

# Analyse Discriminate Linéaire

HUYLENBROECK Florent, DELFOSSE Charly, JOSSE Thomas

7 juin 2019

# Table des matières

## 1 Introduction

## 2 Fonctionnement

- Règle géométrique
- Approche statistique

## 3 Exemple

# Introduction

# Qu'est-ce que l'ADL ?

Cette technique fait partie des techniques d'analyse discriminante prédictive. Le but est de pouvoir expliquer et prédire l'appartenance d'un individu à un groupe prédéfini à partir de caractéristiques qui ont été mesurées au préalable à l'aide de variables prédictives. On peut la comparer à la régression logistique.

# Variables utilisées

$X$  contient les  $j$  variables prédictives [ $X = (X_1, \dots, X_j)$ ],  $n$  est le nombre d'observations réparties dans  $K$  groupe d'effectifs  $n_k$ ,  $Y$  est notre variable à prédire qui prend des valeurs dans l'ensemble  $\{y_1, \dots, y_K\}$ .  $\pi_k = P(Y = y_k)$

## Fonctionnement

# But Final

Cette règle ne prend en compte aucune hypothèse probabiliste. Elle consiste à calculer la distance de  $x$  (vecteur des variables explicatives sur un individu que l'on veut classer) à chacun des  $K$  centres de gravité  $g_1, \dots, g_K$  et affecter  $x$  au groupe le plus proche. Et cette distance du nouvel individu au groupe  $k$  peut être trouvée via la formule :

$$d^2(x, g_k) = (x - g_k)' \mathbf{W}^{-1} (x - g_k)$$

, Où  $\mathbf{W}$  est la matrice des variance-covariance intra-groupe.

# Développement

Pour pouvoir faire cela, nous allons définir notre fonction linéaire discriminante du groupe  $k$  pour savoir si  $x$  appartient au groupe  $k^*$  tel que :

$$k^* = \arg \max_{k=1,\dots,K} d^2(x, g_k)$$

que l'on peut réécrire :

$$k^* = \arg \max_{k=1,\dots,K} L_K(x)$$

où

$$L_K(x) = x' \mathbf{W}^{-1} g_k - \frac{1}{2} g_k' \mathbf{W}^{-1} g_k$$

Et  $L_K(x)$  est notre fonction linéaire discriminante du groupe  $k$ . Chaque  $L_K(x)$  définit une fonction score qui donne une note représentant la probabilité d'appartenance au groupe de la fonction linéaire.  $X$  est donc affecté au groupe dont le score est le plus grand.



# La règle bayesienne

Cette règle consiste à produire une estimation de la probabilité après notre affectation. Cela veut dire que nous devons réaliser une estimation pour une probabilité conditionnelle :

$$P(Y = y_k|X) = \frac{P(Y = y_k) \times P(X|Y = y_k)}{\sum_{i=1}^K P(Y = y_i) \times P(X|Y = y_i)}$$

Nous avons  $P(Y = y_k)$  qui est la probabilité d'appartenance à la classe  $y_k$  et  $P(X|Y = y_k)$  qui est la fonction de densité des  $x$  par rapport à l'appartenance à la classe  $y_k$ .

Et cette règle permet d'affecter une nouvelle observation  $x$  au groupe  $k^*$  tel que :

$$k^* = \arg \max_{k=1,\dots,K} P(Y = y_k|X)$$

$$k^* = \arg \max_{k=1,\dots,K} \pi_k f_k(x)$$

# Hypothèse d'homoscédasticité

L'homoscédasticité est le fait que les variances de chaque groupe soit équivalente (son contraire est l'hétéroscédasticité). Cela veut dire que les données sont réparties de la même manière autour de leur moyenne, ou centre de gravité.

Pour pouvoir effectuer de l'analyse discriminante linéaire c'est la principale hypothèse que l'on doit appliquer à nos données.

# Ce qu'on veut

On veut replacer nos données dans un nouveau repère tel que les points sont le plus possible au sein d'un même groupe, mais ceux-ci doivent être aussi distants que possible entre les 2 groupes.

## Exemple

Mettre package utilisé avec la méthode discrimin + explications

Montrer exemple