

Seriál – Pravděpodobnost a matematická statistika 1

Teorie pravděpodobnosti je nesmírně krásnou vědou. Zatímco většina odvětví matematiky se věnuje pojmům a veličinám, jež jsou „přesně“ dány, umožňuje teorie pravděpodobnosti hovořit i o nejistotě a neurčitosti. A protože v reálném světě je více věcí nejistých nežli jistých, má široké uplatnění v praxi.

Seriál bude mít tři části. V první části zavedeme základní pojmy a probereme to, co někteří z vás ještě nestihli probrat ve škole. Ve druhé si uděláme výlet do světa náhodných veličin a podmíněných nezávislostí. Třetí část se bude věnovat praxi matematické statistiky: normálnímu rozdělení, zákonu velkých čísel a centrální limitní větě.

Anžto jsou znalosti středoškolských studentů v této oblasti značně rozdílné, předem se omlouvám jak těm, kteří se při čtení prvního dílu seriálu budou nudit, tak těm, kteří se v něm zcela ztratí. Pokud nebudete rozumět nějaké definici či větě, podívejte se na osvětlující příklad za ní. Pokud vám schází středoškolské znalosti z kombinatoriky a pojmy jako permutace, kombinace či variace vám nic neříkají, doporučuji si je doplnit. Ať už z nějaké lepší středoškolské učebnice nebo z internetu. Můžete zkusit kupříkladu

<http://www.czech-ware.net/mathes/zobuceb.aspx?zob=kombinatorika>

nebo stručný, ale všezahrnující textík

<http://mks.mff.cuni.cz/knihovna/kombinatorika.pdf>.

Klasická pravděpodobnost – verze pro mateřské školy

Mějme množinu všech možných výsledků nějakého náhodného pokusu. Tuto množinu budeme nazývat *množina elementárních jevů* a značit Ω . O jejích podmnožinách $A \subset \Omega$ budeme hovořit jako o jevech, o jejích prvcích $\omega \in \Omega$ jako o elementárních jevech.

Příklad. Uvažujme výsledek hodu šestistěnnou kostkou¹. Potom množina elementárních jevů jest

$$\Omega = \{1, 2, 3, 4, 5, 6\}.$$

Jevem je kupříkladu to, že padne sudé číslo, $A = \{2, 4, 6\}$. Elementární jevy jsou $\{1\}$, $\{2\}$, $\{3\}$, $\{4\}$, $\{5\}$, $\{6\}$.

Poměrně jednoduchý případ nastává tehdy, je-li zaručeno, že Ω je konečná a že všechny elementární jevy jsou „stejně pravděpodobné“. Potom *pravděpodobností jevu* A rozumíme číslo

$$P(A) = \frac{|A|}{|\Omega|},$$

kde $|X|$ značí počet prvků množiny X .

Jev A nazveme *jistý*, pokud $P(A) = 1$, a *nemožný*, pakliže $P(A) = 0$. Jevem doplňkovým k jevu A míníme jev $A^c = \Omega \setminus A$, zjevně $P(A^c) = 1 - P(A)$.

¹Pokud nebude řečeno jinak, všechny další v textu zmíněné kostky jsou šestistěnné.

Příklad. Tedy v minulém příkladu pravděpodobnost jevu A , že padne sudé číslo, je rovna

$$P(A) = \frac{|A|}{|\Omega|} = \frac{3}{6} = 0,5.$$

Příklad. Na pouti se objevil stánek, kde si návštěvníci mohli vsadit na hod třemi kostkami². Stánkaři zaplatili 10 Kč a pokud hodili v součtu alespoň 17 byla jim vyplacena stokoruna. Jinak vklad propadal. Jaká je šance hráče na výhru?

Uvažujme jako náhodný pokus hod třemi kostkami a zajímejme se, zda je jejich součet alespoň 17. Množinou elementárních jevů Ω tvoří $6 \cdot 6 \cdot 6 = 216$ trojic možných hodů. Z nich čtyři, konkrétně $(6, 6, 6)$, $(5, 6, 6)$, $(6, 5, 6)$, $(6, 6, 5)$, mají součet alespoň 17. Tudíž

$$P(\text{součet je alespoň } 17) = \frac{4}{216} \doteq 0,019.$$

Jako nepříliš těžké cvičení se můžete pokusit přijít na to, jaká byla „průměrná“ výhra hráče, a kolik tedy průměrně na každé sázce vydělával stánkař.

Příklad. Mějme 6 párů (tzn. 12 kusů) střeviců. Náhodně z nich vyberme 5 střeviců. Jaká je pravděpodobnost, že mezi nimi bude alespoň jeden pár?

Množinou elementárních jevů Ω jsou všechny (neuspořádané) pětičky vybrané z množiny všech bot $\{a_1, a_2, b_1, b_2, c_1, c_2, d_1, d_2, e_1, e_2, f_1, f_2\}$, kde stejné písmeno značí příslušnost k témuž páru. Takovýchto pětic je $\binom{12}{5} = 792$.

Nás zajímá pravděpodobnost jevu A , že ve výběru bude alespoň jeden pár. Tu je ale značně nesnadné spočítat. Co je však mnohem snazší, je spočítat pravděpodobnost doplňkového jevu A^c , že tam žádný pár nebude,

$$P(A^c) = \frac{2^5 \cdot \binom{6}{5}}{\binom{12}{5}} = \frac{192}{792} \doteq 0,24,$$

neboť vybraných pět střeviců musí patřit do pěti různých párů (a pět párů z šesti lze vybrat $\binom{6}{5}$ způsoby) a u každého páru máme na výběr ze dvou střeviců.

Odtud už snadno dopočteme pravděpodobnost samotného jevu

$$P(A) = 1 - P(A^c) = \frac{600}{792} \doteq 0,76.$$

Cvičení. „Když se náhodně vyberou dvě mé děti, je stejná pravděpodobnost, že mají stejné pohlaví, jako pravděpodobnost, že mají různé pohlaví,“ pravil sultán.

„Jaká je pravděpodobnost, že to budou dvě dívky?“ otázal se kalif.

„Stejná jako že náhodně vybrané dítě bude chlapec,“ odvětil sultán.

Pokud jste dobře pochopili, co to je pravděpodobnost, a umíte řešit kvadratické rovnice, nemělo by pro vás být těžké určit, kolik má sultán dětí.

Podmínění náhodným jevem, nezávislost jevů

²Příhoda není tak úplně imaginární. Podobný stánek jsem skutečně viděl na Mohelnickém dostavníku.

Nechť máme dva jevy A a B , $P(B) > 0$. *Podmíněnou pravděpodobnost jevu A za podmínky, že nastal jev B* , definujeme jako

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

Dva jevy A a B nazveme *nezávislé*, jestliže

$$P(A \cap B) = P(A) \cdot P(B).$$

Pokud tyto jevy nejsou nemožné, snadno z nezávislosti odvodíme

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = P(A),$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = P(B).$$

Vidíme tedy, že u nezávislých jevů to, že nastal jev A , neovlivní pravděpodobnost, že nastane jev B , a naopak.

Příklad. Uvažujme hod kostkou: Nechť A je jev, že padne sudé číslo. Nechť B je jev, že padne jednička nebo dvojka. Potom si snadno můžete ověřit, že jevy A a B jsou nezávislé.

Rozklad množiny elementárních jevů

Jevy B_1, B_2, \dots, B_m nazveme rozkladem množiny elementárních jevů, jestliže průnik libovolných dvou z nich je prázdná množina a sjednocení všech je Ω .

Jinak řečeno, každý elementární jev $\omega \in \Omega$ náleží do právě jednoho z těchto jevů.

Příklad. U hodu kostkou můžeme uvažovat rozklad na jevy, že padne sudé/liché číslo. Nebo rozklad na jevy $B_1 = \{1\}$, $B_2 = \{3, 4, 6\}$, $B_3 = \{2, 5\}$.

Věta. *Nechť B_1, B_2, \dots, B_m je rozklad množiny elementárních jevů a žádný z těchto jevů není nemožný. Potom pro libovolný jev A platí*

$$P(A) = \sum_{i=1}^m P(B_i) \cdot P(A|B_i).$$

Důkaz. Důkaz je snadný:

$$\sum_{i=1}^m P(B_i) \cdot P(A|B_i) = \sum_{i=1}^m P(A \cap B_i) = P(A).$$

Příklad. K lékaři chodí nemocní i zdraví pacienti. Ví se, že 60% z pacientů, kteří přijdou do ordinace, je nemocných. Pan doktor není tvor neomylný, a tak se u 20% pacientů (ať už zdravých či nemocných) spletě. Kolik procent pacientů označí pan doktor jako nemocné?

Tupě aplikujme předchozí větu s rozkladem pacient je/není nemocný:

$$P(\text{pacient označen za nemocného}) =$$

$$P(\text{pacient označen za nemocného} | \text{pacient je nemocný}) \cdot P(\text{pacient je nemocný}) +$$

$$\begin{aligned} P(\text{pacient označen za nemocného} | \text{pacient je zdravý}) \cdot P(\text{pacient je zdravý}) &= \\ &= 0,8 \cdot 0,6 + 0,2 \cdot 0,4 = 0,56 \end{aligned}$$

Povšimněme si, že nás v tomto případě vůbec nemuselo zajímat, jak vypadá množina Ω .

Následující věta, jež nese jméno po ctihodném otci Bayesovi, slouží k tomu, abychom mohli „prohodit“ jevy v podmiňování:

Věta (Bayesova). *Nechť B_1, B_2, \dots, B_m je rozklad množiny elementárních jevů a žádný z těchto jevů není nemožný. Potom pro libovolný jev A , $P(A) > 0$ platí*

$$P(B_k | A) = \frac{P(B_k) \cdot P(A | B_k)}{\sum_{i=1}^m P(B_i) \cdot P(A | B_i)}, \quad k = 1, \dots, m.$$

Důkaz. Důkaz je opět velmi snadný:

$$P(B_k | A) = \frac{P(A \cap B_k)}{P(A)} = \frac{P(B_k) \cdot P(A | B_k)}{\sum_{i=1}^m P(B_i) \cdot P(A | B_i)}.$$

Příklad. Předchozí věta je značně oblíbená v lékařské diagnostice. Uvažujme testování, zda pacient má AIDS. Máme k dispozici vcelku spolehlivý test, který dá v 95 procentech případů správnou odpověď a v 5 procentech mylnou. Víme, že ve sledované populaci má AIDS 0,5% lidí. Test dopadl pozitivně³. Jaká je pravděpodobnost, že pacient nemoc skutečně má?

Označme N jev, že pacient je nemocný, a Z , že je zdravý. Označme \oplus jev, že test dopadl pozitivně, a tupě dosaďme do předchozí věty:

$$P(N | \oplus) = \frac{P(N) \cdot P(\oplus | N)}{P(N) \cdot P(\oplus | N) + P(Z) \cdot P(\oplus | Z)} = \frac{0,005 \cdot 0,95}{0,005 \cdot 0,95 + 0,995 \cdot 0,05} \doteq 0,087.$$

Odpověď tedy zní: velice nízká, ani ne desetina. Poněkud překvapivý výsledek u tak spolehlivého testu, nezdá se vám? Zkuste přijít na to, proč tomu tak je.

Cvičení. Na procvičení doporučuji následující příklad: Američan, Rus a Čech se hádají, kdo vydrží déle pod vodou. Američan tvrdí, že dvě minuty (pravděpodobnost přežití 0,8). Rus tvrdí, že pět minut (pravděpodobnost přežití 0,5). Čech tvrdí, že deset minut (pravděpodobnost přežití 0,1). Všichni se potopili a jeden z nich to nepřežil. Jaká je pravděpodobnost, že to byl Čech?

Námět k zamyšlení

Jaká je hlavní nevýhoda námi definované pravděpodobnosti? Tušíte správně – v tom, že všechny elementární jevy nemusí být nutně stejně pravděpodobné. Představte si třeba „falešnou kostku“, na které padá šestka o něco častěji než jednička.

Pokud se tedy chceme bavit i o „falešných kostkách“, musíme opustit mateřskou školku. Ale o tom zase až příště ...

³To znamená, že svědčí pro to, že pacient nemoc má.

Seriál – Pravděpodobnost a matematická statistika 2

Motto: Bohužel vaše planeta je jednou z těch, které byly určeny k demolici. Tento proces započne za necelé dvě vaše pozemské minuty. **Nepodléhejte panice!** Děkuji vám.
(z poselství Vogona Jetze, Stopařův průvodce po galaxii, Douglas Adams)

V druhé část seriálu rozšíříme naše znalosti z části první. Místo náhodných jevů budeme pracovat s náhodnými veličinami a krom obvyčejné nezávislosti se budeme zabývat i nezávislostí podmíněnou.

Zkrátka a dobře, všechno bude mnohem zajímavější a dobrodružnější nežli minule.

Podmíněná nezávislost

Minule jsme zavedli pojem „nezávislost jevů“. Rozšíříme jej nyní na podmíněnou nezávislost, tedy nezávislost dvou jevů za podmínky, že nastal jev třetí:

Definice. Jevy A a B nazveme *nezávislé za podmínky, že nastal jev C* ($P(C) > 0$), pakliže platí

$$P(A \cap B|C) = P(A|C)P(B|C).$$

Příklad. Nechť A je jev, že se v náhodně vybraném okrese vyskytuje značně nadprůměrné množství čápů, a B jev, že se tu rodí značně nadprůměrné množství dětí, přičemž význam spojení „značně nadprůměrné“ ponechávám čtenářově fantazii.

Tabulka pravděpodobností A a B by mohla vypadat například takto⁴:

$$\begin{array}{ll} P(A \cap B) = \frac{7}{64} & P(A \cap B^c) = \frac{13}{64} \\ P(A^c \cap B) = \frac{13}{64} & P(A^c \cap B^c) = \frac{31}{64} \end{array}$$

Odtud snadno spočítáme $\frac{7}{64} = P(A \cap B) \neq P(A) \cdot P(B) = \frac{6 \cdot 25}{64}$. Jevy A a B tedy nejsou nezávislé, ke všemu si můžete spočítat, že pravděpodobnost výskytu enormního množství novorozenců je vyšší, pokud nastal výskyt enormního množství čápů, než když tomu tak nebylo.

Teď se možná smějete a myslíte si, že takto by to ve skutečnosti být nemohlo. A to se mýlíte! Skutečně existuje výzkum, který potvrzuje statisticky významnou souvislost mezi počtem čápů a novorozenců. Znamená to tedy, že děti nosí čáp?

Nikoli, znamená to jen, že příčinnou souvislost dvou jevů nelze zkoumat bez ohledu na jevy ostatní. Označme si náš hypotetický vysvětlující jev třeba C . Může odrážet vliv životního prostředí, zeměpisné šířky, ... Jeho pravděpodobnost nechť je $P(C) = \frac{1}{4}$ a pravděpodobnosti A a B podmíněno C jsou

$$\begin{array}{llll} P(A \cap B|C) = \frac{1}{4} & P(A \cap B^c|C) = \frac{1}{4} & P(A \cap B|C^c) = \frac{1}{16} & P(A \cap B^c|C^c) = \frac{3}{16} \\ P(A^c \cap B|C) = \frac{1}{4} & P(A^c \cap B^c|C) = \frac{1}{4} & P(A^c \cap B|C^c) = \frac{3}{16} & P(A^c \cap B^c|C^c) = \frac{9}{16} \end{array}$$

Jako snadné cvičení si můžete ověřit, že A a B jsou podmíněně nezávislé za podmínky C , a obě tabulky jsou navzájem konzistentní. A jako první statistikovo příkázání si запиšte: „Nebudu z popření nezávislosti dvou jevů vyvozovat přímou příčinnou souvislost mezi nimi!“. Kéž by toto příkázání ctili všichni praktičtí výzkumníci.

⁴Připomeňme, že A^c značí doplňkový jev k A , tedy $\Omega \setminus A$.

Klasická pravděpodobnost – rozšíření na různě pravděpodobné elementární jevy

Mějme opět množinu elementárních jevů $\Omega = \{\omega_1, \dots, \omega_n\}$. Každému elementárnímu jevu ω_i přiřadíme pravděpodobnost $p_i \in (0, 1)$, že tento jev nastane, přičemž součet těchto pravděpodobností je nutně jedna:

$$\sum_{i=1}^n p_i = p_1 + p_2 + \dots + p_n = 1.$$

Potom *pravděpodobnost jevu* A je rovna součtu všech pravděpodobností elementárních jevů náležících do A :

$$P(A) = \sum_{i: \omega_i \in A} p_i.$$

Příklad. Uvažujme falešnou kostku, na níž padá šestka s pravděpodobností $\frac{1}{2}$, zatímco ostatní čísla s pravděpodobností $\frac{1}{10}$. Množina elementárních jevů $\Omega = \{1, 2, 3, 4, 5, 6\}$. Jev, že padne sudé číslo, $A = \{2, 4, 6\}$:

$$P(A) = \frac{1}{10} + \frac{1}{10} + \frac{1}{2} = 0,7.$$

Poznámka. Povšimněte si toho, že pokud jsou všechny elementární jevy stejně pravděpodobné $p_1 = p_2 = \dots = p_n = \frac{1}{|\Omega|}$, dostáváme teorii pravděpodobnosti pro mateřskou školku.

Dále si povšimněte toho, že jsme nijak nevyužili faktu, že Ω je konečná, a pokud tedy umíme pracovat s nekonečnými součty, můžeme (a také budeme) mít Ω nekonečnou spočetnou⁵.

Cvičení. Ověřte, že při takto modifikované definici pravděpodobnosti jevu mají smysl všechny definice z minula a platí všechny věty.

Výše uvedenou definici typicky využijeme, když se snažíme popsat nějakou náhodnou věc z praxe, jež je číselné povahy – například počet narozených dětí, počet uzdravených pacientů, ... Zde je použití stejně pravděpodobných elementárních jevů přinejmenším krkolomné.

Nechť tedy Ω je spočetnou podmnožinou reálných čísel⁶ a P pravděpodobnost na Ω . Potom budeme mluvit o náhodné veličině X s rozdělením P a $P(X = k)$ značí pravděpodobnost elementárního jevu k , $P(k)$; resp. $P(X \in A)$ značí pravděpodobnost jevu A , $P(A)$.

Náhodné veličiny budeme obvykle značit velkými písmeny z konce abecedy.

Příklad. U hodu šestistěnnou kostkou můžeme označit číslo, jež padne, jako náhodnou veličinu X . Poté můžeme hovořit o pravděpodobnosti, že X je tři, $P(X = 3) = \frac{1}{6}$; resp. že X je sudé, $P(X \in \{2, 4, 6\}) = \frac{1}{2}$.

Definice. *Střední hodnotou* náhodné veličiny X míníme

$$\sum_{x \in \Omega} x \cdot P(X = x)$$

a značíme EX . Jedná se tedy o vážený průměr možných hodnot X , kde váhy tvoří pravděpodobnosti těchto hodnot.

⁵Pokud této poznámce ani zbla nerozumíte, poněvadž netušíte, co znamená slovíčko „spočetný“, tak si s ní nelamte hlavu. Značně volně řečeno, spočetná je taková množina, která má „stejně nebo méně“ prvků nežli množina celých čísel.

⁶tj. například celá čísla, přirozená čísla, konečně mnoho čísel, ...

V angličtině se místo o střední hodnotě mluví o očekávané hodnotě. A právě takový je její význam – značí, kolik asi tak můžeme očekávat, že bude průměr z mnoha realizací X .

Příklad. Nechť X je náhodná veličina udávající hod spravedlivou kostkou. Potom

$$E X = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6} = 3,5.$$

Nechť Y je náhodná veličina udávající hod falešnou kostkou z minulého příkladu. Potom

$$E Y = 1 \cdot \frac{1}{10} + 2 \cdot \frac{1}{10} + 3 \cdot \frac{1}{10} + 4 \cdot \frac{1}{10} + 5 \cdot \frac{1}{10} + 6 \cdot \frac{1}{2} = 4,5.$$

Definice. *Rozptylem* náhodné veličiny X míníme

$$E X^2 - (E X)^2 = E(X - E X)^2 = \sum_{x \in \Omega} (x - E X)^2 \cdot P(X = x)$$

a značíme $\text{Var } X$.

Rozptyl je mírou toho, jak se X drží kolem své střední hodnoty, resp. jaká je v něm ukryta neurčitost. Čím větší rozptyl, tím větší neurčitost. Jako snadné cvičení můžete ověřit, že $\text{Var } X = 0$ právě tehdy, když $X = E X$ s pravděpodobností jedna, a že rozptyl je vždy nezáporný.

Poznámka. Pozorný čtenář si zajisté povšiml, že předchozí definice není zcela korektní, neboť není zcela jasné, co vlastně míníme výrazem $E X^2$, resp. $E(X - E X)^2$, neboť umíme zatím spočítat střední hodnotu náhodné veličiny X , ale nikoli funkce této veličiny $f(X)$.

Pozornému čtenáři je nutno zcela dát zcela za pravdu a přiznat, že autor zde zneužívá toho, že intuitivně je jasné, co se tím „myslí“. Poctivě by bylo potřeba dokázat, že pokud X je náhodná veličina, potom $X' = f(X)$ je opět náhodná veličina (na nějakém jiné množině Ω' a s jiným rozdělením P'), a tudíž můžeme spočítat její střední hodnotu $E' X'$.

Pro praktický výpočet je možno využít toho, že

$$E f(X) = \sum_{x \in \Omega} f(x) \cdot P(X = x),$$

$$\text{Var } f(X) = E f(X)^2 - (E f(X))^2 = \sum_{x \in \Omega} (f(x) - E f(X))^2 \cdot P(X = x).$$

Důkaz je ponechán pozornému čtenáři.

Příklad. Kdyby nás tedy zajímalo, jaká se střední hodnota třetí mocniny čísla, jež nám padne na spravedlivé kostce (označme jej X), zvětšeného o jedničku, dostali bychom

$$E(X + 1)^3 = 8 \cdot \frac{1}{6} + 27 \cdot \frac{1}{6} + 64 \cdot \frac{1}{6} + 125 \cdot \frac{1}{6} + 216 \cdot \frac{1}{6} + 343 \cdot \frac{1}{6} = 130,5.$$

Příklad. Nechť X je náhodná veličina udávající hod spravedlivou kostkou, potom

$$E X^2 = 1^2 \cdot \frac{1}{6} + 2^2 \cdot \frac{1}{6} + 3^2 \cdot \frac{1}{6} + 4^2 \cdot \frac{1}{6} + 5^2 \cdot \frac{1}{6} + 6^2 \cdot \frac{1}{6} \doteq 15,2,$$

$$\text{Var } X = E X^2 - (E X)^2 \doteq 2,9.$$

Věta. Pro náhodnou veličinu X a libovolné nenáhodné konstanty a a b platí:

$$\begin{aligned} E(aX + b) &= a(E X) + b, \\ \text{Var}(aX + b) &= a^2 \text{Var } X. \end{aligned}$$

Důkaz.

$$E(aX + b) = \sum_{x \in \Omega} (ax + b) \cdot P(X = x) = b \cdot \sum_{x \in \Omega} P(X = x) + a \cdot \sum_{x \in \Omega} x \cdot P(X = x) = b + a \cdot E X.$$

Obdobně

$$\begin{aligned} \text{Var}(aX + b) &= \sum_{x \in \Omega} ((ax + b) - (a E X + b))^2 \cdot P(X = x) = \\ &= a^2 \sum_{x \in \Omega} ((x - E X)^2 \cdot P(X = x)) = a^2 \text{Var } X. \end{aligned}$$

Poznámka. Pozorného čtenáře zajisté napadlo, že sice umíme už pracovat s jednou náhodnou veličinou, ale nikoli se dvěma. Nedokážeme tedy říci, jaká je například pravděpodobnost $P(X = x, Y = y)$ nebo střední hodnota $E XY$.

To je pochopitelně chyba a nelze ji jednoduše napravit. Abychom mohli mluvit o vzájemné interakci více náhodných veličin X_1, X_2, \dots, X_k , musí být nadefinovány na téže množině elementárních jevů Ω , tedy musíme pracovat s jedinou náhodnou veličinou – uspořádanou k -ticí (X_1, X_2, \dots, X_k) a to si žádá opět trochu obecnější definici pravděpodobnosti. Pro názornost tu provedu konstrukci pro případ dvou náhodných veličin X a Y , přičemž pro větší počet se užije obdobného triku.

Mějme dvě náhodné veličiny: X nabývající hodnot z Ω_X a Y nabývající hodnot z Ω_Y . Potom za množinu elementárních jevů Ω budeme brát množinu uspořádaných dvojic $\{(x, y) : x \in \Omega_X, y \in \Omega_Y\}$ a pod $P(X = x, Y = y)$ budeme rozumět $P(x, y)$, kde P je nějaká pravděpodobnost na Ω . V tomto případě $P(X = x)$ značí

$$P(X = x, Y = \text{„cokoli“}) = P(X = x, Y \in \Omega_Y) = \sum_{y \in \Omega_Y} P(X = x, Y = y).$$

Povšimněme si, že pokud se budeme zabývat pouze X, Y či libovolnou nenáhodnou funkcí $f(X, Y)$, dostaneme náhodnou veličinu podle dřívější definice.

Věta. Necht' X a Y jsou náhodné veličiny. Potom

$$E(X + Y) = E X + E Y.$$

Důkaz. Požadované tvrzení dostaneme snadnou úpravou:

$$\begin{aligned} E(X + Y) &= \sum_{x \in \Omega_X, y \in \Omega_Y} (x + y) P(X = x, Y = y) = \\ &= \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} x P(X = x, Y = y) + \sum_{y \in \Omega_Y} \sum_{x \in \Omega_X} y P(X = x, Y = y) = \end{aligned}$$

$$= \sum_{x \in \Omega_X} x P(X = x) + \sum_{y \in \Omega_Y} y P(Y = y) = E X + E Y.$$

Definice. Řekneme, že dvě náhodné veličiny X a Y jsou *nezávislé*, jestliže pro všechna x, y platí

$$P(X = x, Y = y) = P(X = x) \cdot P(Y = y),$$

neboli jevy $\{X = x\}$ a $\{Y = y\}$ jsou *nezávislé*.

Obdobně, jestliže X_1, X_2, \dots, X_n jsou náhodné veličiny nabývající hodnot z $\Omega_1, \Omega_2, \dots, \Omega_n$, potom je nazveme *navzájem nezávislé*, jestliže pro libovolné $x_1 \in \Omega_1, x_2 \in \Omega_2, \dots, x_n \in \Omega_n$ platí

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P(X_1 = x_1) \cdot P(X_2 = x_2) \cdot \dots \cdot P(X_n = x_n).$$

Příklad. Například pokud hodíme n kostkami, tak čísla, jež padnou na jednotlivých kostkách, jsou navzájem nezávislé náhodné veličiny (a to i v případě falešných kostek).

Cvičení. Ověřte, že pokud máme n navzájem nezávislých veličin, pak libovolných k ($k < n$) z nich vybraných jsou opět navzájem nezávislé náhodné veličiny.

Věta. *Nechť X a Y jsou nezávislé náhodné veličiny, potom*

- i) $E(XY) = E X E Y$,
- ii) $\text{Var}(X + Y) = \text{Var } X + \text{Var } Y$.

Cvičení. Dokažte tuto větu: použijte nezávislost a v druhé části definici rozptylu $\text{Var}(X + Y) = E(X + Y)^2 - (E(X + Y))^2$. Pokud vám to přišlo moc jednoduché, formulujte obdobné tvrzení pro více náhodných veličin.

Příklady nejběžnějších rozdělení

Aby nebylo nutné pokaždé znovu odvozovat vlastnosti těch nejběžnějších náhodných veličin, rozlišujeme několik jejich „typů“.

Alternativní rozdělení

Definice. Řekneme, že náhodná veličina X má *alternativní rozdělení* s parametrem p , jestliže X nabývá pouze hodnot 0 a 1 a $P(X = 1) = p$. Značíme $X \sim \text{Alt}(p)$. Pro $X = 1$ budeme mluvit o „úspěchu“, pro $X = 0$ o „neúspěchu“.

Příklad. Typickým příkladem náhodné veličiny s alternativním rozdělením je indikátor nějakého jevu A , tj. náhodná veličina, jež se rovná jedna, pokud jev A nastal (např. padla šestka, padlo sudé číslo, ...), a nula, pokud jev A nenastal.

Střední hodnotu $X \sim \text{Alt}(p)$ spočteme snadno:

$$E X = 1 \cdot p + 0 \cdot (1 - p) = p,$$

rozptyl taktéž:

$$\text{Var } X = E X^2 - (E X)^2 = 1^2 \cdot p + 0^2 \cdot (1 - p) - p^2 = p(1 - p).$$

Binomické rozdělení

Definice. Součet n navzájem nezávislých náhodných veličin s rozdělením $Alt(p)$ nazveme náhodnou veličinou X s *binomickým rozdělením* a značíme $X \sim Binom(n, p)$. Jinak řečeno, X značí počet úspěchů v n nezávislých veličinách s rozdělením $Alt(p)$.

Příklad. Pokud hodíme n spravedlivými kostkami, potom počet šestek, jež nám padne, má binomické rozdělení, konkrétně $Binom(n, \frac{1}{6})$.

Z definice je jasné, že $X \sim Binom(n, p)$ nabývá hodnot $0, 1, \dots, n$. Pravděpodobnost, že nabude hodnoty k , tedy že alternativní náhodné veličiny, jejichž je součet, obsahují právě k úspěchů, je rovna počtu možných výběrů k prvků z n , t.j. $\binom{n}{k}$, krát pravděpodobnost každé z posloupností k úspěchů a $(n - k)$ neúspěchů, t.j. $p^k(1 - p)^{(n-k)}$. Dohromady tedy

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{(n-k)}.$$

Střední hodnotu a rozptyl spočítáme podle věty pro součet náhodných veličin, u rozptylu navíc využijeme nezávislosti alternativně rozdělených veličin Y_1, Y_2, \dots, Y_n , jejichž je X součtem,

$$\begin{aligned} E X &= E \sum_{i=1}^n Y_i = \sum_{i=1}^n E Y_i = np, \\ \text{Var } X &= \text{Var} \sum_{i=1}^n Y_i = \sum_{i=1}^n \text{Var } Y_i = np(1 - p). \end{aligned}$$

Příklad. Známý hráč Samuel Pepys kdysi přišel za jistým Johnem Smithem s následující úlohou: Co si myslíte, že je snazší – hodit 6 kostkami alespoň jednu šestku, hodit 12 kostkami alespoň dvě šestky nebo hodit 18 kostkami alespoň tři šestky?

Ti, co mají dobrou intuici, by měli řešení této úlohy uhádnout i bez počítání. Pro ty ostatní, označme počet šestek při hodu 6, 12 a 18 kostkami po řadě X, Y a Z . Zřejmě $X \sim Binom(6, \frac{1}{6})$, $Y \sim Binom(12, \frac{1}{6})$, $Z \sim Binom(18, \frac{1}{6})$. Tudíž

$$\begin{aligned} P(X \geq 1) &= 1 - P(X = 0) = 1 - \binom{6}{0} \left(\frac{1}{6}\right)^0 \left(\frac{5}{6}\right)^6 \doteq 0,67, \\ P(Y \geq 2) &= 1 - P(Y = 0) - P(Y = 1) = 1 - \binom{12}{0} \left(\frac{1}{6}\right)^0 \left(\frac{5}{6}\right)^{12} - \binom{12}{1} \left(\frac{1}{6}\right)^1 \left(\frac{5}{6}\right)^{11} \doteq 0,62, \\ P(Z \geq 3) &= 1 - P(Z = 0) - P(Z = 1) - P(Z = 2) \doteq 0,60. \end{aligned}$$

Geometrické rozdělení

Definice. Mějme posloupnost nezávislých náhodných veličin s rozdělením $Alt(p)$: Y_1, Y_2, \dots . Náhodnou veličinou s *geometrickým rozdělením* s parametrem p rozumíme počet neúspěchů v předchozí posloupnosti před prvním úspěchem. Značíme $Geom(p)$.

Příklad. Počet hodů kostkou, nežli⁷ nám padne šestka, má rozdělení $Geom(\frac{1}{6})$.

Pravděpodobnost, že $X \sim Geom(p)$ nabude hodnoty k , snadno spočteme jako pravděpodobnost, že nastane k neúspěchů a po nich jeden úspěch:

$$P(X = k) = (1 - p)^k p.$$

⁷Tzn. hod, ve kterém nám padne šestka, se už nepočítá.

Vzorec pro střední hodnotu a rozptyl zde uvedeme bez důkazu:

$$E X = \frac{1-p}{p},$$

$$\text{Var } X = \frac{1-p}{p^2}.$$

Příklad. V každé živýkačce Pedro je ukrytý některý z n obrázků Harry Pottera. Předpokládejme, že každý obrázek se ve živýkačkách vyskytuje se stejnou pravděpodobností. Jaká je střední hodnota počtu živýkaček, které musíme zakoupit, abychom získali všechny Harryho obrázky?

Označme X náhodnou veličinu udávající počet živýkaček, jež musíme zakoupit. Označme Y_i počet živýkaček, zakoupených po dosažení i obrázků do doby, nežli⁸ získáme $(i+1)$ obrázky. Zjevně

$$X = 1 + Y_1 + 1 + Y_2 + 1 + Y_3 + 1 + \dots + Y_{n-1} + 1$$

a veličiny Y_i mají rozdělení $\text{Geom}(\frac{n-i}{n})$. Tudíž

$$E X = n + \sum_{i=1}^{n-1} E Y_i = n + \sum_{i=1}^{n-1} \frac{i}{n-i} = \sum_{i=1}^n \frac{n}{i}.$$

Tedy například pro $n = 20$ vyjde $E X = 72$.

Seriál – Pravděpodobnost a matematická statistika 3

Motto: Po aplikaci preparátu B se 33,3% kuřat uzdravilo, 33,3% kuřat uhynulo a o zbývajících 33,3% kuřat nejsme schopni poskytnout uspokojující informaci. Dosud se nám nepodařilo to třetí kuře chytit.

V závěrečné části seriálu opustíme šedou teorii a vrhneme se do mnohobarevné praxe. Předem varuji, že se zde setkáte s větší mírou neurčitosti a menší mírou preciznosti, než na jakou jste asi zvyklí a než jaká je po vás v korespondenčním semináři vyžadována. Tuto daň za dosažení užitečných a zajímavých výsledků je nutné zaplatit, neboť teorie, jež se za těmito výsledky skrývá, zdaleka přesahuje rámec středoškolské matematiky.

Poissonovo rozdělení

K běžným rozdělením z minula si ještě přidáme tzv. *Poissonovo* rozdělení:

Definice. Řekneme, že náhodná veličina X má Poissonovo rozdělení s parametrem $0 < \lambda < \infty$ (značíme $X \sim \text{Poiss}(\lambda)$), pokud nabývá pouze nezáporných celých čísel a to s pravděpodobností

$$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k = 0, 1, \dots,$$

⁸Ten s $i+1$. obrázkem se už nepočítá.

kde $e \doteq 2,718$ je základ přirozeného logaritmu.

Poissonovo rozdělení se používá k vyjádření počtu náhodných událostí, které nastávají s nějakou intenzitou (tu vyjadřuje parametr λ), tedy například pro počet telefonních hovorů na ústředně v daném časovém intervalu, nebo pro počet lidí, které zvládne obsloužit pokladní v supermarketu v danou hodinu.

Bez důkazu si uveďme střední hodnotu a rozptyl pro $X \sim Poiss(\lambda)$:

$$E X = \text{Var } X = \lambda.$$

Příklad. Pokladní v supermarketu obslouží za pět minut v průměru jednoho zákazníka. Pokud budeme předpokládat, že počet obsloužených zákazníků v pěti minutách má Poissonovo rozdělení s parametrem $\lambda = 1$, spočítejte pravděpodobnost, že pokladní se překoná a obslouží v příštích pěti minutách alespoň tři zákazníky.

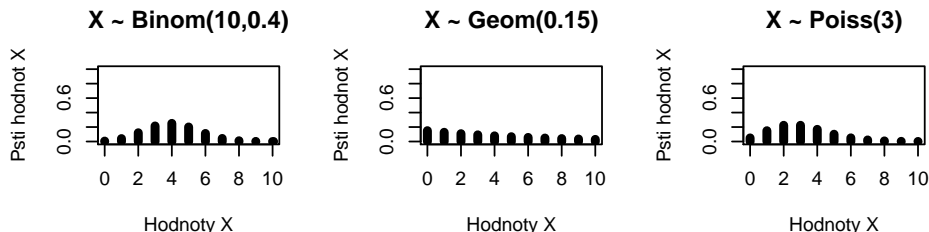
Řešení je jednoduché, spočteme pravděpodobnost doplňkového jevu a odečteme ji od jedničky:

$$\begin{aligned} & P(\text{obslouží alespoň tři}) = \\ &= 1 - P(\text{obslouží dva}) - P(\text{obslouží jednoho}) - P(\text{neobslouží ani jednoho}) = \\ &= 1 - e^{-1} \cdot \frac{1}{2} - e^{-1} \cdot \frac{1}{1} - e^{-1} \cdot \frac{1}{1} \doteq 0,08. \end{aligned}$$

Pravděpodobnost, že pokladní trhne rekord, je tedy přibližně 8%.

Grafické znázornění

Možná už od minula máte pocit, že popis náhodné veličiny jakožto seznamu (či tabulky) pravděpodobností jednotlivých hodnot je sice přesný, ale zato nepřiliš názorný. Nejjednodušším způsobem, jak náhodnou veličinu X graficky znázornit, je vzít jako vodorovnou osu seznam možných hodnot X (či alespoň těch nejvíce pravděpodobných) a ke každé z nich nakreslit sloupec, jehož výška je rovna pravděpodobnosti toho, že X nabude této hodnoty.

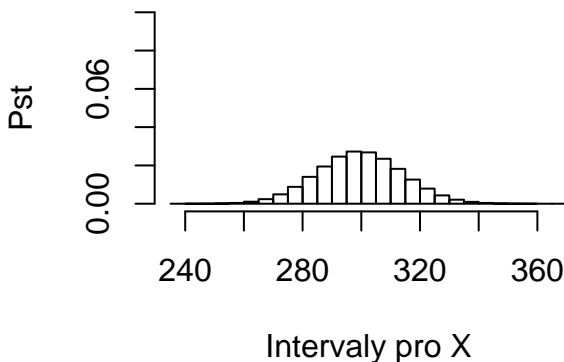


Ukázku na třech různých náhodných veličinách můžete vidět na obrázku výše.

Druhým možným způsobem, jenž oceníte především u veličin s velkým počtem možných hodnot, je použití tzv. *histogramu*. V tomto případě si vodorovnou osu rozdělíte na intervaly stejné délky d a ke každému intervalu I poté nakreslíte sloupec výšky $P(X \in I)/d$.

Histogram je pro $d = 1$ určitou aproximací (odhadem) předchozího obrázku.

Histogram $X \sim \text{Binom}(1000, 0.3)$



Ukázku pro binomické rozdělení s vysokým parametrem n můžete vidět výše.

Posloupnost i.i.d.

Představme si, že máme nějakou náhodnou veličinu X , která nám udává výsledek nějakého pokusu, např. hodu kostkou. Pokud tento pokus budeme několikrát opakovat, budou jeho výsledky X_1, X_2, \dots nezávislé náhodné veličiny, jež budou mít všechny stejné rozdělení jako X .

Definice. Posloupností *nezávislých stejně rozdělených veličin*⁹ neboli posloupností *i.i.d.* (independent, identically distributed) míníme posloupnost náhodných veličin, jež jsou navzájem nezávislé a všechny mají stejné rozdělení X .

Příklad. Posloupností i.i.d. jsou třeba již výše zmiňované hody kostkou. Naopak posloupností i.i.d. není množství srážek v jednotlivé dny (když pršelo dneska, bude asi pršet i zítra, tzn. je porušena nezávislost) nebo vzrůst rostlin, z nichž část jsme hnojili a část ne (porušen předpoklad stejného rozdělení).

V praxi často nemusí být snadné rozhodnout, zda se o posloupnost i.i.d. jedná či nikoli.

Poznámka. V textu o výsledcích experimentů mohou vznikat nejasnosti ohledně toho, zda pod značením X myslíme náhodnou veličinu (t.j. něco neurčitěho) nebo její naměřenou hodnotu (tj. reálné číslo). Domluvme se tedy, že náhodné veličiny budeme značit velkými písmeny, zatímco jejich realizace obvykle písmeny malými.

Otázkou, která nás pochopitelně bude zajímat, jest, kterak z pozorované posloupnosti i.i.d. x_1, x_2, \dots, x_n získat informaci o (neznámém) rozdělení, z něhož pocházejí. Například si představte, že jsme 100 krát hodili falešnou kostkou a zajímá nás, jak je tato kostka falešná (t.j. s jakou pravděpodobností padne které číslo).

⁹Někdy se také mluví o tzv. náhodném výběru.

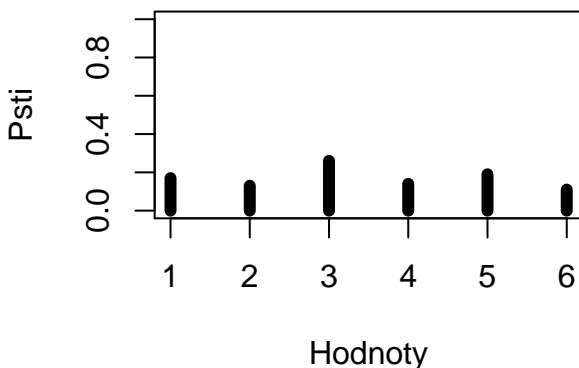
Definice. Empirickým rozdělením¹⁰ k posloupnosti i.i.d. X_1, \dots, X_n budeme rozumět náhodnou veličinu \hat{X} takovou, že

$$P(\hat{X} = x) = \frac{\text{pocet}(x)}{n},$$

kde výrazem $\text{pocet}(x)$ myslíme počet výskytů x v posloupnosti x_1, x_2, \dots, x_n .

Příklad. Při 100 hodech kostkou padlo 17 jedniček, 13 dvojek, 26 trojek, 14 čtyřek, 19 pětěk a 11 šestek. Empirické rozdělení pravděpodobnosti jednotlivých hodnot je vidět na obrázku.

Empirické rozdělení



Průměr

Definice. *Průměrem* posloupnosti i.i.d. X_1, X_2, \dots, X_n budeme rozumět

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Poznámka. Z toho, co víme z minula, můžeme odvodit, že pokud jsou X_1, \dots, X_n realizace X , potom

$$\begin{aligned} E \bar{X}_n &= E \hat{X} = E X, \\ \text{Var } \bar{X}_n &= \frac{1}{n} \text{Var } X. \end{aligned}$$

Průměr je tedy jakýmsi odhadem střední hodnoty.

¹⁰Slovo „empirie“ znamená praxe, zkušenost. Jedná se tedy o odhad rozdělení X na základě zkušenosti z pokusu.

Obdobně můžeme definovat i odhad rozptylu

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Opět bude platit

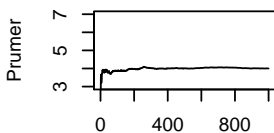
$$E S^2 = \text{Var } X.$$

Tento odhad zde však nepoužijeme.

Zákon velkých čísel

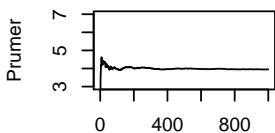
Zmínili jsme, že průměr je odhadem střední hodnoty. Pojdme se přesvědčit, zda tomu tak skutečně je. Pro náhodné veličiny, jejichž rozdělení jsou znázorněna na prvním obrázku, jsme na počítači nagenerovali velké množství jejich realizací. Z těchto jsme pak vždy z prvních i spočítali průměr $\bar{x}_i = \frac{1}{i} \sum_{j=1}^i x_j$ a ten nakreslili do grafu. Pro každou z náhodných veličin jsme tento postup třikrát zopakovali. Výsledky vidíte na obrázcích.

X ~ Binom(10,0.4)



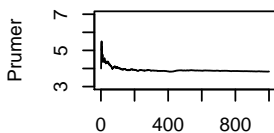
Z kolika velicin je prumer

X ~ Binom(10,0.4)



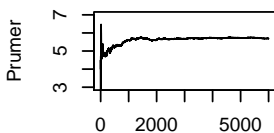
Z kolika velicin je prumer

X ~ Binom(10,0.4)



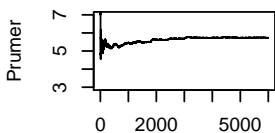
Z kolika velicin je prumer

X ~ Geom(0.15)



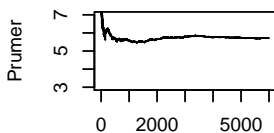
Z kolika velicin je prumer

X ~ Geom(0.15)



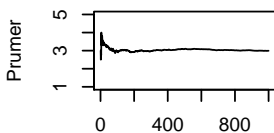
Z kolika velicin je prumer

X ~ Geom(0.15)



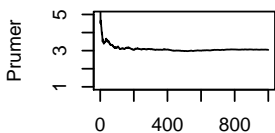
Z kolika velicin je prumer

X ~ Poiss(3)



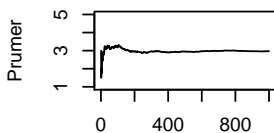
Z kolika velicin je prumer

X ~ Poiss(3)



Z kolika velicin je prumer

X ~ Poiss(3)



Z kolika velicin je prumer

Vidíme, že ačkoli pro malý počet pozorování se průměr od střední hodnoty značně liší, u velkého počtu pozorování je téměř konstantní a jedná se tedy o její dobrý odhad. V statistickém žargonu se tomuto faktu říká zákon velkých čísel (ZVČ):

Věta (ZVČ). *Mějme náhodnou veličinu X s konečnou střední hodnotou EX a i.i.d. posloupnost jejich realizací X_1, X_2, \dots . Potom pro n dostatečně vysoké je \overline{X}_n přibližně EX , značíme*

$$\overline{X}_n \approx EX.$$

Správná otázka pochopitelně je, jak to n má být vysoké či jak blízko se k EX pro dané n dostaneme. Tuto otázku řešit nebudeme či lépe řečeno řešit ji obecně nedokážeme. Všimněte si však, že na předchozích obrázcích máme různé měřítko na vodorovné ose, a že se tedy rychlost přibližování průměru k střední hodnotě může pro různá rozdělení lišit.

Příklad. Student řeší test složený z n stejně obtížných otázek. Každou otázku zodpoví správně s pravděpodobností p , kterou neznáme. Jak tuto pravděpodobnost odhadneme? Jednoduše – dáme mu napsat test a za odhad p zvolíme počet správně zodpovězených otázek podělený n . Pokud bude n dostatečně vysoké, získáme dostatečně přesný odhad p .

Přišlo vám to jednoduché, ba až zřejmé? Dobrá, ukažme si něco obtížnějšího.

Centrální limitní věta

Už víme, že pro velká n se průměr blíží střední hodnotě. To však rozhodně neznamená, že by jí byl roven! Co nás tedy může zajímat, je rozdělení odchylek průměru od střední hodnoty či lépe řečeno výraz

$$Y_n = \sqrt{n} \frac{\overline{X}_n - EX}{\sqrt{\text{Var } X}}.$$

Poznámka. Uvědomme si, že EX a $\text{Var } X$ jsou (nenáhodné) konstanty. Tudiž s použitím vět z minula snadno spočítáme střední hodnotu Y_n

$$E Y_n = \sqrt{\frac{n}{\text{Var } X}} (E \overline{X}_n - EX) = 0,$$

obdobně rozptyl Y_n je roven

$$\text{Var } Y_n = \frac{n}{\text{Var } X} \text{Var } \overline{X}_n = 1.$$

Pro libovolné rozdělení X bude mít tudíž Y_n stejnou střední hodnotu i rozptyl. Za chvíli uvidíme, že pro vysoké n dokonce můžeme říci, že rozdělení Y_n je vždy přibližně stejné, ať bylo původní rozdělení X jakékoli.

Zkusme si pro naše tři staré známé náhodné veličiny počítacem nasimulovat (empirické) rozdělení výše uvedeného výrazu (Y_n). Postupně budeme zkoušet $n = 1$, $n = 10$ a $n = 1\,000$. Pro každý případ nagenenerujeme dostatečně mnoho (10 000) realizací posloupnosti X_1, X_2, \dots, X_n , neboli nagenenerujeme

$$X_{i,1}, X_{i,2}, \dots, X_{i,n} \quad i = 1, \dots, 10\,000,$$

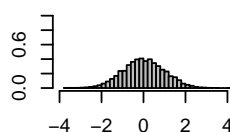
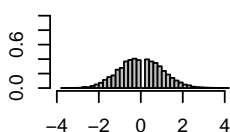
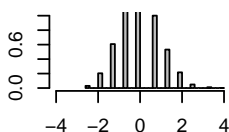
a z každé této i -té posloupnosti spočítáme průměr

$$\overline{X}_{i,n} = \frac{1}{n} \sum_{j=1}^n X_{i,j}.$$

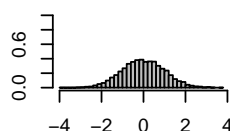
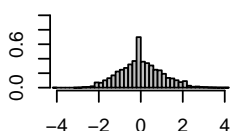
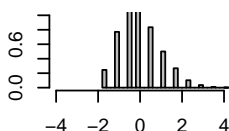
Tak dostaneme empirické rozdělení průměru \widehat{X}_n a z něj i empirické rozdělení výrazu, jenž nás zajímá.

Výsledky (histogramy) můžete vidět na obrázcích:

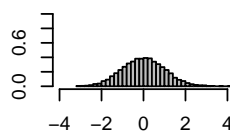
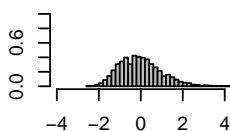
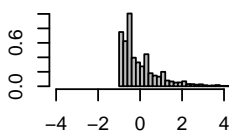
Binomicke rozd., $n = 1$. Binomicke rozd., $n = 10$. Binomicke rozd., $n = 1000$



Poissonovo rozd., $n = 1$. Poissonovo rozd., $n = 10$. Poissonovo rozd., $n = 1000$



Geometricke rozd., $n = 1$. Geometricke rozd., $n = 10$. Geometricke rozd., $n = 100$



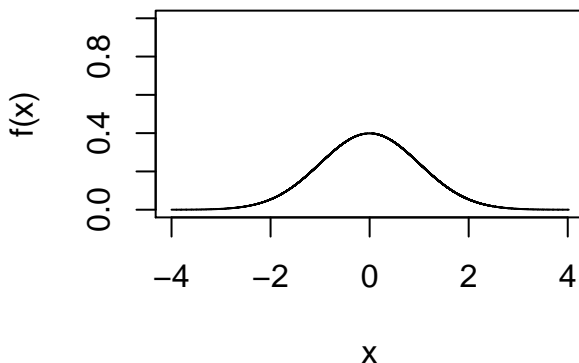
Pokud si výše uvedené obrázky dobře prohlédnete, zjistíte zajímavý fakt, že totiž ačkoli se pro nízké n se jsou obrázky pro různá rozdělení pochopitelně různé, pro n vysoké (zde $n = 1\,000$) jsou si obrázky nápadně podobné.

Kdybychom dále zvyšovali n a počet nagenovaných hodnot a naopak snižovali velikost intervalu v histogramu, začal by se nám (pro libovolné počáteční rozdělení) obrázek histogramu blížit ke grafu funkce

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}},$$

kde $e \doteq 2,718$ je základ základ přirozeného logaritmu. Obrázek (graf) této funkce vidíte dále.

X ~ Norm



Obecně se tomuto faktu říká v statistické hantýrce centrální limitní věta (CLV):

Věta (CLV). *Mějme náhodnou veličinu X s konečným nenulovým rozptylem a i.i.d. posloupnost jejích realizací X_1, X_2, \dots . Potom pro n dostatečně vysoké je rozdělení náhodné veličiny*

$$\sqrt{n} \frac{\bar{X}_n - E X}{\sqrt{\text{Var } X}}$$

vždy přibližně¹¹ stejné. Toto rozdělení budeme nazývat (normované) normální rozdělení. Píšeme rozdělení $\sqrt{n} \frac{\bar{X}_n - E X}{\sqrt{\text{Var } X}} \approx \text{Norm}$.

Poznámka. Názorněji je možná význam věty vidět, když ji napíšeme ve tvaru

$$\text{rozdělení } \bar{X}_n \approx E X + \sqrt{\frac{\text{Var } X}{n}} Z,$$

kde $Z \sim \text{Norm}$. Neboli (přibližné) rozdělení průměru posloupnosti i.i.d. pro dostatečně vysoké n závisí pouze na střední hodnotě a rozptylu X .

Poznámka. Často místo průměru potřebujeme pracovat se sumou, v tomto případě:

$$\sum_{i=1}^n X_i = n \bar{X}_n \approx n E X + \sqrt{n \text{Var } X} Z.$$

¹¹Co toto „přibližně“ znamená by vydalo na samostatný seriál. Jedno z možných zjednodušení říká, že pro libovolně rozdělené X je pravděpodobnost, že Y_n padne do daného intervalu, přibližně tatáž. Neboli že histogramy se příliš neliší.

Náhodnou veličinu $Z \sim Norm$ s normálním rozdělením je za pomoci naší teorie obtížné popsat. Může totiž nabýt libovolné reálné hodnoty, tedy nespočetně mnoha hodnot, každé s nekonečně malou pravděpodobností.

Nelze tedy mluvit o pravděpodobnosti, že Z je rovno z , ale pouze o tom, že $Z \leq z$ nebo že $z_1 \leq Z \leq z_2$. Tyto pravděpodobnosti nelze jednoduše spočítat, lze je však najít v statistických tabulkách nebo vyčíslit za pomoci počítače:

z	$P(Z \leq z)$
-4	0,0000317
-3	0,0013499
-2	0,0227501
-1	0,1586553
0	0,5000000
1	0,8413447
2	0,9772499
3	0,9986501
4	0,9999683

Často potřebujeme také řešit opačnou úlohu. Nalézt pro dané α číslo z takové, že $P(Z \leq z) = \alpha$:

α	$z : P(Z \leq z) = \alpha$
0,5	0
0,75	0,674490
0,9	1,281552
0,95	1,644854
0,975	1,959964
0,99	2,326348
0,995	2,575829

Poznámka. Dlužno podotknout, že normované normální rozdělení je symetrické kolem nuly, neboli platí

$$P(Z \leq -z) = P(Z \geq z) = 1 - P(Z \leq z).$$

Příklad. Stroj na výrobu mikroprocesorů pracuje spolehlivě na 99%. Stroj vyrobil deset tisíc výrobků. Odhadněte počet zmetků. Jaká je pravděpodobnost, že zmetků bude alespoň 120?

Indikátor toho, že je výrobek zmetek, je alternativně rozdělená náhodná veličina (rozdělení $X \sim Alt(0,01)$, $EX = 0,01$, $Var X = 0,0099$). Odhad počtu zmetků v deseti tisících realizacích X je tedy vlastně odhad $10000 \cdot \bar{X}_{10000}$. Veličinu \bar{X}_{10000} odhadneme ze (ZVČ) jakožto 0,01, tedy počet zmetků bude přibližně $10000 \cdot 0,01 = 100$.

Pravděpodobnost, že počet zmetků bude alespoň 120, je stejná jako pravděpodobnost, že $\bar{X}_{10000} - EX \geq 0,002$. Tato pravděpodobnost bude z (CLV) přibližně rovna

$$P\left(\sqrt{\frac{Var X}{n}} Z \geq 0,002\right) = 1 - P\left(Z < \sqrt{\frac{n}{Var X}} 0,002\right) \doteq 1 - P(Z \leq 2) \doteq 0,02275,$$

kde $Z \sim Norm$.

Průměrný počet zmetků tedy bude přibližně 100, pravděpodobnost, že by jich bylo více než 120 je však pouze okolo¹² 2,3%.

Příklad. Obchodník jménem Kolík Aťšepicnu se rozhodl věnovat každému zákazníkovi v jubilejním roce 2000 malého plyšového medvídko. Počet zákazníků Kolíkova krámku má v kterýkoli den v roce Poissonovo rozdělení s parametrem 100. Všichni medvídko musí být objednáni naráz. Kolik jich má Kolík objednat?

Jak správně tušíte není zatím úloha příliš řešitelná, protože když Aťšepicnu objedná z medvídků, je vždy možné (ať už je z vysoké sebevíc), že třeba i v jediný den přijde $z + 1$ zákazníků. Je tedy potřeba určit nějaké riziko, které jsme ochotni podstoupit: Řekněme, že jsme ochotni tolerovat pravděpodobnost nedostatku medvídků nanejvýš na úrovni $\alpha = 5\%$.

Snažíme se tedy najít z takové, že

$$P\left(\sum_{i=1}^{366} X_i > z\right) = 0,05,$$

neboli

$$P\left(\sum_{i=1}^{366} X_i \leq z\right) = 0,95,$$

kde X_1, X_2, \dots, X_{366} jsou i.i.d. s rozdělením *Poiss*(100).

Za použití (CLV) dostáváme

$$\sum_{i=1}^{366} X_i \approx 366 E X + \sqrt{366 \text{Var } X} Z,$$

kde $E X = \text{Var } X = 100$, $Z \sim \text{Norm}$.

Tedy

$$P\left(\sum_{i=1}^{366} X_i \leq z\right) \doteq P(366 \cdot 100 + \sqrt{366 \cdot 100} Z \leq z) = P\left(Z \leq \frac{z - 36600}{\sqrt{36600}}\right) = 0,95.$$

Ovšem z tabulky můžeme vyčíst, že $P(Z \leq x) = 0,95$ právě pro $x = 1,644854$. Dosazením

$$\frac{z - 36600}{\sqrt{36600}} = 1,644854$$

a odtud již snadným výpočtem

$$z = 36914,7.$$

Kolík Aťšepicnu tedy objedná 36 915 medvídků¹³.

Závěr

¹²S pomocí počítače (přesným výpočtem, bez použití CLV) je možné ověřit, že ve skutečnosti je tato pravděpodobnost asi o půl promile nižší, jedná se tedy o vcelku dobrý odhad.

¹³Přesný výpočet na počítači potvrdí, že tentokrát jsme se nespětli ani o jednoho medvídko.

Na závěr tohoto dílu, jakožto i celého seriálu, bych uvedl několik knížek pro ty, kteří by si rádi z tohoto oboru ještě něco přečetli. Řazení je od jednodušších k obtížnějším:

- [1] Disman M.: *Jak se vyrábí sociologická znalost*, pohled na pravděpodobnost a statistiku očima sociologa, vcelku zábavné
- [2] Anděl J.: *Matematika náhody*, nejkrásnější sbírka úloh z diskrétní pravděpodobnosti, jakou znám
- [3] Zvára, K., Štěpán, J.: *Pravděpodobnost a matematická statistika*, pozdější kapitoly předpokládají znalost integrálního počtu
- [4] Anděl J.: *Matematická statistika*, statistikova bible, určeno pro náročnější čtenáře