

Nagy eltérés tétel - bevezetés

1

Legyen X_1, X_2, \dots, X_n független és azonos eloszlású,

$$EX_i = m, \text{Var} X_i = \sigma^2, S_n = X_1 + \dots + X_n.$$

Igy persze $ES_n = nm, \text{Var} S_n = n\sigma^2, DS_n = \sqrt{n}\sigma.$

Kérdés a $P(S_n \leq K)$ vagy $P(S_n \geq K)$ valószínűség,

ahol $K \in \mathbb{R}$.

Persze a kérdés átfejelematható az $\tilde{S}_n := \frac{S_n - nm}{\sqrt{n}\sigma}$

normált val. változóra: $ES_n = 0$ és $\text{Var} \tilde{S}_n = 1,$

így jobban érzelhető, hogy az $\{S_n \leq K\}$ mennyire „extrém” vagy „meglepő”:

$$P(S_n \leq K) = P(\tilde{S}_n \leq x), \text{ ahol } K = nm + \sqrt{n}\sigma x$$

$$P(S_n \geq K) = P(\tilde{S}_n \geq x) \quad \text{avagy} \quad K = ES_n + x DS_n$$

Láttuk: A centrális határelosztás tétel arról az esetről

szól, amikor $x \in \mathbb{R}$ rögzített, vagyis $\underbrace{K - ES_n}_{\text{eltérés a várható értéktől}} \sim \underbrace{\sqrt{n}}_{\text{nagy-számszerű}}$

~~Probléma?~~

Nagy eltérés problémának nevezzük azt, amikor

$$|K - ES_n| \gg \sqrt{n}, \text{ pl. } K - ES_n \sim n$$

~~Az~~ Ilyenkor $P(S_n \geq K)$ becslésére a CHT nem alkalmas:

a keresett valószínűség nagyon kicsi (vagy nagyon nagy), de

a CHT becslés hibája miatt nem lehet CHT-vel megbízhatóan becsülni.

Ilyenkor hasznosak a nagy eltérés tételek.

Hoeffding egyenlőtlenség (Hoeffding, 1963) 3
(ejtsd: "höffding")

Tétel Legyenek X_1, X_2, \dots, X_n független valószínűségi
változók (de nem feltétlen) azonos eloszlásúak, és
legyen $S_n = X_1 + X_2 + \dots + X_n$.

Tegyük fel, hogy mindegyik X_i korlátos: minden i -re
van olyan $a_i < b_i$, hogy $a_i \leq X_i \leq b_i$

Ekkor minden $t > 0$ -ra

$$P(S_n \geq ES_n + t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad \left[\begin{array}{l} \text{nagy eltérés} \\ \text{felé} \end{array} \right]$$

és

$$P(S_n \leq ES_n - t) \leq \text{---} \left[\begin{array}{l} \text{nagy eltérés} \\ \text{lefelé} \end{array} \right]$$

Megjegyzések:

1.) A korlátosság erős feltétel, a tétel legjobb gyengéje.

2.) Ez egy nagy eltérés tétel: Ha az X_i -k azonos
eloszlásúak, akkor $(b_i - a_i)^2$ "szórás négyzet-szerűség", hogy
a nevezőbeli $\sum_{i=1}^n (b_i - a_i)^2 = n(b-a)^2 \sim n \sim \text{Var } S_n$,

Igy a jobboldal akkor lesz kicsi, ha $t \gg \sqrt{n}$. 4

3.) A tétel minden n -re felső becslést ad (nem csak aszimptotikusan $n \rightarrow \infty$ -ben).

4.) A felső becslés általában nem éles: a tényleges érték sokkal kisebb is lehet.

5.) A tétel nagyon könnyen alkalmazható:

- kevés standardissal jár,

és főleg

- keveset kell tudni az X_i -k eloszlásáról:

→ kell az n (darabszám)

→ kellene a korlátok (a_i, b_i)

→ kell az ES_n , vagyis az összeg várható értéke.

Megj: Persze $ES_n = EX_1 + EX_2 + \dots + EX_n$, de nem baj, ha külön-külön nem ismerjük őket.

Példa a Hoeffding egyenlőtlenség alkalmazására

5

Piripécson a kukásautó hűtőnként 1000 háztartásból vисти el a szemetet. A lakók 3féle kukát használnak:

400 háztartásban	60 l-es	kuka van
100 ~	110 ~	~
500 ~	110 120	~

Az egyes háztartások által a kukába tett szemét mennyisége véletlen és független, de persze legfeljebb annyi lehet, amennyi a kukában fér.

Ha nem fér bele az összes szemét a kukásautóba, akkor baj van.

Hány liter szemétnek kell elférni a kukásautóban, hogy ennek valószínűsége egész biztosan legfeljebb 1% legyen?

A kukás edgnek persze fogalma sincs, hogy az egyes emberek külön-külön mennyi szemetet termelnek, de azt tudják, hogy az összes mennyiség átlagosan 40000 liter székelt lenni.

Megoldás: Legyen $n=1000$ és legyen X_1, X_2, \dots, X_n az egyes kukákban a szemét mennyisége, literben. Így

$S_n := X_1 + X_2 + \dots + X_n$ az össz-mennyiség. Olyan K -t keresünk, amire $P(S_n > K) \leq 0.01$. 6

A $P(S_n > K)$ valószínűség becslésére a CHT nem alkalmas, pl. mert nincs okunk feltételezni, hogy az X_i -k azonos eloszlásúak. Ráadásul a szórásukat sem ismerjük.

A Hoeffding egyenlőtlensége viszont alkalmazható, mert az

X_i -k korlátosak: $a_i \leq X_i \leq b_i$, ahol $a_i = 0$,

és $b_i = \begin{cases} 60 & , 400 \text{ db } i\text{-re} & (\text{mondjuk } i=1, \dots, 400) \\ 110 & , 100 \text{ db } i\text{-re} & (\text{mondjuk } i=401, \dots, 500) \\ 120 & , 500 \text{ db } i\text{-re} & (\text{mondjuk } i=501, \dots, 1000) \end{cases}$

Pontos a szumma!

Pontos a négyzet!

Így $\sum_{i=1}^n (b_i - a_i)^2 = 400 \cdot (60 - 0)^2 + 100 \cdot (110 - 0)^2 + 500 \cdot (120 - 0)^2 = 9850000$

Vagyis a Hoeffding egyenlőtlenség szerint minden $t > 0$ -ra

$$P(S_n \geq ES_n + t) \leq \exp\left(-\frac{2t^2}{9850000}\right).$$

Az ebl, hogy a jobb oldali $1\% = 0.01$ legyen:

$$\exp\left(-\frac{2t^2}{9850000}\right) = 0.01 \Leftrightarrow t = \sqrt{\frac{9850000}{2} \cdot (-\ln(0.01))} \approx 4762$$

Vagyis $K := \mathbb{E}S_n + t$ jó lesz $t = 4762$ választással.

ennyivel kell nagyobb kapacitás, mint az $\mathbb{E}S_n = 40000$
várható érték. 7

Válasz: $K = \mathbb{E}S_n + t = 44762$ literes kukásautó elég.