

Analyse et Traitement de l'Information

TP5: Linear Discriminant Analysis (LDA) and Information Theory.

1 Fisher's Discriminant for Two Classes (50%)

Assume we are given a *training set* $\{\mathbf{x}_i, y_i\}_{i=1}^n$ where each \mathbf{x}_i is a high dimensional observation in \mathbb{R}^d and each $y_i = 0, \dots, (L-1)$ corresponds to a label specifying which of the L possible classes \mathbf{x}_i belongs to. For instance, each \mathbf{x}_i may be an image where $y_i = 0, \dots, (L-1)$ corresponds to the image class. The goal of *classification* is to learn a mapping $f : \mathbf{x} \rightarrow y$ based on this labeled training set. Predictions can then be made for any future sample \mathbf{x} with an unknown label by setting its predicted label equal to $f(\mathbf{x})$.

A simple approach to classification is Fisher's linear discriminant analysis (LDA). The basic idea of LDA is to project $\{\mathbf{x}_i\}_{i=1}^n$ into a low dimensional space where different classes are "best" separated. The classification is then performed on this low dimensional space. Here, without loss of generality, we only discuss LDA for two classes, i.e., the number of classes is $L = 2$ where y_i can be either 0 or 1. In this case, the classifier is defined as

$$f(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{a}^T \mathbf{x} + b \geq 0; \\ 1 & \text{if } \mathbf{a}^T \mathbf{x} + b < 0, \end{cases} \quad (1)$$

where $\mathbf{a} \in \mathbb{R}^d$ specifies a linear projection from \mathbb{R}^d to \mathbb{R} and $b \in \mathbb{R}$ is a threshold to determine which of the two classes \mathbf{x} belongs to. The learning task is to find the optimal \mathbf{a} and b based on a given training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$.

Rather than provide the complete derivation, we directly provide the solution of the optimal \mathbf{a} and b and refer the reader to the textbooks [1, 2] for the detailed analysis. The optimal \mathbf{a} with respect to $\{\mathbf{x}_i, y_i\}_{i=1}^n$ is given by

$$\mathbf{a} = S_W^{-1}(\mathbf{g}_0 - \mathbf{g}_1). \quad (2)$$

In eq.(2),

$$\mathbf{g}_0 = \frac{1}{N_0} \sum_{i:y_i=0} \mathbf{x}_i, \quad \mathbf{g}_1 = \frac{1}{N_1} \sum_{i:y_i=1} \mathbf{x}_i,$$

where N_0 and N_1 are the number of samples in class 0 and 1, respectively, and \mathbf{g}_0 and \mathbf{g}_1 are the corresponding mean vectors. S_W is a $d \times d$ *within-class scatter matrix* defined as

$$S_W = \sum_{i:y_i=0} (\mathbf{x}_i - \mathbf{g}_0)(\mathbf{x}_i - \mathbf{g}_0)^T + \sum_{i:y_i=1} (\mathbf{x}_i - \mathbf{g}_1)(\mathbf{x}_i - \mathbf{g}_1)^T.$$

Since the given classes of data are balanced, the optimal b with respect to $\{\mathbf{x}_i, y_i\}_{i=1}^n$ is given by

$$b = -\frac{1}{2} \mathbf{a}^T (\mathbf{g}_0 + \mathbf{g}_1). \quad (3)$$

Consider the binary learning problem to classify images of “3” and “7”. From the MNIST dataset, sample 2000 instances of “3”s and “7”s for training and use all instances for testing. Apply PCA to reduce the dimensionality of the input data to 50 before classification. Then perform the following tasks.

1. Apply LDA to classify the complete set of testing instances. Report the accuracy using the confusion matrix¹.
2. Use the k -NN method that you implemented in TP3 to classify the testing instances with $k = 1$. Report the accuracy using the confusion matrix.
3. Describe how the two classification methods compare.

2 k-means Clustering (50%)

Classification is a *supervised* problem, where labeled (training) data is used to label the unlabeled (testing) data.

Clustering is an *unsupervised* problem, where the objective is to group a given (unlabeled) dataset $\{\mathbf{x}_i\}_{i=1}^n$ into k clusters, so that similar samples appear in the same cluster, and dissimilar samples are in different clusters.

The most basic clustering algorithm, k -means, aims to partition the input data into k clusters in the most compact way, in the sense that the following objective function is minimized

$$J(\mathbf{y}) = \sum_{j=1}^k \sum_{i:y_i=j} \|\mathbf{x}_i - \mathbf{g}_j\|^2, \quad (4)$$

where

- $\mathbf{y} = (y_1, \dots, y_n)$ represents a clustering scheme. Each y_i is an integer taking values from $\{1, 2, \dots, k\}$, indicating which cluster the corresponding sample \mathbf{x}_i belongs to;
- \mathbf{g}_j is the center of cluster j so that

$$\mathbf{g}_j = \frac{1}{N_j} \sum_{i:y_i=j} \mathbf{x}_i \quad (N_j \text{ is the number of samples in cluster } j)$$

Therefore, $J(\mathbf{y})$ means the sum of the square distances from each sample to its associated cluster center. Minimizing $J(\mathbf{y})$ means to take each cluster *compact*, in the sense that a sample should be near to its cluster center. To minimize this $J(\mathbf{y})$, the k -means algorithm is given as follows

- Randomly sample 2000 instances “3”s and “7”s from the MNIST dataset. Perform k -means with $k = 2, 3, \dots, 10$. Compute J defined in eq.(4) for each k . Plot J varying with k . Describe the effect of k over J .

¹https://en.wikipedia.org/wiki/Confusion_matrix

```

1: procedure K-MEANS( $\{\mathbf{x}_i\}_{i=1}^n, k$ )  $\triangleright$  As inputs: a set of samples and the number of clusters
2:    $\triangleright$  As outputs:  $\{y_i\}_{i=1}^n$ .  $y_i \in \{1, 2, \dots, k\}$  specifies  $\mathbf{x}_i$  belongs to which of the  $k$  clusters.
3: Randomly initialize each  $y_i^0$  to one of  $\{1, 2, \dots, k\}$ ;  $t=0$ ;
4: repeat
5:   for  $j \leftarrow 1 \rightarrow k$  do  $\triangleright$  Compute the centroid of each cluster
6:      $\mathbf{g}_j \leftarrow \sum_{i: y_i^t = j} \mathbf{x}_i / N_j$ ;
7:   for  $i \leftarrow 1 \rightarrow n$  do  $\triangleright$  Assign each point to the closest centroid
8:      $y_i^{t+1} \leftarrow \operatorname{argmin}_j \|\mathbf{x}_i - \mathbf{g}_j\|$ ;
9:    $t \leftarrow t + 1$ ;
10: until  $\forall i, y_i^t = y_i^{t-1}$ 
11: Output  $\{y_i\}_{i=1}^n$ 

```

- For the case $k = 2$ in the above experiment, give the confusion matrix.
- Sample 2000 instances of “3”s and “5”s (rather than “3”s and “7”s). Perform k -means with $k = 2$. Build the confusion matrix. Based on your results, explain the difference between these two clustering problems: (1) “3” vs. “7”; (2) “3” vs. “5”.

Assessment

Please archive your report and codes in “Prénom Nom.zip” (replace “Prénom” and “Nom” with your real name), and upload to “Upload TPs/TP4 LDA and k -means” on <https://moodle.unige.ch> before **Monday, November 11 2019, 23:59 PM**. Note, the assessment is mainly based on your report, which should include your answers to all questions and the experimental results.

3 Supplements

1. Explain the k -means algorithm.
2. Explain LDA.

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NY, USA, 2006.
- [2] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT Press, Cambridge, MA, 2012.