

Náhodný výběr

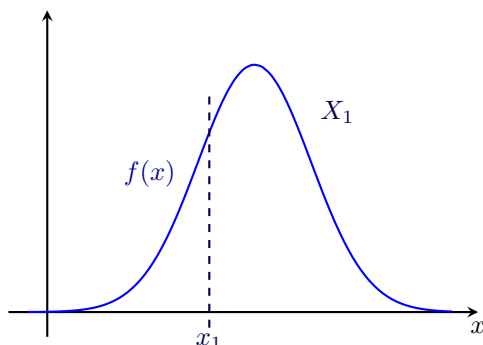
Aplikovaná statistika

16. dubna 2024

JAROSLAV SCHMIDT

Sadě n náhodných veličin X_1, \dots, X_n , které mají všechny shodné rozdělení, říkáme náhodný výběr. Velkým písmenem X značíme, že se stále jedná o náhodnou veličinu. Občas je totiž zaměňována se sadou x_1, \dots, x_n kterou označujeme malými písmeny a míníme tím konkrétní realizaci náhodného výběru.

Představme si to například následovně. Zkoumáme průměrnou výšku populace. Vydáme se do terénu a zeptáme se první oběti na jeho výšku. Než nám odpoví, dá se to vnímat tak, že jeho výška je neznámá veličina X_1 . Jeho výška je tedy prozatím dána jakousi hustotou pravděpodobnosti $f(x)$. To je v obrázku 1 označeno modrou křivkou. Jeho výška je tedy prozatím jen funkcí nabývající nekonečně mnoha hodnot.



Obrázek 1: Náhodná veličina X_1 a její konkrétní realizace x_1 .

Poté co nám prozradí svoji skutečnou výšku, dostáváme to, čemu říkáme realizace náhodné veličiny. Získáváme číslo x_1 . Číslo si zaznamenejme a pokračujeme k další oběti. Opět neznáme její výšku, tedy jedná se opět o náhodnou veličinu X_2 reprezentovanou stejnou funkcí hustoty náhodné veličiny $f(x)$ a poté co nám prozradí svoji výšku, bude se jednat o další realizaci náhodné veličiny x_2 . Předpokládejme na chvíli, že známe funkci $f(x)$, potom pro každou náhodnou veličinu X_i můžeme spočítat střední hodnotu a rozptyl, označme je μ a σ^2 . Co když chceme vyšetřovat průměrnou hodnotu? Tedy hodnotu

$$\bar{X} = \frac{\sum_i X_i}{n}. \quad (1)$$

Je dobré podotknout, že průměrujeme X_i nikoliv x_i . Jedná se tedy o náhodné veličiny, jejichž hodnotu zatím neznáme, může nabývat různých hodnot s různou pravděpodobností. Díky tomu je i veličina \bar{X} náhodnou veličinou se svojí hustotou pravděpodobnosti. Nejedná se tedy o jediné číslo. Naproti tomu $\bar{x} = \sum_i x_i/n$ je prosté číslo a jedná se o pouhý aritmetický průměr naměřených hodnot. Už by nám měly být známé tyto dvě pravidla pro výpočet středních hodnot a rozptylů:

$$\mathbb{E}(aX + bY) = a\mathbb{E}(X) + b\mathbb{E}(Y), \quad (2)$$

$$\text{var}(aX + bY) = a^2 \text{var}(X) + b^2 \text{var}(Y). \quad (3)$$

Ty můžeme ihned použít pro výpočet střední hodnoty a rozptylu naší nové veličiny \bar{X} jako

$$\mathbb{E}\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{\mathbb{E}(X_1) + \dots + \mathbb{E}(X_n)}{n} = \frac{n\mu}{n} = \mu \quad (4)$$

$$\text{var}\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{\text{var}(X_1) + \dots + \text{var}(X_n)}{n^2} = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \quad (5)$$

Hodí se i vědět, jaké má tato náhodná veličina \bar{X} rozdělení. To ovšem není triviální spočítat a navíc k tomu potřebujeme znát funkci $f(x)$, kterou často dopředu neznáme. Naštěstí nás zachraňuje centrální limitní věta. Tu zde nebudeme odvozovat, jen ji uvedeme.

Ta říká, že bez ohledu na to, z jakého rozdělení pochází náhodné veličiny X_i , jejich průměrná hodnota \bar{X} je asymptoticky normálně rozdělená. Asymptoticky v tom smyslu, že plná platnost nastává jen při $n \rightarrow \infty$. Pro menší n je tato platnost jen přibližná. A jelikož už známe střední hodnotu i rozptyl náhodné veličiny \bar{X} , známe její rozdělení

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right). \quad (6)$$

Připomeňme, že toto platí přesně jen asymptoticky. Tento výsledek zaslouží komentář. Povšimněme si výsledného rozptylu. Pokud budeme n zvětšovat nade všechny meze, potom se rozptyl bude blížit nule, tedy

$$n \rightarrow \infty \quad \Rightarrow \quad \bar{X} \sim \mathcal{N}(\mu, 0). \quad (7)$$

Tento výsledek znamená, že pokud budeme mít výběr o velikosti $n = \infty$, potom \bar{X} již nebude náhodná veličina a bude přímo platit $\bar{X} = \mu$