

Glivenko-Cantelli Theorem

Chiara Iannicelli

November 2023

1 Introduction

The Glivenko-Cantelli theorem is a fundamental result in probability theory and mathematical statistics. It provides insights into the convergence of empirical distribution functions to the true underlying distribution as the sample size increases.

The theorem is named after mathematicians Dmitri Glivenko and Francesco Paolo Cantelli and is particularly important in the context of understanding how well the sample distribution function (empirical distribution function) approximates the true distribution function.

2 Statement

Assume that X_1, X_2, \dots are independent and identically distributed random variables in R with common cumulative distribution function $F(x)$. The empirical distribution function for X_1, X_2, \dots, X_n is defined by:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[X_i, \infty)}(x) = \frac{1}{n} |\{i | X_i \leq x, 1 \leq i \leq n\}|$$

where I_C is the indicator function of the set C . For every (fixed) x , $F_n(x)$ is a sequence of random variables which converge to $F(x)$ almost surely by the strong law of large numbers. Glivenko and Cantelli strengthened this result by proving uniform convergence of F_n to F .

3 Theorem

The theorem states as follows:

$$\|F_n - F\|_\infty = \sup_{x \in R} |F_n(x) - F(x)| \longrightarrow 0 \quad \text{almost surely}$$

This theorem originates with Valery Glivenko[4] and Francesco Cantelli,[5] in 1933.

This is an uniform law of large numbers:

$$\begin{aligned}\|F_n - F\|_\infty &= \sup_x |F_n(x) - F(x)| \\ &= \sup_x |P_n(X \leq x) - P[X \leq x]| \\ &\xrightarrow{as} 0,\end{aligned}$$

Figure 1: uniformity

where P_n is the empirical distribution that assigns mass $\frac{1}{n}$ to each X_i .

The law of large numbers says that, for all x , $P_n(X \leq x) \rightarrow P(X \leq x)$. The GC Theorem says that this happens uniformly over x .

4 Simulation

An example of code to simulate the behaviour of the Glivenko-Cantelli Theorem is the following:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 # True distribution function (CDF) for a uniform distribution
5 def true_distribution(x):
6     return np.where(x < 0, 0, np.where(x <= 1, x, 1))
7
8 # Generate random samples from a uniform distribution
9 np.random.seed(42) # Set seed for reproducibility
10 sample_size = 1000
11 samples = np.random.uniform(0, 1, sample_size)
12
13 # Calculate the empirical distribution function (EDF)
14 def empirical_distribution(sample, x):
15     return np.sum(sample <= x) / len(sample)
16
17 # Calculate true distribution values
18 x_values = np.linspace(0, 1, 1000)
19 true_values = true_distribution(x_values)
20
21 # Calculate EDF values
22 edf_values = [empirical_distribution(samples, x) for x in x_values]
```

```

23
24 # Plot the true distribution and EDF
25 plt.plot(x_values, true_values, label='True Distribution')
26 plt.step(x_values, edf_values, label='Empirical Distribution
    ↪ Function', where='post')
27 plt.xlabel('X-axis')
28 plt.ylabel('Cumulative Probability')
29 plt.title('Simulation of Glivenko-Cantelli Theorem')
30 plt.legend()
31 plt.show()
32

```

This generates the following graphic:

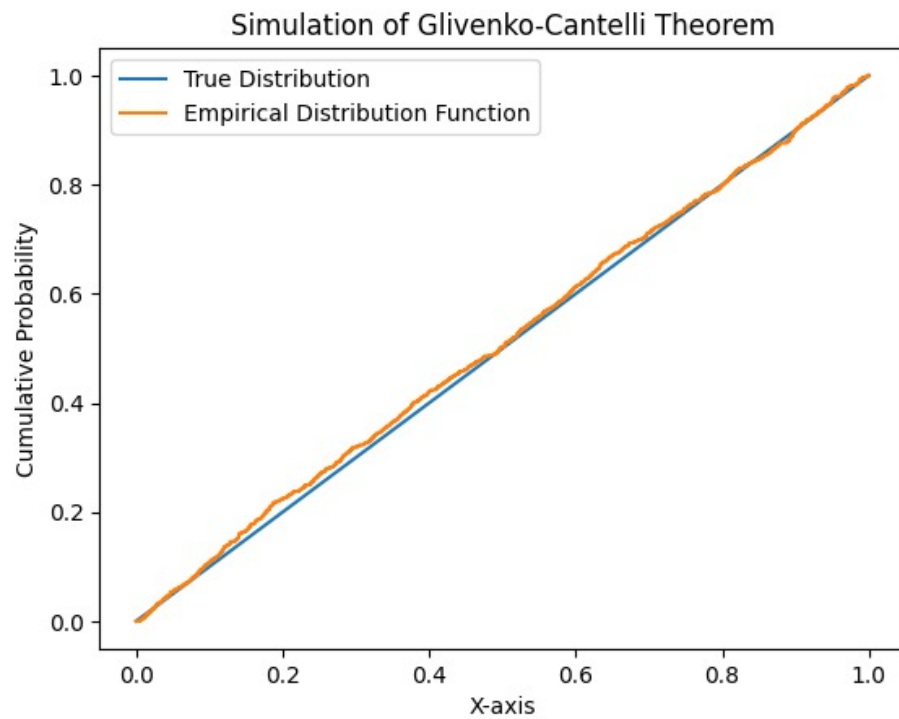


Figure 2: From colab