

Cours 4 – LA METHODE NAIVE BAYES

Présentation

Méthode de classification supervisée

Calcul de probabilités conditionnelles pour la prédiction

Utilisation du théorème de Bayes

Hypothèse de l'indépendance des variables explicatives

=> Méthode qualifiée de naïve

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 ω_i

	X_1	...	X_j	...	X_P	Y
ω_1	$X_1(\omega_1)$...	$X_j(\omega_1)$...	$X_P(\omega_1)$	y_1
ω_2	$X_1(\omega_2)$...	$X_j(\omega_2)$...	$X_P(\omega_2)$	y_2
...
ω_i	$X_1(\omega_i)$...	$X_j(\omega_i)$...	$X_P(\omega_i)$	y_i
...
ω_N	$X_1(\omega_N)$...	$X_j(\omega_N)$...	$X_P(\omega_N)$	y_N

Problème

Pour un nouvel individu ω

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

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

On note

$$X = (X_1, \dots, X_j, \dots, X_P)$$

$$X(\omega) = (X_1(\omega), \dots, X_j(\omega), \dots, X_P(\omega))$$

On veut trouver y , c'est-à-dire $Y(\omega)$ étant donné $X(\omega)$

Solution

Ensemble C des classes distinctes

$$C = \{ y_1, y_2, \dots, y_L \}, \text{ avec } L \leq N$$

La classe y à trouver est la classe $y_k \in C$ qui a la plus grande probabilité d'être la classe de l'individu ω

$$y = \operatorname{argmax}_{y_k \in C} P(y_k | X(\omega))$$

Théorème de Bayes

$$P(y_k | X(\omega)) = \frac{P(X(\omega) | y_k) P(y_k)}{P(X(\omega))}$$

même dénominateur $P(X(\omega))$ quelle que soit la classe y_k

=> Comparer les numérateurs $P(X(\omega) | y_k) P(y_k)$

Calcul de $P(y_k)$

$P(y_k)$ est la probabilité a priori de la classe y_k

= proportion d'individus de la classe y_k

$$= N_k / N$$

avec

N_k le nombre d'individus dans la classe y_k

N le nombre d'individus

Calcul de $P(X(\omega) | y_k)$

Hypothèse : indépendance des variables X_j

$$\Rightarrow P(X(\omega)|y_k)=P(X_1(\omega)|y_k)\times P(X_2(\omega)|y_k) \times \dots \times P(X_P(\omega)|y_k)$$

Pour une variable qualitative X_j

$$P(X_j(\omega) | y_k)$$

= probabilité d'avoir la valeur $X_j(\omega)$ étant donnée la classe y_k

= proportion d'individus de la classe y_k pour lesquelles X_j a pour valeur $X_j(\omega)$

Pour X_j variable quantitative

X_j : variable aléatoire avec une distribution gaussienne

$P(X_j(\omega) | y_k)$ est la fonction de densité de probabilité

$$P(X_j(\omega)|y_k)=\frac{1}{\sqrt{2\pi}\sigma(X_j|y_k)}\exp\left(\frac{-(X_j(\omega)-\mu(X_j|y_k))^2}{2\sigma(X_j|y_k)^2}\right),$$

où :

$\mu(X_j|y_k)$ est la moyenne de X_j dans la classe y_k

$$\mu(X_j|y_k)=\frac{1}{n_k}\sum_{i=1}^{n_k} X_j(\omega_i)$$

$\sigma(X_j | y_k)$ est la déviation standard (écart-type) de X_j dans la classe y_k

$$\sigma(X_j|y_k)=\sqrt{\frac{1}{n_k-1}\sum_{i=1}^{n_k} (X_j(\omega_i)-\mu(X_j|y_k))^2}$$

Exemple d'application

Soit le tableau suivant regroupant des données sur la possibilité de jouer au tennis en fonction de la météo.

Outlook	Temperature	Humidity	Windy	Play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

Prédire la valeur de la variable Play pour un nouvel exemple ω , avec $X(\omega) = (\text{sunny}, 74, 66, \text{true})$ regroupant les valeurs des variables Outlook, Temperature, Humidity et Windy respectivement.

Solution

Comparer $P(\text{yes} | X(\omega))$ et $P(\text{no} | X(\omega))$

Ceci revient à comparer $P(X(\omega) | \text{yes}) \times P(\text{yes})$ et $P(X(\omega) | \text{no}) \times P(\text{no})$

$$P(\text{yes}) = 9/14$$

$$P(\text{no}) = 5/14$$

$$P(X(\omega)|\text{yes}) = P(\text{Outlook}=\text{sunny}, \text{Temperature}=74, \text{Humidity}=66, \text{Windy}=\text{true} | \text{yes})$$

$$\begin{aligned} &= P(\text{Outlook}=\text{sunny} | \text{yes}) \times P(\text{Temperature}=74 | \text{yes}) \\ &\quad \times P(\text{Humidity}=66 | \text{yes}) \times P(\text{Windy}=\text{true} | \text{yes}) \end{aligned}$$

$$P(\text{Outlook}=\text{sunny} | \text{yes}) = 2/9$$

$$\mu(\text{Temperature} | \text{yes}) = 73$$

$$\sigma(\text{Temperature} | \text{yes}) = 6.16$$

$$P(\text{Temperature}=74 | \text{yes}) = \frac{1}{\sqrt{2\pi} 6.16} \exp\left(\frac{-(74-73)^2}{2 \times 6.16^2}\right) = 0.064$$

$$\mu(\text{Humidity} | \text{yes}) = 79.1$$

$$\sigma(\text{Humidity} | \text{yes}) = 10.2$$

$$P(\text{Humidity}=66 | \text{yes}) = \frac{1}{\sqrt{2\pi} 10.2} \exp\left(\frac{-(66-79.1)^2}{2 \times 10.2^2}\right) = 0.045$$

$$P(\text{Windy}=\text{true} | \text{yes}) = 3/9$$

$$\text{Donc } P(X(\omega) | \text{yes}) \times P(\text{yes}) = (2/9) \times 0.064 \times 0.045 \times (3/9) \times (9/14) = 0.64$$

$$P(X(\omega) | \text{yes}) \times P(\text{yes}) = 0.64$$

Les mêmes calculs donnent pour la classe no :

$$P(X(\omega) | \text{no}) \times P(\text{no}) = 0.00027$$

$$P(X(\omega)|\text{yes}) \times P(\text{yes}) \text{ est supérieure à } P(X(\omega)|\text{no}) \times P(\text{no})$$

donc $P(\text{yes} | X(\omega))$ est supérieure à $P(\text{no} | X(\omega))$

donc Play = yes pour $X(\omega) = (\text{sunny}, 74, 66, \text{true})$