

Analyse de données de films

Hakim CHEKIROU

Introduction:

Pour analyser et extraire des informations pertinentes des données fournies, nous prétraitons d'abord les données. Pour décider quelles informations utiliser, nous analysons statistiquement les données avant d'appliquer les algorithmes d'apprentissage.

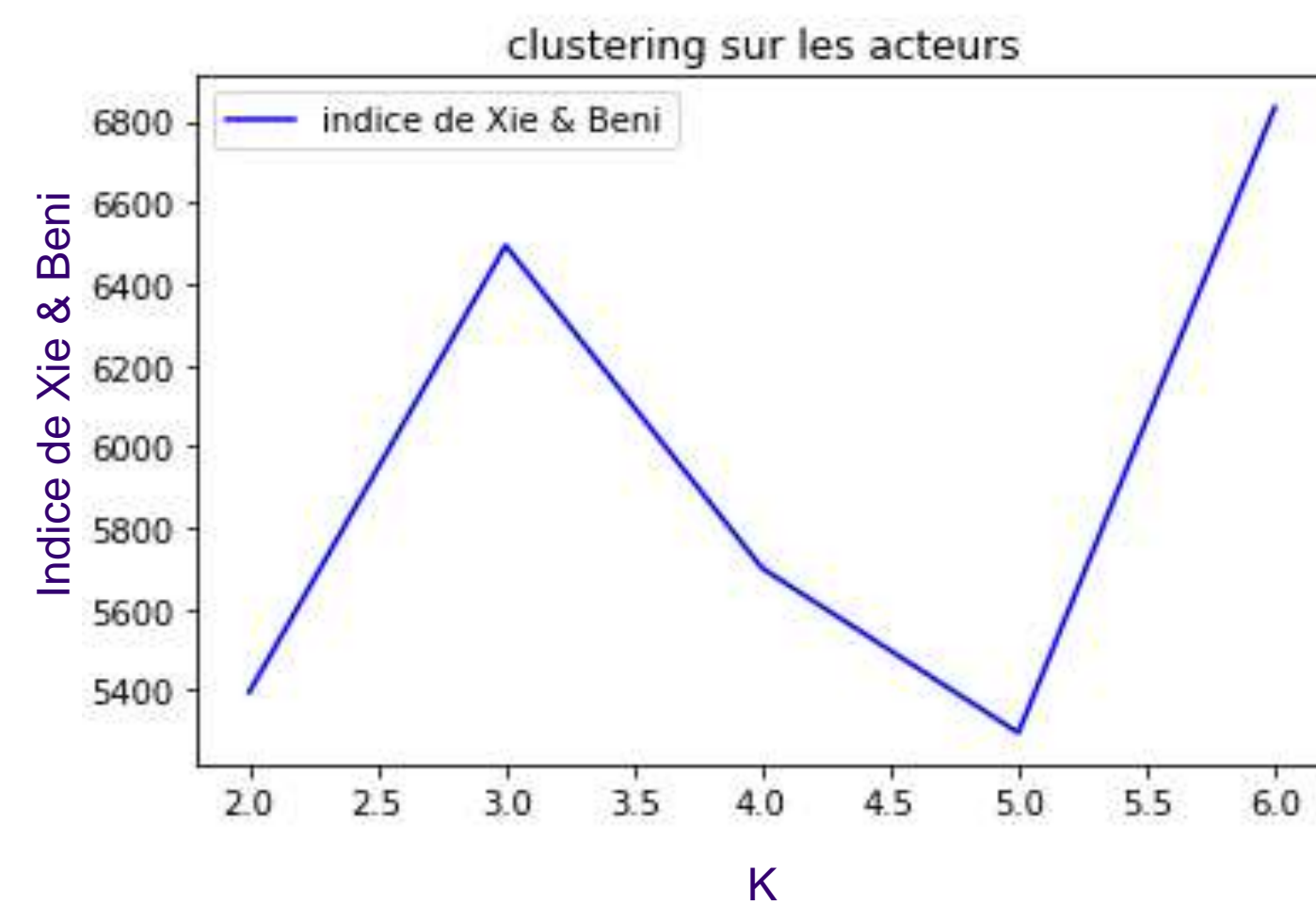
PRE-TRAITEMENT:

- > Réduction de la taille de Ratings
- > Suppression des attributs non significatifs dans toutes les tables.
- > Suppression des valeurs nulles.
- > Rajout de nouvelles colonnes à partir des données brutes.
- > **167** Relier les différentes tables entre eux en utilisant la base Links.

Apprentissage non supervisé : K-moyennes

Nous voulons ici trouver des groupes d'acteurs en se basant sur le type de films dans lesquels ils jouent et leurs performances dans ceux-ci. Les rôles ont été groupés par acteur et de nouvelles colonnes ont été rajoutées comme la qualité moyenne du film par acteur. Du fait du temps de calcul important nous nous sommes limités à **10000 acteurs**

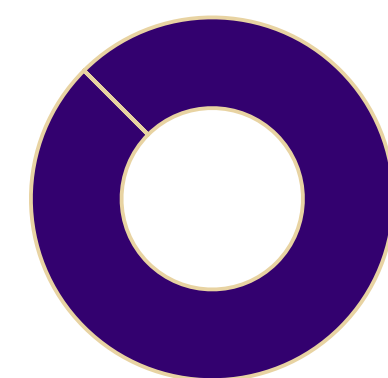
Recherche du meilleur nombre de clusters



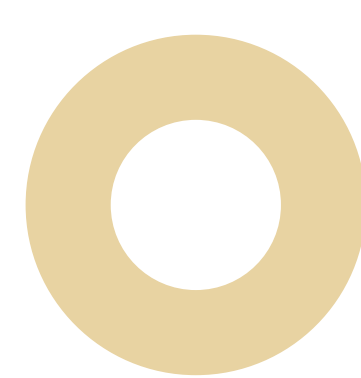
Selon l'indice Xie & Beni, les valeurs de K les plus intéressantes sont 2 et 5

Analyse des clusters trouvés

K = 2



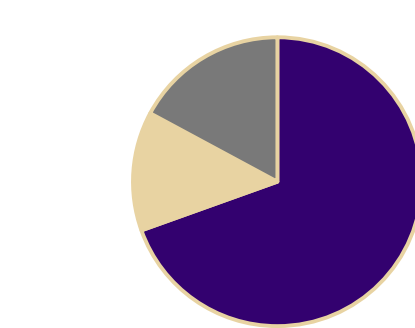
Cluster 1 composé d'hommes



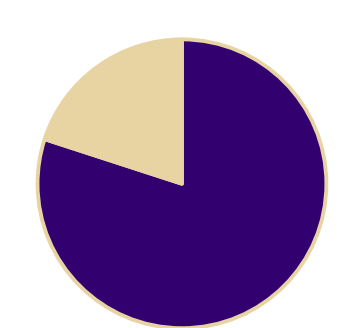
Cluster 2 composé uniquement de femmes.

K = 5

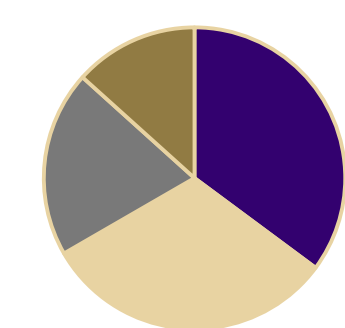
cluster 1 hommes



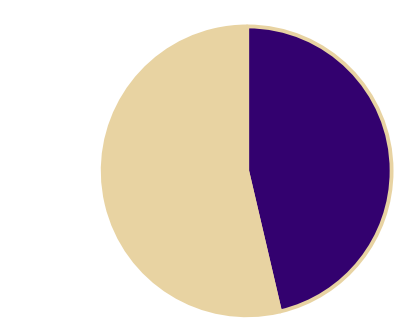
cluster 2 hommes



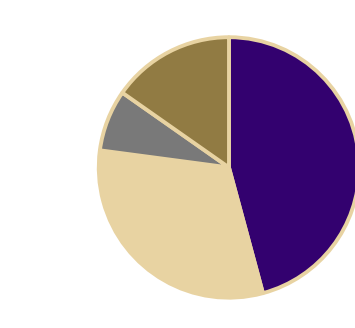
Cluster majorité hommes



cluster 3 mixte



cluster 4 femmes

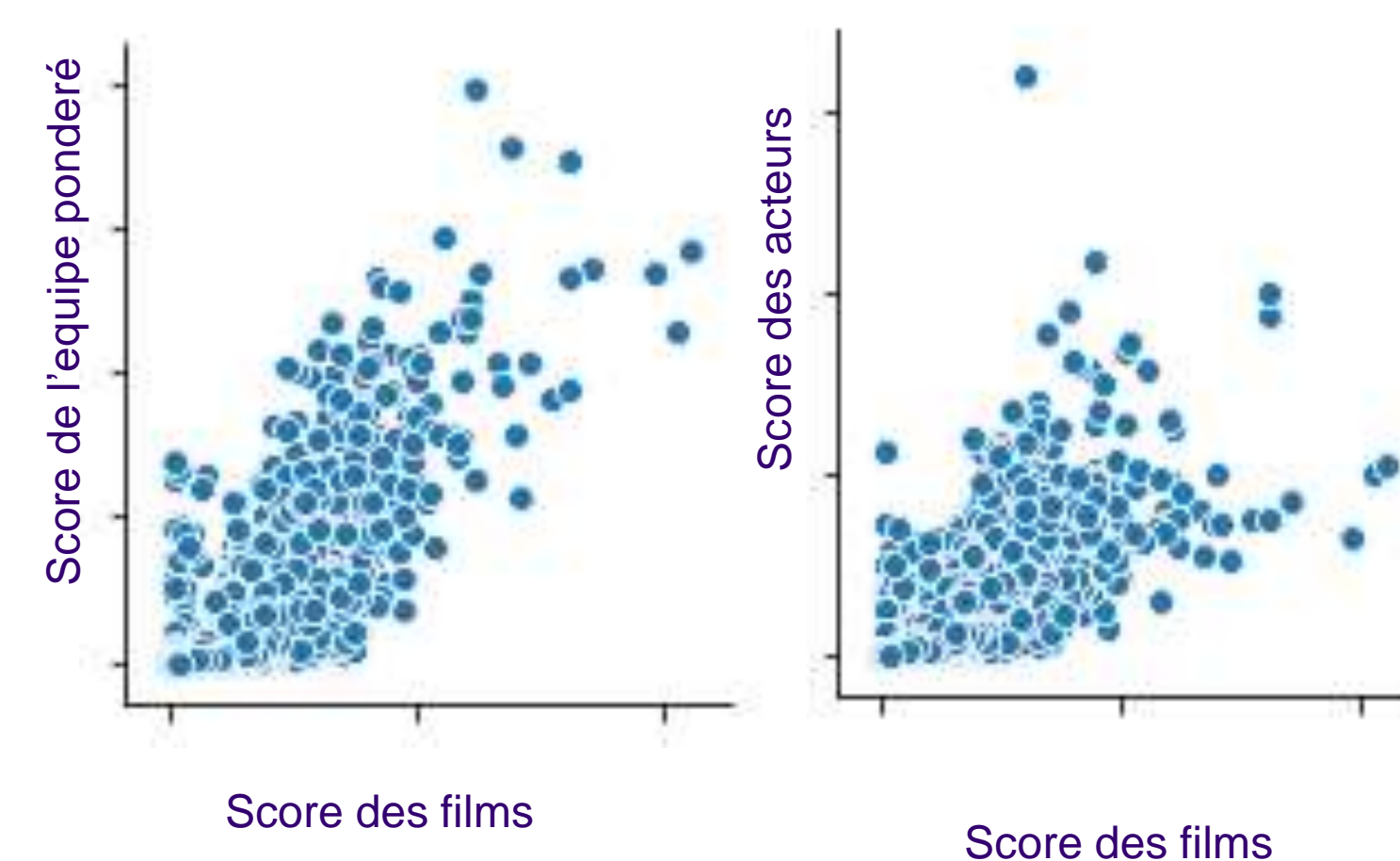


■ dramatique ■ Action ■ Thriller ■ comedies ■ romantique ■ action ■ aventure ■ sci-fi ■ fantastique
■ Horreur ■ Thriller ■ ■ ■ dramatique ■ comédie ■ crime ■ romantique

Regression

Dans cette section on cherche à prédire le score d'un film.

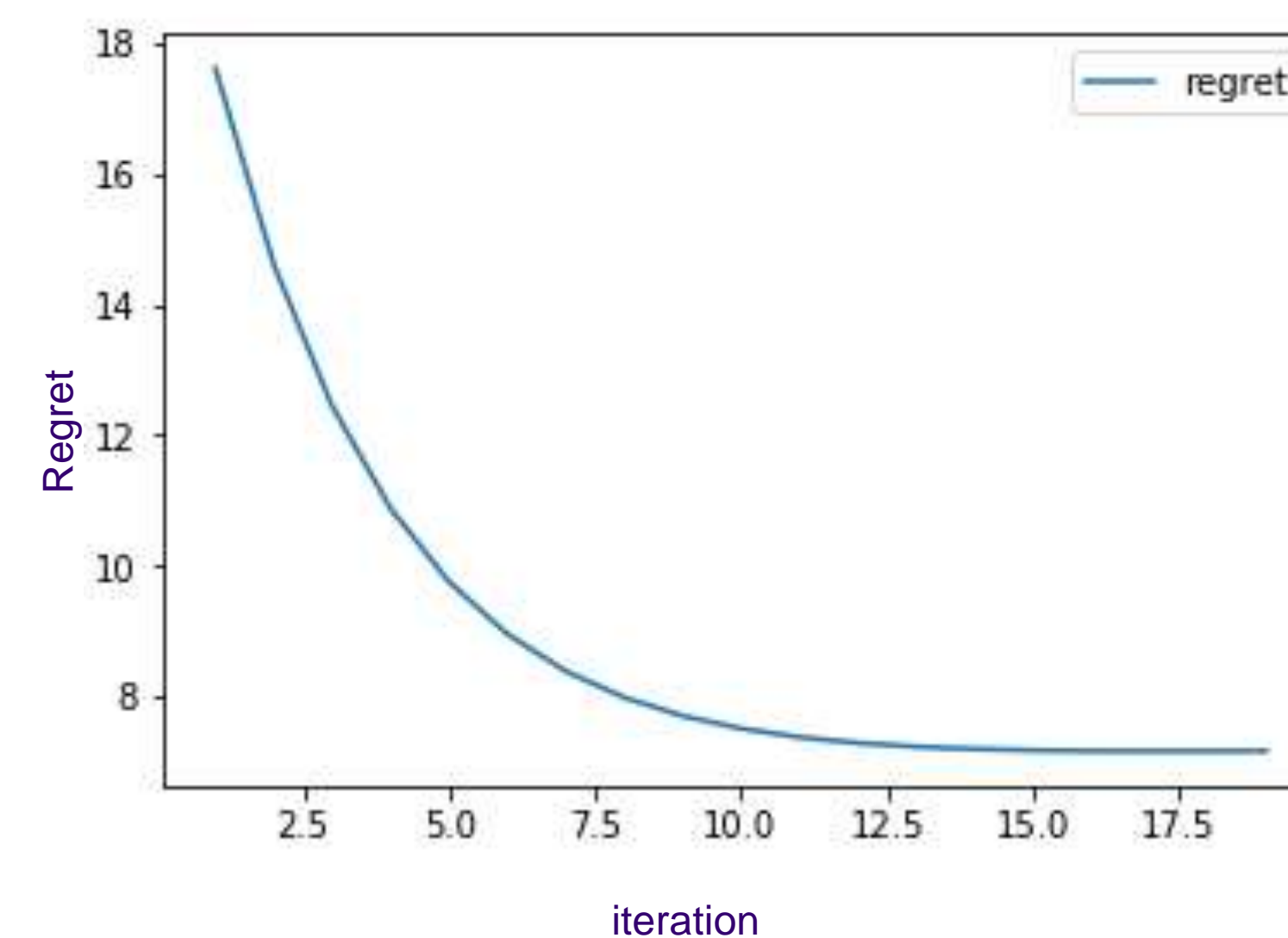
Correlation



Perceptron batch

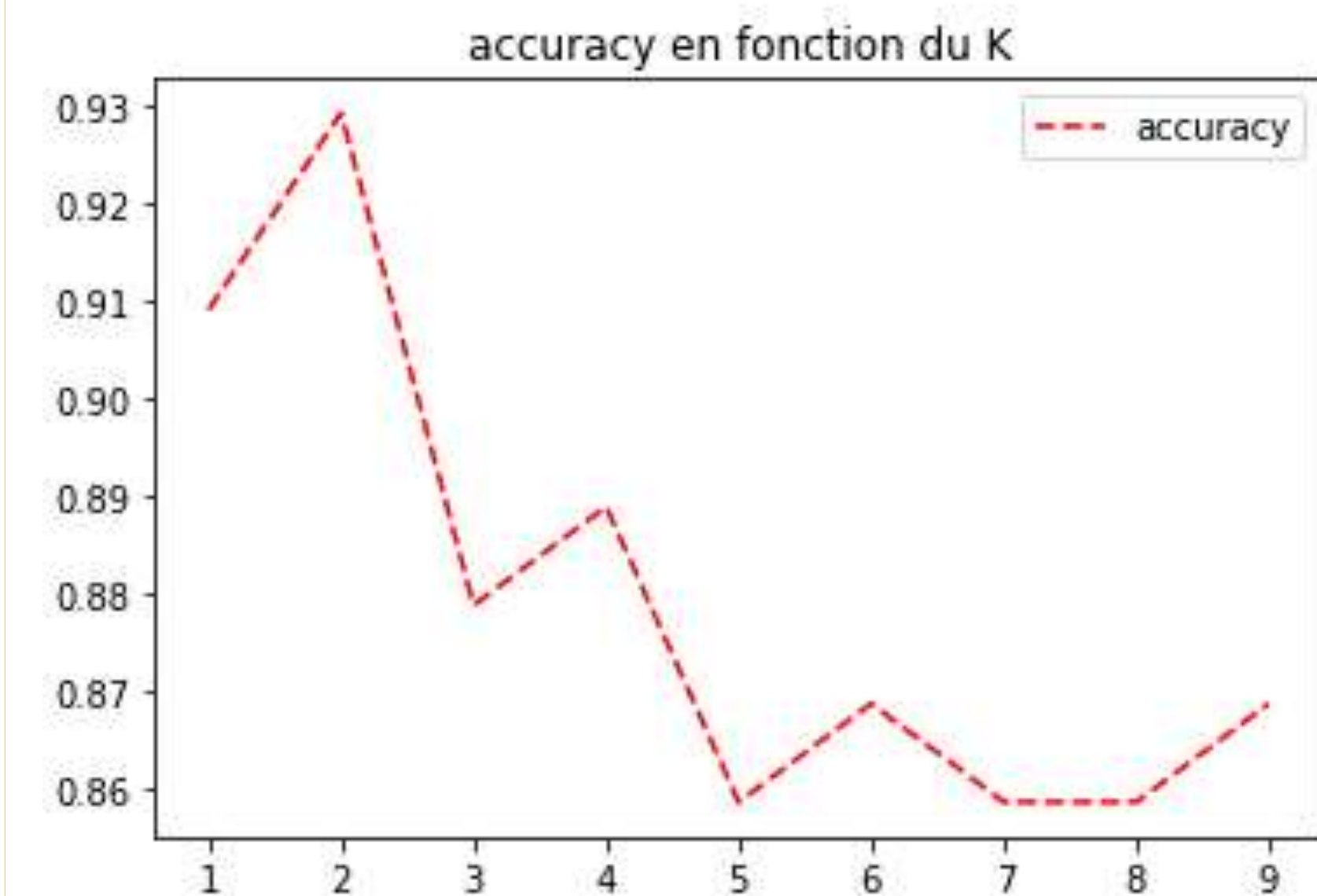
Nous avons entraîné l'algorithme sur 25000 films, avec un taux d'apprentissage de 0,002

Regret en fonction des iterations



Classification Supervisée

Notre but ici est de prédire le type du film.



Arbre de décision

Classifieur qui prédit si la personne est un réalisateur ou non. Epsilon = 0.01,

