

# Fouille de données

Gaëtan Bloch & Maxime Gaudin

février 2011

#### Table des matières

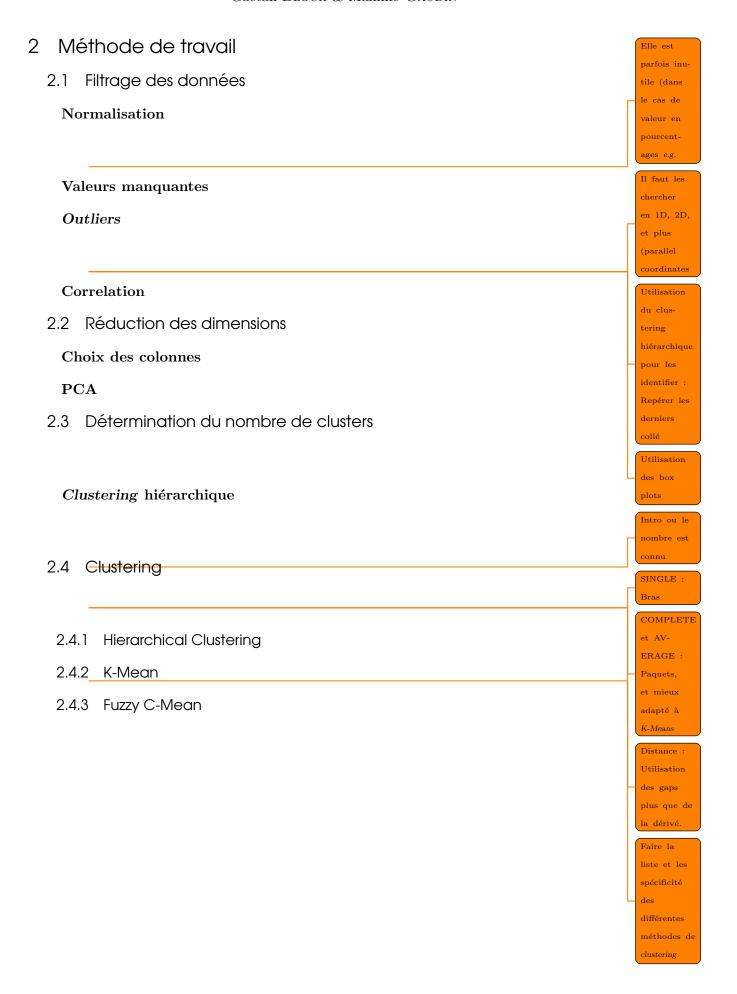
I	Mét	thode (	le travail	
2	2.1	Filtrag	e des données	
2	2.2	Réduc	tion des dimensions	
2	2.3	Déterr	nination du nombre de <i>clusters</i>	
2	2.4	Cluste	ring	
		2.4.1	Hierarchical Clustering	
		2.4.2	K-Mean	
		2.4.2	K-Mean	•
		2.4.2	Fuzzy C-Mean	
c	Zoot	2.4.3	Fuzzy C-Mean	
		2.4.3 teurs d	Fuzzy C-Mean	
	Sect	2.4.3  teurs d	Fuzzy C-Mean	
3		2.4.3  teurs d	Fuzzy C-Mean	
3	3.1	2.4.3  teurs d	Fuzzy C-Mean	
3	3.1	2.4.3  teurs d Sélecti Cluste	Fuzzy C-Mean	
3	3.1	2.4.3  teurs d Sélecti Cluste 3.2.1 3.2.2	Fuzzy C-Mean	

### 1 Description des données

Le jeu de données provient du site http://worldbank.org. Il concerne 210 pays et fournis des informations comme :

- Le pourcentage du Produit Intérieur Brute (PIB) généré par l'agriculture, l'industrie ou les services.
- La quantité de produits importés et exporté.
- Des informations démographiques et géographiques comme la densité et l'age de la population ou la surface du pays.
- Des informations sanitaires comme la quantité de personnes infectés par le Virus d'immunodéficience Humaine (VIH).
- Des centaines d'autres indicateurs.

Par ailleurs, ce jeu de donnée semble être un très bon entrainement pour nous exercer au datamining. En effet, en plus d'exhiber des données très concrètes pour nous, il est évident que des corrélations sont présentes. Enfin, il pourra être très intéressant de voir les corrélations à l'échelle du monde dans un premier temps, puis à l'échelle de l'Europe par exemple.



### 3 Secteurs d'activité des différents pays

nom de la colonne co

Ajouter le

Pour cette première étude de cas, nous avons choisi d'utiliser les données afin de caractériser l'activité d'un pays en fonction de :

- La taille du pays
- La population
- Le PIB
- Le À décrire (RNB)
- Les exports industriels

Notre démarche dans cette analyse est simple : Après avoir établi un *clustering* sur le pourcentage du PIB dû à l'agriculture, l'industrie et les services, nous avons créer un arbre de décision prenant pour critère de sélection les données citées ci-dessus.

Ajouter l'environ-

nement

pour une meilleur

intégration

#### 3.1 Sélection des données

Après avoir sélectionné les colonnes qui nous intéressent (i.e. Agriculture, value added (% of GDP), Industry, value added (% of GDP) et Service, value added (% of GDP)), nous obtenons la matrice suivante :

Le diagramme de est :

Enfin le diagramme parallel coordinates :

boite à chaussure??

Guinée : 95% de son

PIB est industriel

#### 3.2 Clustering

#### 3.2.1 Clustering hiérarchique

Le choix du type de *clustering* hiérarchique est, a *priori*délicat. Voici les résultats obtenus avec un *clustering* hiérarchique :

SINGLE

**COMPLETE** 

#### 3.2.2 Clustering K-Means

Vérification de la stabilité du clustering

Inclure le tableau de vérification de l'entropie

#### 3.3 Interprétation

#### 3.3.1 Arbre de décision

4 Création d'un module de vérification du K-Mean

## Liste des accronymes



Produit Intérieur Brute [Page(s) 1, 3]



À décrire [Page(s) 3]



Virus d'immuno déficience Humaine [Page(s) 1]

# Fin

