

# Project 1

## Statistical Patter Recognition

Girguis Sedky

## 1 ML/Bayes rule estimation

In this method, the image pixel values were assumed to have a gaussian distribution. In addition, the prior probability of all the classes were assumed to be equal during training and testing. Under this assumption, the Bayesian classification rule reduces to assigning a state to the class that maximizes

$$P(\mathbf{x}|\omega_i) = \frac{1}{((2\pi)^d|\Sigma|)^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad (1)$$

where  $d = 784$  is the length of the state vectors. This is equivalent to maximizing

$$\log(P(\mathbf{x}|\omega_i)) = -\frac{1}{2} \log(|\Sigma|) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}). \quad (2)$$

The sample mean and covariances for each class are

$$\hat{\boldsymbol{\mu}}_i = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j, \quad (3)$$

$$\hat{\Sigma}_i = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \hat{\boldsymbol{\mu}}_i)(\mathbf{x}_j - \hat{\boldsymbol{\mu}}_i)^T. \quad (4)$$

A training script was used on the training set comprised of 60,000 images to determine  $\hat{\boldsymbol{\mu}}_i$  and  $\hat{\Sigma}_i$  of every class and a testing script in turn used these values and the Bayesian classification algorithm to classify the test images.

Since the covariance matrices are large and some are singular, special computational techniques were used to ensure good results. A pseudo inverse operation was used to inverse  $|\Sigma_i|$  and small multiple of the identity was added to  $|\Sigma_i|$  before obtaining its determinant. **The classification error on the testing set came to about 35%**

## 2 Nearest Neighborhood classifier

In this method, the image pixel values are arranged in a single column vector that lives in  $\mathbb{R}^{784}$ . In this space, each training image constitutes a prototype. The euclidian distance between each testing image and all the training images is found, and the testing image is assigned to the class of the training image that is closest to it. **The classification error on the testing set came to about 15%**

### 3 PCA/LDA

#### 3.1 PCA

Principal component analysis was used to reduce the dimensionality of the problem. The low-order state vectors were then fed into the Bayesian and Nearest Neighborhood classifier. Figure 1 demonstrates the percentage error in classification based on the number of dimensions kept. **The minimum classification error came to about 20.5% for 50 kept components in the case of the Bayesian classifier and 14.7% for 150 components kept in the case of the Nearest Neighborhood classifier.**

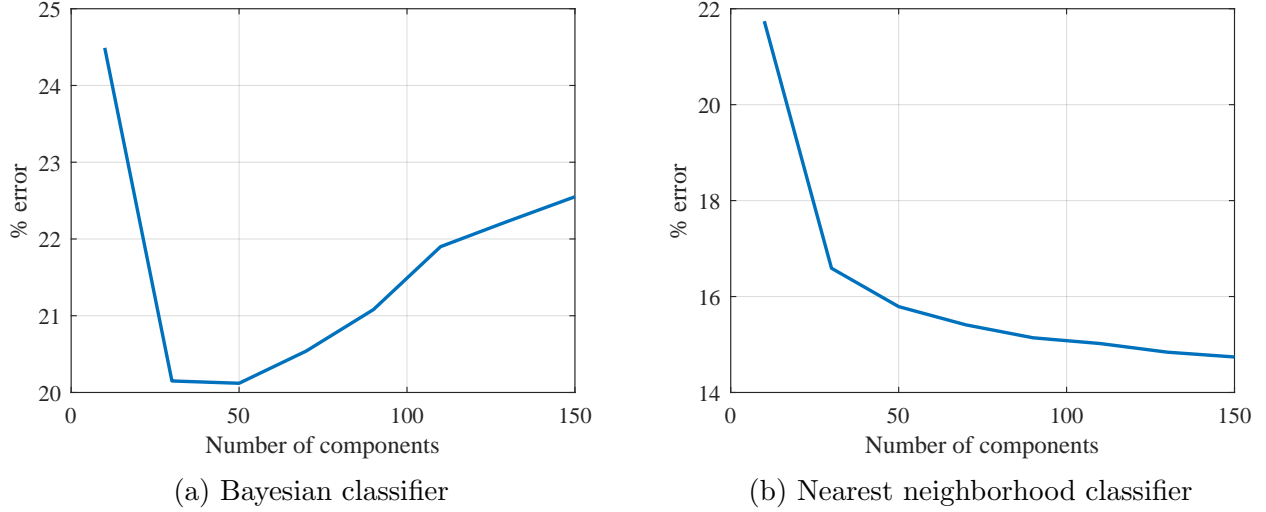


Figure 1: Principal component analysis

#### 3.2 LDA

Linear discriminant analysis was used to reduce the image state vectors from  $\mathbb{R}^{784}$  to  $\mathbb{R}^{M-1}$  where  $M$  is the number of classes. The built-in MATLAB function was used for this purpose. Again, the low-order state vectors were then fed into the Bayesian and Nearest Neighborhood classifier. **The minimum classification error came to about 25.4% in the case of the Bayesian classifier and 28% in the case of the Nearest Neighborhood classifier.**

## 4 Summary

Table 1 shows the classification error of the different methods utilized in this study. The most accurate method of classification was to use principal component analysis to reduce the dimensionality of the state vectors to 50 and then to apply NN to the transformed data.

Table 1: Summary of classification results

Method	Percentage Error
Bayes	35%
NN	15%
PCA, Bayes (50)	20.5%
PCA, NN (150)	14.7%
LDA, Bayes	25.4%
LDA, NN	28%

# Appendix

## ML/Bayes rule estimation

### Training script

```
1 %% training script
2 % find the estimated averages and covariances under gaussian
  assumptions of
3 % the training data and maximum likelihood
4 clear; clc; close all;
5 % Change the filenames if you've saved the files under different names
6 % On some platforms, the files might be saved as
7 % train-images.idx3-ubyte / train-labels.idx1-ubyte
8 images = loadMNISTImages('G:\My Drive\Classes&Books\
  StatisticalPatternRecognition\Project\Data/train-images-idx3-ubyte'
  );
9 labels = loadMNISTLabels('G:\My Drive\Classes&Books\
  StatisticalPatternRecognition\Project\Data/train-labels-idx1-ubyte'
  );
10 mu_hat_mat = [];
11 Epsilon_hat_cell = {};
12 % number of labels
13 labels_num = 10;
14 for jj=1:labels_num
15 images_label = images(:, labels==jj-1);
16 % n = number of images
17 n = size(images_label,2);
18 % L = number of pixels in an image
19 L = size(images_label,1);
20 % Find the maximum likelihoods estimates for the mean and covariance
21 mu_hat = mean(images_label,2);
22 Epsilon_hat = cov(images_label');
23 mu_hat_mat(:,jj) = mu_hat;
24 Epsilon_hat_cell{jj} = Epsilon_hat;
25
26 end
27 save('Estimates.mat','mu_hat_mat','Epsilon_hat_cell');
```

### Testing script

```
1 %% Test script
2 % Test the ML algorithm based on gaussian assumptions of
3 % the training data and maximum likelihood
4 clear; clc; close all;
5 % Change the filenames if you've saved the files under different names
6 % On some platforms, the files might be saved as
7 % train-images.idx3-ubyte / train-labels.idx1-ubyte
```

```

8  images = loadMNISTImages('G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\t10k-images-idx3-ubyte')
    ;
9  labels = loadMNISTLabels('G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\t10k-labels-idx1-ubyte')
    ;
10 load('Estimates.mat');
11 % number of labels
12 labels_num = 10;
13 for ii=1:length(images)
14 % pick an image to classify
15 image = images(:,ii);
16 % n = number of images
17 n = size(image,2);
18 % L = number of pixels in an image
19 L = size(image,1);
20 for jj=1:labels_num
21 % Find the posterior probability estimate of each one
22 cov = Epsilon_hat_cell{jj};
23 cov_det = cov+(10^-10)*eye(size(cov)); % Make covariance matrix full
    rank so that the logdet function can find a determinant
24 mu = mu_hat_mat(:,jj); % average vector
25 % Find the likelihood of this class
26 logP(jj) = -0.5*logdet(cov_det)-0.5*(image-mu)'*pinv(cov)*(image-mu);
27
28 end
29 % Choose the maximum likelihood and assign the image to this class
30 [~,I] = max(logP);
31 labels_est(ii) = I-1;
32 clear P;
33 end
34 labels_est = labels_est';
35 % Test accuracy of classifying algorithm
36 a = find(labels_est==labels);
37 Performance = length(a)*100/length(labels);
38 save('Performance.mat','Performance');

NN

1 %% Nearest Neighborhood
2 clear; clc; close all;
3 % Girguis Sedky
4 %
    %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

5 %% Load data
6 % Prototypes

```

```

7 images_train = loadMNISTImages( 'G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\train-images-idx3-ubyte'
    );
8 labels_train = loadMNISTLabels( 'G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\train-labels-idx1-ubyte'
    );
9 % Test data
10 images_test = loadMNISTImages( 'G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\t10k-images-idx3-ubyte'
    );
11 labels_test = loadMNISTLabels( 'G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\t10k-labels-idx1-ubyte'
    );
12
13 %% Loop through test data
14 for j = 1:length(images_test)
15     % Pick image
16     image = images_test(:,j);
17     % Find eucliden distance to all the training set
18     Dist = vecnorm(images_train-image,2,1);
19     [~,I] = min(Dist);
20     labels_est(j) = labels_train(I);
21 end
22 %% Test accuracy of classifying algorithm
23 a = find(labels_est==labels_test');
24 Performance = length(a)*100/length(labels_test);
25 save('Performance.mat','Performance')

```

## PCA

### 4.1 ML/Bayes rule estimation with PCA

```

1 %% ML with PCA
2 % find the estimated averages and covariances under gaussian
    assumptions of
3 % the training data and maximum likelihood , utilize PCA
4 clear; clc; close all;
5 % Girguis Sedky
6 %
    %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
7 %% Load data
8 images_train = loadMNISTImages( 'G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\train-images-idx3-ubyte'
    );
9 labels_train = loadMNISTLabels( 'G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\train-labels-idx1-ubyte'

```

```

    );
10 images_test = loadMNISTImages('G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\t10k-images-idx3-ubyte')
    ;
11 labels_test = loadMNISTLabels('G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\t10k-labels-idx1-ubyte')
    ;
12 mu_hat_mat = [];
13 Epsilon_hat_cell = {};
14 % number of labels
15 labels_num = 10;
16 % PCA
17 % Find the principal components
18 coeff = pca(images_train');
19 % Number of dimensions kept
20 m = [10:20:150];
21 %% Loop over different values of m, number of dimensions kept
22 for ii=1:length(m)
23 %% Training Stage
24 % reduce the dimensionality of the system to 20
25 images_Transf = coeff(:,1:m(ii))'*images_train;
26 %% Loop over all the pictures in one class inside the training set
27 for jj=1:labels_num
28 images_label = images_Transf(:,labels_train==jj-1);
29 % n = number of images
30 n = size(images_label,2);
31 % L = number of pixels in an image
32 L = size(images_label,1);
33 % Find the maximum likelihoods estimates for the mean and covariance
34 mu_hat = mean(images_label,2);
35 Epsilon_hat = cov(images_label');
36 mu_hat_mat(:,jj) = mu_hat;
37 Epsilon_hat_cell{jj} = Epsilon_hat;
38 end
39 %% Test Stage
40 for kk=1:length(images_test)
41 % pick an image to classify
42 image = images_test(:,kk);
43 % reduce the dimensionality of the system to 20
44 image_Transf = coeff(:,1:m(ii))'*image;
45 for jj=1:labels_num
46 % Find the posterior probability estimate of each one
47 COV = Epsilon_hat_cell{jj};
48 mu = mu_hat_mat(:,jj);
49 P(jj) = (1/(sqrt(det(COV))*(2*pi)^m(ii))))*exp(-0.5*(image_Transf-mu)'*
    inv(COV)*(image_Transf-mu));
50 end

```

```

51 [~,I] = max(P);
52 labels_est(kk) = I-1;
53 clear P;
54 end
55 clear mu_hat_mat Epsilon_hat_cell;
56 %% Test accuracy of classifying algorithm
57 a = find(labels_est==labels_test ');
58 Performance(ii) = length(a)*100/length(labels_test);
59
60 end
61 %% save and plot
62 save('Performance.mat','m','Performance');
63 plot(m,100-Performance,'LineWidth',1.5);
64 xlabel('Number of components');
65 ylabel('% error');
66 grid on;
67 PrintPlot(gcf,'plot.png','-dpng');
68 PrintPlot(gcf,'plot.pdf','-dpdf');

```

## 4.2 NN with PCA

```

1 %% Nearest Neighborhood, with PCA, investigate the effect of different
   components
2 clear; clc; close all;
3 % Girguis Sedky
4 %
   %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
5 %% Load data
6 % Prototypes
7 images_train = loadMNISTImages('G:\My Drive\Classes&Books\
   StatisticalPatternRecognition\Project\Data\train-images-idx3-ubyte'
   );
8 labels_train = loadMNISTLabels('G:\My Drive\Classes&Books\
   StatisticalPatternRecognition\Project\Data\train-labels-idx1-ubyte'
   );
9 % Test data
10 images_test = loadMNISTImages('G:\My Drive\Classes&Books\
   StatisticalPatternRecognition\Project\Data\t10k-images-idx3-ubyte'
   );
11 labels_test = loadMNISTLabels('G:\My Drive\Classes&Books\
   StatisticalPatternRecognition\Project\Data\t10k-labels-idx1-ubyte'
   );
12 % PCA
13 % Find the principal components
14 coeff = pca(images_train);
15 % Number of dimensions kept

```



```

16 m = [10:20:150];
17 %% Loop over different values of m, number of dimensions kept
18 for ii=1:length(m)
19 % reduce the dimensionality of the system to 20
20 images_train_Transf = coeff(:,1:m(ii))'*images_train;
21 images_test_Transf = coeff(:,1:m(ii))'*images_test;
22 %% Loop through test data
23 for j = 1:length(images_test_Transf)
24 % Pick image
25 image = images_test_Transf(:,j);
26 % Find eucliden distance to all the training set
27 Dist = vecnorm(images_train_Transf-image,2,1);
28 [~,I] = min(Dist);
29 labels_est(j) = labels_train(I);
30 end
31 %% Test accuracy of classifying algorithm
32 a = find(labels_est==labels_test ');
33 Performance(ii) = length(a)*100/length(labels_test);
34 end
35 %% save and plot
36 save('Performance.mat','m','Performance');
37 plot(m,100-Performance,'LineWidth',1.5);
38 xlabel('Number of components');
39 ylabel('% error');
40 grid on;
41 PrintPlot(gcf,'plot.png','-dpng');
42 PrintPlot(gcf,'plot.pdf','-dpdf');

```

## LDA

```

1 %% Reduce the data diemsnionality using LDA
2 clear; clc; close all;
3 % Girguis Sedky
4 %
5 %% Load data
6 images_train = loadMNISTImages('G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\train-images-idx3-ubyte'
    );
7 labels_train = loadMNISTLabels('G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\train-labels-idx1-ubyte'
    );
8 images_test = loadMNISTImages('G:\My Drive\Classes&Books\
    StatisticalPatternRecognition\Project\Data\t10k-images-idx3-ubyte'
    );
9 labels_test = loadMNISTLabels('G:\My Drive\Classes&Books\

```

```

        StatisticalPatternRecognition\Project\Data\t10k-labels-idx1-ubyte')
    ;
10 % number of labels
11 A = LDA(images_train', labels_train);
12 A = A(1:end-1, 1:end-1);
13 images_train = A*images_train;
14 images_test = A*images_test;
15 save('LDA_Data.mat', 'images_train', 'images_test', 'labels_train', '
    labels_test');

```