# Data Science Foundation
## Lesson #5 - Exploratory Data Analysis

Ivanovitch Silva
September, 2017

# Agenda

- Case study: unemployment rate
- Tabular vs Visual representation
- Matplotlib
- Line plots
- Case study: movie ratings
- Bar & Scatter plots

# Update the repository

git clone https://github.com/ivanovitchm/EEC2006.git

Or ….

git pull

# Motivation

# Case study: unemployment rate (US)

# Investigating the dataset

| DATE Year-Month-Day | VALUE |
|---|---|
| 1948–01–01 | 3.4 |
| 1948–02–01 | 3.8 |
| 1948–03–01 | 4.0 |
| 1948–04–01 | 3.9 |
| 1948–05–01 | 3.5 |

Conversion of types (Object to Datetime)

```python
import pandas as pd
df['col'] = pd.to_datetime(df['col'])
```

| DATE | VALUE |
|------|-------|
| 1948-01-01 | 3.4 |
| 1948-02-01 | 3.8 |
| 1948-03-01 | 4.0 |
| 1948-04-01 | 3.9 |
| 1948-05-01 | 3.5 |
| 1948-06-01 | 3.6 |
| 1948-07-01 | 3.6 |
| 1948-08-01 | 3.9 |
| 1948-09-01 | 3.8 |
| 1948-10-01 | 3.7 |
| 1948-11-01 | 3.8 |
| 1948-12-01 | 4.0 |

# Observation from the table representation

- What is the minimum value?
- What is the maximum value?
- Is there seasonality?
- What are the trend up periods?
- What are the trend down periods?
- Is the table representation really useful?

# Visual representation

```
import matplotlib.pyplot as plt
plt.plot()
plt.show()
```
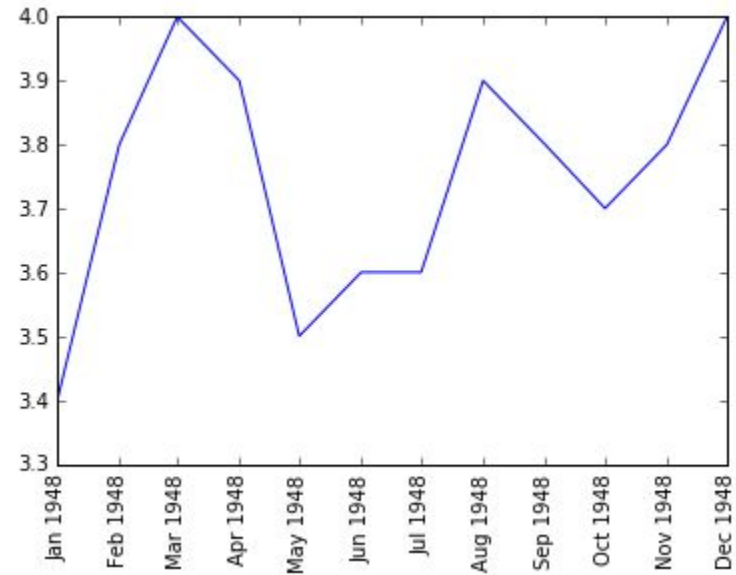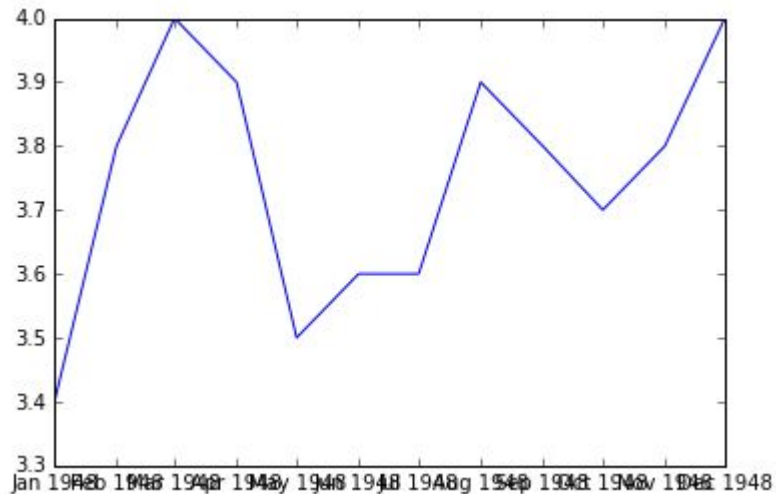
# Adding data

```
plt.plot(x_values, y_values)
```

# Fixing axis ticks
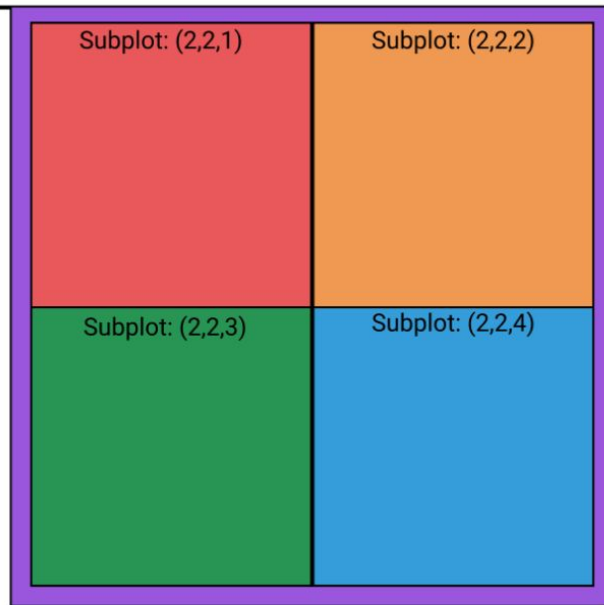


plt.xticks(rotation=90)

# Additional information



Monthly Unemployment Trends, 1948

plt.xlabel("Month")
plt.ylabel("Unemployment Rate")
plt.title("Monthly Unemployment Trends, 1948")

# Grid positioning

```python
import matplotlib.pyplot as plt
fig = plt.figure()
ax1 = fig.add_subplot(2,2,1)
ax2 = fig.add_subplot(2,2,2)
ax3 = fig.add_subplot(2,2,3)
ax4 = fig.add_subplot(2,2,4)
```
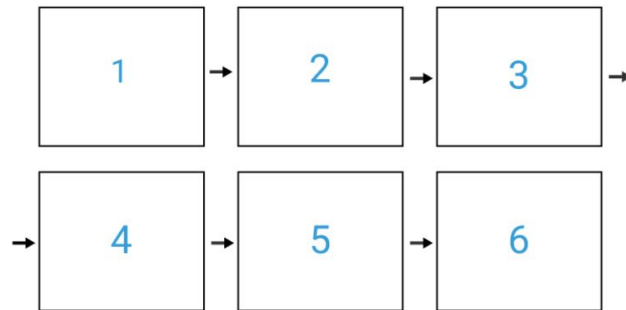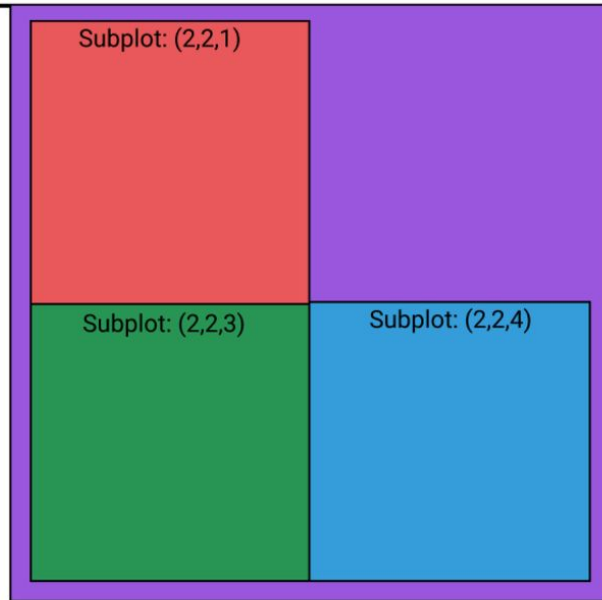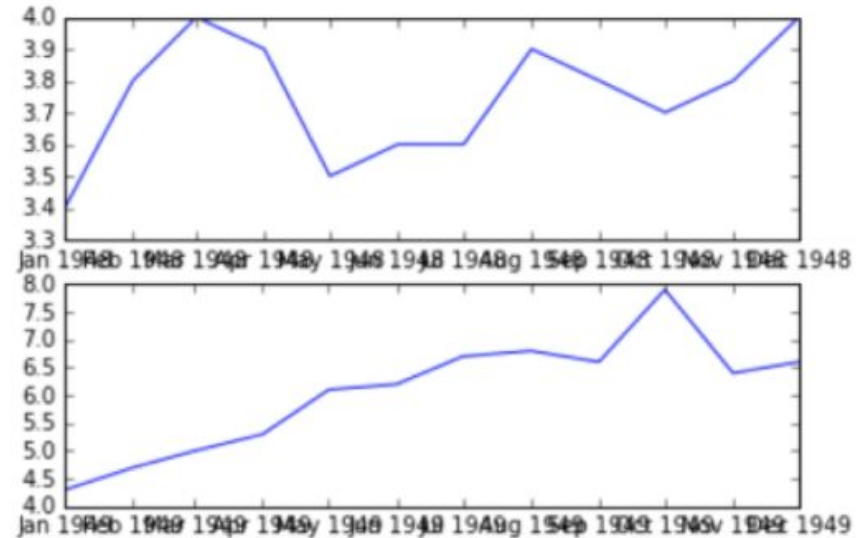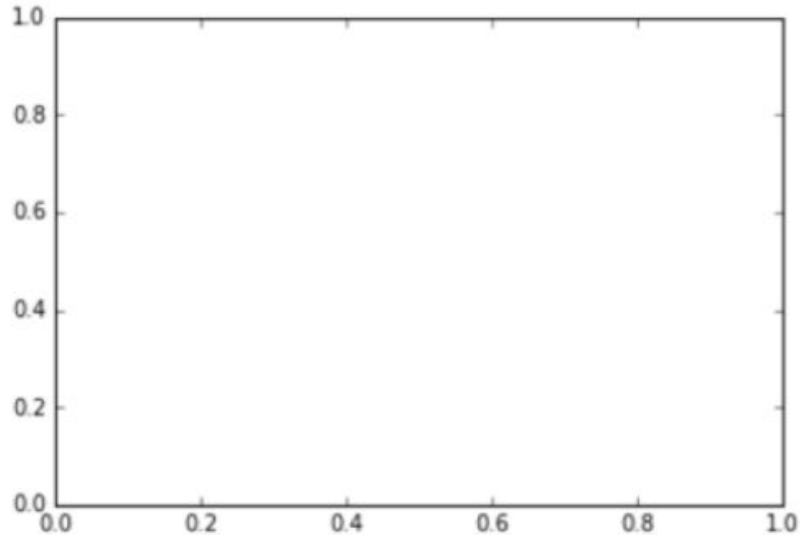
Figure

| | |
|---|---|
| Subplot: (2,2,1) | Subplot: (2,2,2) |
| Subplot: (2,2,3) | Subplot: (2,2,4) |

# Grid positioning

fig.add_subplot(4, 1, x)

```
1 →
→ 2 →
→ 3 →
→ 4 →
```

fig.add_subplot(2, 2, x)

```
1 → 2 →
→ 3 → 4
```

fig.add_subplot(2, 3, x)

```
1 → 2 → 3 →
→ 4 → 5 → 6
```

# Grid positioning

```python
import matplotlib.pyplot as plt
fig = plt.figure()
ax1 = fig.add_subplot(2,2,1)
ax3 = fig.add_subplot(2,2,3)
ax4 = fig.add_subplot(2,2,4)
```
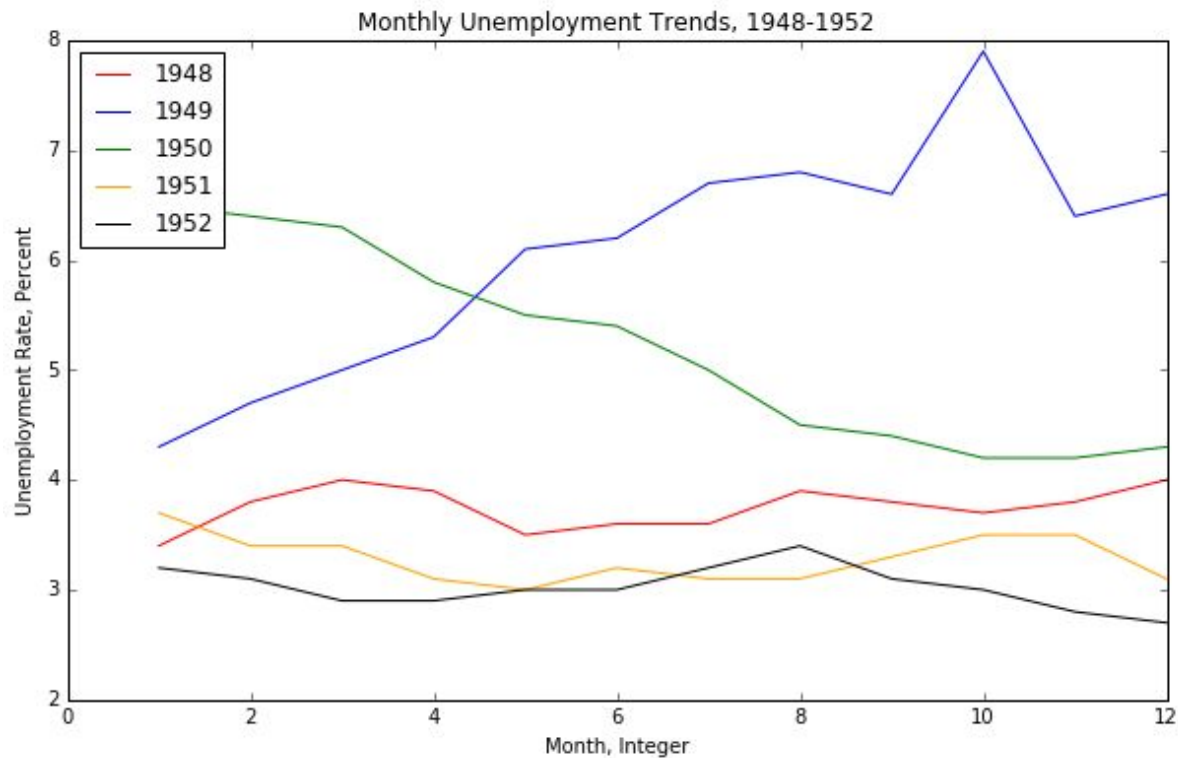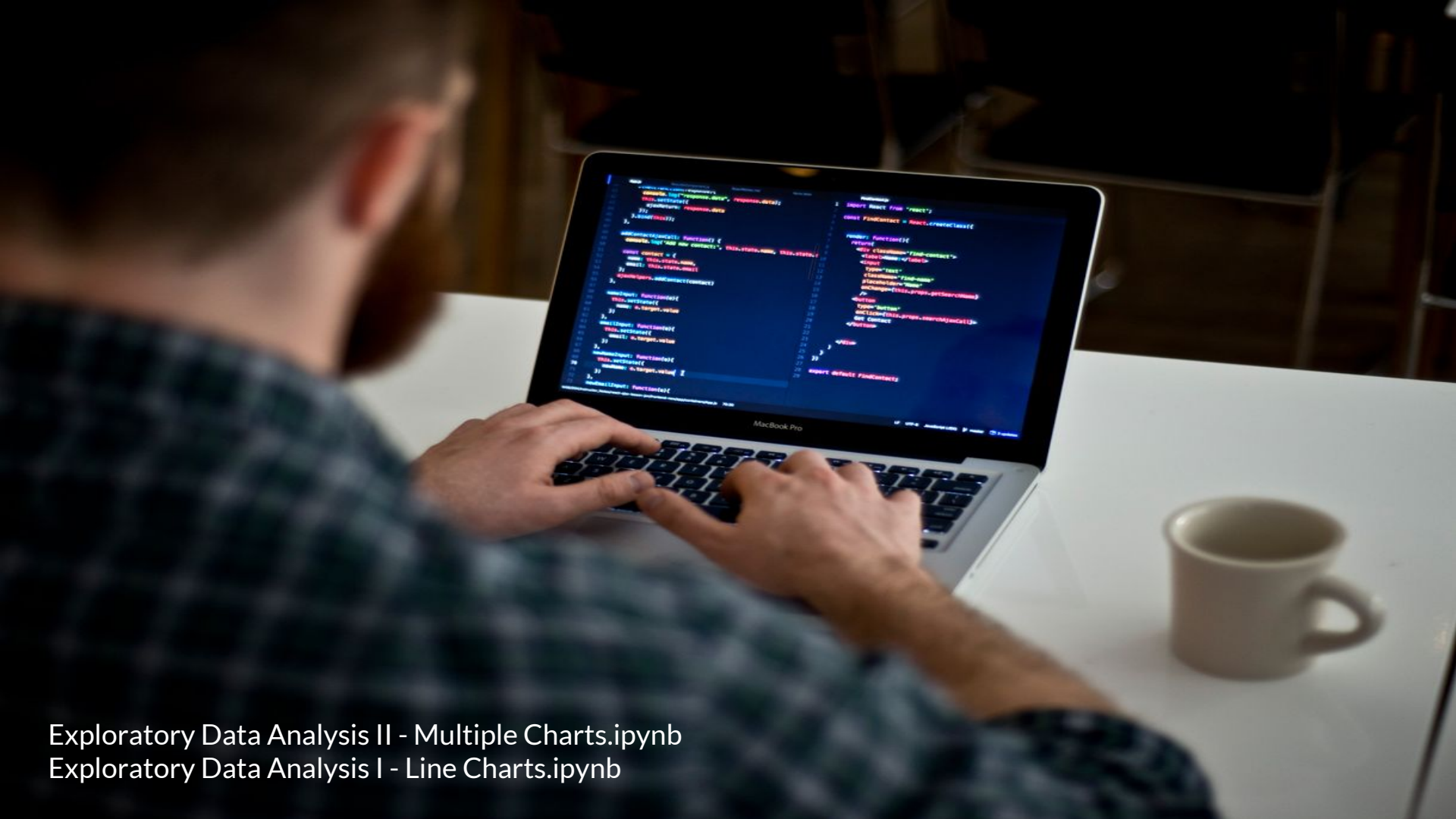
Figure

Subplot: (2,2,1)

Subplot: (2,2,3)

Subplot: (2,2,4)

# Formatting and spacing



```
fig = plt.figure(figsize=(width, height))
```

# Overlaying line charts



Monthly Unemployment Trends, 1948-1952

Exploratory Data Analysis II - Multiple Charts.ipynb
Exploratory Data Analysis I - Line Charts.ipynb

# Case study: movie ratings



Metacritic

Rotten Tomatoes

IMDB

Fandango

# Bias in movie ratings

# Introduction to the data

| | FILM | RT_user_norm | Metacritic_user_nom | IMDB_norm | Fandango_Ratingvalue | Fandango_Stars |
|---|---|---|---|---|---|---|
| **0** | Avengers: Age of Ultron (2015) | 4.3 | 3.55 | 3.90 | 4.5 | 5.0 |
| **1** | Cinderella (2015) | 4.0 | 3.75 | 3.55 | 4.5 | 5.0 |
| **2** | Ant-Man (2015) | 4.5 | 4.05 | 3.90 | 4.5 | 5.0 |
| **3** | Do You Believe? (2015) | 4.2 | 2.35 | 2.70 | 4.5 | 5.0 |
| **4** | Hot Tub Time Machine 2 (2015) | 1.4 | 1.70 | 2.55 | 3.0 | 3.5 |

https://github.com/fivethirtyeight/data/tree/master/fandango

# Bar plot

# Creating bars



Positions of the bars: 0.75, 1.75, 2.75, 3.75, 4.75

Positions of the axis labels: 1.0, 2.0, 3.0, 4.0, 5.0

Width of the bars: 0.5, 0.5, 0.5, 0.5, 0.5

RT_user_norm, Metacritic_user_nom, IMDB_norm, Fandango_Ratingvalue, Fandango_Stars

# Creating bars

```python
fig, ax = plt.subplots()
```

```python
# Positions of the left sides of the 5 bars. [0.75, 1.75, 2.75, 3.75, 4.75]
from numpy import arange
bar_positions = arange(5) + 0.75
# Heights of the bars.  In our case, the average rating for the first movie
 in the dataset.
num_cols = ['RT_user_norm', 'Metacritic_user_nom', 'IMDB_norm', 'Fandango_Ra
tingvalue', 'Fandango_Stars']
bar_heights = norm_reviews[num_cols].iloc[0].values
ax.bar(bar_positions, bar_heights)
```
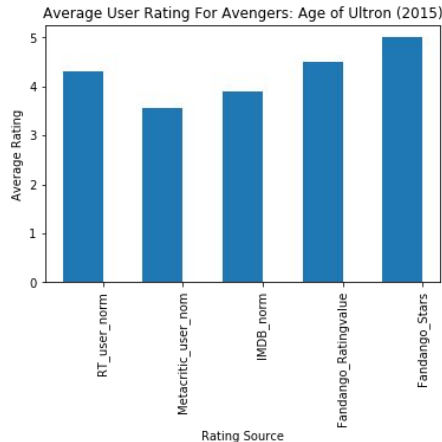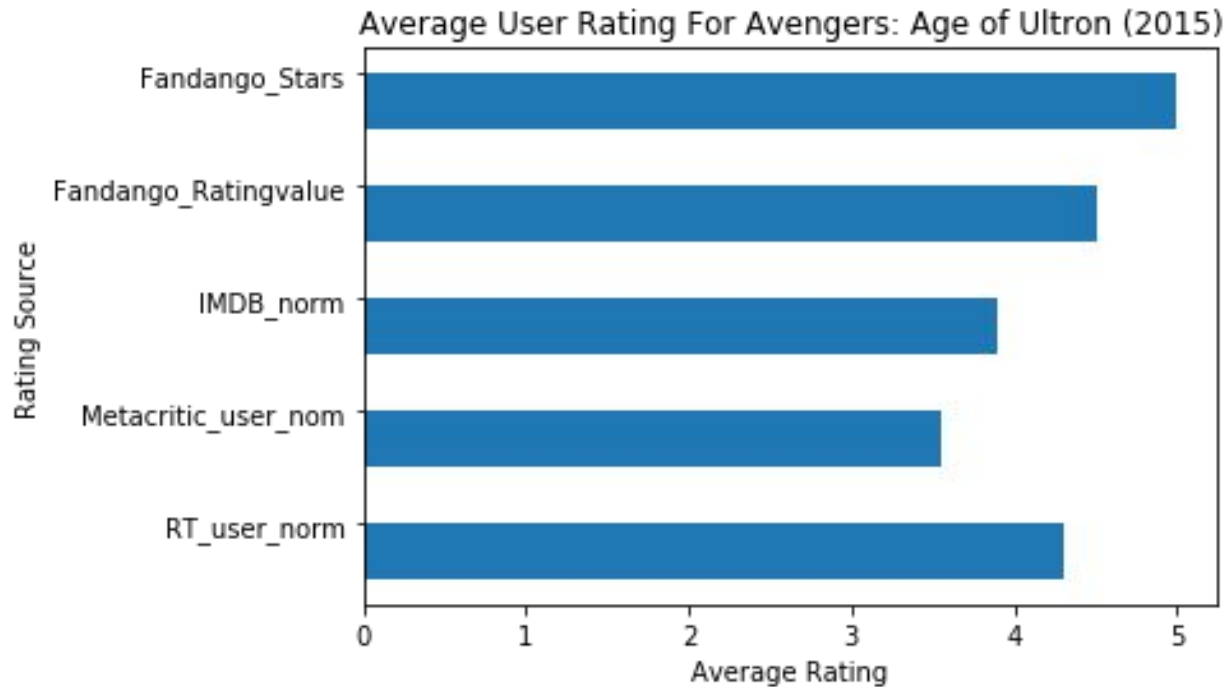
# Aligning axis ticks and labels

```
tick_positions = range(1,6)
ax.set_xticks(tick_positions)

num_cols = ['RT_user_norm', 'Metacritic_user_nom', 'IMDB_norm', 'Fandango_Ra
tingvalue', 'Fandango_Stars']
ax.set_xticklabels(num_cols)

ax.set_xticklabels(num_cols, rotation=90)
```



Average User Rating For Avengers: Age of Ultron (2015)

# Horizontal bar plots



Average User Rating For Avengers: Age of Ultron (2015)
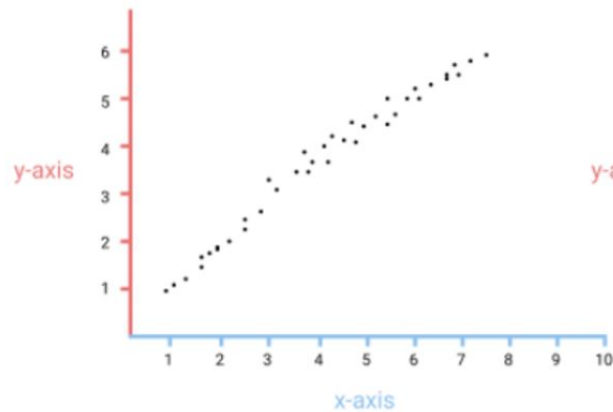
# Horizontal bar plots

```python
bar_widths = norm_reviews[num_cols].iloc[0].values
bar_positions = arange(5) + 0.75
ax.barh(bar_positions, bar_widths, 0.5)
```

```python
tick_positions = range(5) + 1
num_cols = ['RT_user_norm', 'Metacritic_user_nom', 'IMDB_norm', 'Fandango_Ra
tingvalue', 'Fandango_Stars']
ax.set_yticks(tick_positions)
ax.set_yticklabels(num_cols)
```
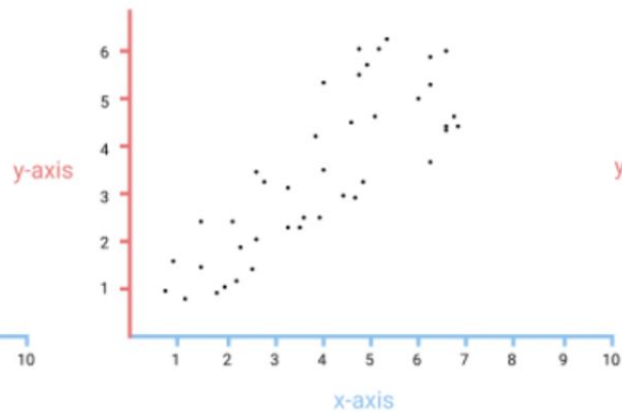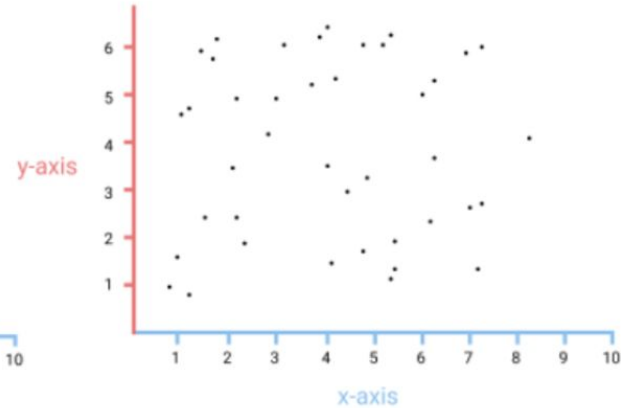
# Scatter plot

# Switching axes

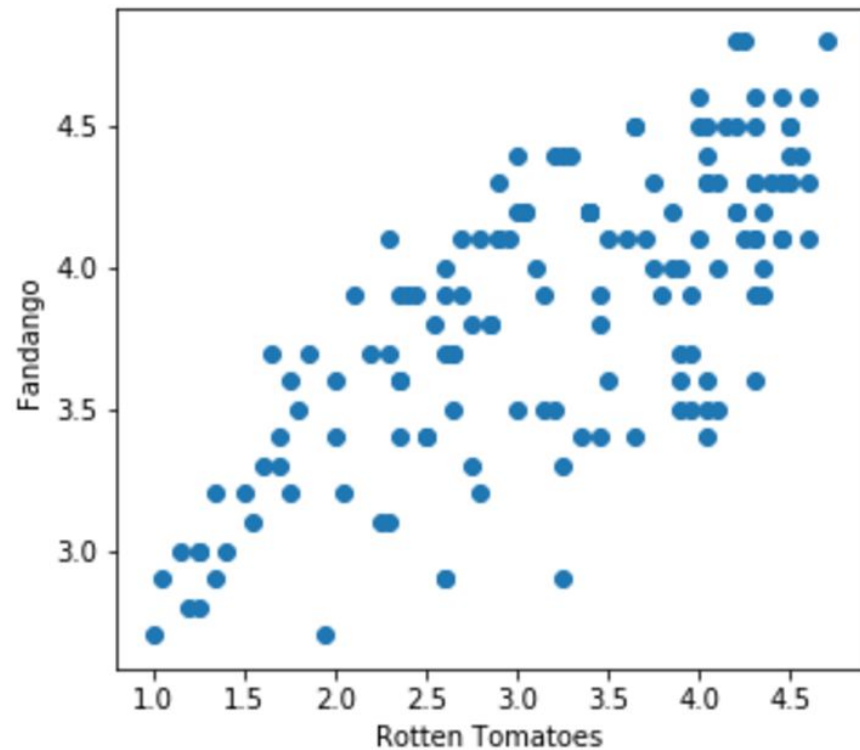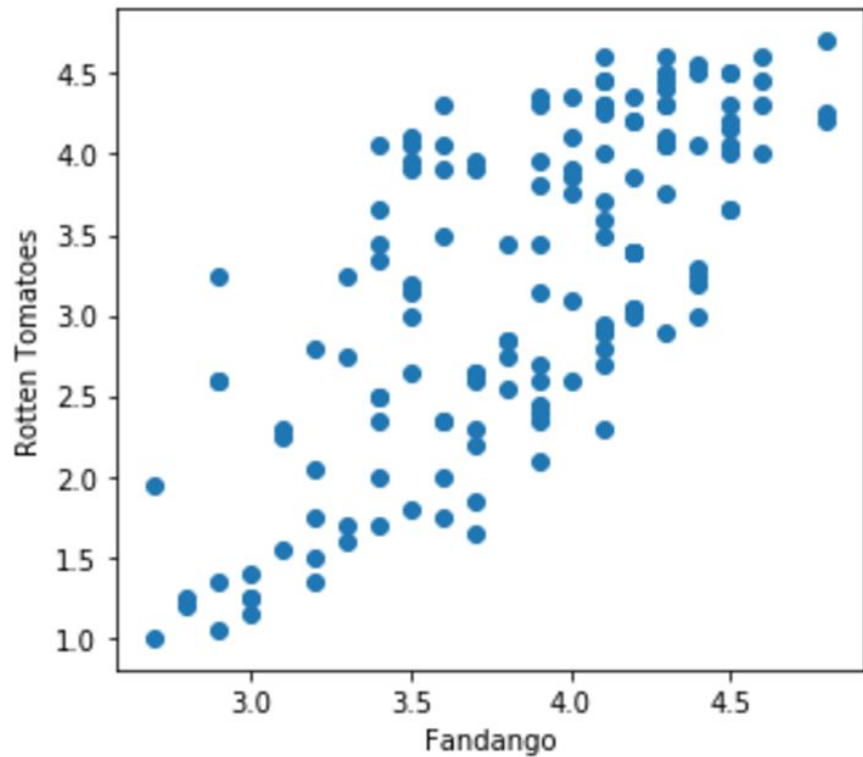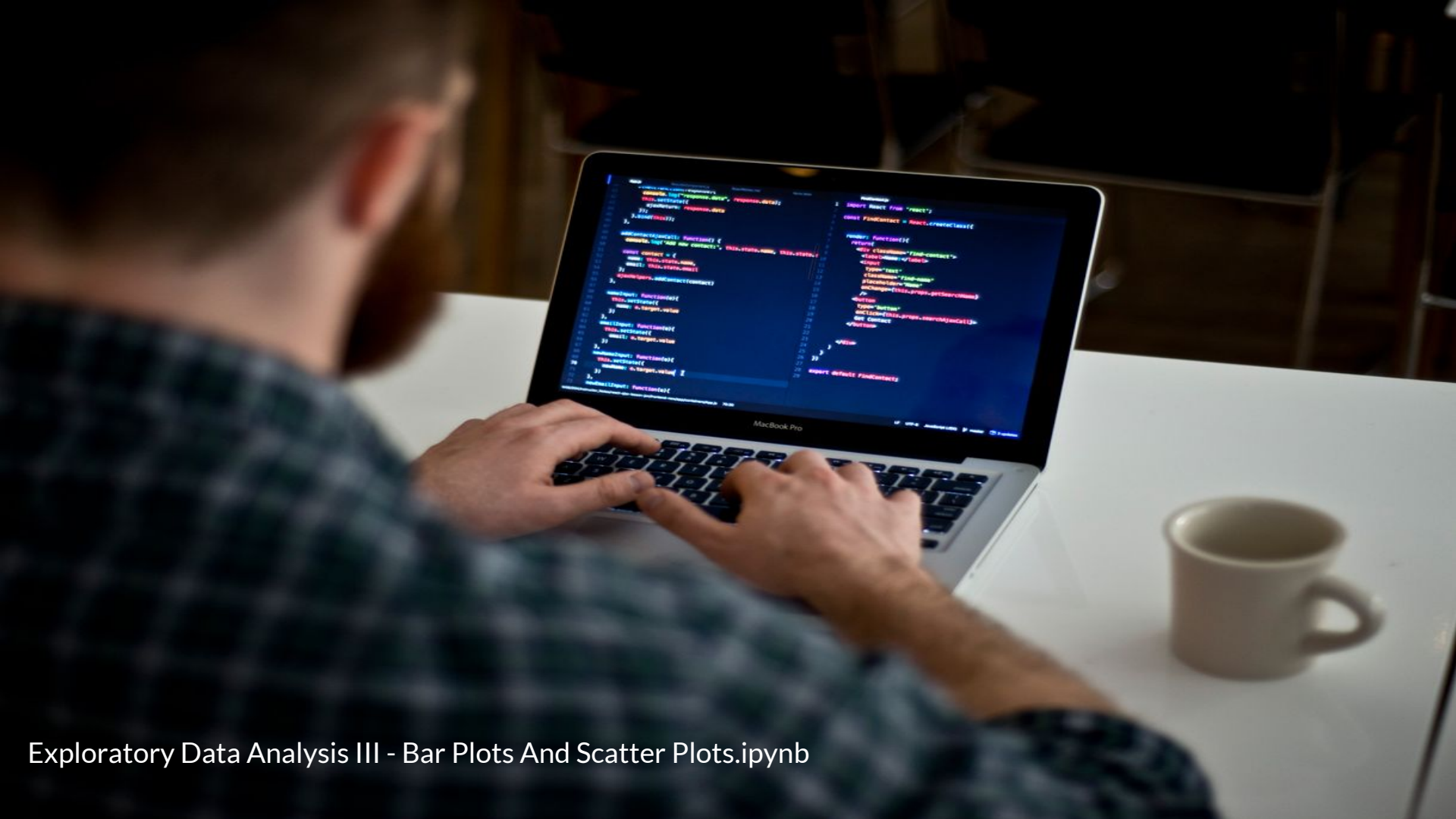Exploratory Data Analysis III - Bar Plots And Scatter Plots.ipynb