

Data Mining and Machine Learning: Fundamental Concepts and Algorithms

dataminingbook.info

Mohammed J. Zaki¹ Wagner Meira Jr.²

¹Department of Computer Science
Rensselaer Polytechnic Institute, Troy, NY, USA

²Department of Computer Science
Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

Chap. 20: Linear Discriminant Analysis

Linear Discriminant Analysis

Given labeled data consisting of d -dimensional points \mathbf{x}_i along with their classes y_i , the goal of linear discriminant analysis (LDA) is to find a vector \mathbf{w} that maximizes the separation between the classes after projection onto \mathbf{w} .

The key difference between principal component analysis and LDA is that the former deals with unlabeled data and tries to maximize variance, whereas the latter deals with labeled data and tries to maximize the discrimination between the classes.

Projection onto a Line

Let \mathbf{D}_i denote the subset of points labeled with class c_i , i.e., $\mathbf{D}_i = \{\mathbf{x}_j | y_j = c_i\}$, and let $|\mathbf{D}_i| = n_i$ denote the number of points with class c_i . We assume that there are only $k = 2$ classes.

The projection of any d -dimensional point \mathbf{x}_i onto a unit vector \mathbf{w} is given as

$$\mathbf{x}'_i = \left(\frac{\mathbf{w}^T \mathbf{x}_i}{\mathbf{w}^T \mathbf{w}} \right) \mathbf{w} = (\mathbf{w}^T \mathbf{x}_i) \mathbf{w} = a_i \mathbf{w}$$

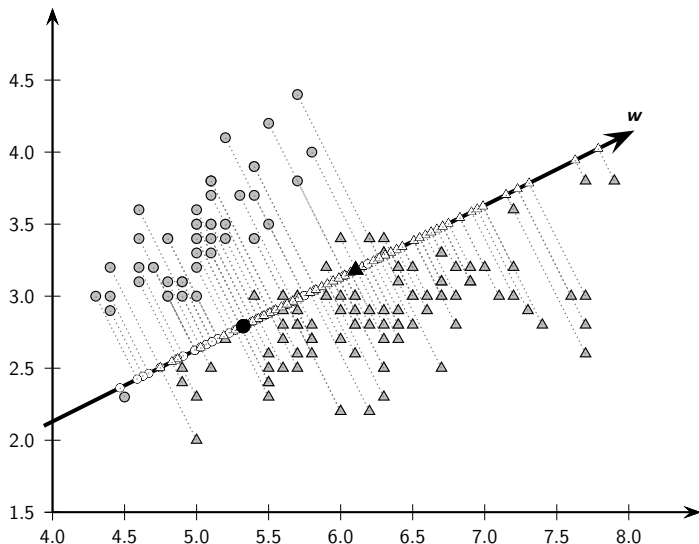
where a_i specifies the offset or coordinate of \mathbf{x}'_i along the line \mathbf{w} :

$$a_i = \mathbf{w}^T \mathbf{x}_i$$

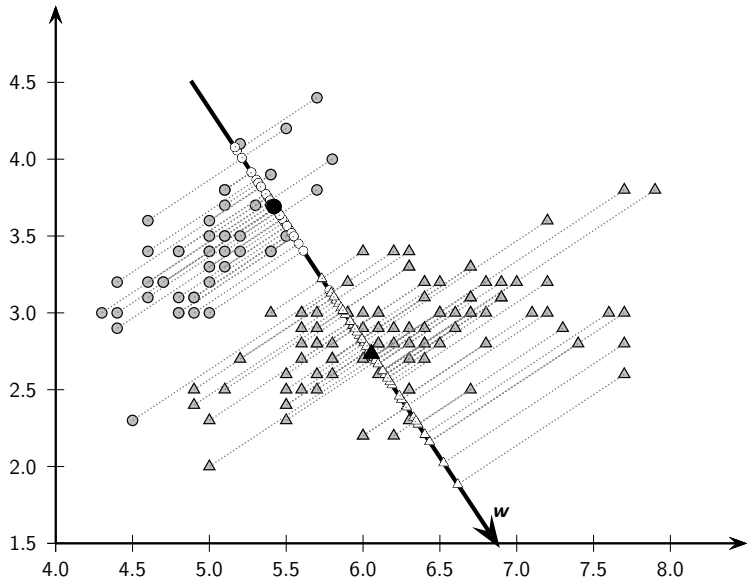
The set of n scalars $\{a_1, a_2, \dots, a_n\}$ represents the mapping from \mathbb{R}^d to \mathbb{R} , that is, from the original d -dimensional space to a 1-dimensional space (along \mathbf{w}).

Projection onto w : Iris 2D Data

iris-setosa as class c_1 (circles), and the other two Iris types as class c_2 (triangles)



Iris 2D Data: Optimal Linear Discriminant Direction



The mean of the projected points is given as:

$$m_1 = \mathbf{w}^T \boldsymbol{\mu}_1 \qquad m_2 = \mathbf{w}^T \boldsymbol{\mu}_2$$

To maximize the separation between the classes, we maximize the difference between the projected means, $|m_1 - m_2|$. However, for good separation, the variance of the projected points for each class should also not be too large. LDA maximizes the separation by ensuring that the *scatter* s_i^2 for the projected points within each class is small, where scatter is defined as

$$s_i^2 = \sum_{\mathbf{x}_j \in \mathcal{D}_i} (a_j - m_i)^2 = n_i \sigma_i^2$$

where σ_i^2 is the variance for class c_i .

Linear Discriminant Analysis: Fisher Objective

We incorporate the two LDA criteria, namely, maximizing the distance between projected means and minimizing the sum of projected scatter, into a single maximization criterion called the *Fisher LDA objective*:

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

In matrix terms, we can rewrite $(m_1 - m_2)^2$ as follows:

$$(m_1 - m_2)^2 = \left(\mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right)^2 = \mathbf{w}^T \mathbf{B} \mathbf{w}$$

where $\mathbf{B} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$ is a $d \times d$ rank-one matrix called the *between-class scatter matrix*.

The projected scatter for class c_i is given as

$$s_i^2 = \sum_{\mathbf{x}_j \in \mathcal{D}_i} (\mathbf{w}^T \mathbf{x}_j - \mathbf{w}^T \boldsymbol{\mu}_i)^2 = \mathbf{w}^T \left(\sum_{\mathbf{x}_j \in \mathcal{D}_i} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T \right) \mathbf{w} = \mathbf{w}^T \mathbf{S}_i \mathbf{w}$$

where \mathbf{S}_i is the *scatter matrix* for \mathcal{D}_i .

Linear Discriminant Analysis: Fisher Objective

The combined scatter for both classes is given as

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_1 \mathbf{w} + \mathbf{w}^T \mathbf{S}_2 \mathbf{w} = \mathbf{w}^T (\mathbf{S}_1 + \mathbf{S}_2) \mathbf{w} = \mathbf{w}^T \mathbf{S} \mathbf{w}$$

where the symmetric positive semidefinite matrix $\mathbf{S} = \mathbf{S}_1 + \mathbf{S}_2$ denotes the *within-class scatter matrix* for the pooled data.

The LDA objective function in matrix form is

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{B} \mathbf{w}}{\mathbf{w}^T \mathbf{S} \mathbf{w}}$$

To solve for the best direction \mathbf{w} , we differentiate the objective function with respect to \mathbf{w} ; after simplification it yields the *generalized eigenvalue problem*

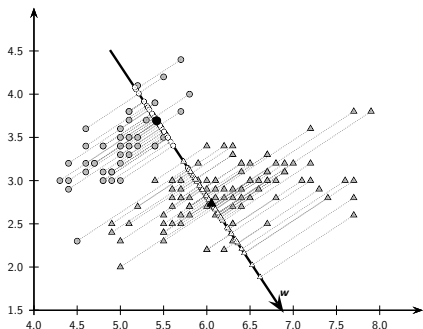
$$\mathbf{B} \mathbf{w} = \lambda \mathbf{S} \mathbf{w}$$

where $\lambda = J(\mathbf{w})$ is a generalized eigenvalue of \mathbf{B} and \mathbf{S} . To maximize the objective λ should be chosen to be the largest generalized eigenvalue, and \mathbf{w} to be the corresponding eigenvector.

LinearDiscriminant ($D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$):

- 1 $D_i \leftarrow \{\mathbf{x}_j \mid y_j = c_i, j = 1, \dots, n\}, i = 1, 2$ // class-specific subsets
- 2 $\mu_i \leftarrow \text{mean}(D_i), i = 1, 2$ // class means
- 3 $B \leftarrow (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$ // between-class scatter matrix
- 4 $Z_i \leftarrow D_i - 1_{n_i} \mu_i^T, i = 1, 2$ // center class matrices
- 5 $S_i \leftarrow Z_i^T Z_i, i = 1, 2$ // class scatter matrices
- 6 $S \leftarrow S_1 + S_2$ // within-class scatter matrix
- 7 $\lambda_1, \mathbf{w} \leftarrow \text{eigen}(S^{-1}B)$ // compute dominant eigenvector

Linear Discriminant Direction: Iris 2D Data



The between-class scatter matrix is

$$B = \begin{pmatrix} 1.587 & -0.693 \\ -0.693 & 0.303 \end{pmatrix}$$

and the within-class scatter matrix is

$$S_1 = \begin{pmatrix} 6.09 & 4.91 \\ 4.91 & 7.11 \end{pmatrix}$$

$$S_2 = \begin{pmatrix} 43.5 & 12.09 \\ 12.09 & 10.96 \end{pmatrix}$$

$$S = \begin{pmatrix} 49.58 & 17.01 \\ 17.01 & 18.08 \end{pmatrix}$$

The direction of most separation between c_1 and c_2 is the dominant eigenvector corresponding to the largest eigenvalue of the matrix $S^{-1}B$. The solution is

$$J(w) = \lambda_1 = 0.11$$

$$w = \begin{pmatrix} 0.551 \\ -0.834 \end{pmatrix}$$

Linear Discriminant Analysis: Two Classes

For the two class scenario, if \mathbf{S} is nonsingular, we can directly solve for \mathbf{w} without computing the eigenvalues and eigenvectors.

The between-class scatter matrix \mathbf{B} points in the same direction as $(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ because

$$\begin{aligned}\mathbf{B}\mathbf{w} &= \left((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \right) \mathbf{w} \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \left((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w} \right) \\ &= b(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\end{aligned}$$

The generalized eigenvectors equation can then be rewritten as

$$\mathbf{w} = \frac{b}{\lambda} \mathbf{S}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

Because $\frac{b}{\lambda}$ is just a scalar, we can solve for the best linear discriminant as

$$\mathbf{w} = \mathbf{S}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

We can finally normalize \mathbf{w} to be a unit vector.

We can directly compute \mathbf{w} as follows:

$$\begin{aligned}\mathbf{w} &= \mathbf{S}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= \begin{pmatrix} 0.066 & -0.029 \\ -0.100 & 0.044 \end{pmatrix} \begin{pmatrix} -1.246 \\ 0.546 \end{pmatrix} = \begin{pmatrix} -0.0527 \\ 0.0798 \end{pmatrix}\end{aligned}$$

After normalizing, we have

$$\mathbf{w} = \frac{\mathbf{w}}{\|\mathbf{w}\|} = \frac{1}{0.0956} \begin{pmatrix} -0.0527 \\ 0.0798 \end{pmatrix} = \begin{pmatrix} -0.551 \\ 0.834 \end{pmatrix}$$

Note that even though the sign is reversed for \mathbf{w} , they represent the same direction; only the scalar multiplier is different.

Kernel Discriminant Analysis

The goal of kernel LDA is to find the direction vector \mathbf{w} in feature space that maximizes

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

It is well known that \mathbf{w} can be expressed as a linear combination of the points in feature space

$$\mathbf{w} = \sum_{j=1}^n a_j \phi(\mathbf{x}_j)$$

The mean for class c_i in feature space is given as

$$\boldsymbol{\mu}_i^\phi = \frac{1}{n_i} \sum_{\mathbf{x}_j \in \mathcal{D}_i} \phi(\mathbf{x}_j)$$

and the covariance matrix for class c_i in feature space is

$$\boldsymbol{\Sigma}_i^\phi = \frac{1}{n_i} \sum_{\mathbf{x}_j \in \mathcal{D}_i} \left(\phi(\mathbf{x}_j) - \boldsymbol{\mu}_i^\phi \right) \left(\phi(\mathbf{x}_j) - \boldsymbol{\mu}_i^\phi \right)^T$$

Kernel Discriminant Analysis

The between-class scatter matrix in feature space is

$$\mathbf{B}_\phi = \left(\boldsymbol{\mu}_1^\phi - \boldsymbol{\mu}_2^\phi \right) \left(\boldsymbol{\mu}_1^\phi - \boldsymbol{\mu}_2^\phi \right)^T$$

and the within-class scatter matrix in feature space is

$$\mathbf{S}_\phi = n_1 \boldsymbol{\Sigma}_1^\phi + n_2 \boldsymbol{\Sigma}_2^\phi$$

\mathbf{S}_ϕ is a $d \times d$ symmetric, positive semidefinite matrix, where d is the dimensionality of the feature space.

The best linear discriminant vector \mathbf{w} in feature space is the dominant eigenvector, which satisfies the expression

$$\left(\mathbf{S}_\phi^{-1} \mathbf{B}_\phi \right) \mathbf{w} = \lambda \mathbf{w}$$

where we assume that \mathbf{S}_ϕ is non-singular.

LDA Objective via Kernel Matrix: Between-class Scatter

The projected mean for class c_i is given as

$$\mathbf{m}_i = \mathbf{w}^T \boldsymbol{\mu}_i^\phi = \frac{1}{n_i} \sum_{j=1}^n \sum_{\mathbf{x}_k \in \mathcal{D}_i} a_j K(\mathbf{x}_j, \mathbf{x}_k) = \mathbf{a}^T \mathbf{m}_i$$

where $\mathbf{a} = (a_1, a_2, \dots, a_n)^T$ is the weight vector, and

$$\mathbf{m}_i = \frac{1}{n_i} \begin{pmatrix} \sum_{\mathbf{x}_k \in \mathcal{D}_i} K(\mathbf{x}_1, \mathbf{x}_k) \\ \sum_{\mathbf{x}_k \in \mathcal{D}_i} K(\mathbf{x}_2, \mathbf{x}_k) \\ \vdots \\ \sum_{\mathbf{x}_k \in \mathcal{D}_i} K(\mathbf{x}_n, \mathbf{x}_k) \end{pmatrix} = \frac{1}{n_i} \mathbf{K}^{c_i} \mathbf{1}_{n_i}$$

where \mathbf{K}^{c_i} is the $n \times n_i$ subset of the kernel matrix, restricted to columns belonging to points only in \mathcal{D}_i , and $\mathbf{1}_{n_i}$ is the n_i -dimensional vector all of whose entries are one.

The separation between the projected means is thus

$$(\mathbf{m}_1 - \mathbf{m}_2)^2 = \left(\mathbf{a}^T \mathbf{m}_1 - \mathbf{a}^T \mathbf{m}_2 \right)^2 = \mathbf{a}^T \mathbf{M} \mathbf{a}$$

where $\mathbf{M} = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$ is the between-class scatter matrix.

LDA Objective via Kernel Matrix: Within-class Scatter

We can compute the projected scatter for each class, s_1^2 and s_2^2 , purely in terms of the kernel function, as follows

$$s_1^2 = \sum_{\mathbf{x}_j \in \mathcal{D}_1} \left\| \mathbf{w}^T \phi(\mathbf{x}_j) - \mathbf{w}^T \boldsymbol{\mu}_1^\phi \right\|^2 = \mathbf{a}^T \left(\left(\sum_{\mathbf{x}_j \in \mathcal{D}_1} \mathbf{K}_j \mathbf{K}_j^T \right) - n_1 \mathbf{m}_1 \mathbf{m}_1^T \right) \mathbf{a} = \mathbf{a}^T \mathbf{N}_1 \mathbf{a}$$

where \mathbf{K}_j is the i th column of the kernel matrix, and \mathbf{N}_1 is the class scatter matrix for c_1 .

The sum of projected scatter values is then given as

$$s_1^2 + s_2^2 = \mathbf{a}^T (\mathbf{N}_1 + \mathbf{N}_2) \mathbf{a} = \mathbf{a}^T \mathbf{N} \mathbf{a}$$

where \mathbf{N} is the $n \times n$ within-class scatter matrix.

The kernel LDA maximization condition is

$$\max_{\mathbf{w}} J(\mathbf{w}) = \max_{\mathbf{a}} J(\mathbf{a}) = \frac{\mathbf{a}^T \mathbf{M} \mathbf{a}}{\mathbf{a}^T \mathbf{N} \mathbf{a}}$$

The weight vector \mathbf{a} is the eigenvector corresponding to the largest eigenvalue of the generalized eigenvalue problem:

$$\mathbf{M} \mathbf{a} = \lambda_1 \mathbf{N} \mathbf{a}$$

When there are only two classes \mathbf{a} can be obtained directly:

$$\mathbf{a} = \mathbf{N}^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

To normalize \mathbf{w} to be a unit vector we scale \mathbf{a} by $\frac{1}{\sqrt{\mathbf{a}^T \mathbf{K} \mathbf{a}}}$.

We can project any point \mathbf{x} onto the discriminant direction as follows:

$$\mathbf{w}^T \phi(\mathbf{x}) = \sum_{j=1}^n a_j \phi(\mathbf{x}_j)^T \phi(\mathbf{x}) = \sum_{j=1}^n a_j K(\mathbf{x}_j, \mathbf{x})$$

Kernel Discriminant Analysis Algorithm

KernelDiscriminant ($D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n, K$):

- 1 $\mathbf{K} \leftarrow \{K(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1,\dots,n}$ // compute $n \times n$ kernel matrix
- 2 $\mathbf{K}^{c_i} \leftarrow \{K(j, k) \mid y_k = c_i, 1 \leq j, k \leq n\}, i = 1, 2$ // class kernel matrix
- 3 $\mathbf{m}_i \leftarrow \frac{1}{n_i} \mathbf{K}^{c_i} \mathbf{1}_{n_i}, i = 1, 2$ // class means
- 4 $\mathbf{M} \leftarrow (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$ // between-class scatter matrix
- 5 $\mathbf{N}_i \leftarrow \mathbf{K}^{c_i} (\mathbf{I}_{n_i} - \frac{1}{n_i} \mathbf{1}_{n_i \times n_i}) (\mathbf{K}^{c_i})^T, i = 1, 2$ // class scatter matrices
- 6 $\mathbf{N} \leftarrow \mathbf{N}_1 + \mathbf{N}_2$ // within-class scatter matrix
- 7 $\lambda_1, \mathbf{a} \leftarrow \text{eigen}(\mathbf{N}^{-1} \mathbf{M})$ // compute weight vector
- 8 $\mathbf{a} \leftarrow \frac{\mathbf{a}}{\sqrt{\mathbf{a}^T \mathbf{K} \mathbf{a}}}$ // normalize \mathbf{w} to be unit vector

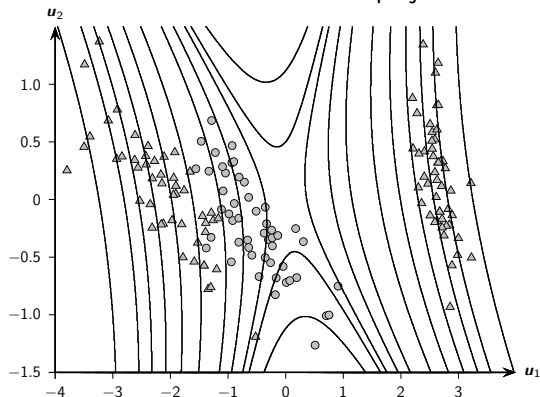
Kernel Discriminant Analysis

Quadratic Homogeneous Kernel

Iris 2D Data: c_1 (circles; iris-virginica) and c_2 (triangles; other two types).

Kernel Function: $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^2$.

Contours of constant projection onto optimal kernel discriminant, i.e., points along both the curves have the same value when projected onto \mathbf{w} .



Kernel Discriminant Analysis

Quadratic Homogeneous Kernel

Projecting $\mathbf{x}_i \in \mathbf{D}$ onto \mathbf{w} , which separates the two classes.

The projected scatters and means for both classes are as follows:

$$m_1 = 0.338$$

$$m_2 = 4.476$$

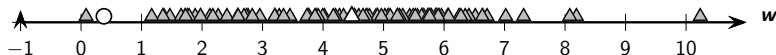
$$s_1^2 = 13.862$$

$$s_2^2 = 320.934$$

The value of $J(\mathbf{w})$ is given as

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} = \frac{(0.338 - 4.476)^2}{13.862 + 320.934} = \frac{17.123}{334.796} = 0.0511$$

which, as expected, matches $\lambda_1 = 0.0511$ from above.



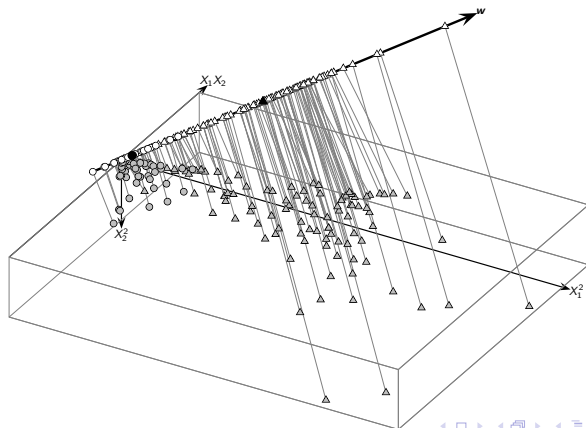
Kernel Feature Space and Optimal Discriminant

It is not desirable or possible to obtain an explicit discriminant vector \mathbf{w} .

$\mathbf{x} = (x_1, x_2)^T \in \mathbb{R}^2$ is mapped to $\phi(\mathbf{x}) = (\sqrt{2}x_1x_2, x_1^2, x_2^2)^T \in \mathbb{R}^3$.

The projection of $\phi(\mathbf{x}_i)$ onto \mathbf{w} is also shown, where

$$\mathbf{w} = 0.511x_1x_2 + 0.761x_1^2 - 0.4x_2^2$$



Data Mining and Machine Learning: Fundamental Concepts and Algorithms

dataminingbook.info

Mohammed J. Zaki¹ Wagner Meira Jr.²

¹Department of Computer Science
Rensselaer Polytechnic Institute, Troy, NY, USA

²Department of Computer Science
Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

Chap. 20: Linear Discriminant Analysis