

U6614: Assignment 3: Subway Fare Evasion Microdata

Sample Solution

2020-09-27

Please submit your knitted .pdf file along with the corresponding R markdown (.rmd) via Courseworks by 11:59pm on Monday, September 28th.

Before knitting your rmd file as a pdf, you will need to install TinyTex for Latex distribution by running the following code:

```
tinytex::install_tinytex()
```

Please visit [this](#) link for more information on TinyTex installation.

1 Load libraries

```
#remember to make sure these packaged are installed before trying to load
library(tidyverse)
library(fastDummies)
```

2 Load and inspect the two public defender client datasets (BDS & LAS)

```
arrests_bds <- read_csv("microdata_BDS_inclass.csv", na = "")
arrests_las <- read_csv("microdata_LAS_inclass.csv", na = "")
```

```
str(arrests_bds)
```

```
## tibble [2,246 x 8] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ client_zip: num [1:2246] 11205 11385 11226 11207 11225 ...
## $ age       : num [1:2246] 25 20 19 17 21 52 59 32 22 19 ...
## $ ethnicity : chr [1:2246] "Hispanic" "Hispanic" "Non-Hispanic" "Non-Hispanic" ...
## $ race      : chr [1:2246] "White" "Black" "Black" "Black" ...
## $ male      : num [1:2246] 1 1 0 1 1 1 1 0 1 ...
## $ loc2      : chr [1:2246] "jefferson st l line station" "myrtle - wyckoff avs station" "winthrop s
## $ st_id     : num [1:2246] 100 119 156 156 156 156 156 156 156 ...
## $ year      : num [1:2246] 2016 2016 2016 2016 2016 ...
## - attr(*, "spec")=
## .. cols(
## ..   client_zip = col_double(),
## ..   age = col_double(),
## ..   ethnicity = col_character(),
## ..   race = col_character(),
## ..   male = col_double(),
```

```
## .. loc2 = col_character(),
## .. st_id = col_double(),
## .. year = col_double()
## .. )

str(arrests_las)

## tibble [1,965 x 9] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ client_zip : num [1:1965] 11222 10016 11236 11236 NA ...
## $ las_race_key : chr [1:1965] "Black" "Asian or Pacific Islander" "Black" "Black" ...
## $ hispanic_flag: chr [1:1965] "N" "N" "N" "N" ...
## $ age : num [1:1965] 32 47 20 64 23 29 26 52 52 22 ...
## $ year : num [1:1965] 2016 2016 2016 2016 2016 ...
## $ male : num [1:1965] 1 0 1 1 1 1 0 1 1 1 ...
## $ dismissal : num [1:1965] 0 1 0 0 0 0 1 0 0 1 ...
## $ loc2 : chr [1:1965] "kingston - throop avs" "avenue h q subway" "nostrand ave and fulton s" ...
## $ st_id : num [1:1965] 106 28 131 150 131 27 68 44 85 31 ...
## - attr(*, "spec")=
## .. cols(
## .. client_zip = col_double(),
## .. las_race_key = col_character(),
## .. hispanic_flag = col_character(),
## .. age = col_double(),
## .. year = col_double(),
## .. male = col_double(),
## .. dismissal = col_double(),
## .. loc2 = col_character(),
## .. st_id = col_double()
## .. )
```

The BDS data includes 2246 observations (client arrest records), and the LAS data includes another 1965 observations. Both datasets include basic demographic information on age, sex, race, ethnicity (coded differently in each dataset), as well as information on the location/subway station where the arrest occurred.

The LAS data also includes information on case dismissal rates.

2.1 For each dataset, what is the unit of observation and representative population?

In each raw dataset, the unit of observation is the arrested individual (client). On the surface, the representative population is all individuals arrested by the NYPD for subway fare evasion in Brooklyn during 2016 who are represented by public defenders. If nearly all individuals arrested for fare evasion are represented by public defenders, then this sample comes close to constituting the universe of subway fare evasion arrests in Brooklyn in 2016. This is difficult to argue convincingly without additional information, but is supported anecdotally by court observers.

2.2 Inspect and describe the coding of race/ethnicity in each dataset.

```
#recode race/ethnicity information from character to factors
arrests_bds <- arrests_bds %>% mutate(race = as.factor(race),
                                     ethnicity = as.factor(ethnicity) )
arrests_las <- arrests_las %>% mutate(race = as.factor(las_race_key),
                                     ethnicity = as.factor(hispanic_flag) )

#compare race coding
summary(arrests_bds$race)
```

```
##           0           Am Indian Asian/Pacific Islander
##           35           1           21
##           Black           Other           Unknown
##           1465           32           2
##           White           NA's
##           533           157
```

```
summary(arrests_las$race)
```

```
## Asian or Pacific Islander           Black           Hispanic
##           11           1247           21
##           Latino           Other           Unknown
##           2           20           10
##           White           NA's
##           426           228
```

```
#compare Hispanic/ethnicity coding
```

```
summary(arrests_bds$ethnicity)
```

```
##           0           Hispanic Non-Hispanic           Other           NA's
##           33           493           1558           5           157
```

```
summary(arrests_las$ethnicity)
```

```
##           N           Y NA's
## 1619  189  157
```

Race information is generally stored in one variable, Hispanic identity in a second variable. To work towards consistent variable names and coding in both datasets, let's first recode the raw race and ethnicity information into two separate columns of data (factors) named `race` and `ethnicity`.

While each dataset refers to similar race and ethnicity categories, there are different category names in each (including some slightly different spellings).

We also note that Hispanic identity factors into both race and Hispanic variables in the Legal Aid Society (LAS) data; in the BDS data, information on Hispanic identity is only included in the ethnicity variable.

Each dataset also contains a different set of values that seem to convey unknown race/ethnicity information, in addition to true missings (e.g. "0" and "Unknown" in addition to blank entries).

2.3 From the outset, are there any data limitations you think are important to note?

It's unclear what processes are used to code race and ethnicity at each public defender group. How much does the information reflect client self-identification rather than identity assigned by police and entered into arrest reports?

It's also important to emphasize what information this data does **not** include that might be relevant to the question of biased fare evasion enforcement:

- fare evasion that resulted in a summons (ticket + fine) rather than an arrest
- fare evasion enforcement on buses

3 Clean BDS race and ethnicity data

3.1 BDS: race data (generate column race_clean).

```
#identify every combination of race-ethnicity in the raw data
table(arrests_bds$race, arrests_bds$ethnicity, useNA = "always")

##
##           0 Hispanic Non-Hispanic Other <NA>
## 0           31         1          3      0    0
## Am Indian    0         0          1      0    0
## Asian/Pacific Islander 0         0         21     0    0
## Black        2        104        1358     1    0
## Other        0         20         11      1    0
## Unknown      0         0          0      2    0
## White        0        368        164     1    0
## <NA>         0         0          0      0   157

#recode as factor in an internally consistent manner (address NAs, specify levels)
arrests_bds.clean <- arrests_bds %>%
  mutate(race_clean = recode(race, "0" = "NA",
                              "Unknown" = "NA",
                              "Am Indian" = "Other" ) ) %>%
  mutate(race_clean = factor(race_clean,
                            levels = c("Black", "White", "Asian/Pacific Islander", "Other")))

#validation: confirm the recode worked as intended
levels(arrests_bds.clean$race_clean)

## [1] "Black"          "White"          "Asian/Pacific Islander"
## [4] "Other"

arrests_bds.clean %>% count(race_clean, sort = TRUE)

## # A tibble: 5 x 2
##   race_clean      n
##   <fct>         <int>
## 1 Black        1465
## 2 White         533
## 3 <NA>         194
## 4 Other         33
## 5 Asian/Pacific Islander 21

table(arrests_bds.clean$race_clean, arrests_bds.clean$race, useNA = "always")

##
##           0 Am Indian Asian/Pacific Islander Black Other
## Black        0         0          0    1465     0
## White        0         0          0      0     0
## Asian/Pacific Islander 0         0         21     0     0
## Other        0         1          0      0    32
## <NA>         35         0          0      0     0
##
##           Unknown White <NA>
## Black        0      0     0
## White        0    533     0
```

```
## Asian/Pacific Islander      0      0      0
## Other                       0      0      0
## <NA>                        2      0    157
```

3.2 BDS: ethnicity data (generate column ethnicity_clean).

```
#ok now let's recode to Hispanic, Non-Hispanic, and NA
arrests_bds.clean <- arrests_bds.clean %>%
  mutate(hispanic = recode(ethnicity, "0" = "NA", "Other" = "Non-Hispanic") ) %>%
  mutate(hispanic = factor(hispanic, levels = c("Hispanic", "Non-Hispanic")))

#validation: confirm the recode worked as intended
summary(arrests_bds.clean$hispanic)
```

```
## Hispanic Non-Hispanic NA's
## 493 1563 190
```

```
table(arrests_bds.clean$race_clean, arrests_bds.clean$hispanic, useNA = "always")
```

```
##
## Hispanic Non-Hispanic <NA>
## Black 104 1359 2
## White 368 165 0
## Asian/Pacific Islander 0 21 0
## Other 20 13 0
## <NA> 1 5 188
```

3.3 Generate a single race/ethnicity factor variable race_eth with mutually exclusive categories.

```
#let's investigate a bit
table(arrests_bds.clean$race_clean, arrests_bds.clean$hispanic, useNA = "always")
```

```
##
## Hispanic Non-Hispanic <NA>
## Black 104 1359 2
## White 368 165 0
## Asian/Pacific Islander 0 21 0
## Other 20 13 0
## <NA> 1 5 188
```

```
prop.table(table(arrests_bds.clean$race_clean, arrests_bds.clean$hispanic), 2) %>% round(2)
```

```
##
## Hispanic Non-Hispanic
## Black 0.21 0.87
## White 0.75 0.11
## Asian/Pacific Islander 0.00 0.01
## Other 0.04 0.01
```

```
#generate race_eth column (as a factor) in steps
arrests_bds.clean <- arrests_bds.clean %>%
  mutate(race_clean_char = as.character(race_clean)) %>% #work with characters
  mutate(hispanic_char = as.character(hispanic)) %>% #work with characters
  mutate(race_eth = ifelse(hispanic_char == "Hispanic",
                           hispanic_char,
```

```

      race_clean_char) ) %>%
mutate(race_eth = as.factor(recode(race_eth, "White" = "Non-Hispanic White"))) %>%
select(-race_clean_char, -hispanic_char)

```

#validate results: joint distribution of race_eth and hispanic

```
table(arrests_bds.clean$race_eth, arrests_bds.clean$hispanic, useNA = "always")
```

```
##
##              Hispanic Non-Hispanic <NA>
## Asian/Pacific Islander      0      21    0
## Black                      0     1359    0
## Hispanic                   493       0    0
## Non-Hispanic White         0      165    0
## Other                      0       13    0
## <NA>                       0       5   190
```

```
summary(arrests_bds.clean$race_eth, useNA = "always")
```

```
## Asian/Pacific Islander      Black      Hispanic
##              21          1359          493
## Non-Hispanic White      Other      NA's
##              165          13          195
```

Note that `race_eth` assigns individuals who identify as both Hispanic and a race other than white as Hispanic. This means, for example, that an individual who identifies as both Black and Hispanic appears as Hispanic in the `race_eth` column.

4 Clean LAS race and ethnicity data

4.1 Follow your own steps to end up at a `race_eth` variable for the LAS data that is coded in a comparable manner as in the BDS data.

NOTE: you may be able to do everything in a single pipe, depending on your approach (but you certainly don't have to).

#inspect LAS data

```
str(arrests_las)
```

```
## tibble [1,965 x 11] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ client_zip : num [1:1965] 11222 10016 11236 11236 NA ...
## $ las_race_key : chr [1:1965] "Black" "Asian or Pacific Islander" "Black" "Black" ...
## $ hispanic_flag: chr [1:1965] "N" "N" "N" "N" ...
## $ age : num [1:1965] 32 47 20 64 23 29 26 52 52 22 ...
## $ year : num [1:1965] 2016 2016 2016 2016 2016 ...
## $ male : num [1:1965] 1 0 1 1 1 1 0 1 1 1 ...
## $ dismissal : num [1:1965] 0 1 0 0 0 0 1 0 0 1 ...
## $ loc2 : chr [1:1965] "kingston - throop avs" "avenue h q subway" "nostrand ave and fulton s" ...
## $ st_id : num [1:1965] 106 28 131 150 131 27 68 44 85 31 ...
## $ race : Factor w/ 7 levels "Asian or Pacific Islander",...: 2 1 2 2 2 2 2 2 7 ...
## $ ethnicity : Factor w/ 2 levels "N","Y": 1 1 1 1 1 1 1 1 1 ...
## - attr(*, "spec")=
## .. cols(
## .. client_zip = col_double(),
## .. las_race_key = col_character(),
```

```
## .. hispanic_flag = col_character(),
## .. age = col_double(),
## .. year = col_double(),
## .. male = col_double(),
## .. dismissal = col_double(),
## .. loc2 = col_character(),
## .. st_id = col_double()
## .. )
```

```
table(arrests_las$las_race_key, arrests_las$hispanic_flag, useNA = "always")
```

```
##
##              N      Y <NA>
## Asian or Pacific Islander    11    0    0
## Black                      1201   46    0
## Hispanic                    20    1    0
## Latino                      2    0    0
## Other                       11    9    0
## Unknown                     10    0    0
## White                       294  132    0
## <NA>                        70    1  157
```

```
#generate race_eth column as a factor
```

```
arrests_las.clean <- arrests_las %>%
```

```
  mutate(race_eth = recode(las_race_key, "Asian or Pacific Islander" = "Asian/Pacific Islander",
                                "Unknown" = "NA",
                                "Latino" = "Hispanic",
                                "White" = "Non-Hispanic White")) %>%
```

```
  mutate(race_eth = ifelse(hispanic_flag == "Y", "Hispanic", race_eth) ) %>%
```

```
  mutate(race_eth = factor(race_eth,
```

```
    levels = c("Black", "Hispanic", "Non-Hispanic White", "Asian/Pacific Islander")
```

```
#validate
```

```
summary(arrests_las.clean$race_eth)
```

```
##              Black              Hispanic      Non-Hispanic White
##              1201              211              294
## Asian/Pacific Islander              Other              NA's
##              11              11              237
```

```
arrests_las.clean %>% count(race_eth, sort = TRUE)
```

```
## # A tibble: 6 x 2
```

```
##   race_eth      n
##   <fct>      <int>
## 1 Black      1201
## 2 Non-Hispanic White    294
## 3 <NA>        237
## 4 Hispanic     211
## 5 Asian/Pacific Islander   11
## 6 Other        11
```

```
table(arrests_las.clean$race_eth, arrests_las.clean$hispanic_flag, useNA = "always")
```

```
##
##              N      Y <NA>
## Black      1201    0    0
```

```
##   Hispanic                22 189    0
##   Non-Hispanic White      294    0    0
##   Asian/Pacific Islander   11    0    0
##   Other                    11    0    0
##   <NA>                     80    0  157
```

5 Combining (appending) the BDS and LAS microdata

5.1 Create a column (pd) to identify public defender data source.

```
arrests_bds.clean <- arrests_bds.clean %>% mutate(pd = "bds")
arrests_las.clean <- arrests_las.clean %>% mutate(pd = "las")
```

5.2 Append arrests_bds.clean and arrests_las.clean using rbind(). Store as new data frame arrests_all and inspect for consistency/accuracy.

```
arrests.clean <- plyr::rbind.fill(arrests_las.clean, arrests_bds.clean) %>%
  mutate(pd = as.factor(pd),
         st_id = as.factor(st_id),
         loc2 = as.factor(loc2)) %>% #station/location info is not continuous
  select(pd, race_eth, age, male, dismissal, st_id, loc2)
str(arrests.clean)
```

```
## 'data.frame':   4211 obs. of  7 variables:
## $ pd          : Factor w/ 2 levels "bds","las": 2 2 2 2 2 2 2 2 2 ...
## $ race_eth    : Factor w/ 5 levels "Black","Hispanic",...: 1 4 1 1 1 1 1 1 3 ...
## $ age         : num  32 47 20 64 23 29 26 52 52 22 ...
## $ male        : num  1 0 1 1 1 1 0 1 1 1 ...
## $ dismissal   : num  0 1 0 0 0 0 1 0 0 1 ...
## $ st_id       : Factor w/ 157 levels "1","2","3","4",...: 106 28 131 150 131 27 68 44 85 31 ...
## $ loc2        : Factor w/ 157 levels "15 st prospect park f g line",...: 104 30 129 149 129 28 68 44 85 ...
```

```
summary(arrests.clean)
```

##	pd	race_eth	age	male	
##	bds:2246	Black	:2560	Min. : 0.00	
##	las:1965	Hispanic	: 704	1st Qu.:20.00	
##		Non-Hispanic White	: 459	Median :26.00	
##		Asian/Pacific Islander	: 32	Mean :29.18	
##		Other	: 24	3rd Qu.:35.00	
##		NA's	: 432	Max. :71.00	
##			NA's	:317	
##				NA's	:314
##	dismissal	st_id		loc2	
##	Min. :0.0000	66 : 223	coney island-stillwell ave	: 223	
##	1st Qu.:0.0000	99 : 198	jay st - metrotech	: 198	
##	Median :1.0000	150 : 143	utica ave and fulton st	: 143	
##	Mean :0.5392	70 : 142	utica ave and eastern parkway	: 142	
##	3rd Qu.:1.0000	114 : 141	marcy ave j m z line	: 141	
##	Max. :1.0000	131 : 141	nostrand ave and fulton st a c station	: 141	
##	NA's :2529	(Other):3223	(Other)	:3223	

5.3 What is the total number of subway fare evasion arrest records?

The total number of subway fare evasion arrest records from both BDS and LAS is 4211.

5.4 Export `arrests_all` as `.csv`, and save as `.rds` file.

```
write_csv(arrests.clean, "arrests_all.csv")
saveRDS(arrests.clean, "../Lecture4/arrests.clean.rds")
```

6 Descriptive statistics by race/ethnicity

6.1 Print the number of arrests for each race/ethnicity category (a frequency table).

```
arrests.clean %>% count(race_eth, sort = TRUE)
```

```
##           race_eth    n
## 1           Black 2560
## 2          Hispanic  704
## 3 Non-Hispanic White  459
## 4              <NA>  432
## 5 Asian/Pacific Islander   32
## 6              Other    24
```

6.2 Print the proportion of total arrests for each race/ethnicity category.

```
#including NAs
prop.table(table(arrests.clean$race_eth, useNA = "always")) %>%
  round(2) %>%
  as.data.frame() %>%
  arrange(desc(Freq)) %>%
  rename(race_eth = Var1)
```

```
##           race_eth Freq
## 1           Black 0.61
## 2          Hispanic 0.17
## 3 Non-Hispanic White 0.11
## 4              <NA> 0.10
## 5 Asian/Pacific Islander 0.01
## 6              Other 0.01
```

```
#excluding NAs
prop.table(table(arrests.clean$race_eth)) %>%
  round(2) %>%
  as.data.frame() %>%
  arrange(desc(Freq)) %>%
  rename(race_eth = Var1)
```

```
##           race_eth Freq
## 1           Black 0.68
## 2          Hispanic 0.19
## 3 Non-Hispanic White 0.12
## 4 Asian/Pacific Islander 0.01
```

```
## 5          Other 0.01
```

6.3 Show the average age, share male, and dismissal rate for each race/ethnicity category. Describe any noteworthy findings.

```
arrests.clean %>%
  group_by(race_eth) %>%
  summarise(n = n(),
            mean_age = mean(age, na.rm = TRUE),
            mean_male = mean(male, na.rm = TRUE),
            mean_dism = mean(dismissal, na.rm = TRUE))
```

```
## # A tibble: 6 x 5
##   race_eth          n mean_age mean_male mean_dism
##   <fct>          <int>   <dbl>   <dbl>   <dbl>
## 1 Black          2560    29.1    0.875    0.514
## 2 Hispanic        704    29.7    0.901    0.538
## 3 Non-Hispanic White 459    29.7    0.898    0.587
## 4 Asian/Pacific Islander 32    28.9    0.938    0.636
## 5 Other           24    28.3    0.833    0.444
## 6 <NA>           432    25.9    0.610    0.75
```

7 Subway-station level analysis

7.1 Create dummy variables for each race/ethnicity category and show summary statistics only for these dummy variables.

```
arrests.clean <- dummy_cols(arrests.clean, select_columns = "race_eth")
summary(arrests.clean[,8:13])
```

```
## race_eth_Black  race_eth_Hispanic race_eth_Non-Hispanic White
## Min. :0.0000   Min. :0.0000   Min. :0.0000
## 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000
## Median :1.0000 Median :0.0000 Median :0.0000
## Mean :0.6774   Mean :0.1863   Mean :0.1215
## 3rd Qu.:1.0000 3rd Qu.:0.0000 3rd Qu.:0.0000
## Max. :1.0000   Max. :1.0000   Max. :1.0000
## NA's :432      NA's :432      NA's :432
## race_eth_Asian/Pacific Islander race_eth_Other  race_eth_NA
## Min. :0.0000                   Min. :0.0000   Min. :0.0000
## 1st Qu.:0.0000                 1st Qu.:0.0000 1st Qu.:0.0000
## Median :0.0000                 Median :0.0000 Median :0.0000
## Mean :0.0085                   Mean :0.0064   Mean :0.1026
## 3rd Qu.:0.0000                 3rd Qu.:0.0000 3rd Qu.:0.0000
## Max. :1.0000                   Max. :1.0000   Max. :1.0000
## NA's :432                     NA's :432
```

7.2 Aggregate to station-level observations (group by loc2), and show a table of stations with at least 50 arrests along with the following information:

- station name (loc2)

- station arrest total
- share of arrests that are Black and Hispanic (excluding race_eth = NA from denominator)
- sorted in ascending order above Black and Hispanic arrest share
- remember to only show stations with at least 50 total arrests

```
arrests_stations_top <- arrests_clean %>%
  group_by(loc2) %>%
  summarise(st_id = first(st_id),
            n = n(),
            n_black = sum(race_eth_Black, na.rm = TRUE),
            n_hisp = sum(race_eth_Hispanic, na.rm = TRUE),
            n_api = sum(`race_eth_Asian/Pacific Islander`, na.rm = TRUE),
            n_nhw = sum(`race_eth_Non-Hispanic White`, na.rm = TRUE),
            n_oth = sum(race_eth_Other, na.rm = TRUE),
            n_bh = sum(race_eth_Black, race_eth_Hispanic, na.rm = TRUE),
            n_na = sum(race_eth_NA)) %>%
  mutate(sh_bh = round(n_bh / (n - n_na), 2)) %>%
  filter(n >= 50) %>%
  arrange(sh_bh)
knitr::kable(arrests_stations_top) #kable functions in knitr package are good for formatted tables
```

loc2	st_id	n	n_black	n_hisp	n_api	n_nhw	n_oth	n_bh	n_na	sh_bh
marcy ave j m z line	114	141	55	42	3	34	0	97	7	0.72
myrtle av and broadway station	117	69	38	15	0	13	0	53	3	0.80
coney island-stillwell ave	66	223	124	48	5	35	1	172	10	0.81
graham ave l line	88	54	28	11	0	9	0	39	6	0.81
broadway and lorimer st j m station	112	70	34	22	0	11	1	56	2	0.82
clinton - washington avs station	64	63	42	6	0	10	0	48	5	0.83
jay st - metrotech	99	198	112	43	3	29	0	155	11	0.83
hozt-schermerhorn a c g line	98	71	46	9	0	10	0	55	6	0.85
myrtle - willoughby avs g line	118	50	27	12	0	5	1	39	5	0.87
canarsie rockaway pkwy	54	133	109	4	1	11	2	113	6	0.89
nevins st 2 3 4 5 lines	123	86	63	11	0	6	1	74	5	0.91
hozt st 2 3	97	77	58	12	0	5	0	70	2	0.93
kingston - throop avs	106	90	69	12	0	6	0	81	3	0.93
nostrand ave and fulton st a c station	131	141	107	20	0	7	1	127	6	0.94
sutter avenue station l line	147	102	79	12	0	6	0	91	5	0.94
utica ave and fulton st	150	143	111	19	0	7	0	130	6	0.95
court st r subway/borough hall 2	68	59	42	11	0	2	0	53	4	0.96
subway 3 subway 4 subway 5 subway										
junius st 3 line	101	75	60	10	1	2	0	70	2	0.96
livonia ave l line	111	75	56	13	0	3	0	69	3	0.96
utica ave and eastern parkway	70	142	118	13	0	5	0	131	6	0.96
rockaway ave c line	141	61	50	7	0	1	0	57	3	0.98
sutter av - rutland rd 3 line	148	68	61	3	0	0	1	64	3	0.98
rockaway ave 3 line	140	61	49	8	0	0	0	57	4	1.00

7.3 Briefly summarize any noteworthy findings from the table you just generated.

At every single high-arrest subway station, the majority of arrested individuals are Black or Hispanic. This isn't surprising, given that 87 percent of *all* arrested individuals with coded race/ethnicity are Black or Hispanic.

7.4 (OPTIONAL) Visualize the distribution of arrests by race/ethnicity at stations with > 100 arrests.

Hint: see R code from class, section 8

#get data frame with obs for every station-race_eth pairings on arrest counts

```
arrests_stations_race <- arrests.clean %>%
  group_by(loc2) %>%
  mutate(st_arrests = n()) %>%
  ungroup() %>%
  group_by(loc2, race_eth) %>%
  summarise(arrests = n(), st_arrests = first(st_arrests)) %>%
  arrange(desc(st_arrests)) %>%
  filter(st_arrests > 100)
arrests_stations_race
```

```
## # A tibble: 39 x 4
## # Groups:   loc2 [8]
##   loc2                race_eth      arrests st_arrests
##   <fct>              <fct>         <int>     <int>
## 1 coney island-stillwell ave Black      124       223
## 2 coney island-stillwell ave Hispanic    48       223
## 3 coney island-stillwell ave Non-Hispanic White 35       223
## 4 coney island-stillwell ave Asian/Pacific Islander 5       223
## 5 coney island-stillwell ave Other        1       223
## 6 coney island-stillwell ave <NA>        10       223
## 7 jay st - metrotech    Black     112       198
## 8 jay st - metrotech    Hispanic    43       198
## 9 jay st - metrotech    Non-Hispanic White 29       198
## 10 jay st - metrotech    Asian/Pacific Islander 3       198
## # ... with 29 more rows
```

```
ggplot(arrests_stations_race,
  aes(x = reorder(loc2, -st_arrests), y = arrests, fill = race_eth)) +
  geom_bar(stat = "identity") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))
```

