Definitions
oooooo

About the data
ooo

Visualizations
oooo

# Implications of Missing Traffic Stop Data

Amber Lee advised by Jo Hardin

Pomona College

November 5, 2021

Definitions
000000

About the data
000

Visualizations
0000

Pomona
College

Table of Contents

**1** Definitions

**2** About the data

**3** Visualizations

Definitions
oooooo

About the data
ooo

Visualizations
oooo

## Motivation

Pomona
College

Public concern over racial profiling has led to the collection of traffic stop
data. As traffic stop data has been made available to the public,
researchers have tried to quantitatively look for evidence of discrimination.

## Motivation

Public concern over racial profiling has led to the collection of traffic stop data. As traffic stop data has been made available to the public, researchers have tried to quantitatively look for evidence of discrimination.

- A generic model may look like search_conducted $\sim$ race + gender + age + day/night.

- 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.

Definitions
000000

About the data
000

Visualizations
0000

## Motivation

Pomona
College

Public concern over racial profiling has led to the collection of traffic stop data. As traffic stop data has been made available to the public, researchers have tried to quantitatively look for evidence of discrimination.

- A generic model may look like search_conducted $\sim$ race + gender + age + day/night.

- 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.

$\Rightarrow$ **What are the trends of missing values in traffic stop data?**

## Row-wise missingness

Pomona
College

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

- Each observation $i = 1, 2, \ldots, n$ represents a single traffic stop.

- We denote a **missingness indicator function** as

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

Definitions
●○○○○○

About the data
○○○

Visualizations
○○○○

# Row-wise missingness

Pomona
College

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

- Each observation $i = 1, 2, \ldots, n$ 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 $\mathrm{SMR}$) is the percentage of missing values for a single traffic stop $i$.

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

**Definitions**
○●○○○○

About the data
○○○

Visualizations
○○○○

## Dataset missingness

Pomona
College

### 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

$$\text{SMR}(X) = 7/12$$

**Definitions**
○○●○○○

About the data
○○○

Visualizations
○○○○

## Missingness by a variable

Pomona College

One limitation of $\mathrm{SMR}_i$ and $\mathrm{SMR}(X)$ is that we can't see how missingness varies *by* a variable, like time of day or the `race` of the driver.

### Definition

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

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

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

Example

**Definitions**
○○○●○○

About the data
○○○

Visualizations
○○○○

## SMR by `race`

Pomona College

Let's quantify how missingness varies by `race`.

Let $r \in \{1, 2, \ldots, k\}$ 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 $(NA) \subseteq \{1, 2, \ldots, n\}$ such that

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

Definitions
○○○○●○

About the data
○○○

Visualizations
○○○○

SMR by `race`

Pomona
College

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

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

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

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

Example

Definitions
○○○○○○●

About the data
○○○

Visualizations
○○○○

## 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 $\mathrm{SMR}$s to be meaningful.

**Definitions**
○○○○○●

About the data
○○○

Visualizations
○○○○

# 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.

Definitions
oooooo

About the data
●oo

Visualizations
oooo

## About the datasets

The Stanford Open Policing Project

- amasses 100 million number of traffic stop observations

- from 21 state patrol agencies and 35 municipal police departments

- from 1999 to 2020.

Definitions
000000

About the data
○●○

Visualizations
0000

# Data pre-processing

Pomona
College

We consider 32 total datasets and $k = 9$ variables.

- driver demographic: race, sex, age;

- situational details: time, date, latitude, longitude; and

- outcomes: search_conducted and arrest_made.

Definitions
○○○○○○
About the data
○○●
Visualizations
○○○○

## Note: data collection is a deeply human process. 🏛 Pomona College

In states like California and New York, officers use *perception*(!!!) to gauge driver race, gender, and age. Evidently, North Carolina definitely uses perception for age, too.

Definitions
oooooo

About the data
ooo

Visualizations
●ooo

# Missingness by race

yeehaw

Definitions
000000

About the data
000

Visualizations
0●00

## Post-stop outcomes by race

Pomona
College

yeehaw

Definitions
oooooo

About the data
ooo

Visualizations
oo●o

Missingness by `date`

graphic

Definitions
oooooo

About the data
ooo

Visualizations
ooo●

Missingness by `time`

Pomona
College

graphic