

Data Analysis and Manifold Learning

Lecture 5: Minimum-error Formulation of PCA and Fisher's Discriminant Analysis

Radu Horaud
INRIA Grenoble Rhone-Alpes, France
Radu.Horaud@inrialpes.fr
<http://perception.inrialpes.fr/>

Outline of Lecture 5

- Minimum-error formulation of PCA
- PCA for high-dimensional spaces
- Fischer's discriminant analysis for two classes and generalization to K classes

Material for This Lecture

- C. M. Bishop. Pattern Recognition and Machine Learning. 2006. (Chapters 4 and 12)
- http://en.wikipedia.org/wiki/Linear_discriminant_analysis
- Numerous textbooks treat PCA and LDA

Projecting the Data

- Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_n) \subset \mathbb{R}^D$,
- Consider an orthonormal basis vector, e.g., the columns of a $D \times D$ orthonormal matrix \mathbf{U} , namely $\mathbf{u}_i^\top \mathbf{u}_j = \delta_{ij}$.
- We can write:

$$\mathbf{x}_j = \sum_{i=1}^D \alpha_{ji} \mathbf{u}_i \text{ with } \alpha_{ji} = \mathbf{x}_j^\top \mathbf{u}_i$$

- Moreover, consider a lower-dimensional subspace of dimension $d < D$. We approximate each data point with:

$$\tilde{\mathbf{x}}_j = \sum_{i=1}^d z_{ji} \mathbf{u}_i + \sum_{i=d+1}^D b_i \mathbf{u}_i$$

Minimizing the Distorsion

- Choose the vectors $\{\mathbf{u}_j\}$ and the scalars $\{z_{ji}\}$ and $\{b_i\}$ that minimize the following distortion error:

$$J = \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j - \tilde{\mathbf{x}}_j\|^2$$

- By substitution of $\tilde{\mathbf{x}}_j$ and by setting the derivatives to 0, $\partial J / \partial z_{ji} = 0$, $\partial J / \partial b_i = 0$ we obtain:

$$z_{ji} = \mathbf{x}_j^\top \mathbf{u}_i, \quad i = 1 \dots d$$

$$b_i = \bar{\mathbf{x}}^\top \mathbf{u}_i, \quad i = d + 1 \dots D$$

Closed-form Expression of the Distorsion

- By substitution we obtain the following distorsion between each data point and its projection onto the *principal subspace*:

$$\mathbf{x}_j - \tilde{\mathbf{x}}_j = \sum_{i=d+1}^D \left((\mathbf{x}_j - \bar{\mathbf{x}})^\top \mathbf{u}_i \right) \mathbf{u}_i$$

This *error-vector* lies in a space perpendicular to the principal space. The distorsion becomes:

$$J = \frac{1}{n} \sum_{j=1}^n \sum_{i=d+1}^D (\mathbf{x}_j^\top \mathbf{u}_i - \bar{\mathbf{x}}^\top \mathbf{u}_i)^2 = \sum_{i=d+1}^D \mathbf{u}_i^\top \Sigma \mathbf{u}_i$$

Minimizing the Distorsion (I)

- Note that in the previous equation,

$$1/n \sum_{j=1}^n (\mathbf{x}_j^\top \mathbf{u}_i - \bar{\mathbf{x}}^\top \mathbf{u}_i)^2$$

corresponds to the variance of the projected data onto \mathbf{u}_i .

Minimizing the distorsion is equivalent to minimizing the variances along the directions perpendicular to the principal directions $\mathbf{u}_1 \dots \mathbf{u}_d$. This can be done by minimizing \tilde{J} with respect to $\mathbf{u}_{d+1} \dots \mathbf{u}_D$:

$$\tilde{J} = \sum_{i=d+1}^D \mathbf{u}_i^\top \Sigma \mathbf{u}_i + \sum_{i=d+1}^D \lambda_i (1 - \mathbf{u}_i^\top \mathbf{u}_i)$$

Minimizing the Distorsion (II)

- By setting the derivatives to 0 we obtain:

$$\frac{\partial \tilde{J}}{\partial \mathbf{u}_i} = 0 \leftrightarrow \Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad i = d + 1 \dots D$$

- The distorsion becomes:

$$\tilde{J} = \lambda_{d+1} + \dots + \lambda_D$$

- The principal directions correspond to the largest d eigenvalue-eigenvector pairs of the covariance matrix Σ

Choosing the Dimension of the Principal Subspace

- The covariance matrix can be written as $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$. The trace of the diagonal matrix $\mathbf{\Lambda}$ can be interpreted as the *total variance*.
- One way to choose the principal subspace is to choose the largest d eigenvalue-eigenvector pairs such that:

$$\alpha(d) = \frac{\lambda_1 + \dots + \lambda_d}{\lambda_1 + \dots + \lambda_D} = \frac{\lambda_1 + \dots + \lambda_d}{\text{tr}(\Sigma)} \approx 0.95$$

High-dimensional Data

- When D is very large, the number of data points n may be smaller than the dimension. In this case it is better to use the $n \times n$ Gram matrix instead of the $D \times D$ covariance matrix.
- For centred data we have:

$$\Sigma = \frac{1}{n} \mathbf{X} \mathbf{X}^\top \text{ with } (\lambda_i, \mathbf{u}_i)$$

$$\mathbf{G} = \mathbf{X}^\top \mathbf{X} \text{ with } (\mu_i, \mathbf{v}_i)$$

- By premultiplication of $\mathbf{X} \mathbf{X}^\top \mathbf{u} = n\lambda \mathbf{u}$ with \mathbf{X}^\top we obtain:

$$\mathbf{v} = \mathbf{X}^\top \mathbf{u} \text{ and } \mu = n\lambda$$

- From which we obtain: $\mathbf{u} = \frac{1}{\mu} \mathbf{X} \mathbf{v}$
- Assuming that the eigenvectors of the Gram matrix are normalized, we obtain:

$$\frac{\mathbf{u}}{\|\mathbf{u}\|} = \frac{1}{\sqrt{\mu}} \mathbf{X} \mathbf{v}$$

Discriminant Analysis

- Project the high-dimensional input vector to one dimension, i.e., along the direction of \mathbf{w} :

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

- This results in a loss of information and well-separated clusters in the initial space may overlap in one dimension.
- With a proper choice of \mathbf{w} one can select a projection that maximizes the class separation.

Two-Class Problem

- Let's assume that the data points belong to two clusters, \mathcal{C}_1 and \mathcal{C}_2 and that the mean vectors of these two clusters are $\bar{\mathbf{x}}_1 = 1/n_1 \sum_{j \in \mathcal{C}_1} \mathbf{x}_j$ and $\bar{\mathbf{x}}_2 = 1/n_2 \sum_{j \in \mathcal{C}_2} \mathbf{x}_j$
- One can choose \mathbf{w} to maximize the distance between the projected means: $\bar{y}_1 - \bar{y}_2 = \mathbf{w}^\top (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$
- We can enforce the constraint $\mathbf{w}^\top \mathbf{w} = 1$ using a Lagrange multiplier and obtain the following solution:

$$\mathbf{w} = \bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2$$

- This solution is optimal when the two clusters are spherical.

Fisher's Linear Discriminant

- The solution consists in enforce small variances within each class. The criterion to be maximized becomes:

$$J(\mathbf{w}) = \frac{(\bar{x}_1 - \bar{x}_2)^2}{\sigma_1^2 + \sigma_2^2}$$

- Where the within cluster projected variance is:

$$\sigma_k^2 = \frac{1}{n_k} \sum_{j \in \mathcal{C}_k} (y_j - \bar{y}_k)^2$$

Maximizing Fisher's Criterion

- The criterion can be rewritten in the form:

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top \Sigma_B \mathbf{w}}{\mathbf{w}^\top \Sigma_W \mathbf{w}}$$

- Where Σ_B is the *between-cluster* covariance and Σ_W the total *within-cluster* covariance. They are given by:

$$\begin{aligned}\Sigma_B &= (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \\ \Sigma_W &= \frac{1}{n_1} \sum_{j \in \mathcal{C}_1} (\mathbf{x}_j - \bar{\mathbf{x}}_1)(\mathbf{x}_j - \bar{\mathbf{x}}_1)^\top \\ &\quad + \frac{1}{n_2} \sum_{j \in \mathcal{C}_2} (\mathbf{x}_j - \bar{\mathbf{x}}_2)(\mathbf{x}_j - \bar{\mathbf{x}}_2)^\top\end{aligned}$$

- Optimal solution: $\mathbf{w} = \Sigma_W^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$

Discriminant Analysis

- The method just described can be applied to a training data set, where each data point belongs to one of the two clusters.
- Once the choice of an optimal direction of projection was performed, the projected data can be used to construct a discriminant.
- The projected training data belongs to a 1D Gaussian mixture with two clusters, \mathcal{C}_1 and \mathcal{C}_2 . The parameters of eachone of these two clusters can be computed with ($k = 1, 2$):

$$\pi_k = n_k/n, \quad \bar{y}_k = \mathbf{w}^\top \bar{\mathbf{x}}_k, \quad \text{and} \quad \sigma_k^2 = \frac{1}{n_k} \sum_{j \in \mathcal{C}_k} (y_j - \bar{y}_k)^2$$

- Classification of a *new data point* y can be done using the class-posterior probabilities:

$$\mathcal{C} = \arg \max_k \pi_k \mathcal{N}(y | \bar{y}_k, \sigma_k^2)$$

Fisher's Discriminant for Multiple Classes

- The two-class discriminant analysis can be extended to $K > 2$ classes.
- The idea is to consider several linear projections, i.e., $y_k = \mathbf{w}_k^\top \mathbf{x}$ with $k = 1 \dots K - 1$.
- A formal derivation can be found in: C. M. Bishop. Pattern Recognition and Machine Learning. 2006. (Chapter 4, pp 191-192).