

# Discriminant Analysis

1. Fisher Linear Discriminant
2. Multiple Discriminant Analysis

# Motivation

Projection that best **separates** the data in a least-squares sense

- PCA finds components that are useful for representing data
- However no reason to assume that components are useful for discriminating between data in different classes
- Pooling data may discard directions that are essential
  - OCR example: Discriminating between O and Q may ignore tail

# Fisher Linear Discriminant

Projecting data from  $d$  dimensions onto a line

Set of  $n$   $d$  - dimensional samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$

$n_1$  in the subset  $D_1$  labelled  $\omega_1$

$n_2$  in the subset  $D_2$  labelled  $\omega_2$

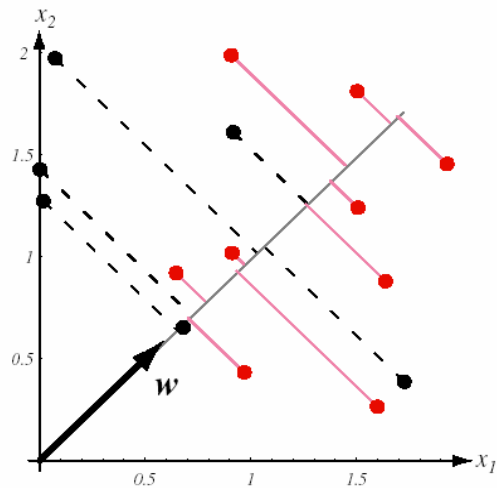
We wish to form a linear combination of the components of  $\mathbf{x}$  as

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

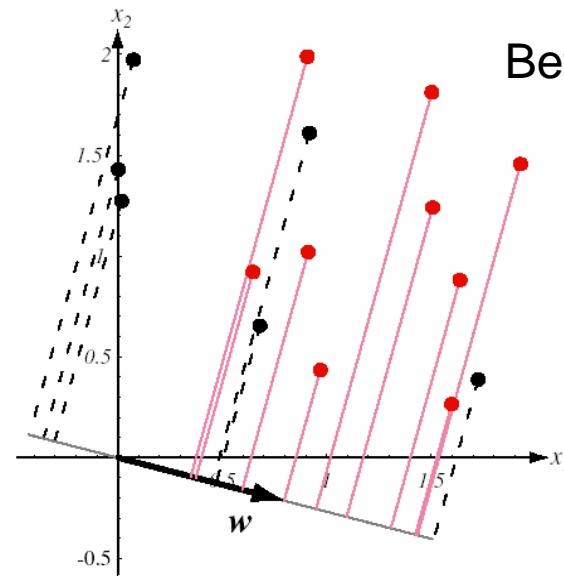
and a corresponding set  $n$  of samples  $y_1, \dots, y_n$

# Fisher Linear Discriminant: two-dimensional example

Classes mixed



Better Separation



Projection of same set of two-class samples onto two different lines in the direction marked  $w$ .

# Finding best direction $w$

Sample Mean in  $d$  - dimensional space :

$$\mathbf{m}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in D_i} \mathbf{x}$$

Sample Mean of projected points :

$$\hat{m}_i = \frac{1}{n_i} \sum_{y \in Y_i} y = \frac{1}{n_i} \sum_{\mathbf{x} \in D_i} \mathbf{w}^t \mathbf{x} = \mathbf{w}^t \mathbf{m}_i$$

Distance between projected means is

$$|\hat{m}_1 - \hat{m}_2| = |\mathbf{w}^t (\mathbf{m}_1 - \mathbf{m}_2)|$$

# Criterion for Fisher Linear Discriminant

Rather than forming sample variances, define scatter for the projected samples

$$s_i^{-2} = \sum_{y \in Y_i} (y - \bar{m}_i)^2$$

Thus  $(1/n)(s_1^{-2} + s_2^{-2})$  is an estimate of the variance of the pooled data

Total within class scatter is  $(s_1^{-2} + s_2^{-2})$

Find that linear function  $w^t x$  for which  $J(w) = \frac{|\bar{m}_1 - \bar{m}_2|^2}{(s_1^{-2} + s_2^{-2})}$

is maximum and independent of  $\|w\|$ .

While maximizing  $J(\cdot)$  leads to best separation between the two projected sets, we will need a threshold criterion to have a true classifier.

# Scatter Matrices

To obtain  $J(\cdot)$  as an explicit function of  $w$ , we define scatter matrices  $S_i$  and  $S_W$

$$S_i = \sum_{x \in \mathcal{D}_i} (x - m_i)(x - m_i)^t$$

and

$$S_W = S_1 + S_2. \quad (98)$$

Then we can write

$$\begin{aligned} \tilde{s}_i^2 &= \sum_{x \in \mathcal{D}_i} (w^t x - w^t m_i)^2 \\ &= \sum_{x \in \mathcal{D}_i} w^t (x - m_i)(x - m_i)^t w \\ &= w^t S_i w; \end{aligned} \quad (99)$$

therefore the sum of these scatters can be written

$$\tilde{s}_1^2 + \tilde{s}_2^2 = w^t S_W w. \quad (100)$$

Similarly, the separations of the projected means obeys

$$\begin{aligned} (\tilde{m}_1 - \tilde{m}_2)^2 &= (w^t m_1 - w^t m_2)^2 \\ &= w^t (m_1 - m_2)(m_1 - m_2)^t w \\ &= w^t S_B w, \end{aligned}$$

where

$$S_B = (m_1 - m_2)(m_1 - m_2)^t. \quad (102)$$

Within Class  
Scatter Matrix

Between Class  
Scatter Matrix

# Scatter Matrices

To obtain  $J(\cdot)$  as an explicit function of  $w$ , we define scatter matrices  $S_i$  and  $S_W$

$$S_i = \sum_{x \in D_i} (x - m_i)(x - m_i)^t$$

Within Class  
Scatter Matrix

$$S_W = S_1 + S_2$$

Between Class  
Scatter Matrix

$$S_B = (m_1 - m_2)(m_1 - m_2)^t$$



# Criterion Function in terms of Scatter Matrices

We call  $\mathbf{S}_W$  the *within-class scatter matrix*. It is proportional to the sample covariance matrix for the pooled  $d$ -dimensional data. It is symmetric and positive semidefinite, and it is usually nonsingular if  $n > d$ . Likewise,  $\mathbf{S}_B$  is called the *between-class scatter matrix*. It is also symmetric and positive semidefinite, but because it is the outer product of two vectors, its rank is at most one. In particular, for any  $\mathbf{w}$ ,  $\mathbf{S}_B \mathbf{w}$  is in the direction of  $\mathbf{m}_1 - \mathbf{m}_2$ , and  $\mathbf{S}_B$  is quite singular.

In terms of  $\mathbf{S}_B$  and  $\mathbf{S}_W$ , the criterion function  $J(\cdot)$  can be written as

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}. \quad (103)$$

This expression is well known in mathematical physics as the generalized Rayleigh quotient. It is easy to show that a vector  $\mathbf{w}$  that maximizes  $J(\cdot)$  must satisfy

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}, \quad (104)$$

for some constant  $\lambda$ , which is a generalized eigenvalue problem (Problem 42). This

# Final form of Fisher Discriminant

can also be seen informally by noting that at an extremum of  $J(\mathbf{w})$  a small change in  $\mathbf{w}$  in Eq. 103 should leave unchanged the ratio of the numerator to the denominator. If  $\mathbf{S}_W$  is nonsingular we can obtain a conventional eigenvalue problem by writing

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}. \quad (105)$$

In our particular case, it is unnecessary to solve for the eigenvalues and eigenvectors of  $\mathbf{S}_W^{-1} \mathbf{S}_B$  due to the fact that  $\mathbf{S}_B \mathbf{w}$  is always in the direction of  $\mathbf{m}_1 - \mathbf{m}_2$ . Because the scale factor for  $\mathbf{w}$  is immaterial, we can immediately write the solution for the  $\mathbf{w}$  that optimizes  $J(\cdot)$ :

$$\mathbf{w} = \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2). \quad (106)$$

Thus, we have obtained  $\mathbf{w}$  for Fisher's linear discriminant—the linear function yielding the maximum ratio of between-class scatter to within-class scatter. (The solution  $\mathbf{w}$  given by Eq. 106 is sometimes called the *canonical variate*.) Thus the classification has been converted from a  $d$ -dimensional problem to a hopefully more manageable one-dimensional one. This mapping is many-to-one, and in theory it cannot possibly reduce the minimum achievable error rate if we have a very large training set. In general, one is willing to sacrifice some of the theoretically attainable performance for the advantages of working in one dimension. All that remains is to find the threshold, that is, the point along the one-dimensional subspace separating the projected points.

# Case of Multivariate Normal pdfs

When the conditional densities  $p(\mathbf{x}|\omega_i)$  are multivariate normal with equal covariance matrices  $\Sigma$ , we can calculate the threshold directly. In that case we recall from Chapter 2 that the optimal decision boundary has the equation

$$\mathbf{w}'\mathbf{x} + w_0 = 0 \quad (107)$$

where

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad (108)$$

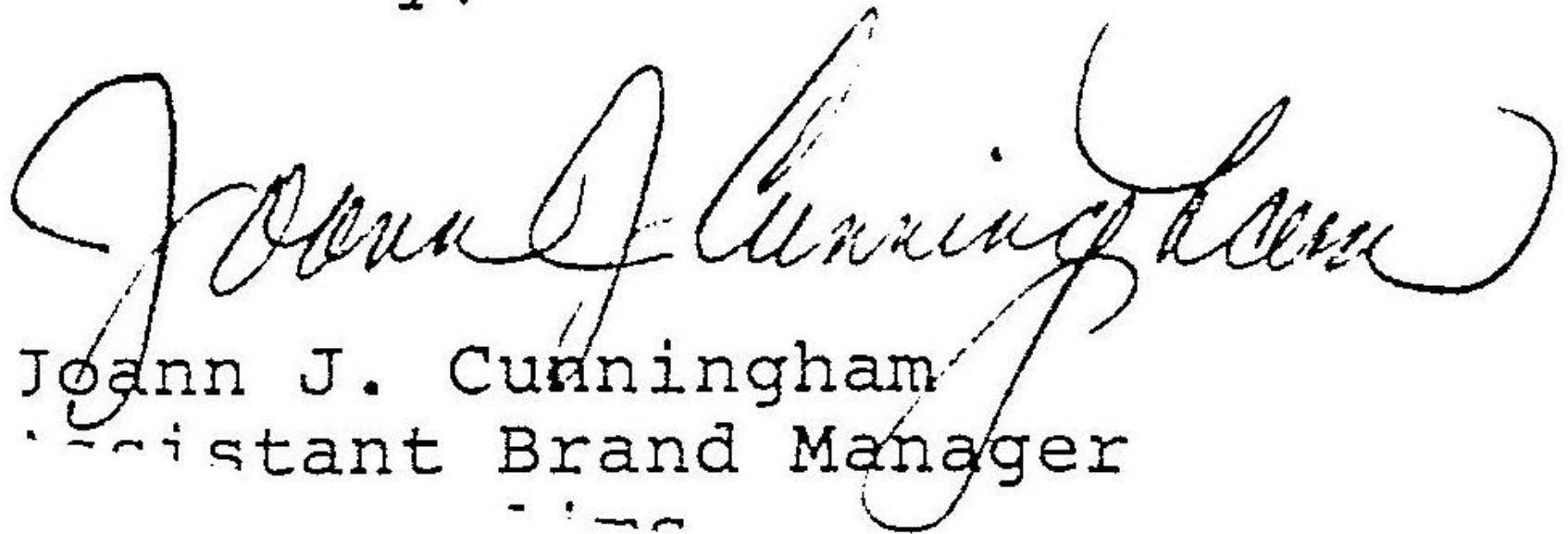
and where  $w_0$  is a constant involving  $\mathbf{w}$  and the prior probabilities. If we use sample means and the sample covariance matrix to estimate  $\boldsymbol{\mu}_i$  and  $\Sigma$ , we obtain a vector in

Bayes Decision rule is to compute Fisher LD and decide  $\omega_1$  if it exceeds a threshold and  $\omega_2$  otherwise

# Fisher's Linear Discriminant Example

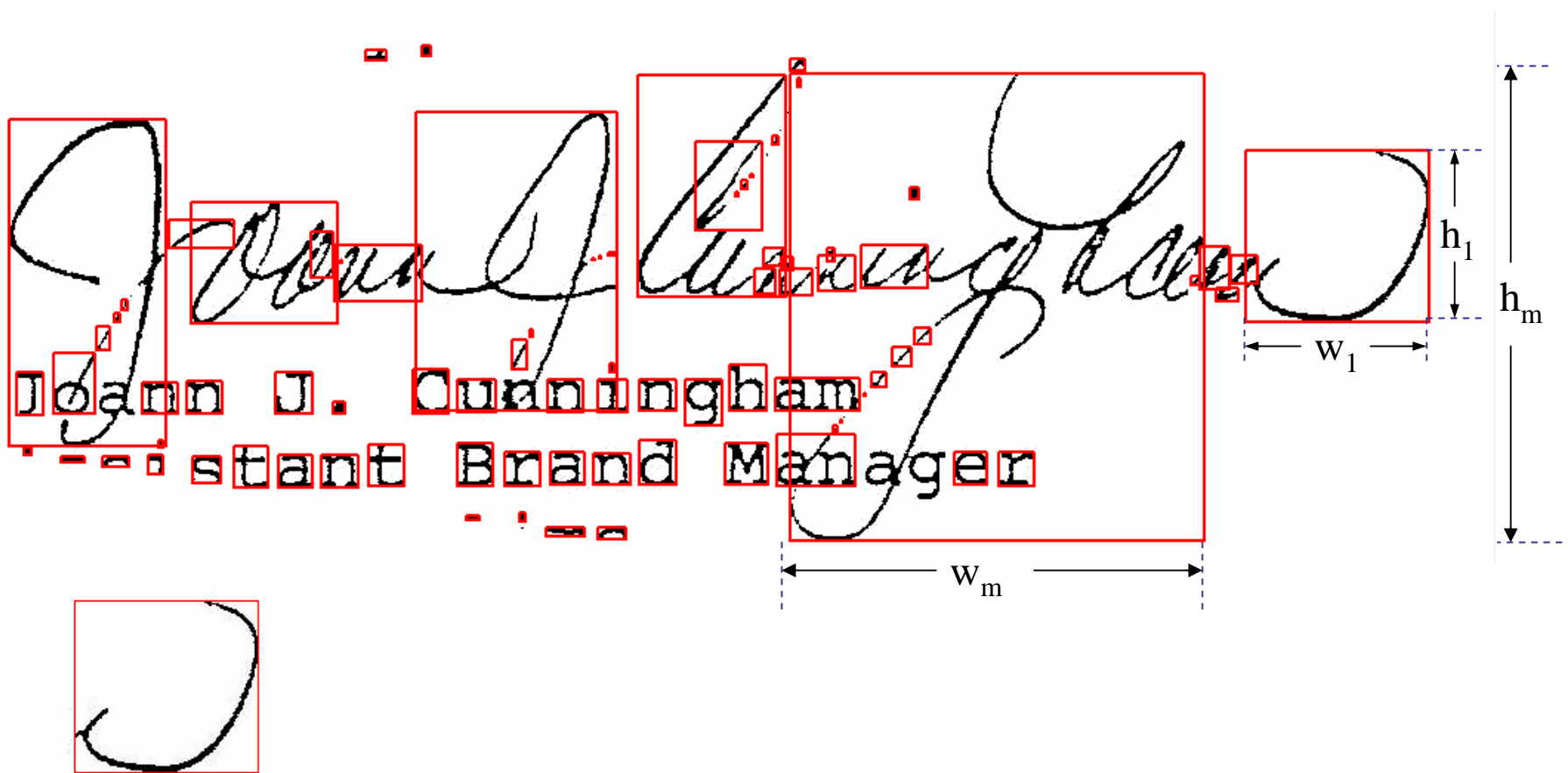
Discriminating between machine-  
print and handwriting

# Cropped signature image



Joann J. Cunningham  
Assistant Brand Manager

# Connected components and their features



$$x_1 = (h_1 + w_1) / (h_m + w_m) = 0.4034$$

$$x_2 = h_1 / w_1 = 0.9355$$

# Feature Values of Labelled Components (from 20 images)

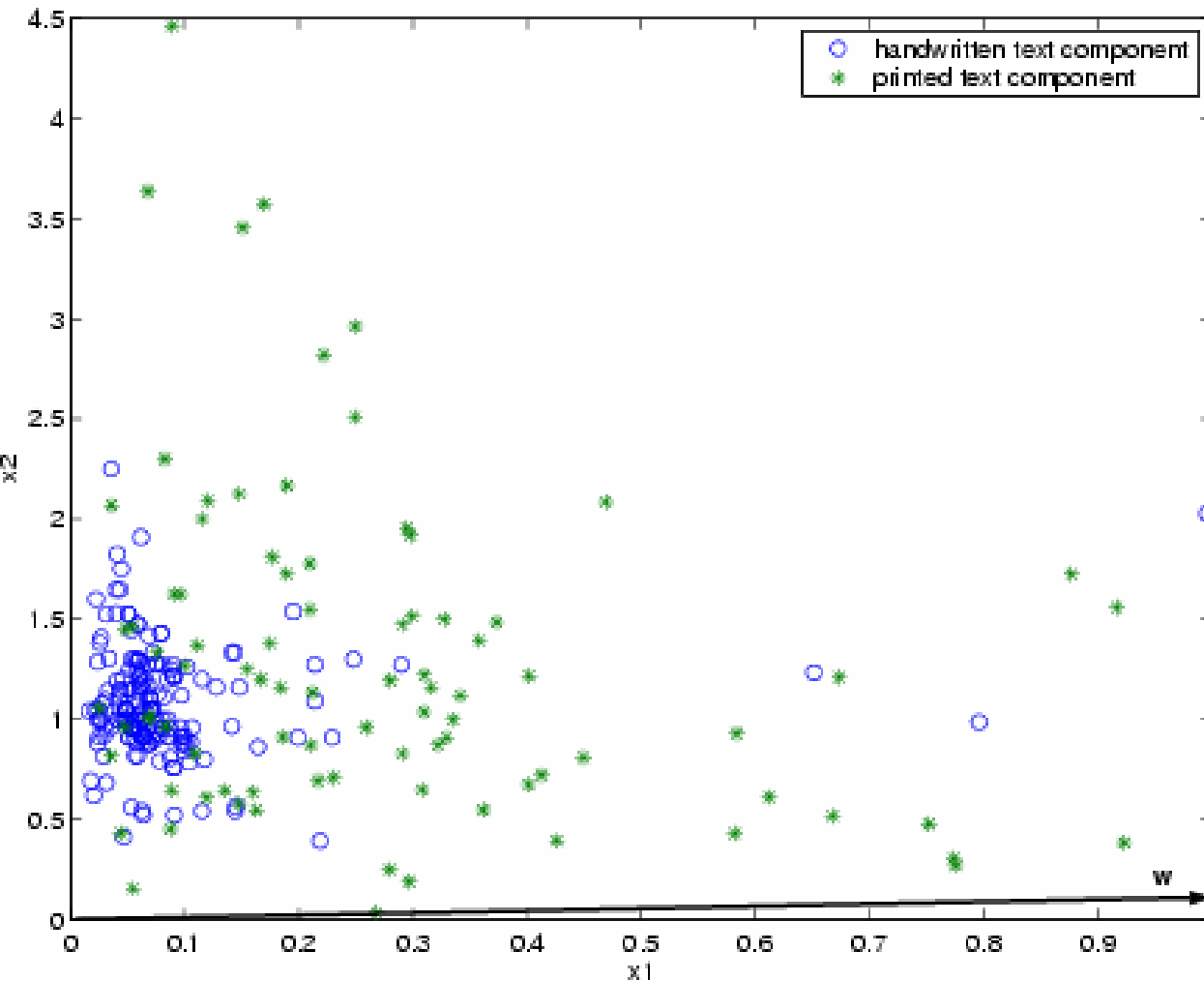
Print components

Component	$X_1$	$X_2$

Handwriting components

Component	$X_1$	$X_2$

# Feature distribution in 2 dimensional space



$$\mathbf{m}_1 = (0.0843 \ 1.0782)^T$$

$$\mathbf{m}_2 = (0.2725 \ 1.2652)^T$$

$$\mathbf{S}_w = \begin{bmatrix} 5.9554 & -3.2349 \\ -3.2349 & 73.3043 \end{bmatrix}$$

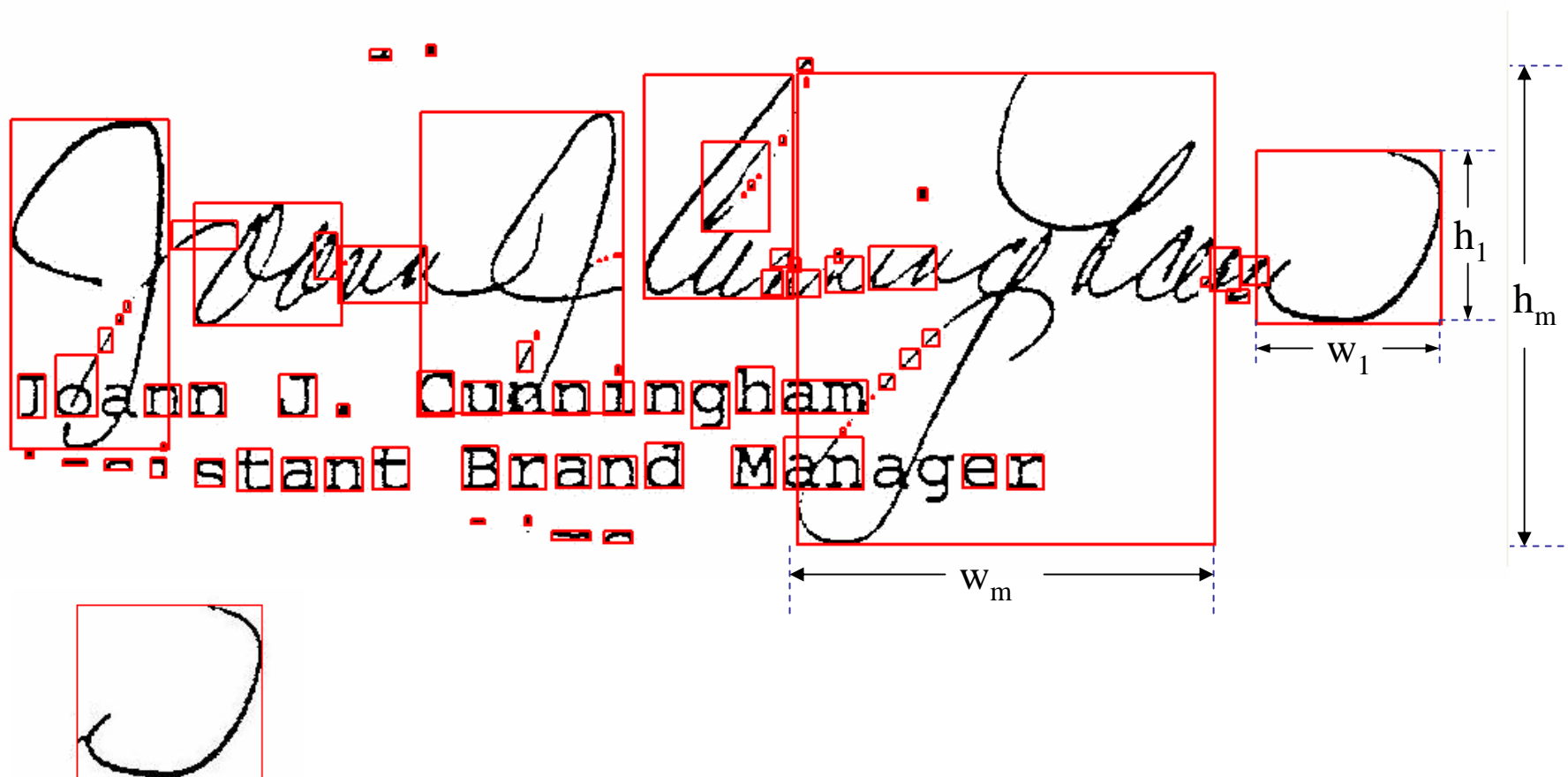
$$\mathbf{S}_w^{-1} = \begin{bmatrix} 0.1720 & 0.0076 \\ 0.0076 & 73.3043 \end{bmatrix}$$

$$\mathbf{w} = \mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2) = (-0.0338 \ -0.0040)^T$$

After projection  
Samples are 1-D  
A Bayes classifier  
is designed for the  
projected samples



# Discrete component analysis



$$x_1 = (h_1 + w_1) / (h_m + w_m) = 0.4034$$

$$x_2 = h_1 / w_1 = 0.9355$$

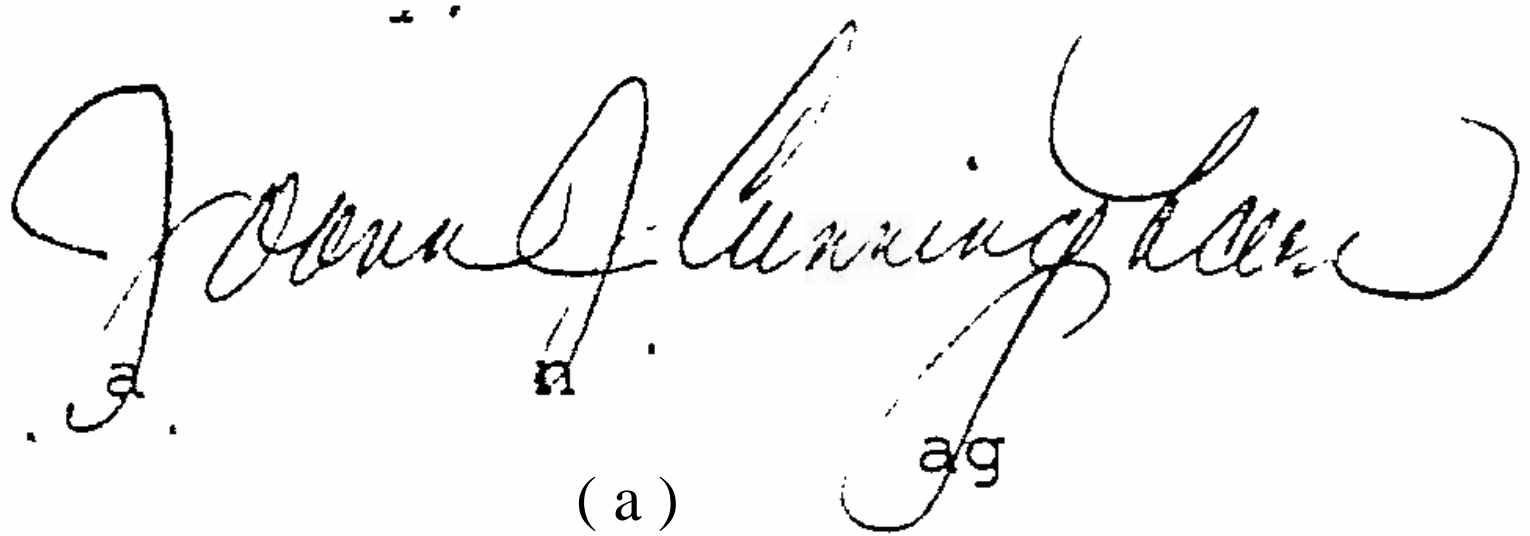
## The obtained printed text component candidates

John J. Cunningham  
Assistant Brand Manager

# Framework of printed text filtering

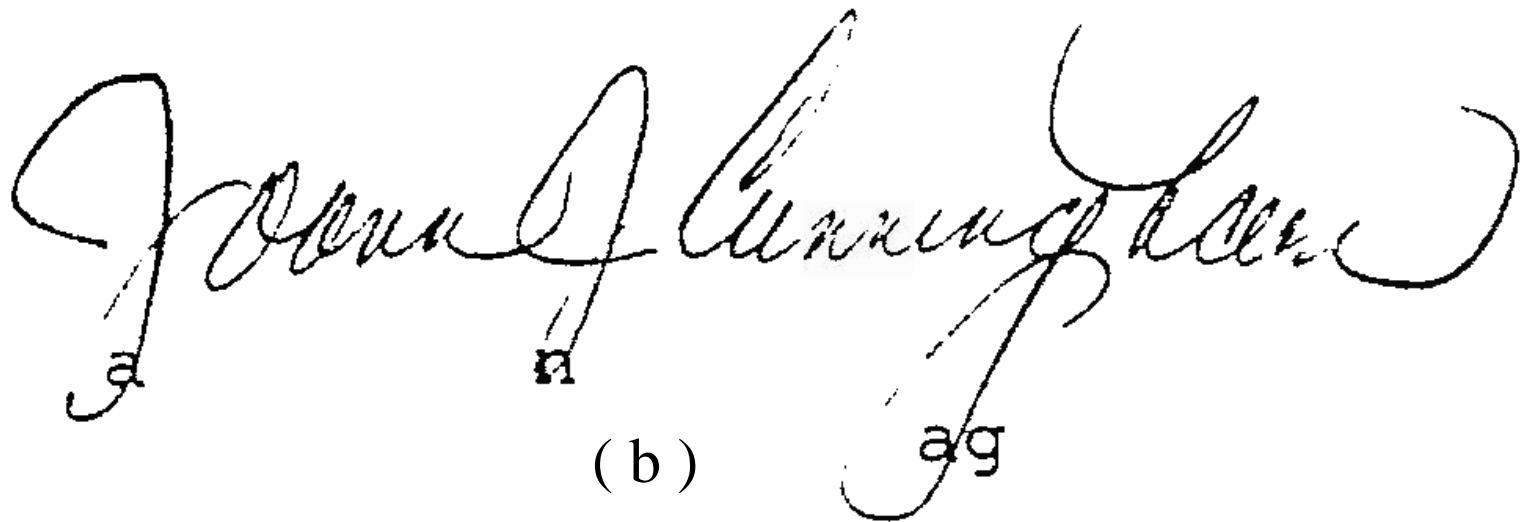
- Step 1:* Extract shape features of bounding box for each discrete component. Detect candidates by Fisher's Linear Discriminant and Bayesian Classification
- Step 2:* Remove candidates that satisfy the spatial relation defined for printed text components
- Step 3:* For candidates surviving from step2, remove isolated and small pieces.

Processed image after ( a ): Step 2, ( b ): Step 3 (final)



( a )

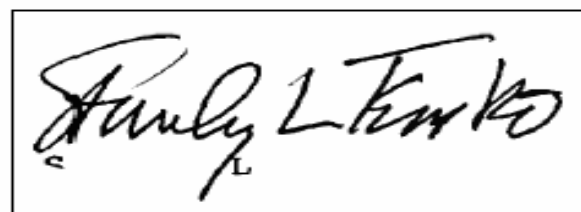
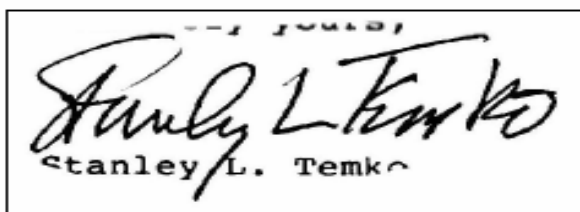
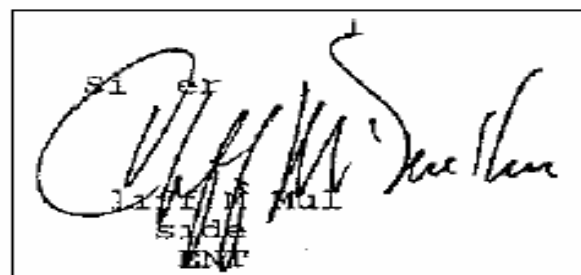
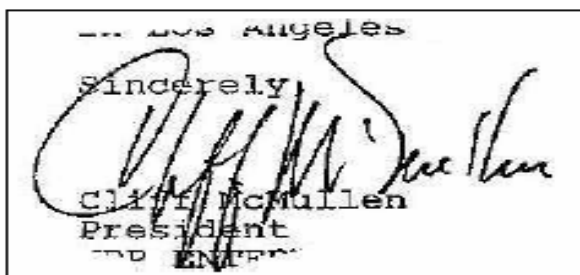
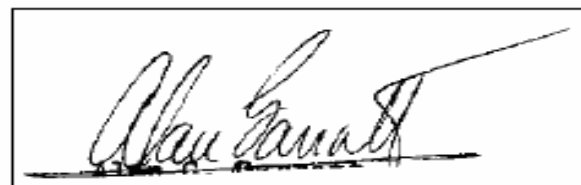
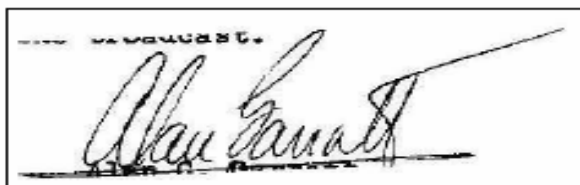
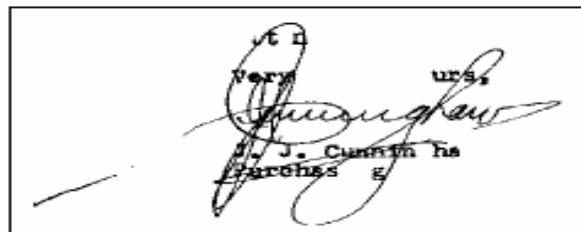
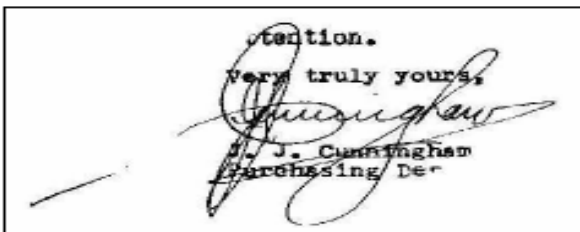
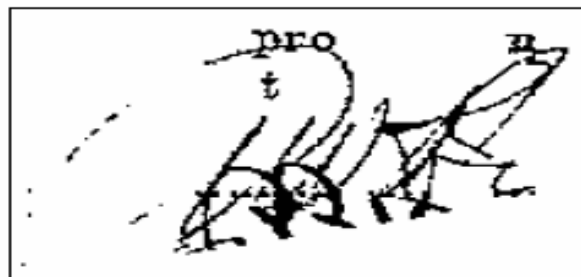
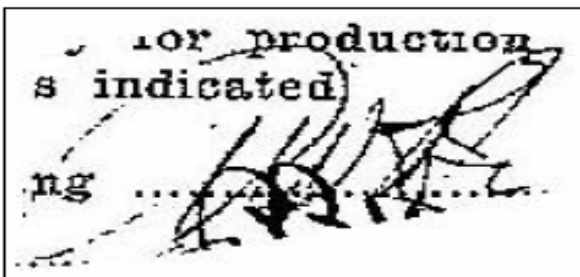
The image shows a handwritten signature 'John J. Cunningham' in cursive. Small black dots are placed above the 'n' and 'g' in 'Cunningham'. Below the signature, the letters 'a', 'n', and 'ag' are printed in a small, black, sans-serif font. 'a' is under the first 'a' in 'John', 'n' is under the 'n' in 'Cunningham', and 'ag' is under the 'g' in 'Cunningham'. The label '( a )' is centered below the signature.



( b )

The image shows the same handwritten signature 'John J. Cunningham' as in (a). Below the signature, the letters 'a', 'n', and 'ag' are printed in a small, black, sans-serif font. 'a' is under the first 'a' in 'John', 'n' is under the 'n' in 'Cunningham', and 'ag' is under the 'g' in 'Cunningham'. The label '( b )' is centered below the signature.

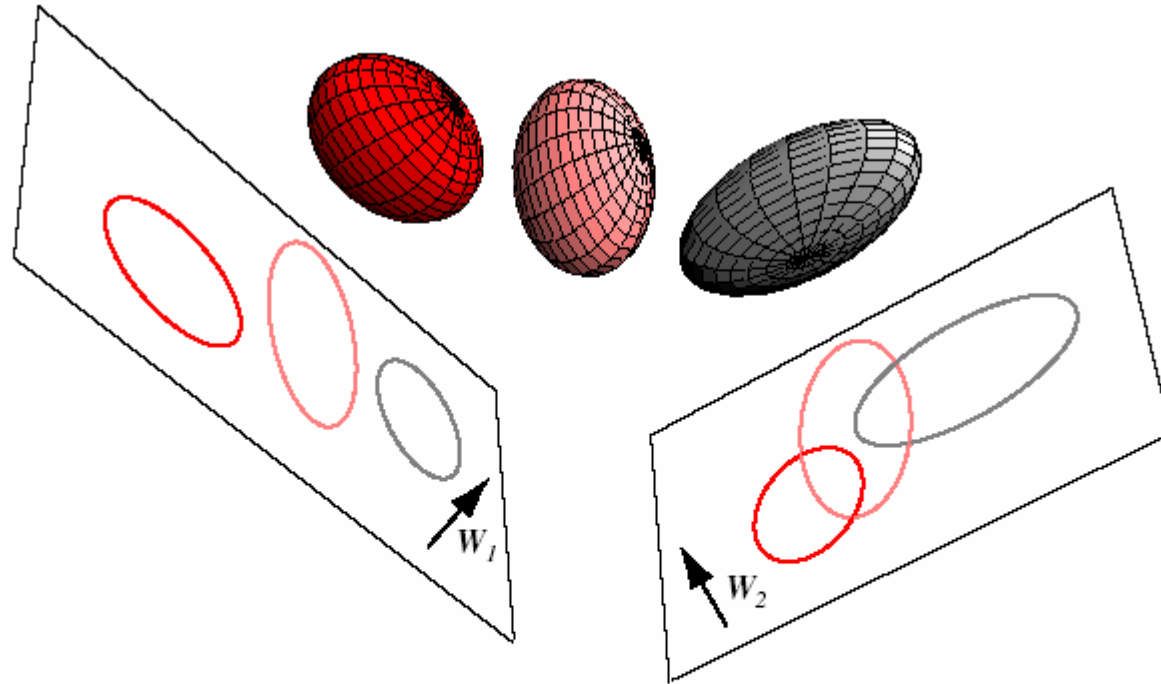
# Sample images before and after enhancement



# Multiple Discriminant Analysis

- c-class problem
- Natural generalization of Fisher's Linear Discriminant function involves  $c-1$  discriminant functions
- Projection is from a  $d$ -dimensional space to a  $c-1$  dimensional space

# Mapping from d-dimensional space to c-dimensional space $d=3, c=3$



Three three-dimensional distributions are projected onto two-dimensional subspaces, described by a normal vectors  $\mathbf{W}_1$  and  $\mathbf{W}_2$ . Informally, multiple discriminant methods seek the optimum such subspace, that is, the one with the greatest separation of the projected distributions for a given total within-scatter matrix, here as associated with  $\mathbf{W}_1$ .