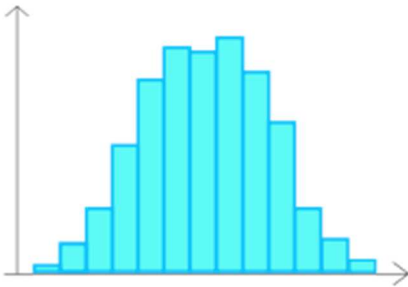


Rappel de quelques éléments simples de statistique

Imaginons que l'on s'intéresse à la taille des hommes de 40 ans en France. On aura alors une population de plusieurs centaines de milliers d'individus ayant cet âge. Si l'on trace un histogramme de la répartition de cette population par taille, on fera concrètement le schéma suivant : en ordonné, c'est-à-dire sur l'axe vertical, le nombre d'hommes mesurant une taille donnée par exemple 1,73 m (+/- 2 cm), ou, ce qui revient au même, le pourcentage des hommes de cette taille dans la population étudiée (ce qu'on peut appeler la fréquence de cette taille) et en abscisse, c'est-à-dire sur l'axe horizontal, les différentes tailles retenues, par tranches. On obtiendra un graphique du type :



Or, un résultat très intéressant de la statistique est résumé par la **loi des grands nombres** : quand on a affaire à une très grande population, la répartition de cette population suivant une variable donnée continue (comme ici la taille) donnera une répartition (en statistique, on parle de distribution) qui prendra la forme d'une courbe en cloche.

La figure 4.8 ci-dessous présente plusieurs courbes en cloche de ce type, également appelées Gaussienne ou **loi normale**.

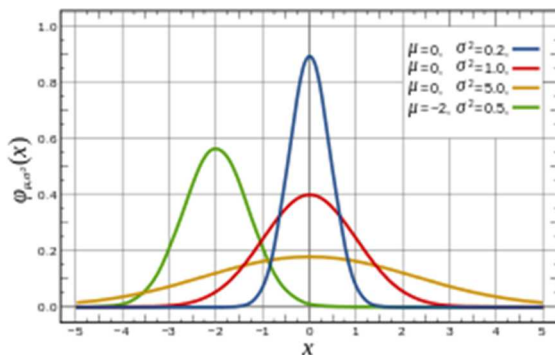


Fig.4.8 – Loi normale

Il est pour nous intéressant de se pencher sur cette loi normale pour deux raisons :

1- Un processus industriel qui se reproduit à l'identique, finit, dans la durée, par générer de grandes quantités, donc des mesures sur une grande population, et donc des quantités qui s'apparentent à des grands nombres... Les distributions que l'on peut construire à partir des mesures observées sur certains paramètres clés d'un processus peuvent ainsi être modélisées par la loi normale. C'est très utile du fait de ce qui suit.

2- Une loi normale peut être caractérisée assez simplement. C'est tout l'intérêt. D'une part, il y a la moyenne de la variable mesurée. (Par exemple 1,75m de taille moyenne pour les hommes adultes en France lors du recensement de 2007). Et puis il y a une mesure de l'aplatissement relatif de la courbe qui décrit la distribution. Sur la figure 4.8, la courbe jaune est très « aplatie », la courbe rouge un peu moins, la courbe verte encore moins et la courbe bleue est plutôt « pointue » : ceci signifie que la population représentée par la courbe jaune est très « dispersée », c'est-à-dire par exemple qu'il y aura peu de briques de lait avec la bonne quantité promise (1litre) et une part importante des briques avec soit trop, soit trop peu de lait. Cela revient à dire que le processus de remplissage des briques de lait, tel qu'il est organisé, ne permet pas d'obtenir ce qui est promis sur l'étiquette, à savoir un litre de lait, même si c'est implicitement « à peu de choses près ». A l'inverse, la population représentée par la courbe bleue, disons celle des sacs de ciment contrôlés en sortie de ligne, est pour l'essentiel proche de la cible de 35kg. Rares sont les « individus », c'est-à-dire les sacs contrôlés, qui sont « loin » de la cible. A ce titre, le processus est satisfaisant.

Pour caractériser l'aplatissement de la courbe vers le bas, ou au contraire, ce qui sera préférable, une forme pointue concentrée autour de la moyenne, on utilise en statistiques un concept appelé « l'écart type ». L'écart type mesure la dispersion des données. Il est traditionnellement symbolisé par la lettre grecque sigma : σ . La courbe jaune, dispersée, a un sigma élevé, la courbe bleue, pointue, a un sigma faible.

Nous ne rentrons pas tout de suite ici dans la définition mathématique formelle du σ . Nous ne nous préoccupons pas non plus du comment il se calcule, autrement que pour rappeler qu'il est calculé à partir des données mesurées. (Voir l'annexe ci-dessous). Mais arrive alors le point essentiel.

L'intérêt de la loi normale et du concept d'écart-type est que la forme connue de cette distribution permet de dire quelle est la part de la population qui est dans un intervalle de 2 ou 4 ou 6 ou n écarts-types autour de la moyenne m.

Ainsi, entre $m - \sigma$ et $m + \sigma$, il y a 68% des individus de la population (c'est donc 68% des individus qui sont concentrés dans un intervalle large de 2σ autour de la moyenne) ; entre $m - 2\sigma$ et $m + 2\sigma$, soit dans un intervalle large de 4σ , c'est 95%. Et entre $m - 3\sigma$ et $m + 3\sigma$, soit dans un intervalle de 6σ , il y a 99,7% de la population, ce qui signifie, en creux, que 3 individus pour mille restent encore hors de cet intervalle. C'est dire que σ devient l'unité de compte de l'écart à la moyenne... D'où son nom, l'écart-type.

Pour mémoire, plutôt que de parler d'écart type, on peut aussi parler de variance (on sent bien l'idée de variation / dispersion véhiculée par le mot) : la variance se trouve être le carré de l'écart type : σ^2

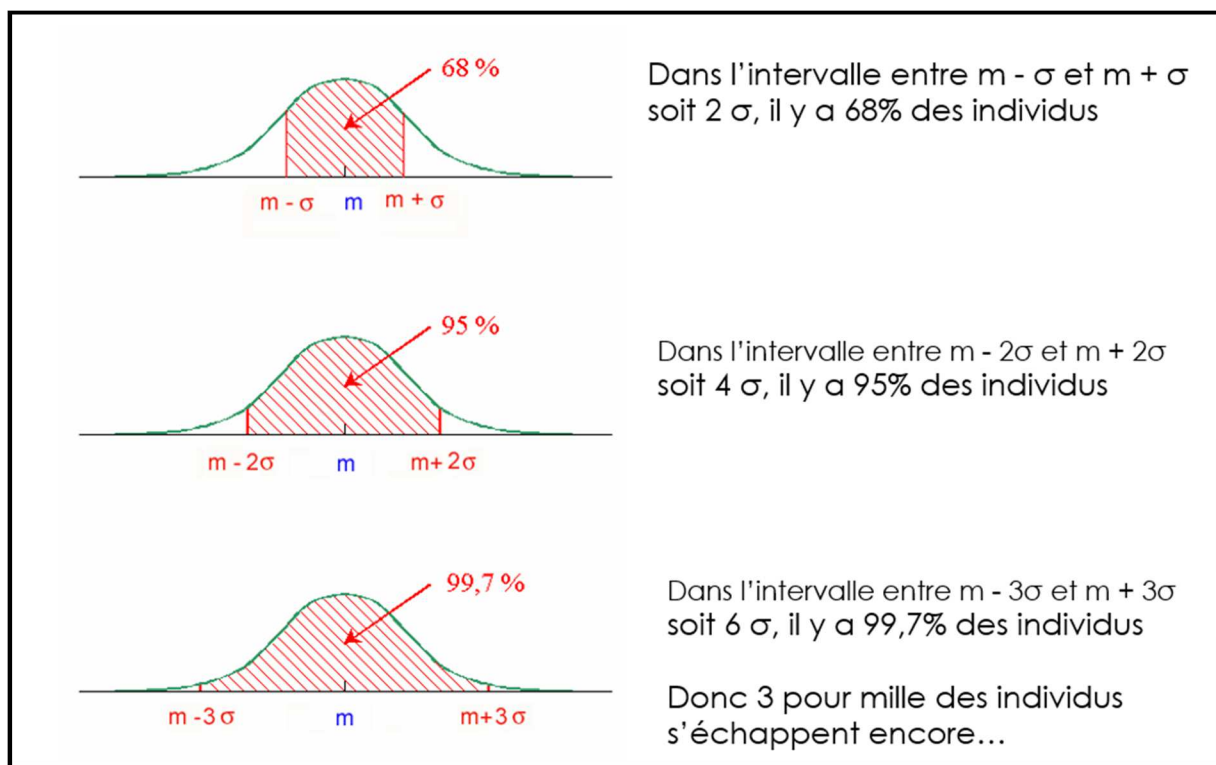


Fig.4.9 – Quelle fraction de la population dans quel intervalle ?

Et ainsi de suite. Mais d'où viennent ces chiffres ? De la table de la loi normale. Voir l'annexe ci-dessous.

Annexe :

Pour mémoire, mais ce n'est pas essentiel pour notre raisonnement, une courbe appelée la densité de la loi normale et représentée ci-dessous est donnée par :

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

où μ est la moyenne (notée m précédemment) et σ est l'écart type. Cette courbe dite de densité est, au sens mathématique, l'intégrale de la loi normale.

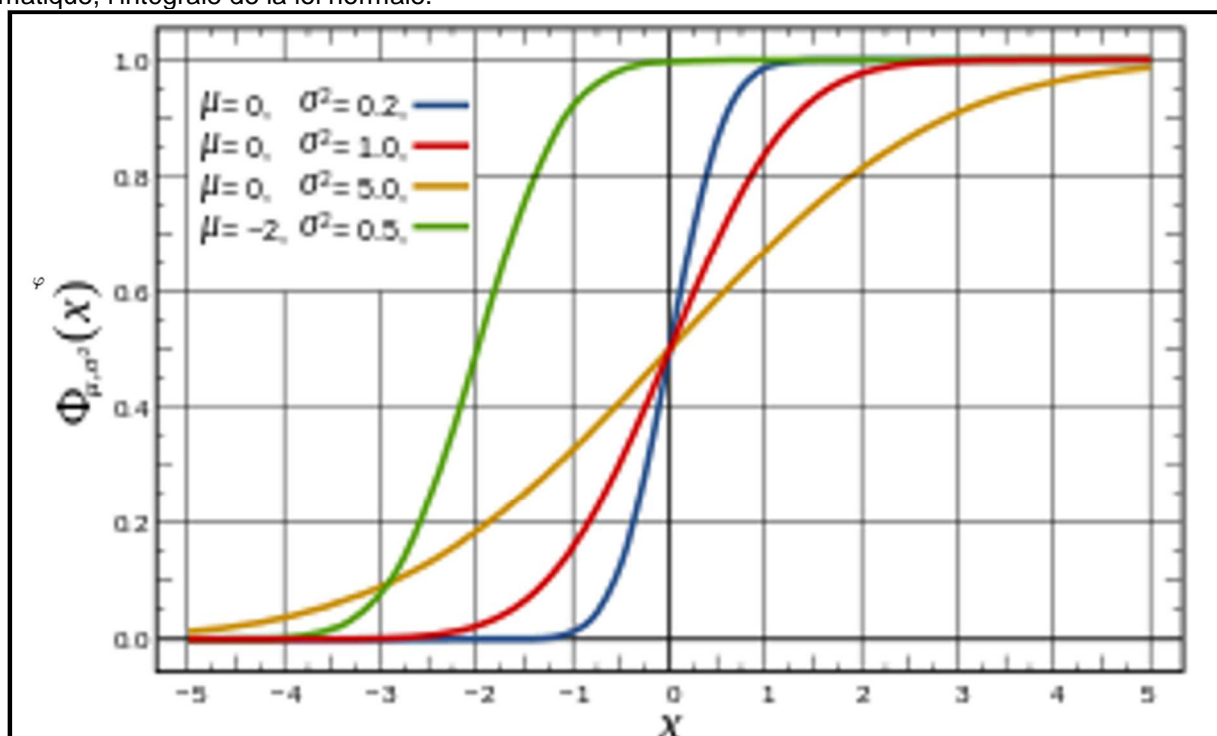


Fig.4.10 – Contenu de la table et Fonction de répartition

On parle de loi normale centrée réduite quand on a ramené la moyenne à 0 (en retirant la valeur moyenne à toutes les mesures, soit $x-\mu$) et qu'on a réduit (normalisé) l'échelle des distances en la divisant par σ .

La table des valeurs de la loi normale est obtenue en calculant l'intégrale de la densité de répartition (le pourcentage de tous les individus de la zone grisée sur la courbe par rapport à la population totale).

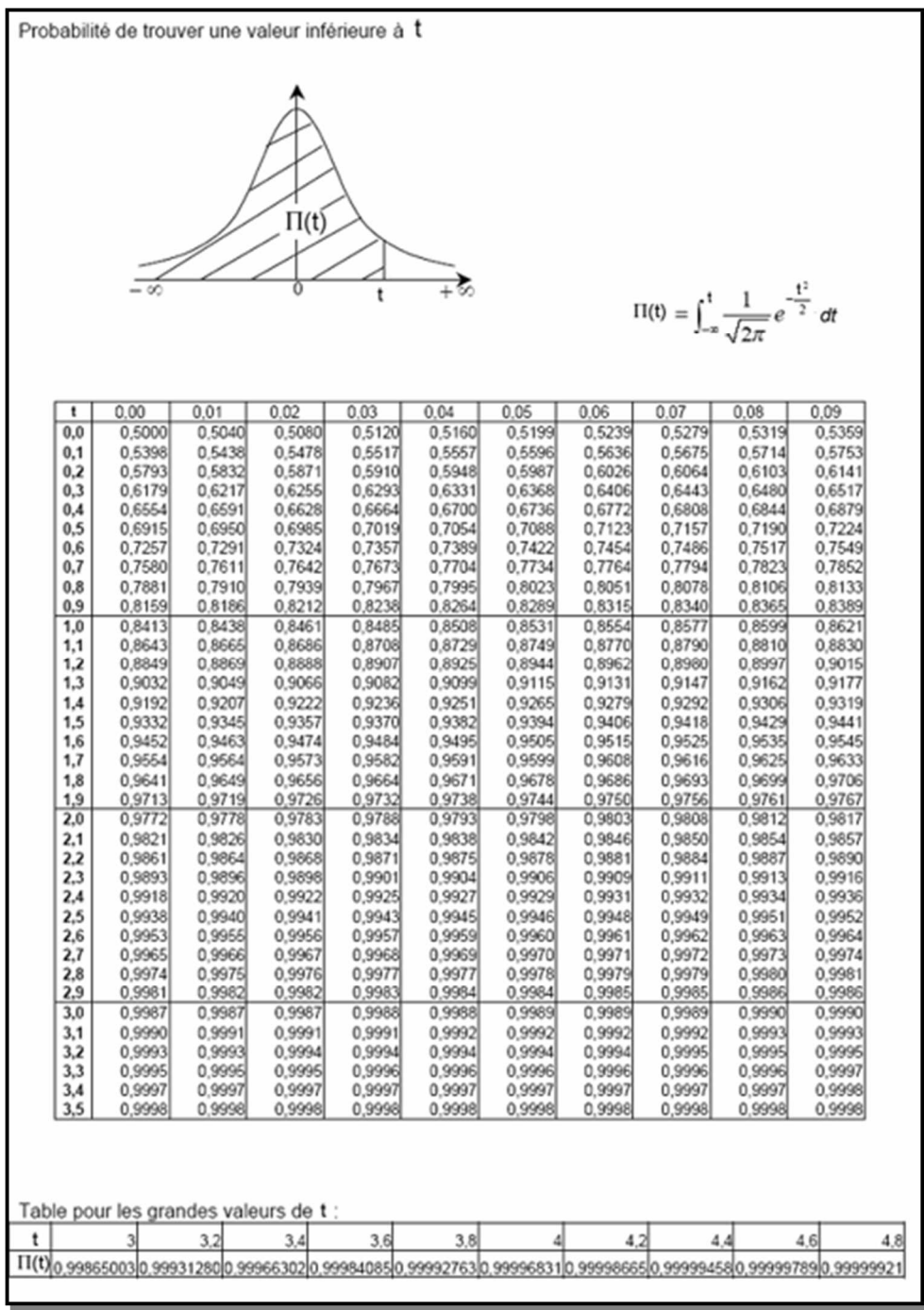


Fig.4.11 – Loi normale centrée réduite