

Micro-module A: Online Urban Data Gathering

A4- Social Media Data Gathering and Processing

The rapid advances in information and communication technology (ICT) are rapidly integrating with the built environment, resulting in the emergence of new urban science manifesting as a new infrastructure of sensing, data gathering, and urbanism analysis. Point of Interests (POIs) is a significant type of social media data, it provides finer-grained information of urban land use and vitality.

In this micromodule, you will learn how to obtain, filter and input POIs into QGIS based on Google Place API.

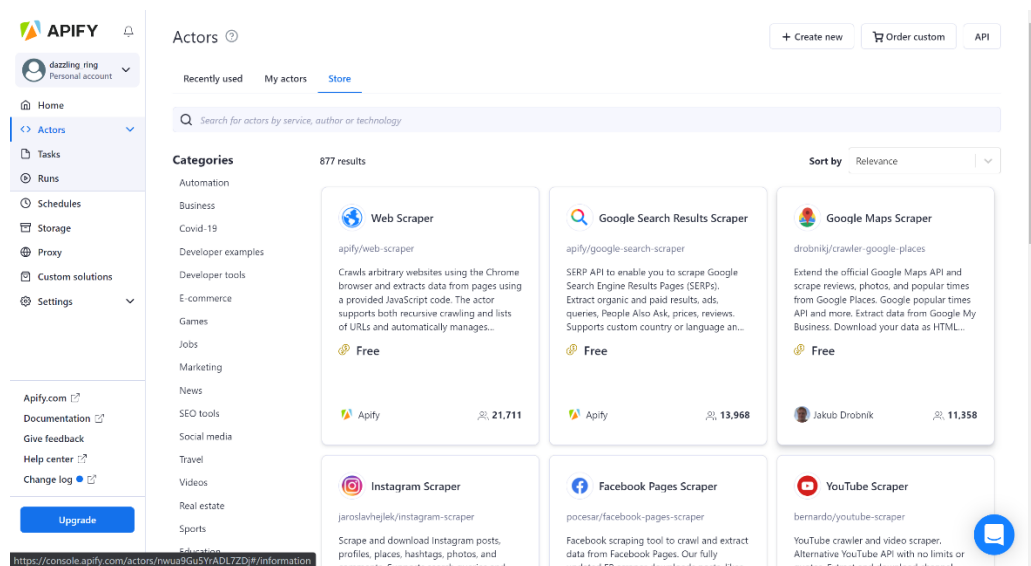
In the first part, this tutorial shows how to employ APIFY to automatically obtained selected functions of POIs in the selected research area, the output .csv file includes geo-location of POIs, review scores distribution, popular time and addresses.

In the second part, this tutorial introduces how to input .csv and filter POIs in the selected area through the 'clip' function.

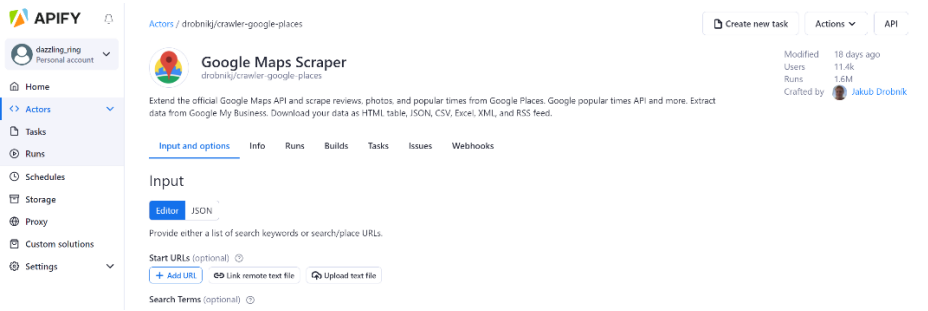
1. Obtaining the POIs

1.1 To Obtain POIs via APIFY

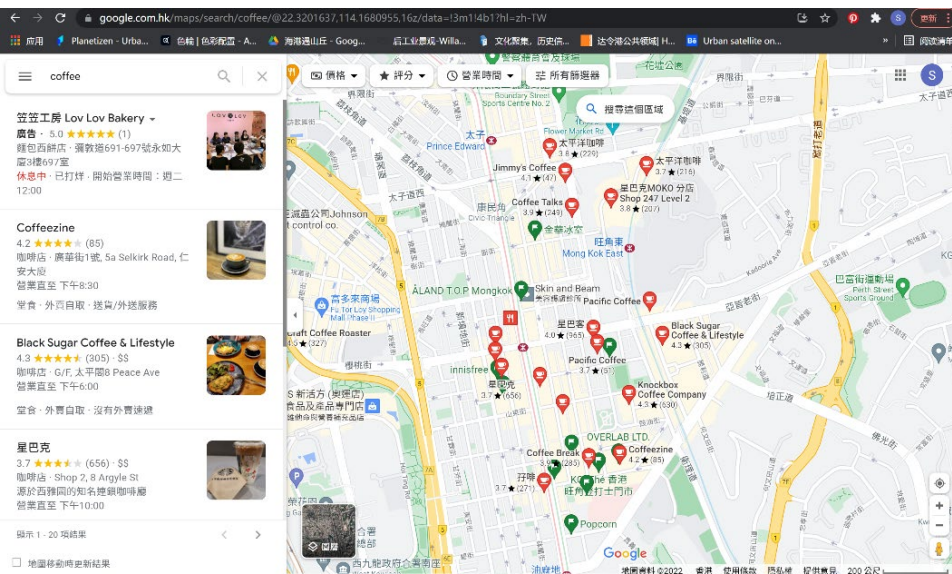
- Sign up to APIFY and find the actor called 'Google Map Scraper'. 'Google Map Scraper' extend the official Google Maps API and scrape reviews, photos, and popular times from Google Places. Google popular times API and more. Extract data from Google My Business. The output data can be downloaded as HTML table, JSON, CSV, Excel, XML, and RSS feed.



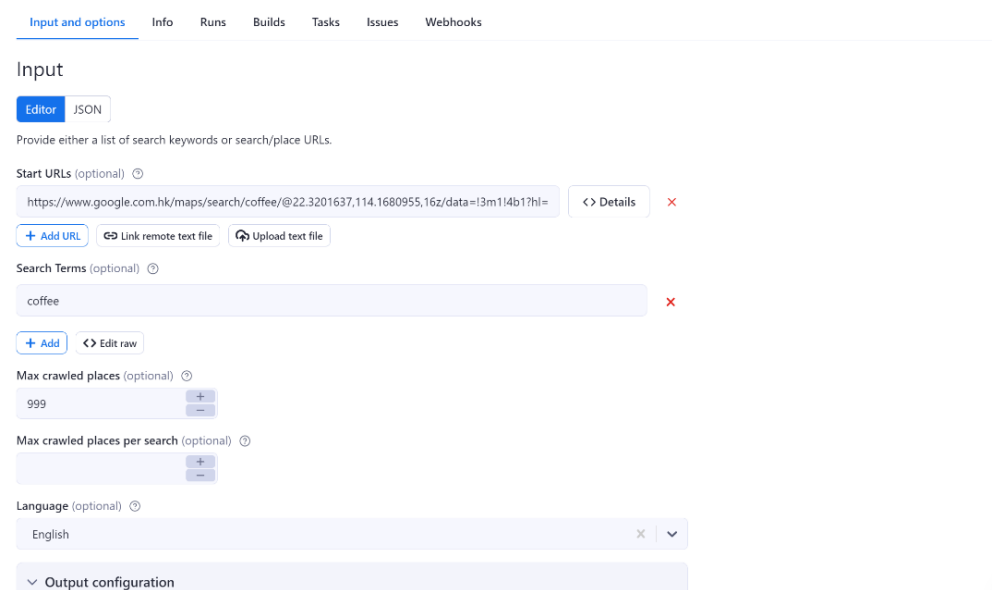
- By inputting the url of selected area and categories of POIs in Google map, the 'google map scraper' can automatically scrap POIs.



- You could zoom in to the specific area you want and search the category, for example, this tutorial search 'coffee' in Mong Kok area. You need to copy the url and paste in Google Map Scraper.



- The 'Search Terms' are array of strings to be searched. It is also possible to fill in Google Place IDs in the format `place_id:ChIjP4JiUCNP0xQR1JaSjpW_Hms`. Setting geolocation input fields is more accurate than including location in the search string. The 'Max crawled places' are the maximum number of places you scrape per whole run. If you want to scrape all available, set this to '9999999'.



- You could change the 'Output configuration' to include popular time, the dataset could be filtered to analysis amenity need pattern for a whole day. Moreover, you could change the 'Personal data' to decide if the output includes reviewers' information.

Output configuration

You can set what extra information you would like extracted. For maximum efficiency, the default setup doesn't include reviews or images. If you need these, just increase the maximum count for them.

Export place URLs only (skips place details) ⓘ
 Include popular times ⓘ
 Include opening hours ⓘ
 Include people also search ⓘ
 Additional Place Info ⓘ

Number of images (Slow for more than 1) (optional) ⓘ

0

Number of reviews (slow) (optional) ⓘ

0

Sort reviews by (optional) ⓘ

Newest

Reviews translation (optional) ⓘ

Original & translated (Google's default)

Personal data (reviews only)

All of these fields contain personal data. Personal data is protected by the GDPR in the European Union and by other regulations around the world. You should not scrape personal data unless you have a legitimate reason to do so. If you're unsure whether your reason is legitimate, consult your lawyers.

Reviewer name ⓘ
 Reviewer ID ⓘ
 Reviewer URL ⓘ
 Review ID ⓘ
 Review URL ⓘ
 Response from owner ⓘ

- Click 'Run', the scraper will start to crawl data. The dashboard shows the status of data-processing.

[+ Add](#) [Edit raw](#)

Max crawled places (optional) ⓘ

999

Max crawled places per search (optional) ⓘ

Run #oRave60UJZhLcsEtU ← Newer run Older run →

Actor drobnijs/crawler-google-places

Started at 2022-01-03 15:25:07.386 Duration 5 minutes ⓘ API

RUNNING

7 results

8 of 409 requests handled

0 of 0 webhooks finished

Log Info Input Key-value store Dataset Request queue Live view Webhooks

Showing only the latest lines of the log.

[Show full log](#) [Download log](#) [Copy log to clipboard](#)

```

2022-01-03T19:25:07.728Z ACTOR: Pulling Docker image from repository.
2022-01-03T19:25:07.818Z ACTOR: Creating Docker container.
2022-01-03T19:25:07.982Z ACTOR: Starting Docker container.
2022-01-03T19:25:11.833Z INFO: System info {"apiVersion":"1.3.4","apiClientVersion":"1.4.2","osType":"Linux","nodeVersion":"v14.18.1"}
2022-01-03T19:25:12.537Z WARN:
2022-01-03T19:25:12.521Z
2022-01-03T19:25:12.523Z -----
2022-01-03T19:25:12.526Z Using Start URLs disables search. You can use either search or Start URLs.
2022-01-03T19:25:12.538Z -----
2022-01-03T19:25:12.532Z
2022-01-03T19:25:12.722Z INFO: Prepared 1 Start URLs (showing max 10):
2022-01-03T19:25:12.726Z [
2022-01-03T19:25:12.729Z   "https://www.google.com.hk/maps/search/coffee/822.3201637,114.1688955,16z/data=!3m1!4m1!3h-TW"
2022-01-03T19:25:12.731Z ]
2022-01-03T19:25:14.452Z INFO: Full list of Start URLs is available on link: https://api.apify.com/v2/key-value-stores/bf7EDbYs7qY9u9Ush/records/START-R
2022-01-03T19:25:14.4512Z INFO: PuppeteerCrawler:AutoscaledPool: state {"currentConcurrency":0,"desiredConcurrency":2,"systemStatus":{"systemId":"tne
2022-01-03T19:25:44.498Z ACTOR: The actor run was aborted by the user.
2022-01-03T19:25:44.778Z
2022-01-03T19:28:36.988Z ACTOR: Pulling Docker image from repository.
2022-01-03T19:28:37.060Z ACTOR: Creating Docker container.
2022-01-03T19:28:37.388Z ACTOR: Starting Docker container.
2022-01-03T19:28:41.893Z INFO: System info {"apiVersion":"1.3.4","apiClientVersion":"1.4.2","osType":"Linux","nodeVersion":"v14.18.1"}
2022-01-03T19:28:41.487Z WARN: Actor was restarted, skipping search step because it was already done...
2022-01-03T19:28:42.164Z INFO: PuppeteerCrawler:AutoscaledPool: state {"currentConcurrency":0,"desiredConcurrency":2,"systemStatus":{"systemId":"tne
                
```

1.2 To Output data into .csv

- Once the data-processing is done, you could click 'Dataset' to choose the type of export data, in this case, we choose 'CSV' for it's convenient to be input into QGIS in the following steps. Click 'Download'.

The screenshot shows a web interface with a top navigation bar containing icons for 'ABORTED', '129 results', '129 of 409 requests handled', and '0 of 0 webhooks finished'. Below this is a 'Log' section with tabs for 'Info', 'Input', 'Key-value store', 'Dataset', 'Request queue', 'Live view', and 'Webhooks'. The 'Dataset' tab is selected, showing 'A storage for tabular results from the actor. Learn more'. The main content area is divided into 'Dataset details' and 'Export'. 'Dataset details' shows 'Clean items: 129', 'Items: 129', 'Dataset ID: My5P1BM2YyV0PsR', and 'Storage size: 492.6 kB'. The 'Export' section has a 'Format' dropdown with options: HTML Table, JSON, CSV (selected), Excel, XML, and RSS. Below this is an 'Advanced options' section with a 'Download' button and a 'Copy link' button.

- The output .csv data includes the address, categories, geo-location, popular time, reviews numbers, reviews score distribution.

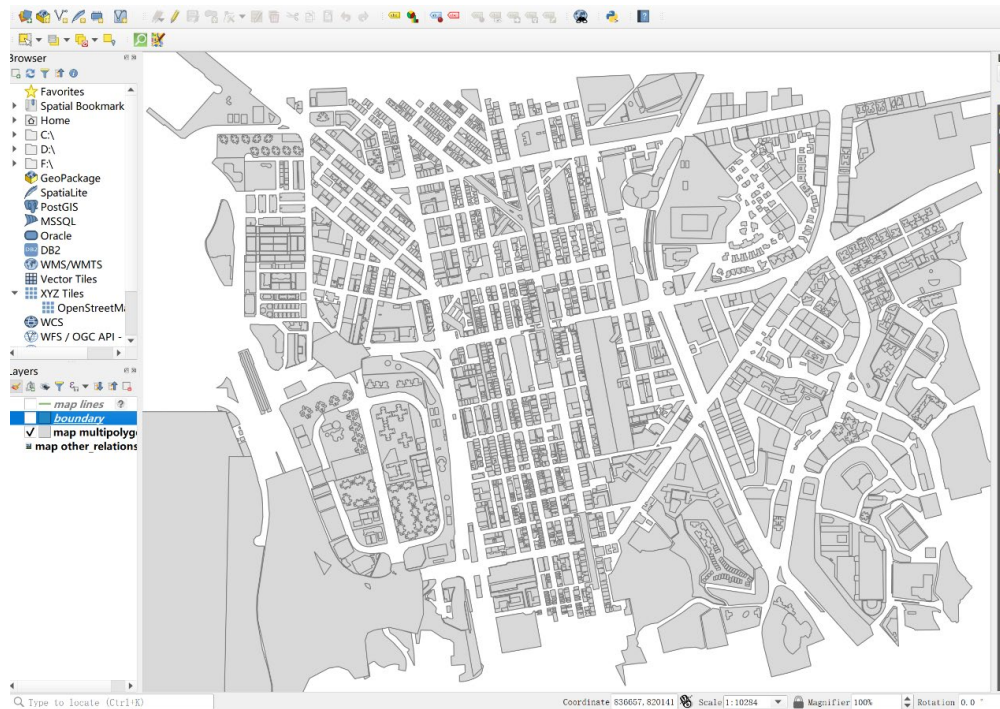
The screenshot shows a CSV data export with columns A through S. The data includes various coffee shops and cafes, such as '29 Tai Kok Cafe', 'Grateful Ti Coffee shc', and 'Shop 4201.Cafe'. The columns contain details like address, categories, geo-location (country, city, location), and reviews (reviewsCount, reviewsDistribution). The data is sorted by reviewsCount in descending order.

JW	JX	JY	JZ	KA	KB	KC	KD	KE	KF	KG	KH
popularTin	popularTin	popularTin	postalCod	price	rank	reviewsCount	reviewsDistribution /fiveStar	reviewsDistribution /fourStar	reviewsDistribution /oneStar	reviewsDistribution /threeStar	reviewsDistribution /twoStar
	41	Less busy than usual	\$\$		3	327	210	87	3	17	10
		Now: Usually not too busy			2	85	40	29	11	2	3
1		Not busy	\$\$		20	271	67	104	18	63	19
		Now: Usually a little b	\$\$\$		80	165	39	57	13	42	14
			10001	\$\$	417	1354	920	323	17	66	28
					360	4	1	1	2	0	0
					359	0	0	0	0	0	0
					358	4	1	1	0	2	0
					357	2	0	0	1	1	0
					320	0	0	0	0	0	0
					356	1	1	0	0	0	0
					355	0	0	0	0	0	0
		Now: Usually a little busy			354	2	0	0	0	1	1
					353	23	6	9	1	5	2
					352	14	5	5	1	3	0
			518010		351	31	14	7	0	9	1
					350	6	3	2	0	1	0
	100	Busier than usual			349	92	46	24	10	8	4
					348	0	0	0	0	0	0
					347	0	0	0	0	0	0
		Now: Usually a little b	\$\$\$		346	1	0	1	0	0	0
	32	Not too busy			319	4272	1738	1726	87	617	104
					318	0	0	0	0	0	0
					317	0	0	0	0	0	0
					316	0	0	0	0	0	0
					380	0	0	0	0	0	0
			20005	\$\$	416	61	42	15	0	3	1
			20007	\$\$	415	755	505	163	21	51	15
			95062	\$\$	414	375	305	52	3	11	4

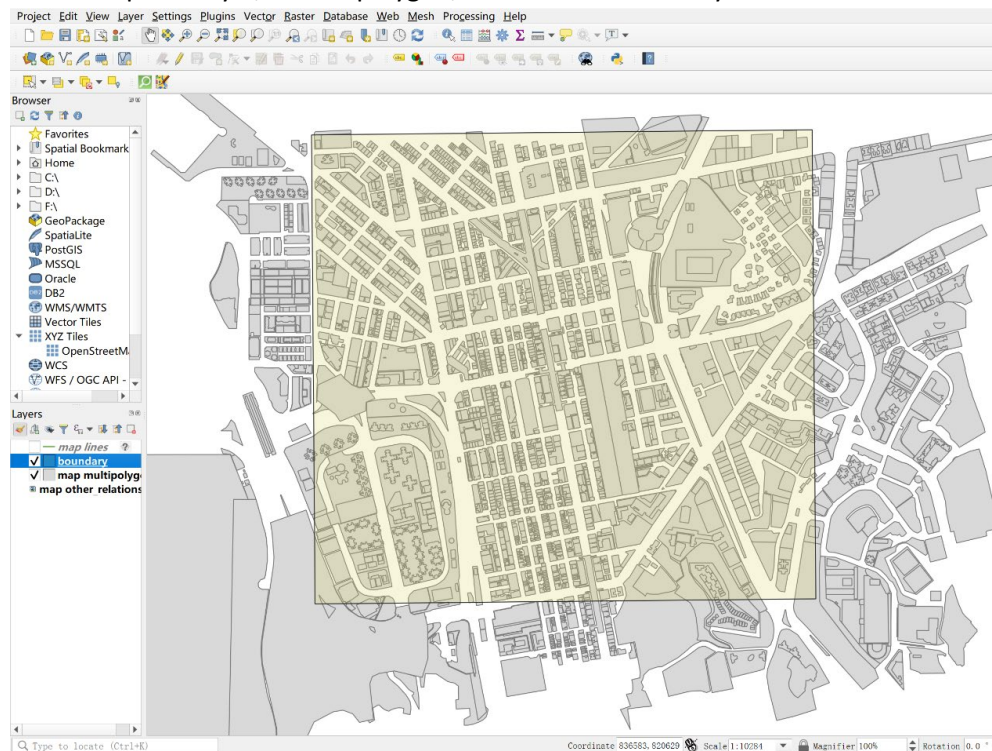
2. Processing the POIs in QGIS

2.1 To input the base map and boundary

- Input the base map through Open Street Map, following the tutorial in A3.



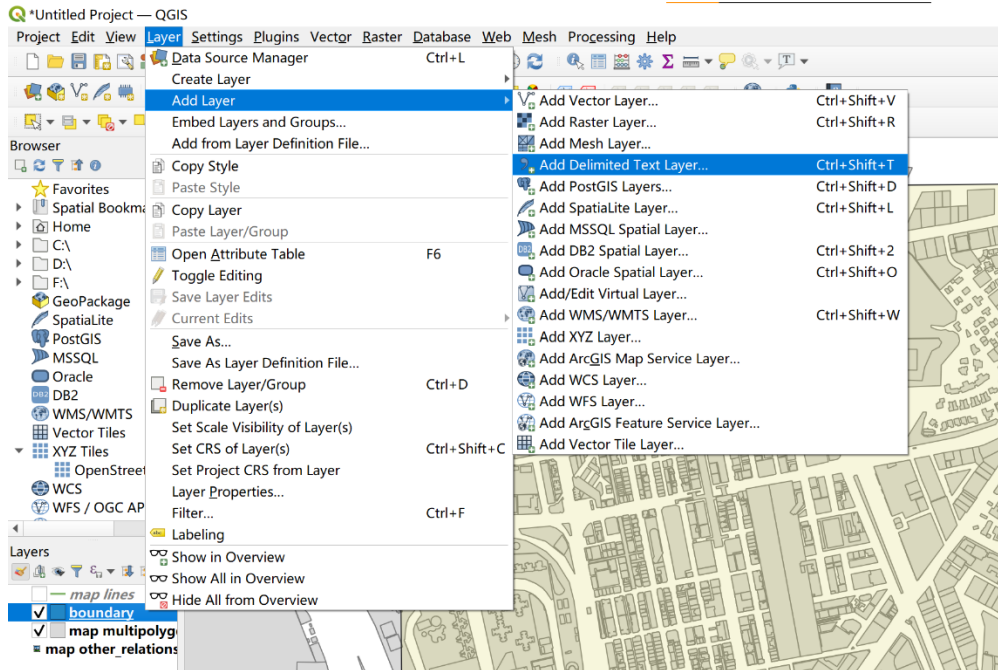
- Add a shapefile layer, draw a polygon, renamed as boundary.



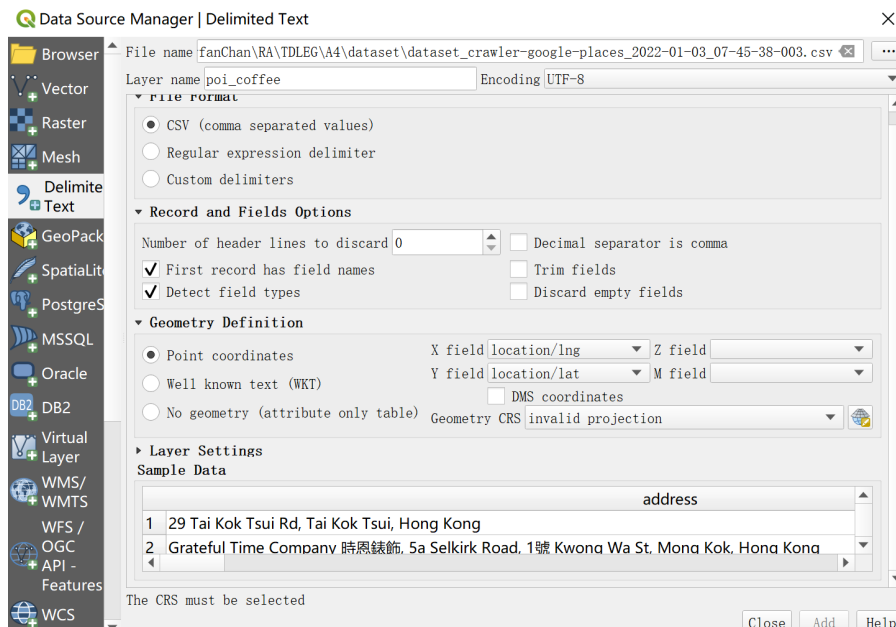
2.2 To input the .csv into QGIS

Once the csv files are formatted properly, you can add them into QGIS.

- Click on the "Layer" menu, mouse-over "Add Layer" and click on "Add Delimited Text Layer..." or click on the "Add Delimited Text Layer" icon in the left column of QGIS.



- The next GUI will have many different options you may need to change depending on the specific data set you have. Here is an outline of the most common fields needed to be changed.
- Browse - click on Browse and find the folder where the csv file is saved and open the file.
- Layer name - the name of the CSV will show up here.
- File format - depending on the version of QGIS you are using, you may need to verify the file format.



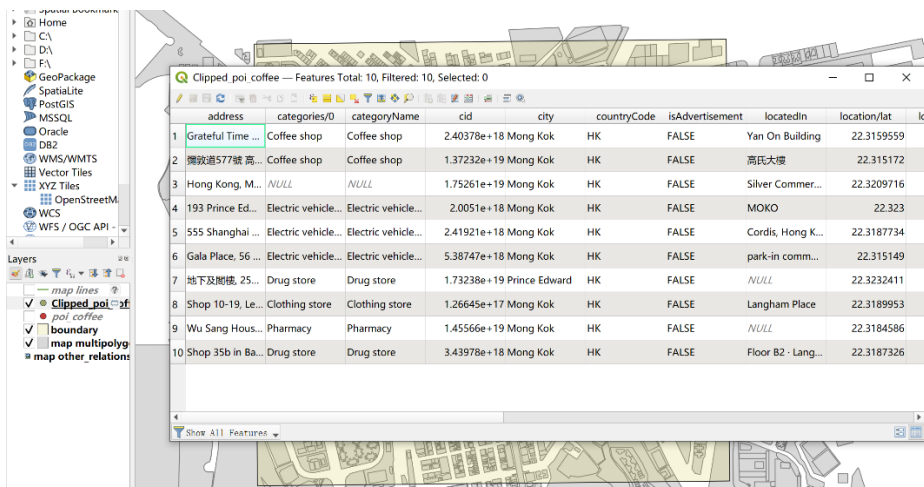
- Geometry definition - If you have x, y coordinates you will choose the "Point coordinates" option. Verify the X Field is pointing to your Longitude field and the Y Field is pointing to your Latitude Field. If you have a table with no x,y coordinates you will choose the "No Geometry" option.

- Layer settings - you will see a preview of the table. Verify that everything looks correct. Click OK.
- If you have a field with no x, y coordinates you are done importing the csv file. However, depending on the version of QGIS you are using, you may be prompted to define the coordinate reference system (CRS) of your x, y coordinates. Longitude/Latitude coordinates are unprojected, and you should choose the CRS of WGS 84 (EPSG:4326). If you have coordinates using something else, like meters in a UTM zone, search for that using the filter box in the CRS Selector Dialog. For instance, type in <utm 17> to get a short list that includes the UTM zone for Durham, 17N (e.g., NAD83 / UTM zone 17N, or EPSG:26917).

2.3 Clip/ Filter the POIs



After inputting the csv file named 'poi_coffee', you could clip the .shp with the boundary to filter POIs in the selected area.



After clipping, you could find the new .shp layer, open the attribute table, you could choose and filter the POIs as values you want.