

CS550: Massive Data Mining and Learning

Homework 2

Due 11:59pm Saturday, March 23, 2019

Only one late period is allowed for this homework (11:59pm Sunday
3/24)

Submission Instructions

Assignment Submission Include a signed agreement to the Honor Code with this assignment. Assignments are due at 11:59pm. All students must submit their homework via Sakai. Students can typeset or scan their homework. Students also need to include their code in the final submission zip file. Put all the code for a single question into a single file.

Late Day Policy Each student will have a total of *two* free late days, and for each homework only one late day can be used. If a late day is used, the due date is 11:59pm on the next day.

Honor Code Students may discuss and work on homework problems in groups. This is encouraged. However, each student must write down their solutions independently to show they understand the solution well enough in order to reconstruct it by themselves. Students should clearly mention the names of all the other students who were part of their discussion group. Using code or solutions obtained from the web is considered an honor code violation. We check all the submissions for plagiarism. We take the honor code seriously and expect students to do the same.

Discussion Group (People with whom you discussed ideas used in your answers):

On-line or hardcopy documents used as part of your answers:

I acknowledge and accept the Honor Code.

(Signed)Xinyang Wang_____

If you are not printing this document out, please type your initials above.

Answer to Question 1(a)

$$(\mathbf{M}\mathbf{M}^T)^T = (\mathbf{M}^T)^T \mathbf{M}^T = \mathbf{M}\mathbf{M}^T$$

$$(\mathbf{M}^T \mathbf{M})^T = \mathbf{M}^T \mathbf{M}$$

So both are symmetric.

Because matrix \mathbf{M} is of size $p \times q$, so \mathbf{M}^T is of size $q \times p$.

So $\mathbf{M}\mathbf{M}^T$ is of size $p \times p$ and

$\mathbf{M}^T \mathbf{M}$ is of size $q \times q$

so both are square

Because \mathbf{M} is real, \mathbf{M}^T is real. So the products of $\mathbf{M}^T \mathbf{M}$ and $\mathbf{M}\mathbf{M}^T$ are real.

Answer to Question 1(b)

Suppose the eigenvalue of $\mathbf{M}\mathbf{M}^T$ is λ and corresponding eigenvector is \mathbf{e} .

$$\mathbf{M}\mathbf{M}^T \mathbf{e} = \lambda \mathbf{e}$$

we multiply \mathbf{M}^T to the equation above from the left:

$$\mathbf{M}^T \mathbf{M}\mathbf{M}^T \mathbf{e} = \mathbf{M}^T \lambda \mathbf{e}$$

we can rewrite it as below:

$$(\mathbf{M}^T \mathbf{M})(\mathbf{M}^T \mathbf{e}) = \lambda(\mathbf{M}^T \mathbf{e})$$

As we can see above, the eigenvalues the same and the eigenvectors are not.

Answer to Question 1(c)

$$\mathbf{M}^T \mathbf{M} = \mathbf{Q} \lambda \mathbf{Q}^T$$

Answer to Question 1(d)

$$\mathbf{M}^T \mathbf{M} = (\mathbf{U} \Sigma \mathbf{V}^T)^T * (\mathbf{U} \Sigma \mathbf{V}^T) = (\mathbf{V} \Sigma \mathbf{U}^T)(\mathbf{U} \Sigma \mathbf{V}^T) = \mathbf{V} \Sigma (\mathbf{U}^T \mathbf{U}) \Sigma \mathbf{V}^T = \mathbf{V} \Sigma^2 \mathbf{V}^T$$

Answer to Question 1(e)(a)

The returned values are as follows

```

1 U is
2 [[-0.27854301  0.5
3  [-0.27854301 -0.5
4  [-0.64993368  0.5
5  [-0.64993368 -0.5
6 s is
7 [ 7.61577311  1.41421356]
8 VT is
9 [[-0.70710678 -0.70710678]
10 [-0.70710678  0.70710678]]

```

Answer to Question 1(e)(b)

The values of Evals and Evecs are as follows(after sorting and re-arranging):

```

1 Evals is
2 [58.0 , 2.0]
3 Evecs is
4 [[ 0.70710678 -0.70710678]
5 [ 0.70710678  0.70710678]]

```

Answer to Question 1(e)(c)

As we can see from the results of our experiments, the V produced by SVD and the matrix of eigenvectors Evecs are equal in value.

Answer to Question 1(e)(d)

As we can see from the results of our experiments, the eigenvalues of $\mathbf{M}^T\mathbf{M}$ are square of the singular values of M.

Answer to Question 2(a)

Target is to prove that: for $r' = Mr$, we have $w(r') = w(r)$

let's consider the first row of r' (r'_1) and the first row of r (r_1). We have $r'_1 = \sum_{i=1}^n (W_{i1}r_1) = (\sum_{i=1}^n W_{i1})r_1 = r_1$. And it is easy to see this relation holds for every row of r . So we can draw the conclusion that $r' = r$

Answer to Question 2(b)

Under the circumstance of $w(r) = 1$ will $w(r') = w(r)$.

Proof is shown below:

Assume $w(r') = w(r)$, and we have $\beta w(r) = \beta w(r')$. So $w(r') = \beta w(r') + (1 - \beta)$, which leads to $w(r') = 1$

Answer to Question 2(c)(a)

$$r'_i = \beta \sum_{livej=1}^n M_{ij}r_j + \beta \sum_{deadj=1}^n M_{ij}r_j + (1 - \beta)/n$$

Answer to Question 2(c)(b)

As we can see in the equation above, we can assume that 'dead' ends are in some sense 'live'.(we can think of it in the following way: each 'dead' end actually has n-1 arc to all the left nodes) Thus we can combine the first two parts in the above equation together, and rewrite it as follows:

$$r'_i = \beta \sum_{allj=1}^n M_{ij}r_j + (1 - \beta)/n$$

Thus if $w(r) = 1$, according to the conclusion of question 2(b), we have $w(r') = 1$ as well.

Answer to Question 3(a)

Before giving out the results, let's first look at the structure of our code: (detailed code is in the question3 file)

```
1 class PageRank():
2
3     def __init__(self,n:int):
4         '''Initialization '''
5
6     def load_file(self,file):
7         '''load file '''
8
9     def generate_matrix(self,file):
10        '''generate matrix M according to file \
11        (during which remove the duplicates) '''
12
13    def Iteration(self, times:int, beta):
14        '''iterate to compute final scores '''
```

And the top 5 ids with scores are shown as follows:

```
1 the top 5 ids with the highest PageRank score
2 id:53, score:[ 0.0357312]
3 id:14, score:[ 0.03417091]
4 id:40, score:[ 0.03363009]
5 id:1, score:[ 0.03000598]
6 id:27, score:[ 0.02972014]
```

Answer to Question 3(b)

The bottom 5 ids with scores are shown as follows:

```
1 the bottom 5 ids with the lowest PageRank score
2 id:85, score:[ 0.00340969]
3 id:59, score:[ 0.00366986]
4 id:81, score:[ 0.00369535]
5 id:37, score:[ 0.0038082]
6 id:89, score:[ 0.00392247]
```

Answer to Question 4(a)

Before giving out the results, let's first look at the structure of our code: (detailed code is in the question4 file)

```
1 def mapper(line):
2     '''splitting the line and then return \
3     in the form of np.array'''
4
5 def closest_center(row, centers):
6     '''for each data, find the closest center \
7     and returned in the form (center,(data,1))'''
8
9 def cost_computing(row, centers):
10    '''compute the distance between data and \
11    corresponding center'''
12
13 def plot_cost(c1_cost, c2_cost):
14    '''plot cost as a function of loop'''
```

some attention 1: if you want to test the code, you should configure the pyspark environment first. Otherwise, you can use the jupyter offered by school server (<https://jupyter.cs.rutgers.edu/hub/login?next=%2Fhub%2Fuser%2Fhz333>).

some attention 2: before running the script, be sure to put the data.txt, c1.txt, c2.txt into the school server HDFS system by using this link: <https://data-services1.cs.rutgers.edu/#/login>. Because pyspark read the file from HDFS system on default.

Below we plot the cost function as a function of iterations.

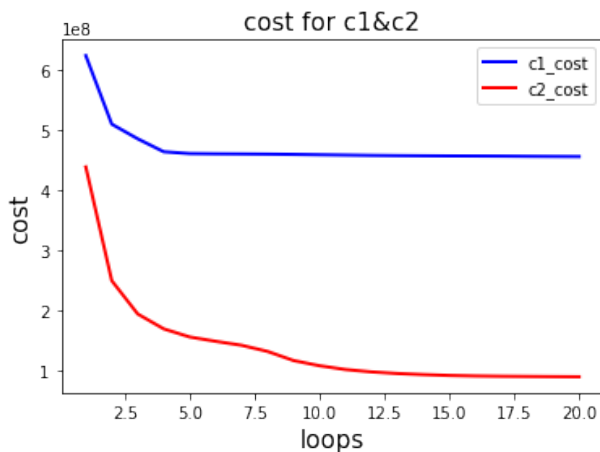


Figure 1: cost

Answer to Question 4(b)

In this part, we compute how much the cost decreases when using c1 and c2. And we also output the final centers in c1 and c2.

```
1 when using c1 as centroids , the percentage change in cost\  
2 after 10 iterations is  
3 0.2648391714456056  
4  
5 when using c2 as centroids , the percentage change in cost\  
6 after 10 iterations is  
7 0.7669795594605946
```

```
centers in c1  
0.170.260.480.420.310.160.080.220.340.630.120.750.150.120.20.250.160.351.720.21.370.310.210.190.410.130.020.020.010.0  
20.010.00.090.00.070.030.110.010.020.030.010.040.010.010.040.030.00.020.070.130.030.690.210.0710.62187.83868.940.75  
  
0.190.180.590.00.50.310.160.140.450.440.180.570.210.270.540.290.180.461.690.480.910.180.430.210.180.040.010.00.00.010  
.00.00.070.00.00.030.080.00.00.040.010.00.050.010.050.060.00.010.020.090.010.370.230.0663.93555.61770.690.85  
  
0.130.120.30.160.30.110.090.130.390.380.10.510.130.110.10.280.120.311.530.240.940.970.150.180.310.150.020.010.00.020.  
030.00.070.00.090.040.150.00.020.030.010.030.010.020.040.030.010.030.040.130.040.280.170.115.06110.281273.630.65  
  
0.060.410.20.00.230.040.050.040.020.140.010.470.060.050.00.190.030.111.670.020.470.00.020.060.580.282.150.120.190.120  
.10.070.090.070.120.110.090.030.10.080.070.290.030.140.570.240.00.040.020.120.020.20.020.011.695.3420.580.12  
  
0.290.130.270.080.070.110.020.130.210.380.140.580.320.590.030.150.260.121.750.040.810.210.420.270.090.050.00.00.00.00  
.00.00.280.00.00.030.160.00.010.020.020.010.00.010.040.050.010.050.070.30.010.260.330.1210.45388.564058.560.77  
  
0.090.140.350.170.450.130.180.130.070.270.080.550.140.040.040.30.310.251.730.140.990.160.160.10.480.210.090.090.050.0  
80.040.030.10.030.070.090.170.00.040.060.010.030.050.040.140.110.00.020.030.150.020.310.110.044.1252.17236.430.61  
  
0.120.090.240.00.250.070.140.130.030.230.070.570.080.030.040.240.110.181.660.040.860.010.060.070.660.360.210.230.090.  
160.070.060.130.060.160.130.190.010.110.060.030.10.080.10.270.240.010.030.020.150.010.250.050.052.6315.9773.250.36  
  
0.110.110.310.00.430.150.140.110.040.210.050.520.090.010.030.310.170.21.660.090.990.030.140.090.520.330.20.160.080.12  
0.050.050.090.050.110.110.150.010.10.070.050.040.080.040.240.150.00.030.030.140.020.250.080.053.0525.79140.450.49  
  
0.160.140.340.050.330.120.090.140.250.350.090.750.120.10.080.270.240.141.470.121.030.320.140.180.740.170.020.020.010.  
020.010.00.090.00.020.070.140.010.020.020.050.080.010.010.060.150.010.020.080.160.010.340.110.076.1776.09582.320.58  
  
0.120.140.330.370.40.110.250.160.170.290.10.530.090.030.040.310.340.211.740.161.00.520.160.10.450.140.090.050.040.030  
.050.020.040.020.050.060.140.020.040.10.020.030.030.030.110.110.010.030.180.140.010.290.120.18.3177.24376.580.59
```

Figure 2: final-center-in-c1

Figure 3: final-center-in-c1

9