

Cover Sheet

Assignment Submission Fill in and include this cover sheet with each of your assignments. Assignments are due at 11:59pm. All students (SCPD and non-SCPD) must submit their homeworks via GradeScope (<http://www.gradescope.com>). Students can typeset or scan their homeworks. Make sure that you answer each question on a separate page. Students also need to upload their code at <http://snap.stanford.edu/submit>. Put all the code for a single question into a single file and upload it. Please do not put any code in your GradeScope submissions.

Late Day Policy Each student will have a total of *two* free late periods. *One late period expires at the start of each class.* (Homeworks are usually due on Thursdays, which means the first late periods expire on the following Tuesday.) Once these late periods are exhausted, any assignments turned in late will be penalized 50% per late period. However, no assignment will be accepted more than *one* late period after its due date.

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 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: MANDI RAVI

Email: mandi1390@gmail.com SUID: mandi1390

CS246: Mining Massive Datasets Homework 2

Answer to Question 1(a)

$$T = R \times R^T$$

$$= \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 2 & 1 & 2 & 1 \\ 1 & 1 & 1 & 0 \\ 2 & 1 & 3 & 0 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

$T_{ii} \Rightarrow$ represents the no. of paths from u_i

$T_{ij} \Rightarrow$ Similarity of u_i paths from u_i and u_j in terms of distribution the paths finally mean no.

The diagonal elements of products transpose(R) and R AND R*transpose(R) contain diagonal elements of Q and P respectively which is also a reflection of the number of paths from ii and ui from the bipartite graph. The manipulation along with the division by square root terms of P and Q inverses will get us the representation of Si and Su

Answer to Question 1(b)

$$R * R^T = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 2 & 1 & 2 & 1 \\ 1 & 1 & 1 & 0 \\ 2 & 1 & 3 & 1 \\ 1 & 0 & 1 & 1 \end{bmatrix} \quad \text{diagonal}$$

$$S_I = R^T * R = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

$$Q = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix}$$

$$S_I = \begin{bmatrix} \frac{3}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{2}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} & \frac{3}{\sqrt{3}} & \frac{1}{\sqrt{3}} \\ \frac{2}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{3}{\sqrt{3}} \end{bmatrix}$$

$$Q^{-1} = \begin{bmatrix} \frac{1}{\sqrt{3}} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \frac{1}{\sqrt{3}} \end{bmatrix}$$

$$(Q^{-1})^{1/2} = \begin{bmatrix} \frac{1}{\sqrt{3}} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \frac{1}{\sqrt{3}} \end{bmatrix}$$

$$= \underline{P^T * R} * ((Q^{-1})^{1/2})$$

$$\therefore S_I = \begin{bmatrix} 1 & \frac{1}{\sqrt{3}} & \frac{2}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} & 1 & \frac{1}{\sqrt{3}} \\ \frac{2}{\sqrt{3}} & \frac{1}{\sqrt{3}} & 1 \end{bmatrix}$$

$$(Q^{-1})^{1/2} * (Q^{-1})^{1/2} = \begin{bmatrix} \frac{1}{\sqrt{3}} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \frac{1}{\sqrt{3}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{3}} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \frac{1}{\sqrt{3}} \end{bmatrix}$$

$$= \underline{P_S} \begin{bmatrix} 1 & \frac{1}{\sqrt{3}} & \frac{2}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} & 1 & \frac{1}{\sqrt{3}} \\ \frac{2}{\sqrt{3}} & \frac{1}{\sqrt{3}} & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 1 & 1 \\ 2 & 1 & 1 \end{bmatrix}$$

$$S_U = \begin{bmatrix} 1 & \frac{1}{\sqrt{2}} & \frac{2}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & 1 & \frac{1}{\sqrt{3}} & 0 \\ \frac{2}{\sqrt{2}} & \frac{1}{\sqrt{3}} & 1 & \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{3}} & 1 \end{bmatrix}$$

$$= \underline{(R^T * R) - P + I} * (P^{-1})^{1/2}$$

$$= \underline{((R^T * R) - Q + I)} * (Q^{-1})^{1/2}$$

The recommendation matrix inherently has the product of R and Si(or)Su matrices per definition. Hence, the definition of a recommendation matrix has variations based upon the type of collaborative filtering and user-user /item-item similarities accordingly. The more similar these are the higher the chances of recommendation for that given item for a given user.

Answer to Question 1(c)

$$R = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

4×3

User - user

$$X_{\text{user}} = \begin{bmatrix} 1 + \frac{1}{\sqrt{2}}(1) + \frac{1}{\sqrt{3}} & \cancel{\frac{\sqrt{2}}{\sqrt{3}}} & 1 + \frac{\sqrt{2}}{\sqrt{3}} + \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} + 1 + \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{3}} \\ \frac{\sqrt{2}}{\sqrt{3}} + \frac{1}{\sqrt{3}} + 1 & 1 & \frac{\sqrt{2}}{\sqrt{3}} + 1 + \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{3}} & 1/\sqrt{3} & \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{3}} + 1 \end{bmatrix}$$

$$= \begin{bmatrix} 2.525 & 0.816 & 2.525 \\ 1.517 + 0.816 & \cancel{0.517} & 1.517 + 0.816 \\ 2.228 & 0.527 & 2.228 \\ 2.29 + 0.527 & 1 & 2.39 \\ 1.28 & 0.527 & 2.28 \end{bmatrix}$$

item - item

$$X_{\text{item}} = \begin{bmatrix} 1 + \frac{2}{\sqrt{3}} & \frac{1}{\sqrt{3}} + \frac{1}{\sqrt{3}} & \frac{2}{3} + 1 \\ 1 + 0 + \cancel{\frac{1}{\sqrt{3}}} + \frac{2}{\sqrt{3}} & \frac{1}{\sqrt{3}} + \frac{1}{\sqrt{3}} & \frac{2}{3} + 1 \\ \cancel{1} & \frac{1}{\sqrt{3}} & \frac{2}{3} \\ 1 + 1/\sqrt{3} + 2/\sqrt{3} & \frac{1}{\sqrt{3}} + 1 + \frac{1}{\sqrt{3}} & \frac{2}{3} + \frac{1}{\sqrt{3}} + 1 \\ 2/\sqrt{3} & 1/\sqrt{3} & 1 \end{bmatrix} = \begin{bmatrix} 2.155 & 1.15 & 1.66 \\ 1 & 0.577 & 0.66 \\ 2.237 & 2.15 & 2.237 \\ 0.66 & 0.577 & 1 \end{bmatrix}$$

$$F_u = \cancel{S_I * R} = [(R^T - R - Q) + \epsilon]^{1/2}$$

$$T_I = R * S_u * \cancel{R} = [(R^T - R - Q) - P + \epsilon]^{1/2}$$

Answer to Question 1(d)

Top 5 user-user collaborative filtering :

Fox 28 News at 10pm

Family Guy

2009 NCAA Basketball Tournament

NBC 4 at Eleven

Two and a Half Men

Top 5 item-item :

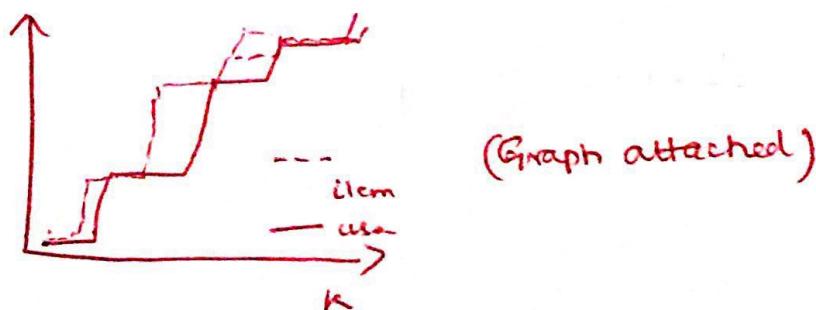
Fox 28 News at 10pm

Family Guy

2009 NCAA Basketball Tournament

NBC 4 at Eleven

Access Hollywood



Answer to Question 2(a)

$M^T M$ & $M M^T$ is symmetric, square and real

$$(M M^T)^T = (M^T)^T M^T = M M^T$$

$$\text{Hence, } (M M^T)^T = (M^T)^T (M^T)^T = M^T M$$

if $M \rightarrow p \times q$ $\Rightarrow (M M^T) \stackrel{\text{dim}}{\Rightarrow} p \times q \times q \times p$
 $M^T \rightarrow q \times p$ $\Rightarrow q \times p$

(or $q \times q$ in the case of $M^T M$)

if M is real $M M^T$ will also be real.

Answer to Question 2(b)

$$\{v \mid Av = \lambda v\}$$

If Eigen value of A is λ then $Av = \lambda v$ non trivial
($v \rightarrow$ non zero vector)
($n \times n$) Subspace

Determinant $|(\lambda I - A)|$ is called characteristic polynomial

Zeros of this polynomial are eigen values of A.

If M is non-singular, $M^{-1}M$ & $M^T M$ are similar

$$A = M^{-1}(M^T M)M$$

For a square matrix A of size $n \times n$

$$\det(A - zI) = (-1)^n (z - \lambda_1)(z - \lambda_2) \dots (z - \lambda_n) \quad (\lambda_1, \lambda_2, \dots \text{ eigen vals})$$

Also a square matrix has same characteristic polynomial as its transpose : $\det(A - zI) = \det((A - zI)^T) = \det(A^T - zI^T)$

$$= \det(A^T - zI)$$

Also similar matrices have same characteristic polynomial
same eigenvalues

If P is non-singular,

$$\begin{aligned} \det(A - zI) &= \det(P^{-1}P - zP^{-1}IP) \det(P) \\ &= \det(P^{-1}AP - zP^{-1}IP) \det(P^{-1}AP - zI) \end{aligned}$$

Hence, A and $P^{-1}AP$ have same characteristic polynomial

Hence, same eigen values. Eigen vectors could be scale multiple of one other and could be transformed into other

Answer to Question 2(c)

$$M^T M = Q \Lambda Q^T$$

Answer to Question 2(d)

$$M = U \Sigma V^T$$

$$MTM = (U \Sigma V^T)^T (U \Sigma V^T)$$

$$MTM = (V \Sigma^T U^T) (U \Sigma V^T)$$

Since $U^T U = I$

$$\boxed{MTM = V \Sigma^2 V^T}$$

Answer to Question 2(e)

$$U = \begin{bmatrix} -0.2785, 0.5 \\ -0.27854, -0.5 \\ -0.6499, 0.5 \\ -0.6499, 0.5 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 7.6157 & 0 \\ 0 & 1.414 \end{bmatrix}$$

$$V^T = \begin{bmatrix} -0.7071 & -0.7071 \\ -0.7071 & 0.7071 \end{bmatrix} \quad V =$$

$$\text{Evals} = [2., 58.]$$

$$\text{Evec} = \begin{bmatrix} 0.7071 & -0.7071 \\ 0.7071 & 0.7071 \end{bmatrix}$$

We can see that Evec and V are same
(and linearly transformable from one other)

Eigenvalues from PCA(MTM) are squares
of singular values of SVD

Answer to Question 3(a)

Centroid based

$$C_1 \quad \frac{0.68 + 0.72 + 0.71}{3} \\ = \frac{2.11}{3} = 0.73$$

$$C_2 \quad \frac{0.87 + 0.92 + 0.91}{3} \\ = 0.91$$

$$C_3 \quad \frac{1.01 + 1.22 + 1.45}{3} \\ = 1.23$$

C_1 and C_2 will be merged

Centroid based merging will form 2 clusters such that one of the clusters correspond to low risk of lung cancer and other represent high risk of lung cancer.

Closest member based

$$C_1 \& C_2 \text{ dist.} = 0.08$$

$$C_2 \& C_3 \text{ dist.} = 0.02$$

C_2 and C_3 will be merged

Answer to Question 3(b)

Centroid based



$$C_1 = 0.73$$

$$C_2 = \frac{0.87 + 0.92}{2} \\ = 0.90$$

$$C_3 = 1.23$$

C_1 and C_2 will be merged

The centroid based merging strategy is stable
without point $p_6 = 0.99$ and existing set of clusters

Closest Member based

$$C_1 \& C_2 \text{ dist.} = 0.08$$

$$C_2 \& C_3 \text{ dist.} = 0.09$$

C_1 & C_2 will be
merged

Answer to Question 3(c)

I would prefer a stable merge strategy over a unstable one because a stable merge strategy would have the chance of being less affected by (or not at all affected) by the % of miscalculated NSR in the dataset and its clustered sets will not be decisive entirely based upon the miscalculated points or not at all dependent on them. Hence, we could get a expect a more reliable merge strategy from a stable one.

Answer to Question 4(a)

1. Graph plot

2. $\xrightarrow{\text{source } C_1}$
 $'C_0 = 6.236603453064109E8 -$
 initial
 $'C_{10} = \underline{4.584906561919809E8}$
 $\downarrow_{10\text{th}}$

using C_2 : (% change)
 Percentage change
 after 10 iterations

$$^2C_0 = 4.3874779002791625E8 -$$

$$^2C_{10} = 1.0223720331799605E8$$

$$\frac{^2C_0 - ^2C_{10}}{^2C_0} = 76.69\% \text{ decrease}$$

using C_1

$$\frac{'C_0 - 'C_{10}}{'C_0}$$

$$= \frac{165169689}{165169689}$$

$$= 26.48\% \text{ decrease}$$

When using Euclidean distance for k-means and calculating the cost function, far-away assignment of initial clusters has a better chance of starting with a relatively lower cost and also in having better algorithm efficiency in terms of percentage decrease of cost function. This is because of basic algorithm of k-means which tends to cluster points closer into a cluster and further away from distance cluster.

Starting with faraway centroids reflects the assumption of far-away clusters which can generally help in good cost reduction unless we have many outliers.

(Percentage change after 10 iterations)

Answer to Question 4(b)

q. Graph plotted 2.

$${}^1 C_0 = 550117.142$$

c₁.txt

$${}^1 C_{10} = 447494.3682$$

$$\rightarrow \frac{{}^1 C_0 - {}^1 C_{10}}{{}^1 C_0} = 18.655\%$$

decrease

$${}^2 C_0 = 1433739.031$$

c₂.txt

$${}^2 C_{10} = 694587.92525$$

$$\rightarrow \frac{{}^2 C_0 - {}^2 C_{10}}{{}^2 C_0} = 51.55\%$$

Since the points in c₂ have been set up based upon farthest distance in L₂ space they need not be farthest ones according to Manhattan distance. So, initializing with lower cost seems more or equally likely with c₁ (random) initialization. But the algorithm

8 c₂

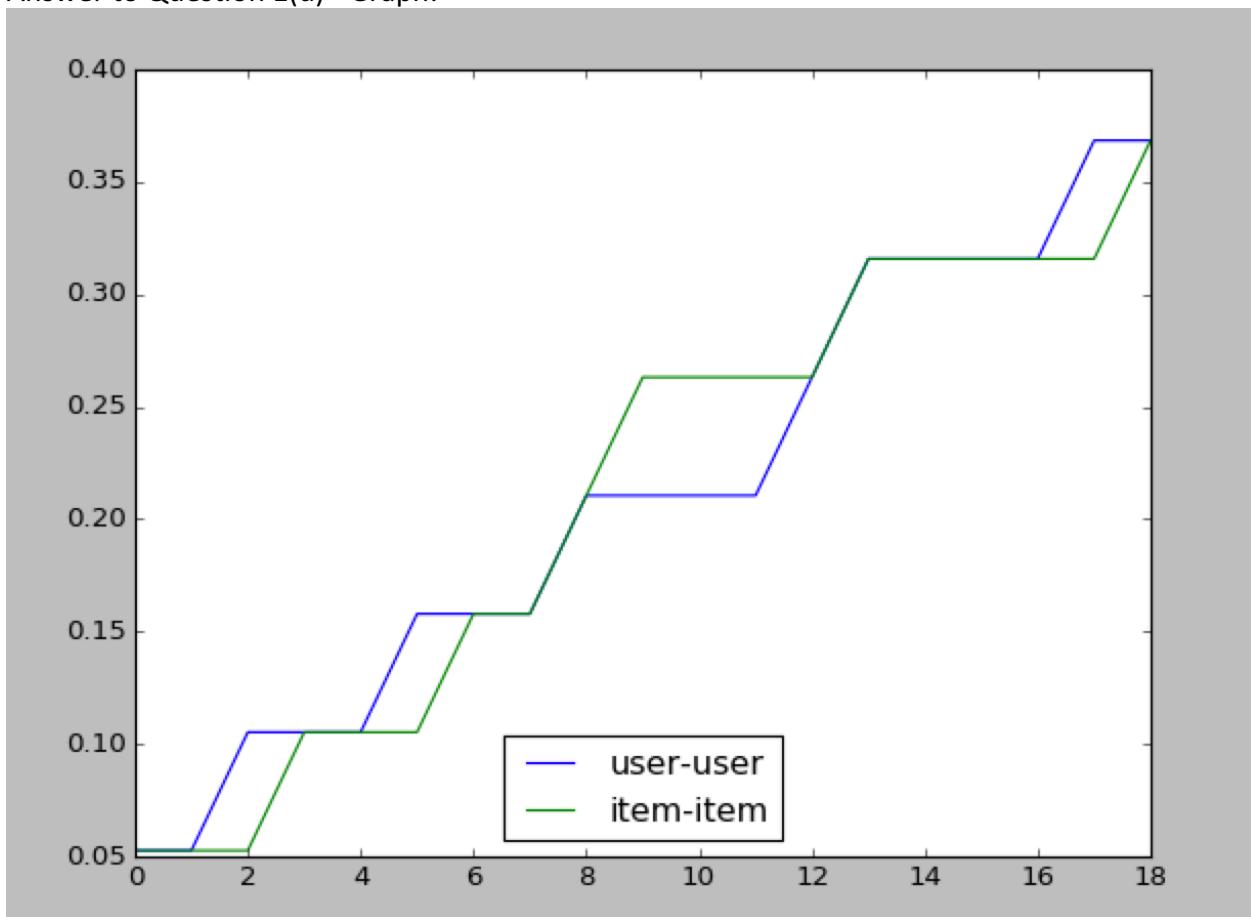
efficiency in terms of reducing the cost still

is better by starting with c₂ (like 4a).

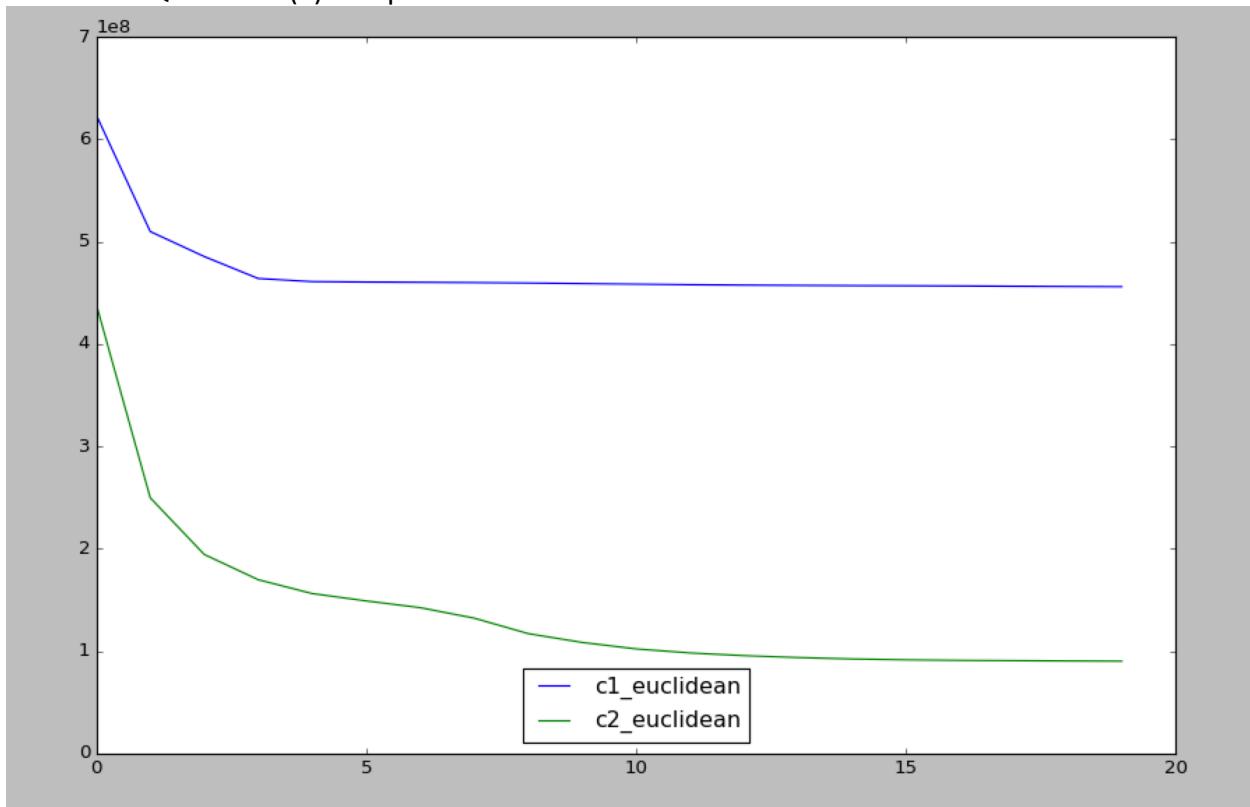
Because, it helps sustain the underlying concept of k-means which clusters closer points together and farther away points in different clusters. unless there is lot of outliers

This should be an ¹⁴ efficient start for max reduction of cost.

Answer to Question 1(d) - Graph:



Answer to Question 4(a)- Graph:



Answer to Question 4(b)- Graph:

