

## HOMEWORK 3

### 1 Task (Visit Prediction)

**Problem 1.** Although we have built a validation set, it only consists of positive samples. For this task we also need examples of user/business pairs that weren't visited. Build such a set by randomly sampling users and businesses until you have 100,000 non-visited user/business pairs. This random sample combined with your 100,000 validation reviews now corresponds to the complete validation set for the visit prediction task. Evaluate the performance (accuracy) of the baseline model on the validation set you have built (1 mark).

**Answer:** First step is to divide 200,000 data into two parts for training and validation. Secondly, in order to generate 100,000 non-visited user/business pairs, we need a user set and a business set containing all recorded users and businesses. Then we randomly pick one user and one business and check whether this pair has ever existed in 200,000 data. If not, we then successfully generate one pair. Keep doing this until we have 100,000 non-visited pairs. Lastly, evaluate the performance (accuracy) of the baseline model on the validation set. The accuracy is 0.648705. Below is the code for this problem.

```
1 import gzip
2 from collections import defaultdict
3 import numpy as np
4 import random
5 import math
6
7 def readGz(f):
8     for l in gzip.open(f):
9         yield eval(l)
10
11 data_set = []
12 u2b = {}
13 users = set([])
14 businesses = set([])
15 for l in readGz("train.json.gz"):
16     user, business = l['userID'], l['businessID']
17     users.add(user)
18     businesses.add(business)
19     if user not in u2b:
20         u2b[user] = set([business])
21     else:
22         u2b[user].add(business);
23     data_set.append([user, business, 1])
24
25 train_set = data_set[:int(len(data_set)/2)]
```

```

26 validate_set = data_set[int(len(data_set)/2):]
27
28 # randomly create 100000 negative samples
29 users = list(users)
30 businesses = list(businesses)
31 neg_samples = []
32 while len(neg_samples) < 100000:
33     selected_user = users[int(np.random.rand()*len(users))]
34     selected_business = businesses[int(np.random.rand()*len(businesses))]
35     if selected_business not in u2b[selected_user]:
36         neg_samples.append([selected_user, selected_business, 0])
37
38 validate_set = validate_set + neg_samples
39
40 # use baseline model to train
41 businessCount = defaultdict(int)
42 totalPurchases = 0
43
44 i = 1;
45 for l in readGz("train.json.gz"):
46     if i <= 100000:
47         user,business = l['userID'],l['businessID']
48         businessCount[business] += 1
49         totalPurchases += 1
50     else:
51         break
52
53 mostPopular = [(businessCount[x], x) for x in businessCount]
54 mostPopular.sort()
55 mostPopular.reverse()
56
57 return1 = set()
58 count = 0
59 for ic, i in mostPopular:
60     count += ic
61     return1.add(i)
62     if count > totalPurchases/2: break
63
64 prediction = []
65 correct_number = 0
66 for elem in validate_set:
67     if elem[1] in return1:
68         prediction.append(1)
69         if elem[2] == 1:
70             correct_number += 1
71     else:
72         prediction.append(0)
73         if elem[2] == 0:
74             correct_number += 1
75
76 accuracy = correct_number/len(validate_set)
77
78 print(accuracy)

```

**Problem 2.** The existing 'visit prediction' baseline just returns True if the business in question is 'popular,' using a threshold of the 50th percentile of popularity ( $\text{totalVisits}/2$ ). Assuming that the 'non-visited' test examples are a random sample of user-visit pairs, is this particular threshold value the best? If not, see if you can find a better one (and report its performance), or if so, explain why it is the best (1 mark).

**Answer:** By trying different percentiles, we can get accuracy data in Table 1.

Table 1: Accuracy on the Validation Set with Different Percentiles

Percentile	1/4	1/3	1/2	2/3	3/4
Accuracy	0.596	0.618	0.648	0.652	0.639

By checking the table, we can know that setting percentile to be approximately 2/3 would be better. This ratio is very close to **golden ratio**: 0.618, and if we try this percentile, we can get a better result: 0.653435. Golden ratio wins.

**Problem 3.** Users may tend to repeatedly visit business of the same type. Build a baseline that returns 'True' if a user has visited a business of the same category before (at least one category in common), or zero otherwise (1 mark).

**Answer:** In this problem, we need to keep track of each user's visiting history to help determine whether he/she has ever visited a business of same category.

First, for each user-business pair in the training set, we add the categories of that business into two dictionaries: u2t(user to type) and b2t(business to type). By doing so, we are storing the categories that each user would visit and each business possesses.

Second, validate the model. For each user-business pair, we retrieve the categories that this user would visit and the categories that this business possess, we check if these two sets have any intersections. If there exists some intersections, then we say this user would visit this business, return 'True'.

Hence, the accuracy of this model on the validation set is: 0.650525.

Below is the code for this problem:

```

1  import gzip
2  from collections import defaultdict
3  import numpy as np
4  import random
5  import math
6
7  def readGz(f):
8      for l in gzip.open(f):
9          yield eval(l)
10
11  data_set = []
12  u2b = {}
13  users = set([])
14  businesses = set([])
15  u2t = {}
16  b2t = {}
17  i = 1
18  for l in readGz("train.json.gz"):
19      user, business, types = l['userID'], l['businessID'], l['categories']
20      users.add(user)
21      businesses.add(business)
22      if user not in u2b:
23          u2b[user] = set([business])
24      else:
25          u2b[user].add(business);
26      data_set.append([user, business, 1])
27      if i <= 100000:
28          i += 1
29          # update user -> types that they visit
30          if user not in u2t:
31              u2t[user] = set(types)
32          else:
33              for t in types:
34                  u2t[user].add(t)
35          # update business -> types that they belong to
36          if business not in b2t:
37              b2t[business] = set(types)
38          else:
39              for t in types:

```

```

40         b2t[business].add(t)
41
42
43 train_set = data_set[:int(len(data_set)/2)]
44 validate_set = data_set[int(len(data_set)/2):]
45
46 # randomly create 100000 negative samples
47 users = list(users)
48 businesses = list(businesses)
49 neg_samples = []
50 while len(neg_samples) < 100000:
51     selected_user = users[int(np.random.rand()*len(users))]
52     selected_business = businesses[int(np.random.rand()*len(businesses))]
53     if selected_business not in u2b[selected_user]:
54         neg_samples.append([selected_user, selected_business, 0])
55
56 validate_set = validate_set + neg_samples
57
58 # predict the results
59 prediction = []
60 correct_number = 0
61 for elem in validate_set:
62     if elem[0] in u2t and elem[1] in b2t:
63         types_user = u2t[elem[0]]
64         types_business = b2t[elem[1]]
65         size1 = len(types_user)
66         size2 = len(types_business)
67         types_user = types_user.union(types_business)
68         size3 = len(types_user)
69         if size1 + size2 != size3:
70             prediction.append(1)
71             if elem[2] == 1:
72                 correct_number += 1
73         else:
74             prediction.append(0)
75             if elem[2] == 0:
76                 correct_number += 1
77     else:
78         prediction.append(0)
79         if elem[2] == 0:
80             correct_number += 1
81
82 accuracy = correct_number/len(validate_set)
83
84 print(accuracy)

```

**Problem 4.** To run our model on the test set, we'll have to use the files 'pairs Visit.txt' to find the userID/businessID pairs about which we have to make predictions. Using that data, run the above model and upload your solution to Kaggle. Tell us your Kaggle user name (1 mark). If you've already uploaded a better solution to Kaggle, that's fine too!

**Answer:** My Kaggle user name is: **jih417**, as is shown in Fig. 1.

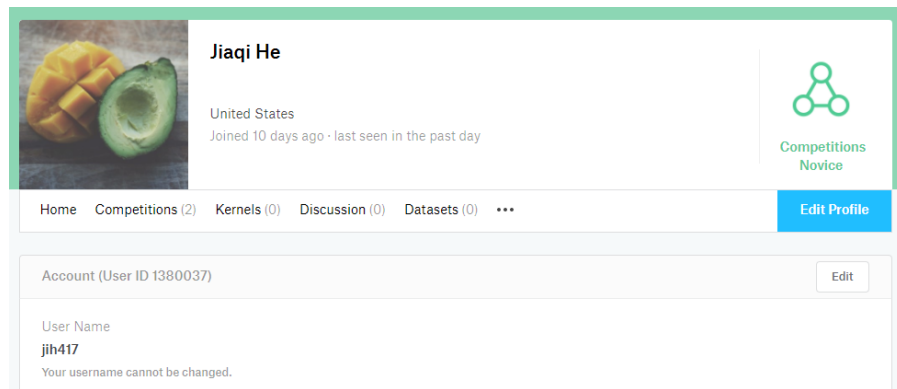


Figure 1: Kaggle profile info

If using the model from problem 3, the score is 0.66490 on Kaggle. We can definitely improve this result by building a collaborative filtering model, which will be solved in assignment 1.

Below is the code for this problem.

```

1  import gzip
2  from collections import defaultdict
3  import numpy as np
4  import random
5  import math
6
7  def readGz(f):
8      for l in gzip.open(f):
9          yield eval(l)
10
11  data_set = []
12  u2b = {}
13  users = set([])
14  businesses = set([])
15  u2t = {}
16  b2t = {}
17  for l in readGz("train.json.gz"):
18      user, business, types = l['userID'], l['businessID'], l['categories']
19      users.add(user)
20      businesses.add(business)
21      if user not in u2b:
22          u2b[user] = set([business])
23      else:

```

```

24     u2b[user].add(business);
25     # update user -> types that they visit
26     if user not in u2t:
27         u2t[user] = set(types)
28     else:
29         for t in types:
30             u2t[user].add(t)
31     # update business -> types that they belong to
32     if business not in b2t:
33         b2t[business] = set(types)
34     else:
35         for t in types:
36             b2t[business].add(t)
37     data_set.append([user, business, 1])
38
39 print (len(users))
40 print (len(businesses))
41 # predict the results
42 predictions = open("predictions_Visit.txt", 'w')
43 for l in open("pairs_Visit.txt"):
44     if l.startswith("userID"):
45         #header
46         predictions.write(l)
47         continue
48     u,i = l.strip().split('-')
49     if u in u2t and i in b2t:
50         types_user = u2t[u]
51         types_business = b2t[i]
52         size1 = len(types_user)
53         size2 = len(types_business)
54         types_user = types_user.union(types_business)
55         size3 = len(types_user)
56         if size1 + size2 != size3:
57             predictions.write(u + '-' + i + ",1\n")
58         else:
59             predictions.write(u + '-' + i + ",0\n")
60     else:
61         predictions.write(u + '-' + i + ",0\n")
62
63 predictions.close()

```

## 2 Task (Rating Prediction)

**Problem 5.** What is the performance of a trivial predictor

$$\text{rating}(\text{user}, \text{item}) = \alpha \quad (1)$$

on the validation set, and what is the value of  $\alpha$  (1 mark)?

**Answer:** This model predicts each user-item pair to be the average rating.

First, we acquire training data and validation data. On the training set, we add up all ratings and calculate the average rating, namely  $\alpha$  in the model.

For validation, we just predict  $\alpha$  ratings and calculate MSE.

Hence,  $\alpha = 4.18703$ , and the MSE on the validation set is 0.7483437444993984.

Below is the code for this problem.

```
1 import gzip
2 from collections import defaultdict
3 import math
4
5 def readGz(f):
6     for l in gzip.open(f):
7         yield eval(l)
8
9 allRatings = []
10 userRatings = defaultdict(list)
11 validate_data = []
12 i = 1
13 for l in readGz("train.json.gz"):
14     if i <= 100000:
15         user,business = l['userID'],l['businessID']
16         allRatings.append(l['rating'])
17         userRatings[user].append(l['rating'])
18     else:
19         user,business,rate = l['userID'],l['businessID'],l['rating']
20         validate_data.append([user, business, rate])
21     i+=1
22
23 globalAverage = sum(allRatings) / len(allRatings)
24
25 print(globalAverage)
26
27 MSE = 0
28 for d in validate_data:
29     MSE += (d[2]-globalAverage)**2
30 MSE /= 100000
31 print(MSE)
```



**Problem 6.** Fit a predictor of the form

$$rating(user, item) \approx \alpha + \beta_{user} - \beta_{item} \quad (2)$$

by fitting the mean and the two bias terms as described in the lecture notes. Use a regularization parameter of  $\lambda = 1$ . Report the MSE on the validation set (1 mark).

**Answer:** This model is more complex than the previous one.  
The optimization problem of this model is

$$\operatorname{argmin}_{\alpha, \beta} \sum_{u, i} (\alpha + \beta_u + \beta_i - R_{u, i})^2 + \lambda [\sum_u \beta_u^2 + \sum_i \beta_i^2] \quad (3)$$

If we differentiate this expression, we can acquire the following:

$$\alpha = \frac{\sum_{u, i \in \text{train}} (R_{u, i} - (\beta_u + \beta_i))}{N_{\text{train}}} \quad (4)$$

$$\beta_u = \frac{\sum_{i \in I_u} (R_{u, i} - (\alpha + \beta_i))}{\lambda + |I_u|} \quad (5)$$

$$\beta_i = \frac{\sum_{u \in U_i} (R_{u, i} - (\alpha + \beta_u))}{\lambda + |U_i|} \quad (6)$$

We need to repeat the updates above until convergence.

After the training, we can try our model on the validation set. The MSE on the validation set is: 0.6456944316928961.

Below is the code for this problem.

```

1  import gzip
2  from collections import defaultdict
3  import math
4
5  def readGz(f):
6      for l in gzip.open(f):
7          yield eval(l)
8
9  alpha = 0
10 beta_u = defaultdict(list)
11 beta_i = defaultdict(list)
12 users = set([])
13 businesses = set([])
14 u2b = {}
15 b2u = {}
16 train_data = []
17 validate_data = []
18 i = 1
19 for l in readGz("train.json.gz"):
20     user, business, rate = l['userID'], l['businessID'], l['rating']

```

```

21     if i <= 100000:
22         train_data.append([user, business, rate])
23         if user not in u2b:
24             u2b[user] = set([(business, rate)])
25             users.add(user)
26         else:
27             u2b[user].add((business, rate))
28         if business not in b2u:
29             b2u[business] = set([(user, rate)])
30             businesses.add(business)
31         else:
32             b2u[business].add((user, rate))
33     else:
34         validate_data.append([user, business, rate])
35     i+=1
36     beta_u[user] = 0
37     beta_i[business] = 0
38
39 users = list(users)
40 businesses = list(businesses)
41 target = businesses[0]
42
43 def iterate(lamda):
44     global alpha
45     # update alpha
46     numerator = 0
47     for d in train_data:
48         numerator += (d[2] - beta_u[d[0]] - beta_i[d[1]])
49     alpha = numerator/len(train_data)
50     # update beta_u
51     for user in users:
52         total = 0
53         for elem in u2b[user]:
54             total += (elem[1] - alpha - beta_i[elem[0]])
55         beta_u[user] = total/(lamda + len(u2b[user]))
56     # update beta_i
57     for b in businesses:
58         total = 0
59         for elem in b2u[b]:
60             total += (elem[1] - alpha - beta_u[elem[0]])
61         beta_i[b] = total/(lamda + len(b2u[b]))
62
63
64 for x in range(20):
65     iterate(1)
66
67 # validate
68 mean_sq_error = 0
69 for d in validate_data:
70     mean_sq_error += (alpha + beta_u[d[0]] + beta_i[d[1]] - d[2])**2
71 mean_sq_error = mean_sq_error/100000
72 print("MSE is equal to ", mean_sq_error)

```

**Problem 7.** Report the user and item IDs that have the largest and smallest values of  $\beta$  (1 mark).

**Answer:** By checking  $\beta_u$  and  $\beta_i$ , we can easily find the user and item with the largest and smallest values.

The user that has the largest value of  $\beta$  is 'U357799541' with  $\beta = 1.1643513853119356$ ;

The user that has the largest value of  $\beta$  is 'U417838537' with  $\beta = -2.8318415410925115$ ;

The item that has the largest value of  $\beta$  is 'B093985406' with  $\beta = 1.1719607471576514$ ;

The item that has the smallest value of  $\beta$  is 'B241777680' with  $\beta = -2.2312453963097525$ .

Below is the code for this problem.

```
76 beta_max_u = 0
77 beta_max_i = 0
78 max_u = []
79 max_i = []
80
81 beta_min_u = 0
82 beta_min_i = 0
83 min_u = []
84 min_i = []
85
86 for user in users:
87     if beta_u[user] > beta_max_u:
88         max_u = [user]
89         beta_max_u = beta_u[user]
90     elif beta_u[user] == beta_max_u:
91         max_u.append(user)
92     if beta_u[user] < beta_min_u:
93         min_u = [user]
94         beta_min_u = beta_u[user]
95     elif beta_u[user] == beta_min_u:
96         min_u.append(user)
97
98 for b in businesses:
99     if beta_i[b] > beta_max_i:
100         max_i = [b]
101         beta_max_i = beta_i[b]
102     elif beta_i[b] == beta_max_i:
103         max_i.append(b)
104     if beta_i[b] < beta_min_i:
105         min_i = [b]
106         beta_min_i = beta_i[b]
107     elif beta_i[b] == beta_min_i:
108         min_i.append(b)
109
110 print(max_u, beta_max_u)
111 print(max_i, beta_max_i)
112 print(min_u, beta_min_u)
113 print(min_i, beta_min_i)
```

**Problem 8.** Find a better value of  $\lambda$  using your validation set. Report the value you chose, its MSE, and upload your solution to Kaggle by running it on the test data (1 mark).

**Answer:** By trying different  $\lambda$ , we can obtain data in Fig. 2.

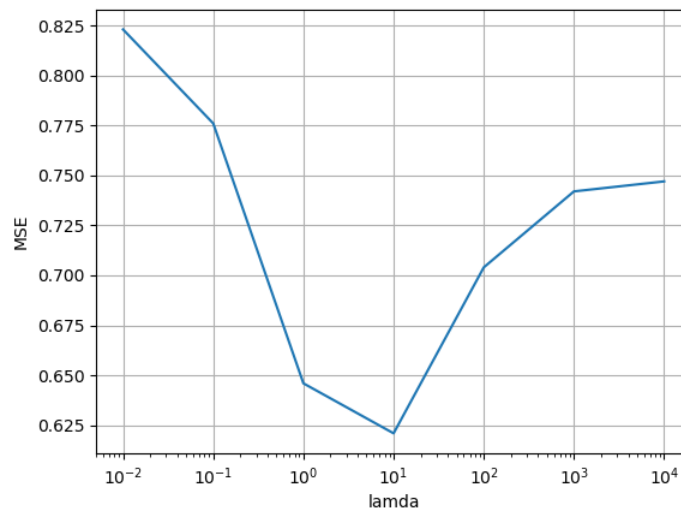


Figure 2: Kaggle profile info

Hence, it is a good choice to select  $\lambda = 10$ . This time the MSE on the validation set is 0.6214475097509096.

Below is the code that has already been modified to make predictions on test file.

```

1 import gzip
2 from collections import defaultdict
3 import math
4
5 def readGz(f):
6     for l in gzip.open(f):
7         yield eval(l)
8
9 alpha = 0
10 beta_u = defaultdict(list)
11 beta_i = defaultdict(list)
12 users = set([])
13 businesses = set([])
14 u2b = {}
15 b2u = {}
16 train_data = []
17 i = 1
18 for l in readGz("train.json.gz"):

```

```

19     user,business,rate = l['userID'],l['businessID'],l['rating']
20     beta_u[user] = 0
21     beta_i[business] = 0
22     train_data.append([user, business, rate])
23     if user not in u2b:
24         u2b[user] = set([(business,rate)])
25         users.add(user)
26     else:
27         u2b[user].add((business, rate))
28     if business not in b2u:
29         b2u[business] = set([(user,rate)])
30         businesses.add(business)
31     else:
32         b2u[business].add((user,rate))
33
34
35 users = list(users)
36 businesses = list(businesses)
37
38 def iterate(lamda):
39     global alpha
40     # update alpha
41     numerator = 0
42     for d in train_data:
43         numerator += (d[2] - beta_u[d[0]] - beta_i[d[1]])
44     alpha = numerator/len(train_data)
45     # update beta_u
46     for user in users:
47         total = 0
48         for elem in u2b[user]:
49             total += (elem[1] - alpha - beta_i[elem[0]])
50         beta_u[user] = total/(lamda + len(u2b[user]))
51     # update beta_i
52     for b in businesses:
53         total = 0
54         for elem in b2u[b]:
55             total += (elem[1] - alpha - beta_u[elem[0]])
56         beta_i[b] = total/(lamda + len(b2u[b]))
57         # if b == target:
58         #     print(beta_i[target])
59
60 for x in range(30):
61     iterate(10)
62 # mean_sq_error = 0
63 # for d in validate_data:
64 #     mean_sq_error += (alpha + beta_u[d[0]] + beta_i[d[1]] - d[2])**2
65 # mean_sq_error = math.sqrt(mean_sq_error)
66 # mean_sq_error = mean_sq_error/100000
67 # print(mean_sq_error)
68
69 idx = 1
70 predictions = open("predictions_Rating.txt", 'w')
71 for l in open("pairs_Rating.txt"):
72     if l.startswith("userID"):
73         predictions.write(l)
74         continue
75     u,i = l.strip().split('-')

```

```

76     if u not in users:
77         if i not in businesses:
78             predictions.write(u + '-' + i + ',' + str(alpha) + '\n')
79         else:
80             predictions.write(u + '-' + i + ',' + str(alpha + beta_i[i]) + '\n')
81     else:
82         if i not in businesses:
83             predictions.write(u + '-' + i + ',' + str(alpha + beta_u[u]) + '\n')
84         else:
85             predictions.write(u + '-' + i + ',' + str(alpha + beta_u[u] + beta_i[i]) +
                                '\n')
86         # if idx < 10:
87         #     print(alpha, " ", beta_u[u], " ", beta_i[i])
88         #     idx+=1
89
90 predictions.close()

```

*Submitted by He, Jiaqi on November 10, 2017.*