

HW: Datatypes and Wrangling

Hector Corrada Bravo

2018-02-04

Data types

1) Provide a URL to the dataset.

I downloaded my dataset from http://www.hcbravo.org/IntroDataSci/misc/BPD_Arrests.csv

2) Explain why you chose this dataset.

I am interested in studying how rates of arrests in different parts of Baltimore are related to demographic statistics.

3) What are the entities in this dataset? How many are there?

Entities are specific arrests. There are 104528.

4) How many attributes are there in this dataset?

There are 15 attributes.

5) What is the datatype of each attribute (categorical -ordered or unordered-, numeric -discrete or continuous-, datetime, geolocation, other)? Write a short sentence stating how you determined the type of each attribute. Do this for at least 5 attributes, if your dataset contains more than 10 attributes, choose 10 of them to describe.

Num	Name	Type	Description
1	arrest	categorical	Identifier of each arrest, takes values from finite set
2	age	numeric continuous	Ages are numeric values measured in time units
3	race	categorical unordered	Can take value from finite set of possible races
4	sex	categorical unordered	Can take value from finite set of possible sexes
5	arrestDate	datetime	Specifies date of arrest
6	arrestTime	datetime	Specifies time of arrest
7	arrestLocation	text - address	Street address of arrest
8	incidentOffense	categorical unordered	Can take value from finite set of possible offenses
9	incidentLocation	text - address	Street address if incident

Num	Name	Type	Description
10	charge	categorical unordered	Can take value from finite set of possible charges

6) Write R code that loads the dataset using function `read_csv`. Were you able to load the data successfully? If no, why not?

```
library(tidyverse)

url <- "http://www.hcbravo.org/IntroDataSci/misc/BPD_Arrests.csv"
arrest_tab <- read_csv(url)
arrest_tab %>% slice(1:10)

## # A tibble: 10 x 15
##   arrest   age sex  race arrestDate arrestTime arrestLocation
##   <int> <int> <chr> <chr> <chr>         <time>      <chr>
## 1 11126858   23 B    M    01/01/2011 00'00"      <NA>
## 2 11127013   37 B    M    01/01/2011 01'00"    2000 Wilkens Ave
## 3 11126887   46 B    M    01/01/2011 01'00"    2800 Mayfield Ave
## 4 11126873   50 B    M    01/01/2011 04'00"    2100 Ashburton St
## 5 11126968   33 B    M    01/01/2011 05'00"    4000 Wilsby Ave
## 6 11127041   41 B    M    01/01/2011 05'00"    2900 Spellman Rd
## 7 11126932   29 B    M    01/01/2011 05'00"    800 N Monroe St
## 8 11126940   20 W    M    01/01/2011 05'00"    5200 Moravia Rd
## 9 11127051   24 B    M    01/01/2011 07'00"    2400 Gainsborough Ct
## 10 11127018   53 B    M    01/01/2011 15'00"    3300 Woodland Ave
## # ... with 8 more variables: incidentOffense <chr>, incidentLocation
## #   <chr>, charge <chr>, chargeDescription <chr>, district <chr>, post
## #   <int>, neighborhood <chr>, `Location 1` <chr>
```

Wrangling

1) My pipeline computes average arrest age (ignoring ages ≤ 0), for each district and writes them in increasing order

```
mean_ages <- arrest_tab %>%
  filter(age > 0) %>%
  select(district, age) %>%
  group_by(district) %>%
  summarize(mean_age=mean(age)) %>%
  arrange(mean_age)
mean_ages
```

```
## # A tibble: 10 x 2
##   district    mean_age
##   <chr>         <dbl>
## 1 NORTHEASTERN    30.4
## 2 SOUTHERN        32.3
## 3 SOUTHWESTERN    32.5
## 4 SOUTHEASTERN    32.5
## 5 CENTRAL         33.1
## 6 NORTHERN        33.1
```

```
## 7 <NA>          33.4
## 8 EASTERN       34.1
## 9 WESTERN       34.4
## 10 NORTHWESTERN 34.6
```

Plotting

1) This barplot shows the average arrest age per district (ignoring ages ≤ 0)

```
mean_ages %>%
  ggplot(aes(x=district, y=mean_age)) +
  geom_bar(stat="identity") +
  coord_flip()
```

