

Cours de statistique inférentielle- STATINF

7 mai 2024

Table des matières

1	Prérequis sur les VA absolument continues	5
2	Échantillonnage et estimation de paramètres	7
I	Tables de la Loi Normale Centrée Réduite.	7
II	Loi du χ^2	9
III	Loi de Student	10
A	Compétences attendues à l'issue de ce cours	13
I	Échantillonnage	13
B	Tables de valeurs	15
I	Tables de la Loi Normale Centrée Réduite.	15
II	Loi du χ^2	17
III	Loi de Student	18

Prérequis sur les VA absolument continues

Ce chapitre constitue des rappels importants de notions vues dans le cours d'ALÉA.

Échantillonnage et estimation de paramètres

I- Construction de nouvelles lois

En combinant des lois normales, on peut créer de nouvelles lois de probabilités utiles en statistique :

1. Somme de lois normales

Proposition I.1

Soient deux variables aléatoires indépendantes $X_1 \sim \mathcal{N}(\mu_1, \sigma_1)$ et $X_2 \sim \mathcal{N}(\mu_2, \sigma_2)$, $a > 0$:

1. $X_1 + X_2$ suit une loi $\mathcal{N}(\mu_1 + \mu_2, \sqrt{\sigma_1^2 + \sigma_2^2})$
2. aX_1 suit une loi $\mathcal{N}(a\mu_1, |a|\sigma_1)$
3. $-X_1$ suit une loi $\mathcal{N}(-\mu_1, \sigma_1)$

Exemple 1.1

Si X_1, X_2, X_3 sont trois variables aléatoires indépendantes suivant une même loi $\mathcal{N}(\mu, \sigma^2)$, alors :

- $X_1 + X_2 + X_3$ suit une loi $\mathcal{N}(3\mu, \sqrt{3}\sigma)$
- en revanche, $3X_1$ suit une loi $\mathcal{N}(3\mu, 3\sigma)$.

Exemple 1.2

Si (X_1, \dots, X_n) est une suite de n variables aléatoires indépendantes suivant chacune une même loi $\mathcal{N}(\mu, \sigma^2)$, alors la variable

$$Y = \frac{1}{n} \sum_{i=1}^n X_i$$

suit une loi $\mathcal{N}(\mu, \frac{\sigma}{\sqrt{n}})$

2. Loi du Chi-deux

Définition I.1

On considère n variables aléatoires Z_1, \dots, Z_n indépendantes de même loi $\mathcal{N}(0, 1)$. Alors la variable

$$\sum_{i=1}^n Z_i^2$$

suit une loi du Chi-deux à n degrés de liberté $\chi^2(n)$.

Proposition I.2

Soit U une variable aléatoire suivant une loi $\chi^2(n)$. Alors U admet une fonction densité

$$f_U(x) = \begin{cases} \frac{1}{\Gamma(n/2)2^{n/2}} x^{n/2-1} e^{-x/2} & \text{si } x > 0 \\ 0 & \text{si } x \leq 0 \end{cases}$$

où la fonction Γ est définie comme suit :

$$\Gamma : z \mapsto \int_0^{+\infty} t^{z-1} e^{-t} dt$$

Pour rappels, la fonction Γ présente les propriétés suivantes :

- La fonction Γ est définie sur le demi plan complexe $\{Re(z) > 0\}$
- $\Gamma(1) = 1$
- $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$
- Pour tout $x > 0$, on a : $\Gamma(x+1) = x\Gamma(x)$

On peut calculer des valeurs à l'aide d'un outil informatique ou d'une table de valeurs. On admet la propriété suivante :

Proposition I.3

Si X suit une loi $\chi^2(n)$ alors

- $\mathbb{E}(X) = n$
- $\sigma^2(X) = 2n$

Proposition I.4

Si X_1, \dots, X_n une suite de n variables aléatoires indépendantes suivant chacune une même loi normale de moyenne μ et d'écart type σ . Alors pour tout entier i compris entre 1 et n , la variable aléatoire

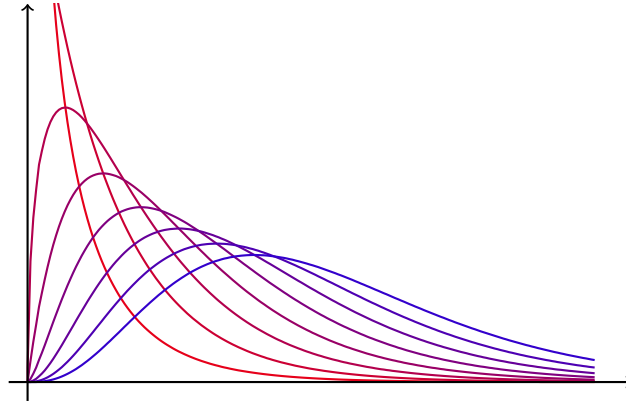
$$Z_i = \frac{X_i - \mu}{\sigma}$$

suit une loi $\mathcal{N}(0, 1)$. Par conséquent, la variable aléatoire

$$U = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2$$

suit une loi $\chi^2(n)$.

Cette loi est utile en statistique pour l'estimation et la théorie des tests statistiques.

FIGURE 2.1 – Courbe de $\chi^2(n)$ pour différentes valeurs n .

En remplaçant la moyenne μ par la moyenne empirique \bar{X} , on obtient le résultat (admis) suivant :

Théorème 2.1 : Fisher

Si X_1, \dots, X_n une suite de n variables aléatoires indépendantes suivant chacune une même loi normale de moyenne μ et d'écart type σ . On note

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Alors la variable aléatoire

$$V = \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma} \right)^2$$

suit une loi $\chi^2(n-1)$. De plus, la variable $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ suit une loi $\mathcal{N}(0, 1)$ et est indépendante de V .

3. Loi de Student

Définition I.2

On considère une variable $Z \sim \mathcal{N}(0, 1)$ et $U \sim \chi^2(n)$. On suppose que Z et U sont indépendantes. Alors la variable

$$T = \frac{Z}{\sqrt{\frac{U}{n}}}$$

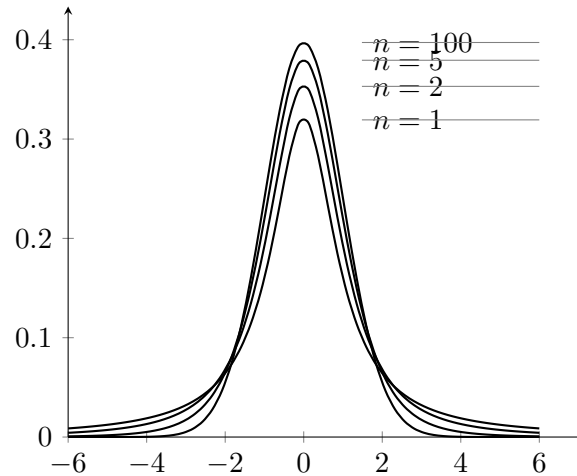
suit une loi de Student à n degrés de liberté $St(n)$.

De même que pour la loi du Chi-deux, on ne retiendra que ses paramètres pour une utilisation en statistiques.

Proposition I.5

Si T suit une loi $St(n)$ alors

- $\mathbb{E}(T) = 0$
- $\sigma^2(T) = \frac{n}{n-2}$ si $n > 2$.

FIGURE 2.2 – Courbe de $St(n)$ pour différentes valeurs de n .**Théorème 2.2**

Soit une suite de variables aléatoires (T_n) telle que pour tout entier $n \geq 1$, T_n suit une loi de Student $St(n)$. Alors (T_n) converge en loi vers une loi normale $\mathcal{N}(0, 1)$.

Démonstration. Soit (Z_n) une suite de variables aléatoires i.i.d. selon une loi normale $\mathcal{N}(0, 1)$. Ainsi, T_n suit la même loi que la variable

$$T'_n = \frac{Z_0}{\sqrt{\frac{\sum_{k=1}^n Z_k^2}{n}}}$$

Or d'après la loi forte des grands nombres, $\frac{1}{n} \sum_{k=1}^n Z_k^2 \xrightarrow[n \rightarrow +\infty]{p.s.} \mathbb{E}(Z_0^2) = 1$.

Donc la suite (T'_n) converge presque sûrement vers Z_0 . D'après la propriété ??, la suite (T'_n) converge en loi vers Z_0 . Or T_n a la même loi que T'_n , d'où la convergence en loi vers la loi de Z_0 . \square

Ce théorème permet de justifier que la table de valeurs fournie donne des résultats pour des lois de Student $St(n)$ avec $n \leq 30$. Au delà, on utilise une approximation par une loi normale centrée réduite. Voir une animation.

4. Loi de Fisher**Définition I.3**

On considère deux variables aléatoires indépendantes U_1, U_2 où $U_1 \sim \chi^2(n_1)$ et $U_2 \sim \chi^2(n_2)$. Alors la variable

$$F = \frac{U_1/n_1}{U_2/n_2}$$

suit une loi de Fisher à n_1 et n_2 degrés de liberté $F(n_1, n_2)$.

II-

Échantillon aléatoire et estimateur ponctuel

Le but d'une étude statistique est d'obtenir des informations sur l'ensemble d'une population. Lorsque celle-ci est de taille trop importante, on étudie un **échantillon** qui doit être prélevé de manière **aléatoire**. Ainsi, un échantillon sera modélisé par des variables aléatoires et la théorie de l'échantillonnage se basera sur la théorie des probabilités.

On considère par la suite qu'un tirage d'un échantillon de taille n se fait **avec remise**, ce qui a pour conséquence que les tirages sont supposés **indépendants**.

1. Définition

Définition II.1

Soit X une variable aléatoire définie sur un espace probabilisé Ω . Un **échantillon** de taille n de X est un n -uplet (X_1, \dots, X_n) de variables aléatoires i.i.d. suivant la même loi que X (appelée loi mère).

Une **réalisation** de cet échantillon est un n -uplet de réels (x_1, \dots, x_n) où pour tout $i \in [1; n]$, $X_i(\omega) = x_i$ avec $\omega \in \Omega$.

Définition II.2

Une **statistique** d'échantillonnage est une fonction de (X_1, \dots, X_n) .

2. Estimateurs ponctuels

Dans une population donnée, on considère un paramètre inconnu θ que l'on souhaite estimer (cela peut être sa moyenne, son écart-type, ou une proportion...). Soit (X_1, \dots, X_n) un échantillon de taille n de cette population.

Définition II.3

Un **estimateur** de θ est une statistique T_n des variables (X_1, \dots, X_n) dont la réalisation (oscillant autour de $\mathbb{E}(T)$) est envisagée comme une « bonne » valeur de θ . Si (x_1, \dots, x_n) est une réalisation de cet échantillon, alors $T_n(x_1, \dots, x_n)$ est un réel appelé **estimation** du paramètre θ .

Il existe une infinité d'estimateurs, mais certains choix d'estimateurs ont de bonnes propriétés que l'on détaillera par la suite.

Exemple 2.1

Un sondage effectué sur 300 votants d'une population de 50000 personnes a montré que 165 personnes sont prêtes à voter pour un candidat C. La proportion observée $p = \frac{165}{300}$ est une estimation de la proportion de votants prêts à voter pour le candidat C de la population de 50000 personnes. Cette estimation est obtenue à l'aide de l'estimateur $F = \frac{1}{300} \sum_{i=1}^{300} X_i$.

3. Propriétés attendues

Définition II.4

Soit T un estimateur d'un paramètre θ .

1. Le **biais** de T est la valeur $\mathbb{E}(T - \theta)$ notée $B(T)$;
2. L'**écart quadratique moyen** est la valeur $\mathbb{E}((T - \theta)^2)$ notée $EQM(T)$

Voici deux caractéristiques souhaitables pour un estimateur :

Définition II.5

Soit T_n un estimateur d'un paramètre θ associé à un échantillon de taille n .

1. L'estimateur T_n est **sans biais** si $B(T_n) = 0$.
2. L'estimateur T_n est **asymptotiquement sans biais** si $\lim_{n \rightarrow +\infty} B(T_n) = 0$.
3. L'estimateur T_n **converge en moyenne quadratique** si $\lim_{n \rightarrow +\infty} EQM(T_n) = 0$.

Proposition II.1

Soit T_n un estimateur d'un paramètre θ associé à un échantillon de taille n .

$$EQM(T_n) = \sigma^2(T_n) + B(T_n)^2$$

Conséquences et interprétations :

1. un estimateur converge en moyenne quadratique si et seulement s'il est asymptotiquement sans biais et sa variance tend vers 0.
2. Plus l'écart quadratique moyen d'un estimateur est petit, meilleure est **l'efficacité** de l'estimateur.

4. Moyenne empirique

On considère par la suite un n -échantillon (X_1, \dots, X_n) de X .

Définition II.6

La **moyenne empirique** de l'échantillon (X_1, \dots, X_n) est la statistique notée \bar{X} définie par

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Proposition II.2

Soit $\mu = \mathbb{E}(X)$ et $\sigma^2 = V(X)$ où X est la variable mère de l'échantillon (X_1, \dots, X_n) et \bar{X} la moyenne empirique de cet échantillon. Alors on a :

- $\mathbb{E}(\bar{X}) = \mu$
- $V(\bar{X}) = \frac{\sigma^2}{n}$

Démonstration. Par linéarité de l'espérance, on a que

$$\mathbb{E}(\bar{X}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X) = \frac{1}{n} \times n \times \mathbb{E}(X) = \mu$$

Par indépendance des variables X_1, \dots, X_n , on a $\sigma^2(X_1 + \dots + X_n) = \sigma^2(X_1) + \dots + \sigma^2(X_n) = n\sigma^2$, d'où le résultat. \square

5. Variance empirique

Définition II.7

La **variance empirique** de l'échantillon (X_1, \dots, X_n) est la statistique notée \tilde{S}^2 définie par

$$\tilde{S}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

Proposition II.3

Soit $\mu = \mathbb{E}(X)$ et $\sigma^2 = V(X)$ où X est la variable mère de l'échantillon (X_1, \dots, X_n) et \tilde{S}^2 la variance empirique de cet échantillon. Alors on a :

$$\mathbb{E}(\tilde{S}^2) = \frac{n-1}{n} \sigma^2$$

Démonstration. Pour faire apparaître $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ dans l'expression de \tilde{S}^2 , on écrit d'abord que pour tout i , $X_i - \bar{X} = X_i - \mu - (\bar{X} - \mu)$. Cela permet d'obtenir que

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X})^2 &= \sum_{i=1}^n (X_i - \mu - (\bar{X} - \mu))^2 \\ &= \sum_{i=1}^n (X_i - \mu)^2 - 2 \sum_{i=1}^n (X_i - \mu)(\bar{X} - \mu) + \sum_{i=1}^n (\bar{X} - \mu)^2 \\ &= \sum_{i=1}^n (X_i - \mu)^2 - 2(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + n(\bar{X} - \mu)^2 \\ &= \sum_{i=1}^n (X_i - \mu)^2 - 2(\bar{X} - \mu) \times n(\bar{X} - \mu) + n(\bar{X} - \mu)^2 \\ &= \sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2 \end{aligned}$$

Par linéarité de l'espérance, on en déduit que

$$\begin{aligned} \mathbb{E}(\tilde{S}^2) &= \frac{1}{n} \mathbb{E} \left(\sum_{i=1}^n (X_i - \mu)^2 \right) - \mathbb{E}((\bar{X} - \mu)^2) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}((X_i - \mu)^2) - \mathbb{E}((\bar{X} - \mu)^2) \\ &= \frac{1}{n} \sum_{i=1}^n \sigma^2(X_i) - \sigma^2(\bar{X}) \\ &= \frac{1}{n} \sum_{i=1}^n \sigma^2 - \frac{\sigma^2}{n} \\ &= \frac{n-1}{n} \sigma^2 \end{aligned}$$

\square

6. Fréquence

Lorsque la variable mère suit une loi de Bernoulli (succès ou échec), on s'intéressera à la fréquence de succès dans l'échantillon à l'aide de la statistique

$$F = \frac{1}{n} \sum_{i=1}^n X_i$$

appelée **fréquence empirique**.

Proposition II.4

Soit p le paramètre de la loi de Bernoulli suivie par la variable mère X de l'échantillon (X_1, \dots, X_n) . Soit F la fréquence empirique de cet échantillon. Alors on a :

- $\mathbb{E}(F) = p$
- $V(F) = \frac{p(1-p)}{n}$

Lors d'un sondage sur 1000 personnes, on peut observer la fréquence d'apparition d'un caractère f . On dit que f est une **estimation ponctuelle** de la proportion p de ce caractère dans la population mère, de taille trop importante pour être étudiée directement. On pourra également estimer la valeur de p à l'aide d'une fourchette $[p_1; p_2]$ appelé intervalle de confiance.

7. Résumé sur les estimateurs usuels

On considère un n -échantillon (X_1, \dots, X_n) de moyenne μ et d'écart-type σ . Les propriétés des statistiques étudiées précédemment permettent les conclusions suivantes :

Moyenne : $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ est un estimateur de la moyenne μ , c'est un estimateur **sans biais** et **convergent**.

Variance : on a 3 estimateurs usuels :

1. si μ est connu, $\Sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ est un estimateur de la variance σ^2 , c'est un estimateur **sans biais** et **convergent**.
2. si μ est inconnu, $\tilde{S}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ est un estimateur de la variance σ^2 , c'est un estimateur **biaisé** mais **convergent**.
3. si μ est inconnu, $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} \tilde{S}^2$ est un estimateur de la variance σ^2 , c'est un estimateur **sans biais** et **convergent**.

Fréquence : $F = \frac{1}{n} \sum_{i=1}^n X_i$, associé à une variable de Bernoulli de paramètre p , est un estimateur de la fréquence de succès p , c'est un estimateur **sans biais** et **convergent**.

III-

Estimation par maximum de vraisemblance

Les estimateurs utilisés précédemment ont été construits de manière empirique. Il existe des méthodes plus globales pour construire des estimateurs, en voici une basée sur le principe suivant : on cherche un estimateur qui rend le plus probable (vraisemblable) l'échantillon observé.

1. Fonction de vraisemblance pour une loi discrète

On considère un échantillon X_1, \dots, X_n de variables aléatoires indépendantes suivant une loi discrète $\mathcal{L}(\theta)$ où θ est le paramètre que l'on cherche à estimer. On note x_1, \dots, x_n une réalisation de l'échantillon. On s'intéresse à l'événement :

$$A = \{X_1 = x_1, \dots, X_n = x_n\}$$

La **fonction de vraisemblance** $L(x_1, \dots, x_n, \theta)$ est la probabilité de cet événement, autrement dit la probabilité que l'échantillon se réalise tel qu'on l'a observé. Par indépendance des variables aléatoires, on obtient :

$$L(x_1, \dots, x_n, \theta) = P_\theta(A) = P_\theta(X_1 = x_1) \times \dots \times P_\theta(X_n = x_n)$$

On notera que le calcul de probabilité dépend du paramètre θ que l'on cherche à estimer, puisque la loi $\mathcal{L}(\theta)$ dépend de ce paramètre.

2. Fonction de vraisemblance pour une loi continue

On considère un échantillon X_1, \dots, X_n de variables aléatoires indépendantes suivant une loi continue $\mathcal{L}(\theta)$ où θ est le paramètre que l'on cherche à estimer. On note x_1, \dots, x_n une réalisation de l'échantillon. L'événement $\{X_i = x_i\}$ étant de probabilité nulle pour une loi continue, on s'intéresse plutôt à l'événement $\{|X_i - x_i| < \varepsilon\}$. Supposons que la loi $\mathcal{L}(\theta)$ admette une densité f_θ : alors $P_\theta(|X_i - x_i| < \varepsilon) = \int_{x_i - \varepsilon}^{x_i + \varepsilon} f_\theta$. D'après le théorème de la moyenne, il existe $c \in]x_i - \varepsilon; x_i + \varepsilon[$ tel que $f_\theta(c) = \frac{1}{2\varepsilon} P_\theta(|X_i - x_i| < \varepsilon)$. Par continuité, on peut écrire $f_\theta(c) = f_\theta(x_i) + \alpha_i(\varepsilon)$ où $\alpha_i(\varepsilon)$ tend vers 0 lorsque $\varepsilon \rightarrow 0$.

De manière analogue au cas discret, on s'intéresse donc à l'événement

$$A_\varepsilon = \{|X_1 - x_1| < \varepsilon, \dots, |X_n - x_n| < \varepsilon\}$$

Or

$$\frac{1}{2\varepsilon} P_\theta(A_\varepsilon) = (f_\theta(x_1) + \alpha_1(\varepsilon)) \times \dots \times (f_\theta(x_n) + \alpha_n(\varepsilon))$$

En faisant tendre $\varepsilon \rightarrow 0$, on définit donc la fonction de vraisemblance

$$L(x_1, \dots, x_n, \theta) = P_\theta(A_\varepsilon) = f_\theta(x_1) \times \dots \times f_\theta(x_n)$$

3. Optimiser la vraisemblance

On cherche une valeur de θ telle que la vraisemblance soit maximale. Lorsque cela est possible, on cherche θ tel que

$$\frac{\partial L(x_1, \dots, x_i, \dots, x_n; \theta)}{\partial \theta} = 0$$

De manière équivalente, on cherche le maximum de la **log-vraisemblance** :

$$\frac{\partial \ln L(x_1, \dots, x_i, \dots, x_n; \theta)}{\partial \theta} = 0$$

ce qui peut conduire à des calculs plus simples.

Définition III.1

L'estimation du maximum de vraisemblance est, lorsqu'elle existe, la valeur $\hat{\theta}$ qui maximise la fonction $\theta \mapsto L(x_1, \dots, x_n, \theta)$.

L'estimation $\hat{\theta}$ dépend des réalisations x_1, \dots, x_n . Pour obtenir l'estimateur Θ recherché, on remplace dans l'expression les réalisations x_1, \dots, x_n par les variables aléatoires X_1, \dots, X_n .

4. Exemple : estimer le paramètre d'une loi de Bernoulli

On considère un échantillon X_1, \dots, X_n de variables aléatoires indépendantes suivant une loi de Bernoulli de paramètre p inconnu. Soit x_1, \dots, x_n une réalisation de cet échantillon :

$$P_p(X_i = x_i) = \begin{cases} p & \text{si } x_i = 1 \\ 1 - p & \text{si } x_i = 0 \end{cases}$$

Par indépendance, la fonction de vraisemblance devient alors :

$$L(x_1, \dots, x_n, p) = p^{\sum_{i=1}^n x_i} \times (1 - p)^{n - \sum_{i=1}^n x_i}$$

Cette fonction de p est strictement positive sur $]0; 1[$, nulle pour $p \in \{0; 1\}$. Elle admet un maximum pour une valeur $\hat{p} \in]0; 1[$. En supposant que $p \notin \{0; 1\}$, il est plus commode de considérer le logarithme de cette expression :

$$\ln L(x_1, \dots, x_n, p) = \left(\sum_{i=1}^n x_i \right) \ln(p) + \left(n - \sum_{i=1}^n x_i \right) \ln(1 - p)$$

puis de chercher à annuler sa dérivée :

$$\frac{\partial \ln L(x_1, \dots, x_i, \dots, x_n; p)}{\partial p} = \frac{\sum_{i=1}^n x_i}{p} - \frac{n - \sum_{i=1}^n x_i}{1 - p} = 0 \iff p = \frac{1}{n} \sum_{i=1}^n x_i$$

On trouve donc une valeur $\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i$. On en déduit un estimateur de p qui est l'estimateur de fréquence déjà connu :

$$F = \frac{1}{n} \sum_{i=1}^n X_i$$

IV- Estimation par intervalle de confiance

On cherche maintenant à estimer un paramètre à l'aide d'un intervalle. Cela permet de contrôler l'erreur lors de l'estimation.

1. Démarche générale

Étant donné un estimateur T d'un paramètre θ et un niveau de confiance $(1 - \alpha)$ où $\alpha \in]0; 1[$, on cherche à déterminer un intervalle $I_{conf}(T)$ tel que

$$P(\theta \in I_{conf}(T)) = 1 - \alpha$$

Le nombre α est le **risque d'erreur**, on choisit souvent $\alpha = 5\%$ ou $\alpha = 1\%$.

On cherche autant que possible à déterminer des intervalles de confiance à **risques symétriques**, c'est-à-dire des intervalles de la forme $[T - \epsilon; T + \epsilon]$ avec

$$P(\theta < T - \epsilon) = P(\theta > T + \epsilon) = \frac{\alpha}{2}$$

Pour calculer ces probabilités, il faut connaître en particulier les lois suivies par les estimateurs usuels (on rappelle qu'un estimateur est avant tout une variable aléatoire).

2. Intervalle de confiance asymptotique

Un intervalle de confiance asymptotique au niveau $1 - \alpha$ est un intervalle $I_n(T)$ tel que

$$\lim_{n \rightarrow +\infty} P(\theta \in I_n(T)) = 1 - \alpha$$

En pratique, on utilise un tel intervalle « quand n est grand ». Cette situation se produit notamment lorsque l'on approche la loi de l'estimateur à l'aide du théorème central limite.

3. Intervalle de confiance asymptotique d'une proportion

Estimateur utilisé : Soit $F = \frac{1}{n} \sum_{i=1}^n X_i$ l'estimateur de fréquence associé à une variable de Bernoulli de paramètre p . Alors la variable nF suit une loi binomiale $\mathcal{B}(n, p)$ et on en déduit que $\mathbb{E}(F) = p$ et $\sigma^2(F) = \frac{p(1-p)}{n}$.

Loi de l'estimateur utilisée : On suppose que $n > 30$, $np > 5$ et $n(1-p) > 5$. D'après le théorème central limite, on sait donc que F peut-être approché par une loi normale $\mathcal{N}\left(p, \frac{p(1-p)}{n}\right)$.

Calcul d'un intervalle de confiance asymptotique : On sait que $\frac{F-p}{\sqrt{\frac{p(1-p)}{n}}}$ suit approximativement une loi normale centrée réduite $\mathcal{N}(0, 1)$. On rappelle qu'il existe un unique réel strictement positif $u_{\alpha/2}$ tel que

$$P(-u_{\alpha/2} < Z < u_{\alpha/2}) = 1 - \alpha$$

où $Z = \frac{F-p}{\sqrt{\frac{p(1-p)}{n}}}$ suit une loi $\mathcal{N}(0, 1)$.

Or

$$\begin{aligned} -u_{\alpha/2} \leq Z \leq u_{\alpha/2} &\iff -u_{\alpha/2} \leq \frac{F-p}{\sqrt{\frac{p(1-p)}{n}}} \leq u_{\alpha/2} \\ &\iff F - u_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq p \leq F + u_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \\ &\iff p \in I_{\text{conf}}(F) \end{aligned}$$

où

$$I_{\text{conf}}(F) = \left[F - u_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} ; F + u_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \right]$$

Cet intervalle aléatoire contient la valeur p recherchée avec une probabilité $1 - \alpha$.

Un intervalle de confiance est aussi une *réalisation* de cet intervalle aléatoire à partir des données observées. La réalisation $F(\omega)$ de la variable F est dite « fréquence observée dans l'échantillon » que l'on note également f_{obs} . La réalisation de cet intervalle est donc :

$$\left[f_{\text{obs}} - u_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} ; f_{\text{obs}} + u_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \right]$$

Modification de l'intervalle avec des données connues : Le problème de l'intervalle de confiance trouvé ici est que ses bornes dépendent du paramètre **inconnu** p . L'écart-type $\tilde{\sigma} = \sqrt{\frac{f(1-f)}{n}}$ associé à la fréquence observée est un estimateur de l'écart-type $\sigma = \sqrt{\frac{p(1-p)}{n}}$ mais il est biaisé, on peut en revanche estimer σ par $\sqrt{\frac{n}{n-1}}\tilde{\sigma}$. On peut donc prendre pour intervalle de confiance :

$$I_{conf}(F(\omega)) = \left[f_{obs} - u_{\alpha/2} \sqrt{\frac{f_{obs}(1-f_{obs})}{n-1}} ; f_{obs} + u_{\alpha/2} \sqrt{\frac{f_{obs}(1-f_{obs})}{n-1}} \right]$$

En pratique :

1. Dans une situation courante : si n est grand, on peut remplacer $n-1$ par n :

$$I_{conf}(F(\omega)) = \left[f_{obs} - u_{\alpha/2} \sqrt{\frac{f_{obs}(1-f_{obs})}{n}} ; f_{obs} + u_{\alpha/2} \sqrt{\frac{f_{obs}(1-f_{obs})}{n}} \right]$$

Cet intervalle, dit « intervalle de Wald », comporte des défauts statistiques mais est très communément utilisé.

2. Si f_{obs} n'est pas connu et qu'on s'intéresse seulement à la longueur de l'intervalle, on peut majorer $p(1-p)$ par 0,25, on obtient ainsi un intervalle de confiance par excès de p :

$$I_{conf}(F(\omega)) = \left[f_{obs} - u_{\alpha/2} \frac{1}{2\sqrt{n}} ; f_{obs} + u_{\alpha/2} \frac{1}{2\sqrt{n}} \right]$$

Ce résultat est d'autant plus pertinent que $p(1-p)$ est proche de son maximum, c'est-à-dire si p est proche de 0.5.

Dans ce cas, la longueur de l'intervalle est $u_{\alpha/2} \frac{1}{\sqrt{n}}$. Cette quantité ne dépend que de la taille de l'échantillon et du niveau de confiance souhaité.

4. Intervalle de confiance d'une moyenne

Estimateur utilisé : Soit $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ l'estimateur de la moyenne μ pour une loi mère de moyenne μ et d'écart-type σ . On a calculé que $\mathbb{E}(\bar{X}) = \mu$ et $\sigma^2(\bar{X}) = \frac{\sigma^2}{n}$.

Loi de l'estimateur utilisé :

1. Si l'échantillon est distribué normalement selon une $\mathcal{N}(\mu, \sigma^2)$ alors \bar{X} suit (exactement) une loi $\mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$.
2. Si on ne sait pas que l'échantillon est distribué normalement et que n est considéré comme grand, on peut considérer que \bar{X} suit approximativement une loi normale $\mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ d'après le théorème central limite.

Calcul d'un intervalle de confiance pour cette loi : On sait que $\frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}}$ suit une loi normale centrée réduite $\mathcal{N}(0, 1)$. De même que précédemment, il existe un unique réel strictement positif $u_{\alpha/2}$ tel que

$$P\left(\bar{X} - u_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + u_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

d'où

$$I_{conf}(\bar{X}) = \left[\bar{X} - u_{\alpha/2} \frac{\sigma}{\sqrt{n}} ; \bar{X} + u_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

On note $x_{obs} = \bar{X}(\omega)$ la moyenne observée dans l'échantillon, c'est une réalisation de la variable \bar{X} .

5. Intervalle de confiance d'une moyenne, l'écart-type est inconnu

Le calcul précédent est possible en connaissant l'écart-type σ de la loi mère. En général, quand on cherche la moyenne μ , l'écart-type σ est également inconnu.

Estimateur utilisé : On fait intervenir un estimateur de σ^2 , ici l'estimateur usuel non biaisé

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Loi de l'estimateur utilisé si la loi mère a une distribution normale : on sait donc que $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ suit une loi normale $\mathcal{N}(0, 1)$. En remplaçant σ par l'estimateur S , d'après la définition d'une loi de Student et le théorème de Fisher (th ??)

$$\frac{Z}{\sqrt{\frac{(n-1)S^2}{\sigma^2}}}$$

suit une loi $St(n-1)$. En remplaçant Z par son expression et en simplifiant, il vient que $U = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}}$ suit une loi $St(n-1)$.

Calcul d'un intervalle de confiance pour cette loi : De même que pour une loi normale, il existe un unique réel strictement positif $t_{\alpha/2}$ tel que

$$P(-t_{\alpha/2} < U < t_{\alpha/2}) = 1 - \alpha$$

On en déduit un intervalle de confiance asymptotique :

$$I_{conf}(\bar{X}) = \left[\bar{X} - t_{\alpha/2} \frac{S}{\sqrt{n}} ; \bar{X} + t_{\alpha/2} \frac{S}{\sqrt{n}} \right]$$

que l'on réalise en utilisant les valeurs observées.

Loi de l'estimateur utilisée si la loi mère a une distribution quelconque et n est grand : d'après le théorème central limite, la variable $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ suit approximativement une loi normale $\mathcal{N}(0, 1)$. On admet qu'alors, en remplaçant σ par son estimateur S , $U = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}}$ suit approximativement une loi $\mathcal{N}(0, 1)$.

On en déduit un intervalle de confiance :

$$I_{conf}(\bar{X}) = \left[\bar{X} - u_{\alpha/2} \frac{S}{\sqrt{n}} ; \bar{X} + u_{\alpha/2} \frac{S}{\sqrt{n}} \right]$$

Si la loi mère a une distribution quelconque et que l'échantillon n'est pas de grande taille n : on ne peut rien dire.

6. Intervalle de confiance d'une variance

Estimateur utilisé : On fait intervenir un estimateur de σ^2 , en l'occurrence

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Loi de l'estimateur utilisé si la loi mère a une distribution normale : on a vu précédemment que la variable $\frac{(n-1)}{\sigma^2} S^2$ suit une loi $\chi^2(n-1)$, on obtient un intervalle de confiance pour estimer une variance.

Remarquons que si la moyenne μ est connue, l'estimateur est

$$\Sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$$

et on sait que $\frac{n}{\sigma^2} \Sigma^2$ suit une loi $\chi^2(n)$.



Compétences attendues à l'issue de ce cours

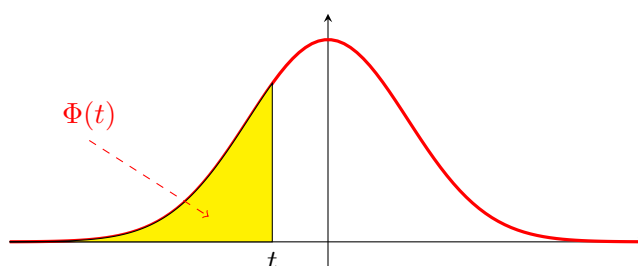
I- Échantillonnage

- Connaître les définitions de population, échantillon, réalisation d'un échantillon ;
- Estimateur : connaître les définitions et les propriétés (biais, convergence, efficacité) ;
- Appliquer la méthode du maximum de vraisemblance pour construire un estimateur pour des lois discrètes et continues ;
- Connaître les estimateurs usuels (moyenne, variance, proportion) et leurs propriétés ;
- Déterminer la loi exacte ou approchée des estimateurs usuels ;
- Construire un intervalle de confiance exact ou asymptotique ;

Tables de valeurs

I- Tables de la Loi Normale Centrée Réduite.

$$\Phi(t) = P(X \leq t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \quad \text{et} \quad \Phi(-t) = 1 - \Phi(t).$$



Voici une table des quantiles les plus couramment utilisées :

p	$\Phi^{-1}(p)$	p	$\Phi^{-1}(p)$	p	$\Phi^{-1}(p)$
0.9	1.2815515655	0.999	3.0902323062	0.99999	4.2648907939
0.95	1.6448536270	0.9995	3.2905267315	0.999995	4.4171734135
0.975	1.9599639845	0.99975	3.4807564043	0.9999975	4.5647877303
0.99	2.3263478740	0.9999	3.7190164855	0.999999	4.7534243088
0.995	2.5758293035	0.99995	3.8905918864	0.9999995	4.8916384757
0.9975	2.8070337683	0.999975	4.0556269811	0.99999975	5.0263128360

La table suivante donne toutes les valeurs de $\Phi(t)$ où t parcourt l'intervalle $[0; 3.99]$ par pas de 10^{-2} .

Exemple 0.1

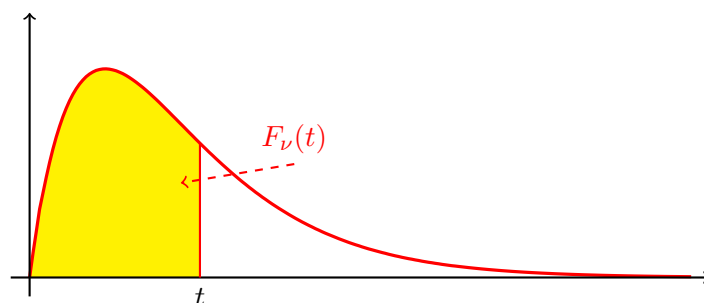
On lira dans la table ci-contre que $P(X \leq 1.96) \approx 97.50\%$.

[illegible]

II- Loi du χ^2

$$F_\nu(t) = P(X \leq t)$$

La table suivante contient les quantiles de la loi χ^2 avec ν degrés de liberté. Pour tout $0 < p < 1$, le quantile est la valeur de t pour laquelle $P\{X \leq t\} = p$, où $X \sim \chi^2(\nu)$. Ainsi $t = F_\nu^{-1}(p)$.



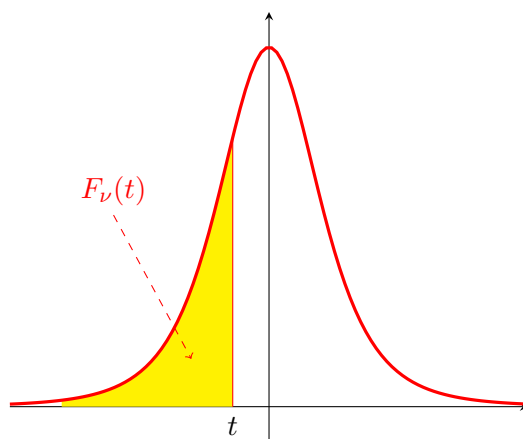
ν	p											
	0.005	0.01	0.025	0.05	0.1	0.5	0.9	0.95	0.975	0.99	0.995	0.999
1	0.0000	0.0002	0.0010	0.0039	0.0158	0.4549	2.7055	3.8415	5.0239	6.6349	7.8794	10.828
2	0.0100	0.0201	0.0506	0.1026	0.2107	1.3863	4.6052	5.9915	7.3778	9.2103	10.597	13.816
3	0.0717	0.1148	0.2158	0.3518	0.5844	2.3660	6.2514	7.8147	9.3484	11.345	12.838	16.266
4	0.2070	0.2971	0.4844	0.7107	1.0636	3.3567	7.7794	9.4877	11.143	13.277	14.860	18.467
5	0.4117	0.5543	0.8312	1.1455	1.6103	4.3515	9.2364	11.070	12.833	15.086	16.750	20.515
6	0.6757	0.8721	1.2373	1.6354	2.2041	5.3481	10.645	12.592	14.449	16.812	18.548	22.458
7	0.9893	1.2390	1.6899	2.1673	2.8331	6.3458	12.017	14.067	16.013	18.475	20.278	24.322
8	1.3444	1.6465	2.1797	2.7326	3.4895	7.3441	13.362	15.507	17.535	20.090	21.955	26.124
9	1.7349	2.0879	2.7004	3.3251	4.1682	8.3428	14.684	16.919	19.023	21.666	23.589	27.877
10	2.1559	2.5582	3.2470	3.9403	4.8652	9.3418	15.987	18.307	20.483	23.209	25.188	29.588
11	2.6032	3.0535	3.8157	4.5748	5.5778	10.341	17.275	19.675	21.920	24.725	26.757	31.264
12	3.0738	3.5706	4.4038	5.2260	6.3038	11.340	18.549	21.026	23.337	26.217	28.300	32.909
13	3.5650	4.1069	5.0088	5.8919	7.0415	12.340	19.812	22.362	24.736	27.688	29.819	34.528
14	4.0747	4.6604	5.6287	6.5706	7.7895	13.339	21.064	23.685	26.119	29.141	31.319	36.123
15	4.6009	5.2293	6.2621	7.2609	8.5468	14.339	22.307	24.996	27.488	30.578	32.801	37.697
16	5.1422	5.8122	6.9077	7.9616	9.3122	15.338	23.542	26.296	28.845	32.000	34.267	39.252
17	5.6972	6.4078	7.5642	8.6718	10.085	16.338	24.769	27.587	30.191	33.409	35.718	40.790
18	6.2648	7.0149	8.2307	9.3905	10.865	17.338	25.989	28.869	31.526	34.805	37.156	42.312
19	6.8440	7.6327	8.9065	10.117	11.651	18.338	27.204	30.144	32.852	36.191	38.582	43.820
20	7.4338	8.2604	9.5908	10.851	12.443	19.337	28.412	31.410	34.170	37.566	39.997	45.315
21	8.0337	8.8972	10.283	11.591	13.240	20.337	29.615	32.671	35.479	38.932	41.401	46.797
22	8.6427	9.5425	10.982	12.338	14.041	21.337	30.813	33.924	36.781	40.289	42.796	48.268
23	9.2604	10.196	11.689	13.091	14.848	22.337	32.007	35.172	38.076	41.638	44.181	49.728
24	9.8862	10.856	12.401	13.848	15.659	23.337	33.196	36.415	39.364	42.980	45.559	51.179
25	10.520	11.524	13.120	14.611	16.473	24.337	34.382	37.652	40.646	44.314	46.928	52.620
26	11.160	12.198	13.844	15.379	17.292	25.336	35.563	38.885	41.923	45.642	48.290	54.052
27	11.808	12.879	14.573	16.151	18.114	26.336	36.741	40.113	43.195	46.963	49.645	55.476
28	12.461	13.565	15.308	16.928	18.939	27.336	37.916	41.337	44.461	48.278	50.993	56.892
29	13.121	14.256	16.047	17.708	19.768	28.336	39.087	42.557	45.722	49.588	52.336	58.301
30	13.787	14.953	16.791	18.493	20.599	29.336	40.256	43.773	46.979	50.892	53.672	59.703

III- Loi de Student

$$F_\nu(t) = P(X \leq t)$$

La table suivante contient les quantiles de la loi de Student avec ν degrés de liberté. Pour tout $0 < p < 1$, le quantile est la valeur de t pour laquelle $P\{X \leq t\} = p$, où $X \sim St(\nu)$. Ainsi $t = F_\nu^{-1}(p)$.

Cette table ne contient que les quantiles pour des valeurs $p \geq \frac{1}{2}$. Si $p < \frac{1}{2}$, les quantiles peuvent être obtenus par symétrie de la loi : $F_\nu^{-1}(p) = -F_\nu^{-1}(1 - p)$.



ν	p											
	0.6	0.7	0.75	0.8	0.85	0.9	0.95	0.975	0.99	0.995	0.999	0.9995
1	0.3249	0.7265	1.0000	1.3764	1.9626	3.0777	6.3138	12.706	31.821	63.657	318.31	636.62
2	0.2887	0.6172	0.8165	1.0607	1.3862	1.8856	2.9200	4.3027	6.9646	9.9248	22.327	31.599
3	0.2767	0.5844	0.7649	0.9785	1.2498	1.6377	2.3534	3.1824	4.5407	5.8409	10.215	12.924
4	0.2707	0.5686	0.7407	0.9410	1.1896	1.5332	2.1318	2.7764	3.7469	4.6041	7.1732	8.6103
5	0.2672	0.5594	0.7267	0.9195	1.1558	1.4759	2.0150	2.5706	3.3649	4.0321	5.8934	6.8688
6	0.2648	0.5534	0.7176	0.9057	1.1342	1.4398	1.9432	2.4469	3.1427	3.7074	5.2076	5.9588
7	0.2632	0.5491	0.7111	0.8960	1.1192	1.4149	1.8946	2.3646	2.9980	3.4995	4.7853	5.4079
8	0.2619	0.5459	0.7064	0.8889	1.1081	1.3968	1.8595	2.3060	2.8965	3.3554	4.5008	5.0413
9	0.2610	0.5435	0.7027	0.8834	1.0997	1.3830	1.8331	2.2622	2.8214	3.2498	4.2968	4.7809
10	0.2602	0.5415	0.6998	0.8791	1.0931	1.3722	1.8125	2.2281	2.7638	3.1693	4.1437	4.5869
11	0.2596	0.5399	0.6974	0.8755	1.0877	1.3634	1.7959	2.2010	2.7181	3.1058	4.0247	4.4370
12	0.2590	0.5386	0.6955	0.8726	1.0832	1.3562	1.7823	2.1788	2.6810	3.0545	3.9296	4.3178
13	0.2586	0.5375	0.6938	0.8702	1.0795	1.3502	1.7709	2.1604	2.6503	3.0123	3.8520	4.2208
14	0.2582	0.5366	0.6924	0.8681	1.0763	1.3450	1.7613	2.1448	2.6245	2.9768	3.7874	4.1405
15	0.2579	0.5357	0.6912	0.8662	1.0735	1.3406	1.7531	2.1314	2.6025	2.9467	3.7328	4.0728
16	0.2576	0.5350	0.6901	0.8647	1.0711	1.3368	1.7459	2.1199	2.5835	2.9208	3.6862	4.0150
17	0.2573	0.5344	0.6892	0.8633	1.0690	1.3334	1.7396	2.1098	2.5669	2.8982	3.6458	3.9651
18	0.2571	0.5338	0.6884	0.8620	1.0672	1.3304	1.7341	2.1009	2.5524	2.8784	3.6105	3.9216
19	0.2569	0.5333	0.6876	0.8610	1.0655	1.3277	1.7291	2.0930	2.5395	2.8609	3.5794	3.8834
20	0.2567	0.5329	0.6870	0.8600	1.0640	1.3253	1.7247	2.0860	2.5280	2.8453	3.5518	3.8495
21	0.2566	0.5325	0.6864	0.8591	1.0627	1.3232	1.7207	2.0796	2.5176	2.8314	3.5272	3.8193
22	0.2564	0.5321	0.6858	0.8583	1.0614	1.3212	1.7171	2.0739	2.5083	2.8188	3.5050	3.7921
23	0.2563	0.5317	0.6853	0.8575	1.0603	1.3195	1.7139	2.0687	2.4999	2.8073	3.4850	3.7676
24	0.2562	0.5314	0.6848	0.8569	1.0593	1.3178	1.7109	2.0639	2.4922	2.7969	3.4668	3.7454
25	0.2561	0.5312	0.6844	0.8562	1.0584	1.3163	1.7081	2.0595	2.4851	2.7874	3.4502	3.7251
26	0.2560	0.5309	0.6840	0.8557	1.0575	1.3150	1.7056	2.0555	2.4786	2.7787	3.4350	3.7066
27	0.2559	0.5306	0.6837	0.8551	1.0567	1.3137	1.7033	2.0518	2.4727	2.7707	3.4210	3.6896
28	0.2558	0.5304	0.6834	0.8546	1.0560	1.3125	1.7011	2.0484	2.4671	2.7633	3.4082	3.6739
29	0.2557	0.5302	0.6830	0.8542	1.0553	1.3114	1.6991	2.0452	2.4620	2.7564	3.3962	3.6594

→

ν	0.6	0.7	0.75	0.8	0.85	0.9	0.95	0.975	0.99	0.995	0.999	0.9995
30	0.2556	0.5300	0.6828	0.8538	1.0547	1.3104	1.6973	2.0423	2.4573	2.7500	3.3852	3.6460