

(I changed some items in the name and officer columns for better search results.)

Column Header Description:

index:

index (original file)

registered_number:

company number retrieved from the website

name:

name (original file)

search_key:

name (original file) if usable.

company name extracted from website (original file) if name (original file) is not usable (e.g.: H22018SUW - CA).

result_name:

name in the search result

name_matched:

whether result_name matches search_key. (None if there's no search result OR either result_name or search_key is missing.)

True if the Jaro distance between result_name and search_key is greater than 2/3.

False otherwise.

(Jaro distance is used to measure the similarity of two strings. The larger the Jaro distance is between two strings, the more similar they are to each other. The range is between 0 and 1. Source: https://rosettacode.org/wiki/Jaro_distance) (I chose Jaro distance over exact match because sometimes there's some redundant information in name (original file) that might cause name (original file) to not match result_name while name (original file) actually matches result_name.)

city: city (original file)

result_address: address in the search result

city_matched:

whether result_address contains city (original file). (None if there's no search result OR either city or result_address is missing.)

officer: officer (original file)

result_officer:

officer found in the search result that matches officer (original officer)

All officers' names concatenated if no officer found in the search result matches officer (original officer)

officer_matched:

whether result_officer matches officer (original file). (None if there's no search result OR either officer or result_officer is missing.)

True if the last names of officer and result_officer match AND the Jaro distance between their first names is greater than 2/3. (This rule strives to minimize false positives and might have caused several false negatives)

False otherwise.

(Jaro distance is used to measure the similarity of two strings. The larger the Jaro distance is between two strings, the more similar they are to each other. The range is between 0 and 1. Source: https://rosettacode.org/wiki/Jaro_distance) (I chose Jaro distance over exact match because some officers' first names are abbreviated in weird ways that might cause officer (original file) to not match result_officer while officer (original file) actually matches result_officer.)

score: (The range is between 0 and 3)

quality of the search result.

starting from 0

+ Jaro distance between result_name and search_key if name_matched is True (The range is between 0 and 1)

+ 1 if city_matched is True

+ the Jaro distance between the first names of officer and result_officer if officer_matched is True (The range is between 0 and 1)

Ranking:

≥ 2 :

Very high quality. Can be considered as an exact match.

≥ 1.5 :

name_matched is True AND officer_matched is True: High quality. Close to an exact match.

Otherwise: Need further confirmation.

< 1.5 : Low quality. (Fuzzy match)

=0: Either name (original file) is missing or there is no search result. (Not found)

Code: <https://github.com/Ziyu-Chen/companies-house-web-scraper/blob/master/index.ipynb>