

Glossaire de statistique pour les professeurs de mathématiques de collège

Glossaire à l'usage des professeurs (ou futurs professeurs) de mathématiques ; c'est sur des exemples que les notions seront introduites auprès des élèves.

Nous définissons tout d'abord le vocabulaire statistique concernant les *données observées sur une population finie* (une classe d'élèves par exemple), puis le vocabulaire supplémentaire concernant les données observées sur un *échantillon* extrait d'une population finie, enfin, nous reprenons l'interprétation des mêmes termes dans le contexte de *données expérimentales* (durée d'un même trajet chaque jour ouvrable de huit semaines consécutives, même moyen de locomotion, départ à la même heure, ...).

Les données observées sur une population finie

<i>étude statistique</i>	C'est dans le but de rechercher de l'information sur une question donnée que l'on peut entreprendre une <i>étude statistique</i> ; il s'agit alors de définir la <i>population</i> et les différents <i>caractères</i> (appelés aussi <i>variables</i>) définis sur cette <i>population</i> qui pourront apporter l'information cherchée.
<i>population (statistique)</i>	ensemble fini et homogène auprès duquel on recherche l'information, noté $E = \{e_1, \dots, e_n\}$ Ce peut être une population humaine ou une population d'entreprises ou de pays, ... Cet ensemble n'est pas ordonné. Lorsque l'ensemble est ordonné, les douze mois d'une année donnée par exemple, les données observées sont appelées <i>séries chronologiques</i> ; le questionnement, et donc le traitement statistique, différent en général de ce qui est présenté ici.
<i>individus (statistiques)</i>	ou <i>unités statistiques</i> : éléments de la <i>population</i>
<i>caractère (ou variable)</i>	application définie sur la <i>population</i> ; si on note X le caractère, il peut être identifié au n -uplet $(X(e_1), \dots, X(e_n))$ des <i>observations</i> de X sur E Dans le cadre d'une enquête par questionnaire, pour une question autorisant plusieurs réponses (sports pratiqués par exemple), on construit autant de <i>caractères dichotomiques</i> (oui, non) que de réponses possibles
<i>série statistique (à 1 variable) de taille n</i>	n -uplet de valeurs pouvant être considéré comme le n -uplet des <i>observations</i> d'un caractère sur une population de taille n
<i>observations</i>	images par un caractère X des individus statistiques, termes de la série statistique
<i>caractère qualitatif (ou variable catégorielle)</i>	<i>caractère</i> à valeurs dans un ensemble fini dont les éléments sont appelés <i>modalités</i>
<i>caractère quantitatif (ou variable réelle)</i>	<i>caractère</i> à valeurs dans l'ensemble des nombres réels

<i>caractère quantitatif discret</i>	<i>caractère quantitatif</i> dont l'ensemble des valeurs possibles est discret
<i>caractère quantitatif continu</i>	<i>caractère quantitatif</i> dont l'ensemble des valeurs possibles est un intervalle ; les valeurs sont alors bien souvent regroupées en <i>classes</i>
$[X \in I]$	notation en statistique et probabilités de $X^{-1}(I)$, image réciproque par l'application X d'un sous-ensemble I de l'ensemble d'arrivée : $[X \in I] = \{e \in E ; X(e) \in I\}$; généralisation $[X = x], [X \leq x], \dots$
<i>codage</i>	numérotation (de 1 à r des r modalités d'un <i>caractère qualitatif</i>) Un caractère qualitatif codé apparaît comme un caractère quantitatif discret mais faire des calculs sur les codes n'a aucun sens. Un caractère quantitatif peut être considéré comme caractère qualitatif tant que l'on ne fait aucun calcul sur les valeurs prises par le caractère. Lorsque le caractère a deux modalités (caractère dichotomique), on préfère souvent les coder 0 et 1.
<i>caractère dichotomique</i>	<i>caractère qualitatif</i> à deux modalités ; lorsque ces modalités sont codées 0 et 1, alors le caractère qualitatif codé est un caractère quantitatif appelé <i>indicateur</i>
<i>indicateur</i>	caractère quantitatif ne prenant que deux valeurs 0 et 1 Il s'agit donc de la fonction caractéristique du sous-ensemble A de la population prenant la valeur 1 (on parle alors de l'indicateur de A). La moyenne de l'indicateur de A est la proportion p de A dans la population et sa variance est égale à $p(1-p)$.
<i>classes</i>	intervalles de l'ensemble des réels ; deux à deux disjoints, la réunion est un intervalle recouvrant l'ensemble des valeurs d'un <i>caractère quantitatif continu</i> . Notation : $[x_{i-1}, x_i[; i \in \{1, \dots, r\}$ avec $x_0 < x_1 < \dots < x_r$ et $X(E) \subset [x_0, x_r]$. on suppose que, pour chaque classe, les observations de la classe sont uniformément réparties dans la classe ; le regroupement des valeurs en classes fait perdre de l'information pour davantage de lisibilité (histogramme)
<i>catégorie</i>	sous-ensemble d'une <i>population</i> de tous les <i>individus</i> prenant une même <i>modalité</i> d'un <i>caractère qualitatif</i> L'ensemble des catégories associées à un caractère qualitatif est la partition de la population engendrée par le caractère qualitatif. On définira de même la partition engendrée par un caractère quantitatif discret ou un caractère quantitatif continu dont les valeurs sont regroupées en classes.
<i>effectif</i>	nombre d'éléments d'un sous-ensemble de la <i>population</i> (effectif ou "fréquence" ou "fréquence absolue" dans quelques ouvrages français "frequency" en anglais, "frecuencia absoluta" en espagnol)
<i>fréquence (relative)</i>	rapport de l' <i>effectif</i> d'un sous-ensemble de la <i>population</i> sur l' <i>effectif</i> de la <i>population</i> ("relative frequency" en anglais, "frecuencia relativa" en espagnol)
<i>distribution d'effectifs (d'un caractère)</i>	ensemble $\{(x_i, n_i), i \in \llbracket 1, r \rrbracket\}$ où $\{x_i, i \in \llbracket 1, r \rrbracket\}$ est l'ensemble des r valeurs distinctes du caractère X et n_i l'effectif du sous-ensemble de la population prenant la valeur x_i ; $n = \sum_{i=1}^r n_i$ Pour un caractère quantitatif continu dont les valeurs sont regroupées en classes, on remplace les valeurs par les classes de valeurs. La distribution d'effectifs est souvent présentée dans un tableau à 2 lignes ou 2 colonnes dont la première est l'ensemble des valeurs (ou des classes) (rangées usuellement dans l'ordre croissant lorsque l'ensemble des valeurs est ordonné) et la deuxième est l'ensemble des effectifs correspondants.

<i>distribution de fréquences (d'un caractère)</i>	<p>ensemble $\{(x_i, f_i); i \in \llbracket 1, r \rrbracket\}$ où $\{x_i; i \in \llbracket 1, r \rrbracket\}$ est l'ensemble des r valeurs distinctes du caractère X et f_i la fréquence du sous-ensemble de la population prenant la valeur x_i; $f_i = n_i / n$ avec $n = \sum_{i=1}^r n_i$</p> <p>Pour un caractère quantitatif continu dont les valeurs sont regroupées en classes, on remplace les valeurs par les classes.</p> <p>La distribution de fréquences est souvent présentée dans un tableau 2 lignes ou 2 colonnes dont la première est l'ensemble des valeurs (rangées usuellement dans l'ordre croissant lorsque l'ensemble des valeurs est ordonné) et la deuxième est l'ensemble des fréquences correspondantes.</p>
<i>effectifs cumulés</i>	<p>$N_0 = 0$, et pour $i = 1, \dots, r$, $N_i = \sum_{j=1}^i n_j$</p> <p>pour un caractère quantitatif (ou un caractère qualitatif dont les modalités sont ordonnées)</p>
<i>fréquences cumulées</i>	<p>$F_0 = 0$, et pour $i = 1, \dots, r$, $F_i = \sum_{j=1}^i f_j$; on a donc : $F_i = N_i / n$</p> <p>pour un caractère quantitatif (ou un caractère qualitatif dont les modalités sont ordonnées)</p>
<i>fonction de répartition (f.d.r.)</i>	<p>d'un caractère quantitatif, définie sur \mathbb{R} par : $F(x) = \text{Freq}([X \leq x])$</p> <p>Cas discret : $F(x) = \sum_{i=1}^r F_i \mathbf{1}_{[x_i, x_{i+1}[}(x)$, avec $x_{r+1} = +\infty$</p> <p>F est une fonction en escalier, croissante, à valeurs dans $[0, 1]$, nulle sur $]-\infty, x_1[$, égale à 1 sur $[x_r, +\infty[$, continue sauf aux points d'abscisse x_i où elle admet une discontinuité de saut f_i et où elle est continue à droite.</p> <p>Cas continu : $F(x) = \sum_{i=1}^r \left[\frac{x - x_{i-1}}{x_i - x_{i-1}} (F_i - F_{i-1}) + F_{i-1} \right] \mathbf{1}_{[x_{i-1}, x_i[}(x) + \mathbf{1}_{[x_r, \infty[}(x)$</p> <p>$F$ est une fonction continue, affine par morceaux, croissante, à valeurs dans $[0, 1]$, nulle sur $]-\infty, x_0]$, égale à 1 sur $[x_r, +\infty[$. Si, pour tout i de $\{1, \dots, r\}$ $f_i > 0$, la restriction de F à $[x_0, x_r]$ est une bijection à valeurs dans $[0, 1]$.</p> <p>On utilise ici à nouveau l'hypothèse d'uniforme répartition des observations dans les classes.</p>
<i>diagramme en secteurs</i>	<p>type de graphique permettant de représenter la <i>distribution d'effectifs</i> (ou de <i>fréquences</i>) d'un <i>caractère qualitatif</i>; il s'agit d'un disque (parfois d'un demi-disque) composé de secteurs angulaires (représentant les modalités) dont les mesures d'angle sont proportionnelles aux effectifs (ou aux fréquences) des modalités</p> <p>Il peut être également utilisé pour représenter, par exemple, la répartition du budget de la commune selon les différents postes budgétaires. Dans ce cas, il ne s'agit pas d'un diagramme d'effectifs ou de fréquences d'un caractère qualitatif. Les euros ne sont pas des effectifs, les différents postes budgétaires ne sont pas les modalités d'un caractère qualitatif.</p>
<i>diagramme en barres</i>	<p>type de graphique permettant de représenter la <i>distribution d'effectifs</i> (ou de <i>fréquences</i>) d'un <i>caractère qualitatif</i>; on place les modalités de la variable sur un axe horizontal (à égales distances les unes des autres) et on élève au dessus de ces modalités des barres de hauteurs proportionnelles aux effectifs (ou aux fréquences) des modalités</p> <p>L'ordre des modalités sur l'axe horizontal peut parfois donner une information trompeuse. (diagramme parfois appelé en "tuyaux d'orgue", les barres ou les tuyaux d'orgue peuvent être présentés horizontalement et non verticalement)</p>
<i>diagramme en bâtons</i>	<p>type de graphique permettant de représenter la <i>distribution d'effectifs</i> (ou de <i>fréquences</i>) d'un <i>caractère quantitatif discret</i>; on place les r valeurs distinctes de la variable sur un axe horizontal représentant la droite des nombres réels et on élève au dessus de ces valeurs des bâtons de hauteurs</p>

	proportionnelles aux effectifs (ou aux fréquences) des valeurs
<i>histogramme</i>	type de graphique permettant de représenter la <i>distribution d'effectifs</i> (ou de <i>fréquences</i>) d'un <i>caractère quantitatif continu</i> dont les valeurs sont regroupées en classes ; on place les classes sur un axe horizontal représentant la droite des nombres réels et on élève au dessus de ces classes des rectangles de mesures d'aires proportionnelles aux effectifs (ou aux fréquences) ; les rectangles sont en cohérence avec l'hypothèse faite selon laquelle les observations d'une même classe sont uniformément réparties dans la classe
<i>graphe des effectifs cumulés</i>	graphe de la fonction définie sur \mathbb{R} par: $N(x) = \text{card}(\{X \leq x\}) = n F(x)$ où F est la <i>fonction de répartition</i> de X
<i>graphe des fréquences cumulées</i>	graphe de la <i>fonction de répartition</i> de X
<i>indice de position ou de tendance centrale</i>	résumé numérique d'une distribution d'effectifs ou de fréquences d'un caractère quantitatif correspondant à une valeur située au "milieu" (selon un critère à préciser) de la distribution, par exemple, le mode, la moyenne, la médiane, ... Il s'agit d'un résumé numérique d'une distribution qui entraîne une grande perte d'information ; aussi, il faut éviter de l'appeler "caractéristique" qui a un autre sens en mathématiques, d'autant plus que ce mot est parfois utilisé à la place de "caractère". Dans le cas où les données sont observées sur un <i>échantillon</i> , on appelle "paramètres" les indices concernant la population (qu'il s'agit d'estimer à partir des indices similaires calculés sur l'échantillon).
<i>indice de dispersion</i>	résumé numérique d'une distribution d'effectifs ou de fréquences d'un caractère quantitatif mesurant le plus ou moins grand "étalement" de la distribution, par exemple, l'étendue, l'écart-type, l'écart interquartile, ...
<i>mode, classe modale</i>	(<i>indice de position d'un caractère quantitatif</i>) le <i>mode</i> est la valeur du caractère correspondant à l'effectif (ou la fréquence) maximal ; dans le cas d'un caractère quantitatif continu dont les données sont regroupées en classes, la <i>classe modale</i> est celle dont l'effectif (ou la fréquence) par unité du caractère est maximal un caractère peut avoir plusieurs modes ou classes modales
<i>étendue</i>	(<i>indice de dispersion d'un caractère quantitatif</i>) l'étendue est la différence $x_r - x_1$ entre la valeur maximale du caractère x_r (appelée <i>max</i>) et la valeur minimale x_1 (appelée <i>min</i>) pour un caractère quantitatif continu dont les valeurs sont regroupées en <i>classes</i> (cf. notation introduite), l'étendue est la différence $x_r - x_0$
<i>moyenne (arithmétique)</i>	(<i>indice de position d'un caractère quantitatif</i>) la <i>moyenne</i> est la somme des n valeurs du caractère divisée par n ; elle correspond à la valeur commune qu'auraient les n individus de la population s'ils se partageaient de façon égale la somme des valeurs positives ou négatives du caractère ; la moyenne a la dimension du caractère X et elle est notée \bar{X} ou \bar{x} (la notation μ est utilisée pour la moyenne de la population dont serait extrait un échantillon) $\bar{x} = \frac{1}{n} \sum_{k=1}^n X(e_k) = \frac{1}{n} \sum_{i=1}^r n_i x_i = \sum_{i=1}^r f_i x_i \text{ avec } n = \sum_{i=1}^r n_i \text{ et } f_i = \frac{n_i}{n}$

<p><i>variance et écart-type</i></p>	<p>(indices de dispersion d'un caractère quantitatif)</p> <p>la <i>variance</i> est la moyenne des carrés des écarts à la moyenne, soit :</p> $\text{var}(X) = \frac{1}{n} \sum_{k=1}^n (X(e_k) - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^r n_i (x_i - \bar{x})^2 = \sum_{i=1}^r f_i (x_i - \bar{x})^2$ <p>la variance est de dimension le carré de la dimension de X ; on vérifie aisément que la variance est égale à la moyenne des carrés moins le carré de la moyenne (2^{ème} formule de la variance utilisée dans les calculs), soit :</p> $\text{var}(X) = \frac{1}{n} \sum_{k=1}^n (X(e_k))^2 - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^r n_i x_i^2 - \bar{x}^2 = \sum_{i=1}^r f_i x_i^2 - \bar{x}^2$ <p>l'<i>écart-type</i> est la racine carrée de la variance (de dimension celle de X) ; il est noté s (la notation σ est utilisée pour l'écart-type de la population dont serait extrait un échantillon).</p>
<p><i>médiane</i></p>	<p>(indice de position d'un caractère quantitatif)</p> <p>intuitivement, la médiane partage la population en deux parties de même effectif ; plus précisément, un réel m est médiane si au moins la moitié de la population prend des valeurs inférieures ou égales à m et au moins la moitié de la population prend des valeurs supérieures ou égales à m, i.e. (définition d'un <i>quantile d'ordre 0.5</i>) : $\text{Freq}([X \leq m]) \geq 0.5$ et $\text{Freq}([X \geq m]) \geq 0.5$;</p> <p>lorsque l'on range les valeurs des n individus de la plus petite à la plus grande, si n est impair, l'unique médiane est la $[(n+1)/2]^{\text{ème}}$ valeur, si n est pair, n'importe quelle valeur comprise au sens large entre la $n^{\text{ème}}$ et la $(n+1)^{\text{ème}}$ est médiane ; par convention, pour avoir unicité, on prend la demi-somme de ces deux valeurs</p>
<p><i>quantile d'ordre p ($p \in]0, 1[$)</i></p>	<p>q est quantile d'ordre p si la proportion d'observations inférieures ou égales à q est supérieure ou égale à p et si la proportion d'observations supérieures ou égales à q est supérieure ou égale à $1 - p$, i.e. si</p> $\text{Freq}([X \leq q]) \geq p \text{ et } \text{Freq}([X \geq q]) \geq 1 - p ;$ <p>l'ensemble des quantiles d'ordre p est un intervalle fermé non vide ; pour avoir unicité, différentes conventions sont utilisées</p>
<p><i>fonction quantile</i></p>	<p>$\forall p \in]0, 1[, q(p) = \inf \{x \in \mathbb{R} ; F(x) \geq p\}$</p> <p>$q(p)$ est la borne inférieure de l'intervalle des quantiles d'ordre p ; $q(p)$ appartient à $X(E)$, c'est la plus petite valeur q de la série des observations ordonnées en croissant telle que la proportion d'observations inférieures ou égales à q soit au moins égale à p</p>
<p><i>médiane, quartiles, déciles</i></p>	<p>On appelle <i>médiane</i> un quantile d'ordre 0.5.</p> <p>On appelle <i>1^{er} quartile</i> (resp. <i>2^{ème} quartile</i>, resp. <i>3^{ème} quartile</i>) un quantile d'ordre 0.25 (resp. 0.50, resp. 0.75).</p> <p>On appelle <i>1^{er} décile</i> (resp. <i>2^{ème} décile</i>, ..., resp. <i>9^{ème} décile</i>) un quantile d'ordre 0.10 (resp. 0.20, ..., resp. 0.90).</p> <p>m est médiane $\Leftrightarrow m$ est 2^{ème} quartile $\Leftrightarrow m$ est 5^{ème} décile</p> <p>dans le secondaire, la définition donnée pour le 1^{er} (resp. 3^{ème}) quartile est $q(0.25)$ (resp. $q(0.75)$), la définition donnée pour le 1^{er} (resp. 9^{ème}) décile est $q(0.1)$ (resp. $q(0.9)$).</p>
<p><i>intervalle interquartile, intervalle interdécile</i></p>	<p>On appelle <i>intervalle interquartile</i> (resp. <i>intervalle interdécile</i>) l'intervalle dont les extrémités sont le 1^{er} et le 3^{ème} quartiles (resp. le 1^{er} et le 9^{ème} déciles) après avoir éventuellement utilisé une convention pour l'unicité.</p>

<i>écart interquartile, écart interdécile</i>	On appelle <i>écart interquartile</i> (resp. <i>écart interdécile</i>) la longueur de l'intervalle interquartile (resp. interdécile).
<i>diagramme en boîte (ou en boîte et moustaches)</i>	(en anglais, "box plot" ou "box and whiskers plot") Diagramme représentant une distribution de fréquences construit sur les cinq indices, <i>min</i> , <i>q₁</i> , <i>m</i> , <i>q₂</i> , <i>max</i> , placés sur un axe représentant la droite des réels : la boîte rectangulaire, de largeur arbitraire, est limitée en longueur par le premier et troisième quartile ; à l'intérieur de la boîte est indiquée par un trait la médiane et, de part et d'autre de la boîte, des segments représentent les valeurs extérieures à l' <i>intervalle interquartile</i> , les extrémités de ces segments indiquent les valeurs extrêmes (min et max) de la variable. Les extrémités des segments peuvent aussi correspondre aux premier et neuvième déciles ; des croix indiquent alors les observations extérieures à l' <i>intervalle interdécile</i> . L'avantage de ce type de diagramme est sa construction rapide et la possibilité de faire une comparaison visuelle de plusieurs distributions.

Les données d'enquête observées sur un *échantillon*

<i>sondage</i>	dans le langage courant, il s'agit d'une enquête d'opinion ("poll" en anglais, "sondeo" en espagnol), en statistique, il s'agit d'une enquête auprès d'une partie de la <i>population</i> appelée <i>échantillon</i> ("sampling" en anglais, "muestreo" en espagnol)
<i>recensement</i>	s'oppose à <i>sondage</i> ; l'enquête est réalisée auprès de la <i>population</i> tout entière
<i>base de sondage</i>	liste numérotée de tous les individus de la <i>population</i> , nécessaire pour sélectionner un échantillon aléatoire, notée $E = \{e_1, \dots, e_N\}$
<i>échantillon</i>	sous-ensemble de la <i>population</i> ; ce peut être un <i>échantillon aléatoire</i> (ou <i>échantillon probabiliste</i>) ou un <i>échantillon empirique</i> On ne parle d'échantillon que lorsque l'on souhaite inférer à la <i>population</i> dont il est issu les indices obtenus sur l'échantillon ; une classe d'élèves est alors une <i>population</i> et non un échantillon.
<i>taille</i>	nombre d'éléments de l' <i>échantillon</i> (noté usuellement <i>n</i>) ou de la <i>population</i> (noté usuellement <i>N</i>) ou de sous-populations ou sous-échantillons ...
<i>taux de sondage</i>	rapport de la <i>taille</i> de l' <i>échantillon</i> sur la <i>taille</i> de la <i>population</i>
<i>échantillon aléatoire (ou probabiliste)</i>	échantillon sélectionné selon une <i>procédure aléatoire</i> (et en utilisant un <i>générateur de nombres pseudo-aléatoires</i> ou une <i>table de nombres au hasard</i>)
<i>échantillon empirique</i>	échantillon non aléatoire ; le plus connu est l'échantillon empirique obtenu par la <i>méthode des quotas</i>
<i>échantillon avec remise</i>	les <i>n</i> individus de l'échantillon sont tirés un à un en remettant dans la <i>population</i> de taille <i>N</i> l'individu sélectionné après chaque tirage ; un même individu de la <i>population</i> peut apparaître plusieurs fois dans l'échantillon ; un tel échantillon est représenté par un élément de l'ensemble E^n de cardinal N^n

<i>échantillon sans remise</i>	<p>les n individus de l'échantillon sont tirés un à un sans remettre dans la population de taille N les individus déjà sélectionnés ; un tel échantillon est représenté par un arrangement de n éléments de E, il y a A_N^n tels échantillons.</p> <p>Dans la mesure où l'on ne tient pas compte de l'ordre de tirage des individus de l'échantillon, on peut considérer un échantillon sans remise comme un sous-ensemble de n éléments de E ; on a $\binom{N}{n}$ tels échantillons.</p>
<i>échantillon aléatoire simple à probabilités égales avec remise</i>	un des N^n échantillons de taille n avec remise d'une population de taille N (tous les échantillons ayant la même probabilité d'être tiré)
<i>échantillon aléatoire simple à probabilités égales sans remise</i>	un des $\binom{N}{n}$ échantillons de taille n sans remise d'une population de taille N (tous les échantillons ayant la même probabilité d'être tiré)
<i>échantillon aléatoire stratifié</i>	on dispose dans la <i>base de sondage</i> d'une (ou plusieurs) variable catégorielle permettant de définir une partition de la population en sous-populations appelées <i>strates</i> ; on tire de façon indépendante, un échantillon aléatoire simple de taille fixée dans chaque strate ; l'échantillon global est la réunion de tous ces échantillons.
<i>échantillonnage aléatoire stratifié proportionnel</i>	échantillon aléatoire stratifié tel que les tailles des échantillons soient proportionnelles aux tailles des strates
<i>méthode des quotas</i>	il s'agit d'une méthode de sélection non aléatoire d'un échantillon tentant d'imiter un échantillon aléatoire stratifié proportionnel ; on connaît la répartition de la population selon deux ou trois variables catégorielles, appelées variables de quotas, l'échantillon doit respecter ces répartitions (appelées quotas)

Au niveau du collège, on travaillera sur des populations et non sur des échantillons.

Au niveau du lycée, on pourra travailler sur des échantillons aléatoires simples à probabilités égales avec remise. Le vocabulaire statistique défini sur la population dans la partie précédente s'applique ici en remplaçant population par échantillon.

Une fois l'échantillon décrit, l'objectif est d'inférer à la population, dont les *paramètres* (proportions, moyenne, écart-type, ...) sont inconnus, les informations obtenues (calculées) sur l'échantillon.

La loi des grands nombres permet d'approcher une proportion (ou fréquence) relative à la population (appelée aussi *fréquence théorique*) par une proportion (ou fréquence) relative à l'échantillon (appelée aussi *fréquence empirique* ou observée) et donc une *distribution de fréquences théoriques* (sur la population) par une *distribution de fréquences empiriques* (sur l'échantillon).

Le théorème central limite permet d'utiliser l'approximation gaussienne pour estimer par intervalle de confiance les proportions ou moyennes de la population à partir des proportions ou moyennes calculées sur l'échantillon.

On notera que l'observation x d'un caractère (qualitatif ou quantitatif) sur *UN individu tiré avec équiprobabilité d'une population finie* est l'observation d'une *variable aléatoire* (catégorielle ou réelle) dont la distribution de probabilité n'est autre que la distribution de fréquences du caractère (le "nombre de cas favorables sur nombre de cas possibles" n'est autre que la fréquence). Nous ne faisons aucune hypothèse sur cette distribution de probabilité, l'aléatoire vient de la procédure de sélection de l'individu (on simule l'équiprobabilité).

Pour un échantillon *aléatoire* de taille n tiré avec équiprobabilité et avec remise d'une population finie, la série statistique de taille n EST l'observation de n variables aléatoires indépendantes et identiquement distribuées (i.i.d.), la distribution de probabilité EST la distribution de fréquences sur la population du caractère considéré (*approche sondage*).

Données expérimentales

<i>étude expérimentale</i>	C'est dans le but de rechercher de l'information sur une question donnée que l'on peut décider de mettre en place un <i>dispositif expérimental</i> ; il s'agit alors de définir les conditions expérimentales permettant de réaliser une série d' <i>expériences</i> dans les mêmes conditions et de définir la ou les variables qui seront observées pour chacune des expériences et qui pourront apporter l'information cherchée.
<i>série statistique (à 1 variable) de taille n</i>	n -uple de valeurs pouvant être considéré comme le n -uple des observations d'une variable sur une suite de n <i>expériences</i>

Si on note $E = \{e_1, \dots, e_n\}$ l'ensemble des expériences, nous sommes dans la situation précédente, observation d'une variable sur un échantillon de taille n . Mais de quelle population serait extrait cet échantillon ? On adaptera donc le vocabulaire décrit dans la première partie pour décrire la série numérique sans référence à des individus statistiques ou à une population.

Si l'expérience est reproduite dans les mêmes conditions, *on fera l'hypothèse* que la série statistique de taille n est l'observation de n variables aléatoires indépendantes et de même distribution de probabilité (i.i.d.) mais on ne connaît en général rien du type de la distribution de probabilité et on pourra être amené à *faire des hypothèses* sur cette distribution de probabilité (*approche modèle*).

C'est par abus de langage que l'on parle encore, dans le cadre de données expérimentales, de *paramètres de la population* ; il s'agit en fait de paramètres de la loi de probabilité du modèle.