

## Cours 4 – LA METHODE DES K PLUS PROCHES VOISINS

### Présentation de la méthode

Méthode des K plus proches voisins (KPPV)

En Anglais, K Nearest Neighbours (KNN)

Méthode d'apprentissage supervisée

classification ou régression

Principe

Se baser sur des cas similaires au cas à résoudre pour faire la prédiction

### Description de la méthode

#### Echantillon d'apprentissage

variables explicatives  $X_j$

variable à expliquer  $Y$

$Y(\omega_i)$  est égal à  $y_i$  pour chaque individu  $\omega_i$

	$X_1$	...	$X_j$	...	$X_P$	$Y$
$\omega_1$	$X_1(\omega_1)$	...	$X_j(\omega_1)$	...	$X_P(\omega_1)$	$y_1$
$\omega_2$	$X_1(\omega_2)$	...	$X_j(\omega_2)$	...	$X_P(\omega_2)$	$y_2$
...	...	...	...	...	...	...
$\omega_i$	$X_1(\omega_i)$	...	$X_j(\omega_i)$	...	$X_P(\omega_i)$	$y_i$
...	...	...	...	...	...	...
$\omega_N$	$X_1(\omega_N)$	...	$X_j(\omega_N)$	...	$X_P(\omega_N)$	$y_N$

### Problème

Pour un nouvel individu  $\omega$

prédire ( trouver )  $y$  , c'est à dire  $Y(\omega)$ .

$\omega$	$X_1(\omega)$	...	$X_j(\omega)$	...	$X_P(\omega)$	?
----------	---------------	-----	---------------	-----	---------------	---

## Solution

- Trouver les K plus proches voisins de  $\omega$

c'est-à-dire les K individus qui lui ressemblent le plus

On peut prendre K égal à la racine carrée de N

=> Calcul des distances entre  $\omega$  et les individus de l'échantillon

- Ensuite calculer y en se basant sur les valeurs  $y_i$  de ces voisins.

## Distance entre deux individus

### Distance euclidienne

$$D(\omega, \omega_i) = \sqrt{\sum_{j=1}^P [d(X_j(\omega), X_j(\omega_i))]^2}$$

$d(X_j(\omega), X_j(\omega_i))$  distance par rapport à  $X_j$

Si  $X_j$  est de type quantitatif

$$d(X_j(\omega), X_j(\omega_i)) = X_j(\omega) - X_j(\omega_i)$$

Si  $X_j$  est de type qualitatif nominal

$$d(X_j(\omega), X_j(\omega_i)) = 0 \text{ si } X_j(\omega) \text{ est égal à } X_j(\omega_i) \\ = 1 \text{ sinon}$$

Si  $X_j$  est de type qualitatif ordinal

$$d(X_j(\omega), X_j(\omega_i)) = 0 \text{ si } X_j(\omega) \text{ est égal à } X_j(\omega_i) \\ = |X_j(\omega) - X_j(\omega_i)| / \text{Card}(D_j) \text{ sinon}$$

Par exemple si  $X_j$  est la variable Mention qui est ordinale

$$D_j = \{\text{Passable}, \text{ABien}, \text{Bien}, \text{TBien}\} = \{0, 1, 2, 3\}$$

Si  $X_j(\omega) = \text{Passable}$  et  $X_j(\omega_i) = \text{Passable}$

$$\Rightarrow d(X_j(\omega), X_j(\omega_i)) = 0$$

Si  $X_j(\omega) = \text{Passable}$  et  $X_j(\omega_i) = \text{Bien}$

$$\Rightarrow d(X_j(\omega), X_j(\omega_i)) = |0 - 2| / 4 = 0.5$$

### **Normalisation des variables quantitatives**

Objectif : éviter que la distance dépende principalement de la variable qui a des valeurs plus élevées que les autres.

#### **Normalisation min-max**

Transformer les  $X_j(\omega_i)$  en des valeurs  $\hat{X}_j(\omega_i)$  entre 0 et 1.

$$\hat{X}_j(\omega_i) = \frac{X_j(\omega_i) - \min(X_j)}{\max(X_j) - \min(X_j)}$$

$\min(X_j)$  et  $\max(X_j)$  : minimum et maximum des valeurs de  $X_j$

### **Cas de la classification**

#### **Vote majoritaire**

la classe  $y$  de  $\omega$  est égal à la classe majoritaire parmi les classes  $\{y_1, y_2, \dots, y_K\}$  de ses plus proches voisins.

#### **Vote pondéré**

vote majoritaire pondéré

le poids  $p_i$  de la classe  $y_i$  d'un voisin  $\omega_i$  est égal à  $1 / D(\omega, \omega_i)$

### **Cas de la régression**

#### **moyenne simple**

$y$  est des valeurs  $\{y_1, y_2, \dots, y_K\}$  de ses voisins.

#### **moyenne pondérée**

$y$  est la moyenne pondérée des valeurs  $\{y_1, y_2, \dots, y_K\}$  de ses voisins.

## Exemple d'application 1

Soit l'ensemble d'apprentissage suivant dans lequel age est une variable quantitative et income, credit\_rating, student sont des variables qualitatives considérées comme non ordonnées.

	<b>age</b>	<b>income</b>	<b>student</b>	<b>credit_rating</b>	<b>buys_computer</b>
$\omega_1$	30	high	no	fair	no
$\omega_2$	30	high	no	excellent	no
$\omega_3$	40	high	no	fair	yes
$\omega_4$	50	medium	no	fair	yes
$\omega_5$	50	low	yes	fair	no
$\omega_6$	50	low	yes	excellent	no
$\omega_7$	40	low	yes	excellent	yes
$\omega_8$	30	medium	no	fair	no
$\omega_9$	30	low	yes	fair	yes
$\omega_{10}$	50	medium	yes	fair	yes
$\omega_{11}$	30	medium	yes	excellent	yes
$\omega_{12}$	40	medium	no	excellent	yes
$\omega_{13}$	40	high	yes	fair	no
$\omega_{14}$	50	medium	no	excellent	no

En utilisant la méthode KNN, avec  $K = 5$ . prédire la valeur de la variable buys\_computer pour un nouvel exemple ayant les caractéristiques suivantes :

age = 35, income = medium, student = yes, credit\_rating = fair

## Solution

Pour le calcul des distances, on utilise la version normalisée de la variable age obtenue avec la normalisation Min-Max.

Il faut donc considérer le nouveau tableau de données suivant :

	age	age*	income	student	credit_rating	buys_computer	D( $\omega, \omega_i$ )	Poids $\pi_i$
$\omega_1$	30	0	high	no	fair	no	1.43	
$\omega_2$	30	0	high	no	excellent	no	1.75	
$\omega_3$	40	0.5	high	no	fair	yes	1.43	
$\omega_4$	50	1	medium	no	fair	yes	1.25	0.8
$\omega_5$	50	1	low	yes	fair	no	1.25	0.8
$\omega_6$	50	1	low	yes	excellent	no	1.60	
$\omega_7$	40	0.5	low	yes	excellent	yes	1.43	
$\omega_8$	30	0	medium	no	fair	no	1.75	
$\omega_9$	30	0	low	yes	fair	yes	1.43	
$\omega_{10}$	50	1	medium	yes	fair	yes	0.75	1.33
$\omega_{11}$	30	0	medium	yes	excellent	yes	1.03	0.97
$\omega_{12}$	40	0.5	medium	no	excellent	yes	1.43	
$\omega_{13}$	40	0.5	high	yes	fair	no	1.03	0.97
$\omega_{14}$	50	1	medium	no	excellent	no	1.60	

$\omega$	35	0.25	Medium	Yes	fair	?		
----------	----	------	--------	-----	------	---	--	--

L'âge normalisé du nouvel exemple à classer est égal à 0.25

$$\begin{aligned} D(\omega, \omega_1) &= ( (0.25 - 0)^2 + (d(\text{medium}, \text{high}))^2 + (d(\text{yes}, \text{no}))^2 \\ &\quad + (d(\text{fair}, \text{fair}))^2 )^{1/2} \\ &= ( 0.25^2 + 1^2 + 1^2 + 0^2 )^{1/2} \\ &= 1.43 \end{aligned}$$

Même calcul pour les autres distances  $D(\omega, \omega_2), \dots, D(\omega, \omega_{14})$

Vote simple :

yes :  $1+1+1 = 3$

no :  $1 + 1=2$

Donc buys\_computer = yes pour le nouvel exemple.

Vote pondéré :

yes :  $0.8+1.33+0.97 = 3.1$

no :  $0.8+0.97 = 1.77$

Donc buys\_computer = yes pour le nouvel exemple

## Exemple d'application 1

Soit l'ensemble d'apprentissage suivant dans lequel  
age est une variable quantitative,  
income est une variable qualitative ordonnée,  
credit\_rating, student, buys\_computer sont des variables qualitatives  
considérées comme non ordonnées.

	<b>income</b>	<b>student</b>	<b>credit rating</b>	<b>buys computer</b>	<b>age</b>
$\omega_1$	high	no	fair	no	30
$\omega_2$	high	no	excellent	no	30
$\omega_3$	high	no	fair	yes	40
$\omega_4$	medium	no	fair	yes	50
$\omega_5$	low	yes	fair	no	50
$\omega_6$	low	yes	excellent	no	50
$\omega_7$	low	yes	excellent	yes	40
$\omega_8$	medium	no	fair	no	30
$\omega_9$	low	yes	fair	yes	30
$\omega_{10}$	medium	yes	fair	yes	50
$\omega_{11}$	medium	yes	excellent	yes	30
$\omega_{12}$	medium	no	excellent	yes	40
$\omega_{13}$	high	yes	fair	no	40
$\omega_{14}$	medium	no	excellent	no	50

En utilisant la méthode KNN, avec  $K = 5$ . prédire la valeur de la variable  
age pour un nouvel exemple ayant les caractéristiques suivantes :

income = medium, student = yes, credit\_rating = fair,

buys\_computer = yes