

Tasks (Read prediction)

Since we don't have access to the test labels, we'll need to simulate validation/test sets of our own. So, let's split the training data ('train Interactions.csv.gz') as follows:

- Reviews 1-190,000 for training
- Reviews 190,001-200,000 for validation
- Upload to Kaggle for testing only when you have a good model on the validation set. This will save you time (since Kaggle can take several minutes to return results), and prevent you from exceeding your daily submission limit.

```
1  # ===== Environment Setup =====
2  # ignore the warnings
3  import warnings
4  warnings.filterwarnings("ignore")
5
6  # import some packages
7  import gzip
8  import random
9  import pandas as pd
10 from tqdm.notebook import tqdm
11 from collections import defaultdict
12
13 # ===== Load and Split data =====
14 def splitDataset(datapath):
15     """
16     Split the training data ('train Interactions.csv.gz') as
17     follows:
18     (1) Reviews 1-190,000 for training
19     (2) Reviews 190,001-200,000 for validation
20     :param datapath:
21     :return: train, valid as dataframe
22     """
23     f = gzip.open(datapath, 'rt')
24     data = pd.read_csv(f)
25     train, valid = data[:190001], data[190001:]
26     return data, train, valid
27
28 data, train, valid =
29 splitDataset("./data/train_Interactions.csv.gz")
```

Question 1

Although we have built a validation set, it only consists of positive samples. For this task we also need examples of user/item pairs that weren't read. For each entry (user,book) in the validation set, sample a negative entry by randomly choosing a book that user hasn't read. Evaluate the performance (accuracy) of the baseline model on the validation set you have built (1 mark).

```
1 def sampleNegative(data, train, valid):
2     """
3     For each entry (user,book) in the validation set, sample a
4     negative entry
5         by randomly choosing a book that user hasn't read.
6     :param train:
7     :param valid:
8     :return:
9     """
10    valid['read'] = 1
11    NegValid = valid
12
13    userBookDict = {}
14    print("Preprocessing Data userBookDict ...")
15    for index, row in tqdm(data.iterrows()):
16        if row['userID'] not in userBookDict:
17            userBookDict[row['userID']] = {row['bookID']}
18        else:
19            userBookDict[row['userID']].add(row['bookID'])
20
21    for index, row in tqdm(valid.iterrows()):
22        samNegBookID = random.sample(set(data['bookID']) -
23        userBookDict[row['userID']], 1)[0]
24        NegValid = NegValid.append({'userID': row['userID'],
25        'bookID': samNegBookID, 'rating': 0, 'read': 0}, ignore_index=True)
26
27    return NegValid, userBookDict
```

```
1 def baselineOnValidation():
2     random.seed(5583)
3     data, train, valid =
4     splitDataset("./data/train_Interactions.csv.gz")
5     print("Sampling Negative samples ...")
6     valid, _ = sampleNegative(data, train, valid)
7
8     bookCount = defaultdict(int)
9     totalRead = 0
10
11    print('Training ...')
```

```

11     for index, row in tqdm(train.iterrows()):
12         bookCount[row['bookID']] += 1
13         totalRead += 1
14
15     mostPopular = [(bookCount[x], x) for x in bookCount]
16     mostPopular.sort()
17     mostPopular.reverse()
18
19     return1 = set()
20     count = 0
21     for ic, i in mostPopular:
22         count += ic
23         return1.add(i)
24         if count > totalRead / 2:
25             break
26
27     print('Evaluating ...')
28     correct = 0
29     for index, row in tqdm(valid.iterrows()):
30         if row['bookID'] in return1:
31             correct += (row['read'] != 0)
32         else:
33             correct += (row['read'] == 0)
34
35     print('Accuracy on Validation set is %.3f' %
36           (correct/len(valid)))
37 baselineOnValidation()

```

```

1 | Sampling Negative samples ...
2 | Preprocessing Data userBookDict ...
3 | Training ...
4 | Evaluating ...

```

```

1 | Accuracy on Validation set is 0.646

```

Question 2

The existing ‘read prediction’ baseline just returns True if the item in question is ‘popular,’ using a threshold of the 50th percentile of popularity ($\text{totalRead}/2$). Assuming that the ‘non-read’ test examples are a random sample of user-book pairs, this threshold may not be the best one. See if you can find a better threshold and report its performance on your validation set (1 mark).

```

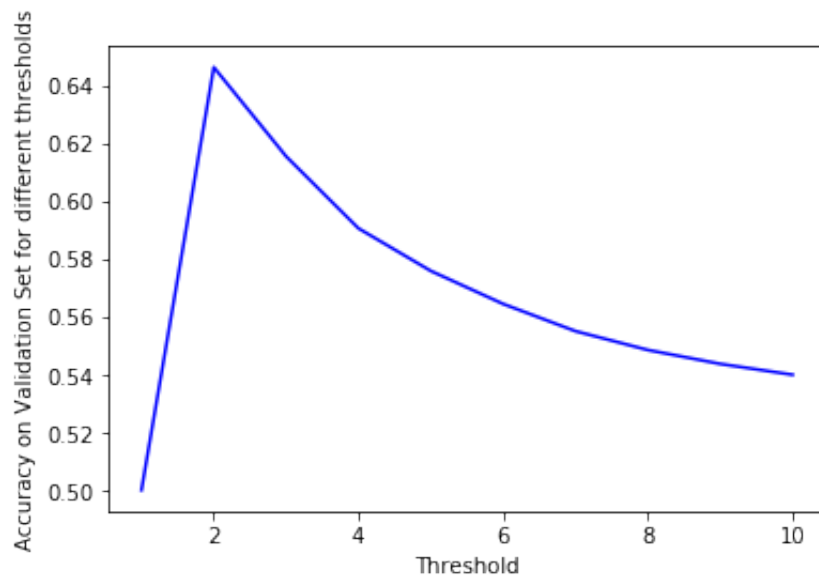
1 import matplotlib.pyplot as plt
2
3 def baselineOnValidation():
4     random.seed(5583)
5     data, train, valid =
splitDataset("./data/train_Interactions.csv.gz")
6     print("Sampling Negative samples ...")
7     valid, _ = sampleNegative(data, train, valid)
8
9     acc = []
10    thresholds = list(range(1, 11))
11
12    bookCount = defaultdict(int)
13    totalRead = 0
14
15    print('Training ...')
16    for index, row in train.iterrows():
17        bookCount[row['bookID']] += 1
18        totalRead += 1
19
20    mostPopular = [(bookCount[x], x) for x in bookCount]
21    mostPopular.sort()
22    mostPopular.reverse()
23
24    for threshold in thresholds:
25        return1 = set()
26        count = 0
27        for ic, i in mostPopular:
28            count += ic
29            return1.add(i)
30            if count > totalRead / threshold:
31                break
32
33        print('Evaluating on threshold %d ...' % threshold)
34        correct = 0
35        for index, row in valid.iterrows():
36            if row['bookID'] in return1:
37                correct += (row['read'] != 0)
38            else:
39                correct += (row['read'] == 0)
40
41        acc.append(correct/len(valid))
42
43    plt.plot(thresholds, acc, 'b-')
44    plt.xlabel('Threshold')
45    plt.ylabel('Accuracy on Validation Set for different
thresholds')
46    plt.show()
47

```

```

1 Sampling Negative samples ...
2 Preprocessing Data userBookDict ...
3
4 Training ...
5 Evaluating on threshold 1 ...
6 Evaluating on threshold 2 ...
7 Evaluating on threshold 3 ...
8 Evaluating on threshold 4 ...
9 Evaluating on threshold 5 ...
10 Evaluating on threshold 6 ...
11 Evaluating on threshold 7 ...
12 Evaluating on threshold 8 ...
13 Evaluating on threshold 9 ...
14 Evaluating on threshold 10 ...

```



As we can see, the best accuracy on validation set was achieved at threshold around 1.9. And the accuracy is around 0.647.

Question 3

A stronger baseline than the one provided might make use of the Jaccard similarity (or another similarity metric). Given a pair (u, b) in the validation set, consider all training items b' that user u has read. For each, compute the Jaccard similarity between b and b' , i.e., users (in the training set) who have read b and users who have read b' . Predict as 'read' if the maximum of these Jaccard

similarities exceeds a threshold (you may choose the threshold that works best). Report the performance on your validation set (1 mark).

```
1 random.seed(5583)
2 data, train, valid =
  splitDataset("./data/train_Interactions.csv.gz")
3 print("Sampling Negative samples ...")
4 valid, _ = sampleNegative(data, train, valid)
5
6 userReadBook, BookwasRead = {}, {}
7 for index, row in tqdm(train.iterrows()):
8     if row['userID'] not in userReadBook:
9         userReadBook[row['userID']] = {row['bookID']}
10    else:
11        userReadBook[row['userID']].add(row['bookID'])
12    if row['bookID'] not in BookwasRead:
13        BookwasRead[row['bookID']] = {row['userID']}
14    else:
15        BookwasRead[row['bookID']].add(row['userID'])
```

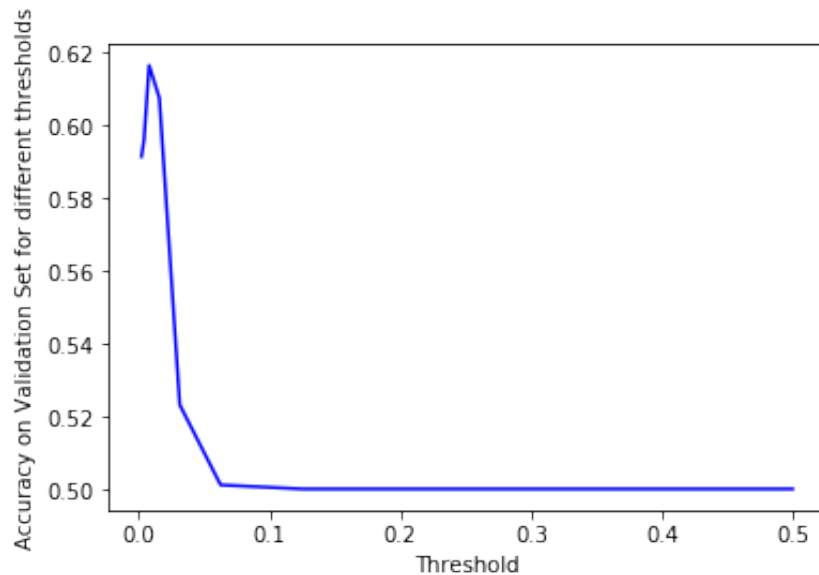
```
1 Sampling Negative samples ...
2 Preprocessing Data userBookDict ...
```

```
1 def Jaccard(s1, s2):
2     numer = len(s1.intersection(s2))
3     denom = len(s1.union(s2))
4     return numer / denom
5
6 thresholds = [1/2**i for i in range(1, 10)]
7 acc = []
8 for threshold in thresholds:
9     print('Evaluating on threshold %.3f ...' % threshold)
10    correct = 0
11    for index, row in tqdm(valid.iterrows()):
12        userReads = userReadBook[row['userID']]
13        jac = []
14        for book in userReads:
15            if row['bookID'] not in BookwasRead:
16                jac.append(0)
17            else:
18                jac.append(Jaccard(BookwasRead[row['bookID']],
19                                   BookwasRead[book]))
20
21        if max(jac) > threshold:
22            correct += (row['read'] != 0)
23        else:
24            correct += (row['read'] == 0)
```

```

24
25     acc.append(correct/len(valid))
26
27 plt.plot(thresholds, acc, 'b-')
28 plt.xlabel('Threshold')
29 plt.ylabel('Accuracy on Validation Set for different thresholds')
30 plt.show()

```



```

1 print('threshold %.3f, accuracy %.3f' %
      (thresholds[acc.index(max(acc))], max(acc)))

```

```

1 threshold 0.008, accuracy 0.616

```

Question 4

Improve the above predictor by incorporating both a Jaccard-based threshold and a popularity based threshold. Report the performance on your validation set (1 mark).

First, we simply Use a "and" operation for ensemble model.

```

1 # preprocessing
2 random.seed(5583)
3 data, train, valid =
  splitDataset("./data/train_Interactions.csv.gz")
4 print("Sampling Negative samples ...")
5 valid, _ = sampleNegative(data, train, valid)

```



```

33         if max(jac) > threshold_jac and row['bookID'] in return1:
34             correct += (row['read'] != 0)
35         else:
36             correct += (row['read'] == 0)
37
38     return correct/len(valid)
39
40 print("Accuracy is %.3f" % ensemble())

```

```

1 Training ...
2 Evaluating ...

```

```

1 Accuracy is 0.650

```

If we use the best thresholds from the single popularity and Jaccard method, the accuracy is 0.650. However, since in the ensemble model, we simply do a "and" operation, we should properly **lower the single thresholds**.

Tune on Popularity Threshold

```

1 threshold_pops = [i/10 for i in range(10, 20)]
2 acc = []
3 for threshold_pop in threshold_pops:
4     acc.append(ensemble(threshold_pop, threshold_jac = 0.008))

```

```

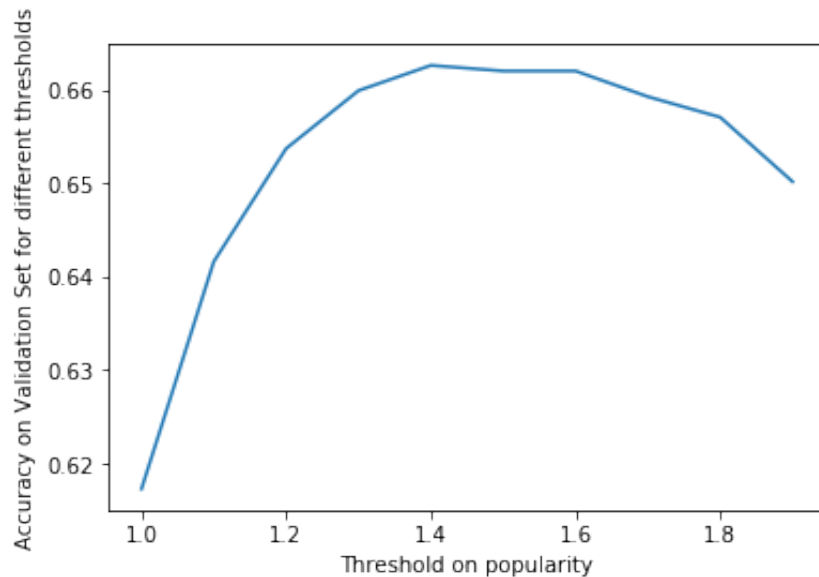
1 Training ...
2 Evaluating ...

```

```

1 plt.plot(threshold_pops, acc)
2 plt.xlabel('Threshold on popularity')
3 plt.ylabel('Accuracy on Validation Set for different thresholds')
4 plt.show()

```



The best accuracy was achieved at threshold 1.4 as 0.663.

```
1 max(acc)
```

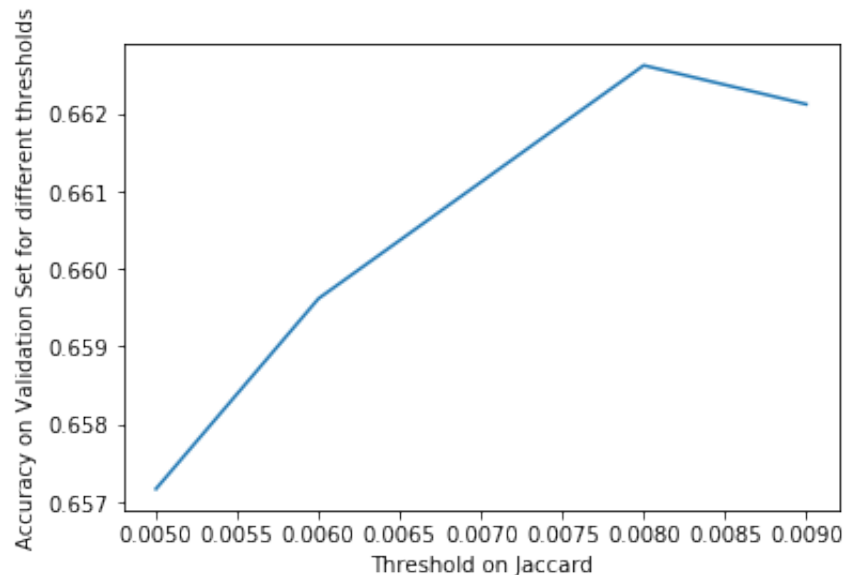
```
1 0.6626162616261626
```

Tune on Jaccard Threshold

```
1 threshold_jacs = [i/10000 for i in range(50, 100, 10)]
2 acc = []
3 for threshold_jac in threshold_jacs:
4     acc.append(ensemble(threshold_pop=1.4,
5                          threshold_jac=threshold_jac))
```

```
1 Training ...
2 Evaluating ...
```

```
1 plt.plot(threshold_jacs, acc)
2 plt.xlabel('Threshold on Jaccard')
3 plt.ylabel('Accuracy on Validation Set for different thresholds')
4 plt.show()
```



According to the previous experiments, the best results 0.663 are achieved at the popularity threshold = 1.4, Jaccard threshold = 0.008.

Question 5

To run our model on the test set, we'll have to use the files 'pairs Read.txt' to find the reviewerID/itemID 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!

```

1 threshold_pop = 1.4
2 threshold_jac = 0.008
3
4 bookCount = defaultdict(int)
5 totalRead = 0
6
7 print('Training ...')
8 for index, row in train.iterrows():
9     bookCount[row['bookID']] += 1
10    totalRead += 1
11
12 mostPopular = [(bookCount[x], x) for x in bookCount]
13 mostPopular.sort()
14 mostPopular.reverse()
15
16 return1 = set()
17 count = 0
18 for ic, i in mostPopular:

```

```

19     count += ic
20     return1.add(i)
21     if count > totalRead / threshold_pop:
22         break
23
24 predictions = open("predictions_Read.txt", 'w')
25 for l in open("./data/pairs_Read.txt"):
26     if l.startswith("userID"):
27         # header
28         predictions.write(l)
29         continue
30     u, b = l.strip().split('-')
31     userReads = userReadBook[u]
32     jac = []
33     for book in userReads:
34         if b not in BookwasRead:
35             jac.append(0)
36         else:
37             jac.append(Jaccard(BookwasRead[b], BookwasRead[book]))
38
39     if max(jac) > threshold_jac and b in return1:
40         predictions.write(u + '-' + b + ",1\n")
41     else:
42         predictions.write(u + '-' + b + ",0\n")
43
44 predictions.close()

```

Kaggle

- User name: shihanran
- Display name: Shihan Ran
- Email: sran@ucsd.edu
- Public Leaderboard Score: 0.67466

Tasks (Rating prediction)

Let's start by building our training/validation sets much as we did for the first task. This time building a validation set is more straightforward: you can simply use part of the data for validation, and do not need to randomly sample non-read users/books.

Question 9

Fit a predictor of the form $\text{rating}(\text{user}, \text{item}) \simeq \alpha + \beta_{\text{user}} + \beta_{\text{item}}$

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).

```
1 import scipy
2 import scipy.optimize
3 import numpy as np
4
5 def prediction(user, item):
6     return alpha + userBiases[user] + itemBiases[item]
7
8 def MSE(predictions, labels):
9     differences = [(x-y)**2 for x,y in zip(predictions,labels)]
10    return sum(differences) / len(differences)
11
12 def unpack(theta):
13     global alpha
14     global userBiases
15     global itemBiases
16     alpha = theta[0]
17     userBiases = dict(zip(users, theta[1:nUsers+1]))
18     itemBiases = dict(zip(items, theta[1+nUsers:]))
19
20 def cost(theta, labels, lamb):
21     unpack(theta)
22     predictions = [prediction(d['userID'], d['bookID']) for index,
23 d in train.iterrows()]
24     cost = MSE(predictions, labels)
25     print("MSE = " + str(cost))
26     for u in userBiases:
27         cost += lamb*userBiases[u]**2
28     for i in itemBiases:
29         cost += lamb*itemBiases[i]**2
30     return cost
31
32 def derivative(theta, labels, lamb):
33     unpack(theta)
34     N = len(train)
35     dalpha = 0
36     dUserBiases = defaultdict(float)
37     dItemBiases = defaultdict(float)
38     for index, d in train.iterrows():
39         u,i = d['userID'], d['bookID']
40         pred = prediction(u, i)
41         diff = pred - d['rating']
42         dalpha += 2/N*diff
43         dUserBiases[u] += 2/N*diff
44         dItemBiases[i] += 2/N*diff
```

```

44     for u in userBiases:
45         dUserBiases[u] += 2*lamb*userBiases[u]
46     for i in itemBiases:
47         dItemBiases[i] += 2*lamb*itemBiases[i]
48     dtheta = [dalpha] + [dUserBiases[u] for u in users] +
[dItemBiases[i] for i in items]
49     return np.array(dtheta)

```

Training on Lambda = 1

```

1  lamb = 1
2  data, train, valid =
splitDataset("./data/train_Interactions.csv.gz")
3
4  ratingMean = train['rating'].mean()
5  alpha = ratingMean
6
7  labels = train['rating']
8
9  userBiases = defaultdict(float)
10 itemBiases = defaultdict(float)
11
12 users = list(set(train['userID']))
13 items = list(set(train['bookID']))
14 nUsers = len(users)
15 nItems = len(items)
16
17 scipy.optimize.fmin_l_bfgs_b(cost, [alpha] + [0.0]*(nUsers+nItems),
18                             derivative, args = (labels, lamb))

```

```

1  MSE = 1.4735439815452338
2  MSE = 1.4560896560389294
3  MSE = 1.4733864367918064
4  MSE = 1.473386434421149

```

```

1  (array([ 3.89708841e+00,  9.39450287e-05,  3.06462740e-05, ...,
2         -1.74937858e-05,  3.89145254e-05, -1.19408114e-04]),
3   1.473465196550704,
4   {'grad': array([-2.47029336e-08, -3.42086976e-09, -1.74216486e-09,
5   ...,
6         -7.99313594e-09, -8.92116090e-09, -2.32554776e-08]),
7    'task': b'CONVERGENCE: NORM_OF_PROJECTED_GRADIENT_<=_PGTOL',
8    'funcalls': 4,
9    'nit': 2,
10   'warnflag': 0})

```

Validating

```

1 predictions = []
2 for index, d in valid.iterrows():
3     u, i = d['userID'], d['bookID']
4     if u in userBiases and i in itemBiases:
5         predictions.append(prediction(u, i))
6     else:
7         predictions.append(0)
8
9 print("MSE %.3f" % MSE(predictions, valid['rating']))

```

```

1 MSE 1.491

```

Question 10

Report the user and book IDs that have the largest and smallest values of β (1 mark).

```

1 print("max user: %s , max value: %f" % (max(userBiases,
2     key=userBiases.get), max(userBiases.values()))
3 print("max book: %s , max value: %f" % (max(itemBiases,
4     key=itemBiases.get), max(itemBiases.values()))
5 print("min user: %s , min value: %f" % (min(userBiases,
6     key=userBiases.get), min(userBiases.values()))
7 print("min book: %s , min value: %f" % (min(itemBiases,
8     key=itemBiases.get), min(itemBiases.values()))

```

```

1 max user: u92864068 , max value: 0.000404
2 max book: b76915592 , max value: 0.000829
3 min user: u11591742 , min value: -0.001580
4 min book: b57299824 , min value: -0.000272

```

Question 11

Find a better value of λ 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).

Tune on Lambda

Training

```

1  lamb = 10**(-5)
2
3  def training(lamb)
4      scipy.optimize.fmin_l_bfgs_b(cost, [alpha] + [0.0]*
      (nUsers+nItems),
5                                     derivative, args = (labels, lamb))
6
7  training(lamb)

```

We tuned on a few values of *lambda* and got the best result at $\lambda = 10^{(-5)}$.

Validating

```

1  predictions = []
2  for index, d in valid.iterrows():
3      u, i = d['userID'], d['bookID']
4      if u in userBiases and i in itemBiases:
5          predictions.append(prediction(u, i))
6      else:
7          predictions.append(0)
8
9  print("MSE %.3f" % MSE(predictions, valid['rating']))

```

```

1  MSE 1.111

```

Predicting on Test

```

1  predictions = open("predictions_Rating.txt", 'w')
2  for l in open("./data/pairs_Rating.txt"):
3      if l.startswith("userID"):
4          #header
5          predictions.write(l)
6          continue
7      u, i = l.strip().split('-')
8      if u in userBiases and i in itemBiases:
9          predictions.write(u + '-' + i + ',' + str(prediction(u, i))
10 + '\n')
11      else:
12          predictions.write(u + '-' + i + ',' + str(0) + '\n')
13  predictions.close()

```

At last, we got $MSE = 1.14702$ on test set.