

Matt McCoy

EC 525

Prof: Eric Zou

5/20/2021

## Empirical Project 2

### **Q1:**

Regression Discontinuity allows us to compare the outcomes of individuals who fall on either side of a cutoff to get an estimate of the causal effect, but this estimate will be biased if individuals can manipulate exactly which side of the cutoff they end up in. A simple comparison of air pollution in the northern cities compared to the southern cities would not measure the causal effect of the Huai River Policy. This is because when we compare all of the North to all of the South, there is a greater number of statistically significant differences, underscoring the value of the RD design in this setting.

The conventional approach, which uses OLS to estimate their model, rests on the assumption that linear adjustment for the limited set of variables available in the census removes all sources of confounding. The issue with this is that previous research has raised considerable concerns about the validity of this assumption.

This paper uses an RD design that exploits the Huai River Policy which provides free or heavily subsidized coal for indoor heating north of the river and no subsidies to the south. They separately tested whether the Huai River Policy caused discontinuous changes in PM10 and life expectancy to the north of the river compared to the south. If the necessary assumptions are that any unobserved determinants of PM10 or mortality change smoothly as they cross the river boundary, and if this assumption is valid then adjusting for a sufficiently flexible polynomial in

distance from the Huai river or local linear regressions on either side of the river will remove all potential sources of bias and allow for causal inference. A simple comparison of air pollution in northern cities versus southern cities would not measure the causal effect of the Huai River Policy because they were not randomly assigned, but they get over this using a quasi-random variation in pollution to estimate the long term impacts

**Q2:**

The outcome variable is  $PM_{10}$  exposure, and the assignment variable is where the individual is located relative to the Huai River, this is measured as degrees north of the Huai River boundary. Where they are located in relation to the river will decide which side of the cutoff they end up in. The river is the cutoff and anyone North will receive treatment while those who are located South of the river will not be treated and will be the control group. This experiment exploits the Huai River Policy that provides free or heavily subsidized coal for indoor heating to those who live north of the river as the treatment, and no subsidies to the south as the control.

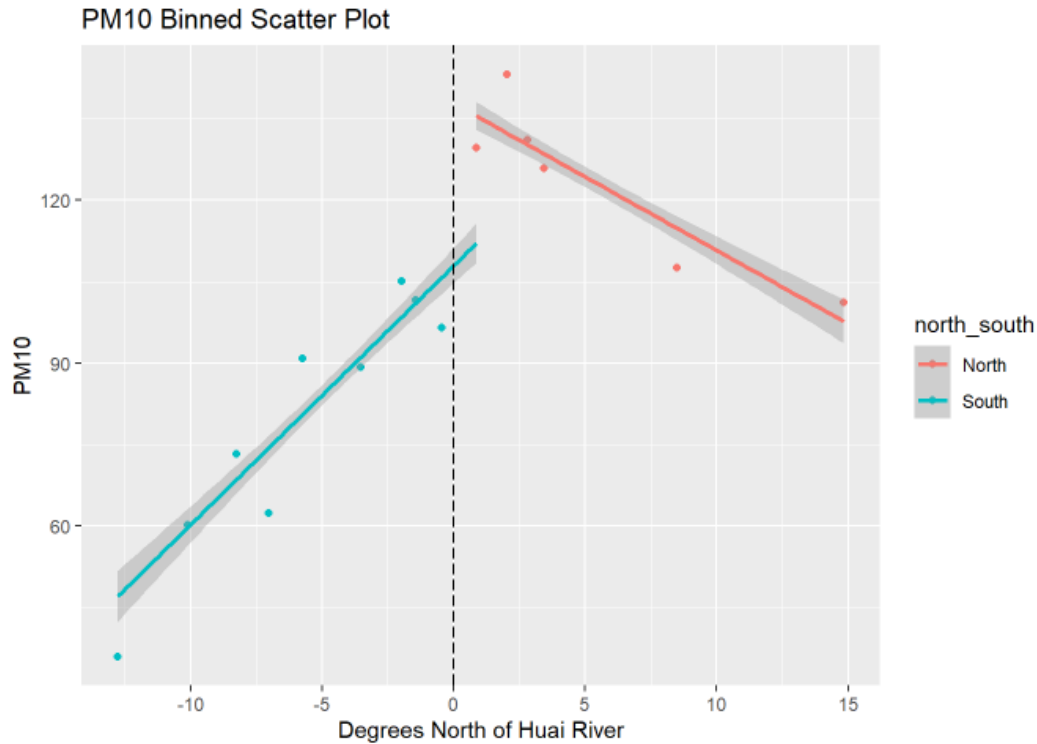
**Q3:**

A Binned scatter plot is a transparent way for us to present our data. It does a great job of displaying to the audience which “parts” of the data is driving the average relationship present. Binscatter allows us to assess the functional form assumption, and is thus known as a “non-parametric” way of getting  $E[Y|X]$ . Binned scatter plots graph the nonparametric relationship between two of our variables, for multiple subgroups, either conditionally or unconditionally on a set of controls.

A binned scatter plot is constructed by taking the raw data points and grouping them into bins, and then an aggregate statistic is used to summarize each bin. We basically just take the raw data and compute averages.

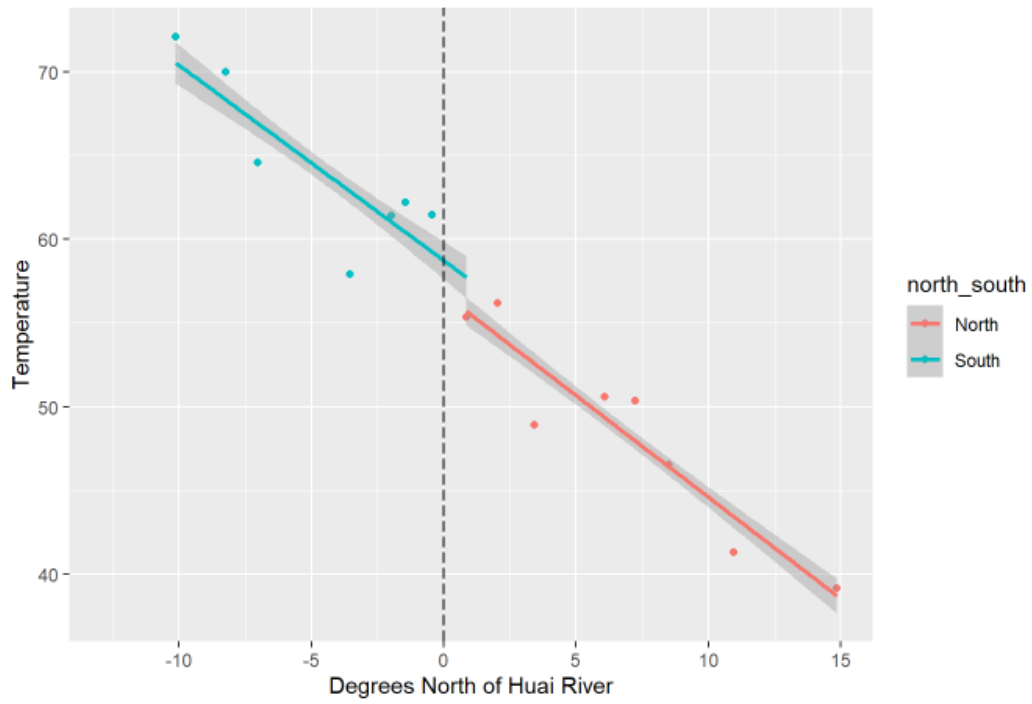
1. Group the x-axis variable into equal-sized bins
2. Compute the mean of the x-axis and y-axis variables within each of our bins
3. Create a scatterplot of these data points
4. Draw the population regression line

**Q4.a:**

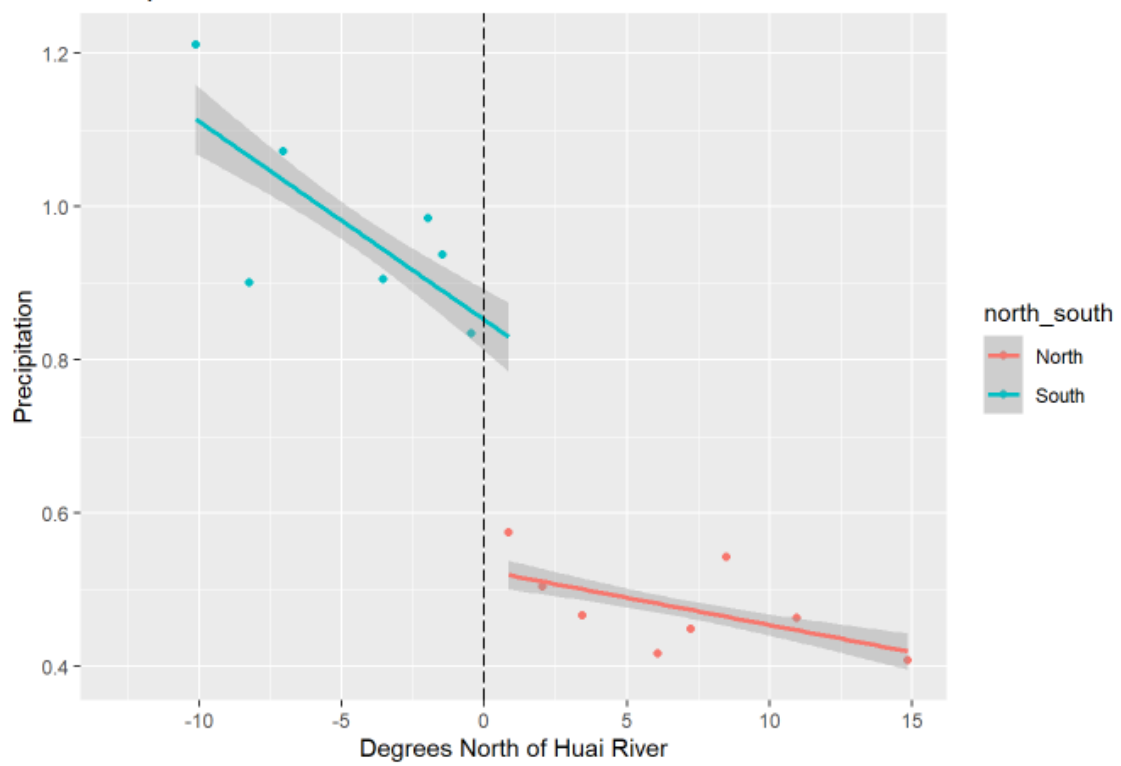


**Q4.b:**

Temperature Binned Scatter Plot



Precipitation Binned Scatter Plot





The 95% confidence interval for the slope is the estimated coefficient  $\pm$  two standard errors. The confidence interval can be calculated using the equation:

$$[Estimate - (2 \cdot SE), Estimate + (2 \cdot SE)] = [.]$$

The confidence interval for the PM10 regression is:

$$[.091 - (2 \cdot 0.67), .091 + (2 \cdot 0.67)] = [-1.25, 1.43]$$

The confidence interval for the temperature regression is:

$$[-1.191 - (2 \cdot 0.107), -1.191 + (2 \cdot 0.107)] = [-1.405, -0.977]$$

The confidence interval for the precipitation regression:

$$[-0.013 - (2 \cdot 0.005), -0.013 + (2 \cdot 0.005)] = [-0.023, -0.003]$$

The confidence interval for the windspeed regression:

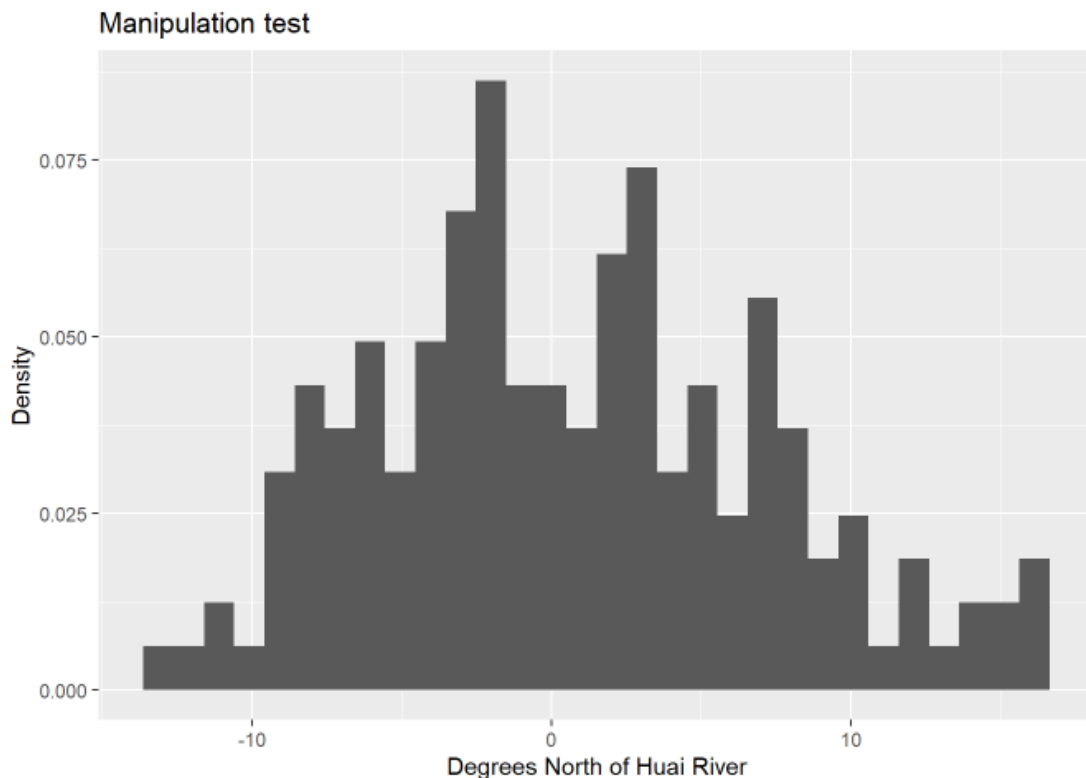
$$[0.046 - (2 \cdot 0.013), 0.046 + (2 \cdot 0.013)] = [0.02, 0.072]$$

#### **Q6:**

For a quasi experiment we must satisfy the Identification Assumption, which states that a change in the treatment variable is the only reason for a discrete jump in the outcome variable around the cutoff. Or, “all observed and unobserved determinants of  $Y_i$  (other than treatment) are smooth around the cutoff”. Some of our graphs in 4b are consistent with this assumption, whereas some are not. We can see that the fitted line for temperature does not have a large jump at the cutoff, but instead the line stays fairly smooth as it crosses showing that it is consistent with our assumption. This is not the case for precipitation. We can see that there is a very large drop in precipitation as you go north of the Huai river. This jump in precipitation was not caused by the treatment, so it fails the Identification Assumption. I do not believe that the wind speed binned scatter supports the Identification Assumption because we appear to have stronger winds south of the cut off, but this one is less clear.

**Q7:**

I believe that a manipulation test here, given the context of the study, would be a good idea because the Huai River policy can cause people to manipulate where they are located relative to the river. If only people living in the North will get free coal during the winter, then it is very possible that people will manipulate which side of the cutoff they end up. A manipulation test is a quick way for us to observe whether or not people are doing something to receive treatment. I primarily believe that we should do a manipulation test because free winter heating is a very big deal for many people who may not be able to afford it otherwise, and even though this study was originally done during a time of restricted movement I still think it's possible for people to try and manipulate the cutoff.



Based on the histogram, it does not seem like there is manipulation at the cutoff. There are actually more people just south of the river rather than just north which is what I would have expected to see. There appears to be no evidence of manipulation around cutoff.

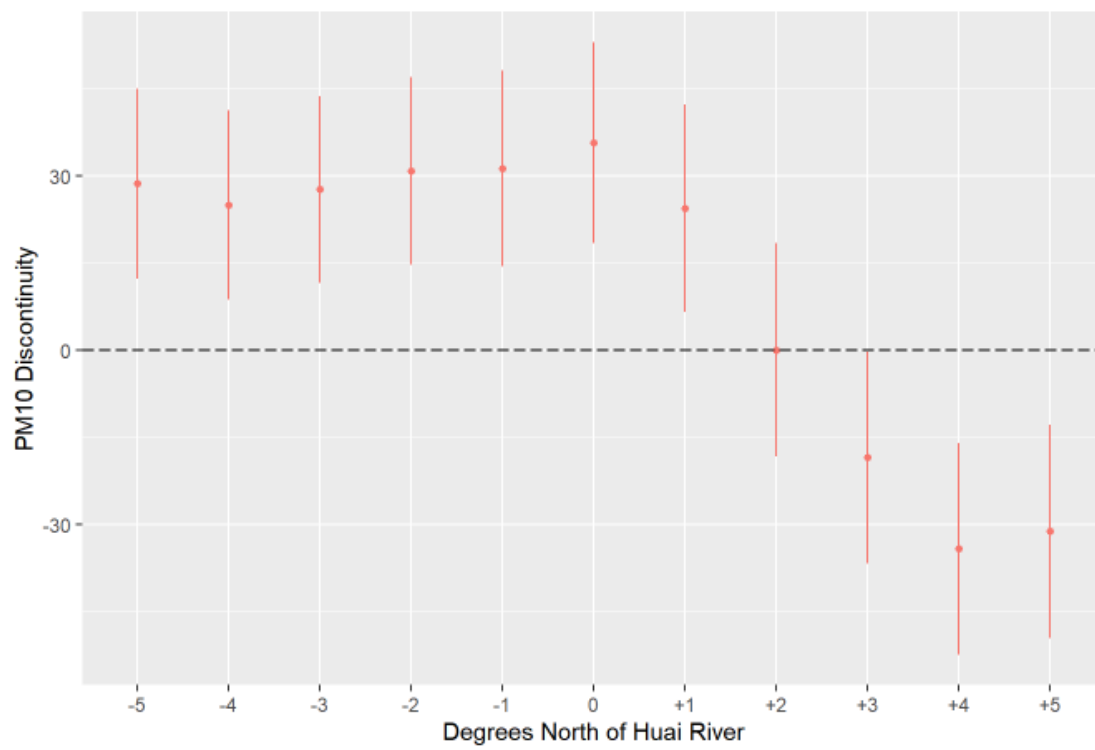
**Q8a:**

A placebo test allows us to probe the soundness of a research design by checking for associations that should be present if the design is flawed but not otherwise. The placebo test in figure 4 estimated discontinuity in pollution and life expectancy at displaced Huai River boundaries using a bandwidth selection method. To do this test they pretended the Huai river was located one degree north or one degree south for example, and they did this “fake” test multiple times with different locations. They tested for whether or not there would be a discrete jump in life expectancy or PM10 at the fake Huai Rivers. Estimating regression discontinuity using false locations of the Huai River allows them to test for these discrete jumps in life expectancy and PM10 which tells them if the research design is flawed. The results of the test provide extra confidence that what they found was causal. The results of this test told them that only at the actual Huai River is there a discontinuity in either life expectancy or PM10 observed, which provides supporting evidence for their overall empirical strategy.

**Q8b:**



Huai River Placebo Test



# Empirical Project 2

Matt McCoy

5/17/2021

```
## Install the pacman package if necessary
if (!require("pacman")) install.packages("pacman")
```

```
## Loading required package: pacman
```

```
## Install other packages using pacman::p_load()
pacman::p_load(tidyverse, haven, sandwich, lmtest, stargazer, dplyr, ggplot2, broom, magrittr)

getwd()
```

```
## [1] "C:/Users/mattm/OneDrive/Desktop"
```

```
river_df <- read_dta(file = "huairiver.dta")
```

```
#create bins for dist
river_df <- river_df %>% mutate(dist_bin= cut(dist_huai, breaks=quantile(dist_huai, probs = seq(
0, 1, by = 0.05), na.rm = TRUE)))

is.factor(river_df$dist_bin)
```

```
## [1] TRUE
```

```
table(river_df$dist_bin)
```

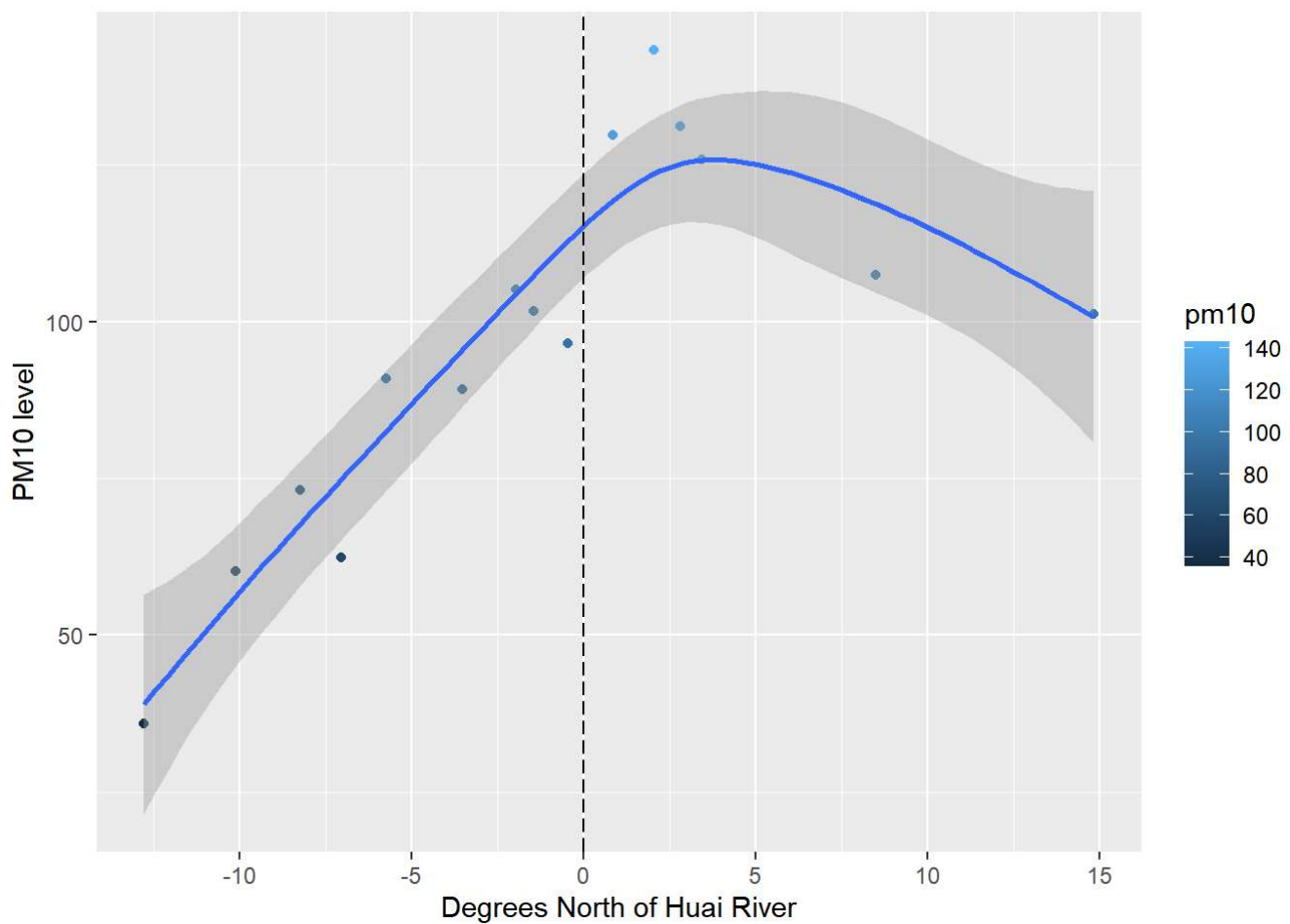
```
##
##  (-12.8,-9.01]  (-9.01,-7.77]  (-7.77,-6.5]  (-6.5,-5.29]  (-5.29,-4.1]
##           8           8           8           8           8
##  (-4.1,-3.24]  (-3.24,-2.17]  (-2.17,-1.78]  (-1.78,-1.07]  (-1.07,-0.0471]
##           8           8           8           8           8
##  (-0.0471,1.4]  (1.4,2.24]  (2.24,3.15]  (3.15,3.89]  (3.89,5.29]
##           8           8           8           8           8
##  (5.29,6.83]  (6.83,7.71]  (7.71,9.29]  (9.29,12.4]  (12.4,16.5]
##           8           8           8           8           8
```

```
river_df %>%
  group_by(dist_bin) %>%
  summarise(pm10 = mean(pm10), dist_huai= mean(dist_huai)) %>%
  ggplot(aes(x=dist_huai,y=pm10,color=pm10)) + geom_point() + geom_smooth(method="gam") + geom_v
line(xintercept=0,linetype="longdash") + ylab("PM10 level") + xlab("Degrees North of Huai River"
)
```

```
## `geom_smooth()` using formula 'y ~ s(x, bs = "cs")'
```

```
## Warning: Removed 6 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 6 rows containing missing values (geom_point).
```



```
river_df %<>% mutate(north_south = ifelse(north_huai == 1, "North", "South"))
```

```
#4a.)
river_df %>%
  group_by(dist_bin) %>%
  summarise(pm10 = mean(pm10), dist_huai= mean(dist_huai), north_south = north_south) %>%
  ggplot(aes(x=dist_huai,y=pm10,color = north_south))+geom_point()+
  geom_smooth(method="lm")+geom_vline(xintercept=0,linetype="longdash") + ggtitle("PM10 Binned S
catter Plot") +
  labs(y="PM10", x = "Degrees North of Huai River")
```

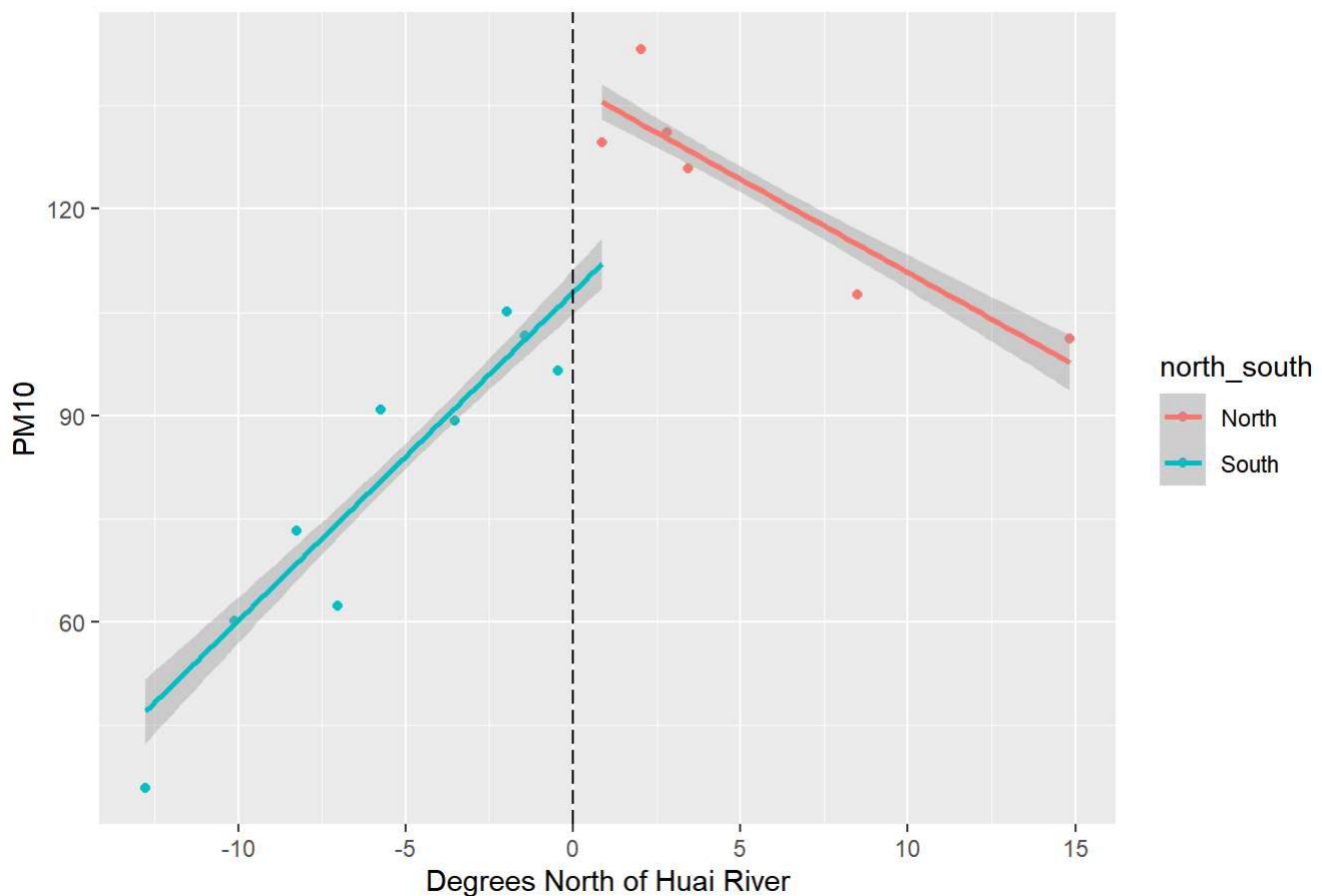
```
## `summarise()` has grouped output by 'dist_bin'. You can override using the `.groups` argumen
t.
```

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

```
## Warning: Removed 48 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 48 rows containing missing values (geom_point).
```

PM10 Binned Scatter Plot



#4b.i.)

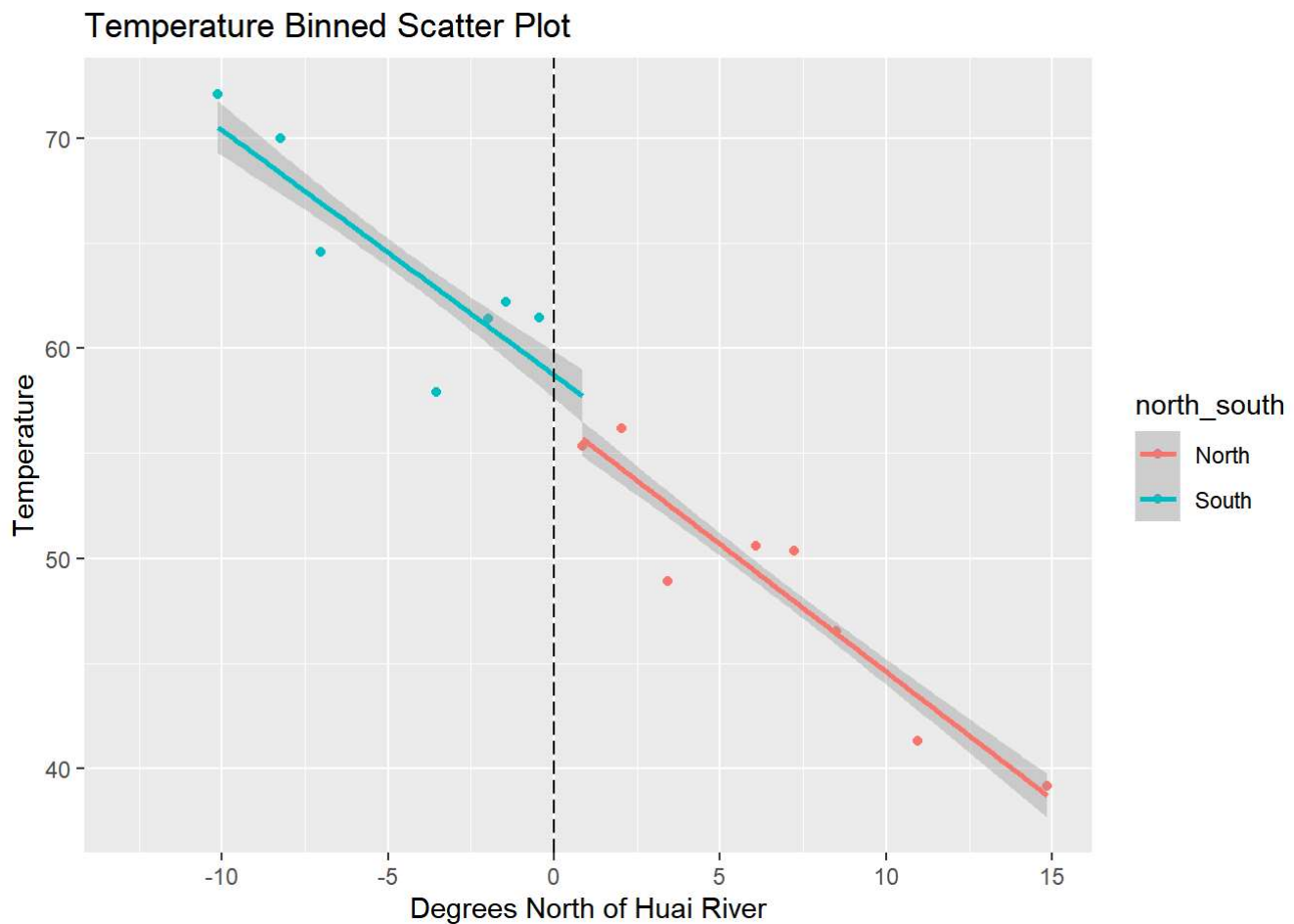
```
river_df %>%
  group_by(dist_bin) %>%
  summarise(temp = mean(temp), dist_huai= mean(dist_huai), north_south = north_south) %>%
  ggplot(aes(x=dist_huai,y=temp,color = north_south))+geom_point()+
  geom_smooth(method="lm")+geom_vline(xintercept=0,linetype="longdash") + ggtitle("Temperature Binned Scatter Plot") +
  labs(y="Temperature", x = "Degrees North of Huai River")
```

## `summarise()` has grouped output by 'dist\_bin'. You can override using the `.groups` argument.

## `geom\_smooth()` using formula 'y ~ x'

## Warning: Removed 41 rows containing non-finite values (stat\_smooth).

## Warning: Removed 41 rows containing missing values (geom\_point).



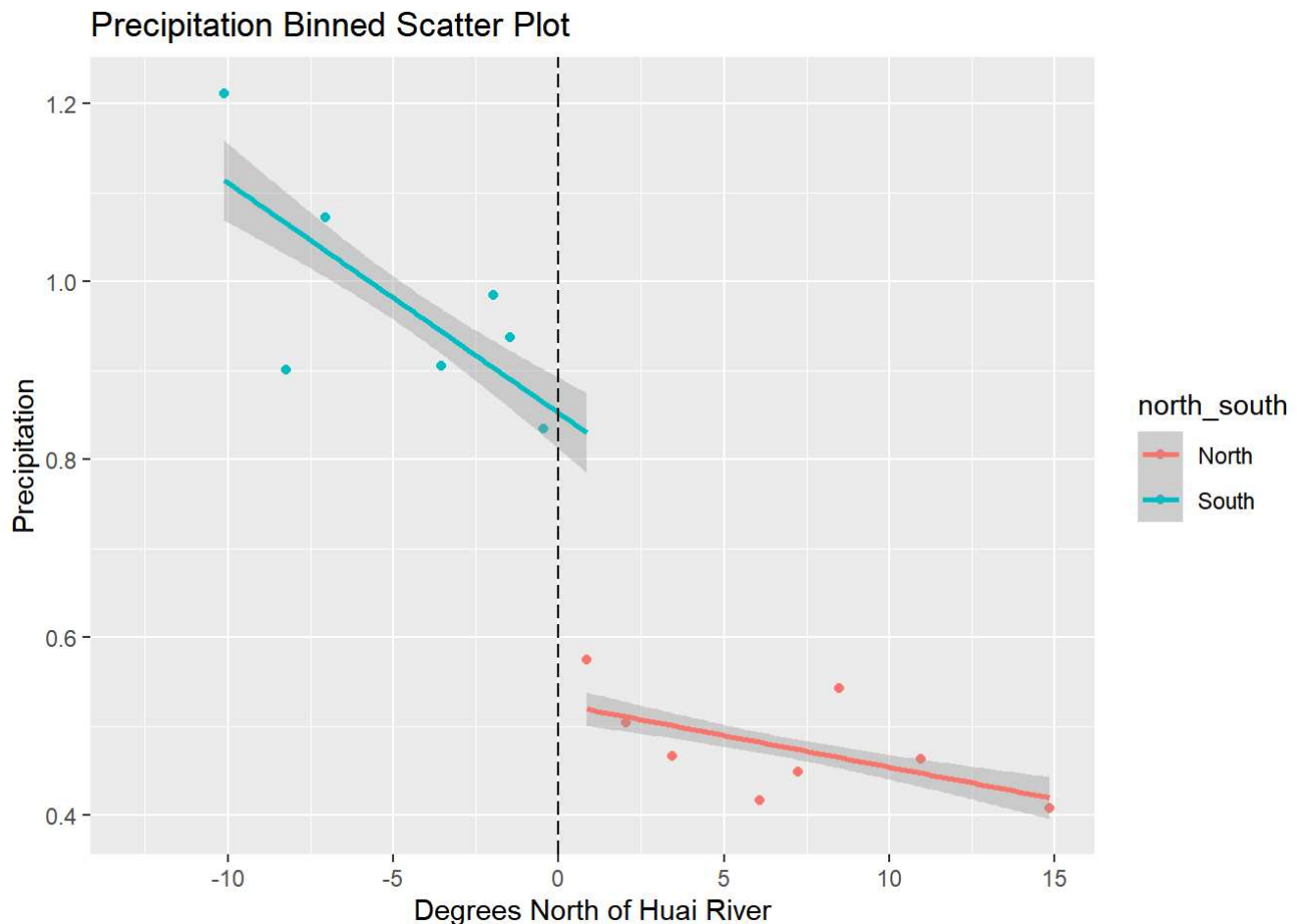
```
river_df %>%
  group_by(dist_bin) %>%
  summarise(prcp = mean(prcp), dist_huai= mean(dist_huai), north_south = north_south) %>%
  ggplot(aes(x=dist_huai,y=prcp,color = north_south))+geom_point()+
  geom_smooth(method="lm")+geom_vline(xintercept=0,linetype="longdash") + ggtitle("Precipitation
Binned Scatter Plot") +
  labs(y="Precipitation", x = "Degrees North of Huai River")
```

## `summarise()` has grouped output by 'dist\_bin'. You can override using the `.groups` argument.

## `geom\_smooth()` using formula 'y ~ x'

## Warning: Removed 41 rows containing non-finite values (stat\_smooth).

## Warning: Removed 41 rows containing missing values (geom\_point).



```
library(ggthemes)
```

## Warning: package 'ggthemes' was built under R version 4.0.5

```
river_df %>%
  group_by(dist_bin) %>%
  summarise(wspd = mean(wspd), dist_huai = mean(dist_huai), north_south = north_south) %>%
  ggplot(aes(x=dist_huai, y=wspd, color = north_south)) + geom_point() +
  geom_smooth(method="lm") + geom_vline(xintercept=0, linetype="longdash") + ggtitle("Windspeed Binned Scatter Plot") +
  labs(y="Windspeed", x = "Degrees North of Huai River")
```

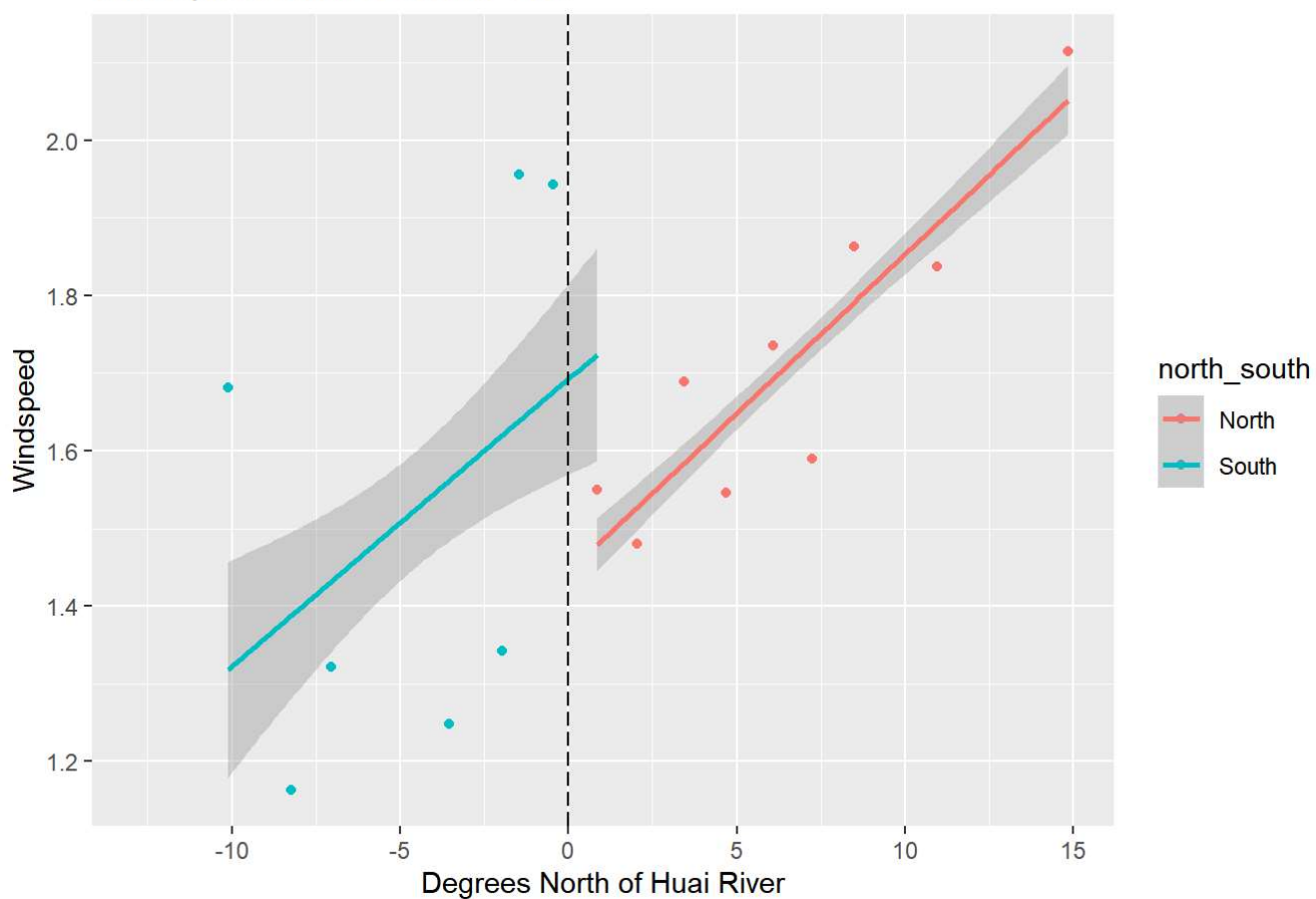
## `summarise()` has grouped output by 'dist\_bin'. You can override using the `.groups` argument.

## `geom\_smooth()` using formula 'y ~ x'

## Warning: Removed 33 rows containing non-finite values (stat\_smooth).

## Warning: Removed 33 rows containing missing values (geom\_point).

Windspeed Binned Scatter Plot



```
pm10_reg <- lm(pm10 ~ dist_huai + north_huai, data = river_df)
stargazer(pm10_reg, type = 'text')
```

```
##
## =====
##                      Dependent variable:
##                      -----
##                      pm10
## -----
## dist_huai           0.091
##                      (0.670)
##
## north_huai          35.728***
##                      (8.823)
##
## Constant            85.876***
##                      (4.666)
##
## -----
## Observations         154
## R2                   0.262
## Adjusted R2          0.252
## Residual Std. Error  31.075 (df = 151)
## F Statistic          26.833*** (df = 2; 151)
## =====
## Note:                *p<0.1; **p<0.05; ***p<0.01
```

```
temp_reg <- lm(temp ~ dist_huai + north_huai , data = river_df)

stargazer(temp_reg, type = 'text')
```

```
##
## =====
##                      Dependent variable:
##                      -----
##                      temp
## -----
## dist_huai           -1.191***
##                      (0.107)
##
## north_huai          -1.329
##                      (1.399)
##
## Constant            57.852***
##                      (0.740)
##
## -----
## Observations         153
## R2                   0.749
## Adjusted R2          0.745
## Residual Std. Error  4.887 (df = 150)
## F Statistic          223.336*** (df = 2; 150)
## =====
## Note:                *p<0.1; **p<0.05; ***p<0.01
```



```
prcp_reg <- lm(prcp ~ dist_huai + north_huai , data = river_df)

stargazer(prcp_reg, type = 'text')
```

```
##
## =====
##                      Dependent variable:
##                      -----
##                      prcp
## -----
## dist_huai            -0.013**
##                      (0.005)
##
## north_huai           -0.367***
##                      (0.070)
##
## Constant             0.912***
##                      (0.037)
## -----
## Observations          153
## R2                    0.532
## Adjusted R2           0.525
## Residual Std. Error   0.243 (df = 150)
## F Statistic           85.143*** (df = 2; 150)
## =====
## Note:                 *p<0.1; **p<0.05; ***p<0.01
```

```
wspd_reg <- lm(wspd ~ dist_huai + north_huai , data = river_df)

stargazer(wspd_reg, type = 'text')
```

```
##
## =====
##                      Dependent variable:
##                      -----
##                      wspd
## -----
## dist_huai           0.046***
##                      (0.013)
##
## north_huai          -0.282*
##                      (0.167)
##
## Constant            1.659***
##                      (0.088)
##
## -----
## Observations         156
## R2                   0.104
## Adjusted R2          0.092
## Residual Std. Error   0.585 (df = 153)
## F Statistic           8.858*** (df = 2; 153)
## =====
## Note:                *p<0.1; **p<0.05; ***p<0.01
```

```
stargazer(pm10_reg, temp_reg, prcp_reg, wspd_reg, type = 'text')
```

```
##
## =====
##
##                                     Dependent variable:
## -----
##          pm10          temp          prcp
wspd
##          (1)          (2)          (3)
(4)
## -----
## dist_huai          0.091          -1.191***          -0.013**
0.046***
##          (0.670)          (0.107)          (0.005)
(0.013)
##
## north_huai          35.728***          -1.329          -0.367***
-0.282*
##          (8.823)          (1.399)          (0.070)
(0.167)
##
## Constant          85.876***          57.852***          0.912***
1.659***
##          (4.666)          (0.740)          (0.037)
(0.088)
##
## -----
## Observations          154          153          153
156
## R2          0.262          0.749          0.532
0.104
## Adjusted R2          0.252          0.745          0.525
0.092
## Residual Std. Error    31.075 (df = 151)    4.887 (df = 150)    0.243 (df = 150)
0.585 (df = 153)
## F Statistic    26.833*** (df = 2; 151) 223.336*** (df = 2; 150) 85.143*** (df = 2; 150)
8.858*** (df = 2; 153)
## =====
## Note:
1; **p<0.05; ***p<0.01
```

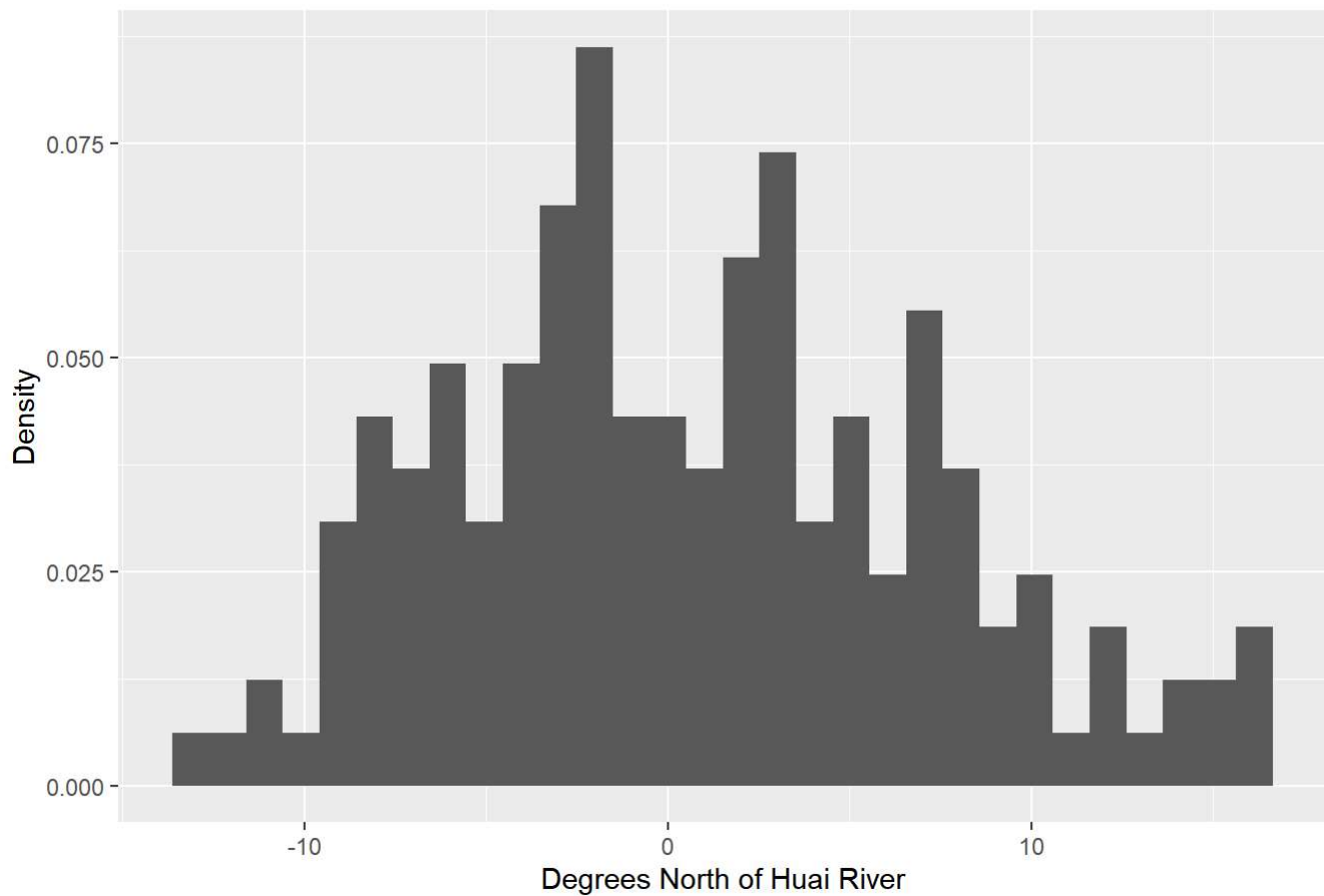
\*p&lt;0.

```
# conf. interval: [1.96 - 2*SE, 1.96 + 2*SE]
```

```
ggplot(river_df, aes(x = dist_huai, after_stat(density))) +
  geom_histogram() + ggtitle("Manipulation test") +
  labs(y="Density", x = "Degrees North of Huai River")
```

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

### Manipulation test



```
river_df$one_degree_north <- (river_df$dist_huai - 1)
river_df$two_degree_north <- (river_df$dist_huai - 2)
river_df$three_degree_north <- (river_df$dist_huai - 3)
river_df$four_degree_north <- (river_df$dist_huai - 4)
river_df$five_degree_north <- (river_df$dist_huai - 5)
river_df$one_degree_south <- (river_df$dist_huai + 1)
river_df$two_degree_south <- (river_df$dist_huai + 2)
river_df$three_degree_south <- (river_df$dist_huai + 3)
river_df$four_degree_south <- (river_df$dist_huai + 4)
river_df$five_degree_south <- (river_df$dist_huai + 5)
```

```
river_df %<>% mutate(one_degree_north_huai = ifelse(one_degree_north > 0, 1, 0))  
river_df %<>% mutate(two_degree_north_huai = ifelse(two_degree_north > 0, 1, 0))  
river_df %<>% mutate(three_degree_north_huai = ifelse(three_degree_north > 0, 1, 0))  
river_df %<>% mutate(four_degree_north_huai = ifelse(four_degree_north > 0, 1, 0))  
river_df %<>% mutate(five_degree_north_huai = ifelse(five_degree_north > 0, 1, 0))  
river_df %<>% mutate(one_degree_south_huai = ifelse(one_degree_south > 0, 1, 0))  
river_df %<>% mutate(two_degree_south_huai = ifelse(two_degree_south > 0, 1, 0))  
river_df %<>% mutate(three_degree_south_huai = ifelse(three_degree_south > 0, 1, 0))  
river_df %<>% mutate(four_degree_south_huai = ifelse(four_degree_south > 0, 1, 0))  
river_df %<>% mutate(five_degree_south_huai = ifelse(five_degree_south > 0, 1, 0))
```

```
new_reg1 <- lm(pm10 ~ one_degree_north_huai + one_degree_north, data = river_df)  
new_reg2 <- lm(pm10 ~ two_degree_north_huai + two_degree_north, data = river_df)  
new_reg3 <- lm(pm10 ~ three_degree_north_huai + three_degree_north , data = river_df)  
new_reg4 <- lm(pm10 ~ four_degree_north_huai + four_degree_north, data = river_df)  
new_reg5 <- lm(pm10 ~ five_degree_north_huai + five_degree_north, data = river_df)  
stargazer(new_reg1, new_reg2, new_reg3, new_reg4, new_reg5, type = 'text')
```

```

##
## =====
##                               Dependent variable:
##                               -----
##                               pm10
##                               (1)      (2)      (3)      (4)      (5)
## -----
## one_degree_north_huai      24.386***
##                               (9.147)
##
## one_degree_north           0.799
##                               (0.693)
##
## two_degree_north_huai      -0.019
##                               (9.356)
##
## two_degree_north           2.325***
##                               (0.699)
##
## three_degree_north_huai    -18.521**
##                               (9.288)
##
## three_degree_north         3.405***
##                               (0.671)
##
## four_degree_north_huai     -34.253***
##                               (9.280)
##
## four_degree_north          4.160***
##                               (0.628)
##
## five_degree_north_huai     -31.251***
##                               (9.385)
##
## five_degree_north          3.961***
##                               (0.626)
##
## Constant                   92.204*** 106.452*** 117.778*** 126.751*** 128.804***
##                               (5.162)  (5.379)  (5.303)  (5.111)  (5.548)
##
## -----
## Observations                154      154      154      154      154
## R2                          0.219      0.182      0.203      0.250      0.238
## Adjusted R2                 0.209      0.171      0.193      0.240      0.228
## Residual Std. Error (df = 151) 31.975      32.718      32.296      31.335      31.580
## F Statistic (df = 2; 151)      21.153*** 16.808*** 19.239*** 25.136*** 23.587***
## =====
## Note:                        *p<0.1; **p<0.05; ***p<0.01

```

```
new_reg6 <- lm(pm10 ~ one_degree_south_huai + one_degree_south , data = river_df)
new_reg7 <- lm(pm10 ~ two_degree_south_huai + two_degree_south, data = river_df)
new_reg8 <- lm(pm10 ~ three_degree_south_huai + three_degree_south, data = river_df)
new_reg9 <- lm(pm10 ~ four_degree_south_huai + four_degree_south, data = river_df)
new_reg10 <- lm(pm10 ~ five_degree_south_huai + five_degree_south, data = river_df)
stargazer(new_reg6, new_reg7, new_reg8, new_reg9, new_reg10, type = 'text')
```

```
##
## =====
##                               Dependent variable:
##                               -----
##                               pm10
##                               (1)      (2)      (3)      (4)      (5)
## -----
## one_degree_south_huai      31.223***
##                               (8.602)
##
## one_degree_south           0.413
##                               (0.652)
##
## two_degree_south_huai      30.868***
##                               (8.241)
##
## two_degree_south           0.542
##                               (0.611)
##
## three_degree_south_huai    27.633***
##                               (8.204)
##
## three_degree_south         0.863
##                               (0.581)
##
## four_degree_south_huai     24.965***
##                               (8.304)
##
## four_degree_south          1.141**
##                               (0.554)
##
## five_degree_south_huai     28.610***
##                               (8.340)
##
## five_degree_south          1.072**
##                               (0.531)
##
## Constant                   85.622*** 82.865*** 81.181*** 79.412*** 75.029***
##                               (4.616) (4.676) (4.976) (5.353) (5.527)
##
## -----
## Observations                154      154      154      154      154
## R2                          0.248      0.252      0.239      0.228      0.241
## Adjusted R2                 0.238      0.242      0.229      0.218      0.231
## Residual Std. Error (df = 151) 31.378 31.297 31.555 31.781 31.514
## F Statistic (df = 2; 151)    24.862*** 25.385*** 23.743*** 22.333*** 24.002***
## =====
## Note:                        *p<0.1; **p<0.05; ***p<0.01
```

```
library(dotwhisker)
```



```
## Warning: package 'dotwhisker' was built under R version 4.0.5
```

```
## Warning in checkMatrixPackageVersion(): Package version inconsistency detected.  
## TMB was built with Matrix version 1.3.2  
## Current Matrix version is 1.2.18  
## Please re-install 'TMB' from source using install.packages('TMB', type = 'source') or ask CRA  
N for a binary version of 'TMB' matching CRAN's 'Matrix' package
```

```
## Warning in readRDS(nsInfoFilePath): error reading the file
```

```
orig_tidy <- tidy(pm10_reg)  
  
n_1_tidy <- tidy(new_reg1)  
n_2_tidy <- tidy(new_reg2)  
n_3_tidy <- tidy(new_reg3)  
n_4_tidy <- tidy(new_reg4)  
n_5_tidy <- tidy(new_reg5)  
  
s_1_tidy <- tidy(new_reg6)  
s_2_tidy <- tidy(new_reg7)  
s_3_tidy <- tidy(new_reg8)  
s_4_tidy <- tidy(new_reg9)  
s_5_tidy <- tidy(new_reg10)  
  
joined_reg <- rbind(s_5_tidy, s_4_tidy, s_3_tidy, s_2_tidy, s_1_tidy, orig_tidy, n_1_tidy, n_2_t  
idy, n_3_tidy, n_4_tidy, n_5_tidy)  
  
joined_reg %<>% mutate(term = ifelse(term == "five_degree_south_huai", "-5", term),  
                        term = ifelse(term == "four_degree_south_huai", "-4", term),  
                        term = ifelse(term == "three_degree_south_huai", "-3", term),  
                        term = ifelse(term == "two_degree_south_huai", "-2", term),  
                        term = ifelse(term == "one_degree_south_huai", "-1", term),  
                        term = ifelse(term == "north_huai", "0", term),  
                        term = ifelse(term == "one_degree_north_huai", "+1", term),  
                        term = ifelse(term == "two_degree_north_huai", "+2", term),  
                        term = ifelse(term == "three_degree_north_huai", "+3", term),  
                        term = ifelse(term == "four_degree_north_huai", "+4", term),  
                        term = ifelse(term == "five_degree_north_huai", "+5", term))  
  
dw_plot <- dwplot(joined_reg, ci=.95) + ylim(breaks=c("-5","-4","-3","-2","-1", "0", "+1", "+2",  
"+3", "+4", "+5")) + coord_flip()  
  
dw_plot + geom_vline(xintercept=0,linetype="longdash") + ggtitle("Huai River Placebo Test") +  
  labs(x="PM10 Discontinuity", y = "Degrees North of Huai River")
```

Huai River Placebo Test

