Stanford Open Policing Project dataset

ANALYZING POLICE ACTIVITY WITH PANDAS

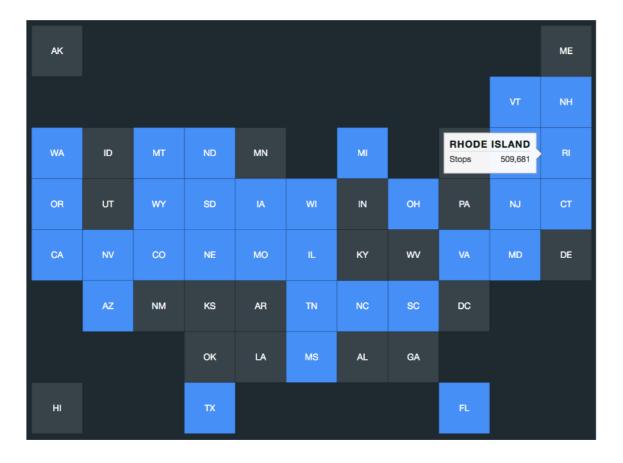


Kevin MarkhamFounder, Data School



Introduction to the dataset

Traffic stops by police officers



Download data for any state:

https://openpolicing.stanford.edu/

Preparing the data

- Examine the data
- Clean the data

```
import pandas as pd
ri = pd.read_csv('police.csv')
ri.head(3)
```

```
stop_date stop_time county_name driver_gender driver_race
 state
        2005-01-04
                        12:55
                                       NaN
                                                                White
                                                        Μ
0
        2005-01-23
                        23:15
                                       NaN
                                                                White
        2005-02-17
                        04:15
                                       NaN
                                                        Μ
                                                                White
```

- Each row represents one traffic stop
- NaN indicates a missing value

Locating missing values (1)

```
ri.isnull()
```

```
state stop_date stop_time county_name driver_gender
False
          False
                    False
                                  True
                                               False
False
          False
                    False
                                  True
                                               False
False
          False
                                               False
                     False
                                  True
```



Locating missing values (2)

```
ri.isnull().sum()
```

```
state0stop_date0stop_time0county_name91741driver_gender5205...
```

- .sum() calculates the sum of each column
- True = 1 , False = 0

Dropping a column

```
ri.isnull().sum()
```

```
state0stop_date0stop_time0county_name91741driver_gender5205driver_race5202...
```

```
    county_name column only contains missing values
```

Drop county_name using the .drop() method

```
ri.drop('county_name',
   axis='columns', inplace=True)
```

```
ri.shape
```

```
(91741, 15)
```

Dropping rows

dropna(): Drop rows based on the presence of missing values

```
ri.head()
```

```
stop_date stop_time driver_gender driver_race
  state
0
     RΙ
        2005-01-04
                                         Μ
                                                  White
                       12:55
        2005-01-23
                       23:15
                                                  White
     RΙ
                                         Μ
     RI
        2005-02-17
                        04:15
                                         M
                                                  White
                                         Μ
        2005-02-20
                                                  White
     RΙ
                       17:15
        2005-02-24
4
     RΙ
                        01:20
                                                  White
```

```
ri.dropna(subset=['stop_date', 'stop_time'], inplace=True)
```

Let's practice!

ANALYZING POLICE ACTIVITY WITH PANDAS



Using proper data types

ANALYZING POLICE ACTIVITY WITH PANDAS



Kevin MarkhamFounder, Data School



Examining the data types

```
stop_date object
stop_time object
driver_gender object
... ...
stop_duration object
drugs_related_stop bool
district object
```

- object: Python strings (or other Python objects)
- bool: True and False values
- Other types: int , float , datetime , category

Why do data types matter?

- Affects which operations you can perform
- Avoid storing data as strings (when possible)
 - int , float : enables mathematical operations
 - datetime: enables date-based attributes and methods
 - category: uses less memory and runs faster
 - o bool: enables logical and mathematical operations



Fixing a data type

```
apple
```

```
date time price
0 2/13/18 16:00 164.34
1 2/14/18 16:00 167.37
2 2/15/18 16:00 172.99
```

```
apple.price.dtype
```

```
dtype('0')
```

```
apple['price'] =
  apple.price.astype('float')
apple.price.dtype
```

```
dtype('float64')
```

- Dot notation: apple.price
- Bracket notation:

```
apple['price']
```

Let's practice!

ANALYZING POLICE ACTIVITY WITH PANDAS



Creating a DatetimeIndex

ANALYZING POLICE ACTIVITY WITH PANDAS



Kevin MarkhamFounder, Data School



Using datetime format

```
ri.head(3)
```

```
      stop_date
      stop_time
      driver_gender
      driver_race

      0
      2005-01-04
      12:55
      M
      White

      1
      2005-01-23
      23:15
      M
      White

      2
      2005-02-17
      04:15
      M
      White
```

```
ri.dtypes
```

```
stop_date object
stop_time object
driver_gender object
driver_race object
...
```

- Combine stop_date and stop_time into one column
- 2. Convert it to datetime format

Combining object columns

```
date time price
0 2/13/18 16:00 164.34
1 2/14/18 16:00 167.37
2 2/15/18 16:00 172.99
```

```
0 2-13-18
1 2-14-18
2 2-15-18
Name: date, dtype: object
```

apple.date.str.replace('/', '-')

```
combined =
  apple.date.str.cat(apple.time, sep=' ')
```

```
combined
```

```
0 2/13/18 16:00
1 2/14/18 16:00
2 2/15/18 16:00
Name: date, dtype: object
```

apple

Converting to datetime format

```
apple['date_and_time'] = pd.to_datetime(combined)
apple
```

```
date time price date_and_time
0 2/13/18 16:00 164.34 2018-02-13 16:00:00
1 2/14/18 16:00 167.37 2018-02-14 16:00:00
2 2/15/18 16:00 172.99 2018-02-15 16:00:00
```

```
apple.dtypes
```

```
date object
time object
price float64
date_and_time datetime64[ns]
```



Setting the index

```
apple.set_index('date_and_time', inplace=True)
apple
```

```
date time price

date_and_time

2018-02-13 16:00:00 2/13/18 16:00 164.34

2018-02-14 16:00:00 2/14/18 16:00 167.37

2018-02-15 16:00:00 2/15/18 16:00 172.99
```

```
apple.index
```



Let's practice!

ANALYZING POLICE ACTIVITY WITH PANDAS

