

Fisher's linear discriminant

Jianxin Wu

LAMDA Group

National Key Lab for Novel Software Technology

Nanjing University, China

wujx2001@gmail.com

February 20, 2019

Contents

| | | |
|----------|--|-----------|
| 1 | FLD for binary classification | 4 |
| 1.1 | Idea: what is far apart? | 4 |
| 1.2 | Translation to mathematics | 5 |
| 1.3 | Scatter vs. covariance matrix | 7 |
| 1.4 | The two scatter matrices and the FLD objective | 8 |
| 1.5 | The optimization | 9 |
| 1.6 | Wait, we have a shortcut! | 9 |
| 1.7 | The FLD method for binary problems | 10 |
| 1.8 | A caveat: what if S_W is not invertible? | 11 |
| 2 | FLD for more classes | 11 |
| 2.1 | A slightly modified notation and S_W | 12 |
| 2.2 | Candidates for S_B | 12 |
| 2.3 | A tale of three scatter matrices | 13 |
| 2.4 | The solution | 14 |
| 2.5 | Finding more projection directions | 14 |
| | Exercises | 15 |

Fisher's Linear Discriminant (FLD), which is also a linear dimensionality reduction method, extracts lower dimensional features utilizing linear relationships among the dimensions of the original input. A natural question is: what makes FLD different from PCA, and why do we still need FLD when we have PCA handy?

A short answer is: FLD is supervised, but PCA is unsupervised. We illustrate their differences using a somehow extreme example in Figure 1.

是的，对于原理相似的两个算法，
我们需要清楚每个算法适用的case
以及各自的优势和不足；

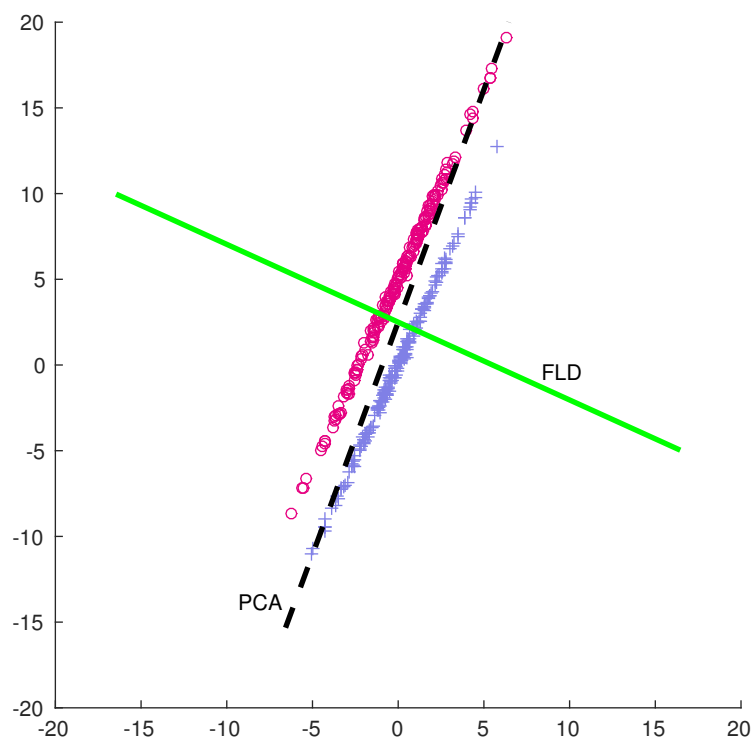


Figure 1: FLD vs. PCA.

Suppose we are to deal with a binary classification problem. We are given a set of input vectors $\mathbf{x} \in \mathbb{R}^D$ (which is two dimensional in Figure 1), and need to classify them into two classes. The class label y can take value in the set $\{1, 2\}$. In Figure 1, we are given 200 examples with $y = 1$ and 200 examples with $y = 2$. These 400 examples form our training set for this binary classification problem. The positive examples (with $y = 1$) are denoted by the red \circ sign, and negative examples ($y = 2$) are denoted by the blue $+$ sign in Figure 1.

In this example, the two classes have special properties: the inherent dimensionality of both seem to be 1. In fact, these data points are generated using the following Matlab command:

```
n1 = 200;
n2 = 200;
rot = [1 2; 2 4.5];
data1 = randn(n1,2) * rot;
data2 = randn(n2,2) * rot + repmat([0 5],n2,1);
```

In other words, we generate two very thin Gaussians, and shift one of them to separate them apart.

Visual inspection tells us that it is easy to classify these examples into two classes if we project any \mathbf{x} into the projection direction labeled “FLD” (the green solid line). After the projection, the positive class examples will clump into a small range of values, so do the negative examples. However, the range of projected values in different classes almost never overlap with each other. That is, a proper linear dimensionality reduction makes our binary classification problem trivial to solve.

But, if we compute the PCA solution using all 400 training examples, the first eigenvector generates the black dashed line (with a label “PCA” in Figure 1). The PCA projection direction is surprisingly bad for the classification. The ranges of projection values are almost identical for the positive and the negative class, that is, this direction is almost useless for classification.

What makes PCA so unsuitable in this example? PCA is an unsupervised method, which does not consider any label information. It tries to eliminate the projection direction that corresponds to small variance (which is the FLD projection direction). However, in this example, the labels tell us that it is this specific direction that separates the two classes!

Hence, when we are given labels of examples in a classification task, it is essential to take these labels into consideration when we perform linear dimensionality reduction. FLD is such a supervised dimensionality reduction method, and we will start the introduction of FLD from the simple binary classification setup again.

The example in Figure 1 is a little bit contrived and extremal. In real applications, both PCA and FLD are useful and they can be combined to produce better linear dimensionality reduction results. PCA+FLD has been a very successful method in face recognition.

1 FLD for binary classification

The example in Figure 1 also gives us some hints on how to incorporate the labels. Given any projection direction \mathbf{w} ($\|\mathbf{w}\| = 1$), we can compute the projected values as $\mathbf{w}^T \mathbf{x}$, and require that projected values from the positive and the negative classes are *far apart* from each other.

1.1 Idea: what is far apart?

But, how do we translate the phrase *two sets of values are far apart from each other* into precise mathematical expressions? That is, we need a way to measure the level of separation between two sets of values. There are many different ways to compare the distance between two sets of values, and we will introduce only the FLD way in this chapter.

Figure 2 shows the histogram (with 10 bins) of the projected values for the data in Figure 1, when the PCA or FLD projection direction is used. A plot that shows all the projected values might be too cluttered to reflect any valuable information. The histogram succinctly captures the distribution properties of the projected values, and hence is a useful tool in various studies.

Let us first look at the FLD figure (Figure 2b). The histograms of both classes are sharp, meaning that the projected values of the positive examples are in a small range, so do those of the negative examples. In this case, we can use the mean of each class to represent that class. And, the level of separation between these two classes can be measured by the distance between the two mean values. In this example, the two mean values are 0.0016 and 2.0654, respectively. The distance between them is 2.0638.

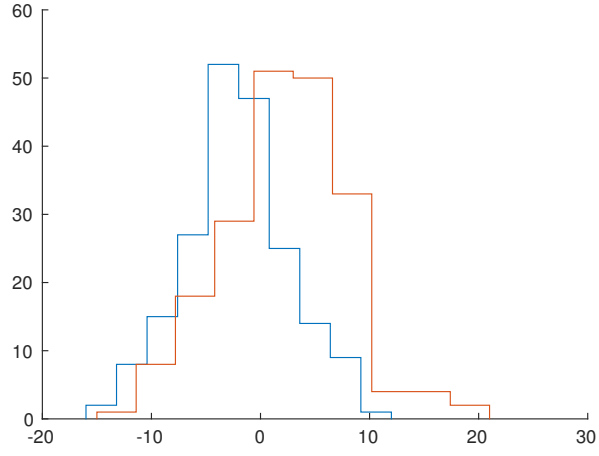
However, if we apply this measure of separation, we will encounter difficulty when we move to the PCA projection direction, whose histogram is shown in Figure 2a. The two mean values are -2.2916 and 2.2916, respectively. The distance, 4.5832, is even larger than that in the FLD case (2.0638). Note that the scale of the x -axis is different in Figure 2b and Figure 2a.

It does not mean the PCA projection direction is better, though. As shown in Figure 2a, the two histograms overlap seriously, meaning that the level of separation is low. This paradox (i.e., PCA has larger distance between projected mean values, but worse separation) appears because the *shape* of the histograms are not taken into account. The FLD histograms are sharp and have small variances (or standard deviations). The distance between mean values is very large compared to *both* standard deviations. Hence, the two histograms do not overlap with each other.

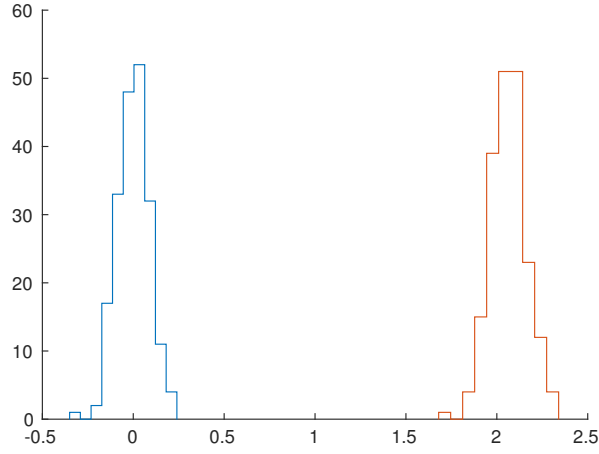
In the PCA case, the distance between mean values is smaller than both standard deviations. Hence, the two histograms significantly overlap with each other, which means low separation ability.

Putting these facts together, it is the *ratio* between

- distance between two mean values, and
- the two standard deviations



(a) PCA



(b) FLD

Figure 2: Histogram of projected values of the dataset in Figure 1 along the PCA (Figure 2a) or FLD (Figure 2b) direction. Please note that the scales in the x -axis are *different* in these figures.

that decides separability, not the distance alone. We want to maximize this ratio.

1.2 Translation to mathematics

Let us define our notations first. One training example for FLD will be a pair (\mathbf{x}_i, y_i) instead of a single vector \mathbf{x}_i , in which $\mathbf{x}_i \in \mathbb{R}^D$ is the original input

vector that describes the i -th example in the training set, and $y_i \in \{1, 2\}$ is a label that specifies to which class \mathbf{x}_i belongs. Hence, a training set X with N such pairs can be written as $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, which is an abbreviation for the set $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$.

Because there are two categories, we can collect all positive examples to form a subset of X , as

$$X_1 = \{\mathbf{x}_i | 1 \leq i \leq N, y_i = 1\},$$

and $N_1 = |X_1|$ is the size of this subset, which is also the number of positive examples in X . Similarly, we define

$$X_2 = \{\mathbf{x}_i | 1 \leq i \leq N, y_i = 2\},$$

and $N_2 = |X_2|$. Note that we do not include the labels y_i in X_1 or X_2 , because labels of examples in them can be easily deduced.

The mean of vectors in these subsets are also easy to define, as

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{\mathbf{x} \in X_1} \mathbf{x}, \quad (1)$$

$$\mathbf{m}_2 = \frac{1}{N_2} \sum_{\mathbf{x} \in X_2} \mathbf{x}, \quad (2)$$

respectively. Later, we will find these two notations very useful, and it is convenient to define them early. Similarly, the covariance matrices are

$$C_1 = \frac{1}{N_1} \sum_{\mathbf{x} \in X_1} (\mathbf{x} - \mathbf{m}_1)(\mathbf{x} - \mathbf{m}_1)^T, \quad (3)$$

$$C_2 = \frac{1}{N_2} \sum_{\mathbf{x} \in X_2} (\mathbf{x} - \mathbf{m}_2)(\mathbf{x} - \mathbf{m}_2)^T, \quad (4)$$

respectively.

Given any projection direction \mathbf{w} , as in the PCA case, we can safely assume $\|\mathbf{w}\| = 1$. Hence, the projected value for \mathbf{x}_i is $\mathbf{x}_i^T \mathbf{w}$. The mean of projected values for the two classes are

$$m_1 = \mathbf{m}_1^T \mathbf{w}, \quad (5)$$

$$m_2 = \mathbf{m}_2^T \mathbf{w}. \quad (6)$$

respectively.

The variance for projected positive examples is

$$\frac{1}{N_1} \sum_{\mathbf{x} \in X_1} (\mathbf{x}^T \mathbf{w} - m_1)^2 \quad (7)$$

$$= \mathbf{w}^T \left(\frac{1}{N_1} \sum_{\mathbf{x} \in X_1} (\mathbf{x} - \mathbf{m}_1)(\mathbf{x} - \mathbf{m}_1)^T \right) \mathbf{w} \quad (8)$$

$$= \mathbf{w}^T C_1 \mathbf{w}, \quad (9)$$

and the standard deviation is

$$\sigma_1 = \sqrt{\mathbf{w}^T C_1 \mathbf{w}}. \quad (10)$$

Similarly, the variance and standard deviation for the negative examples are

$$\mathbf{w}^T C_2 \mathbf{w}, \quad \text{and} \quad (11)$$

$$\sigma_2 = \sqrt{\mathbf{w}^T C_2 \mathbf{w}}, \quad (12)$$

respectively.

The ratio we want to maximize can be translated into different forms, e.g.,

$$\frac{|m_1 - m_2|}{\sigma_1 + \sigma_2}, \quad (13)$$

or

$$\frac{|m_1 - m_2|}{\sqrt{\sigma_1^2 + \sigma_2^2}}. \quad (14)$$

Both equations are sometimes referred to as Fisher's Linear Discriminant by some authors. However, the classic FLD has a slightly different form, which we will introduce soon. We want to point out that the solution of these two equations should also give reasonable projection directions for binary problems, because they share the same intuition with FLD.

1.3 Scatter vs. covariance matrix

Maximizing Equation 14 does not seem convenient. The absolute value and the square root both complicates the optimization. However, we can maximize the alternative objective function

$$\frac{(m_1 - m_2)^2}{\sigma_1^2 + \sigma_2^2}, \quad (15)$$

because it is obviously equivalent to Equation 14.

The projection \mathbf{w} is only implicitly embedded in Equation 15. Plugging in the expressions for m_1 , m_2 , σ_1^2 and σ_2^2 , we get the following objective which is explicit in \mathbf{w} :

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\mathbf{w}^T (C_1 + C_2) \mathbf{w}}. \quad (16)$$

The term $C_1 + C_2$ uses the covariance matrix statistics of X_1 and X_2 to measure how *scattered* they are. Other than covariance matrices, the *scatter matrix* is another interesting statistics to measure how scattered a set of points are. The scatter matrix for a set of points $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$ is defined as

$$S = \sum_{i=1}^n (\mathbf{z}_i - \bar{\mathbf{z}}) (\mathbf{z}_i - \bar{\mathbf{z}})^T, \quad (17)$$

where $\bar{\mathbf{z}}$ is the mean of these points, defined as

$$\bar{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i.$$

It is obvious that these two statistics are related to each other. Let us denote the scatter matrix for X_1 and X_2 as S_1 and S_2 , respectively. Then,

$$S_1 = N_1 C_1, \quad (18)$$

$$S_2 = N_2 C_2. \quad (19)$$

They are equivalent to each other, modulo the slight difference in multiplying the number of points or not.

Traditionally, the scatter matrix is used in FLD, not the covariance matrix. That is, the FLD objective is

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\mathbf{w}^T (S_1 + S_2) \mathbf{w}}. \quad (20)$$

However, we want to point out the solutions for both forms should be similar to each other qualitatively. When the number of examples in all classes are the same, these two should give the same answer.

1.4 The two scatter matrices and the FLD objective

In FLD (for binary problems), two additional symbols are introduced,

$$S_B = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T, \quad (21)$$

$$S_W = S_1 + S_2, \quad (22)$$

in which S_B is called the between-class scatter matrix, and S_W is the within-class scatter matrix (both are $D \times D$ in size). S_W measures how *scattered* is the original input dataset *within* each class. S_B measures the scatter caused by the two class means, which measures how scattered it is *between* different classes.

Now, we can formally define the FLD objective as

$$J = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}. \quad (23)$$

In other words, this optimization aims at finding a projection direction that makes the *between-class scatter much larger than the within-class scatter*.

A side note will be useful soon when we deal with multiclass problems. In Equation 21, one class is approximately represented by its mean (as \mathbf{m}_1 and \mathbf{m}_2 , respectively). This is a reasonable simplification, especially when the scatter within each class is not large.

However, if one mean vector acts as a surrogate for all examples in its respective class, we would expect that the *scatter matrix of the set of mean vectors* to

be the between-class scatter, i.e., we should expect that

$$\sum_{i=1}^2 (\mathbf{m}_i - \bar{\mathbf{m}})(\mathbf{m}_i - \bar{\mathbf{m}})^T \quad (24)$$

is the between-class scatter, rather than S_B . In Equation 24,

$$\bar{\mathbf{m}} = \frac{\mathbf{m}_1 + \mathbf{m}_2}{2}.$$

These two terms (Equation 24 and Equation 21) are not exactly the same, however, it is easy to verify that

$$\frac{1}{2}S_B = \sum_{i=1}^2 (\mathbf{m}_i - \bar{\mathbf{m}})(\mathbf{m}_i - \bar{\mathbf{m}})^T. \quad (25)$$

Hence, it is reasonable to define S_B in the form of Equation 21, because a constant factor $\frac{1}{2}$ will not change the optimal solution for Equation 23.

1.5 The optimization

The right hand side of Equation 23 is called the generalized Rayleigh quotient, or the Rayleigh-Ritz ratio, whose maximization is not only useful in computer science (e.g., in FLD), but also in other subjects (such as physics).

Finding the derivative of J with respect to \mathbf{w} and set that to 0, we get

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{2((\mathbf{w}^T S_W \mathbf{w}) S_B \mathbf{w} - (\mathbf{w}^T S_B \mathbf{w}) S_W \mathbf{w})}{(\mathbf{w}^T S_W \mathbf{w})^2} = 0. \quad (26)$$

Hence, one necessary condition for optimality is

$$S_B \mathbf{w} = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}} S_W \mathbf{w}. \quad (27)$$

Noticing that $\frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$ is a scalar value, this condition is stating that \mathbf{w} should be a *generalized eigenvector* of S_B and S_W , with $\frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$ being its corresponding generalized eigenvalue!

Since $J = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$ is both the generalized eigenvalue and the optimization objective, we should use the generalized eigenvector of S_B and S_W whose corresponding generalized eigenvalue is the largest, which gives the optimal parameter \mathbf{w}^* that maximizes J .

1.6 Wait, we have a shortcut!

Although there are efficient routines to find the generalized eigenvalues and eigenvectors, to solve FLD for a binary problem we have an easier way—if we observe carefully we will find it.

Some small tweaking to Equation 27 tells us that

$$S_W \mathbf{w} = \frac{\mathbf{w}^T S_W \mathbf{w}}{\mathbf{w}^T S_B \mathbf{w}} S_B \mathbf{w} \quad (28)$$

$$= \frac{\mathbf{w}^T S_W \mathbf{w}}{\mathbf{w}^T S_B \mathbf{w}} (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \quad (29)$$

$$= \frac{\mathbf{w}^T S_W \mathbf{w}}{\mathbf{w}^T S_B \mathbf{w}} (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} (\mathbf{m}_1 - \mathbf{m}_2) \quad (30)$$

$$= c(\mathbf{m}_1 - \mathbf{m}_2), \quad (31)$$

in which we define

$$c \triangleq \frac{\mathbf{w}^T S_W \mathbf{w}}{\mathbf{w}^T S_B \mathbf{w}} (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w},$$

and c is a scalar value. Hence, this optimality condition directly leads to the optimal projection direction:

$$S_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2). \quad (32)$$

Then, by normalizing this direction to make its ℓ_2 norm equal 1, we get the optimal \mathbf{w} . In the above equation, we have omitted the factor c , because it will disappear anyway in the normalization process. More details of data normalization will be discussed in Chapter 9.

1.7 The FLD method for binary problems

Algorithm 1 describes the steps to find the FLD projection direction for a binary classification problem.

Algorithm 1 The FLD algorithm for a binary problem

- 1: **Input:** a D -dimensional binary training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$.
- 2: Compute \mathbf{m}_1 , \mathbf{m}_2 , and S_W following the equations in this chapter.
- 3: Compute

$$\mathbf{w} \leftarrow S_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2).$$

- 4: Normalize:

$$\mathbf{w} \leftarrow \frac{\mathbf{w}}{\|\mathbf{w}\|}.$$

The idea of FLD was originally described by Ronald Fisher, a famous statistician and biologist. This method is named after him, as *Fisher's Linear Discriminant*, abbreviated as FLD. This name is sometimes used interchangeably with LDA (which stands for Linear Discriminant Analysis). LDA, however, is often used to refer to any linear discriminant function, including FLD but not limited to it.

1.8 A caveat: what if S_W is not invertible?

There is an obvious shortcoming of Algorithm 1: it is not well-defined if S_W is not invertible.

If there are more positive (or negative) examples than the number of dimensions, i.e., if $N_1 > D$ or $N_2 > D$, S_W will be invertible with high probability. However, if S_W^{-1} does not exist in some cases, we can use the Moore-Penrose (or MP for short) pseudoinverse to replace the inverse. A detailed coverage of the Moore-Penrose pseudoinverse is beyond the scope of this book, we will only briefly explain how to compute it for S_W .

Instead of a $^{-1}$ superscript, the Moore-Penrose pseudoinverse is usually denoted by a $^{+}$ superscript. It is easy to compute the MP pseudoinverse for an 1×1 matrix, i.e., a scalar $x \in \mathbb{R}$, following the rule

$$x^{+} = \begin{cases} 0 & \text{if } x = 0 \\ \frac{1}{x} & \text{otherwise} \end{cases}. \quad (33)$$

Then, a diagonal matrix

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_D)$$

has its MP pseudoinverse as

$$\Lambda^{+} = \text{diag}(\lambda_1^{+}, \lambda_2^{+}, \dots, \lambda_D^{+}). \quad (34)$$

Note that when Λ is invertible (i.e., when $\lambda_i \neq 0$ for $1 \leq i \leq D$), the MP pseudoinverse is the same as the usual matrix inverse.

Since S_W is the sum of two covariance matrices, it is positive semi-definite. Hence, we can find its spectral decomposition as

$$S_W = E \Lambda E^T,$$

where the diagonal matrix Λ contains the eigenvalues of S_W and the columns of the orthogonal matrix E contain the eigenvectors of S_W . The Moore-Penrose pseudoinverse of S_W is then

$$S_W^{+} = E \Lambda^{+} E^T. \quad (35)$$

Note that when S_W is invertible, $S_W^{+} = S_W^{-1}$.

2 FLD for more classes

When our data is multiclass, i.e., when there are more than two possible labels, we need to extend Algorithm 1 to handle this situation.

We still have a training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$. However, y_i now has a different domain. Let $\mathcal{Y} = \{1, 2, \dots, K\}$ denotes the K possible labels in a K -class problem, and $y_i \in \mathcal{Y}$ for all $1 \leq i \leq N$.

The idea to let the projection values of examples in these K classes to be far apart from each other still makes sense in this more complex setting. But it is more difficult to translate this idea to precise mathematics.

In the binary case, we maximize $\frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$ (cf. Equation 23), whose meaning is to make the between-class variance (or scatter) much larger than the within-class variance (or scatter). This intuition may still be valid in multiclass problems. However, the key issue becomes: how to define these two variances (or scatter) when we have $K > 2$ classes (i.e., how do we define S_B and S_W for this new setup?)

2.1 A slightly modified notation and S_W

The within-class scatter matrix S_W can be easily extended to multiclass, and keep its semantic meaning. We just need to modify the notation a bit.

Let X_k , N_k , \mathbf{m}_k , C_k and S_k be the subset of X that belongs to the k -th class, its size, mean, covariance matrix and scatter matrix, respectively, $1 \leq k \leq K$, i.e.,

$$X_k = \{\mathbf{x}_i | y_i = k\}, \quad (36)$$

$$N_k = |X_k|, \quad (37)$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{\mathbf{x} \in X_k} \mathbf{x}, \quad (38)$$

$$C_k = \frac{1}{N_k} \sum_{\mathbf{x} \in X_k} (\mathbf{x} - \mathbf{m}_k)(\mathbf{x} - \mathbf{m}_k)^T, \quad (39)$$

$$S_k = N_k C_k = \sum_{\mathbf{x} \in X_k} (\mathbf{x} - \mathbf{m}_k)(\mathbf{x} - \mathbf{m}_k)^T. \quad (40)$$

The within-class scatter is then simply the sum of scatter matrices of the K subsets, which measures the overall within-class scatter of the multiclass training set,

$$S_W = \sum_{k=1}^K S_k. \quad (41)$$

2.2 Candidates for S_B

It is, however, not trivial to extend the between-class scatter matrix S_B to multiclass in a natural way. A straightforward extension might be

$$\sum_{i=1}^K \sum_{j=1}^K (\mathbf{m}_i - \mathbf{m}_j)(\mathbf{m}_i - \mathbf{m}_j)^T, \quad (42)$$

which mimics Equation 21. But, this equation requires the creation of $K^2 D \times D$ matrices, and also need to sum them up. Hence, it is not very attractive in terms of computational efficiency.

Equation 25 may lead to another possible extension

$$\sum_{k=1}^K (\mathbf{m}_k - \bar{\mathbf{m}})(\mathbf{m}_k - \bar{\mathbf{m}})^T, \quad (43)$$

where

$$\bar{\mathbf{m}} = \frac{1}{K} \sum_{k=1}^K \mathbf{m}_k. \quad (44)$$

Equation 43 is the scatter of the K mean vectors, which fits our intuition of the between-class scatter and is also fast to compute.

However, researchers have found another definition, which makes even more sense.

2.3 A tale of three scatter matrices

Note we can compute the scatter matrix for the entire training set regardless of the class labels. We denote it as S_T ,

$$S_T = \sum_{i=1}^N (\mathbf{x}_i - \mathbf{m})(\mathbf{x}_i - \mathbf{m})^T, \quad (45)$$

where

$$\mathbf{m} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (46)$$

is the mean of all training points. Note that Equation 46 (\mathbf{m}) is different from Equation 44 ($\bar{\mathbf{m}}$).

S_T is called the *total scatter* matrix. It is curious to see what is the difference between the total and within-class scatter, and some algebraic manipulation gives us

$$S_T - S_W = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T. \quad (47)$$

The right hand side of Equation 47 closely resembles the scatter matrix in Equation 43, with only two slight exceptions. First, \mathbf{m} replaces $\bar{\mathbf{m}}$; second, N_k is multiplied before each term in the summation.

In FLD, the between-class scatter is defined as

$$S_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T. \quad (48)$$

We want to point out that when the classes have equal sizes, i.e., when $N_i = N_j$ for any i and j , the above two deviations are trivial. In that case, $\mathbf{m} = \bar{\mathbf{m}}$ and Equation 48 differs from Equation 43 only by a constant multiplier $\frac{K}{N}$.

Furthermore, these changes are effective in imbalanced problems. Consider an imbalanced problem where $N_i \gg N_j$ for some i and j . We expect \mathbf{m}_i to be more important than \mathbf{m}_j . Both Equations 46 and 48 put more weights on \mathbf{m}_i than \mathbf{m}_j , while they are treated equally in Equations 44 and 43.

Finally, it makes sense if the sum of between- and within-class scatter matrices equals the total scatter matrix, in which case we have an easy-to-remember and reasonable equation that describes the tale of three scatter matrices:

$$S_T = S_B + S_W . \quad (49)$$

We also point out that: when $K = 2$, Equation 48 is different from Equation 21. When we are dealing with an imbalanced binary problem (e.g., when the negative class has much more examples than the positive class), it might help if we define S_B following Equation 48 even in the binary setup.

2.4 The solution

Since S_B is not a rank 1 matrix in multiclass problems, Algorithm 1 is not applicable. However, we can still find the optimal projection direction by solving a generalized eigenvalue problem

$$S_B \mathbf{w} = \lambda S_W \mathbf{w} .$$

When S_W is invertible, the generalized eigenvalue problem is equivalent to an eigenvalue problem

$$S_W^{-1} S_B \mathbf{w} = \lambda \mathbf{w} .$$

However, directly solving the generalized eigenvalue problem is more efficient. For example, in Matlab we can obtain the solution using a simple command `eig(A,B)`, where $A = S_B$ and $B = S_W$, and retrieve the generalized eigenvector associated with the largest generalized eigenvalue.

2.5 Finding more projection directions

As in PCA, we can find more projection directions for multiclass FLD: just use the generalized eigenvectors associated with the top $K - 1$ largest generalized eigenvalues.

For a K class problem, we can extract at most $K - 1$ meaningful features (i.e., projected values). We will leave the proof of this statement to the readers as an exercise problem.

The generalized eigenvectors are not necessarily ℓ_2 normalized or perpendicular to each other.

Exercises

1. (Matrix rank) Let A be a real matrix of size $m \times n$. Its rank, denoted as $\text{rank}(A)$, is defined as the maximal number of linearly independent rows in A , which is also called the row rank. Similarly, the column rank is the maximal number of linearly independent columns in A .

(a) Prove that the row rank is equal to the column rank. Hence,

$$\text{rank}(X) = \text{rank}(X^T).$$

(b) Prove that

$$\text{rank}(A) \leq \min(m, n).$$

(c) Let X and Y be two matrices of the same size. Show that

$$\text{rank}(X + Y) \leq \text{rank}(X) + \text{rank}(Y).$$

(d) Show that

$$\text{rank}(XY) \leq \min(\text{rank}(X), \text{rank}(Y))$$

so long as the matrix multiplication is well-defined.

(e) A matrix X of size $m \times n$ is called *full rank* if

$$\text{rank}(X) = \min(m, n).$$

Prove that if X is full rank, then

$$\text{rank}(X) = \text{rank}(XX^T) = \text{rank}(X^TX).$$

(f) Let \mathbf{x} be a vector, what is the rank of $\mathbf{x}\mathbf{x}^T$? Show that the rank of a real symmetric matrix X equals its number of non-zero eigenvalues. This is the way we used to define matrix rank in Chapter 2.

2. (Matrix norm) Similar to vectors norms, various matrix norms can be defined. $f : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$ is a matrix norm if it satisfies the following conditions:

- (a) $f(X) \geq 0$ for any $X \in \mathbb{R}^{m \times n}$;
- (b) $f(X) = 0$ if and only if X is a zero matrix;
- (c) $f(X + Y) \leq f(X) + f(Y)$ for $X, Y \in \mathbb{R}^{m \times n}$;
- (d) $f(cX) = |c|f(X)$ for $c \in \mathbb{R}$ and $X \in \mathbb{R}^{m \times n}$.

(a) Prove that

$$\|X\|_2 = \max_{\|\mathbf{v}\|=1} \|X\mathbf{v}\|$$

is a valid matrix norm (which is called the 2-norm). The vector norm considered here is the vector 2-norm $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}}$.

(b) Let $\sigma_{\max}(X)$ denote the largest singular value of a matrix X . Prove that if X and Y are of the same size, then

$$\sigma_{\max}(X) + \sigma_{\max}(Y) \geq \sigma_{\max}(X + Y).$$

(Hint: what is the relationship between $\sigma_{\max}(X)$ and $\|X\|_2$?)

3. (Condition number) Given any matrix norm $\|\cdot\|$, a corresponding *condition number* can be defined for a non-singular real square matrix. The condition number of a matrix X is defined as

$$\kappa(X) = \|X\| \|X^{-1}\|.$$

One frequently used condition number is the 2-norm condition, as

$$\kappa_2(X) = \|X\|_2 \|X^{-1}\|_2.$$

A matrix is called *ill-conditioned* if its condition number is large.

(a) If you already know the singular values of X as $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$, what is the condition number $\kappa_2(X)$?

(b) Let f be a function from a set \mathbb{X} to \mathbb{Y} . Suppose $f(\mathbf{x}) = \mathbf{y}$ and $f(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{y} + \Delta\mathbf{y}$. We call f as *ill-conditioned* if a small $\Delta\mathbf{x}$ leads to a large $\Delta\mathbf{y}$. Let A be a square full-rank matrix (i.e., invertible matrix) in $\mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$, and we are interested in solving $A\mathbf{x} = \mathbf{b}$. Show that this linear system is ill-conditioned if $\kappa_2(A)$ is large. (You are not required to prove this conclusion. Just give some intuitions why an ill-conditioned matrix A is bad.)

(c) Show that an orthogonal matrix is well-conditioned (i.e., have small condition number).

4. (LU, LDL and Cholesky factorizations) In this problem, we study various factorizations that are useful in computing the generalized eigenvalue problem.

(a) Show that
$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 4 \\ 3 & 4 & 6 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -1 & -2 \\ 0 & -2 & -3 \end{bmatrix}.$$

(b) (Gauss transformation) In a more general form, let \mathbf{v} be a vector in \mathbb{R}^n and $v_k \neq 0$. Then, the *Gauss transformation* M_k is an $n \times n$ lower triangular matrix of the form:

$$M_k = \begin{bmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & -\frac{v_{k+1}}{v_k} & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & -\frac{v_n}{v_k} & 0 & \cdots & 1 \end{bmatrix}. \quad (50)$$

That is, entries in M_k are all zeros except those in the diagonals and those in the k -th column. The diagonals are all ones. In the k -th column, the first $k - 1$ entries are zeros, the k -th entry is 1, and the i -th entry ($i > k$) is $-\frac{v_i}{v_k}$. What is the result of left-multiplying M_k to \mathbf{v} ?

(c) For an $n \times n$ matrix A , what is the effect of running the following algorithm?

```

1:  $A^{(1)} = A$ 
2: for  $k = 1, 2, \dots, n - 1$  do
3:   Let  $\mathbf{v} = A_{:,k}^{(k)}$  (i.e., the  $k$ -th column of  $A^{(k)}$ ), and compute the Gauss
     transformation  $M_k$  based on  $\mathbf{v}$ 
4:    $A^{(k+1)} \leftarrow M_k A^{(k)}$ 
5: end for
```

In the above algorithm, we assume $a_{kk}^{(k)} \neq 0$ for $1 \leq k \leq n - 1$.

(d) Show that $\det(M_k) = 1$ for all $1 \leq k < n$. If $L = M_{n-1}M_{n-2} \dots M_1$, prove that L is lower triangular and its determinant is 1.

(e) (LU Factorization) Let $A \in \mathbb{R}^{n \times n}$. And, for $1 \leq k \leq n - 1$, the leading principal submatrix $A_{1:k, 1:k}$ is non-singular. Prove that there exists a decomposition

$$A = LR,$$

where $L \in \mathbb{R}^{n \times n}$ is lower triangular with its diagonal entries all being 1, and $U \in \mathbb{R}^{n \times n}$ is upper triangular. If A is furthermore non-singular, prove that the LU decomposition is *unique*.

(f) (LDL Factorization) We pay our attention to the LDL factorization of a real symmetric and positive definite matrix. If A is such a matrix, prove that there exist a *unique* lower triangular real matrix L (with its diagonal entries all being 1) and a diagonal matrix D such that

$$A = LDL^T.$$

(g) (Cholesky Factorization) Let A be a real symmetric positive definite matrix. Then, prove that there exists a *unique* real lower triangular matrix G such that

$$A = GG^T,$$

and the diagonal entries in G are all positive. This matrix decomposition is called the *Cholesky* factorization or Cholesky decomposition, named after the French mathematician Andre-Louis Cholesky.

5. (Solving FLD) Consider a C ($C > 2$) class FLD problem, in which there are N examples of D dimensions.

(a) If $N < D$, will S_W be invertible?

(b) Prove that at most $C - 1$ generalized eigenvectors can be obtained. (Hint: What is the rank of S_B ?)

(c) Explain why S_W is very likely invertible if $N > D$, especially if $N \gg$

D ?

(d) The generalized eigenvalue problem $S_B \mathbf{w} = \lambda S_W \mathbf{w}$ is equivalent to $S_W^{-1} S_B \mathbf{w} = \lambda \mathbf{w}$ when S_W is invertible. However, in practice we do not solve $S_W^{-1} S_B \mathbf{w} = \lambda \mathbf{w}$. One algorithm for computing the FLD solution is as follows.

First, assuming S_W is positive definite, the Cholesky factorization finds a G such that $S_W = GG^T$.

Second, compute $C = G^{-1} S_B G^{-T}$.

Third, diagonalize C such that there is an orthogonal matrix Q and $Q^T C Q$ is a diagonal matrix.

Finally, let $X = G^{-T} Q$.

Prove that $X^T S_B X$ is diagonal and $X^T S_W X$ is the identity matrix. Show that the generalized eigenvectors are in the columns of X and the generalized eigenvalues are in the diagonal entries of $X^T S_B X$. Are the generalized eigenvectors computed this way unit norm or orthogonal?

6. (PCA+FLD face recognition) Both PCA and FLD are useful in the face recognition application.

(a) The ORL dataset is an early dataset for face recognition. Download the ORL dataset from <http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>. Read the instructions in that page to understand the format of this dataset.

(b) OpenCV is an open source computer vision library, which provides many functions useful for various computer vision applications. Download the library from <http://opencv.org/>. Learn OpenCV basics from <http://docs.opencv.org/>.

- 10' (c) There is an OpenCV tutorial on face recognition, which is available at http://docs.opencv.org/2.4/modules/contrib/doc/facerec/facerec_tutorial.html. Try to understand every line of code in the tutorial, especially those on Eigenfaces (PCA) and Fisherfaces (FLD). Run experiments on the ORL dataset, analyze the differences between recognition results of these methods. 分析
- 10' (d) In the Eigenface experiment, you can reconstruct an approximate face image using the eigenfaces (i.e., eigenvectors). Modify the source code in the OpenCV tutorial, and observe the visualization using different amount of eigenfaces. How many eigenfaces are required if you want the face reconstructed from eigenfaces are visually indistinguishable from the original input face image?