



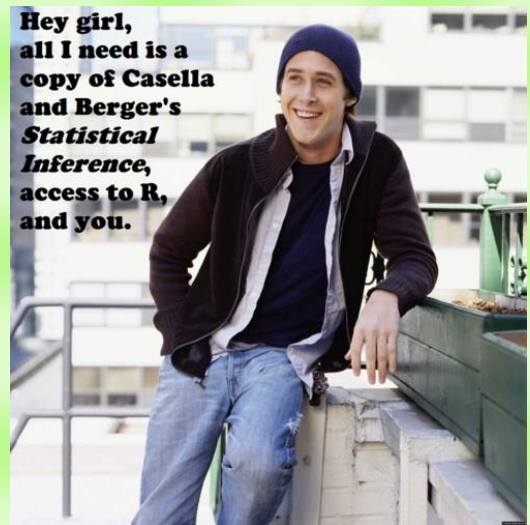
BIOSTATISTIKA

**Ze skript Základy statistiky v prostředí R, přednášek, internetu, a
vlastního úsudku sestavila Majda. Nedoporučuji číst biostatistikům.**

POPISNÁ STATISTIKA

- Měřítko

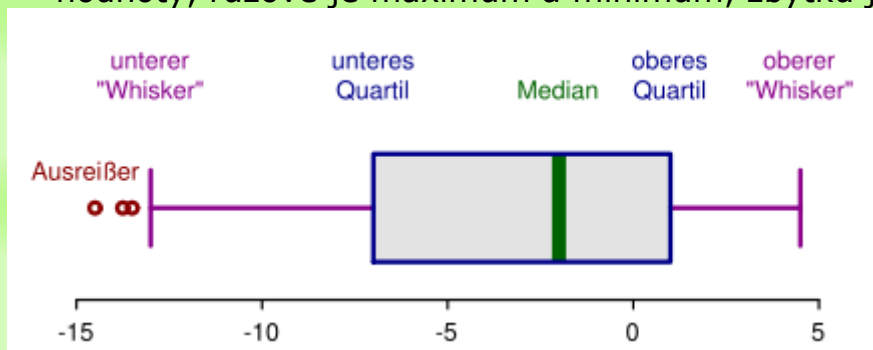
- 0-1 – jen dvě varianty, př. pohlaví, kouří/nekouří
- nominální – více nečíselných možností, tzv. faktor, př. město, odrůda
- ordinální – uspořádaný faktor=uspořádané hodnoty nominálního měřítka, př. stupeň vzdělání
- poměrové – srovnání s jednotkou – kolikrát je x větší než y, př. váha, výška, věk
- intervalové – stejné vzdálenosti sousedních – o kolik je větší, př. rok narození
- hrubší dělení měřítek
 - kvalitativní (diskrétní veličiny) – tam patří: 0-1, nominální, často ordinální – udávají se četnosti hodnot – kolikrát která nastala
 - kvantitativní (spojité veličiny) – tam patří: intervalové, poměrové, někdy ordinální, hodnoty jsou čísla



- Míry polohy – pomocí nich charakteristika vzorku jedním číslem – ukazují na velikost hodnot

- Variační řada - hodnoty uspořádané od nejmenší po největší – z ní odvodíme min, max – ty jsou ale pro popis vzorku dost nevhodné. Variační řada se značí x a číslo v závorce – $X(1)$, $X(2)$,... - z ní odvodíme pořadí – $X(1)$ má pořadí 1 atd... Když jsou čísla vedle sebe stejná, dostanou průměrnou hodnotu pořadí – např. 3,5, 3,5..; v řadě 3, 3, 3 je pořadí 2, 2, 2!
- Průměr, vážený průměr (s využitím četností – vše vynásobíme vahami jednotlivých hodnot) - \bar{x} s čárkou
- Medián – dělí variační řadu na 2 stejné části – větší a menší než medián, když je jich sudý počet, je medián průměr ze středních dvou, - \tilde{x} s vlnovkou
- Minimum, maximum
- Horní kvartil, dolní kvartil – mediány pro $\frac{1}{4}$, značka Q_1 , Q_3 , první decil
- Výběrový kvantil x_p = p -tý percentil – odděluje p -té procento nejmenších hodnot od hodnot větších.. Medián je percentil $x_{0,5}$, kvartily $x_{0,25}$ a $x_{0,75}$!
- Modus – x se stříškou – nejčastěji se vyskytující hodnota, stříškou se ve statistice také značí odhad

- **vlastnosti** – přičítání i násobení konstantou u každého prvku – pak i každá z měř polohy (obecná míra polohy – mí) stejná, jen změněná o konstantu
- graf – box-plot – červené divné slovo znamená odlehlé hodnoty, růžově je maximum a minimum, zbytku je rozumně



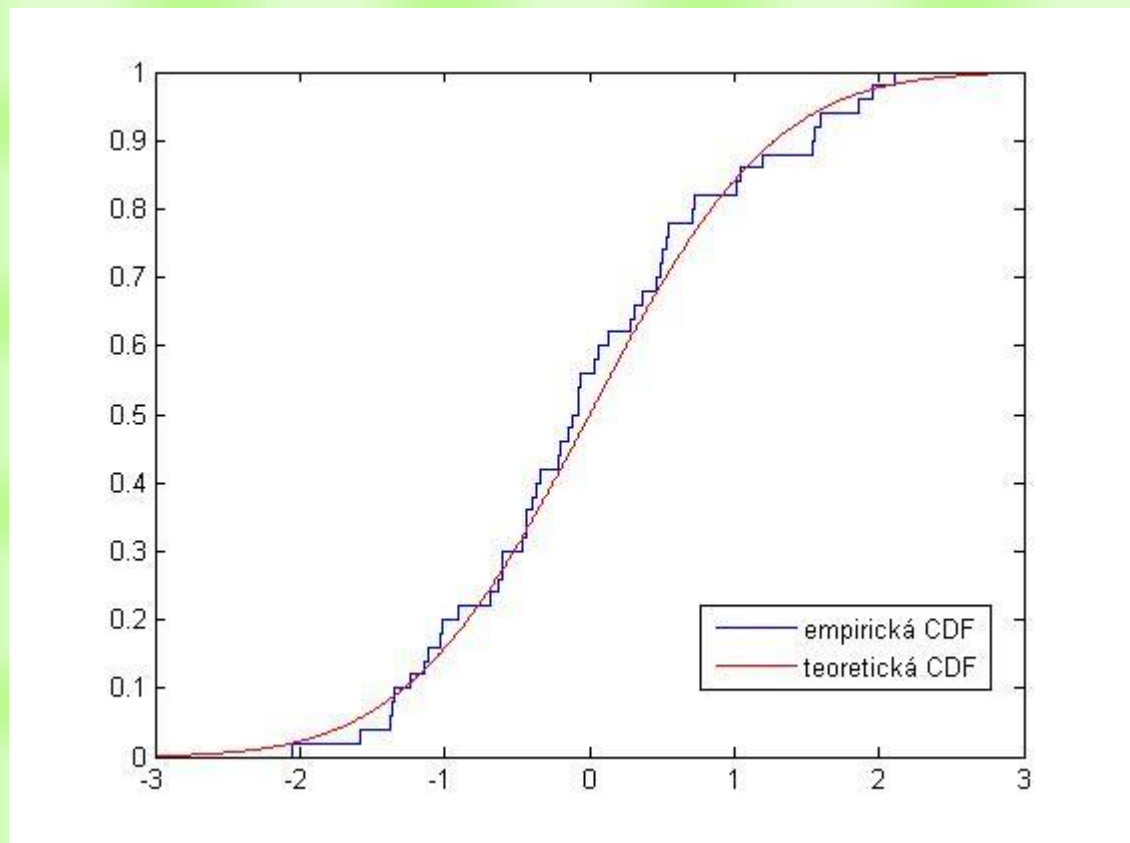
- Míry variability

- číselně charakterizuje jinou vlastnost než míra polohy, ukazuje kolísání, variabilitu
- **Podmínky:** Nesmí být měněny úrovně, ve které počítáme míry polohy, protože přičtu-li ke všem prvkům konstantu a , míry variability se nemění, vynásobím-li všechny prvky konstantou b (větší než 0), směrodatnou odchylku také vynásobím touto konstantou, rozptyl konstantou na druhou (viz níže) a ten zbytek radši ani nevěděť..
- směrodatná odchylka (SD, standard deviation) - s_x – odmocnina rozptylu, platí pro ni obě podmínky.
- (výběrový) rozptyl - s_x^2 – téměř průměr druhých mocnin vzdáleností konkrétních hodnot od průměru – jak moc kolísají jednotlivé hodnoty kolem průměru, také se dá chápat jako míra vzájemné rozdílnosti hodnot x_1-x_n ; platí pro něj pravidla o měřách variability, ale když násobím prvky konstantou b , musím roznásobit rozptyl b^2
- rozpětí - $R = x_{\text{Max}} - x_{\text{min}}$, $x_n - x_1$
- kvartilové rozpětí - $R_q = Q_3 - Q_1$ – lepší, není ovlivněno odlehlými hodnotami – interval, co obsahuje prostřední část dat
- variační koeficient - $V_x = s_x / \bar{x}$ prům. – nesplňuje podmínky nahoře – buď to kazí průměr, nebo se to vykrátí a nezbyde nic, pro co by se podmínka dělala 😊.. Používá se pro porovnání hodně polohou různých souborů dat

- Z-skór

- založený na upravených původních číselných hodnotách, jeho hodnoty nezávisí na míře polohy ani variability
- z-skór - z_i říká, jak daleko je hodnota x_i od průměru všech x . jednotka vzdálenosti je směrodatná odchylka s_x . Kladná hodnota $= x_i$ je větší než průměr, záporná naopak
- $z_i = (x_i - \bar{x}) / s_x$
- používá se při určování normality – šikmost a špičatost

- průměr z-skóru = 0, směrodatná odchylka = 1
- **Šikmost a špičatost**
 - k posouzení, jestli číselná hodnota pochází z normálního rozdělení
 - šikmost – průměr 3.mocnin z-skórů, vyjadřuje symetrii rozložení pozorování kolem průměrné hodnoty
 - špičatost – průměr 4.mocnin z-skórů zmenšený o 3
- **Empirická distribuční funkce**
 - lze ji sestavit pro ordinální a spojitou náhodnou proměnou, pro dané x je dána relativní četností hodnot menších nebo stejných než x
 - ke každému x se zjistí, kolik hodnot je menších nebo stejně velkých jako x , a to číslo se dělí celkovým počtem hodnot.
 - Je odhadem distribuční funkce



TEORIE

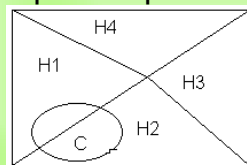
- **Náhodné jevy**
Náhodný pokus – dělá se za přesných podmínek a jeho výsledek je nejistý. S počtem opakování pokusu je ale četnost určitého výsledku stále stabilnější – např. stále stabilněji se rodí 50% holčiček

se ptát, jaká bude pravděpodobnost, že je určený nemocný skutečně nemocný a že zdravý s negativním výsledkem je skutečně zdravý. A to vše chceme zjistit z informace o existenci nebo neexistenci nemoci. Vypočítá nám to Bayesův vzorec: $P(\text{nemocného při pozitivním testu}) = \frac{\text{senzitivita} \times \text{prevalence}}{\text{senzitivita} \times \text{prevalence} + (1 - \text{senzitivita}) \times (1 - \text{prevalence})}$



Obecně:

- $P(A/B) = P(B/A) \times P(A) / (P(B/A) \times P(A) + P(B/\text{non}A) \times P(\text{non}A))$
- H_j – Rozdělení jevu jistého na x neslučitelných jevů, z nichž jeden musí platit. Těmto jevům se říká hypotézy.
- známe apriorní psti $P(H_j)$
- známe pravděpodobnost $P(C|H_j)$



Když spojíme hypotézy, Bayesův vzorec a úplnou pravděpodobnost pro počítání s hypotézami, dostaneme tvar: $P(H_j/B) = P(B/H_j) \times P(H_j) / \text{suma od všech hypotéz } P(B/H_i) \times P(H_i)$

Podmíněnou pravděpodobnost B/H_i známe, apriorní pravděpodobnost H_i také, a protože máme podmínku, že nastal jev B , dopočítáme aposteriorní pravděpodobnost (H_j/B) ☺ ☺ ☺

- Úplná pravděpodobnost

- vzorec pro úplnou pravděpodobnost klíčový pro Bayesův vzorec.
- Pravděpodobnost, s jakou náhodně zvolená osoba pozitivně reaguje na test – sjednocení pozitivních zdravých a pozitivních nemocných – jsou neslučitelné, takže je sečteme. Tento součet je úplná pravděpodobnost

- $P(C)$ je váženým průměrem podmíněných pravděpodobností

$P_{\text{pozit.}} = \text{senzitivita} \times \text{prevalence} + P(\text{zdravého, co reaguje pozitivně}) \times \text{část populace, kde není prevalence}$

Když počítáme s více neslučitelnými jevy – s hypotézami – $P(B)$ úplná pravděpodobnost = suma všech podmíněných pravděpodobností $P(B/H_i) \times$ apriorní pravděpodobnost $P(H_i)$. Je to tedy vážený průměr podmíněných pravděpodobností, který je vážený apriorní pravděpodobností $P(H_i)$.

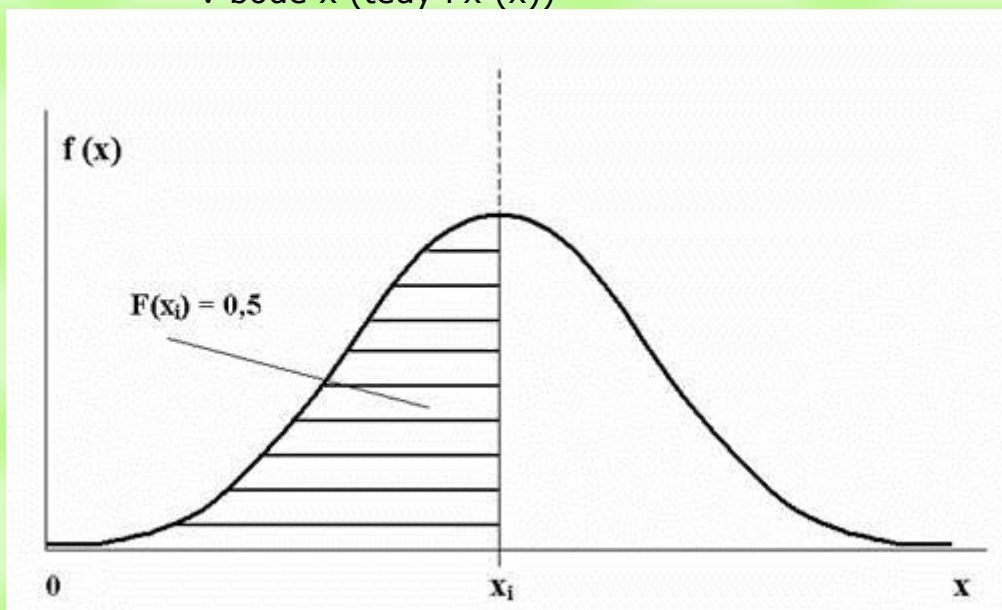
- Náhodná veličina

- číselně vyjádřený výsledek náhodného pokusu – např. jakou váhu má náhodně vybraný muž? Kolikrát padne šestka?
- předem nevím, který výsledek vyjde, známe jen možné hodnoty a jejich pravděpodobnosti – ty musí mít dohromady součet 1
- může mít diskrétní – vyjadřuje počet případů = četnost a pravděpodobnost těchto případů; nebo spojitě rozdělení

- Distribuční funkce

můžeme pomocí ní popsat chování náhodné veličiny

- protějšek empirické distribuční funkce
- distribuční funkce $F(x)$ je nejúplnější popis pravděpodobnostního chování diskrétní nebo spojitě náhodné proměnné X
- $P(x_1 < X \leq x_2) = F(x_2) - F(x_1)$
- $F(x) = P(X \leq x)$ – hodnota distribuční funkce v bodě x = pravděpodobnost, že náhodná veličina nepřekročí x
- Když je náhodná veličina spojitá, pracujeme s hustotou jejího rozdělení (namísto pravděpodobnosti) – pravděpodobnost, že $a < X < b$ je potom rovna velikosti plochy pod grafem na intervalu a až b , celá plocha pod grafem je rovna 1 a plocha od mínus nekonečno k x je rovna hodnotě distribuční funkce v bodě x (tedy $F(x)$)



- Střední hodnota mí náhodné veličiny X s diskrétním rozdělením (označení jako návod na použití – divné ale viz níže - $E X$)

- je to míra polohy, populační průměr (idealizovaný protějšek průměru)
- počítá se z ní populační směrodatná odchylka a populační rozptyl

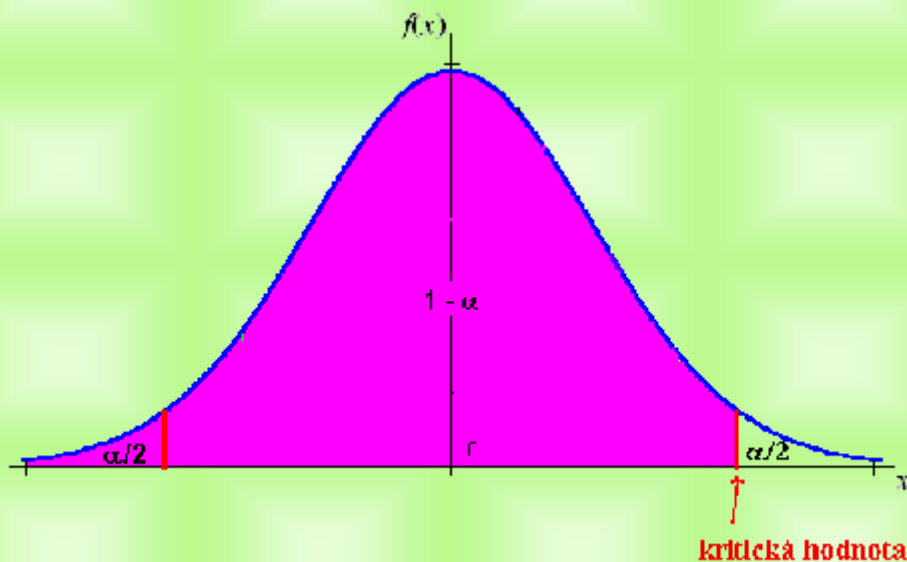
- je to jakési těžiště možných hodnot náhodné veličiny
 $E X = \text{suma všech hodnot vynásobených pravděpodobnostmi vždy dané hodnoty; vážený průměr možných hodnot, kde váhami jsou pravděpodobnosti těchto hodnot}$
 Z toho odvodíme i populační medián – hodnota, pro kterou bude hodnota distribuční funkce F_X 1/2

- **Kvantily a kritické hodnoty**

Kvantil (úplněji 100.p)% kvantil

Zvolíme (podle toho, co chceme otestovat) si číslo p mezi 0 a 1 a pro dané číslo (náhodnou veličinu) X najdeme takové číslo $x(p)$ – kvantil, pro které je distribuční funkce čísla X rovná p – pravděpodobnost, že je náhodná veličina menší nebo rovna kvantilu – taková konstanta, kterou s danou pravděpodobností náhodná veličina nepřekročí

- kritická hodnota- taková konstanta, kterou s danou pravděpodobností náhodná veličina překročí – 90% kvantil=10% kritická hodnota



- **Definice nezávislosti dvou náhodných jevů**
- populační rozptyl - σ^2 – to je značka, označení **var** je návod na to, že máme použít rozptyl, jeho odmocnina je směrodatná odchylka – splňují požadavky na míry variability stejně jako jejich nepopulační obdoby
 - pst jevu D nezávisí na tom, zda platí jev C
- $P(D \cap C) = P(D) P(C)$
 Náhodné veličiny X a Y jsou nezávislé, když jsou navzájem nezávislé všechny dvojice jevů A a B , kdy A je jev, který patří k X a B patří k Y - buď se jevy X a Y náhodou vyskytnou spolu = sdružené rozdělení – udává se pravděpodobností $P=(X=x_i, Y=y_j)$ nebo každý zvlášť – s pravděpodobností $P(X=x_i)$ nebo $P(Y=y_j)$ – tomu se říká marginální rozdělení. Pokud se pravděpodobnost sdruženého rozdělení rovná násobku pravděpodobností marginálních rozdělení,

jsou náhodné veličiny nezávislé ☺ ☺ - tehdy můžeme rekonstruovat sdružené rozdělení z rozdělení marginálních. Pokud jsou náhodné veličiny závislé, už můžeme rekonstruovat jen marginální ze sdružených – to znám z chí kvadrát testu (viz níže) – ty uprostřed sdružené, když je počítám na řádku/sloupci, vzniknou marginální

- **sdružené rozdělení** – popisuje společné chování veličin X a Y
- **marginální rozdělení** – chování jedné bez ohledu na druhou
- **podmíněné rozdělení** – popisuje chování Y při dané hodnotě X

- Kovariance

○ vyjadřuje vzájemnou závislost dvojice náhodných veličin
 $\text{cov}(X,Y) = \text{střední hodnota } ((x_i - \bar{x})(y_i - \bar{y}))$
 platí, že $\text{cov}(X,X) = \text{var}(X)$

Jsou li náhodné veličiny X, Y nezávislé, pak je vždy $\text{cov}(X,Y) = 0$ – kovariance je tedy ukazatelem možné závislosti!

Kovariance nám ale vůbec neukazuje těsnost nějakého vztahu, protože se její hodnoty mění s nastaveným měřítkem.. Pro vyjádření těsnosti vztahu se tedy používá populační korelační koeficient (je idealizovanou populační obdobou výběrového korelačního koeficientu – oba jsou bezrozměrné ukazatele), značí se ρ .. (óóóó, jaká podobnost s rozptylem!!)

Rozdělení

- Hypergeometrické rozdělení

- Náhodná veličina Y, označující počet vybraných jednotek ze všech vykazujících sledovanou vlastnost, se řídí hypergeometrickým rozdělením. Např. 20 otázek na test, 5 z nich student umí (vykazují danou vlastnost), dostane 3, Y= jaká je pravděpodobnost, že si vytáhne alespoň jednu otázku, co umí?
- je to rozdělení náhodné veličiny, kdy při opakování náhodného pokusu je výskyt sledovaného jevu závislý na výsledcích předchozích pokusů, jde tedy o pokusy, které jsou na sobě závislé
- máme ryby v rybníku – je jich tam m, vylovíme jich a, ty označíme, pustíme a vylovíme n ryb. Nějaké z nich budou označené Y– bude jich více než 0, méně než n a méně než a, také musí být neoznačených méně než m-a, tzn $Y = n - m + a$ – rozsah hodnot, jakých to může nabývat, potom s tímto rozsahem počítáme, kolika způsoby můžeme ryby vylovit – z toho vyplyne hypergeometrické rozdělení, z veličin a, m, n lze vypočítat i střední hodnotu a rozptyl Y, ze střední hodnoty pak vznikne odhad – $m = na/Y$

- Binomické rozdělení

- zjišťujeme pouze výskyt či nevýskyt jevu B
- pokusy jsou nezávislé (jejich počet n)

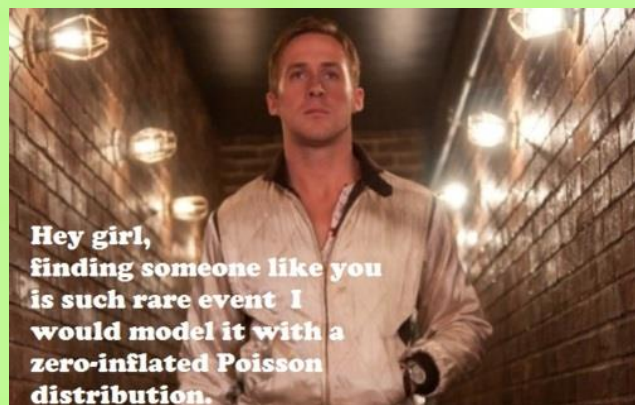
- o pravděpodobnost výskytu B je v každém pokuse stejná
pravděpodobnost $6 - 1/6$, házíme 12x, počet padlých šestek je náhodná veličina Y, který nabývá hodnot od 0 do 12. Určí se pravděpodobnosti všech hodnot – tyto pravděpodobnosti charakterizují binomické rozdělení – Y má binomické rozdělení ($Y = bi(n, p)$) – posloupnost n nezávislých pokusů, kdy každý může skončit zdarem nebo nezdarem. Náhodná veličina Y udává náhodný počet pokusů, v nichž nastal zdar.

- Multinomické rozdělení

- o zobecnění binomického rozdělení (multinomické rozdělení pro $k=2$ (2 možné situace) je binomické, stejně tak, když označíme jednu z marginálních pravděpodobností jedné situace jako zdar a zbytek jako nezdar, máme zase rozdělení binomické)
- o v dílčím pokusu máme k možných výsledků (jsou to jevy neslučitelné, dohromady dávají jev jistý)
Máme n pokusů, v každém nastane jeden z k výsledků ($A_1 - A_k$), pravděpodobnosti toho, že nastane $A_1 - A_k$ se v průběhu opakovaných pokusů nemění. Jako $Y_1 - Y_k$ označujeme počty pokusů, v kterých nastaly jevy $A_1 - A_k$. Hodnoty Y pak mají multinomické rozdělení. V každém případě nastane právě jeden z jevů A, proto: $Y_1 + \dots + Y_k = n$. Y jsou závislé, když známe k-1 hodnot zjistíme všechny (z počtu n). Tzn. K-tice náhodných veličin Y má k-1 stupňů volnosti!
-Velmi časté rozdělení – vždy, když třídíme četnosti znaku v nominálním měřítku – např. počty lidí s krevní skupinou
-Má nějakou tajemnou spojitost s Chí-kvadrát rozdělením – kdo to dočte až dolů, tak mu to možná docvakne.

- Poissonovo rozdělení

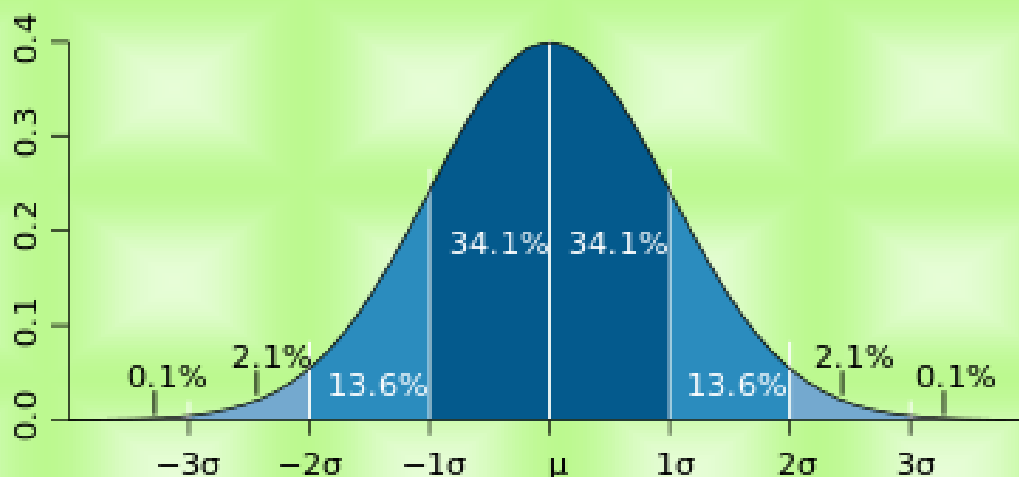
- o =zákon řídkých jevů
– aproximuje (přibližuje) nám binomické rozdělení, když jde o jevy, které se vyskytují jen zřídka – velký počet pokusů s malou pravděpodobností výskytu jevu
- o Střední hodnota i rozptyl náhodné veličiny s Poissonovým rozdělením jsou stejné
- o má ho náhodná veličina, která nabývá hodnot $k = 0, 1, \dots$ (neurčeno), na rozdíl od binomického neznáme maximální hodnotu náhodné veličiny (u binomického to bylo n – počet členů) Y modeluje počet výskytu sledovaného jevu. Jevy zjišťujeme v oddělených časových úsecích, na různých plochách, ... a musí být nezávislé. Parametr lambda, který je



s rozdělením spojený musí být úměrný velikosti úseku, na kterém jevy zkoumáme (např. plochy, času,..)

- Normální rozdělení (Gaussovo)

- o spojité rozdělení, symetrické okolo střední hodnoty



- o Graf – zvon, maximum pro x ve střední hodnotě, výška se rovná převrácené směrodatné odchylce – wooow
- o Plocha pod grafem je rozdělená na 6 částí, jejich základny jsou vždy dlouhé sigma, velikosti těchto částí jsou nezávislé na parametrech – procenta vždy stejná - zase wooow
- o Když veličinu násobíme nebo k ní něco připočítáváme, má stále normální rozdělení
- o Normované normální rozdělení – $\mu=1$, $\sigma=0$. Normovaná normální veličina Z z tohoto rozdělení je s 95% pravděpodobností v intervalu $(-1,96; 1,96)$ – vyplývá z výpočtu kvantilu $z(0,975)=1,96$
- o Normální rozdělení je velmi časté v důsledku centrální limitní věty (viz níže) a je prostě super.

-Náhodný výběr

Už byl náhodný jev a náhodný pokus, né že by se opakovalo pořád to samé dokola

Z x náhodných realizací veličiny X se náhodně vybraná n -tice nazývá náhodný výběr – např. při vážení téhož předmětu

Střední hodnota μ z náhodně vybraných X zastupuje skutečnou hodnotu veličiny X , převrácená hodnota sigma (směrodatné odchylky) zase přesnost vážení. μ v této situaci říkáme populační průměr (průměr X z náhodného výběru je roven parametru μ ; je to tedy nestranný odhad parametru μ).

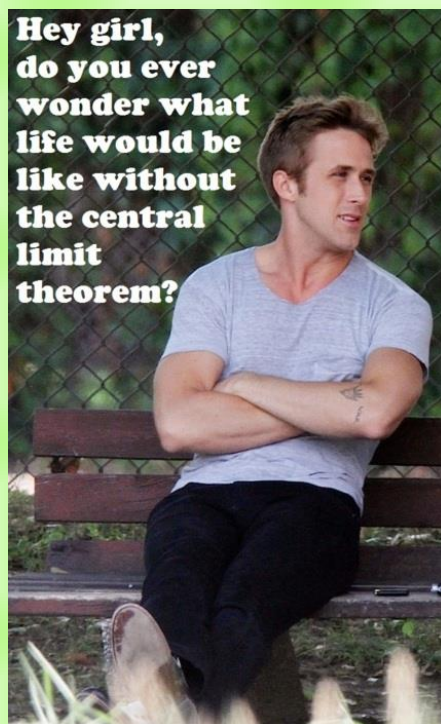
- **S.E.(X) Střední chyba průměru** (To X má čárku nahoře – protože je to průměr a ne něco jiného..)
- Směrodatná odchylka průměru v případě měření s náhodným výběrem z populace. Zjednodušeně řečeno je to číslo, které

označuje, jak moc se asi námi získaný průměr náhodného výběru liší od střední hodnoty základního souboru.

- S.E. ($X_{\text{prům.}}$) = σ / \sqrt{n}
- Průměr tedy bude kolísat za použití náhodného výběru mnohem méně, než by kolísaly jednotlivé náhodné výběry
- Pod tímto pojmem se chápe i odhad střední chyby průměru, který se pak počítá S.E. ($x_{\text{prům.}}$) = Směrodatná odchylka S / \sqrt{n}

- **Centrální limitní věta**

- s rostoucím n se rozdělení náhodné veličiny Z blíží k normovanému normálnímu rozdělení. To znamená, že průměr X má přibližně normální rozdělení
- máme hodnoty, co rozhodně nemají normální rozdělení, mají kladnou šikmost atd.. Když z hodnot náhodně vybereme tisíckrát nějakou a uděláme histogram, trochu se znormální. Když tisíckrát vybereme 10 hodnot, které zprůměrujeme a z těch tisíců průměrů uděláme histogram, bude ještě normálnější, když budeme průměrovat 100 hodnot, tak bude už úplně super.. (Rozptyl se bude neustále zmenšovat)
 - o průměr má pro dost velká n normální rozdělení s rozptylem n -krát menším než jednotlivá pozorování, a to bez ohledu na výchozí rozdělení jednotlivých pozorování
 - o to je často důvodem předpokladu o normálním rozdělení, výsledná hodnota je ovlivněna součtem velkého počtu nahodilých malých vlivů



STATISTICKÁ INDUKCE

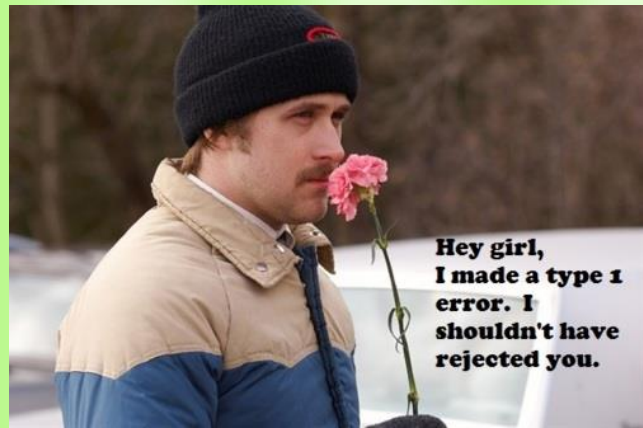
- **95% interval spolehlivosti (konfidenční interval)**
- Rozmezí hodnot, ve kterých leží s 95% pravděpodobností **populační průměr**
- Máme výběrový průměr, vypočtený z X hodnot, ale neznáme směrodatnou odchylku σ , takže si nemůžeme dopočítat z kvantilu normovaného normálního rozdělení, o kolik se hodnota našeho výběrového průměru s 95% pravděpodobností liší od skutečného průměru μ . Ale my si stejně poradíme, směrodatnou odchylku si odhadneme složitým vzorečkem – $S = \sqrt{\frac{1}{n-1} \times \text{suma}(\text{odečtu každé konkrétní hodnoty od průměru})^2}$ na druhou. Tímto jsme si všechno převedli z normálního rozdělení na studentovo, takže použijeme jeho kvantil Studentova rozdělení – pro $z(0,975)$ je

to 2,064 (místo 1,960). Zvýšení kritické hodnoty je náš trest za to, že jsme neznali sigma a nahradili jsme ho jeho odhadem S.

- Ted' už známe všechno, co potřebujeme, a tak konečně vypočítáme 95% interval spolehlivosti. Je to rozmezí, které s 95% pravděpodobností překryje skutečnou hodnotu μ . Intervaly spolehlivosti mohou být i jednostranné – v případě, že nás zajímá jen závislost jedním směrem – např. jestli jsou všichni větší než hodnota 1,5 metru..
- Co určuje velikost intervalu spolehlivosti (jeho délku)? Střední chyba průměru (S.E.X) – přesnost, s jakou výběrový průměr X odhaduje populační průměr μ . Dále pochopitelně požadovaná spolehlivost (1- α) a trošičku i počet stupňů volnosti (n-1)

- **Chyba prvního druhu** – zamítnutí nulové hypotézy, i když platí, tedy falešné potvrzení alternativní hypotézy

Padne nám ze 100 pokusů 21 krát šestka (což se může stát s 10% pravděpodobností) a my řekneme, že je kostka vadná. To ale dost riskujeme, že je kostka v pořádku (platí nulová hypotéza) a ty šestky padly náhodou. Když se chceme této chybě co nejvíce vyhnout, musíme si určit číslo, které náhodu připouští jen s (většinou) 5% pravděpodobností. Největší dovolené pravděpodobnosti chyby 1. druhu se říkáááá...



- **...hladina testu - alfa (hladina pravděpodobnosti, významnosti)**

- o volí se před pokusem, nezávisle na jeho výsledku (ale dost svádí k tomu podvádět..)
- o bývá většinou 5% nebo 1%
- o definuje nám **kritický obor** – vypočítáme si hodnotu (kritickou hodnotu, k_0), kde už začíná být pravděpodobnost, že se jednalo o náhodný úlet hodnot kamsi do dále, menší než alfa. Všechny hodnoty vyšší než kritická nám spadnou do kritického oboru, a to je jako exploze, co zničí (zamítne) nulovou hypotézu, a my můžeme oslavovat ☺
- o Doplnkem kritického oboru je obor přijetí, když tam padnou hodnoty, H_0 nezamítáme



- **Chyba druhého druhu** – nezamítnutí nulové hypotézy, i když neplatí – kostka je pošahaná, ale nám náhodou padne přesně v šestině pokusů šestka. Pravděpodobnost chyby II. druhu se říká beta
- pravděpodobnost této chyby se těžko ovládá – nejlepší způsob je pomocí...
- **...síly testu - 1-beta**
 - o pravděpodobnost zamítnutí neplatné nulové hypotézy
 - o pst, s jakou prokážeme platnou „vědeckou hypotézu“

Nejlepší způsob, jak zvýšit sílu testu, je přidat počet pozorování – je to jasné, čím více se na nějakou věc koukáme, tím větší je šance, že si na ní něčeho zajímavého všimneme.. Se zvýšením počtu pozorování (n) se také pochopitelně zvýší k_0
Síla testu nám přímo řekne, jakou procentuální šanci máme s našimi daty, a při naší hladině významnosti, na to, že nám test něco poví (že správně vyvrátíme nulovou hypotézu)

Tabulka 1. Platnost hypotéz o situaci v základním souboru a možná rozhodnutí na základě testování¹

	Rozhodnutí	
Platí	H ₀	H ₁
H ₀	OK ($P=1-\alpha$)	Chyba prvního druhu ($P=\alpha$)
H ₁	Chyba druhého druhu ($P=\beta$)	OK ($P=1-\beta$) Síla testu

- **Pojmy k populaci**
 - o populace – všichni – prvky celé množiny, základní soubor
 - o výběr – vzorek populace - podmnožina
 - o reprezentativní výběr – podmnožina, pomocí které hodnotíme celý základní soubor, musí být základnímu souboru dostatečně podobná, musí věrně zachycovat poměry platné v celém souboru
 - o parametr – neznámá konstanta, která nějakým způsobem určuje použitý pravděpodobnostní model.. aha..
 - o odhad – statistika použitá k odhadu toho tajemného parametru – lepší definici nemám, ale dá se to vytušit z dalších kapitol..
 - o statistika – je to vhodná funkce naměřených hodnot, která nám zpravidla pomáhá definovat kritický obor – např. výběrový průměr, chí kvadrát,..
 - o populační průměr – μ (při normálním rozdělení) – je to průměr, který získáme měřením náhodného vzorku a bereme ho jako průměr celé populace
- **Dosažená hladina testu – p-hodnota**
 - o nejmenší hladina, na které lze ještě z daných dat zamítnout nulovou hypotézu, pravděpodobnost, s jakou dostaneme výsledek, co bude stejně nebo více odporovat nulové hypotéze, p musí být větší než

hladina významnosti, abychom zamítli H_0 – poslední část je fakt důležitá!

- **možností zkoumání závislosti dvojice znaků**

- kvantitativní – kvantitativní
 - korelace, regrese
 - graficky – rozptylový (bodový) diagram
- kvantitativní – kvalitativní
 - t-test, ANOVA
 - graficky – box-plot (krabicový diagram)
- kvalitativní – kvalitativní
 - chí-kvadrát test, Fischerův exaktní test
 - graficky – kontingenční tabulka

- **výběr testu**

rozdělení	normální	spojité
Jeden výběr	jednovýběrový t-test	jednovýběrový Wilcoxon
výběr dvojic	párový t-test	znaménkový, Wilcoxon
dva nezávislé výběry	dvouvýběrový t-test	Mann-Whitney (=dvouvýběrový Wilcoxon)
k nezávislých výběrů	ANOVA	Kruskal-Wallis

- **neparametrické testy**

- při neparametrickém testu testujeme jinou hypotézu o rozdělení základního souboru než je hypotéza o jeho parametru. Jejich řešení nezávisí na typu rozdělení základního souboru, takže na rozdíl od parametrických testů výsledky nejsou závislé na tom, zda jsme model rozdělení volili správně v souladu se skutečným rozdělením základního souboru.
- Můžeme je pak použít i pro silně nenormální rozdělení, kdy parametrické testy předpokládající normální rozdělení selhávají. Jsou to ty testy, co mají v tabulce nahoře rozdělení spojité.

JEDEN VÝBĚR

- **T-test - jednovýběrový**

- předpoklady – X_1, \dots, X_n jsou nezávislé náhodné veličiny, mají stejné rozdělení (normální) s neznámými parametry μ a σ na druhou (větší než nula)
- srovnáváme μ a μ_0 , což je konstanta, která je daná (víme jí, zajímáme se o ni, porovnáváme s ní data); nulová hypotéza říká, že $\mu = \mu_0$. K ní existují 3 alternativní: oboustranná alternativní hypotéza – $\mu \neq \mu_0$ – a obě jednostranné ($\mu < \mu_0$, $\mu > \mu_0$) – ty ovšem volíme, jen pokud si to žádá náš výzkum nebo zkoumaná situace, ne když se nám k tomu hodí data – při oboustranném testu se počítá

s poloviční alfou, proto nás to svádí podvádět a použít jednostranné testy, kde je alfa celá. Jednostranný musí mít jeden směr, který vede od H_0 k alternativě a žádný druhý směr není.

- kritický obor dostaneme ze srovnání populačního průměru (který předpokládáme, μ_0) a výběrového průměru, kterým se snažíme odhadnout skutečný populační průměr ($T = (X - \mu_0) / S.E.(X)$) – zajímá nás, jestli se liší – velikost rozdílu hodnotíme jeho vydělením $S.E.(X)$ (tj. porovnáváme s přesností výběrového průměru). Vydělením vznikne bezrozměrná statistika T – záleží jen na rozsahu výběru n !!! Chytrá věc..
- úpravou asi milionu vzorečků dojdeme k nezajímavé věci (skripta tvrdí opak), že interval spolehlivosti pro μ je množina všech středních hodnot μ_0 , které bychom nezamítli

Ověření předpokladů: tzv. robustním metodám nedodržení tolik nevadí, ale nemělo by se dít.

T - test má 2 předpoklady: nezávislost jednotlivých měření (např. když testuji výšku, neměla by být ve vzorku jednovaječná dvojčata) a normální rozdělení. To můžeme ověřit graficky pomocí...

- **...pravděpodobnostního diagramu (probability plot, normal q-q plot)**
 - Grafické znázornění variační řady našeho pozorování proti hodnotám, které bychom ve stejném případě měli s normálním rozdělením = porovnání skutečných hodnot s jejich idealizovanými protějšky. Ideál je, aby byly všechny body na přímce se směrnicí sigma.

Test pro ověření normality je **Shapiro-Wilkův test**, dá se to dále udělat např. výběrovou šikmostí a špičatostí.

T-test – párový

Dvojice prvních a druhých měření (U, V), které jsou mezi sebou nezávislé; a rozdíly v obou měřeních ($X = U - V$) u všech hodnot mají normální rozdělení. Z nich testujeme hypotézu o shodě populačních průměrů U a V ($\mu_U = \mu_V$). To je stejné, jako bychom testovali jednovýběrový t -test s $\mu = 0$, který se aplikuje na rozdíly $U - V$.

Znaménkový test

- lze jím prověřovat, zda jsou dva opačné jevy stejně pravděpodobné
- Používá se pro vyhodnocení párových pokusů v případech, kdy studovanou veličinu nemůžeme přesně měřit. V testu nepoužíváme žádné naměřené hodnoty, stačí nám rozhodnutí, zda pokusný zásah „A“ zapůsobil více či méně než pokusný zásah „B“. Pro svou jednoduchost se znaménkový test používá zejména k orientačnímu hodnocení předběžných pokusů, např. v mikrobiologii. Nepoužíváme hodnoty rozdílů, ale pouze jejich znaménka.
- Při platnosti H_0 má jev $X > 0$ pravděpodobnost $\frac{1}{2}$

Párový Wilcoxonův test

Zdokonalení znaménkového testu – řešíme nejen, jaké je znaménko rozdílu, ale také vzdálenost rozdílu od nuly – ne přímo konkrétní číslo, jen pořadí hodnot absolutní hodnota U-V

Yatesova korekce

Využívá se pro malé hodnoty n – např. pár desítek. Čísel se přiblíží o $\frac{1}{2}$ k nule, spolehlivěji se tak dodrží zvolená hladina významnosti.

- řeší problém malých teoretických četností ve čtyřpolní tabulce pro test nezávislosti či homogenity (viz níže)

2 VÝBĚRY

- Srovnání 2 populací pomocí výběru z každé z nich, výběry nezávislé jak uvnitř populace, tak mezi populacemi
- H_0 – oba populační výběry jsou stejné – založená na odhadu obou populačních průměrů
- **Dvouvýběrový t-test**
 - Předpoklady: všechny X a všechny Y nezávislé, s normálním rozdělením, shodné populační rozptyly sigma na druhou (to odhadneme pomocí váženého průměru výběrových rozptylů příšerným vzorečkem)
 - H_0 : $\mu_x = \mu_y$, rozhoduje se statistikou $T = \frac{\bar{x} - \bar{y}}{S.E.(\bar{x} - \bar{y})}$; stupně volnosti $n_x + n_y - 2$
 -
- **Dvouvýběrový Wilcoxonův test (Mannův-Whitneyův)**

pro 2 nezávislé výběry ze spojitého rozdělení, u kterých nepředpokládáme normální rozdělení => porovnáváme 2 populace v kvantitativním znaku

H_0 : rozdělení veličiny je v obou populacích stejné, populace se ve sledovaném znaku neliší = mají stejné populační mediány (rozdělení jsou shodná)

Naměřené hodnoty nahradíme jejich pořadím, bez ohledu na to, jestli jsou v 1. Nebo 2. Populaci, za platnosti H_0 budou pořadí podobná

Sečteme všechny hodnoty pořadí u X a všechny u Y, jejich střední hodnoty budou dávat dohromady celý součet pořadí a to v poměru vůči velikostem souborů x a y – když platí H_0

ANALÝZA ROZPTYLU

- **ANOVA**
 - analýza rozptylu **jednoduchého třídění**
 - Jako dvouvýběrový t-test, ale hodnotíme shodu středních hodnot několika nezávislých výběrů z normálního rozdělení, snaha prokázat závislost mezi znakem Y v kvantitativním měřítku a znakem A v nominálním – hodnoty A se dají zapsat čísly, ale mají slovní vyjádření a označují se jako faktor –

např. závislost mezi hodnotami cholesterolů u různých skupin – faktorem je skupina

- Máme k nezávislých výběrů z normálního rozdělení, k nám tedy říká, kolikrát máme soubor dat Y
- Předpoklad **homoskedasticity** – všechny skupiny mají stejný rozptyl, $\sigma_1 = \dots = \sigma_k$; opak je heteroskedasticita
- Pro $k=2$ je Anova stejná jako t-test, pro více k testujeme, jestli $\mu_1 = \mu_2 = \dots = \mu_k$, testujeme nestejnost většího počtu středních hodnot
- budeme to měřit jejich variabilitou (SS_A), tedy součet druhých mocnin jejich odchylek od celkového průměru (ten zjistíme jako vážený průměr průměrů jednotlivých výběrů), stupňů volnosti bude $k-1$
- problém je, že pokud H_0 neplatí, tzn. střední hodnoty jsou rozdílné, zvětší to automaticky rozptyl. My ho pak nemůžeme správně odhadnout a dále s ním pracovat. Potřebujeme tedy jiný výpočet (odhad) rozptylu, kterému bude jedno, jestli jsou střední hodnoty ve shodě nebo ne. Takovým odhadem je reziduální rozptyl (je to vážený průměr výběrových rozptylů všech k výběrů - blablabla).
- Toto (to dole) je tabulka analýzy rozptylu, ve které jsou prý shrnuty všechny výpočty. Co z ní můžeme dobře odvodit je jev, nazývaný rozklad součtu čtverců. Tedy, že celkovou variabilitu SS_T rozdělíme na 2 složky: SS_A , což je variabilita daná různými středními hodnotami v různých výběrech; a SS_E , což je reziduální variabilita – tedy variabilita uvnitř jednotlivých výběrů, kolísání kolem jejich jednotlivých průměrů.
- Z toho všeho vypočítáme příšernými vzorečky hodnotu statistiky F ; čím je větší v porovnání s průměrným čtvercem (reziduálním rozptylem), tím máme větší šanci, že zamítneme H_0 . Je vzoreček na to, jak zjistit, kdy přesně už jsme překročili kritickou hodnotu a můžeme H_0 zamítnout, ale lepší způsob je nechat si programem vypočítat p -hodnotu a pak se na ní podívat.

Variabilita	S	f	S/f	F	P
Výběry	S_A	$f_A = k-1$	S_A/f_A	F_A	p_A
Reziduální	S_E	$f_E = n-k$	S_E/f_E		
Celková	S_T	$f_T = n-1$			

S – součty čtverců

F – F -statistika

f – počet stupňů volnosti

p – p -hodnota

S/f – průměrné čtverce

- Máme asi milion metod pojmenovaných po milionech pánů, které nám pomůžou ověřit homoskedasticitu. Např. Bartlettův test, Levenův test

- Potom máme moc chytrou metodu pana Tukeya – Tukeyova metoda HSD, která nám poví, jaká byla příčina zamítnuté H_0 (tedy které populační průměry jsou nestejně)

Kruskalův-Wallisův test

Tento test nám řeší stejný problém, jako řešil Wilcoxonův test v případě dvouvýběrového t-testu. Řeší tedy problém s nesplněním předpokladu o normálním rozdělení. Je to vlastně zobecnění Wilcoxonova testu.

Máme k nezávislých výběrů ze spojitého rozdělení (které mají všechny veličiny stejné, mohou se lišit jen posunutím). H_0 říká, že jsou rozdělení v různých výběrech identická; tedy, že posunutí jsou nulová (tedy že jsou stejné populační průměry a mediány). Opět řešíme, jestli jsou průměrná pořadí ve výběrech podobná. Nepodobnosti průměrných pořadí měří statistika Q a H_0 zamítáme podle chí kvadrát rozdělení.

Biostatistickí rozhodně nezahálí, takže vymysleli i Dvojně třídění, kde se řeší nějaké interakce; když je to dvojně třídění, ale žádné interakce se neřeší, tak na to vymysleli Náhodné bloky. Když se to obojí chce počítat jako znaménkový test, jen na to Friedmanův test.

KORELACE A REGRESE

Korelace popisuje sílu vzájemné závislosti dvou veličin, regrese popisuje, jak vysvětlující veličina ovlivní veličinu vysvětlovanou

Výběrový (Pearsonův) korelační koeficient

- Populační korelační koeficient ρ jsme řešili někde nahoře už dávno. Ted' potřebujeme dostat jeho výběrový protějšek (jeho odhad, když počítáme s náhodným výběrem), který se značí r
- $r = \text{výběrová kovariance (S}_{xy}) / \text{rozptyl } X \times \text{rozptyl } Y$ – nezapomenout, že z toho plyne, že je r bezrozměrné!
- Korelační koeficient nám ukazuje na lineární závislost dat, má hodnoty -1 až 1 a tyto hodnoty má v extrémním případě, kdy jsou všechny body na přímce
- Pokud je $\rho = 0$ znamená to, že jsou X a Y nezávislé. Ovšem protože je r jen odhad ρ , vždy to alespoň trochu nenulové vyjde, i když jsou veličiny nezávislé. Musíme proto otestovat H_0 : veličiny X a Y (s normálním rozdělením) jsou nezávislé, tedy r je blízké nule. To prokážeme T statistikou.

- Spearmanův korelační koeficient

Hodnotí závislost spojitých znaků, když není dodrženo normální rozdělení.

Místo hodnot znaků tedy hodnotíme pořadí těchto znaků. Místo

Pearsonova r počítáme Spearmanovo r_s .

Výhody: citlivost na jakoukoliv závislost, nejen lineární; menší citlivost k výskytu odlehlých hodnot.

- Regrese

Slouží k vysvětlení variability (kolísání) odezvy (závislé, vysvětlované proměnné, Y) na regresor (vysvětlující, nezávislá proměnná, x). Jednosměrná závislost Y na x . Velikost písmen x a Y je schválně jiná, aby nám došlo, že jsou jako opravdu zásadně odlišné.

Pomůže nám: prokázat závislost odezvy (její střední hodnoty) na nezávislé proměnné, předpovídat (hmm.. to zní dobře, ale spíš odhadovat) střední hodnotu vysvětlované proměnné podle hodnot nezávislé proměnné.

x známe přesně, není náhodná, Y náhodná je a skládá se ze dvou složek: deterministická složka je ta, co jí ji vnutí její závislost na x ; a náhodná složka, která je daná vlastní vnitřní variabilitou Y . Např. výška synů je závislá na výšce otců (deterministická), ale i stejně vysocí otcové mají různě vysoké syny (stochastická složka).

Regresní přímka

Pro jednu spojitou nezávisle proměnnou.

Střední hodnoty nezávislé náhodné veličiny Y leží na přímce $y = \beta_0 + \beta_1 x$ (absolutní člen = posunutí) + β_1 krát x (β_1 jedna je směrnice). Je to normální směnicový tvar přímky se všemi pravidly jako nás učili na střední. Je dobré si všimnout, že na rozdíl od náhodného výběru tu nepředpokládáme, že budou všechny střední hodnoty stejné, chceme po nich jenom to, aby leželi na té své přímce. Pochopitelně stejně tak jako neznáme rozptyl, neznáme ani hodnoty regresních koeficientů β_0 a β_1 . β_1 je extrémně důležitá hlavně proto, že nám ukazuje citlivost, s jakou reaguje střední hodnota Y na střední hodnotu x .

$Y = \beta_0 + \beta_1 x + E$

$\beta_0 + \beta_1 x$ je vlastní závislost, E je modelem nevysvětlitelná variabilita (něco jako stochastická složka)

Neznámé regresní koeficienty odhadujeme **metodou nejmenších**

čtverců. Je to fakt chytrá metoda, kdy máme body (ty značí naše měření, mají tedy souřadnice x a Y) a ty body posouváme po ose Y na regresní přímku. Délka posunu bodu se rovná délce jedné stany imaginárního čtverce (říkáme jí reziduum). My naší regresní přímku vedeme tak, aby byl součet ploch všech takto vzniklých čtverců (reziduální součet čtverců) co nejmenší.

Z reziduálního součtu čtverců můžeme vypočítat reziduální rozptyl, což je odhad rozptylu sigma na druhou.

Je sice hezké, že jsme si odhadli, jaká je závislost Y na x , ale ještě jsme si tu závislost ani neprokázali. To si prokážeme T statistikou, kdy $H_0: \beta_1 = 0$ (ze střední odvodíme, že by pak přímka byla rovnoběžná s x a konstantní).

Když zvolíme $\beta_1 = 0$, úplně tím vlastně umažeme závislost Y na x . I pro tuto přímku si můžeme spočítat součet čtverců – tentokrát Celkový součet čtverců – a ten nám poví, jaká je variabilita hodnot závisle proměnné.

Když od Celkového součtu čtverců odečteme Reziduální součet,

dostaneme Regresní součet čtverců. Ten nám okomentuje variabilitu Y, kterou vysvětlíme vyšetřovanou závislostí Y na x. Všechny tyto dílčí úvahy nám hezky shrne...

... koeficient determinace

R^2 = regresní součet čtverců/celkový (nebo 1 - reziduální/celkový)

Koeficient determinace nám tedy říká, jakou část variability závisle proměnné vysvětlíme díky její závislosti na x.

KONTINGENČNÍ TABULKY

Hodnotí četnosti nominálních hodnot

- Chí-kvadrát test dobré shody

- při testování kvalitativní – kvalitativní, např. krevní skupiny a jejich četnosti
- Aby byl výběr reprezentativní, musí být empirické četnosti shodné s % podílem v populaci

Očekávané (teoretické) četnosti –

spočítáme je např. pomocí různých pravděpodobností, kombinatorikou a podobnými hrůzami

Empirické četnosti – experimentálně

zjištěné, spočítáme je např. na prstech

H_0 určí očekávané četnosti, říká, že očekávané jsou shodné s empirickými četnostmi

- $\chi^2 = \text{suma} (\text{exper.četnost} - \text{teor.četnost})^2 / \text{teor.četnost}$

Pokud vyjde chí kvadrát statistika velká a překročí kritickou hodnotu pro danou hladinu pravděpodobnosti a stupně volnosti, znamená to, že H_0 zamítáme a výběr považujeme za nerepresentativní.. Krásné, elegantní, jednoduché



- Kontingenční tabulka

Testuje nezávislost dvou znaků v nominálním měřítku.

Zapíšeme četnosti kombinací vždy dvou znaků (sdružené četnosti), jejich součtem doplníme marginální četnosti v řádcích a sloupcích.

V sloupcích jsou pak četnosti druhého znaku bez ohledu na první a

v řádcích naopak. V rohu je celkový počet všech dvojic hodnot – n.

Hodnoty prvního znaku – veličina X, hodnoty 1,...,r, druhého znaku – Y, hodnoty 1,...,s.

- Sdružené četnosti se značí jako n_{RS} – tam se napíší konkrétní čísla, o četnost na jakém místě zrovna půjde (např. B je n_{12}) a marginální četnosti se označí tak, že to, co sčítáme je číslo a to druhé + (např. b+d dole je $+2$). Speciální případ pro $r=c=2$ je čtyřpolní tabulka:

a	B	a+b
c	D	c+d
a+c	b+d	N

Když testujeme hypotézu nezávislosti, normálně chí - kvadrát testem srovnáme očekávané a empirické četnosti. Výpočet jak očekávaných četností, tak testová statistika, jsou formálně shodné s chí - kvadrát testem nezávislosti, ale v tomto případě se mu říká **test homogeneity**. Zajímá nás, zda jsou pravděpodobnosti jednotlivých možných výsledků dílčího pokusu ve všech populacích stejné.

Je to obdoba ANOVY, v případě testování homogeneity dvou multinomických rozdělení dvouvýběrového t-testu (přestože se to počítá jako chí-kvadrát test).

McNemarův test

- obdoba testu párového (ale teď pro nominální hodnoty)

Porovnává výsledky dvojích měření na stejných objektech. Měření se musí týkat nominálního znaku. Vše je zapsáno v kontingenční tabulce (nutně čtvercové) a H_0 říká, že pravděpodobnosti multinomického rozdělení popisujícího četnosti výsledků v prvním a ve druhém pokusu jsou stejné.

