

Statistică - Cursul 9

Olariu E. Florentin

Aprilie, 2016

Table of contents

1 Momentele și forma distribuției

2 Teoremele fundamentale

- Legea numerelor mari

 - Inegalitățile lui Markov și a lui Cebâșev revăzute

 - Teorema lui Cebâșev

 - Legea numerelor mari

- Teorema limită centrală

 - Aproximarea normală a distribuției binomiale

3 Statistică inferențială

- Estimarea parametrilor

 - Estimarea punctuală

 - Estimarea cu intervale

 - Intervale de încredere pentru medie

4 Bibliografie

Momentele

Definition 0.1

Momentul central de ordin k al populației este

$$\mu_k = M \left[(X - \mu)^k \right].$$

Momentul central de ordin k al eșantionului X este

$$m_k = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)^k}{n}.$$

- Momentul central al eșantionului, m_k , este statistica corespunzătoare momentului central μ_k .

Forma distribuției

Definition 0.2

Asimetria^a (skewness) al unei distribuții X este

$$\mu_k = \frac{M \left[(X - \mu)^3 \right]}{\left(M \left[(X - \mu)^2 \right] \right)^{3/2}}.$$

Asimetria eșantionului este $a_3 = \frac{m_3}{s^3}$.

^aSau coeficientul lui Pearson de asimetrie.

- a_3 este o măsură a asimetriei: dacă $a_3 = 0$ avem o distribuție simetrică, dacă $a_3 < 0$, distribuția prezintă o asimetrie la stânga (modul este deplasat către dreapta), iar dacă $a_3 > 0$, distribuția prezintă o asimetrie la dreapta (modul este deplasat către stânga).

Forma distribuției

Definition 0.3

Turtirea (kurtosis) unei distribuții X este

$$\mu_k = \frac{M \left[(X - \mu)^4 \right]}{\left(M \left[(X - \mu)^2 \right] \right)^2}.$$

Turtirea eșantionului este $a_4 = \frac{m_4}{s^4}$.

- Turtirea măsoară înălțimea distribuției: dacă $a_4 = 3$ avem o distribuție mesokurtică (e. g. normal law), dacă $a_4 < 3$, distribuția este platykurtică (comparată cu legea normală este mai degrabă "scundă"), iar dacă $a_4 > 3$, distribuția este leptokurtică (este "înaltă" prin comparație cu legea normală).

Inegalitățile lui Markov și a lui Cebâșev

Proposition 1.1

Fie $X \geq 0$ o variabilă aleatoare. Dacă $a > 0$, atunci

$$P(X \geq a) \leq \frac{M[X]}{a}.$$

proof:

$$\begin{aligned} M[X] &= \int_0^{+\infty} tf(t)dt = \int_0^a tf(t)dt + \int_a^{+\infty} tf(t)dt \geq \\ &\int_a^{+\infty} tf(t)dt \geq a \int_a^{+\infty} f(t)dt = aP(X \geq a). \end{aligned}$$




Inegalitatea lui Cebâșev

Proposition 1.2

Fie X o variabilă aleatoare cu media μ și dispersia σ^2 . Atunci

$$P(|X - \mu| \geq k) \leq \frac{\sigma^2}{k^2}.$$

proof: Considerăm variabila $Y = (X - \mu)^2$ și $a = k^2$ în inegalitatea lui Markov

$$P(|X - \mu| \geq k) = P[(X - \mu)^2 \geq k^2] \leq \frac{\mathbb{E}[(X - \mu)^2]}{k^2} = \frac{\sigma^2}{k^2}.$$


Teorema lui Cebâșev

Statistică

Statistică

Statistică

Statistică

Statistică

Statistică

Theorem 1.1

Fie $(X_n)_{n \geq 1}$ un șir de variabile aleatoare independente având dispersii finite, uniform mărginite, i. e. $D^2[X_n] \leq c$, pentru orice $n \geq 1$. Atunci

$$\lim_{n \rightarrow \infty} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n M[X_i] \right| < \epsilon \right) = 1.$$

proof: Știm că

$$M \left[\frac{1}{n} \sum_{i=1}^n X_i \right] = \frac{1}{n} \sum_{i=1}^n M[X_i] \text{ și}$$

$$D^2 \left[\frac{1}{n} \sum_{i=1}^n X_i \right] = \frac{1}{n^2} \sum_{i=1}^n D^2[X_i] < \frac{c}{n}.$$

Teorema lui Cebâșev

Folosind inegalitatea lui Cebâșev pentru variabila $\frac{1}{n} \sum_{i=1}^n X_i$ obținem

$$1 \geq \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n M[X_i] \right| < \epsilon \right) \geq 1 - \frac{D^2 \left[\frac{1}{n} \sum_{i=1}^n X_i \right]}{\epsilon^2} \geq 1 - \frac{c}{n\epsilon^2}.$$

Trecând la limită,

$$\lim_{n \rightarrow \infty} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n M[X_i] \right| < \epsilon \right) = 1. \blacksquare$$

Legea numerelor mari - varianta slabă

- Legea numerelor mari spune că pe măsură ce crește numărul de variabile independente, identic distribuite, media lor de selecție se apropie de media lor comună.

Theorem 1.2

(Legea slabă numerelor mari , legea lui Khintchine) Fie $(X_n)_{n \geq 1}$ un șir de variabile aleatoare independente și identic distribuite cu media μ și dispersia σ^2 . Atunci

$$\lim_{n \rightarrow \infty} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| < \epsilon \right) = 1 \text{ sau } \lim_{n \rightarrow \infty} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq \epsilon \right) = 0.$$

proof: O consecință a teoremei anterioare, deoarece $\frac{1}{n} \sum_{i=1}^n M[X_i] = \mu$.



Legea numerelor mari - varianta tare

Theorem 1.3

(Legea tare numerelor mari) Fie $(X_n)_{n \geq 1}$ un șir de variabile aleatoare independente și identic distribuite cu media μ și dispersia σ^2 . Atunci

$$P \left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu \right) = 1.$$

proof: Fiind mai complicată este omisă. ■

Un exemplu cu frecvențe

- Bernoulli este primul care a demonstrat legea slabă numerelor mari dar doar pentru distribuții Bernoulli.
- Să presupunem că avem o experiență aleatoare și un eveniment aleator asociat A cu $P(A) = p$.
- Repetăm în mod independent experiența și considerăm următorul șir de variabile aleatoare : $X_i = 1$ dacă A se produce la a i -a repetare și 0 altfel.
- Variabilele sunt independente și distribuite Bernoulli cu media p .

Un exemplu cu frecvențe

- The law of large numbers says that, with probability 1,

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow p.$$

- $\sum_{i=1}^n X_i$ is the number of occurrences of A after n performings.
- In other words the law of large numbers says that the A occurs with frequency p .

Istorie

- James Bernoulli a demonstrat legea slabă a numerelor mari în 1700; Poisson i-a generalizat rezultatul în 1800.
- Cebâșev a descoperit inegalitatea care-i poartă numele în 1866, iar Markov a extins rezultatul lui Bernoulli la variabile aleatoare dependente.
- În 1909 Émile Borel a demonstrat teorema care astăzi este cunoscută sub numele de legea tare a numerelor mari (care generalizează o dată în plus teorema lui Bernoulli).
- În 1926 Kolmogorov a obținut o condiție mai generală, suficientă pentru ca un șir de variabile aleatoare independente să respecte legea numerelor mari. Condiția este

$$\sum_{n \geq 1} \frac{D^2[X_n]}{n^2} < +\infty.$$

Teorema limită centrală

Theorem 2.1

(Teorema limită centrală, Lindeberg-Lévy) Fie $(X_n)_{n \geq 1}$ un șir de variabile aleatoare independente și identic distribuite cu media μ și dispersia σ^2 . Atunci

$$\frac{\frac{1}{n} \sum_{i=1}^n X_i - \mu}{\frac{\sigma}{\sqrt{n}}} \rightarrow N(0, 1) \text{ sau}$$

$$\lim_{n \rightarrow \infty} P \left(\frac{\sum_{i=1}^n X_i - n\mu}{\sigma \sqrt{n}} \leq a \right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a \exp -t^2/2 dt.$$

Teorema limită centrală

- Teorema limită centrală permite estimarea unor probabilități asociate sumelor de variabile (independente și identic distribuite).
- Pe de altă parte, teorema explică de ce atât de multe procese (din științele sociale, biologie, psihologie etc) urmează o lege normală.
- în esența ei teorema limită centrală spune că, pentru eșantioane suficient de mari ($n \geq 30$), variabila

$$\frac{\sum_{i=1}^n X_i - n\mu}{\sigma/\sqrt{n}}$$

urmează o lege normală standard, $N(0, 1)$.

- Teorema limită centrală are lco chiar și pentru variabile dependente, dacă au corelația foarte scăzută.

Teorema lui Bernoulli

Proposition 2.1

Fie α_n numărul de apariții ale unui eveniment A în n repetări independente ale unei experiențe aleatoare. Dacă $f_n = \frac{\alpha_n}{n}$ este frecvența relativă a apariției lui A , atunci șirul $(f_n)_{n \geq 1}$ converge în probabilitate la $p = P(A)$.

proof: $\alpha_n = nf_n$ este o variabilă distribuită binomial, astfel $M[\alpha_n] = np$ și $D^2[\alpha_n] = np(1-p)$. Mai mult,

$$\begin{aligned} P(|f_n - p| < \epsilon) &= P(|\alpha_n - np| < n\epsilon) = P(|\alpha_n - \mathbb{E}[\alpha_n]| < n\epsilon) \geq \\ &\geq 1 - \frac{p(1-p)}{n\epsilon^2}. \end{aligned}$$

Evident, trecând la limită, $\lim_{n \rightarrow \infty} P(|f_n - p| < \epsilon) = 1$, pentru orice $\epsilon > 0$.

Aproximarea normală a distribuției binomiale

- Fie X_n un șir de variabile Bernoulli(p) independente.
- $X = \sum_{i=1}^n X_i$ are o distribuție binomială, $B(n, p)$.
- Folosind teorema limită centrală obținem teorema de Moivre-Laplace care spune că, pentru n suficient de mare, variabila

$$Y = \frac{X - \mathbb{E}[X]}{\sqrt{\text{Var}[X]}} = \frac{X - np}{\sqrt{np(1-p)}}$$

este o variabilă normală standard ($N(0, 1)$).

- Aproximarea este bună pentru $np(1-p) \geq 10$.

Aproximarea normală a distribuției binomiale

- Un alt mod de a vedea acest rezultat este următorul: când k este aproape de np

$$\binom{n}{k} p^k (1-p)^{n-k} \sim \frac{\exp\left(-\frac{(k-np)^2}{2np(1-p)}\right)}{\sqrt{2\pi np(1-p)}}.$$

- Considerăm următorul exemplu: fie X numărul de apariții ale steimei în 40 de aruncări ale unei monede.
- Cât este $P(X = 20)$?

$$\begin{aligned} P(X = 20) &= P(19.5 \leq X \leq 20.5) = \\ &= P\left(\frac{19.5 - 20}{\sqrt{10}} \leq \frac{X - 20}{\sqrt{10}} \leq \frac{20.5 - 20}{\sqrt{10}}\right) = P\left(-0.16 \leq \frac{X - 20}{\sqrt{10}} \leq 0.16\right) \\ &\sim \Phi(0.16) - \Phi(-0.16) = 0.1272, \end{aligned}$$

unde $\Phi(\cdot)$ este funcția de repartiție a variabilei $N(0, 1)$.

Corecția continuă

- Corecția continuă este o ajustare care se face ori de câte ori o distribuție discretă este aproximată printr-una continuă.
- $P(X = 10) = P(9.5 \leq X \leq 10.5)$, $P(X > 15) = P(X \geq 15.5)$,
 $P(X < 13) = P(X \leq 12.5)$.

Statistică inferențială

- Statistica inferențială are scopul să trage concluzii relativ la o populație folosind rezultate ale teoriei probabilităților și statistici obținute din eșantioane.
- Fără utilizarea teoriei probabilităților (legea numerelor mari, teorema limită centrală ș. a.) am putea considera drept sistematic un comportament care este de fapt datorat hazardului sau, din contră, un comportament sistematic ar putea trece neobservat.
- Exemple de inferențe statistice: *intervale de încredere* pentru estimarea parametrilor sau *teste de semnificație*.
- Tehnicile de inferență au la bază distribuțiile eșantioanelor.

Estimarea parametrilor

- Distribuția unei anumite populații poate fi necunoscută în întregime; de aceea ne putem dori să aflăm cel puțin media, dispersia sau alți parametri ai ei.
- Acești parametri ai unei populații pot fi estimați folosind statistici calculate din eșantion.
- Există două tipuri de estimare: *estimare punctuală* și *estimare cu intervale*.

Proposition 1.1

Estimarea punctuală a unui parametru constă în determinarea unui număr, de obicei valoarea unei statistici corespunzătoare, desemnat să estimeze acel parametru.

Estimarea parametrilor

Statistică Statistică Statistică Statistică Statistică Statistică
Statistică Statistică Statistică Statistică Statistică Statistică Statistică
Statistică Statistică Statistică Statistică Statistică Statistică
Statistică Statistică Statistică Statistică Statistică Statistică Statistică

Proposition 1.2

Estimarea cu un interval a unui parametru constă în determinarea unui interval ale cărui limite sunt statistici calculate din eșantioane.

Nivelul de încredere $(1 - \alpha)$ este proporția acelor intervale (de încredere) care conțin parametrul estimat.

Un **interval de încredere** este un interval care are un anumit nivelul de încredere (prescris).

Statistică Statistică Statistică Statistică Statistică Statistică Statistică
Statistică Statistică Statistică Statistică Statistică Statistică Statistică
Statistică Statistică Statistică Statistică Statistică Statistică Statistică

Estimarea punctuală

- Exemple de estimatori punctuali sunt media de selecție \bar{x}_n , dispersia eșantionului s^2 , sau deviația standard a eșantionului s .
- Pentru un parametru dat putem avea mai mulți estimatori punctuali: media populației poate fi estimată prin media de selecție, mediană, mod.
- Se pot ridica anumite întrebări legate de calitatea estimatorilor punctuali.
- Cât de exact este un estimator (i.e., the *încrederea*) - este în mod frecvent mai mare (*supra-estimator*) sau mai mic (*sub-estimator*) în raport cu parametrul estimat?
- Din acest punct de vedere se preferă estimatorii *nedeplasați*.
- Care este variabilitatea unui estimator punctual (văzut ca o variabilă aleatoare) - aceasta este *acuratețea*.

Caracteristici ale estimării punctuale

- De exemplu dispersia mediei de selecție, care este σ/\sqrt{n} : cu cât este mai mare eșantionul cu atât mai mică este variabilitatea mediei de selecție.

Definition 1.1

Fie θ un anumit parametru al unei populații, x un eșantion al acestei populații și $\hat{\theta}_n = \hat{\theta}_n(x)$ un estimator punctual al lui θ .

$\hat{\theta}_n$ este o **statistică nedeplasată** dacă $M[\hat{\theta}_n] = \theta$. Altfel $\hat{\theta}_n$ este numită statistică **deplasată**.

$\hat{\theta}_n$ este o **statistică consistentă** dacă, pentru orice $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| < \epsilon) = 1.$$

Un estimator nedeplasat cu dispersie minimă, dacă există, se numește **statistică eficientă**.

Distribuția mediei de selecție

- Fie \mathcal{P} o populație formată din N indivizi ale căror valori ale atributului sunt a_1, a_2, \dots, a_N .
- Variabila care reprezintă atributul populației este notată cu X .
- Pentru această populație, deci și pentru X , media și dispersia sunt

$$\frac{1}{N} \sum_{i=1}^N a_i = \mu, \text{ respectiv } \frac{1}{N} \sum_{i=1}^N (a_i - \mu)^2 = \sigma^2.$$

- Pentru a estima μ și σ^2 folosim eșantioane de dimensiune n .

Distribuția mediei de selecție

- Pentru un eșantion dat, $x^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})$, media de selecție este

$$\bar{x}_n^{(k)} = \frac{1}{n} \sum_{j=1}^n x_j^{(k)}.$$

- Fiecare $x_j^{(k)}$ este o valoare a unei variabile aleatoare cu aceeași distribuție ca X .

Definition 1.2

Variabila aleatoare care are drept valori toate mediile de selecție posibile, $\bar{x}_n^{(k)}$, pentru eșantioane de dimensiune n , se numește distribuția mediei de selecție.

Distribuția mediei de selecție

- Se poate demonstra următorul rezultat

Theorem 1.1

Media și dispersia mediei de selecție, \bar{x}_n , sunt μ și σ^2/n :

$$M[\bar{x}_n] = \mu, D^2[\bar{x}_n] = \frac{\sigma^2}{n}.$$

În plus, pentru valori mari ale lui n (≥ 30), distribuția mediei de selecție este normală, i. e.

$$\bar{x}_n \sim N(\mu, \sigma^2/n).$$

Estimarea mediei și a dispersiei

- Bineînțeles, un estimator al mediei adevărate a populației (μ) este media de selecție, \bar{x}_n .
- Un estimator pentru adevărata dispersie a populației (σ^2) este dispersia eșantionului s^2 - o statistică nedeplasată.
- Pentru deviația standard a mediei de selecție, σ/\sqrt{n} numită **eroarea standard a mediei** (standard error of the mean - SEM), un estimator este s/\sqrt{n} .

Estimarea cu intervale

- Un estimator de tip interval al mediei populației poate fi obținut cu inegalitatea lui Cebâșev:

$$P(|\bar{x}_n - M[\bar{x}_n]| \geq k \cdot \text{Var}[\bar{x}_n]) \leq \frac{1}{k^2} \Leftrightarrow$$

$$P\left(|\bar{x}_n - \mu| \geq k \cdot \frac{\sigma}{\sqrt{n}}\right) \leq \frac{1}{k^2} \Leftrightarrow$$

$$P\left(|\bar{x}_n - \mu| < k \cdot \frac{\sigma}{\sqrt{n}}\right) \geq 1 - \frac{1}{k^2} \Leftrightarrow$$

$$P\left(\bar{x}_n - k \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x}_n + k \cdot \frac{\sigma}{\sqrt{n}}\right) \geq 1 - \frac{1}{k^2}.$$

- Astfel o estimare cu interval a lui μ este

$$\left(\bar{x}_n - k \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x}_n + k \cdot \frac{\sigma}{\sqrt{n}}\right).$$

Intervale de încredere

Definition 1.3

Un interval de încredere pentru un parametru θ cu $(1 - \alpha)$ nivel de încredere este definit cu ajutorul a două statistici L și U astfel ca

$$P(L \leq \theta \leq U) = (1 - \alpha).$$

- L și U sunt de fapt variabile aleatoare ale căror valori sunt statistici: pentru diferite eșantioane au diferite valori.
- De obicei nivelul de încredere este o probabilitate aproape de 1, cum ar fi 0.90, 0.95, sau 0.99 (care dau $\alpha \in \{0.10, 0.05, 0.01\}$).

Intervale de încredere pentru medie

- Să presupunem că avem un eșantion de dimensiune n și un nivel de încredere $(1 - \alpha)$ și dorim să construim un interval de încredere pentru media μ .
- Știm că media de selecție, \bar{x}_n , urmează o distribuție normală

$$N\left(\mu, \frac{\sigma^2}{n}\right).$$

- Prin standardizare variabila $Z = \frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}}$ este distribuită $N(0, 1)$.
- Căutăm o valoare z^* , numită *valoare critică*, astfel ca o variabilă normal standard să acopere sub grafic o arie de $(1 - \alpha)$ pe intervalul centrat în medie și de lungime $2z^*$ deviații standard.

Intervale de încredere pentru medie

- Fie $Z : N(0, 1)$, valoarea critică este aleasă așa încât

$$P(-z^* \leq Z \leq z^*) = 1 - \alpha.$$

- Definiții echivalente ale lui z^* :

$$P(Z \leq -z^*) = \alpha/2 \text{ sau } P(Z \geq z^*) = \alpha/2.$$

- Notăm cu $\Phi(a) = P(Z \leq a)$, funcția de repartiție normală standard.
- Astfel $z^* = -\Phi^{-1}(\alpha/2)$ - valoare care poate fi găsită în tabele sau aproximată în pachetele de prelucrare statistică uzuale (R, MiniTab, SPSS etc).
- Odată ce am determinat valoarea critică, știm că

$$P\left(-z^* \leq \frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}} \leq z^*\right) = 1 - \alpha.$$

Intervale de încredere pentru medie

- Echiuivalent

$$P\left(\bar{x}_n - z^* \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x}_n + z^* \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

- Am demonstrat astfel

Theorem 1.2

Un interval de încredere cu nivelul de încredere $(1 - \alpha)$ pentru pentru media unei populații cu media cunoscută este

$$\left(\bar{x}_n - z^* \frac{\sigma}{\sqrt{n}}, \bar{x}_n + z^* \frac{\sigma}{\sqrt{n}}\right),$$

unde z^ este valoarea critică asociată cu $\alpha/2$. Mai mult, acest interval este exact pentru o populație distribuită normal și aproximativ altfel, când eșantionul este suficient de mare ($n \geq 30$).*

Intervale de încredere pentru medie

Nivel de încredere	$\alpha/2$	z^*
90%	0.05	1.645
95%	0.025	1.960
99%	0.005	2.576

- $z^* \frac{\sigma}{\sqrt{n}}$ se numește *eroarea marginală*.
- Lungimea unui interval de încredere pentru medie este $2z^* \frac{\sigma}{\sqrt{n}}$.
- Dacă dorim o lungime anume pentru acest interval, w , atunci ne trebuie un eșantion de dimensiune $n = \frac{(2z^*\sigma)^2}{w^2}$. În practică această valoare poate fi nerealistă (dacă n este prea mare).
- Pe măsură ce n crește lungimea intervalului (sau eroarea marginală) scade, ceea ce este util, dar poate fi nepractic.

Intervale de încredere pentru medie

- Să ne amintim că $\frac{\sigma}{\sqrt{n}}$ este eroarea standard a mediei. Un interval de încredere poate fi văzut astfel

$$\text{estimatedMean} \pm z^* \text{estimatedStdDev}$$

- Exemplu.* Un anumit medicament este analizat măsurându-i-se substanța activă de trei ori, rezultatele sunt 0.8403, 0.8636, și 0.8447 g/l. Se știe că această concentrație urmează o lege normală cu deviația standard $\sigma = 0.0068$ g/l. Să se determine un interval de încredere de 99% pentru adevărata medie a concentrației, μ .

Soluție:

$$\bar{x}_3 = 0.8404, \alpha = 0.01, \alpha/2 = 0.005, z^* = 2.576, z^* \frac{\sigma}{\sqrt{n}} = 0.0101$$

Intervalul de încredere este (0.8303, 0.8505).

Bibliography



Freedman, D., R. Pisani, R. Purves, *Statistics*, W. W. Norton & Company, 4th edition, 2007.



Johnson, R., P. Kuby, *Elementary Statistics*, Brooks/Cole, Cengage Learning, 11th edition, 2012.



Shao, J., *Mathematical Statistics*, Springer Verlag, 1998.



Spiegel, M. R., L. J. Stephens, *Theory and Problems of Statistics*, Schaum's Outline Series, McGraw Hill, 3rd edition, 1999.