

Confidence regions

the notion of confidence regions is another way to quantify the precision of an estimation.

I Definition and properties

def: let $\{(\mathcal{Z}, \mathcal{F}), (\mathcal{P}, \mathcal{G}), \{\mathbb{P}_\theta, \theta \in \Theta\}, \mathcal{Z}\}$ be a statistical model. let $g: \Theta \rightarrow \mathbb{R}^p$ and $\alpha \in]0, 1[$. A function $\mathcal{C}: \mathcal{Z} \rightarrow \mathcal{G}(\mathbb{R}^p)$ is called a confidence region at level $1-\alpha$ for $g(\theta)$ if $\forall \theta \in \Theta: \{z \in \mathcal{Z}: g(\theta) \in \mathcal{C}(z)\} \subset \mathcal{Z}$ and

$$\inf_{\theta \in \Theta} \mathbb{P}_\theta(\{z \in \mathcal{Z}: g(\theta) \in \mathcal{C}(z)\}) \geq 1-\alpha$$

If $p=1$ then we talk about confidence intervals, and they take one of the 3 following forms:

- (i) $\mathcal{C}(z) = [m(z); +\infty[$
- (ii) $\mathcal{C}(z) =]-\infty; \pi(z)]$
- (iii) $\mathcal{C}(z) = [m(z); \pi(z)]$

Example: let us go back to the survey.

let X_1, \dots, X_n iid $\mathcal{B}(\theta)$

We have seen that with $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ (the average number of success)

$$\sup_{\theta \in \Theta} \mathbb{E}_\theta[(\hat{\theta}_n - \theta)^2] = \frac{1}{4n}$$

which writes also $\sup_{\theta \in \Theta} \mathbb{E}_\theta[\sqrt{n}(\hat{\theta}_n - \theta)^2] = \frac{1}{4}$

We now construct a confidence interval.

let $\delta > 0$ and $\theta \in \Theta$

thanks to the Bienayme - Tchebychev inequality

$$\mathbb{P}_\theta(|\hat{\theta}_n - \theta| \geq \delta) \leq \frac{1}{\delta^2} \text{Var}_\theta(\hat{\theta}_n) \leq \frac{1}{4n\delta^2}$$

If we choose $\delta_{1-\alpha}$ such that $\frac{1}{4n\delta_{1-\alpha}^2} = \alpha$ i.e. $\delta_{1-\alpha} = \frac{1}{2\sqrt{n\alpha}}$ and write

$$I_{n,\alpha} = \left[\hat{\theta}_n - \frac{1}{2\sqrt{n\alpha}}; \hat{\theta}_n + \frac{1}{2\sqrt{n\alpha}} \right] \text{ then:}$$

$$\forall \theta \in \Theta \quad \mathbb{P}_\theta(\theta \in I_{n,\alpha}) \geq 1-\alpha$$

the quality of this interval is given by its length: $|I_{n,\alpha}| = \frac{1}{\sqrt{n\alpha}}$

lemme de Bienayme - Tchebychev:

Soit Z une v.a. tq $\mathbb{E}[Z^2] < \infty$. Alors $\forall \delta > 0$:

$$\mathbb{P}(|Z - \mathbb{E}[Z]| > \delta) \leq \frac{\text{Var}(Z)}{\delta^2}$$

Note that $|J_{n,n}| \rightarrow +\infty$ when $\alpha \rightarrow 0$. This is a tradeoff that one can easily understand. Indeed, one has to choose between the precision of the estimation ($|J_{n,n}|$ small) and the risk we accept (α small).

II Classical confidence regions:

Let us first introduce some classical distributions.

1) χ_k^2 : the Chi-square with k degrees of freedom.

def: Let (X_1, \dots, X_k) be iid random variables such that $X_i \sim \mathcal{P}(0, 1)$. Then $U = \sum_{i=1}^k X_i^2$ follows a centered χ^2 distribution with k degrees of freedom.

Prop: the pdf of U is $f_k(x) = \frac{1}{2^{k/2} \Gamma(\frac{k}{2})} x^{k/2-1} \exp(-\frac{x}{2}) \mathbb{1}_{x>0}$

$$\bullet \quad \mathbb{E}[U] = k \quad \text{Var}(U) = 2k$$

def: If $X_i \sim \mathcal{P}(\mu_i, 1)$ then $U = \sum_{i=1}^k X_i$ follows a non-centered χ^2 distribution with k degrees of freedom denoted by $\chi_k^2(\gamma)$ where $\gamma = \frac{1}{2} \sum_{i=1}^k \mu_i^2$.

In this case $\mathbb{E}(U) = \gamma + k$ and $\text{Var}(U) = 2k + 4\gamma$

2) Student distribution

def: Let X and Y be 2 independent r.v. such that:

$$\bullet X \sim \mathcal{P}(0, 1)$$

$$\bullet Y \sim \chi_r^2 \text{ (centered)}$$

then $T = \frac{X}{\sqrt{Y/r}} \sim t(r)$ a Student distribution with r degrees of freedom.

Prop: $f_r(t) = \frac{\Gamma(\frac{r+1}{2})}{\Gamma(\frac{r}{2})} \frac{1}{(\pi r)^{1/2}} \left(1 + \frac{t^2}{r}\right)^{-\frac{r+1}{2}}$

Theorem: Let X_1, \dots, X_n iid $\mathcal{P}(\mu, \sigma^2)$ then:

$$(i) \quad \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \text{ (the empirical mean)} \quad \text{and} \quad S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

the variance of the sample are independent

$$(ii) \quad \bar{X}_n \sim \mathcal{P}\left(\mu, \frac{\sigma^2}{n}\right)$$

$$(iii) \quad (n-1) S_n^2 \sim \chi_{n-1}^2$$

3) Fisher's distribution:

def: let $X, Y \in \mathbb{R}^n$ independent r.v such that

- $X \sim \chi_q$
- $Y \sim \chi_x^2$

then $W = \frac{X/q}{Y/r} \sim F(q, r)$ a Fisher's distribution with (q, r) degrees of freedom

Prop: $f_{(q,n)}(w) = \frac{\Gamma\left(\frac{q+2}{2}\right)}{\Gamma\left(\frac{q}{2}\right)\Gamma\left(\frac{n}{2}\right)} \left(\frac{q}{n}\right)^{q/2} \frac{w^{n/2-1}}{\left(1 + \left(\frac{q}{n}\right)w\right)^{\frac{(q+n)}{2}}}$ $1, w > 0$

4) Confidence interval for a normal sample with known variance:

Let (X_1, \dots, X_n) i.i.d samples of $N(\mu, \sigma^2)$, with σ^2 known.

We want to construct a confidence interval for μ of the form $[m(x_1, \dots, x_n); n(x_1, \dots, x_n)]$ of level $1 - \alpha$: i.e. we aim at finding m and n such that

$$P_{\mu}(\mu \in [m(X_1, \dots, X_n); n(X_1, \dots, X_n)]) = 1 - \alpha$$

For all $\mu \in \mathbb{R}$ $G(x_1, \dots, x_n, \mu) = \frac{\sqrt{n}}{\sigma} (\bar{x}_n - \mu) \sim \mathcal{D}(0, 1)$ thanks to the central limit theorem.

Let F the cumulative distribution function of $N(0,1)$ distribution. For all $a \in]0,1[$, let

z_α be the α -quantile of Φ given by $\Phi(z_\alpha) = \alpha$.

In particular $z_{1-\frac{\alpha}{2}} = 1,96$ for $\alpha = 0,05$ and $z_{1-\frac{\alpha}{2}} = 3$ for $\alpha = 0,01$

This yields: $P_{\mu}\left(\frac{\sqrt{n}}{\sigma}(\bar{X}_n - \mu) \in \left[-z_{1-\frac{\alpha}{2}}; z_{1-\frac{\alpha}{2}}\right]\right) = 1 - \alpha$

Note that $\mathcal{I}_n = [-z_{n-\frac{n}{2}}; z_{n-\frac{n}{2}}]$ is independent of μ

therefore
$$J_{n,1} = \left[\bar{X}_n - \frac{\sigma}{\sqrt{n}} Z_{1-\frac{\alpha}{2}} ; \bar{X}_n + \frac{\sigma}{\sqrt{n}} Z_{1-\frac{\alpha}{2}} \right]$$
 is the confidence interval for μ and $|J_{n,1}| \xrightarrow{n \rightarrow \infty} 0$

$$G(x_1, \dots, x_n, \mu) = \frac{\sqrt{n}(\bar{x}_n - \mu)}{S_n^2} \text{ where } S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$$

therefore $\forall A \subset \mathbb{R}$ $\sigma = (\mu, \sigma^2)$

$$P_0(G(x_1, \dots, x_n, y) \in A) = \text{tr}_1(A)$$

Let t_{n-1} the α quantile of t_n . Thanks to the symmetry of the Student's distribution, confidence intervals symmetric as well and this writes

$$P_0(-t_{1-\frac{\alpha}{2}}(n-1) \leq G(X_1, \dots, X_n, \mu) \leq t_{1-\frac{\alpha}{2}}(n-1)) = 1-\alpha$$

this leads to the following interval for μ :

$$m(X_{n-1}, X_n) = \bar{X}_n - \frac{s_0}{\sqrt{n}} t_{1-\frac{\alpha}{2}}(n-1) \quad \text{and} \quad u(X_{n-1}, X_n) = \bar{X}_n + \frac{s_0}{\sqrt{n}} t_{1-\frac{\alpha}{2}}(n-1)$$

Remark: If the sample is small $n \leq 100$ then it is important to consider the Student's distribution. However as soon as $n \geq 100$ then the Student's distribution is very close to a normal one and one can use the normal quantiles.

6) Confidence Interval for the variance:

Let again (X_1, \dots, X_n) iid $\mathcal{D}(\mu, \sigma^2)$ with unknown μ and σ^2 .

$$\text{betr. } G(x_1, \dots, x_n, \sigma^2) = \frac{(n-1)}{\sigma^2} S_n^2$$

We know that $\frac{(n-1)}{s^2} S^2 \sim \chi^2_{n-1}$. let $\chi^2_{\alpha(n-1)}$ the α quantile of χ^2_{n-1} and

fix $d_1 + d_2 = d$ then: $\forall (\mu, \sigma^2)$

$$P_{\text{true}}(X_{\text{dev}}^2(n-1) \leq G(X_n, \dots, X_n, \sigma^2) \leq X_{1-\alpha_2}(n-1)) = 1 - \alpha$$

this yields
$$I_{\alpha, n} = \left[\frac{n-1}{\chi^2_{n-\alpha, (n-1)}} S_0^2 ; \frac{n-1}{\chi^2_{\alpha, (n-1)}} S_n^2 \right]$$
 is a confidence interval of level $1-\alpha$ for σ^2

Note that the length of this interval is random as it depends on S_n^2

