

Problem I: Locality-Sensitive Hashing

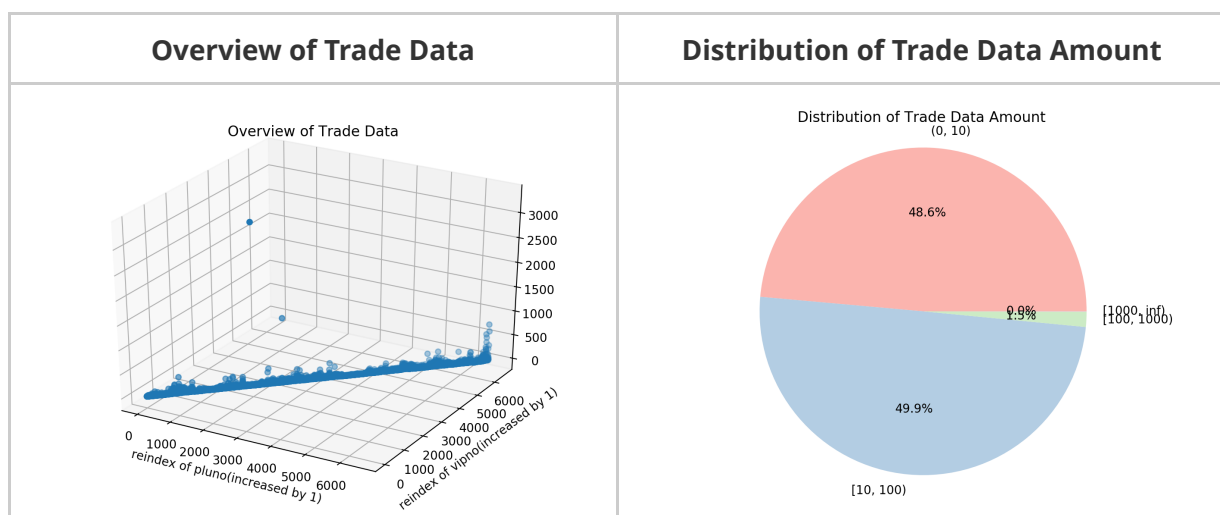
1452669, Yang LI, April 7

Data Overview & Proprocessing

Read data as DataFrame, which indexed by pluno and columned by vipno, basic information shows as follows:

```
1 <class 'pandas.core.frame.DataFrame'>
2 Index: 2635 entries, 10000004 to 40000700
3 Columns: 298 entries, 13205496418 to 6222021615015662822
4 dtypes: float64(298)
5 memory usage: 6.0+ MB
```

To have a basic impression by Data Visualization:



As we can see, it have 2 points above 1k, which lists as follows:

vipno	pluno	amt
2900001356947	14842010	3303.0
2900000350175	15114031	1334.0

Thus, the distribution of the data can be represented by the Pie Chart as above (excluding zero).

Near Neighbor using LSH

Since I use Euclidean Space, following is a brief introduction of Locality-Sensitive Hashing using s-stable distributions.

Let $f_s(t)$ denote the probability density function of the absolute value of the s-stable distribution. For the two vectors p, q , let $u = \|p - q\|_s$ and let $p(u)$ denote the probability (as a function of u) that p, q collide for a hash function uniformly chosen from family. For a random vector a whose entries are drawn from s-stable distribution, $a \cdot p - a \cdot q$ is distributed as cX where X is a random variable drawn from a s-stable distribution. Since b is drawn uniformly from $[0, w]$ it is easy to see that

$$p(u) = Pr_{a,b}[h_{a,b}(p) = h_{a,b}(q)] = \int_0^w \frac{1}{u} f_s\left(\frac{t}{u}\right) \left(1 - \frac{t}{w}\right) dt$$

In practice, all L hash tables use the same primary hash function t_1 (used to determine the index in the hash table) and the same secondary hash function t_2 . These two hash functions have the form

$$t_1(a_1, a_2, \dots, a_k) = ((\sum_{i=1}^k r'_i a_i) \bmod P) \bmod \text{tablesize}$$

$$t_2(a_1, a_2, \dots, a_k) = (\sum_{i=1}^k r''_i a_i) \bmod P$$

where r' and r'' are random integers, tablesize is the size of the hash tables, and P is a prime.

In the current usage of the LSHash module, I use `sample` to generate random vipno. Since it is an Approximate Near Neighbor, I remove the generated random vipno carefully (rather than remove the first), and the whole code shows as follows:

```
1 def knn(df, k, coefficient):
2     hash_size = int(coefficient * df.shape[1])
3     lsh = LSHash(hash_size, input_dim=df.shape[0])
4     for vipno in df:
5         lsh.index(df[vipno], extra_data=vipno)
6     random_column = df[df.columns.to_series().sample(1)]
7     random_vip = random_column.columns.values[0]
8     logging.info('random vipno: {}'.format(random_vip))
9     res = lsh.query(random_column.values.flatten())[0: k + 1]
10    print('vipno in ranked order using kNN(k = {}):'.format(k))
11    knns = []
12    for item in res:
13        if item[0][1] != random_vip:
14            print(item[0][1])
15            knns.append(item[0][1])
16    return random_vip, knns[:5]
```

Performance

Here will analyze the performance both in theory and practice.

Time & Space Complexity in Theory

Given the parameters k and L , the algorithm has the following performance guarantees:

- preprocessing time: $O(nLkt)$, where t is the time to evaluate a function $h \in \mathcal{F}$ on an input point p .
- space: $O(nL)$, plus the space for storing data points.

- query time: $O(L(kt + dnP_2^k)^L)$

Benchmark in Practice

As the statistics shows, most time cost in file I/O. In the algorithms part, `index` and `query` cost almost all time, since number of functions L is constant, it takes linear time, which is acceptable.

Timer unit: 1e-06 s

Total time: 0.073344 s
File: /Users/Yang/Developer/420235DataMining/hw1/q1/knn.py
Function: knn at line 5

Line #	Hits	Time	Per Hit	% Time	Line Contents
=====					
5					@profile
6					def knn(df, k, coefficient):
7	1	13.0	13.0	0.0	hash_size = int(coefficient * df.shape[1])
8	1	223.0	223.0	0.3	lsh = LSHash(hash_size, input_dim=df.shape[0])
9	299	443.0	1.5	0.6	for vipno in df:
10	298	63103.0	211.8	86.0	lsh.index(df[vipno], extra_data=vipno)
11	1	2874.0	2874.0	3.9	random_column = df[df.columns.to_series().sample(1)]
12	1	5.0	5.0	0.0	random_vip = random_column.columns.values[0]
13	1	203.0	203.0	0.3	logging.info('random vipno: {}'.format(random_vip))
14	1	6445.0	6445.0	8.8	res = lsh.query(random_column.values.flatten())[0: k + 1]
15	1	12.0	12.0	0.0	print('vipno in ranked order using kNN(k = {}):'.format(k))
16	1	1.0	1.0	0.0	knns = []
17	7	4.0	0.6	0.0	for item in res:
18	6	4.0	0.7	0.0	if item[0][1] != random_vip:
19	5	7.0	1.4	0.0	print(item[0][1])
20	5	6.0	1.2	0.0	knns.append(item[0][1])
21	1	1.0	1.0	0.0	return random_vip, knns[:5]

Line #	Mem usage	Increment	Line Contents
=====			
5	112.758 MiB	112.758 MiB	@profile
6			def knn(df, k, coefficient):
7	112.758 MiB	0.000 MiB	hash_size = int(coefficient * df.shape[1])
8	112.777 MiB	0.020 MiB	lsh = LSHash(hash_size, input_dim=df.shape[0])
9	130.707 MiB	0.000 MiB	for vipno in df:
10	130.707 MiB	17.930 MiB	lsh.index(df[vipno], extra_data=vipno)
11	130.969 MiB	0.262 MiB	random_column = df[df.columns.to_series().sample(1)]
12	130.969 MiB	0.000 MiB	random_vip = random_column.columns.values[0]
13	130.969 MiB	0.000 MiB	logging.info('random vipno: {}'.format(random_vip))
14	131.016 MiB	0.047 MiB	res = lsh.query(random_column.values.flatten())[0: k + 1]
15	131.016 MiB	0.000 MiB	print('vipno in ranked order using kNN(k = {}):'.format(k))
16	131.016 MiB	0.000 MiB	knns = []
17	131.016 MiB	0.000 MiB	for item in res:
18	131.016 MiB	0.000 MiB	if item[0][1] != random_vip:
19	131.016 MiB	0.000 MiB	print(item[0][1])
20	131.016 MiB	0.000 MiB	knns.append(item[0][1])
21	131.016 MiB	0.000 MiB	return random_vip, knns[:5]

Screenshot

