

# Tests d'hypothèses via une **A**pproche **E**xpérimentale des **P**robabilités

CQLS : cqls@upmf-grenoble.fr  
http://cqls.upmf-grenoble.fr

## 1 Introduction et généralités

### 1.1 Cadre d'estimation à un échantillon

Dans ce cadre, tout problème pratique doit se ramener à l'étude d'une unique variable d'intérêt notée ici  $Y$  (pouvant aussi être vue comme une future unique donnée). En pratique, nous disposerons d'un jeu de  $n$  données  $\mathbf{y} = (y_1, \dots, y_n)$  (i.e. un vecteur ou "paquet" de  $n$  observations "indépendantes" de  $Y$ ) qui peut par conséquent être vu comme un résultat possible d'un futur jeu de  $n$  données  $\mathbf{Y} = (Y_1, \dots, Y_n)$ . Afin d'explicitier le paramètre d'intérêt intimement lié dans les problématiques du cours à la variable d'intérêt  $Y$ , nous imaginerons disposer d'une infinité de données virtuelles  $y_{[1]}, \dots, y_{[m]}, \dots$  dont la notation en indice entre crochet (i.e. " $[\cdot]$ ") nous rappelle qu'il ne faut pas les confondre avec le jeu des  $n$  données  $y_1, \dots, y_n$  qui seront bien réelles. Rappelons aussi que dans ce cours les tailles  $n$  des données **réelles** et  $m$  des données **virtuelles** ont a priori des ordres de grandeur complètement différents, à savoir  $n$  plutôt raisonnablement grand et  $m$  aussi grand que possible voire infini.

#### 1.1.1 Paramètres proportion et moyenne

La moyenne notée  $\mu_Y$  ou plus simplement  $\mu$  (plutôt appelée espérance de  $Y$  dans l'**A**pproche **M**athématique des **P**robabilités et notée  $\mathbb{E}(Y)$ ) s'exprime via l'**A**pproche **E**xpérimentale des **P**robabilités par :

$$\overline{(y_{[\cdot]})}_m = \frac{1}{m} \sum_{k=1}^m y_{[k]} \simeq \overline{(y_{[\cdot]})}_\infty = \mu_Y = \mathbb{E}(Y).$$

Soulignons toutefois que si les données sont exclusivement à valeurs 0 ou 1, la moyenne devient une proportion (ou probabilité) et sera notée  $p$  plutôt que  $\mu$ . Rappelons qu'une future estimation  $\widehat{\mu}_Y(\mathbf{Y})$  de  $\mu_Y$  est tout simplement  $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$  (notée aussi  $\bar{Y}$ ).

#### 1.1.2 Paramètre variance

La variance notée  $\sigma_Y^2$  ou plus simplement  $\sigma^2$  (conservant la même dénomination dans l'**A**.**M**.**P**. et notée  $\mathbb{V}\text{ar}(Y)$ ) s'exprime via l'**A**.**E**.**P**. par :

$$\left( \overleftrightarrow{(y_{[\cdot]})}_m \right)^2 = \frac{1}{m} \sum_{k=1}^m \left( y_{[k]} - \overline{(y_{[\cdot]})}_m \right)^2 \simeq \left( \overleftrightarrow{(y_{[\cdot]})}_\infty \right)^2 = \sigma_Y^2 = \mathbb{V}\text{ar}(Y) = \sigma(Y)^2.$$

Dans le cadre de grands échantillons (voir plus loin), il est plus qu'intéressant de noter que la variance est aussi une moyenne. En effet, nous pouvons écrire  $\sigma_Y^2 = \mu_{\check{Y}}$  puisque  $\mathbb{V}\text{ar}(Y) = \mathbb{E}((Y - \mu_Y)^2) = \mathbb{E}(\check{Y})$  où  $\check{Y} = (Y - \mu_Y)^2$  est le carré de la variable aléatoire  $Y$  préalablement centrée. Le vecteur des futures données  $((Y_1 - \mu_Y)^2, \dots, (Y_n - \mu_Y)^2)$  étant inaccessible puisque  $\mu_Y$  est inconnu, nous le remplacerons par  $\check{\mathbf{Y}} = ((Y_1 - \bar{Y})^2, \dots, (Y_n - \bar{Y})^2)$ . Ainsi, nous pourrions aussi proposer  $\widehat{\mu}_{\check{Y}}(\check{\mathbf{Y}})$  comme future estimation de  $\sigma_Y^2 = \mu_{\check{Y}}$  (plutôt lorsque la taille  $n$  des données sera suffisamment grande).

## 1.2 Cadre d'estimation à deux échantillons (indépendants)

Il y a dans ce cadre deux variables d'intérêts  $Y^{(1)}$  et  $Y^{(2)}$  ("indépendantes") dont on cherche soit à comparer les moyennes soit les variances à partir de deux échantillons, l'un  $\mathbf{y}^{(1)} = (y_1^{(1)}, \dots, y_{n^{(1)}}^{(1)})$  de taille  $n^{(1)}$  et l'autre  $\mathbf{y}^{(2)} = (y_1^{(2)}, \dots, y_{n^{(2)}}^{(2)})$  de taille  $n^{(2)}$ . Il en découle deux futurs jeux de données  $\mathbf{Y}^{(1)} = (Y_1^{(1)}, \dots, Y_{n^{(1)}}^{(1)})$  et  $\mathbf{Y}^{(2)} = (Y_1^{(2)}, \dots, Y_{n^{(2)}}^{(2)})$ . Pour homogénéiser ce cas avec celui à un seul échantillon, nous noterons  $\mathbf{Y}$ , le vecteur agrégé de toutes les futures données  $\mathbf{Y} = (\mathbf{Y}^{(1)}, \mathbf{Y}^{(2)})$  de taille  $n = n^{(1)} + n^{(2)}$ . De manière analogue au cas d'un seul échantillon, nous imaginons disposer de deux infinités de données virtuelles, l'une  $y_{[1]}^{(1)}, \dots, y_{[m]}^{(1)}, \dots$  relative à  $Y^{(1)}$  et l'autre  $y_{[1]}^{(2)}, \dots, y_{[m]}^{(2)}, \dots$  relative à  $Y^{(2)}$ . Pour  $j = 1$  ou  $j = 2$ , on peut alors exprimer :

- la **moyenne**  $\mu_{Y^{(j)}}$  ou plus simplement  $\mu^{(j)}$  par  $\boxed{\overline{\left(y_{[\cdot]}^{(j)}\right)}_m \simeq \overline{\left(y_{[\cdot]}^{(j)}\right)}_\infty = \mu^{(j)}}.$
- la **variance**  $\sigma_{Y^{(j)}}^2$  ou plus simplement  $\sigma_{(j)}^2$  définie par  $\boxed{\left(\overrightarrow{\left(y_{[\cdot]}^{(j)}\right)}_m\right)^2 \simeq \left(\overrightarrow{\left(y_{[\cdot]}^{(j)}\right)}_\infty\right)^2 = \sigma_{(j)}^2}$

Nous sommes alors en mesure d'introduire les paramètres servant à comparer respectivement les moyennes et les variances.

- **comparaison de moyennes** s'étudiant soit à partir de la **différence de moyennes**  $d_\mu = \mu^{(1)} - \mu^{(2)}$  soit à partir du **rapport de moyennes**  $r_\mu = \mu^{(1)}/\mu^{(2)}$  (si  $\mu^{(2)} \neq 0$ ).
- **comparaison de variances** s'étudiant soit à partir de la **différence de variances**  $d_{\sigma^2} = \sigma_{(1)}^2 - \sigma_{(2)}^2$  soit à partir du **rapport de variances**  $r_{\sigma^2} = \sigma_{(1)}^2/\sigma_{(2)}^2$  (si  $\sigma_{(2)}^2 \neq 0$ ).

Insistons sur le fait que les utilisations d'une différence ou d'un rapport ne sont pas anodines puisqu'elles permettent de traiter des assertions d'intérêt différentes.

## 1.3 Les deux cadres usuels : asymptotique et gaussien

- **Cadre asymptotique ou grand(s) échantillon(s)** : par grand échantillon, on entend dans ce cours une taille de données  $n \geq 30$  pour le cas un seul échantillon et des tailles  $n^{(1)} \geq 30$  et  $n^{(2)} \geq 30$  pour celui de deux échantillons.
- **Cadre gaussien** : si une variable d'intérêt est supposée suivre une loi Normale on dit que l'échantillon associé est gaussien. Ce cadre d'étude n'est a priori intéressant que s'il est possible de vérifier (éventuellement à partir d'un outil statistique) cette hypothèse de Normalité de la variable d'intérêt. Alors qu'il faudrait disposer d'un grand échantillon pour cette vérification, l'usage dans la littérature statistique est d'utiliser ce cadre d'étude même pour des petits échantillons. Les résultats reposent alors sur la validité de l'a priori que la variable d'intérêt suit une loi Normale. Cependant, certains phénomènes étudiés peuvent laisser penser que cette hypothèse sur la (ou les) variable(s) d'intérêt ne doit pas être aberrante.

## 1.4 Comparaison entre A.M.P., A.E.P. et Pratique

Dans le tableau suivant, le **jour J** désigne le jour où les données sont réellement récoltées (*Indic* : voir fin de document pour les différentes notations).

Avant le jour J ( $\theta$ fixé éventuellement à une valeur arbitraire pour l'expérimentation)				
Mathématique	$\mathbf{Y}$	$Y$	$\hat{\theta}(\mathbf{Y})$ ou $\hat{\Theta}$	$t(\mathbf{Y})$ ou $T$
Expérimental	$\mathbf{y}_{[1]}$	$\left\{ \begin{array}{c} y_{[1]} \\ \vdots \\ y_{[n]} \end{array} \right.$	$\hat{\theta}(\mathbf{y}_{[1]})$ ou $\hat{\theta}_{[1]}$	$t(\mathbf{y}_{[1]})$ ou $t_{[1]}$
	$\mathbf{y}_{[2]}$	$\left\{ \begin{array}{c} y_{[n+1]} \\ \vdots \\ y_{[2n]} \end{array} \right.$	$\hat{\theta}(\mathbf{y}_{[2]})$ ou $\hat{\theta}_{[2]}$	$t(\mathbf{y}_{[2]})$ ou $t_{[2]}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	$\mathbf{y}_{[m]}$	$\left\{ \begin{array}{c} y_{[(m-1) \times n + 1]} \\ \vdots \\ y_{[m \times n]} \end{array} \right.$	$\hat{\theta}(\mathbf{y}_{[m]})$ ou $\hat{\theta}_{[m]}$	$t(\mathbf{y}_{[m]})$ ou $t_{[m]}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$
Moyenne =		$\mu := \overline{(y_{[\cdot]})}_{\infty} = \mathbb{E}(Y)$	$\overline{(\hat{\theta}(\mathbf{y}_{[\cdot]}))}_{\infty} = \mathbb{E}(\hat{\theta}(\mathbf{Y}))$	$\overline{(t(\mathbf{y}_{[\cdot]}))}_{\infty} = \mathbb{E}(t(\mathbf{Y}))$
Ecart-Type =		$\sigma := \overline{(y_{[\cdot]})}_{\infty}$ = $\sigma(Y)$ = $\sqrt{\text{Var}(Y)}$	$\sigma_{\hat{\theta}} := \overline{(\hat{\theta}(\mathbf{y}_{[\cdot]}))}_{\infty}$ = $\sigma(\hat{\theta}(\mathbf{Y}))$ = $\sqrt{\text{Var}(\hat{\theta}(\mathbf{Y}))}$	$\overline{(t(\mathbf{y}_{[\cdot]}))}_{\infty} = \sigma(t(\mathbf{Y}))$ = $\sqrt{\text{Var}(t(\mathbf{Y}))}$
Proportion dans $[a, b[$ =		$\overline{(y_{[\cdot]} \in [a, b])}_{\infty}$ = $\mathbf{P}(Y \in [a, b])$	$\overline{(\hat{\theta}(\mathbf{y}_{[\cdot]}) \in [a, b])}_{\infty}$ = $\mathbf{P}(\hat{\theta}(\mathbf{Y}) \in [a, b])$	$\overline{(t(\mathbf{y}_{[\cdot]}) \in [a, b])}_{\infty}$ = $\mathbf{P}(t(\mathbf{Y}) \in [a, b])$
Histogramme à pas "zéro" =		$f_Y$	$f_{\hat{\theta}(\mathbf{Y})}$ ou $f_{\hat{\Theta}}$	$f_{g(\mathbf{Y})}$ ou $f_T$
Surface brique ( $m$ fini) =		$\frac{1}{mn}$	$\frac{1}{m}$	$\frac{1}{m}$
Après le jour J ( $\theta$ est égal à $\theta^{\bullet}$ qui est toujours inconnu)				
Pratique	$\mathbf{y}$	$\left\{ \begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \right.$	$\hat{\theta}(\mathbf{y})$ ou $\hat{\theta}$	$t(\mathbf{y})$ ou $t$

Le **jour J** ( $\theta = \theta^{\bullet}$ ), si on essaye d'associer des temps de conjugaison aux différents concepts, nous pouvons dire :

- le **jeu de données réel**  $\mathbf{y}$  représente le *présent*.
- le **jeu de données aléatoire**  $\mathbf{Y}$  représente le *futur* (on pourra alors aussi l'appeler *futur jeu de données*)
- les **jeux de données virtuels**  $\mathbf{y}_{[j]}$  représentent le *conditionnel* (ils représentent une infinité de jeux de données que l'on aurait pu avoir à la place de  $\mathbf{y}$ )

## 2 Test d'hypothèses

De manière générale, la rédaction standard d'un test d'hypothèses s'écrit toujours de la même façon. Elle est décrite ci-dessous pour un paramètre  $\theta$  qui devra être remplacé par  $p$  pour une proportion,  $\mu$  pour une moyenne,  $\sigma^2$  pour une variance,  $d_\mu$  (resp.  $r_\mu$ ) pour une différence (resp. rapport) de moyennes et enfin  $d_{\sigma^2}$  (resp.  $r_{\sigma^2}$ ) pour une différence (resp. rapport) de variances. La valeur de référence  $\theta_0$  et la loi  $\mathcal{L}_0$  devront être adaptée selon la problématique.

### Rédaction standard d'un test d'hypothèses paramétrique

**Hypothèses de test :**

$$\mathbf{H}_0 : \theta = \theta_0 \text{ contre } \mathbf{H}_1 : \begin{cases} \theta > \theta_0 & (\text{cas (a) : test unilatéral droit}) \\ \theta < \theta_0 & (\text{cas (b) : test unilatéral gauche}) \\ \theta \neq \theta_0 & (\text{cas (c) : test bilatéral}) \end{cases}$$

**Statistique de test sous  $\mathbf{H}_0$  :**

$$\widehat{\delta_{\theta, \theta_0}}(\mathbf{Y}) \rightsquigarrow \mathcal{L}_0$$

où  $\mathcal{L}_0$  est une loi standard à préciser (selon la problématique envisagée).

**Règle de décision :**

$$\text{on accepte } \mathbf{H}_1 \text{ si } \begin{cases} \boxed{p\text{-valeur} < \alpha} \\ \text{ou de manière équivalente} \\ \left\{ \begin{array}{ll} \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) > \delta_{\text{lim}, \alpha}^+ & \text{(a)} \\ \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) < \delta_{\text{lim}, \alpha}^- & \text{(b)} \\ \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) < \delta_{\text{lim}, \alpha/2}^- \text{ ou } \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) > \delta_{\text{lim}, \alpha/2}^+ & \text{(c)} \end{array} \right. \end{cases}$$

où  $\delta_{\text{lim}, \alpha}^- = q_\alpha$  et  $\delta_{\text{lim}, \alpha}^+ = q_{1-\alpha}$  désignent respectivement les quantiles d'ordre  $\alpha$  et  $1 - \alpha$  associés à la loi  $\mathcal{L}_0$  et où la  $p$ -valeur est définie mathématiquement par :

$$p\text{-valeur} = \begin{cases} \mathbb{P}_{\theta=\theta_0} \left( \widehat{\delta_{\theta, \theta_0}}(\mathbf{Y}) > \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) \right) & \text{(a) : } p\text{-valeur droite} \\ \mathbb{P}_{\theta=\theta_0} \left( \widehat{\delta_{\theta, \theta_0}}(\mathbf{Y}) < \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) \right) & \text{(b) : } p\text{-valeur gauche} \\ 2 \times \min \left( \mathbb{P}_{\theta=\theta_0} \left( \widehat{\delta_{\theta, \theta_0}}(\mathbf{Y}) < \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) \right), \mathbb{P}_{\theta=\theta_0} \left( \widehat{\delta_{\theta, \theta_0}}(\mathbf{Y}) > \widehat{\delta_{\theta, \theta_0}}(\mathbf{y}) \right) \right) & \text{(c) : } p\text{-valeur bilatérale} \end{cases}$$

**Conclusion :** Application de la règle de décision au vu des données  $\mathbf{y}$ .

Propriétés :

1. La somme des  $p$ -valeur gauche et  $p$ -valeur droite est égale à 1
2. La  $p$ -valeur bilatérale est égale à deux fois la plus petite des  $p$ -valeurs gauche et droite

**Tableaux récapitulatifs :**

Il sera aussi supposé que les données ont été saisies dans le logiciel R soit sous le nom  $\mathbf{y}$  (pour un unique échantillon) soit sous les noms  $\mathbf{y1}$  et  $\mathbf{y2}$  (pour deux échantillons indépendants).

$\theta$	$\widehat{\theta}(\mathbf{Y})$	$\widehat{\theta}(\mathbf{y})$ en R	$\sigma_{\widehat{\theta}}$	$\widehat{\sigma_{\widehat{\theta}}}(\mathbf{Y})$	$\widehat{\sigma_{\widehat{\theta}}}(\mathbf{y})$ en R
$p$	$\widehat{p}(\mathbf{Y}) = \overline{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$	mean(y)	$\sigma_{\widehat{p}} = \sqrt{\frac{p(1-p)}{n}}$	$\sqrt{\frac{\widehat{p}(\mathbf{Y})(1-\widehat{p}(\mathbf{Y}))}{n}}$	seMean(y)
$\mu$	$\widehat{\mu}(\mathbf{Y}) = \overline{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$	mean(y)	$\sigma_{\widehat{\mu}} = \sqrt{\frac{\sigma^2}{n}}$	$\sqrt{\frac{\widehat{\sigma^2}(\mathbf{Y})}{n}}$	seMean(y)
$\sigma^2$	$\widehat{\sigma^2}(\mathbf{Y}) = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \overline{Y})^2$	var(y)	$\sigma_{\widehat{\sigma^2}} = \sqrt{\frac{\sigma_{\widehat{Y}}^2}{n}}$	$\sqrt{\frac{\widehat{\sigma_{\widehat{Y}}^2}(\mathbf{Y})}{n}}$	seVar(y)
$d_{\mu} = \mu^{(1)} - \mu^{(2)}$	$\widehat{d}_{\mu}(\mathbf{Y}) = \widehat{\mu^{(1)}}(\mathbf{Y}^{(1)}) - \widehat{\mu^{(2)}}(\mathbf{Y}^{(2)})$	mean(y1)-mean(y2)	$\sigma_{\widehat{d}_{\mu}} = \sqrt{\frac{\sigma_{\widehat{Y}^{(1)}}^2}{n^{(1)}} + \frac{\sigma_{\widehat{Y}^{(2)}}^2}{n^{(2)}}}$	$\sqrt{\frac{\widehat{\sigma_{\widehat{Y}^{(1)}}^2}(\mathbf{Y}^{(1)})}{n^{(1)}} + \frac{\widehat{\sigma_{\widehat{Y}^{(2)}}^2}(\mathbf{Y}^{(2)})}{n^{(2)}}}$	seDMean(y1,y2)
$d_{\sigma^2} = \sigma_{\widehat{Y}^{(1)}}^2 - \sigma_{\widehat{Y}^{(2)}}^2$	$\widehat{d}_{\sigma^2}(\mathbf{Y}) = \widehat{\sigma_{\widehat{Y}^{(1)}}^2}(\mathbf{Y}^{(1)}) - \widehat{\sigma_{\widehat{Y}^{(2)}}^2}(\mathbf{Y}^{(2)})$	var(y1)-var(y2)	$\sigma_{\widehat{r_{\sigma^2}}} = \sqrt{\frac{\sigma_{\widehat{Y}^{(1)}}^2}{n^{(1)}} + \frac{\sigma_{\widehat{Y}^{(2)}}^2}{n^{(2)}}}$	$\sqrt{\frac{\widehat{\sigma_{\widehat{Y}^{(1)}}^2}(\mathbf{Y}^{(1)})}{n^{(1)}} + \frac{\widehat{\sigma_{\widehat{Y}^{(2)}}^2}(\mathbf{Y}^{(2)})}{n^{(2)}}}$	seDVar(y1,y2)
$r_{\mu} = \frac{\mu^{(1)}}{\mu^{(2)}}$	$\widehat{r}_{\mu}(\mathbf{Y}) = \frac{\widehat{\mu^{(1)}}(\mathbf{Y}^{(1)})}{\widehat{\mu^{(2)}}(\mathbf{Y}^{(2)})}$	mean(y1)/mean(y2)	$\sigma_{\widehat{r_{\mu}}} = \frac{1}{\mu^{(2)}} \sqrt{\frac{\sigma_{\widehat{Y}^{(1)}}^2}{n^{(1)}} + r_{\mu}^2 \times \frac{\sigma_{\widehat{Y}^{(2)}}^2}{n^{(2)}}}$	$\frac{1}{\widehat{\mu^{(2)}}(\mathbf{Y}^{(2)})} \sqrt{\frac{\widehat{\sigma_{\widehat{Y}^{(1)}}^2}(\mathbf{Y}^{(1)})}{n^{(1)}} + \frac{\widehat{\sigma_{\widehat{Y}^{(2)}}^2}(\mathbf{Y}^{(2)})}{n^{(2)}}}$	seRMean(y1,y2)
$r_{\sigma^2} = \frac{\sigma_{\widehat{Y}^{(1)}}^2}{\sigma_{\widehat{Y}^{(2)}}^2}$	$\widehat{r}_{\sigma^2}(\mathbf{Y}) = \frac{\widehat{\sigma_{\widehat{Y}^{(1)}}^2}(\mathbf{Y}^{(1)})}{\widehat{\sigma_{\widehat{Y}^{(2)}}^2}(\mathbf{Y}^{(2)})}$	var(y1)/var(y2)	$\sigma_{\widehat{r_{\sigma^2}}} = \frac{1}{\sigma_{\widehat{Y}^{(2)}}^2} \sqrt{\frac{\sigma_{\widehat{Y}^{(1)}}^2}{n^{(1)}} + r_{\sigma^2}^2 \times \frac{\sigma_{\widehat{Y}^{(2)}}^2}{n^{(2)}}}$	$\frac{1}{\widehat{\sigma_{\widehat{Y}^{(2)}}^2}(\mathbf{Y}^{(2)})} \sqrt{\frac{\widehat{\sigma_{\widehat{Y}^{(1)}}^2}(\mathbf{Y}^{(1)})}{n^{(1)}} + \frac{\widehat{\sigma_{\widehat{Y}^{(2)}}^2}(\mathbf{Y}^{(2)})}{n^{(2)}}}$	seRVar(y1,y2)

Cadre Asymptotique			Cadre Gaussien	
$\theta$	$\theta_0$	$\sigma_{\widehat{\theta}}$ sous $\mathbf{H}_0$	$\delta_{\theta, \theta_0} = (\theta - \theta_0) / \sigma_{\widehat{\theta}}$	$\widehat{\delta}_{\theta, \theta_0}(\mathbf{Y})$ et sa loi sous $\mathbf{H}_0$
$p$	$p_0$	$\sqrt{\frac{p_0(1-p_0)}{n}}$	$\widehat{p}(\mathbf{Y}) - p_0 \underset{\sim}{\approx} \mathcal{N}(0, 1)$	
$\mu$	$\mu_0$	$\sigma_{\widehat{\mu}}$	$\widehat{\mu}(\mathbf{Y}) - \mu_0 \underset{\sim}{\approx} \mathcal{N}(0, 1)$	$\widehat{\mu}(\mathbf{Y}) - \mu_0 \underset{\sim}{\approx} \mathcal{St}(n-1)$
$\sigma^2$	$\sigma_0^2$	$\sigma_{\widehat{\sigma^2}}^2 = \frac{\sigma^2 - \sigma_0^2}{\sigma_{\widehat{\sigma^2}}^2}$	$\widehat{\sigma^2}(\mathbf{Y}) - \sigma_0^2 \underset{\sim}{\approx} \mathcal{N}(0, 1)$	$\widehat{\sigma^2}(\mathbf{Y}) - \sigma_0^2 \underset{\sim}{\approx} \chi^2(n-1)$
$d_{\mu} = \mu^{(1)} - \mu^{(2)}$	$d_0$	$\sigma_{\widehat{d}_{\mu}}$	$\widehat{d}_{\mu}(\mathbf{Y}) - d_0 \underset{\sim}{\approx} \mathcal{N}(0, 1)$	$\widehat{d}_{\mu}(\mathbf{Y}) - d_0 \underset{\sim}{\approx} \mathcal{St}(n^{(1)} + n^{(2)} - 2)$
$d_{\sigma^2} = \sigma_{\widehat{Y}^{(1)}}^2 - \sigma_{\widehat{Y}^{(2)}}^2$	$d_0$	$\sigma_{\widehat{d_{\sigma^2}}}$	$\widehat{d_{\sigma^2}}(\mathbf{Y}) - d_0 \underset{\sim}{\approx} \mathcal{N}(0, 1)$	
$r_{\mu} = \frac{\mu^{(1)}}{\mu^{(2)}}$	$r_0$	$\sigma_{\widehat{r_{\mu}}}$	$\widehat{r_{\mu}}(\mathbf{Y}) - r_0 \underset{\sim}{\approx} \mathcal{N}(0, 1)$	
$r_{\sigma^2} = \frac{\sigma_{\widehat{Y}^{(1)}}^2}{\sigma_{\widehat{Y}^{(2)}}^2}$	$r_0$	$\sigma_{\widehat{r_{\sigma^2}}}$	$\widehat{r_{\sigma^2}}(\mathbf{Y}) - r_0 \underset{\sim}{\approx} \mathcal{N}(0, 1)$	$\widehat{r_{\sigma^2}}(\mathbf{Y}) - r_0 \underset{\sim}{\approx} \mathcal{F}(n^{(1)} - 1, n^{(2)} - 1)$

### 3 Intervalle de confiance

#### 3.1 Généralités

Le concept d'intervalle de confiance d'un paramètre quelconque  $\theta$  consiste à proposer un encadrement (ou une "fourchette") représenté par un intervalle de variables aléatoires  $[\tilde{\theta}_{\inf}(\mathbf{Y}), \tilde{\theta}_{\sup}(\mathbf{Y})]$  de sorte que le paramètre d'intérêt  $\theta$  inconnu ait un niveau de confiance  $1 - \alpha$  (plutôt élevé si  $\alpha$  raisonnablement petit) d'être à l'intérieur de cet intervalle. Mathématiquement cela s'exprime par :

$$\mathbb{P}\left(\theta \in [\tilde{\theta}_{\inf}(\mathbf{Y}), \tilde{\theta}_{\sup}(\mathbf{Y})]\right) = 1 - \alpha$$

Par l'approche expérimentale, si nous pouvions imaginer répéter autant de fois que possible la conception d'intervalles de confiance  $[\tilde{\theta}_{\inf}(\mathbf{y}_{[1]}), \tilde{\theta}_{\sup}(\mathbf{y}_{[1]})], \dots, [\tilde{\theta}_{\inf}(\mathbf{y}_{[m]}), \tilde{\theta}_{\sup}(\mathbf{y}_{[m]})], \dots$  obtenus respectivement sur une infinité de jeux de données virtuels  $\mathbf{y}_{[1]}, \dots, \mathbf{y}_{[m]}, \dots$ , nous constaterions alors qu'il n'y en aurait qu'une proportion  $1 - \alpha$  qui contiendraient le paramètre d'intérêt  $\theta$  inconnu. La construction de ces intervalles dépend en général de la caractérisation du comportement aléatoire de la mesure d'écart standardisée  $\delta_{\hat{\theta}, \theta}(\mathbf{Y})$  proposée dans toutes les problématiques dans le tableau suivant :

$\theta$	Cadre Asymptotique $\delta_{\hat{\theta}, \theta}(\mathbf{Y})$	Cadre Gaussien $\delta_{\hat{\theta}, \theta}(\mathbf{Y})$
$p$	$\delta_{\hat{p}, p}(\mathbf{Y}) = \frac{\hat{p}(\mathbf{Y}) - p}{\widehat{\sigma}_{\hat{p}}(\mathbf{Y})} \overset{approx.}{\rightsquigarrow} \mathcal{N}(0, 1)$	
$\mu$	$\delta_{\hat{\mu}, \mu}(\mathbf{Y}) = \frac{\hat{\mu}(\mathbf{Y}) - \mu}{\widehat{\sigma}_{\hat{\mu}}(\mathbf{Y})} \overset{approx.}{\rightsquigarrow} \mathcal{N}(0, 1)$	$\delta_{\hat{\mu}, \mu}(\mathbf{Y}) = \frac{\hat{\mu}(\mathbf{Y}) - \mu}{\widehat{\sigma}_{\hat{\mu}}(\mathbf{Y})} \rightsquigarrow St(n - 1)$
$\sigma^2$	$\delta_{\widehat{\sigma^2}, \sigma^2}(\mathbf{Y}) = \frac{\widehat{\sigma^2}(\mathbf{Y}) - \sigma^2}{\widehat{\sigma_{\widehat{\sigma^2}}}(\mathbf{Y})} \overset{approx.}{\rightsquigarrow} \mathcal{N}(0, 1)$	$\delta_{\widehat{\sigma^2}, \sigma^2}(\mathbf{Y}) = (n - 1) \frac{\widehat{\sigma^2}(\mathbf{Y})}{\sigma^2} \rightsquigarrow \chi^2(n - 1)$
$d_{\mu} = \mu^{(1)} - \mu^{(2)}$	$\delta_{\widehat{d_{\mu}}, d_{\mu}}(\mathbf{Y}) = \frac{\widehat{d_{\mu}}(\mathbf{Y}) - d_{\mu}}{\widehat{\sigma_{\widehat{d_{\mu}}}}(\mathbf{Y})} \overset{approx.}{\rightsquigarrow} \mathcal{N}(0, 1)$	$\delta_{\widehat{d_{\mu}}, d_{\mu}}(\mathbf{Y}) = \frac{\widehat{d_{\mu}}(\mathbf{Y}) - d_{\mu}}{\widehat{\sigma_{\widehat{d_{\mu}}}}(\mathbf{Y})} \rightsquigarrow St(n^{(1)} + n^{(2)} - 2)$
$d_{\sigma^2} = \sigma_{(1)}^2 - \sigma_{(2)}^2$	$\delta_{\widehat{d_{\sigma^2}}, d_{\sigma^2}}(\mathbf{Y}) = \frac{\widehat{d_{\sigma^2}}(\mathbf{Y}) - d_{\sigma^2}}{\widehat{\sigma_{\widehat{d_{\sigma^2}}}}(\mathbf{Y})} \overset{approx.}{\rightsquigarrow} \mathcal{N}(0, 1)$	
$r_{\mu} = \frac{\mu^{(1)}}{\mu^{(2)}}$	$\delta_{\widehat{r_{\mu}}, r_{\mu}}(\mathbf{Y}) = \frac{\widehat{r_{\mu}}(\mathbf{Y}) - r_{\mu}}{\widehat{\sigma_{\widehat{r_{\mu}}}}(\mathbf{Y})} \overset{approx.}{\rightsquigarrow} \mathcal{N}(0, 1)$	
$r_{\sigma^2} = \frac{\sigma_{(1)}^2}{\sigma_{(2)}^2}$	$\delta_{\widehat{r_{\sigma^2}}, r_{\sigma^2}}(\mathbf{Y}) = \frac{\widehat{r_{\sigma^2}}(\mathbf{Y}) - r_{\sigma^2}}{\widehat{\sigma_{\widehat{r_{\sigma^2}}}}(\mathbf{Y})} \overset{approx.}{\rightsquigarrow} \mathcal{N}(0, 1)$	$\delta_{\widehat{r_{\sigma^2}}, r_{\sigma^2}}(\mathbf{Y}) = \frac{\widehat{r_{\sigma^2}}(\mathbf{Y})}{r_{\sigma^2}} \rightsquigarrow \mathcal{F}(n^{(1)} - 1, n^{(2)} - 1)$

#### 3.2 Cadre asymptotique

Tous les intervalles de confiance relatifs à toutes les problématiques du cours s'obtiennent **dans le cadre de grands échantillons** par la même méthode. Après substitution du paramètre d'intérêt de votre problématique (au choix parmi  $p, \mu, \sigma^2, d_{\mu}, d_{\sigma^2}, r_{\mu}$  et  $r_{\sigma^2}$ ) notée ici de manière générale  $\theta$ , nous allons naturellement utiliser la caractérisation du comportement aléatoire de l'écart entre  $\hat{\theta}(\mathbf{Y})$  et  $\theta$  exprimée via la mesure d'écart standardisée  $\delta_{\hat{\theta}, \theta}(\mathbf{Y})$  suivant approximativement une loi Normale  $\mathcal{N}(0, 1)$ . Très facilement, nous pouvons affirmer que :

$$1 - \alpha \simeq \mathbb{P}\left(\left|\delta_{\hat{\theta}, \theta}(\mathbf{Y})\right| < \delta_{lim, \frac{\alpha}{2}}^+\right) = \mathbb{P}\left(\underbrace{\hat{\theta}(\mathbf{Y}) - \delta_{lim, \frac{\alpha}{2}}^+}_{\tilde{\theta}_{\inf}(\mathbf{Y})} \times \widehat{\sigma_{\hat{\theta}}}(\mathbf{Y}) < \theta < \underbrace{\hat{\theta}(\mathbf{Y}) + \delta_{lim, \frac{\alpha}{2}}^+}_{\tilde{\theta}_{\sup}(\mathbf{Y})} \times \widehat{\sigma_{\hat{\theta}}}(\mathbf{Y})\right)$$

où  $\delta_{lim, \frac{\alpha}{2}}^+$  est le quantile d'ordre  $1 - \frac{\alpha}{2}$  de la loi  $\mathcal{N}(0, 1)$

#### 3.3 Cadre gaussien

Pour construire un intervalle de confiance dans un cadre gaussien du paramètre  $\theta$  (au choix  $\mu, \sigma^2, d_{\mu}$  ou  $r_{\sigma^2}$ ), nous allons naturellement utiliser la caractérisation du comportement aléatoire de

l'écart entre  $\hat{\theta}(\mathbf{Y})$  et  $\theta$  exprimée via la mesure d'écart standardisée  $\delta_{\hat{\theta},\theta}(\mathbf{Y})$ . Il s'agit alors de trouver  $\tilde{\theta}_{\inf}(\mathbf{Y})$  et  $\tilde{\theta}_{\sup}(\mathbf{Y})$  tels que

$$1-\alpha = \mathbb{P}\left(\tilde{\theta}_{\inf}(\mathbf{Y}) < \theta < \tilde{\theta}_{\sup}(\mathbf{Y})\right) \quad \text{en utilisant le fait que} \quad 1-\alpha = \mathbb{P}\left(q_{\frac{\alpha}{2}} < \delta_{\hat{\theta},\theta}(\mathbf{Y}) < q_{1-\frac{\alpha}{2}}\right)$$

où  $q_{1-\frac{\alpha}{2}}$  est le quantile d'ordre  $1 - \frac{\alpha}{2}$  de la loi de la mesure d'écart standardisée  $\delta_{\hat{\theta},\theta}(\mathbf{Y})$ . L'exercice est plus difficile que dans le cadre asymptotique d'une part parce que la mesure d'écart standardisée  $\delta_{\hat{\theta},\theta}(\mathbf{Y})$  ( $\theta$  étant au choix  $\mu$ ,  $\sigma^2$ ,  $d_\mu$  ou  $r_{\sigma^2}$ ) ne se décline pas toujours sur le même schéma de construction et d'autre part parce que la loi de  $\delta_{\hat{\theta},\theta}(\mathbf{Y})$  n'est plus une loi Normale standard. Sans trop nous attarder, voici les différents intervalles de confiance pour les différents choix de  $\theta$  :

- $\theta = \mu$  :

$$\tilde{\mu}_{\inf}(\mathbf{Y}) = \hat{\mu}(\mathbf{Y}) - q_{1-\frac{\alpha}{2}} \times \sqrt{\frac{\widehat{\sigma^2}(\mathbf{Y})}{n}} \quad \text{et} \quad \tilde{\mu}_{\sup}(\mathbf{Y}) = \hat{\mu}(\mathbf{Y}) + q_{1-\frac{\alpha}{2}} \times \sqrt{\frac{\widehat{\sigma^2}(\mathbf{Y})}{n}}$$

où  $q_{1-\frac{\alpha}{2}}$  est le quantile d'ordre  $1 - \frac{\alpha}{2}$  de la loi  $\mathcal{St}(n-1)$ .

- $\theta = \sigma^2$  :

$$\tilde{\sigma^2}_{\inf}(\mathbf{Y}) = \widehat{\sigma^2}(\mathbf{Y}) \times \frac{n-1}{q_{1-\frac{\alpha}{2}}} \quad \text{et} \quad \tilde{\sigma^2}_{\sup}(\mathbf{Y}) = \widehat{\sigma^2}(\mathbf{Y}) \times \frac{n-1}{q_{\frac{\alpha}{2}}}$$

où  $q_{1-\frac{\alpha}{2}}$  (resp.  $q_{\frac{\alpha}{2}}$ ) est le quantile d'ordre  $1 - \frac{\alpha}{2}$  (resp.  $\frac{\alpha}{2}$ ) de la loi  $\chi^2(n-1)$ .

- $\theta = d_\mu$  :

$$\tilde{d}_{\mu\inf}(\mathbf{Y}) = \widehat{d}_\mu(\mathbf{Y}) - q_{1-\frac{\alpha}{2}} \times \sqrt{\widehat{\sigma^2}(\mathbf{Y}) \left( \frac{1}{n^{(1)}} + \frac{1}{n^{(2)}} \right)}$$

et

$$\tilde{d}_{\mu\sup}(\mathbf{Y}) = \widehat{d}_\mu(\mathbf{Y}) + q_{1-\frac{\alpha}{2}} \times \sqrt{\widehat{\sigma^2}(\mathbf{Y}) \left( \frac{1}{n^{(1)}} + \frac{1}{n^{(2)}} \right)}$$

où  $q_{1-\frac{\alpha}{2}}$  est le quantile d'ordre  $1 - \frac{\alpha}{2}$  de la loi  $\mathcal{St}(n^{(1)} + n^{(2)} - 2)$ . La quantité  $\widehat{\sigma^2}(\mathbf{Y})$  est définie comme dans le tableau récapitulatif des tests d'hypothèses (partie cadre gaussien).

- $\theta = r_{\sigma^2}$  :

$$\tilde{r}_{\sigma^2\inf}(\mathbf{Y}) = \widehat{r}_{\sigma^2}(\mathbf{Y}) \times \frac{1}{q_{1-\frac{\alpha}{2}}} \quad \text{et} \quad \tilde{r}_{\sigma^2\sup}(\mathbf{Y}) = \widehat{r}_{\sigma^2}(\mathbf{Y}) \times \frac{1}{q_{\frac{\alpha}{2}}}$$

où  $q_{1-\frac{\alpha}{2}}$  (resp.  $q_{\frac{\alpha}{2}}$ ) est le quantile d'ordre  $1 - \frac{\alpha}{2}$  (resp.  $\frac{\alpha}{2}$ ) de la loi  $\mathcal{F}(n^{(1)} - 1, n^{(2)} - 1)$ .

## 4 Langage mathématique et Systèmes de notation

- Dans ce cours, deux systèmes de notation sont utilisés pour décrire des expressions mathématiques dédiées à la statistique. Le premier, appelé *Norme CQLS* (ou *Norme CQLS Standard*) consiste en un système de notation riche (et peut-être un peu lourde) dont le principal avantage est qu'il est taillé sur mesure pour être traduisible dans le langage littéral. Le deuxième système, appelé *Norme SSE* (ou *Norme CQLS Simplifié*), a pour vocation à être Simple, Synthétique et Explicite (ou du moins le plus possible). Il demande cependant dans son utilisation un meilleur niveau d'expertise essentiellement dû au fait que sa traduction dans le langage littéral est moins explicite que celle pour la *Norme CQLS*.
- Notre conseil est de commencer par l'utilisation de la *Norme CQLS* pour, au fur et à mesure du cours, passer à la *Norme SSE*.
- Conventions communes aux deux Normes CQLS et SSE :

1. Majuscule versus Minuscule : une variable aléatoire (ou susceptible de l'être) est notée en majuscule quand une variable dont on sait qu'elle est déterministe (i.e. non aléatoire) est noté en minuscule.
2. Le Chapeau au dessus d'une quantité (par exemple,  $\hat{\theta}$ ) désigne généralement un remplaçant appelé plus communément estimation dans le cas où la quantité est un paramètre (ici  $\theta$ ).
3. Un vecteur est noté en **caractères gras**.

*Remarque* : une expression écrite sur un document imprimé en caractères gras (ex : "**expression en gras**") est substituée sur un tableau ou sur une feuille papier par sa version soulignée (ex : "expression en gras").

4. “Delta” ( $\delta$  en minuscule et  $\Delta$  en majuscule) est utilisé pour désigner un écart le plus souvent additif (i.e. une soustraction) mais parfois multiplicatif (i.e. une division).
- La *Norme CQLS* a été introduite pour décrire le plus précisément possible l’Approche Expérimentale des Probabilités (A.E.P.). L’A.E.P. s’articulant sur une distinction des différents jeux de données, la *Norme CQLS* repose sur la convention suivante : Toute statistique (i.e. v.a. dépendant d’un jeu de données) s’écrit comme une fonction du jeu de données.
  - Il n’y a pas vraiment de convention propre à la *Norme SSE*. Son objectif est cependant de ne pas respecter la convention spécifique (ci-dessus) à la *Norme CQLS* dans le but de rendre plus synthétique les notations mathématiques.
  - Le tableau ci-dessous exprime plus clairement la spécificité des *Normes CQLS* et *SSE* en proposant les principales expressions utilisées en statistique dans les 2 normes.

Statistique (v.a. fonction de l’échantillon)	Aléatoire ou futur		Réalisé ou présent		Réalizable ou conditionnel	
	<i>CQLS</i>	<i>SSE</i>	<i>CQLS</i>	<i>SSE</i>	<i>CQLS</i>	<i>SSE</i>
Estimation de $\theta$	$\hat{\theta}^\bullet(Y)$	$\hat{\Theta}^\bullet$	$\hat{\theta}^\bullet(y)$	$\hat{\theta}^\bullet$	$\hat{\theta}^\bullet(y_{[k]})$	$\hat{\theta}_{[k]}^\bullet$
Estimation de $p^\bullet$	$\hat{p}^\bullet(Y)$	$\hat{P}^\bullet$	$\hat{p}^\bullet(y)$	$\hat{p}^\bullet$	$\hat{p}^\bullet(y_{[k]})$	$\hat{p}_{[k]}^\bullet$
Estimation de $\mu^\bullet$	$\hat{\mu}^\bullet(Y)$	$\hat{M}^\bullet$	$\hat{\mu}^\bullet(y)$	$\hat{\mu}^\bullet$	$\hat{\mu}^\bullet(y_{[k]})$	$\hat{\mu}_{[k]}^\bullet$
Estimation de $\sigma^2_\bullet$	$\hat{\sigma}^2_\bullet(Y)$	$\hat{\Sigma}^2_\bullet$	$\hat{\sigma}^2_\bullet(y)$	$\hat{\sigma}^2_\bullet$	$\hat{\sigma}^2_\bullet(y_{[k]})$	$\hat{\sigma}^2_{\bullet,[k]}$
Erreur standard de $\hat{\theta}^\bullet$	$\widehat{\sigma_{\theta^\bullet}}(Y)$	$\hat{\Sigma}_{\theta^\bullet}$	$\widehat{\sigma_{\theta^\bullet}}(y)$	$\hat{\sigma}_{\theta^\bullet}$	$\widehat{\sigma_{\theta^\bullet}}(y_{[k]})$	$\hat{\sigma}_{\theta^\bullet,[k]}$
Ecart entre $\hat{\theta}^\bullet$ et $\theta^\bullet$	$\delta_{\hat{\theta}^\bullet,\theta^\bullet}(Y)$	$\Delta_{\theta^\bullet}$	$\delta_{\hat{\theta}^\bullet,\theta^\bullet}(y)$	$\delta_{\theta^\bullet}$	$\delta_{\hat{\theta}^\bullet,\theta^\bullet}(y_{[k]})$	$\delta_{\theta^\bullet,[k]}$
Estimation de $\delta_{\theta^\bullet,\theta_0}$ (ou $\delta_{\theta_0}$ )	$\widehat{\delta_{\theta^\bullet,\theta_0}}(Y)$	$\hat{\Delta}_{\theta_0}$	$\widehat{\delta_{\theta^\bullet,\theta_0}}(y)$	$\hat{\delta}_{\theta_0}$	$\widehat{\delta_{\theta^\bullet,\theta_0}}(y_{[k]})$	$\hat{\delta}_{\theta_0,[k]}$

- Le tableau ci-dessous illustre comment convertir une notation en sa définition littérale ou mathématique pour des concepts de base de la statistique. La conversion dans le langage R y est aussi proposée permettant à l’utilisateur de savoir comment obtenir ces quantités en Pratique :

Notation	Définition littérale	Définition mathématique
$\mathbf{y}$ ou $(y.)_n$	Vecteur des réels $y_1, \dots, y_n$ ( $y_i$ est la $i^{\text{ème}}$ composante de $\mathbf{y}$ )	$(y_1, \dots, y_n)$ (en R : $\mathbf{y} \leftarrow \mathbf{c}(y_1, \dots, y_n)$ )
$\#(\mathbf{y})$	Nombre de composantes de $\mathbf{y}$	$n \stackrel{\text{R}}{=} \text{length}(\mathbf{y})$
$\bar{y}$ ou $\overline{(y.)_n}$  $\frac{\overline{y} = a}{\text{ou } (y. = a)_n}$	Moyenne (empirique) de $\mathbf{y}$  Proportion des $y_1, \dots, y_n$ égaux à $a$	$\frac{1}{n} \sum_{i=1}^n y_i \stackrel{\text{R}}{=} \text{mean}(\mathbf{y})$  $\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{y_i=a} \stackrel{\text{R}}{=} \text{mean}(\mathbf{y}==a)$
$\frac{\overline{a \leq y \leq b}}{\text{ou } (a \leq y. \leq b)_n}$	Proportion des $y_1, \dots, y_n$ dans $[a, b]$ avec $(a \leq b)$	$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[a,b]}(y_i) \stackrel{\text{R}}{=} \text{mean}(a \leq \mathbf{y} \ \& \ \mathbf{y} \leq b)$
$\overleftarrow{y}$ ou $\overleftarrow{(y.)_n}$  $(\overleftarrow{y})^2$ ou $(\overleftarrow{(y.)_n})^2$	Ecart-type (empirique) de $\mathbf{y}$  Variance (empirique) de $\mathbf{y}$	$\sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} \stackrel{\text{R}}{=} \text{sd}(\mathbf{y})$  $\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \stackrel{\text{R}}{=} \text{var}(\mathbf{y})$
$q_\alpha(\mathbf{y})$ ou $q_\alpha((y.)_n)$	Quantile d’ordre $\alpha$ de $\mathbf{y}$ ( $0 < \alpha < 1$ )	$y_{[\alpha n]+1}$ ( $n$ impair) et $\frac{y_{[\alpha n]+1} + y_{[\alpha n]+2}}{2}$ ( $n$ pair) $\stackrel{\text{R}}{=} \text{quantile}(\mathbf{y}, \alpha)$

## 5 Quelques instructions R

**Instructions de base par l’exemple :** des exemples (commentés) valent (peut-être) mieux que de longs discours !

```

1 | > c(-1,1)                # Création du vecteur (-1,1)
2 | [1] -1  1
3 | > 4+2*c(-1,0,1)         # Transformation 4+2*x appliqué pour chaque composante de y
4 | [1]  2  4  6
5 | > y<-c(1,3,2,4,7,6)
```



```

6 | > y
7 | [1] 1 3 2 4 7 6
8 | > 4+2*y
9 | [1] 6 10 8 12 18 16
10 | > mean(y) # Moyenne de y
11 | [1] 3.833333
12 | > sd(y) # Ecart-type de y
13 | [1] 2.316607
14 | > yc <- y-mean(y) # yc correspond au vecteur y centré
15 | > yc
16 | [1] -2.833333 -0.833333 -1.833333 0.166667 3.166667 2.166667
17 | > mean(yc) # Moyenne nulle
18 | [1] -1.480297e-16
19 | > sd(yc) # Idem que l'écart-type de y
20 | [1] 2.316607
21 | > ycr <- (y-mean(y))/sd(y) # ycr correspond au vecteur y centré et réduit
22 | > mean(ycr) # Moyenne nulle
23 | [1] -7.40239e-17
24 | > sd(ycr) # Ecart-type à 1
25 | [1] 1
26 | > var(y) # Variance de y
27 | [1] 5.366667
28 | > sqrt(var(y)) # Ecart-type = racine carrée de variance
29 | [1] 2.316607
30 | > sd(y)^2 # Variance = carré de l'écart-type
31 | [1] 5.366667

```

**Quantiles et fonctions de répartition avec R** : Soit  $p$  un réel appartenant à  $]0, 1[$ , on définit le quantile d'ordre  $p$  associée à une loi de probabilité le réel qui via l'approche expérimentale peut être vu comme le réel qui sépare l'infinité des observations (associée à la loi de probabilité) en deux, une proportion  $p$  à gauche et une proportion  $1 - p$  à droite. On définit également la fonction de répartition en un réel  $q$ , la proportion parmi l'infinité des observations qui se situent avant  $q$ . Ces deux notions sont illustrées dans la figure 1.

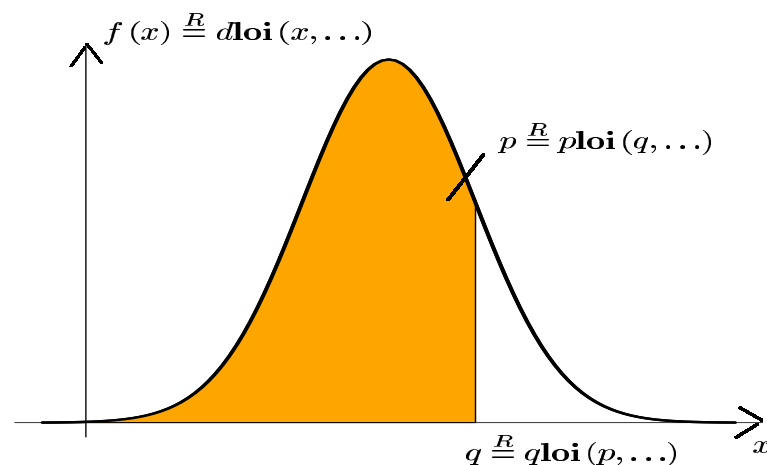


FIGURE 1 – Si  $X \rightsquigarrow \text{loi}(\dots)$  (v.a. continue), alors  $f(x) \stackrel{R}{=} d\text{loi}(x, \dots)$  représente sa densité de probabilité,  $p = F(q) = P(X \leq q) \stackrel{R}{=} p\text{loi}(q, \dots)$  sa fonction de répartition et  $q = F^{-1}(p) \stackrel{R}{=} q\text{loi}(p, \dots)$  son quantile d'ordre  $p$ .

Le tableau suivant résume les différentes lois de probabilités considérées dans ce cours de deuxième année ainsi que les instructions R permettant d'évaluer les quantiles et fonctions de répartition associés à ces lois de probabilités.

lois de probabilités	loi R	quantile d'ordre p	fonction de répartition en q
Normale $\mathcal{N}(\mu, \sigma)$	norm	qnorm(p, $\mu, \sigma$ )	pnorm(q, $\mu, \sigma$ )
Normale $\mathcal{N}(0, 1)$	norm	qnorm(p)	pnorm(q)
Chisquare $\chi^2(n)$	chisq	qchisq(p, n)	pchisq(q, n)
Fisher $\mathcal{F}(n_1, n_2)$	f	qf(p, $n_1, n_2$ )	pf(q, $n_1, n_2$ )
Student $St(n)$	t	qt(p, n)	pt(q, n)

#### Application :

```

1 > pnorm(1.6449) # proba N(0,1) plus petit que 1.6449
2 [1] 0.9500048
3 > qnorm(0.95) # quantile N(0,1) d'ordre 95% proche de 1.6449
4 [1] 1.644854
5 > 1-pnorm(1.96) # proba N(0,1) plus grand que 1.96 proche de 2.5%
6 [1] 0.0249979
7 > qnorm(c(.95,.975,.99)) # quantiles N(0,1) d'ordre 95%, 97.5% et 99%
8 [1] 1.644854 1.959964 2.326348
9 > qt(c(.95,.975,.99),10) # quantiles St(10) d'ordre 95%, 97.5% et 99%
10 [1] 1.812461 2.228139 2.763769
11 > pt(c(1.812461,2.228139,2.763769),10) # les probas correspondantes
12 [1] 0.950 0.975 0.990
13 > qchisq(c(.95,.975,.99),10) # quantiles Khi2(10) d'ordre 95%, 97.5% et 99%
14 [1] 18.30704 20.48318 23.20925
15 > pchisq(c(18.30704,20.48318,23.20925),10) # les probas correspondantes
16 [1] 0.950 0.975 0.990
17 > qf(c(.95,.975,.99),10,20) # quantiles F(10,20) d'ordre 95%, 97.5% et 99%
18 [1] 2.347878 2.773671 3.368186
19 > pf(c(2.347878,2.773671,3.368186),10,20) # les probas correspondantes
20 [1] 0.950 0.975 0.990

```

**Illustration du lien entre A.E.P. et A.M.P. :** Une instruction `rloi(n,...)` (du même type que les intructions `ploi(q,...)` et `qlloi(p,...)` présentées précédemment) permet de générer simultanément  $n$  réalisations  $\mathbf{y} := (y_1, \dots, y_n)$  d'une v.a.  $Y$  ayant pour loi `loi(...)`. Illustrons-le sur une vérification expérimentale (A.E.P.) d'obtention de probabilité, quantile, moyenne et variance relatifs à une loi  $\mathcal{N}(1, 2)$ .

```

1 > yy<-rnorm(10000,1,2) # les m=10000 réalisations ont stockées dans le vecteur yy
2 > yy # les 10 premières et 10 dernières composantes de yy
3 [1] 4.279724e+00 8.447115e-01 -1.098879e+00 2.826055e+00 -1.356146e+00
4 [6] 1.540536e+00 2.304664e+00 -3.084724e+00 -1.098613e+00 9.650271e-01
5 ...
6 [9991] 2.483222e+00 3.517996e-01 -1.382401e-02 2.162169e+00 4.103853e-01
7 [9996] -1.998113e+00 5.178801e+00 6.135185e-01 -3.672471e-01 9.240147e-01
8 > mean(yy<0.5) # proportion des m=10000 composantes strictement inférieur à 0.5
9 [1] 0.4053
10 > pnorm(0.5,1,2) # idem si m=infini
11 [1] 0.4012937
12 > mean(yy==0.5) # proportion des m=10000 composantes égale à 0.5 (=0 si m=infini)
13 [1] 0
14 > mean(0.5<=yy && yy<=3) # proportion des m=10000 composantes compris entre 0.5 et 3
15 [1] 0
16 > pnorm(3,1,2)-pnorm(.5,1,2) # idem si m=infini
17 [1] 0.4400511
18 > quantile(yy,.95) # quantile d'ordre 95% des m=10000 composantes
19 95%
20 4.316004
21 > qnorm(.95,1,2) # idem si m=infini
22 [1] 4.289707
23 > mean(yy) # moyenne des m=10000 composantes (=1 si m=infini)
24 [1] 0.9720361

```

```
25 |> var(yy)                                # variance des m=10000 composantes (=2^2=4 si m=infini)
26 | [1] 4.134311
```

Tableaux de lois usuelles de variables aléatoires continues (pour la statistique)			
Nom	Graphique	Densité de probabilité	Esperance et Variance
Uniforme $X \rightsquigarrow \mathcal{U}([a, b])$ $a < b$	<p><math>\mathcal{U}([0, 1])</math> et <math>\mathcal{U}([2, 4])</math>.</p>	$f(x) = \begin{cases} \frac{1}{b-a} & \text{si } x \in [a, b] \\ 0 & \text{sinon} \end{cases}$	$E(X) = \frac{a+b}{2}$ et $Var(X) = \frac{(b-a)^2}{12}$ <p>La densité de probabilité d'une loi uniforme est un histogramme à 1 classe.</p>
Normale $X \rightsquigarrow \mathcal{N}(\mu, \sigma)$ $\mu$ réel et $\sigma$ réel $> 0$	<p><math>\mathcal{N}(-2, 0.5), \mathcal{N}(0, 1)</math> puis <math>\mathcal{N}(4, 2)</math>.</p>	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$	1) Si $X \rightsquigarrow \mathcal{N}(\mu, \sigma)$ alors $\frac{X-\mu}{\sigma} \rightsquigarrow \mathcal{N}(0, 1)$ 2) Si $X \rightsquigarrow \mathcal{N}(\mu_X, \sigma_X)$ et $Y \rightsquigarrow \mathcal{N}(\mu_Y, \sigma_Y)$ sont des v.a. indépendantes alors $X + Y \rightsquigarrow \mathcal{N}\left(\mu_X + \mu_Y, \sqrt{\sigma_X^2 + \sigma_Y^2}\right).$
Chisquare $X \rightsquigarrow \chi^2(\nu)$ $\nu$ entier $> 0$	<p><math>\nu = 3, \nu = 6</math> puis <math>\nu = 9</math>.</p>	$f(x) = \begin{cases} \frac{e^{-\frac{x}{2}} x^{\frac{(\nu-2)}{2}}}{2^{\frac{\nu}{2}} \Gamma(\frac{\nu}{2})} & \text{si } x > 0 \\ 0 & \text{sinon} \end{cases}$	Si $X_1, \dots, X_n$ sont $n$ lois $\mathcal{N}(0, 1)$ indépendantes alors $Y = \sum_{i=1}^n X_i^2 \rightsquigarrow \chi^2(n)$
Student $X \rightsquigarrow \mathcal{St}(\nu)$ $\nu$ entier $> 0$	<p><math>\nu = 2</math> et <math>\nu = 30</math>.</p>	$f(x) = \frac{\left(1 + \frac{x^2}{\nu}\right)^{-\frac{1}{2}(\nu+1)}}{\beta\left(\frac{1}{2}, \frac{\nu}{2}\right) \sqrt{\nu}}$	Si $X \rightsquigarrow \mathcal{N}(0, 1)$ et $Y \rightsquigarrow \chi^2(\nu)$ sont indépendantes alors $Z = \frac{X}{\sqrt{\frac{Y}{\nu}}} \rightsquigarrow \mathcal{St}(\nu)$
Fisher $X \rightsquigarrow \mathcal{F}(\nu_1, \nu_2)$ $\nu_1, \nu_2$ entiers $> 0$	<p><math>\mathcal{F}(5, 200), \mathcal{F}(200, 5)</math> puis <math>\mathcal{F}(30, 30)</math>.</p>	$f(x) = \begin{cases} \frac{\frac{1}{2}\nu_1 \frac{1}{2}\nu_2 x^{\frac{\nu_1}{2}-1}}{(\nu_1 x + \nu_2)^{\frac{1}{2}(\nu_1 + \nu_2)} \beta\left(\frac{\nu_1}{2}, \frac{\nu_2}{2}\right)} & \text{si } x > 0 \\ 0 & \text{sinon} \end{cases}$	Si $X_1 \rightsquigarrow \chi^2(\nu_1)$ et $X_2 \rightsquigarrow \chi^2(\nu_2)$ sont indépendantes alors $Y = \frac{X_1/\nu_1}{X_2/\nu_2} \rightsquigarrow \mathcal{F}(\nu_1, \nu_2)$