

$$P_{obs} = P_{H_0}(|T| > |t_{obs}|)$$

$$= P_{H_0}(T > |t_{obs}|) + P_{H_0}(T < -|t_{obs}|)$$

sous H_0 , $T \sim t_{n-1}$ qui est symétrique

$$= 2P_{H_0}(T > |t_{obs}|) = 2(1 - P_{H_0}(T \leq |t_{obs}|))$$

$$= 2(1 - F_{t_{n-1}}(|t_{obs}|))$$

Cadre statistique : [4] La p -valeur

On peut utiliser la p -valeur pour faire un test d'hypothèse :

$$\boxed{\text{rejetter } H_0 \iff p_{obs} < \alpha}$$

Exemple bilatérale (203) $p_{obs} = 2(1 - F_{t_{n-1}}(t_{obs}))$ donc

$$p_{obs} < \alpha \iff F_{t_{n-1}}(|t_{obs}|) > 1 - \alpha/2 \iff |t_{obs}| > t_{n-1, 1-\alpha/2}$$

De manière générale, l'approche de la p -valeur est équivalente à l'approche des valeurs critiques. Cependant, la p -valeur p_{obs} fournit une information plus facilement interprétable que la valeur t_{obs} . Il s'agit d'une mesure de l'évidence contre H_0 contenue dans les données.

Attention : la p -valeur **n'est pas** la probabilité que H_0 soit vraie.

Résumé : les éléments d'un test

- A Une **hypothèse nulle** H_0 à tester contre une hypothèse alternative H_1 .
- B Une **statistique de test** T , choisie de telle sorte que des valeurs "extrêmes" de T (en direction de H_1) suggèrent que H_0 est fausse. La valeur observée de T est t_{obs} .
- C Un **niveau de significativité** α , qui la probabilité d'erreur de type I (rejet de H_0 quand H_0 est vraie) maximale que nous allons tolérer.
- D1 Des **valeurs critiques**, telles que quand T tombe au-delà de ces valeurs, nous rejetons H_0 en faveur de H_1 . Les valeurs critiques sont choisies pour respecter le niveau de significativité α .

Au lieu de D1, nous pouvons utiliser l'approche équivalente D2 :
- D2 Une **valeur** p_{obs} donnant la probabilité d'observer une valeur de T aussi élevée que t_{obs} sous H_0 . On rejette alors H_0 en faveur de H_1 quand $p_{\text{obs}} \leq \alpha$.

Exemple

Exemple On a contrôlé 10 compteurs d'électricité nouvellement fabriqués et obtenu les valeurs suivantes (en MW) :

983 1002 998 996 1002 983 994 991 1005 986.

On suppose qu'il s'agit de réalisation d'un échantillon iid d'une loi normale. On aimerait savoir s'il y a un écart entre la moyenne attendue de 1000 MW et la moyenne réelle des compteurs qui sortent de la fabrication. Nous avons obtenu $\bar{x} = 994 < 1000$. S'agit-il d'un hasard ou une faute de production ?

On va prendre $\alpha = 5\%$.

Solution Exemple 208

Supposons que nos observations x_1, \dots, x_n soient des réalisations de variables aléatoires $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, avec σ^2 inconnu. On veut tester : $H_0 : \mu = \mu_0$ contre $H_1 : \mu \neq \mu_0$, où $\mu_0 = 1000$. On prend comme statistique de test

$$T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim t_{n-1} \text{ sous } H_0 : \mu = \mu_0.$$

Dans notre cas $n = 10$, $\mu_0 = 1000$, $\bar{x} = 994$, et

$$s^2 = \frac{1}{9} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{9} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) = 64.88,$$

donc $t_{\text{obs}} = -2.35$. $\quad = \frac{994 - 1000}{\sqrt{64.88/10}}$

On rejette H_0 si et seulement si $|t_{\text{obs}}| > t_{n-1, 1-\alpha/2}$ (cas bilatéral). , $t_{n-1, 1-\alpha/2} = 2.262$ (voir les tables), et comme $t_{\text{obs}} = -2.35 < -2.262$, on rejette l'hypothèse H_0 .

$$p_{\text{obs}} \approx 0.046 < 0.05 \Rightarrow \text{rejet}$$

Intervalles de confiance et tests

- Soient $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$ avec μ, σ inconnus
- Considérons $H_0 : \mu = 0$ et $H_1 : \mu \neq 0$
- On rejette H_0 au niveau α si et seulement si

$$|T_n| = \left| \sqrt{n} \frac{\bar{Y}_n}{S_n} \right| > t_{n-1, 1-\alpha/2}$$

- Intervalle de confiance (IC) de niveau $1 - \alpha$

$$\left[\bar{Y}_n - \frac{t_{n-1, 1-\alpha/2}}{\sqrt{n}} S_n, \bar{Y}_n + \frac{t_{n-1, 1-\alpha/2}}{\sqrt{n}} S_n \right]$$

- Cet intervalle contient zéro si et seulement si

$$|\bar{Y}_n| \leq t_{n-1, 1-\alpha/2} S_n / \sqrt{n} \text{ et donc}$$

on rejette $H_0 \iff 0$ n'est pas dans l'intervalle de confiance

- De manière générale, rejet de $H_0^{\theta_0} : \theta = \theta_0$ est équivalent à {l'IC ne contient pas θ_0 } si on se base sur la même statistique de test
- α petit \iff difficile à rejeter \iff IC de niveau $1 - \alpha$ large

Intervalle de confiance (IC) et tests

$$\text{on rejette } p_{\text{obs}} \iff \sqrt{n} \left| \frac{\bar{Y}_n}{s_n} \right| \leq t_{n-1, 1-\frac{\alpha}{2}}$$

$$\iff -t_{n-1, 1-\frac{\alpha}{2}} \leq -\sqrt{n} \frac{\bar{Y}_n}{s_n} \leq t_{n-1, 1-\frac{\alpha}{2}}$$

$$\iff \bar{Y}_n - \frac{t_{n-1, 1-\frac{\alpha}{2}} s_n}{\sqrt{n}} \leq 0 \leq \bar{Y}_n + \frac{t_{n-1, 1-\frac{\alpha}{2}} s_n}{\sqrt{n}}$$

$$\iff 0 \in \text{IC}(1-\alpha)$$

- Rejet $\iff p_{\text{obs}} \leq \alpha \iff \theta_0 \notin \text{IC de niveau } 1-\alpha$
- **IC** : α fixé, pour quels $\theta_0 \in \mathbb{R}$ $H_0^{\theta_0}$ n'est pas rejetée?
- **p-valeur** : θ_0 fixé, pour quels $\alpha \in (0, 1)$ $H_0^{\theta_0}$ est rejetée?

$$\theta = \theta_0$$

3.4 Tests khi-deux

Le test khi-deux

- On se pose la question de l'adéquation d'une distribution théorique à des données
- Supposons que dans une expérience on observe n résultats différents avec des
 - fréquences observées** dans k classes disjointes o_1, \dots, o_k , alors que les
 - fréquences théoriques** correspondantes sont e_1, \dots, e_k ,
 - où $\sum_{i=1}^k o_i = \sum_{i=1}^k e_i = n$
- On a H_0 : "les observations proviennent de la loi théorique spécifiée"
- Une mesure de l'écart entre les o_j et les e_j est donnée par la **statistique khi-deux**

$$T_n = \sum_{i=1}^k \frac{(o_i - e_i)^2}{e_i}$$

- Si n est grand et les e_i ne sont pas trop petites (règle de pouce : $e_i \geq 5$ pour la plupart), alors $T_n \stackrel{\text{app}}{\sim} \chi_\nu^2$ sous H_0 , où

- χ_ν^2 est la loi **khi-carré**, la loi de $\sum_{i=1}^\nu Z_i^2$ où $Z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$
- $\nu = k - 1$ si les e_i peuvent être calculés sans devoir estimer des paramètres inconnus
- $\nu = k - 1 - c$ si les e_i sont calculés après avoir estimé c paramètres

Exemple

P valeur $W \sim \chi^2_5$ $P_{obs} = Pr(W > t_{obs}) \approx 0,131 \Rightarrow$ pas de rejet

rejet de H_0 au niveau $\alpha \iff t_{obs} = \sum_{i=1}^k \frac{(o_i - e_i)^2}{e_i} > \chi^2_{v, 1-\alpha} \approx 11,4$

Exemple $n = 60$ jets d'un dé ont donné la répartition suivante :

Valeur x_i	1	2	3	4	5	6	
Valeur o_i	8	10	9	16	13	4	60

es 10 10 10 10 10 10

$\alpha = 0,05$

$v = k - 1 = 5$

$\chi^2_{5, 0,95}$

Ici $k = 6$ et H_0 : "équilibre du dé" est équivalente au modèle

même $\sigma: \mathcal{L} = 0,1$
 $Pr(X = x) = 1/6, \quad x \in \{1, \dots, 6\}$

$$t_{obs} = \frac{(8-10)^2}{10} + \frac{(10-10)^2}{10} + \dots + \frac{(4-10)^2}{10} = 8,5$$

$t_{obs} < \chi^2_{5, 0,95} \Rightarrow$ pas de rejet

Exemple

Exemple L'intelligence (QI) X de $n = 1000$ personnes est testée :

$$a_i = 0 \quad b_i = 70$$

Score X	[0, 70)	[70, 85)	[85, 100)	[100, 115)	[115, 130)	[130, ∞)
Nombre o_i	34	114	360	344	120	28

$$e_i \quad 22.75 \quad 285.91 \quad 391.34 \quad 341.34 \quad 235.91 \quad 22.75$$

Est-il plausible que $X \sim \mathcal{N}(100, 15^2)$?

On a

$$e_i = n \Pr_{\mathcal{N}(100, 15^2)}(a_i \leq X \leq b_i) = n \Pr \left(\frac{a_i - 100}{15} \leq \frac{X - 100}{15} \leq \frac{b_i - 100}{15} \right)$$

$$= n \Phi \left(\frac{b_i - 100}{15} \right) - n \Phi \left(\frac{a_i - 100}{15} \right) \quad V = 6 - 1 = 5$$

$$\alpha = 0.05$$

$$t_{obs} = \frac{(34 - 22.75)^2}{22.75} + \dots + t \approx 13.27 \quad n \approx \chi_{5, 0.9}^2$$

$$p_{obs} \approx 0.027$$

Tableaux de contingence

Un **tableau de contingence** est une classification de n objets ou individus selon plusieurs critères

- Une question fondamentale concerne **l'indépendance** des critères
- Supposons qu'on observe deux caractères A (h classes) et B (k classes) sur chacun de n individus, donnant le **tableau de contingence** suivant :

	B						
A	1	2	...	j	...	k	Σ
1	n_{11}	n_{12}	...	n_{1j}	...	n_{1k}	$n_{1\cdot}$
2	n_{21}	n_{22}	...	n_{2j}	...	n_{2k}	$n_{2\cdot}$
...
i	n_{i1}	n_{i2}	...	n_{ij}	...	n_{ik}	$n_{i\cdot}$
...
h	n_{h1}	n_{h2}	...	n_{hj}	...	n_{hk}	$n_{h\cdot}$
Σ	$n_{\cdot 1}$	$n_{\cdot 2}$...	$n_{\cdot j}$...	$n_{\cdot k}$	$n_{\cdot\cdot} = n$

sommes des lignes

Sommes de colonnes

Indépendance

- Soit n_{ij} le nombre de personnes tombant dans la classe i du caractère A et dans la classe j du caractère B , et soit

$$n_{i.} = \sum_{j=1}^k n_{ij}, \quad n_{.j} = \sum_{i=1}^h n_{ij}, \quad i \in \{1, \dots, h\}, j \in \{1, \dots, k\}$$

- On désire tester H_0 : **A et B sont indépendants**. Dans ce cas

$$\Pr(A = i, B = j) = \Pr(A = i) \times \Pr(B = j), \quad i \in \{1, \dots, h\}, j \in \{1, \dots, k\},$$

et les probabilités empiriques sont

$$\widehat{\Pr(A = i)} = \frac{n_{i.}}{n}, \quad \widehat{\Pr(B = j)} = \frac{n_{.j}}{n}$$

Ainsi, sous H_0 , l' (i, j) ième élément du tableau de contingence est

$$E_{ij} = n \times \Pr(A = i, B = j) = n \times \Pr(A = i) \times \Pr(B = j)$$

qu'on va estimer par

$$e_{ij} = n \frac{n_{i.}}{n} \frac{n_{.j}}{n} = \frac{n_{i.} n_{.j}}{n}$$

valeur attendue

Calcul sous H_0

- Sous H_0 et pour n grand, la statistique T_n suit une distribution χ^2_ν avec $\nu = (h-1)(k-1)$, car on a dû estimer $c = (h-1) + (k-1)$ probabilités, et
$$hk - 1 - c = kh - 1 - (h-1) - (k-1) = (h-1)(k-1)$$
- Pour tester à un niveau de significativité α fixé, on rejette H_0 si $t_{\text{obs}} > \chi^2_{(h-1)(k-1), 1-\alpha}$, sinon on ne rejette pas H_0

Exemple On a relevé parmi 95 personnes la couleur de leurs yeux (caractère A) et celle de leurs cheveux (caractère B) et on a obtenu les résultats suivants :

A	B		
	Cheveux clairs	Cheveux sombres	
Yeux bleus	$n_{11} = 32$	$n_{12} = 12$	$n_{1.} = 44$
Yeux bruns	$n_{21} = 14$	$n_{22} = 22$	$n_{2.} = 36$
Autres	$n_{31} = 6$	$n_{32} = 9$	$n_{3.} = 15$
Σ	$n_{.1} = 52$	$n_{.2} = 43$	$n = 95$

On désire tester si la couleur des cheveux est indépendante de celle des yeux

Donc on a H_0 : indépendance entre couleur des cheveux et couleur des yeux

Pobs
 ≈ 0.0048

Q: 32	$\frac{44 \times 52}{95}$	12	}	44	
L: 14		22		$\frac{36 \times 43}{95}$	36
	$\frac{15 \times 52}{95}$	9			15
52		43			95

$$t_{obs} = \frac{\left(32 - \frac{44 \times 52}{95}\right)^2}{\frac{44 \times 52}{95}} + \dots + \frac{\left(9 - \frac{15 \times 43}{95}\right)^2}{\frac{15 \times 43}{95}}$$

≈ 10.67 . $v = (h-1)(k-1) = 2$ $\chi^2_{2,0.05} \approx 5.99 < t_{obs}$
 reject ²¹⁹

Régularité (non-examinable)

Les conditions de régularité sont compliquées. Elles sont fausses le plus souvent quand

- un des paramètres est discret
- le support de $f(y; \theta)$ dépend de θ
- le vrai θ se trouve sur une borne des valeurs possibles

Elles sont satisfaites dans la grande majorité des cas rencontrés en pratique

Exemple Soient $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} U(0, \theta)$, trouver la vraisemblance $L(\theta)$ et le $\hat{\theta}_{\text{ML}}$. Montrer que la loi limite de $n(\theta - \hat{\theta}_{\text{ML}})/\theta$ quand $n \rightarrow \infty$ est $\exp(1)$.
Discuter.

Preuve (non-examinable)

- Les fonctions de densité et de répartition de y_j sont

$$f(y; \theta) = \theta^{-1} I(0 \leq y \leq \theta), \quad F(y) = y/\theta, \quad 0 \leq y \leq \theta.$$

- L'indépendance donne

$$L(\theta) = \prod \theta^{-1} I(Y_j < \theta) = \theta^{-n} I(\max Y_j \leq \theta), \quad \theta > 0,$$

qui est maximisée au point $\hat{\theta}_{\text{ML}} = M_n = \max Y_j$

- On a

$$\Pr(M_n \leq x) = \prod_{j=1}^n \Pr(Y_j \leq x) = (x/\theta)^n, \quad 0 \leq x \leq \theta$$

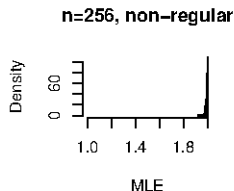
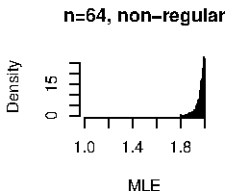
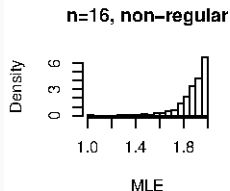
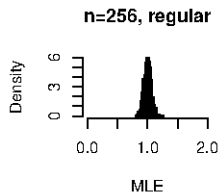
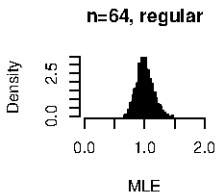
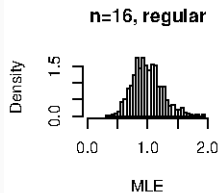
- Donc pour $x \geq 1$ et $n \geq x$,

$$\begin{aligned} \Pr \left\{ n(\theta - \hat{\theta}_{\text{ML}})/\theta \leq x \right\} &= \Pr(\hat{\theta}_{\text{ML}} \geq \theta - x\theta/n) = 1 - \{(\theta - x\theta/n)/\theta\}^n \\ &\rightarrow 1 - \exp(-x), \end{aligned}$$

comme requis

Exemple (non examinable)

Comparaison des lois de $\hat{\theta}$ dans un cas régulier (en haut, avec écart-type $\propto n^{-1/2}$ et loi limite normale) et dans un cas non-régulier (en bas, avec écart-type $\propto n^{-1}$ et loi limite non-normale)



4. Régression

4.1 Introduction

Motivation

La **régression** concerne la relation entre une variable d'intérêt et d'autres variables. On note

- la variable d'intérêt, la **variable de réponse**, y , et on la considère comme variable aléatoire
- les autres variables, les **covariables** (variables explicatives) sont notées $x^{(1)}, \dots, x^{(p)}$, on les considère comme fixées

On peut s'intéresser à

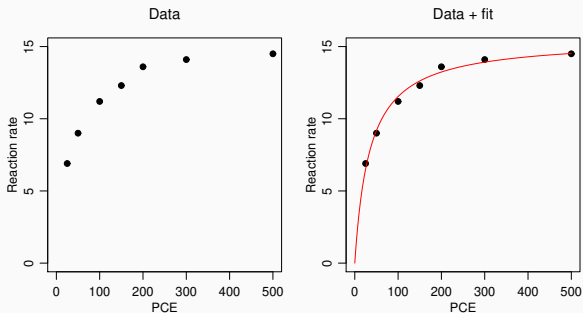
- l'**estimation** d'une relation éventuelle entre y et les $x^{(j)}$, ou
- la **prévision** des valeurs futures/manquantes de y sur la base des $x^{(j)}$ correspondantes

Réaction chimique

Professeur Christophe Holliger (SIE) : on essaye de déterminer les paramètres cinétiques d'une 'reductive dehalogenase dechlorinating tetrachloroethene (PCE)'. Ceci dépend de la concentration du substrat, et la vitesse de la réaction peut être exprimé par l'équation de Michaelis-Menten

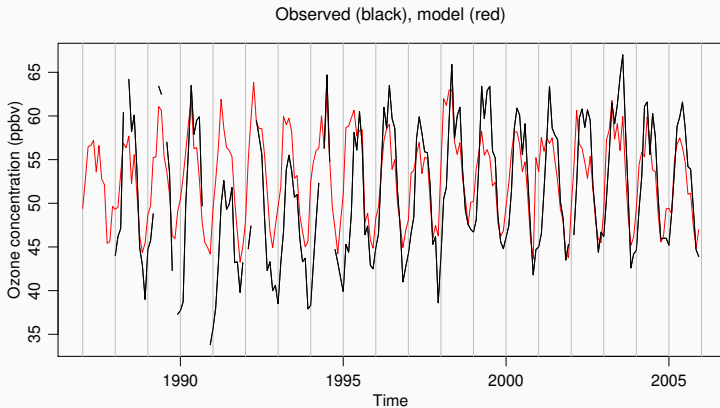
$$y = \frac{\gamma_0 x}{\gamma_1 + x},$$

où x est la concentration de PCE, γ_0 est la vitesse maximale, et γ_1 est la concentration quand $y = \gamma_0/2$. Comment estimer γ_0 et γ_1 ? Quelles sont leurs incertitudes?



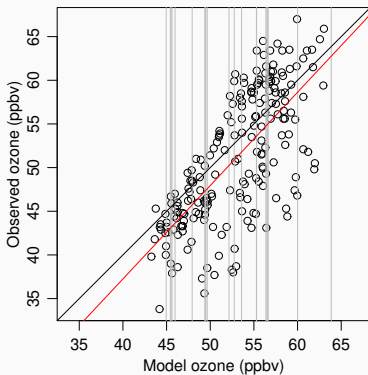
Ozone atmosphérique

Observations de la concentration de l'ozone au Jungfraujoch, de janvier 1987 à décembre 2005 (quelques valeurs manquantes), et résultats d'une modélisation



Soient y les données réelles et x les résultats du modèle

Relation linéaire ?



Les lignes verticales grises montrent des x dont les y sont manquants. La ligne noire montre la relation $y = x$, et la ligne rouge montre la meilleure estimation d'une relation linéaire entre y et x .

Comment utiliser la relation entre les résultats du modèle x et les données observées y pour estimer les y manquants ?

Problème d'ajustement

une seule covariable

- On considère une variable de réponse y que l'on cherche à expliquer par une covariable x
- On dispose d'un ensemble de points

$$\begin{pmatrix} x_1 \\ y_1 \end{pmatrix}, \dots, \begin{pmatrix} x_n \\ y_n \end{pmatrix}$$

qu'on peut représenter par un nuage de points (scatterplot) comme ceux d'auparavant

- D'une manière générale, le **problème d'ajustement** consiste à trouver une courbe $y = \mu(x)$ qui résume "le mieux possible" le nuage de points. La fonction $\mu(x)$ dépend de paramètres qu'il faut estimer
- S'il y a une **relation linéaire**, on peut utiliser la corrélation pour mesurer la dépendance linéaire entre les y et x . La régression linéaire permet de résumer cette dépendance par une droite

