



Opakování základních pojmů statistiky

Petr Pošík

Části dokumentu jsou převzaty (i doslovně) z
Mirko Navara: Pravděpodobnost a matematická statistika,
https://cw.felk.cvut.cz/lib/exe/fetch.php/courses/a6m33ssl/pms_print.pdf
s laskavým svolením autora.



Úvod do statistiky



Pravděpodobnost a statistika

Teorie pravděpodobnosti: Nástroj pro rozhodování v systémech, jejichž popis známe, ale jejichž *budoucí* stav a chování závisí na okolnostech, které neznáme. Deduktivní uvažování.

Statistika: Nástroj pro hledání a ověřování pravděpodobnostního popisu reálných systémů na základě jejich pozorování. Induktivní uvažování.

Úvod do statistiky

- Pravděpodobnost a statistika

- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Pravděpodobnost a statistika

Teorie pravděpodobnosti: Nástroj pro rozhodování v systémech, jejichž popis známe, ale jejichž *budoucí* stav a chování závisí na okolnostech, které neznáme. Deduktivní uvažování.

Statistika: Nástroj pro hledání a ověřování pravděpodobnostního popisu reálných systémů na základě jejich pozorování. Induktivní uvažování.

Dedukce: Ze znalosti „*obecného*“ usuzujeme na vlastnosti „*konkrétního*“. *Specializace* obecných znalostí, zde využití *pravděpodobnosti*.

Indukce: Ze znalosti „*konkrétního*“ usuzujeme na vlastnosti „*obecného*“. *Generalizace* poznatků, zde využití *statistického usuzování*.

Úvod do statistiky

- Pravděpodobnost a statistika

- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Povaha statistiky

Úvod do statistiky

- Pravděpodobnost a statistika
 - Účel
 - Analýza dat
 - Základní pojmy
 - Druhy veličin
 - Šetření vs. exp.
 - Randomizace
 - Statistika
 - Náhodný výběr
 - Histogram a empirické rozdělení
 - Odhady
 - Výběrový průměr
 - Cent. lim. věta
 - Výběrový rozptyl
 - Rozdělení výběrového rozptylu
 - Chí-kvadrát
 - Výběrová sm.odch.
 - Výběrový medián
 - Míry polohy

Nemusím sníst celého vola, abych poznal, že je tuhý.

Samuel Johnson



Statistika: účel a členění

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Statistika jako matematická disciplína:

- Zkoumá *společné vlastnosti* velkého počtu obdobných jevů.
- Využívá jen *vybraný vzorek* jevů, nikoli všechny.
- Zabývá se sběrem, prezentací, analýzou a interpretací dat popisujících jevy či vlastnosti pozorovaných objektů.
- Typické úlohy:
 - Odhad parametrů pravděpodobnostního modelu
 - Testování hypotéz

Matematická (teoretická) statistika: výzkum a popis nových metod.

Aplikovaná statistika: použití stat. metod v konkrétních problémech různých oborů, např. například v přírodních či společenských vědách, v politice nebo v lékařství.

Deskriptivní statistika se zabývá numerickým nebo grafickým popisem získaných dat

Inferenční (induktivní) statistika se zabývá *vyhledáváním zákonitostí* v datech naměřených na vzorku jedinců nebo objektů a *zobecnováním* těchto zákonitostí na skupinu, z níž byl vzorek vybrán. Inferenční statistika vychází z počtu pravděpodobnosti.



Analýza dat

Z <http://www.quora.com/What-is-the-difference-between-statistics-and-machine-learning>:

- **Machine Learning** is AI people doing data analysis.
- **Data Mining** is database people doing data analysis.
- **Applied Statistics** is statisticians doing data analysis (mathematical statistics is mathematicians doing statistics, a meta-activity analogous to theoretical computer scientists investigating, say, PAC and convergence - which sometimes gets put into machine learning).
- **Infographics** is Graphic Designers doing data analysis.
- **Data Journalism** is Journalists doing data analysis.
- **Econometrics** is Economists doing data analysis (and here you can win a Nobel Prize).
- **Psychometrics** is Psychologists doing data analysis.
- **Chemometrics** and **Cheminformatics** are Chemists doing data analysis.
- **Bioinformatics** is Biologists doing data analysis.

:-)

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- **Analýza dat**
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Základní pojmy

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- **Základní pojmy**
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Statistické jednotky: Objekty, jejichž vlastnosti zkoumáme. (Lidé, buňky, hřídele, ...)

Statistický soubor: Specificky vymezená množina statistických jednotek; o libovolném prvku musíme být schopni rozhodnout, zda do statistického souboru patří či nikoliv. (Např. učitelé FEL, kteří měli na FEL v roce 2013 nadpoloviční úvazek.)

Základní soubor (populace): *Úplný* statistický soubor (soubor všech jednotek). Může být i nekonečný.

Výběrový soubor (výběr, vzorek) rozsahu n : Konečný soubor obsahující jen těch n prvků, které skutečně pozorujeme nebo měříme. Jejich výběr ze základního souboru musí být proveden *náhodně* (s rovnoměrným rozdělením), nebo podle znaku, který se studovanými znaky nespojuje.

Proč výběr?

1. Omezené zdroje
Nejsou prostředky nebo čas na zkoumání celé populace.
2. Destruktivní zkoušky
Můžeme použít *všechny* červené krvinky pacienta, abychom zjistili jejich skutečnou průměrnou velikost?
3. Vzorek bývá přesnější
Sběr dat pro menší vzorek lze provést menší skupinou lépe proškolených lidí.



Druhy veličin

Znak	Škála	Možné operace	Příklady
Kval.	Nominální	Popsat příslušnost	
	Ordinální	Seřadit	
Kvant.	Intervalová	Porovnat vzdálenosti	
	Poměrová	Porovnat velikosti	

- Spojité vs. diskrétní
- Nezávislé (vstupy) vs. závislé (výstupy)

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- **Druhy veličin**
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Druhy veličin

Znak	Škála	Možné operace	Příklady
Kval.	Nominální	Popsat příslušnost	Barva očí, národnost, pohlaví, místo narození
	Ordinální	Seřadit	
Kvant.	Intervalová	Porovnat vzdálenosti	
	Poměrová	Porovnat velikosti	

- Spojité vs. diskrétní
- Nezávislé (vstupy) vs. závislé (výstupy)

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- **Druhy veličin**
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Druhy veličin

Znak	Škála	Možné operace	Příklady
Kval.	Nominální	Popsat příslušnost	Barva očí, národnost, pohlaví, místo narození
	Ordinální	Seřadit	Popis velikosti (S,M,L,XL,XXL), vzdělání (ZŠ, SŠ, VŠ)
Kvant.	Intervalová	Porovnat vzdálenosti	
	Poměrová	Porovnat velikosti	

- Spojité vs. diskrétní
- Nezávislé (vstupy) vs. závislé (výstupy)

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- **Druhy veličin**
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Druhy veličin

Znak	Škála	Možné operace	Příklady
Kval.	Nominální	Popsat příslušnost	Barva očí, národnost, pohlaví, místo narození
	Ordinální	Seřadit	Popis velikosti (S,M,L,XL,XXL), vzdělání (ZŠ, SŠ, VŠ)
Kvant.	Intervalová	Porovnat vzdálenosti	Kalendářní datum, teplota, úhel
	Poměrová	Porovnat velikosti	

- Spojité vs. diskrétní
- Nezávislé (vstupy) vs. závislé (výstupy)

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- **Druhy veličin**
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Druhy veličin

Znak	Škála	Možné operace	Příklady
Kval.	Nominální	Popsat příslušnost	Barva očí, národnost, pohlaví, místo narození
	Ordinální	Seřadit	Popis velikosti (S,M,L,XL,XXL), vzdělání (ZŠ, SŠ, VŠ)
Kvant.	Intervalová	Porovnat vzdálenosti	Kalendářní datum, teplota, úhel
	Poměrová	Porovnat velikosti	Objem prodeje, průměr hřídle, hmotnost, teplota v Kelvinech, úhel vzhledem k ...

- Spojité vs. diskrétní
- Nezávislé (vstupy) vs. závislé (výstupy)

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- **Druhy veličin**
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Šetření vs. experiment

Šetření Průzkumy veřejného mínění, potvrzovací studie prováděné mezi obyvatelstvem...

- Pouze sledujeme, nijak nezasahujeme.
- Nemůžeme ovlivnit rozdělení subjektů do skupin.
- Rozdíl mezi skupinami bývá často ovlivněn tzv. matoucí (confounding) veličinou.
- Částečně lze vyřešit zahrnutím matoucích veličin do modelu; problém je v tom, že nevíme, co může být matoucí veličinou.

Experiment Klinické studie, laboratorní testy

- Aktivně zasahujeme a ovlivňujeme podmínky.
- Snažíme se vyloučit vlivy, které nejsou předmětem výzkumu. Náhodným rozdělením subjektů se tyto vlivy vyruší.
- Zkoumané veličiny aktivně nastavujeme tak, aby vzorek nebyl vychýlený.

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Randomizovaný experiment

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- **Randomizace**
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Vzorek zkoumaných subjektů

- se *náhodně rozdělí* do skupin,
- ke kterým se snažíme *chovat* naprosto *shodným způsobem*.
- Pokud zkoumané subjekty neví, ve které skupině jsou zařazeny (což je obvyklé), mluvíme o **slepém experimentu**.
- Pokud navíc ani lidé, kteří subjekty hodnotí, příp. se o ně starají, neví, do které skupiny jsou subjekty zařazeny, mluvíme o **dvojitě slepém experimentu**.

Randomizované experimenty dávají přesnější představu o sledovaném jevu než výběrová šetření, ale:

- **Randomizace nemožná**, např. souvislost pohlaví s výší platů: pohlaví nelze přiřadit náhodně.
- **Randomizace možná, ale nepraktická**, např. studie kriminality ve městech a na vesnici: lidé se nepřestěhují jen kvůli průzkumu.
- **Randomizace možná, praktická, přesto se nedělá**, např. studie výhod předškolních vzdělávacích programů poskytovaných zdarma, kdy poptávka převyšuje nabídku: náhodné přidělení míst je férový postup, ale lidé odmítají připustit, že generátor náhodných čísel udělá jejich práci lépe než oni.
- **Etické problémy**: Experimenty na lidech a na zvířatech ano či ne? Správné otázky by měly znít: *Experimentujeme s rozmyslem nebo hazardujeme? Provádíme experimenty, abychom se z nich dozvěděli maximum, nebo jsou naše experimenty chabé, poskytují špatné informace a poškozují lidi?*



Etika: ilustrace

Jednoho dne naši školu navštívil jeden významný chirurg z Bostonu a udělal nám skvělou přednášku o velké skupině pacientů, na kterých vyzkoušel svou novou metodu vaskulární rekonstrukce. Na konci přednášky se zeptal jeden ze studentů.

- *Student:* „Měl jste nějakou kontrolní skupinu?“

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- **Randomizace**
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Etika: ilustrace

Jednoho dne naši školu navštívil jeden významný chirurg z Bostonu a udělal nám skvělou přednášku o velké skupině pacientů, na kterých vyzkoušel svou novou metodu vaskulární rekonstrukce. Na konci přednášky se zeptal jeden ze studentů.

- *Student*: „Měl jste nějakou kontrolní skupinu?“
- *Chirurg*: „Myslíte tím, zda jsem svou novou metodou operoval jen polovinu pacientů?“

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- **Randomizace**
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Etika: ilustrace

Jednoho dne naši školu navštívil jeden významný chirurg z Bostonu a udělal nám skvělou přednášku o velké skupině pacientů, na kterých vyzkoušel svou novou metodu vaskulární rekonstrukce. Na konci přednášky se zeptal jeden ze studentů.

- *Student*: „Měl jste nějakou kontrolní skupinu?“
- *Chirurg*: „Myslíte tím, zda jsem svou novou metodou operoval jen polovinu pacientů?“
- *Student*: „Ano, to je přesně to, co mám na mysli.“

Lékař praštil pěstí do stolu a zahřměl:

- *Chirurg*: „Samozřejmě, že ne! Tím bych odsoudil polovinu z nich k smrti!“

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- **Randomizace**
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Etika: ilustrace

Jednoho dne naši školu navštívil jeden významný chirurg z Bostonu a udělal nám skvělou přednášku o velké skupině pacientů, na kterých vyzkoušel svou novou metodu vaskulární rekonstrukce. Na konci přednášky se zeptal jeden ze studentů.

- *Student*: „Měl jste nějakou kontrolní skupinu?“
- *Chirurg*: „Myslíte tím, zda jsem svou novou metodou operoval jen polovinu pacientů?“
- *Student*: „Ano, to je přesně to, co mám na mysli.“

Lékař praštil pěstí do stolu a zahřměl:

- *Chirurg*: „Samozřejmě, že ne! Tím bych odsoudil polovinu z nich k smrti!“
- *Student*: „A kterou polovinu?“

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- **Randomizace**
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Kontrolní skupina

Statistikovi a jeho ženě se narodila dvojčata. Ihned po návratu z porodnice volá muž do kostela a oznamuje tu skvělou zprávu. Kněz má samozřejmě radost:

- „To je skvělé! Tak je co nejdříve přivezte a pokřtíme je!“
- „Ne,“ řekl statistik, „pokřtíme jen jedno. To druhé si necháme jako kontrolní skupinu.“

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- **Randomizace**
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Statistika

Statistika je

- matematická disciplína ... — to už víme. Ale také
- každá *měřitelná*¹ funkce G definovaná na náhodném výběru libovolného (dostatečného) rozsahu, tj. počítá se z náhodných veličin výběru, a tudíž *sama je náhodnou veličinou*.
- Obvykle se používá jako *odhad parametrů rozdělení* (které nám zůstávají skryty).

Značení:

θ ... jakákoli hodnota parametru (reálné číslo)

θ^* ... skutečná (správná) hodnota parametru (reálné číslo)

$\hat{\Theta}, \hat{\Theta}_n$... odhad parametru založený na náhodném výběru rozsahu n (náhodná veličina)

$\hat{\theta}, \hat{\theta}_n$... realizace odhadu (reálné číslo)

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- **Statistika**
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

¹ V praxi se setkáte jen s měřitelnými funkcemi. **Měřitelná funkce** G je taková funkce, že pro každé $t \in \mathbb{R}$ je definována pravděpodobnost

$$P[G(X_1, \dots, X_n) \leq t] = F_{G(X_1, \dots, X_n)}(t).$$



Náhodný výběr

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- **Náhodný výběr**
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Náhodný výběr $X = (X_1, \dots, X_n)$ je vektor náhodných veličin, které jsou *nezávislé* a mají *stejné rozdělení* (independent and identically distributed, i.i.d., IID).

Realizace $x = (x_1, \dots, x_n)$ **náhodného výběru** X je výsledkem konkrétního pokusu.

- Popisuje ji *empirické rozdělení*: Vybereme $j \in \{1, \dots, n\}$ s rovnoměrným rozdělením, výsledkem je x_j .
- Je to diskrétní rozdělení, směs Diracových: $\text{Mix}_{(1/n, \dots, 1/n)}(x_1, \dots, x_n)$.

funkce $f : D \rightarrow \mathbb{R}$	funkční hodnota $f(x) \in \mathbb{R}, \quad x \in D$
náhodná veličina $X : \Omega \rightarrow \mathbb{R}$	realizace náhodné veličiny $x := X(\omega) \in \mathbb{R}, \quad \omega \in \Omega$
náhodný vektor/výběr $\mathbf{X} = (X_1, \dots, X_n) : \Omega \rightarrow \mathbb{R}^n$	realizace náhodného vektoru/výběru $\mathbf{x} = (x_1, \dots, x_n) := \mathbf{X}(\omega) \in \mathbb{R}^n, \quad \omega \in \Omega$

Realizace náhodného výběru může mít význam *trénovací množiny*: neznámé parametry odhadujeme tak, aby na trénovací množině byly optimální.



Histogram a empirické rozdělení

V realizaci náhodného výběru $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ *nezáleží na pořadí* hodnot, ale *záleží na jejich četnostech*. Náhodný výběr tak lze popsat

1. množinou (nejvýše n) hodnot $H = \{x_1, \dots, x_n\}$ a
2. jejich četnostmi $n_t, t \in H$,

které se obvykle prezentují ve formě **tabulky četností** nebo **histogramu**.

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- **Histogram a empirické rozdělení**
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Histogram a empirické rozdělení

V realizaci náhodného výběru $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ *nezáleží na pořadí* hodnot, ale *záleží na jejich četnostech*. Náhodný výběr tak lze popsat

1. množinou (nejvýše n) hodnot $H = \{x_1, \dots, x_n\}$ a
2. jejich četnostmi $n_t, t \in H$,

které se obvykle prezentují ve formě **tabulky četností** nebo **histogramu**.

Empirické rozdělení $\text{Emp}(x)$, přesněji jeho psaní funkce $p_{\text{Emp}(x)}$, vznikne normováním četností: $r_t := \frac{n_t}{n} = p_{\text{Emp}(x)}(t)$. Jak uvidíme později:

- Obecné momenty empirického rozdělení jsou rovny výběrovým momentům původního rozdělení.

$$E(\text{Emp}(x))^k = \sum_{t \in H} t^k \cdot r_t = \frac{1}{n} \sum_{t \in H} t^k \cdot n_t = \frac{1}{n} \sum_{j=1}^n x_j^k = m_x^k, \text{ z čehož plyne } E \text{Emp}(x) = \bar{x}.$$

- Rozptyl empirického rozdělení odpovídá odhadu $\widehat{\sigma}_X^2 = \frac{n-1}{n} S_X^2$ rozptylu původního rozdělení, ale odlišnému od S_X^2 .

$$D \text{Emp}(x) = \sum_{t \in H} (t - \bar{x})^2 \cdot r_t = \frac{1}{n} \sum_{t \in H} (t - \bar{x})^2 \cdot n_t = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 = \widehat{\sigma}_x^2 = \frac{n-1}{n} s_x^2.$$

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- **Histogram a empirické rozdělení**
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- **Histogram a empirické rozdělení**
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Jednou jsem potkal hezkou a milou statističku. Tak jsem ji hned požádal o telefonní číslo.

Ale ona mi dala jenom odhad.

Anonym



Odhady

Cílem odhadování je

- najít takovou statistiku (funkci náhodného výběru), jejíž hodnoty *co nejlépe* odpovídají odhadovanému parametru populace.

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- **Odhady**
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Odhady

Cílem odhadování je

- najít takovou statistiku (funkci náhodného výběru), jejíž hodnoty *co nejlépe* odpovídají odhadovanému parametru populace.

Žádoucí vlastnosti:

- $E \hat{\Theta}_n = \theta^*$ **nestranný** (opak: **vychýlený**)
- $\lim_{n \rightarrow \infty} E \hat{\Theta}_n = \theta^*$ **asymptoticky nestranný**
- více, resp. méně **eficientní** = s menším, resp. větším rozptylem, což posuzujeme podle $E \left((\hat{\Theta}_n - \theta^*)^2 \right) = D \hat{\Theta}_n + \left(E \hat{\Theta}_n - \theta^* \right)^2$.
Pro nestranný odhad se redukuje na $D \hat{\Theta}_n$
- **nejlepší nestranný** odhad je ze všech nestranných ten, který je nejvíce eficientní (mohou však existovat více eficientní vychýlené odhady)
- $\lim_{n \rightarrow \infty} E \hat{\Theta}_n = \theta^*$, $\lim_{n \rightarrow \infty} \sigma_{\hat{\Theta}_n} = 0$ **konzistentní**
- **robustní**, tj. odolný vůči šumu („i při zašuměných datech dostáváme dobrý výsledek“) – přesné kritérium chybí, ale je to velmi praktická vlastnost

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- **Odhady**
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Odhady

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- **Odhady**
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Cílem odhadování je

- najít takovou statistiku (funkci náhodného výběru), jejíž hodnoty *co nejlépe* odpovídají odhadovanému parametru populace.

Žádoucí vlastnosti:

- $E \hat{\Theta}_n = \theta^*$ **nestranný** (opak: **vychýlený**)
- $\lim_{n \rightarrow \infty} E \hat{\Theta}_n = \theta^*$ **asymptoticky nestranný**
- více, resp. méně **eficientní** = s menším, resp. větším rozptylem, což posuzujeme podle $E \left((\hat{\Theta}_n - \theta^*)^2 \right) = D \hat{\Theta}_n + \left(E \hat{\Theta}_n - \theta^* \right)^2$.
Pro nestranný odhad se redukuje na $D \hat{\Theta}_n$
- **nejlepší nestranný** odhad je ze všech nestranných ten, který je nejvíce eficientní (mohou však existovat více eficientní vychýlené odhady)
- $\lim_{n \rightarrow \infty} E \hat{\Theta}_n = \theta^*$, $\lim_{n \rightarrow \infty} \sigma_{\hat{\Theta}_n} = 0$ **konzistentní**
- **robustní**, tj. odolný vůči šumu („i při zašuměných datech dostáváme dobrý výsledek“) – přesné kritérium chybí, ale je to velmi praktická vlastnost

Rozlišujeme odhady

- **bodové** (výsledkem je hodnota aproximující skutečnou hodnotu optimálně ve smyslu jistého kritéria) a
- **intervalové** (výsledkem je interval, v němž se skutečná hodnota nachází s danou pravděpodobností).



Výběrový průměr

Výběrový průměr \bar{X} je statistika (náhodná veličina) definovaná jako aritmetický průměr náhodného výběru:

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n X_j$$

Realizace výběrového průměru je rovna aritmetickému průměru realizace náhodného výběru a také střední hodnotě empirického rozdělení:

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j = E \text{Emp}(x)$$

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- **Výběrový průměr**
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Výběrový průměr

Výběrový průměr \bar{X} je statistika (náhodná veličina) definovaná jako aritmetický průměr náhodného výběru:

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n X_j$$

Realizace výběrového průměru je rovna aritmetickému průměru realizace náhodného výběru a také střední hodnotě empirického rozdělení:

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j = E \text{Emp}(x)$$

Platí:

$$E \bar{X}_n = \frac{1}{n} \sum_{j=1}^n E X_j = E X,$$

$$D \bar{X}_n = \frac{1}{n^2} \sum_{j=1}^n D X_j = \frac{1}{n} D X,$$

$$\sigma_{\bar{X}_n} = \sqrt{\frac{1}{n} D X} = \frac{1}{\sqrt{n}} \sigma_X, \text{ pokud existují. (Zde } E X = E X_j \text{ atd.)}$$

Důsledek: Výběrový průměr je *nestranný konzistentní* odhad střední hodnoty (nezávisle na typu rozdělení).

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- **Výběrový průměr**
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Centrální limitní věta

Věta: Výběrový prům. z *normálního* rozdělení $N(\mu, \sigma^2)$ má normální rozdělení $N(\mu, \frac{1}{n}\sigma^2)$.

Centrální limitní věta

Věta: Výběrový prům. z *normálního* rozdělení $N(\mu, \sigma^2)$ má normální rozdělení $N(\mu, \frac{1}{n}\sigma^2)$.

Podobná věta platí i pro jiná rozdělení alespoň asymptoticky.

Centrální limitní věta

Věta: Výběrový prům. z *normálního* rozdělení $N(\mu, \sigma^2)$ má normální rozdělení $N(\mu, \frac{1}{n}\sigma^2)$.

Podobná věta platí i pro jiná rozdělení alespoň asymptoticky.

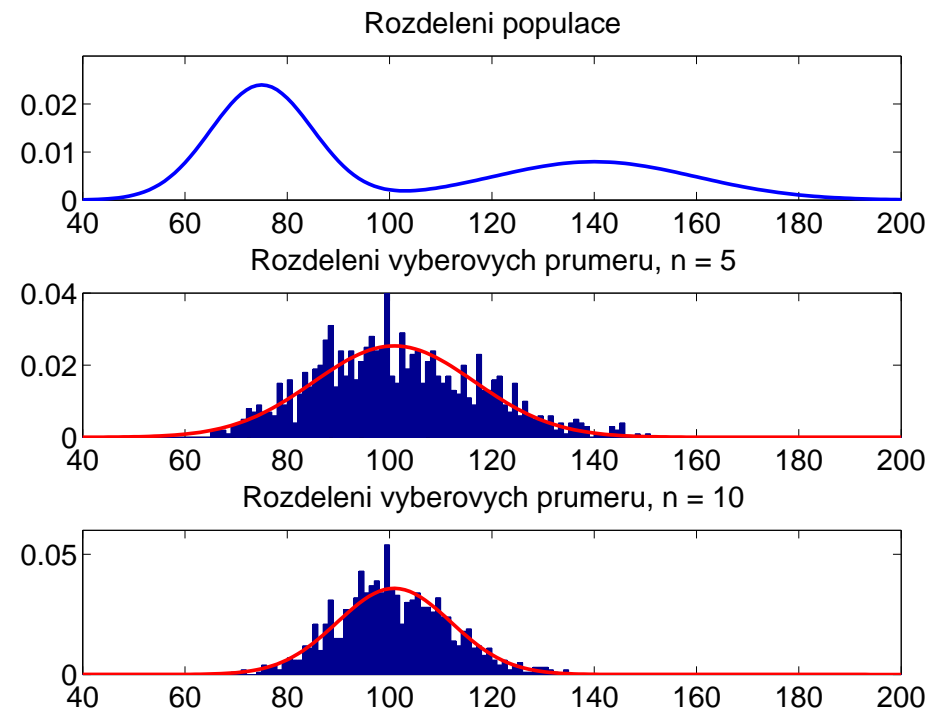
Centrální limitní věta: Necht $X_j, j \in \mathbb{N}$, jsou nezávislé stejně rozdělené náhodné veličiny se střední hodnotou $E X$ a směrodatnou odchylkou $\sigma_X \neq 0$. Pak normované náhodné veličiny

$$Y_n = \text{norm } \bar{X}_n = \frac{\sqrt{n}}{\sigma_X} (\bar{X}_n - E X)$$

konvergují k normovanému normálnímu rozdělení v následujícím smyslu:

$$\forall t \in \mathbb{R} : \lim_{n \rightarrow \infty} F_{Y_n}(t) = \lim_{n \rightarrow \infty} F_{\text{norm } \bar{X}_n}(t) = \Phi(t).$$

Ilustrace:





Výběrový rozptyl

Výběrový rozptyl S_X^2 je statistika (náhodná veličina) definovaná jako

$$S_X^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X}_n)^2$$

a jeho **realizace**:

$$s_x^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x}_n)^2.$$

Věta: Výběrový rozptyl je *nestranný* ($E S_X^2 = D X$) *konzistentní* odhad rozptylu původního rozdělení (má-li původní rozdělení rozptyl a 4. centrální moment).

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- **Výběrový rozptyl**
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Výběrový rozptyl

Výběrový rozptyl S_X^2 je statistika (náhodná veličina) definovaná jako

$$S_X^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X}_n)^2$$

a jeho **realizace**:

$$s_x^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x}_n)^2.$$

Věta: Výběrový rozptyl je *nestranný* ($E S_X^2 = D X$) *konzistentní* odhad rozptylu původního rozdělení (má-li původní rozdělení rozptyl a 4. centrální moment).

POZOR: Odhad rozptylu pomocí rozptylu empirického rozdělení

$$\widehat{\sigma_x^2} = D \text{Emp}(x) = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2,$$

který je realizací odhadu

$$\widehat{\sigma_X^2} = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X})^2,$$

je *vychýlený* (pouze *asymptoticky nestranný*) odhad rozptylu!

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- **Výběrový rozptyl**
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Rozdělení výběrového rozptylu

Speciální případ: pro $N(0, 1)$ a $n = 2$:

$$\bar{X} = \frac{X_1 + X_2}{2}, \quad X_1 - \bar{X} = -(X_2 - \bar{X}) = \frac{X_1 - X_2}{2} \text{ má rozdělení } N\left(0, \frac{1}{2}\right),$$

$$S_X^2 = (X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 = 2 \left(\frac{X_1 - X_2}{2}\right)^2 = \left(\frac{X_1 - X_2}{\sqrt{2}}\right)^2 = U^2,$$

kde $U = \frac{X_1 - X_2}{\sqrt{2}}$ má rozdělení $N(0, 1)$.

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- **Rozdělení výběrového rozptylu**
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Rozdělení výběrového rozptylu

Speciální případ: pro $N(0, 1)$ a $n = 2$:

$$\bar{X} = \frac{X_1 + X_2}{2}, \quad X_1 - \bar{X} = -(X_2 - \bar{X}) = \frac{X_1 - X_2}{2} \text{ má rozdělení } N\left(0, \frac{1}{2}\right),$$

$$S_X^2 = (X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 = 2 \left(\frac{X_1 - X_2}{2}\right)^2 = \left(\frac{X_1 - X_2}{\sqrt{2}}\right)^2 = U^2,$$

kde $U = \frac{X_1 - X_2}{\sqrt{2}}$ má rozdělení $N(0, 1)$. Tomu říkáme:

Rozdělení χ^2 s 1 stupněm volnosti, $\chi^2(1)$, je rozdělení náhodné veličiny $V = U^2$, kde U má normované normální rozdělení $N(0, 1)$.

Vlastnosti:

$$E V = E U^2 = D U + (E U)^2 = 1 \quad (\text{protože } E U = 0, D U = 1),$$

$$D V = 2.$$

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- **Rozdělení výběrového rozptylu**
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Rozdělení χ^2

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- **Chí-kvadrát**
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Rozdělení χ^2 s η stupni volnosti, $\chi^2(\eta)$:

- rozdělení náhodné veličiny $Y = \sum_{j=1}^{\eta} V_j$, kde V_j jsou *nezávislé* náhodné veličiny s rozdělením $\chi^2(1)$.
- rozdělení náhodné veličiny $Y = \sum_{j=1}^{\eta} U_j^2$, kde U_j jsou *nezávislé* náhodné veličiny s *normovaným normálním* rozdělením $N(0, 1)$.

Vlastnosti:

$$E Y = E \sum_{j=1}^{\eta} V_j = \sum_{j=1}^{\eta} E V_j = \eta$$

$$D Y = D \sum_{j=1}^{\eta} V_j = \sum_{j=1}^{\eta} D V_j = 2\eta$$

Věta: Nechť X, Y jsou *nezávislé* náhodné veličiny s rozdělením $\chi^2(\eta)$, resp. $\chi^2(\xi)$. Pak $X + Y$ má rozdělení $\chi^2(\eta + \xi)$.



Výběrový rozptyl z normálního rozdělení

Pro výběrový rozptyl z *normálního* rozdělení $N(E X, D X)$ platí:

$$\frac{(n-1)S_X^2}{D X} = \frac{n\widehat{\sigma}_X^2}{D X} \text{ má rozdělení } \chi^2(n-1).$$

Z vlastností rozdělení χ^2 plyne

- pro střední hodnotu výběrového rozptylu

$$E \frac{(n-1)S_X^2}{D X} = n-1 \text{ takže}$$

$$E S_X^2 = D X, \quad \text{což potvrzuje nestrannost odhadu.}$$

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- **Chí-kvadrát**
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Výběrový rozptyl z normálního rozdělení

Pro výběrový rozptyl z *normálního* rozdělení $N(E X, D X)$ platí:

$$\frac{(n-1)S_X^2}{D X} = \frac{n\widehat{\sigma}_X^2}{D X} \text{ má rozdělení } \chi^2(n-1).$$

Z vlastností rozdělení χ^2 plyne

- pro střední hodnotu výběrového rozptylu

$$E \frac{(n-1)S_X^2}{D X} = n-1 \text{ takže}$$

$$E S_X^2 = D X, \quad \text{což potvrzuje nestrannost odhadu.}$$

- pro rozptyl výběrového rozptylu

$$D \frac{(n-1)S_X^2}{D X} = 2(n-1),$$

$$\frac{(n-1)^2 D S_X^2}{(D X)^2} = 2(n-1), \quad \text{takže}$$

$$D S_X^2 = \frac{2}{n-1} (D X)^2.$$

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- **Chí-kvadrát**
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Výběrový rozptyl z normálního rozdělení

Pro výběrový rozptyl z *normálního* rozdělení $N(E X, D X)$ platí:

$$\frac{(n-1)S_X^2}{D X} = \frac{n\widehat{\sigma}_X^2}{D X} \text{ má rozdělení } \chi^2(n-1).$$

Z vlastností rozdělení χ^2 plyne

- pro střední hodnotu výběrového rozptylu

$$E \frac{(n-1)S_X^2}{D X} = n-1 \text{ takže}$$
$$E S_X^2 = D X, \quad \text{což potvrzuje nestrannost odhadu.}$$

- pro rozptyl výběrového rozptylu

$$D \frac{(n-1)S_X^2}{D X} = 2(n-1),$$
$$\frac{(n-1)^2 D S_X^2}{(D X)^2} = 2(n-1), \quad \text{takže}$$
$$D S_X^2 = \frac{2}{n-1} (D X)^2.$$

Věta: Pro náhodný výběr X_n z *normálního* rozdělení je \bar{X} nejlepší nestranný odhad střední hodnoty, S_X^2 je nejlepší nestranný odhad rozptylu a statistiky \bar{X} a S_X^2 jsou konzistentní a *nezávislé*.

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- **Chí-kvadrát**
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Výběrová směrodatná odchylka

Výběrová směrodatná odchylka S_X je statistika (náhodná veličina) definovaná jako

$$S_X = \sqrt{S_X^2} = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X}_n)^2}$$

a její **realizace**:

$$s_x = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x}_n)^2}.$$

Věta: Výběrová směrodatná odchylka je *vychýleným* ($E S_X \leq \sigma_X$) *konzistentním* odhadem směrodatné odchylky původního rozdělení (má-li původní rozdělení rozptyl a 4. centrální moment).

Důkaz: $D X = E S_X^2 = (E S_X)^2 + D S_X$, a protože $D S_X \geq 0$, tak $D X \geq (E S_X)^2$, takže $\sigma_X \geq E S_X$.

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Výběrová směrodatná odchylka

Výběrová směrodatná odchylka S_X je statistika (náhodná veličina) definovaná jako

$$S_X = \sqrt{S_X^2} = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X}_n)^2}$$

a její **realizace**:

$$s_x = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x}_n)^2}$$

Věta: Výběrová směrodatná odchylka je *vychýleným* ($E S_X \leq \sigma_X$) *konzistentním* odhadem směrodatné odchylky původního rozdělení (má-li původní rozdělení rozptyl a 4. centrální moment).

Důkaz: $D X = E S_X^2 = (E S_X)^2 + D S_X$, a protože $D S_X \geq 0$, tak $D X \geq (E S_X)^2$, takže $\sigma_X \geq E S_X$.

POZOR: Odhad směrodatné odchylky pomocí sm. odch. empirického rozdělení

$$\hat{\sigma}_x = \sigma_{\text{Emp}(x)} = \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2}$$

je taktéž *vychýlený* odhad směrodatné odchylky!

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy



Výběrový medián

Výběrový medián je statistika (náhodná veličina), která se používá jako odhad mediánu původního rozdělení. Je to 50% kvantil empirického rozdělení, $q_{\text{Emp}(x)}\left(\frac{1}{2}\right)$.

- Je *robustnější* než výběrový průměr (odolnější vůči vlivu odlehlých hodnot).
- Víme, jak se změní po transformaci monotónní funkcí.
- Má vyšší výpočetní náročnost než výběrový průměr: seřazení hodnot má náročnost $\mathcal{O}(n \log n)$, výpočet průměru jen n .
- Má vyšší paměťovou náročnost než výběrový průměr: musíme si pamatovat všech n čísel, u průměru stačí 2 registry.
- Špatně se decentralizuje / paralelizuje.

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- **Výběrový medián**
- Míry polohy

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Která míra polohy je vhodná? Modus, medián nebo průměr?

Příklad: Rozdělení platů v republice. Je poměrně velká část populace, která nedostává žádný plat (děti, důchodci, nezaměstnaní, lidé, co pracovat nechtějí).

- Modus je

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Která míra polohy je vhodná? Modus, medián nebo průměr?

Příklad: Rozdělení platů v republice. Je poměrně velká část populace, která nedostává žádný plat (děti, důchodci, nezaměstnaní, lidé, co pracovat nechtějí).

- Modus je 0. Největší část populace nedostává plat. Nestane-li se s platy něco opravdu převratného, zůstane modus nulový. (Což nemusí platit pro jiné veličiny.)
- Medián je

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Která míra polohy je vhodná? Modus, medián nebo průměr?

Příklad: Rozdělení platů v republice. Je poměrně velká část populace, která nedostává žádný plat (děti, důchodci, nezaměstnaní, lidé, co pracovat nechtějí).

- Modus je 0. Největší část populace nedostává plat. Nestane-li se s platy něco opravdu převratného, zůstane modus nulový. (Což nemusí platit pro jiné veličiny.)
- Medián je asi nejlepší ukazatel toho, co si vydělá typický občan. I když se vysoké platy 10x zvýší, zůstane stejný.
- Průměr je

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Která míra polohy je vhodná? Modus, medián nebo průměr?

Příklad: Rozdělení platů v republice. Je poměrně velká část populace, která nedostává žádný plat (děti, důchodci, nezaměstnaní, lidé, co pracovat nechtějí).

- Modus je 0. Největší část populace nedostává plat. Nestane-li se s platy něco opravdu převratného, zůstane modus nulový. (Což nemusí platit pro jiné veličiny.)
- Medián je asi nejlepší ukazatel toho, co si vydělá typický občan. I když se vysoké platy 10x zvýší, zůstane stejný.
- Průměr je dobrá míra pro ekonomiku a daňové úřady, protože z něj mohou spočítat celkový plat. Není to ale dobrá míra typického platu, je snadno ovlivnitelný (především vysokými platy).

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Která míra polohy je vhodná? Modus, medián nebo průměr?

Příklad: Rozdělení platů v republice. Je poměrně velká část populace, která nedostává žádný plat (děti, důchodci, nezaměstnaní, lidé, co pracovat nechtějí).

- Modus je 0. Největší část populace nedostává plat. Nestane-li se s platy něco opravdu převratného, zůstane modus nulový. (Což nemusí platit pro jiné veličiny.)
- Medián je asi nejlepší ukazatel toho, co si vydělá typický občan. I když se vysoké platy 10x zvýší, zůstane stejný.
- Průměr je dobrá míra pro ekonomiku a daňové úřady, protože z něj mohou spočítat celkový plat. Není to ale dobrá míra typického platu, je snadno ovlivnitelný (především vysokými platy).

Teď už přesně víme, co je průměr a co je medián. Takže nás vůbec nepřekvapí, že:

- Naprostá většina lidí má nadprůměrný počet nohou. *Je to tak?*

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Která míra polohy je vhodná? Modus, medián nebo průměr?

Příklad: Rozdělení platů v republice. Je poměrně velká část populace, která nedostává žádný plat (děti, důchodci, nezaměstnaní, lidé, co pracovat nechtějí).

- Modus je 0. Největší část populace nedostává plat. Nestane-li se s platy něco opravdu převratného, zůstane modus nulový. (Což nemusí platit pro jiné veličiny.)
- Medián je asi nejlepší ukazatel toho, co si vydělá typický občan. I když se vysoké platy 10x zvýší, zůstane stejný.
- Průměr je dobrá míra pro ekonomiku a daňové úřady, protože z něj mohou spočítat celkový plat. Není to ale dobrá míra typického platu, je snadno ovlivnitelný (především vysokými platy).

Teď už přesně víme, co je průměr a co je medián. Takže nás vůbec nepřekvapí, že:

- Naprostá většina lidí má nadprůměrný počet nohou. *Je to tak?*
- Polovina populace má inteligenci nižší, než medián. Většina populace má ovšem nadprůměrnou inteligenci. To je jen důsledek toho, že lidská inteligence je shora omezená. Pro lidskou hloupost ovšem žádné limity neexistují. *Jak by muselo vypadat rozdělení inteligence v populaci, aby to pravda byla?*

Míry polohy

1. *Výběrový modus* je nejčastější hodnota ve výběru. Lze jej stanovit i pro celou populaci.
2. *Výběrový medián* je hodnota, pod níž (i nad níž) leží 50% hodnot. 50%-ní kvantil, 50. percentil, 5. decil.
3. *Výběrový průměr* je takový díl, že naskládáme-li jich za sebe stejný počet, jako je původních hodnot, dostaneme stejný součet, jako dávají původní hodnoty.

Která míra polohy je vhodná? Modus, medián nebo průměr?

Příklad: Rozdělení platů v republice. Je poměrně velká část populace, která nedostává žádný plat (děti, důchodci, nezaměstnaní, lidé, co pracovat nechtějí).

- Modus je 0. Největší část populace nedostává plat. Nestane-li se s platy něco opravdu převratného, zůstane modus nulový. (Což nemusí platit pro jiné veličiny.)
- Medián je asi nejlepší ukazatel toho, co si vydělá typický občan. I když se vysoké platy 10x zvýší, zůstane stejný.
- Průměr je dobrá míra pro ekonomiku a daňové úřady, protože z něj mohou spočítat celkový plat. Není to ale dobrá míra typického platu, je snadno ovlivnitelný (především vysokými platy).

Teď už přesně víme, co je průměr a co je medián. Takže nás vůbec nepřekvapí, že:

- Naprostá většina lidí má nadprůměrný počet nohou. *Je to tak?*
- Polovina populace má inteligenci nižší, než medián. Většina populace má ovšem nadprůměrnou inteligenci. To je jen důsledek toho, že lidská inteligence je shora omezená. Pro lidskou hloupost ovšem žádné limity neexistují. *Jak by muselo vypadat rozdělení inteligence v populaci, aby to pravda byla?*

Už jste slyšeli o tom politikovi, který ve své předvolební kampani sliboval, že se po svém zvolení zasadí o to, aby měl každý občan nadprůměrný příjem?



Místo závěru

Úvod do statistiky

- Pravděpodobnost a statistika
- Účel
- Analýza dat
- Základní pojmy
- Druhy veličin
- Šetření vs. exp.
- Randomizace
- Statistika
- Náhodný výběr
- Histogram a empirické rozdělení
- Odhady
- Výběrový průměr
- Cent. lim. věta
- Výběrový rozptyl
- Rozdělení výběrového rozptylu
- Chí-kvadrát
- Výběrová sm.odch.
- Výběrový medián
- Míry polohy

Existují tři druhy lži: lži, naprosté lži a statistiky.

Benjamin Disraeli