

# CPSC 340: Assignment2

Yiwei Hou(84435156)  
Qiong Zhang(85896158)

## 1 K-Nearest Neighbours

### 1.1 KNN Prediction

1. The code for prediction is given as follows:

```

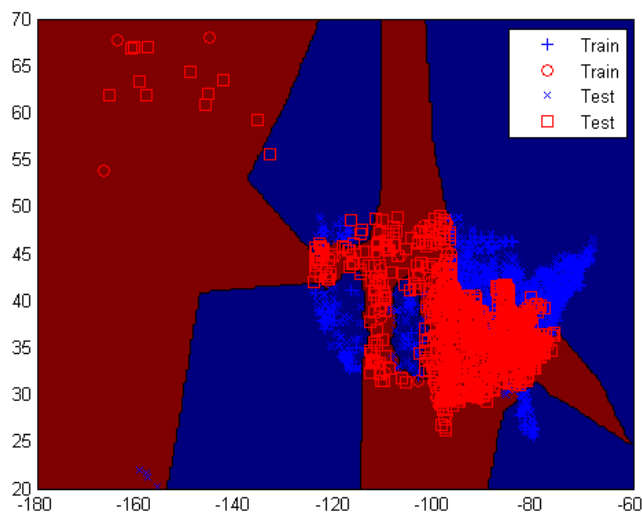
1 function [yhat] = predict(model,Xtest)
2 [n,d]=size(model.X);
3 [t,d]=size(Xtest);
4 D=(model.X).^2*ones(d,t)+ones(n,d)*(Xtest').^2-2*(model.X)*Xtest';
5 [~,pos]=sort(D,1);
6 if model.k==1
7     yhat=model.y(pos(1:model.k,:));
8 else
9     yhat=mode(model.y(pos(1:model.k,:)),1);
10 end
11 end

```

2. The training and test error obtained on *citiesSmall.mat* dataset with different  $k$  values is given by the table below:

k	1	3	10
Training error	0	0.03	0.07
Test error	0.06	0.07	0.1

3. The plot generated by classifier2Dplot when  $k = 1$  is given by:



4. When  $k = 1$ , the classification is done by finding out the label of the object that is closest to the object which we want to classify. The closest point to a point in training set is itself, thus by KNN, we classify all the objects in training set correctly.
5. We can use cross validation for training dataset to find out a  $k$ . The  $k$  should be chosen as the one that minimizes error.

## 1.2 Condensed Nearest Neighbours

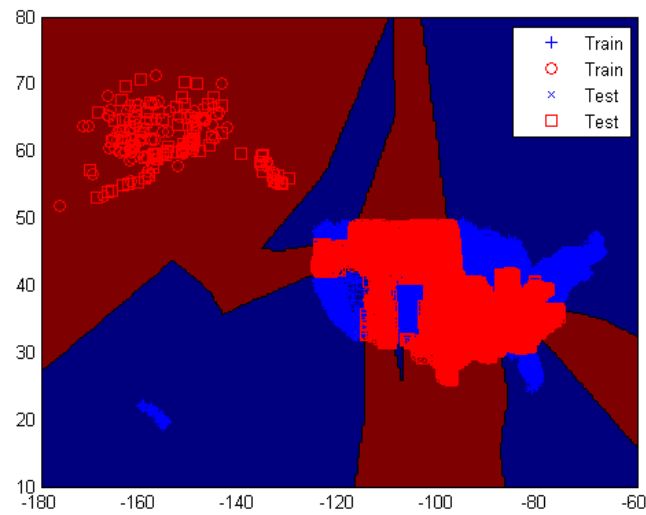
1. The code for *cnn.m*:

```

1  function [model] = cnn(X,y,k)
2  % [model] = cnn(X,y,k)
3  %
4  % condensed nearest neighbour
5  [n d]=size(X);
6  Xsubset=X(1,:);
7  ysubset=y(1);
8  for i=2:n;
9      model=knn(Xsubset,ysubset,k);
10     yhat=model.predict(model,X(i,:));
11     if yhat~=y(i);
12         Xsubset=[Xsubset;X(i,:)];
13         ysubset=[ysubset;y(i)];
14     end
15 end
16
17 model.X = Xsubset;
18 model.y = ysubset;
19 model.k = k;
20 model.c = max(y);
21 model.predict = @predict;
22 end
23
24 function [yhat] = predict(model,Xtest)
25 % Write me!
26 [n,d]=size(model.X);
27 [t,d]=size(Xtest);
28 D=(model.X).^2*ones(d,t)+ones(n,d)*(Xtest').^2-2*(model.X)*Xtest';
29 [~,pos]=sort(D,1);
30 if model.k==1
31     yhat=model.y(pos(1:model.k,:));
32 else
33     yhat=mode(model.y(pos(1:model.k,:),:),1)';
34 end
35 end

```

2. For dataset *citiesBig1.mat* with  $k = 1$ , the training and testing errors are 0.01 and 0.02 respectively. The number of objects in subset is 457.
3. The plot generated by classifier2Dplot for *citiesBig1.mat* when  $k = 1$  is given by:



4. In getting the subset of the training set, we go through the data by order. When a new example comes in, those examples who are classified correctly by a former classifier might be classified wrongly by the classifier built on the current subset. We are not using the whole data set when making predictions. Thus, the training error is greater than 0.
5. It costs  $O(sd)$  to predict one test object, since we have  $t$  test objects, the total cost for prediction will be  $O(std)$ .
6. For dataset *citiesBig2.mat* with  $k = 1$ , the training and testing errors are 0.14 and 0.21 respectively. The number of subjects in subset is 30. The high error is caused by the small size of the subset of the training set. The final classifier is built on the training set of size 30 compared to a training set of size 457 in question 2, both training and test errors increase. It might be because the order of the data set influences the size of the subset that we get from going through each data once.

## 2 Random Forests

### 2.1 Random Trees

1. The training error and test error for the *decisionTree* model based on *vowel.mat* is given as follows(Figure 1):

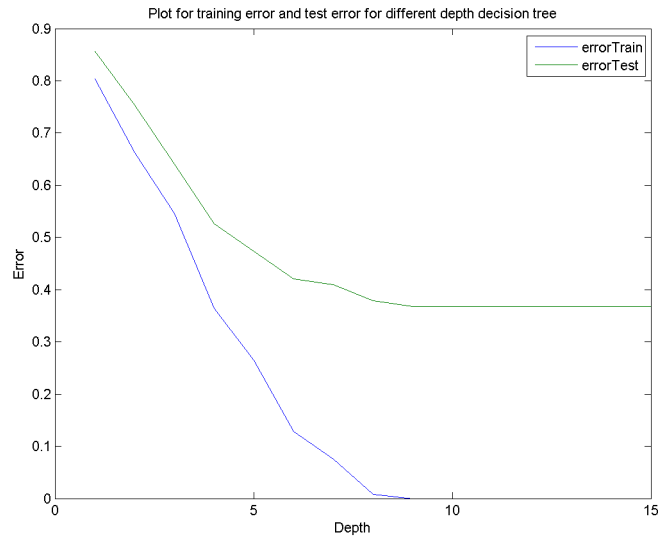


Figure 1: Training and test error for decisionTree model

2. The training set have 264 objects, a decision tree with depth 9 has  $2^9 - 1 = 512$  nodes which is larger than the size of the training set. The decision tree stop growing after the depth larger than 9 since there's no more node to split.
3. New training function is:

```

1 function [model] = randomStump(X,y)
2 % [model] = randomStump(X,y)
3 %
4 % Fits a decision stump that splits on a single variable.
5
6 % Choose random features.
7 [n,d] = size(X);
8 m=floor(sqrt(d));
9 Xrandomfeature=zeros(n,m);
10 for i=1:m;
11     j=ceil(rand*d);
12     Xrandomfeature(:,i)=X(:,j);
13 end
14
15
16 % Computer number of class lables
17 k = max(y);
18
19 % Address the trivial case where we do not split
20 count = accumarray(y,ones(size(y)),[k 1]); % Counts the number of occurrences
    of each class
21 [maxCount,maxLabel] = max(count);
22
23 % Compute total entropy (needed for information gain)
24 p = count/sum(count); % Convert counts to probabilities
25 entropyTotal = -sum(p.*log0(p));
26

```

```

27 maxGain = 0;
28 splitVariable = [];
29 splitThreshold = [];
30 splitLabel0 = maxLabel;
31 splitLabel1 = [];
32
33 % Loop over features looking for the best split
34 if any(y ~= y(1))
35     for j = 1:m
36         thresholds = sort(unique(Xrandomfeature(:,j)));
37
38         for t = thresholds'
39
40             % Count number of class labels where the feature is greater than
              threshold
41             yVals = y(X(:,j) > t);
42             count1 = accumarray(yVals, ones(size(yVals)), [k 1]);
43             count0 = count - count1;
44
45             % Compute infogain
46             p1 = count1/sum(count1);
47             p0 = count0/sum(count0);
48             H1 = -sum(p1.*log0(p1));
49             H0 = -sum(p0.*log0(p0));
50             prob1 = sum(X(:,j) > t)/n;
51             prob0 = 1-prob1;
52             infoGain = entropyTotal - prob1*H1 - prob0*H0;
53
54             % Compare to minimum error so far
55             if infoGain > maxGain
56                 % This is the lowest error, store this value
57                 maxGain = infoGain;
58                 splitVariable = j;
59                 splitThreshold = t;
60
61                 % Compute majority class
62                 [maxCount, splitLabel1] = max(count1);
63                 [maxCount, splitLabel0] = max(count0);
64             end
65         end
66     end
67 end
68 model.splitVariable = splitVariable;
69 model.splitThreshold = splitThreshold;
70 model.label1 = splitLabel1;
71 model.label0 = splitLabel0;
72 model.predict = @predict;
73 end
74
75 function [y] = predict(model,X)
76 [t,d] = size(X);
77
78 if isempty(model.splitVariable)
79     y = model.label0*ones(t,1);
80 else

```

```

81     y = zeros(t,1);
82     for n = 1:t
83         if X(n,model.splitVariable) > model.splitThreshold
84             y(n,1) = model.label1;
85         else
86             y(n,1) = model.label0;
87         end
88     end
89 end
90 end

```

When *randomStump* only considers  $\lfloor \sqrt{d} \rfloor$  randomly-chosen features. One run of the method gives the training error and test error as follows(Figure 2):

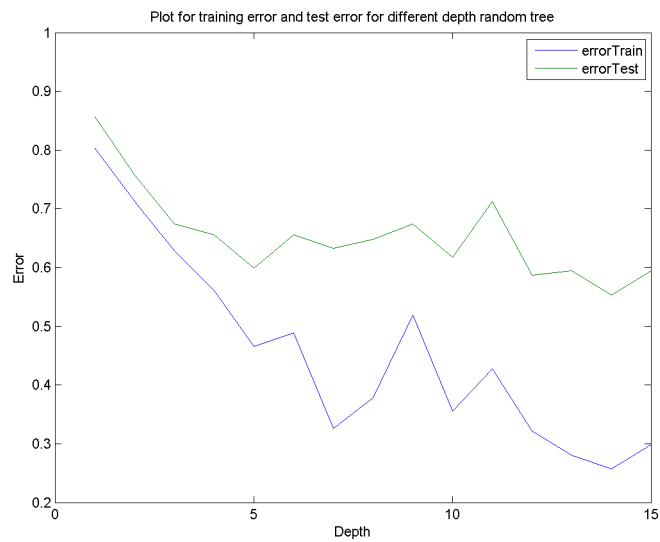


Figure 2: Training and test error for randomTree model

- For each depth, train *decisionTree* on a different bootstrap sample, then the training error and test error become(Figure 3):

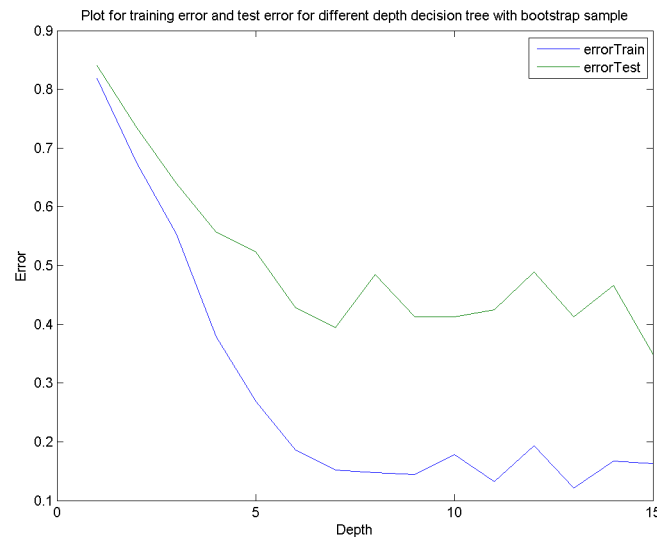


Figure 3: Training and test error of a decisionTree trained on a bootstrap sample

## 2.2 Random Decision Forests

When set depth equals 1000, one run gives the output as follows:

1. The test error of the *decisionForest* classifier with a depth of  $\infty$  and 50 trees is **0.3674**.
2. When each tree is trained on a different bootstrap sample, the test error of the *decisionForest* classifier with a depth of  $\infty$  and 50 trees becomes **0.2652**.
3. When repeatedly calls *randomTree* instead of *decisionTree* on the original dataset, the test error of the *decisionForest* classifier with a depth of  $\infty$  and 50 trees becomes **0.4091**.
4. When use *randomTree* and train each tree on a different bootstrap sample, the test error becomes **0.4053**. The code for *decisionForest.m*:

```

1 function [model] = decisionForestmodified(X,y,depth,nBootstraps)
2 [n,d]=size(X);
3
4 % Fit model to each bootstrap sample of data
5 for m = 1:nBootstraps
6     for i=1:n;
7         j=ceil(rand*n);
8         Xbootstrap(i,:)=X(j,:);
9         ybootstrap(i)=y(j);
10    end
11    model.subModel{m} = randomTree(Xbootstrap,ybootstrap',depth);
12 end
13
14 model.predict = @predict;
15
16 end
17
18 function [y] = predict(model,X)

```

```

19
20 % Predict using each model
21 for m = 1:length(model.subModel)
22     y(:,m) = model.subModel{m}.predict(model.subModel{m},X);
23 end
24
25 % Take the most common label
26 y = mode(y,2);
27 end

```

5. Bootstrapping can decrease the test error as it allows visits to different samples. Random split increases the test error.

## 3 K-Means Clustering

### 3.1 Selecting among Initializations

1. The code for *error* is given as:

```

1 function [error] = error(model,X)
2 [n,d] = size(X);
3 W = model.W;
4 y = model.y;
5 k = size(W,1);
6
7 % Compute the sum of squared distances between samples and their means
8 X2 = X.^2*ones(d,k);
9 distance2 = X2 + ones(n,d)*(W').^2 - 2*X*W';
10 error=0;
11 for i=1:n;
12     error=error+distance2(i,y(i));
13 end
14 end

```

2. The plot of clustering that minimizes the objective function by running k-means 50 times is give by(Figure 4):



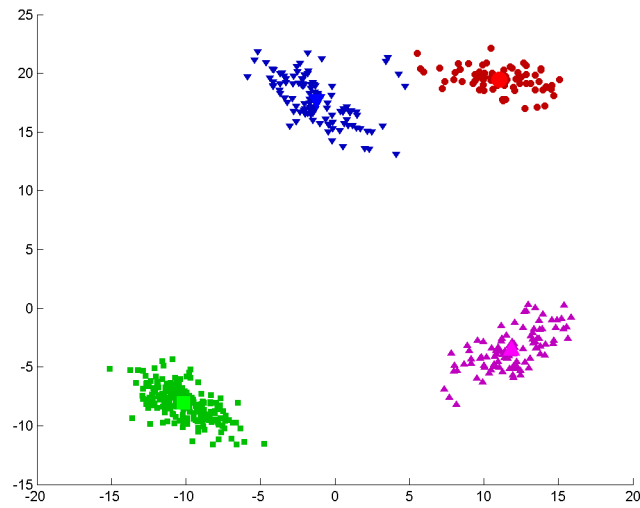


Figure 4: K-Means Clustering with the minimum error over 50 runs when  $k = 4$

### 3.2 Selecting $k$

1. When  $k$  equals the number of objects in the training set, the objective function is minimized and equals to 0. This choice of  $k$  does not make any sense for clustering, only one object forms a cluster.
2. Clustering is an unsupervised learning, and we usually don't split the dataset in to training data and test data.
3. Figure 5 shows the minimum error when  $k$  varies from 1 to 10.

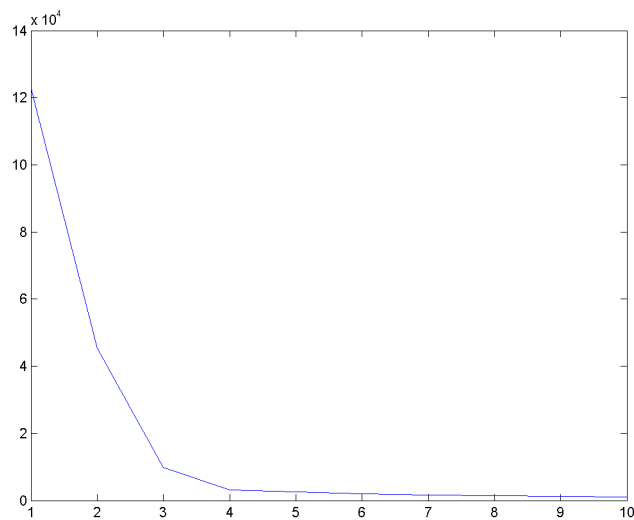


Figure 5: Error curve of Elbow Method for K-Means clustering without outlier

- The plot for elbow method is given in Figure 5. According to the plot,  $k = 3$  gives the sharpest elbow.

### 3.3 $k$ -Medians

- The clustering with the lowest error obtained by running  $k$ -means 50 times(Figure 6):

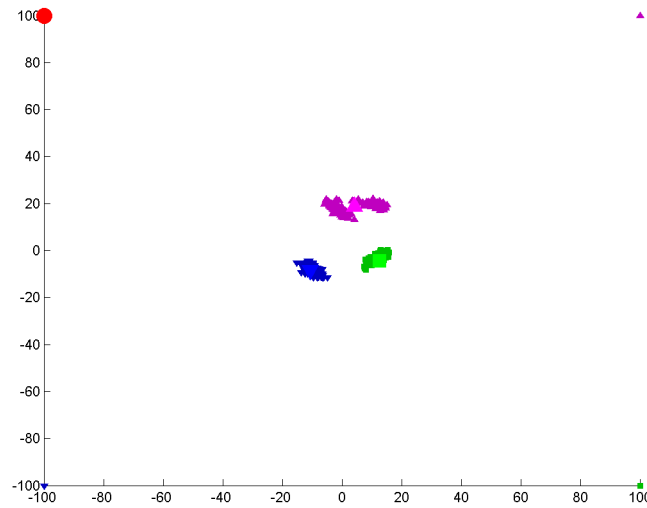


Figure 6: K-Means Clustering for dataset with an outlier

- By elbow method, the error curve is given by(Figure 7):

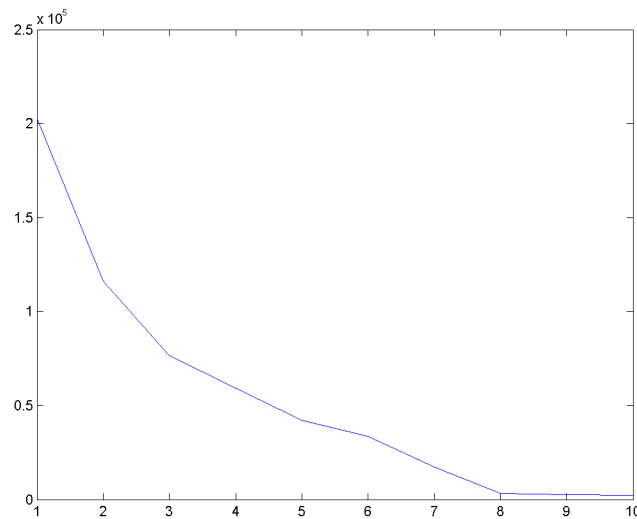


Figure 7: Error curve of Elbow Method for K-Medians clustering without outlier

The  $k$  should be chosen as 3.

3. The clustering by K-Medians is given by(Figure 8):

```

1     function [e] = error (medmodel,X)
2     [t,d] = size(X);
3     W = medmodel.W;
4     k = size(W,1);
5
6     % Compute Euclidean distance between each data point and each mean
7     distance=zeros(t,k);
8     % Compute L1 norm between each data point and each mean
9     for j=1:k
10        for i=1:t
11            distance(i,j)=norm((X(i,:) - W(j,:)),1);
12        end
13    end
14    sortD = sort(distance',1);
15    e =sum(sortD(1,:));
16 end

```

4. The error curve of Elbow method for K-Medians is given by(Figure 9):

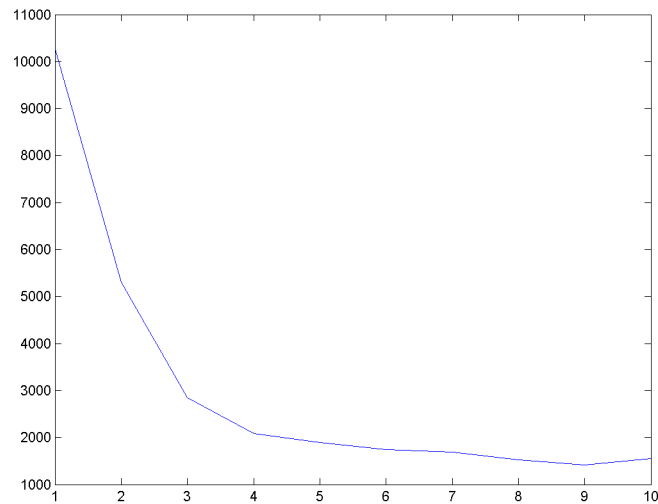


Figure 8: Error cure of Elbow method for K-Medians clustering

According to the error curve, the  $k$  should be chosen as 3.

## 4 Vector Quantization and Density-Based Clustering

### 4.1 Image Colour-Sapce Compression

1. The code for vector quantization:

```

1 function [ Iquant ] = quantizeImage( I,b )
2 % function for vector quantization

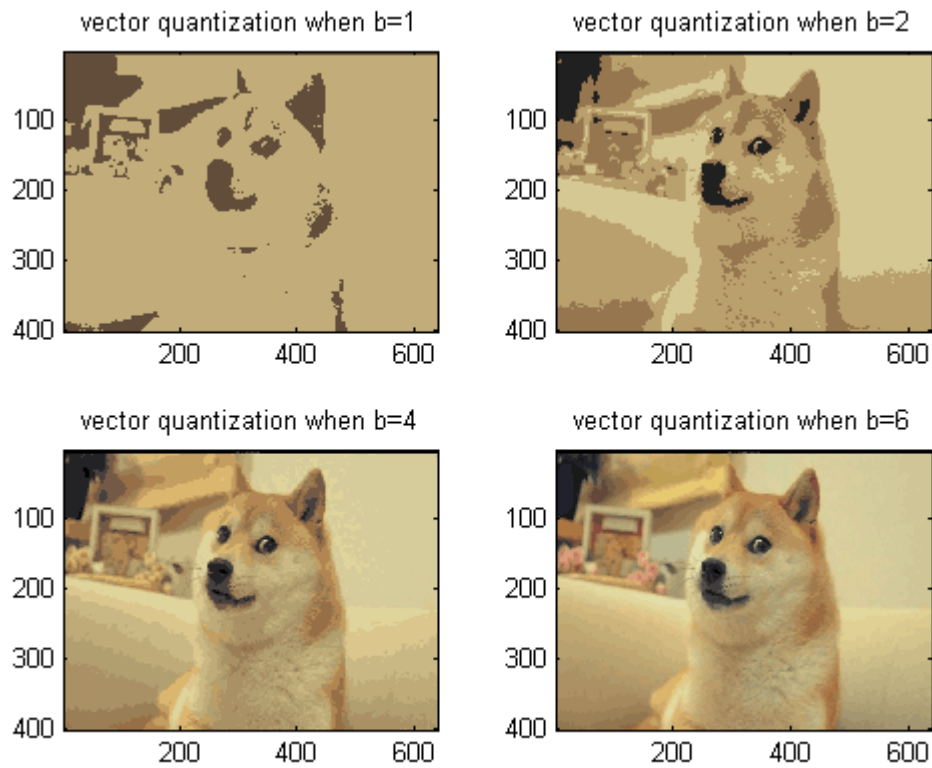
```

```

3  [n d t]=size(I);
4  I=reshape(I,n*d,t);
5  % number of clusters
6  k=2^b;
7
8  % find the prototype
9  doPlot=0;
10 model=clusterKmeans(I,k,doPlot);
11
12 % replace the original pixel's colour with the nearest prototype
13 W=model.W;
14 y=model.y;
15 Iquant=W(y,:);
16 Iquant=reshape(Iquant,[n,d,t]);
17 end

```

2. When  $b = 1, 2, 4, 6$ , the corresponding output images are



## 4.2 Effect of Parameters on DBSCAN

Table 1: Values for parameters for density-based clustering method

nclusters	1	2	3	4
radius	300	200	50	10
minPts	10	10	10	10

### 4.3 K-Means vs. DBSCAN Clustering

1.  $k$ -Means clustering heavily depends on the initialization of  $k$  centers. If there's one initial center  $w_1$  such that the closest center to each objects is not  $w_1$ , then there exists an empty cluster. The center of this empty cluster won't change in the following steps, this cluster returns NA when the algorithm converges.
2. Set the radius to 16 and minPoints to 3, and the clusters are given by:  
Cluster 1: grizzly+bear killer+whale beaver bat otter giant+panda polar+bear raccoon  
Cluster 2: antelope horse hippopotamus moose elephant ox sheep rhinoceros giraffe buffalo zebra deer pig cow  
Cluster 3: dalmatian persian+cat german+shepherd siamese+cat skunk mole tiger leopard fox hamster squirrel rabbit wolf chihuahua rat weasel bobcat lion mouse collie  
Cluster 4: blue+whale humpback+whale seal walrus dolphin  
Cluster 5: spider+monkey gorilla chimpanzee