

~~Jupyter~~

HackerRank

Jupyter comes with abilities to convert to other formats

`jupyter nbconvert --to FORMAT mynotebook.ipynb`

Slideshow

`jupyter nbconvert notebook.ipynb --to slides`

Numpy → arrays / work with numerical values

Pandas → manipulate / analyze data

Matplotlib → visualize data

→ `/c/users/abahn/AppData/Roaming/SPB_Data`

Process

1) Question 2) Wrangle make sure that data is in good quality

A) Gather B) Assess C) Clean

3) Explore patterns & relationships exploratory data analysis

4) Conclusions Descriptive Statistics 5) Communications

Packages

Questions

CSV Files

`df = pd.read_csv('csv file', separator)`

if another separator other than comma is used

→ header pick a line to be header

→ index\_col specify 1st column

replace column names → `header = 0, names = labels`

`df.info()` summary of dataframe

• `nunique()` number of unique values in each column

• `describe()` descriptive analysis for each column

### Select multiple columns

```
df = df.iloc[:, np.r_[ : , single line ]]
```

### View index number & label for each column

```
for i, v in enumerate(df.columns):
```

```
    print(i, v)
```

if diagnosis refers to an element as object

```
type(df[' '][0]) for further investigation
```

### Clean Data

1) missing data

2) Duplicates

3) incorrect datatypes

• `fillna(mean, inplace=True)`

`pd.to_datetime()`

`df.duplicated()`

`df.drop_duplicates(inplace=True)`

for large datasets

`sum(df.duplicated())`

### Plotting % matplotlib inline

quick histogram view

• `hist(figsize=( , ))`

• `plot(kind = ' ')` bar, scatter, pie, box

`pd.plotting.scatter_matrix( , figsize=(15, 15))`

show all possible relations

`pd.plot(x= , y= )`



## Conclusions descriptive statistics & visualizations

`df_M = df[df['diagnosis'] == 'M']` boolean indexing

## Communication

to sort graphs according to index

`ind = df_a['education'].value_counts().index`

`df_a['education'].value_counts()[ind].plot(kind = "bar")`

## → Gathering

- **TSV** tab separated files
- **CSV** comma separated values

## File Types

**Flat Files** → Simple, human readable, software readable

data in plain text format with one data record per line

**Web Scraping** get readable html using code

`span` → groups text together

to know which encoding

to use

`<meta charset = " " >`

`from bs4 import BeautifulSoup`

`with open('rt.html', 'r') as file:`

`soup = BeautifulSoup(file, 'lxml')`

`soup.find('tag').contents[0]`

returns a list of tags' children

• **Download Files** using codes

`import requests` creates request

`url = ' ' ,`

`response = requests.get(url)`

creates folder

`import os`

`folder_name = ' ' ,`

`if not os.path.exists(' '):`

`os.makedirs(folder_name)`

access content & write to a file

`with open(os.path.join(folder_name, url.split('/')[-1]), mode = 'wb') as file:`

`file.write(response.content)`

glob

```
for ebert_review in glob.glob('Data/*_*.txt'):  
    with open(ebert_review, encoding = ' ') as file:
```

file.read

file.readline

after read the "file"

put the variables in a dictionary

put dictionary in a list

add list to dataframe

```
with open (.....) as f:
```

```
    for line in f:
```

```
        # Do something with 'line'
```

## API (application programming interface)

→ choose API over scrapping

we use wikipedia API

MediaWiki

access library ease work for programmers

MediaWiki

uptools

- page ("last bit of wikipedia code")
- get() gets all data

JSON

great for representing complicated data hierarchies



```
df.query(' _ > _ ')
```

greater than

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

more  
visualization  
controls

```
plt.bar(x-coord, bar heights)
```

→ specify names for x-coordinates

```
plt.xticks([1, 2, 3], [' ', ' ', ' '])
```

- title(' \_ ')

- xlabel(' \_ ')

- ylabel(' \_ ');

- tick\_label

normalize  
produces  
proportions

```
pd.cut(df[' _ '], , labels = )
```

column  
name

values

strings that represents