# MATH444 Midterm

## Jiasen Zhang

## 1 Problem 1

Before processing the dataset "cities.mat", notice that the ratings of different attributes have different orders of magnitudes. For example, the average rating of "climate" is about 500 while the average rating of "housing" is about 8000. So instead of centering the data, I normalize the data so that the mean value is 0 and the standard variance is 1.

Computing the first two principal components of the normalized data and plot the absolute values of them in Figure 1, we can clearly indentify the largest components. The first two singular values are 33.5 and 20.0, that means the first PCA vector plays much more important role than the second one. So we focus on the first PCA vector.

According to Figure 1, the two largest components of the first PCA vetor correspond to the attributes "health" and "arts". So we can conclude that health and arts are the predominant factor when differentiating the quality of life in the data.
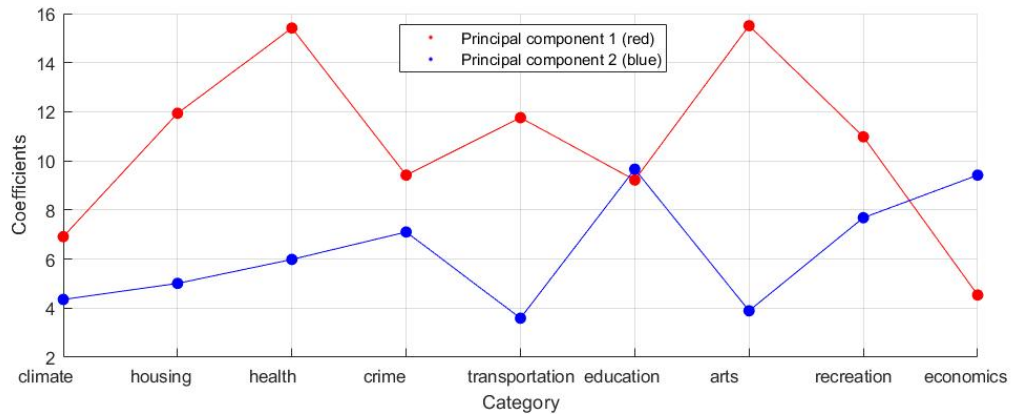


Figure 1: Absolute value of PCA components

## 2    Problem 2

To check if there are any natural clusters, I plot the first two principal components of the transpose of the data. As shown in Figure 2, there is a natural cluster: most data points tend to cluster together on the left. In fact, the data points corresponding to cities with lower ratings are closer to the left dense area. The data point at the rightmost corresponds to New York City, which has the highest average rating.
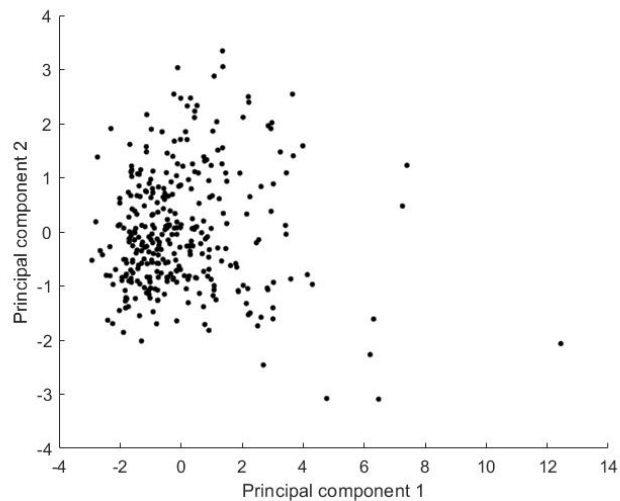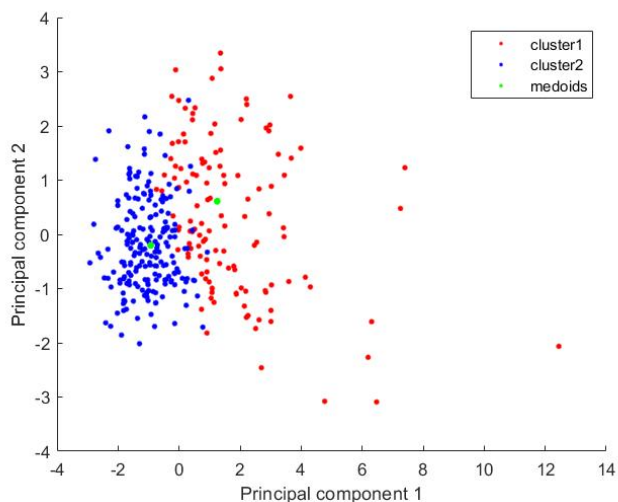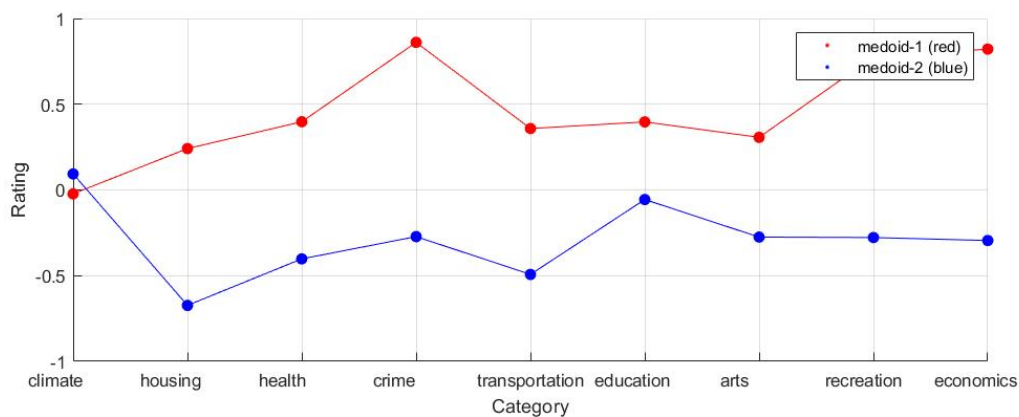
Figure 2: PCA of cities.mat

Running k-medoids with k=2, the results are:

|           | color | medoid city      |
|-----------|-------|------------------|
| cluster 1 | red   | Phoenix, AZ      |
| cluster 2 | blue  | Evansville, IN-KY |

As shown in Figure 3 A, when k=2, there are two clusters: red and blue. As we mentioned above, the red cluster should have higher ratings generally because it's farther from the dense area. In Figure 3 B, as we expected, the rating of medoid 1 corresponding to "Phoenix, AZ" has higher overall rating than the rating of medoid 2 corresponding to the less famous "Evansville, IN-KY".
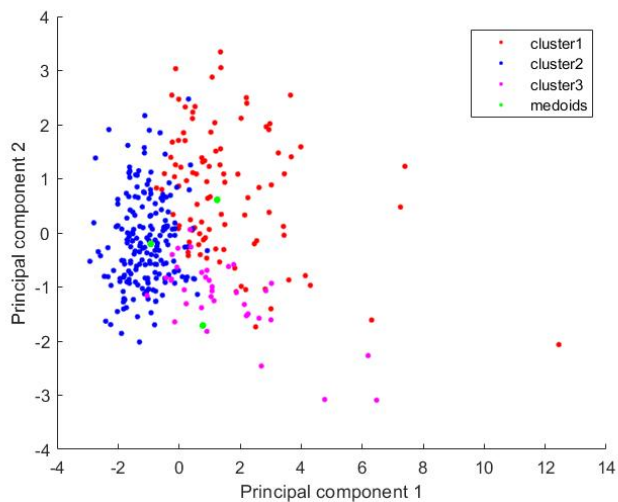
A : k-medoids with k=2. Green points are medoids.
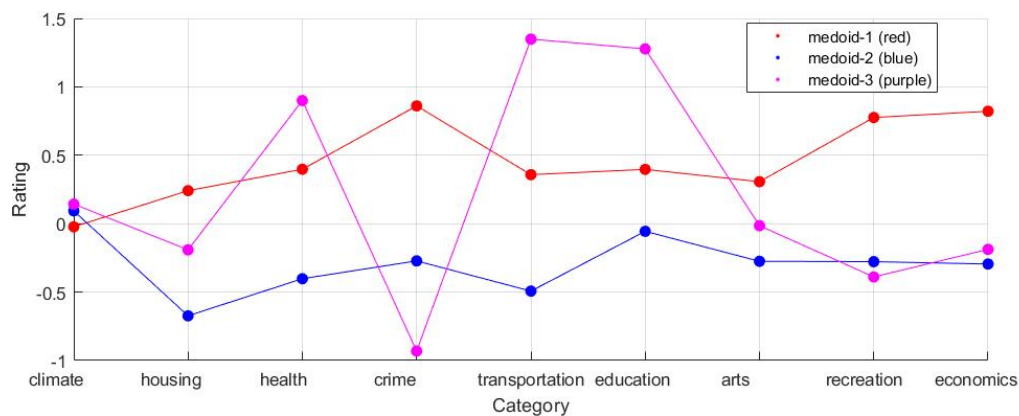
B : ratings of the medoids

Figure 3

Running k-medoids with k=3, the results are:

|          | color  | medoid city                     |
|----------|--------|---------------------------------|
| cluster 1 | red    | Phoenix, AZ                     |
| cluster 2 | blue   | Evansville, IN-KY               |
| cluster 3 | purple | Harrisburg-Lebanon-Carlisle, PA |

As shown in Figure 4 A, when k=3, the third purple cluster is between the blue and red clusters. In Figure 4 B, most attribute ratings of purple cluster is lower than red cluster and higher than blue cluster. Generally the ratings of purple cluster is between the red and blue clusters, which is corresponding to the result in Figure 4 A.
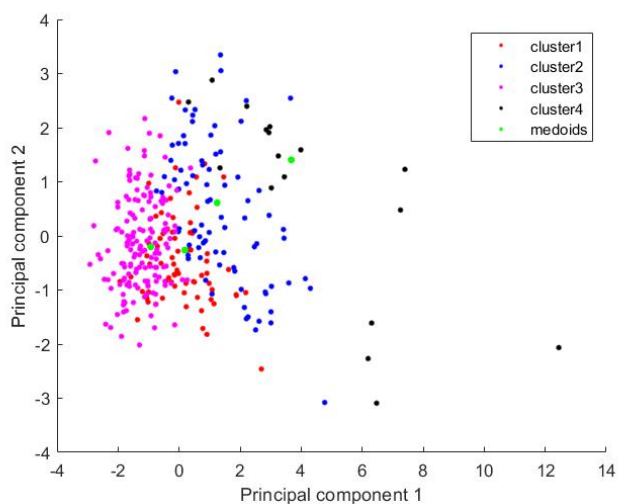


A : k-medoids with k=3. Green points are medoids.
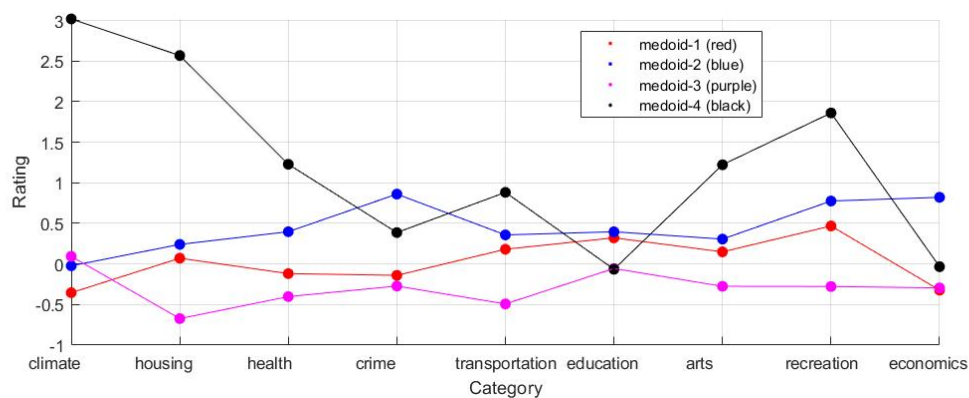


B : ratings of the medoids

Figure 4

Running k-medoids with k=4, the results are:

|          | color  | medoid city       |
|----------|--------|-------------------|
| cluster 1 | red    | Kenosha, WI       |
| cluster 2 | blue   | Phoenix, AZ       |
| cluster 3 | purple | Evansville, IN-KY |
| cluster 4 | black  | San Diego, CA     |

With the results above, in Figure 5 A we can easily know the four clusters from low to high ratings are purple, red, blue and black, and it's proved by Figure 5 B. The result is reasonable: we know "Kenosha, WI" and "Evansville, IN-KY" are not famous and therefore shouldn't have high ratings. Both "San Diego, CA" and "Phoenix, AZ" are famous and generally Dan Diego is more popular than Phoenix mainly for better climate.



A : k-medoids with k=4. Green points are medoids.



B : ratings of the medoids

Figure 5

# 3 Problem 3

When setting the rank as 2, for the most frequent result, the two cities that are closest to the feature vectors are Pheonix and Janesville. I think they are representative. As shown in Figure 6, Pheonix is better than Janesville in all attributes, it can represent the big cities with high quality of life. Janesville is a small city in the Wisconsin State with bad climate. It represents small and remote cities with low quality of life.

|        | color | city                  |
|--------|-------|-----------------------|
| city 1 | red   | Phoenix, AZ           |
| city 2 | blue  | Janesville-Beloit, WI |



Figure 6: NMF: rank = 2

When the rank increases, the situation will be complicated. The result cities will represent different levels of quality of life. For example, when the rank is 3, the three cities are Phoenix, Sheboygan and Birmingham. According to Figure 7, the cities from high to low levels are Phoenix, Birmingham and Sheboygan.

|        | color  | city            |
|--------|--------|-----------------|
| city 1 | red    | Phoenix, AZ     |
| city 2 | blue   | Sheboygan, WI   |
| city 3 | purple | Birmingham, AL  |



Figure 7: NMF: rank = 3

# 4 Problem 4

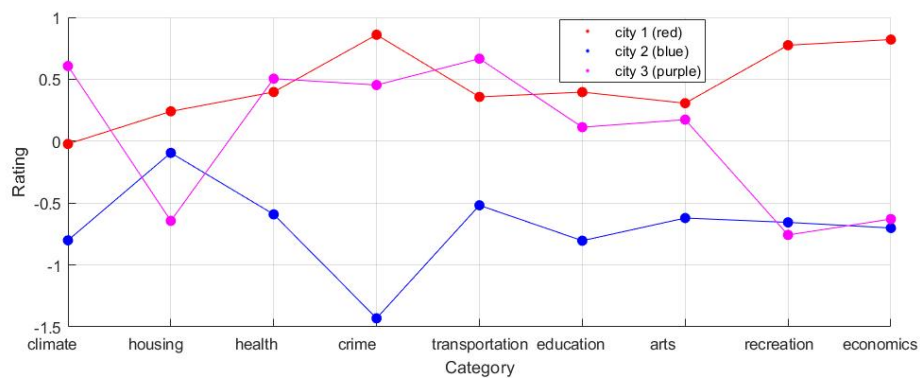I run SOM algorithm with a 10 by 10 lattice and get 100 prototypes, then put the name of the cities that are closest to the prototypes in a 10 by 10 lattice. In addition, to help illustrate I compute the average ratings of the cities and put them in a lattice. The average rating can partly reflect the overall quality of a city.

The result in Figure 8 looks reasonable. As we can see, the cities on the upper left are big cities with higher ratings and quality of life. The lower right cities are less famous with lower quality of life. From upper left to lower right, the average ratings of cities decrease.

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| NewYork, NY | Boston, MA | Baltimore, MD | St.Louis, MO-IL | Richmond-Petersburg, VA | Springfield, IL | KansasCity, MO | Birmingham, AL | Fayetteville, NC | Chico, CA |
| LosAngeles, LongBeach, CA | Newark, NJ | Minneapolis-St.Paul, MN-WI | Providence, RI | Albany-Troy, NY | Champaign-Urbana-Rantoul, IL | Kalamazoo, MI | Nashville, TN | Shreveport, LA | NewBritain, CT |
| Newark, NJ | Bridgeport-Milford, CT | Milwaukee, WI | Buffalo, NY | Trenton, NJ | Kalamazoo, MI | Kalamazoo, MI | DesMoines, IA | Topeka, KS | Rochester, MN |
| Norwalk, CT | Monmouth-Ocean, NJ | Norfolk-VirginiaBeach-NewportNews, | Norfolk-VirginiaBeach-NewportNews, | Bellingham, WA | Spokane, WA | Pawtucket-Woonsocket-Attleboro, RI- | Alton, GraniteCity, IL | Lincoln, NE | St.Cloud, MN |
| SanDiego, CA | SanDiego, CA | Tuscon, AZ | Tacoma, WA | Redding, CA | NiagaraFalls, NY | Lynchburg, VA | Lancaster, PA | Lafayette, IN | Wausau, WI |
| SantaRosa-Petaluma, CA | Tuscon, AZ | Tuscon, AZ | Fresno, CA | FallRiver, MA-RI | BattleCreek, MI | Lima, OH | Canton, OH | Dubuque, IA | Waterloo-CedarFalls, IA |
| Salem, MA | Portland, ME | Orlando, FL | Sacramento, CA | LittleRock, NorthLittleRock, AR | Pueblo, CO | Lawrence, KS | Lima, OH | TerreHaute, IN | Dubuque, IA |
| Poughkeepsie, NY | ColoradoSprings, CO | FortWorth-Arlington, TX | Bakersfield, CA | LasCruces, NM | Yakima, WA | Chattanooga, TN-GA | Columbus, GA-AL | Columbus, GA-AL | Kokomo, IN |
| FortWaltonBeach, FL | Lawrence-Haverhill, MA-NH | Lafayette, LA | BatonRouge, LA | Beaumont-PortArthur, TX | Springfield, MO | Sherman-Denison, TX | Columbus, GA-AL | Bloomington, IN | York, PA |
| Visalia-Tulare-Porterville, | FortCollins-Loveland, CO | Sarasota, FL | CorpusChristi, TX | SanAngelo, TX | WichitaFalls, TX | Macon, WarnerRobbins, GA | Albany, GA | Athens, GA | Fitchburg-Leominster, MA |

A : name of the cities
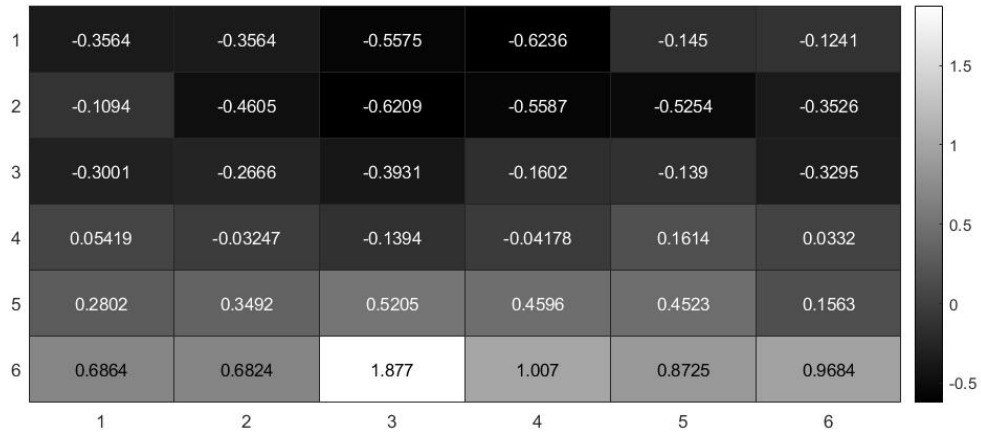


B : average ratings

Figure 8

# 5   Problem 5

I run SOM algorithm with a 6 by 6 lattice and get 36 prototypes, then put the name of the cities that are closest to the prototypes in a 6 by 6 lattice.

The result in Figure 9 looks reasonable too. This time the better cities gather at the edge rather than corner. Cities with similar quality are still closer to each other.

When running SOM algorithm, I found that when the size of lattice is smaller (less nodes), there will be less prototypes. As a result, the dissimilarity of the prototypes and their closest cities will be larger, and the map will be organized better. That's reasonable becasue it's easier to cluster less data.

| | | | | | |
|---|---|---|---|---|---|
| Hagerstown,MD | Hagerstown,MD | Columbus,GA-AL | Albany,GA | LasCruces,NM | Vineland-Millville-Bridgeton,NJ |
| Roanoke,VA | Lima,OH | Montgomery,AL | Macon,WarnerRobbins,GA | Monroe,LA | NewBedsford,MA |
| NiagaraFalls,NY | Hamilton-Middletown,OH | GreenBay,WI | Bradenton,FL | Lowell,MA-NH | Visalia-Tulare-Porterville,CA |
| Birmingham,AL | Lansing-EastLansing,MI | Lincoln,NE | Wichita,KS | ColoradoSprings,CO | Sarasota,FL |
| Springfield,IL | Wilmington,DE-NJ-MD | Buffalo,NY | Phoenix,AZ | Tuscon,AZ | LakeCounty,IL |
| Providence,RI | Middlesex-Somerset,Hunterdon,NJ | Boston,MA | Cleveland,OH | SantaBarbara-SantaMaria-Lompoc,CA | Norwalk,CT |

A : name of the cities



B : average ratings

Figure 9

# 6 Matlab Code

The code of PCA, k-medoids and SOM are completly the same as the previous homeworks.

## 6.1 NMF function

```matlab
function [W, H] = NMF(X, k)
%Jiasen Zhang: NMF

% input:
% k = rank of W and H
% X = nonnegative matrix

% output: W[n,k], H[k,p]

tau = 0.001; % tolerance
tmax = 100000;
[n,p]=size(X);

% initialization
W = abs(rand(n,k));
H = abs(rand(k,p));
for j=1:k
    W(:,j)=W(:,j)/norm(W(:,j),inf);
end

% iteration
for t=1:tmax
    Xc = W*H;
    Hn = (W'*X)./(W'*Xc).*H; % new H
    Xc = W*Hn; % new Xc
    Wn = (X*Hn')./(Xc*Hn').*W; % new W
    for j=1:k
        Wn(:,j)=Wn(:,j)/norm(Wn(:,j),inf);
    end
    dWH = norm(Wn-W,'fro')/norm(W,'fro')+norm(Hn-H,'fro')/norm(H,'fro');
    W=Wn;
    H=Hn;
    fprintf('%d, %e \n',t,dWH);
    if dWH<tau
        break;
    end
end
end
```

## 6.2 Problem 1

```matlab
clear all;clc;

load cities;
X = ratings;
% % % normalize X
stdX = std(X,0,1);
rateBar = mean(X,1);
for i=1:9
    X(:,i)=(X(:,i)-rateBar(i))/stdX(i);
end

[u,d,v] = svd(X);
diag(d) % singular values
d=diag(d.^2);
proportion=d/sum(d) % importance of PCA components
```

```
16  z = u(:,1:2)'*X; % PCA components
17
18  % plot PCA components
19  absz=abs(z);
20  figure();hold on;grid on;
21  scatter((1:9),absz(1,:),'r','k.','SizeData',500);
22  scatter((1:9),absz(2,:),'b','k.','SizeData',500);
23  plot(absz(1,:),'r');
24  plot(absz(2,:),'b');
25  xlabel('Category');
26  ylabel('Coefficients');
27  set(gca,'xticklabel',categories(1:9,:));
28  legend('Principal component 1 (red)','Principal component 2 (blue)');
```

## 6.3   Problem 2

```
1   clear all;clc;
2
3   load cities;
4   X = ratings;
5   %normalize X
6   stdX = std(X,0,1);
7   rateBar = mean(X,1);
8   for i=1:9
9       X(:,i)=(X(:,i)-rateBar(i))/stdX(i);
10  end
11  X=X';
12  [n,p]=size(X);
13
14  % Get distance matrix D
15  D = zeros(p,p);
16  for i=1:p
17      for j=1:i
18          D(i,j)=norm(X(:,i)-X(:,j),2);
19          D(j,i)=D(i,j);
20      end
21  end
22
23  k = 4;
24  tau = 1e-6;
25  ninit=20;
26  [I_assign, I_bar] = my_k_medoids(k, D, tau, ninit);
27
28  [u,d,v] = svd(X);
29  z = u(:,1:2)'*X;
30
31  figure();hold on;
32  % plot the two clusters
33  scatter(z(1,I_assign==1),z(2,I_assign==1),'r','k.','SizeData',100);
34  scatter(z(1,I_assign==2),z(2,I_assign==2),'b','k.','SizeData',100);
35  scatter(z(1,I_assign==3),z(2,I_assign==3),'m','k.','SizeData',100);
36  scatter(z(1,I_assign==4),z(2,I_assign==4),'k','k.','SizeData',100);
37  for L=1:k
38      scatter(z(1,I_bar(L)),z(2,I_bar(L)),'g','k.','SizeData',200);
39  end
40  legend('cluster1','cluster2','cluster3','cluster4','medoids');
41  xlabel('Principal component 1');
42  ylabel('Principal component 2');
43
44  % plot the average rating of two clusters
45  figure();hold on;grid on;
46  scatter((1:9),X(:,I_bar(1)),'r','k.','SizeData',500);
47  scatter((1:9),X(:,I_bar(2)),'b','k.','SizeData',500);
48  scatter((1:9),X(:,I_bar(3)),'m','k.','SizeData',500);
49  scatter((1:9),X(:,I_bar(4)),'k','k.','SizeData',500);
```

```
50   plot(X(:,I_bar(1)),'r');
51   plot(X(:,I_bar(2)),'b');
52   plot(X(:,I_bar(3)),'m');
53   plot(X(:,I_bar(4)),'k');
54   legend('medoid-1 (red)','medoid-2 (blue)','medoid-3 (purple)','medoid-4 (black)');
55   xlabel('Category');
56   ylabel('Rating');
57   set(gca,'xticklabel',categories(1:9,:));
58
59   % show the medoid cities
60   for L=1:k
61       fprintf(names(I_bar(L),:));
62       fprintf('\n');
63   end
```

## 6.4   Problem 3

```
1    clear all;clc;
2
3    load cities;
4    X = ratings;
5    %normalize X
6    stdX = std(X,0,1);
7    rateBar = mean(X,1);
8    for i=1:9
9        X(:,i)=(X(:,i)-rateBar(i))/stdX(i);
10   end
11   X=X';
12   [n,p]=size(X);
13
14   k=3;
15   [W,H] = NMF(X,k);
16
17   %normalize W
18   stdW = std(W,0,2);
19   meanW = mean(W,2);
20   for i=1:9
21       W(i,:)=(W(i,:)-meanW(i))/stdW(i);
22   end
23
24
25   % W contains feature vectors
26   % find closest cities
27   num = zeros(k,1);
28   for t=1:k
29       num(t)=1;
30       for tt=2:p
31           Xt = X(:,tt);% city to be compared
32           Xnum = X(:,num(t));% current chosen city
33           if norm(Xt-W(:,t))<norm(Xnum-W(:,t))
34               num(t) = tt;
35           end
36       end
37   end
38   names(num(:),:) % show the cities
39
40   figure();hold on;grid on;
41   scatter((1:9),X(:,num(1)),'r','k.','SizeData',500);
42   scatter((1:9),X(:,num(2)),'b','k.','SizeData',500);
43   scatter((1:9),X(:,num(3)),'m','k.','SizeData',500);
44   %scatter((1:9),X(:,num(4)),'k','k.','SizeData',500);
45   plot(X(:,num(1)),'r');
46   plot(X(:,num(2)),'b');
47   plot(X(:,num(3)),'m');
48   %plot(X(:,num(4)),'k');
```

```
49  %legend('city 1 (red)','city 2 (blue)');
50  legend('city 1 (red)','city 2 (blue)','city 3 (purple)');
51  %legend('city 1 (red)','city 2 (blue)','city 3 (purple)','city 4 (black)');
52  xlabel('Category');
53  ylabel('Rating');
54  set(gca,'xticklabel',categories(1:9,:));
```

## 6.5   Problem 4 and 5

```
1   clear all;clc;
2
3   load cities;
4   X = ratings;
5   %normalize X
6   stdX = std(X,0,1);
7   rateBar = mean(X,1);
8   for i=1:9
9       X(:,i)=(X(:,i)-rateBar(i))/stdX(i);
10  end
11  X=X';
12  [n,p]=size(X);
13
14  % run SOM
15  K=36;
16  N = round(sqrt(K));
17  Tmax = 600*K;
18  T0 = 2000;
19  M = SOM(K, X, Tmax, T0);
20
21  name0 = char(K,43);
22  rate0 = zeros(N,N);
23  for i=1:N
24      for j=1:N
25          L = N*(i-1)+j;
26          % find data closest to M(:,L)
27          numForL = 1;
28          for tt=1:p
29              if norm((M(:,L)-X(:,tt)),2)<norm((M(:,L)-X(:,numForL)),2)
30                  numForL = tt;
31              end
32          end
33          name0(L,1:43)=names(numForL,:);
34          rate0(i,j)=mean(X(:,numForL));
35      end
36  end
37  name0
38  %contourf(rate0')
39  heatmap(rate0');colormap(gray);
```