

# team13\_capstone\_project

June 1, 2020

## 1 Team13: Capstone project of Python Bootcamp

This is [the Capstone project for Team 13 of the Python Data Analysis Bootcamp](#). We are trying, more or less, to follow the structure of [jupytertemplate](#).

### 1.1 Purpose

State the purpose of the notebook.

### 1.2 Methology

Quickly describe assumptions and processing steps.

### 1.3 TODO / Improvements

- ☒ Find a dataset that has at least 2 CSV files
- ☐ Come up with 5 questions that you want to answer while exploring the dataset
- ☐ Perform EDA (Exploratory Data Analysis) on your dataset with basic visualisations

### 1.4 Results

### 1.5 Setup

```
[16]: # install system dependencies
import sys
import os

!conda install -c conda-forge --yes --prefix {sys.prefix} pandas jupyterthemes
↪ seaborn jupyter_contrib_nbextensions pandoc
```

```
Collecting package metadata (current_repodata.json): done
Solving environment: done
```

```
# All requested packages already installed.
```

### 1.5.1 Library Import

```
[96]: # load libraries and setup environment
# mandatory
import pandas as pd

%matplotlib inline
import matplotlib.pyplot as plt

# optional
import numpy as np
import seaborn as sns
from jupyterthemes import jtplot
from IPython.core.display import HTML
jtplot.style(theme='monokai', context='notebook', ticks=True, grid=False)
```

## 1.6 Parameter definition

We set all relevant parameters for our notebook. By convention, parameters are uppercase, while all the other variables follow Python's guidelines.

## 1.7 Data import

We retrieve all the required data for the analysis.

```
[83]: cost_of_living = pd.read_csv('../data/andytran1996_cost-of-living/
↳ datasets_73059_162758_cost-of-living-2018.csv')

# we are droppping the Rank column because it's entirely empty
cost_of_living = cost_of_living.drop(columns = 'Rank')
```

## 1.8 Data processing

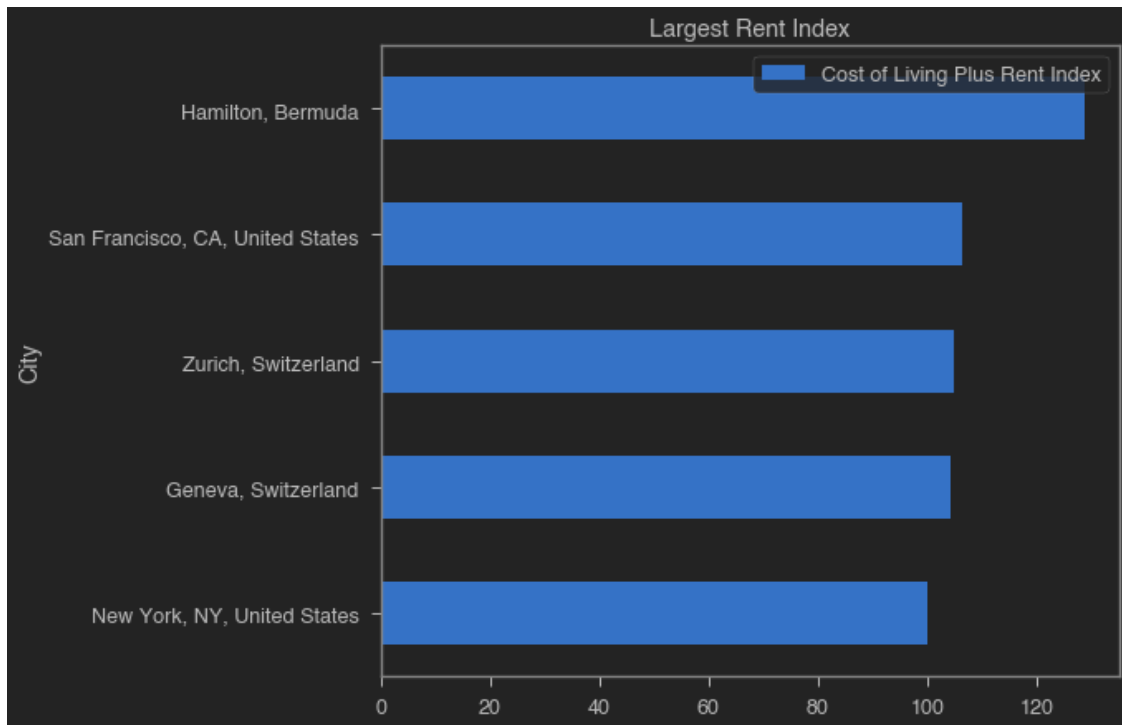
### 1.8.1 1. What are the five cities with the highest/lowest cost of living (incl. rent)?

```
[153]: caption_column = 'City'
index_column = 'Cost of Living Plus Rent Index'

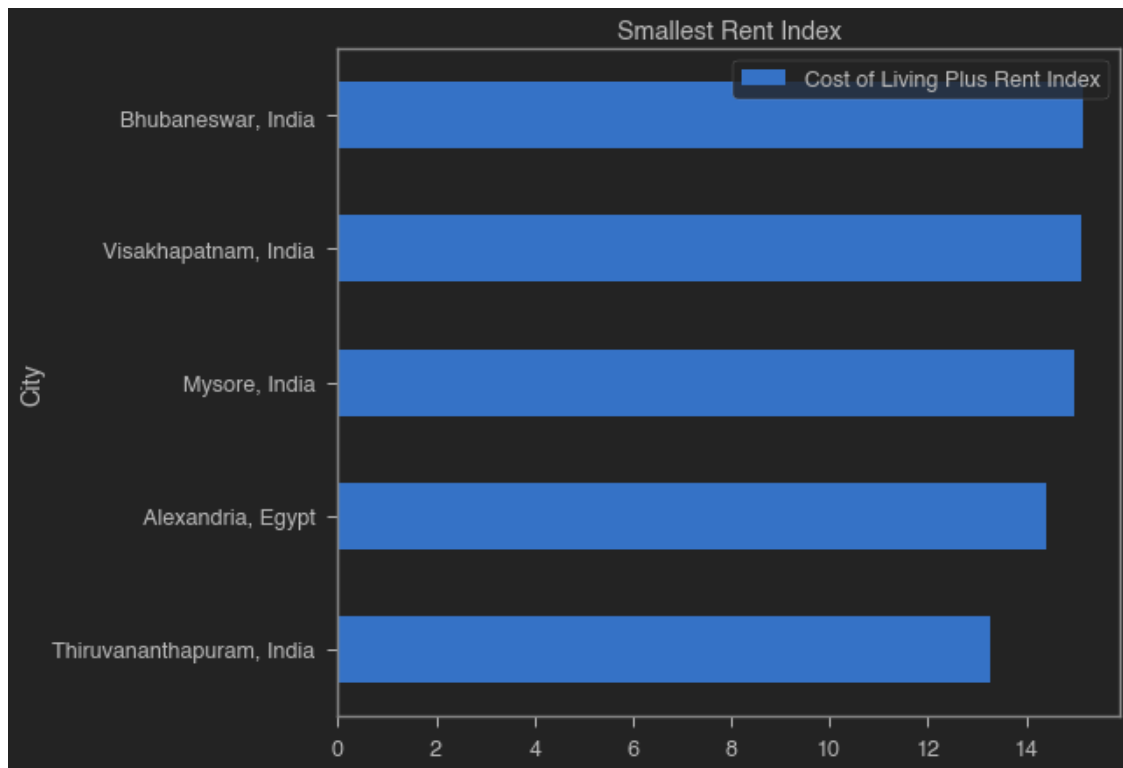
def display_cost_of_living(costs, title):
    filtered_costs = costs[[caption_column, index_column]].
    ↳sort_values(index_column)
    filtered_costs.plot.barh(title = title, x = caption_column, y =
    ↳index_column)
    plt.show()
    display(filtered_costs.style.hide_index())

# print the ten most expensive cities in the database in 2018
```

```
display_cost_of_living(cost_of_living.nlargest(5, index_column), 'Largest Rent_  
↪Index')  
display_cost_of_living(cost_of_living.nsmallest(5, index_column), 'Smallest_  
↪Rent Index')
```



<pandas.io.formats.style.Styler at 0x7f3fd402a210>



<pandas.io.formats.style.Styler at 0x7f3fd40acad0>

## 1.9 References

- [data for the cost of living](#)
- [base data for countries of the world](#)
- [data for life expectancy from the WHO](#)
- [roshansharma\\_europe-datasets](#)