Modèles statistiques

On étudie une **population** (ensemble d'individus ou d'éléments) à partir d'un **échantillon** (sous-ensemble).

- modèle statistique : X = la quantité étudiée (variable aléatoire); la loi F
 de X est supposée connue sauf un nombre fini de paramètres θ
- échantillon (doit être représentatif de la population) : "données" x_1, \ldots, x_n , souvent supposées comme étant une réalisation de $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} F$ (indépendantes et identiquement distribuées) ou $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} f$
- **statistique** : une fonction $T_n = h(X_1, ..., X_n)$ des variables aléatoires $X_1, ..., X_n$
- estimateur : une statistique utilisée pour estimer certains paramètres de F
- Notations :

$$T_n = h(X_1, \dots, X_n)$$
 est la statistique (variable aléatoire)
 $t_n = h(x_1, \dots, x_n)$ est la réalisation de T_n au moyen des x_i
 $\widehat{\theta}$ ou $\widehat{\theta}_n$ est un **estimateur** d'un paramètre θ

151

Commentaires

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$ et x_1, \ldots, x_n une réalisation correspondante. Alors

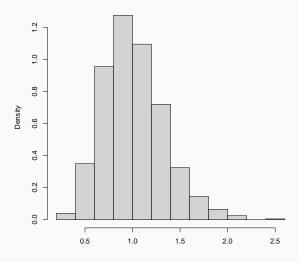
- $\widehat{\mu}_n = \overline{X}_n$ est une estimateur de μ , dont la valeur observée est \overline{x}_n
- $\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (X_i \overline{X}_n)^2$, est un estimateur de σ^2 , dont la valeur observée est $n^{-1} \sum_{i=1}^n (x_i \overline{X}_n)^2$

Remarques:

- une statistique T_n étant fonction des variables aléatoires X_1, \ldots, X_n , c'est elle-même une variable aléatoire!
- La loi de T_n dépend de la loi des X_i, et est appelée distribution d'échantillonnage de T_n
- Si on ne peut pas déduire la loi de T_n de celle des X_i , on doit se contenter parfois de connaître $\mathbb{E}(T_n)$ et $\text{var}(T_n)$, ou, si T_n est liée à \overline{X}_n , l'approximer à l'aide de $\mathbb{E}(X)$ et var(X) et le théorème central limite

Loi d'échantillonnage

Histogramme de 1000 réalisations de $\overline{X}_n = \frac{1}{10}(X_1 + \ldots + X_{10})$ où les $X_i \stackrel{\text{iid}}{\sim} \exp(1)$



Problèmes attaqués par la statistique

On suppose un **modèle** (c'est à dire une famille de distributions $f(x; \theta)$) et on souhaite

- estimer les paramètres θ de ce modèle
- poser des questions au sujet de la valeur de ces paramètres, par exemple tester si $\theta=0$
- prédire les valeurs des observations futures

Il existe plusieurs méthodes pour estimer les paramètres d'un modèle. On va décrire les suivantes :

- méthode des moments (simple)
- méthode du maximum de vraisemblance (souvent utilisée car optimale dans beaucoup de situations)

3.1 Estimation de paramètres

Méthode des moments

- Supposons que l'échantillon tiré soit représentatif de la population
- Pour obtenir des estimateurs pour les paramètres inconnus de la population, on égalise les "moments" de l'échantillon ("empirique") à ceux de la population ("théorique")
- kème moment
 - Population ("théorique") : $m_k = \mathbb{E}_{\theta}(X^k) = m_k(\theta)$. Comme la loi de X dépend de θ , $m_k = m_k(\theta)$
 - Echantillon ("empirique") : $\widehat{m}_k = \frac{1}{n} \sum_{i=1}^n X_i^k$
 - L'estimateur des moments s'obtient on égalisant m_k et \widehat{m}_k , ce qui donne des équation(s) pour $\theta \in \mathbb{R}^p$
 - Par exemple $m_1(\theta) = \mathbb{E}_{\theta}(X) = \sum_{i=1}^n X_i/n$
- On a donc besoin d'autant de moments (supposés finies!) que de paramètres inconnus

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} U(0, \theta)$. Estimer θ .

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$. Estimer μ et σ^2 .

Méthode des moments

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} U(0, \theta)$. Estimer θ .

Méthode des moments

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$. Estimer μ et σ^2 .

Méthode du maximum de vraisemblance

Définition: Soient x_1, \ldots, x_n des données supposées être une réalisation d'un échantillon aléatoire $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} f(x; \theta)$. La **vraisemblance** pour θ est la fonction

$$L(\theta) = f(x_1; \theta) \times f(x_2; \theta) \times \cdots \times f(x_n; \theta) = \prod_{i=1}^n f(x_i; \theta).$$

f est la fonction de densité ou de masse (on suppose qu'elle existe)

Définition: L'estimateur du maximum de vraisemblance (maximum likelihood) $\widehat{\theta}_{\mathrm{ML}}$ d'un paramètre θ est celui qui maximise la fonction de vraisemblance parmi tous les θ possibles :

$$L(\widehat{\theta}_{\mathrm{ML}}) \geq L(\theta)$$
 pour tout θ

Il est plus facile de maximiser $\ell(\theta) := \log L(\theta)$, souvent en résolvant $d\ell(\theta)/d\theta = 0$, et vérifiant qu'il s'agit bien d'un maximum (par exemple si la deuxime dérivée est négative $d^2\ell(\theta)/d\theta^2 < 0$)

Exemple x_1, \ldots, x_n réalisations d'une loi $\exp(\lambda)$ avec $\lambda > 0$. Estimer λ

Exemple : maximum de vraisemblance

Exemple Supposons que x_1, \ldots, x_n soient des réalisations i.i.d. d'une loi exponentielle,

$$f(x; \lambda) = \lambda e^{-\lambda x}, \ x \ge 0, \quad \lambda > 0.$$

Trouver $\widehat{\lambda}_{ML}$.

Erreur quadratique moyenne

Définition: L'erreur quadratique moyenne de l'estimateur $\widehat{\theta}$ de θ est

$$\mathrm{EQM}_{\theta}(\widehat{\theta}) = \mathbb{E}_{\theta}\{(\widehat{\theta} - \theta)^2\} = \cdots = \mathrm{Var}_{\theta}(\widehat{\theta}) + [b_{\theta}(\widehat{\theta})]^2,$$

où $b_{ heta}(\widehat{ heta}) = \mathbb{E}_{ heta}(\widehat{ heta}) - heta$ est le **biais** de $\widehat{ heta}$

- La distribution de $\widehat{\theta}$ dépend de celle des X_i et donc de θ
- Si $\widehat{\theta}$ et $\widehat{\theta}'$ sont deux estimateurs du même paramètre θ et $\mathsf{EQM}_{\theta}(\widehat{\theta}) \leq \mathsf{EQM}_{\theta}(\widehat{\theta}')$, on préfère $\widehat{\theta}$
- si $b_{\theta}(\theta) < 0$, alors $\widehat{\theta}$ sous-estime θ
- si $b_{\theta}(\theta) > 0$, alors $\widehat{\theta}$ sur-estime θ
- si $b_{\theta}(\theta) \equiv 0$, alors $\widehat{\theta}$ est **non biaisé**, et $\mathsf{EQM}_{\theta}(\widehat{\theta}) = \mathsf{var}(\widehat{\theta})$

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$. La médiane M_n et la moyenne \overline{X}_n sont (à peu près) non-biaisés pour μ mais $\text{var}(M_n) > \text{var}(\overline{X}_n)$. Lequel des estimateurs \overline{X}_n et M_n de μ est préférable? Et si des valeurs aberrantes peuvent apparaître?

Biais et variance

High bias, low variability



Low bias, high variability



High bias, high variability



The ideal: low bias, low variability



- θ = "bulle centrale", supposée être la vraie valeur
- fléchettes rouges = réalisations de $\widehat{\theta}$ qui estime θ

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, $\widehat{\mu} = \overline{X}_n$, et $\widehat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (X_i - \overline{X})^2$. On admet que $\text{var}_{\mu, \sigma^2}(\sum_{i=1}^n (X_i - \overline{X})^2) = 2\sigma^4(n-1)$. Trouver le biais et la variance de $\widehat{\mu}$. Trouver les valeurs de a qui minimisent le biais, la variance, et la EQM, pour $\widehat{\sigma}^a := a\widehat{\sigma}_n^2$.

Biais et variance : exemple

Exemple Soient $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, $\widehat{\mu} = \overline{X}_n$, et $\widehat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (X_i - \overline{X})^2$. On admet que $\text{var}_{\mu, \sigma^2}(\sum_{i=1}^n (X_i - \overline{X})^2) = 2\sigma^4(n-1)$. Trouver le biais et la variance de $\widehat{\mu}$. Trouver les valeurs de a qui minimisent le biais, la variance, et la EQM, pour $\widehat{\sigma}^a := a\widehat{\sigma}_n^2$.

Retour sur le maximum de vraisemblance

• Soient x_1, \ldots, x_n des données supposées être une réalisation d'un échantillon aléatoire $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} f(x; \theta)$ (densité / masse). La **vraisemblance** pour θ est la fonction

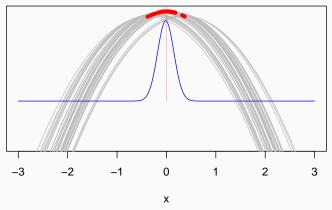
$$L(\theta) = f(x_1; \theta) \times f(x_2; \theta) \times \cdots \times f(x_n; \theta) = \prod_{i=1}^n f(x_i; \theta).$$

• $\widehat{\theta}_{\mathrm{ML}}$ satisfait

$$L(\widehat{\theta}_{\mathrm{ML}}) \geq L(\theta)$$
 pour tout θ

- Interprétation : dans le cas discret, on maximise $\Pr_{\theta}(X_1 = x_1, \dots, X_n = x_n)$
- Il est plus facile de maximiser $\ell(\theta) := \log L(\theta)$, souvent en résolvant $d\ell(\theta)/d\theta = 0$, et vérifiant qu'il s'agit bien d'un maximum

La vraisemblance est une fonction aléatoire



Fonctions de vraisemblances correspondantes à 50 échantillons de taille n=40. La vraie valeur est $\theta=0$ et les 50 maxima sont en rouge.

Distribution asymptotique de $\widehat{ heta}_{\mathrm{ML}}$

- Dans des modèles "réguliers" (normal, exponentiel, Poisson, \dots) on a $\ell'(\widehat{\theta}_{\mathrm{ML}})=0$
- Comme $\ell(\theta)$ est une somme de variables aléatoires iid, le théorème central limite s'applique
- Grâce à la méthode delta $\widehat{\theta}_{\mathrm{ML}} \stackrel{\mathrm{app}}{\sim} \mathcal{N}(\theta, 1/I_{n}(\theta))$ avec **l'information** de Fisher

$$I_n(\theta) = -\mathbb{E}_{\theta}(\ell''(\theta)) = \mathbb{E}_{\theta}(J_n(\theta)), \qquad J_n(\theta) = -\ell''(\theta)$$

 $1/I_n(\theta)$ est la **variance asymptotique** de $\widehat{\theta}_{\mathrm{ML}}$, $I_n(\theta)$ est la courbure

- L'information observée est $J_n(\widehat{\theta}_{\mathrm{ML}})$
- Attention! Certains modèles, tel que $U[0,\theta]$, ne sont pas réguliers

Exemple x_1, \ldots, x_n réalisations d'une loi $\exp(\lambda)$ avec densité $\lambda e^{-\lambda x} I(x \ge 0), \ \lambda > 0$

Exemple

Exemple x_1, \ldots, x_n réalisations d'une loi $\exp(\lambda)$ avec densité $\lambda e^{-\lambda x} I(x \ge 0), \ \lambda > 0$

3.2 Estimation par intervalle

Les intervalles de confiance

Un élément clé de la statistique est de donner une idée de l'incertitude d'un constat Soit θ un paramètre inconnu, et soit $\tilde{\theta}_n=1$ une estimation de θ basée sur y_1,\ldots,y_n :

- si $n=10^5$ on est beaucoup plus sûr que $heta pprox ilde{ heta}_n$ que si n=10
- pour exprimer ceci on aimerait donner un intervalle qui serait plus large quand n=10 que quand $n=10^5$, pour expliciter l'incertitude liée à $\tilde{\theta}_n$

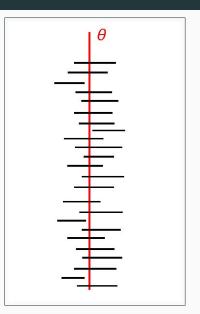
Définition: Soient $Y \equiv Y_1, \ldots, Y_n$ des données issues d'une loi F de paramètre $\theta \in \mathbb{R}$. Un **intervalle de confiance** (IC, 'confidence interval' en anglais) (L_n, U_n) pour θ est une statistique sous forme d'intervalle qui contient θ avec une probabilité spécifiée $1 - \alpha$

$$\Pr_{\theta}(L_n \leq \theta \leq U_n) = 1 - \alpha \quad \forall \theta$$

- $1-\alpha \in (0,1)$ est le **niveau**, souvent $\alpha \in \{0.05, 0.01, 0.1\}$
- Les bornes $L_n = L_n(Y)$, $U_n = U_n(Y)$ sont des statistiques et non pas des inconnus
- Un IC bilatéral, de la forme (L_n, U_n) , est le plus souvent utilisé
- Un IC unilatéral à droite est $(-\infty, U_n]$ tel que $\Pr_{\theta}(U_n \ge \theta) = 1 \alpha$
- Un IC unilatéral à gauche est $[L_n, \infty)$ tel que $\Pr_{\theta}(L_n \leq \theta) = 1 \alpha$

Interprétation d'un intervalle de confiance

- (L_n, U_n) est un intervalle aléatoire qui contient θ avec probabilité $(1 - \alpha)$
- Si on répète l'expérience avec d'autres données, on aura un autre intervalle de confiance
- Nous ne savons pas si notre IC contient θ , mais cette procédure nous fournit une garantie statistique que cet événement a une probabilité $(1-\alpha)$



Constuire un IC dans le cas normal

• Soient $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, 1)$, et

$$\overline{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$$

la moyenne de l'échantillon. On peut montrer que $\overline{Y}_n \sim \mathcal{N}(\mu, 1/n)$. Donc

$$Z_n = n^{1/2}(\overline{Y}_n - \mu) \sim \mathcal{N}(0, 1)$$

a une distribution qui ne dépend pas de μ . On obtient

$$\Pr_{\mu}\left(\overline{Y}_n - \frac{z_{1-\alpha/2}}{\sqrt{n}} \le \mu \le \overline{Y}_n + \frac{z_{1-\alpha/2}}{\sqrt{n}}\right) = 1 - \alpha$$

avec $z_{\beta} = \Phi^{-1}(\beta)$ les quantiles de la loi $\mathcal{N}(0,1)$

• Si $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, et on connaît σ^2 , alors

$$\frac{Z_n}{\sigma} = \frac{n^{1/2}(\overline{Y}_n - \mu)}{\sigma} \sim \mathcal{N}(0, 1)$$

et l'intervalle de confiance pour μ est

$$\left[\overline{Y}_n - \frac{z_{1-\alpha/2}}{\sqrt{n}}\sigma, \overline{Y}_n + \frac{z_{1-\alpha/2}}{\sqrt{n}}\sigma\right]$$

171

La largeur \downarrow avec n, \uparrow avec σ et avec le niveau $(1-\alpha)$, car $\Pr_{\theta}(L_n \leq \theta \leq U_n) \uparrow$

Exemple

Exemple On suppose que la résistance Y d'un certain type d'équipements électriques est distribuée selon une loi normale avec $\sigma = 0.12$ ohm.

Un échantillon de taille n=64 a donné comme moyenne la valeur $\bar{y}_n=5.34$ ohm.

Trouver un intervalle de confiance pour μ au niveau 95%.

Intervalle de Student

- Soient $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$ avec σ^2 inconnue
- $\qquad \qquad \boxed{\overline{Y}_n \frac{z_{1-\alpha/2}}{\sqrt{n}}\sigma, \overline{Y}_n + \frac{z_{1-\alpha/2}}{\sqrt{n}}\sigma} \text{ n'est plus un intervalle de confiance}$
- On suppose n > 1 et estime σ^2 par

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \overline{Y}_n)^2$$

• On remplace σ^2 par S_n : soit

$$T_n = \frac{Z_n}{S_n} = n^{1/2} \frac{\overline{Y}_n - \mu}{S_n} = \frac{Z_n/\sigma}{\sqrt{S_n^2/\sigma^2}}$$

- Nominateur et dénominateur indép, leurs lois ne dépendent pas de (μ,σ^2)
- T_n suit une loi de Student avec n-1 degrés de liberté, $T_n \sim t_{n-1}$, qui ne dépend de $\theta = (\mu, \sigma^2)$. C'est une loi symétrique, on dénote les quantiles par $t_{n-1,\beta}$
- L'intervalle de confiance qui en résulte est

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

■ La largeur \downarrow avec n, \uparrow avec S_n et avec le niveau $(1 - \alpha)$

Exemple

Exemple Pour déterminer le point de fusion μ d'un certain alliage, on a procédé à n=9 observations qui ont donné une moyenne $\overline{y}_n=1040^\circ$ avec $s_n=16^\circ$.

Trouver un intervalle de confiance pour μ à niveau 95%.

Intervalles de confiance pour μ dans le cas normal : résumé

Soient $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$ avec μ inconnu

• Pour σ connue $n^{1/2} \frac{\overline{Y}_n - \mu}{\sigma} \sim \mathcal{N}(0, 1)$ et on a les IC de niveau $1 - \alpha$

$$\left[\overline{Y}_n - \frac{z_{1-\alpha/2}}{\sqrt{n}} \sigma, \overline{Y}_n + \frac{z_{1-\alpha/2}}{\sqrt{n}} \sigma \right]$$

$$\left[\overline{Y}_n - \frac{z_{1-\alpha}}{\sqrt{n}} \sigma, \infty \right] \qquad \left[-\infty, \overline{Y}_n + \frac{z_{1-\alpha}}{\sqrt{n}} \sigma \right]$$

avec $z_{\beta} = \Phi^{-1}(\beta)$ les quantiles de la loi $\mathcal{N}(0,1)$

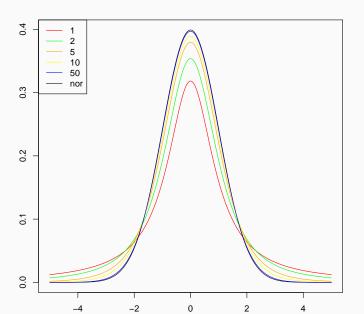
• Pour σ inconnue (et n>1) $n^{1/2} \frac{Y_n-\mu}{S_n} \sim t_{n-1}$ et on a les IC de niveau $1-\alpha$

$$\begin{bmatrix} \overline{Y}_n - \frac{t_{n-1,1-\alpha/2}}{\sqrt{n}} S_n, \overline{Y}_n + \frac{t_{n-1,1-\alpha/2}}{\sqrt{n}} S_n \end{bmatrix}$$
$$\begin{bmatrix} \overline{Y}_n - \frac{t_{n-1,1-\alpha}}{\sqrt{n}} S_n, \infty \end{bmatrix} \quad \text{ou} \quad \begin{bmatrix} -\infty, \overline{Y}_n + \frac{t_{n-1,1-\alpha}}{\sqrt{n}} S_n \end{bmatrix}$$

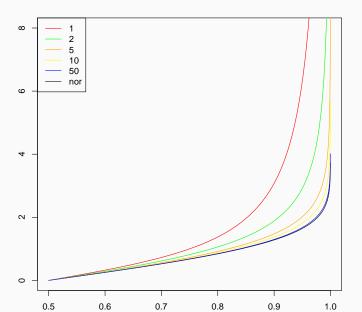
avec t_{n-1} la **loi de Student avec** n-1 **degrés de liberté**, $t_{n-1,\beta}$ sont les quantiles de cette distribution, S_n définie à la diapositive 174.

 t_{n-1} symétrique comme $\mathcal{N}(0,1)$, mais quantiles plus grands : $|t_{n-1,\beta}| > |z_{\beta}|, \ \beta \neq \frac{1}{2}$

Densités de t_k **et** $\mathcal{N}(0,1)$



Quantiles de t_k et $\mathcal{N}(0,1)$



Pivots

- Une fonction $T(Y_1, \ldots, Y_n, \theta)$ dont la loi est connue et ne dépend pas de θ s'appelle un **pivot**
- **Attention!** Ce n'est pas une statistique car c'est une fonction de θ
- Si a ≤ b,

$$\alpha_1 = \Pr(T < a), \quad \alpha_2 = \Pr(T > b), \quad \alpha_1 + \alpha_2 = \alpha \in [0, 1]$$

(souvent $\alpha_1 = \alpha_2 = \alpha/2$) alors

$$\Pr_{\theta}(a \leq T \leq b) = \Pr_{\theta}(T \leq b) - \Pr_{\theta}(T < a) = (1 - \alpha_2) - \alpha_1 = 1 - \alpha_2$$

• Si on peut isoler θ , on peut trouver des variables aléatoires L et U telles que

$$\Pr_{\theta}(L \leq \theta \leq U) = 1 - \alpha$$

Exemples T_n (diapositive 174), Z_n (diapositive 171).

Pivot : cas uniforme

Exemple Soient $Y_1, \ldots, Y_n \sim U(0, \theta)$, $M_n = \max(Y_i)$. Alors $T(Y_1, \ldots, Y_n, \theta) = M_n/\theta$ est un pivot.

Intervalles de confiance approximatifs

 En pratique la plupart des intervalles de confiance se basent sur des approximations fournies par le théorème central limite, étant de la forme

$$(L_n, U_n) = (\widehat{\theta}_n - \sqrt{V_n} z_{1-\alpha/2}, \widehat{\theta}_n + \sqrt{V_n} z_{1-\alpha/2}),$$

où V_n est une estimateur de $\text{var}_{\theta}(\widehat{\theta}_n)$ dont la racine s'appelle **erreur type** (standard error) de $\widehat{\theta}_n$. Sa réalisation $v_n^{1/2}$ est aussi appelée erreur type

• L'intervalle de confiance est approximatif dans le sens que

$$\Pr_{\theta}(L_n \leq \theta \leq U_n) \to 1 - \alpha, \quad n \to \infty$$

■ Dans des modèles réguliers, la variance asymptotique de $\widehat{\theta}_{\mathrm{ML}}$ (voir diapositive 166) est $1/I_n(\theta)$, estimée par $1/J_n(\widehat{\theta}_{\mathrm{ML}})$ et donc

$$(L_n, U_n) = \left(\widehat{\theta}_{\mathrm{ML}} - z_{1-\alpha/2} / \sqrt{J_n(\widehat{\theta}_{\mathrm{ML}})}, \widehat{\theta}_{\mathrm{ML}} + z_{1-\alpha/2} / \sqrt{J_n(\widehat{\theta}_{\mathrm{ML}})}\right)$$

Exemples

Exemple Soient $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \operatorname{Poiss}(\lambda)$, $\lambda > 0$ inconnu. Trouver une erreur type pour $\widehat{\lambda}_{\mathrm{ML}}$, et ainsi donner un intervalle de confiance approximatif de niveau 90% pour λ .

Exemples

Exemple Soient $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} U(0, \theta)$ avec θ inconnu et $\overline{Y}_n = n^{-1} \sum_{j=1}^n Y_j$. Utiliser le théorème central limite pour trouver un intervalle de confiance approximatif de niveau 95% pour θ .

3.3 Tests statistiques

Démarche scientifique

Toute **démarche scientifique** s'effectue selon le même schéma. Afin d'analyser la plausibilité d'une théorie, on itère les étapes suivantes :

- Enoncé d'une hypothèse (théorie) pouvant être contredite par des données.
- Récolte de données
- Comparaison des données avec les prédictions/implications de l'hypothèse.
- Non-rejet, rejet ou modification éventuelle de l'hypothèse.

Dans un cadre statistique, en supposant que l'on dispose d'un modèle pour le phénomène étudié, on itère les étapes suivantes :

- Enoncé d'une hypothèse (typiquement sur les paramètres du modèle statistique). Cette hypothèse peut être contredite par des données (via une statistique, appelée statistique de test).
- Récolte de données
- Rejet (ou non) de l'hypothèse à partir de la comparaison entre les données et les implications de l'hypothèse. En cas d'écart, à partir de quel seuil juge-t-on cet écart significatif, i.e., suffisamment important pour justifier le rejet de l'hypothèse?

Exemple

Exemple Question: L'alcool ralentit-il les réflexes?

Afin d'étudier l'effet de l'alcool sur les réflexes, on fait passer à 14 sujets un test de dextérité avant et après qu'ils aient consommé 100 ml de vin. Leurs temps de réaction (en ms) avant et après sont donnés dans le tableau suivant :

Sujet														
Avant														
Après	55	60	68	69	70	73	74	74	75	76	76	78	81	90

Cadre statistique : [1] Hypothèse nulle et alternative

Etant donné un modèle statistique (de densité $f(x;\theta)$), nous voulons choisir entre deux théories concurrentes à propos du paramètre θ . Ces dernières forment une paire d'hypothèses :

 H_0 : l'hypothèse <u>nulle</u> vs H_1 : l'hypothèse <u>alternative</u>.

Exemple. Dans une population décrite par la loi $\mathcal{N}(\mu, \sigma^2)$, nous pouvons former des hypothèses sur μ comme suit :

$$\underbrace{\left\{ \begin{array}{l} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{array} \right\}}_{\text{paire bilatérale}} \quad \text{ou} \quad \underbrace{\left\{ \begin{array}{l} H_0: \mu = \mu_0 \\ H_1: \mu > \mu_0 \end{array} \right\}}_{\text{paires unilatérales}} \quad \text{ou} \quad \underbrace{\left\{ \begin{array}{l} H_0: \mu = \mu_0 \\ H_1: \mu < \mu_0 \end{array} \right\}}_{\text{paires unilatérales}}.$$

Cadre statistique : [2] Statistique de test

Comment choisir entre les deux hypothèses?

- Nous tirons un échantillon $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} f(x; \theta)$ tiré de la population. Comment l'utiliser pour prendre notre décision?
- Nous choisissons une statistique $T = T_n = g(X_1, ..., X_n)$ qui a tendance à prendre des valeurs "typiques" sous l'hypothèse nulle H_0 (i.e., si H_0 est vraie) et "extrêmes" (dans la direction de l'hypothèse alternative H_1) sous H_1
- Ainsi, si on observe une valeur plutôt "extrême" ("extrême" dans la direction de l'hypothèse alternative H_1) de T, nous avons de l'évidence contre H_0 .

Notre règle de décision est donc :

- Rejeter H₀ si la valeur observée de T est assez extrême (au-delà d'une valeur critique à déterminer).
- Ne pas rejeter H_0 si la valeur observée de T n'est pas assez extrême.

Exemple : paire bilatérale

Soient $X_1,...,X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu,\sigma^2)$, où σ^2 est inconnue, et considérons la paire d'hypothèses :

$$\left\{ \begin{array}{l} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{array} \right\}.$$

On parle de paire bilatérale car $\mu \neq \mu_0$ est équivalent à $\mu < \mu_0$ ou $\mu > \mu_0$.

Considérons la statistique de test $T = \frac{X - \mu_0}{S/\sqrt{n}}$.

- Si H_0 est vraie, alors $T \sim t_{n-1}$ (donc si H_0 est vraie, T prend typiquement des valeurs proches de 0).
- Compte tenu de H₁, nous considérons donc les valeurs de T comme "extrêmes" si elles sont "éloignées" de 0. Notons qu'ici, la notion d'"extrême" dans la direction de l'hypothèse alternative H₁ signifie une valeur "extrême" de la valeur absolue de T.
- Nous allons rejeter H_0 si |T| est suffisamment élevée, i.e., $|T| > v^*$, où $v^* > 0$ est une valeur critique à déterminer.

Exemple : paire unilatérale

Soient $X_1,...,X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu,\sigma^2)$, où σ^2 est inconnu, et considérons la paire d'hypothèses :

$$\left\{ \begin{array}{l} H_0: \mu = \mu_0 \\ H_1: \mu < \mu_0 \end{array} \right\}.$$

Considérons la statistique de test $T_n = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$.

- Si H_0 est vraie, alors $T \sim t_{n-1}$.
- Compte tenu de H₁, nous considérons donc les valeurs de T comme "extrêmes" si elles sont fortement négatives. Donc ici, la notion d'"extrême" dans la direction de l'hypothèse alternative H₁ signifie une valeur "extrême" de | min(T,0)| et non de |T|.
- Nous allons donc rejeter H_0 si T est suffisamment négative, i.e., $T < v_*$, où $v_* < 0$ est la valeur critique à déterminer.

Choix de la valeur critique (par exemple v^* et v_*) : Comment définir suffisamment élevée ou suffisamment négative. En d'autres termes, quelle ampleur est considérée comme significative?

Pour répondre à cette question, il faut considérer les deux types d'erreurs que l'on peut commettre lorsque l'on se décide en faveur de l'une des hypothèses :

Décision \ Vérité	H₀ vraie	H₀ fausse
Non-rejet de <i>H</i> ₀	<u>"</u>	Erreur de type II
Rejet de H ₀	Erreur de type I	<u></u>

Erreur de type I (faux positif) considérée plus grave que l'erreur de type II (faux négatif) — filtre de spam

On ne peut pas contrôler les deux erreurs à la fois :

Pr(erreur type I) petite

 \updownarrow

rejet uniquement pour des valeurs très extrêmes

1

difficile de rejetter

⇕

Pr(erreur type II) grande

- ullet L'asymétrie entre la gravité des erreurs nous aide à choisir H_0 et H_1
- On va contrôler l'erreur de type I, qui est plus grave
- Parfois on choisit H₀ par convenience mathématique (de sorte que mathématiquement il serait plus facile de contrôler l'erreur de type I)

- Nous choisissons la valeur maximale que l'on tolère pour la probabilité d'erreur de type I (éventuellement en tenant compte de l'avis d'un spécialiste). Cette quantité est notée α et appelée **niveau/seuil de significativité du test**; $\alpha \in (0,1)$. On choisit généralement une valeur faible pour α . Typiquement, $\alpha = 0.1, 0.05, 0.01, 0.001$; le plus souvent, $\alpha = 0.05$.
- La valeur critique est déterminée de manière à ce que

$$\Pr_{H_0}[\text{Rejet de } H_0] = \alpha.$$

Ainsi, la valeur critique est telle que

$$\Pr_{H_0}[|T| > \text{valeur critique}] = \alpha \quad \text{(cas bilatéral)},$$

$$\Pr_{H_0}[T < \text{valeur critique}] = \alpha \quad \text{ou}$$

$$\Pr_{H_0}[T > \text{valeur critique}] = \alpha \quad \text{(cas unilatéral)}.$$

Les probabilités sont sous l'hypothèse que H₀ est vraie!

Exemple, paire bilatérale : Soient $X_1,...,X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu,\sigma^2)$, où σ^2 est inconnu, et considérons la paire $H_0: \mu = \mu_0$ contre $H_1: \mu \neq \mu_0$.

Nous allons rejeter
$$H_0$$
 si $|T| = \left| \frac{\bar{X} - \mu_0}{S / \sqrt{n}} \right|$ est assez large, c'est à dire $|T| > v^*$.

Soit α le niveau de significativité. La valeur critique v^* satisfait

$$\Pr_{H_0}[|T| > \mathbf{v}^*] = \alpha,$$

i.e.,

$$\Pr_{H_0}[T < -\mathbf{v}^* \text{ ou } T > \mathbf{v}^*] = \alpha.$$

Quand H_0 **est vraie** $T \sim t_{n-1}$. On doit donc choisir

$$v^*=t_{n-1,1-\alpha/2},$$

où $t_{n-1,1-\alpha/2}$ est le $(1-\alpha/2)$ quantile de la loi de Student t_{n-1} .

Cadre statistique : [4] La p-valeur

Au lieu d'utiliser des valeurs critiques pour choisir entre H_0 et H_1 , nous pouvons utiliser une autre approche, basée sur la notion de p-valeur.

- La p-valeur (notée pobs) est la probabilité d'obtenir une valeur de la statistique de test au moins aussi élevée (élevée dans la direction de H₁) que celle que nous avons observée si H₀ était vraie.
- Supposons que la réalisation de la statistique de test sur nos données est $T=t_{
 m obs}.$ Alors :
 - Cas bilatéral : $p_{\mathrm{obs}} = \mathrm{Pr}_{H_0}[|T| > |t_{\mathrm{obs}}|]$,
 - Cas unilatéral à gauche : $p_{
 m obs} = \Pr_{H_0}[T < t_{
 m obs}]$,
 - Cas unilatéral à droite : $p_{\text{obs}} = \Pr_{H_0}[T > t_{\text{obs}}]$.
- Petite valeurs de $p_{\rm obs}$ s'opposent à H_0 car elles démontrent que la realité observée serait très improbable si l'hypothèse nulle H_0 était vraie.
- Cas bilatéral (195) : $p_{obs} = 2(1 F_{t_{n-1}}(|t_{obs}|))$, où $F_{t_{n-1}}$ est la fonction de répartition de la loi de Student t_{n-1} .

Cadre statistique : [4] La p-valeur

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

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

Exemple bilatérale (195) $p_{obs} = 2(1 - F_{t_{n-1}}(t_{obs}))$ 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 $p_{\rm obs}$ fournit une information plus facilement interprétable que la valeur $t_{\rm obs}$. 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 $t_{\rm obs}$.
- 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 $p_{\rm obs}$ donnant la probabilité d'observer une valeur de T aussi élevée que $t_{\rm obs}$ sous H_0 . On rejette alors H_0 en faveur de H_1 quand $p_{\rm obs} < \alpha$.

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 200

Supposons que nos observations x_1,\ldots,x_n soient des réalisations de variables aléatoires $X_1,\ldots,X_n\stackrel{\mathrm{iid}}{\sim}\mathcal{N}(\mu,\sigma^2)$, 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 $t_{obs} = -2.35$.

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

Intervalles de confiance et tests

- Soient $Y_1, \ldots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$ 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

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

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

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

• Cet intervalle contient zéro si et seulement si

$$|\overline{Y}_n| \leq t_{n-1,1-\alpha/2} S_n / \sqrt{n}$$
 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 rejetter \iff IC de niveau $1-\alpha$ large

Intervalles de confiance (IC) et tests

- Rejet $\iff p_{obs} \leq \alpha \iff \theta_0 \notin IC$ de niveau 1α
- 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 ?

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 o₁,..., o_k, alors que les
 - **fréquences théoriques** correspondantes sont e_1, \ldots, 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 : $e_i \ge 5$ pour la plupart), alors $T_n \stackrel{\text{app}}{\sim} \chi^2_{\nu}$ sous H_0 , où
 - χ^2_{ν} 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

rejet de
$$H_0$$
 au niveau $\alpha \iff t_{\text{obs}} = \sum_{i=1}^{\kappa} \frac{(o_i - e_i)^2}{e_i} > \chi^2_{\nu, 1-\alpha}$

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

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

$$Pr(X = x) = 1/6, \quad x \in \{1, \dots, 6\}.$$

Exemple

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

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

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

On a

$$\begin{split} e_i &= n \mathrm{Pr}_{\mathcal{N}(100,15^2)} \big(a_i \leq X \leq b_i \big) = n \mathrm{Pr} \left(\frac{a_i - 100}{15} \leq \frac{X_i - 100}{15} \leq \frac{b_i - 100}{15} \right) \\ &= n \Phi \left(\frac{b_i - 100}{15} \right) - n \Phi \left(\frac{a_i - 100}{15} \right) \end{split}$$

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 :

	В						
Α	1	2		j		k	Σ
1	n ₁₁	n_{12}		n_{1j}		n_{1k}	n_1 .
2	n ₂₁	n_{22}		n_{2j}	• • •	n_{2k}	n_2 .
:	:	:	:	:	:	:	:
i	n _{i1}	n_{i2}		n _{ij}		n _{ik}	n_i .
:	:	:	:	:	:	:	:
h	n_{h1}	n_{h2}		n_{hj}		n_{hk}	n _h .
Σ	n.1	n. ₂		n.j		$n_{\cdot k}$	n = n

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

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

On désire tester H₀: A et B sont indépendants. Dans ce cas

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

et les probabilités empiriques sont

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

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

$$e_{ij} = n \frac{n_i}{n} \frac{n_{ij}}{n} = \frac{n_i \cdot n_{ij}}{n}$$

Calcul sous H₀

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 $t_{\rm obs} > \chi^2_{(h-1)(k-1),1-\alpha}$, 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 :

	I		
Α	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	$n_{\cdot 1} = 52$	$n_{.2} = 43$	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

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

Preuve (non-examinable)

Les fonctions de densité et de répartition de y_j sont

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

L'indépendance donne

$$L(\theta) = \prod \theta^{-1} I(Y_j < \theta) = \theta^{-n} I(\max Y_j \le \theta), \quad \theta > 0,$$
 qui est maximisée au point $\widehat{\theta}_{\mathrm{ML}} = M_n = \max Y_i$

On a

$$\Pr(M_n \le x) = \prod_{i=1}^n \Pr(Y_i \le x) = (x/\theta)^n, \quad 0 \le x \le \theta$$

• Donc pour $x \ge 1$ et $n \ge x$,

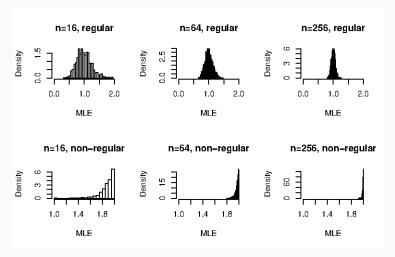
$$\Pr\left\{n(\theta-\widehat{\theta}_{\mathrm{ML}})/\theta \leq x\right\} = \Pr(\widehat{\theta}_{\mathrm{ML}} \geq \theta - x\theta/n) = 1 - \{(\theta - x\theta/n)/\theta\}^n$$

$$ightarrow 1 - \exp(-x),$$

comme requis

Exemple (non examinable)

Comparaison des lois de $\widehat{\theta}$ 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 explanatoires) sont notées $x^{(1)}, \ldots, x^{(p)}$, on les considère comme fixées

On peut s'interesser à

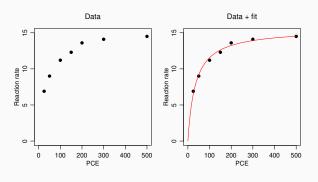
- l'estimation d'une relation eventuelle entre y et les $x^{(j)}$, ou
- la prévision des valeurs futures/manquantes de y sur la base des x^(j) correspondantes

Réaction chimique

Professeur Christophe Holliger (SIE) : on essaye de déterminer les paramètres kinétiques 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

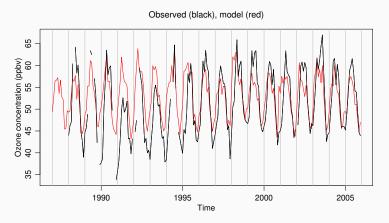
$$y = \frac{\gamma_0 x}{\gamma_1 + x},$$

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?



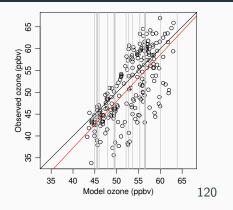
Ozone atmosphérique

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



Soient y les données réeles et x les resultats 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 resultats du modèle x et les données observées y pour estimer les y manquants?

Problème d'ajustement

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

$$\left(\begin{array}{c}x_1\\y_1\end{array}\right),\ldots,\left(\begin{array}{c}x_n\\y_n\end{array}\right)$$

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

Moindres carrés

Les écarts verticaux entre les données y_j et la courbe $\mu(x_j)$ sont

$$y_j - \mu(x_j), \quad j = 1, \ldots, n$$

 On cherche les paramètres de la fonction μ(x) pour minimiser la somme des carrés des écarts verticaux

$$\sum_{j=1}^{n} \{y_j - \mu(x_j)\}^2$$

• L'ajustement est dit linéaire si $\mu(x) = a + \beta x$. Dans ce cas, il faut minimiser

$$S(a,\beta) = \sum_{j=1}^{n} \{y_j - \mu(x_j)\}^2 = \sum_{j=1}^{n} \{y_j - (a + \beta x_j)\}^2$$

Estimateurs de moindres carrés

Théorème Soient $(x_1, y_1), \ldots, (x_n, y_n)$ issues d'un relation $y = a + \beta x$ et telles que pas tous les x_j sont égaux. Alors les **estimateurs de moindres carrés** de a et β sont

$$\widehat{a}_n = \overline{y}_n - \widehat{\beta}_n \overline{x}_n, \quad \widehat{\beta}_n = \frac{\sum_{j=1}^n (x_j - \overline{x}_n) y_j}{\sum_{j=1}^n (x_j - \overline{x}_n)^2}$$

Définition: La droite

$$\widehat{a}_n + \widehat{\beta}_n x$$

s'appelle la **droite des moindres carrés**, la **valeur ajustée** qui correspond à (x_i, y_i) est

$$\widehat{y}_j = \widehat{a}_n + \widehat{\beta}_n x_j,$$

et la différence

$$r_j = y_j - \widehat{y}_j = y_j - (\widehat{a}_n + \widehat{\beta}_n x_j)$$

s'appelle un résidu

Preuve

Il faut minimiser

$$S(a,\beta) = \sum_{i=1}^{n} (y_i - a - \beta x_i)^2$$

en a et β . On calcule

$$\frac{dS}{da}(a,\beta) = -2\sum_{i=1}^{n} (y_i - a - \beta x_i) = 2na + 2n\beta \overline{x}_n - 2n\overline{y}_i$$

$$\frac{dS}{d\beta}(a,\beta) = -2\sum_{i=1}^{n} x_i (y_i - a - \beta x_i) = 2na\overline{x}_i + 2\beta \sum_{i=1}^{n} x_i^2 - 2\sum_{i=1}^{n} x_i y_i$$

$$\frac{d^2S}{da^2}(a,\beta) = 2n > 0 \qquad \frac{d^2S}{d\beta^2}(a,\beta) = 2\sum_{i=1}^{n} x_i^2 \qquad \frac{d^2S}{d\beta da}(a,\beta) = 2n\overline{x}_n$$

La matrice hessienne est donc

$$H = \begin{pmatrix} 2n & 2n\overline{x}_n \\ 2n\overline{x}_n & 2\sum_{i=1}^n x_i^2 \end{pmatrix} \qquad \text{d\'efinie positive}$$

car 2n > 0 et $\det(H) = 4n[(\sum_{i=1}^{n} x_i^2) - n\overline{x}_n] = 4n \sum_{i=1}^{n} (x_i - \overline{x}_n)^2 > 0$ 22

Propriétés

- La droite de moindres carrés passe par $(\overline{x}_n, \overline{y}_n)$
- $\sum_{j=1}^{n} x_j r_j = \sum_{j=1}^{n} x_j (y_j \widehat{y}_j) = 0$

(voir série d'exercices). Donc

$$\sum_{j=1}^{n} (y_j - \bar{y}_n)^2 = \sum_{j=1}^{n} \left(\underbrace{y_j - \hat{y}_j}_{r_j} + \widehat{y}_j - \bar{y}_n \right)^2 = \cdots = \sum_{j=1}^{n} (\widehat{y}_j - \bar{y}_n)^2 + \sum_{j=1}^{n} r_j^2,$$

nous donnant la décomposition de la somme des carrés total

$$SC_{Total} = SC_R + SC_E$$

en une partie due à la régression (variation expliquée par le modèle) et une partie due à l'erreur (variation non-expliquée par le modèle)

Ozone atmosphérique

- Il y a n = 207 paires (Observée, Modèle) = (y_j, x_j) , et en plus 21 valeurs de x sans valeur observée
- Avec les n paires complètes on trouve comme droite des moindres carrés

$$\widehat{y} = \widehat{a}_n + \widehat{\beta}_n x = -5.511 + 1.069x$$

avec décomposition de la somme des carrés

$$\mathrm{SC}_{\mathrm{Total}} = \mathrm{SC}_{\mathrm{R}} + \mathrm{SC}_{\mathrm{E}} = 5813 + 5832.$$

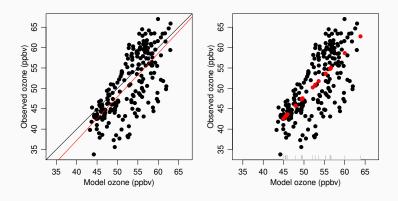
La régression explique donc une moitié de la somme des carrés totale

 Pour un paire (Observée, Modèle) = (?, x₊) dont la valeur observée manque, on peut la remplacer par la valeur ajustée correspondante,

$$\widehat{y}_{+}=\widehat{a}_{n}+\widehat{\beta}_{n}x_{+}.$$

On parle d'imputation de donnée

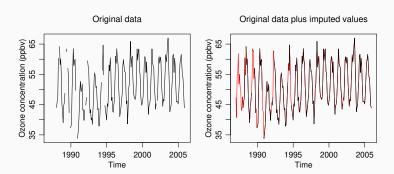
Modèle ajusté



Gauche : droite y=x et droite ajustée $\widehat{y}=\widehat{a}_n+\widehat{\beta}_n x=-5.511+1.069x$

Droite : valeurs ajustées pour des valeurs manquantes de x

Données imputées



Gauche : données originales

Droite : données originales (noir) avec valeurs imputées (rouge). Comparer avec la diapositive 230

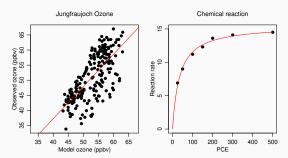
4.2 Modèle statistique

Modèle normale

- On observe une version perturbée d'une relation $y = \mu(x)$
- Pour modéliser ceci, on peut souvent supposer que

$$y_j \overset{\mathrm{ind}}{\sim} \mathcal{N}\left\{\mu(x_j), \sigma^2\right\} \quad \text{ou bien} \quad y_j = \mu(x_j) + \epsilon_j, \quad \epsilon_j \overset{\mathrm{ind}}{\sim} \mathcal{N}(0, \sigma^2)$$

- Ainsi la dépendance entre la réponse y et la variable explicative x est donnée par $\mathbb{E}(y) = \mu(x)$, alors que le bruit dépend de σ^2
- À gauche : $\mu(x)$ linéaire, σ^2 grand, donc beaucoup de bruit
- À droite : $\mu(x)$ non-linéaire, σ^2 petite, donc peu de bruit



Linéarité

La linéarité du modèle concerne les paramètres :

$$y = a + \beta x + \epsilon,$$

où $\epsilon \sim \mathcal{N}(0, \sigma^2)$ est la différence entre y et la droite $a + \beta x$

■ Le modèle

$$y = a + \beta x + \gamma x^2 + \delta x^3 + \epsilon$$

est linéaire en $(a, \beta, \gamma, \delta)$.

■ Le modèle

$$y = \gamma_0 x^{\gamma_1} \eta, \quad \eta \sim \exp(1),$$

devient linéaire après transformation logarithmique :

$$\log y = \log \gamma_0 + \gamma_1 \log x + \log \eta = a + \beta x' + \log \eta$$

Le modèle

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

n'est pas linéaire en les paramètres γ_0, γ_1

Estimation des paramètres

- Dans le cas $\mu(x) = a + \beta x$ il y a trois paramètres inconnus : (intercepte, pente, bruit), $\theta = (a, \beta, \sigma^2) \in \mathbb{R} \times \mathbb{R} \times \mathbb{R}_+$
- Nous utilisons la méthode de maximum de vraisemblance pour les estimer
- La log vraisemblance est

$$\ell(a, \beta, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^{n} \{y_j - (a + \beta x_j)\}^2 - \frac{n}{2} \log(2\pi),$$

et en maximisant celle-ci par rapport à heta nous trouvons

$$\widehat{a}_n = \overline{y}_n - \widehat{\beta}_n \overline{x}_n, \quad \widehat{\beta}_n = \frac{\sum_{j=1}^n (x_j - \overline{x}_n) y_j}{\sum_{j=1}^n (x_j - \overline{x}_n)^2}, \quad \widehat{\sigma}_n^2 = n^{-1} \sum_{j=1}^n r_j^2$$

avec $r_j = y_j - \widehat{y}_j$ les **résidus** et $\widehat{y}_j = \widehat{a}_n + \widehat{\beta}_n x_j$ les **valeurs ajustées**

Les estimateurs \widehat{a}_n et $\widehat{\beta}_n$ sont les estimateurs de moindres carrés et sont sans biais, mais $\mathbb{E}(\widehat{\sigma}_n^2) < \sigma^2$, et on utilise souvent l'estimateur non-biaisé (comparer avec 174)

$$S_n^2 = \frac{1}{n-2} \sum_{i=1}^n r_i^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Inférence pour les paramètres du modèle linéaire simple

- Le coefficient β (pente) est plus intéressant que a (ordonnée à l'origine). On se concentre donc ici sur le premier
- On peut montrer que

$$\operatorname{Var}(\widehat{\beta}_n) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \overline{x}_n)^2}$$

• On estime σ^2 par S^2 pour estimer cette variance. En prenant la racine carrée on obtient **l'erreur type** (standard error)

$$\widehat{\mathrm{sd}}(\widehat{\beta}_n) = \frac{S_n}{\sqrt{\sum_{i=1}^n (x_i - \overline{x}_n)^2}}$$

On peut montrer que

$$\frac{\widehat{\beta_n} - \beta}{\widehat{\mathrm{sd}}(\widehat{\beta_n})} \sim t_{n-2}$$

On a donc un pivot. On peut construire des intervalles de confiance et tester des hypothèses

Intervalles de confiance pour β

On en déduit des intervalles de confiance pour β au niveau de confiance $1-\alpha$, pour $\alpha \in (0,1)$:

• Intervalle de confiance bilatéral symétrique :

$$\left[\widehat{\beta}_n - t_{n-2,1-\alpha/2} \frac{S_n}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}, \widehat{\beta}_n + t_{n-2,1-\alpha/2} \frac{S_n}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}\right].$$

Intervalle de confiance unilatéral à gauche :

$$\left(-\infty,\widehat{\beta}_n+t_{n-2,1-\alpha}\frac{S_n}{\sqrt{\sum_{i=1}^n(x_i-\bar{x})^2}}\right].$$

Intervalle de confiance unilatéral à droite :

$$\left[\widehat{\beta}_n - t_{n-2,1-\alpha} \frac{S_n}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}, \infty\right).$$

Comparer avec diapositive 186 : $[\widehat{\theta} \pm t_{k,1-\alpha/2}\widehat{sd}(\widehat{\theta})]$: en 186, $k=n-1,\widehat{\theta}=\overline{Y}_n$, ici $k=n-2,\widehat{\theta}=\widehat{\beta}_n$

Tests pour β

On peut effectuer les tests statistiques classiques au niveau de significativité α , pour $\alpha \in (0,1)$:

- Test bilatéral $H_0: \beta = \beta_0$ contre $H_1: \beta \neq \beta_0$. On rejette H_0 si et seulement si $|t_{\text{obs}}| > t_{n-2,1-\alpha/2}$.
- Test unilatéral à gauche $H_0: \beta = \beta_0$ contre $H_1: \beta < \beta_0$. On rejette H_0 si et seulement si $t_{\text{obs}} < t_{n-2,1-\alpha}$.
- Test unilatéral à droite $H_0: \beta = \beta_0$ contre $H_1: \beta > \beta_0$. On rejette H_0 si et seulement si $t_{\text{obs}} > t_{n-2,1-\alpha}$.

La statistique de test est

$$T = \frac{\widehat{\beta_n} - \beta_0}{\widehat{\mathrm{sd}}(\widehat{\beta_n})} = \frac{\widehat{\beta_n} - \beta_0}{S_n / \sqrt{\sum_{i=1}^n (x_i - \overline{x}_n)^2}}$$

qui suit la loi t_{n-2} quand H_0 est vraie

Nos données

```
> JungOzone
  Observed Model
1
        NA 49.42
      40.7 52.79
3
      NA 56.49
      NA 56.61
      61.8 57.22
     NA 53.59
7
      NA 56.61
8
       NA 52.75
        NA 52.15
10
        NA 45.43
> MM <- data.frame(
+ x=c(25, 50, 100, 150, 200, 300, 500),
+ y=c(6.9, 9.0, 11.2, 12.3, 13.6, 14.1, 14.5))
> MM
   х
       У
1 25 6.9
2 50 9.0
3 100 11.2
4 150 12.3
5 200 13.6
6 300 14.1
7 500 14.5
```

Inférence

Voici le résultat de l'ajustement du modèle linéaire aux données d'ozone :

> fit <- lm(Observed~Model,data=JungOzone)</pre>

```
> summary(fit)
. . .
Coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) -5.51072 3.98014 -1.385 0.168
Model 1.06903 0.07479 14.294 <2e-16 ***
Signif. codes: 0 "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1 " " 1
Residual standard error: 5.334 on 205 degrees of freedom
  (21 observations deleted due to missingness)
Multiple R-Squared: 0.4992, Adjusted R-squared: 0.4967
F-statistic: 204.3 on 1 and 205 DF, p-value: < 2.2e-16
```

Exemple : données d'ozone (inférence)

- On sait d'après les slides précédentes que l'intervalle de confiance bilatéral symétrique pour β au niveau de confiance $1-\alpha$ est $\left[\hat{\beta}_n-t_{n-2,1-\alpha/2}\hat{\mathrm{sd}}(\hat{\beta}_n),\hat{\beta}_n+t_{n-2,1-\alpha/2}\hat{\mathrm{sd}}(\hat{\beta}_n)\right]$.
- Ainsi, en lisant les sorties du logiciel, on obtient qu'une réalisation de l'IC précédent pour β au niveau de confiance 95% est donnée par $1.06903 \pm t_{205,0.975} \times 0.07479 \approx 1.07 \pm 1.97 \times 0.07 = [0.93, 1.21].$
- Souvent, on veut tester si le terme impliquant la covariable est significatif. Cela revient à tester H_0 : $\beta=0$.
- Ici, le scatter plot semble clairement indiquer que β est différent de 0 et on effectue donc plutôt le test $H_0: \beta=1$. On choisit comme niveau de significativité $\alpha=0.05$. On rejette H_0 si et seulement si la valeur absolue de la réalisation $t_{\rm obs}$ de $\hat{\beta}=1$

$$T = \frac{\beta_n - 1}{\hat{\mathrm{sd}}(\hat{\beta}_n)}$$

est strictement supérieure à $t_{n-2,1-\alpha/2}=t_{205,0.975}\approx 1.97$. On a $t_{\rm obs}\approx 0.92$ et on ne rejette donc pas H_0 .

Modèle nonlinéaire (non-examinable)

- Les mêmes idées s'appliquent aux modèles nonlinéaires, mais comme approximations
- Il faut donner des valeurs initiales pour γ_0 et γ_1 , en principe il faut en essayer plusieurs, car il est possible que la vraisemblance ait des maximas locaux
- Pour ajuster le modèle $\mu(x)=\gamma_0x/(\gamma_1+x)$ aux données chimiques :

```
> fit <- nls(y~g0*x/(g1+x),data=MM, start=c(g0=1,g1=1))
> summary(fit)
```

```
Formula: y \sim g0 * x/(g1 + x)
```

```
Estimate Std. Error t value Pr(>|t|)
g0 15.5269    0.2876    53.99 4.12e-08 ***
g1 34.5990    2.8777    12.02 7.02e-05 ***
---
Signif. codes: 0 "***" 0.001 "**" 0.05 "." 0.1 " " 1
```

Residual standard error: 0.3341 on 5 degrees of freedom

Coefficient de détermination

Nous avons dêjà vu la décomposition de la somme des carrés total

$$\sum_{j=1}^{n}(y_j-\bar{y})^2=\sum_{j=1}^{n}(\widehat{y}_j-\bar{y})^2+\sum_{j=1}^{n}r_j^2,\quad\text{soit}\quad\mathrm{SC}_{\mathrm{Total}}=\mathrm{SC}_{\mathrm{R}}+\mathrm{SC}_{\mathrm{E}},$$

en une partie SC_R due à la régression et une partie SC_E due à l'erreur

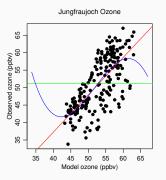
• Le proportion (ou pourcentage) de la variation totale expliquée par le modèle

$$R^2 = rac{\mathrm{SC_R}}{\mathrm{SC_{Total}}} = rac{\mathrm{SC_{Total}} - \mathrm{SC_E}}{\mathrm{SC_{Total}}}$$

est appelé coefficient de détermination ; $0 \le R^2 \le 1$

- Si $R^2 \approx 1$, alors $y_j \approx \widehat{y}_j$ pour tout j et donc tous les $r_j \approx 0$, et donc le modèle explique les données presque parfaitement
- Si $R^2 \approx 0$, alors l'inclusion de x n'explique presque rien de la variation totale
- Pour les données d'ozone, $R^2=0.5$, donc la moitié de la variance est expliquée par le modèle
- Pour les données chimiques, $R^2 = 0.99$, donc le modèle explique presque toute la variation

Comparaison des modèles



Voici trois modèles :

constant (vert) : $y = a + \epsilon$,

linéaire (rouge) : $y = a + \beta x + \epsilon$,

cubique (bleu) : $y = a + \beta x + \gamma x^2 + \delta x^3 + \epsilon$?

Le rouge semble être bien meilleur que le vert, mais que le rouge et le bleu semblent être semblables. Comment tester ces constats?
241

Décomposition de la variance

- Comparons le modèle constante $y = a + \epsilon$ et le modèle linéaire $y = a + \beta x + \epsilon$
- Pour tester s'il vaut la peine d'ajouter βx , on calcule

$$F = \frac{\mathrm{SC_R}/1}{\mathrm{SC_E}/(n-2)} \sim F_{1,n-2}$$

si l'hypothèse nulle H_0 : $\beta=0$ que le modèle est constant est vraie

- F_{d_1,d_2} est la **loi de Fisher(-Snedecor)** avec d_1 et d_2 degrés de liberté
- Pour un niveau de significativité $\alpha \in (0,1)$ donné, il faut comparer la valeur observée de F avec le $1-\alpha$ quantile $F_{1,n-2,1-\alpha}$ (rejet pour grandes valeurs de F)
- Pour les données d'ozone, on trouve $f_{obs}=204.32$, à comparer avec $F_{1,205,0.95}=3.887$
- Ce test est équivalent au t-test pour $H_0: \beta = 0$ vu précédemmant car : $T \sim t_{\nu} \Longrightarrow T^2 \sim F_{1,\nu}$

Statistique *F*

• Pour tester $H_0: \beta_{q+1} = \cdots = \beta_p = 0$ dans le modèle linéaire

$$y_i = \beta_0 + \beta_1 x_i^{(1)} + \dots + \beta_q x_i^{(q)} + \beta_{q+1} x_i^{(q+1)} + \dots + \beta_p x_i^{(p)} + \epsilon,$$

on a deux sommes des carrés, l'un $SC_{E,p}$ qui correspond au modèle avec $x^{(1)}, \ldots, x^{(p)}$ et l'autre $SC_{E,q}$ qui correspond au modèle réduit avec $x^{(1)}, \ldots, x^{(q)}$, q < p. On a $SC_{E,p} \leq SC_{E,q}$, et pour tester H_0 on calcule

$$F = \frac{(\mathrm{SC}_{E,q} - \mathrm{SC}_{E,p})/(p-q)}{\mathrm{SC}_{E,p}/(n-p-1)} \sim F_{p-q,n-p-1}$$

si H_0 : $\beta = 0$ est vraie

- On rejette H_0 au niveau lpha si $f_{obs} > F_{p-q,n-p-1,1-lpha}$
- Pour les données d'ozone, pour tester $\gamma=\delta=0$ dans le modèle cubique

$$y = a + \beta x + \gamma x^2 + \delta x^3 + \epsilon,$$

on a n = 207, p = 3, q = 1, et

$$F = \frac{(5831.9 - 5712.2)/(3 - 1)}{5712/(207 - 3 - 1)} = 2.13 \sim F_{3-1,207-3-1} = F_{2,203},$$

dont le 0.95 quantile est $F_{2,203,0.95} = 3.04$.

Validation du modèle de régression linéaire (non-examinable)

• Le modèle normale $y \sim \mathcal{N} \{\mu(x), \sigma^2\}$ implique que

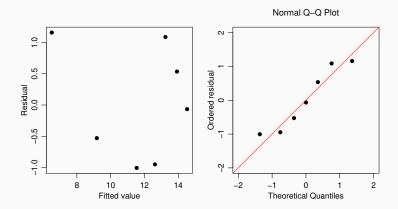
$$\frac{y - \mu(x)}{\sigma} \sim \mathcal{N}(0, 1),$$

et donc que le résidu standardisé

$$r_j^S = \frac{r_j}{s_n} = \frac{y_j - \widehat{y}_j}{s_n} = \frac{r_j}{s_n} = \frac{y_j - (\widehat{a}_n + \widehat{\beta}_n x_j)}{s_n} \stackrel{\text{app}}{\sim} \mathcal{N}(0, 1)$$

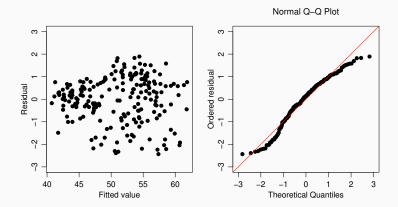
- On teste cela graphiquement avec un quantile-quantile plot (Q-Q plot) normal. C'est un graphique des quantiles empiriques des données (ici les résidus standardisés) contre les quantiles théoriques d'une loi $\mathcal{N}(0,1)$. Si les r_i^S suivent effectivement la loi $\mathcal{N}(0,1)$, alors les points du Q-Q plot doivent se trouver (plus ou moins) sur la diagonale y=x. Des écarts trop importants par rapport à la diagonale indiquent une violation de l'hypothèse de normalité des erreurs.
- Par ailleurs, il faut qu'il n'y ait pas de relation entre les r_j^S et les valeurs ajustées \hat{y}_i

Données chimiques (non-examinable)



- À gauche : r_j^S contre \widehat{y}_j
- À droite : QQplot des r_j^S
- Avec n = 7, il est presque impossible de contradire le modèle

Données d'ozone (non-examinable)



• À gauche : r_j^S contre \hat{y}_j

• À droite : QQplot des r_j^S

■ La loi des erreurs n'est pas normale, mais asymétrique, et la variance semble changer avec $\mathbb{E}(y)$