

Empirical cdf:

Assume that X_1, X_2, \dots, X_n are randomly sampled from a population whose cdf is the continuous function $F(x)$

The estimate of $F(x)$ is called the empirical cdf:

$\hat{F}(x)$ = fraction of observations $\leq x$

$$\hat{F}(x) = P(X \leq x)$$

The empirical cdf is a step function that takes a step at each observed data value.

If the data points are distinct, the size of each step is $1/n$, and if k data points have the value x , the step size is k/n .

Example 1.2.1

We are interested in the number of on-off cycles of a mechanical device before the device fails. The hypothetical data in Table 1.2.1 are the number of cycles (in thousands) that it takes for 20 door latches to fail, and the empirical cdf for the data is in Figure 1.2.1

If the data points are distinct, the size of each step is $1/n$, and if k data points have the value x , the step size is k/n .

Standard Deviation

If $\hat{F}(x)$ is the fraction of observations $\leq x$, then what is $SD(\hat{F}(x))$?

$$P(X \leq x) = \frac{a}{b} \quad SD(\hat{F}(x)) = \sqrt{\frac{\frac{a}{b}(1 - \frac{a}{b})}{n}}$$

Hint: What is the distribution of the number of observations for which $X_i \leq x$? Let's call that $S = \sum I(X_i \leq x)$ What is the $SD(S)$?

Hint: $p = P(X_i \leq x) = F(x)$

Hint: Remember your rules of variances

$$S \sim \text{Bin}(n, p)$$

$$X \Rightarrow \sigma_x^2 = \sigma^2$$
$$bX \Rightarrow b^2 \sigma_x^2$$

$$\mu_S = n \cdot p$$
$$SD_S = \sqrt{n p (1-p)}$$
$$= \sqrt{\frac{n F(x) (1-F(x))}{n^2}}$$

Utilizing approximate normality of the binomial distribution

An approximate $100(1-\alpha)$ confidence interval for $F(x)$ is given by:

$$\hat{F}(x) \pm z_{(1-\alpha/2)} \sqrt{\frac{\hat{F}(x)(1-\hat{F}(x))}{n}}$$

point estimate critical value standard error

margin of error

Example 1.2.2

From the door latch data, 14 out of the 20 observations are less than or equal to 45.

What is the approximate 90% confidence interval for $F(45)$?

Is this interval particularly helpful? No. Want to have as narrow CI as possible, without losing

How do we make inferences for the $(100p)$ th percentile θ_p ?

Confidence interval is of the form: $X_{(a)} < \theta_p < X_{(b)}$

A and b are chosen so that: $P(X_{(a)} < \theta_p < X_{(b)}) = 1 - \alpha$

where $1 - \alpha$ is the desired probability that the interval captures θ_p .

Since observations are independent, the probability that at least a and at most $b - 1$ of the observations fall less than θ_p is given by the binomial probability with $p = p$:

$$\sum_{i=a}^{b-1} \binom{n}{i} p^i (1-p)^{n-i}$$

Without using the continuity correction for this approximation, we may obtain a and b by solving for them in the equations:

$$\frac{a - np}{\sqrt{np(1-p)}} = -z_{(1-\alpha/2)} \quad \frac{b - 1 - np}{\sqrt{np(1-p)}} = z_{(1-\alpha/2)}$$

and rounding to the nearest integers. ^ symmetric

CLT Test vs. Binomial Test: How to Compare?

Type I Error: Occurs if H_0 is rejected when it is true (think $\alpha = .05$)

Power: The probability of rejecting H_0 if it is false.

If both tests have the correct Type I error (yield the Type I error that it claims), then the one with the greater power is the preferred test.

In this setting, let's consider the hypothesis test:

$H_0 : \mu = \mu_0$, $H_a : \mu > \mu_0$

If σ is known (or we have a large sample size and known that σ is finite), then:.....

CLT Test: If conditions are met to apply the CLT, then Type I errors found by using the standard normal distribution will be approximately correct. Many studies have shown that the approximation is quite good even for moderate sample sizes and a range of population distributions.

UMP (Uniformly Most Powerful)