

EECE 5698 Assignment 4: Matrix Factorization

Zifeng Wang

Problem 1

1(a)

We can calculate the gradients as below

$$\begin{aligned}\nabla_{u_i} \text{RSE}(U, V) &= \nabla_{u_i} \left(\sum_{(i,j,r_{ij}) \in \mathcal{D}} (u_i^\top v_j - r_{ij})^2 + \lambda \sum_{i=1}^n \|u_i\|_2^2 + \mu \sum_{j=1}^n \|v_j\|_2^2 \right) \\ &= \sum_{j=1}^m 2 (u_i^\top v_j - r_{ij}) v_j + 2\lambda u_i\end{aligned}$$

$$\begin{aligned}\nabla_{v_j} \text{RSE}(U, V) &= \nabla_{v_j} \left(\sum_{(i,j,r_{ij}) \in \mathcal{D}} (u_i^\top v_j - r_{ij})^2 + \lambda \sum_{i=1}^n \|u_i\|_2^2 + \mu \sum_{j=1}^n \|v_j\|_2^2 \right) \\ &= \sum_{i=1}^n 2 (u_i^\top v_j - r_{ij}) u_i + 2\mu v_j\end{aligned}$$

1(b)

$$\begin{aligned}\nabla^2 \ell(u, v) &= \nabla^2 ((u^\top v - r)^2) \\ &= \nabla^2 ((uv - r)^2) \\ &= \begin{bmatrix} 2v^2 & 2(2uv - r) \\ 2(2uv - r) & 2u^2 \end{bmatrix} \\ \nabla^2 \ell(0, 0) &= \begin{bmatrix} 0 & -2r \\ -2r & 0 \end{bmatrix}\end{aligned}$$

It's eigenvalues are $2r, -2r$, due to the noise, $r \neq 0$, so the eigenvalues could not be all non-negative. So $\nabla^2 \ell(0, 0)$ is not positive semi-definite. Thus $\ell(u, v)$ is not convex.

1(c)

No. Because RSE is not convex. Gradient equals to zero doesn't mean global minimum.

1(d) Derive from results in 1(a)

$$\begin{aligned}
\nabla_{u_i} \text{RSE}(U, V) &= \sum_{j=1}^m 2 (u_i^\top v_j - r_{ij}) v_j + 2\lambda u_i \\
&= \sum_{j=1}^m 2 (0^\top 0 - r_{ij}) 0 + 2\lambda 0 \\
&= 0
\end{aligned}$$

$$\begin{aligned}
\nabla_{v_j} \text{RSE}(U, V) &= \sum_{i=1}^n 2 (u_i^\top v_j - r_{ij}) u_i + 2\mu v_j \\
&= \sum_{i=1}^n 2 (0^\top 0 - r_{ij}) 0 + 2\mu 0 \\
&= 0
\end{aligned}$$

Problem 2

2(a)

$$\begin{aligned}
\text{RSE}(U^k, V^k) &= \sum_{(i,j,r_{ij}) \in \mathcal{D}} \left(u_i^{k\top} v_j^k - r_{ij}^k \right)^2 + \lambda \sum_{i=1}^n \|u_i^k\|_2^2 + \mu \sum_{j=1}^n \|v_j^k\|_2^2 \\
&= \sum_{i=1}^n \sum_{j=1}^m \delta_{ij}^k{}^2 + \lambda \sum_{i=1}^n \|u_i^k\|_2^2 + \mu \sum_{j=1}^n \|v_j^k\|_2^2
\end{aligned}$$

$$\begin{aligned}
\nabla_{u_i} \text{RSE}(U^k, V^k) &= \sum_{j=1}^m 2 \left(u_i^{k\top} v_j^k - r_{ij}^k \right) v_j^k + 2\lambda u_i^k \\
&= \sum_{j=1}^m 2\delta_{ij}^k v_j^k + 2\lambda u_i^k
\end{aligned}$$

$$\begin{aligned}
\nabla_{v_j} \text{RSE}(U^k, V^k) &= \sum_{i=1}^n 2 \left(u_i^{k\top} v_j^k - r_{ij}^k \right) u_i^k + 2\mu v_j^k \\
&= \sum_{i=1}^n 2\delta_{ij}^k u_i^k + 2\mu v_j^k
\end{aligned}$$

2(b) The implementation is as follow:

```
def predict(u,v):
    """ Given a user profile uprof and an item profile vprof, predict the rating given by the user to this item

    Inputs are:
        -u: user profile, in the form of a numpy array
        -v: item profile, in the form of a numpy array

    The return value is
        - the inner product <u,v>
    """
    return np.dot(u, v)

def pred_diff(r,u,v):
    """ Given a rating, a user profile u and an item profile v, compute the difference between the prediction and actual rating

    Inputs are:
        -r: the rating a user gave to an item
        -u: user profile, in the form of a numpy array
        -v: item profile, in the form of a numpy array

    The return value is the difference
        -  $\delta = \langle u, v \rangle - r$ 
    """
    return np.dot(u, v) - r
```

```
def gradient_u(delta,u,v):
    """ Given a user profile u and an item profile v, and the difference in rating predictions  $\delta$ , compute the gradient

         $\nabla_u l(u,v) = 2 (\langle u, v \rangle - r) v$ 

    of the square error loss:

         $l(u,v) = (\langle u, v \rangle - r)^2$ 

    Inputs are:
        - $\delta$ : the difference  $\langle u, v \rangle - r$ 
        -u: user profile, in the form of a numpy array
        -v: item profile, in the form of a numpy array

    The return value is
        - The gradient w.r.t. u
    """
    return 2 * delta * v
```

```
def gradient_v(delta,u,v):
    """ Given a user profile u and an item profile v, and the difference in rating predictions  $\delta$ , compute the gradient

         $\nabla_v l(u,v) = 2 (\langle u, v \rangle - r) u$ 

    of the square error loss:

         $l(u,v) = (\langle u, v \rangle - r)^2$ 

    Inputs are:
        - $\delta$ : the difference  $\langle u, v \rangle - r$ 
        -u: user profile, in the form of a numpy array
        -v: item profile, in the form of a numpy array

    The return value is
        - the gradient w.r.t. v
    """
    return 2 * delta * u
```

2(c)

```
def generateItemProfiles(R,d,seed,sparkContext,N):
    """ Generate the item profiles from rdd R and store them in an RDD containing tuples of the form
        (j,vj)
    where v is a random np.array of dimension d.

    The random uis are generated using normalVectorRDD(), a function in RandomRDDs.

    Inputs are:
        - R: an RDD that contains the ratings in (user, item, rating) form
        - d: the dimension of the user profiles
        - seed: a seed to be used for in generating the random vectors
        - sparkContext: a spark context
        - N: the number of partitions to be used during joins, etc.

    The return value is an RDD containing the item profiles
    """
    V = R.map(lambda (i,j,rij):i).distinct(numPartitions = N)
    numUsers = V.count()
    randRDD = RandomRDDs.normalVectorRDD(sparkContext, numUsers, d,numPartitions=N, seed=seed)
    V = V.zipWithIndex().map(swap)
    randRDD = randRDD.zipWithIndex().map(swap)
    return V.join(randRDD,numPartitions = N).values()
```

We initialized the user and item profiles to random values instead of zero because we want to get rid of stuck at saddle point. If we initialize all profiles to be zero vectors, the RSE and its gradients will stay zero no matter how many iterations go, we will not get closer to our optimal solution.

Problem 4

4(a)

```
def joinAndPredictAll(R,U,V,N):
    """ Receives as inputs the ratings R, the user profiles U, and the items V, and constructs a joined RDD.

    Inputs are:
    - R: an RDD containing tuples of the form (i,j,rij)
    - U: an RDD containing tuples of the form (i,ui)
    - V: an RDD containing tuples of the form (j,vj)
    - N: the number of partitions to be used during joins, etc.

    The output is a joined RDD containing tuples of the form:

    (i,j,δij,ui,vj)

    where
    δij = <ui,vj>-rij
    is the prediction difference.

    """
    U_V = U.cartesian(V).map(lambda (i_u, j_v): ((i_u[0], j_v[0]), (i_u[1], j_v[1])))
    R = R.map(lambda (i, j, rij): ((i, j), rij))
    joined = R.join(U_V, numPartitions = N).map(lambda (ij, v) : (ij[0], ij[1], pred_diff(v[0],v[1][0],v[1][1]), v[1][0], v[1][1]))
    return joined
```

4(b)

```
def SE(joinedRDD):
    """ Receives as input a joined RDD as well as a λ and a μ and computes the MSE:

    SE(R,U,V) = Σ{i,j in data} (<ui,vj>-rij)^2

    The input is
    -joinedRDD: an RDD with tuples of the form (i,j,δij,ui,vj), where δij = <ui,vj> - rij is the prediction difference.

    The output is the SE.

    """
    SE = joinedRDD.map(lambda (i,j,sig,ui,vj): sig**2).reduce(add)
    return SE

def normSqRDD(profileRDD,param):
    """ Receives as input an RDD of profiles (e.g., U) as well as a parameter (e.g., λ) and computes the square of norms:
    λ Σi ||ui||2^2

    The input is:
    -profileRDD: an RDD of the form (i,u), where i is an index and u is a numpy array
    -param: a scalar λ>0

    The return value is:
    λ Σi ||ui||2^2

    """
    sum_u_norm = profileRDD.map(lambda (i, u): np.dot(u, u)).reduce(add)
    return param * sum_u_norm
```

In `__main__`, they are used as below to calculate the regularized objective function.

```
obj = SE(joinedRDD) + normSqRDD(U,args.lam) + normSqRDD(V,args.lam)
```

4(c)

```

def adaptU(joinedRDD, gamma, lam, N):
    """
    Receives as input a joined RDD
    as well as a gain  $\gamma$ , and regularization parameters  $\lambda$  and  $\mu$ , and constructs a new RDD of user profiles of the form

    
$$u_i = u_i - \gamma \nabla_{u_i} \text{RegSE}(R, U, V)$$


    where

    
$$\text{RegSE}(R, U, V) = \sum_{\{i, j \text{ in } R\}} (\langle u_i, v_j \rangle - r_{ij})^2 + \lambda \sum_i ||u_i||_2^2 + \mu \sum_j ||v_j||_2^2$$


    Inputs are
    -joinedRDD: an RDD with tuples of the form (i,j, $\delta_{ij}$ ,ui,vj), where  $\delta_{ij} = \langle u_i, v_j \rangle - r_{ij}$ 
    -gamma: the gain  $\gamma$ 
    -lam: the regularization parameter  $\lambda$ 
    -N: the number of partitions to be used in reduceByKey operations

    The return value is an RDD with tuples of the form (i,ui). The returned rdd contains exactly N partitions.
    """
    grad_u = joinedRDD.map(lambda (i,j,sig,ui,vj): (i, gradient_u(sig,ui,vj))).reduceByKey(add, numPartitions=N)
    U = joinedRDD.map(lambda (i,j,sig,ui,vj): (i, ui)).reduceByKey(lambda x, y: (x+y)/2, numPartitions=N)
    new_U = U.join(grad_u, numPartitions=N).mapValues(lambda (ui, grad_u): ui-gamma*(grad_u+2*lam*ui))
    return new_U

def adaptV(joinedRDD, gamma, mu, N):
    """
    Receives as input a joined RDD
    as well as a gain  $\gamma$ , and regularization parameters  $\lambda$  and  $\mu$ , and constructs a new RDD of user profiles of the form

    
$$u_i = u_i - \gamma \nabla_{u_i} \text{RegSE}(R, U, V)$$


    where

    
$$\text{RegSE}(R, U, V) = \sum_{\{i, j \text{ in } R\}} (\langle u_i, v_j \rangle - r_{ij})^2 + \lambda \sum_i ||u_i||_2^2 + \mu \sum_j ||v_j||_2^2$$


    Inputs are
    -joinedRDD: an RDD with tuples of the form (i,j, $\delta_{ij}$ ,ui,vj), where  $\delta_{ij} = \langle u_i, v_j \rangle - r_{ij}$ 
    -gamma: the gain  $\gamma$ 
    -mu: the regularization parameter  $\mu$ 
    -N: the number of partitions to be used in reduceByKey operations

    The return value is an RDD with tuples of the form (j,vj). The returned rdd contains exactly N partitions.
    """
    grad_v = joinedRDD.map(lambda (i,j,sig,ui,vj): (j, gradient_v(sig,ui,vj))).reduceByKey(add, numPartitions=N)
    V = joinedRDD.map(lambda (i,j,sig,ui,vj): (j, vj)).reduceByKey(lambda x, y: (x+y)/2, numPartitions=N)
    new_V = V.join(grad_v, numPartitions=N).mapValues(lambda (vj, grad_v): vj-gamma*(grad_v+2*mu*vj))
    return new_V

```

Problem 5

5(a)

This part is for reading data for each fold if **args.folds** is specified, otherwise, read all data in a single fold.

```

folds = {}

if args.output is None:
    for k in range(args.folds):
        folds[k] = readRatings(args.data+"/fold"+str(k),sc)
else:
    folds[0] = readRatings(args.data,sc)

```

This part constructs the training and testing set for each code.

```

train_folds = [folds[j] for j in folds if j is not k ]

if len(train_folds)>0:
    train = train_folds[0]
    for fold in train_folds[1:]:
        train=train.union(fold)
    train.repartition(args.N).cache()
    test = folds[k].repartition(args.N).cache()
    Mtrain=train.count()
    Mtest=test.count()

    print("Initiating fold %d with %d train samples and %d test samples" % (k,Mtrain,Mtest) )
else:
    train = folds[k].repartition(args.N).cache()
    test = train
    Mtrain=train.count()
    Mtest=test.count()
    print("Running single training over training set with %d train samples. Test RMSE computes RMSE on training set" % Mtrain )

```

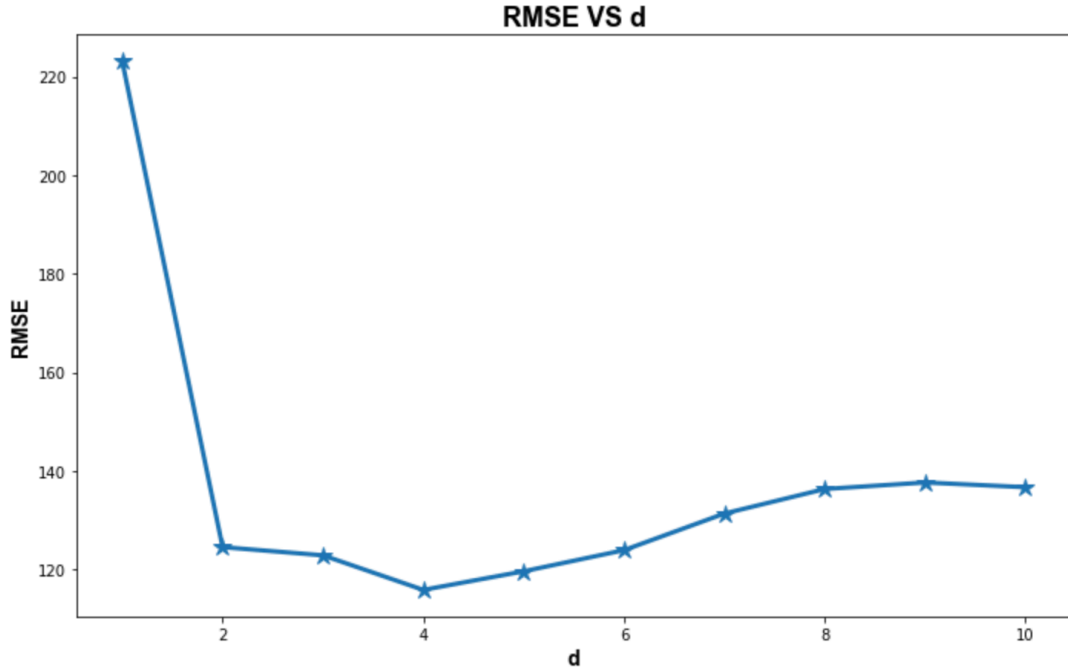
Explanation of what each operation does is in the comments in red below:

```
# specify training folds as all folds except the test fold
train_folds = [folds[j] for j in folds if j is not k ]
if len(train_folds)>0:
    # the next 3 lines merge the folds inside train_folds into one
    single list
    train = train_folds[0]
    for fold in train_folds[1:]:
        train=train.union(fold)
    # repartition the train data and cache it
    train.repartition(args.N).cache()
    # repartition the train data and cache it
    test = folds[k].repartition(args.N).cache()
    # calculate number of training samples
    Mtrain=train.count()
    # calculate number of test samples
    Mtest=test.count()
    # print info
    print("Initiating fold %d with %d train samples and %d test
samples" % (k,Mtrain,Mtest) )
else:
    # use all data as training set and repartition & cache it
    train = folds[k].repartition(args.N).cache()
    # use training set as test set
    test = train
    # calculate number of training samples
    Mtrain=train.count()
    # calculate number of test samples
    Mtest=test.count()
    # print info
    print("Running single training over training set with %d train
samples. Test RMSE computes RMSE on training set" % Mtrain )
```

5(b) step size is set as below, i is the number of iteration, **args.gain** and **args.power** are user specified parameters:

```
gamma = args.gain / i**args.power
```

5(c) Plot the RMSE w.r.t. d as below:



We can see that $b = 4$ is the optimal value

5(d)

For **--gain 0.1 --pow 0.0 --maxiter 5**, we have the observation below

```
Iteration: 1 Time: 9.074272 Objective: 1013500.042176 TestRMSE: 165.467171
Iteration: 2 Time: 20.283857 Objective: 192387332.578143 TestRMSE: 2309.034806
Iteration: 3 Time: 30.292711 Objective: 313778766684237948233728.000000 TestRMSE: 237055796220.879547
Iteration: 4 Time: 42.134367 Objective: 3156851951687905134926714567158335603860742213578240926818096813965312.000000 TestRMSE: 8905796323923928549709324216172544.000000
Iteration: 5 Time: 52.153823 Objective:
4411177661468997000228206321331790297785343502071231801627761775796236689393956133128258666500963174959472179665456704812477888591540742279500909207131195841344404951401330271814904206454042487575980720
259072.000000 TestRMSE: 64747973450158910312374268541083223719470910226173344799834268784317605305368308542511972305191567360.000000
Initiating fold 1 with 40 train samples and 10 test samples
Training set contains 10 users and 10 items
Iteration: 1 Time: 8.352871 Objective: 934540.007357 TestRMSE: 187.817433
Iteration: 2 Time: 18.836488 Objective: 2025824149.155913 TestRMSE: 5114.998322
Iteration: 3 Time: 30.346551 Objective: 7689176559326446978334720.000000 TestRMSE: 92883984755.931299
Iteration: 4 Time: 40.665933 Objective: 64351920497074678652348626426728149959055198200128328669628322806776832.000000 TestRMSE: 18449534623344717054295035346944000.000000
Iteration: 5 Time: 51.119073 Objective:
4110586007484096398190993888283767174445697683320855752712471429068983197891487097454398899065647855064891800006650905829761618380662720232098107438025201648584521017945124307685266449722576692868784400
3906846720.000000 TestRMSE: 34690998838395647095791246085143714607040170343558604333816324662343417879708829296026842928459806670848.000000
Initiating fold 2 with 40 train samples and 10 test samples
Training set contains 10 users and 10 items
Iteration: 1 Time: 8.361678 Objective: 1135618.639875 TestRMSE: 123.156506
Iteration: 2 Time: 19.509996 Objective: 2594881440.592070 TestRMSE: 6211.310434
Iteration: 3 Time: 30.373434 Objective: 9463139258607852767412224.000000 TestRMSE: 109927540214.046951
Iteration: 4 Time: 41.694653 Objective: 870112131350369054872560550861952904857232391701940177941329561123815424.000000 TestRMSE: 17426197321977737591445268541210624.000000
Iteration: 5 Time: 52.736168 Objective:
830507043334298176312731598233197905396429502611739161249731927241242273201427650858796182863283486978749477344611910685999667787654805578572136319403385138478599693492569538355053834491467812754386283
2197009408.000000 TestRMSE: 2670576818966763400917666158243583181701288248604892648550415800410556283935027958663309538302925733888.000000
Initiating fold 3 with 40 train samples and 10 test samples
Training set contains 10 users and 10 items
Iteration: 1 Time: 8.356851 Objective: 1007655.537988 TestRMSE: 167.223007
Iteration: 2 Time: 18.924245 Objective: 5111830477.907906 TestRMSE: 8403.427325
Iteration: 3 Time: 29.763850 Objective: 10845051638235119901261824.000000 TestRMSE: 568305429493.300415
Iteration: 4 Time: 40.451817 Objective: 1740187604225424373122004013670634517765614800494006407063606648700487598080.000000 TestRMSE: 1213967623127571945345702013519790800.000000
Iteration: 5 Time: 51.705304 Objective:
805116425814381958907246263442593127036580996403095511436697094108327089122971104424263143399408338646929720472335641750401310839599953587635790891470873974804804205429871170047370265072641450167840
83749236277262129152.000000 TestRMSE: 53189708078947679431312571394175848358057137793717832849848355496505490416484841767825912477605215338496.000000
Initiating fold 4 with 40 train samples and 10 test samples
Training set contains 10 users and 10 items
Iteration: 1 Time: 9.126647 Objective: 1057955.040949 TestRMSE: 151.439377
Iteration: 2 Time: 19.599344 Objective: 5019020205.414135 TestRMSE: 1720.935899
Iteration: 3 Time: 30.767882 Objective: 1083297297436629723293616.000000 TestRMSE: 56454004136.552391
Iteration: 4 Time: 42.049717 Objective: 1763793273321029049365631547114704388262179050215704960830067954420423852032.000000 TestRMSE: 58010308585005927224029553389297664.000000
Iteration: 5 Time: 53.351197 Objective:
83968283242638761700985123080478035319853487996236573647131730543314886150179722410881654794547126536705468854910239093120699546270284515785073267884474064426190867621493085364536839198654602980285337964
64671852733141090304.000000 TestRMSE: 126074452646853528135059866378918250239727206316147873638192972974581113998784825110719490414584132869816320.000000
5 fold cross validation error is: 108002135993501844158096394785638024330005507160964376223271725469822390176278818412185099717058358719068432.000000
```

We can see that the step size is too big, so the RMSE is ‘exploding’.

For **--gain 0.0001 --pow 1.0 --maxiter 20**, we have the observation below (for simplicity I just capture screen for 1 fold, the rest are similar):

```

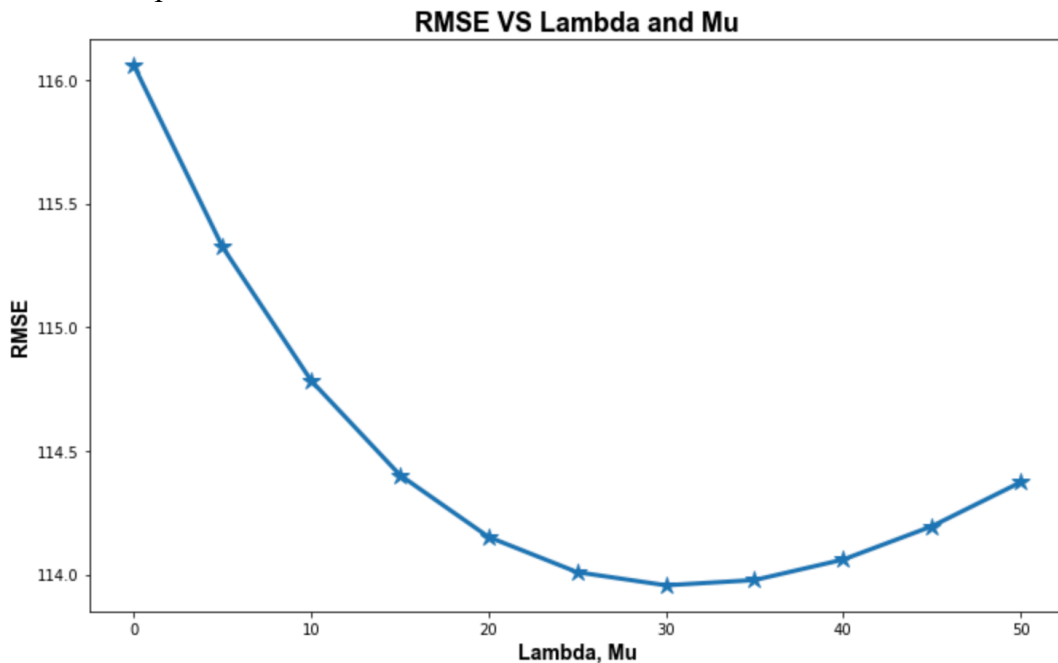
Training set contains 10 users and 10 items
Iteration: 1   Time: 9.031763  Objective: 1013500.042176  TestRMSE: 165.467171
Iteration: 2   Time: 20.364278 Objective: 1010755.181194  TestRMSE: 165.464027
Iteration: 3   Time: 31.279617 Objective: 1009408.319779  TestRMSE: 165.453452
Iteration: 4   Time: 43.331643 Objective: 1008501.947583  TestRMSE: 165.443073
Iteration: 5   Time: 53.420159 Objective: 1007813.094133  TestRMSE: 165.433514
Iteration: 6   Time: 63.803859 Objective: 1007254.467580  TestRMSE: 165.424752
Iteration: 7   Time: 76.152883 Objective: 1006782.808386  TestRMSE: 165.416680
Iteration: 8   Time: 86.259500 Objective: 1006373.486731  TestRMSE: 165.409194
Iteration: 9   Time: 97.272816 Objective: 1006011.123790  TestRMSE: 165.402207
Iteration: 10  Time: 107.417636 Objective: 1005685.461192  TestRMSE: 165.395651
Iteration: 11  Time: 117.996144 Objective: 1005389.307289  TestRMSE: 165.389468
Iteration: 12  Time: 129.741703 Objective: 1005117.420821  TestRMSE: 165.383611
Iteration: 13  Time: 139.845720 Objective: 1004865.861072  TestRMSE: 165.378044
Iteration: 14  Time: 150.207647 Objective: 1004631.588413  TestRMSE: 165.372734
Iteration: 15  Time: 160.763873 Objective: 1004412.207617  TestRMSE: 165.367655
Iteration: 16  Time: 170.849454 Objective: 1004205.796721  TestRMSE: 165.362785
Iteration: 17  Time: 181.421509 Objective: 1004010.789316  TestRMSE: 165.358104
Iteration: 18  Time: 192.244272 Objective: 1003825.891393  TestRMSE: 165.353595
Iteration: 19  Time: 202.592768 Objective: 1003650.021234  TestRMSE: 165.349245
Iteration: 20  Time: 214.696185 Objective: 1003482.265060  TestRMSE: 165.345041

```

We can see that the step size is too small, so the optimization process is too slow.

5(e)

Here as the searching space is continuous and so big, I just let $\lambda = \mu = 0 : 5 : 50$ to search for the optimal value. However, a finer search may be conducted further around the optimal value I find.



From the figure, we can see that $\lambda = \mu = 30$ is approximately the optimal solution.

Bonus:

The method is similar as above.

- Set $\lambda = \mu = 0$, and set $d = 1:10$ to evaluate the least RMSE. If we can observe a first decreasing and then increasing curve, we take the minimum value of d , either do finer search or fix that value. Otherwise, give d a wider range to search.
- For investigating λ and μ , I use a greedy search. I search λ in range $[1,50]$ first, find the optimal value for λ . Then fix the optimal value for λ , search μ in range $[1,50]$
- As it is too time consuming to do all this experiments, I just offer the general method here, more work needs to be done to find the exact optimal triplet.