

Génie Biologique - Statistiques S3

Résumé de cours

Table des matières

1 Lois de Variables aléatoires (1h)	7
2 Estimation et intervalles de confiance (2h)	13
3 Généralités sur les tests statistiques	21
4 Tests statistiques de base (2h)	23
5 Comparaison de distributions (3h)	29
6 Validation de la normalité (2h)	35
7 Comparaison de $m > 2$ populations normalement distribuées (4h + 3 TP)	43
8 Cartes de contrôle (2h + 2TP)	55

- Les 5h d'amphi ne sont pas comptées dans le prévisionnel
- La dernière séance de TP servira pour un TP noté

Résumé des fonctions Excel / Calc

Fonctions de base

- max, min, moyenne, racine : rien à signaler
- var, ecartype : empirique ou estimé ?

Lois statistiques

- loi.normale(.standard) : cumulative = vrai pour avoir l'aire au lieu de la valeur de la fonction
- loi.normale(.standard).inverse
- loi.f :
- inverse.loi.f
- loi.student
- loi.student.inverse
- loi.khideux
- khideux.inverse

A chacune de ces fonctions se posent les questions fondamentales suivantes :

- Probas ou quantiles ?
- Probas d'être inférieur ou d'être supérieur au paramètre donné ?
- Quantiles unilatéraux ou bilatéraux ?

Fonctions techniques

- somme.ecarts.carres(plage)
- sommeprod(colonne1 ; colonne2) et sommeprod(colonne1 ; colonne1 ; colonne1) pour la moyenne et la variance avec effectifs
- sommeprod((cellule < plage)*1 ; (plage<=cellule)*1) pour les effectifs de classe (le *1 sert à convertir des "vrai / faux" en 0 / 1)

1 Lois de Variables aléatoires (1h)

Modèles mathématiques basiques pour des comportements aléatoires.

1.1 Probabilités

- $\mathbb{P} : \{\text{événements}\} \mapsto [0; 1]$
- $\mathbb{P}(\bar{A}) = 1 - \mathbb{P}(A)$
- $\mathbb{P}(A \text{ ou } B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \text{ et } B)$
- $\mathbb{P}(A \text{ et } B) = \mathbb{P}(A|B)\mathbb{P}(B) = \begin{cases} 0 & \text{si } A \text{ et } B \text{ sont incompatibles} \\ \mathbb{P}(A)\mathbb{P}(B) & \text{si } A \text{ et } B \text{ sont indépendants} \end{cases}$

1.2 Variables aléatoires (V.A.)

- $X : \{\text{événements}\} \mapsto$ modalités de X (souvent notées x , ou même x_i si on peut les dénombrer)
- \mathbb{P} (événements tels que X prenne certaines valeurs) est généralement abrégé en $\mathbb{P}(X = \text{valeurs})$
- Espérance : $\mu(X) = \int x * \text{densite}(x)$ ou $\sum x * \mathbb{P}(X = x)$
- Variance : $\sigma^2(X) = \mu((X - \mu(X))^2) = \mu(X^2) - (\mu(X))^2$

Remarques

- Pour les variances et la covariance il faut calculer les moyennes avec 4 chiffres après la virgule
- X n'est *pas* un événement - $\mathbb{P}(X)$ n'a *pas* de sens
- X est une *fonction* qui transforme une expérience en modalité, qu'on appelle réalisation de X . Par exemple, on peut faire l'expérience de jeter une pièce en l'air : la variable aléatoire peut porter sur la prédiction du côté de la pièce qui finira sur le dessus ; mais une fois que la pièce a atterri on a une réalisation (c'est soit pile soit face, mais il n'y a plus rien à prédire)
- Espérance et variance d'une V.A. utilisent l'ensemble des modalités possibles pour X : ce sont la moyenne et la variance sur l'ensemble de la *population*. Par opposition, moyenne empirique et variance empirique sont calculées sur un échantillon : un ensemble de réalisations.

1.3 Lois de variables aléatoires discrètes

Une variable aléatoire est dite discrète si ses modalités sont dénombrables. La loi d'une V.A. discrète X est la donnée des valeurs $\mathbb{P}(X = x_i)$ pour toutes les modalités x_i de X .

1.3.1 Loi uniforme sur $\{1, \dots, n\}$: $X \sim \mathcal{U}(n)$

Cette loi modélise l'équiprobabilité sur n évènements.

$$\begin{aligned} - \mathbb{P}(X = k) &= \begin{cases} \frac{1}{n} & \text{si } 1 \leq k \leq n \\ 0 & \text{sinon} \end{cases} \\ - \mu(X) &= \frac{n+1}{2} \\ - \sigma^2 &= \frac{n^2-1}{12} \end{aligned}$$

1.3.2 Loi de Bernoulli : $X \sim \mathcal{B}(n; p)$

Cette loi modélise la somme des succès après n tentatives ayant chacune p chances de réussite.

$$\begin{aligned} - \mathbb{P}(X = k) &= C_n^k p^k (1-p)^{n-k} \\ - \mu(X) &= np \\ - \sigma^2 &= np(1-p) \end{aligned}$$

1.3.3 Loi de Poisson : $X \sim \mathcal{P}(\lambda)$

Cette loi modélise la somme des succès après un temps assez long en sachant qu'en moyenne on a λ succès sur cette durée.

$$\begin{aligned} - \mathbb{P}(X = k) &= \frac{\lambda^k}{e^\lambda k!} \\ - \mu(X) &= \lambda \\ - \sigma^2 &= \lambda \end{aligned}$$

1.4 Lois de variables aléatoires continues

Une variable aléatoire est dite continue si ses modalités sont des grandeurs physiques (le temps, le poids, ...). La loi d'une V.A. continue X est la donnée des valeurs $\mathbb{P}(X < x)$ pour toutes les modalités x de X .

1.4.1 Loi uniforme sur $[a; b]$: $X \sim \mathcal{U}[a; b]$

Cette loi modélise l'équiprobabilité sur un intervalle.

$$\begin{aligned} - \mathbb{P}(X < x) &= \begin{cases} 0 & \text{si } x \leq a \\ \frac{x-a}{b-a} & \text{si } a \leq x \leq b \\ 1 & \text{si } x \geq b \end{cases} \\ - \mu(X) &= \frac{a+b}{2} \\ - \sigma^2 &= \frac{(b-a)^2}{12} \end{aligned}$$

1.4.2 Loi exponentielle de paramètre λ

Cette loi modélise la durée de vie d'un phénomène sans mémoire : à un instant donné, la probabilité de sa durée de vie ne dépend pas de son âge.

$$\begin{aligned} - \mathbb{P}(X < x) &= \begin{cases} 0 & \text{si } x \leq 0 \\ 1 - e^{-\lambda x} & \text{sinon} \end{cases} \\ - \mu(X) &= \frac{1}{\lambda} \\ - \sigma^2 &= \frac{1}{\lambda^2} \end{aligned}$$

1.4.3 Loi normale standard : $X \sim \mathcal{N}(0; 1)$

Cette loi modélise quasiment tout le reste !

- $\mathbb{P}(X < x)$ sera donné par un tableau pré-calculé ou l'ordinateur
- $\mu = 0$
- $\sigma^2 = 1^2$
- Cette loi est symétrique : $\mathbb{P}(X < -x) = \mathbb{P}(X > x)$

1.4.4 Loi normale générale : $X \sim \mathcal{N}(\mu; \sigma)$

Si $X \sim \mathcal{N}(\mu; \sigma)$, alors on va utiliser une variable aléatoire qui servira d'intermédiaire de calcul : $Y = \left(\frac{X-\mu}{\sigma}\right)$. On a alors $Y \sim \mathcal{N}(0; 1)$ et

$$\mathbb{P}(X < x) = \mathbb{P}\left(Y < \frac{x - \mu}{\sigma}\right)$$

Approximations classiques :

- $\mathbb{P}(\mu - \sigma < X < \mu + \sigma) \approx 68\%$
- $\mathbb{P}(\mu - 2\sigma < X < \mu + 2\sigma) \approx 95\%$

Lois de Variables aléatoires - Exercices

Exercice 1 Loi exponentielle

1. Soit X une variable aléatoire qui suit la loi exponentielle de paramètre $\frac{1}{2}$. Calculer $\mathbb{P}(1 \leq X \leq 2)$.
2. Trouver λ , sachant que la variable X suit une loi exponentielle et que $\mathbb{P}(X \leq 50) = 1 - \frac{1}{e}$.
3. Avec ce paramètre, que vaut $\mathbb{P}(X \geq 25)$
4. De même, que vaut $\mathbb{P}(X \geq 50)$

Exercice 2 Loi normale

1. Soit X une variable aléatoire dont la loi est $\mathcal{N}(20, 5)$.
 - a) Que vaut $\mathbb{P}(X \leq 20)$?
 - b) Si $\mathbb{P}(X \leq 15) = 0.16$, que valent $\mathbb{P}(X \geq 15)$ et $\mathbb{P}(X \geq 25)$?
2. Soit X une variable aléatoire dont la loi est $\mathcal{N}(20, 5)$. Quelle est la valeur de a telle que :
 - a) $\mathbb{P}(X \leq a) = 0.975$,
 - b) $\mathbb{P}(X \geq a) = 0.005$,
 - c) $\mathbb{P}(20 - a \leq X \leq 20 + a) = 0.95$.
3. Soit X une variable aléatoire dont la loi est $\mathcal{N}(176, 6)$. Déterminer a tel que $\mathbb{P}(176 - a \leq X \leq 176 + a) = 0.68$.
4. Soit X une variable aléatoire dont la loi est $\mathcal{N}(-5, 1)$. Déterminer a tel que $\mathbb{P}(-5 - a \leq X \leq -5 + a) = 0.95$.
5. On considère une V.A. X normalement distribuée, de moyenne 3 et d'écart-type inconnu σ . Déterminer σ , sachant que $\mathbb{P}(X > 0) = 0.9$.
6. (Plus difficile) Soit X une (autre) V.A. distribuée suivant une loi normale dont la moyenne et l'écart-type sont inconnus. Sachant que $\mathbb{P}(X > 0) = 0.3$ et $\mathbb{P}(X < -1) = 0.4$, déterminer les paramètres de la loi de X .

Exercice 3 Loi du χ^2 et loi T Soient X une V.A. suivant une loi du χ^2 à 24 degrés de liberté et S une V.A. suivant une loi de Student (T) à 10 degrés de liberté.

1. Que vaut $\mathbb{P}(X < 13.848)$? et $\mathbb{P}(X < 36.415)$? En déduire ce que calcule l'ordinateur quand on appelle la fonction "loi.khideux(36.415;24)" (cette question est fondamentale!).
2. Trouver a tel que $\mathbb{P}(X < a) = 1\%$.
3. Trouver b tel que $\mathbb{P}(X > b) = 1\%$.
4. Que vaut $\mathbb{P}(S < 1.812)$? $\mathbb{P}(S > 1.812)$
5. Trouver a tel que $\mathbb{P}(-a < S < a) = 99\%$
6. Trouver b tel que $\mathbb{P}(S < b) = 2.5\%$

Réponse 1 Loi exponentielle

1. $-\frac{1}{e} + \frac{1}{\sqrt{e}}$.
2. $\lambda = \frac{1}{50}$.
3. $\frac{1}{\sqrt{e}}$.
4. $\frac{1}{e}$.

Réponse 2 Loi normale

1. a) 0.5
b) 0.84 et 0.16
2. a) $a = 29.8$
b) $a = 32.75$ (en réalité un peu plus, ceci est dû au manque de précision des tableaux)
c) $a = 9.8$
3. $a = 6$ (en réalité un peu moins, ceci est une approximation rapide)
4. $a = 2$ (en réalité un peu moins, ceci est une approximation rapide)
5. $\sigma = 2.34$
6. $\sigma = 1.28$ et $\mu = -0.67$.

Réponse 3 Loi du χ^2 et loi T

1. 5% et 95%. De fait, la fonction "loi.khideux(36.415 ;24)" calcule la valeur $\mathbb{P}(X > 36.415)$ et il faut bien faire attention à la convention utilisée par le logiciel !
2. $a = 10,856$.
3. $b = 42,980$.
4. 95% et 5%
5. $a = 3.169$.
6. $b = -2.228$

2 Estimation et intervalles de confiance (2h)

Inférence des paramètres d'une population trop grande pour être recensée intégralement.

2.1 Estimation

Approximation sans garantie

2.1.1 Définitions

- L'estimation d'un paramètre d'une *population* est une valeur calculée sur un *échantillon* et censée être "proche" du paramètre estimé.
- L'estimateur est la fonction qui donne l'estimation à partir des valeurs de l'échantillon ; c'est une variable aléatoire.

Remarque "proche" a une définition rigoureuse dont voici l'idée :

- Si la taille de l'échantillon grandit, les estimations ainsi obtenues doivent se rapprocher du paramètre estimé.
- Si on prend de plus en plus d'échantillons différents, la moyenne des estimations ainsi obtenues doit se rapprocher du paramètre estimé.

2.1.2 Estimations principales

- Une "bonne" estimation de la proportion π de la *population* est $\pi_{est} = \pi_E$, où π_E est la proportion empirique.
- Une "bonne" estimation de l'espérance μ de la *population* est $\mu_{est} = \mu_E$, où μ_E est la moyenne empirique.
- Une "bonne" estimation de la variance σ^2 de la *population* est $\sigma_{est}^2 = \frac{n}{n-1} \sigma_E^2$, où σ_E^2 est la variance empirique.
- Une "bonne" estimation de l'écart-type σ de la *population* est $\sigma_{est} = \sqrt{\frac{n}{n-1}} \sigma_E$, où σ_E est l'écart-type empirique.

Remarques

- $\sigma^2(X) = \mu(X^2) - (\mu(X))^2$
- $\sigma_E^2 = \frac{1}{n} \sum (x_i - x.)^2$
- $\sigma_{est}^2 = \frac{1}{n-1} \sum (x_i - x.)^2$

2.2 Intervalles de confiance (IDC) pour une V.A. $X \sim \mathcal{N}(\mu; \sigma)$

Encadrement avec garantie

2.2.1 Définitions

- Un intervalle de confiance au risque α est un *intervalle* construit à partir des valeurs d'un échantillon tel que, sur l'ensemble des réalisations de la V.A. étudiée, exactement $\alpha\%$ des IDC qui découlent ne contiennent *pas* le paramètre étudié.
- $u(\alpha)$ est le quantile de la loi normale standard tel que $\mathbb{P}(U > u) = \alpha$, $U \sim \mathcal{N}(0; 1)$
- $t(\alpha)$ est le quantile de la loi de Student à $(n - 1)$ degrés de liberté tel que $\mathbb{P}(T > t) = \alpha$, $T \sim Student(n - 1)$
- $k(\alpha)$ est le quantile de la loi du Khi-deux à $(n - 1)$ degrés de liberté tel que $\mathbb{P}(K > k) = \alpha$, $K \sim \chi^2(n - 1)$

2.2.2 Interprétation

- Le niveau de confiance $(1 - \alpha)$ est le pourcentage d'intervalles qui contiendront le paramètre en prenant des échantillons différents; et non pas la probabilité pour que le paramètre étudié soit dans l'intervalle calculé. C'est encore la différence entre la V.A. qui permet de calculer l'intervalle (la formule) et réalisation de cette V.A. (résultat de la formule avec l'échantillon particulier dont on dispose).
Par exemple : le jeu des trois gobelets. On place un haricot sous un des gobelets et on les mélange. On a alors 33% de chances que le haricot soit sous le premier gobelet (niveau de confiance de la prédiction). Mais une fois le gobelet choisi, même si on ne regarde pas dessous, il n'y a plus de hasard : le haricot s'y trouve ou ne s'y trouve pas. De même, une fois les mesures faites et l'intervalle construit, il n'y a plus de hasard : la paramètre étudié s'y trouve ou ne s'y trouve pas, mais c'est définitif. Mais si on fait 100 fois la même expérience, 33% de nos choix devraient être bons.
- Un intervalle de confiance est un compromis entre faible risque d'erreur, pertinence de l'intervalle et petite taille de l'échantillon. On ne peut généralement pas avoir les trois en même temps.
- Les IDC qui vont suivre sont donnés pour quand on sait que X suit une loi normale. Ils sont aussi valables quand on ne le sait pas mais que n est "assez grand" : au delà de 30 c'est acceptable, au delà de 50 c'est correct et au delà de 100 c'est indiscernable.

2.2.3 IDC pour les valeurs de X

Aussi appelé intervalle de fluctuation.

- σ connu : $\left[\mu_{(est)} - u\left(\frac{\alpha}{2}\right)\sigma ; \mu_{(est)} + u\left(\frac{\alpha}{2}\right)\sigma \right]$
- σ estimé : $\left[\mu_{(est)} - t\left(\frac{\alpha}{2}\right)\sigma_{est} ; \mu_{(est)} + t\left(\frac{\alpha}{2}\right)\sigma_{est} \right]$

2.2.4 IDC pour la moyenne de X

- σ connu : $\left[\mu_{est} - u\left(\frac{\alpha}{2}\right)\frac{\sigma}{\sqrt{n}} ; \mu_{est} + u\left(\frac{\alpha}{2}\right)\frac{\sigma}{\sqrt{n}} \right]$

$$- \sigma \text{ estimé : } \left[\mu_{est} - t\left(\frac{\alpha}{2}\right) \frac{\sigma_{est}}{\sqrt{n}} ; \mu_{est} + t\left(\frac{\alpha}{2}\right) \frac{\sigma_{est}}{\sqrt{n}} \right]$$

Remarque On voit que la variance de la moyenne est n fois plus petite que la variance des valeurs

2.2.5 IDC pour la variance de X

$$- \left[\frac{(n-1)\sigma_{est}^2}{k\left(\frac{\alpha}{2}\right)} ; \frac{(n-1)\sigma_{est}^2}{k\left(1-\frac{\alpha}{2}\right)} \right]$$

Remarque On obtient celui pour l'écart-type en prenant la racine des bornes.

Estimation et intervalles de confiance - Exercices

Exercice 1 Estimation ponctuelle de la moyenne et de l'écart-type Lors d'un concours radiophonique, on note X le nombre de réponses reçues chaque jour. Durant les 10 premiers jours, on a obtenu :

$$\begin{array}{cccccc} x_1 = 200 & x_2 = 240 & x_3 = 190 & x_4 = 150 & x_5 = 220 \\ x_6 = 180 & x_7 = 170 & x_8 = 230 & x_9 = 210 & x_{10} = 210 \end{array}$$

Déterminer une estimation ponctuelle de la moyenne et de l'écart-type.

Exercice 2 Estimation de variance Lors d'un contrôle d'une chaîne de médicaments, on s'intéresse au nombre de comprimés défectueux dans un lot. L'étude de 200 lots a donné les résultats suivants :

# de comprimés défectueux	0	1	2	3	4	5
Nombre de lots	75	53	39	23	9	1

1. Calculer la moyenne, le mode et les quartiles du nombre de comprimés défectueux pour cet échantillon de 200 lots.
2. Calculer la variance, l'écart-type et le coefficient de variation du nombre de comprimés défectueux pour cet échantillon de 200 lots.
3. Donner une estimation de la variance et de l'écart-type du nombre de comprimés défectueux sur l'ensemble de la production.

Exercice 3 Estimation par IDC de la moyenne Dans une station service, on suppose que le montant des chèques essence suit une loi normale de paramètres μ et σ . On considère un échantillon de taille $n = 51$ et on obtient une moyenne de 13€ et un écart-type (empirique) de 4€.

1. Donner une estimation de μ par un intervalle de confiance au niveau de confiance 95%.
2. Donner une estimation de σ par un intervalle de confiance au niveau de confiance 95%.

Exercice 4 Exemple pratique – suite Voici les indicateurs statistiques de base pour un contrôle de statistiques :

- effectif = 49
- max = 20
- moyenne = 10.4551
- médiane = 10
- min = 1.15
- var = 16.8648

2 Estimation et intervalles de confiance (2h)

- écart-type = 4.1067
- coeff. variation = 39.28 %

1. Estimez la moyenne et la variance des notes en général.
2. Calculez l'intervalle de confiance à 95 % pour les valeurs (aussi appelé intervalle de fluctuation) des notes.
3. Calculez l'intervalle de confiance à 98 % pour les valeurs des notes (oui, on peut le faire avec seulement les tableaux).
4. Comparez avec le min et le max de l'échantillon.
5. [Avec Excel] Quel serait le pourcentage de risque exact pour que l'IDC aille précisément du min au max de l'échantillon ?
6. A l'inverse, quelle serait l'écart-type estimé pour la population si l'IDC à 95% pour les valeurs allait précisément du min au max de l'échantillon ?

Exercice 5 Estimation et IDC 1 On considère la variable X masse d'un ressort provenant d'une certaine fabrication. Cette variable suit une loi normale de moyenne μ et d'écart-type σ . On donne la répartition des masses de 219 ressorts :

X Masses (g)	[8,2;8,4[[8,4;8,6[[8,6;8,8[[8,8;9,0[[9,0;9,2[[9,2;9,4[[9,4;9,6[
Nb de ressorts	9	21	39	63	45	27	15

1. Donner une estimation ponctuelle de μ et σ .
2. Que dire de la qualité d'une estimation par Intervalle de Confiance dans ce cas ?

Exercice 6 Estimation et IDC 2 On veut estimer l'espérance mathématique μ d'une variable aléatoire Gaussienne X dont on connaît l'écart-type $\sigma = 2,3$. Quelle est la taille minimale de l'échantillon de X qui est à prendre si l'on veut obtenir pour μ un intervalle de confiance de seuil 0,95 et dont la longueur ne dépasse pas 0,1.

Réponse 1 Estimation ponctuelle de la moyenne et de l'écart-type $\mu_{est} = 200$ et $\sigma_{est} = \sqrt{777.78} = 27.89$

Réponse 2 Estimation de variance

1. $\mu_E = 1.205$, mode = 0, effectifs cumulés = 75, 128, 167, 190, 199, 200 : $Q_1 = 0$, $Q_2 = 1$ et $Q_3 = 2$.
2. $\sigma_E = 1.214$, $CV = 1.007 (> 100\%)$
3. $\sigma_{est} = \sqrt{1.4803} = 1.217$

Réponse 3 Estimation par IDC de la moyenne $\mu_{est} = 13$, $\sigma_{est} = 4.040$

1. $11.864 \leq \mu \leq 14.136$
2. $3.380 \leq \sigma \leq 5.022$

Réponse 4 Exemple pratique

1. $\mu = 10.4551, \sigma_{est}^2 = 17.2159$.
2. Le quantile bilatéral de Student à 5% vaut 2.011 et l'IDC est alors [2.1127; 18.7975].
3. Le quantile bilatéral de Student à 2% vaut 2.407 et l'IDC est alors [0.469620.4406].
4. On voit que l'IDC à 95% n'englobe pas toutes les valeurs de l'échantillon, mais que l'IDC à 98% les dépasse : un IDC à 95% est donc assez parlant pour décrire les valeurs communément mesurées, et ne nécessite que trois informations (moyenne, écart-type et effectif).
5. $max - min = 18.85 = 2t_{\alpha/2}\sigma_{est}$. En gardant $\sigma_{est} = 4.149$, on trouve $\alpha = 2.7638\%$.
6. En gardant $t_{\alpha/2} = 2.011$, la même formule donne $\sigma_{est} = 4.6877$.

Réponse 5 Estimation et IDC 1

1. $\mu_{est} = 8.9329, \sigma_{est} = 0.2986$
2. L'écart-type (numérateur) est petit et la taille de l'échantillon (dénominateur) est grande : la largeur de l'intervalle de confiance sera faible et l'encadrement sera donc précis.

Réponse 6 Estimation et IDC 2 Il faut considérer au moins $n = 8129$ individus.

3 Généralités sur les tests statistiques

- Ce sont des aides à la décision pour trancher systématiquement entre "oui" et "non" quand la réponse n'est pas évidente.
- La question porte sur la *population* ! On n'a pas besoin de faire de longues études pour savoir si 3.14 est plus grand que 3 ou pas.
- Comme tout résultat statistique, on n'obtiendra pas des vérités absolues, mais des éventualités assez probables en se basant sur les informations incomplètes disponibles

3.1 Préliminaires

- \mathcal{H}_0 : hypothèse privilégiée — imposée par le test
- \mathcal{H}_1 : hypothèse alternative — parfois orientable en fonction du résultat recherché

3.2 Estimation contextuelle

Sous l'hypothèse privilégiée \mathcal{H}_0 , on aura alors une certaine V.A. qui suivra une certaine loi

- On peut ainsi calculer l'intervalle de confiance pour cette V.A., généralement avec un niveau de risque $\alpha = 5\%$: c'est la zone de validité pour \mathcal{H}_0
- En parallèle, on calcule une estimation de la V.A. à partir de l'échantillon disponible

3.3 Prise de décision

- Si l'estimation ne rentre pas dans l'intervalle de confiance, c'est donc que les observations ne sont pas compatibles avec la l'hypothèse privilégiée : on réfute \mathcal{H}_0 pour accepter \mathcal{H}_1
- Si l'estimation rentre dans l'intervalle de confiance, les observations sont donc compatibles avec \mathcal{H}_0 : on accepte \mathcal{H}_0

3.4 Risques d'erreur

- On réfute \mathcal{H}_0 pour accepter \mathcal{H}_1 quand l'estimation est dans la zone de risque de l'IDC : α , appelé risque de première espèce, est donc précisément le risque d'erreur à réfuter \mathcal{H}_0 — On dit alors que le test est significatif au seuil α et idéalement on précise le risque le plus bas possible permettant de réfuter \mathcal{H}_0
- Quand l'estimation est dans l'IDC, on ne peut pas réfuter \mathcal{H}_0 et α n'est donc plus le risque d'erreur mais plutôt une sorte d'indicateur de qualité — On dit que le test est non significatif au seuil α et on précise idéalement le risque le plus haut possible ne réfutant pas \mathcal{H}_0

3 Généralités sur les tests statistiques

- Si on ne peut pas réfuter \mathcal{H}_0 , le risque d'erreur, appelé risque de deuxième espèce, est noté β
- Il est généralement difficile à calculer mais c'est lui qui fait la différence d'efficacité (appelée puissance) entre deux tests différents pour la même question

Remarque Si on ne peut pas réfuter \mathcal{H}_0 , cela ne veut pas dire que l'hypothèse privilégiée soit vraie. Cela veut juste dire qu'on n'a aucune preuve de contraire. Si on la pensait vraie pour d'autres raisons, alors on peut continuer d'y croire. Sinon on ne peut pas dire grand chose, d'où la non-significativité du test

4 Tests statistiques de base (2h)

Ces tests sont valables quand la V.A. étudiée suit une loi normale, ou quand la taille de l'échantillon est supérieure à 30.

4.1 Comparaison d'un paramètre observé à un paramètre théorique

Un seul échantillon

4.1.1 Comparaison de variance

- $\mathcal{H}_0 : (\sigma(X) = \sigma_0)$
- $\mathcal{H}_1 : (\sigma(X) \neq \sigma_0)$ ou $(\sigma(X) > \sigma_0)$ ou $(\sigma(X) < \sigma_0)$
- Sous l'hypothèse \mathcal{H}_0 , $K = \frac{(n-1)\sigma^2(X)}{\sigma_0^2} \sim \chi^2(n-1)$

4.1.2 Comparaison de moyenne

- $\mathcal{H}_0 : (\mu(X) = \mu_0)$
- $\mathcal{H}_1 : (\mu(X) \neq \mu_0)$ ou $(\mu(X) > \mu_0)$ ou $(\mu(X) < \mu_0)$
- Sous l'hypothèse \mathcal{H}_0 :
 - Si $\sigma(X)$ est connu : $U = \frac{(\mu(X) - \mu_0)}{\sqrt{\frac{\sigma^2}{n}}} \sim \mathcal{N}(0; 1)$
 - Si $\sigma(X)$ est estimé : $T = \frac{(\mu(X) - \mu_0)}{\sqrt{\frac{\sigma^2(X)}{n}}} \sim Student(n-1)$

4.2 Comparaison de deux paramètres observés - A

Deux échantillons *indépendants*

4.2.1 Comparaison de variances

- $\mathcal{H}_0 : (\sigma(X_1) = \sigma(X_2))$

4 Tests statistiques de base (2h)

- $\mathcal{H}_1 : (\sigma(X_1) \neq \sigma(X_2))$ ou $(\sigma(X_1) > \sigma(X_2))$
- Sous l'hypothèse \mathcal{H}_0 , $F = \frac{\sigma^2(X_1)}{\sigma^2(X_2)} \sim Fisher(n_1 - 1 ; n_2 - 1)$

4.2.2 Comparaison de moyennes

- $\mathcal{H}_0 : (\mu(X_1) = \mu(X_2))$
 - $\mathcal{H}_1 : (\mu(X_1) \neq \mu(X_2))$ ou $(\mu(X_1) > \mu(X_2))$ ou $(\mu(X_1) < \mu(X_2))$
 - Sous l'hypothèse \mathcal{H}_0 :
 - Si $\sigma(X_1)$ et $\sigma(X_2)$ sont connus : $U = \frac{(\mu(X_1) - \mu(X_2))}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim \mathcal{N}(0; 1)$
 - Si $\sigma(X_1)$ et $\sigma(X_2)$ sont estimés, il faut commencer par tester leur égalité.
 - Si on admet que $\sigma(X_1) = \sigma(X_2)$: $T = \frac{(\mu(X_1) - \mu(X_2))}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}} \sim Student(n_1 + n_2 - 2)$, où

$$s^2 = \frac{(n_1 - 1)\sigma^2(X_1) + (n_2 - 1)\sigma^2(X_2)}{n_1 + n_2 - 2}$$
 (la moyenne pondérée des variances)
 - Si on admet que $\sigma(X_1) \neq \sigma(X_2)$: $T = \frac{(\mu(X_1) - \mu(X_2))}{\sqrt{\frac{\sigma^2(X_1)}{n_1} + \frac{\sigma^2(X_2)}{n_2}}} \sim Student(k)$, où

$$k = \text{Partie entière de } \left(\frac{\left(\frac{\sigma_{1est}^2}{n_1} + \frac{\sigma_{2est}^2}{n_2} \right)^2}{\left(\frac{\sigma_{1est}^2}{n_1} \right)^2 + \left(\frac{\sigma_{2est}^2}{n_2} \right)^2} \right)$$
. Exemple : partie entière de 4.87 = 4.
- Fonction Excel : ENT().
Version simplifiée : $k = \text{minimum entre } n_1 - 1 \text{ et } n_2 - 1$.

4.3 Comparaison de deux paramètres observés - B

Deux échantillons *appariés*

4.3.1 Comparaison de moyennes

On considère les valeurs de la V.A. $D = X_1 - X_2$: la différence entre les deux groupes.

- $\mathcal{H}_0 : (\mu(D) = 0)$
- $\mathcal{H}_1 : (\mu(D) \neq 0)$ ou $(\mu(D) > 0)$ ou $(\mu(D) < 0)$
- Sous l'hypothèse \mathcal{H}_0 , $T = \frac{\mu(D)}{\sqrt{\frac{\sigma^2(D)}{n}}} \sim Student(n - 1)$

Tests statistiques de base - Exercices

Exercice 1 La teneur en hémoglobine du sang des femmes non malades a pour valeur moyenne 14.5 g/100 mL et pour écart-type 1.1 g/100 mL qu'on supposera constant quelle que soit la population étudiée. Ce paramètre biologique suit une loi normale. Sur un échantillon de 20 femmes on trouve une teneur moyenne en hémoglobine de 13.8 g/100 mL. Au risque de 5%, peut-on conclure que la population de femmes dont est extrait cet échantillon présente une teneur en hémoglobine trop faible ?

Exercice 2 Dans un laboratoire pharmaceutique sont fabriquées des ampoules de soluté injectable dont le volume doit être de 1 mL. Afin de contrôler sa production, l'industriel effectue régulièrement des prélèvements de lots de 50 ampoules. Au cours d'un tel prélèvement, les résultats sont les suivants :

$$\sum_{i=1}^{50} x_i = 51.34 \text{ mL} \quad \text{et} \quad \sum_{i=1}^{50} x_i^2 = 53.64 \text{ mL}^2.$$

1. Calculez moyenne, variance et écart-type pour cet échantillon.
2. Déduisez-en une estimation ponctuelle de ces mêmes paramètres sur la population.
3. Au vu des ces résultats, peut-on conclure que le prélèvement est conforme à la valeur de référence ?

Exercice 3 Dans un laboratoire d'analyses médicales, on effectue le dosage de calcium sérique par une méthode colorimétrie dont l'écart-type sur la mesure est de 1.1 mg/L. Ce dosage suit une loi normale. Après une remise à niveau de l'appareil de mesure, le directeur de ce laboratoire veut vérifier que l'écart-type sur la mesure n'a pas augmenté. Il fait refaire par un même technicien, dans les mêmes conditions, 32 dosages d'un même prélèvement sérique et obtient les résultats suivants :

<i>Teneur en Ca mg/L</i>	[96.5;97.5[[97.5;98.5[[98.5;99.5[[99.5;100.5[[100.5;101.5[
<i>Nombre</i>	3	7	11	6	5

Exercice 4 Dans une usine sont fabriquées des micro-pipettes de laboratoire de volume calibré à 100 μ L. En cours de fabrication sur deux machines différentes, on prélève deux lots indépendants de 25 micro-pipettes et on obtient les valeurs suivantes pour les variances calculées sur les deux échantillons : $\sigma_{E1}^2 = 0.02 \mu\text{L}^2$ et $\sigma_{E2}^2 = 0.015 \mu\text{L}^2$. On suppose que la variable *volume d'une micro-pipette* suit une loi normale. Peut-on dire que les variances des volumes de micro-pipettes fabriquées par ces deux machines sont différentes (au risque de 5%) ?

Exercice 5 Dans une unité de fabrication, deux machines *A* et *B* produisent des ampoules pour soluté injectable. On sait que les variances du volume de ces ampoules dans toute la production de l'usine valent pour la machine *A* : 0.04 mL² et pour la machine *B* : 0.03 mL². Deux échantillons indépendants de 20 ampoules sont prélevés et on obtient les valeurs suivantes pour les volumes moyens : 4.9 mL pour l'échantillon issu de la production par la machine *A* et 5.3 mL pour l'échantillon issu de la production par la machine *B*.

4 Tests statistiques de base (2h)

1. Peut-on dire que la machine A produit des ampoules de volume moyen différent de celles produites par la machine B ?
2. On a donc conclu à une différence des volumes des machines à partir d'une différence des volumes des échantillons. A partir de quelle différence entre les volumes moyens des échantillons peut-on conclure, au risque de 5%, à une différence significative entre les volumes moyens de toutes les ampoules produites par les deux machines ?

Exercice 6 Le dosage des transaminases (exprimés en UI/L) informe sur le fonctionnement hépatique. Lors de la mise au point d'un traitement qui peut perturber la fonction hépatique, deux variantes du protocole thérapeutique sont appliquées à deux groupes indépendants de malades :

- la variante 1 à un groupe de 12 malades,
- la variante 2 à un groupe de 15 malades.

On suit la sensibilité au traitement par dosage des transaminases dont les valeurs sont reportées dans le tableau suivant :

<i>Variante 1</i>	22	18	25	14	16	22	17	19	30	17	23	17			
<i>Variante 2</i>	31	35	32	30	28	26	27	19	25	20	18	31	27	29	24

La variable *dosage des transaminases* est supposée suivre une loi normale. Peut-on dire, au risque de 5%, que les deux variantes du protocole modifient de façons différentes la fonction hépatique ?

Exercice 7 Dans un laboratoire pharmaceutique, deux machines assurent la fabrication des comprimés, l'une est rotative, l'autre est verticale. La production est contrôlée par pesée des comprimés et ce poids est supposé suivre une loi normale. On extrait de la production deux échantillons indépendants, de même taille $n = 25$, sur lesquels les comprimés sont pesés. On obtient les résultats suivants (en g et en g^2) :

- Pour la machine rotative, $\sum_{i=1}^{25} x_{i,1} = 2541$ et $\sum_{i=1}^{25} x_{i,1}^2 = 258\,505$,
- Pour la machine verticale, $\sum_{i=1}^{25} x_{i,2} = 2555$ et $\sum_{i=1}^{25} x_{i,2}^2 = 262\,547$.

Peut-on dire que la machine rotative produit des comprimés de poids inférieur à ceux produits par la machine verticale ?

Exercice 8 On étudie l'effet d'une nouvelle forme médicamenteuse sur la tension artérielle. On constitue un groupe de 12 personnes et on mesure leur tension artérielle avant et après ingestion de la substance. Les résultats sont les mesures de la tension systolique exprimées en mm de mercure.

<i>Avant</i>	120	130	132	125	140	145	135	125	133	140	138	137
<i>Après</i>	122	125	118	135	142	138	125	115	123	135	122	126

Peut-on dire que cette forme médicamenteuse a fait baisser la tension artérielle ?

Réponse 1 Comparaison à une moyenne théorique d'un petit échantillon suivant une loi normale d'écart type connu : avec $U_{est} = -2.85$, $U_{5\%} = -1.6449$ (test unilatéral), au risque de 5% on peut conclure que la population dont est extrait l'échantillon est constituée de femmes dont la teneur en hémoglobine est trop faible.

Réponse 2

1. $\mu_E = 1.027, \sigma_E^2 = 0.018482, \sigma_E = 0.13595$.
2. $\mu_{est} = 1.027, \sigma_{est}^2 = 0.018859, \sigma_{est} = 0.13733$.
3. Comparaison à une moyenne théorique avec un échantillon suffisamment grand pour être supposé suivre une loi normale d'écart type inconnu : avec $T_{est} = 1.39$ et $T_{2.5\%} = \pm 2.009$ (test bilatéral), au risque de 5% on peut conclure que le volume moyen semble conforme à la valeur de référence.

Réponse 3 Ecart type estimé $\sigma_{est} = 1.2011\text{mg/L}$. Avec $K_{est} = 36.96$ et $\chi_{5\%}^2 = 44.9$ (test unilatéral), au risque de 5% on ne peut pas conclure que l'écart type après intervention soit supérieur à l'écart type initial. Remarque : Sur les tableaux, on n'aura que la ligne 30 DDL, au lieu de 31, et on prendra donc la limite $\chi_{5\%}^2 = 43.7$.

Réponse 4 $\sigma_{est.1}^2 = 0.02083, \sigma_{est.2}^2 = 0.015625$. Avec $F_{est} = 1.33$ et $F_{2.5\%} = 2.27$, au risque de 5% (test bilatéral – avec les tableaux, comme on n'a pas (24;24), on prendra la valeur 2.3) on ne peut pas conclure que les variances associées à chacune des deux machines soient différentes.

Réponse 5 Comme on a un échantillon petit, il faut ajouter l'hypothèse de normalité pour la V.A. étudiée (le volume des ampoules) – les conditions énoncées ne permettant pas de la garantir.

1. Avec $U_{est} = -6.76$ et $U_{5\%} = -1.96$ (test bilatéral), au risque de 5% on peut dire que les ampoules A sont d'un volume moyen significativement inférieur aux B (conclusion valable même avec un risque d'erreur de 0.1%).
2. On conclut à une différence, au risque de 5%, entre les productions dès que $|U_{est}| > 1.96$, ce qui donne $|\mu_A - \mu_B| > 0.116 \text{ mL}$.

Réponse 6 Comparaison de moyennes observées sur des échantillons indépendants de variances inconnues.

1. Comparaison des variances : avec des variances estimées $\sigma_{est.1}^2 = 20.54545, \sigma_{est.2}^2 = 24.45714$, une variable de décision $F_{est} = 1.19039$ et une valeur seuil $F_{2.5\%} = 3.36$, on peut conclure au risque de 5% (test bilatéral) que les variances sont égales pour les deux protocoles. On conserve alors comme variance commune $s^2 = 22.73$.
2. Comparaison des moyennes : avec les moyennes $\mu_{E1} = 20, \mu_{E2} = 26.8$, la variable de décision $T_{est} = -3.68$ et la valeur seuil $T_{2.5\%} = 2.06$ (test bilatéral), au risque de 5% on peut conclure que les moyennes sont significativement différentes selon le protocole utilisé. Cette conclusion est même valable avec un risque d'erreur de 1%.

Réponse 7

1. Comparaison des variances : avec des variances estimées $\sigma_{est.1}^2 = 9.91, \sigma_{est.2}^2 = 59.42$, une variable de décision $F_{est} = 5.998$ et une valeur seuil $F_{2.5\%} = 2.27$ (test bilatéral – avec les tableaux, comme on n'a pas (24;24), on prendra la valeur 2.3), on peut conclure au risque $\alpha = 5\%$ que les variances sont différentes. On ne peut donc pas calculer de variance commune.
2. Comparaison des moyennes : avec les moyennes $\mu_{E1} = 101.64, \mu_{E2} = 102.2$, la variable de décision $T_{est} = -0.336$ et la valeur seuil $T_{2.5\%} = \pm 2.04$ (test unilatéral avec 31 DDL – valeur calculée avec la formule adaptée au cas de variances non connues mais supposés différentes), au risque de $\alpha = 5\%$ on ne peut pas conclure que les poids moyens diffèrent.

4 Tests statistiques de base (2h)

Réponse 8 $\mu_{est}(D) = 6.16667, \sigma_{est}^2(D) = 56.69697, T_{est} = 2.837, T_{5\%} = 1.796$ (test unilatéral) : on peut conclure que la tension artérielle diminue après le traitement, au risque de 5% et même au risque de 1%.

5 Comparaison de distributions (3h)

- Le test du χ^2 est un test générique permettant de comparer les valeurs observées à des valeurs théoriques.
- Dans toute cette partie, la variable de test suit la loi du χ^2
- Remarque essentielle : les réalisations du χ^2 se calculent toujours sur des effectifs.

5.1 Le test du χ^2

Le test du χ^2 compare des effectifs observés (o_i) à des effectifs théoriques (t_i).

Le test

- \mathcal{H}_0 : Variable (induit le calcul des t_i)
- \mathcal{H}_1 : Les observations ne concordent pas avec l'hypothèse privilégiée.
- La réalisation de la V.A. de test pour l'ajustement est donnée par $k = \sum_{\text{modalités}} \frac{(o_i - t_i)^2}{t_i}$
- La zone de rejet se calcule par la loi du χ^2 avec un nombre de degrés de liberté variable

Remarques

- La réalisation est une somme d'écart positifs : même si les hypothèses peuvent faire penser le contraire, ce test est unilatéral
- Ce test n'est valable que si tous les effectifs théoriques t_i sont ≥ 5 (quantité arbitraire mais communément admise). Si un des t_i est < 5 , il faut regrouper des classes pour augmenter les t_i trop faibles ou bien utiliser la correction de Yates

Correction de Yates Si un des t_i est < 5 et que le nombre de degrés de liberté $\nu = 1$, on peut utiliser le test du χ^2 avec la réalisation de V.A. de test suivante : $k = \sum_{\text{modalités}} \frac{(|o_i - t_i| - \frac{1}{2})^2}{t_i}$

5.2 Comparaison d'une distribution observée à une distribution théorique : test d'ajustement

On teste l'adéquation entre observations et loi théorique pour la V.A étudiée.

Tableau de contingence

Modalités \ Effectifs	Observés	Théoriques
m_1	o_1	t_1
m_2	o_2	t_2
m_3	o_3	t_3
\vdots	\vdots	\vdots

— Effectifs théoriques : $t_i = \mathbb{P}(m_1) \cdot n$

Le test

- \mathcal{H}_0 : La V.A. X suit la loi prescrite
- \mathcal{H}_1 : Ce n'est pas le cas
- Degrés de liberté pour le $\chi^2 =$ nombre de lignes -1 -nombre de paramètres estimés dans la loi prescrite

5.3 Comparaison de plusieurs distributions observées : test d'indépendance

On teste l'indépendance entre deux V.A.

Tableau de contingence

$X \setminus Y$	y_1	y_2	y_3	\dots	Sommes
x_1	$o_{1,1}$	$o_{1,2}$	$o_{1,3}$	\dots	$o_{1,\cdot}$
x_2	$o_{2,1}$	\ddots			$o_{2,\cdot}$
x_3	$o_{3,1}$				$o_{3,\cdot}$
\vdots	\vdots				\vdots
Sommes	$o_{\cdot,1}$	$o_{\cdot,2}$	$o_{\cdot,3}$	\dots	n

— Effectifs théoriques : $t_{i,j} = \mathbb{P}(X = x_i \text{ et } Y = y_j) \cdot n = \mathbb{P}(X = x_i) \cdot \mathbb{P}(Y = y_j) \cdot n = \frac{o_{i,\cdot}}{n} \cdot \frac{o_{\cdot,j}}{n} \cdot n = \frac{o_{i,\cdot} \cdot o_{\cdot,j}}{n}$

Le test

- \mathcal{H}_0 : Les deux V.A. sont indépendantes
- \mathcal{H}_1 : Ce n'est pas le cas
- Degrés de liberté pour le $\chi^2 =$ (nombre de lignes -1)(nombre de colonnes -1)

5.4 Comparaison de plusieurs distributions observées : test d'homogénéité

On teste l'égalité de loi entre plusieurs V.A.

Tableau de contingence

Modalités \ Variables	X_1	X_2	X_3	...	Sommes
m_1	$o_{1,1}$	$o_{1,2}$	$o_{1,3}$...	$o_{1,\cdot}$
m_2	$o_{2,1}$	\ddots			$o_{2,\cdot}$
m_3	$o_{3,1}$				$o_{3,\cdot}$
\vdots	\vdots				\vdots
Sommes	$o_{\cdot,1}$	$o_{\cdot,2}$	$o_{\cdot,3}$...	n

— Effectifs théoriques : $t_{i,j} = (\text{proportion des } m_i) \cdot (\text{proportion des } X_j) \cdot n = \frac{o_{i,\cdot}}{n} \cdot \frac{o_{\cdot,j}}{n} \cdot n = \frac{o_{i,\cdot} \cdot o_{\cdot,j}}{n}$

Le test

- \mathcal{H}_0 : Les V.A. suivent la même loi
- \mathcal{H}_1 : Ce n'est pas le cas
- Degrés de liberté pour le $\chi^2 = (\text{nombre de lignes} - 1)(\text{nombre de colonnes} - 1)$

Comparaison de distributions - Exercices

Exercice 1 Test d'adéquation On étudie les notes obtenues par une promotion d'étudiants Génie Biologique en deuxième année. Les résultats, regroupés par classes, étaient les suivants :

Classes]0-4]]4-6]]6-8]]8-10]]10-16]
Effectifs	14	12	6	5	9

La moyenne précise de cet échantillon de 46 notes (avant regroupement par classes) est de 6.6, et l'écart-type de 3.7.

1. Estimer les paramètres de la population dont est issu cet échantillon.
2. Utilisez le test du χ^2 pour vérifier si cette population suit une loi normale avec ces paramètres.

Exercice 2 Test d'homogénéité On se propose de comparer les réactions produites par deux vaccins BCG désignés par les lettres A et B. Un groupe d'enfants a été divisé, par tirage au sort, en deux séries qui ont été vaccinées, l'une avec le vaccin A, l'autre avec le B. La réaction a été appréciée ensuite par une personne ignorant le vaccin utilisée.

- Dans le cas du vaccin A, on a noté (dans l'ordre d'importance des réactions) 12 réactions légères, 156 réactions moyennes, 8 ulcérations et 1 abcès.
- Avec le vaccin B, on a obtenu 29 réactions légères, 135 réactions moyennes, 6 ulcérations et 1 abcès.

Que pouvez vous conclure de ces résultats ?

Exercice 3 Test d'homogénéité Lors d'un contrôle de la mise en suspension de pénicilline retard en flacons siliconés, le mode opératoire consiste classer, selon des critères visuels, les flacons en trois catégories : A (suspension bonne), B (suspension correcte) et C (suspension mauvaise). Pour déterminer la subjectivité des critères, on soumet à trois contrôleurs (C1, C2 et C3), trois échantillons de 200 flacons tirés au hasard dans un même lot. Les résultats sont les suivants :

	A	B	C
C1	125	53	22
C2	112	59	29
C3	128	57	15

Peut-on dire que ces critères sont objectifs ?

Exercice 4 Test d'indépendance Afin d'étudier l'influence des rayons X sur la spermatogenèse des chenilles, on a irradié des mâles au 2^e jour du 4^e stade larvaire. Ces mâles ont été ensuite accouplés avec des femelles non irradiées. Un lot de mâles normaux a servi de contrôle. On a compté le nombre d'œufs fertiles dans les pontes de femelles. On a ainsi obtenu les résultats suivants :

	Nombre total d'œufs	Nombre d'œufs fertiles
Mâles irradiés	5646	4998
Mâles normaux	6351	6236

L'irradiation a-t-elle une influence sur la fertilité des œufs ?

Exercice 5 Test d'indépendance On a étudié, sur un ensemble de 100 individus pris au hasard, la présence ou l'absence de deux caractères A et B (V.A. binaires). On a trouvé :

- 40 individus possédant à la fois le caractères A et le caractères B.
- 50 individus possédant le caractère A.
- 60 individus possédant le caractère B.

Les deux caractères A et B sont-ils indépendants ?

Réponse 1 Test d'adéquation En utilisant la loi normale de moyenne 6.6 et d'écart-type 3.7409, on trouve les arrondis suivants :

Effectifs théoriques	9.42	8.87	9.64	7.93	8.08
Termes somme (k)	2.23	1.11	1.38	1.08	0.1

Enfin, avec $k = 5.9$ et $k_{5\%} = 5.991$ (DDL=2), on peut conclure de justesse que, au risque $\alpha = 5\%$, la distribution des notes pour ce module suit bien une loi normale de moyenne 6.6 et d'écart-type 3.74.

Réponse 2 Test d'homogénéité $k = 8.706$. Avec 2 DDL, les 3 zones de confiance disponibles sur la table du χ^2 permettent de conclure que les réactions sont différentes au risque de 5%.

Réponse 3 Test d'homogénéité $k = 5.976$. Avec 4 DDL, les 3 zones de confiance disponibles sur la table du χ^2 permettent de dire que les critères ne sont pas subjectifs au seuil de 5%.

Réponse 4 Test d'indépendance $k = 468.955$. Avec 1 DDL, les 3 zones de confiance disponibles sur la table du χ^2 permettent de dire que l'irradiation a une influence sur la fertilité au risque de 0.1%.

Réponse 5 Test d'indépendance $k = 16.667$. Avec 1 DDL, les 3 zones de confiance disponibles sur la table du χ^2 permettent de dire que les deux caractères sont liés au risque de 0.1%.

6 Validation de la normalité (2h)

De nombreux tests demandent à ce que la V.A. étudiée suive une loi normale. Comment s'en assurer ?

6.1 Vérification rapide de la symétrie

Avant une vérification sérieuse, on peut rapidement évaluer la symétrie et réfuter la normalité si $\frac{\text{moyenne}}{\text{mediane}} > 2$ (ou, ce qui est la même chose, si $\frac{\text{moyenne}}{\text{mediane}} < 0.5$).

Attention : comme cette vérification ne tient pas compte de la taille de l'échantillon, elle n'est que peu fiable avec des petits échantillons.

6.2 Graphiques

- Histogramme : graphe des (c_i, n_i) où les c_i sont les centres de classes qu'on choisira et les n_i sont les effectifs associés. L'histogramme doit ressembler à une cloche symétrique.
- Boxplot : on reporte min, Q_1 , Q_2 , Q_3 et max sur une droite puis on fait un rectangle allant de Q_1 à Q_3 . Les valeurs Q_1 , Q_2 et Q_3 doivent découper l'intervalle [min; max] en quatre parties semblables et Q_2 doit être proche de la moyenne.
- QQplot : graphe des (x_i, y_i) où les x_i sont les mesures de la V.A. X étudiée et les y_i sont les quantiles théoriques tels que $\mathbb{P}(Y \leq y_i) = \frac{\text{nombre de mesures} \leq x_i}{\text{nombre de mesures}}$ où Y suit la loi normale $\mathcal{N}(0, 1)$. Ce graphe doit ressembler à une droite.

6.3 Tests statistiques

Pour trancher, quand les graphiques ne sont pas clairement disqualifiants

6.3.1 Pour des petits échantillons ($n < 50$) : Shapiro-Wilk

- \mathcal{H}_0 : les mesures correspondent à une variable aléatoire suivant une loi normale
- \mathcal{H}_1 : non

$$SW = \frac{\left(\sum_{j=1 \dots k} a_j d_j \right)^2}{\sum_{i=1 \dots n} (x_i - \bar{x})^2} \text{ où}$$

6 Validation de la normalité (2h)

- n est l'effectif de l'échantillon
- $k = n/2$ arrondi à l'inférieur
- a_j est lu sur une table
- x_j sont les mesures ordonnées de la plus petite à la plus grande
- $d_j = x_{n+1-j} - x_j$: on liste les mesures décroissantes et on soustrait la liste des mesures croissantes
- \bar{x} : moyenne des mesures
- On rejette \mathcal{H}_0 si la valeur SW est inférieure à la valeur lue sur la table de Shapiro

Remarque $\sum (x_i - \bar{x})^2 = (n-1)\sigma_{est}^2$

6.3.2 Pour des grands échantillons ($n > 50$) : Kolmogorov-Smirnov

Remarques

- Ce test est aussi valable pour des petits échantillons, mais alors il faut d'autres valeurs critiques que celles présentées ici
- Ce test permet de vérifier si on suit n'importe quelle loi, pas seulement la normale
- De fait, ceci est un test *non paramétrique* : on n'a pas besoin que X suive une loi normale

Le test

- \mathcal{H}_0 : les mesures correspondent à une variable aléatoire suivant la loi $\mathbb{P}_{the}(X < x)$
- \mathcal{H}_1 : non
- $KS = \max \left| \mathbb{P}_{the}(X \leq x_i) - \mathbb{P}_{obs}(X \leq x_i) \right|$ pour toutes les mesures x_i
- On rejette \mathcal{H}_0 si la valeur KS est supérieure à $\begin{cases} \frac{1.358}{\sqrt{n}} & \text{au risque } \alpha = 5\% \\ \frac{1.518}{\sqrt{n}} & \text{au risque } \alpha = 2\% \\ \frac{1.629}{\sqrt{n}} & \text{au risque } \alpha = 1\% \end{cases}$

6.4 Valeurs aberrantes sous l'hypothèse de normalité : Grubs

Remarques

- Les valeurs aberrantes peuvent provenir d'erreurs de mesure, d'individus étrangers,...
- En général, on ne peut détecter que des écarts évidents. L'intérêt est le côté automatique de la procédure.
- Cette détection, et l'élimination éventuelle de valeurs aberrantes (*outliers*) peut servir deux buts :
 - Rapprocher le comportement des mesures de celui d'une vraie loi normale
 - Réduire la variance de l'échantillon et/ou d'avoir des variances égales entre plusieurs échantillons

Le test

- \mathcal{H}_0 : les valeurs suspectes ne sont pas aberrantes
- \mathcal{H}_1 : si
- $G = \frac{SCE^*}{SCE}$ où

6.4 Valeurs aberrantes sous l'hypothèse de normalité : Grubs

- $SCE = \sum_{x_i \in E} (x_i - x.)^2$, $x.$ étant la moyenne de E (l'échantillon complet)
- $SCE^* = \sum_{x_i \in E^*} (x_i - x.)^2$, $x.$ étant la moyenne de $E^* = E \setminus \{\text{valeurs suspectes}\}$
- On rejette \mathcal{H}_0 si la valeur G est inférieure à la valeur lue sur la table de Grubs

Validation de la Normalité - Exercices

Exercice 1 Etude d'un petit échantillon

6.70	6.87	6.92	6.95	6.95	6.96	6.96	6.99
7.02	7.03	7.05	7.05	7.06	7.08	7.10	7.13

1. Etudier la normalité de cette série de données (histogramme, boîte à moustaches, QQPlot, test de Shapiro).
2. La valeur 6.70 peut-elle être considérée comme aberrante au risque 5% ?
3. Dans l'affirmative, comparer la variance avec et sans la valeur 6.7, et refaire l'étude de la normalité sans la valeur 6.70

Exercice 2 Exemple pratique Etudions la normalité des données de l'exercice 4, page 17

1. Divisez moyenne par médiane pour avoir un premier avis sur la symétrie des données.
2. Commentez l'histogramme et le QQ-plot en Figure 6.1.
3. Confirmez par un test statistique ($\sum a_j d_j = 28.3245$ pour Shapiro, et $KS = 0.1515$ pour Kolmogorov)
4. Pour cette même valeur $KS = 0.1515$, à partir de quelle taille d'échantillon aurait-on conclu que la loi de ces données n'est pas la loi normale de moyenne 10.4551 et d'écart-type 4.1492 ?

Exercice 3 Etude d'un grand échantillon

6.0	8.9	9.2	9.6	9.9	10.2	10.5	10.9	11.4	12.0
7.1	9.0	9.2	9.6	9.9	10.2	10.5	10.9	11.5	12.0
7.7	9.0	9.2	9.7	9.9	10.3	10.5	10.9	11.5	12.1
8.0	9.0	9.2	9.7	9.9	10.3	10.6	10.9	11.5	12.1
8.0	9.0	9.3	9.7	10.0	10.3	10.6	10.9	11.8	12.2
8.2	9.1	9.4	9.8	10.0	10.3	10.6	11.0	11.8	12.3
8.3	9.1	9.4	9.8	10.0	10.4	10.7	11.0	11.8	12.5
8.3	9.1	9.5	9.8	10.0	10.4	10.7	11.0	11.9	12.7
8.6	9.1	9.5	9.8	10.1	10.4	10.8	11.1	11.9	12.7
8.7	9.1	9.6	9.9	10.1	10.4	10.8	11.4	12.0	14.0

1. En se basant sur les graphiques de la Figure 6.2, peut-on s'attendre à une distribution normale ?
2. Valider cette impression par un test de Kolmogorov. ($\mu_{\text{est}} = 10.2120$, $\sigma_{\text{est}} = 1.3073$, $KS = 0.05$).
3. Peut-on considérer les valeurs extrêmes (6 et 14) comme aberrantes au risque 5% ? Comparez les estimations de la variance avec et sans ces deux valeurs. ($SCE = 169.2056$, $SCE^* = 137.1139$)

6 Validation de la normalité (2h)

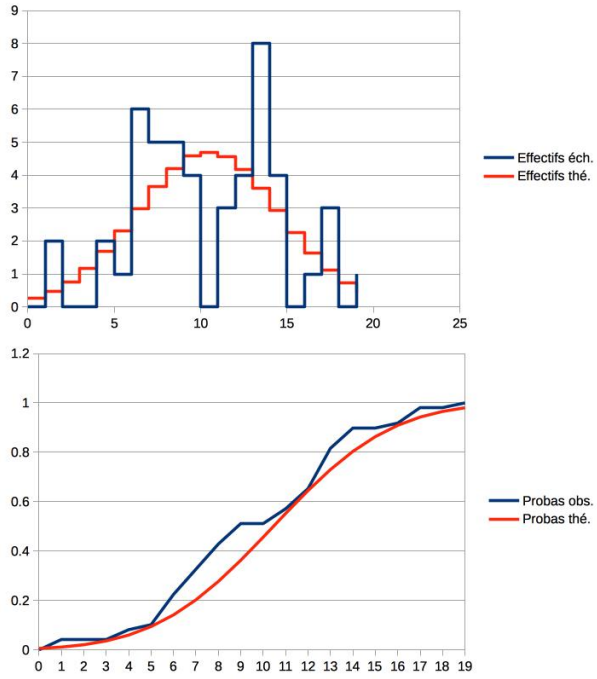


FIGURE 6.1 – Histogramme (gauche) et QQ-plot (droite) des notes de Statistiques

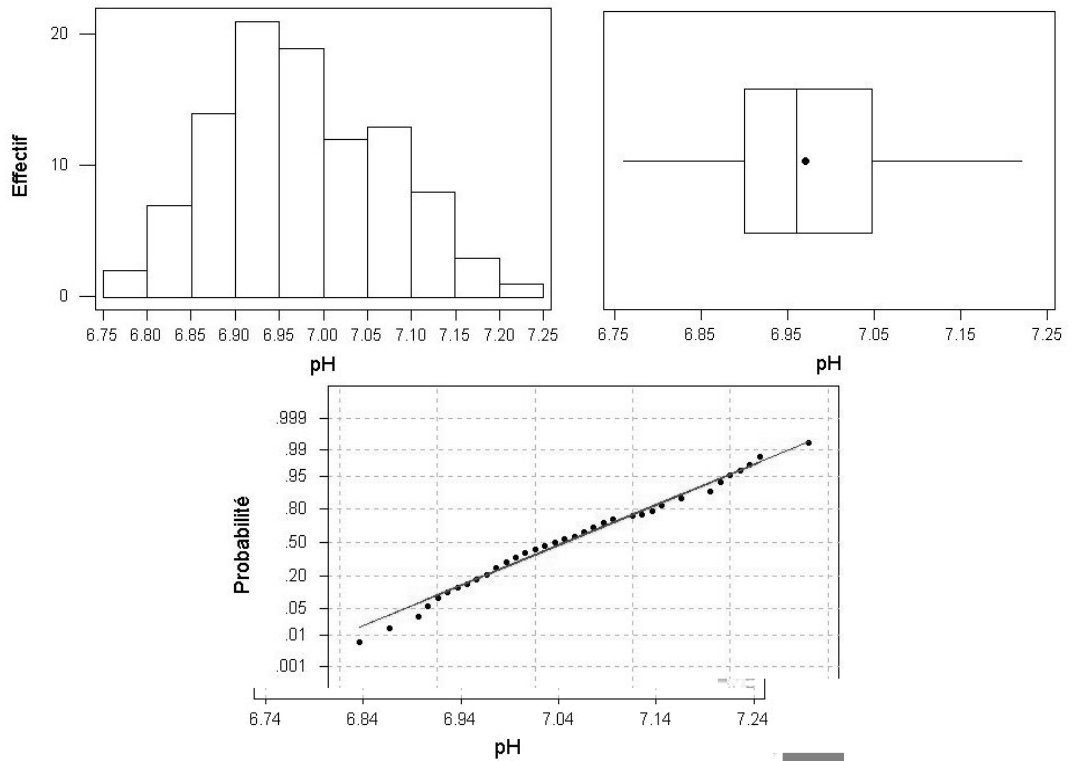


FIGURE 6.2 – Histogramme, boxplot et qqplot pour les données de l'exercice 2.

Réponse 1 Etude d'un petit échantillon

1. Graphiques : Figure 6.3.

Somme des $a_j d_j = 0.38$, $\mu_{\text{éch}} = 6.988$, $SCE = 0.16$, $SW = 0,901$, $SW_{5\%} = 0,837$. On ne peut donc pas dire au risque $\alpha = 5\%$ que la population, donc l'échantillon est tirée, n'est pas issue d'une Loi Normale.

2. $\mu_{\text{éch}}^* = 7.008$, $SCE^* = 0.07$, $G = 0.453$, $G_{5\%} = 0.5755$, $V_{1\%} = 0.4634$. On peut donc conclure avec un risque $\alpha = 1\%$ que la valeur 6.7 est aberrante.
3. $\sigma_{\text{est}} = 0.0108$ et $\sigma_{\text{est}}^* = 0.0052$: deux fois moins.

Moyenne et SCE sont les version étoilées de la question précédente, seule la somme des $a_j d_j$ doit être calculée : 0.26. Ainsi on trouve $SW = 0.975$ et on lit $SW_{5\%} = 0.881$. Sans surprise, on conserve donc encore la conclusion que la population étudiée suit une loi normale au risque $\alpha = 5\%$. Cependant, SW est plus grand qu'auparavant. La conclusion s'en trouve donc renforcée.

Réponse 2 Exemple pratique

1. moy/méd = 1.046, ce qui est très proche de 1 : la symétrie semble très forte.
2. Comme dans beaucoup de cas réels, c'est pas très clair : Le trou en 10 et le pic en 13 donnent une impression de profil en M, au lieu d'une cloche. Cependant, vu la largeur du trou, il pourrait seulement s'agir de hasard ponctuel d'échantillonnage. En effet, l'histogramme suit quand même assez bien la courbe d'une loi normale avec la moyenne et l'écart-type estimés. Pour un faible effectif (49), cet histogramme semble ainsi assez correct.
De même, le QQ-plot est un peu chahuté en son centre, mais ne présente pas de coude franc et semble donc aller en faveur de données normalement distribuées.
3. Par pur hasard (véridique) la taille de l'échantillon est à la limite entre Shapiro et Kolmogorov. On pourra donc faire l'un ou l'autre. Pour Shapiro, on trouve $SW_{5\%} = 0,947$ et $SW = 0.9709$: on accepte donc la normalité des données ; De justesse. Pour Kolmogorov, le seuil étant de $KS_{5\%} = 0.194$, on conclut au risque de 5% que les données suivent une loi normale de moyenne 10.4551 et d'écart-type 4.1492. L'acceptation plus franche de Kolmogorov est à mettre sur la petitesse de l'échantillon : moins on a de données, plus les tests généraux comme celui-ci sont permissifs.
4. La formule $0.1515 = 1.358/\sqrt{n}$ donne $n > 80$.

Réponse 3 Etude d'un grand échantillon

1. L'histogramme présente une forme en cloche caractéristique. De plus, le boxplot indique des données fortement symétriques. Enfin, le qqplot nous montre des points fortement alignés, même dans les valeurs extrêmes. On peut donc raisonnablement s'attendre à ce que la V.A. suive une loi normale.
2. $KS_{5\%} = 0.1358$. Au risque $\alpha = 5\%$, rien ne vient donc infirmer l'hypothèse que la V.A. suive une loi normale avec la moyenne et l'écart-type estimé.
3. $G = 0.810$, $G_{5\%} = 0.821$ et $G_{1\%} = 0.794$. On peut donc affirmer au risque $\alpha = 5\%$ que les valeurs extrêmes sont aberrantes, mais pas au risque de 1%.
 $\sigma_{\text{est}} = \frac{SCE}{n-1} = 1.70$ et $\sigma_{\text{est}}^* = \frac{SCE^*}{n-3} = 1.41$: c'est moins, mais pas dramatiquement.

6 Validation de la normalité (2h)

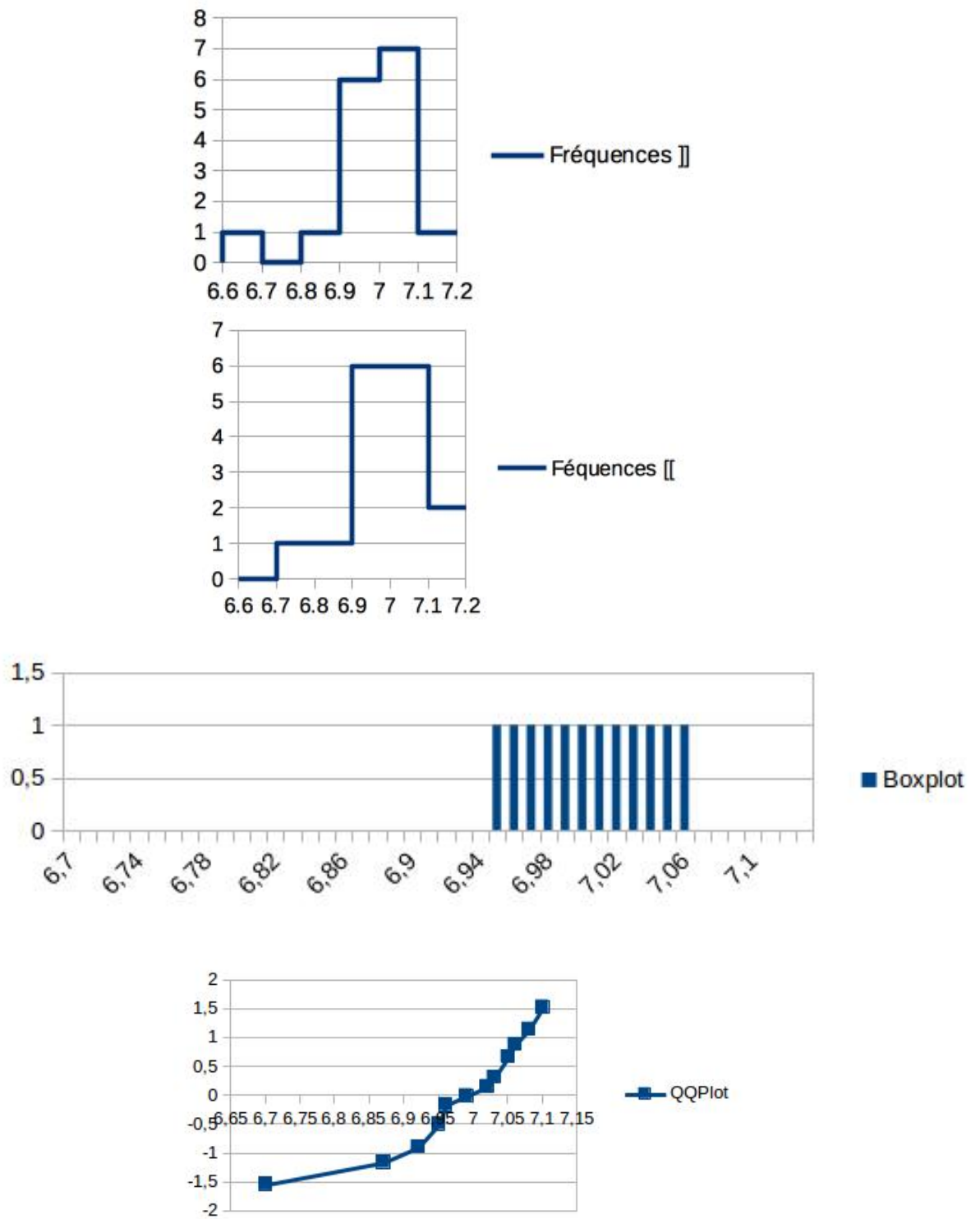


FIGURE 6.3 – Histogrammes, boxplot et qqplot pour les données de l'exercice 1

7 Comparaison de $m > 2$ populations normalement distribuées (4h + 3 TP)

- On dispose de m échantillons, de tailles respectives n_j (et $n = \sum n_j$) identifiés par un (seul) facteur
- En faisant les comparaisons deux à deux, les risques se cumulent. Il faut donc un test global.

Exemple Si on compare les données issues de différents laboratoires, mais avec le même protocole pour chacun, chaque jeu de mesures n'est donc identifié que par son laboratoire d'origine (un seul facteur) :

Valeurs \ Échantillons	E_1	E_2	...	E_m
m_1	$x_{1,1}$	$x_{1,2}$...	
m_2	$x_{2,1}$	$x_{2,2}$		
m_3	$x_{n_1,1}$	$x_{3,2}$	\ddots	
m_4		$x_{n_2,2}$	\ddots	
				\vdots
Moyennes	$x_{.,1}$	$x_{.,2}$...	$x_{.,m}$
Variances	$\sigma_{1,est}^2$	$\sigma_{2,est}^2$...	$\sigma_{m,est}^2$

7.1 Égalité des variances

Il existe différents tests, chacun ayant ses avantages, dont on en présente deux

Remarque Avant de se lancer dans une longue vérification, on peut gagner du temps en vérifiant rapidement si $\frac{\max(\sigma_{est,j}^2)}{\min(\sigma_{est,j}^2)} > 2$ (auquel cas, il est fort probable que les variances seront déclarées comme différentes)

7.1.1 Test de Bartlett

Pour des échantillons de tailles distinctes

Le test

- \mathcal{H}_0 : les variances des m populations sont égales
- \mathcal{H}_1 : au moins une diffère significativement des autres

7 Comparaison de $m > 2$ populations normalement distribuées (4h + 3 TP)

- $B = \frac{1}{c} \left((n - m) \ln(d) - \sum_{j=1 \dots m} (n_j - 1) \ln(\sigma_{est,j}^2) \right)$ avec la fonction Excel SOMMEPROD des ddl et des LN des variances) où
 - $n = \sum_{j=1 \dots m} n_j$
 - $c = 1 + \frac{1}{3(m - 1)} \sum_{j=1 \dots m} \left(\frac{1}{n_j - 1} - \frac{1}{n - m} \right)$ (fonction Excel SOMMEPROD des 1/ddl)
 - $d = \frac{1}{n - m} \sum_{j=1 \dots m} (n_j - 1) \sigma_{est,j}^2$ (fonction Excel SOMMEPROD des ddl et variances)
- On rejette \mathcal{H}_0 si la valeur B est supérieure à $k(\alpha)$, le quantile de la loi du $\chi^2(m - 1)$ tel que $\mathbb{P}(K > k) = \alpha$, $K \sim \chi^2(m - 1)$

Remarques

- C'est un test paramétrique qui demande à ce que chaque échantillon soit normalement distribué – dans le cas contraire, le test de Levene est mieux adapté.
- De plus, une condition d'application de ce test est que tous les $n_j > 6$.
- Contrairement à ce que les hypothèses peuvent laisser penser, c'est un test unilatéral

7.1.2 Test de Cochran

Pour des échantillons de même taille

Le test

- \mathcal{H}_0 : les variances des m populations sont égales
- \mathcal{H}_1 : au moins une diffère significativement des autres
- $C = \frac{\max_{j=1 \dots m} (\sigma_{est,j}^2)}{\sum_{j=1 \dots m} (\sigma_{est,j}^2)}$
- On rejette \mathcal{H}_0 si la valeur C est supérieure à la valeur lue sur la table de Cochran

Remarques

- C'est un test paramétrique qui demande à ce que chaque échantillon soit normalement distribué – dans le cas contraire, le test de Levene est mieux adapté.
- Contrairement à ce que les hypothèses peuvent laisser penser, c'est un test unilatéral

7.2 Egalité des moyennes – ANOVA

Conditions d'application et remarques

1. Chaque échantillon doit être issu d'une V.A. suivant une loi normale – ce critère est théoriquement nécessaire, mais le test est relativement souple à cet égard et on peut se contenter d'accepter la normalité à 2%, au lieu du traditionnel 5%.
2. Les échantillons doivent être indépendants

3. Les variances des populations pour chaque échantillon doivent être égales – ceci est par contre très important et doit être validé à 5% au minimum

4. Comme précédemment, le test est unilatéral

5. Décomposition de la quantité de variation :
$$\underbrace{\sum_{i,j} (x_{i,j} - x_{\cdot,\cdot})^2}_{Q_T} = \underbrace{\sum_j \left(\sum_i (x_{i,j} - x_{\cdot,j})^2 \right)}_{Q_I} + \underbrace{\sum_j n_j (x_{\cdot,j} - x_{\cdot,\cdot})^2}_{Q_E}$$

6. Décomposition des degrés de liberté :
$$\underbrace{n-1}_{Total} = \underbrace{n-m}_{Interne} + \underbrace{m-1}_{Externe}$$

7. Si l'égalité des variances n'est pas assez franche, on peut essayer de détecter et supprimer des valeurs aberrantes (test de Grubs)

Le test

- \mathcal{H}_0 : les moyennes des m populations sont égales
- \mathcal{H}_1 : au moins une diffère significativement des autres
- $F = \frac{V_E}{V_I} \sim Fisher(m-1; n-m)$ (test unilatéral) où
 - $V_I = \frac{Q_I}{n-m}$ est la variance interne.
Cas particulier : si tous les n_j sont égaux, $V_I =$ moyenne des variances estimées.
 - $V_E = \frac{Q_E}{m-1}$ est la variance externe.
Cas particulier : si tous les n_j sont égaux, $V_E = n_j$ fois la variance estimée des moyennes.

Remarques

- $Q_T = (n-1)\sigma_{est}^2(\{x_{i,j}\})$ (la variance estimée de toutes les valeurs confondues : fonction Excel VAR)
- $Q_I = \sum (n_j - 1)\sigma_{est,j}^2$ (les variances estimées de chaque variable : fonction Excel SOMMEPROD)
- $Q_E = Q_T - Q_I$

7.3 Répétabilité et reproductibilité

- L'analyse de la variance peut aussi être utilisée pour étudier la fidélité d'un protocole de mesure sur plusieurs laboratoires
- Si la moyenne exacte des mesures est connue, une méthode de mesure peut être validée sur plusieurs laboratoires

Répétabilité La variance de répétabilité V_r indique la fidélité des mesures dans des conditions identiques. C'est la moyenne pondérée des variances obtenues dans chaque laboratoire :

$$V_r = V_I$$

- Ecart-type de répétabilité : $\sqrt{V_r}$
- Coefficient de variation de répétabilité : $\frac{\sqrt{V_r}}{x_{\cdot,\cdot}}$

Reproductibilité La variance externe V_E est égale à la somme de la variance interne et de toutes les variances inter-laboratoire : $V_E = V_I + nV_L$, où $n = \frac{1}{m-1} \left(n - \frac{\sum (n_j^2)}{n} \right)$. Cas particulier : si tous les n_j

7 Comparaison de $m > 2$ populations normalement distribuées (4h + 3 TP)

sont égaux, alors $n = n_j$. Ainsi, la variance inter-laboratoire est

$$V_L = \frac{V_E - V_I}{n}$$

Remarque Cette valeur n'est pas une vraie variance à proprement parler et pourra parfois être négative

La variance de reproductibilité V_R indique la fidélité des mesures d'un laboratoire à un autre. C'est la somme de la variance interne et de la variance inter-laboratoire :

$$V_R = V_I + V_L$$

- Ecart-type de reproductibilité : $\sqrt{V_R}$
- Coefficient de variation de reproductibilité : $\frac{\sqrt{V_R}}{x_{.,}}$

Le test On compare la moyenne générale $\mu(X)$, estimée par $x_{.,}$, et μ_0 , la moyenne attendue :

- \mathcal{H}_0 : La méthode de mesure est exacte ($\mu(X) = \mu_0$)
- \mathcal{H}_1 : Elle ne l'est pas ($\mu(X) \neq \mu_0$)
- $T = \frac{\mu(X) - \mu_0}{\sqrt{\frac{V_R}{n}}} \sim Student(n - 1)$ (cette fois-ci, le test est bilatéral comme il se doit)

Comparaison de plus de 2 moyennes - Exercices

Exercice 1 Analyse de la variance à un facteur contrôlé et Application aux circuits interlabos 5 laboratoires participent à une campagne pour évaluer les limites de répétabilité et de reproductibilité d'une dosage.

Ces 5 labos dosent (3 fois) une solution en aveugle (le résultat est inconnu) afin de quantifier les limites de répétabilité et de reproductibilité.

Voici les résultats obtenus :

Participants	Répétitions			Moyenne	Ecart-type (est)	Variance (est)	SCE
Labo 1	7,77	7,74	7,52	7,6767	0,1365	0,0186	0,0373
Labo 2	7,14	7,46	7,21	7,2700	0,1682	0,0283	0,0566
Labo 3	7,42	7,27	7,22	7,3033	0,1041	0,0108	0,0217
Labo 4	7,64	7,64	7,64	7,6400	0,0000	0,0000	0,0000
Labo 5	6,22	6,28	6,29	6,2633	0,0379	0,0014	0,0029

- SCE des 5 moyennes : 1.3090 (moyenne des 5 : 7.2307)
- SCE des moyennes 1-4 : 0.1394 (moyenne des 4 : 7.4725)
- SCE des 12 premières répétitions : 0.5336

1. Peut-on dire que le labo 5 fournit des résultats *en moyenne* aberrants? Si oui, le supprimer de l'étude.
2. On suppose que, par labo, les résultats sont normalement distribués et que les labos sont indépendants, au moins au sens statistique du terme. Par conséquent, pour faire une analyse de la variance, il nous suffit de vérifier l'égalité statistique des variances.
3. Isoler les quantités de variation interne et externe et conclure en ce qui concerne la différence des valeurs obtenues en moyenne par les laboratoires étudiés.
4. Calculer les variances de répétabilité et de reproductibilité.
5. Admettons que la moyenne générale qu'on aurait du avoir soit $\mu_0 = 7.42$. Peut alors dire que le protocole utilisé pour ces mesures (toujours en excluant le labo 5) ait été fiable?
6. Et si on avait conservé les résultats du labo 5? (il faut recalculer V_R et les degrés de liberté changent)

Réponse 1

1. $G = 0.1065$ et $G_{5\%} = 0.127$: La moyenne du labo 5 est aberrante au risque $\alpha = 5\%$ (mais on ne peut pas l'affirmer au risque $\alpha = 1\%$).
2. Cochran : 0.4899, valeur critique : 0.7679. On ne peut donc pas conclure à la différence des variances au risque $\alpha = 5\%$.

7 Comparaison de $m > 2$ populations normalement distribuées (4h + 3 TP)

3. $V_E = 0.1394$, $V_I = 0.0144$, $F = 9.65$, $F_{1\%} = 7.59$. On peut donc conclure à la différence des moyennes au risque $\alpha = 1\%$: il y a une différence entre les laboratoires, même après avoir éliminé le Labo 5.
4. $V_r = 0.0144$. Tous les échantillons ont la même taille, ainsi $n = 3$, $V_L = 0.0416$ et $V_R = 0.0561$.
5. $T = 0.7680$ et $T_{5\%} = 1.7959$: au risque $\alpha = 10\%$, le protocole de mesure utilisé est fiable.
6. Les différentes variances : $Q_E = 3.9271$, $V_E = 0.9818$, $V_L = 0.3233$, $Q_I = 0.1184$, $V_I = 0.0118$, $V_R = 0.3352$.

Ainsi : $T = -1.2666$ et $T_{5\%} = 1.7613$: on voit que la valeur de la variable de test s'écarte franchement plus de 0 que sans le labo 5, mais que même au risque $\alpha = 10\%$ on ne peut pas dire que le protocole de mesure utilisé n'est pas fiable.

Indicateurs de base

Bactérie	Eubacteria	Enterobacteria	Bacteroides	Firmicutes	Enterococcus	Clostridium	Bacteroides	Lactobacillus
Max	12.81210847	10.4845809	11.37883048	13.31213684	13.81210847	12.64180098	10.3673	12.19376759
Moyenne	11.0996261659	8.9105892531	9.5733451013	11.4464190855	12.0996261659	11.0870052944	9.099597	8.787623345
Médiane	11.15711114	9.032087816	9.78635422	11.406615255	12.15711114	11.09833924	9.096666	8.349684072
Min	7.64208771	6.577288143	7.055028284	9.856283889	8.64208771	9.639566317	8.008355	7.094285163
Variance	0.8234285433	0.9060647377	1.0265131883	0.7100586245	0.8234285433	0.522049417	0.41534	1.3527379895
Ecart type	0.9074296354	0.9518743287	1.0131698714	0.8426497638	0.9074296354	0.7225298728	0.644468	1.1630726501
Effectif	97	99	96	94	97	100	87	93
DDL	96	98	95	93	96	99	86	92

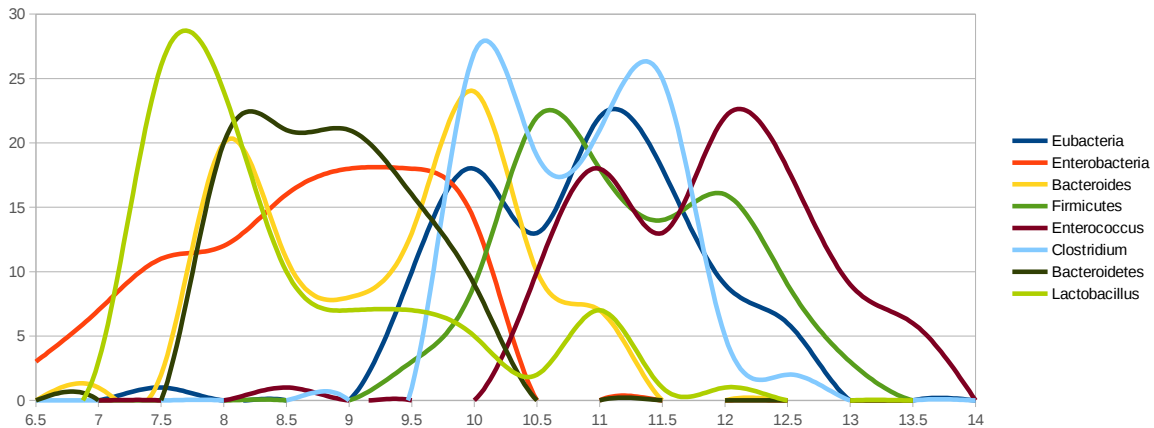
Validation de la normalité des échantillons

Vérifier rapidement la symétrie de chaque échantillon

Divisions	0.9948476829	0.9865481198	0.9782340682	1.0034895391	0.9952714939	0.998978771	1.000322	1.0524498016
Moyenne sur médiane	Semble correct, il est donc pertinent de faire des graphiques pour une vérification plus poussée							

Tracer sur le même diagramme tous les histogrammes pour des classes d'effectifs de largeur 0.5

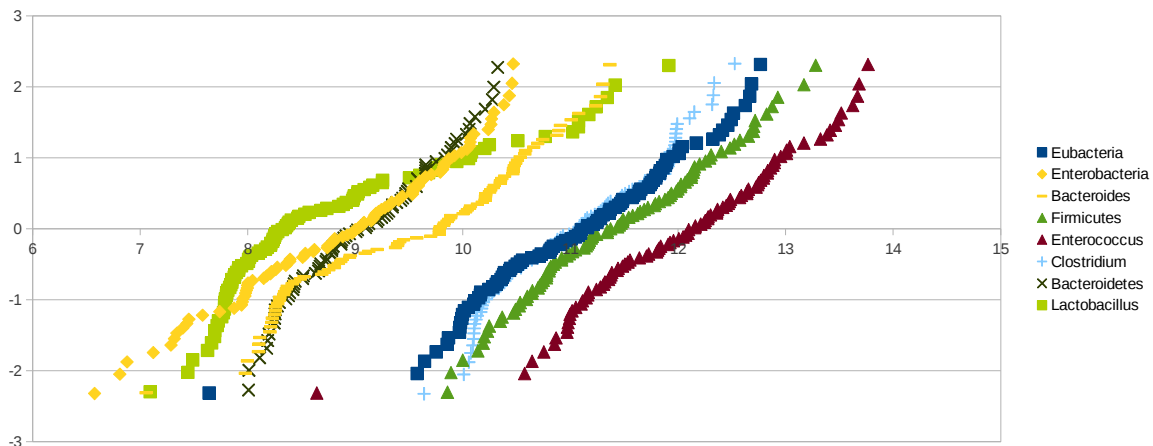
- Calculer le min et le max global pour établir les bornes des classes
- Calculer les cumulés avec la fonction `sommeprod` et vérifier l'effectif total



Remarque : on devrait faire des courbes en escalier, mais ici elles se chevaucheraient trop et on ne verrait rien – d'où ces versions « lissées »
On voit que Lactobacillus et Bacteroides n'ont pas franchement une forme de cloche

Tracer sur un même diagramme tous les QQPlots

- Calculer les probas observées à partir d'un tableau des effectifs cumulés pour chaque valeur x_i de chaque bactérie
- Calculer les quantiles théoriques à partir de la loi normale standard pour ces probas observées
- Ajouter les séries de données une par une à un diagramme de type XY



(ne pas faire relier les points, pour éviter les aller-retours des x_i qui ne sont pas dans un ordre différent pour chaque série de données)

Première remarque : il ne faut pas tenir compte des valeurs aux bouts. A cause de l'échantillonnage, elles sont généralement tordues.

On voit que Lactobacillus (encore) et Bacteroides (encore) sont franchement courbés au centre. Les autres sont « acceptables ».

Effectuer un test de Kolmogorov pour chaque échantillon

- Calculer les probas théoriques en utilisant les moyennes et écarts-type estimés pour chaque bactérie
- Calculer le tableau des valeurs absolues (abs) des différences entre les deux tableaux et remplacer manuellement les #valeur par des NA
- Calculer le max de chaque colonne des valeurs absolues des différences

	Eubacteria	Enterobacteria	Bacteroides	Firmicutes	Enterococcus	Clostridium	Bacteroides	Lactobacillus
KS	0.0560998025	0.0559216209	0.1123450995	0.0764370457	0.0560998025	0.0897089959	0.084196	0.1770771422
KS 5 %	0.1378840092	0.1364841353	0.1386002946	0.1400670052	0.1378840092	0.1358	0.145593	0.1408180401
Conclusion :	Acceptation	Acceptation	Acceptation	Acceptation	Acceptation	Acceptation	Acceptatio	Rejet

On constate que seule Lactobacillus est rejetée. Ce test est donc assez permissif et montre à quel point on peut être tolérant avec les graphiques.

Egalité des Variances

Validation rapide de l'égalité des variances, en écartant Lactobacillus

max sur min 2.4715032924 C'est fort, et on voit que les valeurs sont en effet assez différentes.

Validation rapide de l'égalité des variances, en écartant Lactobacillus et aussi les Bacteroidetes

max sur min 1.9663142124 C'est mieux, mais vraiment limite

Test de Bartlett (en écartant Lactobacillus et Bacteroidetes)

m 6
N-m 577
d 0.8009180621
c 1.003467921
B 12.5403788042
Khi2(5/100, 5 DDL) 11.0704976935

Conclusion : le test de Bartlett réfute l'égalité des variances au risque de 5%, mais guère mieux. Ecarter Clostridium de l'étude devrait suffire.

Test de Bartlett (en écartant Lactobacillus, Bacteroidetes et Clostridium)

m 5
N-m 478
d 0.8586753757
c 1.003488014
B 3.484307735
Khi2(40/100, 4 DDL) 4.0446264906

Conclusion : le test de Bartlett accepte l'égalité des variances au risque $\alpha=40\%$, ce qui est très fort.

Egalité des moyennes – ANOVA – en écartant Lactobacillus, Bacteroidetes et Clostridium

QT	1103.7663622936	VI	0.8586753757
QI	410.4468295669	VE	173.3298831817
QE	693.3195327266		

m 5
N-m 478
F 201.8572886731
Valeur seuil à 5 % 2.390591367
Valeur seuil à 0,1 % 4.6970712879

On rejette donc l'hypothèse d'égalité des moyennes avec presque sûrement aucune chance de se tromper

Pour la suite, admettons que les mesures ne concernent qu'une seule bactérie et que les noms des colonnes désignent en fait des noms de labos
On vérifie alors si le protocole de mesure pour cette bactérie fictive est fiable

Répétabilité – en se basant encore sur les 5 colonnes de la partie précédente

Moyenne générale	10.6159013168
Ecart type de répétabilité	0.9266473847
Coefficient de variation de répétabilité (%)	8.73%

Reproductibilité

n	483
n gothique	96.5931677019
VL	1.7855425172
VR	2.6442178929
Ecart type de reproductibilité	1.6261051297
Coefficient de variation de reproductibilité (%)	15.32%

Fiabilité de la mesure avec une moyenne attendue de : 10

T	-8.3240767937
T 0.1 %	3.3110398369

On constate que les mesures ne seraient certainement pas fiables (au risque de 0.1 %).

Ce qui est normal puisque les valeurs proviennent de mesures qui n'ont rien à voir entre elles et que les moyennes ont été déclarées hyper significativement différentes

8 Cartes de contrôle (2h + 2TP)

- Détecte les dérives lors d'un suivi
- Processus discontinus – dérives brusques : cartes de Shewhart (ex : chaîne de production d'objets bien distincts)
- Processus continu – dérives lentes : cartes de Cumsum (ex : production à grande vitesse, ou de fluides) – Non présenté ici –
- Il existe plusieurs sortes de suivi : respecter des bornes, respecter une moyenne, ...

8.1 Préliminaires

On étudie d'abord m échantillons, mesurés à des intervalles de temps assez rapprochés pour déterminer la loi de la production qu'on va suivre

- Etablir les paramètres de la loi de la production
 - Si les paramètres de la production (moyenne et écart-type) sont spécifiés, on prendra bien entendu ces valeurs pour μ_* et σ_* .
 - Si on spécifie un min (a) et un max (b) à respecter, on prendra $\mu_* = \frac{a+b}{2}$.
 - Enfin, en l'absence de données, on prendra comme paramètres pour la production $\mu_* = \mu_{est}$ et $\sigma_*^2 = V_I$, la variance interne.
- Valider la normalité de la production
 - Idéalement, il faudrait valider la normalité de chaque échantillon – cependant, pour des échantillons trop petits ($n_j < 20$?) ces validations n'ont plus vraiment de sens et il est préférable de valider la normalité sur l'ensemble des valeurs d'un coup
- Valider la stabilité de la production (égalité des variances)
 - Si les variances ne sont pas égales, la production n'est pas fiable dès le début et effectuer un suivi n'a pas de sens
 - Si les variances sont égales, alors une bonne estimation de la variance pour la production est la variance interne V_I .

Une fois la normalité et les paramètres μ_* et σ_* de la production globale validés, on peut donc admettre que la production suit la loi $\mathcal{N}(\mu_* ; \sigma_*)$.

8.2 Carte de contrôle

On prépare ensuite les graphiques sur lesquels on pourra effectuer le suivi de la production, avec des échantillons successifs de taille n (qui ne sera donc pas le même n que dans la phase préliminaire)

8.2.1 Carte de contrôle pour la moyenne

- Intervalle de fluctuation pour les estimations de la moyenne : $\left[\mu_* - u\left(\frac{\alpha}{2}\right) \frac{\sigma_*}{\sqrt{n}} ; \mu_* + u\left(\frac{\alpha}{2}\right) \frac{\sigma_*}{\sqrt{n}} \right]$

8 Cartes de contrôle (2h + 2TP)

- Les recommandations AFNOR sont $\alpha = 5\%$ pour les limites de surveillance et $\alpha = 0.2\%$ pour les limites de contrôle : 4 bornes d'intervalle + le centre = 5 droites horizontales
- On trouve parfois dans la littérature une version simplifiée avec la valeur 2 au lieu de $u(2.5\%)$ et la valeur 3 au lieu de $u(0.1\%)$
- Cas particulier : si on a spécifié un min (a) et un max (b) à respecter, les bornes sont donc déjà données. On calcule alors la bonne taille d'échantillon n pour que les bornes de l'intervalle ci-dessus correspondent à a et b avec $\alpha = 5\%$ (l'intervalle avec $\alpha = 0.2\%$ se calcule ensuite en suivant la formule classique)

8.2.2 Carte de contrôle pour la variance

- Intervalle de fluctuation pour les estimations de la variance : $\left[\frac{\sigma_*^2 k (1 - \frac{\alpha}{2})}{n - 1} ; \frac{\sigma_*^2 k (\frac{\alpha}{2})}{n - 1} \right]$
- Les recommandations AFNOR sont $\alpha = 5\%$ pour les limites de surveillance et $\alpha = 0.2\%$ pour les limites de contrôle : 4 bornes d'intervalle = 4 droites horizontales

8.3 Suivi de la production

- On reporte enfin, sur leur carte respective, la moyenne et la variance des échantillons prélevés à intervalles réguliers
- Si un point sort de la zone de surveillance, il faut refaire immédiatement un prélèvement pour confirmer ou infirmer l'éventuelle dérive détectée (seulement 5% de chances qu'elle provienne du hasard de l'échantillonnage)
- Si un point sort de la zone de contrôle, il faut arrêter la production et vérifier la machine (seulement 0,2% de chances que la dérive provienne du hasard de l'échantillonnage)
- En plus de ces règles de base, chaque spécialiste et chaque logiciel a ses propres règles additionnelles basées sur le même principe : avec la loi établie, on calcule le degré de bizarrerie à partir duquel la probabilité associée à un événement est trop faible et demande une vérification. Exemple : 11 points d'affilée au dessus de la moyenne, ou 7 points croissants d'affilée *etc.*

Cartes de Contrôle - Exercices

Exercice 1 Un laboratoire produit des capsules contenant un traitement pour un virus. On souhaite contrôler cette production dans le temps. Pour cela, on commence par établir sa loi de production. Ensuite, on vérifiera si cette production est (suffisamment) constante.

1. Le volume contenu dans les capsules produites par la machine est réputée suivre une loi normale de moyenne $\mu = 0,5mL$ et d'écart-type $\sigma = 0,2mL$. Pour déterminer si cette calibration est fiable, on prélève un échantillon de 5 individus par jour, pendant 5 jours. Les résultats synthétiques sont listées sur la Table 8.1. Effectuez alors un test de Kolmogorov (en se contentant des cinq

Jour	Moyenne	Ecart-type estimé	Borne supérieure de classe	Effectifs cumulés
1	0,53	0,19	0,20	1,00
2	0,53	0,22	0,40	8,00
3	0,48	0,14	0,60	17,00
4	0,43	0,23	0,80	22,00
5	0,57	0,19	1,00	25,00
(Résultats quotidiens)			(Histogramme)	

TABLE 8.1 – Phase préliminaire - Constat de la calibration

probabilités observées qu'on pourra déduire des effectifs disponibles dans la Table 8.1) et un test de Cochran pour établir si la machine suit bien la loi spécifiée et que la calibration est stable.

2. A titre informatif, estimez la moyenne et l'écart-type de la production.
3. On lance alors la production de masse et on contrôle pendant un an le bon fonctionnement de la machine, en analysant chaque mois un échantillon de 5 capsules (voir Table 8.2).

Mois	Moyenne	Ecart type estimé
1	0,48	0,17
2	0,43	0,46
3	0,50	0,13
4	0,51	0,11
5	0,40	0,27
6	0,40	0,24
7	0,44	0,21
8	0,56	0,21
9	0,60	0,24
10	0,51	0,15
11	0,68	0,22
12	0,46	0,13

TABLE 8.2 – Phase dynamique - Suivi de la production

Donnez les limites de surveillance et de contrôle (intervalles à 5% et 0.2%) pour la moyenne et l'écart-type. On utilisera pour cela les moyenne et écart type spécifiés.

8 Cartes de contrôle (2h + 2TP)

4. Quelle taille d'échantillon mensuel faudrait-il pour avoir une plage de surveillance correspondant à $\pm 10\%$ de la moyenne spécifiée (ce qui est déjà assez laxiste) ?
5. Contrôlez la production en traçant une carte de contrôle et en interprétant le résultat.

Réponse 1

1. Valeur calculée pour le test de Cochran : 0.27, valeur critique à 5% : 0.5441. On conclut donc à la stabilité de la production au risque $\alpha = 5\%$.

$\mathbb{P}_{\text{thé}}$	\mathbb{P}_{obs}
0,07	0,04
0,31	0,32
0,69	0,68
0,93	0,88
0,99	1,00

Valeur calculée pour le test de Kolmogorov : 0.053, valeur critique à 5% : 0.2716. On conclut donc à la normalité de l'échantillon, au risque α de 5%.

2. $\mu_{est} = 0.5080, \sigma_{est} = 0.1965$.

3.

Seuils	Moyenne	Ecart-type
Contrôle sup	0,7764	0,4297
Surveillance sup	0,6753	0,3338
Référence	0,50	0,20
Surveillance inf	0,3247	0,0696
Contrôle inf	0,2236	0,0301

4. $n = 62$.

5. Voir Figure 8.1

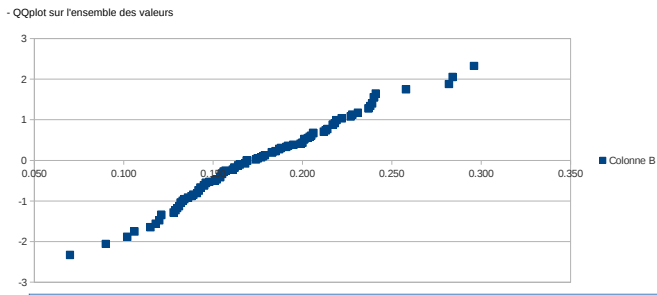
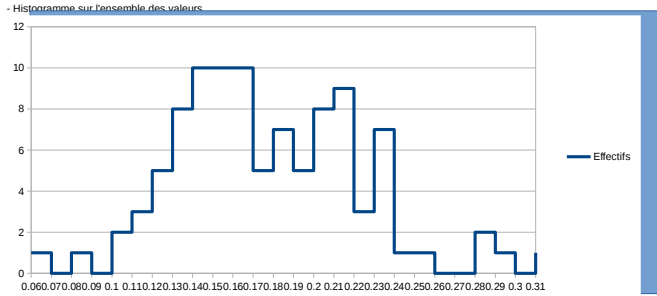
Phase Préliminaire

Date	Echantillon				
01-mai	0.200	0.183	0.228	0.142	0.121
02-mai	0.213	0.154	0.258	0.187	0.136
03-mai	0.128	0.241	0.162	0.239	0.205
04-mai	0.132	0.204	0.222	0.132	0.217
05-mai	0.155	0.237	0.155	0.148	0.217
06-mai	0.157	0.238	0.201	0.192	0.178
07-mai	0.161	0.201	0.183	0.169	0.130
08-mai	0.195	0.174	0.090	0.134	0.219
09-mai	0.203	0.185	0.102	0.284	0.070
10-mai	0.156	0.138	0.151	0.146	0.314
11-mai	0.169	0.212	0.191	0.169	0.188
12-mai	0.240	0.164	0.296	0.217	0.133
13-mai	0.168	0.145	0.145	0.141	0.139
14-mai	0.115	0.219	0.187	0.146	0.201
15-mai	0.165	0.164	0.179	0.219	0.152
16-mai	0.131	0.206	0.177	0.231	0.199
17-mai	0.143	0.142	0.183	0.106	0.154
18-mai	0.121	0.237	0.129	0.143	0.282
19-mai	0.155	0.120	0.214	0.227	0.206
20-mai	0.162	0.118	0.175	0.156	0.240

Estimations

Moyenne	Ecart-type	Variance	Effectif
0.1748	0.043309352	0.0018757	5
0.1896	0.048407644	0.0023433	5
0.195	0.049320381	0.0024325	5
0.1814	0.045571921	0.0020768	5
0.1824	0.041422216	0.0017158	5
0.1932	0.030044966	0.0009027	5
0.1688	0.026480181	0.0007012	5
0.1622	0.050835027	0.0025842	5
0.1688	0.085039403	0.0072317	5
0.181	0.07464583	0.005572	5
0.1858	0.017908099	0.0003207	5
0.21	0.064011718	0.0040975	5
0.1476	0.011696153	0.0001368	5
0.1736	0.042388678	0.0017968	5
0.1758	0.025974988	0.0006747	5
0.1888	0.037619144	0.0014152	5
0.1456	0.027645976	0.0007643	5
0.1824	0.0725865	0.0052688	5
0.1844	0.045191813	0.0020423	5
0.1702	0.044409458	0.0019722	5

Vérification de la normalité



Vérification de la loi de la production

- Test de Cochran

H0 : Les écarts-type sont égaux
 H1 : Au moins un des écarts-type est différent

Valeur calculée : 0.1574669245
 Valeur seuil à 5 % : 0.1921 (on doit ici recourir à la table de Cochran)

Conclusion : On accepte l'hypothèse privilégiée, au risque de 5 %

Ecart-type estimé : 0.0479193072

- Test de Kolmogorov sur l'ensemble des valeurs

Moyenne estimée : 0.17807
 Ecart-type estimé : 0.0479193072

H0 : Les fréquences observées correspondent à une loi normale avec la moyenne et l'écart type ci-dessus
 H1 : Les fréquences observées ne correspondent pas à une loi normale avec la moyenne et l'écart type ci-dessus

Valeur calculée : 0.0750619539
 Valeur seuil à 5 % : 0.1358

Conclusion : On accepte l'hypothèse privilégiée, au risque de 5 %

- Tests de Shapiro pour chaque jour

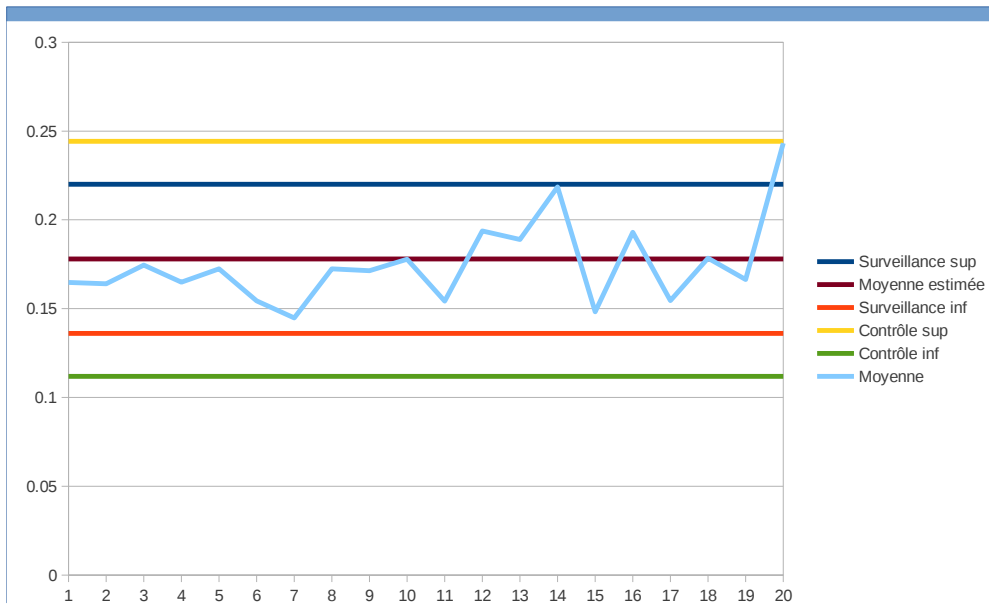
Copier les valeurs puis trier préalablement chaque ligne (menu données, trier, options, de gauche à droite)
 Attention à ne pas étendre la sélection : ligne par ligne
 Puis re-copier ce tableau trié et le trier d'un bloc dans l'ordre décroissant (toujours de gauche à droite)
 Enfin, calculer SW avec les fonctions sommeprod et somme.carres.ecarts

aj	0.6646	0.2413	
Jour	SW	SW 5 %	Conclusion
1	0.9654133894	0.762	Normalité
2	0.9693063266	0.762	Normalité
3	0.9019448781	0.762	Normalité
4	0.7766787004	0.762	Normalité
5	0.800252375	0.762	Normalité
6	0.9765928066	0.762	Normalité
7	0.9825107042	0.762	Normalité
8	0.9633218115	0.762	Normalité
9	0.9594606821	0.762	Normalité
10	0.6394564422	0.762	Anormalité
11	0.895141959	0.762	Normalité
12	0.9789465666	0.762	Normalité
13	0.7485397112	0.762	Anormalité
14	0.9444701166	0.762	Normalité
15	0.8589747352	0.762	Normalité
16	0.9532281107	0.762	Normalité
17	0.9562813267	0.762	Normalité
18	0.8400978269	0.762	Normalité
19	0.8916949923	0.762	Normalité
20	0.9302614368	0.762	Normalité

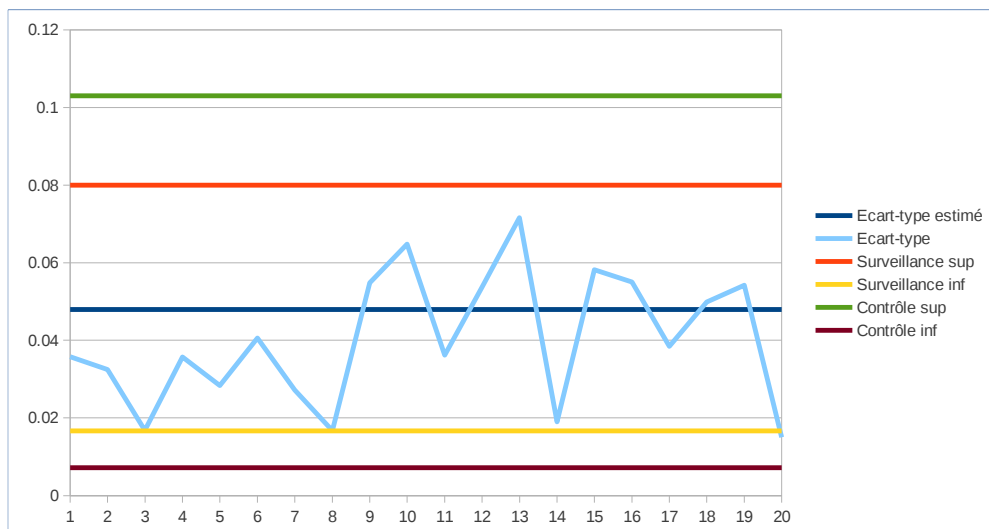
Phase Dynamique

Date	Echantillon					Moyenne	Ecart-type
21-mai	0.162	0.115	0.168	0.163	0.216	0.165	0.036
22-mai	0.162	0.152	0.119	0.206	0.181	0.164	0.033
23-mai	0.194	0.182	0.163	0.182	0.152	0.175	0.017
24-mai	0.204	0.128	0.191	0.127	0.175	0.165	0.036
25-mai	0.189	0.213	0.158	0.142	0.160	0.172	0.028
26-mai	0.198	0.180	0.169	0.125	0.100	0.154	0.041
27-mai	0.118	0.126	0.132	0.176	0.172	0.145	0.027
28-mai	0.164	0.150	0.175	0.195	0.178	0.172	0.017
29-mai	0.122	0.110	0.188	0.242	0.195	0.171	0.055
30-mai	0.250	0.105	0.121	0.233	0.180	0.178	0.065
31-mai	0.149	0.190	0.108	0.133	0.191	0.154	0.036
01-juin	0.223	0.224	0.228	0.193	0.101	0.194	0.054
02-juin	0.242	0.117	0.249	0.105	0.232	0.189	0.072
03-juin	0.214	0.224	0.213	0.247	0.195	0.219	0.019
04-juin	0.122	0.108	0.132	0.251	0.128	0.148	0.058
05-juin	0.238	0.232	0.127	0.139	0.229	0.193	0.055
06-juin	0.117	0.202	0.146	0.121	0.187	0.155	0.038
07-juin	0.104	0.170	0.243	0.185	0.190	0.178	0.050
08-juin	0.244	0.186	0.124	0.171	0.107	0.166	0.054
09-juin	0.270	0.236	0.236	0.238	0.236	0.243	0.015

Carte de contrôle de la moyenne



Carte de contrôle de la variance



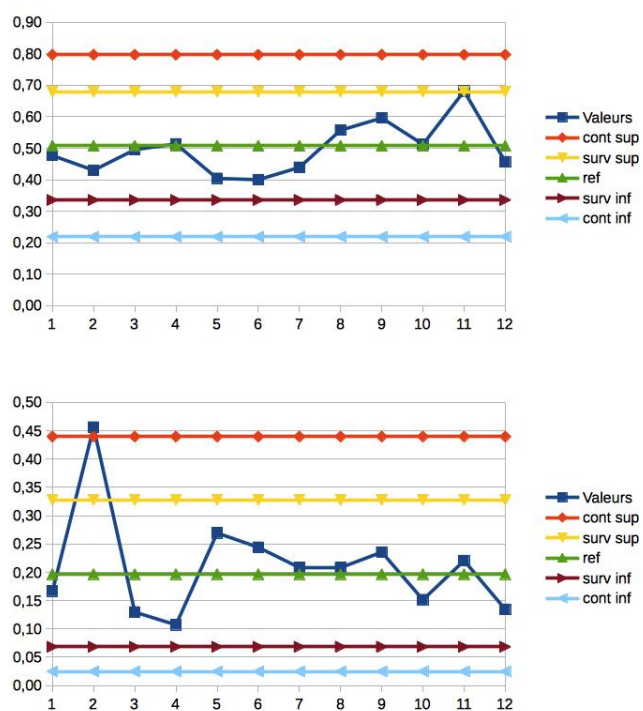


FIGURE 8.1 – Cartes de contrôle de la moyenne et de l'écart-type