

Th4: The Glivenko–Cantelli theorem, Proof, Simulations Glivenko–Cantelli

In the theory of probability, the Glivenko–Cantelli theorem (sometimes referred to as the Fundamental Theorem of Statistics), named after Valery Ivanovich Glivenko and Francesco Paolo Cantelli, determines the asymptotic behavior of the empirical distribution function as the number of independent and identically distributed observations grows.

The uniform convergence of more general empirical measures becomes an important property of the Glivenko–Cantelli classes of functions or sets.

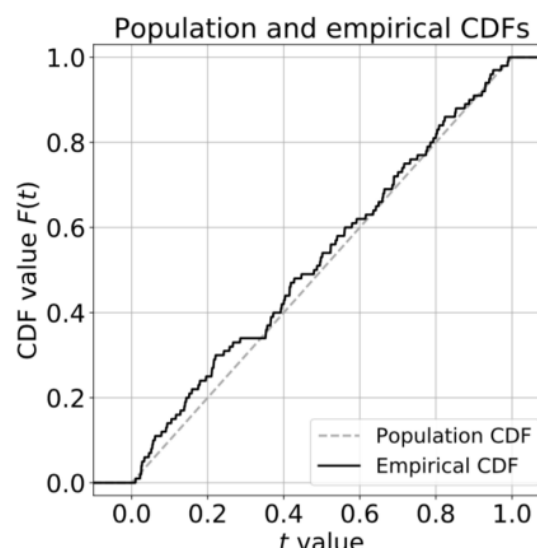
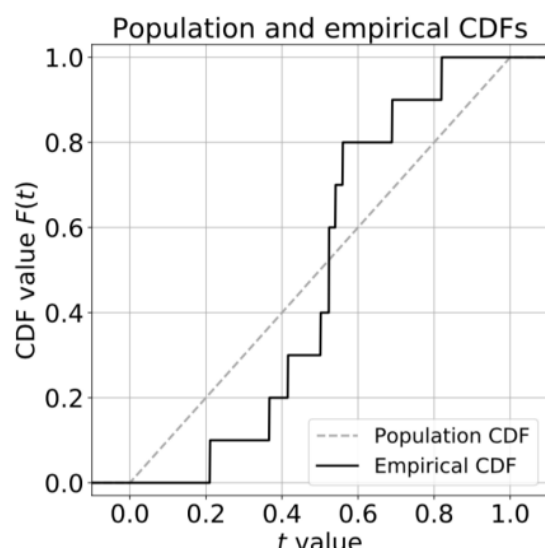
In the most general statistical setup, we have some unknown probability distribution μ , and we try to estimate properties of it by taking independent, identically distributed samples X_1, X_2, \dots, X_n with distribution μ (written $X_i \sim \mu$). The general hope is that as we take enough samples ($n \rightarrow \infty$), we can recover properties of the distribution μ .

The hope and general assumption of statistics is that it is possible to estimate any property of the distribution μ consistently. A distribution μ on \mathbb{R} specifies a cumulative distribution function (CDF).

$$F_\mu(t) := \mathbb{P}_\mu(X \leq t).$$

The CDF uniquely specifies the distribution it arises from, as well.

We can see this approximation in the following graph



We have $\mu((a, b]) = F_\mu(b) - F_\mu(a)$, which determines the value of μ by the Caratheodory outer measure construction. So if we can consistently estimate F_μ with the data X_1, \dots, X_n as $n \rightarrow \infty$, we should be able to estimate any property of the distribution μ . The Glivenko-Cantelli theorem says that this estimation of the entire distribution is indeed possible.

So in simpler words we can say that the Glivenko-Cantelli (GC) Theorem gives the relationship between the empirical ("data-driven") CDF and the true CDF. A CDF gives the probability that a random variable takes value less than a specified number. The empirical CDF estimates this

probability with the proportion of data less than a specified number. The GC theorem says that the difference between the empirical ("data driven") CDF and the true CDF goes away as you get a larger sample size.

the Glivenko-Cantelli theorem states the following:

With probability one

$$D_n := \sup_{-\infty < t < \infty} |F_n(t) - F(t)| = \sup_{t \text{ rational}} |F_n(t) - F(t)| \longrightarrow 0$$

The proof of the Glivenko-Cantelli theorem is the following:

If $-\infty \leq s \leq t < u \leq \infty$, then

$$F_n(s) - F(s) - (F(u-) - F(s)) \leq F_n(t) - F(t) \leq F_n(u-) - F(s) = (F_n(u-) - F(u-)) + (F(u-) - F(s)).$$

Thus if S is a finite subset of \mathbb{R} , say, $S = \{s_1, \dots, s_k\}$ with $s_1 < s_2 < \dots < s_k$, with associated partition

$$\pi_S := \{(-\infty =: s_0, s_1), [s_1, s_2), \dots, [s_{k-1}, s_k), [s_k, s_{k+1} := \infty)\},$$

we have

$$1) \quad D_n = \sup_t |F_n(t) - F(t)| \leq \max_{s \in S} (|F_n(s) - F(s)| \vee |F_n(s-) - F(s-)|) + w_{\pi_S}(F)$$

where

$$2) \quad w_{\pi_S}(F) := \max_{1 \leq j \leq k+1} |F(s_j-) - F(s_{j-1})|.$$

Now observe that there exists a sequence $(S_m)_{m \geq 1}$ of finite subsets $S_m \subset \mathbb{R}$ such that $w_{\pi_{S_m}}(F) \rightarrow 0$ as $m \rightarrow \infty$;

This is clear from inspection of a diagram I'll draw in class. Put $S_\infty := \bigcup_{m=1}^\infty S_m$. S_∞ is countable, and countably many applications of the SLLN imply that there exists a P -null set N off of which

$$3) \quad |F_n(s) - F(s)| \vee |F_n(s-) - F(s-)| \rightarrow 0 \text{ for all } s \in S_\infty$$

Now from (1), (2), (3) we get $D_n \rightarrow 0$ off of N , as desired.

A simulation of the Glivenko-Cantelli theorem can be found on this [page](#).

Bibliography

-https://en.wikipedia.org/wiki/Glivenko-Cantelli_theorem

-<https://pillowmath.github.io/Expository%20Notes/VC-Dimension-and-Glivenko-Cantelli-Notes.pdf>

-https://www.ams.jhu.edu/~fill/550.621/prob_two/GCthm.pdf