

Science des données

Quentin Fortier

November 28, 2023

La **science des données** (*data science*) a pour objectif d'extraire de l'information à partir de données brutes.

La **science des données** (*data science*) a pour objectif d'extraire de l'information à partir de données brutes.

Exemples :

- Données sur des fleurs : longueur et largeur des pétales et des sépales.

La **science des données** (*data science*) a pour objectif d'extraire de l'information à partir de données brutes.

Exemples :

- Données sur des fleurs : longueur et largeur des pétales et des sépales.
- Données sur les clients d'une banque : âge, sexe, épargne, ...

Représentation des données

Pour pouvoir avoir une notion de distance entre deux données, **on représente chaque donnée comme un vecteur** de \mathbb{R}^p .

Représentation des données

Pour pouvoir avoir une notion de distance entre deux données, **on représente chaque donnée comme un vecteur** de \mathbb{R}^p .

Exemple : chaque donnée de fleur peut être représentée par un quadruplet de \mathbb{R}^4 correspondant à la longueur et largeur des pétales et des sépales.

Les composantes de ce vecteur sont appelées les **attributs**.

Parfois il est moins évident de représenter une donnée par un vecteur :

Parfois il est moins évident de représenter une donnée par un vecteur :

- Variable catégorielle (non numérique : genre, couleur, etc.) : on utilise souvent un vecteur avec un 1 et que des 0 (*one-hot vector*).

Parfois il est moins évident de représenter une donnée par un vecteur :

- Variable catégorielle (non numérique : genre, couleur, etc.) : on utilise souvent un vecteur avec un 1 et que des 0 (*one-hot vector*).
- Image : On passe d'une matrice de pixels avec n lignes, p colonnes à un vecteur de taille np .

Parfois il est moins évident de représenter une donnée par un vecteur :

- Variable catégorielle (non numérique : genre, couleur, etc.) : on utilise souvent un vecteur avec un 1 et que des 0 (*one-hot vector*).
- Image : On passe d'une matrice de pixels avec n lignes, p colonnes à un vecteur de taille np .
- Son : Transformée de Fourier discrète.

Représentation des données

On représente classiquement l'ensemble des données (donc de vecteurs de \mathbb{R}^p) par une matrice X dont chaque ligne est une donnée et chaque colonne est un attribut.

Représentation des données

On représente classiquement l'ensemble des données (donc de vecteurs de \mathbb{R}^p) par une matrice X dont chaque ligne est une donnée et chaque colonne est un attribut.

Python	Matrice	Données
<code>X[i]</code>	i ème ligne	i ème donnée
<code>len(X)</code>	nombre de lignes	nombre de données
<code>X[i][j]</code>	élément ligne i , colonne j	j ème attribut de la i ème donnée
<code>len(X[0])</code>	nombre de colonnes	nombre d'attributs

Distance

On a besoin de savoir si deux données sont « proches » l'une de l'autre. Pour cela, on utilise une **distance sur les données**, c'est-à-dire sur \mathbb{R}^p .

Distance

On a besoin de savoir si deux données sont « proches » l'une de l'autre. Pour cela, on utilise une **distance sur les données**, c'est-à-dire sur \mathbb{R}^p .

On utilise souvent la **distance euclidienne** :

$$d(x, y) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$$

Question

Écrire une fonction $d(x, y)$ renvoyant la distance euclidienne entre deux vecteurs x et y .

Distance

On a besoin de savoir si deux données sont « proches » l'une de l'autre. Pour cela, on utilise une **distance sur les données**, c'est-à-dire sur \mathbb{R}^p .

On utilise souvent la **distance euclidienne** :

$$d(x, y) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$$

Question

Écrire une fonction $d(x, y)$ renvoyant la distance euclidienne entre deux vecteurs x et y .

On peut utiliser d'autres distances, par exemple la distance de Manhattan :

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

Standardisation

Quand les attributs n'ont pas la même échelle (par exemple, l'argent d'un client d'une banque peut être beaucoup plus élevé que son âge), un attribut peut avoir beaucoup plus d'importance qu'un autre dans les calculs de distance.

Standardisation

Quand les attributs n'ont pas la même échelle (par exemple, l'argent d'un client d'une banque peut être beaucoup plus élevé que son âge), un attribut peut avoir beaucoup plus d'importance qu'un autre dans les calculs de distance.

Pour que les attributs aient la même importance, on peut **standardiser** (ou : **normaliser**) les données, c'est-à-dire les modifier pour avoir une moyenne de 0 et un écart-type de 1.

Standardisation

Quand les attributs n'ont pas la même échelle (par exemple, l'argent d'un client d'une banque peut être beaucoup plus élevé que son âge), un attribut peut avoir beaucoup plus d'importance qu'un autre dans les calculs de distance.

Pour que les attributs aient la même importance, on peut **standardiser** (ou : **normaliser**) les données, c'est-à-dire les modifier pour avoir une moyenne de 0 et un écart-type de 1.

La plupart des algorithmes de science des données fonctionnent mieux avec des données standardisées.

Standardisation

Si les données sont x_1, \dots, x_n , on calcule la moyenne μ et l'écart-type σ de chaque attribut, puis on remplace chaque x_i par $\frac{x_i - \mu}{\sigma}$.

Question

Écrire une fonction `standardiser(X)` qui renvoie la matrice obtenue en standardisant les données de la matrice `X`.