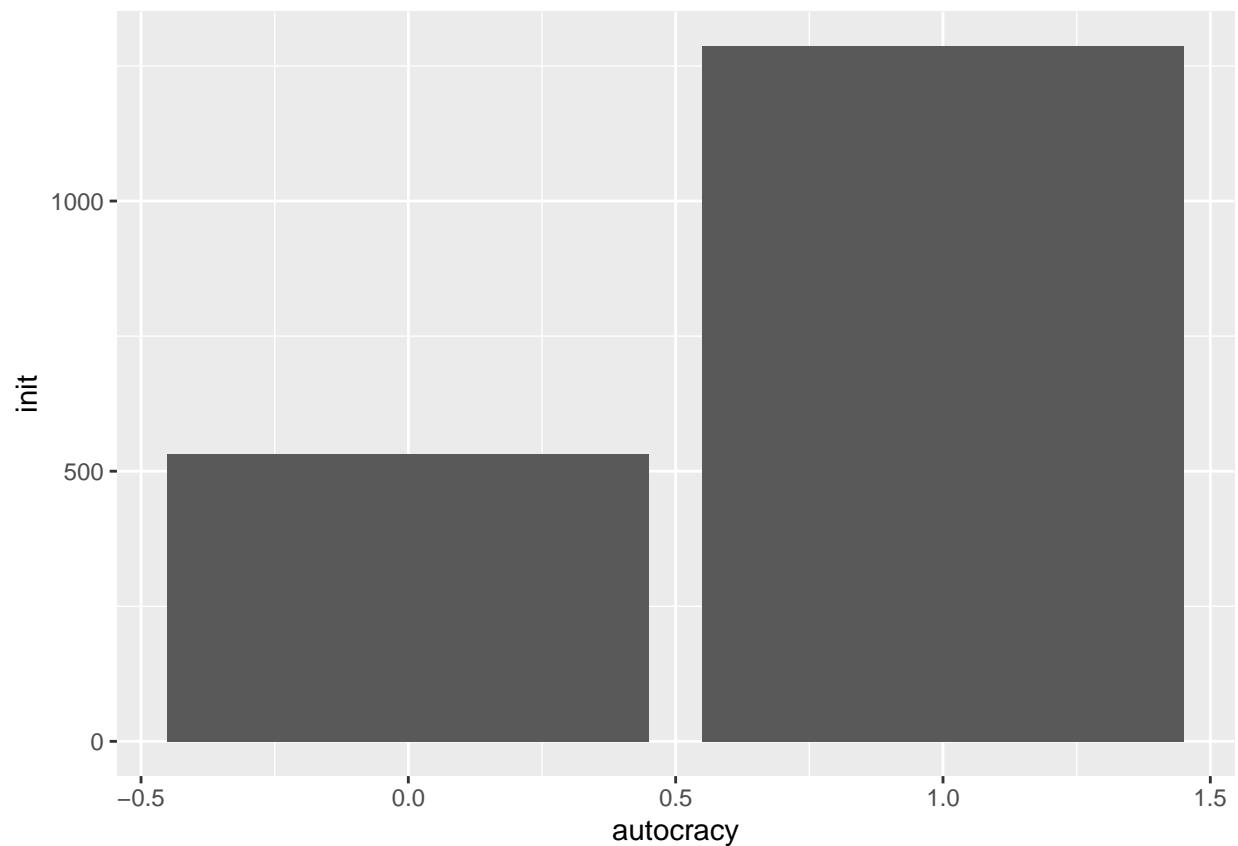
```
print("Democracy initating with Democracy")
```

```
## [1] "Democracy initating with Democracy"
```

```
mean(dems$init_democ)
```

```
## [1] 0.04016854
```

```
s_data %>%  
  group_by(autocracy) %>%  
  ggplot(aes(x=autocracy, y=init)) + geom_bar(stat = "identity")
```



#### Question 4

Now that we've taken a look at the classic democratic peace theory, let's take an initial look at how women's suffrage is related to initiating conflict. **Conduct a bivariate regression, modeling the number of conflicts initiated with women's suffrage (i.e.  $\text{init} \sim \text{suffrage}$ ).** This will help inform you about how the number of conflict initiated in a year depends on women's suffrage. Report the coefficient on suffrage. Interpret your results. If you like, extend the problem by reporting the 95% confidence interval for the suffrage coefficient. Is the relationship statistically significant?

The `lm()` function is used to calculate regressions in R. [Here](#) is a guide to linear regression in R that may be helpful.

```
model <- lm(formula = init ~ suffrage, data = s_data)
summary(model)
```

```
##
## Call:
## lm(formula = init ~ suffrage, data = s_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.2045 -0.2045 -0.2045 -0.1538  23.7955
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.204515   0.008588  23.815 < 2e-16 ***
## suffrage    -0.050727   0.013532  -3.749 0.000179 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6592 on 9863 degrees of freedom
## Multiple R-squared:  0.001423,    Adjusted R-squared:  0.001321
## F-statistic: 14.05 on 1 and 9863 DF,  p-value: 0.0001788
```

Yes, the relationship is statistically significant, as seen by the p-value above!

## Question 5

The model in the previous question was very simple; we modeled initiation only as a function of suffrage. In reality, the relationship is probably more complicated - conflict initiation probably depends on more than just women's suffrage. **Look at the other variables available in the data and find one or more that you think may also be related to conflict initiation. Explain why you think so, then add the variable(s) to the right side of the regression (as explanatory variables) in question 4. Interpret what you find.**

```
model_2 <- lm(formula = init ~ suffrage + nuclear + autocracy, data = s_data)
summary(model_2)
```

```
##
## Call:
## lm(formula = init ~ suffrage + nuclear + autocracy, data = s_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.8931 -0.2010 -0.2010 -0.0969  23.7990
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.14263    0.02275   6.270 3.75e-10 ***
## suffrage    -0.04569    0.02175  -2.101 0.03565 *
## nuclear      0.69205    0.03730  18.555 < 2e-16 ***
## autocracy    0.05841    0.02220   2.632 0.00851 **
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6479 on 9861 degrees of freedom
## Multiple R-squared:  0.03538,    Adjusted R-squared:  0.03509
## F-statistic: 120.6 on 3 and 9861 DF,  p-value: < 2.2e-16
```

## Question 6: Data Science Question

Estimate a regression of the following form:  $\text{init} \sim \text{suffrage} + \text{polity} + \text{polity} * \text{suffrage}$ , where  $\text{polity} * \text{suffrage}$  is the interaction between polity score and women's suffrage. Compare this to the same model but without the interaction term. Interpret your results.

```
model_3 <- lm(formula = init ~ suffrage + polity + polity*suffrage, data = s_data)
summary(model_3)
```

```
##
## Call:
## lm(formula = init ~ suffrage + polity + polity * suffrage, data = s_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.2441 -0.2093 -0.1923 -0.1275  23.7861
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.193154   0.012553   15.387 <2e-16 ***
## suffrage        0.063886   0.038459    1.661  0.0967 .
## polity         -0.002309   0.001861   -1.240  0.2149
## suffrage:polity -0.010643   0.004748   -2.242  0.0250 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6589 on 9861 degrees of freedom
## Multiple R-squared:  0.002468,    Adjusted R-squared:  0.002165
## F-statistic: 8.132 on 3 and 9861 DF,  p-value: 2.093e-05
```

```
model_3 <- lm(formula = init ~ suffrage + polity, data = s_data)
summary(model_3)
```

```
##
## Call:
## lm(formula = init ~ suffrage + polity, data = s_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.2245 -0.2127 -0.1733 -0.1458  23.7794
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.1851035   0.0120309   15.386 <2e-16 ***
## suffrage        0.0001303   0.0258959    0.005  0.9960
## polity         -0.0039446   0.0017127   -2.303  0.0213 *
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.659 on 9862 degrees of freedom
## Multiple R-squared:  0.00196,    Adjusted R-squared:  0.001757
## F-statistic: 9.682 on 2 and 9862 DF,  p-value: 6.301e-05
```

In the social sciences, we use interaction terms in regressions to capture heterogeneous effects. As an example of how to implement and interpret this type of model, suppose we wanted to understand the relationship between education on the one hand (as the outcome variable), and age and gender on the other hand (as explanatory variables). We might think that the effect of age on education depends on whether you're talking about men or women. Maybe for men, age has no effect on education, but for women, there is a negative effect, as older women were discouraged or barred from seeking higher education. To assess whether this is true, we can use an interaction between gender and age. You can model this in R using this formula in the `lm()` function: `education ~ age + female + age*female` (supposing gender is coded into a binary variable `female`). Here, `age*female` is what creates the interaction.

Lets say that we ran this regression in R and found that the model looks like this:  $education = 1.5 + .005 * age + .01 * female + -.4 * age * female$ . Here, the coefficient on `age` is .005, .01 on `female`, and -.4 on the interaction between the two. Without an interaction, to interpret the coefficient on `age`, we would say the effect of `age` on `education` is .005. However, the interaction term modifies that relationship - the effect of `age` on `education` now depends on gender.

To see this, we must plug in values for `female` and `age`. When `female = 0`, then the interaction term vanishes, and then the effect of `age` on `education` is .005. In other words, for non-women, there is a very small relationship between age and education. Now plugging in `female = 1`, the effect of `age` on `education` becomes .005 (the coefficient on age) + -.4 (the coefficient on the interaction) = -.395. In other words, the effect of age on education among women is negative.

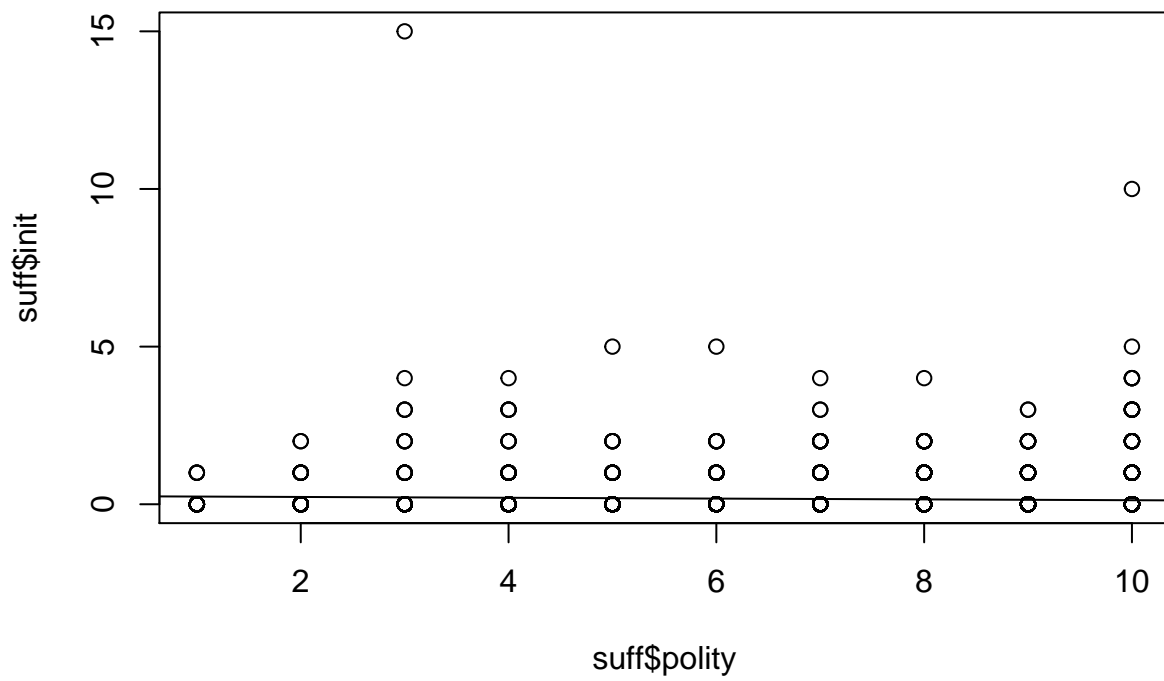
Interpreting the effect of gender is a bit more complicated, and in this case is nonsensical. To do so, we set `age = 0` (which doesn't make a ton of sense) to find that the effect of `female` on `education` is .01 when `age = 0`. Always pay attention to whether the coefficients you're focusing on are even substantively meaningful.

## Question 7: Data Science Question

When using regression, especially with interactions, sometimes it is useful to visualize the results. **Create two plots of the predicted number of conflicts per year on the y-axis and Polity score on the x-axis (among countries with a Polity score greater than or equal to one only), split by suffrage.** That is, one plot should plot the predicted number of conflict per year among suffrage democracies, and the other among non-suffrage democracies. This way you will be able to visualize the interaction between suffrage and Polity score that we saw in the previous question. [This guide](#) may be helpful in doing so - it uses a different type of regression model (binary logit), but the principle of prediction is the same. Make sure to hold the suffrage variable at 0 or 1. Comment on what you find.

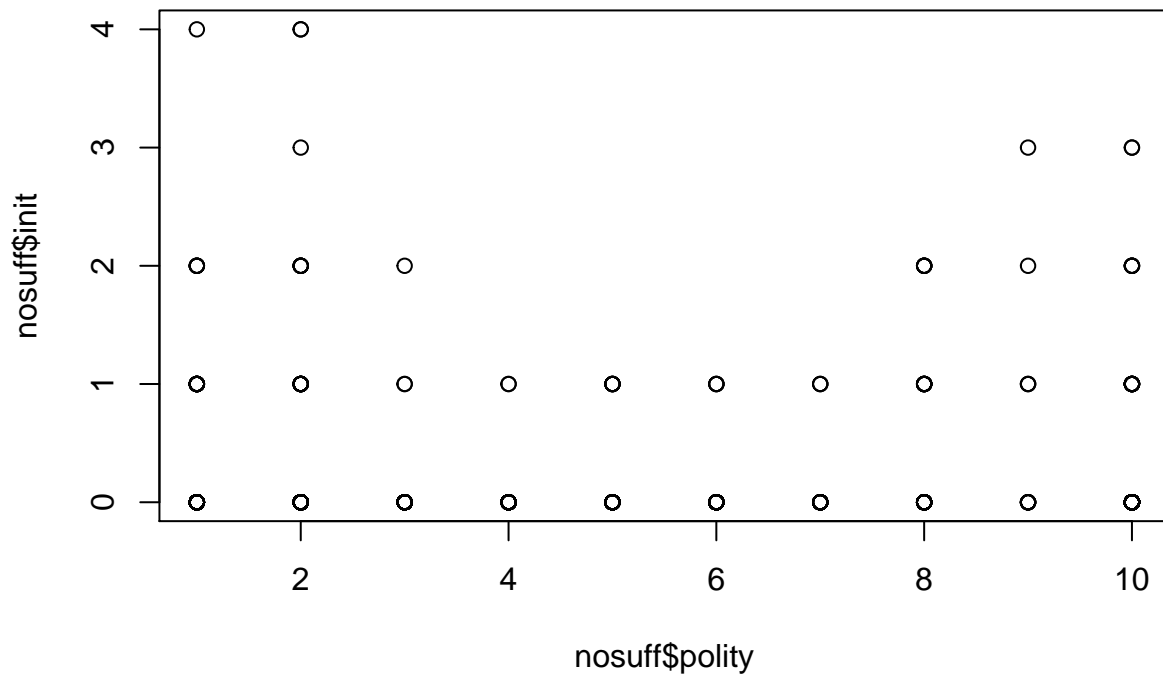
```
suff <- filter(s_data, suffrage == 1 & polity >= 1)
nosuff <- filter(s_data, suffrage == 0 & polity >= 1)

lm_suff <- lm(init ~ polity, data = suff)
plot(suff$polity, suff$init)
abline(lm_suff)
```



```
plot(nosuff$polity, nosuff$init)
```

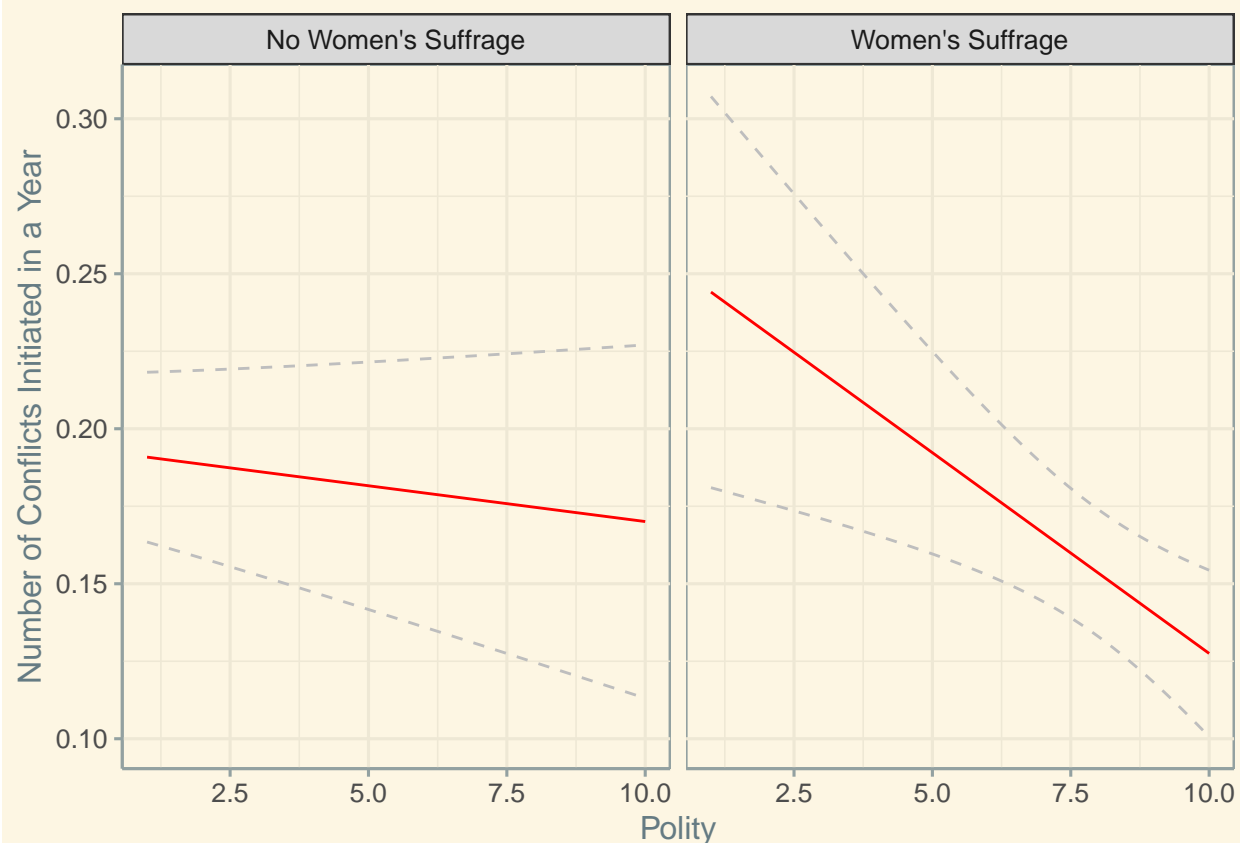




```
m1 <- lm(init ~ polity*suffrage, data = s_data)
pred_data1 <- data.frame(polity = seq(1, 10, 0.01), suffrage = 0)
pred_data2 <- data.frame(polity = seq(1, 10, 0.01), suffrage = 1)
pred_data <- bind_rows(pred_data1, pred_data2)

preds <- predict(m1, pred_data, se.fit = T)
preds_df <- data.frame(fit = preds$fit, se = preds$se.fit, polity = pred_data$polity, suffrage = pred_data$suffrage)
```

```
preds_df %>%
  mutate(suffrage = recode(suffrage, '0' = "No Women's Suffrage", '1' = "Women's Suffrage")) %>%
  ggplot(aes(x=polity)) +
  geom_line(aes(y=fit), col = 'red') +
  geom_line(aes(y=fit +1.96*se), col="gray", lty="dashed") +
  geom_line(aes(y=fit -1.96*se), col="gray", lty="dashed") +
  facet_wrap(~ suffrage) +
  labs(x="Polity", y="Number of Conflicts Initiated in a Year") +
  theme_solarized() + scale_fill_solarized(accent = "blue")
```



## Question 8

One of the advantages of the data we have is that we can plot trends over time. **Group countries into those with and without suffrage and plot the average number of disputes initiated by those countries in each year covered by the data. Comment on what you find.**

## Question 9

With geographical data like we are working with, you may want to make a map. For example, you may want to be able to visualize which countries had suffrage and which did not in a given year. As an example of how to create maps using the `ggplot2` and `sf` packages, below is a map of suffrage in North America in 1960. **Use the `map` data to plot a map of the number of conflicts initiated by each country in the Americas (that is, countries in North and South America), in 1960. You can modify the example code given below.**

Note that if you encounter a country that is missing from your map, you should check how the country name is spelled in each of the data sets (`world` and `s_data`). We merge the two data sets together to allow for mapping based on country name, so if they country names don't match exactly then the merge will return `NA` and the mapping will fail. To see an example of how to change the country names when need be, see below. For example, in the `world` data set, the US is called "United States", whereas in the `s_data` data set, it is called "United States of America".

```
# if you don't have these libraries already, download them using install.packages()
library(sf) # this is for plotting maps in ggplot
```

```
## Linking to GEOS 3.8.1, GDAL 3.2.1, PROJ 7.2.1
```

```
library(spData) # this is for the 'world' data set
```

```
## To access larger datasets in this package, install the spDataLarge  
## package with: 'install.packages('spDataLarge',  
## repos='https://nowosad.github.io/drat/', type='source')'
```

```
# edit this to change country names when the two data sets don't exactly match
```

```
s_data <- s_data %>%
```

```
  mutate(country_name = case_when(  
    country_name == "United States of America" ~ "United States",  
    country_name == "Russia" ~ "Russian Federation",  
    T ~ country_name  
  ))
```

```
map <- left_join(s_data, world[, c("continent", "geom", "name_long")], by = c("country_name" = "name_lo
```

```
map %>% filter(year == 1960, continent == "North America") %>%
```

```
  ggplot() +  
  geom_sf(aes(fill = as.factor(suffrage), geometry = geom)) +  
  labs(fill = "Suffrage", title = "Women's Suffrage in North America, 1950")
```

Women's Suffrage in North America, 1950

