

Methods that assume a particular form for the distribution of a r.v. and try to estimate the parameters of that particular distribution

from a "random sample,"
i.e. i.i.d. r.v.'s X_1, \dots, X_n
can be called parametric methods.

Here, we focus on the so-called non-parametric methods, that are generic and do not assume parametric forms for probability

distributions.

Empirical Distribution

Definition: Assume that X_1, X_2, \dots, X_n are i.i.d. samples with common CDF F_X . The empirical distribution function of X_1, X_2, \dots, X_n is

defined as

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[X_i, +\infty)}(x)$$

where I_A is the indicator function of set A.

Informally, the empirical Cdf

is increased by $\frac{1}{n}$ at each

X_i .

Goodness -of- fit

To measure how well a fitted distribution resembles

the sample data, we can

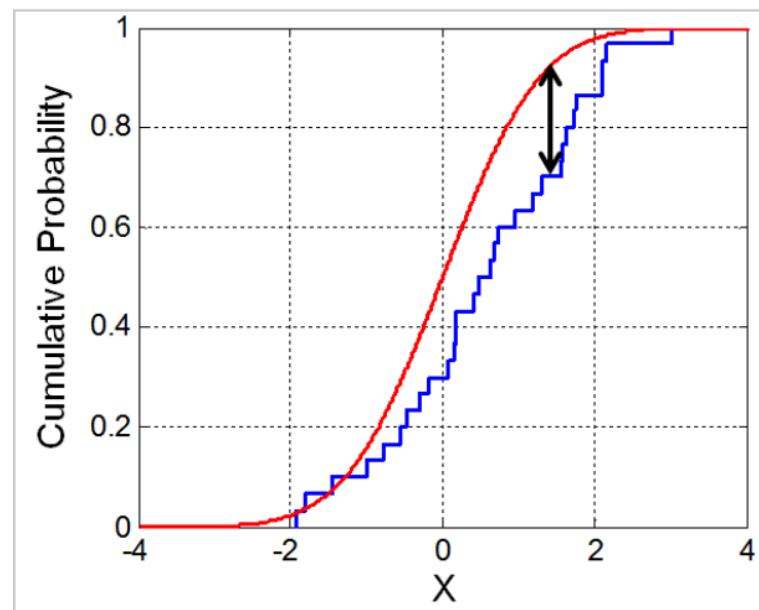
use the Kolmogorov-Smirnov (K-S) test statistic :

$$D_n = \sup_x | \hat{F}_n(x) - F(x) |$$

empirical
distribution

fitted
distribution

$$D_n = \sup_x |F_n(x) - F(x)|$$



The Golivenko - Cantelli Theorem

Theorem: Assume that X_1, X_2, \dots, X_n are i.i.d. samples with common cdf $F_X(x)$ and the empirical cdf $\hat{F}_n(x)$. Then the

Kolmogorov-Smirnov Statistic

$$D_n = \sup_x |\hat{F}_n(x) - F(x)|$$

Converges to Zero almost surely

for all x , i.e.

$$\text{P}(\lim_{n \rightarrow \infty} D_n \neq 0) = 0.$$

Remark: $\hat{F}_n(x)$ is a random variable at each x . D_n is also a r.v.

Remark: The G-K Theorem is sometimes called the

Fundamental theorem of statistics, and is a basis for the Kolmogorov-Smirnov test.

Assume that multiple distributions are estimated from i.i.d samples using different (parametric or non-parametric) methods. Which distribution is "the best?"

The Kolmogorov-Smirnov Test Criterion says that the KS statistic has to be computed for each distribution and the distribution with the minimum D_n has to be selected.

Goodness of fit as hypothesis

Testing

H_0 : All samples X_1, \dots, X_n

Come from the cdf F

H_1 : At least one of the

samples X_1, \dots, X_n does not

Come from the cdf F

Under the null hypothesis

(assuming H_0 is true)

the K-S test statistic

D_n follows a K-S

distribution, which has

a complicated form.

Therefore, to reject the null,

we have to observe a

sufficiently large D_n so that

we are in the rejection

region.

K-S distribution has a complicated form.

$$P\left(D_n < \frac{1}{2n} + v\right) = n! \int_{\frac{1}{2n}-v}^{\frac{1}{2n}+v} \int_{\frac{3}{2n}-v}^{\frac{3}{2n}+v} \cdots \int_{\frac{2n-1}{2n}-v}^{\frac{2n-1}{2n}+v} g(u_1, u_2, \dots, u_n) du_n \dots du_2 du_1 \quad 0 \leq v \leq \frac{2n-1}{2n}$$

$g(u_1, u_2, \dots, u_n) = 1$ over $0 \leq u_1 \leq u_2 \leq \dots \leq u_n \leq 1$,

And 0 otherwise

Birnbaum, Z.W. (1952), "Numerical Tabulation of the Distribution of Kolmogorov's Statistic for Finite Sample Size," Journal of the American Statistical Association, 47, 425–441

$$F_{D_6}(t) = \begin{cases} 0 & t < \frac{1}{12} \\ 46080t^6 - 23040t^5 + 4800t^4 - \frac{1600}{3}t^3 + \frac{100}{3}t^2 - \frac{10}{9}t + \frac{5}{324} & \frac{1}{12} \leq t < \frac{1}{6} \\ 2880t^6 - 4800t^5 + 2360t^4 - \frac{1280}{3}t^3 + \frac{235}{9}t^2 + \frac{10}{27}t - \frac{5}{81} & \frac{1}{6} \leq t < \frac{1}{4} \\ 320t^6 + 320t^5 - \frac{2600}{3}t^4 + \frac{4240}{9}t^3 - \frac{785}{9}t^2 + \frac{145}{27}t - \frac{35}{1296} & \frac{1}{4} \leq t < \frac{1}{3} \\ -280t^6 + 560t^5 - \frac{1115}{3}t^4 + \frac{515}{9}t^3 + \frac{1525}{54}t^2 - \frac{565}{81}t + \frac{5}{16} & \frac{1}{3} \leq t < \frac{5}{12} \\ 104t^6 - 240t^5 + 295t^4 - \frac{1985}{9}t^3 + \frac{775}{9}t^2 - \frac{7645}{648}t + \frac{5}{16} & \frac{5}{12} \leq t < \frac{1}{2} \\ -20t^6 + 32t^5 - \frac{185}{9}t^3 + \frac{175}{36}t^2 + \frac{3371}{648}t - 1 & \frac{1}{2} \leq t < \frac{2}{3} \\ 10t^6 - 38t^5 + \frac{160}{3}t^4 - \frac{265}{9}t^3 - \frac{115}{108}t^2 + \frac{4651}{648}t - 1 & \frac{2}{3} \leq t < \frac{5}{6} \\ -2t^6 + 12t^5 - 30t^4 + 40t^3 - 30t^2 + 12t - 1 & \frac{5}{6} \leq t < 1 \\ 1 & t \geq 1. \end{cases}$$



Fortunately, there are tables

for the critical values of

the distribution of K-S

test.

K-S is distribution free, i.e. it does not assume any particular distribution for the data.

Drawback: Only for one dimensional distributions

There are other goodness of fit tests including the Anderson-Darling test.

The Chi-Squared Goodness of Fit Test

This test is based on the histogram of data. The idea is that if the pmf or

pdf ~~of~~ estimated for a sample is too different from its histogram, it probably is not a good fit for the data.

In particular, if

$P(e_j \leq X_i \leq e_{j+1})$ calculated

according to the estimated

pmf/pdf is not close to

the histogram estimate

of $P(e_j \leq X_i \leq e_{j+1})$,

which is $\frac{h_j}{n}$, the

estimated pmf/pdf is

not a good fit.

It can be shown that

the quantity

$$\chi^2 = \sum_{j=1}^m \left(\frac{|H_j - np_j|^2}{np_j} \right)$$

has a Chi-squared distribution

with $m-1$ degrees of freedom

when n is large.

This is because

$$H_j - np_j$$

$$\sqrt{np_j}$$

has a standard normal distribution

when n is large.

Therefore, if

$$\chi^2 > \chi^2_{m-1, \alpha}$$

(i.e. if the discrepancy is too large), we reject the null

that the pdf approx we

found is a good fit for

the data.

