

Algèbre linéaire et applications aux sciences des données

28 novembre 2024

Cette semaine

CM : Applications

TD : (demain) : Connexion CC + TD Espaces euclidiens

Vecteurs de \mathbb{R}^n , normes et produits scalaires

Cadre de travail : \mathbb{R}^n en tant qu'espace euclidien

Produit scalaire: $x^T y = \sum_{i=1}^n x_i y_i$ $\forall x \in \mathbb{R}^n$
 $\forall y \in \mathbb{R}^n$

Forme quadratique associée $x^T x = \sum_{i=1}^n x_i^2$

Norme associée : $\|x\| = \sqrt{x^T x} = \sqrt{\sum_{i=1}^n x_i^2}$

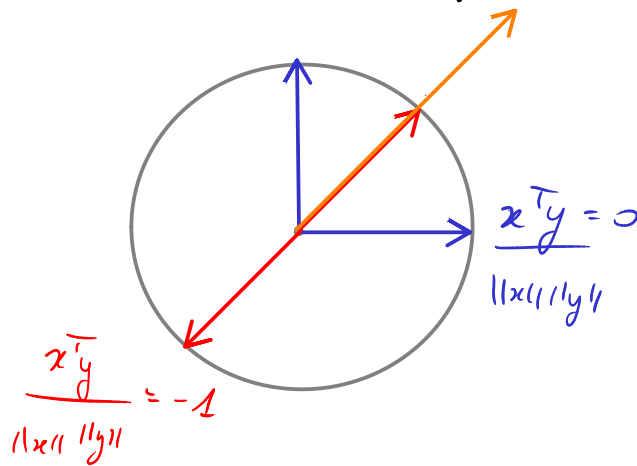
Etant donné deux vecteurs x et y de \mathbb{R}^n , on peut définir pour les comparer

a) $\|x - y\| = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \geq 0$ avec égalité ssi $x = y$

b) $\frac{x^T y}{\|x\| \|y\|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}} \in [-1, 1]$

(= 0 ssi x et y sont orthogonaux
= ± 1 ssi x et y sont colinéaires)

$\left(\frac{x}{\|x\|} \right)^T \left(\frac{y}{\|y\|} \right)$
↑
vecteurs x
'normalisés'



Def: La similarité cosinus ou la distance angulaire entre deux vecteurs normés x et y de \mathbb{R}^n est définie par

$$\frac{x^T y}{\|x\| \|y\|} \in [-1, 1]$$

On dit que la valeur $\arccos\left(\frac{x^T y}{\|x\| \|y\|}\right)$ est l'angle entre les vecteurs x et y . $\in [0, \pi]$ par convention

• Comparer des vecteurs en utilisant la distance euclidienne $\|x-y\|$

↳ Approche dite des plus proches voisins

(x_1, \dots, x_N) N points de \mathbb{R}^m

But: Regrouper en k groupes (clustering)

Processus: Partant de (z_1, \dots, z_k) k points de \mathbb{R}^m

• Pour chaque x_i , définir le plus proche voisin comme $z_{c_i} \in \{z_1, \dots, z_k\}$ tel que

$$\|x_i - z_{c_i}\| \leq \|x_i - z_j\| \quad \forall j=1..k$$

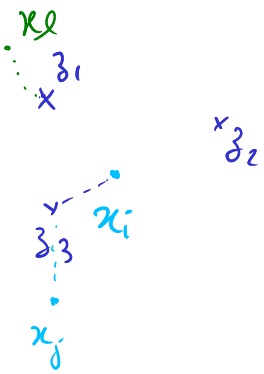
• Mettre à jour les z_j en fonction des x_i dont ils sont les plus proches voisins

$$\forall j=1..k, \quad z_j \leftarrow \frac{1}{n_j} \sum_{i \in I_j} x_i$$

$$\text{avec } I_j = \{i \in \{1, \dots, N\} \mid z_{c_i} = z_j\} \quad |I_j| = n_j$$

• Comparer des vecteurs en utilisant la distance angulaire

→ Technique de base en analyse de texte



Cadre général: Corpus de p documents (textes
- article Wikipédia
- syllabus de cours)
Dictionnaire de n mots

Pour chaque document, on définit le vecteur $x_i \in \mathbb{R}^m$ qui
contient les occurrences des mots du dictionnaire dans le
document i ($x_i \in \mathbb{N}^m$)

- Pour évaluer la similarité entre des documents,
on regarde $\frac{x_i^T x_j}{\|x_i\| \|x_j\|} \in [0, 1]$ car les x_i
sont à coefficients positifs

$$\frac{x_i^T x_j}{\|x_i\| \|x_j\|} \approx 1 \Rightarrow \text{Documents similaires}$$

$$\frac{x_i^T x_j}{\|x_i\| \|x_j\|} \approx 0 \Rightarrow \text{Documents sont "orthogonaux"} \\ (\text{traitent de sujets différents au sein du dictionnaire})$$

- On peut aussi effectuer des requêtes

Étant donné $q \in \mathbb{N}^m$ et $\tau \in (0, 1]$, les documents
du corpus les plus similaires à celui représenté par q
sont les vecteurs x_i tels que $\frac{x_i^T q}{\|x_i\| \|q\|} \geq \tau$.

Exemples: Soient 4 résumés de cours disponibles sur le site de Dauphine

Machine Learning

② Fondements du ML : (L3)

① Algèbre linéaire et applications aux Sciences des données (DL2)

③ Mathématiques pour les sciences des données (M1)

④ Programmation stochastique (M2 ~~DL2~~)

$p=4$

Dictionnaire : { "données", "algèbre linéaire", "statistiques", "optimisation" }

$m=4$

→ Pour construire les vecteurs x_i , on compte le nombre d'occurrences de chacun des mots du dictionnaire dans le résumé du cours

	① Algèbre linéaire et applications aux SD	② Fondements du ML	③ Maths pour les SD	④ Programmation stochastique
"données"	2	3	3	0
"algèbre linéaire"	2	2	1	0
"statistiques"	0	2	3	1
"optimisation"	0	1	4	3

$x_1 = \begin{bmatrix} 2 \\ 2 \\ 0 \\ 0 \end{bmatrix}$

\Rightarrow on obtient $x_1 = \begin{bmatrix} 2 \\ 2 \\ 0 \end{bmatrix}$, $x_2 = \begin{bmatrix} 3 \\ 2 \\ 2 \\ 1 \end{bmatrix}$, $x_3 = \begin{bmatrix} 3 \\ 1 \\ 3 \\ 4 \end{bmatrix}$, $x_4 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 3 \end{bmatrix}$

On calcule l'ensemble des $\frac{x_i^T x_j}{\|x_i\| \|x_j\|}$ pour $\{i, j\} \in \{1, 2, 3, 4\}^2$

\Rightarrow On présente en général cela sous la forme d'une matrice

$\frac{x_i^T x_j}{\ x_i\ \ x_j\ }$	x_1	x_2	x_3	x_4	$\frac{x_i^T x_j}{\ x_i\ \ x_j\ }$
x_1	1	0.83	0.48	0	
x_2	0.83	1	0.84	0.37	
x_3	0.48	0.84	1	0.80	
x_4	0	0.37	0.80	1	

NB: Par symétrie du produit scalaire, $\frac{x_i^T x_j}{\|x_i\| \|x_j\|} = \frac{x_j^T x_i}{\|x_j\| \|x_i\|}$

Interprétation de la 1^{re} ligne

$\rightarrow x_2$ est le vecteur le plus "semblable" (au sens de la similitude cosinus) à x_1 . On interprète cela comme

"Le (descriptif du) cours Fondements du ML est le plus semblable à celui du cours Algèbre linéaire et applications aux sciences des données" (C'est vrai!)

$\rightarrow x_4$ est orthogonal à x_1 . On en déduit que "Le (descriptif du) cours Programmation stochastique ne

Couvre pas les mêmes sujets que celui du cours
d'Algèbre linéaire " (c'est vrai!)

NB

$$\|x_1 - x_2\| = 2.45$$

$$\|x_1 - x_3\| = 5.2$$

$$\|x_1 - x_4\| = 4.24$$

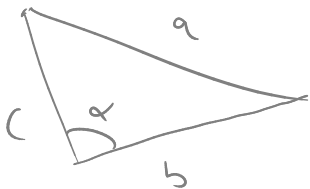
→ Au sein de la
norme, x_2 est le
vecteur le plus proche de
 x_1 mais x_4 est plus
proche que x_3
⇒ Classement moins
pertinent
⇒ Perte de l'information
d'orthogonalité

A retenir

La distance euclidienne et la similarité cosinus
ne permettent pas toujours les mêmes interprétations
et donc il peut arriver que l'une de ces mesures
soit plus utile que l'autre.

⇒ Il y a quand même entre les deux mesures

$$\forall (x, y) \in \mathbb{R}^m, \|x - y\|^2 = \|x\|^2 - 2\|x\|\|y\| \frac{x^T y}{\|x\|\|y\|} + \|y\|^2$$



$$a^2 = b^2 + c^2 - 2bc \cos(\alpha)$$

(Al-Kashi)

Remarques

• La similitude cosinus est invariante par rapport à la norme

$$\forall (x, y) \in \mathbb{R}^n, \quad \forall \lambda > 0,$$

$$\frac{(\lambda x)^T (\lambda y)}{\|\lambda x\| \|\lambda y\|} = \frac{x^T y}{\|x\| \|y\|}$$

$$(\text{alors que } \|\lambda x - \lambda y\| = \lambda \|x - y\|)$$

A venir : Pour $x \in \mathbb{R}^n$ fixe, comment projeter x sur un sous-espace de \mathbb{R}^n ?

Réponse : Trouver le y du sous-espace qui minimise $\|x - y\|$