

Fall 2022 MAT265 Random Samplings

- Shaokang Li

Introduction:

Random Sampling is a method to estimate characteristics of the whole population by sample a subset within the whole population randomly. For this week's assignment, I am interested in questions below:

- By using random sampling technique, I try to find if there is any pattern within the sampled data? Is the pattern any different from the whole population?
- If the pattern sampled in CD is any different from that in DVD or Books?
 - For each media type, sample 3 times in a row, compare their differences.
- Note that the bibNumber and itemNumber follows a linear pattern within a lower itemNumber range. If we randomly sample the bibNumber and itemNumber within the same range, will the same pattern appear?

Question 01:

For the first question. I first tried out the general query to random sampling in the whole `inraw` table by using query below.

Initial Query:

This query should randomly return 100 records from the `inraw` table. However, a error message was returned instead.

```
1 SELECT *
2 FROM spl_2016.inraw
3 order by rand(42) LIMIT 100;
```

Message:

```
SQL: SELECT * FROM spl_2016.inraw order by rand(42) LIMIT 100
```

```
The table '/tmp/#sql_8a1_20' is full
```

After searching on [stackoverflow](https://stackoverflow.com), the error is likely due to a limit in the size of temporary table. If we don't apply any restriction to the query, the temporary table will exceeds the size limit.

Thus, I decided to do the random sampling in the CD category and set up a time scope from 2022 to today.

Updated Query

Query below is used to find the checkout times of a single title, with a randomly selection size of 100.

```
1 SELECT
2     count(cout) as checkouts,
3     bibNumber,
4     title
5 FROM spl_2016.inraw
6 where
7     itemtype in (
8         'arcd',
9         'nacd',
10        'jrzd',
11        'accd',
12        'cacz',
13        'cccd',
14        'jccd',
15        'nccd'
16    )
17     and cout > '2022-01-01'
18 group by bibNumber,title
19 order by rand(42)
20 LIMIT 100;
```

Result:

See attached CD_100_#.csv for results.

After getting the result, I applied the deviation-based outlier method from last week to find out the ratio of outliers in this randomly sampled data. The outliers are the CDs with extremely low or high checkout times. And I'm interested in the difference between sampled data and the whole data.

```
1 import pandas as pd
2 import numpy as np
3
4 df = pd.read_csv('CD_100_01.csv')
5 df.head()
6 print(df.shape)
7
8 diff = df.checkouts.quantile(0.75)-df.checkouts.quantile(0.25)
9 upper_limit = df.checkouts.quantile(0.75) + 1.5*diff
10 lower_limit = df.checkouts.quantile(0.25) - 1.5*diff
11 print(upper_limit, lower_limit)
12
13 filtered_data = df[df.checkouts>upper_limit].sort_values(by=
14 ['checkouts'],ascending=False)
15 ratio = filtered_data.shape[0]/df.shape[0]
16 print(ratio)
```

```
16 filtered_data.head()
```

Results are listed below:

Sample Times	Upper Limit for Non-outlier	Lower Limit for Non-outliers	Outlier Ratio
01	9.125	-3.875	7%
02	11.0	-5.0	5%
03	11.0	-5.0	10%
Whole data	11.0	-5.0	7.224%

Even though we only perform the random sample 3 times, the average outlier ratio is $(7 + 5 + 10)/3 \approx 7.3$, which is very close to the result on whole data (7.224%). In 2 out of 3 queries, the upper limit and lower limit for non-outliers are the same as those on the whole data.

Problem 02:

I would like to perform the procedure above on DVD and Books, trying to find out if the randomly sampling also produce a correct estimation of the outlier ratio.

Query:

```
1  SELECT
2      count(cout) as checkouts,
3      bibNumber,
4      title
5  FROM spl_2016.inraw
6  where
7      itemtype in (
8          'cadvdf',
9          'ncdvdnf',
10         'nalndvd',
11         'jrdvd',
12         'jcdvd',
13         'acdvd',
14         'nadvd',
15         'ardvd',
16         'ccdvdnf',
17         'ncdvd',
18         'nadvdnf'
19     )
20     and cout > '2022-01-01'
21  group by bibNumber, title
22  order by rand(42)
23  LIMIT 100;
```

```

24
25 SELECT
26     count(cout) as checkouts,
27     bibNumber,
28     title
29 FROM spl_2016.inraw
30 where
31     itemtype in (
32         'cabknf',
33         'ncbknf',
34         'nalnbk',
35         'jrbk',
36         'jcbk',
37         'acbknf',
38         'nabk',
39         'arbk',
40         'ccbknf',
41         'ncbk',
42         'nabknf',
43         'jcbknf'
44     )
45     and cout > '2022-01-01'
46 group by bibNumber, title
47 order by rand(44)
48 LIMIT 100;

```

RESULT:

See attached DVD_100_#.csv / Book_100_#.csv for results.

Sample Times	Upper Limit for Non-outlier	Lower Limit for Non-outliers	Outlier Ratio
Book_100_01	16.0	-8.0	10%
Book_100_02	20.125	-8.875	14%
Book_100_03	18.5	-9.5	14%
Book_all	18.5	-9.5	11.9%
DVD_100_01	38.0	-18.0	5%
DVD_100_02	41.75	-20.25	5%
DVD_100_03	59.625	-29.375	4%
DVD_all	44.5	-22.5	8.9%

The data on books performed similarity to CDs. The randomly sampled average outlier ratio is approximately the same as the overall outlier ratio.

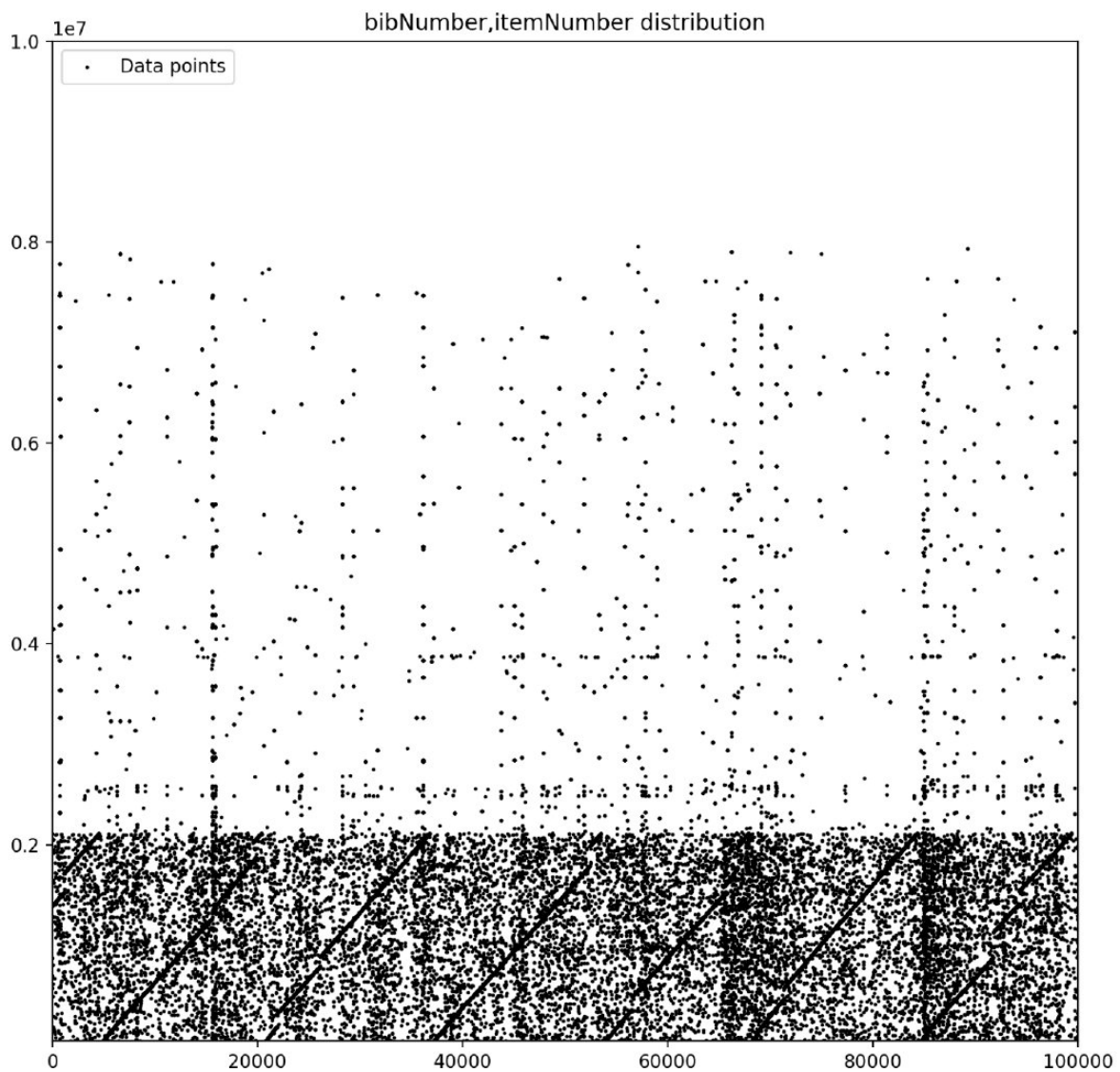
$(38\%/3 \approx 12.7)$

However, this doesn't apply to the DVD category. For the DVD category, the randomly sampled average outlier ratio is 4.6%, whereas the overall outlier ratio is 8.9%, nearly 2 times difference. This is can be the sampling error or the dataset for DVD category is skewed.

If we compare the limits for each media types, we find that DVD has a higher limit, then it is Book and then CD. The limits difference indicate the deviation of the items' checkout times. In this case, DVD has a much higher deviation than the rest category.

Question 03:

From last week's assignment, there is a linear pattern in the distribution of itemNumber and bibNumber. (As shown in the graph below)



If we randomly sampled some datapoints within the same range of bibNumber, will the sampled data follow the same pattern or produce a different distribution pattern?

Query:

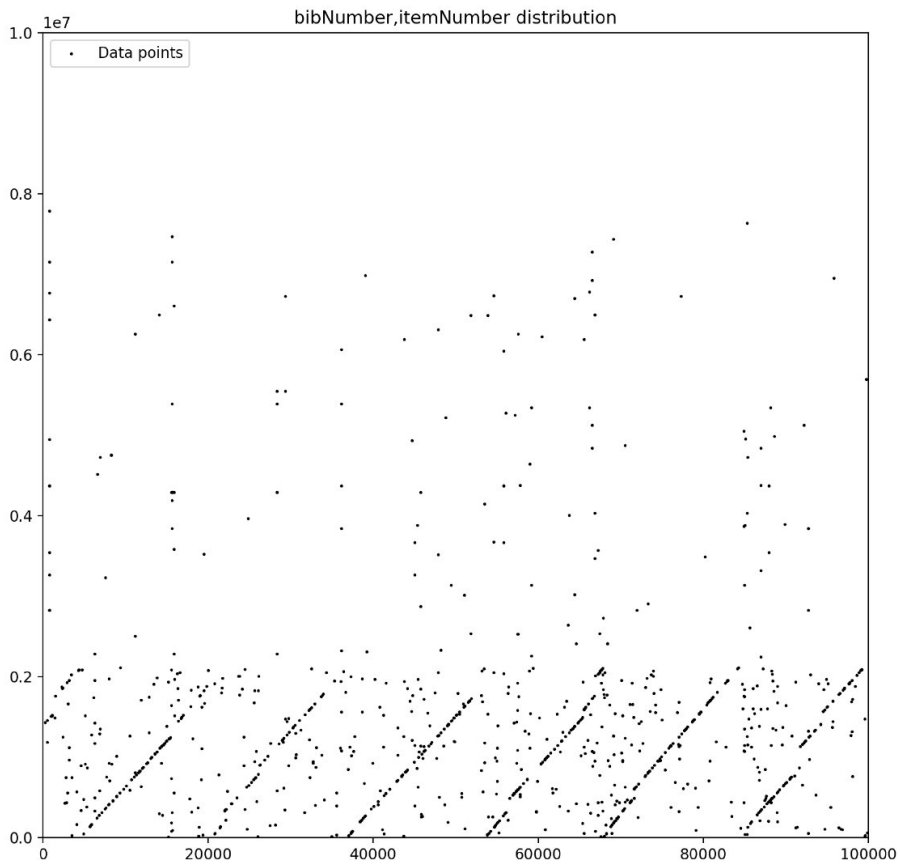
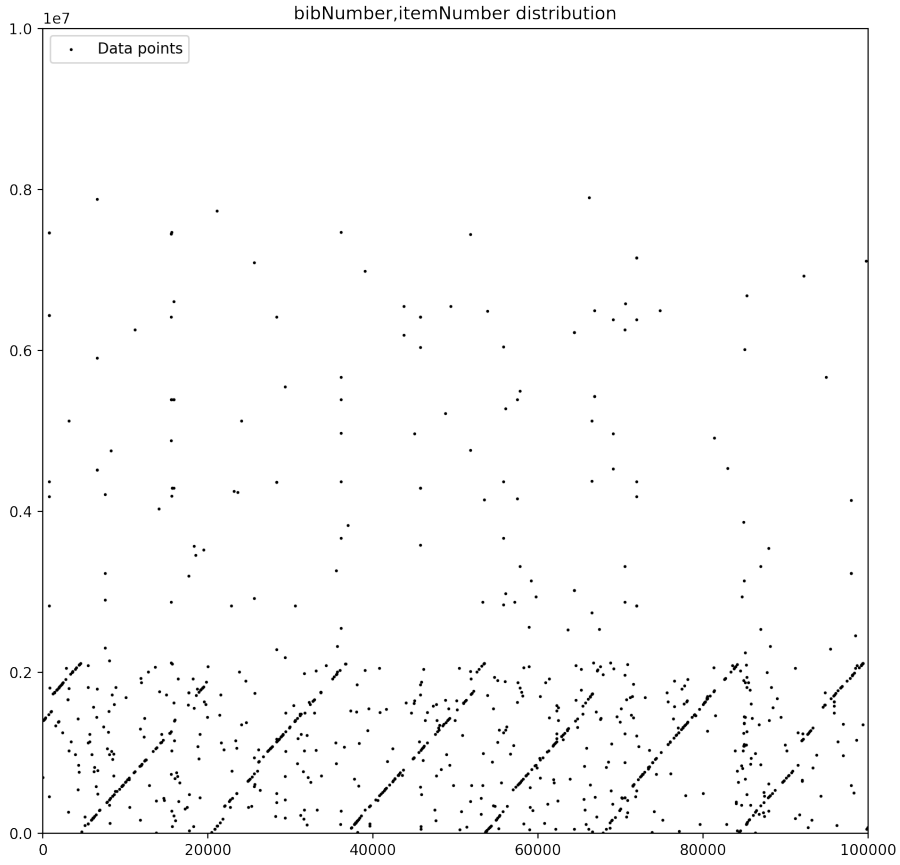
```
1 select
2     distinct bibNumber,
3     itemNumber
4 from spl_2016.inraw
5 where `bibNumber` < 100000
6 order by rand(44),bibNumber
7 LIMIT 1000;
```

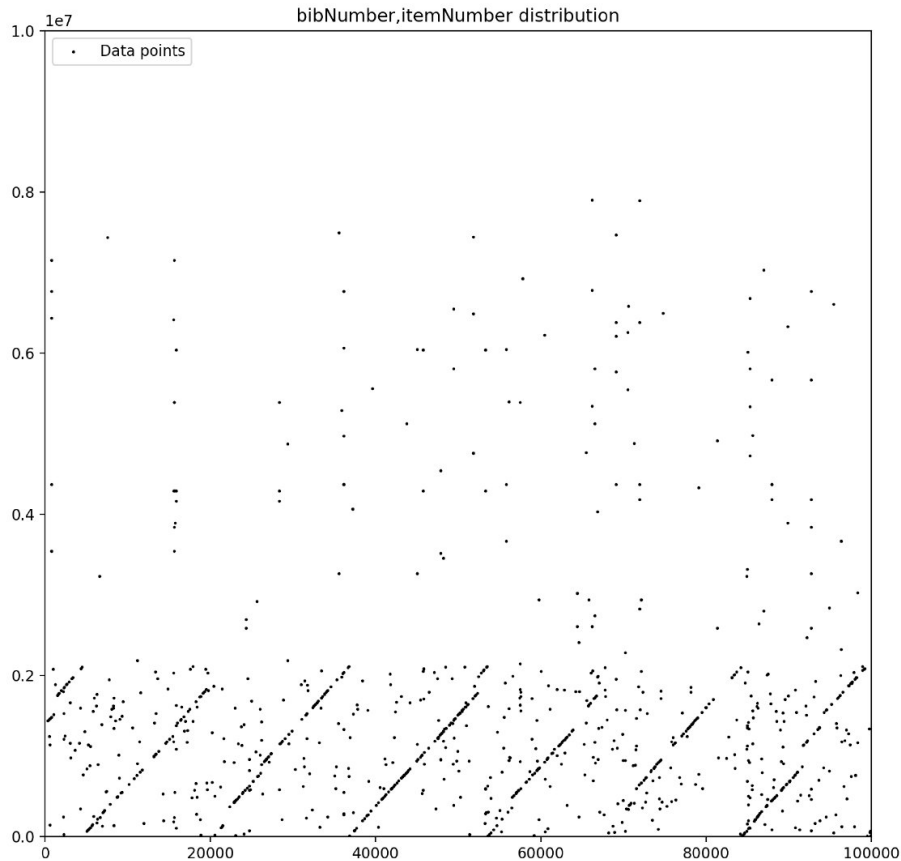
See bibDist_#.csv for results

The random sampling is performed 3 times. Each time the data will be feed to below visualization code to generate the distribution of the bibNumber and itemNumber:

```
1 arr_dist = np.array(df_dist[['bibNumber','itemNumber']])
2 figure(figsize=(10, 10), dpi=200)
3 # Generate train data
4 X = arr_dist
5
6 plt.title("bibNumber,itemNumber distribution")
7 plt.scatter(X[:, 0], X[:, 1], color="k", s=1.0, label="Data points")
8
9 plt.axis("tight")
10 plt.xlim((0, 100000))
11 plt.ylim((50000, 10000000))
12 legend = plt.legend(loc="upper left")
13 # legend.legendHandles[0]._sizes = [10]
14 # legend.legendHandles[1]._sizes = [20]
15 plt.show()
16 plt.savefig('img.png')
```

Result:





The randomly sampled results are similar to the overall figure, 6 lines all appears in the 3 graphs and they share similar starting points and period, which means some of the itemNumbers start over after 18000 bibNumbers. The linear pattern is much more evident in the random sampled datapoints than the overall one.

Conclusion:

Comparing the data between random sampling and the whole data, here's some of my observations.

- The random sampled results are similar to the overall dataset on CD and Books, whereas for DVD, this doesn't hold. This may indicate that DVD checkout has a more uneven distribution.
- The result on the overall dataset is approximately the average of the randomly sampled results.
- For the bibNumber / itemNumber distribution, the random sample result follows the similar pattern from the overall dataset. But the pattern is more prominent on the sampled result.