

Notes on Fisher's Linear Discriminant

Zhihan Yang

March 27, 2020

1 Binary classification problem

There are two classes, C_1 and C_2 , each containing a non-zero number of data vectors, \mathbf{x}^n . For each class, C_k , we define some properties that will be used later on in this document.

The number of vectors in class k is given by N_k .

The mean vector of class k is given by:

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}^n$$

The projected mean of class k is given by:

$$m_k = \mathbf{w}^T \mathbf{m}_k$$

The within-class variance of class k is given by:

$$s_k = \sum_{n \in C_k} (y^n - m_k)^2$$

2 Fisher's criterion

A linear combination between input variables, x_i , and weights w_i , can be written in vector form as:

$$y = \mathbf{w}^T \mathbf{x}$$

where y can be interpreted as the projection of \mathbf{x} onto \mathbf{w} .

Given two classes, C_1 and C_2 , we would like to select a \mathbf{w} such that the projection of \mathbf{x} maximizes class separation. Intuitively, this means maximizing between-class separation while minimizing within-class separation. Mathematically, we define between-class separation as $(m_2 - m_1)^2$, and define within-class separation as $s_1^2 + s_2^2$. We can complete the two optimization tasks simultaneously by maximizing the following ratio, known as the Fisher criterion:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

Now, the question remains: how do we find the optimal \mathbf{w} to maximize $J(\mathbf{w})$?

3 The Fisher criterion in matrix form

It turns out that the Fisher criterion is more conveniently optimized in matrix form. In this section, we convert the algebraic form derived in the last section into its matrix form.

Expand the numerator:

$$\text{numerator}(J(\mathbf{w})) = (m_2 - m_1)^2 \tag{1}$$

$$= (m_2 - m_1)(m_2 - m_1) \tag{2}$$

$$= (\vec{w}^T \vec{x} - \vec{w}^T \vec{m}_1)(\vec{w}^T \vec{x} - \vec{w}^T \vec{m}_2) \tag{3}$$

where both factors are scalars, so the order of multiplication does not matter within each factor,

$$= \vec{w}^T (\vec{x} - \vec{m}_1)(\vec{x} - \vec{m}_2)^T \vec{w} \tag{4}$$

Since matrix multiplication is associative, i.e., $A(BC) = (AB)C$, we can rewrite the numerator as:

$$\text{numerator}(J(\mathbf{w})) = \vec{w}^T S_B \vec{w}$$

where $S_B = (\vec{x} - \vec{m}_1)(\vec{x} - \vec{m}_2)^T$ and is called the between-class covariance matrix.

Expand the denominator (somewhat tedious):

$$\text{denominator}(J(\vec{w})) = s_1^2 + s_2^2 \quad (5)$$

$$= \sum_{n \in C_1} (y^n - m_1)^2 + \sum_{n \in C_2} (y^n - m_2)^2 \quad (6)$$

$$= \sum_{n \in C_1} (\vec{w}^T \vec{x} - \vec{w}^T \vec{m}_1)^2 + \sum_{n \in C_2} (\vec{w}^T \vec{x} - \vec{x}^T \vec{m}_2)^2 \quad (7)$$

$$= \sum_{n \in C_1} \vec{w}^T (\vec{x} - \vec{m}_1) (\vec{x} - \vec{m}_1)^T \vec{w} + \sum_{n \in C_2} \vec{w}^T (\vec{x} - \vec{m}_2) (\vec{x} - \vec{m}_2)^T \vec{w} \quad (8)$$

$$= \vec{w}^T \left\{ \sum_{n \in C_1} (\vec{x} - \vec{m}_1) (\vec{x} - \vec{m}_1)^T \right\} \vec{w} + \vec{w}^T \left\{ \sum_{n \in C_2} (\vec{x} - \vec{m}_2) (\vec{x} - \vec{m}_2)^T \right\} \vec{w} \quad (9)$$

$$= \vec{w}^T \left\{ \sum_{n \in C_1} (\vec{x} - \vec{m}_1) (\vec{x} - \vec{m}_1)^T + \sum_{n \in C_2} (\vec{x} - \vec{m}_2) (\vec{x} - \vec{m}_2)^T \right\} \vec{w} \quad (10)$$

$$= \vec{w}^T S_W \vec{w} \quad (11)$$

where $S_W = \sum_{n \in C_1} (\vec{x} - \vec{m}_1) (\vec{x} - \vec{m}_1)^T + \sum_{n \in C_2} (\vec{x} - \vec{m}_2) (\vec{x} - \vec{m}_2)^T$ and is called the within-class covariance matrix.

Putting the numerator and the denominator together, obtain:

$$J(\vec{w}) = \frac{\vec{w}^T S_B \vec{w}}{\vec{w}^T S_W \vec{w}}$$

4 Fisher's linear discriminant

Since $J(\vec{w})$ involves \vec{w} in a quadratic fashion, $\frac{\partial J(\vec{w})}{\partial \vec{w}}$ involves \vec{w} in a linear fashion and hence $\frac{\partial J(\vec{w})}{\partial \vec{w}} = 0$ has a unique solution for \vec{w} . First, let's find $\frac{\partial J(\vec{w})}{\partial \vec{w}}$.

$$J(\vec{w}) = \frac{\vec{w}^T S_B \vec{w}}{\vec{w}^T S_W \vec{w}} \quad (12)$$

In order to obtain $\frac{\partial J(\vec{w})}{\partial \vec{w}}$, we need the quotient rule for matrix calculus, which requires us to define the numerator, u , the denominator, v , the derivative of numerator, u' , and the derivative of denominator, v' . Let's define them below:

$$u = \vec{w}^T S_B \vec{w}$$

$$v = \vec{w}^T S_W \vec{w}$$

According to the Matrix Cookbook,

$$u' = \frac{\partial u}{\partial \vec{w}} = (S_B + S_B^T) \vec{w}$$

Derivation of $(S_B + S_B^T) \vec{w} = 2S_B \vec{w}$:

$$S_B = (\vec{m}_2 - \vec{m}_1) (\vec{m}_2 - \vec{m}_1)^T$$

$$S_B^T = \left((\vec{m}_2 - \vec{m}_1) (\vec{m}_2 - \vec{m}_1)^T \right)^T$$

$$= (\vec{m}_2 - \vec{m}_1)^{TT} (\vec{m}_2 - \vec{m}_1)^T$$

$$= (\vec{m}_2 - \vec{m}_1) (\vec{m}_2 - \vec{m}_1)^T$$

Therefore, $S_B^T = S_B$ and $u' = 2S_B \vec{w}$.

Again, according to the Matrix Cookbook,

$$v' = \frac{\partial v}{\partial \vec{w}} = (S_W + S_W^T) \vec{w}$$

Derivation of $(S_W + S_W^T) \vec{w} = 2S_W \vec{w}$:

$$S_W = \sum_{n \in C_1} (\vec{x}^n - \vec{m}_1) (\vec{x}^n - \vec{m}_1)^T + \sum_{n \in C_2} (\vec{x}^n - \vec{m}_2) (\vec{x}^n - \vec{m}_2)^T$$

Since a transpose of sums is a sum of transposes,

$$S_W^T = \sum_{n \in C_1} \left[(\vec{x}^n - \vec{m}_1) (\vec{x}^n - \vec{m}_1)^T \right]^T + \sum_{n \in C_2} \left[(\vec{x}^n - \vec{m}_2) (\vec{x}^n - \vec{m}_2)^T \right]^T$$

$$= \sum_{n \in C_1} (\vec{x}^n - \vec{m}_1) (\vec{x}^n - \vec{m}_1)^T + \sum_{n \in C_2} (\vec{x}^n - \vec{m}_2) (\vec{x}^n - \vec{m}_2)^T$$

Therefore, $S_W^T = S_W$ and $v' = 2S_W \vec{w}$.

Using u , u' , v and v' and the quotient rule to find $\frac{\partial J(\vec{w})}{\partial \vec{w}}$:

$$\frac{\partial J(\vec{w})}{\partial \vec{w}} = \frac{u'v - uv'}{v^2} \quad (13)$$

$$= \frac{(2S_B \vec{w}) (\vec{w}^\top S_W \vec{w}) - (\vec{w}^\top S_B \vec{w}) (2S_w \vec{w})}{(\vec{w}^\top S_W \vec{w}) (\vec{w}^\top S_W \vec{w})} \quad (14)$$

$$= \frac{2S_B \vec{w}}{\vec{w}^\top S_W \vec{w}} - \frac{(\vec{w}^\top S_B \vec{w}) (2S_w \vec{w})}{(\vec{w}^\top S_W \vec{w}) (\vec{w}^\top S_W \vec{w})} \quad (15)$$

Setting $\frac{\partial J(\vec{w})}{\partial \vec{w}} = 0$ and solving for this equation:

$$\frac{\partial J(\vec{w})}{\partial \vec{w}} = 0 \quad (16)$$

$$\frac{2S_B \vec{w}}{\vec{w}^\top S_w \vec{w}} = \frac{(\vec{w}^\top S_B \vec{w}) (2S_w \vec{w})}{(\vec{w}^\top S_w \vec{w}) (\vec{w}^\top S_w \vec{w})} \quad (17)$$

$$\underbrace{2 (\vec{w}^\top S_w \vec{w})}_{\text{scalar}} (S_B \vec{w}) = \underbrace{2 (\vec{w}^\top S_B \vec{w})}_{\text{scalar}} (S_w \vec{w}) \quad (18)$$

Since we do not care about the magnitude of \vec{w} , only its direction, we remove the scalars:

$$S_B \vec{w} \propto S_w \vec{w}.$$

Since $S_B = (\vec{m}_2 - \vec{m}_1)(\vec{m}_2 - \vec{m}_1)^T$ and $(\vec{m}_2 - \vec{m}_1)^T \vec{w}$ is a scalar, $S_B \vec{w}$ is always in the direction of $\vec{m}_2 - \vec{m}_1$:

$$(\vec{m}_2 - \vec{m}_1) \propto S_w \vec{w}$$

Multiplying both sides by S_w^{-1} :

$$\vec{w} \propto S_w^{-1}(\vec{m}_2 - \vec{m}_1)$$

"This choice of \vec{w} is known as Fisher's linear discriminant, although it is strictly not a discriminant but rather a specific choice of direction for projection of data down to one dimension." (Bishop, 1995)