

Module 1.3: CLASSIFICATION

LINEAR (FISHER'S) DISCRIMINANT ANALYSIS

Suppose we have vectors $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N$ and we know that N_1 of these vectors belong to class D_1 , N_2 of them belong to class D_2 , and so forth for c classes. Assume $\vec{x} \in \mathbb{R}^n$

We want to find a vector in \mathbb{R}^n (and in fact set of vectors) such that when we project all $\vec{x} \in D_i$, they will be distant from the projection of other $\vec{x} \in D_j, j \neq i$ onto the same vector.

A good way to do this is by having the mean values of the projections of all $\vec{x} \in D_i$ distant from the projected mean value of the entire data set.

Consider a candidate vector w , the projection of $\vec{x} \in D_i$ onto w is

$$\Rightarrow y_i = \vec{w}^T \vec{x}_i \quad \text{where } X_i \text{ is a matrix whose columns are } \vec{x} \in D_i \text{ row vector}$$

The mean of the projections of $\vec{x} \in D_i$ onto w is

$$\tilde{\mu}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} y_j \vec{1}_i^T = \vec{w}^T \left(\frac{1}{N_i} \vec{x}_i \vec{1}_i \right) \quad \text{where } \vec{1}_i = \underbrace{[1, 1, \dots, 1]}_{N_i}^T$$

$$= \vec{w}^T \vec{\mu}_i$$

where $\vec{\mu}_i$ is the mean of the vectors of $\vec{x} \in D_i$

$$\vec{\mu}_i = \frac{1}{N_i} \sum_{\vec{x} \in D_i} \vec{x}$$

Now, the mean of all the samples is $\vec{\mu} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i$, and its projection to w is

$$\tilde{\mu} = \vec{w}^T \vec{\mu}$$

We want

$$\begin{aligned} \sum_{i=1}^c N_i (\tilde{\mu}_i - \tilde{\mu})^2 &= \sum_{i=1}^c N_i \vec{w}^\top (\vec{\mu}_i - \vec{\mu}) (\vec{\mu}_i - \vec{\mu})^\top \vec{w} \\ &= \vec{w}^\top \underbrace{\left(\sum_{i=1}^c N_i (\vec{\mu}_i - \vec{\mu}) (\vec{\mu}_i - \vec{\mu})^\top \right)}_{S_B} \vec{w} \end{aligned}$$

$S_B \leftarrow$ Between class scatter matrix

to be as large as possible.

Meanwhile we want to variance of the projections, within a class,

$$\begin{aligned} \sum_{y \in \mathcal{Y}_i} (y - \tilde{\mu}_i)^2 &= \sum_{x \in \mathcal{D}_i} (\vec{w}^\top \vec{x} - \vec{w}^\top \vec{\mu}_i)^2 \\ &= \sum_{x \in \mathcal{D}_i} \vec{w}^\top (\vec{x} - \vec{\mu}_i) (\vec{x} - \vec{\mu}_i)^\top \vec{w} \\ &= \vec{w}^\top \underbrace{\left(\sum_{x \in \mathcal{D}_i} (\vec{x} - \vec{\mu}_i) (\vec{x} - \vec{\mu}_i)^\top \right)}_{S_i} \vec{w} \end{aligned}$$

$S_i \leftarrow$ Scatter matrix for class i

to be small as possible.

Let $S_W = \sum_{i=1}^c S_i$ be the total within-class scatter matrix.

We want

$$J(\vec{w}) = \frac{\vec{w}^\top S_B \vec{w}}{\vec{w}^\top \left(\sum_{i=1}^c S_i \right) \vec{w}} = \frac{\vec{w}^\top S_B \vec{w}}{\vec{w}^\top S_W \vec{w}}$$

to be maximized (with $\|\vec{w}\| = 1$). We need \vec{w}^* that accomplishes this

$$\omega^* = \underset{\omega}{\operatorname{argmax}} \frac{\omega^T S_B \omega}{\omega^T S_W \omega}$$

$$\text{or } \omega^* = \underset{\omega}{\operatorname{argmax}} \omega^T S_B \omega \text{ subject to } \omega^T S_W \omega = K$$

The Lagrangian is given as

$$L = \omega^T S_B \omega - \lambda (\omega^T S_W \omega - K) = \omega^T (S_B - \lambda S_W) \omega + 2K$$

$$\nabla_{\omega} L = 2(S_B - \lambda S_W)\omega = 0 \Rightarrow S_B \omega^* = \lambda S_W \omega^*$$

which is a generalized eigenvalue problem.

If S_W is full-rank (or invertible), we can write

$$S_W^{-1} S_B \omega^* = \lambda \omega^*$$

From our definition, $S_B = \sum_{i=1}^c N_i (\vec{\mu}_i - \vec{\mu}) (\vec{\mu}_i - \vec{\mu})^T \in \mathbb{R}^{n \times n}$ has, at most,

rank $c-1$. Meanwhile, $S_W = \sum_{i=1}^c \sum_{x_i \in D_i} (x - \mu_i)(x - \mu_i)^T \in \mathbb{R}^{n \times n}$ has, at

$$\text{most rank } c(N-1) = N - c$$

If the data size $N > c+n$, and $c > n$, S_B and S_W are generally likely to be full rank

For projection onto more vectors, in general, the maximization problem is

$$W^* = \underset{W}{\operatorname{argmax}} \frac{|W^T S_B W|}{|W^T S_W W|} = \underset{W}{\operatorname{argmax}} \frac{|\tilde{S}_B|}{|S_W|}$$

Which yields an extension of the previous solution

$$S_B w_k^* = \lambda S_W w_k^* \quad k=1, \dots, m$$

But the issue of S_W being singular still persists. We can address this the following ways:

① S_W is very likely to have a rank of $N-c$, we can reduce the

dimension of our data from n to $N-c$. And what better way can we preserve the information in our data while reducing its dimension than by using PCA? So we pick out the first $N-c$, components or less of our PCA transformation

$$y_i' = y_{(1:N-c)} = (T(:, 1:N-c))^T x_i \quad v=1, \dots, N$$

We now carry out finding S_B' & S_W' using y_i'
& find $W' \in \mathbb{R}^{N-c \times m}$

$$\begin{matrix} N \\ \overbrace{\quad \quad \quad}^{N-c} [y'] \end{matrix}$$

② We can find the the Moore-Penrose Inverse (Pseudoinverse)

$$S_W' S_B' w_k^* = \lambda w_k^*$$

NOTE: In ①

We have taken two steps which can be expressed as

$$y_{FLD} = \underbrace{W_{FLD}^T T_{1:N-C}^T}_{W_{opt}^T} x$$

$m \times N-C \times N-C \times n$
 $= m \times n$

Note that because S_B os rank $C-1$, m should be, at most, $C-1$

The columns of W_{opt} are called Fisherfaces.

③ Regularize S_w :

$$S_w' + \epsilon I_d \quad \text{where } \epsilon \text{ is a small constant}$$

You can show that for a 2-class case:

$$S_w = \frac{1}{\lambda} S_w^{-1} S_B w = S^{-1} (\mu_1 - \mu_2)$$