

Fondements du ML

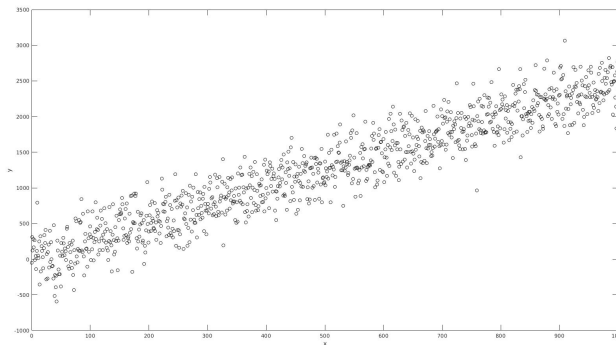
29/10/2024

Aujourd'hui : Régression linéaire (1/3)

TD (Fin modèle linéaire + Début régression)

Projet : A venir (en binôme, sur de la classification "linéaire")

RÉGRESSION LINÉAIRE ET LIEN AVEC LES MOINDRES CARRÉS



$$y_i = x_i^T \beta^* + \varepsilon_i$$

avec $\varepsilon_i \sim \mathcal{N}(0, 1)$

Contexte : $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times m}$, $y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$

avec $y_i = x_i^T \beta^* + \varepsilon_i$

β^* : Vérité terrain ("ground truth")

ε_i : bruit (variable aléatoire)

But : Trouver un modèle linéaire $x \mapsto x^T \beta$ tel que $y_i \approx x_i^T \beta$

→ On a vu comment calculer une solution de ce problème au sens des moindres carrés

→ Mais en général, β^* n'est pas une solution de ce problème

Q) Quel est le lien entre la solution au sens des moindres carrés et β^* ?

⇒ Besoin d'outils d'estimation statistique

① Régression linéaire simple

$$n=1$$

$$X \in \mathbb{R}^{m \times 1}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

$$y_i = \underbrace{x_i}_{1 \times 1} \underbrace{\beta}_{1 \times 1} + \varepsilon_i$$

ε_i variable aléatoire

Approche par moindres carrés linéaires

On cherche β solution de $\min_{\beta \in \mathbb{R}} \frac{1}{2} \|X\beta - y\|^2 = \frac{1}{2} \sum_{i=1}^m (x_i \beta - y_i)^2$

→ Le vecteur $X^T y$ est toujours une solution, et c'est même la solution de norme minimale. On l'appelle l'estimateur des moindres carrés ordinaires, et on le note $\hat{\beta}^{OLS} = X^T y$

(OLS = Ordinary Least Squares)

→ Dans le cas qui nous intéresse ($n=1$), on a 2 possibilités:

- Soit $\text{rang}(X) = 0$: alors $X = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \in \mathbb{R}^{m \times 1}$, $X^T = [0 \ \dots \ 0] \in \mathbb{R}^{1 \times m}$

$$\text{et donc } \hat{\beta}^{OLS} = 0$$

- Soit $\text{rang}(X) = 1 \leq \min(1, m)$, alors on sait que

$$\hat{\beta}^{OLS} = X^T y = \underbrace{(X^T X)^{-1}}_{\substack{1 \times m \times m \times 1 \\ 1 \times 1}} \underbrace{X^T y}_{\substack{1 \times m \quad m \times 1 \\ 1 \times 1}} = \frac{\sum_{i=1}^m x_i y_i}{\sum_{i=1}^m x_i^2}$$

⊕ Formule explicite

⊖ A priori indépendante de la distribution du bruit $\{\varepsilon_i\}_{i=1..m}$

Approche par maximum de vraisemblance

Hypothèse: $y_i = x_i \beta^* + \varepsilon_i$ où $\varepsilon_1, \dots, \varepsilon_m$ sont des variables aléatoires iid suivant une loi gaussienne $N(0, 1)$ (moyenne 0, variance 1)
 x_i fixés/courus
(Valeurs avec $\varepsilon_i \sim N(0, \sigma^2)$)
 $\varepsilon_i \sim N(0, 1) \Rightarrow y_i \sim N(x_i \beta^*, 1)$

→ Calculer un modèle linéaire, c'est trouver un estimateur de β^*
→ Principe du maximum de vraisemblance: Trouver la valeur la plus probable étant données les observations y_1, \dots, y_m

Def: L'estimateur du maximum de vraisemblance (sous l'hypothèse ci-dessus) est défini comme la/une solution, notée $\hat{\beta}^{MLE}$, du problème d'optimisation

$$\max_{\beta \in \mathbb{R}} \mathcal{L}(y_1, \dots, y_m | \beta)$$

↑ MLE
Maximum Likelihood Estimator

↑
Vraisemblance: Loi jointe de $y_1 - x_1 \beta, \dots, y_m - x_m \beta$

Interprétation

Si β définit un bon modèle,
alors $y_i - x_i^T \beta$ suit une loi gaussienne $N(0, 1)$ (ou $y_i \sim N(x_i \beta, 1)$)

$$\text{Vraisemblance: } \mathcal{L}(y_1, \dots, y_m | \beta) = P\left(Y_1 = y_1, \dots, Y_m = y_m \mid Y_i - x_i \beta \sim N(0, 1) \forall i\right)$$

↳ On aurait pu vouloir maximiser $P(B = \beta \mid y_1, \dots, y_m \sim N(0,1) \text{ iid})$
 mais cette probabilité n'est pas calculable de manière explicite,
 alors on utilise $P(B \mid y_1, \dots, y_m) \propto P(y_1, \dots, y_m \mid B)$
 (formule de Bayes)

↳ Formule de la vraisemblance dans notre cas:

$$\mathcal{L}(y_1, \dots, y_m \mid \beta) = P(y_1 = y_1, \dots, y_m = y_m \mid y_i - x_i \beta \sim N(0,1))$$

$$y_i - x_i^T \beta = \varepsilon_i \text{ iid} \rightarrow \left(= \prod_{i=1}^m P(y_i = y_i \mid y_i - x_i \beta \sim N(0,1)) \right)$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y_i - x_i \beta)^2\right)$$

$$\varepsilon \sim N(0,1) \quad P(\varepsilon = t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right)$$

$$= \left(\frac{1}{\sqrt{2\pi}}\right)^m \exp\left(-\frac{1}{2} \sum_{i=1}^m (y_i - x_i \beta)^2\right)$$

$$= \left(\frac{1}{\sqrt{2\pi}}\right)^m \exp\left(-\frac{1}{2} \|X\beta - y\|^2\right)$$

Théorème: Les problèmes $\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \|X\beta - y\|^2$

et $\max_{\beta \in \mathbb{R}^n} \mathcal{L}(y_1, \dots, y_m \mid \beta)$

ont le même ensemble de solutions.

⚠ Sur l'hypothèse $y_i - x_i \beta \sim N(0,1)$

Remarque: $\max_{\beta \in \mathbb{R}} \mathcal{L}(y_1, \dots, y_m | \beta)$ et $\max_{\beta \in \mathbb{R}} \ln(\mathcal{L}(y_1, \dots, y_m | \beta))$

ont le même ensemble de solutions

$$\ln(\mathcal{L}(y_1, \dots, y_m | \beta)) = -\frac{1}{2} \|X\beta - y\|^2 + \underbrace{\ln\left(\frac{1}{\sqrt{2\pi}}\right)^m}_{\text{constante}}$$

$\max_{\beta \in \mathbb{R}} \ln(\mathcal{L}(y_1, \dots, y_m | \beta))$ revient à maximiser $\beta \mapsto -\frac{1}{2} \|X\beta - y\|^2$

$\max_{\beta \in \mathbb{R}} -\frac{1}{2} \|X\beta - y\|^2$ et $\min_{\beta \in \mathbb{R}} \frac{1}{2} \|X\beta - y\|^2$ ont le même ensemble de solutions

↳ Dans notre cas (hypothèse ε_i iid $\sim N(0,1)$), on peut donc poser

$$\hat{\beta}^{MLE} = \hat{\beta}^{OLS} = X^T y \in \mathbb{R} \quad (m=1)$$

Propriétés statistiques de $\hat{\beta}^{MLE}$

• $\hat{\beta}^{MLE}$ est un estimateur sans biais: $\mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} [\hat{\beta}^{MLE}] = \beta^*$

(si $X = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, on considère que $\beta^* = 0$)

$$\text{Si } \text{rang}(X) = 1, \quad \hat{\beta}^{MLE} = \frac{\sum_{i=1}^m x_i y_i}{\sum_{i=1}^m x_i^2} = \frac{\sum_{i=1}^m x_i (x_i \beta^* + \varepsilon_i)}{\sum_{i=1}^m x_i^2}$$

et donc
$$\mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} [\hat{\beta}^{MLE}] = \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} \left[\frac{\sum_{i=1}^m (x_i^2 \beta^* + x_i \varepsilon_i)}{\sum_{i=1}^m x_i^2} \right]$$

linéarité de l'espérance $\left(\right) = \frac{1}{\sum_{i=1}^m x_i^2} \sum_{i=1}^m (x_i^2 \beta^* + x_i \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} [\varepsilon_i])$

$\varepsilon_i \text{ iid } \sim \mathcal{N}(0, 1)$ $\left(\right) = \frac{1}{\sum_{i=1}^m x_i^2} \sum_{i=1}^m (x_i^2 \beta^* + 0)$

$$= \frac{1}{\sum_{i=1}^m x_i^2} \sum_{i=1}^m x_i^2 \beta^* = \beta^*$$

• $\hat{\beta}^{MLE}$ est un estimateur convergent

$$\text{Var}_{\varepsilon_1, \dots, \varepsilon_m} [\hat{\beta}^{MLE}] = \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} \left[(\hat{\beta}^{MLE} - \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} [\hat{\beta}^{MLE}])^2 \right]$$

$$= \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} \left[(\hat{\beta}^{MLE} - \beta^*)^2 \right]$$

$\hat{\beta}^{MLE}$ convergent $(\Leftrightarrow) \text{Var}_{\varepsilon_1, \dots, \varepsilon_m} [\hat{\beta}^{MLE}] \xrightarrow{m \rightarrow \infty} 0$

lorsque $\text{rang}(X) = 1$, on a

$$\text{Var}_{\varepsilon_1, \dots, \varepsilon_m} [\hat{\beta}^{MLE}] = \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} \left[\left(\frac{\sum_{i=1}^m x_i y_i}{\sum_{i=1}^m x_i^2} - \beta^* \right)^2 \right]$$

$$= \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} \left[\left(\frac{\sum_{i=1}^m x_i^2 \beta^* + x_i \varepsilon_i}{\sum_{i=1}^m x_i^2} - \frac{\sum_{i=1}^m x_i^2 \beta^*}{\sum_{i=1}^m x_i^2} \right)^2 \right]$$

$$\varepsilon_i \sim N(0,1) \text{ iid}$$

$$E\left[\sum_{i=1}^m x_i \varepsilon_i\right] = 0$$

$$\text{Var}[X] = E\left[\left(X - E[X]\right)^2\right]$$

$$\begin{matrix} \varepsilon_i \text{ iid} \\ \varepsilon_1, \dots, \varepsilon_m \end{matrix} \quad \text{Var}\left[\sum_{i=1}^m x_i \varepsilon_i\right] = \sum_{i=1}^m x_i^2 \text{Var}_{\varepsilon_i}[\varepsilon_i]$$

$$= E_{\varepsilon_1, \dots, \varepsilon_m} \left[\left(\frac{\sum_{i=1}^m x_i \varepsilon_i}{\sum_{i=1}^m x_i^2} \right)^2 \right]$$

$$= \frac{1}{\left(\sum_{i=1}^m x_i^2\right)^2} E_{\varepsilon_1, \dots, \varepsilon_m} \left[\left(\sum_{i=1}^m x_i \varepsilon_i \right)^2 \right]$$

$$= \frac{1}{\left(\sum_{i=1}^m x_i^2\right)^2} \text{Var}_{\varepsilon_1, \dots, \varepsilon_m} \left(\sum_{i=1}^m x_i \varepsilon_i \right)$$

$$= \frac{1}{\left(\sum_{i=1}^m x_i^2\right)^2} \sum_{i=1}^m x_i^2 \text{Var}_{\varepsilon_i}[\varepsilon_i]$$

$$= \frac{\sum_{i=1}^m x_i^2}{\left(\sum_{i=1}^m x_i^2\right)^2} = \frac{1}{\sum_{i=1}^m x_i^2} \rightarrow 0 \quad m \rightarrow \infty$$

si on mesure en
une infinité de $x_i \neq 0$

Bilan: Dans le cas $m=1$ et sous l'hypothèse de bruit gaussien, l'approche par maximum de vraisemblance est équivalente à l'approche par moindres carrés, dans le sens où elle conduit au même estimateur.

En connaissant la loi du bruit, on peut en plus donner des garanties statistiques pour cet estimateur (non biaisé $E[\hat{\beta}^{MLE}] = \beta^*$ convergent $\text{Var}[\hat{\beta}^{MLE}] \rightarrow 0$ $m \rightarrow \infty$)

(2) Régression linéaire multiple

$$m \geq 1$$

$$y_i = x_i^T \beta^* + \varepsilon_i, \quad \varepsilon_i \text{ iid } \sim \mathcal{N}(0, 1)$$

Les résultats précédents se généralisent à ce cas :

- Equivalence entre la maximisation de la vraisemblance et les moindres carrés $(\hat{\beta}^{MLE} = \hat{\beta}^{OLS} = X^+ y \in \mathbb{R}^m)$

- Estimateur sans biais

$$\mathbb{E}_{\varepsilon_1, \dots, \varepsilon_m} [\hat{\beta}^{MLE}] = \beta^* \in \mathbb{R}^m$$

$v \in \mathbb{R}^m$ aléatoire

$$\mathbb{E}[v] = \begin{bmatrix} \mathbb{E}[v_1] \\ \vdots \\ \mathbb{E}[v_m] \end{bmatrix}$$

$$\begin{aligned} \text{Cov}[v] &= \mathbb{E} \left[\overbrace{(v - \mathbb{E}[v])}^{m \times 1} \overbrace{(v - \mathbb{E}[v])^T}^{1 \times m} \right] \\ &= \left[\mathbb{E}[(v_i - \mathbb{E}[v_i])(v_j - \mathbb{E}[v_j])] \right]_{ij} \end{aligned}$$

- Estimateur convergent

$$\text{trace} \left(\text{Cov}[\hat{\beta}^{MLE}] \right) \xrightarrow{m \rightarrow \infty} 0$$

Pour la suite

→ ε_i suit une loi non gaussienne ?

→ A priori sur β^* ?