




# Scrape and Analytics Workshop

V1.0



# Agenda

## Scrape and Analytics

- Preparation
  - Scraping
  - Data Analysis
  - Assignments
- 



# **Part 1**

# **Workshop Preparation**

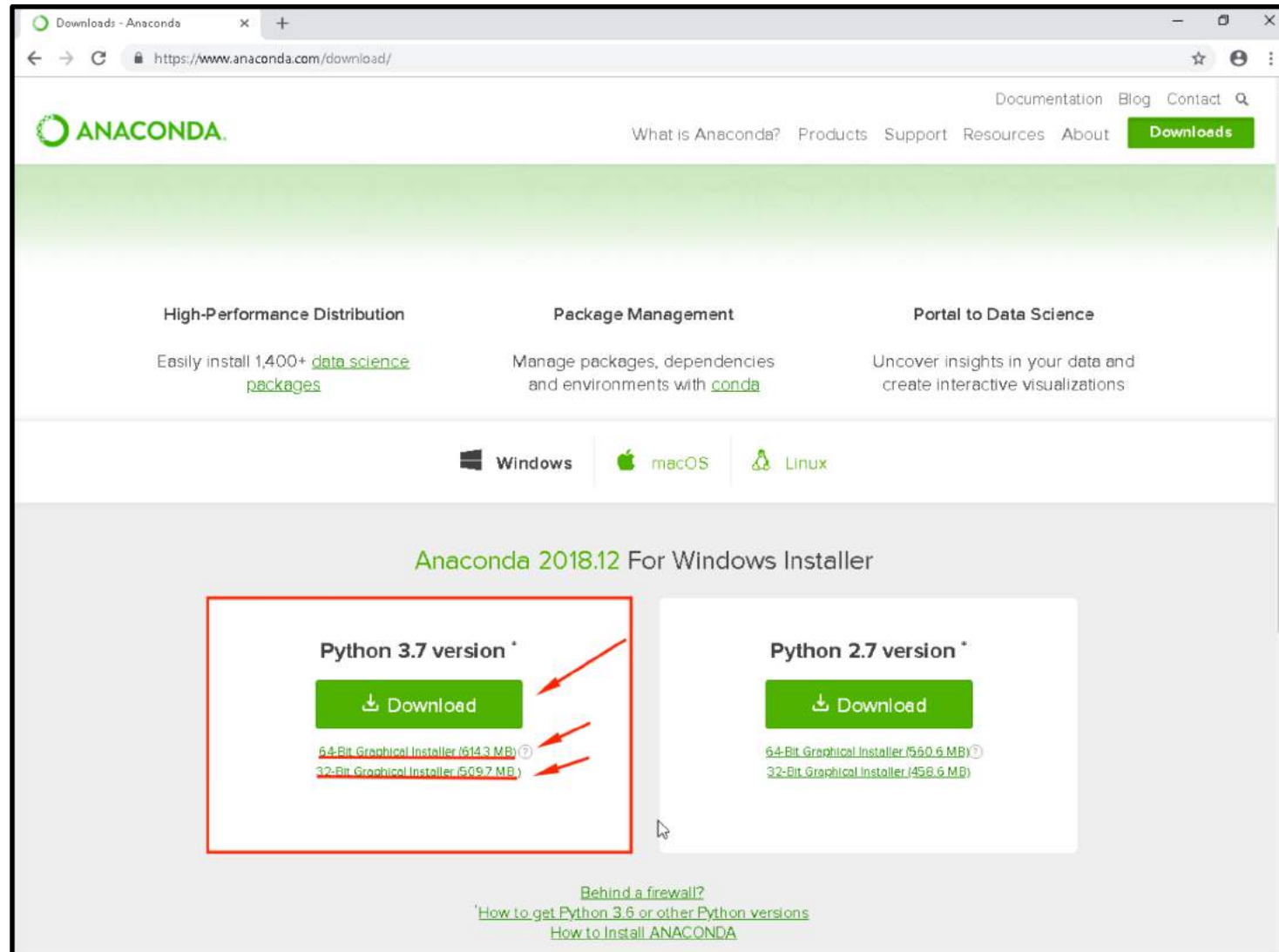


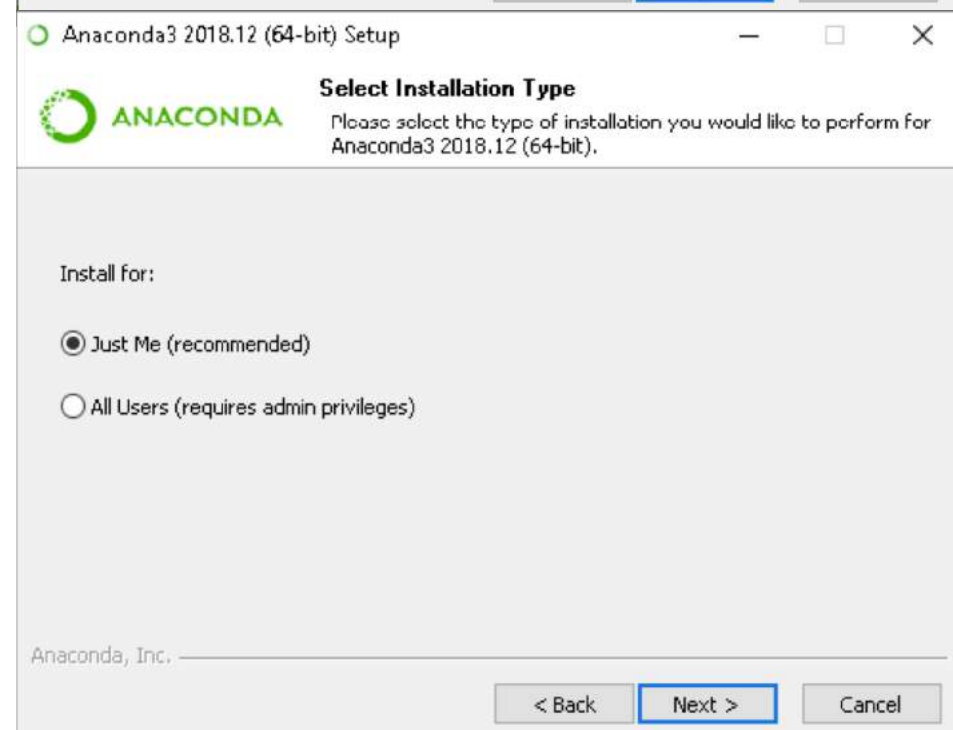
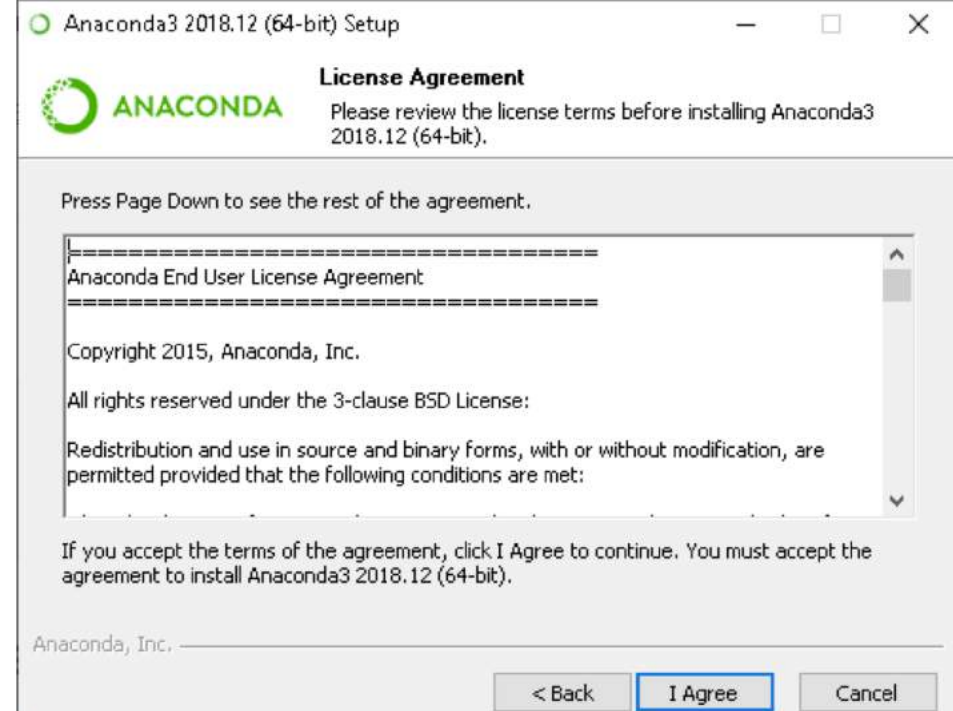
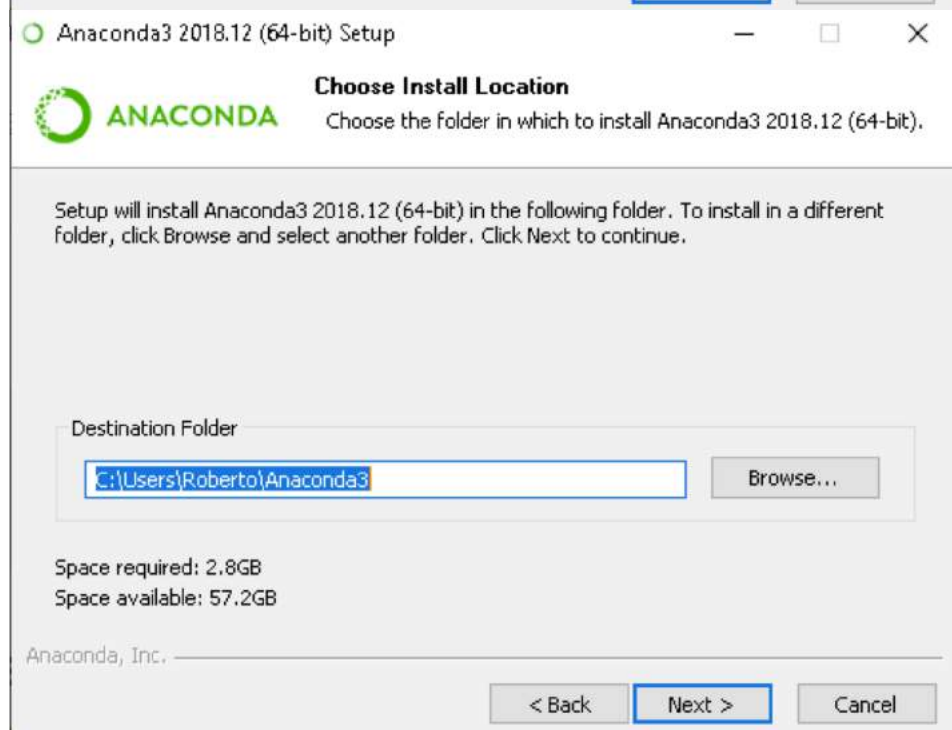


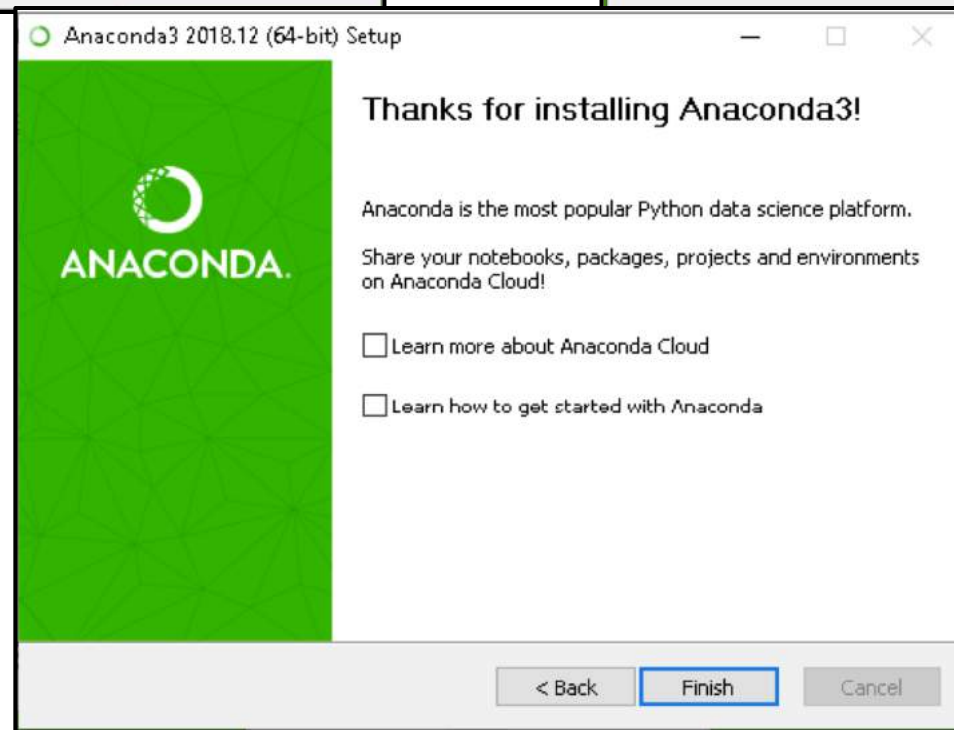
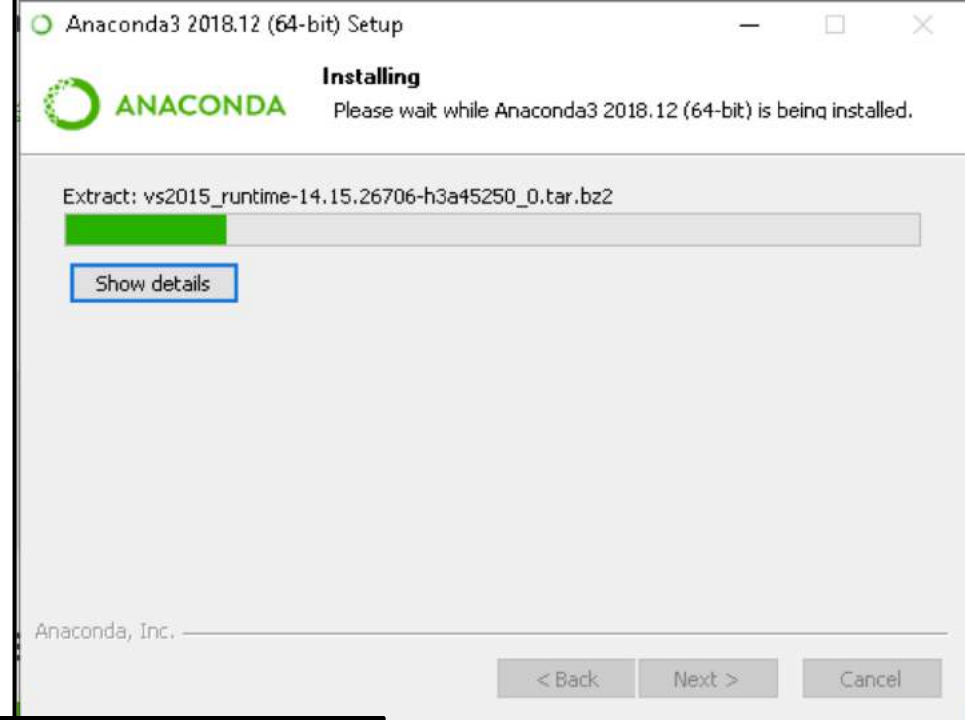
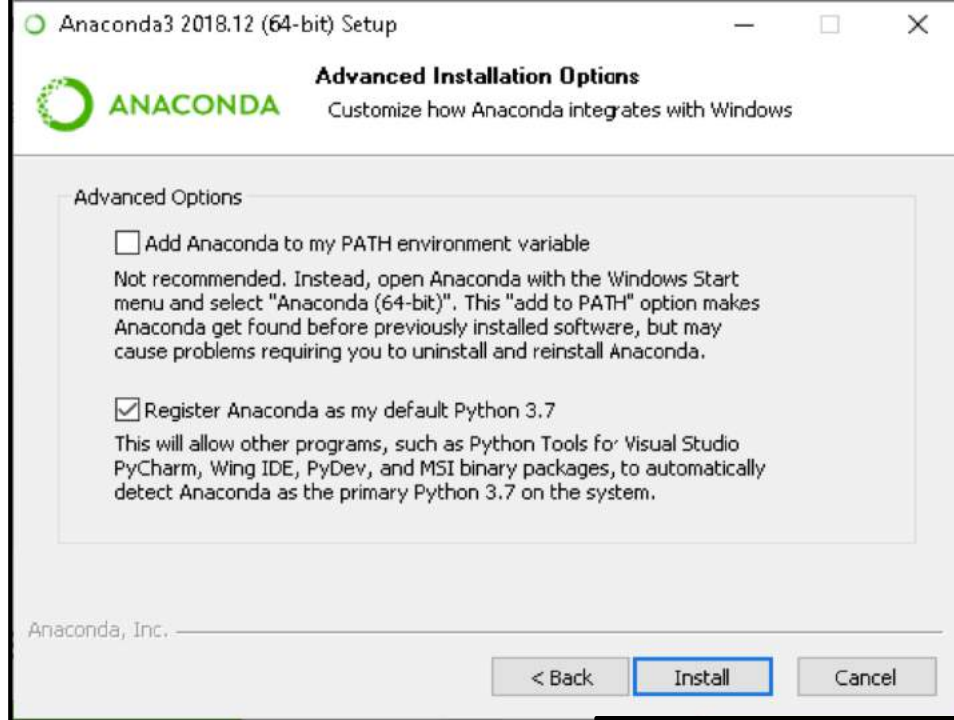
# Environment

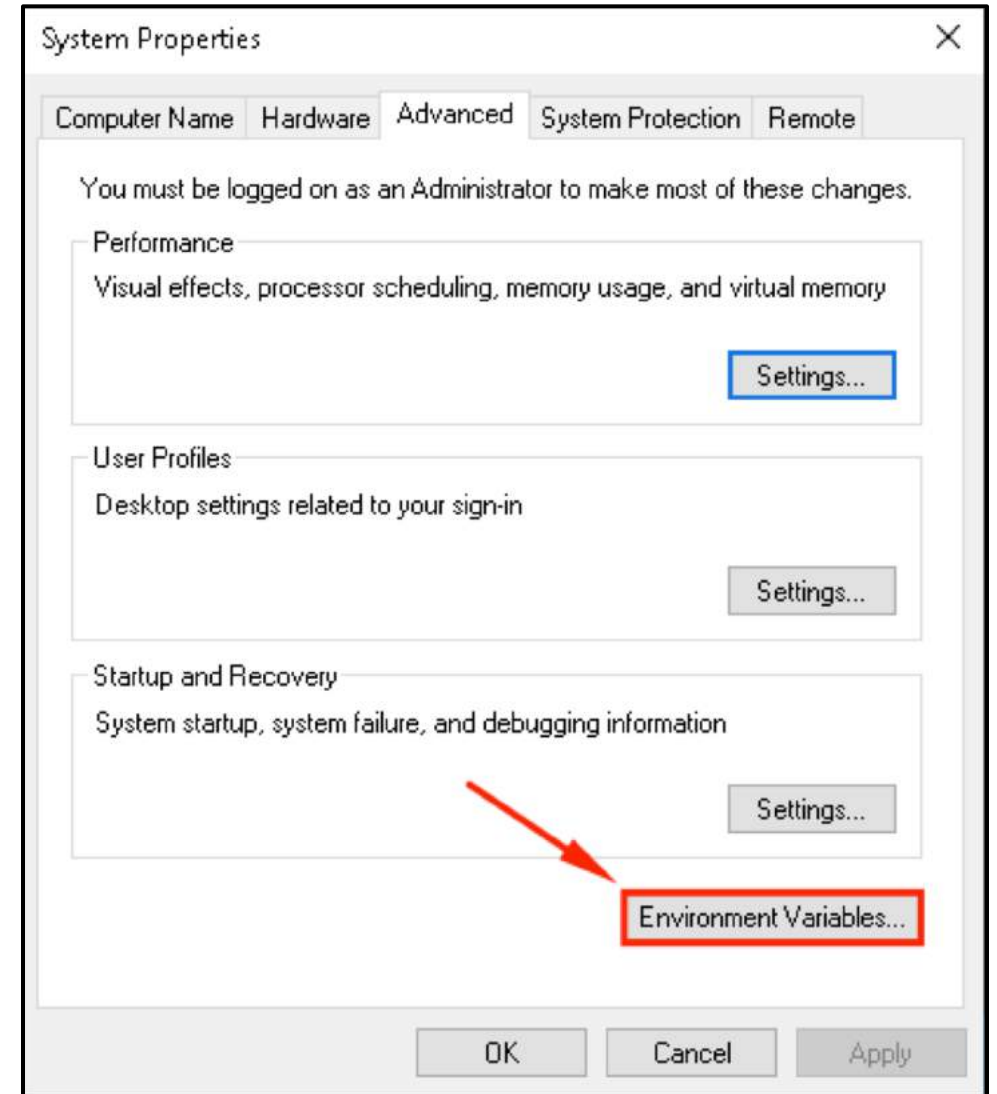
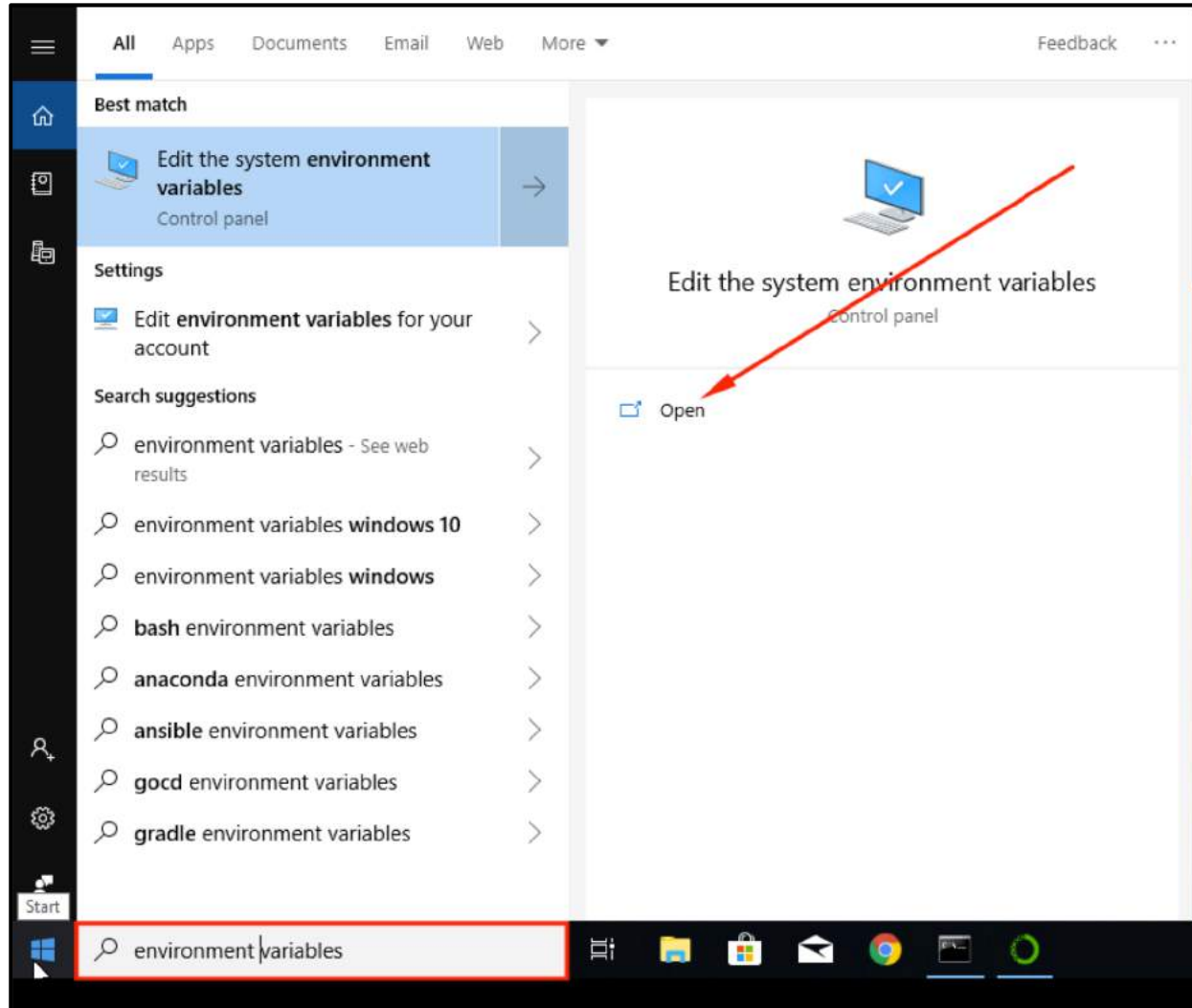


# www.anaconda.com

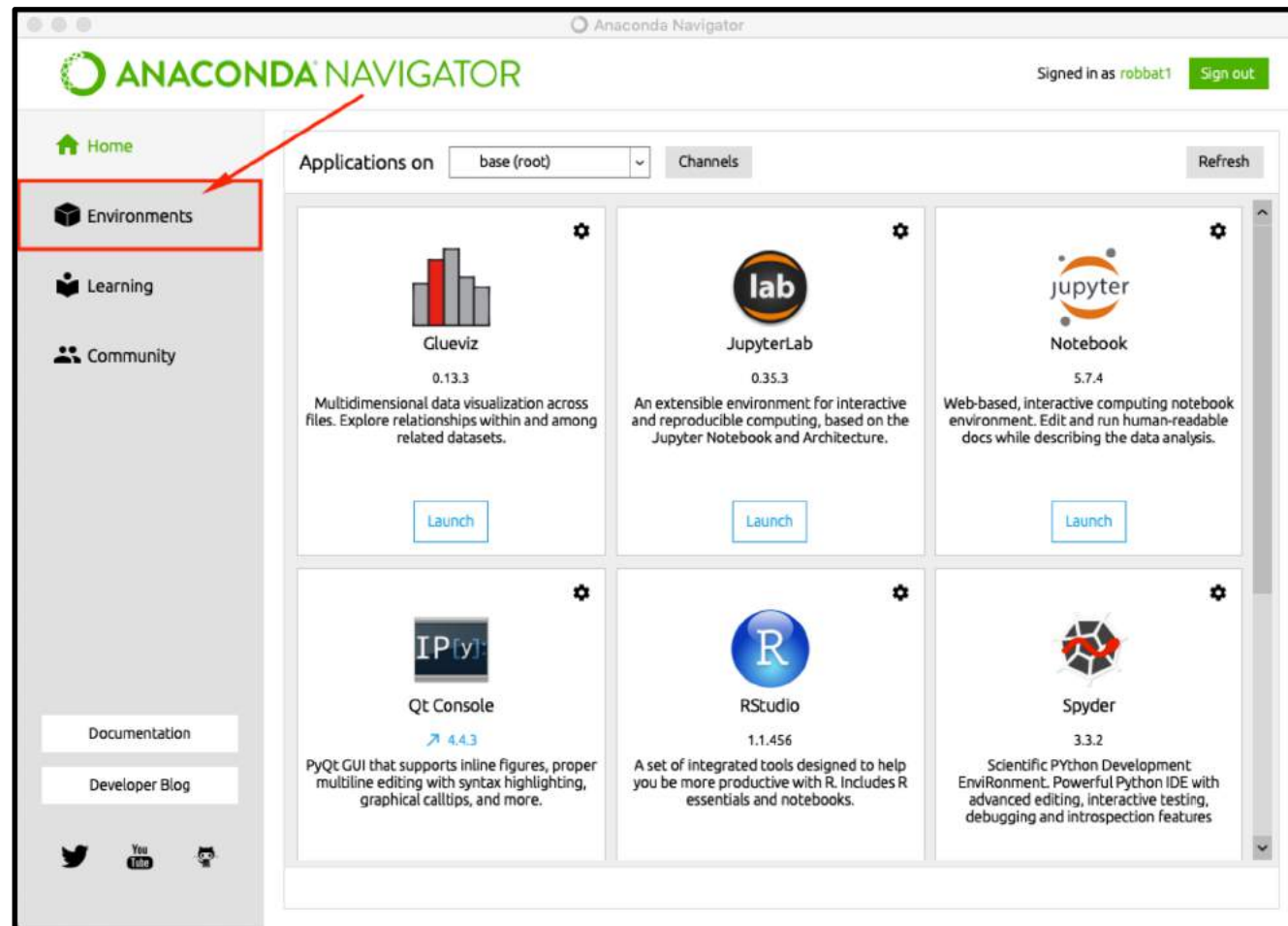
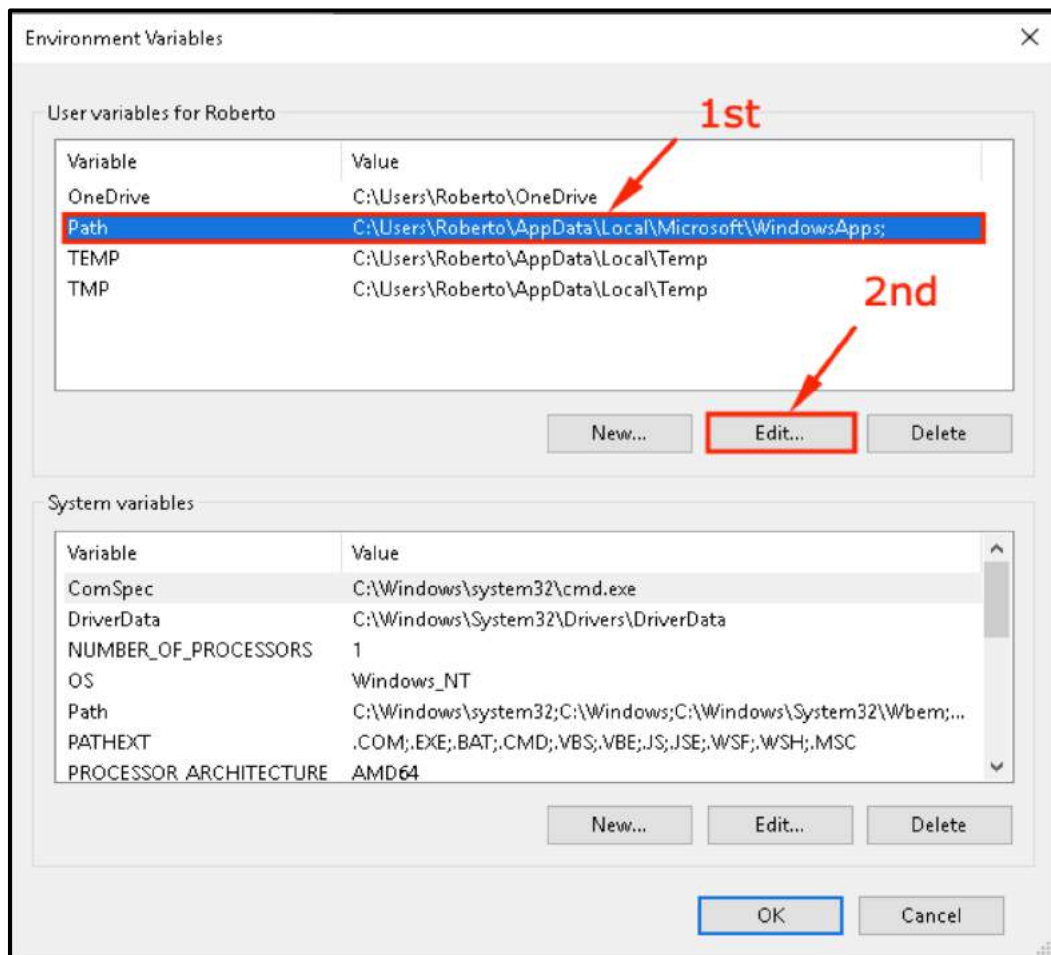












Anaconda Navigator

Signed in as robbat1 [Sign out](#)

Home

Environments

Learning

Community

Documentation

Developer Blog

Twitter YouTube GitHub

Create Clone Import Remove

Search Environments

base (root)

Installed

Channels

Update index...

Search Packages

Create new environment

Name: Scrapy

Location: /anaconda3/envs/Scrapy

Packages: ☒ Python ☐ R

3.7

3.7

3.6

3.5

2.7

Cancel Create

Name	Description	Version
<input checked="" type="checkbox"/> _ipyw_jlab_nb_ex...	A configuration metapackage for enabling anaconda-bundled jupyter extensions	0.1.0
<input checked="" type="checkbox"/> astroid	A abstract syntax tree for python with inference support.	2.1.0
<input checked="" type="checkbox"/> astropy	Community-developed python library for astronomy	3.0.5
<input checked="" type="checkbox"/> atomicwrites	Atomic file writes.	1.2.1
<input checked="" type="checkbox"/> attrs	Attrs is the python package that will b...	18.2.0
<input checked="" type="checkbox"/> automat	Self-service finite-state machines for the programmer on the go	0.7.0

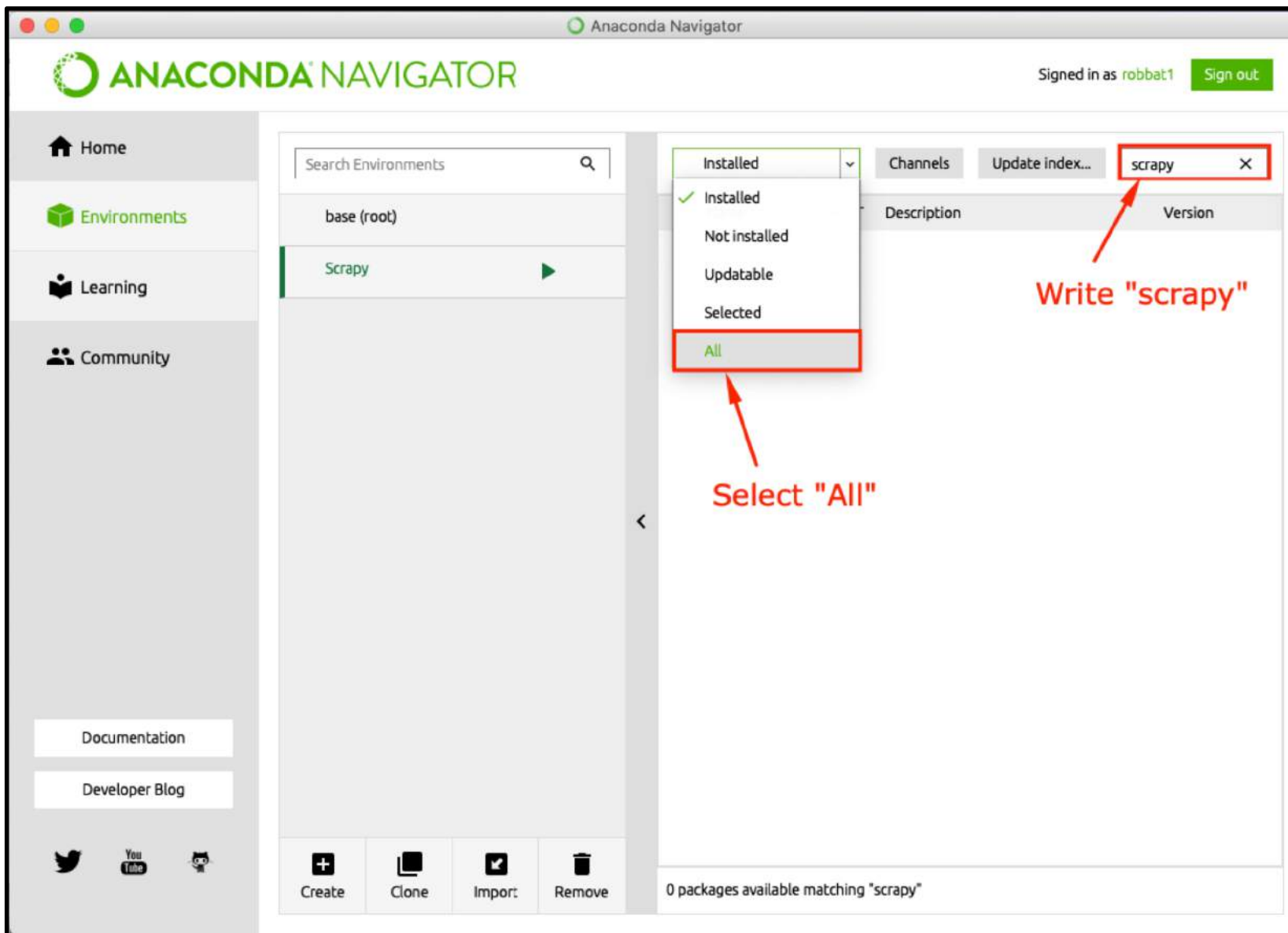
491 packages available

1

2

3

4



Home

Environments

Learning

Community

Documentation

Developer Blog

base (root)

Scrapy

CreateCloneImportRemove

Search Environments

All

Channels

Update index...

scrapy

Name	Description	Version
<input checked="" type="checkbox"/> scrapy	A high-level python screen scraping framework	1.5.1

1 package available matching "scrapy" 1 package selected

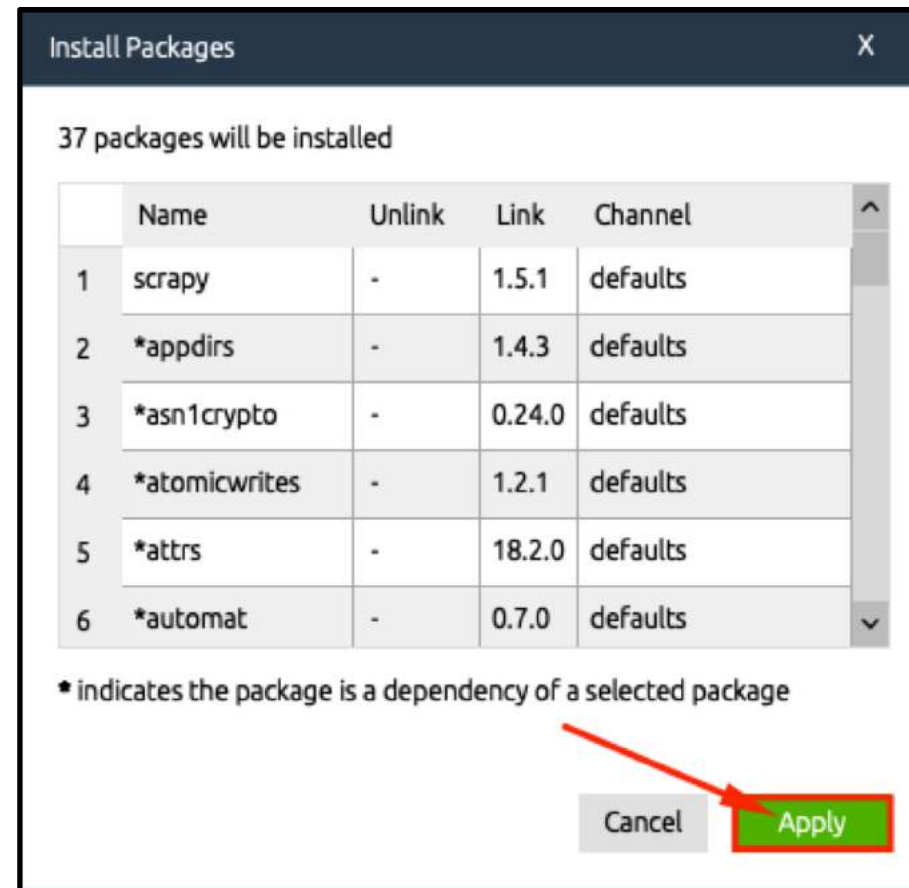
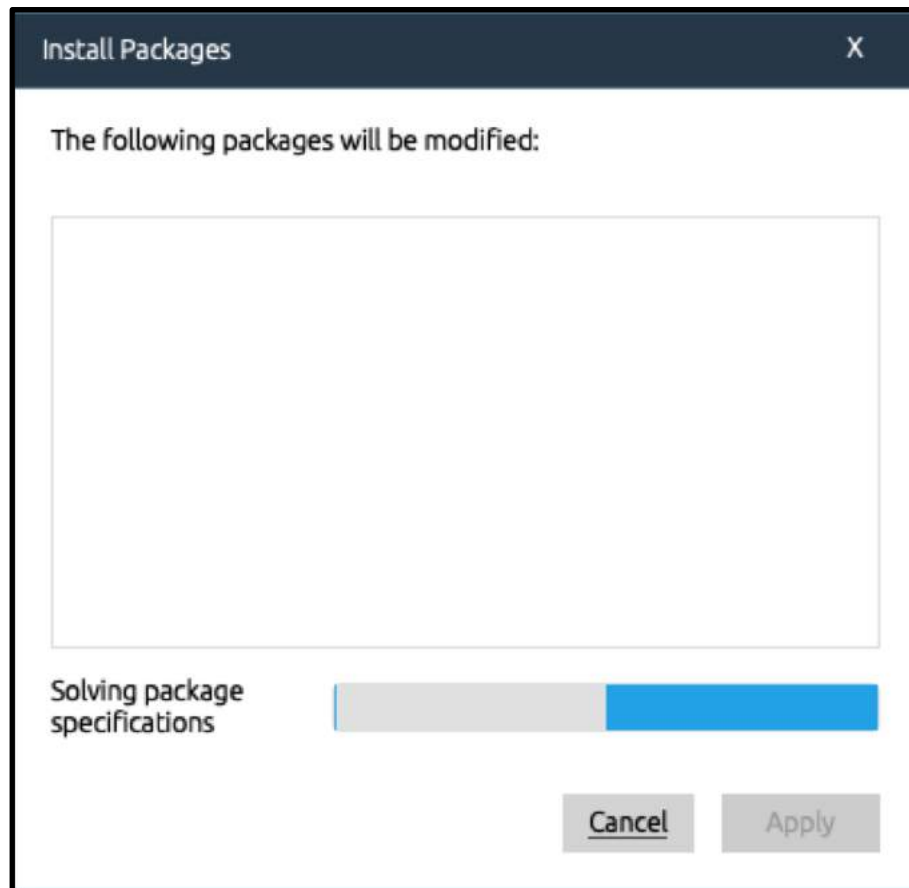
ApplyClear

Signed in as robbat1

Sign out

Check the box.

Then, click on "Apply"



Home




Environments

Learning

Community

Documentation

Developer Blog



Search Environments

base (root)

Scrapy

CreateCloneImportRemove

Anaconda Navigator

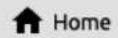
Signed in as robbat1

Sign out

AllChannelsUpdate index...scrapy

Name	T	Description	Version
scrapy		A high-level python screen scraping framework	1.5.1

Installing packages on /anaconda: CancelApplyClear



Home



Environments



Learning



Community

Documentation

Developer Blog



Search Environments



base (root)

Scrapy



Create



Clone



Import



Remove

All



Channels

Update index...

scrapy



Name



T

Description

Version



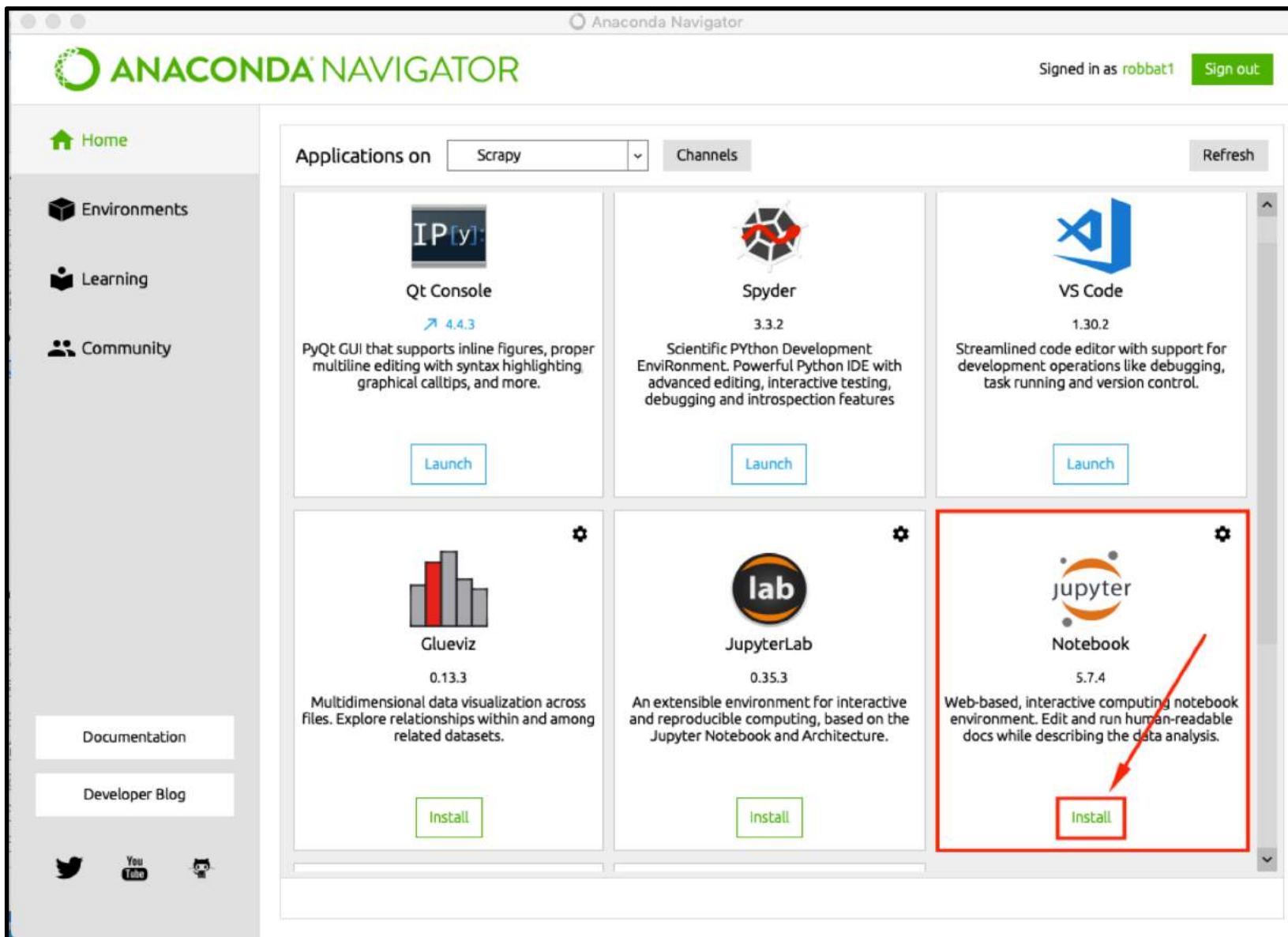
scrapy



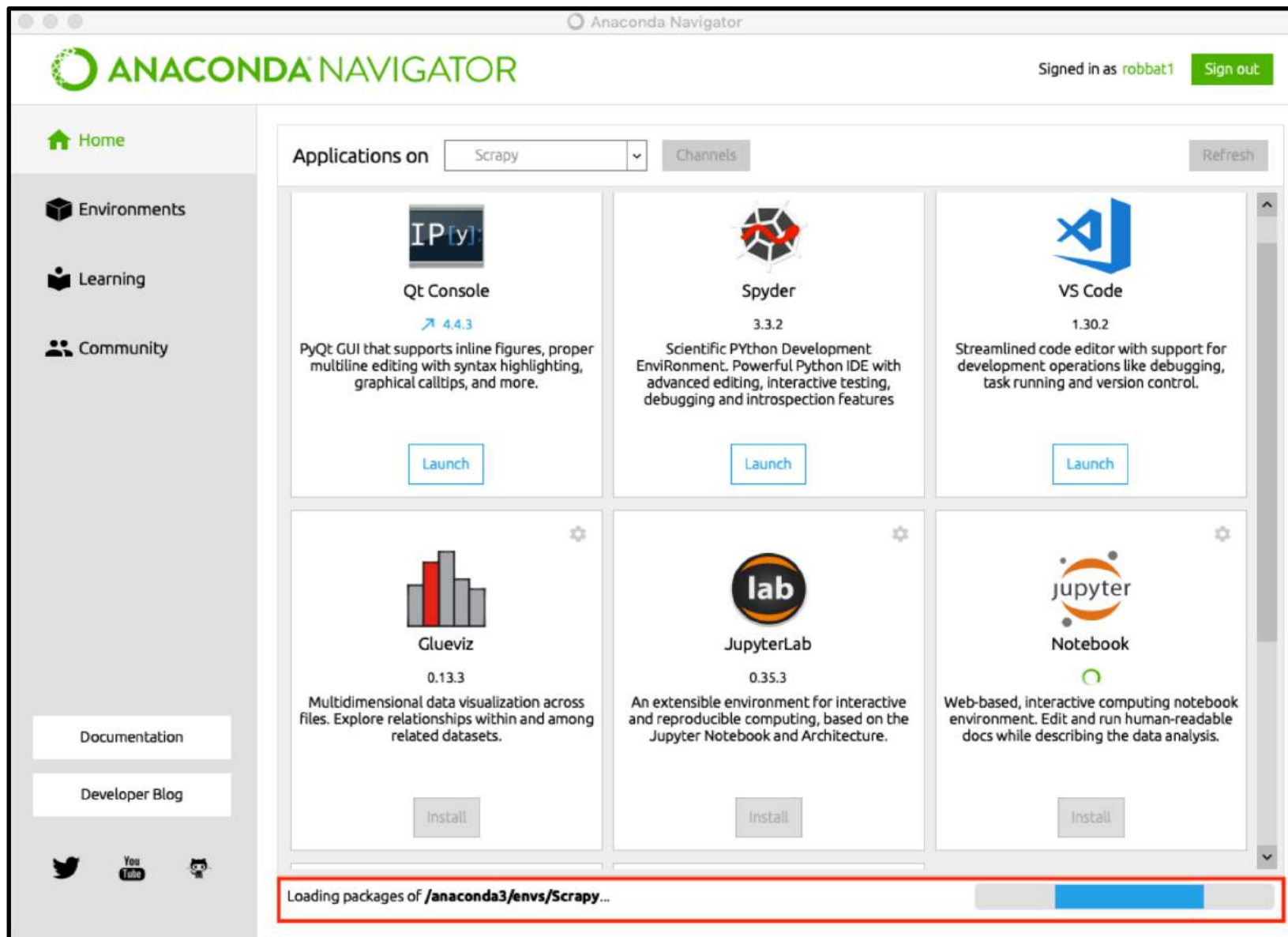
A high-level python screen scraping framework

1.5.1

1 package available matching "scrapy"





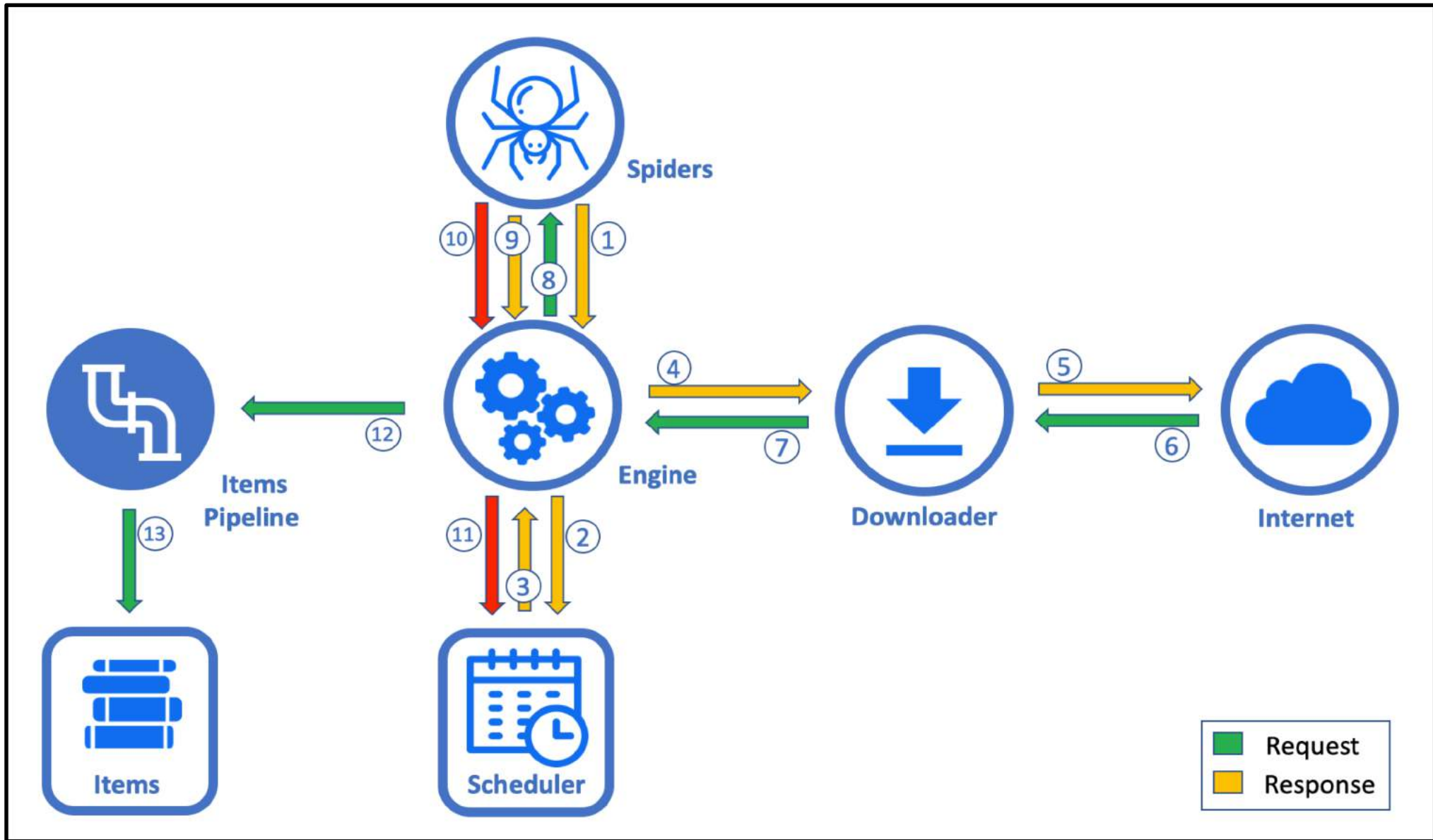




# **Part 2**

# **Scraping Introduction**








# Web site investigation





```
└─ smartcities
  └─ smartcities
    └─ __pycache__
      ├── __init__.cpython-36.pyc
      ├── items.cpython-36.pyc
      └── settings.cpython-36.pyc
    └─ spiders
      └─ __pycache__
        ├── __init__.cpython-36.pyc
        └── sc_job.cpython-36.pyc
      ├── __init__.py
      ├── sc_job.py
      ├── __init__.py
      ├── items.py
      ├── middlewares.py
      ├── pipelines.py
      ├── settings.py
      └── scrapy.cfg
```



```
$ scrapy shell
```

After execute the command above the final data showed in the terminal is the available Scrapy objects in the shell environment:


```
[s] Available Scrapy objects:
[s] scrapy scrapy module (contains scrapy.Request, scrapy.Selector, etc)
[s] crawler <scrapy.crawler.Crawler object at 0x1092933c8>
[s] item {}
[s] settings <scrapy.settings.Settings object at 0x10b825898>
[s] Useful shortcuts:
[s] fetch(url[, redirect=True]) Fetch URL and update local objects (by
    default, redirects are followed)
[s] fetch(req) Fetch a scrapy.Request and update local objects
[s] shelp() Shell help (print this help)
[s] view(response) View response in a browser
```

## Step One - Fetching the page

Once on the Scrapy Shell ">>>" we will use the 'fetch' command (short cut) to fetch the target web site. Which in this case will be the following search from www.ebay.com.

```
'https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&LH_ItemCondition=3&rt=nc&LH_BIN=1&_ipg=200'
```

```
>>>
fetch("https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&LH_ItemCondition=3&rt=nc&LH_BIN=1&_ipg=200")
2019-01-14 07:56:13 [scrapy.core.engine] INFO: Spider opened
2019-01-14 07:56:15 [scrapy.core.engine] DEBUG: Crawled (200) <GET
https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&LH_ItemCondition=3&rt=nc&LH_BIN=1&_ipg=200> (referer: None)
```



---

```
>>> response.url
'https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&LH_ItemCondition=3&rt=nc&LH_BIN=1&_ipg=200'
```

We could identify that the URL used by fetch was correct. The next verification is to retrieve the web page content.

```
>>> view(response)
True
```

---



Hi! Sign in My eBay

heart rate monitor

Related: heart rate monitor watch heart rate monitor polar heart rate monitor garmin...

All Listings Accepts Offers Auction Buy It Now

Best Match

71,569 results Save this search

New

**h3.s-item\_\_title | 300.39 x 38**

**New GARMIN HRM1G Run Sport Heart Rate Monitor With Chest Strap ANT+ USA Sent**

Brand New

★★★★★ 12 product ratings

**\$14.80**

or Best Offer

+\$2.80 shipping

**638 Sold**

2 new & refurbished from \$14.80

**F1 Fitness Blood Pressure Tracker Heart Rate Monitor Smart Wrist Bracelet Watch**

Brand New

**\$14.80**

Elements Console Sources Network Performance Memory 33

**<h3 class="s-item\_\_title" role="text">New GARMIN HRM1G Run Sport Heart Rate Monitor With Chest Strap ANT+ USA Sent</h3>**

Styles Event Listeners DOM Breakpoints Properties

Filter

element.style {

.s-item\_\_link:visited .s-item\_\_title {

color: #6a29b9;

.s-item\_\_link:visited>.s-item\_\_title {

color: #6a29b9;

.s-page .s-item\_\_link>.s-item\_\_title {

color: inherit;

@media only screen and (min-width: 480px)


.srp-list .s-item\_\_title {

max-height: none;

Console What's New

Failed to load resource: the server responded with a status of 404 ()

choices.truste.com/c...ext/javascript%22:1




---

The expression copied will be used in the Scrapy Shell to test if we in fact can extract the content of interest. Expression:

```
#srp-river-results-listing1 > div > div.s-item__info.clearfix > a > h3
```


Using the command response, we will test the expression above:

```
>>> response.css('#srp-river-results-listing1 > div > div.s-  
item__info.clearfix > a > h3')  
[<Selector xpath="descendant-or-self::*[@id = 'srp-river-results-  
listing1']/div/div[@class and contains(concat(' ', normalize-space(@class), '  
' ), ' s-item__info ' ) and (@class and contains(concat(' ', normalize-  
space(@class), ' '), ' clearfix '))]/a/h3" data='<h3 class="s-item__title"  
role="text">Ne'>]
```



The result given is not generalized to all the titles, but to the specific title we copy the selector. In order to make it generalized it is necessary to remove the parent node which is “#srp-river-results-listing1 >”. Let's execute without this part.

```
>>> response.css('div > div.s-item__info.clearfix > a > h3')
[<Selector xpath="descendant-or-self::div/div[@class and contains(concat(' ',
normalize-space(@class), ' '), ' s-item__info ') and (@class and
contains(concat(' ', normalize-space(@class), ' '), ' clearfix '))]/a/h3"
data='<h3 class="s-item__title" role="text">Ne'>,
...]
```




---

```
/" s-item__info ' ) and (@class and contains(concat(' ', normalize-  
space(@class), ' '), ' clearfix '))]/a/h3" data='<h3 class="s-item__title s-  
item__title--">]
```

The result above just show the first and last line of response. We can identify which the response potentially covers all the titles present in the document. The next step is to add the option `::text` in the expression and a function called `extract()` in order to clean-up the data retrieved.

---




---

```
>>> response.css('div > div.s-item__info.clearfix > a > h3::text').extract()
['New GARMIN HRM1G Run Sport Heart Rate Monitor With Chest Strap ANT+ USA
Sent', 'Sports Blood Pressure/Heart Rate Monitor Fitness Smart Watch Wrist
Band Bracelet', 'F1 Fitness Blood Pressure Tracker Heart Rate Monitor Smart
Wrist Bracelet Watch', 'Sports Blood Pressure/Heart Rate Monitor Fitness
Smart Watch Wrist Band Bracelet',
```

The response above shows few lines that the text titles extracted from the page. The next step is to collect the price from each item in the web page. Similar to the process of get the selector of the title we will proceed with the price.

---





```
>>> response.css('#srp-river-results-listing1 > div > div.s-  
item__info.clearfix > div.s-item__details.clearfix > div:nth-child(1) >  
span')  
[<Selector xpath="descendant-or-self::*[@id = 'srp-river-results-  
listing1']/div/div[@class and contains(concat(' ', normalize-space(@class), '  
' ), ' s-item__info ' ) and (@class and contains(concat(' ', normalize-  
space(@class), ' '), ' clearfix '))]/div[@class and contains(concat(' ',  
normalize-space(@class), ' '), ' s-item__details ' ) and (@class and  
contains(concat(' ', normalize-space(@class), ' '), ' clearfix  
' ) )]/div[count(preceding-sibling::*) = 0]/span" data='<span class="s-  
item__price">$14.80</span'>]
```

Once again the response was specific to the selected price, now we will generalize the expression removing the parent node, and add the entire command with `::text` and the function `extract()`.

```
>>> response.css('div > div.s-item__info.clearfix > div.s-  
item__details.clearfix > div:nth-child(1) > span::text').extract()  
['$14.80', '$15.99', '$14.81', '$14.92', '$15.99', '$13.45', '$22.95',  
'$13.45', '$13.50', '$49.99', '$22.78', '$56.49', '$8.90', '$10.99',  
'$13.49', '$16.18', '$13.58', '$14.24',
```



# Scraping



## Step 1 - Create a Project Directory

After we test our selector, we will create our Scrapy project in a new folder that we will create, for this example was created the folder 'SCRAPY\_PROJECT' in '/Users/username/'.

## Step 2 - Create a Project

To create a new Scrapy project we will use the command `scrapy startproject <project_name> [project_dir]`, in this case we will use 'ebay' as the 'project\_name', and the default directory by omitting the argument 'project\_dir'.

```
$ scrapy startproject ebay
New Scrapy project 'ebay', using template directory
'/anaconda3/lib/python3.6/site-packages/scrapy/templates/project', created
in:
  /Users/user_name/SCRAPY_PROJECTS/ebay
```

You can start your first spider with:

```
cd ebay
scrapy genspider example example.com
```



### Step 3 – Create a Spider

After create the project we will create our first spider inside of the directory defined (/Users/user\_name/SCRAPY\_PROJECTS/ebay) using the command scrapy genspider [options]  
<name> <domain>

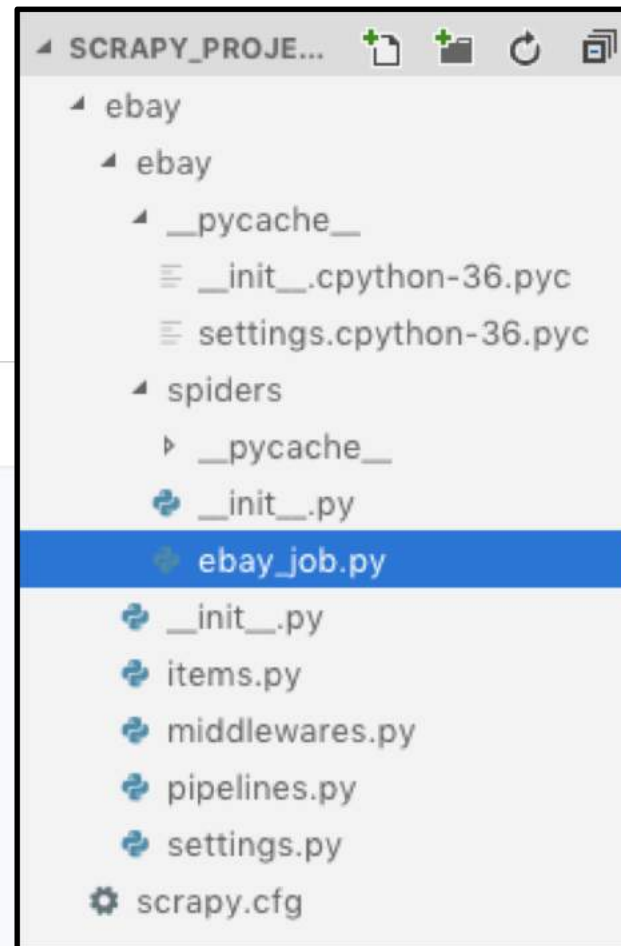
```
$ scrapy genspider ebay_job ebay.com
Created spider 'ebay_job' using template 'basic' in module:
ebay.spiders.ebay_job
```

Opening the file 'ebay\_job.py' we can see the default content created:

```
# -*- coding: utf-8 -*-
import scrapy

class EbayJobSpider(scrapy.Spider):
    name = 'ebay_job'
    allowed_domains = ['ebay.com']
    start_urls = ['http://ebay.com/']

    def parse(self, response):
        pass
```



We will edit few components in this file. First it is necessary to substitute the 'start url' by the link we investigated previously:

```
'https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&LH_ItemCondition=3&rt=nc&LH_BIN=1&ipg=200'
```

```
# -*- coding: utf-8 -*-
import scrapy
from ebay.items import EbayItem

class EbayJobSpider(scrapy.Spider):
    name = 'ebay_job'
    allowed_domains = ['ebay.com']
    start_urls =
["https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&
LH_ItemCondition=3&rt=nc&LH_BIN=1&_ipg=200"]

    def parse(self, response):
        titles = response.css("div > div.s-item__info.clearfix > a >
h3::text").extract()
        prices = response.css("div > div.s-item__info.clearfix > div.s-
item__details.clearfix > div:nth-child(1) > span").extract()

        for item in zip(titles, prices):
            new_item = EbayItem()

            new_item['titles'] = item[0]
            new_item['prices'] = item[1]

            yield new_item
```

```
# -*- coding: utf-8 -*-
import scrapy
from ebay.items import EbayItem

class EbayJobSpider(scrapy.Spider):
    name = 'ebay_job'
    allowed_domains = ['ebay.com']
    start_urls =
    ['https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&
    LH_ItemCondition=3&rt=nc&LH_BIN=1&ipg=200']


    def parse(self, response):

        titles = response.css('div > div.s-item__info.clearfix > a >
        h3::text').extract()
        prices = response.css('div > div.s-item__info.clearfix > div.s-
        item__details.clearfix > div:nth-child(1) > span').extract()

        for item in zip(titles, prices):
            new_item = EbayItem()

            new_item['titles'] = item[0]
            new_item['prices'] = item[1]

            yield new_item
```



## Step 6 – Crawling

After we set the scrapy framework structure it is time to run the Crawling function. To do this we will execute the following command:

```
$ scrapy crawl ebay_job
```





```
[...]
2019-01-20 20:12:32 [scrapy.core.scrapers] DEBUG: Scraped from <200
https://www.ebay.com/sch/i.html?_from=R40&_nkw=heart+rate+monitor&_sacat=0&LH
_ItemCondition=3&rt=nc&LH_BIN=1&ipg=200>
{'prices': '<span class="s-item__price">$13.39<span class="DEFAULT"> to '
          '</span>$13.79</span>',
 'titles': 'Original Xiaomi Mi Band 2 Smart Wristband Bracelet Heart Rate '
          'Monitor Smartwatch'}
2019-01-20 20:12:32 [scrapy.core.engine] INFO: Closing spider (finished)
2019-01-20 20:12:32 [scrapy.statscollectors] INFO: Dumping Scrapy stats:
{'downloader/request_bytes': 759,
 'downloader/request_count': 2,
 'downloader/request_method_count/GET': 2,
 'downloader/response_bytes': 103540,
 'downloader/response_count': 2,
 'downloader/response_status_count/200': 2,
 'finish_reason': 'finished',
 'finish_time': datetime.datetime(2019, 1, 21, 1, 12, 32, 735435),
 'item_scraped_count': 206,
 'log_count/DEBUG': 209,
 'log_count/INFO': 7,
 'memusage/max': 51384320,
 'memusage/startup': 51376128,
 'response_received_count': 2,
 'scheduler/dequeued': 1,
 'scheduler/dequeued/memory': 1,
 'scheduler/enqueued': 1,
 'scheduler/enqueued/memory': 1,
 'start_time': datetime.datetime(2019, 1, 21, 1, 12, 29, 486444)}
```



---

### **Step 7 – Exporting the extracted data for Data Analytics**

It is not enough the extraction of data, it is necessary to export it in a tabular dataset for further analyses. In order to do it is necessary to add arguments in the crawled command used:

```
$ scrapy crawl ebay_job -o ebay_extraction.csv -t csv
```



---



# Data Analytics







**Thank you**