

# Supervised Dimension Reduction

David Barber

University College London

# Supervised Linear Projections

- In cases where class information is available, and our ultimate interest is to reduce dimensionality for improved classification, it makes sense to use the available class information in forming the projection.
- We consider data from two different classes. For class 1, we have a set of  $N_1$  datapoints

$$\mathcal{X}_1 = \{\mathbf{x}_1^1, \dots, \mathbf{x}_1^{N_1}\}$$

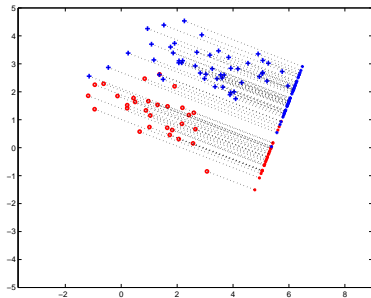
and similarly for class 2, we have a set of  $N_2$  datapoints

$$\mathcal{X}_2 = \{\mathbf{x}_2^1, \dots, \mathbf{x}_2^{N_2}\}$$

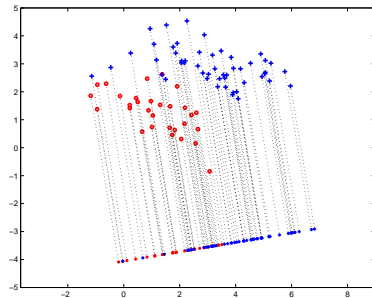
Our interest is then to find a linear projection,

$$\mathbf{y} = \mathbf{W}^\top \mathbf{x}$$

where  $\dim \mathbf{W} = D \times L$ ,  $L < D$ , such that for two datapoints  $\mathbf{x}^i$  and  $\mathbf{x}^j$  in the same class, the distance between their projections  $\mathbf{y}^i$  and  $\mathbf{y}^j$  should be small.



(a)



(b)

**Figure :** The large crosses represent data from class 1, and the large circles from class 2. Their projections onto 1 dimension are represented by their small counterparts. **(a):** Fisher's Linear Discriminant Analysis. Here there is little class overlap in the projections. **(b):** Unsupervised dimension reduction using Principal Components Analysis for comparison. There is considerable class overlap in the projection. In both (a) and (b) the one dimensional projection is the distance along the line, measured from an arbitrary chosen fixed point on the line.

# Fisher's Linear Discriminant

We model the data from each class with a Gaussian. That is

$$p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_1, \mathbf{S}_1), \quad p(\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_2 | \mathbf{m}_2, \mathbf{S}_2)$$

where  $\mathbf{m}_1$  is the sample mean of class 1 data, and  $\mathbf{S}_1$  the sample covariance; similarly for class 2. The projections of the points from the two classes are then given by

$$y_1^n = \mathbf{w}^T \mathbf{x}_1^n, \quad y_2^n = \mathbf{w}^T \mathbf{x}_2^n$$

Because the projections are linear, the projected distributions are also Gaussian,

$$\begin{aligned} p(y_1) &= \mathcal{N}(y_1 | \mu_1, \sigma_1^2), & \mu_1 &= \mathbf{w}^T \mathbf{m}_1, & \sigma_1^2 &= \mathbf{w}^T \mathbf{S}_1 \mathbf{w} \\ p(y_2) &= \mathcal{N}(y_2 | \mu_2, \sigma_2^2), & \mu_2 &= \mathbf{w}^T \mathbf{m}_2, & \sigma_2^2 &= \mathbf{w}^T \mathbf{S}_2 \mathbf{w} \end{aligned}$$

We search for a projection  $\mathbf{w}$  such that the projected distributions have minimal overlap. A useful objective function therefore is

$$\frac{(\mu_1 - \mu_2)^2}{\pi_1 \sigma_1^2 + \pi_2 \sigma_2^2}$$

where  $\pi_i$  represents the fraction of the dataset in class  $i$ .

In terms of the projection  $\mathbf{w}$ , the objective is

$$F(\mathbf{w}) = \frac{\mathbf{w}^\top (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{w}}{\mathbf{w}^\top (\pi_1 \mathbf{S}_1 + \pi_2 \mathbf{S}_2) \mathbf{w}} = \frac{\mathbf{w}^\top \mathbf{A} \mathbf{w}}{\mathbf{w}^\top \mathbf{B} \mathbf{w}}$$

where

$$\mathbf{A} = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^\top, \quad \mathbf{B} = \pi_1 \mathbf{S}_1 + \pi_2 \mathbf{S}_2$$

The optimal  $\mathbf{w}$  can be found by differentiating:

$$\frac{\partial}{\partial \mathbf{w}} \frac{\mathbf{w}^\top \mathbf{A} \mathbf{w}}{\mathbf{w}^\top \mathbf{B} \mathbf{w}} = \frac{2}{(\mathbf{w}^\top \mathbf{B} \mathbf{w})^2} [(\mathbf{w}^\top \mathbf{B} \mathbf{w}) \mathbf{A} \mathbf{w} - (\mathbf{w}^\top \mathbf{A} \mathbf{w}) \mathbf{B} \mathbf{w}]$$

and therefore the zero derivative requirement is

$$(\mathbf{w}^\top \mathbf{B} \mathbf{w}) \mathbf{A} \mathbf{w} = (\mathbf{w}^\top \mathbf{A} \mathbf{w}) \mathbf{B} \mathbf{w}$$

Multiplying by the inverse of  $\mathbf{B}$  we have

$$\mathbf{B}^{-1} (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{w} = \frac{\mathbf{w}^\top \mathbf{A} \mathbf{w}}{\mathbf{w}^\top \mathbf{B} \mathbf{w}} \mathbf{w}$$

Since  $(\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{w}$  is a scalar, the optimal projection is explicitly given by

$$\mathbf{w} \propto \mathbf{B}^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

## When the naive method breaks down

- The above derivation relied on the existence of the inverse of  $\mathbf{B}$ .
- A case where  $\mathbf{B}$  is not invertible is when there are fewer datapoints  $N_1 + N_2$  than dimensions  $D$ .
- A related problematic case is when there are elements of the input vectors that never vary. For example, in the hand-written digits case, the pixels at the corner edges are actually always zero. Let's call such a pixel  $z$ . The matrix  $\mathbf{B}$  will then have a zero entry for  $[B]_{z,z}$  (indeed the whole  $z^{th}$  row and column of  $\mathbf{B}$  will be zero) so that for any vector of the form

$$\mathbf{w}^T = (0, 0, \dots, w_z, 0, 0, \dots, 0) \Rightarrow \mathbf{w}^T \mathbf{B} \mathbf{w} = 0$$

This shows that the denominator of Fisher's objective can become zero, and the objective ill defined.

# Canonical Variates

Canonical Variates generalises Fisher's method to projections of more than one dimension and more than two classes. The projection of any point is given by

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

where  $\mathbf{W}$  is a  $D \times L$  matrix. Assuming that the data  $\mathbf{x}$  from class  $c$  is Gaussian distributed,

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{m}_c, \mathbf{S}_c)$$

the projections  $\mathbf{y}$  are also Gaussian

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{W}^T \mathbf{m}_c, \mathbf{W}^T \mathbf{S}_c \mathbf{W})$$

Find the mean  $\mathbf{m}$  of the whole dataset and  $\mathbf{m}_c$ , the mean of the each class  $c$ . Form

$$\mathbf{A} \equiv \sum_{c=1}^C N_c (\mathbf{m}_c - \mathbf{m}) (\mathbf{m}_c - \mathbf{m})^T$$

where  $N_c$  is the number of datapoints in class  $c$ .

Compute the covariance matrix  $\mathbf{S}_c$  of the data for each class  $c$ . Define

$$\mathbf{B} \equiv \sum_{c=1}^C N_c \mathbf{S}_c$$

Assuming  $\mathbf{B}$  is invertible we can define the Cholesky factor  $\tilde{\mathbf{B}}$ , with

$$\tilde{\mathbf{B}}^T \tilde{\mathbf{B}} = \mathbf{B}$$

A natural objective is then to maximise

$$F(\mathbf{W}) \equiv \frac{\text{trace} \left( \mathbf{W}^T \tilde{\mathbf{B}}^{-T} \mathbf{A} \tilde{\mathbf{B}}^{-1} \mathbf{W} \right)}{\text{trace} \left( \mathbf{W}^T \mathbf{W} \right)}$$

If we assume an orthonormality constraint on  $\mathbf{W}$ , then we equivalently require the maximisation of

$$F(\mathbf{W}) \equiv \text{trace} \left( \mathbf{W}^T \mathbf{C} \mathbf{W} \right), \text{ subject to } \mathbf{W}^T \mathbf{W} = \mathbf{I}$$

where

$$\mathbf{C} \equiv \frac{1}{D} \tilde{\mathbf{B}}^{-T} \mathbf{A} \tilde{\mathbf{B}}^{-1}$$



Since  $\mathbf{C}$  is symmetric and positive semidefinite, it has a real eigen-decomposition

$$\mathbf{C} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^T$$

where  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_D)$  is diagonal with non-negative entries containing the eigenvalues, sorted by decreasing order,  $\lambda_1 \geq \lambda_2 \geq \dots$  and  $\mathbf{E}^T\mathbf{E} = \mathbf{I}$ . Hence

$$F(\mathbf{W}) = \text{trace}(\mathbf{W}^T\mathbf{E}\mathbf{\Lambda}\mathbf{E}^T\mathbf{W})$$

- By setting  $\mathbf{W} = [\mathbf{e}_1, \dots, \mathbf{e}_L]$ , where  $\mathbf{e}_l$  is the  $l^{th}$  eigenvector, the objective  $F(\mathbf{W})$  becomes the sum of the first  $L$  eigenvalues.
- This setting maximises the objective function since forming  $\mathbf{W}$  from any other columns of  $\mathbf{E}$  would give a lower sum.
- Note that since  $\mathbf{A}$  has rank  $C$ , there can be no more than  $C - 1$  non-zero eigenvalues and corresponding directions.

# Canonical Variates Algorithm

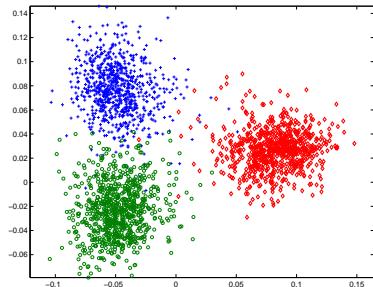
---

**Algorithm 1** Canonical Variates

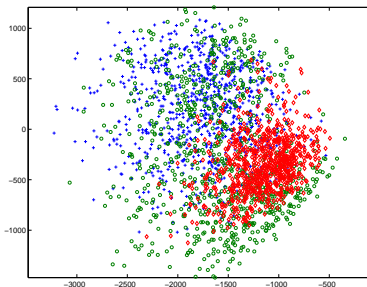
---

- 1: Compute the between and within class scatter matrices  $\mathbf{A}$ , and  $\mathbf{B}$ .
  - 2: Compute the Cholesky factor  $\tilde{\mathbf{B}}$  of  $\mathbf{B}$ .
  - 3: Compute the  $L$  principal eigenvectors  $[\mathbf{e}_1, \dots, \mathbf{e}_L]$  of  $\tilde{\mathbf{B}}^{-\top} \mathbf{A} \tilde{\mathbf{B}}^{-1}$ .
  - 4: Return  $\mathbf{W} = [\mathbf{e}_1, \dots, \mathbf{e}_L]$  as the projection matrix.
-

# Using canonical variates on the digit data



(a)



(b)

**Figure :** **(a):** Canonical Variates projection of examples of handwritten digits 3('+'), 5('o') and 7(diamond). There are 800 examples from each digit class. Plotted are the projections down to 2 dimensions. **(b):** PCA projections for comparison.

# Dealing with the nullspace

- One may encounter situations where  $\mathbf{B}$  is not invertible.
- A solution is to require that  $\mathbf{W}$  lies only in the subspace spanned by the data (that is there can be no contribution from the nullspace).
- To do this we first concatenate the training data from all classes into one large matrix  $\mathbf{X}$ .
- A basis for  $\mathbf{X}$  can be found using, for example, the thin-SVD technique which returns an orthonormal non-square basis matrix  $\mathbf{Q}$ .
- We then require the solution  $\mathbf{W}$  to be expressed in this basis

$$\mathbf{W} = \mathbf{Q}\mathbf{W}'$$

for some matrix  $\mathbf{W}'$ .

Substituting this in the Canonical Variates objective we obtain

$$F(\mathbf{W}') \equiv \frac{\text{trace}(\mathbf{W}'^T \mathbf{Q}^T \mathbf{A} \mathbf{Q} \mathbf{W}')}{\text{trace}(\mathbf{W}'^T \mathbf{Q}^T \mathbf{B} \mathbf{Q} \mathbf{W}')}$$

This is of the same form as the standard quotient on replacing the between-scatter  $\mathbf{A}$  with

$$\mathbf{A}' \equiv \mathbf{Q}^T \mathbf{A} \mathbf{Q}$$

and the within-scatter  $\mathbf{B}$  with

$$\mathbf{B}' \equiv \mathbf{Q}^T \mathbf{B} \mathbf{Q}$$

In this case  $\mathbf{B}'$  is guaranteed invertible since  $\mathbf{B}$  is projected down to the basis that spans the data. One may then carry out Canonical Variates, as above, which returns the matrix  $\mathbf{W}'$ . Transforming, back,  $\mathbf{W}$  is then given by  $\mathbf{W} = \mathbf{Q} \mathbf{W}'$ .