

CPSC 340: Assignment5

Yiwei Hou(84435156)

Qiong Zhang(85896158)

1 Principal Component Analysis

1.1 PCA by Hand

1. The original data set is $X = \begin{pmatrix} -2 & -1 \\ -1 & 0 \\ 0 & 1 \\ 1 & 2 \\ 2 & 3 \end{pmatrix}$. After centralization, $\tilde{X} = \begin{pmatrix} -2 & -2 \\ -1 & -1 \\ 0 & 0 \\ 1 & 1 \\ 2 & 2 \end{pmatrix}$, which

means after centralization, the observation for the first feature and the second feature are the same. When plot the observations where x_1 is the x-axis and x_2 is the y-axis, then the observations fall on the line $x_2 = x_1$. Thus, we have $w_1 = w_2$ and the first principal component is $W = (\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2})$ since $\|W\| = 1$.

2. Since point (3,3) also falls on the line $x_2 = x_1$, thus the reconstruction error is 0.
3. The reconstruction error for point (3,4) is the distance of point (3,4) to the line $x_2 = x_1$, thus the reconstruction error equals $\frac{\sqrt{2}}{2}$.

1.2 Data Visualization

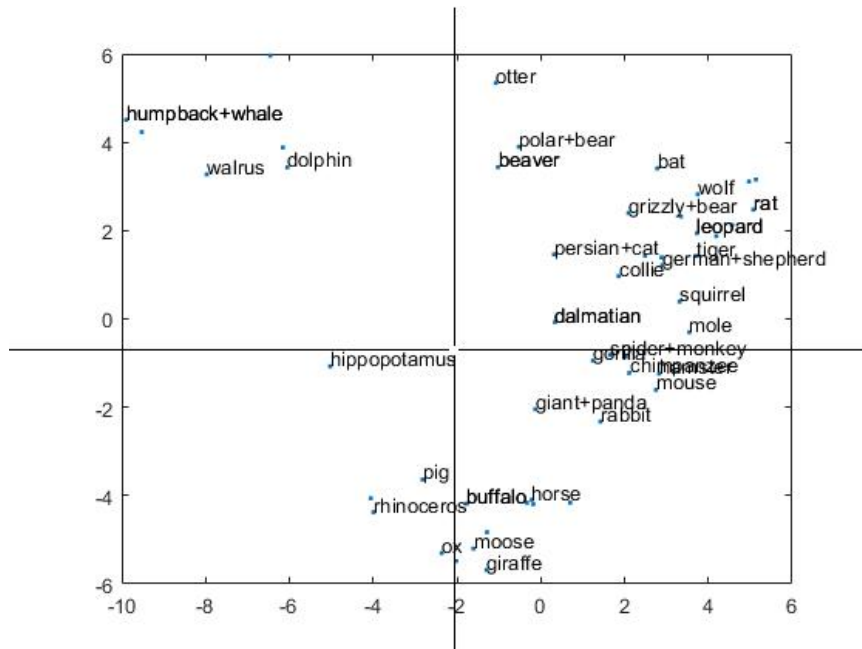


Figure 1: Scatter plot for the weight on the first two principal components

The demo:

```

1 load animals.mat
2 [n,d] = size(X);
3 X = standardizeCols(X);
4 model=dimRedPCA(X,2);
5 Z=model.compress(model,X);
6 plot(Z(:,1),Z(:,2),'.' );
7 gname(animals);

```

1.3 Data Compression

1. The proportion of the variances can be explained by two-dimensional representation is calculated by $1 - \frac{\|ZW-X\|_F^2}{\|X\|_F^2} = 16.45\%$.
2. By taking different values for k , we find out 14 PCS are required to explain 50% of the variance in the data.

2 PCA Generalizations

2.1 Robust PCA

```

1 function [model] = dimRedRPCA(X,k)
2
3 [n,d] = size(X);
4 epsilon=0.0001;
5
6 % Subtract mean
7 mu = mean(X);
8 X = X - repmat(mu,[n 1]);
9
10 % Initialize W and Z
11 W = randn(k,d);
12 Z = randn(n,k);
13
14 R = Z*W-X;
15 f = sum(sum(sqrt(R.^2+epsilon)));
16 for iter = 1:50
17     fOld = f;
18
19     % Update Z
20     Z(:) = findMin(@funObjZ,Z(:),10,0,X,W);
21
22     % Update W
23     W(:) = findMin(@funObjW,W(:),10,0,X,Z);
24
25     R = Z*W-X;
26     f = sum(sum(sqrt(R.^2+epsilon)));
27     fprintf('Iteration %d, loss = %.5e\n',iter,f);
28
29     if fOld - f < 1e-4
30         break;
31     end

```

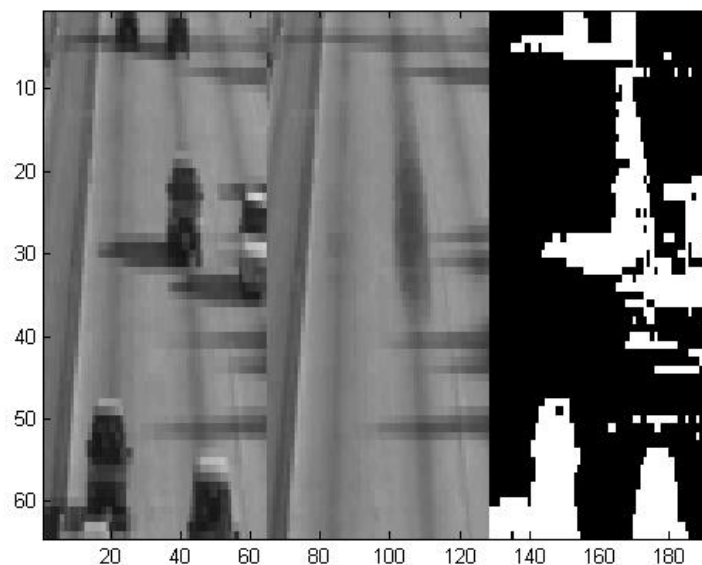
```
32 end
33
34 model.mu = mu;
35 model.W = W;
36 model.compress = @compress;
37 model.expand = @expand;
38 end
39
40 function [Z] = compress(model,X)
41 [t,d] = size(X);
42 mu = model.mu;
43 W = model.W;
44 k = size(W,1);
45
46 X = X - repmat(mu,[t 1]);
47 % We didn't enforce that W was orthogonal so we need to optimize to find Z
48 Z = zeros(t,k);
49 Z(:) = findMin(@funObjZ,Z(:),100,0,X,W);
50 end
51
52 function [X] = expand(model,Z)
53 [t,d] = size(Z);
54 mu = model.mu;
55 W = model.W;
56
57 X = Z*W + repmat(mu,[t 1]);
58 end
59
60 function [f,g] = funObjW(W,X,Z)
61 % Resize vector of parameters into matrix
62 d = size(X,2);
63 k = size(Z,2);
64 W = reshape(W,[k d]);
65 epsilon=0.0001;
66 % Compute function value
67 R = Z*W-X;
68 f = sum(sum(sqrt(R.^2+epsilon)));
69
70 % Compute derivative with respect to each residual
71 dR = R./sqrt(R.^2+epsilon);
72
73 % Multiply by Z' to get elements of gradient
74 g = Z'*dR;
75
76 % Return a vector
77 g = g(:);
78 end
79
80 function [f,g] = funObjZ(Z,X,W)
81 % Resize vector of parameters into matrix
82 n = size(X,1);
83 k = size(W,1);
84 epsilon=0.0001;
85 Z = reshape(Z,[n k]);
86
```

```

87 % Compute function value
88 R = Z*W-X;
89 f = sum(sum(sqrt(R.^2+epsilon)));
90
91 % Compute derivative with respect to each residual
92 dR = R./sqrt(R.^2+epsilon);
93
94 % Multiply by W' to get elements of gradient
95 g = dR*W';
96
97 % Return a vector
98 g = g(:);
99 end

```

The output of the highway example by robust PCA is given by



2.2 L1-Regularized and Binary Latent-Factor Models

1. Higher value of w will encourage the sparsity of W and does not directly affect the sparsity of Z . λ_Z controls the shrinkage of Z but does not directly affect the sparsity of W .
2. The fundamental trade-off in machine learning is a trade-off between how small you can make the training error and how well training error approximates the test error. As mentioned in 1, the higher the λ_Z , the result is more robust which means the training error approximates the test error well, but the training error becomes larger. In terms of the effect of the number of principal component k on fundamental trade, the higher the k value, the smaller the training error. When $k = d$, the error is smallest. However, the large k values may lead to over-fitting which means the test error would be high.
3. If λ_W is zero, the answer in (2) will not change as λ_W only affect the sparsity of W and with zero w , we don't penalize large coefficients. But the model fit is determined by the product of Z and W . It would not directly affect the training and test error.

4.

$$f(Z, W) = \sum_{(i,j) \in \mathcal{R}} \frac{1}{1 + \exp(-x_{ij} w_j^T z_i)} + \lambda_W \sum_{j=1}^d \|w_j\|_1 + \frac{\lambda_Z}{2} \sum_{i=1}^n \|z_i\|^2$$

3 Multi-Dimensional Scaling

3.1 ISOMAP

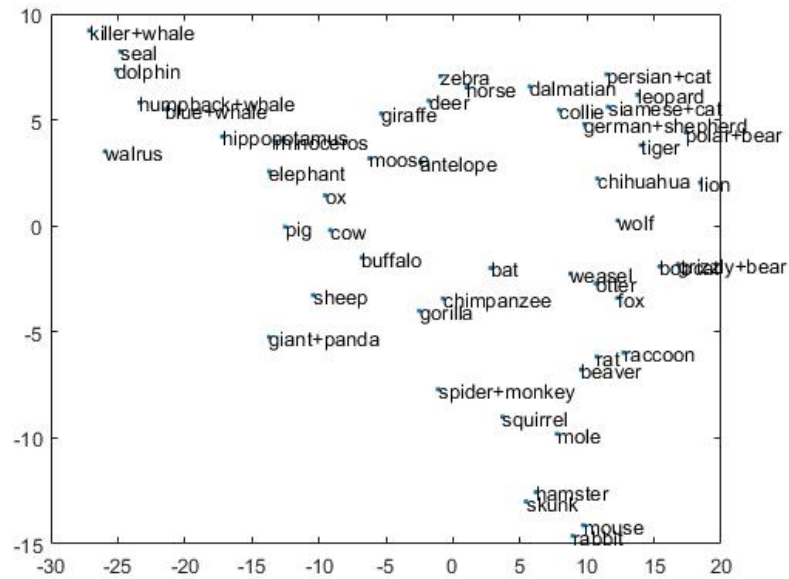
We can implement the *visualizeMDS* function by a different distance function given by

```

1 % Compute all distances
2 D = X.^2*ones(d,n) + ones(n,d)*(X') .^2 - 2*X*X';
3 D = sqrt(abs(D));
4
5 % Intialize the KNN graph
6 G=zeros(n,n);
7 % Find the 3 nearest neighbours
8 [~,pos]=sort(D,1);
9 neighbour=pos(2:4,:);
10 % include an edge i to j if i is a KNN of j or if j is a KNN of i
11 for i=1:n;
12     for j=neighbour(:,i);
13         G(i,j)=D(i,j);
14         G(j,i)=D(j,i);
15     end
16 end
17 %computes the approximate geodesic distance
18 D=zeros(n);
19 for i=1:n;
20     for j=1:n;
21         D(i,j)=dijkstra(G,i,j);
22     end
23 end

```

The final result is shown below.



3.2 ISOMAP with Disconnected Graph

```

1 %computes the approximate geodesic distance
2 D=zeros(n);
3 for i=1:n;
4     for j=1:n;
5         D(i,j)=dijkstra(G,i,j);
6     end
7 end
8 D(isinf(D))=max(max(~isinf(D)));

```

The final result is shown below.

