

# Implications of Missing Traffic Stop Data

Amber Lee advised by Jo Hardin  
Pomona College

November 5, 2021



# Motivation



Concern over racial profiling has led to mass collection of traffic stop data. As the data has been made available to the public, researchers have modeled the data to look for evidence of discrimination.

# Motivation



Concern over racial profiling has led to mass collection of traffic stop data. As the data has been made available to the public, researchers have modeled the data to look for evidence of discrimination.

- Most statistical models do not handle missing values, so most studies simply use complete-case analysis (`na.rm = True`).
- Chanin and Welsh (2020) draw attention to the issue of (non-random) missing values.

# Motivation



Concern over racial profiling has led to mass collection of traffic stop data. As the data has been made available to the public, researchers have modeled the data to look for evidence of discrimination.

- Most statistical models do not handle missing values, so most studies simply use complete-case analysis (`na.rm = True`).
- Chanin and Welsh (2020) draw attention to the issue of (non-random) missing values.

⇒ **What are the trends of missing values in traffic stop data?**

# Row-wise missingness

Let  $X = (x_{ij})$  be a **dataset** with  $n$  observations and  $k$  variables.

- Each observation  $i = 1, 2, \dots, n$  represents a single traffic stop.
- We denote a **missingness indicator function** as

$$\mathbb{1}_M(x_{ij}) = \begin{cases} 1 & \text{if } x_{ij} \text{ is missing} \\ 0 & \text{if } x_{ij} \text{ is observed} \end{cases}$$

# Row-wise missingness

Let  $X = (x_{ij})$  be a **dataset** with  $n$  observations and  $k$  variables.

- Each observation  $i = 1, 2, \dots, n$  represents a single traffic stop.
- We denote a **missingness indicator function** as

$$\mathbb{1}_M(x_{ij}) = \begin{cases} 1 & \text{if } x_{ij} \text{ is missing} \\ 0 & \text{if } x_{ij} \text{ is observed} \end{cases}$$

## Definition

The **stop missingness rate for observation**  $i$  (or row-wise SMR) is the percentage of missing values for a single traffic stop  $i$ .

$$\text{SMR}_i = \frac{1}{k} \sum_{j=1}^k \mathbb{1}_M(x_{ij})$$

# Dataset missingness



## Definition

The **stop missingness rate for dataset**  $X$  (or dataset SMR) is the percentage of missing values in the entire dataset. Equivalently, it is the average row-wise missingness.

$$\text{SMR}(X) = \frac{1}{n} \sum_{i=1}^n \text{SMR}_i$$

# Dataset missingness



## Definition

The **stop missingness rate for dataset**  $X$  (or dataset SMR) is the percentage of missing values in the entire dataset. Equivalently, it is the average row-wise missingness.

$$\text{SMR}(X) = \frac{1}{n} \sum_{i=1}^n \text{SMR}_i$$

## Example

$$\text{SMR}(X) = 2/16$$

subject_race	subject_sex	subject_age	search_conducted
black	male	NA	0
other	male	22	0
hispanic	male	NA	0
black	male	35	0



# Missingness by a variable



## Definition

The SMR **for observation  $i$  restricted by variable  $j'$**  is the percentage of missing values for a traffic stop  $i$ , excluding variable  $j'$ .

Let  $j' \in \{1, 2, \dots, k\}$  be the column index for a variable of interest. The SMR for observation  $i$  restricted by  $j'$  is given by

$$\text{SMR}_{i,j'} = \frac{1}{k-1} \sum_{j \neq j'} \mathbb{1}_M(x_{ij}).$$

# Missingness by a variable



## Definition

The SMR **for observation  $i$  restricted by variable  $j'$**  is the percentage of missing values for a traffic stop  $i$ , excluding variable  $j'$ .

Let  $j' \in \{1, 2, \dots, k\}$  be the column index for a variable of interest. The SMR for observation  $i$  restricted by  $j'$  is given by

$$\text{SMR}_{i,j'} = \frac{1}{k-1} \sum_{j \neq j'} \mathbb{1}_M(x_{ij}).$$

## Example

$$\text{SMR}_{1,\text{race}} = 1/3$$

$$\text{SMR}_{2,\text{race}} = 0$$

subject_race	subject_sex	subject_age	search_conducted
black	male	NA	0
other	male	22	0
hispanic	male	NA	0
black	male	35	0

# SMR by race



Let's quantify how missingness varies by race.

Let  $r \in \{1, 2, \dots, k\}$  be the column index corresponding to the race variable in  $X$ .

Assume that  $(x_{i,r})$  has only three levels: "White", "Other", and NA (missing). We can partition the observations into index sets  $(W)$ ,  $(O)$ , and  $(NA) \subseteq \{1, 2, \dots, n\}$  such that

$$x_{i,r} = \begin{cases} \text{"White"} & i \in (W) \\ \text{"Other"} & i \in (O) \\ \text{NA} & i \in (NA). \end{cases}$$

# SMR by race



The White-SMR, Other-SMR, and NA-SMR are given by:

$$\text{SMR}_{(W)} = \frac{1}{|(W)|} \sum_{i \in (W)} \text{SMR}_{i,r}$$

$$\text{SMR}_{(O)} = \frac{1}{|(O)|} \sum_{i \in (O)} \text{SMR}_{i,r}$$

$$\text{SMR}_{(NA)} = \frac{1}{|(NA)|} \sum_{i \in (NA)} \text{SMR}_{i,r}.$$

## Example

$$\text{SMR}_{(B)} = 1/6$$

$$\text{SMR}_{(H)} = 1/3$$

$$\text{SMR}_{(O)} = 0$$

subject_race	subject_sex	subject_age	search_conducted
black	male	NA	0
other	male	22	0
hispanic	male	NA	0
black	male	35	0

# SMR by date and time



We apply a similar method to continuous variables.

We need to be careful with partitioning the indices – the partitions need enough observations for the average restricted SMRs to be meaningful.

# SMR by date and time



We apply a similar method to continuous variables.

We need to be careful with partitioning the indices – the partitions need enough observations for the average restricted SMRs to be meaningful.

- For date, we partition observations by the week.
- For time, we partition observations by the month and day/night.

# About the datasets



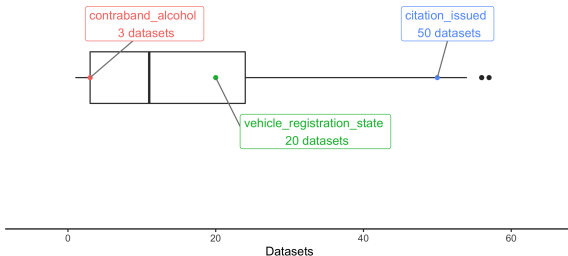
The Stanford Open Policing Project has 100 million traffic stop observations from 56 police departments (so, 56 datasets) from 1999 to 2020. Each datasets records a different set of variables.

# About the datasets



The Stanford Open Policing Project has 100 million traffic stop observations from 56 police departments (so, 56 datasets) from 1999 to 2020. Each dataset records a different set of variables.

The variables recorded during a traffic stop depends on the police department.





# Data pre-processing



We select  $k = 9$  variables about

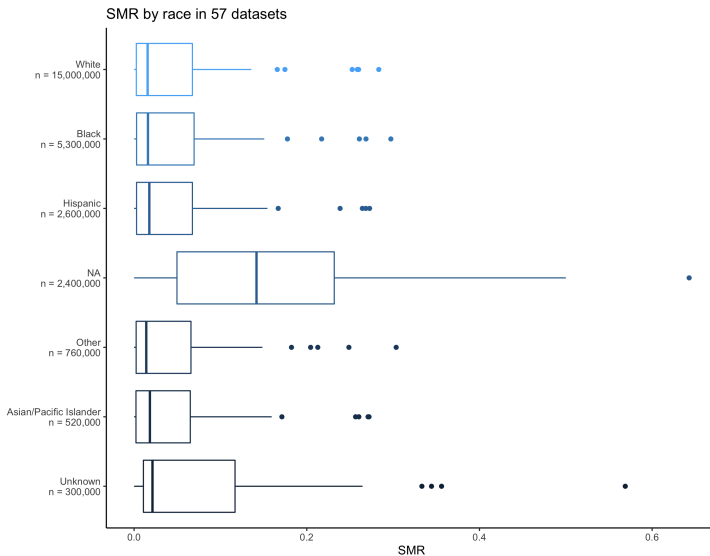
- driver demographic: race, sex, age;
- situational details: time, date, latitude, longitude; and
- outcomes: `search_conducted` and `arrest_made`.

And, for computational reasons, I take a 30% random sample of each dataset.

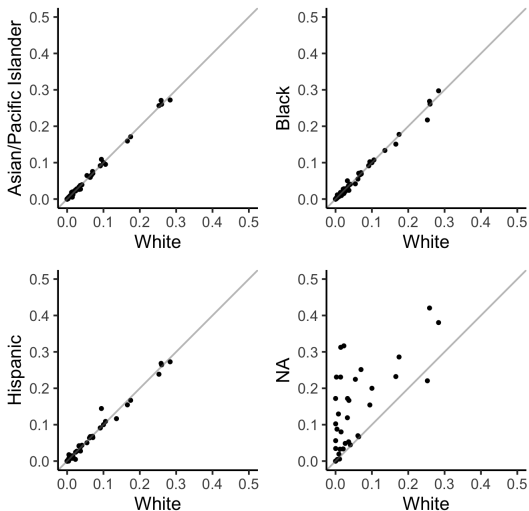
## Section 3

# Visualizations

# SMR by race



# SMR by race



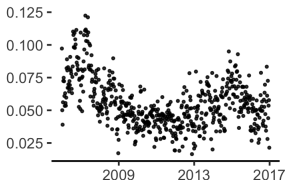
# SMR by date



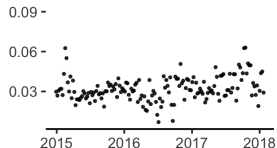
Iowa (statewide)  
n = 620,000



Wichita, KS  
n = 280,000



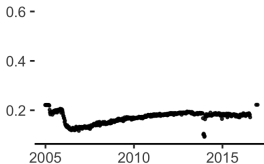
Louisville, KY  
n = 33,000



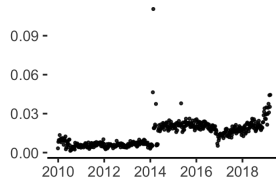
Saint Paul, MN  
n = 200,000



South Carolina (statewide)  
n = 2,700,000

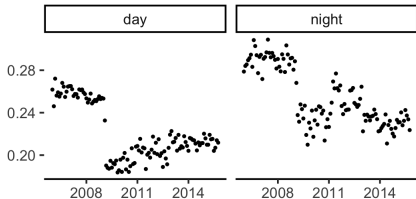


Nashville, TN  
n = 930,000

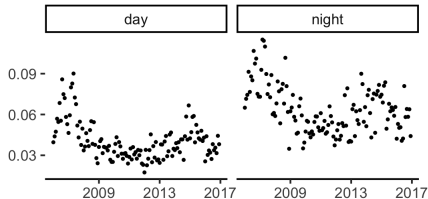


# SMR by time

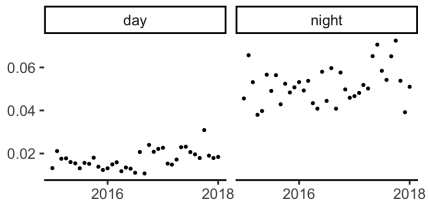
Iowa (statewide)  
n = 510,000



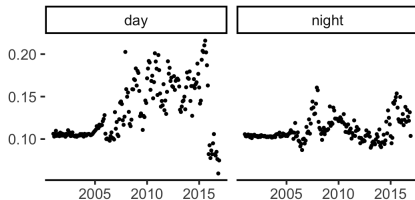
Wichita, KS  
n = 260,000



Louisville, KY  
n = 32,000



Saint Paul, MN  
n = 200,000



# Discussion



- Regarding race, there isn't much evidence that Asian/Pacific Islander, Black, and Hispanic traffic stops are recorded differently than White traffic stops.
- Traffic stops with missing race have higher rates of missing values in other variables, too.
- Traffic stops are recorded differently based on the date and time of the traffic stop. Exactly *how* the missingness is different depends on the dataset.
- Future research could investigate the missingness mechanism (missing at random and missing not at random).

Thank you!

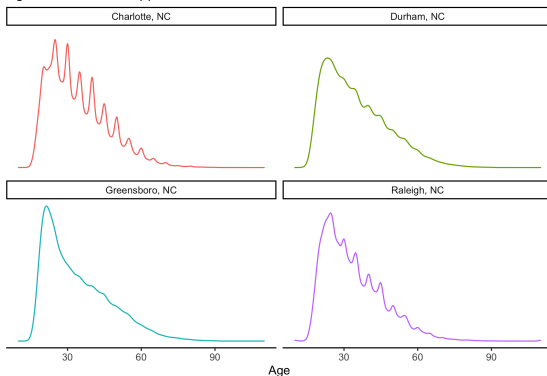


# Note: Data collection is a deeply human process.



In states like California and New York, officers use *perception* to record driver race, gender, and age.

Age densities for stopped drivers in North Carolina



Evidently, North Carolina uses perception for age, too.