

Réduction de dimensionnalité et clustering

Module 7

Objectifs

- utiliser les bonnes méthodes de réduction de dimensionnalité
- faire du clustering avec k-means
- projeter des variables catégorielles dans des embeddings

PCA

- réduire les dimensions tout en conservant la variance
- trouve une nouvelle base aux données et ordonne les axes par variance
- sélectionner les axes dans l'ordre pour conserver le plus de variance

<http://setosa.io/ev/principal-component-analysis/>

Réduire le nombre de dimensions en gardant $x\%$ de la variance. x souvent égal à 99.

t-SNE

- approche probabiliste
- map au mieux des distributions de distances dans l'espace haut-dimensionnel et dans l'espace bas-dimensionnel
- très utile pour visualiser des données en très hautes dimensions

- approche probabiliste
- map au mieux des distributions de distances dans l'espace haut-dimensionnel et dans l'espace bas-dimensionnel
- très utile pour visualiser des données en très hautes dimensions

- calcul d'une distribution P des similarités des objets en haute dimension
- calcul d'une distribution Q de similarités des objets en basse dimension
- minimisation de la KL divergence de Q par rapport à P

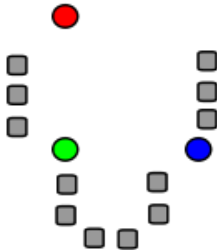
Attention, t-SNE est très sensible au choix de ses paramètres :

<https://distill.pub/2016/misread-tsne/>

k-means

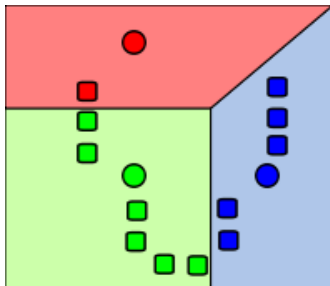
- algo incontournable en clustering
- sert aussi de base à d'autres algorithmes
- sert aussi de préprocessing

k-means — 1^{re} étape



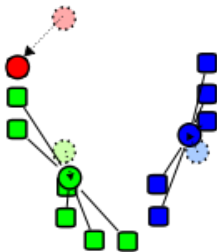
Décider de l'emplacement de k centroïds.

k-means — 2^e étape



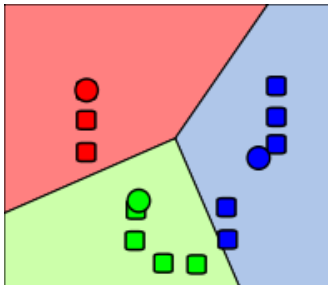
À chaque point, attribuer le cluster le plus proche.

k-means — 3^e étape



Ajuster les centroïdes pour qu'ils soient les barycentres des points.

k-means — 4^e étape



Répéter les étapes 2 et 3 jusqu'à convergence.

- k : nombre de clusters.
 - centroïdes de départ : principalement deux stratégies :
 - hasard
 - points existants
 - k-means++ :
 1. choisir un point dans les données avec une distribution uniforme
 2. calculer une distance des autres points à ce point
 3. utiliser une distribution biaisée par la distance
 4. recommencer 2 et 3 jusqu'à avoir k points
- Choisir par redémarrages multiples la meilleure option
- iterations : moins important

Embeddings

Projection des catégories de variables catégorielles dans un espace à n dimensions.

Pour les mots, $n \in \llbracket 50, 1000 \rrbracket$

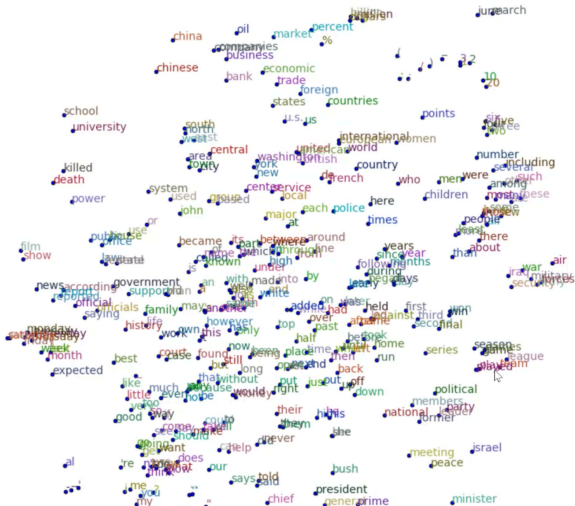
- utilisé en surcouche des one-hot encodings dans les réseaux
- initialisés aléatoirement
- entraînés comme le reste du réseau
- converge vers un espace riche

Problèmes de l'encodage one-hot

- très volumineux
- deux mots proches ne partageront rien de commun dans l'input

Pour pallier ces problèmes, projection dans un espace restreint.

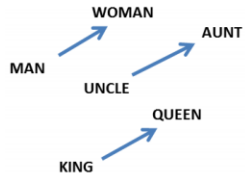
Embeddings



Espace appris riche et requêtable :

- trouver les symboles les plus proches dans l'espace
- résoudre des analogies :
woman - man = aunt - ?

→ Énorme aide pour le réseau



Applicable à toutes les variables catégorielles :

- utilisateurs
- pays
- types de carte SIM
- types d'appels
- ...

Pas limité aux mots !

Conclusion

- méthodes globale et « plus locale » de réduction de dimensionnalité
- clustering avec k-means, initialisé correctement et cross-validé sur son nombre de clusters
- embeddings pour gérer les variables catégorielles