

Data Exploration: Gender and World View

Yao Yu

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> with autocracies 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> with democracies 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).

```
summary(s_data)
```

```
##           X1           ccode1      country_name           year
## Min.      : 1      Min.      : 2.0      Length:9865      Min.      :1900
## 1st Qu.:2467      1st Qu.:200.0      Class :character      1st Qu.:1949
## Median :4933      Median :390.0      Mode  :character      Median :1973
## Mean    :4933      Mean    :415.6                      Mean    :1968
## 3rd Qu.:7399      3rd Qu.:640.0                      3rd Qu.:1991
## Max.    :9865      Max.    :950.0                      Max.    :2007
##
##           init           init_autoc      init_democ      democracynosuff
## Min.      : 0.0000      Min.      :0.0000      Min.      : 0.00000      Min.      :0.0000
## 1st Qu.: 0.0000      1st Qu.:0.0000      1st Qu.: 0.00000      1st Qu.:0.0000
## Median : 0.0000      Median :0.0000      Median : 0.00000      Median :0.0000
## Mean    : 0.1841      Mean    :0.1001      Mean    : 0.05646      Mean    :0.0294
## 3rd Qu.: 0.0000      3rd Qu.:0.0000      3rd Qu.: 0.00000      3rd Qu.:0.0000
## Max.    :24.0000      Max.    :9.0000      Max.    :14.00000      Max.    :1.0000
##
##           suffrage           autocracy           polity           wcivillibs
## Min.      :0.0000      Min.      :0.0000      Min.      : -10.0000      Min.      :0.0010
## 1st Qu.:0.0000      1st Qu.:0.0000      1st Qu.: -7.0000      1st Qu.:0.3220
## Median :0.0000      Median :1.0000      Median : -1.0000      Median :0.5650
## Mean    :0.4027      Mean    :0.6391      Mean    : 0.2715      Mean    :0.5536
## 3rd Qu.:1.0000      3rd Qu.:1.0000      3rd Qu.: 8.0000      3rd Qu.:0.8150
## Max.    :1.0000      Max.    :1.0000      Max.    :10.0000      Max.    :0.9810
##                                     NA's      :642
##
##           nuclear
## Min.      :0.00000
## 1st Qu.:0.00000
## Median :0.00000
## Mean    :0.03254
## 3rd Qu.:0.00000
## Max.    :1.00000
##
```

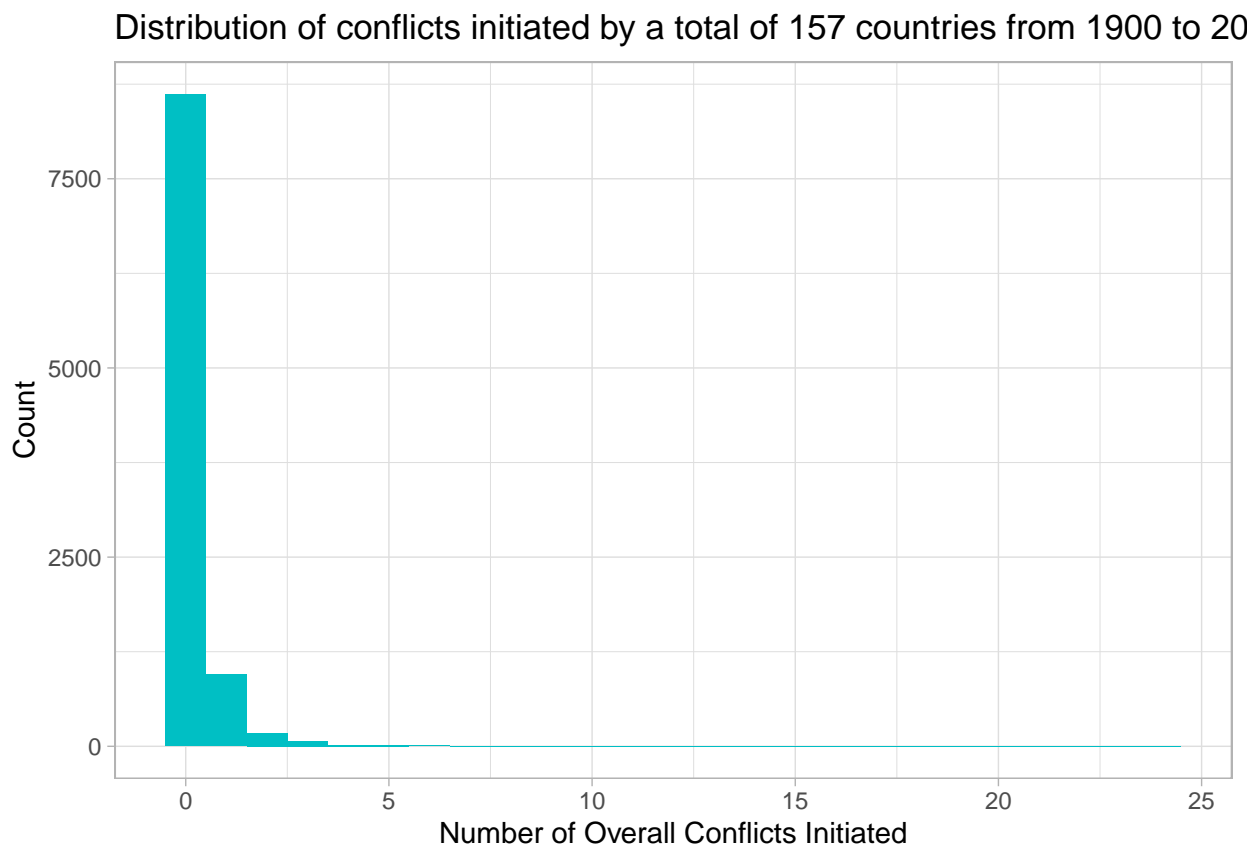
The unit of observation would be different countries and which year the data was collected.

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?

```
init_distribution_plot <- s_data %>%
  ggplot(aes(x = init)) +
  geom_histogram(binwidth = 1, fill = "#00bfc4") +
  theme_light() +
  labs(
    title = "Distribution of conflicts initiated by a total of 157 countries from 1900 to 2007",
    x = "Number of Overall Conflicts Initiated",
    y = "Count"
  )

init_distribution_plot
```



```
# png("init_distribution_plot.png", units="in", width=8, height=5, res=300)
# print(init_distribution_plot)
# dev.off()
```

From the histogram, we see that the distribution in number of conflicts initiated is heavily right-skewed. This aligns with the results I was expecting because most countries are not constantly initialing conflicts with other nations every year (that would be concerning). So, the frequency of initiating conflict is fairly low.

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.)

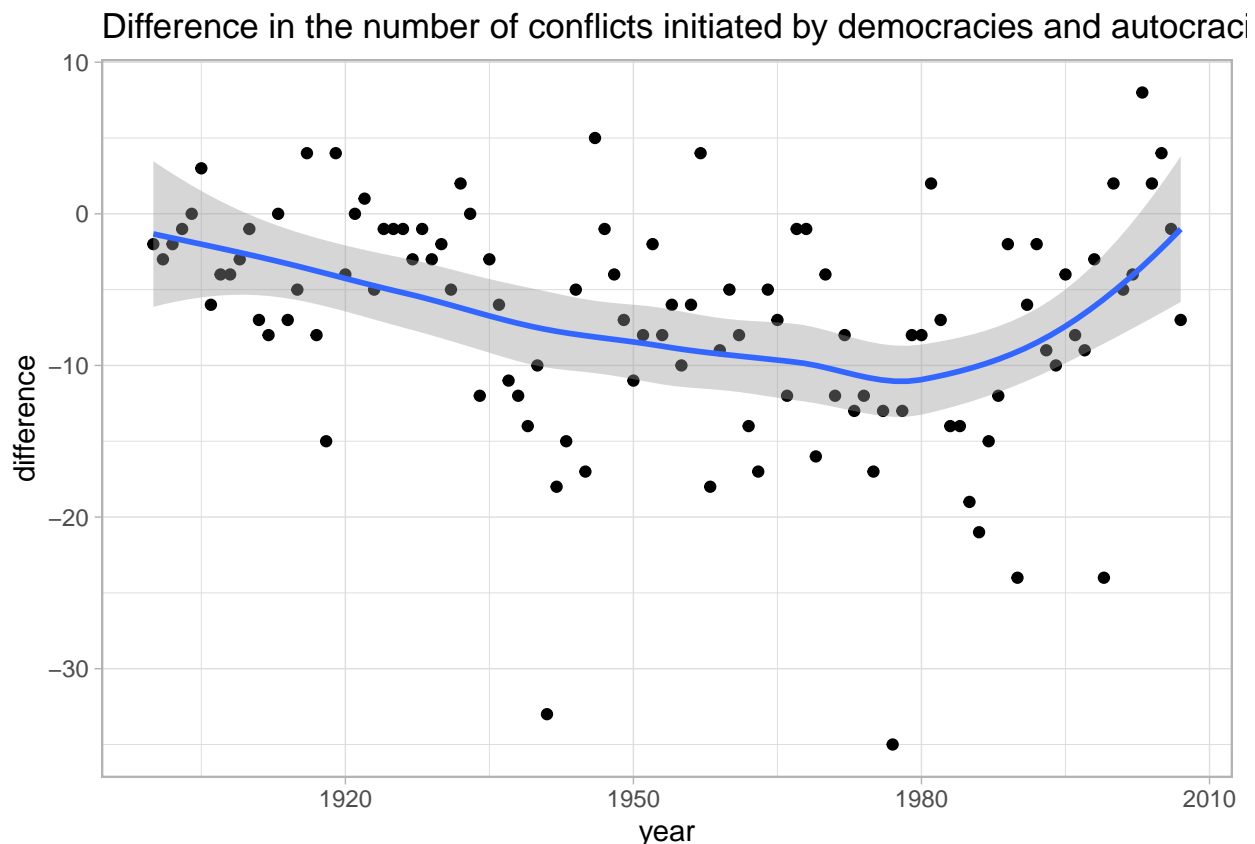
suffrage in this dataset is coded based on two components. First, if any women are able to vote in national elections. Second, if state's Polity score is 1 or higher. Otherwise, **suffrage** is coded as a 0. **autocracy** is coded as a 1 if a state's Polity score is 5 or lower. Otherwise, it is coded a 0.

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?

```
# Plotting difference
init_diff_plot <- s_data %>%
  group_by(autocracy, year) %>%
  summarize(total_init = sum(init),
            .groups = "drop") %>%
  group_by(year) %>%
  pivot_wider(names_from = autocracy, values_from = total_init) %>%
  mutate(difference = `0` - `1`) %>%
  ggplot(aes(x = year, y = difference)) +
  geom_point() +
  geom_smooth(method = "loess", formula = "y ~ x") +
  theme_light() +
  labs(
    title = "Difference in the number of conflicts initiated by democracies and autocracies"
  )
```

init_diff_plot



```
# png("init_diff_plot.png", units="in", width=8, height=5, res=300)
# print(init_diff_plot)
```

```

# dev.off()

t.test(init ~ autocracy, data = s_data)

##
## Welch Two Sample t-test
##
## data:  init by autocracy
## t = -4.4086, df = 9501.9, p-value = 1.052e-05
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
##  -0.07958308 -0.03059433
## sample estimates:
## mean in group 0 mean in group 1
##      0.1488764      0.2039651

s_data %>%
  group_by(autocracy) %>%
  summarize(total_init_autoc = sum(init_autoc),
            total_init_democ = sum(init_democ)) %>%
  mutate(autocracy = case_when(
    autocracy == 0 ~ "Democracy",
    autocracy == 1 ~ "Autocracy",
  )) %>%
  rename(
    num_init_autocracy = total_init_autoc,
    num_init_democracy = total_init_democ
  ) %>%
  gt() %>%
  tab_header(title = "Who Initiates Conflict with who?") %>%
  cols_label(
    autocracy = "Type",
    num_init_autocracy = "Autocracy",
    num_init_democracy = "Democracy"
  )

```

Who Initiates Conflict with who?

Type	Autocracy	Democracy
Democracy	280	143
Autocracy	707	414

Based on the plot, we can see that the data does suggest that modern democracies initiate fewer conflicts than autocracies. This can also be confirmed with a t-test showing that this difference is statistically significant at the 95% confidence level. Democracies on average initiate between 0.03 and 0.08 fewer conflicts than autocracies at the 95% confidence level. From the table, we can see that both democracies and autocracies tend to initiate conflict more with other autocracies than other democracies.

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.

```
fit_1 <- lm(init ~ suffrage, data = s_data)

stargazer(fit_1, header = FALSE,
  dep.var.labels = c("Number of Conflicts Initiated"),
  covariate.labels = c("Suffrage", "Constant"),
  title = "Number of Conflicts Initiated as a Function of Suffrage")
```

Number of Conflicts Initiated as a Function of Suffrage

<i>Dependent variable:</i>	
Number of Conflicts Initiated	
Suffrage	-0.051*** (0.014)
Constant	0.205*** (0.009)
Observations	9,865
R ²	0.001
Adjusted R ²	0.001
Residual Std. Error	0.659 (df = 9863)
F Statistic	14.052*** (df = 1; 9863)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01	

```
# Confidence Interval
print(confint(fit_1, 'suffrage', level=0.95))
```

```
2.5 %      97.5 %
suffrage -0.07725183 -0.02420122
```

From the regression model, we can see that countries with women's suffrage on average initiation on average 0.05 fewer conflicts than countries without women's suffrage. This finding is statistically significant with a p-value of about 0.002 and a confidence interval of -0.077 and -0.024.

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.**

```
fit_2 <- lm(init ~ nuclear + suffrage*autocracy, data = s_data)

stargazer(fit_2, header = FALSE,
  dep.var.labels = c("Number of Conflicts Initiated"),
  covariate.labels = c("Nuclear", "Suffrage", "Autocracy", "Suffrage:Autocracy"),
  title = "Number of Conflicts Initiated as a Function of Suffrage and Other Variables of Interest")
```

Number of Conflicts Initiated as a Function of Suffrage and Other Variables of Interest

	Dependent variable:
	Number of Conflicts Initiated
Nuclear	0.697*** (0.037)
Suffrage	-0.139*** (0.040)
Autocracy	-0.031 (0.039)
Suffrage:Autocracy	0.132*** (0.047)
Constant	0.228*** (0.038)
Observations	9,865
R ²	0.036
Adjusted R ²	0.036
Residual Std. Error	0.648 (df = 9860)
F Statistic	92.435*** (df = 4; 9860)
Note:	*p<0.1; **p<0.05; ***p<0.01

From this model, we can see that if a country is a nuclear power, then they are on average estimated to have 0.697 more conflicts than countries that are not nuclear powers. This finding is statistically significant even while controlling for suffrage, autocracy, and the interaction between women's suffrage and autocracy.

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.

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: $\text{education} = 1.5 + .005 * \text{age} + .01 * \text{female} + -.4 * \text{age} * \text{female}$. 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.

```
fit_3_no_interaction <- lm(init ~ polity + suffrage, data = s_data)
fit_3_interaction <- lm(init ~ polity*suffrage, data = s_data)

stargazer(fit_3_no_interaction, fit_3_interaction, header = FALSE,
  dep.var.labels = c("Number of Conflicts Initiated"),
  covariate.labels = c("Polity", "Suffrage", "Polity:Suffrage"),
  title = "Number of Conflicts Initiated as a Function of Suffrage and Other Variables of Interest")
```

Comparing these two models, we can see that without considering the interaction between polity score and women's suffrage, an increase in polity score by 1 (more democratic) is associated on average with a decrease in the number of conflicts initiated by about 0.004. However, a country with women's suffrage is associated on average with an increase in the number of conflicts initiated by about 0.0001.

However, by considering the interaction between polity score and women's suffrage, we see a big shift in that a country with women's suffrage is associated on average with an increase in the number of conflicts initiated by about 0.064. However, an increase in polity score by 1 in a country with women's suffrage will be associated with on average a decrease in the number of conflicts initiated by about 0.011 compared to an increase in polity score by 1 in a country without women's suffrage. This makes sense as a higher polity score (up to 10) makes a country more democratic and these countries tend to mostly all have women's suffrage.

Number of Conflicts Initiated as a Function of Suffrage and Other Variables of Interest

	<i>Dependent variable:</i>	
	Number of Conflicts Initiated	
	(1)	(2)
Polity	−0.004** (0.002)	−0.002 (0.002)
Suffrage	0.0001 (0.026)	0.064* (0.038)
Polity:Suffrage		−0.011** (0.005)
Constant	0.185*** (0.012)	0.193*** (0.013)
Observations	9,865	9,865
R ²	0.002	0.002
Adjusted R ²	0.002	0.002
Residual Std. Error	0.659 (df = 9862)	0.659 (df = 9861)
F Statistic	9.682*** (df = 2; 9862)	8.132*** (df = 3; 9861)

Note:

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

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.

```
fake_data <- tibble(
  suffrage = c(rep(0, 10), rep(1, 10)),
  polity = rep(seq(1, 10), 2)
)

fake_data_preds <- predict(fit_3_interaction, newdata = fake_data, se.fit = TRUE)

fake_data <- fake_data %>%
  mutate(preds = fake_data_preds$fit,
         se = fake_data_preds$se.fit,
         se_lower = preds - 1.96*se,
         se_upper = preds + 1.96*se)

suffrage_polity_interaction_plot <- fake_data %>%
  mutate(suffrage = case_when(
    suffrage == 0 ~ "No Women's Suffrage",
```

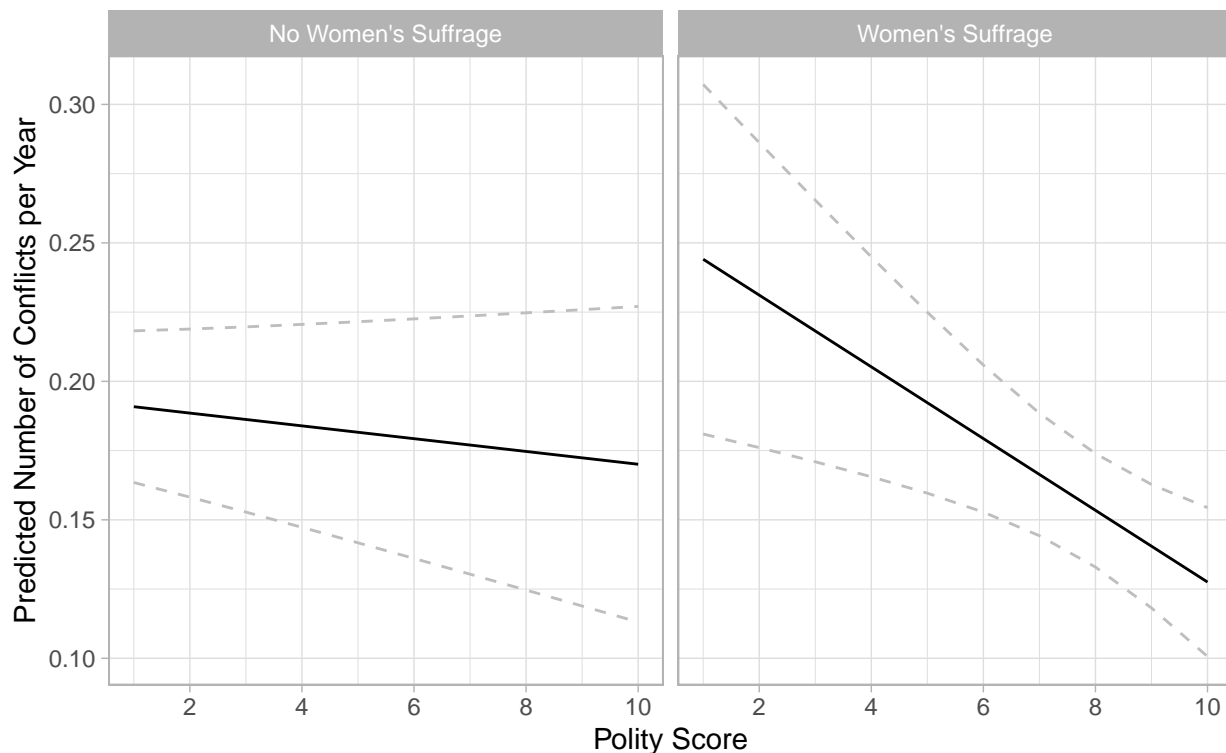
```

    TRUE ~ "Women's Suffrage"
  )) %>%
  ggplot(aes(x = polity, preds)) +
  geom_line() +
  # 95% confidence interval
  geom_line(aes(x = polity, se_lower), linetype = "dashed", color = "grey") +
  geom_line(aes(x = polity, se_upper), linetype = "dashed", color = "grey") +
  facet_wrap(~suffrage) +
  scale_x_continuous(breaks = seq(0, 10, by = 2)) +
  theme_light() +
  labs(
    title = "Predicted Number of Conflicts per Year by Polity Score",
    subtitle = "with the 95% confidence interval",
    x = "Polity Score",
    y = "Predicted Number of Conflicts per Year"
  )
)

```

suffrage_polity_interaction_plot

Predicted Number of Conflicts per Year by Polity Score
with the 95% confidence interval



```

png("suffrage_polity_interaction_plot.png", units="in", width=8, height=5, res=300)
print(suffrage_polity_interaction_plot)
dev.off()

```

```

## pdf
## 2

```

From this plot, we can see that as polity score increases, countries with women's suffrage tend to have much lower predicted number of conflicts on average than countries without women's suffrage. This visualizes the

trend we found in the previous question.

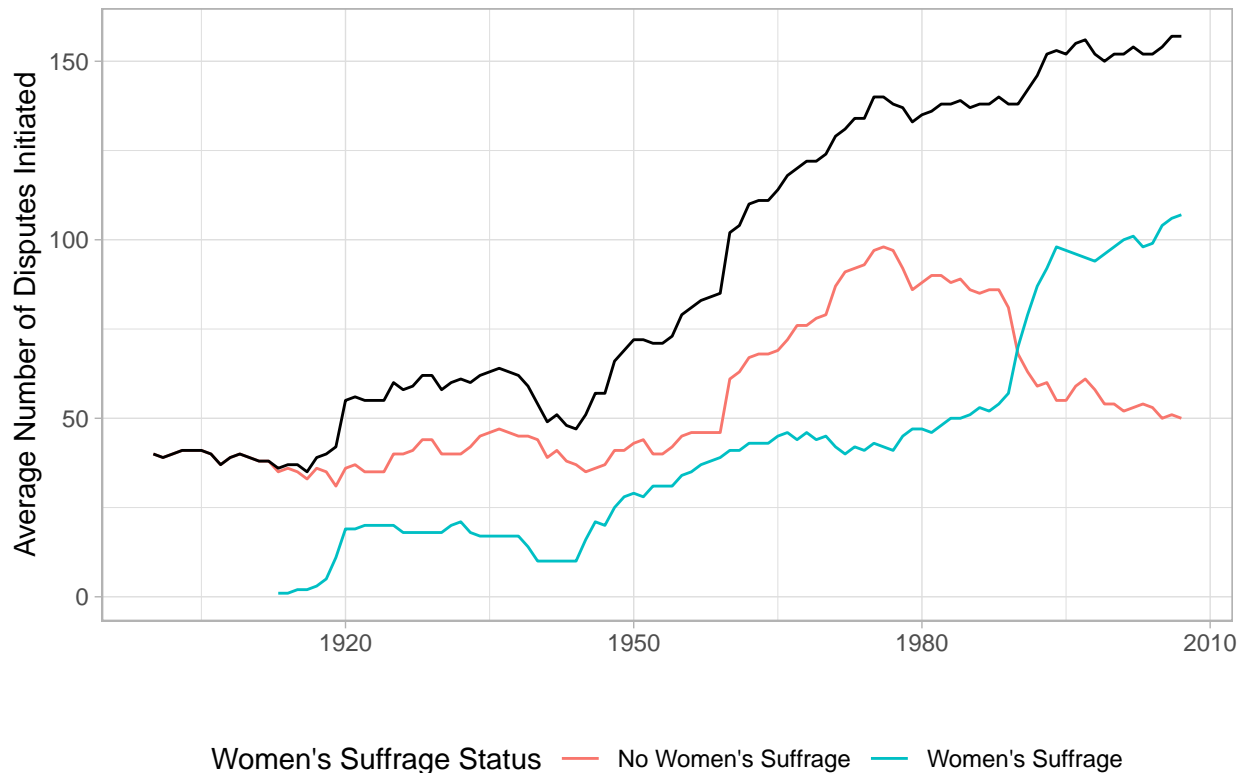
Question 8

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.**

```
# Number of countries with suffrage over time
suffrage_total_year <- s_data %>%
  count(suffrage, year) %>%
  mutate(
    suffrage = case_when(
      suffrage == "0" ~ "No Women's Suffrage",
      suffrage == "1" ~ "Women's Suffrage"
    )
  ) %>%
  group_by(year) %>%
  summarize(suffrage = suffrage,
            n = n,
            total_n = sum(n),
            .groups = "drop") %>%
  ggplot() +
  geom_line(aes(x = year, y = n, color = suffrage)) +
  geom_line(aes(x = year, y = total_n)) +
  theme_light() +
  theme(legend.position = "bottom") +
  labs(
    title = "Total Number of Countries with and without Women's Suffrage",
    x = "",
    y = "Average Number of Disputes Initiated",
    color = "Women's Suffrage Status"
  )

suffrage_total_year
```

Total Number of Countries with and without Women's Suffrage



```
# png("suffrage_total_year.png", units="in", width=8, height=5, res=300)
# print(suffrage_total_year)
# dev.off()
```

```
# Averages
init_suffrage_avg_plot <- s_data %>%
  mutate(suffrage = as.character(suffrage)) %>%
  group_by(suffrage, year) %>%
  summarize(avg_init = mean(init),
            .groups = "drop") %>%
  mutate(
    suffrage = case_when(
      suffrage == "0" ~ "No Women's Suffrage",
      suffrage == "1" ~ "Women's Suffrage"
    )
  ) %>%
  ggplot(aes(x = year, y = avg_init, color = suffrage)) +
  geom_point() +
  geom_smooth(method = "loess", formula = "y ~ x") +
  theme_light() +
  labs(
    title = "Average Number of Disputes Initiated by Countries",
    subtitle = "with and without Women's Suffrage",
    x = "",
    y = "Average Number of Disputes Initiated",
    color = "Women's Suffrage Status"
  )
```

```

# Totals
init_suffrage_total_plot <- s_data %>%
  mutate(suffrage = as.character(suffrage)) %>%
  group_by(suffrage, year) %>%
  summarize(total_init = sum(init),
            .groups = "drop") %>%
  mutate(
    suffrage = case_when(
      suffrage == "0" ~ "No Women's Suffrage",
      suffrage == "1" ~ "Women's Suffrage"
    )
  ) %>%
  ggplot(aes(x = year, y = total_init, color = suffrage)) +
  geom_point() +
  geom_smooth(method = "loess", formula = "y ~ x") +
  theme_light() +
  labs(
    title = "Total Number of Disputes Initiated by Countries",
    subtitle = "with and without Women's Suffrage",
    x = "",
    y = "Total Number of Disputes Initiated",
    color = "Women's Suffrage Status"
  )

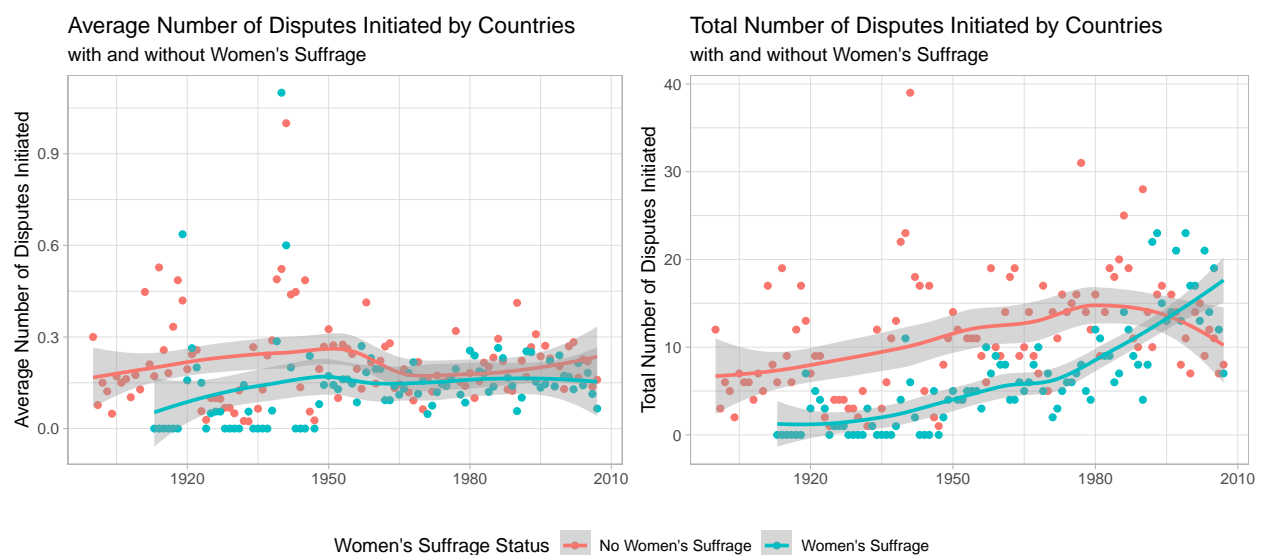
init_suffrage_grid <- plot_grid(init_suffrage_avg_plot +
  theme(legend.position = "none"),
  init_suffrage_total_plot +
  theme(legend.position = "none"))

legend <- get_legend(init_suffrage_avg_plot +
  theme(legend.position = "bottom"))

init_suffrage_grid_legend <- plot_grid(init_suffrage_grid, legend, ncol = 1, rel_heights = c(1, .1))

init_suffrage_grid_legend

```



```
# png("init_suffrage_grid.png", units="in", width=10, height=5, res=300)
# print(init_suffrage_grid_legend)
# dev.off()
```

From the two plots, we can see two unique time-series trends related to the average and total number of disputes initiated by countries with and without women's suffrage. The first trend looking at the average number of disputes shows that countries without women's suffrage have on average always initiated more conflicts than countries with women's suffrage since 1900.

The second trend looking at the total number of disputes shows that countries without women's suffrage have on average initiated more conflicts than countries with women's suffrage up to the late 1990s. Then, there was a switch and countries with women's suffrage started having on average initiated more conflicts than countries without women's suffrage. I think one potential cause of this is that more countries began having women's suffrage over time, so this second plot doesn't show the full story.

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 %in% c("North America", "South America")) %>%
  ggplot() +
  geom_sf(aes(fill = init, geometry = geom)) +
  labs(fill = "init", title = "Number of Conflicts initited in the Americas, 1960") +
  theme_void()
```

Number of Conflicts initiated in the Americas, 1960

