

Practical n°3

stoehr@ceremade.dauphine.fr

Notations and Objectives. Let X_1, \dots, X_n be a sequence of *i.i.d.* random variables with an **unknown** cumulative distribution function F and denote x_1, \dots, x_n realisations of these random variables. The practical focuses on the empirical distribution of data \hat{F}_n . The first objective is to study the properties of \hat{F}_n as an estimate of F for a given sample x_1, \dots, x_n . The second objective is to sample from it.

Exercise. A marketing brand studies the time spent by customers on its website. The file "time.csv" contains the duration (in seconds) for $n = 100$ customers.

1. Download "time.csv" in a folder "TP3" and import the data in R using the command `read.table` or `read.csv`.

The empirical distribution function. We are first interested in getting some information about the underlying unknown distribution F the data are coming from.

Definition

The empirical distribution function is a random function (as function of the random variables X_1, \dots, X_n) associated to the empirical measure and given for $y \in \mathbb{R}$ by

$$\hat{F}_n(y) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq y\}} = \frac{\text{Card}\{X_i \mid X_i \leq y\}}{n}.$$

Its evaluation on a dataset x_1, \dots, x_n brings information about the distribution contains in the data itself.

2. Write a function `cdf_hat(y, x)` which returns the empirical distribution at point y given a n -sample $\mathbf{x} = (x_1, \dots, x_n)$.

\hat{F}_n is an estimate of the cumulative distribution F that generated the data. The first major result is that it converges pointwise.

Theorem. Strong Law of Large Numbers (SLLN)

For all y , $\hat{F}_n(y)$ converges pointwise to the true cumulative distribution function $F(y)$ as almost surely (*a.s.*), *i.e.*,

$$\hat{F}_n(y) \xrightarrow[n \rightarrow +\infty]{} F(y) = \mathbb{P}[X_1 \leq y] \quad \textit{a.s.}$$

3. Let assume the brand is interested in the probability that a customer spend less than 50 seconds on its website.
- Compute `cdf_hat(50, x)`.
 - Write a function `cv_cdf_hat(y, x)` which returns $(\hat{F}_1(y), \hat{F}_2(y), \dots, \hat{F}_n(y))$ for a given point y and a given sample $x = (x_1, \dots, x_n)$.
 - Using a scatterplot, check if $\hat{F}_n(50)$ converged or not with our sample size n .

\hat{F}_n is a random function and hence has a variability depending on the data it is based on, *i.e.*, \hat{F}_n changes when we use other data. The randomness of \hat{F}_n can be asymptotically characterised pointwise.

Theorem. Central Limit Theorem (CLT)

For all $y \in \mathbb{R}$, $\hat{F}_n(y)$ has asymptotically normal distribution with the standard \sqrt{n} rate of convergence, *i.e.*,

$$\sqrt{n} [\hat{F}_n(t) - F(t)] \xrightarrow[n \rightarrow +\infty]{d} \mathcal{N}(0, F(t)\{1 - F(t)\})$$

Remark. This result is a simple application of the Central Limit theorem to the sequence of *i.i.d.* random variables $(\mathbb{1}_{\{x_n \leq t\}})_{n \geq 1}$, that is a sequence of Bernoulli variables with parameter equal to $F(t)$.

4. (a) Show that

$$\sqrt{\frac{n}{\hat{F}_n(t)(1 - \hat{F}_n(t))}} [\hat{F}_n(t) - F(t)] \xrightarrow[n \rightarrow +\infty]{d} \mathcal{N}(0, 1).$$

- (b) Assume that the asymptotic regime is valid. Compute the asymptotic confidence interval for $F(50)$ for 0.95 confidence level. Comment the result.

We actually have more than the pointwise convergence of the empirical distribution to the true cumulative distribution function.

Theorem. Glivenko-Cantelli theorem

\hat{F}_n converges uniformly to F over \mathbb{R} , *i.e.*,

$$\|\hat{F}_n - F\|_{\infty} \xrightarrow[n \rightarrow +\infty]{} 0 \quad a.s.$$

5. (a) Write a function `cdf_hat_vec(y, x)` which returns the empirical distribution function at multiple points $y := (y_1, \dots, y_k)$ given a n -sample $x = (x_1, \dots, x_n)$.
- (b) Check if \hat{F}_n converges uniformly to the Gamma distribution with shape parameter $\alpha = 1.9$ and rate parameter $\beta = 0.04$.

We can further characterized the asymptotic behaviour of the sup-norm between \hat{F}_n and F . We also have a rate of convergence of \sqrt{n} and the Dvoretzky–Kiefer–Wolfowitz inequality provides bound on the tail probabilities.

Theorem. Dvoretzky–Kiefer–Wolfowitz inequality

For all $\varepsilon \geq \sqrt{\frac{1}{2n} \log 2}$,

$$\mathbb{P}[\|\widehat{F}_n - F\|_\infty > \varepsilon] \leq \exp(-2n\varepsilon^2).$$

As per this inequality, we have another confidence interval for the true cumulative distribution function value $F(y)$, $y \in \mathbb{R}$. Given $\alpha \in]0, 1[$,

$$\mathbb{P}[\max\{0, \widehat{F}_n(y) - \varepsilon\} \leq F(y) \leq \min\{\widehat{F}_n(y) + \varepsilon, 1\}] \geq 1 - \alpha, \quad \text{where } \varepsilon = \sqrt{\frac{1}{2n} \log \frac{2}{\alpha}}.$$

6. Compute the confidence interval based on the Dvoretzky–Kiefer–Wolfowitz inequality for $F(50)$ and 0.95 confidence level. Comment the result.

Sampling from the empirical distribution function. Another key feature of the empirical distribution is that it is easy to sample from. The empirical distribution defines a discrete probability distribution on sample space $\{x_1, \dots, x_n\}$. Put in other words, if X^* follows \widehat{F}_n , then X^* is taking values in $\{x_1, \dots, x_n\}$ with probabilities

$$\mathbb{P}[X^* = x_i] = \frac{1}{n}, \quad \text{for } i \in \llbracket 1, n \rrbracket.$$

Remark. Above we did not take into account possible repetitions of values in x_1, \dots, x_n . We consider that observations are ordered elements and that indices for the probability vector match indices of the observation. For instance, if we have observations $(x_1, x_2, x_3) = (1, 2, 1)$, the probability vector is $(p_1, p_2, p_3) = (1/3, 1/3, 1/3)$.

Property

In order to get a m -sample from \widehat{F}_n , we sample m times from $\{x_1, \dots, x_n\}$ with replacement according to the vector of probabilities (p_1, \dots, p_n) . The re-sampled sample of size m is denoted (x_1^*, \dots, x_m^*) and the associated empirical distribution function \widehat{F}_n^* .

7. (a) For each of the following values, $m = 10$, $m = 50$ and $m = 100$, generate 100 samples (x_1^*, \dots, x_m^*) from \widehat{F}_n .
- (b) Compute the percentile of each sample.
- (c) Why is it important to have a resample of the same size as the original sample?