

MÔ TẢ QUÁ TRÌNH CÀI ĐẶT VÀ CHẠY CODE

I. Cài đặt

1. Hệ điều hành sử dụng: Window 8 64-bit
2. Editor sử dụng: Sublime Text 3
3. Cài đặt và thiết lập Virtualenvironment:

- Cài đặt thư viện virtualenv

Sử dụng virtualenv để tạo 1 môi trường isolated Python và có thể cài đặt các packages vào trong môi trường này. Tham khảo cách cài đặt tại đây: <http://python-guide-pt-br.readthedocs.io/en/latest/dev/virtualenvs/>

Sau khi cài đặt

```
C:\Users\duyen>virtualenv --version
15.1.0
```

- Tạo môi trường ảo trong thư mục **env** bên trong thư mục **Code**

```
$ cd Code
```

```
$ virtualenv env
```

```
C:\Users\duyen\Desktop\Code>virtualenv env
New python executable in C:\Users\duyen\Desktop\Code\env\Scripts\python.exe
Installing setuptools, pip, wheel...done.
C:\Users\duyen\Desktop\Code>
```

- Cài đặt packages

```
$ pip install -r requirements.txt
```

```
C:\Users\duyen\Desktop\Code>pip install -r requirements.txt
Collecting amqp==1.4.9 (from -r requirements.txt (line 1))
  Using cached amqp-1.4.9-py2.py3-none-any.whl
Collecting anyjson==0.3.3 (from -r requirements.txt (line 2))
Collecting Babel==2.4.0 (from -r requirements.txt (line 3))
  Using cached Babel-2.4.0-py2.py3-none-any.whl
Collecting backports.ssl-match-hostname==3.5.0.1 (from -r requirements.txt (line 4))
Requirement already satisfied: beautifulsoup4==4.6.0 in c:\users\duyen\anaconda2\lib\site-packages (from -r requirements.txt (line 5))
Collecting billiard==3.3.0.23 (from -r requirements.txt (line 6))
  Using cached billiard-3.3.0.23-cp27-none-win_amd64.whl
Requirement already satisfied: bs4==0.0.1 in c:\users\duyen\anaconda2\lib\site-packages (from -r requirements.txt (line 7))
Collecting celery==3.1.25 (from -r requirements.txt (line 8))
  Using cached celery-3.1.25-py2.py3-none-any.whl
Collecting certifi==2017.4.17 (from -r requirements.txt (line 9))
  Using cached certifi-2017.4.17-py2.py3-none-any.whl
Requirement already satisfied: cyclr==0.10.0 in c:\users\duyen\anaconda2\lib\site-packages (from -r requirements.txt (line 10))
Requirement already satisfied: enum-compat==0.0.2 in c:\users\duyen\anaconda2\lib\site-packages (from -r requirements.txt (line 11))
Requirement already satisfied: enum34==1.1.6 in c:\users\duyen\anaconda2\lib\site-packages (from -r requirements.txt (line 12))
Requirement already satisfied: eventlet==0.21.0 in c:\users\duyen\anaconda2\lib\
```

- Cài đặt redis để lưu trữ dữ liệu trên ram (kết quả của quá trình crawl dữ liệu từ Amazon)

Vào trang <https://github.com/dmajkic/redis/downloads> và tải file **redis-2.4.5-win32-win64.zip**

Tạo thư mục redis trong thư mục Code đặt file để chứa file redis đã tải. Giải nén file **redis-2.4.5-win32-win64.zip**

Code > redis

Name	Date modified	Type	Size
32bit	10/17/2011 2:42 PM	File folder	
64bit	10/17/2011 2:33 PM	File folder	
00-RELEASENOTES	12/28/2011 5:32 PM	File	11 KB
BUGS	9/27/2011 9:02 PM	File	1 KB
CONTRIBUTING	9/27/2011 9:02 PM	File	1 KB
COPYING	9/27/2011 9:02 PM	File	2 KB
README	10/28/2011 2:32 PM	File	5 KB
README-Windows.txt	2/22/2011 11:14 PM	Text Document	3 KB
redis-2.4.5-win32-win64.zip	7/8/2017 7:21 PM	WinRAR ZIP archive	602 KB

Do hệ điều hành em sử dụng là win 64bit nên em chỉ cần giữ thư mục 64bit còn lại có thể xóa cho gọn

Code > redis >

Name	Date modified	Type	Size
64bit	10/17/2011 2:33 PM	File folder	

- Cài đặt RabbitMQ để phân tán task tới các celery worker

Vào trang <https://www.rabbitmq.com/install-windows.html> và tải file **rabbitmq-server-3.6.10.exe**

his PC > Downloads > Programs

Name	Date modified	Type	Size
ZoomInstaller.exe	12/12/2016 11:47 ...	Application	9,040 KB
ViberSetup.exe	12/25/2016 11:05 ...	Application	66,506 KB
vcredist_x64_2.exe	1/11/2017 4:45 PM	Application	7,026 KB
SkypeSetupFull.exe	12/12/2016 8:01 PM	Application	42,851 KB
rabbitmq-server-3.6.10.exe	6/11/2017 10:48 AM	Application	4,688 KB

Cài đặt bằng cách click file exe và chạy

II. Thực hiện chạy code

Thư mục Code gồm:

- Thư mục source chứa source code.
- Thư mục DATA chứa các thư mục dữ liệu thu thập tại các thời gian khác nhau và chưa được tiền xử lý.
- File DATA.csv là file chứa toàn bộ dữ liệu sau khi được tổng hợp và tiền dữ liệu được sử dụng để thực hiện huấn luyện.
- File requirements.txt chứa tên các package cần thiết khi chạy.

PC > Desktop > Code

Name	Date modified	Type	Size
DATA	7/8/2017 9:02 PM	File folder	
source	7/8/2017 9:07 PM	File folder	
DATA.csv	6/30/2017 11:45 PM	Microsoft Excel C...	1,789 KB
requirements.txt	7/8/2017 8:28 PM	Text Document	1 KB

1. Thu thập dữ liệu

- Start redis-server

Vào thư mục **redis/64bit**, click file **redis-server.exe** để start redis-server

PC > Desktop > Code > redis > 64bit

Name	Date modified	Type	Size
libhiredis.dll	12/28/2011 8:51 PM	Application extens...	65 KB
redis.conf	11/6/2011 10:43 PM	CONF File	21 KB
redis-benchmark.exe	12/28/2011 8:51 PM	Application	78 KB
redis-check-aof.exe	12/28/2011 8:51 PM	Application	44 KB
redis-check-dump.exe	12/28/2011 8:51 PM	Application	78 KB
redis-cli.exe	12/28/2011 8:51 PM	Application	94 KB
redis-server.exe	12/28/2011 8:51 PM	Application	278 KB

Start thành công, giữ nguyên cmd này mà mở cmd khác để start celery worker

```

C:\Users\duyen\Desktop\Code\redis\64bit\redis-server.exe
[7636] 08 Jul 20:01:04 # Warning: no config file specified, using the default co
nfig. In order to specify a config file use 'redis-server /path/to/redis.conf'
[7636] 08 Jul 20:01:04 * Server started. Redis version 2.4.5
[7636] 08 Jul 20:01:04 # Open data file dump.rdb: No such file or directory
[7636] 08 Jul 20:01:04 * The server is now ready to accept connections on port 6
379
[7636] 08 Jul 20:01:05 - 0 clients connected (0 slaves), 1179896 bytes in use
[7636] 08 Jul 20:01:10 - 0 clients connected (0 slaves), 1179896 bytes in use

```

- Vào thư mục **Code** sau đó activate môi trường ảo

```
$ cd Code
```

```
$ env\Scripts\activate
```

```
C:\Users\duyen\Desktop\Code>env\Scripts\activate
(env) C:\Users\duyen\Desktop\Code>_
```

- Start celery worker với task **Crawler_Worker** và **concurrency = 15**

```
(env) C:\Users\duyen\Desktop\Code>cd source
(env) C:\Users\duyen\Desktop\Code\source>celery -A Crawler_Worker worker -c 15
```

Start thành công

```
(env) C:\Users\duyen\Desktop\Code\source>celery -A Crawler_Worker worker -c 15
[2017-07-08 20:09:44,115: WARNING/MainProcess] c:\users\duyen\anaconda2\lib\site
-packages\celery\apps\worker.py:161: CDeprecationWarning:
Starting from version 3.2 Celery will refuse to accept pickle by default.

The pickle serializer is a security concern as it may give attackers
the ability to execute any command. It's important to secure
your broker from unauthorized access when using pickle, so we think
that enabling pickle should require a deliberate action and not be
the default choice.

If you depend on pickle then you should set a setting to disable this
warning and to be sure that everything will continue working
when you upgrade to Celery 3.2::

    CELERY_ACCEPT_CONTENT = ['pickle', 'json', 'msgpack', 'yaml']

You must only enable the serializers that you will actually use.

warnings.warn(CDeprecationWarning(W_PICKLE_DEPRECATED))

----- celery@DUYEN v3.1.25 (Cipater)
-----
* *** * -- Windows-8.1-6.3.9600
* - ****
** ----- [config]
** .> app: Crawler_Worker:0x1daa630
** .> transport: amqp://guest:**@localhost:5672//
** .> results: redis://localhost:6379/0
** *** -- * .> concurrency: 15 (prefork)
** -----
** ***** [queues]
** ----- .> celery exchange=celery<direct> key=celery

[2017-07-08 20:10:15,334: WARNING/MainProcess] celery@DUYEN ready.
```

- Mở cmd khác để chạy code thu thập dữ liệu

```
C:\Users\duyen\Desktop\Code>env\Scripts\activate
(env) C:\Users\duyen\Desktop\Code>cd source
(env) C:\Users\duyen\Desktop\Code\source>python Crawler_Worker.py
```

Sau đó quan sát bên cmd celery sẽ thấy bắt đầu crawl dữ liệu

```
[2017-07-08 20:33:21,744: WARNING/MainProcess] celery@DUYEN ready.
[2017-07-08 20:33:36,924: WARNING/Worker-1] https://www.amazon.com/dp/B00B81XR1Y
[2017-07-08 20:33:37,128: WARNING/Worker-1] https://www.amazon.com/ask/questions
/asin/B00B81XR1Y/1/ref=ask_q1_psf_q1_hza
[2017-07-08 20:33:38,305: WARNING/Worker-1] https://www.amazon.com/dp/B0
0J5104FC
[2017-07-08 20:33:38,374: WARNING/Worker-1] https://www.amazon.com/dp/B01FLO5914
[2017-07-08 20:33:38,358: WARNING/Worker-1] https://www.amazon.com/dp/B011QHFB4A
[2017-07-08 20:33:38,372: WARNING/Worker-1] https://www.amazon.com/dp/B010FHR44G
[2017-07-08 20:33:38,375: WARNING/Worker-1] https://www.amazon.com/dp/B00P07GKLM
[2017-07-08 20:33:38,394: WARNING/Worker-1] https://www.amazon.com/dp/B01M0PB8DZ
[2017-07-08 20:33:38,344: WARNING/Worker-1] https://www.amazon.com/dp/B004X8JUT
6
```

Hoàn tất thu thập dữ liệu tại 1 thời điểm

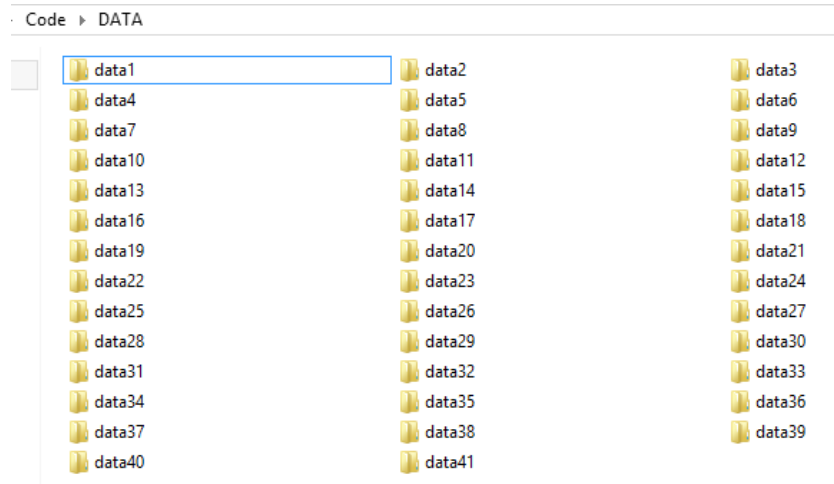
```
C:\Windows\system32\cmd.exe - celery -A Crawler_Worker worker -c 15
[2017-07-08 20:55:35,302: WARNING/Worker-1] [99, 'B00CTQ0EUY', 3.8, 789, 52, 7,
0, 5.99, 0, 11]
[2017-07-08 20:55:35,693: WARNING/Worker-1] https://www.amazon.com/ask/questions
/asin/B0009EILKS/3/ref=ask_q1_psf_q1_hza
[2017-07-08 20:55:38,799: WARNING/Worker-1] [100, 'B0009EILKS', 4.0, 1354, 16, 6
3, 7.49, 17, 11]
[2017-07-08 20:55:39,309: WARNING/Worker-1] Test [194, 'B01772051W', 4.5, 736, 6
2, 8, 0, 11.99, 33, 11, [83, 'B01M1H1IUT', 4.6, 650, 76, 6, 0, 19.99, 50, 11, [9
2, 'B0009D701U', 4.1, 401, 12, 3, 1, 18.95, 5, 91, [88, 'B01FLO5914', 4.3, 1565,
84, 21, 0, 18.95, 5, 11, [91, 'B00D9NU2D4', 4.7, 6273, 246, 56, 0, 15.25, 0, 11
[90, 'B00J5KDC02', 4.5, 805, 46, 14, 0, 29.0, 0, 11, [86, 'B00YJISUI8', 4.1, 6
46, 67, 6, 0, 18.95, 0, 11, [95, 'B00NR1YQK4', 4.3, 1710, 126, 49, 0, 17.09, 5,
131, [93, 'B071WL889Q', 4.6, 457, 0, 0, 0, 14.99, 0, 11, [89, 'B00LU6VDG2', 4.5,
2174, 71, 14, 0, 16.99, 59, 21, [98, 'B011YRJ310', 3.0, 2087, 68, 43, 0, 9.3, 0
511, [81, 'B0017209NO', 4.5, 2272, 86, 9, 0, 14.0, 0, 11, [82, 'B01AYC74DK', 4
.8, 1016, 49, 8, 0, 11.99, 5, 11, [85, 'B007SAC6XE', 4.2, 506, 42, 17, 0, 12.15,
5, 391, [87, 'B00S2G7AUC', 4.3, 696, 64, 20, 0, 10.75, 82, 11, [84, 'B010SS5UXXC
', 4.7, 1564, 44, 11, 0, 13.49, 0, 11, [97, 'B00120UWJ0', 4.8, 2037, 119, 100, 0
0, 0, 141, [96, 'B004D2C4Q4', 4.1, 985, 49, 17, 0, 16.98, 15, 11, [99, 'B00CTQ
0EUY', 3.8, 789, 52, 7, 0, 5.99, 0, 11, [100, 'B0009EILKS', 4.0, 1354, 16, 6, 3,
7.49, 17, 11]
[2017-07-08 20:55:39,311: WARNING/Worker-1] Page number: 5 of cate Skin-Care-Pro
ducts
[2017-07-08 20:55:39,315: WARNING/Worker-1] Category -----Skin-Care-Products--
-----
is done
```

Sau khi hoàn tất crawl trong thư mục **Code** đã được tạo thư mục **DATA** chứa thư mục con data<time> chứa dữ liệu đã crawl

This PC > Desktop > Code > DATA > data2017_07_08 20_33_33

Name	Date modified	Type	Size
Bath-Bathing-Accessories.csv	7/8/2017 8:55 PM	Microsoft Excel C...	5 KB
Gift-Sets.csv	7/8/2017 8:55 PM	Microsoft Excel C...	5 KB
Hair-Care-Products.csv	7/8/2017 8:55 PM	Microsoft Excel C...	5 KB
Makeup.csv	7/8/2017 8:55 PM	Microsoft Excel C...	5 KB
Perfumes-Fragrances.csv	7/8/2017 8:55 PM	Microsoft Excel C...	5 KB
Skin-Care-Products.csv	7/8/2017 8:55 PM	Microsoft Excel C...	5 KB
Tools-Accessories.csv	7/8/2017 8:55 PM	Microsoft Excel C...	5 KB

Em đã đánh số các thư mục 1-41 từ thứ tự thu thập dữ liệu theo thời gian thu thập trong thư mục **DATA**



2. Tổng hợp dữ liệu và tiền xử lý dữ liệu

- Chạy file **Pre_Data.py** để thực hiện quá trình tổng hợp dữ liệu và tiền xử lý.

```
(env) C:\Users\duyen\Desktop\Code\source>python Pre_Data.py
***** count_rank_freq Folder 1*****
*****Folder 2*****
***** count_rank_freq Folder 2*****
```

Sau khi hoàn tất file tổng hợp DATA.csv sẽ được tạo.

3. Huấn luyện

- Chạy file **Train_Data.py** để thực hiện huấn luyện

```
(env) C:\Users\duyen\Desktop\Code\source>python Train_Data.py
Epoch: 0, training err: 73.590476, val err: 74.028571, patience: 20000
Epoch: 100, training err: 35.314286, val err: 36.228571, patience: 19999
```

Sau khi hoàn tất training thư mục **RESULT** sẽ được tạo chứa cái file kết quả độ lỗi tập test (*_TEST.txt), độ lỗi tập train và val (*_Train.txt), ghi lại tất cả độ lỗi trên tập train và val trong quá trình train (*_ERR.csv) và xuất ra file đồ thị biểu diễn sự biến thiên của train error và val error (*_Plot.png). Tên file đặt dựa vào các tham số huấn luyện

Name	Date modified	Type	Size
5_0.1_20000_ERR.csv	7/7/2017 12:39 PM	Microsoft Excel C...	1,837 KB
5_0.1_20000_Plot.png	7/7/2017 12:39 PM	PNG image	46 KB
5_0.1_20000_Test.txt	7/8/2017 4:36 PM	Text Document	1 KB
5_0.1_20000_Train.txt	7/7/2017 12:39 PM	Text Document	1 KB
10_0.1_20000_ERR.csv	7/7/2017 1:10 PM	Microsoft Excel C...	1,694 KB
10_0.1_20000_Plot.png	7/7/2017 1:10 PM	PNG image	47 KB
10_0.1_20000_Test.txt	7/7/2017 1:10 PM	Text Document	1 KB
10_0.1_20000_Train.txt	7/7/2017 1:10 PM	Text Document	1 KB