# Python

## Web Scrapping

# Outline

- 3 Popular Tools and Libraries used for Web Scraping in Python
- Basic Data Scrapping
- Web Scraping
    - Crawl
    - Parse and Transform
    - Store
- Scraping Images

# Three populars library:

- **BeautifulSoup**
  - *BeautifulSoup* is an amazing parsing library in Python that enables us to extract data from HTML and XML documents.
- **Scrapy**
  - Scrapy is a [Python](#) framework for large scale web scraping.
- **Selenium**
  - Selenium is another popular tool for automating browsers. It's primarily used for testing in the industry but is also very handy for web scraping.

# Basic Data Scrapping

# Data Scrapping

- Create code like this.

```
dokumen = '''
<html>
<head>
    <title>Tutorial BeautifulSoup</title>
</head>

<body>
    <p class="judul">Judul Dokumen</p>

    <p class="paragraf">Ini adalah contoh paragraf</p>

    <a href="https://www.ifa.com" class="url">ITERA</a>
</body>

</html>
'''
```

# Print Data using HTML Parser

- To use HTML parser, we need BeautifulSoup library. By default, BeautifulSoup has been availabled in Anaconda (Jupyter Notebook)

```python
from bs4 import BeautifulSoup
html_soup = BeautifulSoup(dokumen, 'html.parser')

print(html_soup)
```

```html
<html>
<head>
<title>Tutorial BeautifulSoup</title>
</head>
<body>
<p class="judul">Judul Dokumen</p>
<p class="paragraf">Ini adalah contoh paragraf</p>
<a class="url" href="https://www.ifa.com">ITERA</a>
</body>
</html>
```

# Get Data From Web Document

- Get HTML data From Element )

```python
judul = html_soup.find('p')
print(judul)
```

```
<p class="judul">Judul Dokumen</p>
```

- Get HTML data from element using **class**

```python
judul = html_soup.find('p', class_='judul')
paragraf = html_soup.find('p', class_='paragraf')
print(judul)
print(paragraf)
```

```
<p class="judul">Judul Dokumen</p>
<p class="paragraf">Ini adalah contoh paragraf</p>
```

- Get text in class 'judul' with element 'p'

```python
judul_saja = html_soup.find('p', class_='judul').text
print(judul_saja)
```

```
Judul Dokumen
```

# Find Data or element or class in web site

- The find () function can only extract one output while usually many of the same HTML tags that all of them want to retrieve.

- To retrieve HTML content with the same tag use the find_all () function

```python
all_paragraf = html_soup.find_all('p')
print(all_paragraf)
```

[<p class="judul">Judul Dokumen</p>, <p class="paragraf">Ini adalah contoh paragraf</p>]
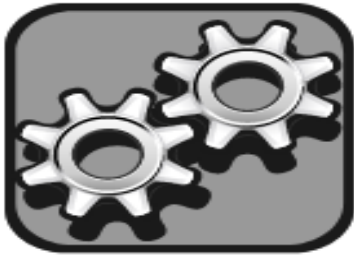
# Web Scrapping

# Components of Web Scraping

## 1. Crawl

The first step is always navigate to the target website by making an HTTP Request and download the response you get.

## 2. Parse and Transform

Once you have received the response, Now its time to parse this downloaded data into a HTML Parser like Beautiful Soup and Extract the Required Data.
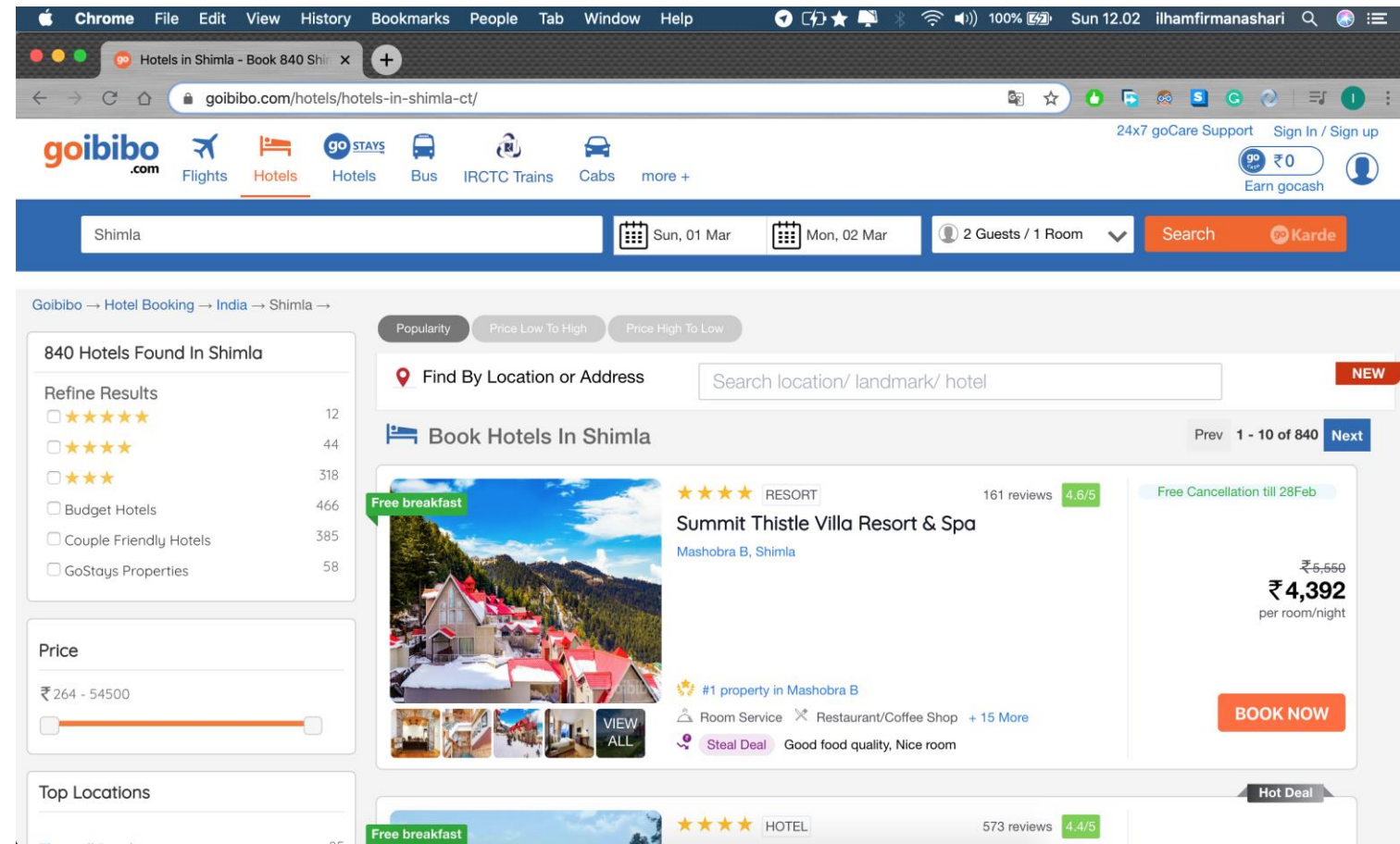
## 3. Store

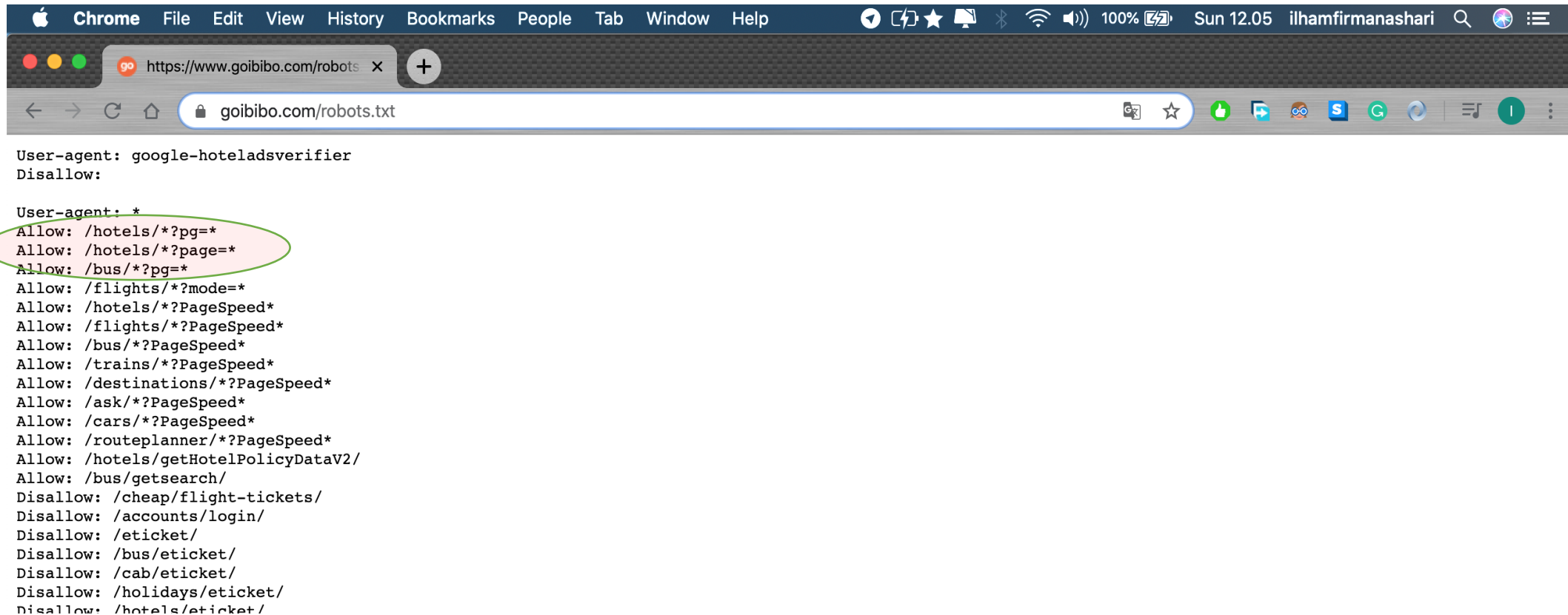Now that you have extracted the required data. You can simply store this as JSON or CSV file or directly into the the DataBase like MongoDB

# Get-Started

- Let's understand these components in detail. We'll do this by scraping hotel details like the name of the hotel and price per room from the **goibibo** website: *goibibo.com/hotels/hotels-in-shimla-ct/*

# Analysis before scrapping using (**robots.txt**)

- So, looks like we are allowed to scrape the data from our targeted URL. We are good to go and write the script of our web robot. Let's begin!

# Step 1 : Start to Scrapping
# Define libraries

- To scrapping the data from website, we need 3 libraries

```python
"""
Web Scraping using Beautiful Soup
"""

import requests
from bs4 import BeautifulSoup
import pandas as pd
```

# Step 1 : Start to Scrapping
## URL Web that will be scrapped

```python
url = "https://www.goibibo.com/hotels/hotels-in-shimla-ct/"

headers = {
    'User-Agent': "Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHT
    }
|
response = requests.request("GET", url, headers=headers)
```
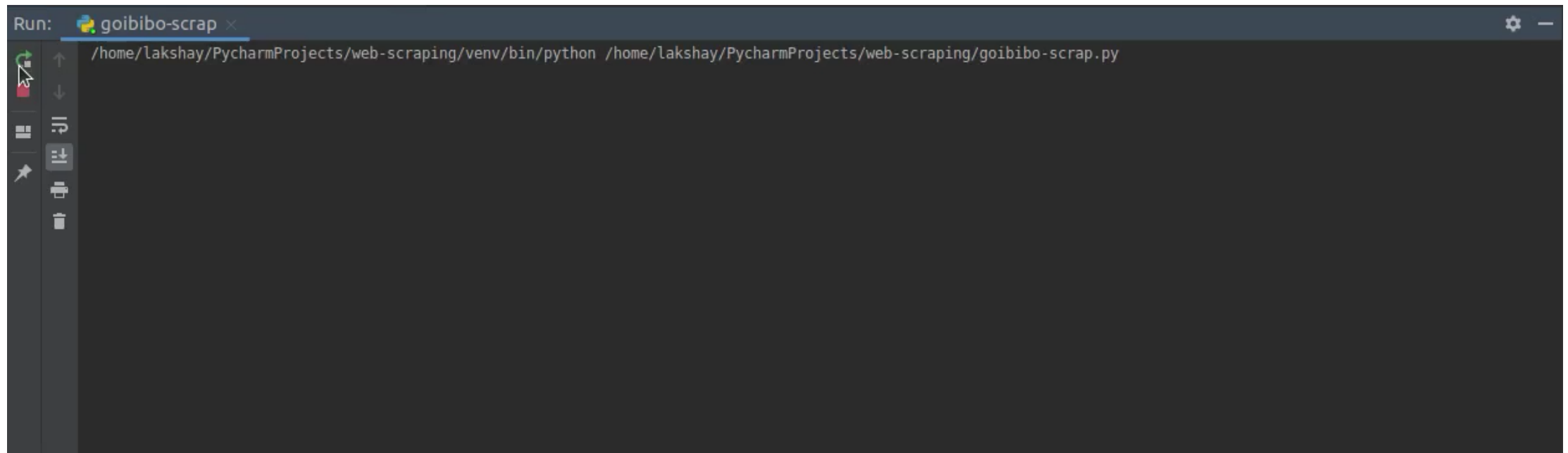
Note: to know the headers of the website, inspect element.
*optional*

# Step 2 : Parsing and Transform
# Parsing html (web-document)

```python
data = BeautifulSoup(response.text, 'html.parser')
print(data)
```

Run:  goibibo-scrap ✕

/home/lakshay/PycharmProjects/web-scraping/venv/bin/python /home/lakshay/PycharmProjects/web-scraping/goibibo-scrap.py

# Step 2: Parsing and Transform Inspect element

- The next step is to parse this data into an HTML Parser and for that, we will use the *BeautifulSoup* library.

- Next, we will select the card and click on the 'Inspect Element' option to get the source code of that particular card. You will get something like this:



```
▼<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-
container-template new-htl-design-tile-main-block" id="htl-
8487828901284192936" onclick="window.open('https://www.goibibo.com/hotels/
saffronstays-thanedhar-estate-shimla-hotel-in-shimla-8487828901284192936/');
event.stopPropagation();"> == $0
```

# Step 2: Parsing and Transform
## Get root element of the data

- The class name of all the cards would be the same and we can get a list of those cards by just passing the tag name and attributes like the *<class>* tag with its name like I've shown below:

```python
cards_data = data.find_all('div', attrs={'class', 'width100 fl htlListSeo
print('Total Number of Cards Found : ', len(cards_data))

for card in cards_data:
    print(card)
```

- **Note:** for more clearly, do inspect element in your web browser to get full class name

# Step 2: Parsing and Transform
Result



```
/home/lakshay/PycharmProjects/web-scraping/venv/bin/python /home/lakshay/PycharmProjects/web-scraping/goibibo-scrap.py
Total Number of Cards Found :  10
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-7071577211857429914"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-7035113582935215619"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-7543907359980632827"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-5109554460899404317"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-2485666429068449291"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-1686119541396289166"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-2866704395127249073"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-2668930289000053794"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-4132340955783784970"
<div class="width100 fl htlListSeo hotel-tile-srp-container hotel-tile-srp-container-template new-htl-design-tile-main-block" id="htl-7393074725561129529"
```

# Step 2: Parsing and Transform
## Get Hotel name and Room Price

- Select only the Hotel Name, perform the Inspect Element step, and do the same with the Room Price:
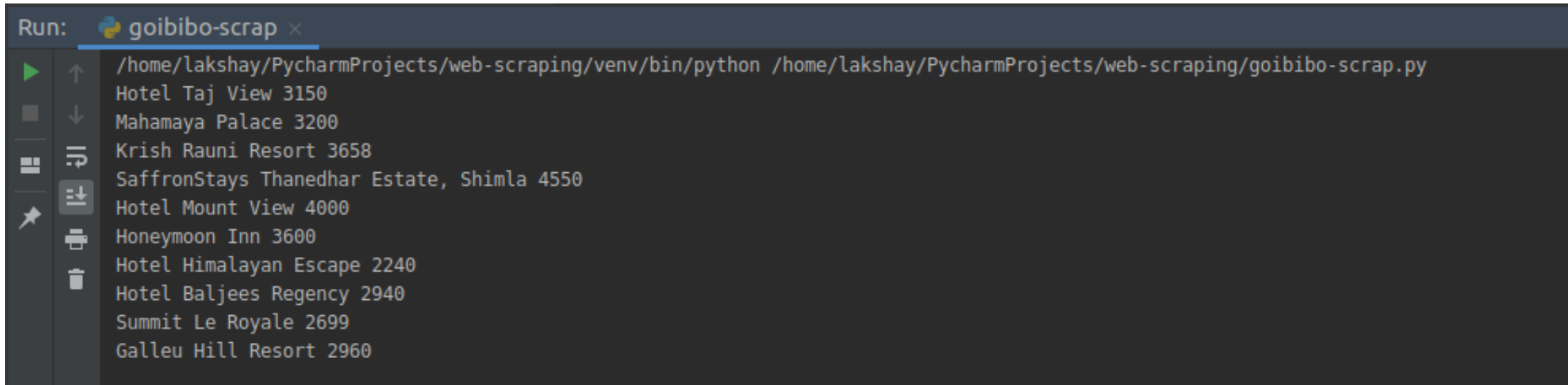
# Step 2: Parsing and Transform
## Get element in web for the data

- Now, for each card, we have to find the Hotel Name which can be extracted from the *<p>* tag .

- **This is because there is only one *<p>* tag for each card and Room Price by <li> tag along with the <class> tag and class name:**

```python
for card in cards_data:

    hotel_name = card.find('p')
    room_price = card.find('li', attrs={'class': 'htl-tile-discount-prc'})
    print(hotel_name.text, room_price.text)
```

# Step 2: Parsing and Transform
result



Run: 🐍 goibibo-scrap ✕

```
/home/lakshay/PycharmProjects/web-scraping/venv/bin/python /home/lakshay/PycharmProjects/web-scraping/goibibo-scrap.py
Hotel Taj View 3150
Mahamaya Palace 3200
Krish Rauni Resort 3658
SaffronStays Thanedhar Estate, Shimla 4550
Hotel Mount View 4000
Honeymoon Inn 3600
Hotel Himalayan Escape 2240
Hotel Baljees Regency 2940
Summit Le Royale 2699
Galleu Hill Resort 2960
```
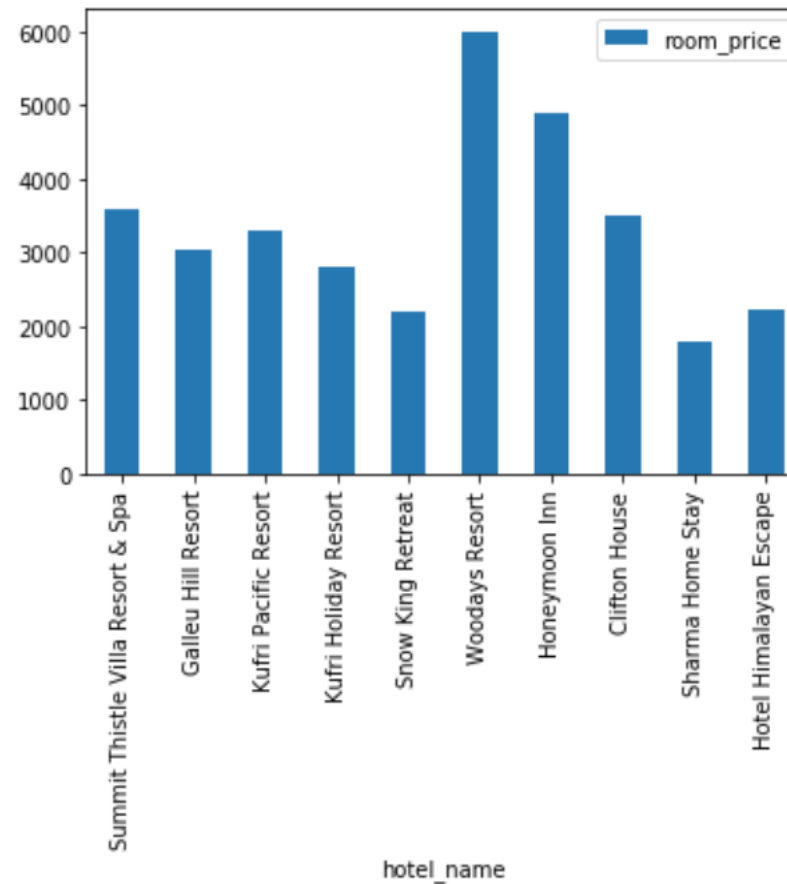
# Step 3 : Store the Data

- Next, let's go ahead and transform this list to a Pandas dataframe as it allows us to convert the dataframe into CSV or JSON files:

```python
scraped_data = []

for card in cards_data:

    card_details = {}

    hotel_name = card.find('p')
    room_price = card.find('li', attrs={'class': 'htl-tile-discount-prc'})

    card_details['hotel_name'] = hotel_name.text
    card_details['room_price'] = room_price.text

    scraped_data.append(card_details)

dataFrame = pd.DataFrame.from_dict(scraped_data)
dataFrame.to_csv('hotels_data.csv', index=False)
```

# Result

You can use this file to data analysis(create images plot, etc)

| A | B |
|---|---|
| hotel_name | room_price |
| Hotel Taj View | 3150 |
| Mahamaya Palace | 3200 |
| Krish Rauni Resort | 3658 |
| SaffronStays Thanedhar Estate, Shimla | 4550 |
| Hotel Mount View | 4000 |
| Honeymoon Inn | 3600 |
| Hotel Himalayan Escape | 2240 |
| Hotel Baljees Regency | 2940 |
| Summit Le Royale | 2699 |
| Galleu Hill Resort | 2960 |

# Data Visualization

- Show the room price on each hotel

# Data Visualization

```python
# harga dari range 2000-3000 ruppe
fd = df.loc[(df.room_price >= 2000) & (df.room_price <= 3000)]
print(fd)
```

```
                               hotel_name  room_price room_rating
0   OYO 1706 Hotel The Alpine Heritage Residency        2718         79%
1                      OYO 14170 Kufri Star Inn        2209         91%
2                       OYO 10692 Hotel Shubham        2987         84%
6                 OYO 37261 Shirdi Sai Residency        2379         88%
```
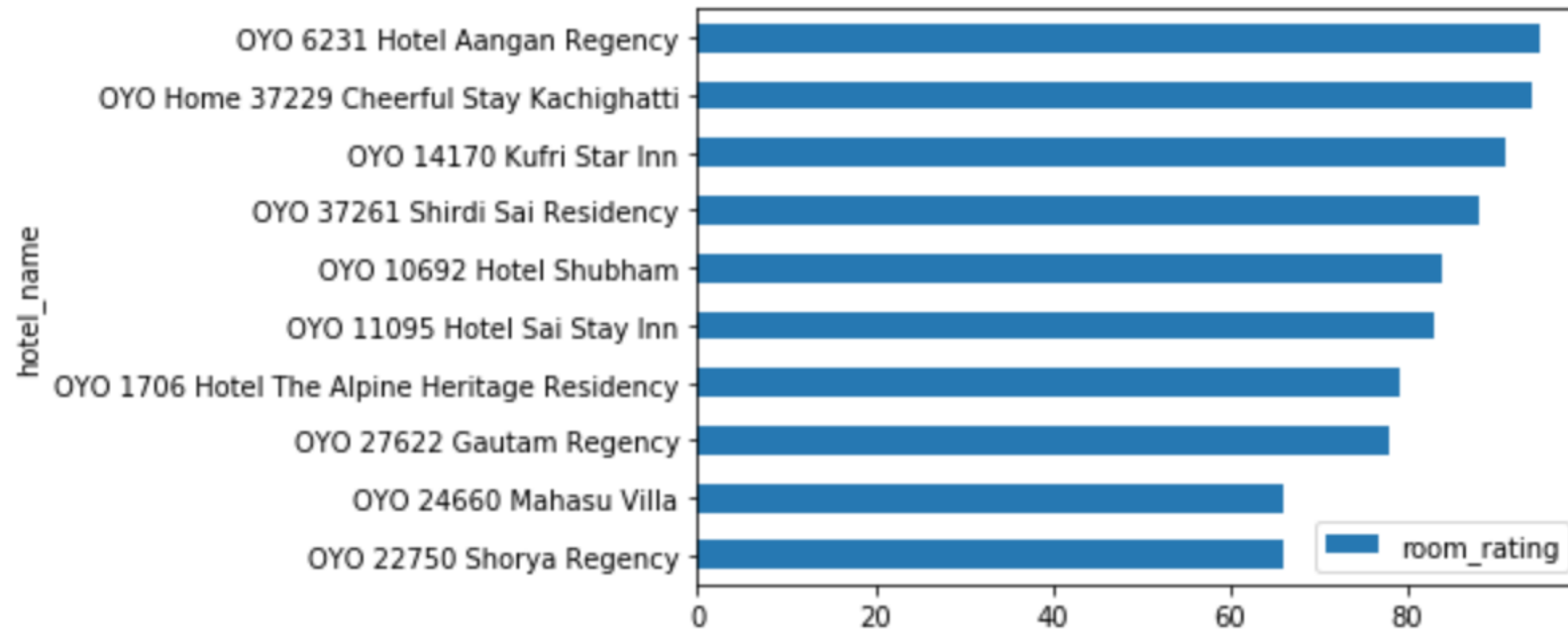
```python
fd.plot(kind="line", x="hotel_name", y="room_price", figsize=(12,5), color="green")
plt.show()
```

# Data Visualization

- Show the highest hotel rating

```
fgd = df.loc[:len(df)].sort_values(["room_rating"], ascending = True)
fgd.plot(kind="barh", x = "hotel_name", y="room_rating")
plt.show()
```

# Scrapping Images

# Scrape Images in Python

- In this section, we will scrape all the images from the same *goibibo* webpage. *goibibo.com/hotels/hotels-in-shimla-ct/*

- The first step would be same to navigate to the target website and download the source code.

- Next, we will find all the images using the **<img>** tag.

- To find all the images, we can use find_all() method

# Find All Images

```python
"""
Web Scraping - Scrap Images
"""

# importing required libraries
import requests
from bs4 import BeautifulSoup

# target URL
url = "https://www.goibibo.com/hotels/hotels-in-shimla-ct/"

headers = {
    'User-Agent': "Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHT
    }

response = requests.request("GET", url, headers=headers)

data = BeautifulSoup(response.text, 'html.parser')

# find all with the image tag
images = data.find_all('img', src=True)

print('Number of Images: ', len(images))

for image in images:
    print(image)
```

# Result

# Get data image source

- From all the image tags, select only the **src** part. Also, notice that the hotel images are available in **jpg** format. So we will select only those

```python
image_src = [x['src'] for x in images]

image_src = [x for x in image_src if x.endswith('.jpg')]

for image in image_src:
    print(image)
```

# Result



```
/home/lakshay/PycharmProjects/web-scraping/venv/bin/python /home/lakshay/PycharmProjects/web-scraping/images.py
Number of Images:  52
https://cdn1.goibibo.com/t_g/hotel-taj-view-shimla-_dsc4695-170163271449-orijgp.jpg
https://cdn1.goibibo.com/t_r/hotel-taj-view-shimla-img-20190421-wa0013-169761665138-orijgp.jpg
https://cdn1.goibibo.com/t_r/hotel-taj-view-shimla-_dsc4732-171263452301-orijgp.jpg
https://cdn1.goibibo.com/t_r/hotel-taj-view-shimla-_dsc4693-170163231235-orijgp.jpg
https://cdn1.goibibo.com/t_r/hotel-taj-view-shimla-_dsc4695-170163271449-orijgp.jpg
https://cdn1.goibibo.com/t_g/mahamaya-palace-shimla-1499008497093jpg-157737283046-orijgp.jpg
https://cdn1.goibibo.com/t_r/mahamaya-palace-shimla-1499008497093jpg-157737283046-orijgp.jpg
https://cdn1.goibibo.com/mahamaya-palace-shimla-1499008314001jpg-113437550491-jpeg-r.jpg
https://cdn1.goibibo.com/t_r/mahamaya-palace-shimla-mm8jpg-157737286146-orijgp.jpg
```
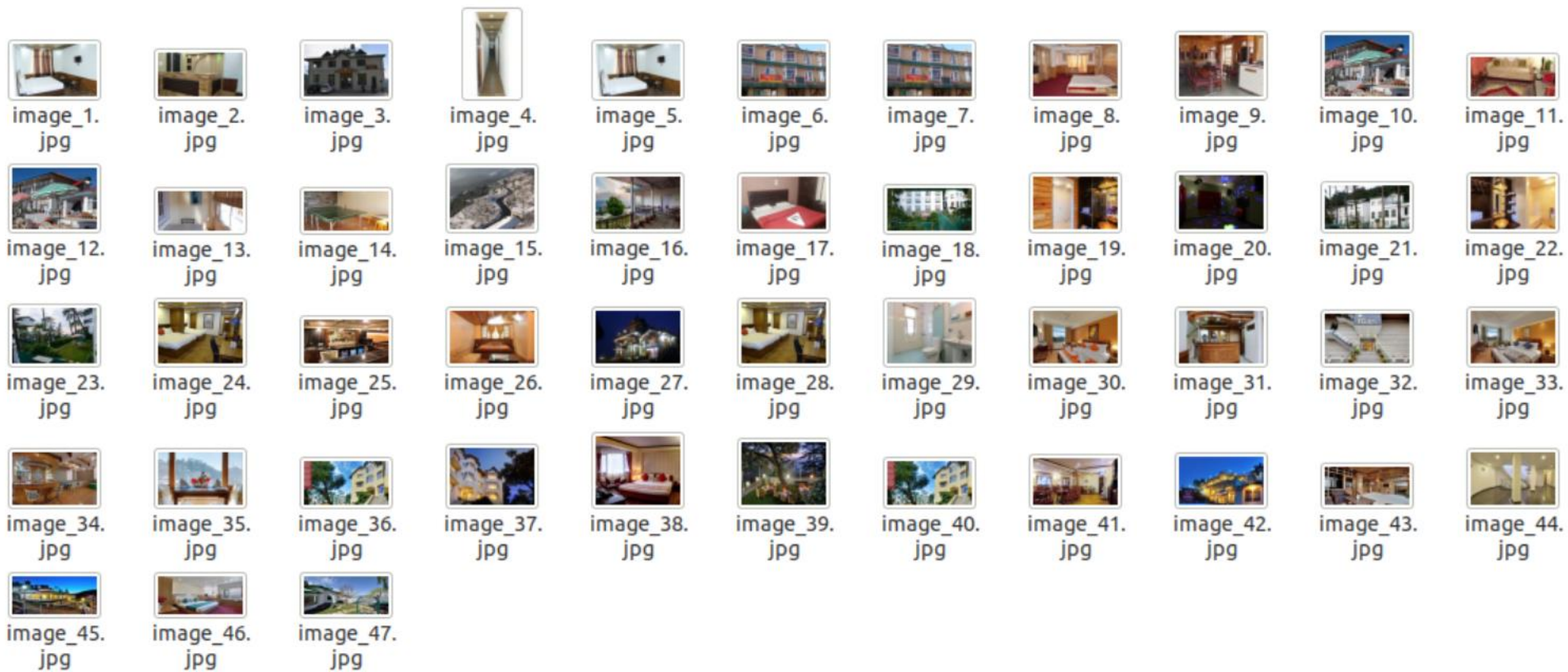
# Store Data

- Now that we have a list of image URLs, all we have to do is request the image content and write it in a file. Make sure that you open the file *'wb'* (write binary) form:

```python
image_count = 1
for image in image_src:
    with open('image_'+str(image_count)+'.jpg', 'wb') as f:
        res = requests.get(image)
        f.write(res.content)
    image_count = image_count+1
```

# Result

# THANK YOU