



MAP433 Statistique

PC 6: Introduction aux tests statistiques

2 octobre 2015

1 Tests d'égalité et d'inégalité

Soient (X_1, \dots, X_n) des variables gaussiennes indépendantes respectivement de loi $\mathcal{N}(\mu, \sigma^2)$.

1. On suppose σ^2 connu. Soient $\mu_0 \neq \mu_1 \in \mathbb{R}$. On considère les hypothèses

(H-0) $\mu = \mu_0$.

(H-1) $\mu = \mu_1$.

i Déterminer le test δ_α de Neyman-Pearson de niveau $\alpha \in (0, 1)$.

ii Dans le cas $\mu_0 < \mu_1$, en donner une forme simplifiée consistant à comparer la statistique

$$\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$$

à un seuil s (appelé aussi *valeur critique*) à déterminer.

iii Calculer sa puissance. Quel est son comportement quand α et n varient ?

2. On conserve la même forme de test qu'à la question précédente mais pour les hypothèses

(H'-0) $\mu \leq \mu_0$.

(H'-1) $\mu > \mu_0$.

i Déterminer le niveau α' et la puissance du test δ_α pour ces nouvelles hypothèses.

ii Soit $\delta'_{\alpha'}$ un test de niveau α' pour (H'-0) contre (H'-1). En utilisant le fait que ce test peut servir aussi pour tester (H-0) contre (H-1) avec $\mu_1 > \mu_0$ arbitraire, comparer la puissance de $\delta'_{\alpha'}$ à celle de δ_α .

iii Conclure.

3. On considère à présent les hypothèses

(H''-0) $\mu = \mu_0$.

(H''-1) $\mu \neq \mu_0$.

i Construire un test δ''_α de niveau α pour tester ces hypothèses, basé sur la même statistique de test que précédemment.

ii Calculer la puissance du test δ''_α et le comparer au test δ_α . Ce test est-il UPP ?

4. Comment adapter ces différents tests au cas où σ^2 est inconnu ?

Application interactive : <http://webpopix.org:8080/X/TestMean>

2 Test pour une certification bio

Pour avoir la certification "bio", un fabricant de produits "bios" doit garantir pour chaque lot un pourcentage d'OGM inférieur à 1%. Il prélève donc $n = 25$ produits par lot et teste si le pourcentage d'OGM est inférieur à 1%. On note X_i le logarithme du pourcentage d'OGM du paquet numéro i .

Modèle : On suppose que les X_i sont indépendants et suivent une loi gaussienne $\mathcal{N}(\theta, 1)$.

1. Pour le fabricant, le pourcentage d'OGM est inférieur à 1% sauf preuve du contraire. Il veut tester l'hypothèse $\mathbf{H}_0 : \theta \leq 0$ contre $\mathbf{H}_1 : \theta > 0$ et il souhaite que pour $\theta \leq 0$ le test se trompe avec une probabilité inférieure à 5%. Calculez un seuil $t_{25,5}$ tel que

$$\sup_{\theta \leq 0} P_{\theta}(\bar{X}_{25} > t_{25,5}) = 5\%.$$

On pourra utiliser que $P(Z > 1.645) \approx 5\%$, pour $Z \sim \mathcal{N}(0, 1)$.

2. Une association "anti-OGM" veut s'assurer qu'il n'y a effectivement pas plus de 1% d'OGM dans les produits labélisés "bio". En particulier, elle s'inquiète de savoir si le test parvient à éliminer les produits pour lesquels le pourcentage d'OGM dépasse de 50% le maximum autorisé. Quelle est la probabilité que le test ne rejette pas \mathbf{H}_0 lorsque le pourcentage d'OGM est de 1.5% ?
3. Scandalisée par le résultat précédent, l'association milite pour que le test du fabricant prouve effectivement que le pourcentage d'OGM est inférieur à 1%. Pour elle, le pourcentage d'OGM est supérieur à 1% sauf preuve du contraire, donc \mathbf{H}_0 est $\theta > 0$ et \mathbf{H}_1 est $\theta \leq 0$. Proposez un test de \mathbf{H}_0 contre \mathbf{H}_1 tel que la probabilité que le test rejette à tort \mathbf{H}_0 soit inférieure à 5%.
4. Les agences de régulation accepte le test statistique proposé par le fabricant ($\mathbf{H}_0 : \theta \leq 0$ contre $\mathbf{H}_1 : \theta > 0$) mais exige qu'un dépassement de 10% du maximum autorisé soit détecté dans 80% des cas. Le niveau du test restant fixé à 5% que doit faire le fabricant pour se soumettre à cette législation ?.

3 Des financiers sans scrupules

Roger a lu dans le journal que 20% des professionnels de la finance pensent qu'il faut enfreindre la loi pour réussir :

<http://www.slate.fr/story/101785/wall-street-enfreindre-loi-reussir>

Il pense que ce pourcentage est sous-estimé et se propose d'enquêter pour vérifier son hypothèse. Il a dans son carnet d'adresses seulement deux noms de financiers qu'il peut interroger et à qui il va poser la question suivante : "*Pensez-vous qu'il faut enfreindre la loi pour réussir ?*".

1. En supposant que la conclusion de Roger ne dépende que des réponses à cette question, quelles règles de décision peut-il mettre en place, et quelles sont alors les risques pour Roger de conclure à tort ?
2. Roger souhaite que la probabilité de conclure à tort que ce pourcentage est supérieur à 20% soit exactement de 5%. Il réalise alors qu'il n'a pas besoin d'interroger qui que ce soit pour que cette contrainte soit satisfaite : il lui suffit de mettre dans sa poche 19 jetons noir et un rouge. Pourquoi ? Quel est le défaut de cette méthode ?

3. Roger choisit maintenant une carte dans un jeu de 32 cartes et note soigneusement le résultat avant de faire son enquête auprès de ses deux connaissances. Pourquoi a-t-il adopté cette nouvelle stratégie ? Quelles sont alors les risques pour Roger de conclure à tort ?
4. Roger sent bien qu'une telle décision ne peut être prise sur la base d'un tirage aléatoire. . . Que devrait-il faire alors ?

4 Piscine

Un constructeur de piscine veut comparer 2 produits à dissoudre dans l'eau qui permettent de tuer les bactéries présentes. Les deux produits garantissent que 95 % des bactéries seront éliminées. Par contre il se peut que le pH soit plus basique et donc plus nocif pour les utilisateurs. Il mène une expérience dans son magasin où 2 piscines remplies sont exposées. Il met le produit A dans la piscine 1 et le produit B dans la piscine 2. Il a un pH-mètre vendu dans le commerce et il se doute que la mesure n'est pas fiable. Il prend donc 10 mesures dans chaque piscine.

Dans la piscine 1, il trouve x_1, \dots, x_{10} qui valent respectivement

7.33 ; 6.17 ; 7.46 ; 8.13 ; 6.68 ; 6.76 ; 7.97 ; 6.76 ; 6.81 ; 8.40

La moyenne empirique vaut $\bar{x} = 7.247$ et la variance empirique (celle divisée par " $n - 1$ ") vérifie $\sqrt{\widehat{\sigma_x^2}} = 0.73$.

Dans la piscine 2, il trouve y_1, \dots, y_{10} qui valent respectivement

10.40 ; 7.27 ; 8.99 ; 7.28 ; 9.18 ; 9.10 ; 7.96 ; 7.71 ; 9.59 ; 9.61

Le moyennes et variance empiriques valent $\bar{y} = 8.709$ et $\sqrt{\widehat{\sigma_y^2}} = 1.08$.

Le constructeur n'a aucun parti pris entre les 2 produits et veut juste conseiller au mieux ses clients. Il voudrait pouvoir leur dire clairement "préférez A", "préférez B" ou "faites ce que vous voulez, les deux produits sont indiscernables vu mes mesures".

On sait que $pH = 7$ est neutre pour la peau tandis que $pH = 9$ est basique et donc nocif pour la peau.

1. Supposons que les observations sont les réalisations de variables aléatoires indépendantes gaussiennes de même variance σ^2 et de moyenne m_1 quand la mesure est prise dans la piscine 1, respectivement m_2 quand la moyenne est prise dans la piscine 2, m_1 et m_2 . Préciser le modèle statistique sous la forme d'un seul vecteur gaussien dont on donnera la moyenne et la matrice de covariance.
2. Supposons que σ^2 est connue, égale à 1.
 - i Donner l'expression la plus simple possible du test du rapport de vraisemblance de $\mathbf{H}_0 : "m_1 = m_2 = 7"$ contre $\mathbf{H}_1 : "m_1 = 7, \quad m_2 = 9"$ au niveau α .
 - ii Expliquer pourquoi ce test est uniformément le plus puissant parmi tous les tests de niveau α de ces mêmes hypothèses.
 - iii Calculer la p-valeur du test et conclure.
3. Toujours dans le cas où σ^2 est connu, construire un test de $\mathbf{H}_0 : "m_1 = m_2"$ contre $\mathbf{H}_1 : "m_1 < m_2"$ au niveau α .

4. Revenons au cas général (plus réaliste vu l'énoncé) où σ^2 est inconnue.

- i Donner les lois de $Z = \bar{Y} - \bar{X}$ et de $W = 9\widehat{\sigma_X^2} + 9\widehat{\sigma_Y^2}$. Expliquer pourquoi W est indépendante de Z .
- ii Toujours dans le cas où σ^2 est inconnue, construire un test de $\mathbf{H}_0 : "m_1 = m_2"$ contre $\mathbf{H}_1 : "m_1 < m_2"$ au niveau α . Calculer la p-valeur du test et conclure par l'une des 3 phrases "préférez A", "préférez B" ou "faites ce que vous voulez, les deux produits sont indiscernables vu les mesures".

Pour vous aider voici un tableau de valeurs qui donne t tel que $P(U > t) = a$

	a=0.05	a=0.025	a=0.01	a=0.005	a=0.0025	a=1.10 ⁻⁴
$U \sim \mathcal{N}(0, 1)$	1.645	1.960	2.326	2.576	2.807	3.719
$U \sim T(9)$	1.833	2.262	2.821	3.250	3.670	6.010
$U \sim T(10)$	1.812	2.228	2.763	3.169	3.581	5.694
$U \sim T(18)$	1.734	2.100	2.552	2.878	3.196	4.648
$U \sim T(20)$	1.724	2.085	2.528	2.845	3.153	4.538

où $T(k)$ est la loi du Student à k degrés de liberté.

5 BONUS : Test à rapport de vraisemblance monotone

Dans cet exercice, on se donne un modèle statistique dominé paramétré par un sous-ensemble $\Theta \subset \mathbb{R}$. On note $L_n(\theta)$ la vraisemblance de l'échantillon en θ . On suppose qu'il existe une statistique $\hat{S} = S(X_1, \dots, X_n)$ et, pour tout $(\theta_0, \theta_1) \in \Theta^2$ tel que $\theta_0 < \theta_1$, une fonction g_{θ_0, θ_1} strictement croissante telles que $\frac{L_n(\theta_1)}{L_n(\theta_0)} = g_{\theta_0, \theta_1}(\hat{S})$. On fixe $\theta_* \in \Theta$ et, pour tester $\mathbf{H}_0 : \theta \leq \theta_*$ contre $\mathbf{H}_1 : \theta > \theta_*$, on propose de rejeter \mathbf{H}_0 lorsque $r_n(\theta_*) \geq t_\alpha$ où

$$r_n(\theta_*) = \frac{\sup_{\theta_1 > \theta_*} L_n(\theta_1)}{\sup_{\theta_0 \leq \theta_*} L_n(\theta_0)}, \quad t_\alpha = \inf \left\{ t \in \mathbb{R}, \text{ tels que } \sup_{\theta_0 \leq \theta_*} P_\theta(r_n(\theta_*) \geq t) \leq \alpha \right\}.$$

1. Montrer qu'il existe $s_\alpha \in \mathbb{R}$ tel que le test rejette \mathbf{H}_0 si $\hat{S} \geq s_\alpha$.

On suppose maintenant qu'il existe s_α tel que $P_{\theta_*}(\hat{S} \geq s_\alpha) = \alpha$.

2. Soit $\theta_0 \leq \theta_*$. Montrer que $P_{\theta_0}(\hat{S} \geq s_\alpha) \leq \alpha$ et donc que le test est de niveau α .

Indications : On pourra utiliser le test de Neyman-Pearson de $\mathbf{H}_0 : \theta = \theta_0$ contre $\mathbf{H}_1 : \theta = \theta_*$ de niveau $\alpha' = P_{\theta_0}(\hat{S} \geq s_\alpha)$ et comparer sa puissance au test trivial rejetant cette hypothèse lorsque $T = 1$ où T est une variable aléatoire de Bernoulli de paramètre α' indépendante des observations.

3. Soit maintenant un autre test de niveau α défini par sa région de rejet R et notons $x_{1:n} = (x_1, \dots, x_n)$, $T : (x_{1:n}) \mapsto 1_{S(x_{1:n}) \geq s_\alpha}$ et $\phi : (x_{1:n}) \mapsto 1_{x_{1:n} \in R}$. Soit $\theta_1 > \theta_*$. Montrer que

$$\begin{aligned} P_{\theta_1}(\hat{S} \geq s_\alpha) - P_{\theta_1}(R) &= g_{\theta_*, \theta_1}(s_\alpha) \left(P_{\theta_*}(\hat{S} \geq s_\alpha) - P_{\theta_*}(R) \right) \\ &\quad + \mathbb{E}_{\theta_*} \left[(T(X_{1:n}) - \phi(X_{1:n})) \left(\frac{L_n(\theta_1)}{L_n(\theta_*)} - g_{\theta_*, \theta_1}(s_\alpha) \right) \right] \\ &\quad + \mathbb{E}_{\theta_1} [(T(X_{1:n}) - \phi(X_{1:n})) 1_{L_n(\theta_*)=0}] \end{aligned}$$

4. Montrer que chacun des trois termes du membre de gauche est positif ou nul et conclure.