

Get started

Open in app



Follow

569K Followers



How can I have a duplicate of my hotel removed from trivago?



trivago Hotel Manager
10 May 2017 05:34

Follow

A duplication of your property can appear on trivago because we received double information about it from another booking site.

Our Hotelier Care team can easily resolve this technical issue for you, and you can inform them about it by [clicking here](#).

Meanwhile, continue optimizing the fuller of the two hotel profiles with your **Hotel Information**, which includes **Hotel Details, Images, Room Types, and Descriptions**.

Photo credit: Trivago

De-duplicate the Duplicate Records from Scratch

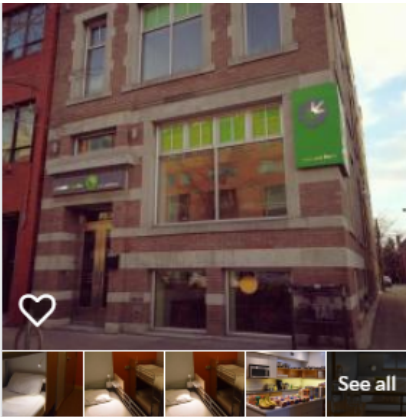
Identify similar records, Sparse matrix multiplication



Susan Li Oct 7, 2019 · 4 min read

Online world is full of duplicate listings. In particular, if you are an online travel agency, and you accept different suppliers that provide you information for the same property.

Sometimes the duplicate records are obvious that makes you think: How is it possible?



Planet Traveler Hostel

📍 Old Town Toronto, Toronto (ON) - [View on map](#)

Exceptional location

City center

Breakfast

Free cancellation

🏆 Best price for 8+ rated properties

Popular! Last booked 5 hours ago

Exceptional

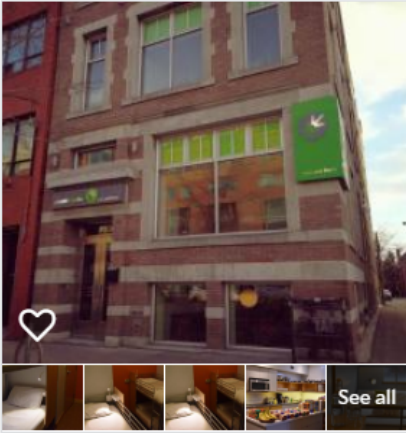
256 reviews

9.2

ONLY 2 LEFT

Price per night as low as

CAD115



Planet Traveler Hostel

📍 Old Town Toronto, Toronto (ON) - [View on map](#)

Exceptional location

City center

Breakfast

Free cancellation

🏆 Best price for 8+ rated properties

Popular! Last booked 5 hours ago

Exceptional

256 reviews

9.2

ONLY 2 LEFT

Price per night as low as


CAD115

Photo credit: agoda

Another time, the two records look like they are duplicates, but we were not sure.

Noel Suites - York and Bremner

York st and Bremner Blvd., Toronto, ON, M5J 0B1, Canada, 888-734-8514



Toronto Entertainment District

- 1.5 km to City center
- 19 km to Toronto, ON (YYZ-Pearson Intl.)

Collect nights

- Parking available • Pool • Gym
- Kitchen • Air Conditioning

Top Hotel

Exceptional 9.6

37 Hotels.com guest reviews

We have 1 left at

\$215 CAD

nightly price per unit

(\$645 CAD for 3 nights)

✓ free cancellation

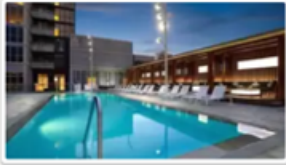
Continue

Secret Price available

Unlock now

Whitehall Suites - Front Street West

300 Front Street, Toronto, ON, M5V0E9, Canada, 888-734-8514



Toronto Entertainment District

- 1.5 km to City center
- 18 km to Toronto, ON (YYZ-Pearson Intl.)

Collect nights

- Parking available • Pool • Gym
- Kitchen • Air Conditioning

Good 7.2

53 Hotels.com guest reviews

\$209 CAD

nightly price per unit

(\$627 CAD for 3 nights)

✓ free cancellation

✓ pay now or at property


Continue

Secret Price available

Unlock now

Noel Suites-York St and Lakeshore Blvd

York Street and Lakeshore Boulevard W, Toronto, ON, M5J, Canada, 888-734-8514



Toronto Entertainment District

- 1.5 km to City center
- 19 km to Toronto, ON (YYZ-Pearson Intl.)

Collect nights

- Parking available • Pool • Gym

Loved by guests

Superb 9.0

33 Hotels.com guest reviews

We have 1 left at

\$205 CAD

nightly price per unit

(\$615 CAD for 3 nights)

✓ free cancellation

Continue

Secret Price available


[Unlock now](#)

Photo credit: expedia

Or, if you work for a company that has significant amount of data about companies or customers, but because the data comes from different source systems, in which are often written in different ways. Then you will have to deal with duplicate records.

first name	last name	address	phone
bob	roberts	1600 pennsylvania ave.	555-0123
Robert	Roberts	1600 Pensylvannia Avenue	

Photo credit: dedupe.io

The Data

I think the best data set is to use my own. Using the Seattle Hotel data set that I created a while ago. I removed hotel description feature, kept hotel name and address features, and added duplicate records purposely, and the data set can be found [here](#).

An example on how two hotels are duplicates:

```

1 from scipy.sparse import csr_matrix
2 import sparse_dot_topn.sparse_dot_topn as ct
3 from sklearn.feature_extraction.text import TfidfVectorizer
4
5 df = pd.read_csv('data/Seattle_Hotels_Duplicates.csv', encoding="latin-1")
6 df.loc[df['name'] == 'Roy Street Commons']

```

duplicate.py hosted with ❤ by GitHub

[view raw](#)

duplicate.py

	name	address
82	Roy Street Commons	621 12th Ave E, Seattle, WA 98102
90	Roy Street Commons	621 12th Avenue East, Seattle, Washington 98102

Table 1

The most common way of duplication is how the street address is input. Some are using the abbreviations and others are not. For the human reader it is obvious that the above two listings are the same thing. And we will write a program to determine and remove the duplicate records and keep one only.

TF-IDF + N-gram

- We will use name and address for input features.
- We all familiar with tfidf and n-gram methods.
- The result we get is a sparse matrix that each row is a document(name_address), each column is a n-gram. The tfidf score is computed for each n-gram in each document.

```
1 df['name_address'] = df['name'] + ' ' + df['address']
2 name_address = df['name_address']
3 vectorizer = TfidfVectorizer("char", ngram_range=(1, 4), sublinear_tf=True)
4 tf_idf_matrix = vectorizer.fit_transform(name_address)
```

tfidf_ngram.py hosted with ❤ by GitHub

[view raw](#)

tfidf_ngram.py

Sparse_dot_topn

I discovered an excellent [library that developed by ING Wholesale Banking](#), [sparse_dot_topn](#) which stores only the top N highest matches for each item, and we can choose to show the top similarities above a threshold.

It claims that it provides faster way to perform a sparse matrix multiplication followed by top-n multiplication result selection.

The function takes the following things as input:

- A and B: two CSR matrix
- ntop: n top results

- `lower_bound`: a threshold that the element of $A*B$ must greater than output

The output is a resulting matrix.

```
1  def awesome_cossim_top(A, B, ntop, lower_bound=0):
2
3      A = A.tocsr()
4      B = B.tocsr()
5      M, _ = A.shape
6      _, N = B.shape
7
8      idx_dtype = np.int32
9
10     nnz_max = M*ntop
11
12     indptr = np.zeros(M+1, dtype=idx_dtype)
13     indices = np.zeros(nnz_max, dtype=idx_dtype)
14     data = np.zeros(nnz_max, dtype=A.dtype)
15
16     ct.sparse_dot_topn(
17         M, N, np.asarray(A.indptr, dtype=idx_dtype),
18         np.asarray(A.indices, dtype=idx_dtype),
19         A.data,
20         np.asarray(B.indptr, dtype=idx_dtype),
21         np.asarray(B.indices, dtype=idx_dtype),
22         B.data,
23         ntop,
24         lower_bound,
25         indptr, indices, data)
26
27     return csr_matrix((data, indices, indptr), shape=(M, N))
28
29 matches = awesome_cossim_top(tf_idf_matrix, tf_idf_matrix.transpose(), 5)
```

awesome_cossim_top.py hosted with ❤ by GitHub

[view raw](#)

awesome_cossim_top.py

After running the function. The matrix only stores the top 5 most similar hotels.

The following code unpacks the resulting sparse matrix, the result is a table where each hotel will match to every hotel in the data(include itself), and cosine similarity score is computed for each pair.

```

1  def get_matches_df(sparse_matrix, name_vector, top=840):
2      non_zeros = sparse_matrix.nonzero()
3
4      sparserows = non_zeros[0]
5      sparsecols = non_zeros[1]
6
7      if top:
8          nr_matches = top
9      else:
10         nr_matches = sparsecols.size
11
12         left_side = np.empty([nr_matches], dtype=object)
13         right_side = np.empty([nr_matches], dtype=object)
14         similairity = np.zeros(nr_matches)
15
16         for index in range(0, nr_matches):
17             left_side[index] = name_vector[sparserows[index]]
18             right_side[index] = name_vector[sparsecols[index]]
19             similairity[index] = sparse_matrix.data[index]
20
21         return pd.DataFrame({'left_side': left_side,
22                             'right_side': right_side,
23                             'similarity': similairity})
24
25  matches_df = get_matches_df(matches, name_address)

```

get_matches_df.py hosted with ❤ by GitHub

[view raw](#)

get_matches_df.py

We are only interested in the top matches except itself. So we are going to visual examine the resulting table sort by similarity scores, in which we determine a threshold a pair is the same property.

```
matches_df[matches_df['similarity'] < 0.99999].sort_values(by=
['similarity'], ascending=False).head(30)
```

	left_side	right_side	similarity
826	Pike's Place Lux Suites by Barsala 2nd Ave and...	Pike's Place Lux Suites by Barsala 2nd Ave and...	0.715406
831	Pike's Place Lux Suites by Barsala 2nd Ave and...	Pike's Place Lux Suites by Barsala 2nd Ave and...	0.715406
206	Holiday Inn Express & Suites Seattle-City Cent...	Holiday Inn Express & Suites Seattle City Cent...	0.712321

256	Holiday Inn Express & Suites Seattle City Cent...	Holiday Inn Express & Suites Seattle-City Cent...	0.712321
181	Travelodge Seattle by The Space Needle 200 6th...	Travelodge Seattle by The Space Needle 200 6th...	0.669974
211	Travelodge Seattle by The Space Needle 200 6th...	Travelodge Seattle by The Space Needle 200 6th...	0.669974
791	citizenM Seattle South Lake Union hotel 201 We...	citizenM Seattle South Lake Union hotel 201 We...	0.651961
836	citizenM Seattle South Lake Union hotel 201 We...	citizenM Seattle South Lake Union hotel 201 We...	0.651961
586	Quality Inn & Suites Seattle Center 618 John S...	Quality Inn & Suites Seattle Center 618 John S...	0.627400
551	Quality Inn & Suites Seattle Center 618 John S...	Quality Inn & Suites Seattle Center 618 John S...	0.627400
46	Hilton Garden Inn Seattle Downtown 1821 Boren ...	Hilton Garden Inn Seattle Downtown 1821 Boren ...	0.617412
1	Hilton Garden Inn Seattle Downtown 1821 Boren ...	Hilton Garden Inn Seattle Downtown 1821 Boren ...	0.617412
781	Hyatt Regency Lake Washington At SeattleS Sout...	Hyatt Regency Lake Washington At SeattleS Sout...	0.614371
746	Hyatt Regency Lake Washington At SeattleS Sout...	Hyatt Regency Lake Washington At SeattleS Sout...	0.614371
346	Home2 Suites by Hilton Seattle Airport 380 Upl...	Home2 Suites by Hilton Seattle Airport 380 Upl...	0.582787
341	Home2 Suites by Hilton Seattle Airport 380 Upl...	Home2 Suites by Hilton Seattle Airport 380 Upl...	0.582787
561	Renaissance Seattle Hotel 515 Madison Street, ...	Renaissance Seattle Hotel 515 Madison St, Seat...	0.531174

Table 2

I decided my safe bet is to remove any pairs where the similarity score is higher than or equal to 0.50.

```
matches_df[matches_df['similarity'] < 0.50].right_side.nunique()
```

152

After that, we now have 152 properties left. If you remember, in our original data set, we did have 152 properties.

Jupyter notebook and the dataset can be found on Github. Have a productive week!

References:

<https://github.com/ing->

[bank/sparse dot topn/blob/master/sparse dot topn/awesome_cossim_topn.py](https://github.com/ing-bank/sparse/blob/master/sparse_dot_topn/awesome_cossim_topn.py)

Super Fast String Matching in Python

Traditional approaches to string matching such as the Jaro-Winkler or Levenshtein distance measure are too slow for...

bergvca.github.io

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

Get this newsletter

[Data Science](#) [Nlp Tutorial](#) [NLP](#) [Naturallanguageprocessing](#) [Cosine Similarity](#)

[About](#) [Write](#) [Help](#) [Legal](#)

Get the Medium app

