

## Exercice 1 : Analyse Discriminante Probabiliste

On considère un échantillon de  $n$  couples de données observées  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  issue de  $K$  groupes (classes) homogènes où chaque individu  $\mathbf{x}_i$ ,  $i = 1, \dots, n$  est décrit par  $d$  variables  $x_{ij}$ ,  $j = 1, \dots, d$ . La variable  $y_i$  représente la classe du  $i$ ème individu et est à valeurs dans  $\{1, \dots, K\}$  :  $y_i = k \in \{1, \dots, K\}$ . On suppose que les données de chaque classe  $k$  sont générées selon une loi normale  $(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$  où la matrice  $\boldsymbol{\Sigma}$  est supposée la même pour toutes les classes.

1. Quel est le modèle conditionnel  $\mathbb{P}(Y|\mathbf{X} = \mathbf{x})$  induit par ce modèle génératif?
2. Montrer qu'il correspond à un modèle probabiliste de classification linéaire (frontières de décision linéaires entre les classes)
3. Quelles frontières de décision obtient-on si l'on autorise chaque classe à avoir sa propre matrice de covariance?
4. Appliquez le modèle aux données en utilisant la méthode du maximum de vraisemblance pour estimer ses paramètres

## Exercice 2 : Analyse Factorielle Discriminante

On considère les données des départements français. Les données sont disponibles dans le fichier `depart.data`.

On se propose de mettre en évidence les plus grandes disparités inter-régionales et donc de rechercher les variables ou combinaisons de variables expliquant au mieux le découpage régional. Pour simplifier, nous procédons à des regroupements afin de construire des régions moins nombreuses comprenant des nombres de départements plus semblables. D'autre part, la région "Ile de France", trop particulière et donc trop facilement discriminée, est laissée à part.

Les données des départements français Les données proviennent du Groupe d'Etude et de Reflexion Inter-Regional (GERI). Elles décrivent, quatre grands thèmes : la démographie, l'emploi, la fiscalité directe locale, la criminalité, de chacun des départements français métropolitains et de la Corse. Les indicateurs ont été observés pendant l'année 1990, ils sont, pour la plupart, des taux calculés relativement à la population totale du département concerné. Voici leur liste :

- identificateur : numéro du département,
- identificateur : code du département,
- identificateur : code de la région
- URBR : indicateur de concentration de la population mesurant le caractère urbain ou rural du département,
- TXCR : taux de croissance de la population sur la période intercensitaire 1982-1990,
- JEUN : part des 0-19 ans dans la population totale,
- AGE : part des plus de 65 ans dans la population totale,
- FE90 : taux de fécondité (pour 1000) égal au nombre de naissances rapportées au nombre de femmes âgées de 15 à 49 ans en moyenne triennale,
- ETRA : part des étrangers dans la population totale,
- CHOM : taux de chômage,
- CRIM : taux de criminalité : nombre de délits par habitant,
- FISC : produit, en francs constants 1990 et par habitant des quatre taxes directes locales (professionnelle, habitation, foncier bâti, foncier non bâti). Suivent ensuite les parts de chaque profession et catégorie socioprofessionnelle (PCS) dans la population active occupée du département :
- AGRI : agriculteurs,
- ARTI : artisans,
- CADR : cadres supérieurs,

- EMPL : employ es,
- OUVR : ouvriers,
- PROF : professions interm ediaires
- Charger les données
- Afficher les
- Regrouper les en régions : Nd (Nord) ; Es (Est) ; Ws (West) ; CN (Centre Nord) ; Centre West (CW) ; Centre Est (CE) ; Sud West (SW) ; Sud Est (SE). (vous pouvez utiliser `select`)
- Appliquer une analyse factorielle discriminante et analyser les résultats

## Exercice 2 : Analyse Discriminante décisionnelle

On considère d'infarctus disponible dans le fichier `infarctu.txt`.

Ce jeu de données provient d'une étude épidémiologique où le but de l'étude était d'évaluer l'existence d'un risque plus élevé de survenue d'un infarctus du myocarde chez les femmes qui utilisent ou ont utilisé des contraceptifs oraux.

101 observations et 8 variables. FRCAR Frequence Cardiaque (i.e. heart rate)

INCAR Index Cardiaque (cardiac index)

INSYS Index Systolique (systolic index)

PRDIA Pression Diastolique (diastolic pressure)

PAPUL Pression Arterielle Pulmonaire (pulmonary artery pressure)

PVENT Pression Ventriculaire (ventricular pressure)

REPUL Resistance Pulmonaire (pulmonary resistance)

PRONO Pronostic (prognosis) : a factor with levels dead and survive

- Charger les données
- Afficher les
- Appliquer une analyse discriminante linéaire et analyser les résultats
- Appliquer une analyse discriminante quadratique et analyser les résultats
- Dans chacun des deux cas précédents, considérer le cas où les proportions des classes sont égales ou non
- utiliser une validation croisée (cross-validation en anglais) pour évaluer les capacités en généralisation de l'analyse discriminante quadratique pour ce jeux de données.