

Mesures de similarités, dissimilarités et distances

I. Les étapes préliminaires à une analyse de données

1) Analyse des données

Un ensemble de méthodes dont l'objectif essentiel est la mise en relief des relations existantes entre les objets, entre les paramètres qui les caractérisent et entre les objets et les paramètres

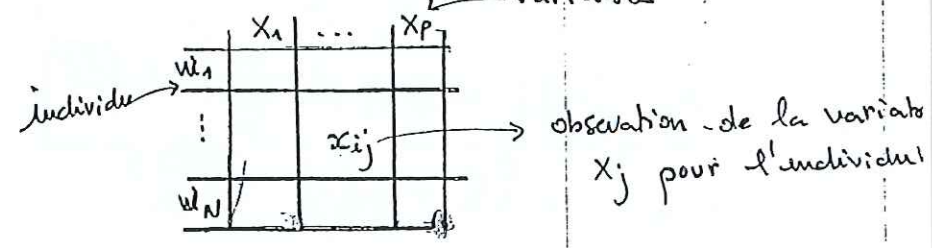
Objet \Leftrightarrow item \Leftrightarrow sujet \Leftrightarrow Forme \Leftrightarrow individu ...
 paramètre \Leftrightarrow variable \Leftrightarrow description ...

2) Modélisation des individus / variables.

- \mathcal{O} : ensemble des individus observés
 - \mathcal{O} : espace des observations.
 - \mathcal{S} : une structure algébrique définie sur \mathcal{O}
- Une variable $v : \mathcal{O} \rightarrow \mathcal{O}$ muni de la structure \mathcal{S} .

| Cardinal Structure \mathcal{S} | Continu | Fini ou Dénombrable | |
|-------------------------------------|--------------------|------------------------|-----------|
| $= \#$ | | CSP | Nominale |
| \leq | Âge Température | Rang Ressemblance | Ordinale |
| $\leq + \times$ | Revenu | | Mesurable |
| | Quantitative | Qualitative | Variable |

II. Type de Tableaux de données à analyser.



a) Tableaux quantitatifs

$\forall j \in \{1..p\}$ X_j est une variable quantitative.

b) Tableaux de contingence (fréquence)

$\forall i \in \{1..N\}, j \in \{1..p\}$ x_{ij} est le nb d'occurrences ou la fréquence d'apparition de la variable X_j pour l'ind w_i .

- exp:
- w_i : un utilisateur identifié par son IDP
 - X_j : un thème spécifique d'un site
 - x_{ij} : la fréquence d'accès de l'utilisateur w_i au thème X_j du site

c) Tableaux de préférence

$j \in \{1..p\}$ X_j : les modalités d'une variable qualitative
 x_{ij} : une note exprimant un ordre de préférence de la modalité X_j par l'individu w_i .

d) Tableaux binaires.

Origine : données d'enquêtes
 X_j : une modalité d'une variable qualitative.
 x_{ij} : la valeur 0, ou 1 exprimant la présence ou non de la modalité X_j chez l'ind w_i .

III e) Tableaux de proximité (individu x individu)

| | | | |
|----------|-------|----------|-------|
| | w_1 | ... | w_N |
| w_1 | | | |
| \vdots | | d_{ij} | |
| w_N | | | |

mesure la proximité (similarité, Dissimilarité, distance, ...)
entre w_i et w_j

f) Tableaux Hétérogène

Tableaux combinant des variables de # types.

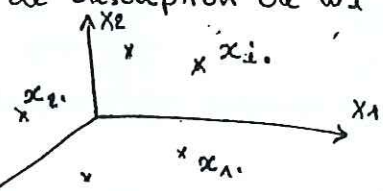
3) Résumés numériques et espaces associés

3.1 Notations:

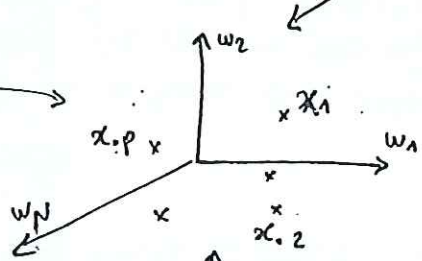
$$X = \begin{matrix} & 1 & \dots & p \\ \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} & \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{iN} \end{bmatrix} \end{matrix}$$

$$x_{i\cdot} = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$$

vecteur $(p \times 1)$ de description de w_i



$$x_{\cdot j} = \begin{pmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{Nj} \end{pmatrix}$$



l'espace des individus - de dimension p

l'espace des variables de

IV 3.2 La matrice des poids

$$D = \begin{pmatrix} p_1 & & 0 \\ & \ddots & \\ 0 & & p_N \end{pmatrix}$$

avec $\sum_{i=1}^N p_i = 1$
 $p_i \geq 0$
 p_i : poids associé à l'individu w_i

En général, on considère $D = \frac{1}{N} I$ (I : matrice identité)

3.3 Centre de gravité g

$$g = \mathbf{1}' D X = (\bar{x}_{\cdot 1}, \dots, \bar{x}_{\cdot p})$$

$$\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\bar{x}_{\cdot j} = \sum_{i=1}^N p_i x_{ij}$$

3.4 Centrage et réduction des données

On note la matrice $Y (y_{ij})$ / $y_{ij} = x_{ij} - \bar{x}_{\cdot j}$

la matrice centrée associée à X ($Y = X - \mathbf{1}g'$)

3.5 Matrice de Variance-Covariance / Corrélation

$$V = X' D X - g'g = Y' D Y$$

$$v_{jj'} = \sum_{i=1}^N p_i (x_{ij} - \bar{x}_{\cdot j})(x_{ij'} - \bar{x}_{\cdot j'})$$

covariance entre X_j et $X_{j'}$

On note $s_j = \sqrt{v_{jj}}$: écart-type associé à x_j

$$D_{1/s} = \begin{pmatrix} 1/s_1 & & 0 \\ & \ddots & \\ 0 & & 1/s_p \end{pmatrix} \quad D_{1/s^2} = \begin{pmatrix} 1/v_{11} & & 0 \\ & \ddots & \\ 0 & & 1/v_{pp} \end{pmatrix}$$

(p x p) (p x p)

La matrice Z centrée & réduite associée à X est :

$$Z = Y D_{1/s} \quad \boxed{z_{ij} = \frac{x_{ij} - \bar{x}_{.j}}{s_j}}$$

(N x p)

On note $R_{(p \times p)}$ la matrice des coefficients de corrélation linéaire entre les p variables :

$$R = D_{1/s} V D_{1/s} = Z' D Z, \quad \boxed{r_{jj'} = \sum_{i=1}^N p_i z_{ij} z_{ij'}}$$

4 Mesures de Similarité / Dissimilarité / Distance

4.1 Indice de similarité : un indice de similarité est une application qui vérifie les trois propriétés :

- 1) s est une application $\mathcal{O} \times \mathcal{O} \rightarrow \mathbb{R}^+$
- 2) s est symétrique : $\forall (w_i, w_j) \in \mathcal{O} \times \mathcal{O} \quad s(w_i, w_j) = s(w_j, w_i)$
- 3) $\forall (w_i, w_j) \in \mathcal{O} \times \mathcal{O}$ avec $w_i \neq w_j$
 $s(w_i, w_i) = s(w_j, w_j) > s(w_i, w_j)$

4.2 Indice de dissimilarité : un indice de dissimilarité est une application qui satisfait aux conditions 1) et 2) d'un indice de similarité et à :

$$s'(w_i, w_i) = 0 \quad \forall w_i \in \mathcal{O}$$

4.3 Une distance

une distance est un indice de dissimilarité qui vérifie en plus les deux propriétés suivantes :

- 1) $s'(w_i, w_j) = 0 \Leftrightarrow w_i = w_j$
- 2) $s'(w_i, w_j) \leq s'(w_i, w_k) + s'(w_k, w_j) \quad \forall w_i, w_k, w_j \in \mathcal{O}$
 (Inégalité triangulaire)

Remarque : un indice de dissimilarité vérifiant uniquement la propriété 1) est un indice de distance. Uniquement la propriété 2) est un "écart".

4.4 Une ultramétrie

un indice de dissimilarité qui vérifie en plus la propriété 2) est un indice de distance. Uniquement la propriété 1) est un "écart".

- 1) $s'(w_i, w_j) = 0 \Leftrightarrow w_i = w_j$
- 2) $s'(w_i, w_j) \leq \max(s'(w_i, w_k), s'(w_k, w_j))$
 $\forall (w_i, w_k, w_j) \in \mathcal{O}^3$

5 Mesures de proximité. Standard. entre individus

5.1 Cas de variables quantitatives.

$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$ donnant la description de w_i

5.1.1 distance de Minkowski

$$d(w_i, w_{i'}) = \left(\sum_{j=1}^p |x_{ij} - x_{i'j}|^r \right)^{1/r} \quad r \geq 1$$

$r=1$ (distance Manhattan, taxicab, city block)

$$d(w_i, w_{i'}) = \sum_{j=1}^p |x_{ij} - x_{i'j}|$$

$r=2$ (distance euclidienne)

$$d(w_i, w_{i'}) = \sqrt{\sum_{j=1}^p (x_{ij} - x_{i'j})^2}$$

$r \rightarrow \infty$ $d(w_i, w_{i'}) = \max_{1 \leq j \leq p} (|x_{ij} - x_{i'j}|)$ (Chebychev)

5.1.2 Distance de Mahalanobis

distance souvent utilisée en analyse discriminante

$$d^2(w_i, w_{i'}) = (x_i - x_{i'})' V^{-1} (x_i - x_{i'})$$

5.1.3 Distance de χ^2 (principe d'équivalence distributionnelle)

souvent utilisée sur des Tableaux de fréquence (Analyse factorielle des correspondances)

$$d^2(w_i, w_{i'}) = \sum_{j=1}^p \frac{1}{x_i} \left(\frac{x_{ij}}{x_i} - \frac{x_{i'j}}{x_{i'}} \right)^2$$

x_{ij} : nb d'occurrence de la modalité x_j pour w_i

5.2 Cas de variables binaires

| | | | | |
|----------|-----------|-----|-----------|----------------|
| | x_1 | ... | x_p | |
| w_i | x_{i1} | ... | x_{ip} | $x_{ij} = 0/1$ |
| $w_{i'}$ | $x_{i'1}$ | ... | $x_{i'p}$ | |

| | | |
|-------------------------|---|---|
| $w_i \backslash w_{i'}$ | 1 | 0 |
| 1 | a | b |
| 0 | c | d |

- a : nb de fois où $x_{ij} = x_{i'j} = 1$
- b : " " $x_{ij} = 0$ et $x_{i'j} = 1$
- c : " " $x_{ij} = 1$ et $x_{i'j} = 0$
- d : " " $x_{ij} = 0 = x_{i'j}$

Mesures de ressemblances usuelles

- Russel et Rao $D_1(w_i, w_{i'}) = \frac{a}{a+b+c+d} \quad (p = a+b+c+d)$
- Jaccard & Needham $D_2(w_i, w_{i'}) = \frac{a}{a+b+c}$
- Dice $D_3(w_i, w_{i'}) = \frac{2a}{2a+b+c}$
- Sokal & Sneath $D_4(w_i, w_{i'}) = \frac{a}{a+2(b+c)}$
- Sokal & Michener $D_5(w_i, w_{i'}) = \frac{a+d}{a+b+c+d}$
- Kulzinsky $D_6(w_i, w_{i'}) = \frac{a}{b+c}$
- Roger & Tanimoto $D_7(w_i, w_{i'}) = \frac{a+d}{a+d+2(b+c)}$
- Yule $D_8(w_i, w_{i'}) = \frac{ad-bc}{ad+bc}$

$D_5, D_7, (D_8+1), (D_9+1)$ sont des mesures de similarités

5.3 Cas de variables qualitatives nominales

On transforme les variables qualitatives en variables binaires par codage disjonctif complet, puis applique les mesures de ressemblances vu2 en 5.2

exp de codage disjonctif complet

| | Contenu | Forme |
|----------|---------|------------|
| w_i | rouge | Ellipsoïde |
| $w_{i'}$ | jaune | circulaire |

Contenu $\in \{ \text{rouge, jaune, bleue} \}$
 Forme $\in \{ \text{Ellipsoïde, circulaire, } \}$

↓ codage

| | Rouge | jaune | bleue | Ellipsoïde | circulaire |
|----------|-------|-------|-------|------------|------------|
| w_i | 1 | 0 | 0 | 1 | 0 |
| $w_{i'}$ | 0 | 1 | 0 | 0 | 1 |

5.4 Cas des variables qualitatives ordinales

Idem cas des variables qualitatives nominales.

6) Mesures de proximité entre variables

6.1 Cas de deux variables quantitatives

$$r_{ij} = d^2(x_j, x_{j'}) = \sum_{i=1}^N p_i \left(\frac{x_{ij} - \bar{x}_j}{s_j} \right) \left(\frac{x_{ij'} - \bar{x}_{j'}}{s_{j'}} \right)$$

coefficient de corrélation linéaire entre $x_j, x_{j'}$

6.2 Cas de deux variables qualitatives nominales

soit X_1, X_2 deux variables nominales à q, r modalités

$$\begin{matrix}
 & \begin{matrix} x_1 & 1 & 2 & \dots & r \end{matrix} \\
 \begin{matrix} x_1 \\ 1 \\ 2 \\ \vdots \\ q \end{matrix} & \left(\begin{matrix} \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{matrix} \right)
 \end{matrix}$$

n_{ij} = nb d'individus ayant la modalité i de X_1 et j de X_2

$n_{.j} = \sum_{i=1}^q n_{ij}$ $n_{i.} = \sum_{j=1}^r n_{ij}$

6.3.2 Coefficients

On associe à chaque variable ordinale X_i une nouvelle var $r_i : \mathcal{O} \rightarrow \{1 \dots r\}$ $n = \text{Card}(\mathcal{O})$

En présence d'ex-aequo on utilise le codage par rang moyen.

exemple:

| | x_j | $x_{j'}$ |
|-------|-------|----------|
| w_1 | a | f |
| w_2 | b | e |
| w_3 | b | e |
| w_4 | d | f |

avec $a < b < c$ et $e < f < g$

| | x_j | $x_{j'}$ |
|-------|-------|----------|
| w_1 | 1 | 2 |
| w_2 | 2 | 1 |
| w_3 | 3 | 4 |
| w_4 | 4 | 3 |

→

| | x_j | $x_{j'}$ |
|-------|-------|----------|
| w_1 | 1 | 2,5 |
| w_2 | 2,5 | 1 |
| w_3 | 2,5 | 4 |
| w_4 | 4 | 3 |

↑

codage par rang moyen

On note r_{ij} le rang de w_i pour la variable x_j

$$\bar{r}_{.j} = \frac{1}{n} \sum_{i=1}^n r_{ij}$$

Le coefficient de corrélation des rangs de Spearman R_s :

$$R_s = \frac{\sum_{i=1}^n (r_{ij} - \bar{r}_{.j}) (r_{ij'} - \bar{r}_{.j'})}{\sqrt{\sum_{i=1}^n (r_{ij} - \bar{r}_{.j})^2 \sum_{i=1}^n (r_{ij'} - \bar{r}_{.j'})^2}}$$

La mesure standard entre deux variables nominales est

la mesure du χ^2 ,

$$\chi^2 = D_A(X_1, X_2) = n_{..} \sum_{i=1}^q \sum_{j=1}^q \left(n_{ij} - \frac{n_{i.} \cdot n_{.j}}{n_{..}} \right)^2 \times \frac{1}{n_{ij}}$$

6.3 Cas de deux variables qualitatives ordinales.

La distance entre deux variables qualitatives ordinales revient à mesurer la distance entre ordres. \exists trois mesures :

- les coefficients de corrélation des rangs de Kendall
- " " " de Spearman
- " " " de Guttman

6.3.1 Coefficients de corrélation des rangs de Kendall

Soit X_1, X_2 deux variables ordinales. On associe à chaque variable une nouvelle variable Y_j définie comme suit

$$Y_j : \mathcal{O} \times \mathcal{O} \longrightarrow \{-1, 0, 1\}$$

$$\begin{cases} Y_j(w_i, w_{i'}) = -1 & \text{si } x_{i1} < x_{i2} \\ Y_j(w_i, w_{i'}) = 0 & \text{si } x_{i1} = x_{i2} \\ Y_j(w_i, w_{i'}) = +1 & \text{si } x_{i1} > x_{i2} \end{cases}$$

le rang de Kendall (Taux de Kendall) τ_{X_1, X_2} :

$$\tau_{X_1, X_2} = \text{cor}(Y_1, Y_2)$$

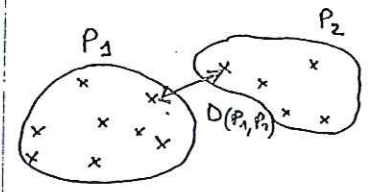
Remarque Y_j porte sur $\frac{n(n-1)}{2}$ couples d'individus.

7) Mesures de proximité entre groupes

soit P_1, P_2 deux populations d'individus de \mathcal{O} .
il existe trois mesures classiques de ressemblance entre P_1, P_2

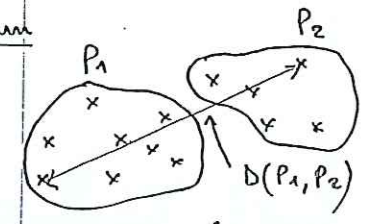
7.1 Distance du lieu minimum

$$D(P_1, P_2) = \min_{\substack{w_i \in P_1 \\ w_{i'} \in P_2}} d(w_i, w_{i'})$$



7.2 Distance du lieu maximum

$$D(P_1, P_2) = \max_{\substack{w_i \in P_1 \\ w_{i'} \in P_2}} d(w_i, w_{i'})$$



7.3 Distance des moyaux

Un sous ensemble N_j de P_j constitue un moyau de P_j à $n_i = \text{Card}(N_j)$ individus ssi il minimise l'expression :

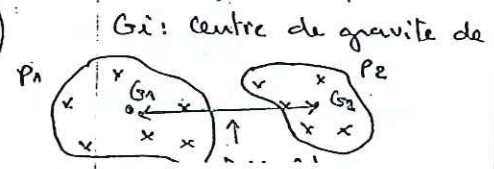
$$\sum_{w \in N_j} \sum_{w_i \in P_j} d(w, w_i)$$

La distance des moyaux de P_1, P_2 :

$$D(P_1, P_2) = \sum_{\substack{w_i \in N_1 \\ w_{i'} \in N_2}} d(w_i, w_{i'})$$

Cas particuliers :
- $\text{Card}(N_1) = \text{Card}(N_2)$
- espace d'observation est 1 espace vectoriel muni d'une distance Quadratique

$$D(P_1, P_2) = d(G_1, G_2)$$



On se situe dans le cas où l'espace d'observation est un espace vectoriel muni d'une distance quadratique.

le barycentre de P_j : $G(P_j) = \frac{1}{P(P_j)} \sum_{w_i \in P_j} p(w_i) x_i$

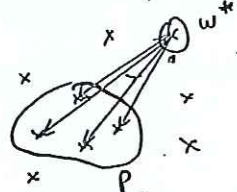
↑
vecteur de description de w_i

$p(w_i)$: poids de l'individu w_i

$$P(P_j) = \sum_{w_i \in P_j} p(w_i)$$

l'inertie de P_j par rapport à w^* de Ω est:

$$I_{w^*}(P_j) = \sum_{w_i \in P_j} p(w_i) d(w^*, w_i)$$



Cas particulier où $w^* = G(P_j)$, l'inertie de P_j est alors:

$$I(P_j) = \sum_{w_i \in P_j} p(w_i) d(G(P_j), w_i)$$

d'après le théorème de Huygens⁽¹⁾, l'inertie $I_{w^*}(P_j)$ est minimale pour $w^* = G(P_j)$.

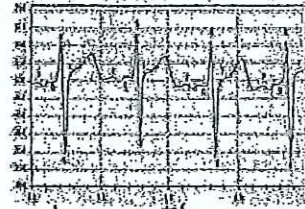
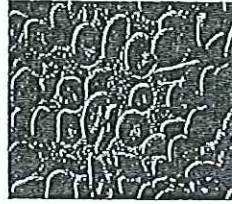
$$(1) \quad I_{w^*}(P_j) = I(P_j) + P(P_j) \cdot d(G(P_j), w^*)^2$$

La mesure de proximité fondée sur l'inertie proposée par Ward (1963)

$$d(P_1, P_2) = I(P_1 \cup P_2) - I(P_1) - I(P_2)$$

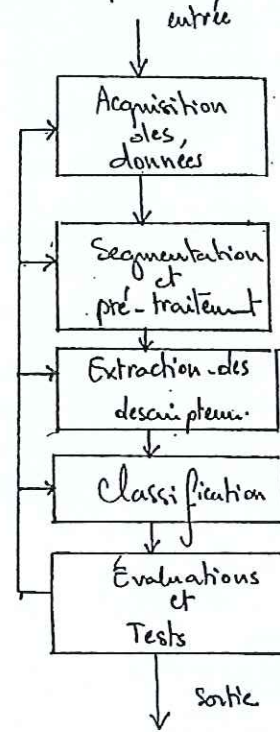
$$= \frac{P(P_1) P(P_2)}{P(P_1) + P(P_2)} d(G(P_1), G(P_2))^2$$

1) Type de données à analyser



| Prétraitement par paramètre | N | n | % | DS 551C |
|-------------------------------------|-------|-----|------|-------------|
| Type machines | 3 575 | 845 | 24.0 | [22,8-25,2] |
| Type propriétaires par le field | 3 516 | 572 | 16.3 | [14,1-16,1] |
| Objets machines | 3 533 | 372 | 10.6 | [2,7-1,5] |
| Système de données par le field | 3 530 | 340 | 10.2 | [1,1-1,3] |
| Système de données par le field | 3 530 | 202 | 5.7 | [10,6-5] |
| Dispositifs d'écriture | 3 525 | 473 | 13.4 | [12,4-14,4] |
| Dispositifs d'écriture par le field | 3 525 | 160 | 4.5 | [4,5-5,7] |
| Autres dispositifs d'écriture | 3 511 | 213 | 6.1 | [5,4-6,8] |

2) Les composantes d'un système de reconnaissance vocale Forme.



- choix de la technique de segmentation
- choix du modèle de pré-traitement de données
- choix des descripteurs
- choix du modèle de classification
- conception du classifieur par apprentissage
- Evaluation du classifieur.

2.1 Acquisition des données

- Constitution d'un échantillon d'apprentissage et d'un échantillon de Test exhaustifs

- Segmentation
 - Extraire l'objet à analyser de son environnement (isoler les cellules; Extraire les sous-séquences ORS, ...)
- Pré-traitement
 - Élimination du bruit dû aux instruments de collecte (bruitage d'1 image, translation d'1 signal, déviation d'1 capteur)

2.3 Sélection des descripteurs caractérisant les objets à analyser

L'objectif est de choisir des descripteurs dont les mesures sont proches pour des objets d'1 même classe et loins pour des objets de \neq classes

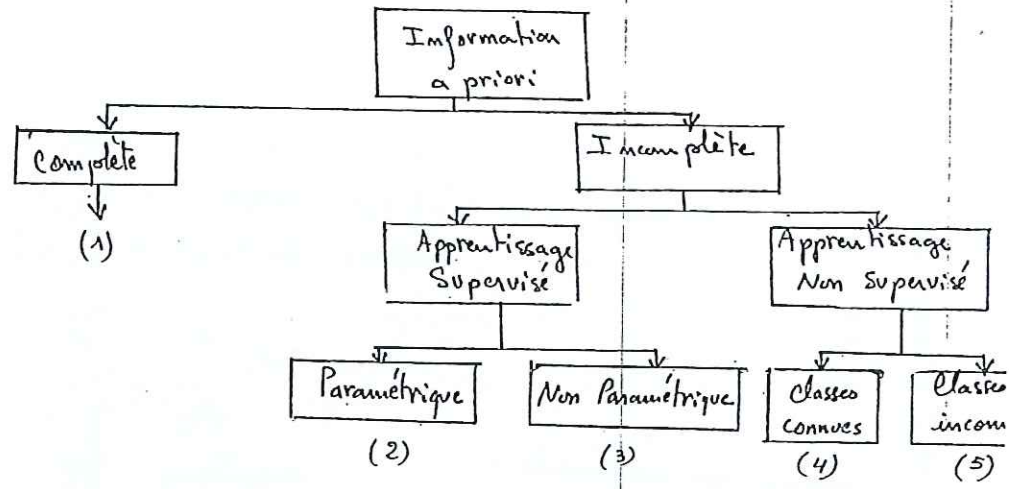
critères de choix:

- Invariance en translation / rotation
(Exp: - position absolue d'1 cellule est non significative
- l'instant où début d'1 cycle ...)
- Invariance par rapport à l'échelle
(Exp: - Échelle d'acquisition de l'image)
- Invariance aux distorsions des projections
(Exp: - distance entre l'objectif et la caméra change)

2.4 Choix du modèle de classification.

- un ensemble (échantillon d'apprentissage) d'objets décrits dans un espace à K dimensions (K descripteurs)
- l'objectif consiste à déterminer les frontières séparatrices des différentes classes.

3) Les différentes approches classificatoires en KT Statistique



| | Connu a priori | Inconnu | à Estimer | Méthodes Standards |
|-----|---|---|-------------------------------------|--|
| (1) | $C(x)$: classe de x $P(C_j)$: proba de C_j $P(x/C_j)$ | $P(C_j/x)$ | $P(C_j/x)$ | • Théorie de la décision bayésienne |
| (2) | $C(x)$ $P(C_j)$ $P(x/C_j) \sim L(\theta)$ $L(\theta)$: loi de la densité de proba de $P(x/C_j)$ | θ $P(C_j/x)$ | θ $P(C_j/x)$ | • Maximum de vraisemblance • Estimation Bayésienne |
| (3) | $C(x)$ $P(C_j)$ | • $\int dp$ de $P(x/C_j)$ • $P(C_j/x)$ • $P(x/C_j)$ | $P(C_j/x)$ | • Estimation $\int dp$ - fenêtre de Parzen - KPPV |
| (4) | Nb classe $P(C_j)$ $P(x/C_j) \sim L(\theta)$ | $C(x)$ $P(C_j/x)$ | • θ $P(C_j/x)$ | Modèles de densité mixtes |
| (5) | | • Nb classe $P(C_j)$, $C(x)$ $P(x/C_j)$ | Extraire les classe C_j $R(x)$ | • classification - Hiérarchique - Partitionnement - K-MEANS |