

PROJECT 1

BÁO CÁO ĐỒ ÁN VỀ CRAWLER DỮ LIỆU

I/ THÔNG TIN

Author	version	Description
20424030 - Đặng Trung Hiếu	1.0	Tham gia viết script crawl của trang Foody
20424025 - Nguyễn Thị Thúy Hằng	1.0	Tham gia viết script crawl của trang Google review
20424051 - Nguyễn Thành Long	1.0	Xây dựng core,viết script crawl của trang ShopeeFood, Login foody và hỗ trợ các plugin khác

II/ KẾ HOẠCH PHÂN CÔNG

Tuần	20424030	20424025	20424051
Tuần 1	Họp về vấn đề chọn chủ đề và các thông tin cần crawl, chọn trang web để crawl. Chọn ngôn ngữ lập trình và thống nhất cấu trúc Project	Họp về vấn đề chọn chủ đề và các thông tin cần crawl, chọn trang web để crawl. Chọn ngôn ngữ lập trình và thống nhất cấu trúc Project	Họp về vấn đề chọn chủ đề và các thông tin cần crawl, chọn trang web để crawl. Chọn ngôn ngữ lập trình và thống nhất cấu trúc Project Tạo source code trên github
Tuần 2	Bắt đầu tiến hành nghiên cứu và code để crawl dữ liệu ở trang Foody	Bắt đầu tiến hành nghiên cứu và code để crawl dữ liệu ở trang Google Review	Bắt đầu tiến hành nghiên cứu và code để crawl dữ liệu ở trang Shopee Food và kiểm tra tiến độ thành viên để support nếu cần
Tuần 3	Họp thống nhất kế hoạch và nêu các khó khăn với nhau để cùng	Họp thống nhất kế hoạch và nêu các khó khăn với nhau để	Họp thống nhất kế hoạch và nêu các khó khăn với nhau

	nhau nghĩ cách giải quyết. Sau khi thống nhất tiếp tục code xử lý những khó khăn gặp phải	cùng nhau nghĩ cách giải quyết. Sau khi thống nhất tiếp tục code xử lý những khó khăn gặp phải	để cùng nhau nghĩ cách giải quyết. Sau khi thống nhất tiếp tục code xử lý những khó khăn gặp phải
Tuần 4	Chạy kiểm kết quả lần cuối và bắt đầu hoàn thành báo cáo	Chạy kiểm kết quả lần cuối và hỗ trợ viết và hoàn thành báo cáo	Tổng hợp 3 bản crawl vào Project, đảm bảo project hoạt động ổn nhất.
Tuần 5	Hoàn thiện báo cáo và đóng gói	Hoàn thiện báo cáo và đóng gói	Thực hiện crawl data demo và hoàn thiện báo cáo và đóng gói.

III/ ĐÁNH GIÁ BÁO CÁO

STT	Tiêu chí	Giải quyết vấn đề	Tỉ lệ hoàn thành
1	Trình bày chủ đề, lý do chọn chủ đề, trang web lấy hoặc trang web là hạt giống.	<ul style="list-style-type: none"> - Thống nhất chọn chủ đề “ Ăn uống ” và các trang web dùng để crawl là: Foody, Shopee Food, Google Review. - Sử dụng chủ đề ăn uống vì sau khi nghiên cứu đây là sở thích của các thành viên và là chủ đề được quen thuộc của người dùng, các trang web về ăn uống có rất nhiều và đa dạng, những trang web lớn và có tính ổn định, dữ liệu để crawl về để xử lý là nhiều và đa dạng. 	10%
2	Mô tả thuật toán, cấu trúc mã nguồn, các thành phần hệ thống	<p>Cấu trúc mã nguồn chia làm các phần chính:</p> <ol style="list-style-type: none"> 1. config: chứa các file cấu hình chung như môi trường, database, logger, launcher,... 2. database: chứa dữ liệu crawler 3. utils: chứa các hàm hỗ trợ: db, lock, metric, signal, tracert 4. model: mô tả mô hình lưu trữ <p><code>website_id_quan: str</code></p>	30%

		<pre> website: str url: str district: str rate: float active_time: str full_name: str = None phone: str = None rate_count: float = None favorite: float = None price_from: float = None price_to: float = None other_service: str = None </pre> <p>5. plugin: chứa thuật toán chính của từng website: Foody, GoogleReview và ShopeeFood</p> <p>Các thành phần trong hệ thống:</p> <ol style="list-style-type: none"> 1. Hệ thống menu lựa chọn 2. Hệ thống logger 3. Hệ thống đo lường <p>Thuật toán chính sẽ trình bày phía dưới.</p>	
3	Các vấn đề xảy ra đối với crawler và phương pháp xử lý	<p>- Vấn đề lấy trang: kích thước frontier tăng lên khá nhanh, nên sẽ giới hạn lại trước khi lấy tiếp; với tốc độ mạng chậm, thì sẽ xử lý bằng cách bỏ qua hoặc chờ 1 khoảng thời gian nếu không tải thành công thì bỏ qua</p> <p>- Vấn đề phân tích trang: Với những trang SPA, hoặc có sử dụng ajax, webpack sử dụng phát sinh mã id, class rất khó, => chưa có phương án xử lý => có thể bỏ qua nếu lỗi.</p> <p>- Trích xuất liên kết: khi sử dụng thư viện python thì đã xử lý các vấn đề url relative</p> <p>- Lưu trữ data: sử dụng 1 field định danh làm khóa chính lưu vào database SQLite</p> <p>- Vấn đề xử lý song song: áp dụng đa luồng trong python để giải quyết - áp dụng kĩ thuật lock để tránh tranh chấp xử lý</p> <p>- Trang Shopee Food phải khai thác api thông qua phân tích</p> <p>- Trang Foody phải Login mới có thể crawl được tất cả dữ liệu.</p> <p>- Trang Google Review phải sử dụng</p>	10%

		selenium để fetch và extract data	
4	Các tính năng phức tạp của crawler	<ul style="list-style-type: none"> - Giới hạn lượng thông tin cần lấy - Xử lý tốc độ mạng không ổn định - Thao tác qua database, độ ổn định cao - Áp dụng đa luồng để crawl, độc lập nhau - Cho phép lựa chọn chủ đề ĐỘNG - Cho phép lựa chọn website crawl cùng lúc - Sử dụng api đã được phân tích cho kết quả nhanh - Áp dụng kĩ thuật vào selenium để tăng hiệu suất crawl data (headless, network, cookie,...) - Có nhiều màn hình để theo dõi từng luồng, thống kê - Có chức năng thống kê tốc độ lấy dữ liệu 	10%
5	Đánh giá hiệu năng crawler	<ul style="list-style-type: none"> - Đa luồng nên khả năng khai thác dữ liệu tối đa - Sử dụng api được phân tích trước đó đem lại tốc độ cao - Khả năng chịu lỗi cao không gây crash ứng dụng - Dữ liệu thống kê lấy về có giá trị rất cao - Số lượng thông tin lấy về cao hơn những gì team mong đợi 	5%
6	Mô tả và đánh giá dữ liệu thu thập được.	<ul style="list-style-type: none"> - Tổng quan từ 3 trang Foody, Shopee Food, Google Review - Dữ liệu đề ra để lấy và lưu vào DB: (website_id, website, full_name, district, rate, active_time, price_from, price_to, favorite) - Các dữ liệu sau khi crawl và xử lý về, thu được hợp lệ so với đề ra, do cấu trúc trang web các dữ liệu lấy về đôi khi có vài trang bị lỗi và phải catch xử lý - Dữ liệu rất đa dạng 	5%
7	Tiền xử lý dữ liệu thu thập	<p>Format lại dữ liệu quận lowercase Format lại để lấy được url đầy đủ Format lại số Rate của từ 3 trang web theo 1 định dạng để lưu vào DB Format lại số phone của quán theo 1 định dạng thống nhất Format lại giờ hoạt động của quán</p> <p>Mẫu vd chưa xử lý tiền dữ liệu:</p> <ul style="list-style-type: none"> - website: Foody - website_url: /ho-chi-minh/sang-trang/foods/san-fu-lu - name: San Fu Lu - phone: (+84) 824110244 - time: 10h30 - 20h30 - rate: 3.7 	10%

		- district: quan-1 - time: 10h30 - 20h00 chưa cập nhật - price: 500.000vnd - 1.000.000vnd Mẫu vd sau khi xử lý tiền dữ liệu: - website: Foody - website_url: https://www.foody.com.vn/ho-chi-minh/sang-trang/foods/san-fu-lu - full_name: San Fu Lu - phone: 0824110244 - time: 10h30 - 20h30 - rate: 7.4 - district: 1 - time: 10h30 - 20h00 None - price_to: 500.000 - price_from: 1.000.000	
8	Báo cáo rõ ràng các mục đã thực hiện, có thể hiện mức độ hoàn thiện của từng công việc	- Báo cáo đã trình bày đủ các mục đề ra, trong đó có cả giải thích các bước đầy đủ và tỉ lệ hoàn thành	20%
Tổng			100%

IV/ GIẢI THUẬT:

1. Crawl theo chiều sâu: tải các trang từ cùng 1 mức trước khi chuyển sang trang kế tiếp
2. Trong quá trình crawl sẽ có thống kê và thể hiện tốc độ lấy dữ liệu.
3. Foody:
 - a. Login: người dùng sẽ sử dụng tài khoản của mình để thực hiện login, từ đây, chúng ta export ra cookie và chia sẻ cho lần chạy sau.
 - b. Thực hiện lấy chủ đề lớn: request lên trang template lấy ra data và cho người dùng chọn lựa
 - c. Thực hiện crawl data url từ trang số 1 ứng với chủ đề người dùng chọn cho vào frontier
 - d. Lập trên frontier
 - e. Sẽ có 2 trường hợp:
 - i. Trường hợp là link danh sách chi nhánh: ta thực hiện fetch data url chi nhánh và lấy tất cả link cửa hàng cho vào frontier

- ii. Trường hợp là link của hàng, ta thực hiện áp dụng request-html để lấy data và trích xuất lưu vào db
- f. Kiểm tra có tồn tại trang thứ 2 không, nếu có ta lặp lại bước C với trang thứ 2. Ngược lại ta kết thúc quá trình crawl và thông báo

4. ShopeeFood:

- a. Gọi api lấy ra danh sách quận và danh sách chủ đề lớn cho người dùng chọn lựa
- b. Thực hiện lấy tất cả id cửa hàng ứng với chủ đề người dùng đã chọn
- c. Chia nhỏ danh sách các cửa hàng thực hiện thành mỗi cụm 25 cửa hàng
- d. Lặp trên từng id cửa hàng
 - i. Truy vấn api lấy data từ id cửa hàng
 - ii. Thực hiện chuẩn hóa
- e. Kết thúc quá trình crawl

5. Google Review:

- a. Cho phép người dùng tự lựa chọn chủ đề và nhập vào cho chương trình
- b. Từ keyword người dùng, chương trình sẽ tiến hành thiết lập selenium và khai thác dữ liệu từ trang google review
- c. Kiểm tra trên trang có tồn tại cửa hàng nào không:
 - i. Nếu có thì thực hiện lấy data thông qua tương tác click và bóc tách dữ liệu
 - ii. Nếu không thì kết thúc quá trình crawl
- d. Nếu chưa kết thúc quá trình, ta kiểm tra xem có nút qua trang hay không?
 - i. Nếu có thì ta tương tác click vào và đợi dữ liệu load xong quay lại **bước c**
 - ii. Nếu không thì kết thúc quá trình crawl

VI/ MÃ NGUỒN CHƯƠNG TRÌNH:

Link drive: [Tại đây](#)

Link github (dự phòng): [Tại đây](#)

VII/ BIÊN DỊCH và CHẠY CHƯƠNG TRÌNH

- Project sử dụng python 3, hệ thống lưu trữ là SQLite
- Hướng dẫn chạy chương trình:
 - + Đọc hướng dẫn cài đặt file "README.md"

WebScience-TopicalCrawler

Thu Thập Dữ Liệu Từ Web

- Cài đặt các package:

```
python -m pip install -r requirements.txt
```

- Tải webdriver giải nén và đưa vào thư mục **webdriver**
 - Chrome: <https://sites.google.com/a/chromium.org/chromedriver/downloads>
 - Edge: <https://developer.microsoft.com/en-us/microsoft-edge/tools/webdriver/>
 - Firefox: <https://github.com/mozilla/geckodriver/releases>
 - Safari: <https://webkit.org/blog/6900/webdriver-support-in-safari-10/>

- File cấu hình môi trường:

```
/config/env.py
```

- python_launcher: cài đặt theo môi trường (python, python3)
- external_log: quyết định có mở console để theo dõi log hay không

- Chạy chương trình:

```
python main.py
```

+

- + Sử dụng chrome webdriver:

<https://sites.google.com/a/chromium.org/chromedriver/downloads>

- + Tiến hành chạy:

- + Thực hiện mở command line và cd tới thư mục Source

C:\Windows\system32\cmd.exe

Microsoft Windows [Version 10.0.19043.1165]

(c) Microsoft Corporation. All rights reserved.

(dev-env) C:\Users\thafnhlong\Desktop\HK3\WebScience-TopicalCrawler\Source>python main.py

- + Đánh lệnh khởi chạy

```
(dev-env) C:\Users\thafnhlong\Desktop\HK3\WebScience-TopicalCrawler\Source>python main.py
Project: TopicalCrawler
Web of Science
```

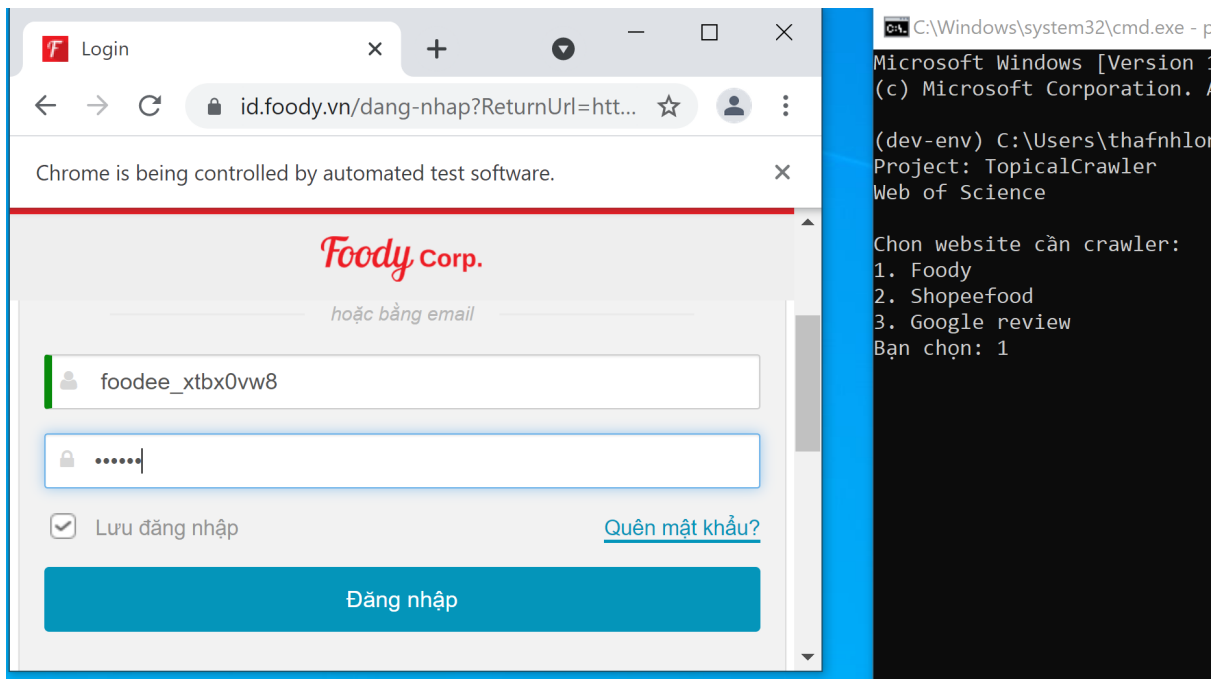
Chon website cần crawler:

1. Foody
2. Shopeefood
3. Google review

Bạn chọn:

- + Lựa chọn 1 plugin

Ví dụ: Foody, thì chương trình sẽ hiện lên cửa sổ, chúng ta sẽ điền tài khoản và mật khẩu



- + Sau khi bấm đăng nhập, chương trình sẽ tự kiểm tra xem đăng nhập thành công hay không và kết quả trả về :

```
C:\Windows\system32\cmd.exe - python main.py
(c) Microsoft Corporation. All rights reserved.

(dev-env) C:\Users\thafnhlong\Desktop\HK3\WebScience-TopicalCrawler\Source>python main.py
Project: TopicalCrawler
Web of Science

Chon website cần crawler:
1. Foody
2. Shopeefood
3. Google review
Bạn chọn: 1
Đang truy xuất lấy data...ok
Foody có những chủ đề ẩm thực hot như sau:
1. Sang trọng
2. Buffet
3. Nhà hàng
4. Ăn vặt/vía hè
5. Ăn chay
6. Café/Dessert
7. Quán ăn
8. Bar/Pub
9. Quán nhậu
10. Beer club
11. Tiệm bánh
12. Tiệc tận nơi
13. Shop Online
14. Giao cơm văn phòng
15. Khu Ẩm Thực
16. Quay lại
Bạn chọn: 12
```

- + Lựa chọn 1 chủ đề lớn: ví dụ 12: tiệc tận nơi


```
C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Monitor] Số trang xử lý: 0.00 / giây
[Monitor] Lượng thông tin xử lý: 0.0B / giây
[Monitor]
[Monitor] Số trang xử lý: 0.19 / giây
[Monitor] Lượng thông tin xử lý: 54.9KiB / giây
[Monitor]

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Foody 1] Lấy chủ đề: Tiệc tận nơi
[Foody 1] Xử lý các cửa hàng trên trang số 1...

C:\Windows\system32\cmd.exe - python main.py
7. Quán ăn
8. Bar/Pub
9. Quán nhậu
10. Beer club
11. Tiệm bánh
12. Tiệc tận nơi
13. Shop Online
14. Giao cơm văn phòng
15. Khu Ẩm Thực
16. Quay lại
Bạn chọn: 12
Đang tiến hành crawl data theo chủ đề: Tiệc tận nơi
Nhập bất kỳ để tiếp tục chọn website khác
Nhập exit để thoát chương trình
_
```

- + Sẽ có 2 cửa sổ hiện lên, 1 là monitor thống kê trên / giây; cửa sổ còn lại là thông tin của luồng xử lý crawl Foody (# Foody 1)
- + Trong lúc đó ta cũng có thể chọn lựa 1 chủ đề khác hoặc kể cả 1 website(plugin) khác
- + Bấm enter để tiếp, hoặc exit để thoát

```
C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Monitor] Lượng thông tin xử lý: 112.9KiB / giây
[Monitor]
[Monitor] Số trang xử lý: 0.53 / giây
[Monitor] Lượng thông tin xử lý: 109.7KiB / giây
[Monitor]
[Monitor] Số trang xử lý: 0.52 / giây
[Monitor] Lượng thông tin xử lý: 106.7KiB / giây
[Monitor]
[Monitor] Số trang xử lý: 0.50 / giây
[Monitor] Lượng thông tin xử lý: 103.8KiB / giây
[Monitor]
[Monitor] Số trang xử lý: 0.49 / giây
[Monitor] Lượng thông tin xử lý: 101.1KiB / giây
[Monitor]

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Foody 1] Lấy chủ đề: Tiệc tận nơi
[Foody 1] Xử lý các cửa hàng trên trang số 1...
[Foody 1] Thời gian xử lý: 12.47
[Foody 1] Xử lý các cửa hàng trên trang số 2...
[Foody 1] Thời gian xử lý: 6.67
[Foody 1] Xử lý các cửa hàng trên trang số 3...
[Foody 1] Thời gian xử lý: 13.38
[Foody 1] Xử lý các cửa hàng trên trang số 4...
[Foody 1] Thời gian xử lý: 5.21
[Foody 1] Xử lý các cửa hàng trên trang số 5...
[Foody 1] Thời gian xử lý: 12.50
[Foody 1] Xử lý các cửa hàng trên trang số 6...
[Foody 1] Thời gian xử lý: 5.36
[Foody 1] Xử lý các cửa hàng trên trang số 7...
[Foody 1] Thời gian xử lý: 5.63
[Foody 1] Xử lý các cửa hàng trên trang số 8...
[Foody 1] Thời gian xử lý: 2.62
[Foody 1] Xử lý các cửa hàng trên trang số 9...
[Foody 1] Thời gian xử lý: 0.04
[Foody 1] Đã crawler hết tất cả cửa hàng của chủ đề này.
Bấm enter để kết thúc

C:\Windows\system32\cmd.exe - python main.py
Bạn chọn: 12
Đang tiến hành crawl data theo chủ đề: Tiệc tận nơi
Nhập bất kỳ để tiếp tục chọn website khác
Nhập exit để thoát chương trình

Chon website cần crawler:
1. Foody
2. ShopeeFood
3. Google review
Bạn chọn: 2
Đang truy xuất lấy data...ok
ShopeeFood có những chủ đề hot như sau:
1. Đồ ăn-Đồ ăn
2. Đồ ăn-Đồ uống
3. Đồ ăn-Đồ chay
4. Đồ ăn-Bánh kem
5. Đồ ăn-Tráng miệng
6. Đồ ăn-Homemade
7. Đồ ăn-Via hè
8. Đồ ăn-Pizza/Burger
9. Đồ ăn-Món gà
10. Đồ ăn-Món lẩu
11. Đồ ăn-Sushi
12. Đồ ăn-Mì phở
13. Đồ ăn-Cơm hộp
14. Đồ ăn-Tất cả
15. Thực phẩm-Đồ chay
16. Thực phẩm-Trái cây
17. Thực phẩm-Thịt / Trứng
18. Thực phẩm-Thủy hải sản
```

- + Giả sử lựa chọn tiếp Plugin ShopeeFood, chọn 1 chủ đề

```
C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Monitor] Lượng thông tin xử lý: 77.0KiB / giây
[Monitor] Số trang xử lý: 0.43 / giây
[Monitor] Lượng thông tin xử lý: 75.5KiB / giây
[Monitor] Số trang xử lý: 0.42 / giây
[Monitor] Lượng thông tin xử lý: 74.1KiB / giây
[Monitor] Số trang xử lý: 0.41 / giây
[Monitor] Lượng thông tin xử lý: 72.7KiB / giây
[Monitor] Số trang xử lý: 0.41 / giây
[Monitor] Lượng thông tin xử lý: 71.4KiB / giây
[Monitor]

C:\Windows\system32\cmd.exe - python main.py
30. Siêu thị-Mỹ phẩm
31. Siêu thị-Mẹ & Bé
32. Siêu thị-Đồ chơi
33. Siêu thị-Quần áo / Giày dép
34. Siêu thị-Điện tử / Điện gia dụng
35. Siêu thị-Trang sức / Phụ kiện
36. Siêu thị-Tất cả
37. Thuốc-Nhà thuốc
38. Thuốc-Tất cả
39. Thú cưng-Thú cưng
40. Thú cưng-Tất cả
41. Quay lại
Bạn chọn: 1
Đang tiến hành crawl data theo chủ đề: Đồ ăn-Đồ ăn
Nhập bất kỳ để tiếp tục chọn website khác
Nhập exit để thoát chương trình

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[ShopeeFood 1] Lấy chủ đề: Đồ ăn-Đồ ăn
[ShopeeFood 1] Tổng số trang tìm được: 1 ứng với 18 cửa hàng
[ShopeeFood 1] Xử lý các cửa hàng trên trang số 1...
[ShopeeFood 1] Thời gian xử lý: 5.25
[ShopeeFood 1] Đã crawler hết tất cả cửa hàng của chủ đề này.
Bấm enter để kết thúc

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Foody 1] Lấy chủ đề: Tiệc tận nơi
[Foody 1] Xử lý các cửa hàng trên trang số 1...
[Foody 1] Thời gian xử lý: 12.47
[Foody 1] Xử lý các cửa hàng trên trang số 2...
[Foody 1] Thời gian xử lý: 6.67
[Foody 1] Xử lý các cửa hàng trên trang số 3...
[Foody 1] Thời gian xử lý: 13.38
[Foody 1] Xử lý các cửa hàng trên trang số 4...
[Foody 1] Thời gian xử lý: 5.21
[Foody 1] Xử lý các cửa hàng trên trang số 5...
[Foody 1] Thời gian xử lý: 12.50
[Foody 1] Xử lý các cửa hàng trên trang số 6...
[Foody 1] Thời gian xử lý: 5.36
[Foody 1] Xử lý các cửa hàng trên trang số 7...
[Foody 1] Thời gian xử lý: 5.63
[Foody 1] Xử lý các cửa hàng trên trang số 8...
[Foody 1] Thời gian xử lý: 2.62
[Foody 1] Xử lý các cửa hàng trên trang số 9...
[Foody 1] Thời gian xử lý: 0.04
[Foody 1] Đã crawler hết tất cả cửa hàng của chủ đề
Bấm enter để kết thúc
```

+ Chúng ta có thể ngưng crawl bằng cách nhập **exit** ở console quản lý

```
C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Monitor] Lượng thông tin xử lý: 180.6B / giây
[Monitor] Số trang xử lý: 0.63 / giây
[Monitor] Lượng thông tin xử lý: 663.7B / giây
[Monitor] Số trang xử lý: 0.75 / giây
[Monitor] Lượng thông tin xử lý: 793.2B / giây
[Monitor] Số trang xử lý: 0.85 / giây
[Monitor] Lượng thông tin xử lý: 891.1B / giây
[Monitor] Số trang xử lý: 0.84 / giây
[Monitor] Lượng thông tin xử lý: 876.4B / giây
[Monitor]
Bấm enter để kết thúc

C:\Windows\system32\cmd.exe - python main.py
(dev-env) C:\Users\thafnhlong\Desktop\HK3\WebScience-TopicalCrawler\Source>python m
Project: TopicalCrawler
Web of Science

Chon website cần crawler:
1. Foody
2. Shopeefood
3. Google review
Bạn chọn: 3
GoogleReview plugin
Nhập vào chủ đề cần crawl (exit để quay lại): Nhà hàng quận 5
Đang tiến hành crawl data theo chủ đề: Nhà hàng quận 5
Nhập bất kỳ để tiếp tục chọn website khác
Nhập exit để thoát chương trình
exit
Chương trình đang tắt .....

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Foody 1] Xử lý các cửa hàng trên trang số 5...
[Foody 1] Thời gian xử lý: 12.50
[Foody 1] Xử lý các cửa hàng trên trang số 6...
[Foody 1] Thời gian xử lý: 5.36
[Foody 1] Xử lý các cửa hàng trên trang số 7...
[Foody 1] Thời gian xử lý: 5.63
[Foody 1] Xử lý các cửa hàng trên trang số 8...
[Foody 1] Thời gian xử lý: 2.62
[Foody 1] Xử lý các cửa hàng trên trang số 9...
[Foody 1] Thời gian xử lý: 0.04
[Foody 1] Đã crawler hết tất cả cửa hàng của chủ đề
Bấm enter để kết thúc

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[GoogleReview 1] Lấy chủ đề: Nhà hàng quận 5
[GoogleReview 1] Xử lý các cửa hàng trên trang số 1...
[GoogleReview 1] Thời gian xử lý: 18.10
[GoogleReview 1] Xử lý các cửa hàng trên trang số 2...
```

+ Luồng chính sẽ quản lý xóa bộ nhớ thừa và tự tắt khi hoàn thành

```

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[Monitor] Lượng thông tin xử lý: 180.6B / giây
[Monitor] Số trang xử lý: 0.63 / giây
[Monitor] Lượng thông tin xử lý: 663.7B / giây
[Monitor] Số trang xử lý: 0.75 / giây
[Monitor] Lượng thông tin xử lý: 793.2B / giây
[Monitor] Số trang xử lý: 0.85 / giây
[Monitor] Lượng thông tin xử lý: 891.1B / giây
[Monitor] Số trang xử lý: 0.84 / giây
[Monitor] Lượng thông tin xử lý: 876.4B / giây
[Monitor] Bấm enter để kết thúc

C:\Windows\system32\cmd.exe
Project: TopicalCrawler
Web of Science

Chon website cần crawler:
1. Foody
2. ShopeeFood
3. Google review
Bạn chọn: 3
GoogleReview plugin
Nhập vào chủ đề cần crawl (exit để quay lại): Nhà hàng quận 5
Đang tiến hành crawl data theo chủ đề: Nhà hàng quận 5
Nhập bất kỳ để tiếp tục chọn website khác
Nhập exit để thoát chương trình
exit
Chương trình đang tắt .....ok!

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[dev-env] C:\Users\thafnhlong\Desktop\HK3\WebScience-TopicalCrawler\Source>

C:\Users\thafnhlong\AppData\Local\Programs\Python\Python39\python.exe
[GoogleReview 1] Lấy chủ đề: Nhà hàng quận 5
[GoogleReview 1] Xử lý các cửa hàng trên trang số 1...
[GoogleReview 1] Thời gian xử lý: 18.10
[GoogleReview 1] Xử lý các cửa hàng trên trang số 2...
[GoogleReview 1] Thời gian xử lý: 17.16
[GoogleReview 1] Bấm enter để kết thúc

```

+ Kiểm tra data đã crawl

baiviet (137 rows)

```
SELECT * FROM 'baiviet' WHERE website='Foody' LIMIT 0,30
```

website_id	website	url	full_name
Foody_215046	Foody	https://www.foody.vn/ho-chi...	Tiệc Lăng Mạn Cho 2 Ngườ...
Foody_3832	Foody	https://www.foody.vn/ho-chi...	Nhà Hàng Cù Đất - Âm Thụ...
Foody_200912	Foody	https://www.foody.vn/ho-chi...	The Pizza Company - Vinco...
Foody_694795	Foody	https://www.foody.vn/ho-chi...	Nhóm Nấu Chef Minh - Dìc...
Foody_294284	Foody	https://w	
Foody_1062686	Foody	https://w	
Foody_1021478	Foody	https://w	
Foody_143987	Foody	https://w	

baiviet (137 rows)

```
SELECT * FROM 'baiviet' WHERE website='GoogleReview' LIMIT 0,30
```

website_id	website	url	full_name
GoogleReview_48fadf2737...	GoogleReview	https://www.google.com/sea...	Nhà Hàng Bắ
GoogleReview_85105e938...	GoogleReview	https://www.google.com/sea...	Quán Lầu Nư
GoogleReview_2d45637a7...	GoogleReview	https://www.google.com/sea...	Nhà hàng Bi
Review	https://www.google.com/sea...		Quán Ăn Gia
Review	https://www.google.com/sea...		Nhà hàng há
Review	https://www.google.com/sea...		Nhà Hàng Cì
Review	https://www.google.com/sea...		Nhà Hàng Hì

baiviet (137 rows)

```
SELECT * FROM 'baiviet' WHERE website='ShopeeFood' LIMIT 0,30
```

website_id	website	url	full_name
ShopeeFood_66931	ShopeeFood	https://shopeefood.vn/ho-ch...	GS25 - Cửa Hàng Tiệ
ShopeeFood_194040	ShopeeFood	https://shopeefood.vn/ho-ch...	Mom Market - Thực Ph
ShopeeFood_51329	ShopeeFood	https://shopeefood.vn/ho-ch...	Yến Sào Đồng Nam Á
ShopeeFood_15228	ShopeeFood	https://shopeefood.vn/ho-ch...	Bakers' Mart - Bánh & M
ShopeeFood_23088	ShopeeFood	https://shopeefood.vn/ho-ch...	Cửa Hàng Thực Phẩm

VII/ TÀI LIỆU THAM KHẢO

<https://selenium-python.readthedocs.io/>
<https://docs.python-requests.org/projects/requests-html/en/latest/>
<https://stackoverflow.com/>
<https://codelearn.io>
<https://simple-deeplearning.medium.com/>