

Rappel: Modèle de régression linéaire simple Nupt

Base: The social production of Criminal Homicide
13 variables 20 observations sociétés

Les VI: City / Pop / Homic / Poor %
ville Arm nb de pu on millia victime / 100 hab

Obj: expliquer pourquoi le taux d'Homicide est élevé dans certaines villes que dans d'autres

• m^{ême} nombre d'homicide mais nb pop $\neq \Rightarrow$ Ça p^{eut} une tient compte de la criminalité

Le taux d'homicide standardisé, rationalisé (normalisé) par rapport à la population, a permis de comparer des villes de population différente.

Histo de la distribution d'Homicide: montre que la variable endogène Homicide ne suit pas une dist normale.

Estimation Kernel densité / densité du noyau: Approche érigée dans la méthodologie de l'histogramme: Estimer la fonction de densité en un point x en utilisant les observi

Soit la fonction de poids: $K(x) = \begin{cases} 1/2 & \text{si } |x| < 1 \\ 0 & \text{sinon} \end{cases}$
de noyau de Kernel

l'estimateur densité Kernel:
fonction additif de Kernel / noyau

$$f(x) = \frac{1}{n \cdot h} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

Certaines fonctions de Kernel:

Uniform: $\frac{1}{2} \mathbb{1}(|x| < 1)$ Epanechnikov: $\frac{3}{4} (1-x^2) \mathbb{1}(|x| < 1)$

Gaussian: $\frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2)$ Biweight: $\frac{15}{16} (1-x^2)^2 \mathbb{1}(|x| < 1)$

La mom normalisée: valeur élevée de min et max

Asymétrie de la dist: faible mesure de centrage alternance

Mediane = Moyenne: dist sym Skewness = 0

Moyenne > Mediane: dist affectée par des val. ext. ≥ 0 skew > 0

Moyenne < Mediane: dist affectée par des outliers < 0 skew < 0

Test de normalité: JB-test: S: Skewness EK: Excess de Kurtosis

$H_0: S = EK = 0$ ou $EK = \text{Kurtosis} - 3 \Leftrightarrow \text{Kurtosis} = 3$

\Leftrightarrow dist normale

La stat: $JB = N \left(\frac{\hat{S}^2}{6} + \frac{\hat{EK}^2}{24} \right) \sim \chi^2(2)$

avec: $\hat{S}^2 = \frac{\frac{1}{N} \sum (y_i - \bar{y})^3}{\left(\frac{1}{N} \sum (y_i - \bar{y})^2 \right)^{3/2}}$ $\hat{EK} = \frac{\frac{1}{N} \sum (y_i - \bar{y})^4}{\left(\frac{1}{N} \sum (y_i - \bar{y})^2 \right)^2} - 3$

$\rightarrow p\text{-value} < \alpha\% \Rightarrow$ distribution non normale

Box Plot:

○ v. observés

\rightarrow à cause du pom taux d'homicide exceptionnellement élevés dans l'échantillon \rightarrow a un très grand résidu

\Rightarrow E de la régression ne sont pas normalement distribuées

\Rightarrow On ne peut pas appliquer les distributions t et f

regression homie poor: homie = 0,9438 Poor = 0,8152
 Bo 11% Cas pas de ville sans pauvreté ← Pas de / } Le taux dans les villes 0%
 La hom 0,9

RNS = NES: l'estimation de la variance d'erreur

$$RNS = S^2 = \frac{\sum \hat{e}_i^2}{n-p} (= 29,53)$$

Root NSE: déviation standard: mesure de diagnostic de la regression

$$\text{Root NSE} = \sqrt{S^2} (= 5,43)$$
 Erreur stand de la regression

Root NSE élevée \Rightarrow résultats de regression non fiables

homie hat: valeur prédite res: sont les résidus de l'estim
 \Rightarrow Pas de corrélation entre les val. prédites et résidus

La distance de Cook, effet levier / leverage, les résidus standardisés: outils qui détectent les valeurs aberrantes

• Effet de levier: l'éloignement d'un point % à \bar{Y}
 Le levier $h_i \gg 2 \cdot \frac{p+1}{n}$: règle de détection: valeur de levier à partir de laquelle on commence à s'inquiéter

Horiz: les valeurs moy d'effet de levier
 Vertical: les valeurs moy des résidus normalisés
 v. aberrante \rightarrow n'a pas trop de levier

Distance de Cook: influence globale des y_i sur les \hat{y}_i
 se calcule à l'aide des valeurs de levier et résidus standardisés
 Comparer à la F-stat à $(n-p)$ ddl

<20% : peu d'inf 50% : grande inf

région critique: $D_i > \frac{4}{m-p-2}$

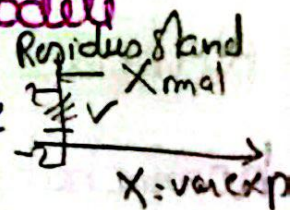
• Résidus Standardisés: Comparer y_i à \hat{y}_i

$$h_{ii} = \frac{1}{m} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \quad E(2.2)$$
$$r_i = \frac{e_i}{\sqrt{1-h_{ii}}} \quad \forall i=1, \dots, m$$

résidus stand \sim les résidus normalisés sauf $V(r_i)=1$

Résidu élevé \Rightarrow point mal reconstitué par le modèle

Pour détecter les valeurs atypiques : on trace le nuage

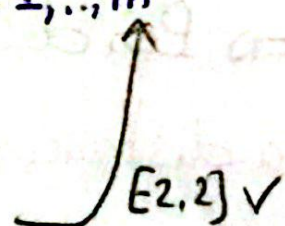


• Résidus Studentisés: $es_{(i)} = \frac{e_i}{s_{(i)} \sqrt{1-h_{ii}}} \quad \forall i=1, \dots, m$

↳ permettent de mieux évaluer

l'importance d'un résidu e_i

pour détecter les
valeurs atypiques



R Student élevé \Rightarrow point mal modélisé et R stand faible

point exagérément influent
altère les résultats visuels

Régression homoc pour pop:

$$\text{homoc} = 0,656 \text{ pour} + 0,022 \text{ pop} - 3,899$$

La pauvreté exerce un effort sur

homoc plus fort que pop

pauv-poor en % et population-pop en millions : On calcule

les beta pour rendre comparable la force des coeff

beta : les coeff après une 1^{re} conversion Z-scores

↳ les coeff beta sont directement comparables

reg hommic pour pop, beta

$$\text{homic} = 0,656 \text{ poor} + 0,022 \text{ pop} - 3,899$$

Les β sont inversés
pop est plus colinéaire
avec la var hommic

$$\tilde{\beta}_K = \hat{\beta}_K \frac{s_{\text{exg}}}{s_y} \quad R=1, \dots, K-1$$

\downarrow
 $\tilde{\beta}_K \in [-1, 1]$: coeff standardisé de la Kème var exp

$\hat{\beta}_K$: coeff OLS habituel

s : déviation standard

Comparaison des modèles simples et multiples :

coeff de la var. exogène poor : $0,94 \rightarrow 0,65$

hautement biaisé car
la var inclue poor est ≥ 0 et fortement
corrélée à pop

F-test = 36,3 et p-value = 0 \Rightarrow Rejette l'hypothèse de base $\neq 0$

\Rightarrow Les 2 paramètres sont conjointement significatifs $\neq 0$
dans la population

$$y_i = \alpha + \varepsilon_i \quad \forall i=1, \dots, m$$

Le bruit imprévisible

\downarrow tout regression est apte
à prédire la variable
endog hommic

Les étapes pour le calcul du F-test :

Estimer le modèle (OLS) \rightarrow $SCR_{mc} = 135,276$ \xrightarrow{RSS}

Estimer le modèle avec seulement l'intercept $\rightarrow SCR_e = 712,94$

Reg hommic : $\text{homic} = 6,9$ Estimation OLS de l'intercept $= \bar{y}$

$$\text{Stat du test : } \hat{F} = \frac{SCR_e - SCR_{mc}}{SCR_{mc}} \times \frac{N-K}{9} \sim F(9, N-K)$$

N : nb Obs

K : nb des param

nb des param restreints $9 = K-1$

On peut tester la significativité des 2 paramètres par une autre approche:

Calcul des résidus OLS restreints: $\hat{\epsilon}_i = y_i - \hat{a}$ définis sous l'hypothèse $H_0: \beta_1 = \beta_2 = 0$, X_{i1} et X_{i2} doivent être non corrélés

Régression par OLS $\hat{\epsilon}_i$ sur X_{i1} et X_{i2}

Calcul LN-stat = $NR^2 \sim \chi^2(q)$ q : nbr de contraintes à tester

$R^2 \uparrow \Rightarrow \hat{\epsilon}_i$ est signif corrélé avec X_{i1} et $X_{i2} \Rightarrow$ rejet $H_0 \Rightarrow$ significative

il faut inclure X_{i1} et X_{i2} dans la régression

RNSE non proche de 0 \Rightarrow le modèle évalué n'est pas meilleur en terme d'exactitude