

# STAT108 Project - Exploratory Data Analysis

Yijia Sun

2022-11-19

## Research Question

We want to examine the impact of state-level abortion restrictions on abortions rate and cross-state movement to obtain abortion care in the United States, 2017 - 2019. Our modeling objective is to infer the overall impact on abortion outcomes in the United States.

## Data

Our data is collected from Centers for Disease Control and Prevention Abortion Surveillance System, Guttmacher Institute, and Advancing New Standards in Reproductive Health (ANSIRH) from University of California San Francisco.

We have the following variables in the dataset:

- **State** - 50 states in the United States and the District of Columbia
- **year** - 2017, 2018, 2019
- **policy\_index** - Score of each states based on whether they had policies in effect in six categories of abortion restrictions and six categories of measures that protect or expand abortion rights and access (6 is most supportive, -6 is most restrictive)
- **abortion\_rate\_residence** - Number of abortions per 1,000 women aged 15 - 44, by state of residence
- **percentage\_leaving** - Percentage of residents obtaining abortions who traveled out of state for care
- **facility\_density** - The number of abortion-providing facility per 100,000 women of reproductive age 15 - 49
- **facilities\_only\_procedural** - Percentage of facilities offering only procedural abortion
- **facilities\_only\_medication** - Percentage of facilities offering only medication abortion
- **Facilities\_both** - Percentage of facilities offering both procedural and medication abortion
- **gestational\_limit\_medication** - Mean gestational limit for medication abortion
- **gestational\_limit\_procedural** - Mean gestational limit for procedural abortion
- **cost\_medication** - The median self-pay costs for abortion services, in U.S. dollars
- **cost\_first\_trimester** - The median self-pay costs for first trimester procedural abortion services, in U.S. dollars
- **cost\_second\_trimester** - The median self-pay costs for second trimester procedural abortion services, in U.S. dollars
- **accepts\_insurance** - Percentage of abortion facilities accepting insurance
- **independent\_clinics** - Percentage of independent clinics
- **planned\_parenthoods** - Percentage of Planned Parenthoods
- **poverty** - Average percentage of people in poverty, 2019 - 2020

Our response variables are **abortion\_rate\_residence** and **percentage\_leaving**.

## Section 1 - Data Cleaning

### Import Data

```
df <- read_csv("data/new_abortion_data.csv")
df <- df[,-c(1,2)]
```

### Categorize state by policy index

We generate a new variable `policy_catog` to categorize states by their policy index. States with scores ranging from -6 to -2 are reported by Guttmacher to be hostile, -1 to +1 are neutral, and +2 to +6 are supportive.

```
df <- df %>% mutate(policy_catog =
  case_when(policy_index <= -2 ~ "hostile",
            policy_index <= 1 & policy_index >= -1 ~ "neutral",
            policy_index >= 2 ~ "supportive")
)
```

### Data Cleaning

Since states including California, Maryland, New Hampshire, and Wyoming don't report their data or only reported by occurrence, their abortion rate by state of residence aren't accurate so we remove these states from analysis. We also remove the District of Columbia from analysis, which was not included in the Guttmacher abortion policy report.

```
# filter out California, Maryland, New Hampshire, Wyoming, District of Columbia
df_filtered <- df[is.na(df$percentage_leaving) == FALSE & is.na(df$policy_index) == FALSE, ]
```

## Section 2 - Exploratory data analysis

Before building the model, we first conduct exploratory data analysis.

### Univariate

#### (1) Pie Chart and Map for categories of state policy

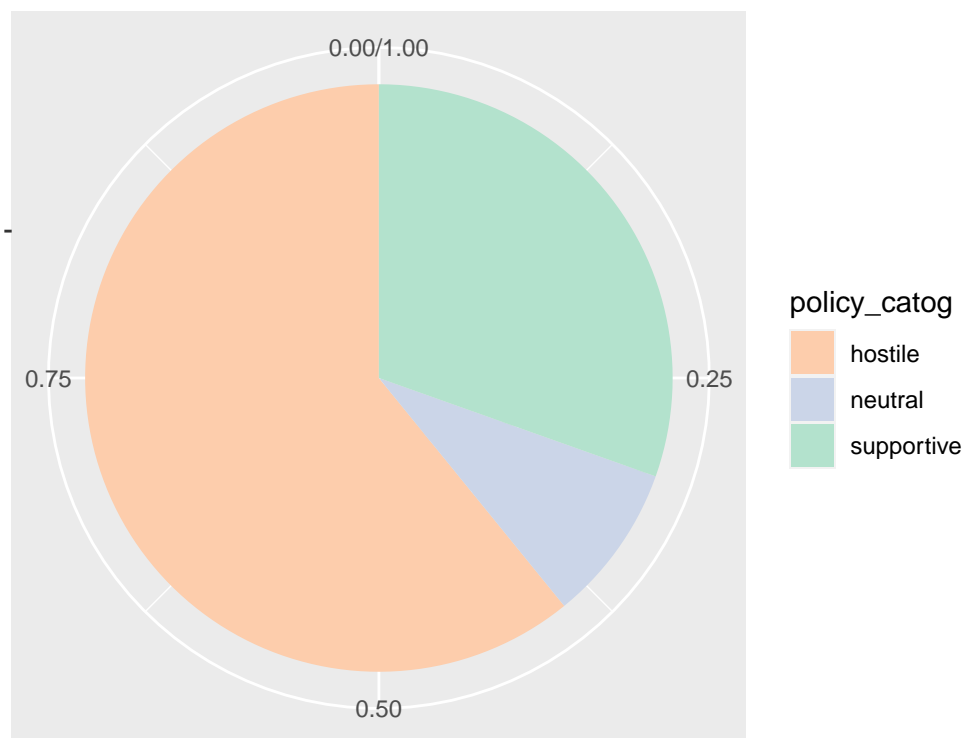
We create a pie chart and map for the categorical variable `policy_catog`, in order to see the distribution of policy categories.

```
#pie chart
df_filtered %>%
  count(policy_catog) %>%
  mutate(proportion = n / sum(n))
```

```
## # A tibble: 3 x 3
##   policy_catog      n proportion
##   <chr>          <int>      <dbl>
## 1 hostile         84      0.609
## 2 neutral         12      0.0870
## 3 supportive      42      0.304

df_filtered %>%
  ggplot(aes(x = "", fill = policy_catog)) +
  geom_bar(position = "fill", width = 1) +
  coord_polar(theta = "y") +
  labs(
    title = "Pie Chart for Abortion Policy Category",
    x = "",
    y = ""
  ) +
  scale_fill_manual(values=c("#fdcdac", "#cbd5e8", "#b3e2cd"))
```

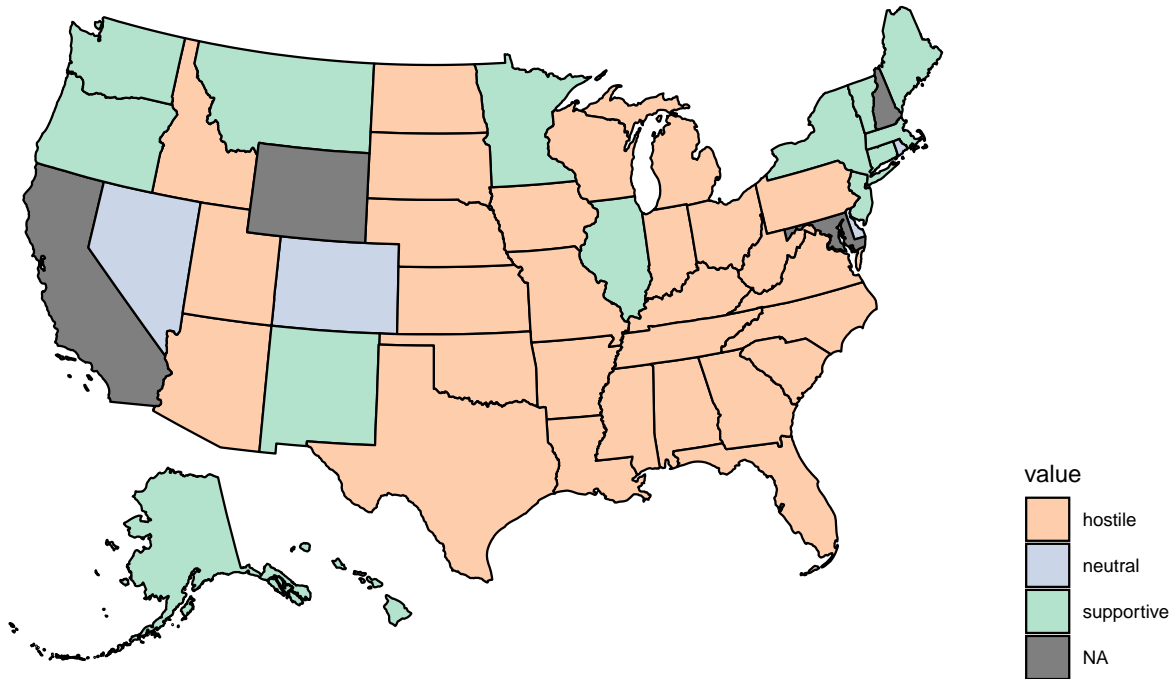
Pie Chart for Abortion Policy Category



```
#USMAP
#organzing data for state, transfer into abbr
map <- df_filtered[c(1,19)]
map <- map[order(df_filtered$State),]
map <- merge(map, statepop, by.x = "State", by.y = "full", all = TRUE)
map <- map[, -c(5)]
colnames(map) <- c('state', 'value', 'fips', 'abbr')
```

```
#plot usmap
plot_usmap(data = map, values = "value") + ggtitle("Abortion policy category") + scale_fill_manual(values = c("hostile", "neutral", "supportive", "NA")) +
  theme(legend.position = "right")
```

Abortion policy category



From pie chart and map above, we find:

- Most states have hostile abortion policy.
- Supportive states mostly are in the northeast and west parts.

## (2) Histograms for outcome variables

We create histograms to check the distribution of outcome variables, `abortion_rate_residence` and `percentage_leaving`.

```
his1 <- ggplot(df_filtered, aes(x = abortion_rate_residence)) +
  geom_histogram(aes(y=..density..), colour="black", fill="white", binwidth = 0.5) +
  geom_density(alpha=.2, fill="#FF6666") +
  #facet_grid(year ~ ., scales = "free") +
  xlab("abortion rate by state of residence") + ggtitle("Histogram for Abortion rate by State of Residence")

his2 <- ggplot(df_filtered, aes(x = percentage_leaving)) +
  geom_histogram(aes(y=..density..), colour="black", fill="white", binwidth = 3) +
  geom_density(alpha=.2, fill="#FF6666") +
  #facet_grid(year ~ ., scales = "free") +
```

```

xlab("percentage of leaving") + ggtitle("Histogram for Percentage of Leaving")

his1 + his2

```



From the histograms above, we find:

- Distribution for abortion rate is slightly right-skewed.
- Distribution for percentage of leaving is strongly right-skewed.

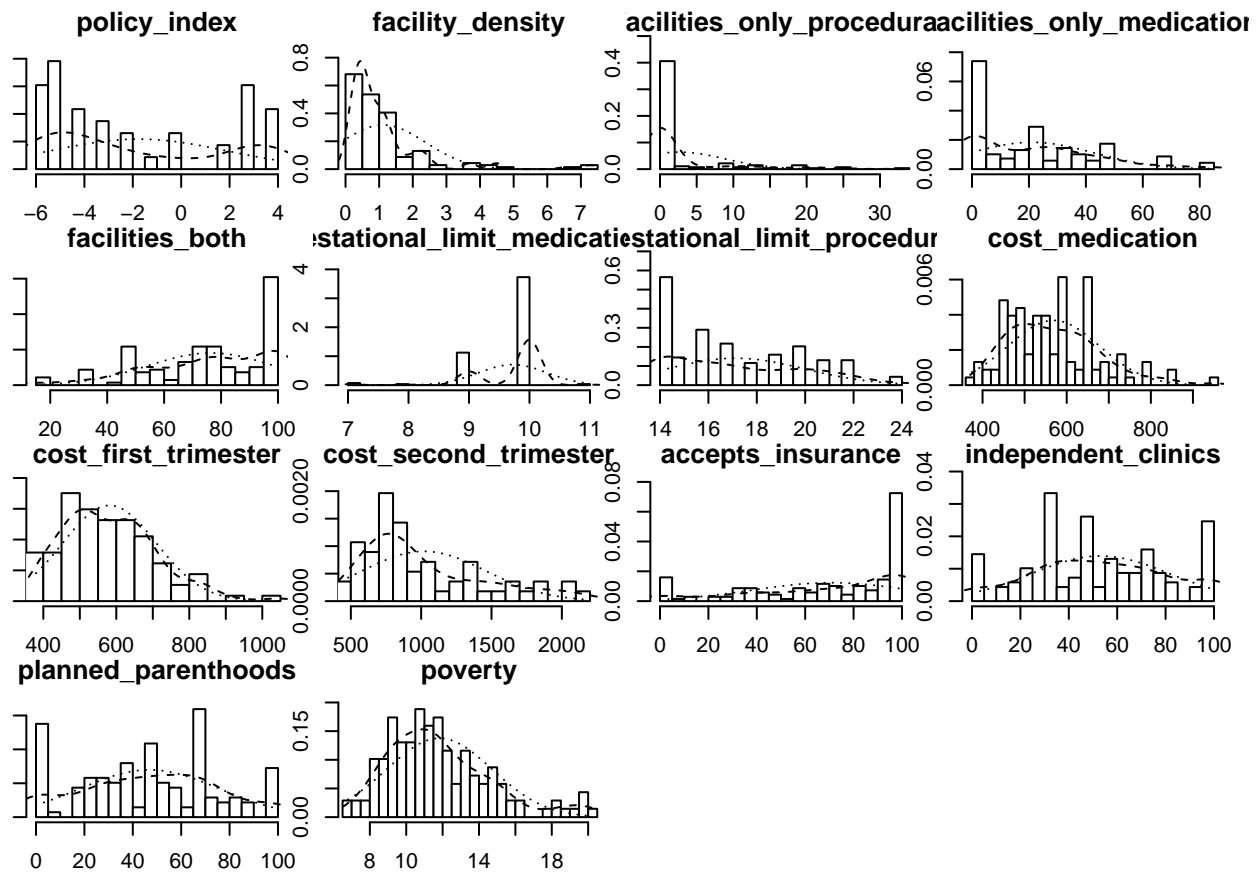
### (3) Histograms for covariates

Then, for all other variables, we create histograms to check their distributions.

```

cov <- df_filtered[,-c(1,2,4,5,19)]
multi.hist(cov, global = FALSE, )

```



From the histograms above, we find:

- `facility_density` and `cost_second_trimester` are more right-skewed.
- `facility_both` and `accepts_insurance` are more left-skewed.
- Most states have 0 facility with only procedural or medication.

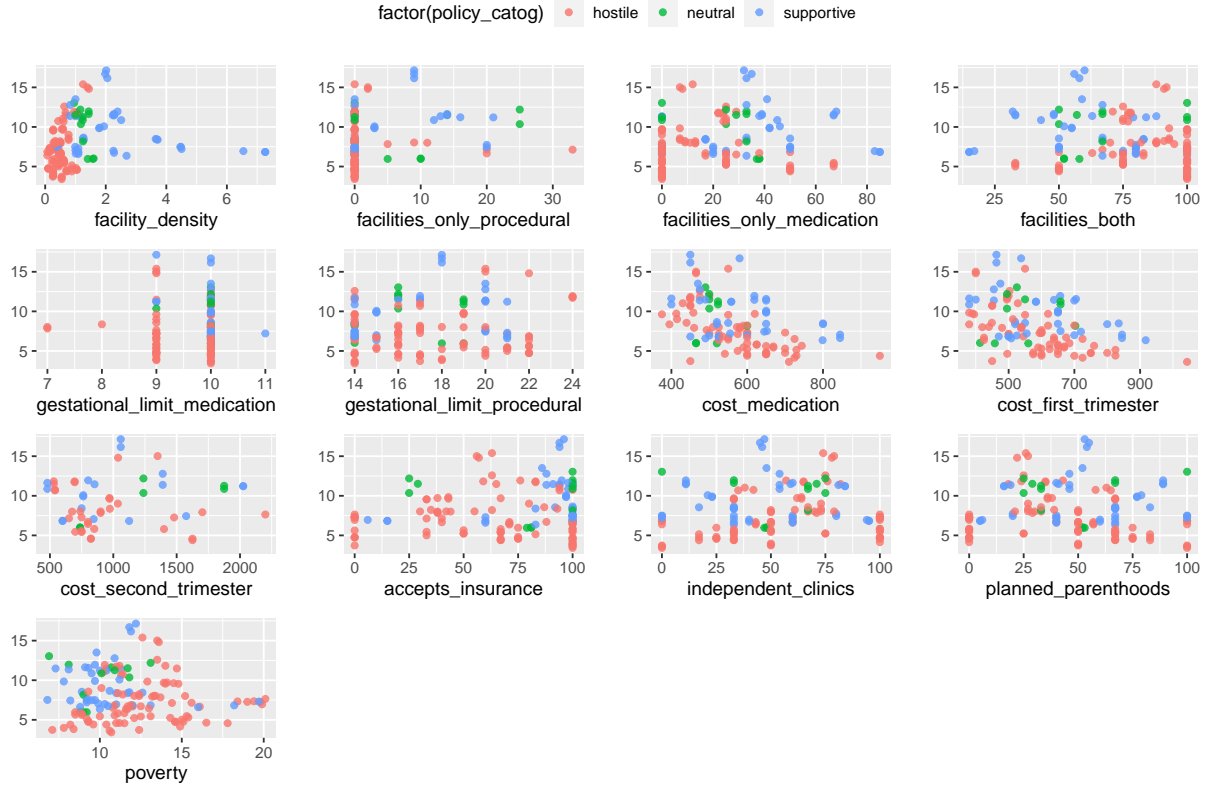
## Bivariate

Then, we use dotplots to analyze the relationship between all variables with each outcome variables, grouped by the policy category.

### (1) Dotplots for abortion rate

*(Code is too long so gets hidden.)*

## Dotplots between explanatory variables and abortion rate



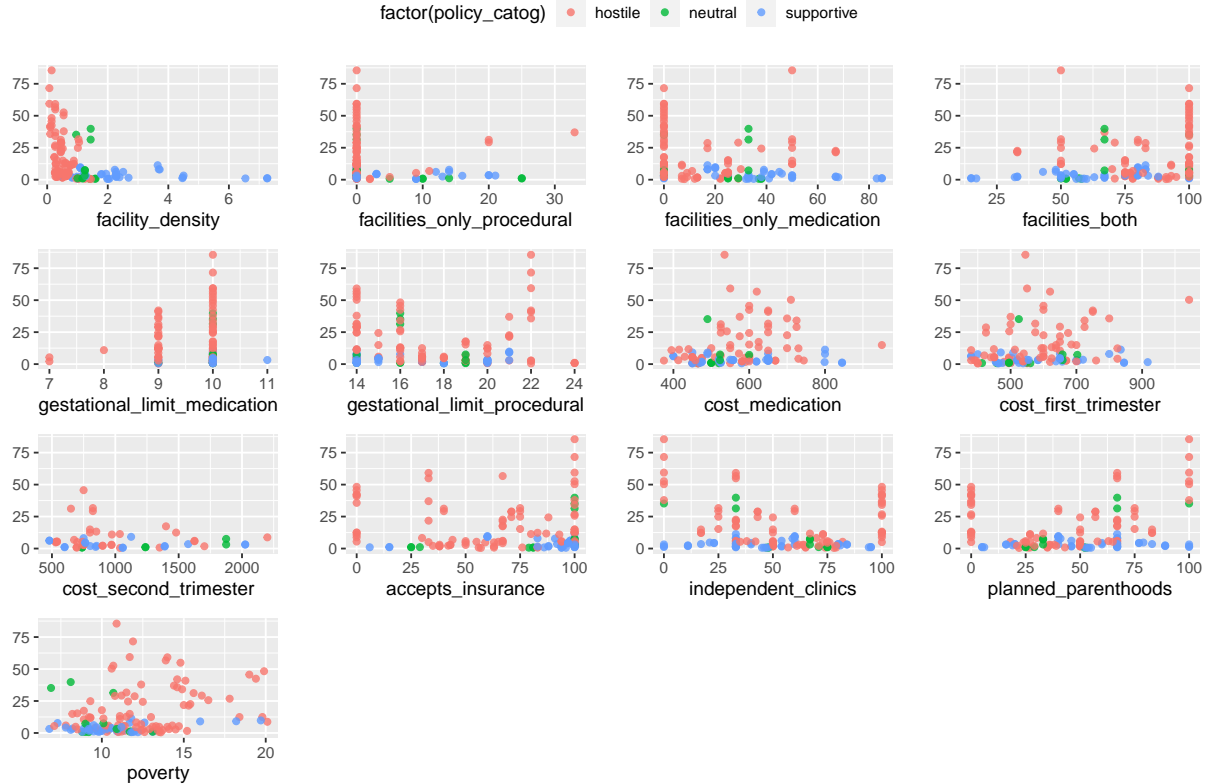
From the dotplots above, we find:

- Hostile states with higher `facility_density` tend to have relatively lower abortion rate.
- `cost_medication` and `cost_fisrt_trimester` both have negative association with abortion rate for all types of states.

## (2) Dotplots for percentage of leaving

(Code is too long so get hidden.)

Dotplots between explanatory variables and percentage of leaving



From the dotplots above, we find:

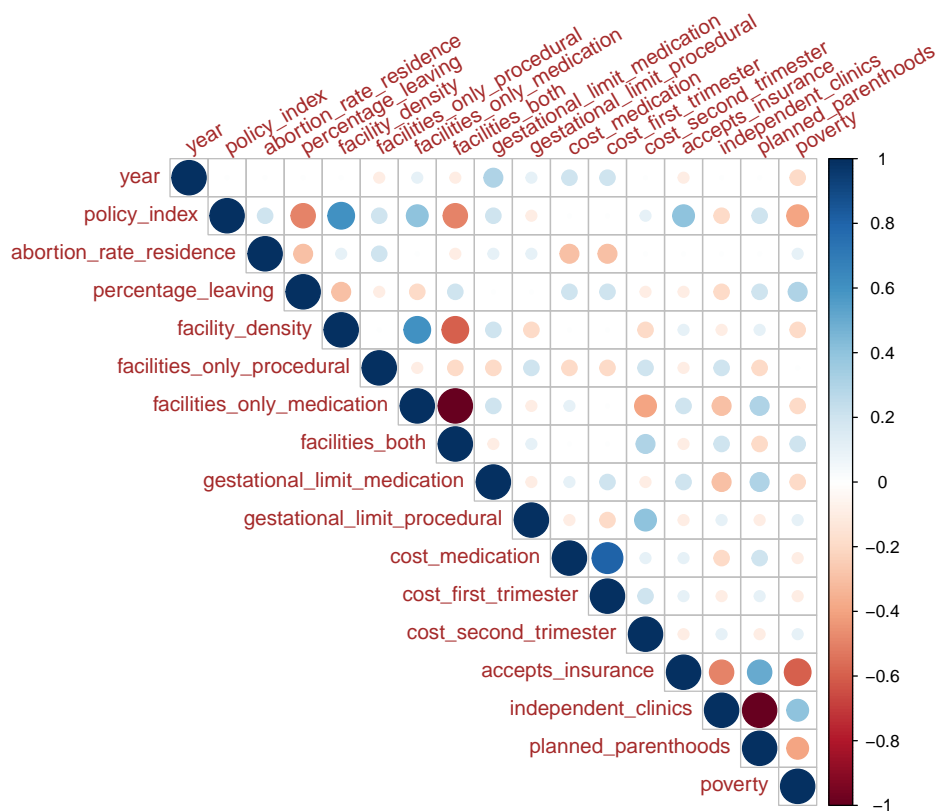
- Hostile states with higher `facility_density` tend to have relatively higher percentage of leaving.
- Hostile states with `gestational_limit_medication` of 9 or 10 tend to have relatively higher percentage of leaving.

## Multivariate

Then, to analyze potential interaction terms or multicollinearity, we use correlation matrix to see the correlation coefficient between all variables.

```
correlation <- df[, -c(1,19)]
corr <- round(cor(correlation,use="pairwise.complete.obs"), 1)
corrplot(corr, tl.col = "brown", bg = "White", tl.srt=30, tl.cex =1,type = "upper")
```





From the correlation matrix above, we find:

- `policy_index` & `facility_density`, `accepts_insurance` & `poverty` have relatively strong negative correlation.
- `facilities_only_procedural` & `facilities_both`, `independent_clinics` & `planned_parenthoods` have very strong negative correlation.
- `cost_medication` & `cost_first_trimester` have very strong positive correlation.