

In [1]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip

def readGz(path):
    for l in gzip.open(path, 'rt'):
        yield eval(l)

def readCSV(path):
    f = gzip.open(path, 'rt')
    f.readline()
    for l in f:
        yield l.strip().split(',')

def findBook(user, userReadBook, bookAllUser):
    l3 = [x for x in list(bookAllUser) if x not in userReadBook[user]]
    proxy = random.choice(l3)
    return proxy
#cut the train set and init valid set
f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")

header = f.readline()
header = header.strip().split(',')

datatrain = []
datavalid = []
count=0
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < 190000 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(set)
bookAllUser = defaultdict(set)

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    userReadBook[user].add(book)
    bookAllUser[book].add(user)

#create the new valid set
i=0
for d in datavalid:
    if i<10000:
        dd = dict(zip(header, fields))
        dd['userID'] = d['userID']
        dd['bookID'] = findBook(d['userID'],userReadBook,bookAllUser)
        dd['rating'] = 0
        datavalid.append(dd)
        i=i+1
    else:
        break

#use train set to train the model
bookCount = defaultdict(int)
totalRead = 0
for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    bookCount[book] += 1
    totalRead += 1
```

```

mostPopular = [(bookCount[x], x) for x in bookCount]
mostPopular.sort()
mostPopular.reverse()

return1 = set()
count = 0
for ic, i in mostPopular:
    count += ic
    return1.add(i)
    if count > totalRead/2: break

#get the return1 set
count = 0
prediction = []
for d in datavalid:
    user,book,r =d['userID'],d['bookID'],d['rating']
    if book in return1 :
        prediction.append(1)
    else:
        prediction.append(0)
    count=count+1

#cal the accuracy
count =0
Tcount=0
for d in datavalid:
    if prediction[count] >0 and int(d['rating'])>0:
        Tcount+=1
    if prediction[count] ==0 and int(d['rating'])==0:
        Tcount+=1
    count+=1
accuracy = Tcount/len(prediction)

print("accuracy of valid set is ",accuracy)

```

accuracy of valid set is 0.6526

In [ ]:

In [5]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict # Dictionaries with default values
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip

def readGz(path):
    for l in gzip.open(path, 'rt'):
        yield eval(l)

def readCSV(path):
    f = gzip.open(path, 'rt')
    f.readline()
    for l in f:
        yield l.strip().split(',')

def findBook(user, userReadBook, bookAllUser):
    l3 = [x for x in list(bookAllUser) if x not in userReadBook[user]]
    proxy = random.choice(l3)
    return proxy

f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")

header = f.readline()
header = header.strip().split(',')

datatrain = []
datavalid = []
count=0
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < 190000 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(set)
bookAllUser = defaultdict(set)

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    userReadBook[user].add(book)
    bookAllUser[book].add(user)

i=0
for d in datavalid:
    if i<10000:
        dd = dict(zip(header, fields))
        dd['userID'] = d['userID']
        dd['bookID'] = findBook(d['userID'],userReadBook,bookAllUser)
        dd['rating'] = 0
        datavalid.append(dd)
        i=i+1
    else:
        break

bookCount = defaultdict(int)
totalRead = 0
for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
```

```

        bookCount[book] += 1
        totalRead += 1

mostPopular = [(bookCount[x], x) for x in bookCount]
mostPopular.sort()
mostPopular.reverse()

return1 = set()
count = 0
threshold = 1.8
for ic, i in mostPopular:
    count += ic
    return1.add(i)
    if count > totalRead/threshold: break

prediction = []
for d in datavalid:
    user, book, r = d['userID'], d['bookID'], d['rating']
    if book in return1 :
        prediction.append(1)
    else:
        prediction.append(0)

count = 0
Tcount = 0
for d in datavalid:
    if prediction[count] > 0 and int(d['rating']) > 0:
        Tcount += 1
    if prediction[count] == 0 and int(d['rating']) == 0:
        Tcount += 1
    count += 1
accuracy = Tcount/len(prediction)
print("the threshold is ", 1/1.8, "the accuracy is ", accuracy)

```

the threshold is 0.5263157894736842 the accuracy is 0.6543

In [ ]:

In [1]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict # Dictionaries with default values
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip

def readGz(path):
    for l in gzip.open(path, 'rt'):
        yield eval(l)

def readCSV(path):
    f = gzip.open(path, 'rt')
    f.readline()
    for l in f:
        yield l.strip().split(',')

def findBook(user, userReadBook, bookAllUser):
    l3 = [x for x in list(bookAllUser) if x not in userReadBook[user]]
    proxy = random.choice(l3)
    return proxy

def Jaccard(book1, book2, bookAllUser):
    s1 = bookAllUser[book1]
    s2 = bookAllUser[book2]
    numer = len(s1.intersection(s2))
    denom = len(s1.union(s2))
    return numer / denom

f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")

header = f.readline()
header = header.strip().split(',')

datatrain = []
datavalid = []
count=0
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < 190000 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(set)
bookAllUser = defaultdict(set)

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    userReadBook[user].add(book)
    bookAllUser[book].add(user)

i=0
for d in datavalid:
    if i<10000:
        dd = dict(zip(header, fields))
        dd['userID'] = d['userID']
        dd['bookID'] = findBook(d['userID'],userReadBook,bookAllUser)
        dd['rating'] = 0
        datavalid.append(dd)
        i=i+1
    else:
        break

prediction =[]
threshold =0.006
```

```
#cal the Jac by compare two books' users set's similiarity
```

```
for d in datavalid:
    user,book,r =d['userID'],d['bookID'],d['rating']
    flag =0
    for b in userReadBook[user] :
        similarJ = Jaccard(b,book,bookAllUser)
        if similarJ > threshold:
            flag =1
            break
    prediction.append(flag)

count =0
Tcount=0
for d in datavalid:
    if prediction[count] >0 and int(d['rating'])>0:
        Tcount+=1
    if prediction[count] ==0 and int(d['rating'])==0:
        Tcount+=1
    count+=1
accuracy = Tcount/len(prediction)
print(accuracy)
```

0.5833

In [ ]:

In [1]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict # Dictionaries with default values
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip

def findBook(user, userReadBook, bookAllUser):
    l3 = [x for x in list(bookAllUser) if x not in userReadBook[user]]
    proxy = random.choice(l3)
    return proxy
def Jaccard(book1, book2, bookAllUser):
    s1 = bookAllUser[book1]
    s2 = bookAllUser[book2]
    numer = len(s1.intersection(s2))
    denom = len(s1.union(s2))
    return numer / denom

f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")

header = f.readline()
header = header.strip().split(',')

datatrain = []
datavalid = []
count=0
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < 190000 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(set)
bookAllUser = defaultdict(set)

for d in datatrain:
    user, book, r = d['userID'], d['bookID'], d['rating']
    userReadBook[user].add(book)
    bookAllUser[book].add(user)

i=0
for d in datavalid:
    if i<10000:
        dd = dict(zip(header, fields))
        dd['userID'] = d['userID']
        dd['bookID'] = findBook(d['userID'], userReadBook, bookAllUser)
        dd['rating'] = 0
        datavalid.append(dd)
        i=i+1
    else:
        break

bookCount = defaultdict(int)
totalRead = 0
for d in datatrain:
    user, book, r = d['userID'], d['bookID'], d['rating']
    bookCount[book] += 1
    totalRead += 1

mostPopular = [(bookCount[x], x) for x in bookCount]
```

```

mostPopular.sort()
mostPopular.reverse()

return1 = set()
count = 0

for ic, i in mostPopular:
    count += ic
    return1.add(i)
    if count > totalRead/1.7: break

prediction = []
thresholdJ = 0.006
#when the book is similar and it is popular, it is true
for d in datavalid:
    user, book, r = d['userID'], d['bookID'], d['rating']
    flag = 0
    for b in userReadBook[user] :
        similarJ = Jaccard(b, book, bookAllUser)
        if similarJ > thresholdJ:
            if book in return1:
                flag = 1
                break
    prediction.append(flag)

count = 0
Tcount = 0
for d in datavalid:
    if prediction[count] > 0 and int(d['rating']) > 0:
        Tcount += 1
    if prediction[count] == 0 and int(d['rating']) == 0:
        Tcount += 1
    count += 1
accuracy = Tcount/len(prediction)
print(accuracy)

```

0.65985

In [ ]:



In [ ]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict # Dictionaries with default values
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip

def readGz(path):
    for l in gzip.open(path, 'rt'):
        yield eval(l)

def readCSV(path):
    f = gzip.open(path, 'rt')
    f.readline()
    for l in f:
        yield l.strip().split(',')

def findBook(user, userReadBook, bookAllUser):
    l3 = [x for x in list(bookAllUser) if x not in userReadBook[user]]
    proxy = random.choice(l3)
    return proxy

def Jaccard(book1, book2, bookAllUser):
    s1 = bookAllUser[book1]
    s2 = bookAllUser[book2]
    numer = len(s1.intersection(s2))
    denom = len(s1.union(s2))
    return numer / denom

f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")

header = f.readline()
header = header.strip().split(',')

datatrain = []
datavalid = []
count=0
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < 190000 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(set)
bookAllUser = defaultdict(set)

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    userReadBook[user].add(book)
    bookAllUser[book].add(user)

i=0
for d in datavalid:
    if i<10000:
        dd = dict(zip(header, fields))
        dd['userID'] = d['userID']
        dd['bookID'] = findBook(d['userID'],userReadBook,bookAllUser)
        dd['rating'] = 0
        datavalid.append(dd)
        i=i+1
    else:
        break

bookCount = defaultdict(int)
totalRead = 0
```

```

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    bookCount[book] += 1
    totalRead += 1

mostPopular = [(bookCount[x], x) for x in bookCount]
mostPopular.sort()
mostPopular.reverse()

return1 = set()
count = 0
for ic, i in mostPopular:
    count += ic
    return1.add(i)
    if count > totalRead/1.7: break

predictions = open("predictions_Read.txt", 'w')
thresholdJ=0.006
for l in open("pairs_Read.txt"):
    if l.startswith("userID"):
        #header
        predictions.write(l)
        continue
    user,book = l.strip().split('-')
    flag='0'
    for b in userReadBook[user] :
        similarJ = Jaccard(b,book,bookAllUser)
        if similarJ > thresholdJ:
            if book in return1 :
                flag='1'
                break
    predictions.write(user + '-' + book + "," +flag+"\n")
predictions.close()
#my kaggle name is HumphreySD, and my score is 0.67200, my email address is haz013@eng.ucsd.edu

```

In [1]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict # Dictionaries with default values
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip

def Jaccard(book1,book2,bookAllUser):
    s1=[]
    for d in bookAllUser[book1]:
        s1.append(d['userID'])
    s2=[]
    for d in bookAllUser[book2]:
        s2.append(d['userID'])
    s1=set(s1)
    s2=set(s2)
    numer = len(s1.intersection(s2))
    denom = len(s1.union(s2))
    return numer / denom

f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")
header = f.readline()
header = header.strip().split(',')
datatrain = []
datavalid = []
count=0
lenAll = 200000
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < lenAll*0.95 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(list)
bookAllUser = defaultdict(list)

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    userReadBook[user].append(d)
    bookAllUser[book].append(d)

ratingMean = sum([float(d['rating']) for d in datatrain]) / len(datatrain)

N = len(datatrain)
nUsers = len(userReadBook)
nBooks = len(bookAllUser)
users = list(userReadBook)
books = list(bookAllUser)

alpha = ratingMean

userBiases = defaultdict(float)
bookBiases = defaultdict(float)

def prediction(user, book):
    return alpha + userBiases[user] + bookBiases[book]

def predictRating(user,item):
    ratings = []
    similarities = []
    for d in userReadBook[user]:
        i2 = d['bookID']
        if i2 == item: continue
```

```

        ratings.append(d['rating'])
        similarities.append(Jaccard(item,i2,bookAllUser))
    if (sum(similarities) > 0):
        weightedRatings = [(x*y) for x,y in zip(ratings,similarities)]
        return sum(weightedRatings) / sum(similarities)
    else:
        # User hasn't rated any similar items
        return ratingMean

def unpack(theta):
    global alpha
    global userBiases
    global bookBiases
    alpha = theta[0]
    userBiases = dict(zip(users, theta[1:nUsers+1]))
    bookBiases = dict(zip(books, theta[1+nUsers:]))

def cost(theta, labels, lamb):
    unpack(theta)
    predictions = [prediction(d['userID'], d['bookID']) for d in datatrain]
    cost = MSE(predictions, labels)
    print("MSE = " + str(cost))
    for u in userBiases:
        cost += lamb*userBiases[u]**2
    for i in bookBiases:
        cost += lamb*bookBiases[i]**2
    return cost

def derivative(theta, labels, lamb):
    unpack(theta)
    N = len(datatrain)
    dalpha = 0
    dUserBiases = defaultdict(float)
    dbookBiases = defaultdict(float)
    for d in datatrain:
        u,i = d['userID'], d['bookID']
        pred = prediction(u, i)
        diff = pred - float(d['rating'])
        dalpha += 2/N*diff
        dUserBiases[u] += 2/N*diff
        dbookBiases[i] += 2/N*diff
    for u in userBiases:
        dUserBiases[u] += 2*lamb*userBiases[u]
    for i in bookBiases:
        dbookBiases[i] += 2*lamb*bookBiases[i]
    dtheta = [dalpha] + [dUserBiases[u] for u in users] + [dbookBiases[i] for i in books]
    return numpy.array(dtheta)

def MSE(predictions, labels):
    differences = [(x-y)**2 for x,y in zip(predictions,labels)]
    return sum(differences) / len(differences)

labels = [float(d['rating']) for d in datatrain]
scipy.optimize.fmin_l_bfgs_b(cost, [alpha] + [0.0]*(nUsers+nBooks),derivative, args = (labels, 1))

#If this user not in the beta map, the model will become simple model with jac
prediction_valid = []
for d in datavalid:
    if d['userID'] in users and d['bookID'] in books:
        prediction_valid.append(prediction(d['userID'],d['bookID']))
    else:
        #退化成相似度模型
        prediction_valid.append(predictRating(d['userID'],d['bookID']))
#求出了预测
labels = [float(d['rating']) for d in datavalid]
print(MSE(prediction_valid,labels))

```

```

MSE = 1.4735475011336192
MSE = 1.4560931393014562
MSE = 1.473389955772163
MSE = 1.4733899534013817
1.4907800984682955

```

In [ ]:



In [1]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict # Dictionaries with default values
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip
f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")
header = f.readline()
header = header.strip().split(',')
datatrain = []
datavalid = []
count=0
lenAll = 200000
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < lenAll*0.95 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(set)
bookAllUser = defaultdict(set)

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    userReadBook[user].add(book)
    bookAllUser[book].add(user)
ratingMean = sum([float(d['rating']) for d in datatrain]) / len(datatrain)

N = len(datatrain)
nUsers = len(userReadBook)
nBooks = len(bookAllUser)
users = list(userReadBook)
books = list(bookAllUser)

alpha = ratingMean

userBiases = defaultdict(float)
bookBiases = defaultdict(float)

def prediction(user, book):
    return alpha + userBiases[user] + bookBiases[book]

def unpack(theta):
    global alpha
    global userBiases
    global bookBiases
    alpha = theta[0]
    userBiases = dict(zip(users, theta[1:nUsers+1]))
    bookBiases = dict(zip(books, theta[1+nUsers:]))

def cost(theta, labels, lamb):
    unpack(theta)
    predictions = [prediction(d['userID'], d['bookID']) for d in datatrain]
    cost = MSE(predictions, labels)
    print("MSE = " + str(cost))
    for u in userBiases:
        cost += lamb*userBiases[u]**2
    for i in bookBiases:
        cost += lamb*bookBiases[i]**2
    return cost

def derivative(theta, labels, lamb):
    unpack(theta)
    N = len(datatrain)
```

```

dalpha = 0
dUserBiases = defaultdict(float)
dbookBiases = defaultdict(float)
for d in datatrain:
    u,i = d['userID'], d['bookID']
    pred = prediction(u, i)
    diff = pred - float(d['rating'])
    dalpha += 2/N*diff
    dUserBiases[u] += 2/N*diff
    dbookBiases[i] += 2/N*diff
for u in userBiases:
    dUserBiases[u] += 2*lamb*userBiases[u]
for i in bookBiases:
    dbookBiases[i] += 2*lamb*bookBiases[i]
dtheta = [dalpha] + [dUserBiases[u] for u in users] + [dbookBiases[i] for i in books]
return numpy.array(dtheta)

def MSE(predictions, labels):
    differences = [(x-y)**2 for x,y in zip(predictions,labels)]
    return sum(differences) / len(differences)

alwaysPredictMean = [ratingMean for d in datatrain]
labels = [float(d['rating']) for d in datatrain]

MSE(alwaysPredictMean, labels)
scipy.optimize.fmin_l_bfgs_b(cost, [alpha] + [0.0]*(nUsers+nBooks),derivative, args = (labels, 1))

user_largest_beta = -100
user_largest_id = ""
user_smallest_beta = 100
user_smallest_id = ""
book_largest_beta = -100
book_largest_id = ""
book_smallest_beta = 100
book_smallest_id = ""
for user in userBiases:
    if userBiases[user]>user_largest_beta:
        user_largest_beta = userBiases[user]
        user_largest_id = user

for user in userBiases:
    if userBiases[user]<user_smallest_beta:
        user_smallest_beta = userBiases[user]
        user_smallest_id = user

for book in bookBiases:
    if bookBiases[book]>book_largest_beta:
        book_largest_beta = bookBiases[book]
        book_largest_id = book

for book in bookBiases:
    if bookBiases[book]<book_smallest_beta:
        book_smallest_beta = bookBiases[book]
        book_smallest_id = book
print("user_largest_beta is ", user_largest_beta,"user_largest_id is ",user_largest_id)
print("user_smallest_beta is",user_smallest_beta,"user_smallest_id",user_smallest_id)
print("book_largest_beta is",book_largest_beta,"book_largest_id is ",book_largest_id)
print("book_smallest_beta is ",book_smallest_beta,"book_smallest_id is ",book_smallest_id)

```

```

MSE = 1.4735475011336192
MSE = 1.4560931393014562
MSE = 1.473389955772163
MSE = 1.4733899534013817
user_largest_beta is  0.00040413237874470305 user_largest_id is  u92864068
user_smallest_beta is -0.0015796730337471908 user_smallest_id u11591742
book_largest_beta is 0.0008292191795822705 book_largest_id is  b76915592
book_smallest_beta is -0.0002721486787445039 book_smallest_id is  b57299824

```

In [ ]:

In [1]:

```
import numpy
import urllib
import scipy.optimize
import random
from collections import defaultdict # Dictionaries with default values
import nltk
import string
from nltk.stem.porter import *
from sklearn import linear_model
import ast
import gzip
def Jaccard(book1,book2,bookAllUser):
    s1=[]
    for d in bookAllUser[book1]:
        s1.append(d['userID'])
    s2=[]
    for d in bookAllUser[book2]:
        s2.append(d['userID'])
    s1=set(s1)
    s2=set(s2)
    numer = len(s1.intersection(s2))
    denom = len(s1.union(s2))
    return numer / denom
f = gzip.open("train_Interactions.csv.gz", 'rt', encoding="utf8")
header = f.readline()
header = header.strip().split(',')
datatrain = []
datavalid = []
count=0
lenAll = 200000
for line in f:
    fields = line.strip().split(',')
    d = dict(zip(header, fields))
    if count < lenAll*0.95 :
        datatrain.append(d)
    else:
        datavalid.append(d)
    count=count+1

userReadBook = defaultdict(list)
bookAllUser = defaultdict(list)

for d in datatrain:
    user,book,r =d['userID'],d['bookID'],d['rating']
    userReadBook[user].append(d)
    bookAllUser[book].append(d)
ratingMean = sum([float(d['rating']) for d in datatrain]) / len(datatrain)

N = len(datatrain)
nUsers = len(userReadBook)
nBooks = len(bookAllUser)
users = list(userReadBook)
books = list(bookAllUser)

alpha = ratingMean

userBiases = defaultdict(float)
bookBiases = defaultdict(float)

def prediction(user, book):
    return alpha + userBiases[user] + bookBiases[book]

def predictRating(user,item):
    ratings = []
    similarities = []
    for d in userReadBook[user]:
        i2 = d['bookID']
        if i2 == item: continue
        ratings.append(d['rating'])
        similarities.append(Jaccard(item,i2,bookAllUser))
    if (sum(similarities) > 0):
        weightedRatings = [(x*y) for x,y in zip(ratings,similarities)]
```



```

        return sum(weightedRatings) / sum(similarities)
    else:
        # User hasn't rated any similar items
        return ratingMean

def unpack(theta):
    global alpha
    global userBiases
    global bookBiases
    alpha = theta[0]
    userBiases = dict(zip(users, theta[1:nUsers+1]))
    bookBiases = dict(zip(books, theta[1+nUsers:]))

def cost(theta, labels, lamb):
    unpack(theta)
    predictions = [prediction(d['userID'], d['bookID']) for d in datatrain]
    cost = MSE(predictions, labels)
    print("MSE = " + str(cost))
    for u in userBiases:
        cost += lamb*userBiases[u]**2
    for i in bookBiases:
        cost += lamb*bookBiases[i]**2
    return cost

def derivative(theta, labels, lamb):
    unpack(theta)
    N = len(datatrain)
    dalpha = 0
    dUserBiases = defaultdict(float)
    dbookBiases = defaultdict(float)
    for d in datatrain:
        u,i = d['userID'], d['bookID']
        pred = prediction(u, i)
        diff = pred - float(d['rating'])
        dalpha += 2/N*diff
        dUserBiases[u] += 2/N*diff
        dbookBiases[i] += 2/N*diff
    for u in userBiases:
        dUserBiases[u] += 2*lamb*userBiases[u]
    for i in bookBiases:
        dbookBiases[i] += 2*lamb*bookBiases[i]
    dtheta = [dalpha] + [dUserBiases[u] for u in users] + [dbookBiases[i] for i in books]
    return numpy.array(dtheta)

def MSE(predictions, labels):
    differences = [(x-y)**2 for x,y in zip(predictions,labels)]
    return sum(differences) / len(differences)

labels = [float(d['rating']) for d in datatrain]
scipy.optimize.fmin_l_bfgs_b(cost, [alpha] + [0.0]*(nUsers+nBooks),derivative, args = (labels,2e-5))
predictions = open("predictions_Rating.txt", 'w')
for l in open("pairs_Rating.txt"):
    if l.startswith("userID"):
        #header
        predictions.write(l)
        continue
    user,book = l.strip().split('-')
    if user in users and book in books:
        prediction_valid=prediction(user,book)
    else:
        prediction_valid=predictRating(user,book)
    predictions.write(user + '-' + book + ", "+str(prediction_valid)+"\n")
predictions.close()
#my kaggle name is HumphreySD, and my score is 1.14656, my email address is haz013@eng.ucsd.edu
#I choose lamuda as 2e-5,MSE is 0.92566

```

```

MSE = 1.4735475011336192
MSE = 1.4560931393014562
MSE = 1.392284346401447
MSE = 8.08588938870108
MSE = 1.3695427812836598
MSE = 1.202888290365505
MSE = 1.201333071387938
MSE = 1.195213051215806
MSE = 1.1723466623360814
MSE = 1.0566656605572056

```

MSE = 1.05888586055/2056  
MSE = 1.0109866093480053  
MSE = 0.9792061553479576  
MSE = 0.9641404735702934  
MSE = 0.952691455209057  
MSE = 0.9382161957927877  
MSE = 0.9324830107667381  
MSE = 0.9310949475578871  
MSE = 0.9306016200670771  
MSE = 0.9519124072854905  
MSE = 0.9305585881131831  
MSE = 0.9297759227007407  
MSE = 0.9285850480621985  
MSE = 0.9270178994992827  
MSE = 0.9263009232898817  
MSE = 0.9257156932469571  
MSE = 0.9249547536300587  
MSE = 0.9256445353993737  
MSE = 0.9260333750758728  
MSE = 0.9262261853216237  
MSE = 0.9262792851467087  
MSE = 0.9319734642514613  
MSE = 0.9262812564150997  
MSE = 0.9261590468682027  
MSE = 0.9258838932958376  
MSE = 0.9257510823640084  
MSE = 0.9252038589117146  
MSE = 0.9255926057178915  
MSE = 0.9256206677681702  
MSE = 0.9256531516906248  
MSE = 0.9256757115957949  
MSE = 0.9256977359736026  
MSE = 0.9256076963007648  
MSE = 0.925651556577265  
MSE = 0.9256600242778392

In [ ]: