

Éléments de compréhension des statistiques

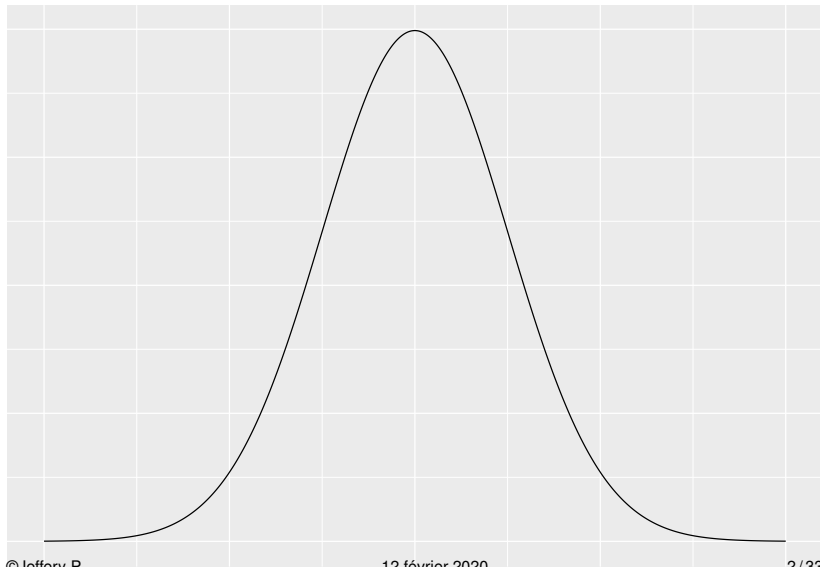
Jeffery P.

Doctorant au Laboratoire des Sciences du Numérique de Nantes (LS2N)

2019

Quelques mots sur la loi normale

S'il y a bien une loi populaire en statistique, il s'agit de la loi normale. . .la célèbre courbe en cloche !

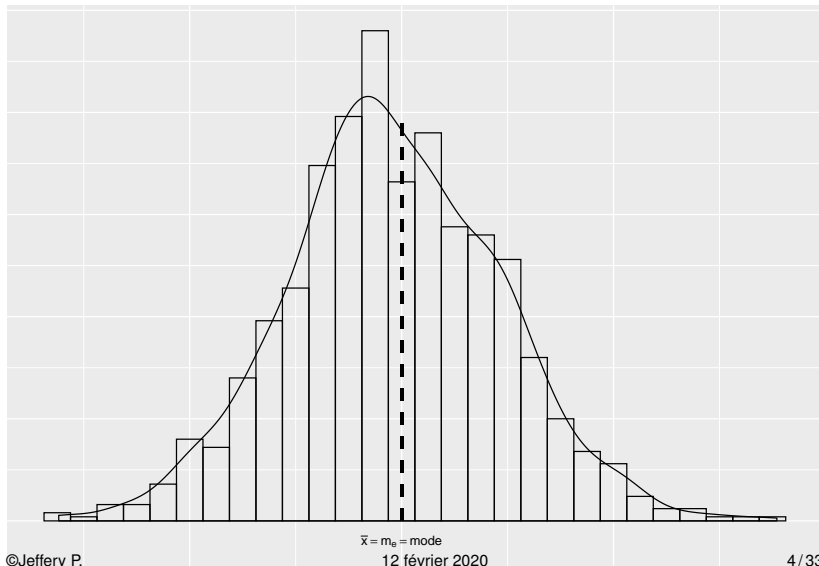


Caractéristiques de la loi normale

- ▶ La loi normale est une loi **symétrique, centrée autour de sa moyenne**
- ▶ La symétrie de la distribution implique que **la médiane est égale à la moyenne**
- ▶ C'est une loi unimodale, **son mode est égale à la moyenne**

Reconnaître une loi normale

En pratique, on peut supposer une distribution normale d'après l'histogramme, ou le diagramme en barre



Quelques mots sur la loi normale

- ▶ La loi normale est définie par deux paramètres :
- 1. Sa moyenne souvent notée μ (où *mean*)
- 2. Son écart-type notée σ (où *sd* pour “Standard Deviation”)

Sa densité de probabilité est la suivante :

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \times \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

avec une variable X qui suit une loi normale de moyenne μ et d'écart-type σ

On note $X \sim \mathcal{N}(\mu, \sigma^2)$

Illustration de la loi normale

Changer la moyenne d'une loi normale revient à translater la distribution vers la droite ou la gauche

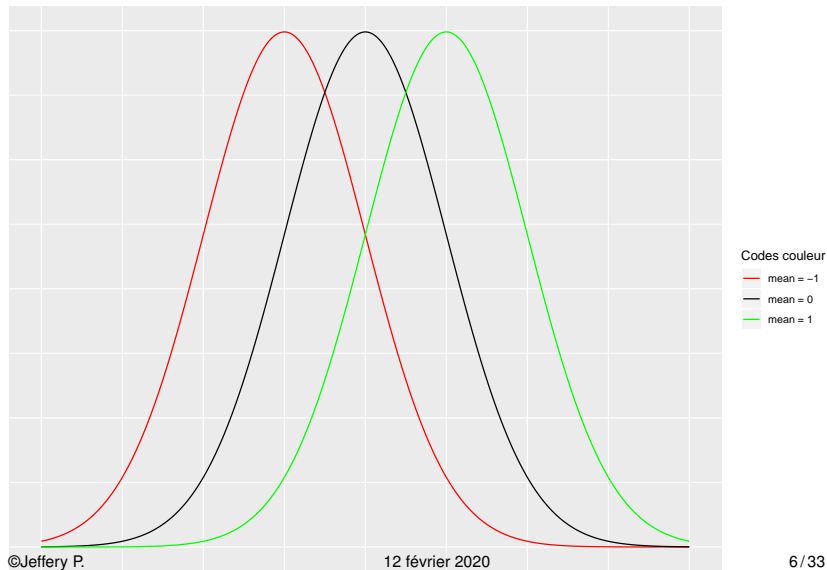
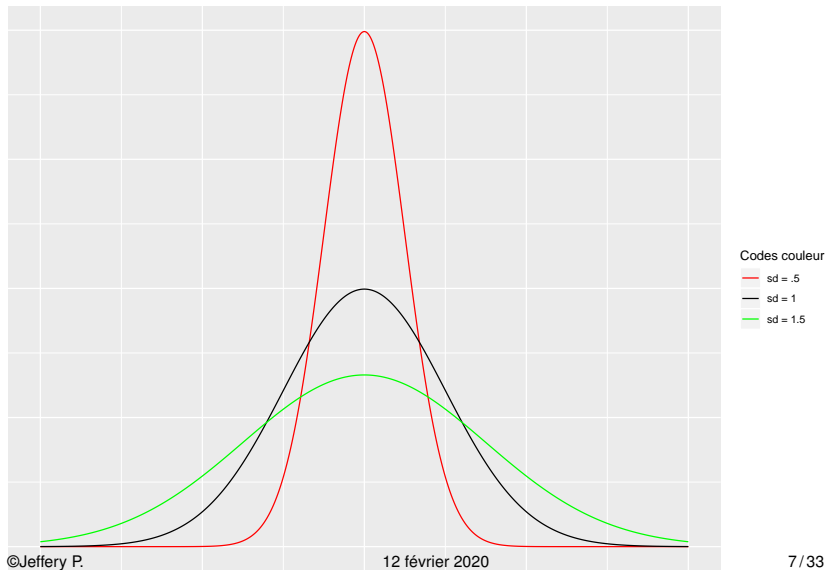


Illustration de la loi normale

Changer l'écart-type d'une loi normale revient à aplatir ou resserrer sa distribution autour de sa moyenne



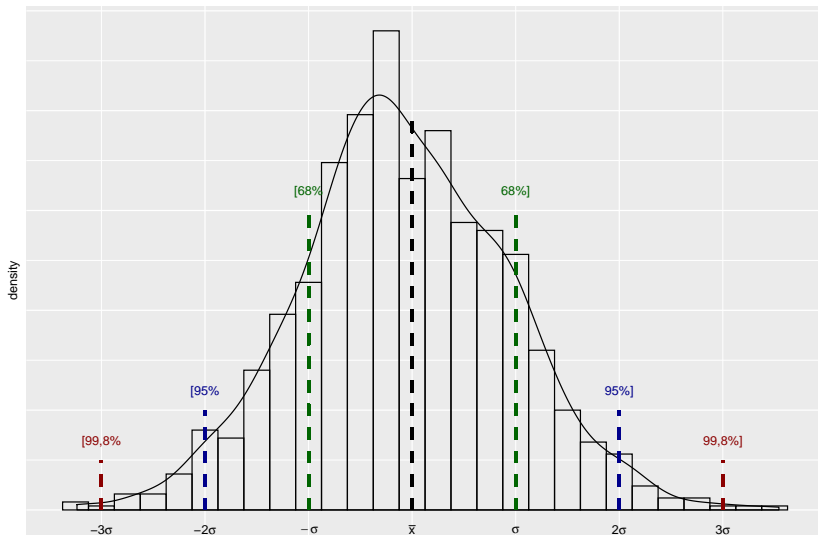
Propriété intéressante

Si l'on dispose d'observations d'une distribution $\mathcal{N}(\mu, \sigma^2)$ alors :

- ▶ 68% des observations sont comprises dans l'intervalle $[\mu - \sigma; \mu + \sigma]$
- ▶ 95% des observations sont comprises dans l'intervalle $[\mu - 2\sigma; \mu + 2\sigma]$
- ▶ 99,8% des observations sont comprises dans l'intervalle $[\mu - 3\sigma; \mu + 3\sigma]$

Estimation graphique de μ et σ

- En pratique on estime μ en prenant le centre de la distribution, et on déduit σ en estimant l'intervalle vert



Calcul de probabilité : généralité

La probabilité pour qu'une variable aléatoire X soit inférieure à une quelconque valeur x s'écrit $P(X \leq x)$

- ▶ **Une probabilité est toujours positive et inférieure à 1**

Remarque

- ▶ Pour une variable aléatoire réelle on a :
 $P(-\infty < X < +\infty) = 1$
- ▶ Pour une variable continue, on a $P(X = x) = 0$, d'où :

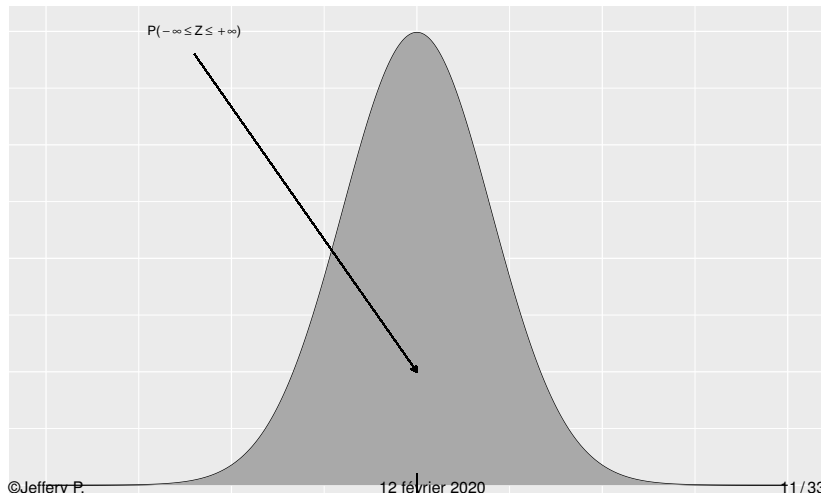
$$P(X \leq x) = P(X < x) \text{ où encore } P(X \geq x) = P(X > x)$$

Ce résultat s'explique avec un peu de théorie mathématique que l'on ne détaillera pas ici, mais il peut être utile d'avoir ces propriétés en tête pour les exercices. . .

Calcul de probabilité : lien avec l'aire sous la courbe de densité

Cas 0 : soit $Z \sim \mathcal{N}(0, 1)$, l'aire sous la courbe est égale à 1

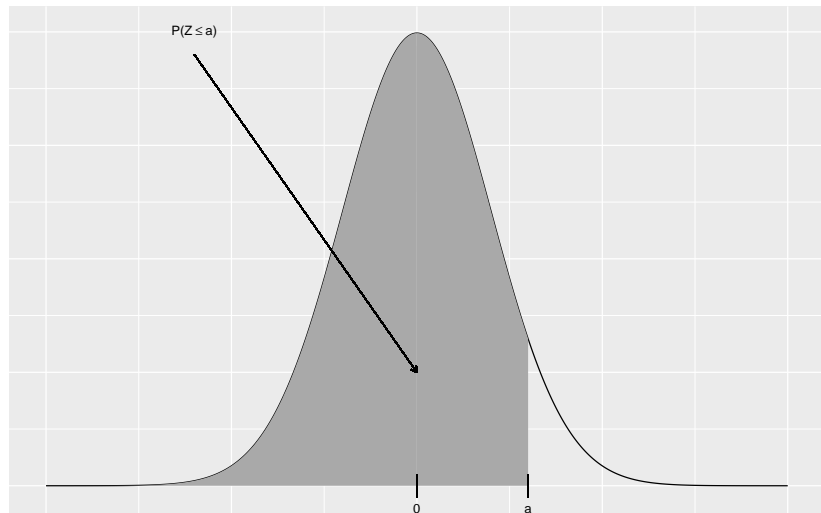
$$P(-\infty < Z < +\infty) = 1$$



Calcul de probabilité : valeurs de la table

Cas I : soient $Z \sim \mathcal{N}(0, 1)$ et a un nombre réel **positif**

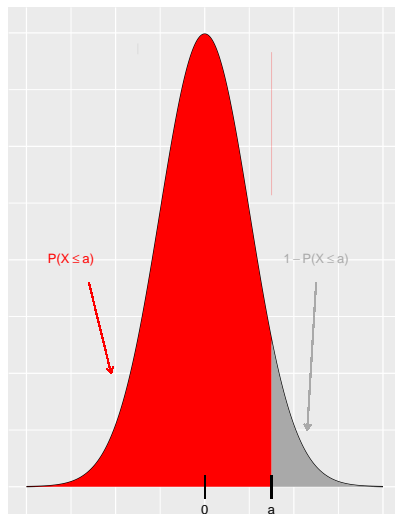
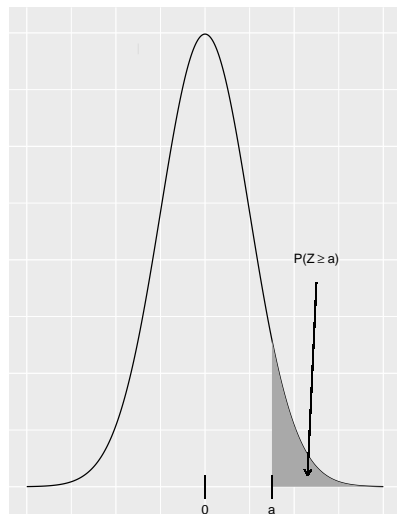
$P(Z \leq a) \longrightarrow$ le résultat se trouve dans la table !



Calcul de probabilité avec la loi normale

Cas II : soient $Z \sim \mathcal{N}(0, 1)$ et a un nombre réel **positif**

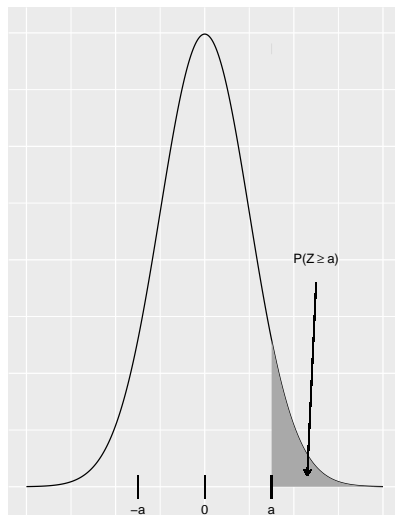
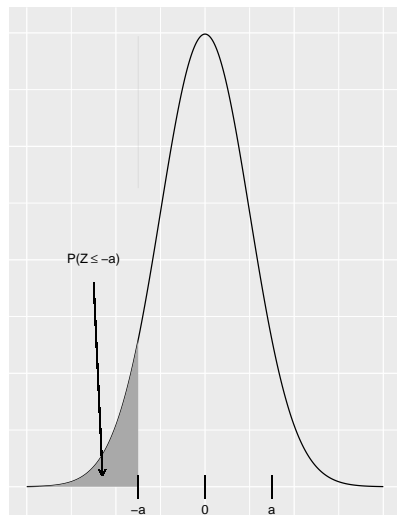
$P(Z \geq a) = 1 - P(Z \leq a) \rightarrow$ on se ramène au **cas I**



Calcul de probabilité avec la loi normale

Cas III : soient $Z \sim \mathcal{N}(0, 1)$ et a un nombre réel **positif**

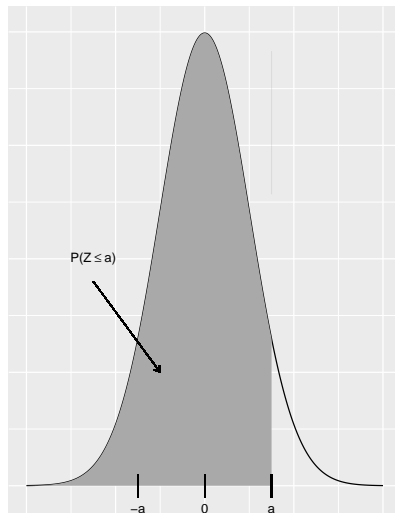
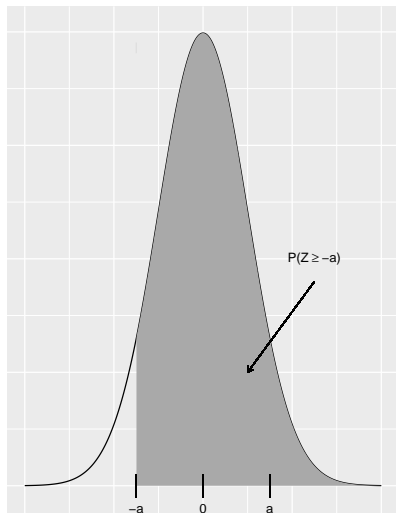
$P(Z \leq -a) = P(Z \geq a) \rightarrow$ on se ramène au **cas II**



Calcul de probabilité avec la loi normale

Cas IV : soient $Z \sim \mathcal{N}(0, 1)$ et a un nombre réel **positif**

$P(Z \geq -a) = P(Z \leq a) \longrightarrow$ on se ramène au **cas I**



Résumé








I	$\mathbb{P}(X \leq a)$		\Rightarrow table
II	$\mathbb{P}(X \geq a)$	 $= 1 -$ 	\Rightarrow cas I
III	$\mathbb{P}(X \leq -a)$	 $=$ 	\Rightarrow cas II
IV	$\mathbb{P}(X \geq -a)$	 $=$ 	\Rightarrow cas I

FIGURE 1 – Calcul de probabilité (*Extrait du cours de M. Gérin, Paris Ouest 2012-2013*)

Formule de calcul avec une intervalle quelconque

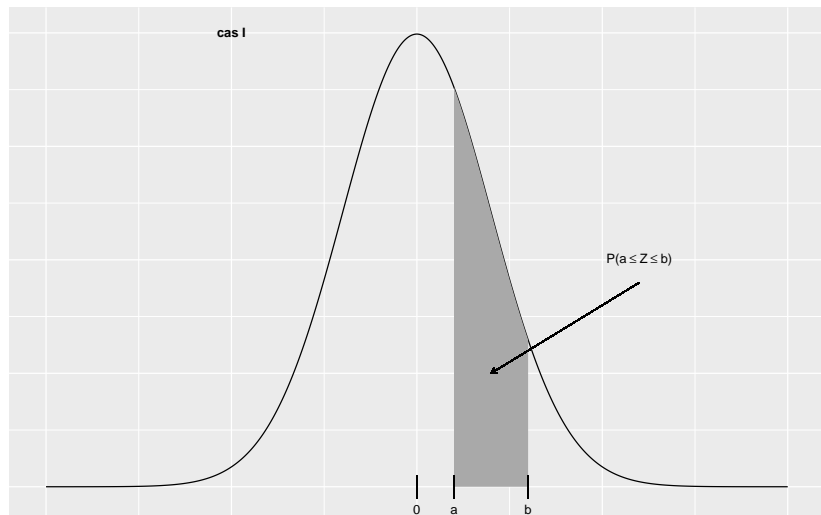
Soient $Z \sim \mathcal{N}(0, 1)$ et u, v deux nombres réels avec $u \leq v$

$$P(u \leq Z \leq v) = P(Z \leq v) - P(Z \leq u)$$

Calcul de probabilité avec la loi normale

Cas V.1 : soient $Z \sim \mathcal{N}(0, 1)$ et $0 \leq a \leq b$

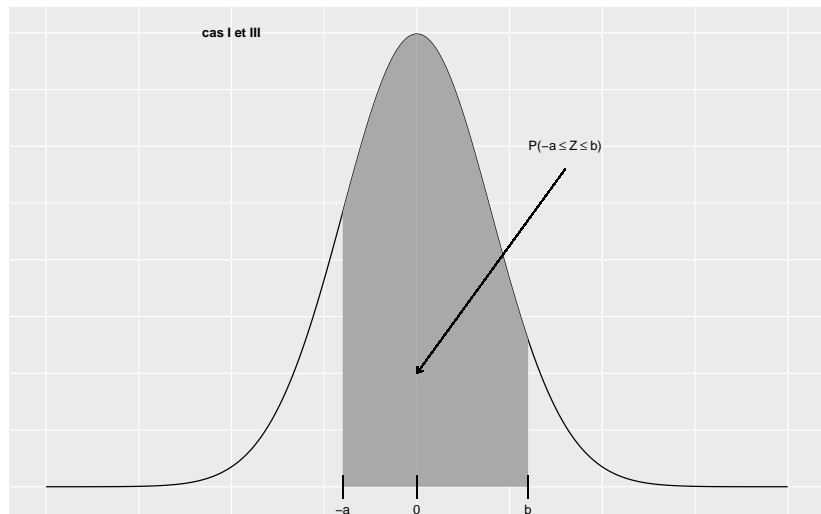
$$P(a \leq Z \leq b) = P(Z \leq b) - P(Z \leq a)$$



Calcul de probabilité avec la loi normale

Cas V.2 : soient $Z \sim \text{Norm}(0, 1)$ et $0 \leq a \leq b$

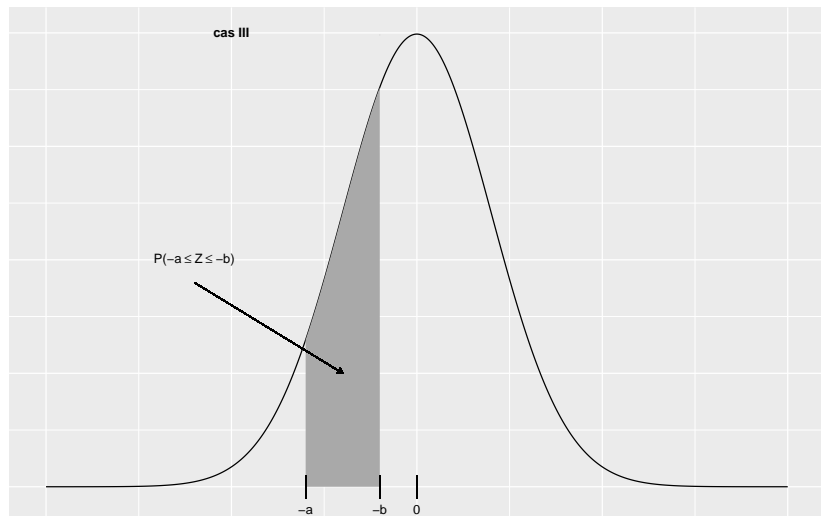
$$P(-a \leq Z \leq b) = P(Z \leq b) - P(Z \leq -a) = P(Z \leq b) + P(Z \leq a) - 1$$



Calcul de probabilité avec la loi normale

Cas V.3 : soient $Z \sim \mathcal{N}(0, 1)$ et $-a \leq -b \leq 0$

$$P(-a \leq Z \leq -b) = P(Z \leq -b) - P(Z \leq -a) = P(Z \leq a) - P(Z \leq b)$$



Transformation de la loi normale

Toute transformation *affine* d'une loi normale est encore une loi normale i.e., quelque soit les nombre réels a et $b \neq 0$ on a :

$$X \sim \mathcal{N}(\mu, \sigma^2) \Rightarrow aX + b \sim \mathcal{N}(\mu + b, a^2 \times \sigma^2)$$

Par conséquent :

- ▶ si $X \sim \mathcal{N}(\mu, \sigma^2)$ alors $X - \mu \sim \mathcal{N}(0, \sigma^2)$ (**centrage**)
- ▶ si $X \sim \mathcal{N}(\mu, \sigma^2)$ alors $\frac{X}{\sigma} \sim \mathcal{N}(\mu, 1)$ (**réduction**)

En pratique, on fait souvent les deux en même temps :

- ▶ si $X \sim \mathcal{N}(\mu, \sigma^2)$ alors $\frac{X-\mu}{\sigma} \sim \mathcal{N}(0, 1)$
(**normalisation=centrage+réduction**)

Calcul de probabilité : en pratique

Si la variable X ne suit pas une loi normale centrée réduite, on peut toujours se ramener aux cas précédents, par exemple :

$$P(X \leq x) = P\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) = P\left(Z \leq \frac{x - \mu}{\sigma}\right)$$

Où $Z \sim \mathcal{N}(0, 1)$

→ Rappel : il s'agit de la **normalisation**, cette petite transformation est très utile et beaucoup utilisée !

Cas pratique : calcul de probabilité

68% des observations sont comprises dans l'intervalle $[\mu - \sigma; \mu + \sigma]$. . . Est-ce bien vrai ?

$$\begin{aligned} P(\mu - \sigma \leq X \leq \mu + \sigma) &= P(-\sigma \leq X - \mu \leq \sigma), \quad \text{on centre} \\ &= P(-1 \leq \frac{X - \mu}{\sigma} \leq 1), \quad \text{on réduit} \\ &= P(-1 \leq Z \leq 1) \end{aligned}$$

Où $Z \sim \mathcal{N}(0, 1)$

Il faut maintenant trouver la valeur de cette probabilité sur les tables. . .

Cas pratique

$$\begin{aligned}P(\mu - \sigma \leq X \leq \mu + \sigma) &= P(-1 \leq Z \leq 1) \\&= P(Z \leq 1) - P(Z \leq -1) \\&= P(Z \leq 1) - (1 - P(Z \leq 1)) \\&= 2P(Z \leq 1) - 1\end{aligned}$$

Où $Z \sim \mathcal{N}(0, 1)$

Remarque

Plus généralement, on a la formule suivante quelque soit k nombre entier positif :

$$P(\mu - k \times \sigma \leq X \leq \mu + k \times \sigma) = P(-k \leq Z \leq k)$$

Calcul des quantiles

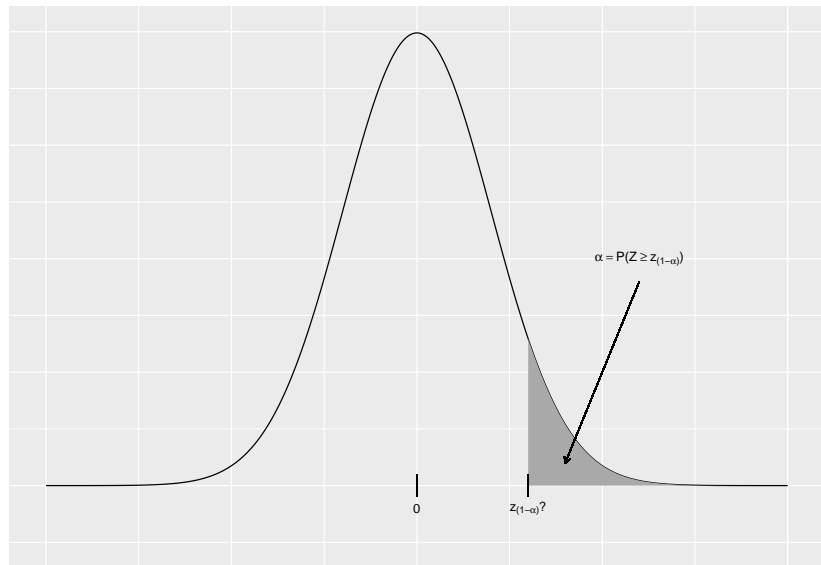
Dans certains cas, on ne cherche pas à calculer une probabilité pour un b donné (e.g., $P(Z \leq b)$) mais on cherche b telle que la probabilité soit égale à une certaine valeur :

On veut trouver b tel que $P(X \geq b) = \alpha$ où $P(|X| \geq b) = \alpha$

- Dans ce cas, on s'assure que la probabilité concerne une variable suivant une **loi normale centrée réduite**, puis on regarde dans une des tables

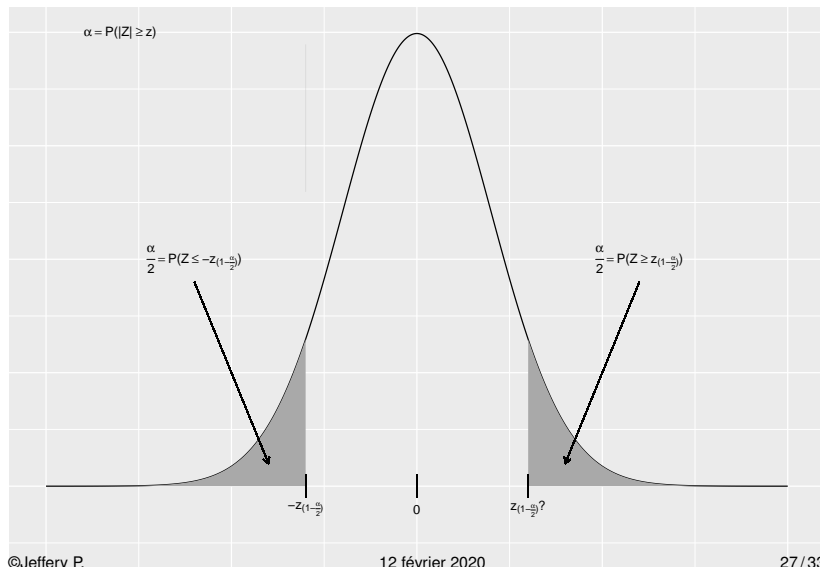
Calcul des quantiles : table 1

On recherche une valeur z telle que $P(Z > z)$ avec $Z \sim \mathcal{N}(0, 1)$:



Calcul des quantiles : table 2

On recherche une valeur z telle que $P(|Z| > z)$ avec $Z \sim \mathcal{N}(0, 1)$:



Calcul des quantiles

Attention, cette problématique est faussement facile. . . dans la majorité des cas elle nécessite une bonne maîtrise du calcul des probabilités de la loi normale !

Cas fréquents

- ▶ On ne demande pas de trouver z tel que $P(Z \geq z) = \alpha$ mais plutôt $P(Z \leq z) = 1 - \alpha$.
- ▶ On ne demande pas de trouver z tel que $P(|Z| \geq z) = \alpha$ mais plutôt $P(|Z| \leq z) = 1 - \alpha$.
- ▶ La probabilité porte sur une variable $X \sim \mathcal{N}(\mu, \sigma^2)$. Il faut bien penser à centrer et réduire l'évènement (ce qui se trouve dans la probabilité) e.g., $P(X > x) = P(Z > \frac{x-\mu}{\sigma})$ où $Z \sim \mathcal{N}(0, 1)$

Calcul des quantiles

Exemple général (théorique mais utile !)

Soit $X \sim \mathcal{N}(\mu, \sigma)$, on suppose $\alpha \in]0; 0.5[$ connu.

On cherche b tel que $P(X > b) = \alpha$

Voici les étapes qui nous permettent de trouver la solution :

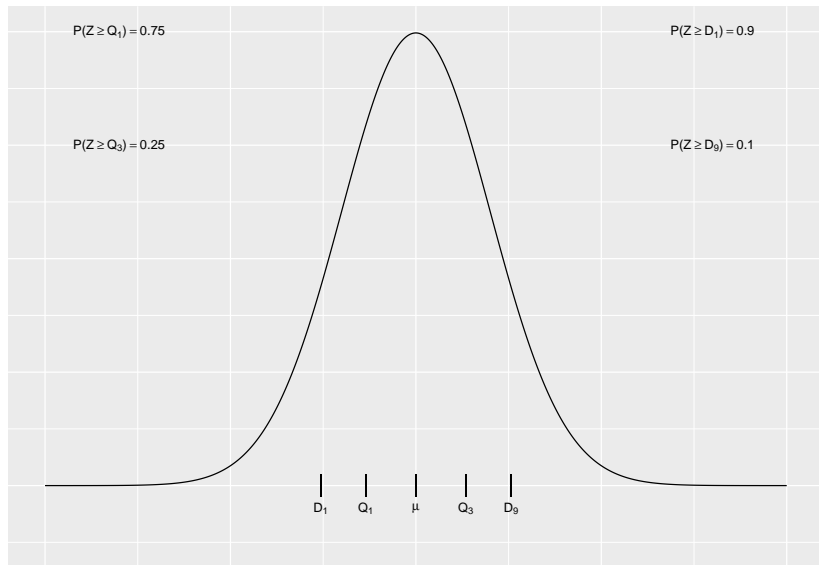
$$\begin{aligned}P(X > b) &= \alpha \\ \Leftrightarrow P(Z > \frac{b-\mu}{\sigma}) &= \alpha \\ \Leftrightarrow P(Z \leq \frac{b-\mu}{\sigma}) &= 1 - \alpha \\ \Rightarrow \frac{b-\mu}{\sigma} &= z_{(1-\alpha)} \\ \text{Donc } b &= \mu + \sigma \times z_{(1-\alpha)}\end{aligned}$$

Avec :

- ▶ $Z \sim \mathcal{N}(0, 1)$
- ▶ $z_{(1-\alpha)}$ est le **quantile** d'ordre $(1 - \alpha)$ de la loi normale centrée réduite

Calcul des quantiles

► Rappel sur la position des déciles et quantiles extrêmes



Calcul des quantiles

Soit $X \sim \mathcal{N}(\mu, \sigma)$

Détermination des déciles et quartiles extrêmes

En utilisant la formule précédente, il vient :

$$\begin{aligned} P(X > Q_3) &= 0.25 \\ \Rightarrow Q_3 &= \mu + \sigma \times z_{(0.75)} \end{aligned}$$

$$\begin{aligned} P(X > D_9) &= 0.1 \\ \Rightarrow D_9 &= \mu + \sigma \times z_{(0.9)} \end{aligned}$$

Pour rappel :

- ▶ Q_3 : troisième quartile
- ▶ D_9 : neuvième décile

Calcul des quantiles

Soit $X \sim \mathcal{N}(\mu, \sigma)$

Détermination des déciles et quartiles extrêmes

Pour D_1 et Q_1 , le calcul est un peu **sioux**. Par exemple pour D_1 , on a :

$$\begin{aligned} P(X > D_1) &= 0.9 \\ \Leftrightarrow P(Z > \frac{D_1 - \mu}{\sigma}) &= 0.9 \\ \Leftrightarrow P(Z \leq \frac{D_1 - \mu}{\sigma}) &= 0.1 \\ \Rightarrow \frac{D_1 - \mu}{\sigma} &= z_{(0.1)} \quad \text{or } z_{(0.1)} = -z_{(1-0.1)} = -z_{(0.9)} \\ \text{Donc } D_1 &= \mu - \sigma \times z_{(0.9)} \end{aligned}$$

Calcul des quantiles

Soit $X \sim \mathcal{N}(\mu, \sigma)$

Détermination des déciles et quartiles extrêmes

On peut utiliser les formules suivantes pour D_1 et Q_1 :

$$Q_1 = \mu - \sigma \times z_{(0.75)}$$

$$D_1 = \mu - \sigma \times z_{(0.9)}$$

Avec :

- ▶ Q_1 : premier quartile i.e., $P(X \leq Q_1) = 0.25$
- ▶ D_1 : premier décile i.e., $P(X \leq D_1) = 0.1$