# U6614: Subway Fare Evasion Arrests and Racial Bias

## Sample Solution

### 2022-10-10

*Please submit your knitted .pdf file along with the corresponding R markdown (.rmd) via Courseworks by 11:59pm on Friday, Oct 7 (extension until Sunday, Oct 9).*

# 1 Load libraries

```
library(tidyverse)
library(weights)
library(lmtest)
library(sandwich)
library(knitr)
```

# 2 Aggregating to subway station-level arrest totals

**2a) Load full set of cleaned arrest microdata (arrests.clean.rdata).**
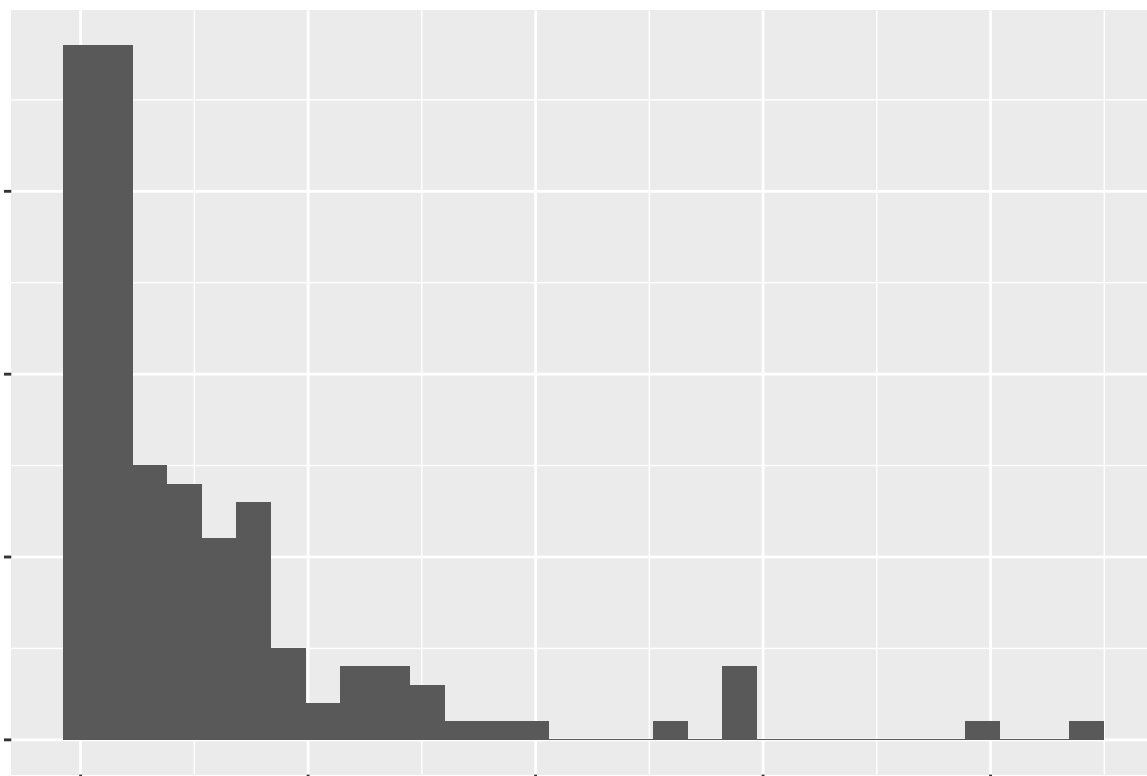
```
load("arrests.clean.RData")
```

**2b) Create new data frame (`st_arrests`) that aggregates microdata to station-level observations including the following information:**

- *st_id, loc2, total arrests*

```
st_arrests <- arrests.clean %>%
  group_by(st_id, loc2) %>%
  summarise(arrests_all = n() ) %>%
  arrange(desc(arrests_all))
```

**2c) Plot histogram of arrests and briefly describe the distribution of arrests across stations.**

```
ggplot(data = st_arrests, aes(x = arrests_all)) +
  geom_histogram()
```

This histogram shows that the majority of subway stations had a relatively small number of fare evasion arrests. The median station arrest total is 13 compared to a mean of 26.82, with 8 stations home to more than 100 arrests.

# 3 Joining subway ridership and neighborhood demographic data

**3a) Import and inspect poverty and ridership data files (`station_povdataclean_2016.csv` and `Subway Ridership by Station - BK.csv`).**

```
st_poverty <- read.csv("station_povdataclean_2016.csv",
                       stringsAsFactors = TRUE)

st_ridership <- read.csv("Subway Ridership by Station - BK.csv",
                         stringsAsFactors = TRUE)

str(st_poverty)
```

```
## 'data.frame':    157 obs. of  10 variables:
##  $ x              : num  -74 -74 -74 -74 -74 ...
##  $ y              : num  40.7 40.6 40.6 40.6 40.6 ...
##  $ st_id          : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ mta_name       : Factor w/ 157 levels "15 St-Prospect Park F subway G subway",..: 1 2 3 4 5 6 7 8 9
##  $ pop_black_2016 : int  306 184 232 202 104 33 237 192 125 841 ...
##  $ pov_black_2016 : int  51 110 161 2 71 1 84 24 41 211 ...
##  $ pop_all_2016   : int  14094 17846 17528 15560 14788 14811 18356 8719 8090 11657 ...
##  $ pov_all_2016   : int  1060 3325 5685 4125 2738 3341 3342 1674 2600 1390 ...
##  $ povrt_all_2016 : num  0.0752 0.1863 0.3243 0.2651 0.1852 ...
##  $ shareblack     : num  0.02171 0.01031 0.01324 0.01298 0.00703 ...
```

```
str(st_ridership)
```

```
## 'data.frame':    157 obs. of  8 variables:
##  $ st_id     : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ mta_name  : Factor w/ 157 levels "15 St-Prospect Park F subway G subway",..: 1 2 3 4 5 6 7 8 9 10
##  $ swipes2011: int  1449241 1654305 1306839 1828462 1370001 1220065 1266146 1243737 3659139 4020665
##  $ swipes2012: int  1852346 1661770 1294568 1819255 1377648 1245145 1206526 1198226 3772420 4300400
##  $ swipes2013: int  2005424 1954729 1313344 1641928 1547029 1286961 1421090 1218842 3826400 4210873
##  $ swipes2014: int  2017347 1996205 1277024 1775695 1564060 1305487 1488192 1274939 3970318 4120946
##  $ swipes2015: int  2011237 1925239 1392552 1784992 1592525 1279525 1508412 1221949 4138929 4199208
##  $ swipes2016: int  1958444 2003777 1344286 723187 1656437 409416 1992945 1188884 4235509 4138758 ..
```

**3b) Join ridership and demographic data to `st_arrests` and inspect results (store new data frame as `st_joined`).**

```
st_joined <- st_arrests %>%
    mutate(st_id = as.integer(st_id)) %>%
    inner_join(st_poverty, by = c("st_id")) %>%
    inner_join(st_ridership, by = c("st_id" = "st_id",
                                    "mta_name" = "mta_name"))
```

- Inspect results from joins, drop unnecessary columns from the ridership data, and group `st_joined` by `st_id` and `mta_name`.
- Only display ungrouped version of `st_joined` for compactness.

```
drop_vars <- c("swipes2011", "swipes2012", "swipes2013", "swipes2014", "swipes2015")

st_joined <- st_joined %>%
  select(-all_of(drop_vars)) %>%
  group_by(st_id, mta_name)

#display structure of ungrouped data frame to avoid lengthy output listing every group
st_joined %>%
  ungroup() %>%
  str()
```

```
## tibble [157 x 13] (S3: tbl_df/tbl/data.frame)
##  $ st_id        : int [1:157] 66 99 150 70 114 131 54 147 106 123 ...
##  $ loc2         : Factor w/ 157 levels "15 st prospect park f g line",..: 66 100 149 148 110 129 54
##  $ arrests_all  : int [1:157] 223 198 143 142 141 141 133 102 90 86 ...
##  $ x            : num [1:157] -74 -74 -73.9 -73.9 -74 ...
##  $ y            : num [1:157] 40.6 40.7 40.7 40.7 40.7 ...
##  $ mta_name     : Factor w/ 157 levels "15 St-Prospect Park F subway G subway",..: 66 99 150 70 114
##  $ pop_black_2016: int [1:157] 36 1939 14825 13135 1542 10311 5624 11804 16176 2698 ...
##  $ pov_black_2016: int [1:157] 2 677 4592 3796 483 2437 900 6706 3832 306 ...
##  $ pop_all_2016 : int [1:157] 5186 12437 18556 17561 23711 15934 6753 15751 20610 13654 ...
##  $ pov_all_2016 : int [1:157] 1329 1939 6149 5565 9182 3511 1156 9104 4809 1221 ...
##  $ povrt_all_2016: num [1:157] 0.256 0.156 0.331 0.317 0.387 ...
##  $ shareblack   : num [1:157] 0.00694 0.15591 0.79893 0.74796 0.06503 ...
##  $ swipes2016   : int [1:157] 5025598 13091255 5152649 9051970 4272443 5861658 3897784 1435112 20317
```

**3c) Print the top 10 stations by total arrest counts**

- Only display `st_id`, `mta_name`, `arrests_all`, `shareblack`, `povrt_all_2016` (no other columns)

```
st_joined %>%
  arrange(desc(arrests_all)) %>%
  select(st_id, mta_name, arrests_all, shareblack, povrt_all_2016) %>%
  mutate(shareblack = round(shareblack, 2),
         povrt_all_2016 = round(povrt_all_2016, 2)) %>%
  head(n = 10)
```

```
## # A tibble: 10 x 5
## # Groups:   st_id, mta_name [10]
##     st_id mta_name                        arrests_all shareblack povrt_all_2016
##     <int> <fct>                                 <int>      <dbl>          <dbl>
## 1     66 "Coney Island-Stillwell Av D sub~        223       0.01           0.26
## 2     99 "Jay St-MetroTech A subway C sub~        198       0.16           0.16
## 3    150 "Utica Av  A subway C subway "           143       0.8            0.33
## 4     70 "Crown Heights-Utica Av 3 subway~        142       0.75           0.32
## 5    114 "Marcy Av  J subway M subway Z s~        141       0.07           0.39
## 6    131 "Nostrand Av  A subway C subway"         141       0.65           0.22
## 7     54 "Canarsie-Rockaway Pkwy L subway"        133       0.83           0.17
## 8    147 "Sutter Av  L subway"                    102       0.75           0.58
## 9    106 "Kingston-Throop Avs C subway"            90       0.78           0.23
## 10   123 "Nevins St  2 subway 3 subway 4 ~         86       0.2            0.09
```

# 4 Explore relationship between arrest intensity and poverty rates across subway station (areas)

**4a) Compute arrest intensity and other explanatory variables for analysis.**

- Drop the observation for the Coney Island station and very briefly explain your logic
- Create new column of data for the following:
  - fare evasion arrest intensity: `arrperswipe_2016` = arrests per 100,000 swipes
  - a dummy indicating if a station is high poverty: `highpov` = 1 if pov rate is > median pov rate across all Brooklyn station areas
  - a dummy for majority Black station areas: `nblack` = 1 if `shareblack` > 0.5
- Coerce new dummy variables into factors with category labels
- Assign results to new data frame called `stations`
- Display top 10 station areas by arrest intensity using `kable()` in the `knitr` package

```
stations <- st_joined %>%
  filter(st_id != 66) %>% #Coney Island
  mutate(arrperswipe = round(arrests_all / (swipes2016 / 100000), 2),
         highpov = as.numeric(povrt_all_2016 > median(st_joined$povrt_all_2016)),
         nblack = as.numeric(shareblack > .5),
         highpov = factor(highpov, levels = c(0,1),
                          labels = c("Not high poverty", "High poverty")),
         nblack = factor(nblack, levels = c(0,1),
                         labels = c("Majority non-Black", "Majority Black")),
         shareblack = round(shareblack, 2),
```

Table 1: Top 10 Stations by Arrest Intensity

| st_id | mta_name | arrperswipe | arrests_all | shareblack | povrt_all_2016 | highpov | nblack |
|---|---|---|---|---|---|---|---|
| 101 | Junius St 3 subway | 11.00 | 75 | 0.78 | 0.48 | High poverty | Majority Black |
| 26 | Atlantic Av L subway | 8.48 | 37 | 0.66 | 0.51 | High poverty | Majority Black |
| 111 | Livonia Av L subway | 7.17 | 75 | 0.83 | 0.45 | High poverty | Majority Black |
| 147 | Sutter Av L subway | 7.11 | 102 | 0.75 | 0.58 | High poverty | Majority Black |
| 106 | Kingston-Throop Avs C subway | 4.43 | 90 | 0.78 | 0.23 | High poverty | Majority Black |
| 112 | Lorimer St J subway M subway | 4.39 | 70 | 0.15 | 0.34 | High poverty | Majority non-Black |
| 140 | Rockaway Av 3 subway | 3.97 | 61 | 0.78 | 0.40 | High poverty | Majority Black |
| 54 | Canarsie-Rockaway Pkwy L subway | 3.41 | 133 | 0.83 | 0.17 | Not high poverty | Majority Black |
| 141 | Rockaway Av C subway | 3.41 | 61 | 0.80 | 0.22 | Not high poverty | Majority Black |
| 144 | Shepherd Av C subway | 3.40 | 36 | 0.61 | 0.30 | High poverty | Majority Black |

```r
        povrt_all_2016 = round(povrt_all_2016, 2))

#display top 10 stations by arrest intensity
#show st_id, mta_name, arrests_all and new variables
  stations_top10 <- stations %>%
    arrange(desc(arrperswipe)) %>%
    select(st_id, mta_name, arrperswipe, arrests_all, shareblack,
           povrt_all_2016, highpov, nblack) %>%
    head(n = 10)

  kable(stations_top10,
        booktabs=TRUE,
        caption = "Top 10 Stations by Arrest Intensity") %>%
    kableExtra::kable_styling(latex_options="scale_down") #scale to fit
```
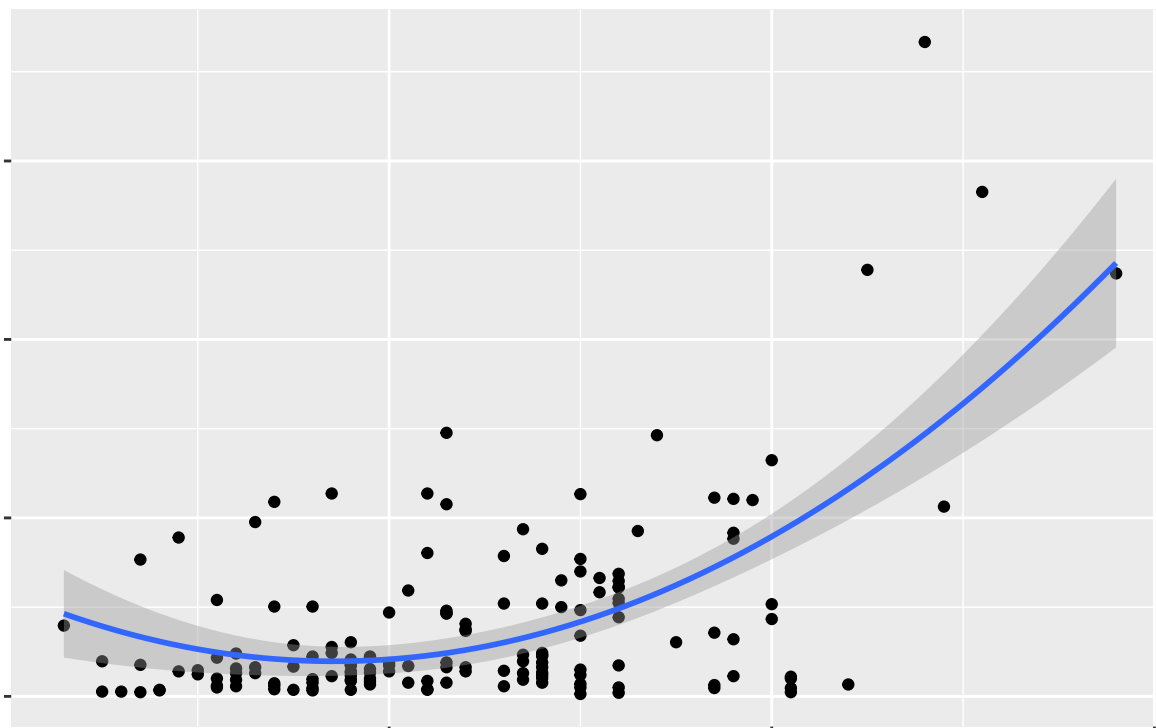
**4b) Examine the relationship between arrest intensity and poverty rates**

- Show a scatterplot of arrest intensity vs. poverty rates along with the regression line you think best fits this relationship.
- Which regression specification do you prefer: linear or quadratic? Be clear about your logic and if applicable cite statistical evidence to support your decision.
- Explain your logic about whether to weight observations or not.

```r
#quadratic
ggplot(stations, #specify data frame to use
       aes(x = povrt_all_2016, y = arrperswipe)) + #specify columns to use
  geom_point() + #specify plot geometry
  ggtitle('Fare evasion arrest intensity vs. poverty rate') + #add title
  labs(x = 'poverty rate',
       y = 'arrests per 100,000 MetroCard swipes') + #change axis labels
  geom_smooth(method = 'lm', formula = y ~ x + I(x^2)) #add regression line
```

```
#linear model (all stations)
ols1l <- lm(arrperswipe ~ povrt_all_2016, data = stations)
summary(ols1l)
coeftest(ols1l, vcov = vcovHC(ols1l, type="HC1")) #get robust SEs

#quadratic model(all stations)
ols1q <- lm(arrperswipe ~ povrt_all_2016 + I(povrt_all_2016^2),
            data = stations)
summary(ols1q)
coeftest(ols1q, vcov = vcovHC(ols1q, type="HC1"))
```

Based on visual inspection, both the linear and quadratic models appear to fit the relationship between fare evasion arrest intensity and poverty rates across all stations fairly well. We prefer the quadratic model because it explains more of the variation in arrest intensity than the linear model; the quadratic model has an adjusted R-squared of 0.36 compared to 0.23 for the linear model. Here we choose not to weight station observations by the number of MetroCard swipes, so that each station area iq equally weighted in the regression analysis. When computing statistics for groups of stations in the next section, we do weight by swipes so that statistics are representative of ridership in each group.

If you prefer the linear specification because it is a bit simpler to interpret without changing the substantive conclusions, that is a reasonable justification.

**4c) Estimate and test difference in mean arrest intensity between high/low poverty areas**

- Report difference and assess statistical significance
- Weight observations by the number of MetroCard swipes

```
stations %>%
  ungroup() %>%
  group_by(highpov) %>%
```

```
  summarise(n = n(),
           mean_pov = weighted.mean(povrt_all_2016, swipes2016),
           mean_arrper = weighted.mean(arrperswipe, swipes2016))
```

```
## # A tibble: 2 x 4
##   highpov              n mean_pov mean_arrper
##   <fct>            <int>    <dbl>       <dbl>
## 1 Not high poverty    79    0.146       0.783
## 2 High poverty        77    0.319        1.42
```

```
#regress arrest intensity on highpov dummy to implement diff in means test
#weighted, robust SEs

ols_diff1 <- lm(formula = arrperswipe ~ highpov, data = stations,
               weights = swipes2016)
ols_diff1_robSE <- coeftest(ols_diff1, vcov = vcovHC(ols_diff1, type="HC1"))
```

The difference in average fare evasion arrest intensity between high- and low-poverty subway stations (weighted by MetroCard swipes) is 0.63 with a p-value of 0.0018. Thus we can conclude that this difference is statistically significant beyond the 1% level.

# 5 How does neighborhood racial composition mediate the relationship between poverty and arrest intensity?

- In this section, you will examine the relationship between arrest intensity & poverty by Black vs. non-Black station area (`nblack`).

**5a) Present a table showing the difference in mean arrests per swipe for each group in a 2x2 table of `highpov` vs `nblack`.**

- Remember to weight by the number of MetroCard swipes at each station
- Could the difference in arrest intensity be explained by differences in poverty rate?

```
t1_arrper_wtd <- with(stations,
                  tapply(arrperswipe * swipes2016,
                        list("High Poverty" = highpov,
                            "Predominantly Black" = nblack),
                        mean)/
                  tapply(swipes2016,
                        list("High Poverty" = highpov,
                            "Predominantly Black" = nblack),
                        mean))

t1_povrt_wtd <- with(stations,
                  tapply(povrt_all_2016 * swipes2016,
                        list("High Poverty" = highpov,
                            "Predominantly Black" = nblack),
                        mean) /
                  tapply(swipes2016,
                        list("High Poverty" = highpov,
                            "Predominantly Black" = nblack),
```

```
                      mean))

round(t1_arrper_wtd, 2)
```

```
##                  Predominantly Black
## High Poverty     Majority non-Black Majority Black
##   Not high poverty               0.66           1.19
##   High poverty                   0.82           2.49
```

```
round(t1_povrt_wtd, 2)
```

```
##                  Predominantly Black
## High Poverty     Majority non-Black Majority Black
##   Not high poverty               0.13           0.19
##   High poverty                   0.32           0.32
```

The above tables show that mean arrests per 100,000 MetroCard swipes are more than 3 times as high at subway stations in majority Black areas compared to non-Black areas. Poverty rates, on the other hand, are very similar between majority-Black and non-Black high-poverty subway station areas, suggesting this is not a likely explanation for the difference in fare evasion arrest intensity (but we can use regression analysis to explore how the relationship between poverty rates and fare evasion differs based on neighborhood racial composition).
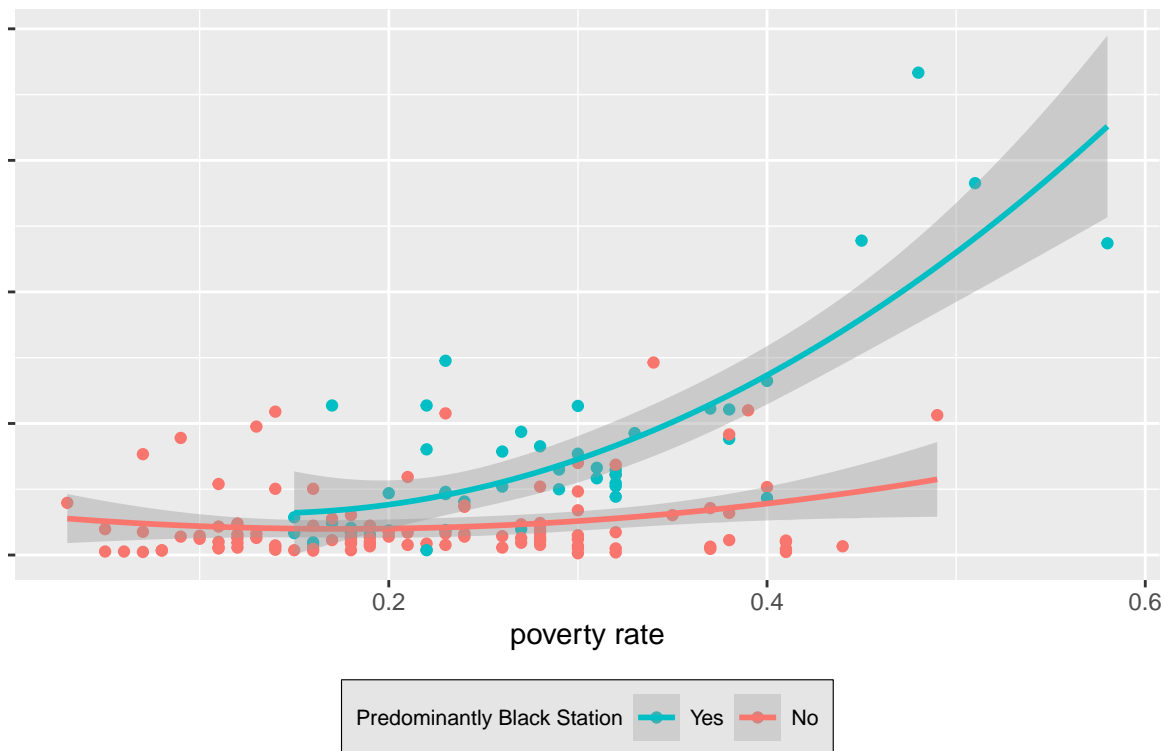
**5b) Show a scatterplot of arrest intensity vs. poverty rates along with the regression line you think best fits this relationship.**

```r
#quadratic
ggplot(stations, aes(x = povrt_all_2016, y = arrperswipe, color = nblack)) +
    geom_point()  +
    geom_smooth(method = 'lm', formula = y ~ x + I(x^2)) +
    ylab("arrests per 100,000 MetroCard swipes") + xlab("poverty rate") +
    ggtitle("Fare evasion arrest intensity vs poverty by race",
            subtitle = "Subway stations in Brooklyn (2016)") +
    scale_color_discrete(name = "Predominantly Black Station",
                         labels=c("No", "Yes"),
                         guide = guide_legend(reverse=TRUE)) +
    theme(legend.position = "bottom",
          legend.background = element_rect(color = "black", fill = "grey90",
                                           size = .2, linetype = "solid"),
          legend.direction = "horizontal",
          legend.text = element_text(size = 8),
          legend.title = element_text(size = 8))
```

8

poverty rate

Predominantly Black Station ──●── Yes ──●── No

```r
#get separate data frames by predominantly Black stations to estimate separate models
stations_black <- stations %>% filter(nblack == "Majority Black")
stations_nonblack <- stations %>% filter(nblack == "Majority non-Black")

#nblack == 1: linear model with station observations
ols_b_l <- lm(arrperswipe ~ povrt_all_2016,
              data = stations_black)

#nblack == 1: quadratic model with station observations
ols_b_q <- lm(arrperswipe ~ povrt_all_2016 + I(povrt_all_2016^2),
              data = stations_black)

#nblack == 0: linear model with station observations
ols_nb_l <- lm(arrperswipe ~ povrt_all_2016,
               data = stations_nonblack)

#nblack == 0: quadratic model with station observations
ols_nb_q <- lm(arrperswipe ~ povrt_all_2016 + I(povrt_all_2016^2),
               data = stations_nonblack)
```

**5c) Which regression specification do you prefer: linear or quadratic? Be clear about your logic and if applicable cite statistical evidence to support your decision.**

Quadratic results are shown here because it explains a greater share of the variation in fare evasion arrest intensity for predominantly Black station areas than the linear model (0.63 compared to 0.58), but the same substantive conclusion holds regardless of functional form.

Visual inspection of the fitted regression lines reveal a clear pattern for both the linear and quadratic specifications: fare evasion arrest intensity increases (at an increasing rate) along with poverty rates at subway stations in predominantly Black areas, but not at other stations. Said another way, the result

suggest that a predominantly Black station area tends to experience significantly higher arrest intensity than a non-Black station with a similarly high poverty rate.

Note that the above interpretation is qualitative in nature: it's a bit more straightforward to provide a numerical interpretation of coefficient estimates with a linear model. Alternatively, it would be informative to compare predicted fare evasion arrest intensity for a predominantly Black station area with a specified poverty rate (say, 40%) compared to a non-Black station area with the same poverty rate. If you prefer the linear specification because it is a bit simpler to interpret without changing the substantive conclusions, that is a reasonable justification.

**5d) Interpret your preferred regression specification (caerfully)!**

For both quadratic and linear models, poverty rates explain very little of the variation in arrest intensity among non-Black station areas in Brooklyn (0.04 and 0.02, respectively).

Regardless of functional form, poverty is only a statistically significant determinant of fare evasion arrest intensity at subway stations in predominantly Black station areas.

# 6 Examine the relationship between arrest intensity and crime

**6a) Read in `nypd_criminalcomplaints_2016.csv`.**

```
st_crime <- read.csv("nypd_criminalcomplaints_2016.csv")
```

**6b) Join stations dataframe to subway station area crime data**

- join on st_id
- exclude the stations with the 4 highest counts of criminal complaints, since they do not face comparable neighbourhood policing conditions

```
stations_wcrime <- stations %>%
  inner_join(st_crime) %>%
  arrange(desc(crimes))

cutoffs <- stations_wcrime %>%
  select(crimes)

#exclude the stations with the 4 highest counts of criminal complaints
stations_wcrime <- stations_wcrime %>%
  filter(crimes < cutoffs$crimes[4])
```

**6c) i. Examine the overall relationship between arrest intensity and crime (without taking neighborhood racial composition or poverty into account) (comparable to Section 4.2).**
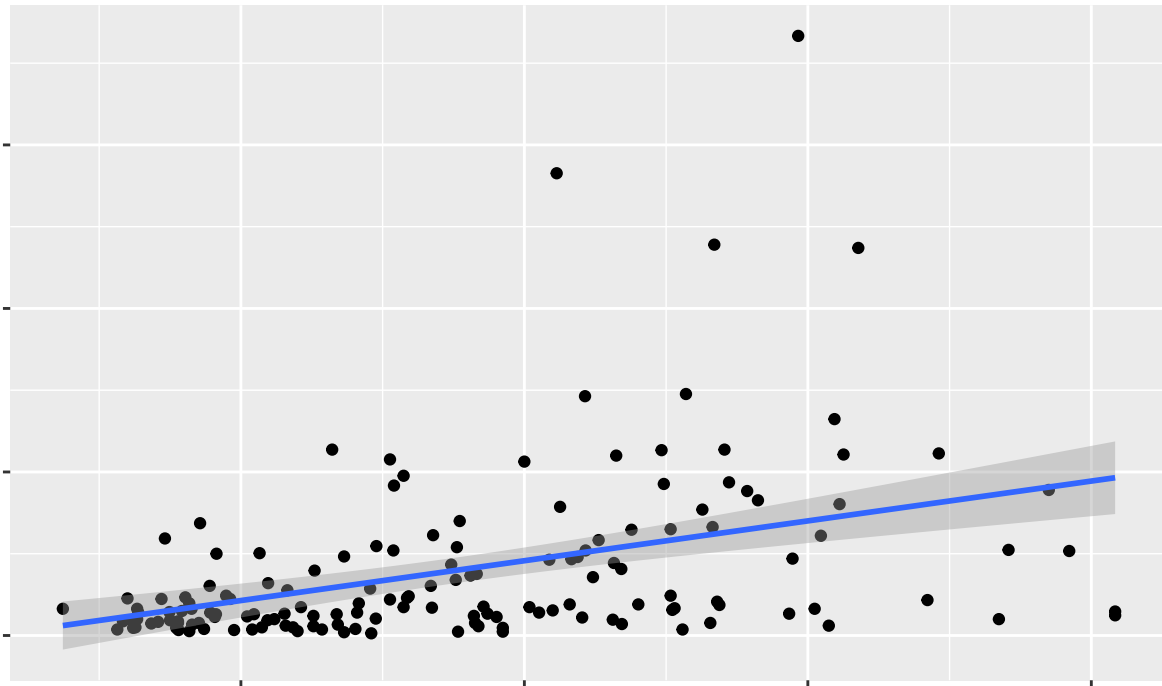
```
#linear
ggplot(stations_wcrime, aes(x = crimes, y = arrperswipe)) +
    geom_point()  +
    geom_smooth(method = 'lm', formula = y ~ x) +
    ylab("arrests per 100,000 MetroCard swipes") + xlab("criminal complaints") +
    ggtitle("Fare evasion arrest intensity vs criminal complaints",
            subtitle = "subway stations in Brooklyn (2016)") +
    scale_color_discrete(name = "Predominantly Black Station",
```

```
                    labels=c("No", "Yes"),
                    guide = guide_legend(reverse=TRUE)) +
    theme(legend.position = "bottom",
        legend.background = element_rect(color = "black", fill = "grey90",
                                    size = .2, linetype = "solid"),
        legend.direction = "horizontal",
        legend.text = element_text(size = 8),
        legend.title = element_text(size = 8))
```
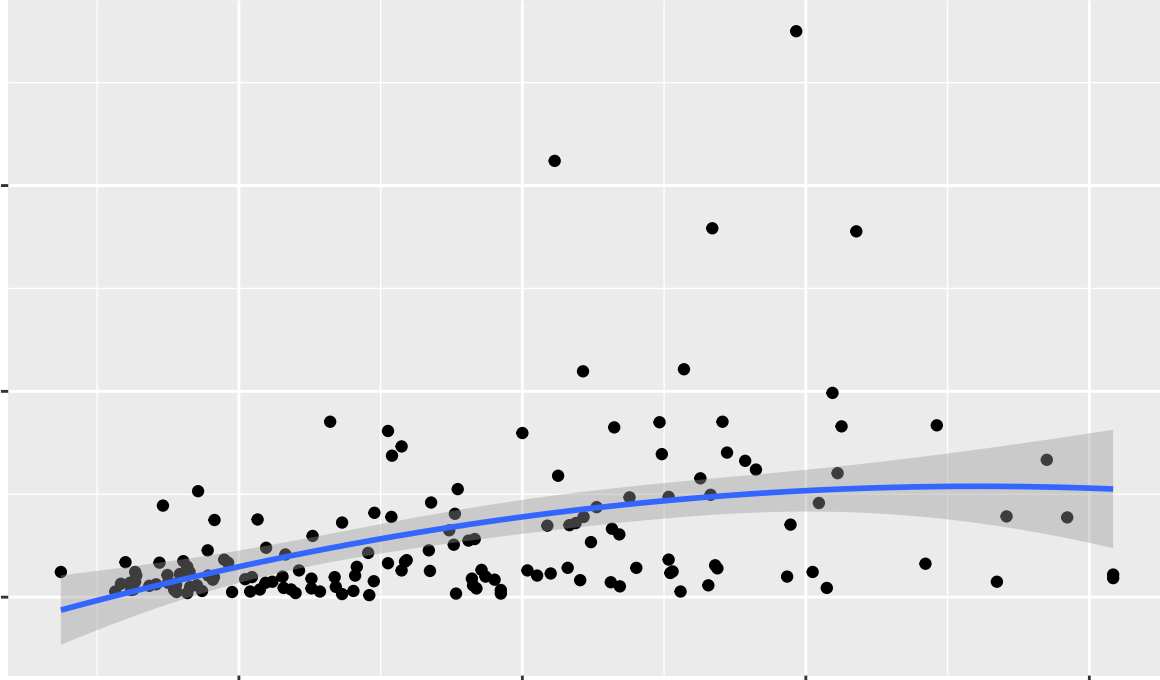


```
#quadratic
ggplot(stations_wcrime, aes(x = crimes, y = arrperswipe)) +
    geom_point()  +
    geom_smooth(method = 'lm', formula = y ~ x + I(x^2)) +
    ylab("arrests per 100,000 MetroCard swipes") + xlab("criminal complaints") +
    ggtitle("Fare evasion arrest intensity vs criminal complaints",
            subtitle = "Subway stations in Brooklyn (2016)") +
    scale_color_discrete(name = "Predominantly Black Station",
                    labels=c("No", "Yes"),
                    guide = guide_legend(reverse=TRUE)) +
    theme(legend.position = "bottom",
        legend.background = element_rect(color = "black", fill = "grey90",
                                    size = .2, linetype = "solid"),
        legend.direction = "horizontal",
        legend.text = element_text(size = 8),
        legend.title = element_text(size = 8))
```

```
ols_c_l <- lm(arrperswipe ~ crimes, data = stations_wcrime)
ols_c_l_robSE <- coeftest(ols_c_l, vcov = vcovHC(ols_c_l, type="HC1")) #get robust SEs

ols_c_q <- lm(arrperswipe ~ crimes + I(crimes^2), data = stations_wcrime)
ols_c_q_robSE <- coeftest(ols_c_q, vcov = vcovHC(ols_c_q, type="HC1")) #get robust SEs
```
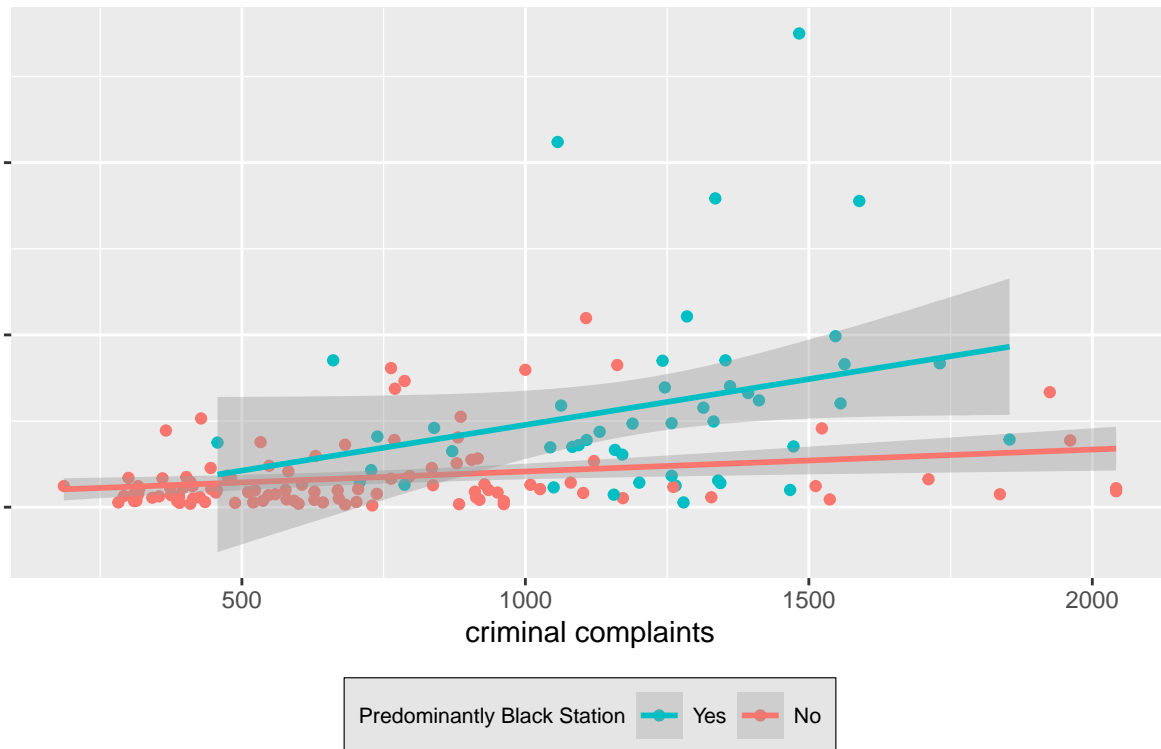
Regardless of the functional form, criminal complaints explain about 16% of the variation in fare evasion arrest intensity across subway stations in Brooklyn (0.166 and 0.156 for quadratic and linear models, respectively).

From the linear model, we can see that the effect of criminal complaints on arrest intensity (0.0015) is statistically significant beyond the 1% level (p-value = 0).

**6c) ii. Examine how neighborhood racial composition mediates the relationship between arrest intensity and crime (comparable to Section 5.2).**

```
#linear
ggplot(stations_wcrime, aes(x = crimes, y = arrperswipe, color = nblack)) +
    geom_point()  +
    geom_smooth(method = 'lm', formula = y ~ x) +
    ylab("arrests per 100,000 MetroCard swipes") + xlab("criminal complaints") +
    ggtitle("Fare evasion arrest intensity vs criminal complaints",
            subtitle = "Subway stations in Brooklyn (2016)") +
    scale_color_discrete(name = "Predominantly Black Station",
                         labels=c("No", "Yes"),
                         guide = guide_legend(reverse=TRUE)) +
    theme(legend.position = "bottom",
          legend.background = element_rect(color = "black", fill = "grey90",
                                           size = .2, linetype = "solid"),
          legend.direction = "horizontal",
          legend.text = element_text(size = 8),
          legend.title = element_text(size = 8))
```

```
#get separate data frames by predominantly Black stations to estimate separate models
stations_wcrime_black <- stations_wcrime %>%
  filter(nblack == "Majority Black")
stations_wcrime_nonblack <- stations_wcrime %>%
  filter(nblack == "Majority non-Black")

#nblack == 1: linear model with station observations
ols_c_b_l <- lm(arrperswipe ~ crimes, data = stations_wcrime_black)
ols_c_b_l_robSE <- coeftest(ols_c_b_l, vcov = vcovHC(ols_c_b_l, type="HC1"))

#nblack == 1: quadratic model with station observations
ols_c_b_q <- lm(arrperswipe ~ crimes + I(crimes^2),
                data = stations_wcrime_black)

#nblack == 0: linear model with station observations
ols_c_nb_l <- lm(arrperswipe ~ crimes, data = stations_wcrime_nonblack)
ols_c_nb_l_robSE <- coeftest(ols_c_nb_l, vcov = vcovHC(ols_c_nb_l, type="HC1"))

#nblack == 0: quadratic model with station observations
ols_c_nb_q <- lm(arrperswipe ~ crimes + I(crimes^2),
                 data = stations_wcrime_nonblack)
```

Estimating separate linear models for the relationship between criminal complaints and arrest intensity for predominantly Black and non-Black station areas reveals a similar pattern as with poverty rates, but less pronounced differences.

Focusing on the linear model for ease of interpretation: the linear relationship between criminal complaints and arrest intensity explains under 6% of the variation regardless of neighborhood racial composition, but the estimated positive effect is four times as large in predominantly Black station areas (0.002 compared to 0.001) and statistically significant at the 5% level (p-value = 0.0127).

# 7   Summarize and interpret your findings with respect to subway fare evasion enforcement bias based on race

- Is there any additional analysis you'd like to explore with the data at hand?
- Are there any key limitations to the data and/or analysis affecting your ability to assess enforcement bias based on race?
- Is there any additional data you'd like to see that would help strengthen your analysis and interpretation?
- For this question, try to be specific and avoid vaguely worded concerns.

The results presented here are consistent with race-based enforcement of fare evasion at subway stations in Brooklyn. As the poverty rate for a subway station area increases, fare evasion arrest intensity tends to increase in predominantly Black station areas (and the association is statistically significant) but not in non-Black station areas.

A similar pattern holds for criminal complaints and fare evasion arrest intensity, though the disparities based on neighborhood racial composition are far less pronounced.

One additional test worth doing is to confirm that the positive association between poverty rates and fare evasion arrest intensity in predominantly Black neighborhoods is still statistically significant when simultaneously controlling for criminal complaints (but not in non-Black neighborhoods). This test confirms that, regardless of where the NYPD enforcement of other crimes is more prevalent, higher poverty Black neighborhoods face considerably higher fare evasion arrests than similarly higher poverty neighborhoods that are not predominantly Black.

The results of this analysis are consistent with disproportionately enforcing fare evasion as a crime of poverty in Black communities; the totality of NYPD policing decisions result in heightened enforcement of fare evasion in higher-poverty, predominantly Black neighborhoods. This analysis does not, however, inform the relative importance of different mechanisms driving these patterns: do police deployment decision explain these disparities, implicit and/or explicit bias in who is stopped and what enforcement action is taken (arrest vs summons), or some combination of these mechanisms? There may also be other differences in subway rider characteristics and behavior that could explain the observed relationship between neighborhood racial composition and fare evasion enforcement intensity, but disparate impact by race is clear even if the all of the underlying mechanisms are not.

Analyzing differences in fare evasion summonses compared to arrests would also be informative: are there significant differences in the demographics of individuals who are stopped for fare evasion, in addition to differences in the enforcement action taken once they are stopped?