









Image Scrapper from scratch to proudction



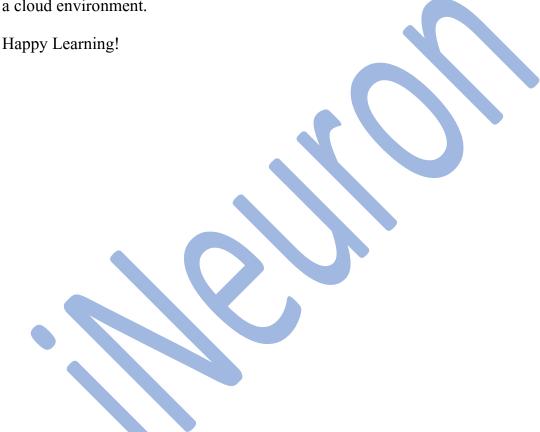
Table of Contents

Preface	 2
Introduction	 3
Prerequesites	4
PyCharm Installation	 4
MongoDB Installation	7
MongoDB Installation	7
Starting MongoDB	10
Application Architecture	 12
Implementation in Python	12
Heroku	17
Heroku Basics	17
Steps before cloud deployment	19
Heroku app creation and deploying the app to cloud	



Preface

This book is intended to help all the data scientists out there. It is a step by step guide for creating a web scraper, in this case, an image scrapper right from scratch and then deploying it to the heroku cloud platform. Image scrappers are extensively used in the industry today for collecting a huge number of images that are used as inputs for training the object detection/classification/identification models. This book takes a simple example of an online image search and tries to explain the concepts simply, extensively, and thoroughly to create a review scrapper right from scratch and then its deployment to a cloud environment.





Web Scraping(Images)

1. Introduction:

Web scraping is a technique using which the webpages from the internet are fetched and parsed to understand and extract specific information similar to a human being. Web scrapping consists of two parts:

- Web Crawling \rightarrow Accessing the webpages over the internet and pulling data from them.
- HTML Parsing Parsing the HTML content of the webpages obtained through web crawling and then extracting specific information from it.

Hence, web scrappers are applications/bots, which automatically send requests to websites and then extract the desired information from the website output.

Let's take an example:

how do we search for images online?

- We go to a website that shows images, and then we enter the keyword to search for the photos.
- The website shows us the images.
- We decide to download the image(s).

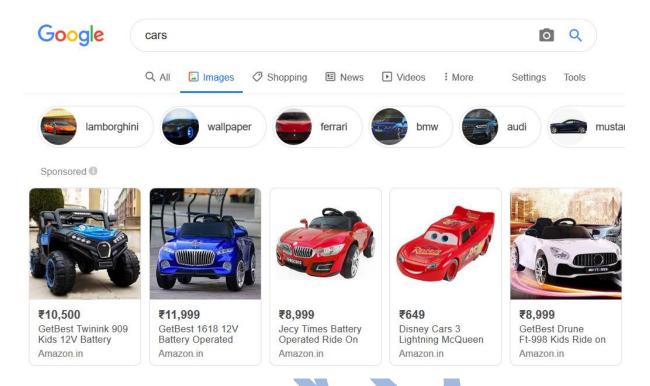
What if there is a computer program that can do all of these for us? That's what web scrappers necessarily do. They try to understand the webpage content as a human would do.

Let's take an example where we need thousands of images to train our object detection/classification model. It is an excellent use case for implementing an image scraper.

In this document, we'll take the example of searching for images and downloading them and take it further.

For example, if we open google.com and search for 'cars, the search result will be as follows:





Our end goal is to build a web scraper that collects the images for a keyword from the internet.

2. Prerequisites:

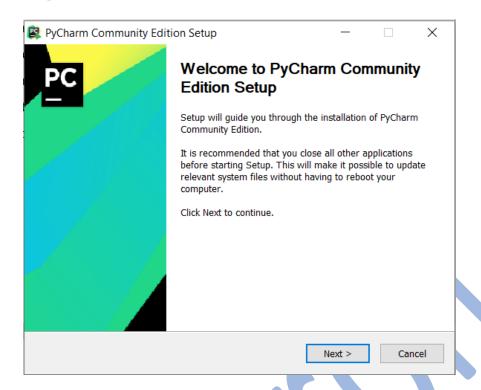
The things needed before we start building a python based web scraper are:

- Python installed.
- A Python IDE (Integrated Development Environment): like PyCharm, Spyder, or any other IDE of choice (Explained Later)
- Flask Installed. (A simple command: pip install flask)
- MongoDB installed (Explained Later).
- Basic understanding of Python and HTML.
- Basic understanding of Git (download Git CLI from https://gitforwindows.org/)

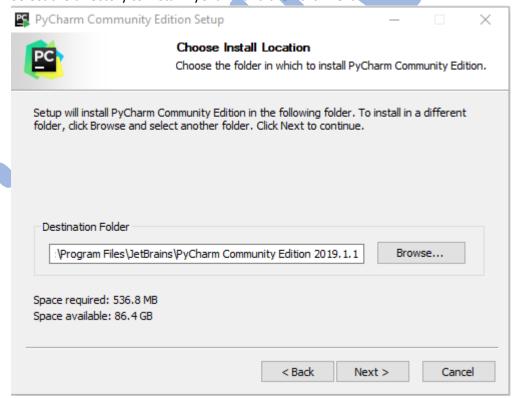
2.1 PyCharm Installation:

- a) Go to the link https://www.jetbrains.com/pycharm/download/#section=windows and download the community edition.
- b) Double click on the installation file to start the installation process and click next to continue.



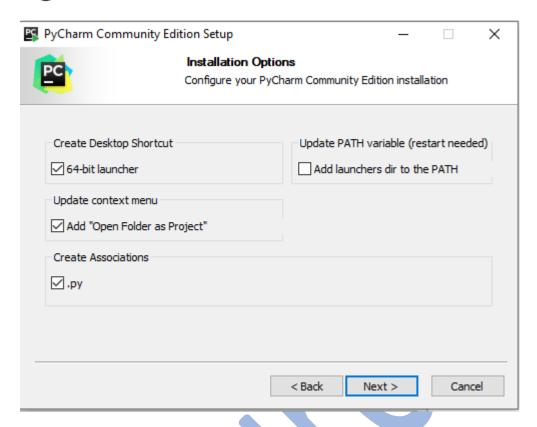


c) Select the directory to install PyCharm and then click next.

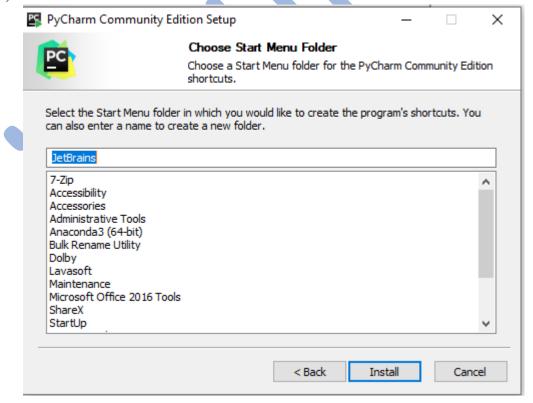


d) Check the appropriate checkboxes and then click next.





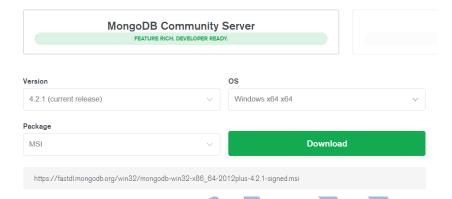
e) Choose the name of the start menu folder and then click on install to finish the installation.





2.2 MongoDB Installation:

1. Go to the page: https://www.mongodb.com/download-center/community and select the MongoDB installation to download based on your operating system.



2. After the installer gets downloaded, double click on the installer file to start installing the application.

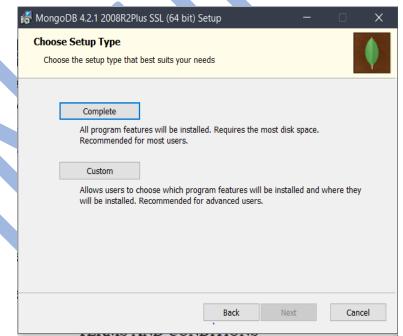




3. Click on the next button to move to the next step and accept the agreement.

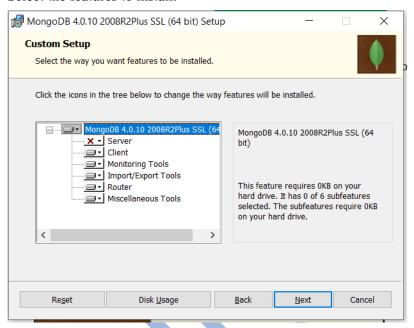


4. Select the type of installation:

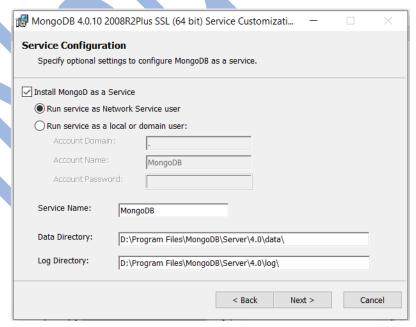




5. Select the features to install.

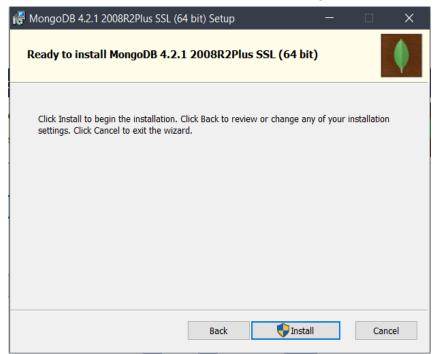


6. Click on next and then configure/customize the way you want the application to be installed.



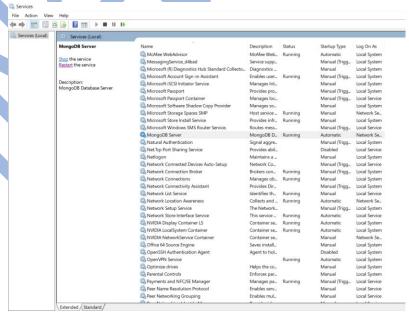


7. Click next and then click on install to start the MongoDB installation.



2.3 Starting MongoDB:

1. Go the services section and then start the MongoDB service if not already started.



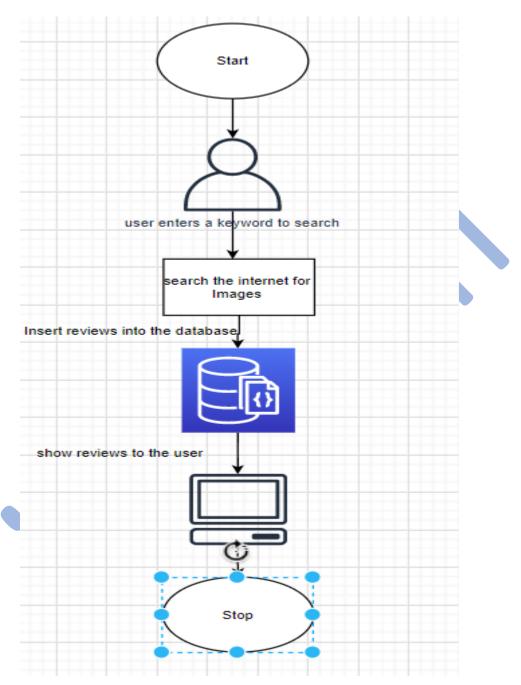
2. Now, to check whether the database server is up or not, go to the bin directory of the MongoDB installation and run the 'mongo' command as shown. If the command runs successfully, it means that the server is up and running, and we can proceed.





3. Application Architecture:

The architecture of the application is:



4. Implementation in Python:

Note: I have used PyCharm as an IDE for this documentation

- 1. Let's create a folder called 'ImageScrapper' on our local machines.
- 2. Inside that folder, let's create two more folders called 'imagescrapper', 'imagescrapperservice', 'imagescrapperutils', 'static', and 'templates'



to hold the code for the UI of our application. Inside 'static', let's create a folder called 'css' for keeping the stylesheets for our UI.

3. Inside 'imagescrapper', 'imagescrapperservice', and 'imagescrapperutils' create the files 'ImageScrapper.py', 'ImageScrapperService.py' and 'ImageScrapperUtils.py' respectively. The files are attached here for reference.



4. Let's create a file called 'app.py' inside the 'ImageScrapper' folder. The file is attached here for reference.



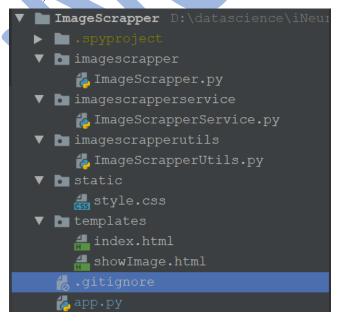
5. Inside the folder 'css', create the file: 'style.css'. The file is attached here for reference.



6. Inside the folder 'templates', create three HTML files called: 'index.html', and 'showImage.html'. The files are attached here for reference.



- showImage.html > Page to show the images for the searched keyword.
- 7. Now, the folder structure should look like:



8. Now, let's understand the flow:



- a) When the application starts, the user sees the page called 'index.html'.
- b) The user enters the search keyword into the search box and presses the submit button.
- c) The application now searches for images and shows the result on the 'showImage.html' page.
- 9. Understanding 'app.py'.
 - a) Import the necessary libraries:

```
from flask_cors import CORS,cross_origin
from imagescrapperservice.ImageScrapperService import
ImageScrapperService
from imagescrapper.ImageScrapper import ImageScrapper
from flask import Flask, render_template, request,jsonify
```

b) Initialize the flask app

```
app = Flask(__name__) # initialising the flask app with the name
'app'
```

- c) Creating the routes to redirect the control inside the application itself. Based on the route path, the control gets transferred inside the application. Let's understand the various routes:
 - i. The route for redirecting to home page

ii. Showing the images on the screen once our parser successfully gives the list of images.

```
@app.route('/showImages') # route to show the images on a
@cross_origin()
                                            show_images():
    scraper_object=ImageScrapper()
object
                             class
                                             ImageScrapper
list_of_jpg_files=scraper_object.list_only_jpg_files('st
atic') # obtaining the list of image files from the static
    print(list_of_jpg files)
        if(len(list_of_jpg_files)>0):
                                                there
          present,
            return
render_template('showImage.html',user_images
list_of_jpg_files)
                                                    folder
static
                      Exception
```



```
The route to actually scrape the web for images and then
iii.
      preparing the list of images and then calling the method in step
      (ii) to show the images to the user.
@app.route('/searchImages',
                                 methods=['GET','POST'
                                         searchImages():
                                                  'POST':
               request.method
       print("entered
        keyWord = request.form['keyword'] # assigning the
    else:
                                                  post")
       print("did
                                                keyWord)
    scraper object = ImageScrapper() # instantiating the
class
    list of jpg files
scraper_object.list_only_jpg_files('static') # obtaining
scraper_object.delete_existing_image(list_of_jpg_files)
# deleting the old image files stored from the previous
search
                                        keyWord.split()
    image_name
                                    '+'.join(image_name)
    image_name
    header
Chrome/43.0.2357.134
                                         Safari/537.36"}
    service = ImageScrapperService # instantiating the
                                   ImageScrapperService
obiect
        (imageURLList, header,
                                   keyWord, fileLoc)
   masterListOfImages = service.downloadImages(keyWord,
header) # getting the master list from keyword
    imageList = masterListOfImages[0] # extracting the
list of images from the master list
    imageTypeList = masterListOfImages[1] # extracting
    response = "We have downloaded ", len(imageList),
              " + image_name
    return show_images() # redirect the control to the
      The route to send the list of image URLs when the API is not
      called from a web browser.
  @app.route('/api/showImages', methods=['GET','POST
  # route to return the list of file locations for API
```



```
calls
@cross_origin()
                                      get_image_url():
              request.method
                                                post")
        print("entered
        keyWord = request.json['keyword'] # assigning
the value of the input keyword to the variable keyword
   else:
        print("Did
                                                post")
                                              keyWord)
                                                 words
                                       keyWord.split()
   image_name
                                  '+'.join(image_name)
                                              metadata
   header
                                        like Gecko)
                                       Safari/537.36"}
   service = ImageScrapperService # instantiating the
   url_enum = service.get_image_urls(keyWord, header)
                                           enumeration
   url_list=[] # initializing and empty url list
   for i, (img, Type) in enumerate(url_enum[0:5]):
       dict={'image_url':img}
       url_list.append(dict)
   return jsonify(url_list) # send the url list in JSON
```

d) After this, we'll just run our python app on our local system, and it'll start scraping for images as shown below:

Home Page:



Search Results:





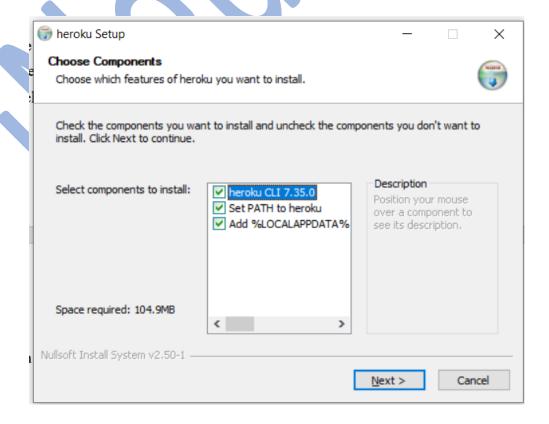
5. Heroku:

The Python app that we have developed is residing on our local machine. But to make it available to end-users, we need to deploy it to either an on-premise server or to a cloud service. Heroku is one such cloud service provider. It is free to use(till 5 applications).

We'll deploy this application to the Heroku cloud, and then anybody with the URL can then consume our app.

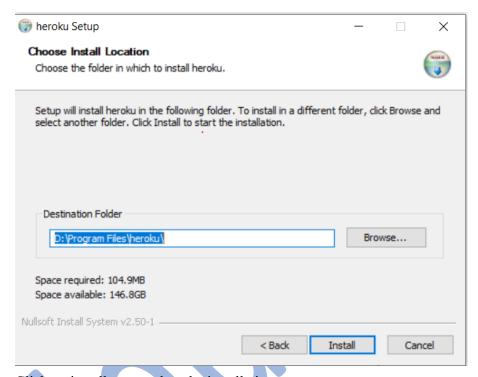
6. Heroku Basics:

- We'll first go to heroku.com, and we'll create a new account if we already don't have one.
- We'll download and install the Heroku CLI from the Heroku website: https://devcenter.heroku.com/articles/heroku-cli.
- Double-click the installation file and the following window shall appear:





• Click on next and select the installation directory for the CLI.



• Click on install to complete the installation.

7. Steps before cloud deployment:

We need to change our code a bit so that it works unhindered on the cloud, as well.

a) Add a file called 'gitignore' inside the 'reviewScrapper' folder. This folder contains the list of the files which we don't want to include in the git repository. My gitignore file looks like:

.idea
As I am using PyCharm as an IDE, and it's provided by the Intellij Idea community, it automatically adds the .idea folder containing some metadata. We need not include them in our cloud app.

b) Add a file called 'Procfile' inside the 'reviewScrapper' folder. This folder contains the command to run the flask application once deployed on the server:

web: gunicorn app:app

Here, the keyword 'web' specifies that the application is a web application. And the part 'app:app' instructs the program to look for a flask application called 'app' inside



the 'app.py' file. Gunicorn is a Web Server Gateway Interface (WSGI) HTTP server for Python.

c) Open a command prompt window and navigate to your 'reviewScrapper' folder. Enter the command 'pip freeze > requirements.txt'. This command generates the 'requirements.txt' file. My requirements.txt looks like:

```
beautifulsoup4==4.8.1
bs4==0.0.1
certifi==2019.9.11
Click==7.0
Flask==1.1.1
Flask-Cors==3.0.8
gunicorn==20.0.4
itsdangerous==1.1.0
Jinja2 = 2.10.3
MarkupSafe==1.1.1
numpy = 1.17.4
opency-python==4.1.2.30
Pillow==6.2.1
pymongo == 3.9.0
requests==2.21.0
requests-oauthlib==1.2.0
six = 1.13.0
soupsieve==1.9.5
Werkzeug==0.16.0
requirements.txt helps the Heroku cloud app to install all the dependencies before starting
the webserver.
```

d) A change has been made to app.py. The part where the port number for local machine was provided has been commented and the part without port number has been uncommented to run on the cloud.

```
if __name__ == "__main__":
    #app.run(host='127.0.0.1', port=8000) # port to run on local machine
    app.run(debug=True) # to run on cloud
```

8. Heroku app creation and deploying the app to cloud:

- a. After installing the Heroku CLI, Open a command prompt window and navigate to your 'reviewScrapper' folder.
- b. Type the command 'heroku login' to login to your heroku account as shown below:

```
D:\datascience\iNeuron\docs\reviewScrapper>heroku login
heroku: Press any key to open up the browser to login or q to exit:
```

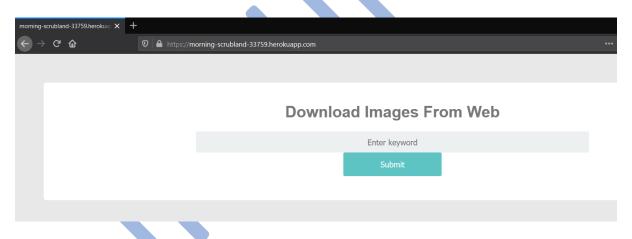
- c. After logging in to Heroku, enter the command 'heroku create' to create a heroku app. It will give you the URL of your Heroku app after successful creation.
- d. Before deploying the code to the Heroku cloud, we need to commit the changes to the local git repository.
- e. Type the command 'git init to initialize a local git repository as shown below:



D:\datascience\iNeuron\docs\reviewScrapper>git init

- f. Enter the command 'git status' to see the uncommitted changes
- g. Enter the command 'git add .' to add the uncommitted changes to the local repository.
- h. Enter the command 'git commit -am "make it better"' to commit the changes to the local repository.
- i. Enter the command 'git push heroku master' to push the code to the heroku cloud.
- j. After deployment, heroku gives you the URL to hit the web API.
- k. Once your application is deployed successfully, enter the command 'heroku logs --tail' to see the logs.

Final Result:



Thank You!