

Lab Assignment- PDV

Scraped the website - <https://shekouwoman.com/collections/home>

To run Scrapy, install the libraries and create a folder to store the files needed for the scrapy project. The command `'scrapy startproject'` creates the required .py files and spider to scrape the website.

```
Anaconda Powershell Prompt (anaconda3)

(base) PS C:\Users\aditi> cd C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment
(base) PS C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment> scrapy startproject clothshopping
New Scrapy project 'clothshopping', using template directory 'C:\Users\aditi\anaconda3\lib\site-packages\scrapy\template\project', created in:
  C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment\clothshopping

You can start your first spider with:
  cd clothshopping
  scrapy genspider example example.com
(base) PS C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment> cd clothshopping
(base) PS C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment\clothshopping> cd clothshopping
(base) PS C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment\clothshopping\clothshopping> ls

Directory: C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment\clothshopping\clothshopping

Mode                LastWriteTime         Length Name
----                -
a---l             23-12-2021         16:19         spiders
a---l             28-01-2022         23:49          269 items.py
a---l             28-01-2022         23:49       3662 middlewares.py
a---l             28-01-2022         23:49          367 pipelines.py
a---l             28-01-2022         23:49       3134 settings.py
a---l             12-10-2020         17:32           0 __init__.py

(base) PS C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment\clothshopping\clothshopping>
```


Scraping the website directly in the shell:

The `'scrapy shell'` command is used to scrape the website on the shell.

```
Anaconda Powershell Prompt (anaconda3)

(base) PS C:\Users\aditi\OneDrive\Desktop\ScrapyAssignment\clothshopping\clothshopping> scrapy shell "https://shekouwoman.com/collections/home"
2022-01-28 23:51:00 [scrapy.utils.log] INFO: Scrapy 2.4.0 started (bot: clothshopping)
2022-01-28 23:51:00 [scrapy.utils.log] INFO: Versions: lxml 4.6.3.0, libxml2 2.9.10, cssselect 1.1.0, parsel 1.5.2, w3lib 1.21.0, Twisted 20.3.0, Python 3.8.8 (default, Apr 13 2021, 15:08:03) [MSC v.1916 64 bit (AMD64)], pyOpenSSL 20.0.1 (OpenSSL 1.1.1k 25 Mar 2021), cryptography 3.4.7, Platform Windows-10-10.0.22000-SP0
2022-01-28 23:51:00 [scrapy.utils.log] DEBUG: Using reactor: twisted.internet.selectreactor.SelectReactor
2022-01-28 23:51:00 [scrapy.crawler] INFO: Overridden settings:
{'BOT_NAME': 'clothshopping',
 'DUPEFILTER_CLASS': 'scrapy.dupefilters.BaseDupeFilter',
 'LOGSTATS_INTERVAL': 0,
 'NEWSPIDER_MODULE': 'clothshopping.spiders',
 'ROBOTSTXT_OBEY': True,
 'SPIDER_MODULES': ['clothshopping.spiders']}
2022-01-28 23:51:00 [scrapy.extensions.telnet] INFO: Telnet Password: 296616c9baa65c3e
2022-01-28 23:51:00 [scrapy.middleware] INFO: Enabled extensions:
['scrapy.extensions.corestats.CoreStats',
 'scrapy.extensions.telnet.TelnetConsole']
2022-01-28 23:51:00 [scrapy.middleware] INFO: Enabled downloader middlewares:
['scrapy.downloadermiddlewares.robotstxt.RobotsTxtMiddleware',
 'scrapy.downloadermiddlewares.httpauth.HttpAuthMiddleware',
 'scrapy.downloadermiddlewares.downloadtimeout.DownloadTimeoutMiddleware',
 'scrapy.downloadermiddlewares.defaultheaders.DefaultHeadersMiddleware',
 'scrapy.downloadermiddlewares.useragent.UserAgentMiddleware',
 'scrapy.downloadermiddlewares.retry.RetryMiddleware',
 'scrapy.downloadermiddlewares.redirect.MetaRefreshMiddleware',
 'scrapy.downloadermiddlewares.httpcompression.HttpCompressionMiddleware',
 'scrapy.downloadermiddlewares.redirect.RedirectMiddleware',
 'scrapy.downloadermiddlewares.cookies.CookiesMiddleware',
 'scrapy.downloadermiddlewares.httpproxy.HttpProxyMiddleware',
 'scrapy.downloadermiddlewares.stats.DownloaderStats']
2022-01-28 23:51:00 [scrapy.middleware] INFO: Enabled spider middlewares:
['scrapy.spidermiddlewares.httperror.HttpErrorMiddleware',
 'scrapy.spidermiddlewares.offsite.OffsiteMiddleware',
 'scrapy.spidermiddlewares.referrer.RefererMiddleware',
 'scrapy.spidermiddlewares.urllength.UrlLengthMiddleware',
 'scrapy.spidermiddlewares.depth.DepthMiddleware']
2022-01-28 23:51:00 [scrapy.middleware] INFO: Enabled item pipelines:
[]
2022-01-28 23:51:00 [scrapy.extensions.telnet] INFO: Telnet console listening on 127.0.0.1:6023
2022-01-28 23:51:00 [scrapy.core.engine] INFO: Spider opened
2022-01-28 23:51:00 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://shekouwoman.com/robots.txt> (referer: None)
```

The `'response.css'` command is used to along with the `extract()` to extract the tag we want and the relevant information.



```
Anaconda Powershell Prompt (anaconda3)
In [1]: response.css('title::text').extract()
Out[1]:
['New Arrivals | Shekou Woman',
 'American Express',
 'Apple Pay',
 'Google Pay',
 'JCB',
 'Mastercard',
 'Shop Pay',
 'Visa']
```

```
Anaconda Powershell Prompt (anaconda3)
In [22]: product.css('p.product--title::text').getall()
Out[22]:
['\n      PRE-ORDER: On Your Mark Sweatpants- Mens\n      Jeans Bundle (4 ITEMS)\n      My Baby Blue Sweatpants\n      Bikini Bundle (8 Separates)\n      Game Player Jeans\n      Kyoto Jeans\n      My Baby Blue Hoodie\n      Tropicana Motel Sweatpants\n      Tropicana Motel Sweatshirt\n      PRE-ORDER: In the Woodlands Sweatpants\n      In the Woodlands Hoodie\n      Jeans Bundle (4 ITEMS)\n      PRE-ORDER: Fleece Lined Hoodie\n      Sweatpants + Hoodie Bundle - Grey (2 ITEMS)\n      I'm Not In The Mood Bundle (4 ITEMS)\n      Hoodie & Sweatpants Bundle (2 ITEMS)\n      Sweatpants Mega Bundle (9 ITEMS)\n      Soft Girl Bundle (4 ITEMS)\n      Beachy Bundle (5 Items)\n      Indie Bundle (5 ITEMS)\n      Astrology Bundle (4 ITEMS)\n      Unlock Your Dreams Phone Chain\n      Running On Time Necklace\n      90's Bundle (4 ITEMS)\n      Centre Of Attention Necklace\n      Key West Kitten Bundle (5 ITEMS)\n      Danish Pastel Bundle (5 ITEMS)\n      Sweatpants Bundle (3 ITEMS)\n      PRE-ORDER: Sweatpants Bundle (3 ITEMS)\n      PRE-ORDER: Sweatpants Bundle (3 ITEMS)\n      Sweatpants Bundle (3 ITEMS)\n      PRE-ORDER: Angel Aura Halter Top\n      PRE-ORDER: Midsummer Night's Dream Dress\n      PRE-ORDER: Don't @ Me Hoodie\n      PRE-ORDER: There's No Better Time Hoodie\n      PRE-ORDER: Collect Moments Hoodie\n      PRE-ORDER: Take What You Need Hoodie\n']
In [23]:
```

```
Anaconda Powershell Prompt (anaconda3)
In [21]: product.css('span.product--price.money::text').getall()
Out[21]:
['\n      Rs. 1,600.00\n      Rs. 5,500.00\n      Rs. 1,600.00\n      Rs. 5,000.00\n      Rs. 1,400.00\n      Rs. 1,400.00\n      Rs. 2,700.00\n      Rs. 1,600.00\n      Rs. 2,700.00\n      Rs. 1,600.00\n      Rs. 2,700.00\n      Rs. 5,400.00\n      Rs. 2,700.00\n      Rs. 3,900.00\n      Rs. 4,900.00\n      Rs. 3,900.00\n      Rs. 13,100.00\n      Rs. 3,100.00\n      Rs. 3,500.00\n      Rs. 7,300.00\n      Rs. 5,400.00\n      Rs. 200.00\n      Rs. 300.00\n      Rs. 3,500.00\n      Rs. 300.00\n      Rs. 3,100.00\n      Rs. 5,400.00\n      Rs. 4,300.00\n      Rs. 4,300.00\n      Rs. 4,300.00\n      Rs. 4,300.00\n      Rs. 1,400.00\n      Rs. 1,400.00\n      Rs. 2,700.00\n      Rs. 2,700.00\n      Rs. 2,700.00\n      Rs. 2,700.00\n      Rs. 1,600.00\n']
```



```
Anaconda Powershell Prompt (anaconda3)
In [26]: products.css('a').attrib['href']
Out[26]: '/collections/home/products/pre-order-on-your-mark-sweatpants-mens'

In [27]:
```

First, we need to create a .py file in the spider's folder and then this folder can be used to scrape the website.

```
Shopcloth.py ×
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Jan 29 01:19:35 2022
4
5  @author: aditi
6  """
7
8  import scrapy
9
10 class ShopclothSpider(scrapy.Spider):
11     name='Clothes'
12     start_urls=['https://shekouwoman.com/collections/home']
13
14     def parse(self, response):
15         title=response.css('title::text').extract()
16         yield{'titleset':title}
17         for products in response.css('div.product--details-container'):
18             yield {
19                 'name': products.css('p::text').get(),
20                 'real_price':products.css('span.product--compare-price.money::text').get(),
21                 'discount_price':products.css('span.product--price.money::text').get()
22             }
23
```

```
1 [{"titleset": ["New Arrivals | Shekou Woman", "American Express", "Apple Pay", "Google Pay", "JCB", "Mastercard", "Sho  
2 {"name": "\n PRE-ORDER: On Your Mark Sweatpants- Mens\n ", "real_price": "\n Rs. 3,100.0  
3 {"name": "\n Jeans Bundle (4 ITEMS)\n ", "real_price": "\n Rs. 10,700.00\n ", "di  
4 {"name": "\n My Baby Blue Sweatpants\n ", "real_price": null, "discount_price": "\n Rs. 1  
5 {"name": "\n Bikini Bundle (8 Separates)\n ", "real_price": "\n Rs. 10,100.00\n "  
6 {"name": "\n Game Player Jeans\n ", "real_price": "\n Rs. 2,700.00\n ", "discount  
7 {"name": "\n Kyoto Jeans\n ", "real_price": "\n Rs. 2,700.00\n ", "discount price  
8 {"name": "\n My Baby Blue Hoodie\n ", "real_price": "\n Rs. 5,400.00\n ", "discou  
9 {"name": "\n Tropicana Motel Sweatpants\n ", "real_price": "\n Rs. 3,100.00\n ",  
10 {"name": "\n Tropicana Model Sweatshirt\n ", "real_price": "\n Rs. 5,400.00\n ",  
11 {"name": "\n PRE-ORDER: In the Woodlands Sweatpants\n ", "real_price": "\n Rs. 3,100.00  
12 {"name": "\n In the Woodlands Hoodie\n ", "real_price": "\n Rs. 5,400.00\n ", "di  
13 {"name": "\n Jeans Bundle (4 ITEMS)\n ", "real_price": "\n Rs. 10,700.00\n ", "di  
14 {"name": "\n PRE-ORDER: Fleece Lined Hoodie\n ", "real_price": "\n Rs. 5,400.00\n "  
15 {"name": "\n Sweatpants + Hoodie Bundle - Grey (2 ITEMS)\n ", "real_price": "\n Rs. 8,40  
16 {"name": "\n I'm Not In The Mood Bundle (4 ITEMS)\n ", "real_price": "\n Rs. 10,100.00\n "  
17 {"name": "\n Hoodie & Sweatpants Bundle (2 ITEMS)\n ", "real_price": "\n Rs. 8,400.00\n "  
18 {"name": "\n Sweatpants Mega Bundle (9 ITEMS)\n ", "real_price": "\n Rs. 24,100.00\n "  
19 {"name": "\n Soft Girl Bundle (4 ITEMS)\n ", "real_price": "\n Rs. 7,400.00\n ",  
20 {"name": "\n Beachy Bundle (5 Items)\n ", "real_price": "\n Rs. 7,400.00\n ", "di  
21 {"name": "\n Indie Bundle (5 ITEMS)\n ", "real_price": "\n Rs. 15,200.00\n ", "di  
22 {"name": "\n Astrology Bundle (4 ITEMS)\n ", "real_price": "\n Rs. 11,600.00\n "  
23 {"name": "\n Unlock Your Dreams Phone Chain\n ", "real_price": "\n Rs. 300.00\n "  
24 {"name": "\n Running On Time Necklace\n ", "real_price": null, "discount_price": "\n Rs. 3  
25 {"name": "\n 90's Bundle (4 ITEMS)\n ", "real_price": "\n Rs. 7,600.00\n ", "disc  
26 {"name": "\n Centre Of Attention Necklace\n ", "real_price": null, "discount_price": "\n R  
27 {"name": "\n Key West Kitten Bundle (5 ITEMS)\n ", "real_price": "\n Rs. 6,400.00\n "  
28 {"name": "\n Danish Pastel Bundle (5 ITEMS)\n ", "real_price": "\n Rs. 11,300.00\n "  
29 {"name": "\n Sweatpants Bundle (3 ITEMS)\n ", "real_price": "\n Rs. 9,200.00\n ",  
30 {"name": "\n PRE-ORDER: Sweatpants Bundle (3 ITEMS)\n ", "real_price": "\n Rs. 9,200.00  
31 {"name": "\n PRE-ORDER: Sweatpants Bundle (3 ITEMS)\n ", "real_price": "\n Rs. 9,200.00  
32 {"name": "\n Sweatpants Bundle (3 ITEMS)\n ", "real_price": "\n Rs. 9,200.00\n ",  
33 {"name": "\n PRE-ORDER: Annel Aura Halter Top\n ", "real_price": "\n Rs. 2,700.00\n "
```



```
Shopcloth.py × items.py × pipelines.py × settings.py ×
6 """
7
8 import scrapy
9
10 from ..items import ClothshoppingItem
11
12 class ShopclothSpider(scrapy.Spider):
13
14     name='Clothes'
15     start_urls=['https://shekouwoman.com/collections/home']
16
17     def parse(self, response):
18         #title=response.css('title::text').extract()
19         #yield{'titleset':title}
20         item=ClothshoppingItem()
21
22         for products in response.css('div.product--details-container'):
23
24             item['name']=products.css('p::text').get(),
25             item['real_price']=products.css('span.product--compare-price.money::text').get(),
26             item['discount_price']=products.css('span.product--price.money::text').get()
27
28         yield item
29
```

Using Pipelines-

We use pipeline to push the scraped item into the output file/database or to perform any pre-processing steps by making use of middleware.

In the settings.py file, uncomment the code for “Item_Pipeline”-

```
Shopcloth.py × items.py × pipelines.py × settings.py ×
56
57 # Enable or disable extensions
58 # See https://docs.scrapy.org/en/latest/topics/extensions.html
59 #EXTENSIONS = {
60 #     'scrapy.extensions.telnet.TelnetConsole': None,
61 #}
62
63 # Configure item pipelines
64 # See https://docs.scrapy.org/en/latest/topics/item-pipeline.html
65 ITEM_PIPELINES = {
66     'clothshopping.pipelines.ClothshoppingPipeline': 300,
67 }
```

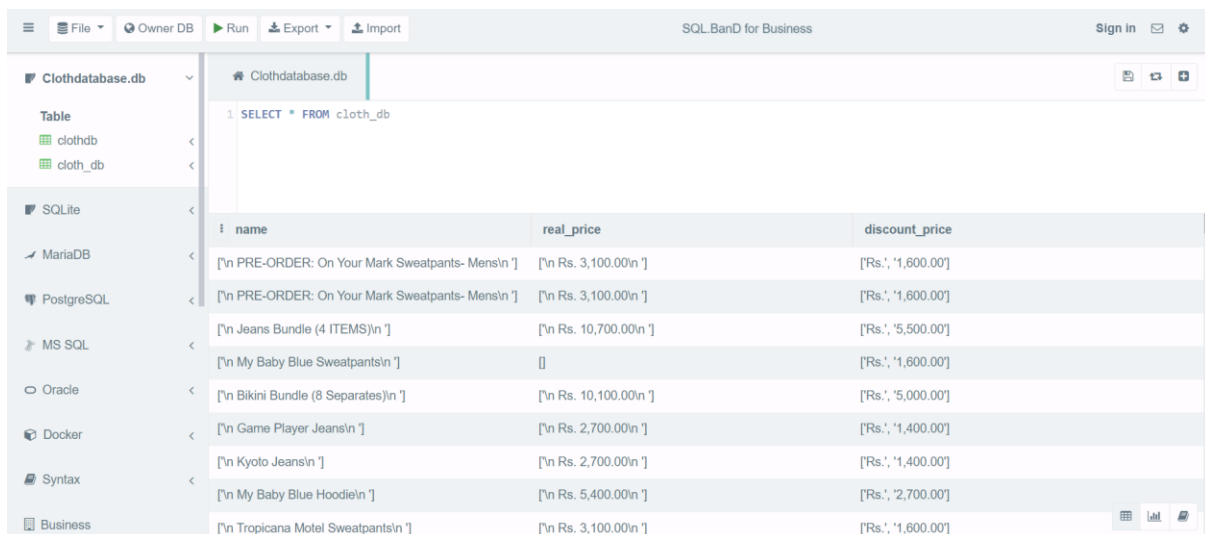
In the pipelines.py file, we define classes to create the database –


```

Shopcloth.py × items.py × pipelines.py × settings.py ×
4 # See: https://docs.scrapy.org/en/latest/topics/item-pipeline.html
5
6
7 # useful for handling different item types with a single interface
8 import sqlite3
9 from itemadapter import ItemAdapter
10
11
12 class ClothshoppingPipeline:
13     def __init__(self):
14         self.create_connection()
15         self.create_table()
16
17     def create_connection(self):
18         self.con = sqlite3.connect('Clothdatabase.db')
19         self.cur = self.con.cursor()
20
21     def create_table(self):
22         self.cur.execute(""" DROP TABLE IF EXISTS Clothdb""")
23         self.cur.execute(""" CREATE TABLE Clothdb(
24             name TEXT
25             real_price REAL
26             discount_price REAL)""")
27
28     def process_item(self, item, spider):
29         self.store_db(item)
30         return item
31
32     def store_db(self, item):
33         self.cur.execute("""INSERT INTO Clothdb VALUES(?, ?, ?)""",
34             (item['name'], item['real_price'], item['discount_price']))
35         self.con.commit()
36         return item

```

Using sqliteonline.com to run and display the database



| name | real_price | discount_price |
|--|-----------------------|---------------------|
| ["n PRE-ORDER: On Your Mark Sweatpants- Mens\n] | ["n Rs. 3,100.00\n] | ["Rs.", "1,600.00"] |
| ["n PRE-ORDER: On Your Mark Sweatpants- Mens\n] | ["n Rs. 3,100.00\n] | ["Rs.", "1,600.00"] |
| ["n Jeans Bundle (4 ITEMS)\n] | ["n Rs. 10,700.00\n] | ["Rs.", "5,500.00"] |
| ["n My Baby Blue Sweatpants\n] | [] | ["Rs.", "1,600.00"] |
| ["n Bikini Bundle (8 Separates)\n] | ["n Rs. 10,100.00\n] | ["Rs.", "5,000.00"] |
| ["n Game Player Jeans\n] | ["n Rs. 2,700.00\n] | ["Rs.", "1,400.00"] |
| ["n Kyoto Jeans\n] | ["n Rs. 2,700.00\n] | ["Rs.", "1,400.00"] |
| ["n My Baby Blue Hoodie\n] | ["n Rs. 5,400.00\n] | ["Rs.", "2,700.00"] |
| ["n Tropicana Motel Sweatpants\n] | ["n Rs. 3,100.00\n] | ["Rs.", "1,600.00"] |

The database cloth_db is created with the columns -name, real_price and discount_price. when we use pipelines to scrape the data, [Sqliteonline.com](https://sqliteonline.com) can then be used to view the scraped data.

To create an offline database sqlitebrowser.com is used and the database can be downloaded directly from the website.