

ERASMUS B4 Course

Winter Semester 2020/2021, 21 JANUARY

Subject: topics in statistical inference for given general population with the feature X .

Introduction

Let the general population P with the feature X be given.

We know that the main goal of M.S. is to establish the probability distribution of X .

To do this, we need the result of observations of P . It means, that we have to choose from P its representation \tilde{P} and examine it.

In M.S. we assume that the result of the above problem has a form

$$\mathbb{R}^n \rightarrow (x_1, x_2, \dots, x_n) = (X_1, X_2, \dots, X_n)(v_0),$$

where $d(X_i) = d(X)$, $i = 1 \dots n$,

and X_1, \dots, X_n are stoch. indep.

Then we say that we have the SIMPLE SAMPLE.

The last lecture showed us that the presentation of S.S. is very important in statistical inference.

It is clear that S.S. describes P_0 CIP. We are going to explain below that the information for P_0 can be extrapolated to the whole IP.

The concept of Empirical Cumulative Prob. Dist. (ECPD)

Fix the S.S.

$$(x_1, x_2, \dots, x_n) = (X_1, X_2, \dots, X_n)(v_0).$$

The function

$$R \rightarrow t \longrightarrow F_n(t, v_0) = \frac{1}{n} \# \{ i : x_i < t \}$$

is called ECPD corresponding to P with X for given S.S.

We know, that ECPD we can obtain from histogram (or frequency diagram) and vice versa. Moreover, the graph of ECPD is a "step function".

Why ECPD is so important in M.S.?

It is explained by the following theorem

Theorem (on ECPD) (GNIEDENKO)

For given G.P. P with the feature X let us consider the functions

$$\mathbb{R} \times \Omega \rightarrow (t, \omega) \rightarrow F_m(t, \omega) \stackrel{\text{def}}{=} \frac{1}{m} \#\{i : X(\omega) < t\},$$

when for $n \geq 2$

X_1, X_2, \dots, X_n are st. independent with the same distribution X.

If F_X is (unknown) a cumulative prob. dist. f. corresponding to X, then for the events

$$A_+ = \left\{ \omega \in \Omega : F_n(t, \omega) \rightarrow F_X(t) \right\}, \quad t \in \mathbb{R}$$

We have

$$\underline{P(A_+)} = 1, \quad t \in \mathbb{R}.$$

Remarks

1^v. Each member $\bar{F}_m(t, u)$ of the above sequence is a generalisation of ECDF, namely

$$\bar{F}_n(t, w_0) = \bar{F}_m(t, u) \Big|_{u=w_0}.$$

So if we take $(\bar{F}_n(t, u))_{n \geq 2}$ it means we consider (from theoretical point of view!) all S.S. (with respect of m and $w \in \Omega$).

2^v. Since for given t , $P(A_t) = 1$ we can assume Ω cl. from the definition of S.S. belongs to A_t .

Therefore we have an approximation rule

$$\bar{F}_m(t, w_0) \approx F_X(t) \quad |$$

which is the basic argument in extrapolation problem.

The concept of estimation

The stage of the statistical inference using the concept of ECDF generally does not solve the main problem of M.S.

On the output of this stage the typical solution is as follows :

$F_X(t, \theta)$ - we know the "shape" of F_X , but the value of the parameter θ is unknown.

So, we need the next stage of S.I., to ~~estimate~~ establish the value θ .

It is called the ESTIMATION PROCEDURE.

The main idea of E.P. says :

For given $\alpha \in (0, 1)$ (signifiant level),
for instance $\alpha = 0.05$ or 0.1

We try to construct of two r.v. Z_1, Z_2
such that:

(i) Z_1, Z_2 is generated by using
the random vector (X_1, X_2) which
defines the S.S.

(ii) $Z_1(v) < Z_2(v)$ with probability 1

→ (iii) $P(\text{Real: } Z_1(v) < \theta < Z_2(v)) = 1-\alpha$.

Then we say we have a random interval (Z_1, Z_2) .

Now, if we use $\omega_{0,1}^n$ from the S.S.

We obtain the real interval $(Z_1(v_0), Z_2(v_0)) \subset \mathbb{R}$,

and finally we have the following conclusion

$\theta \in (Z_1(v_0), Z_2(v_0))$ with prob. $1-\alpha$

Example

Suppose that we know (for instance from analysis of histogram) from the S.I.

$(X_{n-}, X_n) = (X_{n-}, X_n)(U_0)$, if
 $X \in N(m, \sigma^2)$, and the value σ^2 is fixed.

If we fix $\alpha \in (0, 1)$ and we take

$$Z_1 \stackrel{df}{=} \bar{X}_n - n\alpha \frac{\sigma}{\sqrt{n}}, \quad Z_2 = \bar{X}_n + n\alpha \frac{\sigma}{\sqrt{n}},$$

where $\bar{X}_n = \frac{1}{n}(X_{n-} + X_n)$, we can find $n\alpha$ if we resolve the equation

$$(\#) \quad P(\text{event: } Z_1(u) < m < Z_2(u)) = 1 - \alpha.$$

Namely,

$$(\#) \Leftrightarrow P\left(\{v \mid \bar{X}_n(v) - n\alpha \frac{\sigma}{\sqrt{n}} < m < \bar{X}_n(v) + n\alpha \frac{\sigma}{\sqrt{n}}\}\right) = 1 - \alpha$$

\Leftrightarrow

$$\text{Pf (contd): } -\alpha < \frac{\bar{X}_n(\alpha) - m}{\frac{\sigma}{\sqrt{n}}} < \alpha \Leftrightarrow 1 - \alpha.$$

But, according to our assumption of X ,

$$\frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}} \in N(0, 1), \text{ so}$$

$$(\#1) \Leftrightarrow \Phi(\alpha) - \Phi(-\alpha) = 1 - \alpha \Leftrightarrow$$

$$2\Phi(\alpha) - 1 = 1 - \alpha \Leftrightarrow$$

$$\boxed{\Phi(\alpha) = 1 - \frac{\alpha}{2}}$$

$$\text{For instance, for } \alpha = 0.05, \underline{m_\alpha = 1.96}.$$

The above ~~process~~ shows, that by using Z_1, Z_2 given as above, we can construct the ~~standard~~ confidence interval and ~~finally~~ finally determine E.P.

Let's realize the following numerical simulation.

$$\rightarrow (x_1, x_2, x_3, x_4) = (2, 0), 2, 12, 1, 97, 2, 07)$$

for $X \in N(m, (0, 0.25)^2)$ and $\alpha = 0, 0.5$.

Since for $\alpha = 0.05$, $n\alpha = 1.96$, we have

$$Z_1(v) = \bar{X}_h(v) - 1.96 \frac{0.25}{2}$$

$$Z_2(v) = \bar{X}_h(v) + 1.96 \frac{0.25}{2}, v \in V.$$

so, for $v = v_0$

$$\bar{X}_h(v_0) = 2, 03 \text{ and}$$

$$Z_1(v_0) = 2, 03 - 1.96 \frac{0.25}{2} = 1, 838$$

$$Z_2(v_0) = 2, 03 + 1.96 \frac{0.25}{2} = 2, 228.$$

Therefore, in this case

$m \in (1, 838, 2, 228)$ with prob. $0, 95$

Conclusion

The ideas presented above belong to many others that make up the stochastic inference procedure (S.I.P.).

The main goal of S.I.P. is to establish the probability distribution of the feature X for given general population P .

But to present this method, you need another course that I invite you to take in the future.

RK