

Statistiques

En statistiques, je fais des sondages, des essais sur un échantillon ou la totalité d'une population. Je traduis alors ces données et les décris par les nombres.

Je ne cherche pas à calculer les résultats que j'aurais le plus de « chance » d'avoir en théorie, sinon je fais des probabilités. Si l'échantillon, le nombre de sondages ou les essais tendent vers l'infini, ces résultats statistiques doivent correspondre aux valeurs théoriques déterminées par les lois de probabilité.

Vocabulaire - Définitions

- ◆ **x :** **valeur** obtenue au cours d'une étude statistique
Ex. : somme d'argent, note d'élève, score, etc. ...
 - ◆ **($x_1 ; x_2 ; x_3 ; \dots ; x_{k-1} ; x_k$) :** **série statistique de taille k** regroupant toutes les valeurs x_i obtenues (elles sont généralement toutes différentes deux à deux, donc $k \leq N$).
 - **n_i :** **effectif associé à la valeur x_i** , c'est-à-dire le nombre de fois où cette valeur x_i est « apparue » dans la série statistique.
 - **N :** **effectif total** (ou population totale), somme de tous les effectifs $n_i \rightarrow N = \sum_{i=1}^k n_i$
- Remarque : Si la série statistique est $(x_1 ; x_2 ; x_3 ; \dots ; x_{N-1} ; x_N)$, alors $k = N$, les valeurs x_j n'étant pas forcément toutes différentes, donc on aura : $n_j = 1 \quad \forall j$. $\rightarrow N = \sum_{j=1}^N n_j = \sum_{j=1}^N 1$
- ◆ **Événement :** il désigne le fait qu'une valeur x_i apparaisse ou qu'une action se produise

Ex. : « Voici le tableau de notes d'une classe après un devoir.

x = Notes (/20)	7,5	9	11	11,5	13	15	16	19
n = nombre d'élèves	2	4	6	5	4	2	3	2

Soit A l'événement : 'j'ai entre 10/20 et 15/20 (notes incluses)'.

Calcule l'effectif total de cette classe N et l'effectif n_A correspondant à l'événement A.

Précise quelle valeur x_i a le plus grand effectif associé. »

Nous travaillons avec la série **statistique de taille 8** : (7,5 ; 9 ; 11 ; 11,5 ; 13 ; 15 ; 16 ; 19).

→ Effectif total de la classe : $N = \sum_{i=1}^k n_i = 2 + 4 + 6 + 5 + 4 + 2 + 3 + 2 = 28$

Il y a **N = 28 élèves** dans cette classe.

→ $n_A = 6 + 5 + 4 + 2 = 17$ **élèves ont eu entre 10/20 et 15/20 (notes incluses).**

→ la valeur $x_3 = 11/20$ **est la note la plus fréquemment obtenue par les élèves** ($n_3 = 6$).

- **Fréquence** : $f_i = \frac{n_i}{N} \rightarrow$ indique la **proportion de la valeur** x_i dans la série statistique.
NB : Si $N \rightarrow +\infty$, alors : $f_i \rightarrow p_i$ (probabilité associée à x_i)
- **Moyenne** : $\bar{x} = \frac{1}{N} \cdot \sum_{i=1}^k n_i \cdot x_i = \sum_{i=1}^k \frac{n_i}{N} \cdot x_i = \sum_{i=1}^k f_i \cdot x_i \rightarrow$ donne le « centre » de la répartition des valeurs x_i relevées.
NB : $\overline{x^2} = \frac{1}{N} \cdot \sum_{i=1}^k n_i \cdot x_i^2 = \sum_{i=1}^k \frac{n_i}{N} \cdot x_i^2 = \sum_{i=1}^k f_i \cdot x_i^2 \rightarrow$ (moyenne des carrés des valeurs x_i)
- **Variance** : $V = \frac{1}{N} \cdot \sum_{i=1}^k n_i \cdot (x_i - \bar{x})^2 = \sum_{i=1}^k f_i \cdot (x_i - \bar{x})^2$
$$V = \overline{(x_i - \bar{x})^2}$$

 $\rightarrow V$ est le nombre égal à la moyenne des valeurs $(x_i - \bar{x})^2$, c'est-à-dire la **moyenne des carrés des écarts des valeurs x_i à la moyenne \bar{x}** . Il est très abstrait mais permet de mesurer la « dispersion » des valeurs x_i autour de la valeur moyenne. **Il permet surtout de calculer ensuite l'écart-type σ .**

Remarque : Il existe une autre formule de la variance qui peut être plus pratique à calculer sur tableur ou calculatrice :

$$\begin{aligned} V &= \frac{1}{N} \cdot \sum_{i=1}^k n_i \cdot (x_i - \bar{x})^2 = \frac{1}{N} \cdot \sum_{i=1}^k n_i \cdot (x_i^2 - 2 \cdot x_i \cdot \bar{x} + \bar{x}^2) \\ &= \frac{1}{N} \cdot \sum_{i=1}^k n_i \cdot x_i^2 - 2 \cdot \bar{x} \cdot \frac{1}{N} \cdot \sum_{i=1}^k n_i \cdot x_i + \bar{x}^2 \cdot \frac{1}{N} \sum_{i=1}^k n_i \\ &= \overline{x^2} - 2 \cdot \bar{x} \cdot \bar{x} + \bar{x}^2 \cdot 1 \\ V &= \overline{x^2} - \bar{x}^2 \end{aligned}$$

Remarque : Si la série statistique relève toutes les valeurs x_j , même si celles-ci ne sont pas toutes différentes deux à deux, alors $k = N$ et $n_j = 1 \quad \forall j$. On peut alors aussi écrire que :

$$V = \frac{1}{N} \cdot \sum_{j=1}^N (x_j - \bar{x})^2 = \frac{1}{N} \cdot \sum_{j=1}^N x_j^2 - \bar{x}^2 = \overline{x^2} - \bar{x}^2$$

- **Écart-type** : $\sigma = \sqrt{V}$
 $\rightarrow \sigma$ mesure la **dispersion des valeurs x_i autour de la valeur \bar{x} moyenne**, mais il a la même unité de mesure que x_i .

Ex. : Si l'on revient sur les notes de la classe de 28 élèves précédente, nous trouvons :

$x = \text{Notes (/20)}$	7,5	9	11	11,5	13	15	16	19
$n = \text{nombre d'élèves}$	2	4	6	5	4	2	3	2
$(x_i - \bar{x})^2$	22,39	10,45	1,52	0,54	0,59	7,66	14,20	45,80

\rightarrow Moyenne de la classe (vivent les tableurs ou calculatrices !) :

$$\begin{aligned} \bar{x} &= \frac{1}{28} \cdot \sum_{i=1}^8 n_i \cdot x_i = \frac{1}{28} \cdot (2 \times 7,5 + 4 \times 9 + 6 \times 11 + 5 \times 11,5 + 4 \times 13 + 2 \times 15 + 3 \times 16 + 2 \times 19) = \frac{342,5}{28} \\ \bar{x} &\approx 12,23 \end{aligned}$$

→ Variance des notes de la classe :

$$V \approx \frac{1}{28} \cdot \sum_{i=1}^8 n_i \cdot (x_i - 12,23)^2$$

$$V \approx \frac{1}{28} \cdot (2 \times 22,39 + 4 \times 10,45 + 6 \times 1,52 + 5 \times 0,54 + 4 \times 0,59 + 2 \times 7,66 + 3 \times 14,2 + 2 \times 45,8)$$

$$V \approx \frac{250,24}{28} \approx 8,94 \quad \text{J'en déduis l'écart-type } \sigma :$$

→ Écart-type des notes de la classe :

$\sigma = \sqrt{8,94} \approx 2,99$ Cela signifie qu'une grande majorité des élèves a une note qui ne s'écarte pas plus de 3 points de la note moyenne 12,23 / 20.

Probabilités

En probabilités, je cherche à calculer les résultats que j'ai le plus de « chance » d'avoir en théorie. Si l'échantillon, le nombre de sondages ou les essais tendent vers l'infini, les résultats statistiques doivent rejoindre les valeurs théoriques déterminées par les lois de probabilité.

Vocabulaire - Définitions

- ◆ Ω : Univers, ensemble des événements possibles.
- ◆ Événement A_i : chaque événement A_i peut être associé à :
 - un **nombre de cas possibles** n_i
 - une **valeur** (ou un gain) x_i
- ◆ Événements indépendants : événements qui n'ont **aucun lien** (la réalisation de l'un n'a aucun effet sur la réalisation de l'autre)
- ◆ Tirages successifs : les événements ne sont **pas simultanés**, ils se suivent
- ◆ Tirages successifs avec remise : à chaque tirage, **le nombre d'événements possibles ne diminue pas** (événements indépendants)
- ◆ Tirages successifs sans remise : à chaque tirage, **le nombre d'événements possibles diminue**

Ex. : Si je place dans une urne 2 boules noires et 3 boules rouges, je peux effectuer 3 tirages successifs d'une boule avec remise.

Soit R l'événement : « une boule rouge est tirée »

Soit N l'événement : « une boule noire est tirée »

L'univers Ω sera alors :

$\Omega = \{ (R ; R ; R) ; (R ; R ; N) ; (R ; N ; R) ; (N ; R ; R) ; (R ; N ; N) ; (N ; R ; N) ; (N ; N ; R) \}$,
soient 7 tirages possibles.

- **Probabilité** : $p(A) = \frac{n_A}{n_\Omega} \rightarrow$ « chances », **probabilité que l'événement A se réalise.**

NB : $p(\Omega) = \sum_{i=1}^k p_i = 1 \rightarrow$ La somme des probabilités qu'au moins un des événements de l'univers Ω se réalise vaut toujours 1 (soient 100 % de chances de réalisation).

Ex. : dans le cas précédent, $p(2 \text{ rouges tirées}) = \frac{3}{7} \rightarrow$ 3 chances sur 7 de tirer 2 boules rouges.

- **Espérance** : $E = \frac{1}{n_\Omega} \cdot \sum_{i=1}^k n_i \cdot x_i = \sum_{i=1}^k \frac{n_i}{n_\Omega} \cdot x_i = \sum_{i=1}^k p_i \cdot x_i \rightarrow$ Avec $k \leq n_\Omega$, chaque événement A_i étant associé à un nombre de cas possibles n_i et une valeur (ou gain) x_i , **E mesure la valeur (le gain) que l'on peut espérer gagner (moyenne théorique).**

NB : en statistiques, si $N \rightarrow +\infty$, alors : $f_i \rightarrow p_i$ (probabilité associée à x_i). Si l'on compare les deux formules, on constate donc que :
si $N \rightarrow +\infty$, alors : $\bar{x} \rightarrow E$ (les statistiques rejoignent les probabilités sur de grands échantillons). **L'espérance E est donc (en probabilités) la valeur théorique vers laquelle tend la moyenne \bar{x} (en statistiques).**

Ex. : Dans le cas précédent, si je perds 3 € par boule rouge et gagne 5 € par boule noire, alors j'obtiens le tableau suivant :

i	1	2	3
Événements A_i	3 fois R	2 fois R ; 1 fois N	1 fois R ; 2 fois N
Gains x_i (en €)	- 9 €	-1 €	7 €
Nombre de cas favorables n_i	1	3	3
Probabilité p_i	$\frac{1}{7}$	$\frac{3}{7}$	$\frac{3}{7}$

Je peux donc calculer l'espérance : $E = \sum_{i=1}^3 p_i \cdot x_i = \frac{-9}{7} + \frac{-3}{7} + \frac{21}{7} = \frac{9}{7} \approx 1,29 \text{ €}$
 (ce jeu n'est donc pas une arnaque !).

Statistiques et probabilités avancées

Statistiques et probabilités avancées

(Du discret au continu...)

- ✓ Quand la valeur x n'est plus discrète, on parle de **variable aléatoire continue** :

$$X = k \in I \text{ (X est un réel)}$$

Celle-ci est **associée à un événement** par la notation suivante :

$$\{X = k \in [a ; b]\}, \text{ avec } [a ; b] \subset I.$$

- ✓ La probabilité $p(A_i)$ est remplacée par une probabilité **$P(X \in [a ; b])$** ou **$P(a \leq X \leq b)$** ou **$P([a ; b])$** .
- ✓ Pour calculer une probabilité $P(X \in [a ; b])$, il faut connaître sa **fonction de densité** (ou **densité de probabilité**) f :
- f est **définie et continue sur $I = (x_{\min} ; x_{\max})$**
 - f est **positive sur I**
 - l'**intégrale** de f sur I doit être **égale à 1** (la somme des probabilités vaut toujours 1...)
 - Je peux alors calculer la probabilité $P(X \in [a ; b])$ à partir de sa densité de probabilité f :

$$P(X \in [a ; b]) = \int_a^b f(t).dt$$

NB : ➤ $[a ; b] \subset I$

➤ t est la variable pouvant prendre **toutes les valeurs de X sur $[a ; b]$**

➤ $f(t).dt$ est la probabilité associée à un **intervalle de largeur dt infinitésimale** autour de t .

➤ **une intégrale (fonctions continues) étant l'équivalent d'une somme (de termes discrets)**, on retrouve donc ici une correspondance entre les formules :

Valeurs discrètes (en nombre fini)	Valeurs continues (en nombre infini)
Valeurs $\{x_m ; \dots ; x_{m+q}\} \in [a ; b]$	Valeurs de $X = t \in [a ; b]$
Probabilités $\{p_m = p(x_m) ; \dots ; p_{m+q} = p(x_{m+q})\}$	Infinité de probabilités de $f(a).dt$ à $f(b).dt$
Événement $A_i = \{x_i \in [a ; b]\}$	Événement $\{X \in [a ; b]\}$
Probabilité $p(A_i) = \sum_{i=m}^{m+q} p_i$	Probabilité $P([a ; b]) = \int_a^b f(t).dt$

NB : Que les inégalités ci-dessus soient **strictes ou larges, cela ne change rien !**

$$\text{Preuve : } P(a \leq X) = P(X < a) \quad \text{car : } P([a ; a]) = \int_a^a f(t).dt = F(a) - F(a) = 0$$

- ✓ Si $I = [x_{\min} ; x_{\max}]$, alors : $E(X) = \int_{x_{\min}}^{x_{\max}} t.f(t).dt$

NB : De même que ci-dessus, on retrouve l'équivalent avec la formule discrète : $E = \sum_{i=1}^k p_i \cdot x_i$

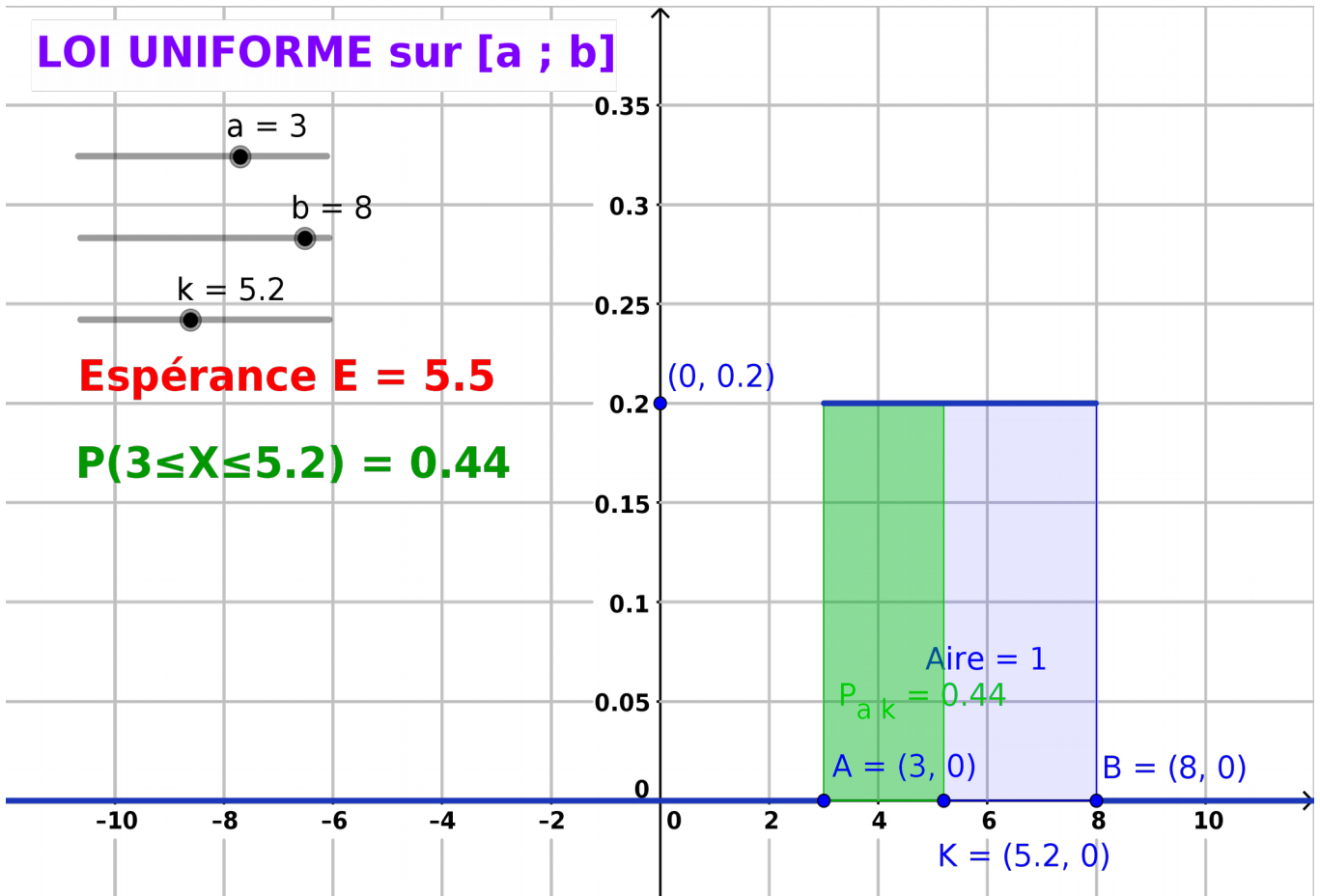
x **Loi uniforme** (notée : $U([a ; b])$) :

$$I = [a ; b]$$

$$f(x) = \frac{1}{b - a} \quad (\text{fonction constante})$$

$$P(a \leq X \leq x) = \frac{x - a}{b - a} \quad \text{avec } x \in [a ; b]$$

$$E(X) = \frac{a + b}{2}$$



x **Durée de vie et lois exponentielles de paramètre λ** :

$$I = [0 ; +\infty[$$

$$f(x) = \lambda \cdot e^{-\lambda \cdot x} \quad \text{avec } \lambda \in \mathbb{R}$$

$$P(X \leq t) = 1 - e^{-\lambda \cdot t} \quad \text{pour tout } t \in [0 ; +\infty[\quad \text{donc :}$$

$$P(X \geq t) = e^{-\lambda \cdot t} \quad \text{pour tout } t \in [0 ; +\infty[$$

$$P([a ; b]) = e^{-\lambda \cdot a} - e^{-\lambda \cdot b}$$

➤ **Durée de vie sans vieillissement :**

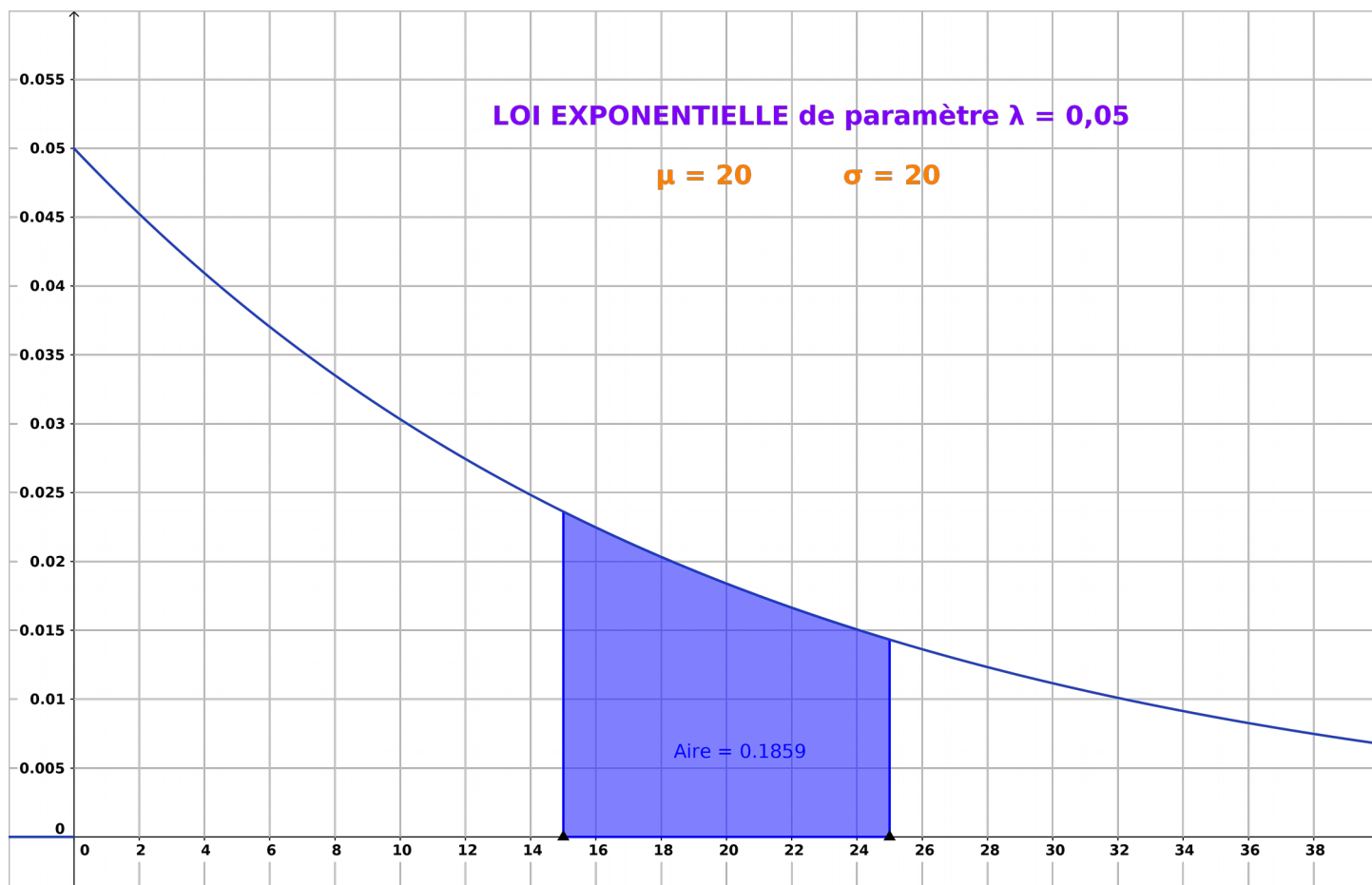
➔ « $X \geq t$ » est l'événement « la durée de vie dépasse t années »

➔ On dit que la durée de vie est « sans vieillissement » lorsque la probabilité qu'il fonctionne encore h heures supplémentaires, sachant qu'il fonctionne à l'instant t , ne dépend pas de t .

La durée de vie sur une période h ne dépend pas de l'âge t à partir duquel on considère cet événement.

➔ Pour tout t et h positifs : **$P_{X \geq t}(X \geq t+h) = P(X \geq h)$**

➔ **$E(X) = \frac{1}{\lambda}$**

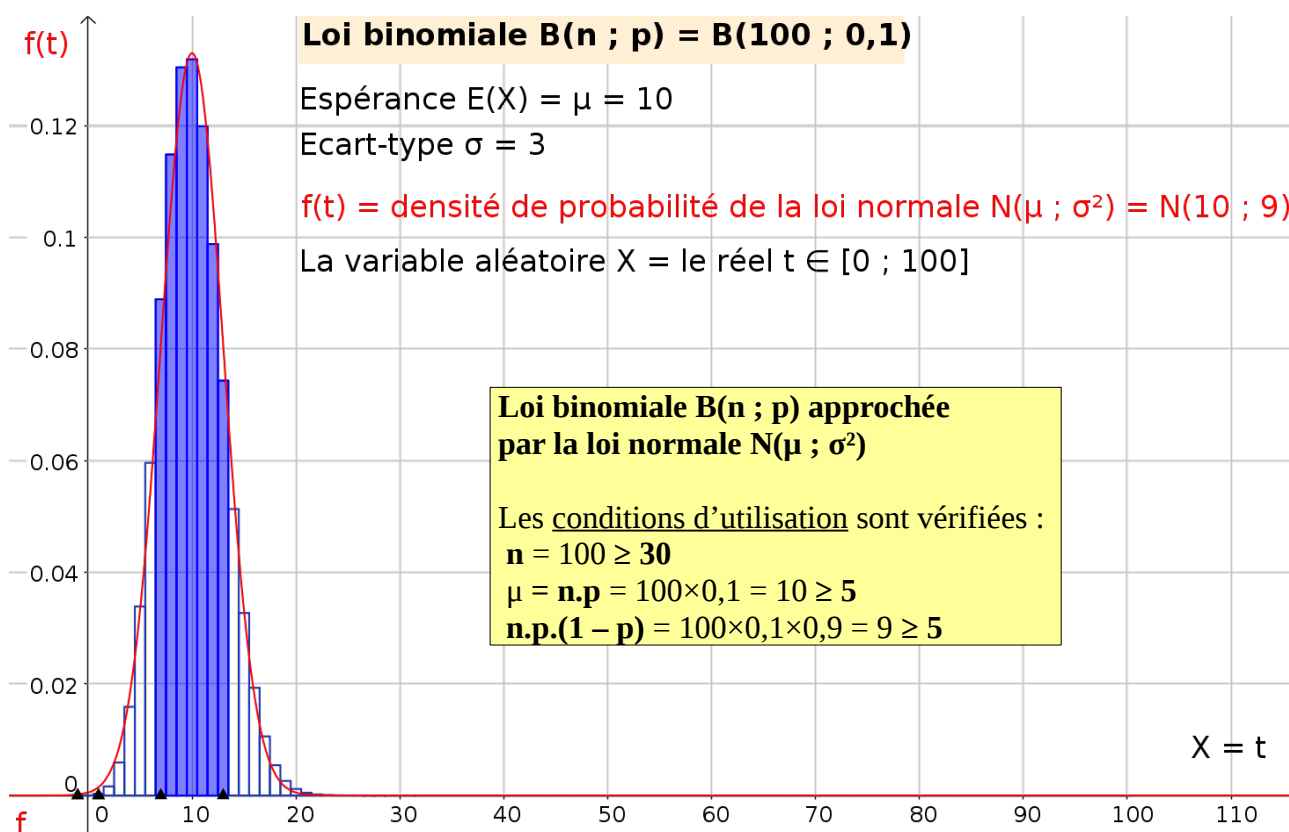


✓ **Loi de BERNOULLI B_p** de paramètre p ,

Loi BINOMIALE $B(n ; p)$ (soient n répétitions d'une épreuve de Bernoulli B_p),

Loi NORMALE $N(\mu ; \sigma^2)$ (espérance μ ; variance σ^2) :

- NB : la loi normale **$N(\mu ; \sigma^2)$** est aussi appelée « **loi de LAPLACE-GAUSS** » (d'espérance μ et de variance σ^2)



Loi :	de Bernoulli B_p ou B_α	Binomiale $B(n; p)$ ou $B(n; \alpha)$	Passage au continu	Normale $N(\mu; \sigma^2)$ ou $N(\mu; \sigma)$
Variable DISCRÈTE (ou « finie ») : calculs exacts			Variable CONTINUE : calculs approchés	
Contexte :	Définie sur un univers à deux éventualités (ou deux issues) : $x \in \{0; 1\}$ avec $p(1)=p \in]0; 1[$	Expérience suivant la loi de Bernoulli B_p et répétée n fois (expériences indépendantes)	Contexte :	Conditions de validité de l'approximation : $n \geq 30$ $n \cdot p \geq 5$ $n \cdot (1 - p) \geq 5$ (p n'est alors pas trop proche de 0 ou 1) Si X suit une loi binomiale $B(n; p)$, je peux alors utiliser les fonctions continues pour faire des calculs approchés.
Exemples :	Pile ou face (un seul lancé) : pile $\rightarrow 0$ face $\rightarrow 1$ $p = 0,5 = p(\text{face})$	Pile ou face (n lancers) : pile $\rightarrow 0$ face $\rightarrow 1$ $p = 0,5 = p(\text{face})$ Combien de fois face ($x=1$) sur n lancers ?	Exemple :	Un appareil en sortie d'usine a 2 issues : bon état $\rightarrow x = 0$ en panne $\rightarrow x = 1$ Production : $n = 1000$. Probabilité de tomber en panne : $p = 0,08$.
Nombre d'expériences	1 seule	n (indépendantes)		n ≥ 30
$n_\Omega =$	2	n+1 \rightarrow tirage 1 : $X \in \{0; 1\}$ (2 cas pour face) \rightarrow tirage 2 : $X \in \{0; 1; 2\}$ (3 cas pour face) \rightarrow tirage n : $X \in \{0; 1; \dots; n\}$ soient n+1 issues face		Grand
Valeurs x :	$x \in \{0; 1\}$	$x \in \{0; 1\}$		$x \in \{0; 1\}$
Variable aléatoire $X_n = k$ = nombre de fois ou ($x = 1$) se réalise	k $\in [0; 1]$ $\rightarrow (X = 1)$ ou $\rightarrow (X = 0)$	k $\in [0; n]$ \rightarrow Face peut sortir 1 fois ($X_n = 1$), ou 4 fois ($X_n = 4$), etc...	Variable aléatoire $\rightarrow X = k$	Réel k $\in I$ <u>NB</u> : $[a; b] \subset I \subseteq \mathbb{R}$
$P(X_n = k) =$	$p^k \cdot (1 - p)^{1-k}$ $\rightarrow P(X=0) = 1 - p$ $\rightarrow P(X=1) = p$	$\binom{n}{k} p^k \cdot (1 - p)^{n-k}$ (Cf. binôme de Newton) $\rightarrow P(X_n = 0) = 1 - p$ $\rightarrow P(X_n = n) = p^n$	$\rightarrow P(a_x \leq X \leq b_x)$ $= \int_{a_x}^{b_x} f(x) \cdot dx \approx$	$\int_{a_x}^{b_x} \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{\frac{-(x-\mu)^2}{2\sigma^2}} \cdot dx$ (admis) avec : $f(x) = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{\frac{-(x-\mu)^2}{2\sigma^2}}$

	→ NB : $P(\Omega) = \sum_{k=0}^1 P(X=k) = 1$	→ NB : $P(\Omega) = \sum_{k=0}^n P(X_n=k) = (p + (1-p))^n = 1$ (binôme de Newton)		
$E(X_n) =$	p $= p.1 + (1-p).0$	$n.p$ (admis)	→ $\mu =$	$n.p$ (admis)
$V(X_n) =$	$p.(1-p)$ (Cf ci-dessous *)	$n.p.(1-p)$ (admis)	→ $\sigma^2 =$	$n.p.(1-p)$ (admis)
$\sigma(X_n) =$	$\sqrt{p.(1-p)}$	$\sqrt{n.p.(1-p)}$	→ $\sigma =$	$\sqrt{n.p.(1-p)}$

➤ Démonstration (*) : selon la loi binomiale B_p , en choisissant que $x_1 = 0$ et $x_2 = 1$:

$$V = \sum_{i=1}^{n_0} p_i \cdot (x_i - \bar{x})^2 = (1-p)(0-p)^2 + p(1-p)^2 = [(-p)^2 + p(1-p)] \cdot (1-p) = p \cdot (1-p)$$

➤ Théorème de MOIVRE-LAPLACE :

Si X est une variable aléatoire qui suit la loi binomiale $B(n; p)$, et $n \rightarrow +\infty$;

soit la variable aléatoire : $Z = \frac{X - \mu}{\sigma} = \frac{X - n.p}{\sqrt{n.p(1-p)}}$.

→ Pour tous réels a_z et b_z tels que $a_z \leq b_z$, on a alors : $\lim_{n \rightarrow +\infty} P(a_z \leq Z \leq b_z) = \int_{a_z}^{b_z} \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{t^2}{2}} \cdot dt$

• Conditions de validité de l'approximation (comme celles de la loi normale $N(\mu; \sigma^2)$) :

✓ $n \geq 30$

✓ $n.p \geq 5$

✓ $n.(1-p) \geq 5$

(p n'est alors pas trop proche de 0 ou 1)

• Ce changement de variable aléatoire permet de calculer plus facilement $P(a_x \leq X \leq b_x)$ car, quand $n \rightarrow +\infty$: $P(a_x \leq X \leq b_x) \approx \lim_{n \rightarrow +\infty} P(a_x \leq X \leq b_x) = \lim_{n \rightarrow +\infty} P(a_z \leq Z \leq b_z)$.

Le calcul est effectivement plus facile car la variable aléatoire Z suit alors la loi normale

centrée réduite $N(0; 1)$ (cf. ci-dessous), de densité de probabilité : $f(t) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{t^2}{2}}$.

On peut donc dire en résumé que **si je vérifie les conditions ci-dessus, alors je peux écrire que :**

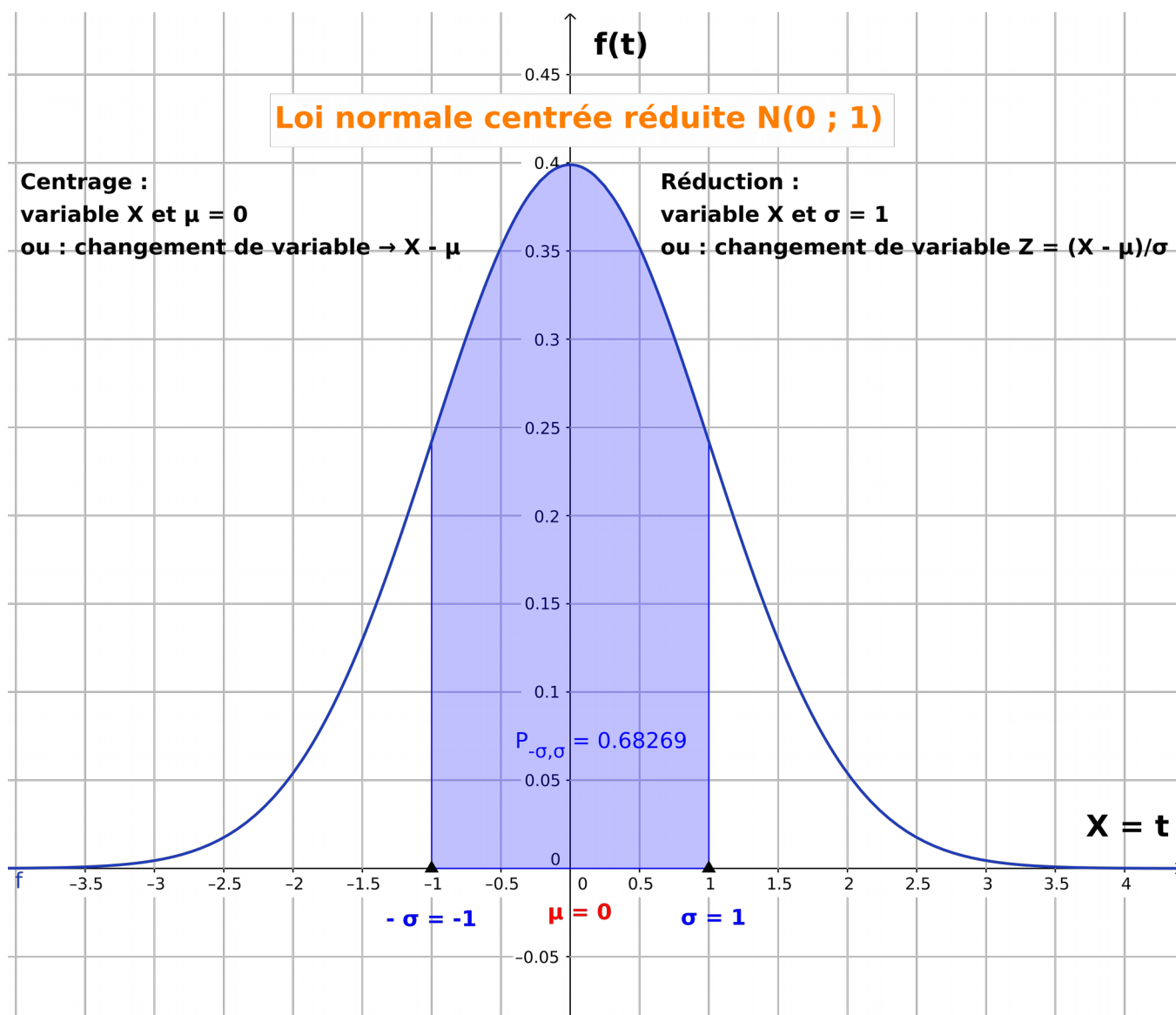
$$P(a_x \leq X \leq b_x) = P(a_z \leq Z \leq b_z) \approx \int_{a_z}^{b_z} \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{t^2}{2}} \cdot dt$$

× **Loi normale centrée réduite** notée : $N(0; 1)$:

➤ Densité de probabilité de **Laplace-Gauss** :

$$f(t) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{t^2}{2}} \quad \text{avec } t \in \mathbb{R}$$

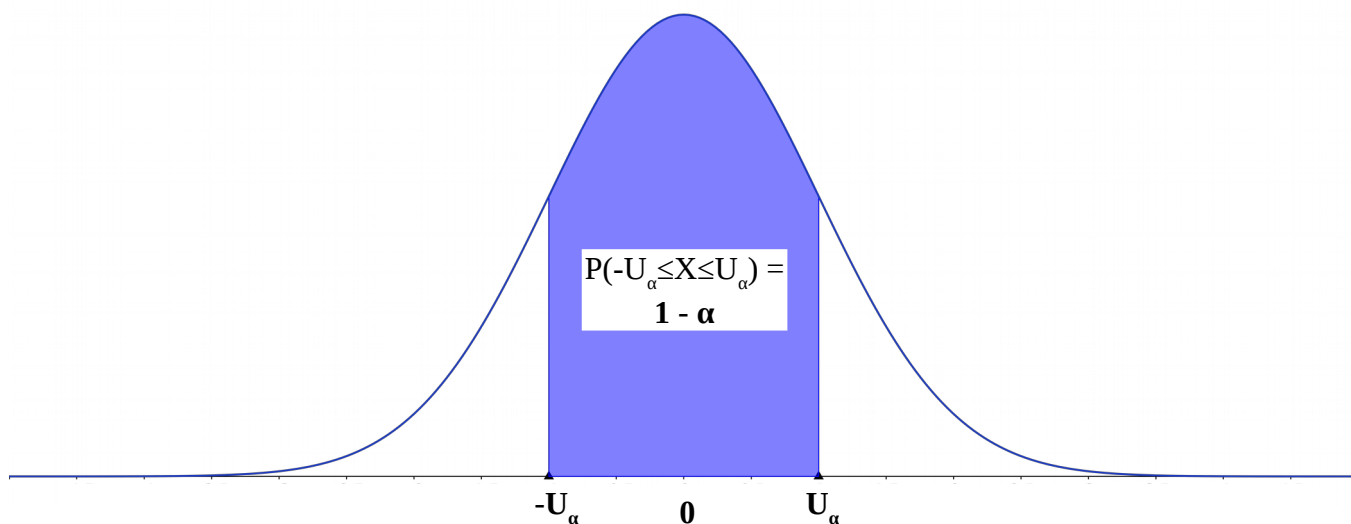
➤ f(t) est représentée par une « **courbe en cloche** » appelée aussi « **courbe de Gauss** » :



➤ Cette fonction est paire (symétrie par rapport à $(0y)$), d'où le **théorème :**

X est une variable aléatoire qui suit la loi normale centrée réduite $N(0 ; 1)$.

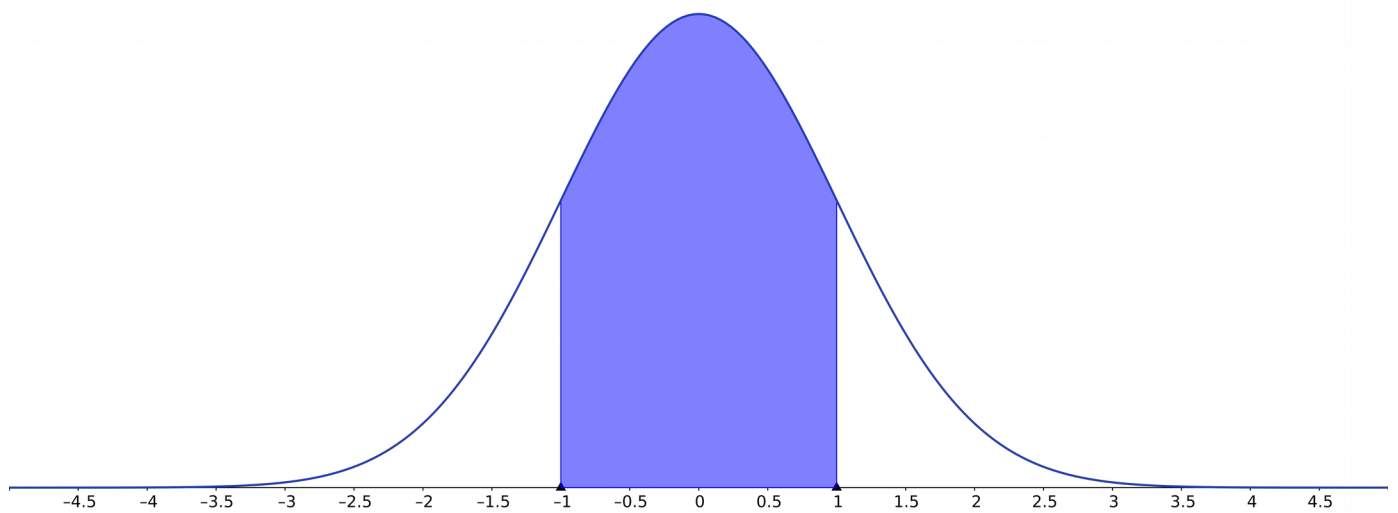
Pour tout $\alpha \in]0 ; 1[$, il existe un **unique réel** $U_\alpha > 0$ tel que : $P(-U_\alpha \leq X \leq U_\alpha) = 1 - \alpha$



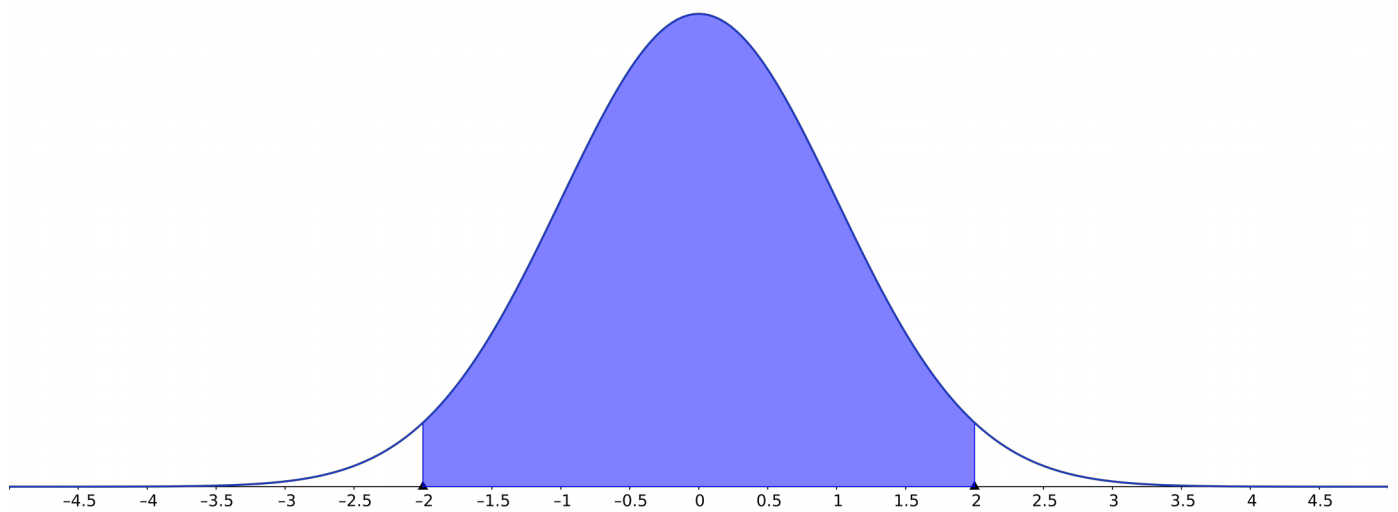
Cas particuliers : $u_{0,05} \approx 1,96$ (95%)

$u_{0,01} \approx 2,58$ (99%)

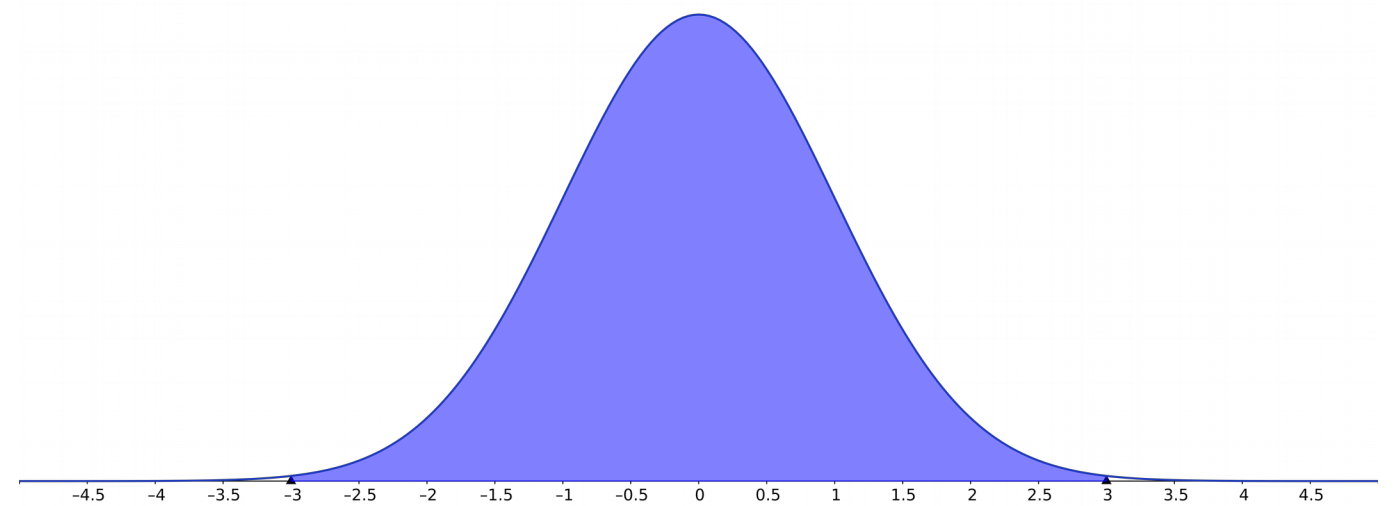
➤ $P(\mu - \sigma \leq X \leq \mu + \sigma) \approx 0,683$



➤ $P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 0,954$



➤ $P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 0,997$



FLUCTUATIONS ET ESTIMATIONS

- x **Conditions de validité générales** des approximations utilisées ci-dessous (celles de la loi normale $N(\mu ; \sigma^2)$) :

- $n \geq 30$
- $n \cdot p \geq 5$
- $n \cdot (1 - p) \geq 5$
(p n'est alors pas trop proche de 0 ou 1)

- x **Contexte d'utilisation :**

- Statistiques appliquées à la médecine, l'astronomie, la météorologie...
- Exemple : une urne contient un **très grand nombre de boules blanches et de boules noires**.
Proportion de boules blanches : p
Proportion de boules noire : $1 - p$
Tirage avec remise (\rightarrow événements indépendants) : **n boules (échantillon)**
Deux issues : blanche ($x = 1$) ou noire ($x = 0$)

Je peux alors pratiquer :

➔ ou bien **un échantillonnage** :

- **Si je connais p** , je calcule l'**intervalle de fluctuation asymptotique (centré sur p)** de la fréquence, au **seuil $1 - \alpha$** ($= 0,95$ généralement). Concrètement, je peux donc calculer de façon théorique l'**intervalle dans lequel 95 % des fréquences doivent se trouver**.

Je peux en parallèle mesurer la **fréquence de succès « apparition des boules blanches »** pour comparer.
- **Si je ne connais pas p mais peux faire une hypothèse sur sa valeur**, permettant de calculer de même un **intervalle de fluctuation asymptotique (centré sur la valeur hypothétique de p)** au **seuil $1 - \alpha$** ($= 0,95$ généralement), je vais alors pouvoir le comparer à la fréquence réelle. Je dois donc mesurer ou connaître la **fréquence de succès « apparition des boules blanches »**. Suivant la correspondance ou non entre l'intervalle et la fréquence connue, je ferai une **prise de décision** sur la **validité de l'hypothèse sur p** .

➔ ou bien **une estimation** :

Je ne connais ni ne peux faire d'hypothèse sur la valeur de p . Je dois d'abord mesurer ou connaître la **fréquence de succès « apparition des boules blanches »**. Je choisis ensuite le **niveau de confiance $1 - \alpha$** ($= 0,95$ généralement), puis je calcule à l'aide de la fréquence connue un **intervalle de confiance** (encadrant p) afin d'**estimer la valeur de p** .

➔ **NB** : **plus n augmente, plus les intervalles** de fluctuation ou de confiance **se resserrent**.

- **Variable aléatoire X_n** = k dénombrant le nombre de succès k (exemple : « tirer une boule blanche ») parmi un échantillon de n essais, et qui suit une **loi binomiale $B(n ; p)$** .
- **Variable aléatoire fréquence de succès $F_n = \frac{X_n}{n}$** pour un « **schéma de Bernoulli (?)** de paramètres n et p ».

→ **Théorème :**

Soit $\alpha \in]0 ; 1 [$.

Soit X_n une variable aléatoire qui suit une loi binomiale $B(n ; p)$.

La **probabilité** que la **fréquence** F_n prenne ses valeurs dans l'intervalle

$I_n = \left[p - u_\alpha \cdot \sqrt{\frac{p(1-p)}{n}} ; p + u_\alpha \cdot \sqrt{\frac{p(1-p)}{n}} \right]$ se rapproche de $1 - \alpha$ quand la taille de l'échantillon n devient grande.

On note : $\lim_{n \rightarrow +\infty} P(F_n \in I_n) = 1 - \alpha$

I_n est appelé « **intervalle de fluctuation asymptotique de la fréquence F_n au seuil $1 - \alpha$** ».

NB :

- en général, $\alpha = 0,05$, donc $1 - \alpha = 0,95$ (seuil de 95%) et $u_\alpha \approx 1,96$.
- « **n devient grand** » revient simplement à vérifier les conditions sur n et p validant la loi normale
- Ici, on ne cherche pas à calculer F_n mais à trouver un intervalle I_n dans lequel elle a $1 - \alpha$ chances de se trouver.

Exemple : « Un élève a manqué les cours et n'a rien révisé à l'issue d'un chapitre en histoire. L'interrogation comportera heureusement un QCM de 40 questions ! Mais quelles sont ses « chances » de bien répondre en se reposant sur le « pur » hasard, sachant que chaque question comporte 1 seule bonne réponse et 4 propositions de réponse ? »

✓ Probabilité de bien répondre à une question : $p = \frac{1}{4} = 0,25$.

✓ **Conditions de validité des approximations :**

- échantillon de taille $n = 40 \geq 30$
- $\mu = n.p = 40 \times \frac{1}{4} = 10 \geq 5$
- $n.(1 - p) = 40 \times \frac{3}{4} = 30 \geq 5$

Les conditions sont vérifiées (nous pouvons donc utiliser les théorèmes de ce chapitre)

✓ Cherchons maintenant, non pas la fréquence F_n , mais l'intervalle I_n dans lequel F_n a une probabilité $1 - \alpha = 1 - 0,05 = 0,95$ de se trouver, soit 95 % de chances :

$$I_{40} \approx \left[0,25 - 1,96 \cdot \sqrt{\frac{0,25 \times 0,75}{40}} ; 0,25 + 1,96 \cdot \sqrt{\frac{0,25 \times 0,75}{40}} \right]$$

$$I_{40} \approx [0,116 ; 0,384]$$

✓ **Conclusion :**

$$P(F_{40} \in I_{40} = [0,116 ; 0,384]) \approx \lim_{n \rightarrow +\infty} P(F_n \in I_n) = 1 - \alpha = 0,95$$

En clair, il a 95 % de chances d'avoir entre 11,6 % et 38,4 % de bonnes réponses...

→ **Propriété (Règle de prise de décision) :**

Soit f la fréquence du caractère étudié d'un échantillon de **taille n** .

Soit l'**hypothèse** : « **la proportion de ce caractère dans la population est p** ».

Soit I_n l'intervalle de fluctuation asymptotique au **seuil 0,95**.

- Si $f \in I_n$, alors on **accepte l'hypothèse** faite sur la proportion p .
- Si $f \notin I_n$, alors on **rejette l'hypothèse** faite sur la proportion p .

NB :

- On parle de « prise de décision » parce que l'on est en mesure de faire une hypothèse sur la proportion p , ce qui permet de calculer I_n .
- On note ici la fréquence « f » (et non « F_n ») car f peut être obtenue par un étude d'un échantillon restreint sur l'échantillon total de taille n (prélèvements de contrôle sur une chaîne de production, ...).

Exemple : Reprenons l'exemple précédent, mais supposons maintenant que le QCM de 40 questions comporte par question aléatoirement 3, 4 ou 5 propositions de réponse. Sachant que cette variété dans les proposition semble être plutôt homogène, on formule l'**hypothèse suivante** : $p = \frac{1}{4} = 0,25$. Or nous savons cette fois que l'élève a finalement obtenu une fréquence $f = \frac{2}{5}$ de bonnes réponses, soit 2 fois sur 5 en moyenne.

✓ **Conditions de validité des approximations :** vérifiées comme précédemment.

✓ La fréquence f est connue, mais cherchons l'intervalle I_n dans lequel f a une probabilité $1 - \alpha = 1 - 0,05 = 0,95$ de se trouver, soit 95 % de chances :

$$I_{40} \approx \left[0,25 - 1,96 \cdot \sqrt{\frac{0,25 \times 0,75}{40}} ; 0,25 + 1,96 \cdot \sqrt{\frac{0,25 \times 0,75}{40}} \right]$$
$$I_{40} \approx [0,116 ; 0,384]$$

✓ **Prise de décision :**

$$f = 0,4 \notin I_n \approx [0,116 ; 0,384]$$

L'hypothèse $p = 0,25$ est donc **rejetée**.



x **Estimation :**

→ **Théorème :**

Soit X_n une variable aléatoire qui suit une loi binomiale **$B(n ; p)$** .

$F_n = \frac{X_n}{n}$ est la fréquence associée à X_n .

Pour **n suffisamment grand**, p appartient à l'intervalle : $J_n = \left[F_n - \frac{1}{\sqrt{n}} ; F_n + \frac{1}{\sqrt{n}} \right]$ avec une probabilité supérieure ou égale à **0,95**.

Remarque :

- L'intervalle J_n se rapproche de l'intervalle de fluctuation I_n , à la différence près qu'on considère que :
 - ✓ $U_\alpha = U_{0,05} \approx 2$
 - ✓ $p(1-p) \approx 0,25$ (dans les conditions évoquées ci-dessus)
 - ✓ la proportion p (issue des probabilités) est remplacée par la fréquence F_n (issue des statistiques).
- Il est plus facile de s'appuyer sur la définition ci-dessous pour faire une estimation →

→ **Définition :**

Soit **f** une **fréquence observée** du caractère étudié sur un **échantillon de taille n** (dans le cadre de la loi binomiale $B(n ; p)$).

On appelle « **intervalle de confiance** » de la proportion p au **niveau de confiance 0,95**

l'intervalle : $J_n = \left[f - \frac{1}{\sqrt{n}} ; f + \frac{1}{\sqrt{n}} \right]$.

Remarque :

- p étant parfaitement inconnu ici, il n'est pas possible de vérifier les conditions énoncées sur n et p rappelées précédemment. Cependant, il faudra les vérifier en remplaçant la proportion p par la fréquence observée f :
 - $n \geq 30$
 - $n \cdot f \geq 5$
 - $n \cdot (1 - f) \geq 5$
- Un niveau de confiance 0,95 signifie qu'on affirme à juste titre que $p \in J_n$ dans 95 % des cas.
- Il n'est pas possible de situer p dans J_n , donc on ne peut affirmer ici que J_n est centré sur p .

Exemple :

« Un institut de sondage interroge 1 052 personnes entre les deux tours de l'élection présidentielle sur leur intention de vote.
614 déclarent avoir l'intention de voter pour Michelle Coluchi.
En supposant que les votes seront conformes aux intentions, la candidate a-t-elle raison de croire qu'elle sera élue ? »

- ✓ Fréquence de succès observée dans les intentions : $f = \frac{614}{1052} \approx 0,584$
- ✓ Échantillon de taille $n \geq 1052$ (65 millions d'habitants en France, dont plus de la moitié ont droit de vote ; mais il y a aussi les abstentions...).
- ✓ La proportion p finale des votes favorables n'est pas connue en revanche.
Je vais donc calculer son **intervalle de confiance J_n au seuil 0,95**.
- ✓ Mais je dois vérifier d'abord le **respect des conditions sur n et f** :
 - $n \geq 1052 \geq 30$
 - $n \cdot f \geq \frac{1052 \times 614}{1052} = 614 \geq 5$
 - $n \cdot (1-f) \geq \frac{1052 \times 438}{1052} = 438 \geq 5 \rightarrow$ Conditions **vérifiées**.
- ✓ D'où : $J_n =_{1052} = \left[\frac{614}{1052} - \frac{1}{\sqrt{1052}} ; \frac{614}{1052} + \frac{1}{\sqrt{1052}} \right] \approx [0,553 ; 0,614]$
 $p \in J_n =_{1052}$ donc p a 95 % de chances d'être supérieur à 0,50 (d'autant que J_n se resserre quand n augmente !) : la candidate peut croire en ses chances !

Références en ligne :

- **Yvan MONKA – Math-et-tiques :**
<https://www.youtube.com/watch?v=4Y12jMMYyVM>
- **jaicompris Maths :**
<https://www.youtube.com/watch?v=TnkHUF1Kxao>
- (...)