

Data Transformation

GAHAN SARAIYA (18MCEC10), PRIYANKA BHATI (18MCEC02)

18mcec10@nirmauni.ac.in, 18mcec02@nirmauni.ac.in

I. INTRODUCTION

Data transformation is the process of converting/transforming data from one form¹ or structure into another form or structure. Data transformation is critical to activities such as data integration and data management. Data transformation can include a range of activities such as: converting data types, data cleaning by removing null values or duplicate data, enrich the data, or perform aggregations, depending on the requirements of the project.

A typical data transformation can be classified in two stages:

- Stage 1
 - Discovering the data and identify various sources and data types.
 - Find out the structure of data and data transformations that need to be performed.
 - Map data to define how individual fields are mapped, modified, joined, filtered, and aggregated.
- Stage 2
 - Extract data from the original source.
The range of sources can vary from structured sources like databases to non-structured sources such as streaming sources.
 - Do the transformation
i.e. aggregating values, changing date formats, editing text or joining attributes.
 - Migrate this transformed data to data warehouse.

II. TRANSFORMATION OF DATA

I. Data Discovery

We have taken a dataset as described in Section IV.I for which we need to perform knowledge gain of user mobility with correlation between locations.

II. Data Structure

Data is a raw data of multiple text files (one file per taxi id) containing records of booking of taxi with timestamp and latitude and longitude of pickup location.

III. Extracting Data

To perform operation we will first combine these multiple data files into one single file containing all the record of taxi-id, latitude, longitude and timestamp

¹i.e. format for different data mining tools

IV. Performing Transformation

The two attributes in dataset [IV.I](#), latitude and longitude are concatenate together to form unique area with method known as geohash described in Section [IV.II](#).

This will generate a file (as in [IV.II.2](#)) containing `taxi-id`, `geohash` and `timestamp`

III. DATA VISUALIZATION

Now we need to compute correlation between locations from our transformed data [IV.II.2](#).

To compute correlation we simply create dissimilarity between two records (which is considered as the timestamp).

Visualization result is obtained as in Section [IV.IV.1](#) which contains `geohash`, `geohash` and `dissimilarity`.

IV. EXPERIMENT

I. Dataset

For the experiments we have taken a dataset of user mobility which contains data attributes as described below:

Attribute	Data-type	Description
Taxi-id	Numerical (Integer)	unique id of each taxi
Latitude	Discrete	latitude of pickup location
Longitude	Discrete	longitude of pickup location
Timestamp	date-time	time-stamp when taxi booked

Table 1: Description of Dataset

I.1 Chunk 1

```
2, 116.36422, 39.88781, 2008-02-02 13:33:52
2, 116.37481, 39.88782, 2008-02-02 13:37:16
2, 116.37677, 39.88791, 2008-02-02 13:38:53
2, 116.38033, 39.88795, 2008-02-02 13:42:18
2, 116.39392, 39.89014, 2008-02-02 13:43:55
2, 116.41149, 39.89152, 2008-02-02 13:47:20
2, 116.42105, 39.89194, 2008-02-02 13:48:57
2, 116.4215, 39.89823, 2008-02-02 13:52:21
2, 116.4215, 39.89902, 2008-02-02 13:53:58
2, 116.42972, 39.90726, 2008-02-02 13:57:23
```

I.2 Chunk 2

```
3, 116.50029, 39.91458, 2008-02-03 23:18:57
3, 116.51617, 39.91126, 2008-02-03 23:23:57
3, 116.55026, 39.90787, 2008-02-03 23:28:56
3, 116.57449, 39.90895, 2008-02-03 23:33:56
3, 116.59587, 39.9006, 2008-02-03 23:43:56
3, 116.61065, 39.91218, 2008-02-03 23:48:56
3, 116.6347, 39.92558, 2008-02-03 23:53:56
3, 116.652, 39.93273, 2008-02-03 23:58:56
3, 116.65518, 39.95388, 2008-02-04 00:03:56
3, 116.69494, 39.98554, 2008-02-04 00:08:56
3, 116.75385, 40.02126, 2008-02-04 00:18:57
```

II. Geohash

The Geohash is the method for encoding regions with specific precision of latitude and longitude.

Geohashes offer properties like arbitrary precision and the possibility of gradually removing characters from the end of the code to reduce its size (and gradually lose precision).

For example, the coordinate pair 57.64911, 10.40744 (near the tip of the peninsula of Jutland, Denmark) produces a slightly shorter hash of u4pruydqvj.

Implementation of Geohash is described below:

II.1 Implementation

```
1 from math import log10
2
3 # Note: the alphabet in geohash differs from the common base32
4 # alphabet described in IETF's RFC 4648
5 # (http://tools.ietf.org/html/rfc4648)
6 __base32 = '0123456789bcdefghjkmnpqrstuvwxyz'
7 __decodemap = {}
8
9
10 for i in range(len(__base32)):
11     __decodemap[__base32[i]] = i
12 del i
13
14
15 def decode_exactly(geohash):
16     """
```

```
17 Decode the geohash to its exact values, including the error
18 margins of the result. Returns four float values: latitude,
19 longitude, the plus/minus error for latitude (as a positive
20 number) and the plus/minus error for longitude (as a positive
21 number).
22 """
23 lat_interval, lon_interval = (-90.0, 90.0), (-180.0, 180.0)
24 lat_err, lon_err = 90.0, 180.0
25 is_even = True
26 for c in geohash:
27     cd = __decodemap[c]
28     for mask in [16, 8, 4, 2, 1]:
29         if is_even: # adds longitude info
30             lon_err /= 2
31             if cd & mask:
32                 lon_interval = ((lon_interval[0] + lon_interval[1]) / 2,
33                               ↪ lon_interval[1])
34             else:
35                 lon_interval = (lon_interval[0], (lon_interval[0] +
36                               ↪ lon_interval[1]) / 2)
37         else: # adds latitude info
38             lat_err /= 2
39             if cd & mask:
40                 lat_interval = ((lat_interval[0] + lat_interval[1]) / 2,
41                               ↪ lat_interval[1])
42             else:
43                 lat_interval = (lat_interval[0], (lat_interval[0] +
44                               ↪ lat_interval[1]) / 2)
45         is_even = not is_even
46     lat = (lat_interval[0] + lat_interval[1]) / 2
47     lon = (lon_interval[0] + lon_interval[1]) / 2
48     return lat, lon, lat_err, lon_err
49
50 def decode(geohash):
51     """
52     Decode geohash, returning two strings with latitude and longitude
53     containing only relevant digits and with trailing zeroes removed.
54     """
55     lat, lon, lat_err, lon_err = decode_exactly(geohash)
56     # Format to the number of decimals that are known
57     lats = "%. *f" % (max(1, int(round(-log10(lat_err)))) - 1, lat)
```

```
55     lons = "%. *f" % (max(1, int(round(-log10(lon_err)))) - 1, lon)
56     if '.' in lats:
57         lats = lats.rstrip('0')
58     if '.' in lons:
59         lons = lons.rstrip('0')
60     return lats, lons
61
62
63 def encode(latitude, longitude, precision=10):
64     """
65     Encode a position given in float arguments latitude, longitude to
66     a geohash which will have the character count precision.
67     """
68     lat_interval, lon_interval = (-90.0, 90.0), (-180.0, 180.0)
69     geohash = []
70     bits = [16, 8, 4, 2, 1]
71     bit = 0
72     ch = 0
73     even = True
74     while len(geohash) < precision:
75         if even:
76             mid = (lon_interval[0] + lon_interval[1]) / 2
77             if longitude > mid:
78                 ch |= bits[bit]
79                 lon_interval = (mid, lon_interval[1])
80             else:
81                 lon_interval = (lon_interval[0], mid)
82         else:
83             mid = (lat_interval[0] + lat_interval[1]) / 2
84             if latitude > mid:
85                 ch |= bits[bit]
86                 lat_interval = (mid, lat_interval[1])
87             else:
88                 lat_interval = (lat_interval[0], mid)
89         even = not even
90         if bit < 4:
91             bit += 1
92         else:
93             geohash += __base32[ch]
94             bit = 0
95             ch = 0
96     return ''.join(geohash)
```

II.2 Data File after applying geohash

II.3 Chunk 1

```
2, uzurg, 2008-02-02 13:33:52
2, uzurg, 2008-02-02 13:37:16
2, uzurg, 2008-02-02 13:38:53
2, uzurg, 2008-02-02 13:42:18
2, uzurg, 2008-02-02 13:43:55
2, uzurg, 2008-02-02 13:47:20
2, uzurg, 2008-02-02 13:48:57
2, uzurg, 2008-02-02 13:52:21
2, uzurg, 2008-02-02 13:53:58
2, uzuru, 2008-02-02 13:57:23
```

II.4 Chunk 2

```
3, uzuru, 2008-02-03 23:18:57
3, uzuru, 2008-02-03 23:23:57
3, uzuru, 2008-02-03 23:28:56
3, uzuru, 2008-02-03 23:33:56
3, uzurg, 2008-02-03 23:43:56
3, uzuru, 2008-02-03 23:48:56
3, uzuru, 2008-02-03 23:53:56
3, uzuru, 2008-02-03 23:58:56
3, uzurv, 2008-02-04 00:03:56
3, uzurv, 2008-02-04 00:08:56
3, uzury, 2008-02-04 00:18:57
```

III. Implementation of Utilities for parsing data

```
1 import os
2 import geohash
3
4 datasets_dir = os.path.abspath(os.path.join(os.path.dirname(__file__),
    ↳ 'datasets'))
5 dataset_files = [os.path.join(datasets_dir, file) for file in
    ↳ os.listdir(datasets_dir)]
```

```
6
7
8 def parseline(line, precision=5):
9     if line[-1] == '\n':
10         line = line[:-1]
11     taxi_id, timestamp, latitude, longitude = line.split(',')
12     latitude, longitude = float(latitude), float(longitude)
13     gh = geohash.encode(latitude, longitude, precision=precision)
14     return {
15         'taxi_id': taxi_id,
16         'geohash': gh,
17         'latitude': latitude,
18         'longitude': longitude,
19         'timestamp': timestamp
20     }
21
22
23 def nextline(dataset_file, precision):
24     taxi_id = dataset_file.split('.')[0]
25     taxi_id = taxi_id.split('/')[1]
26     with open(dataset_file, mode='r') as fp:
27         for line in fp.readlines():
28             yield parseline(line, precision)
29     yield None
30
31
32 def get_random_data(precision):
33     file_pointers = [nextline(dataset_file, precision) for dataset_file in
34                      ↪ dataset_files]
35     while len(file_pointers) > 0:
36         for fp in file_pointers:
37             line = fp.__next__()
38             if line is None:
39                 file_pointers.remove(fp)
40             else:
41                 # print(line)
42                 yield line
43
44 def get_sequential_data(precision):
45     file_pointers = [nextline(dataset_file, precision) for dataset_file in
46                      ↪ dataset_files]
```

```
46     for fp in file_pointers:
47         line = fp.__next__()
48         while line is not None:
49             yield line
50             line = fp.__next__()
51
52     # Print iterations progress
53     def printProgressBar (iteration, total, prefix = '', suffix = '', decimals = 1,
54         ↪ length = 100, fill = '='):
55         """
56         Call in a loop to create terminal progress bar
57         @params:
58             iteration - Required : current iteration (Int)
59             total     - Required : total iterations (Int)
60             prefix    - Optional : prefix string (Str)
61             suffix    - Optional : suffix string (Str)
62             decimals  - Optional : positive number of decimals in percent complete
63             ↪ (Int)
64             length    - Optional : character length of bar (Int)
65             fill      - Optional : bar fill character (Str)
66         """
67         percent = ("{0:." + str(decimals) + "f}").format(100 * (iteration /
68             ↪ float(total)))
69         filledLength = int(length * iteration // total)
70         bar = fill * filledLength + '-' * (length - filledLength)
71         print('\r%s |%s| %s%% %s' % (prefix, bar, percent, suffix), end = '\r')
72         # Print New Line on Complete
73         if iteration == total:
74             print()
```

IV. Implementation of storing and transforming data

```
1  import os
2  import simulator
3  import itertools
4  from time import sleep
5  from datetime import datetime
6
7  PRECISION = 5
8
9  def generate_file_1():
```

```
10     with open('file1.txt', 'w') as fp:
11         for df in simulator.get_sequential_data(precision=PRECISION):
12             fp.write("{} , {} , {} , {}\n".format(
13                 df['taxi_id'], df['latitude'], df['longitude'], df['timestamp']))
14
15 def generate_file_2():
16     with open('file2.txt', 'w') as fp:
17         for df in simulator.get_sequential_data(precision=PRECISION):
18             fp.write("{} , {} , {}\n".format(
19                 df['taxi_id'], df['geohash'], df['timestamp']))
20
21
22 def parseline(line):
23     if line[-1] == '\n':
24         line = line[:-1]
25     taxi_id, geohash, timestamp = line.split(',')
26     epoch_time = datetime.strptime(timestamp, ' %Y-%m-%d %H:%M:%S').timestamp()
27     return geohash, epoch_time, taxi_id
28
29 def generate_file_3(prevent_skip=False):
30     records_to_process = int(os.popen("wc -l
31     ↪ file2.txt").read().strip().split()[0])
32     print("Processing Dissimilarity computation for {}
33     ↪ records..".format(records_to_process))
34     # Initial call to print 0% progress
35     simulator.printProgressBar(0, records_to_process-1, prefix = 'Progress:',
36     ↪ suffix = 'Complete', length = 50)
37     with open('file3.txt', 'w') as write_file:
38         skip_n = 0
39         with open('file2.txt', 'r') as file1:
40             count = 0
41             for file1_line in file1.readlines():
42                 skip_n += 1
43                 count += 1
44                 dataset1 = parseline(file1_line)
45                 simulator.printProgressBar(count, records_to_process, prefix =
46                 ↪ 'Progress:', suffix = 'Complete', length = 50)
47                 with open('file2.txt', 'r') as file2:
48                     line_no = 0
49                     for file2_line in file2.readlines():
50                         line_no += 1
51                         if line_no > skip_n or prevent_skip:
```

```

48         dataset2 = parseline(file2_line)
49         if dataset1[2] == dataset1[2] and dataset1 !=
           ↳ dataset2 and dataset1[0] != dataset2[0]:
50             cmp_str = "{}-{}".format(dataset1[2],
           ↳ dataset2[2])
51             write_file.write("{} , {} , {} , {}
           ↳ \n".format(cmp_str, dataset1[0], dataset2[0],
           ↳ dataset1[1] - dataset2[1]))
52
53
54 if __name__ == '__main__':
55     generate_file_1()
56     generate_file_2()
57     generate_file_3()

```

IV.1 Result

```

2-2,  uzurg,  uzuru,  -1411.0
2-2,  uzurg,  uzuru,  -1508.0
2-2,  uzurg,  uzuru,  -1713.0
2-2,  uzurg,  uzuru,  -1810.0
2-2,  uzurg,  uzuru,  -1810.0
2-2,  uzurg,  uzuru,  -2111.0
2-2,  uzurg,  uzurf,  -10466.0
2-2,  uzurg,  uzurf,  -10723.0
2-2,  uzurg,  uzurf,  -10768.0
2-2,  uzurg,  uzurf,  -11025.0

```

```

:
:

```

```

2-4,  uzurg,  uzuru,  -397746.0
2-4,  uzurg,  uzuru,  -398347.0
2-4,  uzurg,  uzuru,  -398947.0
2-4,  uzurg,  uzuru,  -399546.0
2-4,  uzurg,  uzuru,  -400146.0
2-4,  uzurg,  uzuru,  -400746.0
2-4,  uzurg,  uzuru,  -401946.0
2-4,  uzurg,  uzuru,  -402546.0
2-4,  uzurg,  uzuru,  -403146.0
2-4,  uzurg,  uzuru,  -403745.0
2-4,  uzurg,  uzuru,  -404345.0

```

```

2-4,  uzurg,  uzuru,  -404945.0
2-4,  uzurg,  uzuru,  -405546.0
2-4,  uzurg,  uzuru,  -406146.0
2-4,  uzurg,  uzuru,  -406745.0

```

```

:

```

```

5-5,  uzuru,  uzurf,  -229832.0
5-5,  uzuru,  uzurf,  -230432.0
5-5,  uzuru,  uzurf,  -231032.0
5-5,  uzuru,  uzurf,  -231632.0
5-5,  uzuru,  uzurf,  -232232.0
5-5,  uzuru,  uzurf,  -232832.0
5-5,  uzuru,  uzurf,  -233431.0
5-5,  uzuru,  uzurf,  -234031.0
5-5,  uzuru,  uzurf,  -234631.0
5-5,  uzuru,  uzurf,  -235231.0

```

Attribute	Data-type	Description
taxi-id-taxi-id	string	ids of two taxis separated by hyphen
Area	geohash	hashed location value
Area	geohash	hashed location value
Dissimilarity	Float value	Float value containing difference of epoch timestamp

Table 2: Description of final result: dissimilarity between location

V. CONCLUSION

The technique geohash which we have used for data transformation gives the dissimilarity with greater precision. But this precision leads to higher computation cost which is the major drawback. So to reduce the computation cost the value of precision needs to be decreased to reduce computational cost.