

Apprentissage non-supervisé

1

• Typologie de la réduction de dimension

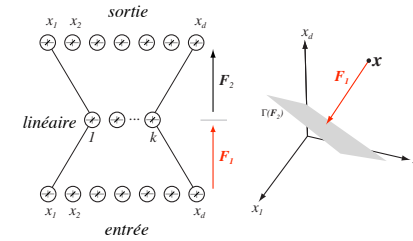
- méthode de base: **ACP**
- “groupement (clustering) des dimensions”
- extensions:
 - **ACP non-linéaire (NLPCA)**
 - **échelonnement multidimensionnel (multidimensional scaling – MDS)**
 - **cartes auto-organisatrices (self-organizing maps – SOM)**
 - **local linear embedding (LLE)**
 - **ISOMAP**
 - **courbes principales (principal curves)**

Apprentissage non-supervisé

2

• ACP non-linéaire – auto-encodage

- modèle de réseau de ACP

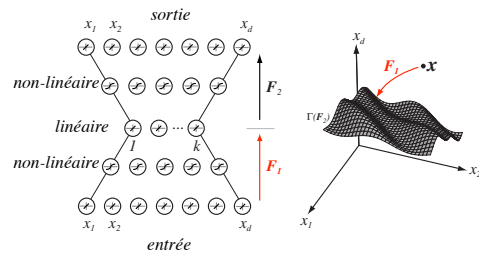


Apprentissage non-supervisé

3

• ACP non-linéaire – auto-encodage

- extension **non-linéaire**

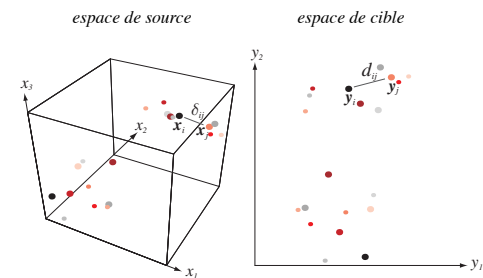


Apprentissage non-supervisé

4

• Échelonnement multidimensionnel (MDS)

- représentation de dimension réduite qui **préserve les distances**



- Échelonnement multidimensionnel (MDS)

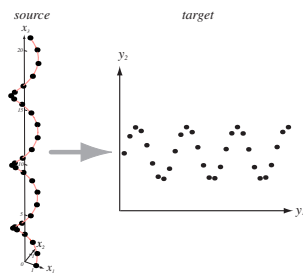
- fonctions d'erreur

- $J_{ee} = \frac{\sum_{i < j} (d_{ij} - \delta_{ij})^2}{\sum_{i < j} \delta_{ij}^2}$

- $J_{ff} = \sum_{i < j} \left(\frac{d_{ij} - \delta_{ij}}{\delta_{ij}} \right)^2$

- $J_{ef} = \frac{1}{\sum_{i < j} \delta_{ij}} \sum_{i < j} \frac{(d_{ij} - \delta_{ij})^2}{\delta_{ij}}$

- Échelonnement multidimensionnel (MDS)



- Échelonnement multidimensionnel (MDS)

- minimisation

- descente de gradient standard

- initialisation

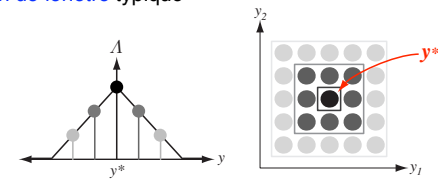
- les d' coordonnées avec les variances plus grandes
 - ACP avec d' composantes

- Cartes auto-organisatrices (SOM)

- x_i appartient à V_ℓ avec un poids $W_{i,\ell}$

- $W_{i,\ell}$ ne dépend que de la distance entre v_ℓ et $v(x_i)$

- fonction de fenêtre typique



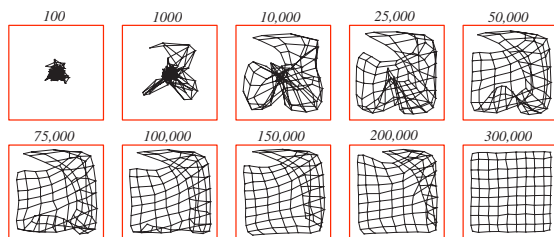
• Cartes auto-organisatrices (SOM)

```

SOM( $X_n$ )
1  $C^{(0)} \leftarrow \{v_1^{(0)}, \dots, v_k^{(0)}\}$ 
2  $j \leftarrow 0$ 
3 faire
4   recalculer  $W^{(j)}$ 
5   pour  $\ell \leftarrow 1$  à  $k$  faire
6      $v_\ell^{(j+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n W_{i,\ell}^{(j)} x_i$ 
7    $j \leftarrow j + 1$ 
8 jusqu'à changement > seuil
    
```

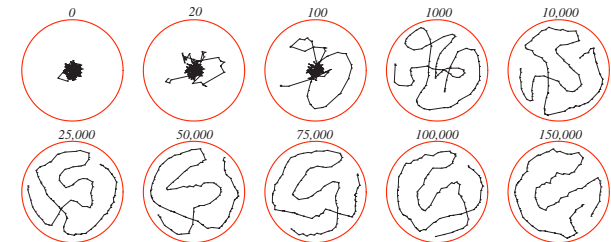
• Cartes auto-organisatrices (SOM)

• 2 dimensions \rightarrow 2 dimensions



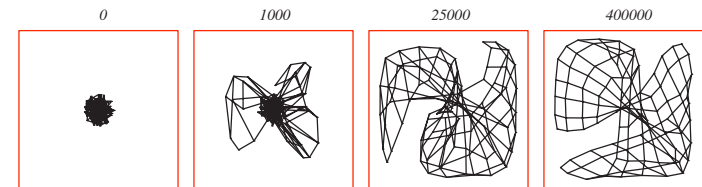
• Cartes auto-organisatrices (SOM)

• 2 dimensions \rightarrow 1 dimension



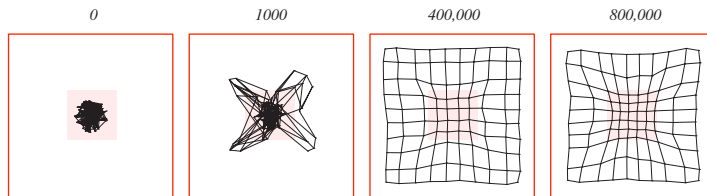
• Cartes auto-organisatrices (SOM)

• problème: **minimum local**



- Cartes auto-organisatrices (SOM)

- estimation de densité



- Cartes auto-organis. (SOM) – [théorie de communication](#)

- Codage de [source](#) – quantification vectorielle:

- [fonction d'erreur](#): $J_s = \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{v}_{(\mathbf{x}_i)}\|^2$

- Codage de [canal](#) – correction d'erreur:

- [probabilité d'erreur](#) d'un bit: p
- [distance de Hamming](#) entre des mots de code: $d_{i,j} = d_H(c(\mathbf{v}_i), c(\mathbf{v}_j))$
- probabilité d'erreur de [code](#): $p_{i,j} = p^{d_{i,j}}(1-p)^{d-d_{i,j}}$
- [fonction d'erreur](#): $J_c = \sum_{i=1}^n \sum_{j=1}^c \|\mathbf{v}_{(\mathbf{x}_i)} - \mathbf{v}_j\|^2 p_{i,j}$

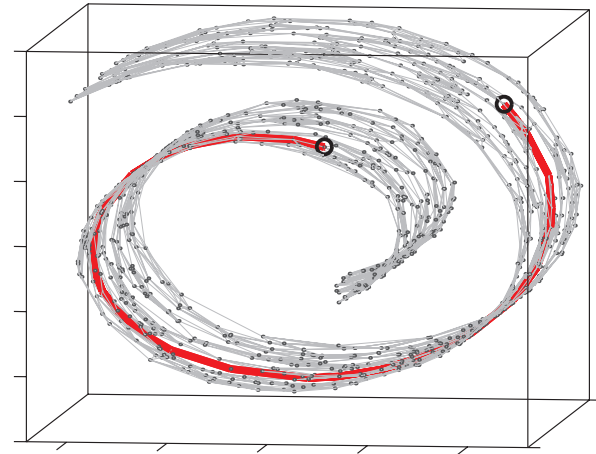
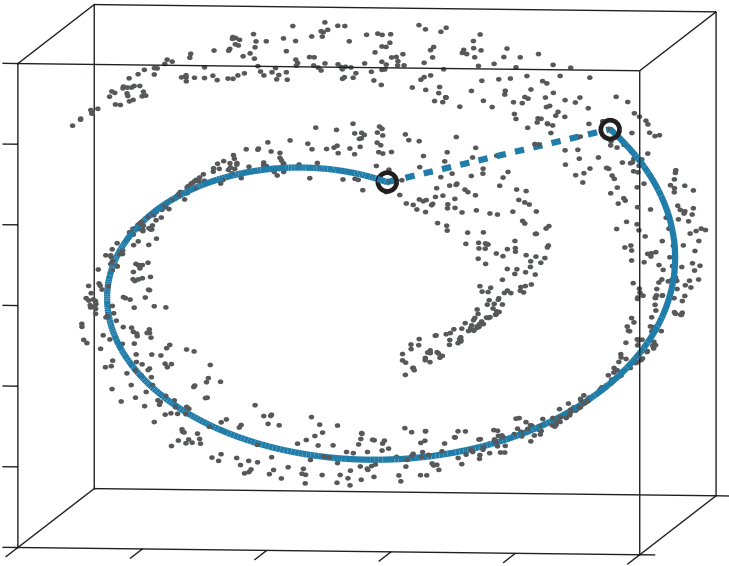
- Codage [conjoint de canal-source](#)

- [fonction d'erreur](#): $J_{s+c} = \sum_{i=1}^n \sum_{j=1}^c \|\mathbf{x}_i - \mathbf{v}_j\|^2 p_{i,j}$

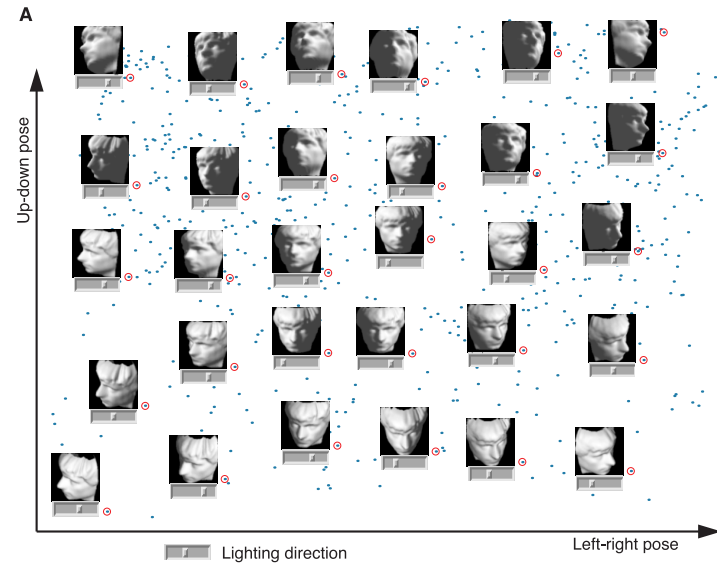
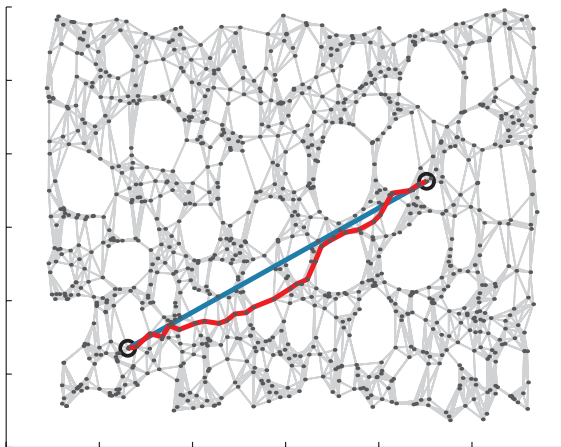
- Problème générale: surfaces compliquées → [minima local](#)

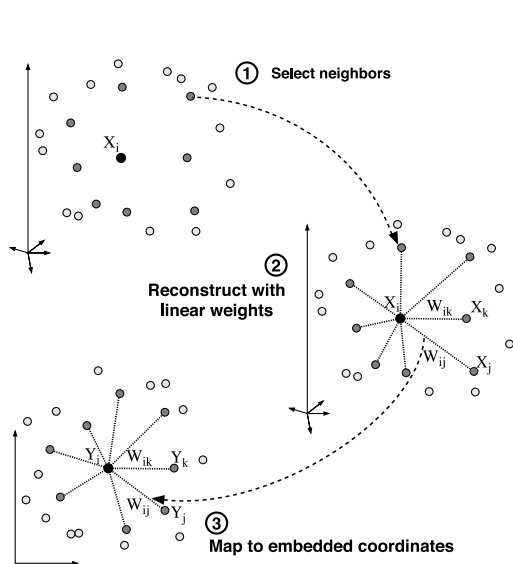
- Solution 1: ISOMAP

- [distance géodésique](#): chemins plus courts dans le graphe de similarité
- MDS standard sur les distances géodésiques



C





The con-
²³
 these recons
 symmetry: 1
 are invaria
 translations
 bors. By syi
 structure we
 metric prop
 opposed to
 ticular fram
 invariance t
 forced by t
 rows of the
 Suppose
 nonlinear m
 $\ll D$. To
 exists a lin
 translation,
 maps the l
 each neighb
 mates on the
 structure we
 ric propertie
 exactly suc
 expect their
 trv in the o

Apprentissage non-supervisé

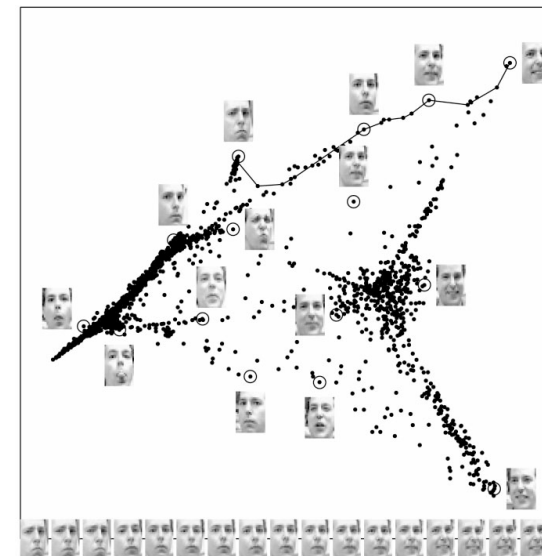
• Solution 2: Local linear embedding (LLE)

- Étape 1: trouver l'ensemble des voisins V_{x_i}
- Étape 2: approximer les points avec une combinaison linéaire de leurs plus proches voisins:

$$\min_W \sum_{i=1}^n \left\| x_i - \sum_{x_j \in V_{x_i}} w_{i,j} x_j \right\|^2$$

- Étape 3: reconstruire les points dans l'espace de projection en utilisant les mêmes poids:

$$\min_Y \sum_{i=1}^n \left\| y_i - \sum_{x_j \in V_{x_i}} w_{i,j} y_j \right\|^2$$



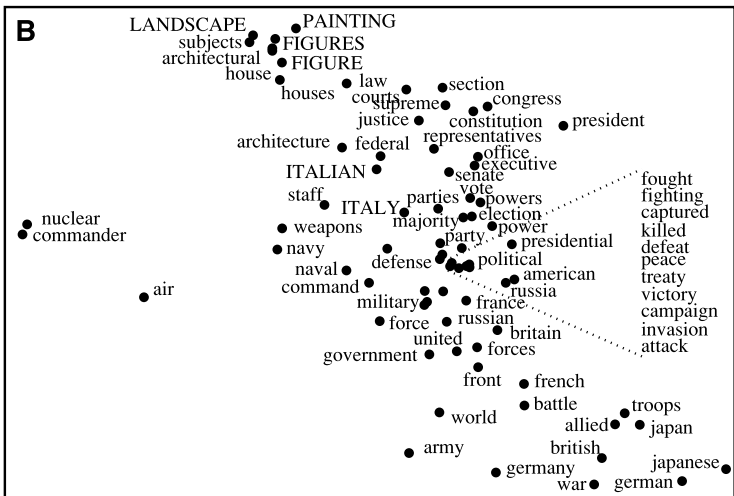
valid for
 particular
 struct the
 should al
 fold coor

LLE c
 mapping
 step of th
 observati
 vector Y_i
 nates on th
 d -dimensi
 embeddin

Φ

This cost
 based on
 but here
 mizing th
 cost in E
 vectors Y
 the probl
 by solvin
 lem (9),
 tors prov
 coordinat

Imple
 straightfo
 points we
 est neigh
 tance or

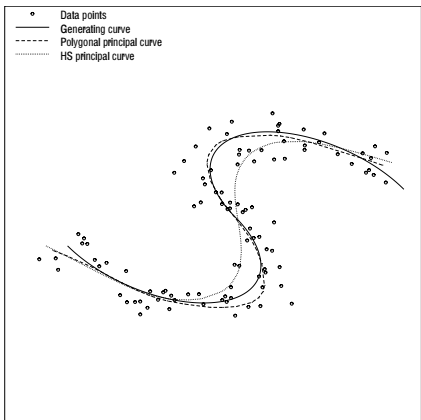


Apprentissage non-supervisé

- désavantage d'ISOMAP:
 - temps d'exécution: $O(n^3)$
- projeter des nouveaux points
 - construire la fonction de projection explicitement
 - problème d'interpolation
 - problème d'apprentissage supervisé (régression multidimensionnelle)

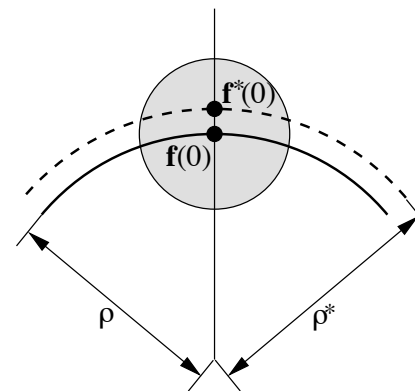
Apprentissage non-supervisé

- Problème: bruit

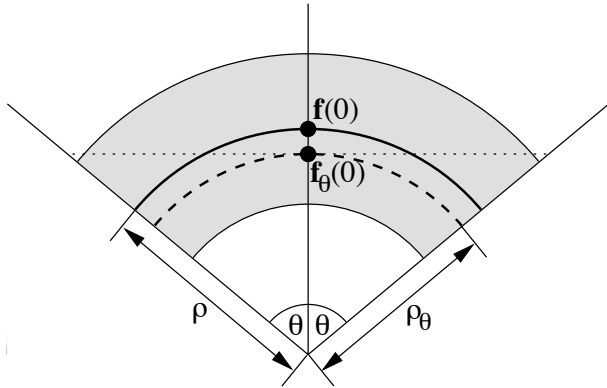


Apprentissage non-supervisé

- Le biais du modèle

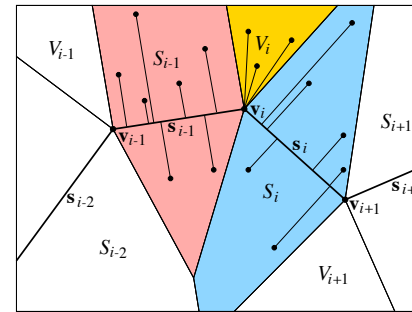


- Le biais de l'estimation

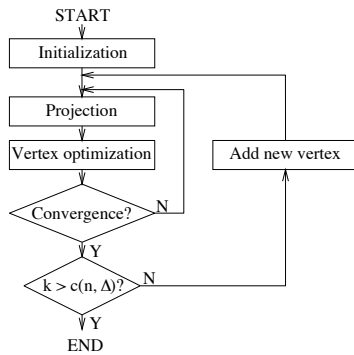


- Solution: courbes principales polygonales

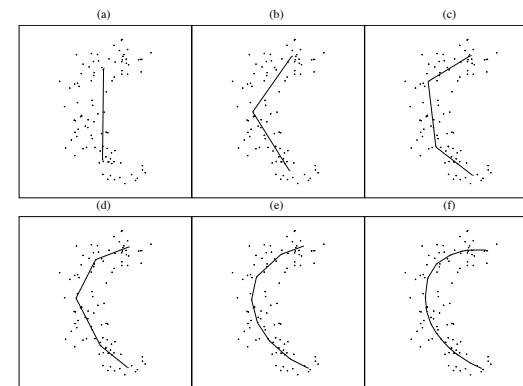
- Mesurer la distance de la courbe au lieu des sommets



- Courbes principales polygonales

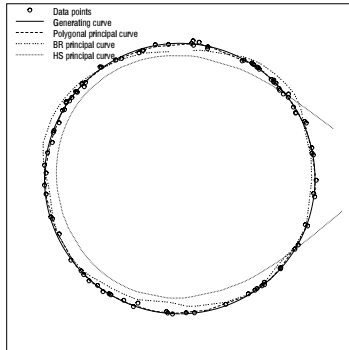


- Courbes principales polygonales



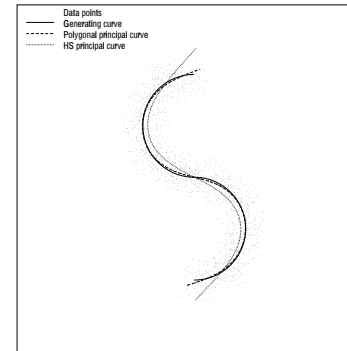
• Courbes principales polygonales

• bruit réduit



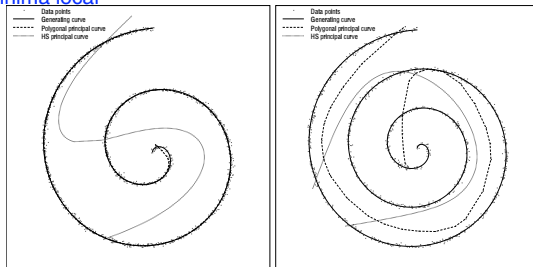
• Courbes principales polygonales

• beaucoup de points



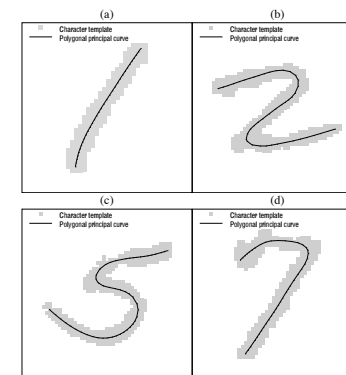
• désavantages des courbes principales:

• minima local

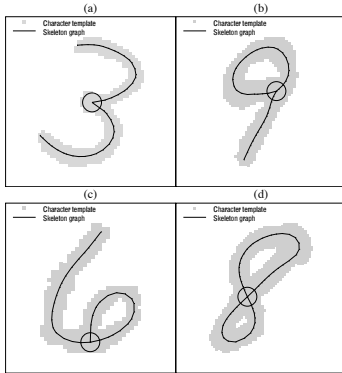


- extension aux surfaces n'est pas évident
→ la plupart des applications sont dans le traitement d'image

• Skeletisation des caractères



• Skeletisation des caractères



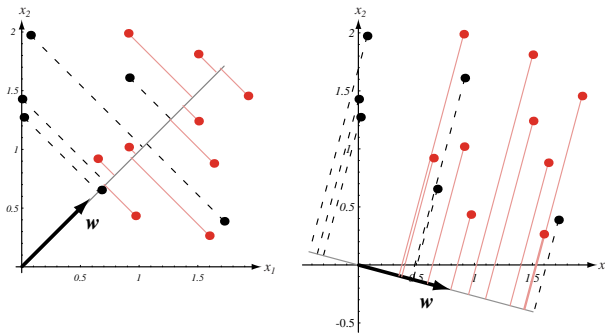
• Apprentissage non-supervisé pour la classification: **analyse discriminante**

- but: trouver la meilleure projection qui **préserv**e l'information **discriminante**

• Discriminante de Fisher

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

• Analyse discriminante



• Idée 1: **séparer les moyennes projetées**

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

- $\tilde{m}_i = \frac{1}{n_i} \sum_{y \in I_i} y = \frac{1}{n_i} \sum_{\mathbf{x} \in D_i} \mathbf{w}^t \mathbf{x}$

- trouver \mathbf{w} qui maximise $|\tilde{m}_1 - \tilde{m}_2| = |\mathbf{w}^t (\mathbf{x}_1 - \mathbf{x}_2)|$

• Idée 2: **séparer les moyennes projetées normalisées par les variances par classe**

- $\hat{s}_i^2 = \sum_{y \in I_i} (y - \tilde{m}_i)^2$

- $J(\mathbf{w}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\hat{s}_1^2 + \hat{s}_2^2}$

- Maximiser $J(\mathbf{w})$:

- $\mathbf{S}_i = \sum_{\mathbf{x} \in D_i} (\mathbf{x} - i)(\mathbf{x} - i)^t$

- $\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$

- $\hat{s}_i^2 = \sum_{\mathbf{x} \in D_i} (\mathbf{w}'\mathbf{x} - \mathbf{w}'_i)^2 = \sum_{\mathbf{x} \in D_i} \mathbf{w}'(\mathbf{x} - i)(\mathbf{x} - i)^t \mathbf{w} = \mathbf{w}'\mathbf{S}_i \mathbf{w}$

- $\hat{s}_1^2 + \hat{s}_2^2 = \mathbf{w}'\mathbf{S}_W \mathbf{w}$

- $\mathbf{S}_B = (1-2)(1-2)^t$

- $(\hat{m}_1 - \hat{m}_2)^2 = (\mathbf{w}'_1 - \mathbf{w}'_2)^2 = \mathbf{w}'(1-2)(1-2)^t \mathbf{w} = \mathbf{w}'\mathbf{S}_B \mathbf{w}$

- $J(\mathbf{w}) = \frac{\mathbf{w}'\mathbf{S}_B \mathbf{w}}{\mathbf{w}'\mathbf{S}_W \mathbf{w}}$

- $\mathbf{w}_{max} = \mathbf{S}_W^{-1}(1-2)$