

**Question 1(a), Homework 2, CS246**

---

Both matrices are symmetric square and real.

Symmetric:  $(M^T M)^T = M^T (M^T)^T = M^T M$ . And  $(M M^T)^T = (M^T)^T M^T = M M^T$

Square:  $M^T M$  is size  $q \times q$ .  $M M^T$  is size  $p \times p$ . If we multiply matrix size  $p \times q$  by its transpose, which is size  $q \times p$  we get square matrix size  $p \times p$ .

Real: If  $M$  is real, the product of  $M$  and its transpose will also be real.

For each eigenvalue  $MM^T$ , we can write  $MM^T v = v\lambda$ . By multiplying both sides by prefix of  $M^T$ , we get  $M^T M(M^T v) = (M^T v)\lambda$ , where the parentheses can be added because of associativity rule.  $\lambda$  is also an eigenvalue of  $M^T M$ , with a corresponding eigenvector being  $M^T v$ . The eigenvectors are usually not the same. An exception is  $M^T = I$ .

We can write  $M^T M = Q \lambda Q^T$ , because  $M^T M$  is symmetric, square and real by 1(a).

$$M^T M = V \sum U^T U \sum V^T = V \sum^2 V^T, \text{ which matches 1(c).}$$

```
U
[[-0.27854301  0.5      ]
 [-0.27854301 -0.5     ]
 [-0.64993368  0.5      ]
 [-0.64993368 -0.5     ]]
S
[7.61577311  1.41421356]
Vh
[[-0.70710678 -0.70710678]
 [-0.70710678  0.70710678]]
Evals
[ 2. 58.]
Evecs
[[-0.70710678  0.70710678]
 [ 0.70710678  0.70710678]]
```

Figure 1: Output

On figure 1 we can see output of assignment.

V is equivalent to the matrix of the eigenvectors if we reorder the columns based on ordering of their singular values.

Each and every singular value of M is square root of an eigenvalues of  $M^T M$ .

Initialization	Percentage change after 10 iterations
c1.txt	26.5%
c2.txt	76.7%

Table 1: Euclidian distance

C2 is better than c1, because it distributes the initial clusters far apart. Because there is less overlap true clusters are split less often which leads to a better final set of clusters.

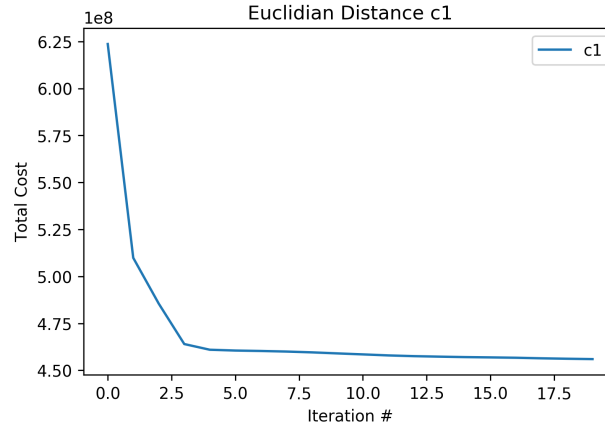


Figure 2: Euclidian graph c1

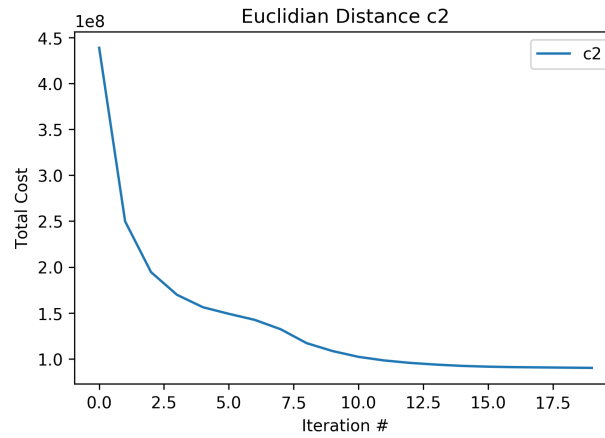


Figure 3: Euclidian graph c2

Initialization	Percentage change after 10 iterations
c1.txt	18.7%
c2.txt	51.6%

Table 2: Manhattan distance

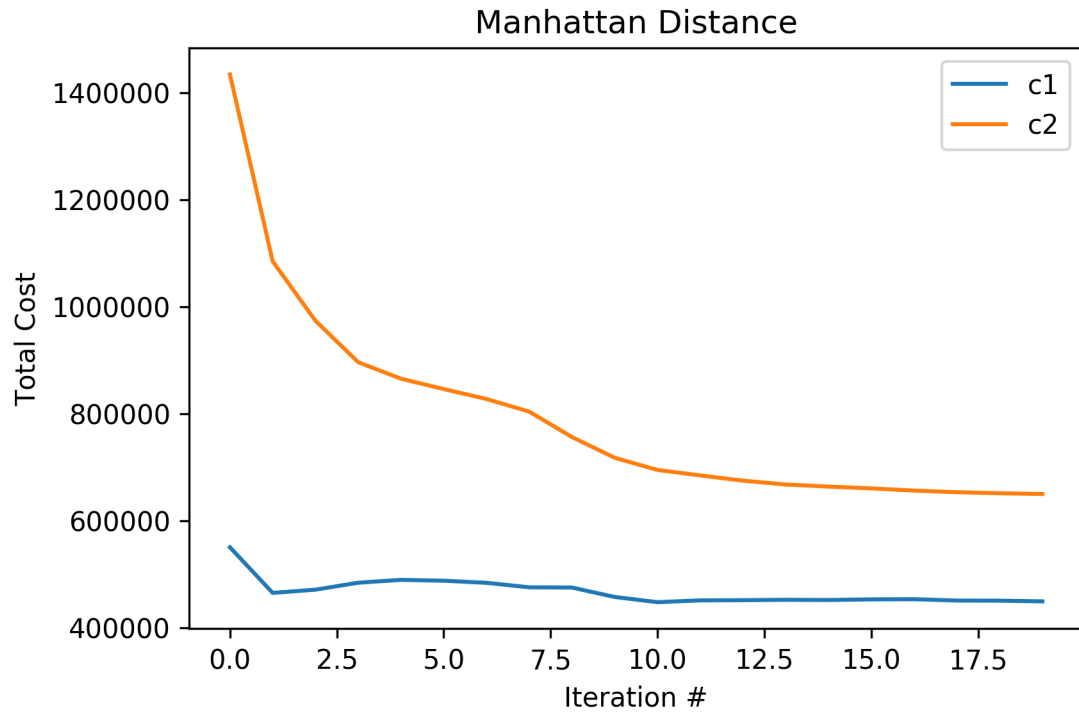


Figure 4: Manhattan graph

For the Manhattan distance, random initialization doesn't appear to be better, because points in c2 are far apart in Euclidean distance, they are not necessarily apart in the Manhattan distance.

$$\epsilon_{iu} = 2 * (R_{iu} - q_i * p_u^T)$$

For SGD we then denote:

$$\frac{\partial E}{\partial q_i} = \left[ \sum_{u(iu) \in ratings} 2(R_{iu} - q_i p_u^T)(-p_u) \right] + 2\lambda q_i$$

$$\frac{\partial E}{\partial p_u} = \left[ \sum_{u(iu) \in ratings} 2(R_{iu} - q_i p_u^T)(-q_i) \right] + 2\lambda p_u$$

$$q_i := q_i + \eta * (\epsilon_{iu} * p_u - 2 * \lambda * q_i)$$

$$p_u := p_u + \eta * (\epsilon_{iu} * q_i - 2 * \lambda * p_u)$$



We got best result for  $\eta = 0.015$ .

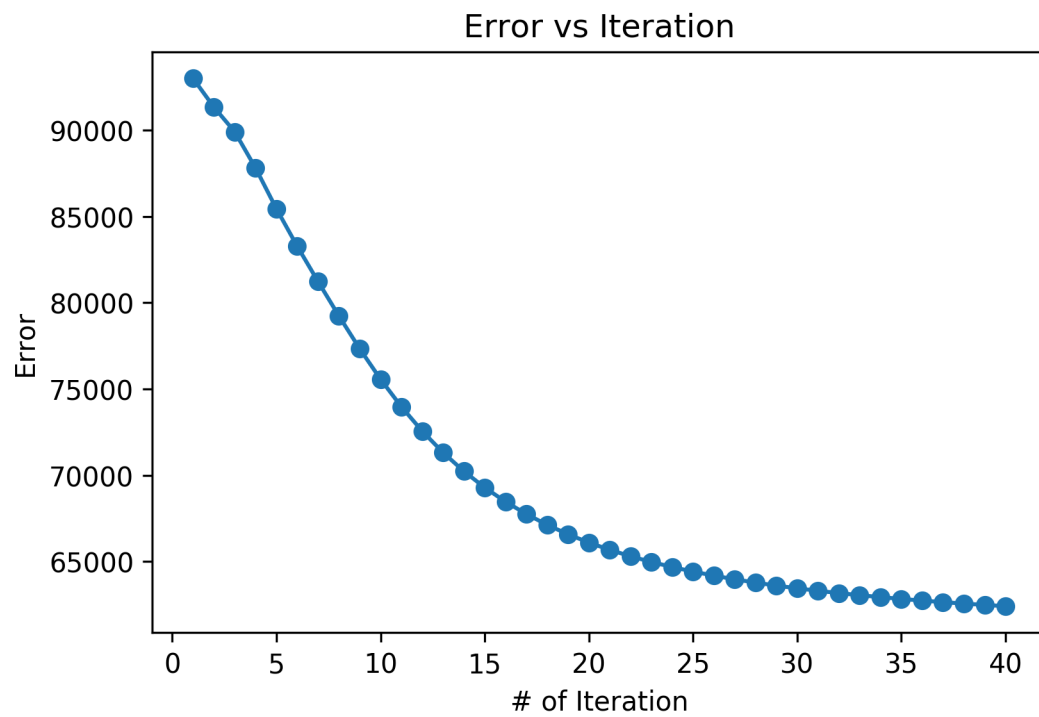


Figure 5: Error vs Iterations

**Question 4(a), Homework 2, CS246**

---

$T_{ii}$  is the number of items that user  $i$  likes, so out-degree of user node  $i$ .

Some math behind this:  $T_{ii} = \sum_{j=1}^n R_{ij} * (R^T)_{ji} = \sum_{j=1}^n R_{ij}^2 = \sum_{j=1}^n R_{ij}$ .  $R_{ij} = 0$  or  $1$ , so  $T_{ii} = \text{degree}(\text{user})$ .

$T_{ij}$  is the number of paths between user  $i$  and user  $j$  with the length of 2. It also represents the number of items they both like. I.e. Number of common neighbours of user  $i$  and  $j$ .

Again some math behind this:  $T_{ij} = \sum_{k=1}^n R_{ik} * R_{jk}$ ,  $R_{ik} * R_{jk} (\neq 0)$  this means it exists a 2 step path started from user  $i$  ended in user  $j$  via item  $k$ .

A diagonal matrix can scale the columns of its left multiplier, so to normalize items we use  $RQ^{-\frac{1}{2}}$ . So, similar to 4(a) the item similarity matrix  $S_I = (RQ^{-\frac{1}{2}})^T RQ^{-\frac{1}{2}} = Q^{-\frac{1}{2}} R^T R Q^{-\frac{1}{2}}$ .

Similar to  $S_I$ ,  $S_U = (R^T P^{-\frac{1}{2}})^T R^T P^{-\frac{1}{2}} = P^{-\frac{1}{2}} R R^T P^{-\frac{1}{2}}$ .

User-user collaborative filtering:  $\Gamma_{us} = \sum_{x \in \text{users}} \cos - \text{sim}(x, u) * R_{x,s}$

So  $\Gamma = S_U R = P^{-\frac{1}{2}} R R^T P^{-\frac{1}{2}} R$ .

Item-item collaborative filtering:  $\Gamma_{us} = \sum_{x \in \text{items}} R_{ux} * \cos - \text{sim}(x, s)$

So  $\Gamma = R S_I = R Q^{-\frac{1}{2}} R^T R Q^{-\frac{1}{2}}$ .

Note: We take  $i_{alex} = 499$  due to 0-indexing.

User\_user shows:

- FOX 28 News at 10 pm
- Family Guy
- 2009 NCAA Basketball Tournament
- NBC 4 at Eleven
- Two and a Half Men

item\_item shows:

- FOX 28 News at 10 pm
- Family Guy
- NBC 4 at Eleven
- 2009 NCAA Basketball Tournament
- Access Hollywood

There is no significant advantage to any of used methods. Precision also decreases for both datasets as k increases.

# Information sheet

## CS246: Mining Massive Data Sets

**Assignment Submission** Fill in and include this information sheet with each of your assignments. This page should be the last page of your submission. Assignments are due at 11:59pm and are always due on a Thursday. All students (SCPD and non-SCPD) must submit their homework via Gradescope (<http://www.gradescope.com>). Students can typeset or scan their homework. Make sure that you answer each (sub-)question on a separate page. That is, one answer per page regardless of the answer length. Students also need to upload their code on Gradescope. Put all the code for a single question into a single file and upload it.

**Late Homework Policy** Each student will have a total of *two* late periods. *Homework are due on Thursdays at 11:59pm PT and one late period expires on the following Monday at 11:59pm PT.* Only one late period may be used for an assignment. Any homework received after 11:59pm PT on the Monday following the homework due date will receive no credit. Once these late periods are exhausted, any assignments turned in late will receive no credit.

**Honor Code** We strongly encourage students to form study groups. Students may discuss and work on homework problems in groups. However, each student must write down their solutions independently, i.e., each student must understand the solution well enough in order to reconstruct it by him/herself. 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 (GitHub/Google/previous year's solutions etc.) is considered an honor code violation. We check all the submissions for plagiarism. We take the honor code very seriously and expect students to do the same.

**Your name:** \_\_\_\_\_

**Email:** \_\_\_\_\_ **SUID:** \_\_\_\_\_

Discussion Group: \_\_\_\_\_

I acknowledge and accept the Honor Code.

(Signed) \_\_\_\_\_