

## 9. cvičení z PST

11. - 15. dubna 2022

**Poznámky ke kovarianci a korelaci:** Náhodné veličiny (jako funkce na pravděpodobnostním prostoru  $\Omega$ ) tvoří přirozeně (reálný) vektorový prostor (kde ještě navíc dvě veličiny budeme pokládat za totožné, pokud se rovnají s pravděpodobností 1). Na vektorovém *pod*prostoru veličin s konečnou střední hodnotou a konečným rozptylem pak můžeme přirozeným způsobem zavést skalární součin jako

$$\langle X|Y \rangle := E(X \cdot Y)$$

Díky němu můžeme přirozeně zavést *normu*  $\|X\|$  (neboli "délku" vektoru  $X$ ) jako

$$\|X\| := \sqrt{\langle X|X \rangle} = \sqrt{E(X^2)}.$$

Mimo jiné si všimněme, že pro  $X$  je  $\text{var}(X) = \|X - E(X)\|^2$ , takže platí

$$\|norm(X)\| = \left\| \frac{X - EX}{\sqrt{\text{var}(X)}} \right\| = \frac{\|X - EX\|}{\sqrt{\text{var}(X)}} = 1$$

neboli  $norm(X)$  má délku skutečně znormovanou na hodnotu 1.

Skalární součin nám dále umožňuje měřit také úhel mezi dvěma vektory. Pro veličiny  $X$  a  $Y$  je užitečné znát, jestli jejich výchyly vůči středním hodnotám (tj. veličiny  $X - EX$  a  $Y - EY$ ) mají podobné chování (tj. jestli korelují). Zavádíme proto korelaci mezi veličinami  $X$  a  $Y$  jako kosinus úhlu  $\alpha$  mezi vektory  $X - EX$  a  $Y - EY$ , tedy

$$\text{corr}(X, Y) := \frac{\langle X - EX | Y - EY \rangle}{\|X - EX\| \cdot \|Y - EY\|} = \dots = \frac{E(XY) - E(X) \cdot E(Y)}{\sqrt{\text{var}(X) \cdot \text{var}(Y)}}.$$

(Veličiny  $X - EX$  a  $Y - EY$  mají nulovou střední hodnotu).

A kromě toho máme:

$$\text{cov}(X, Y) := \langle X - EX | Y - EY \rangle = \dots = E(XY) - E(X) \cdot E(Y).$$

### Praktický důsledek korelace:

Pokud máme dvě veličiny  $X$  a  $Y$  takové, že

$$X - EX \geq 0 \Leftrightarrow Y - EY \geq 0 \quad \left( \text{což implikuje, že } (X - EX)(Y - EY) \geq 0 \right)$$

pak je  $\text{corr}(X, Y) \geq 0$ .

Obdobně platí: Jestliže

$$X - E(X) \geq 0 \Leftrightarrow Y - E(Y) \leq 0$$

pak je  $\text{corr}(X, Y) \leq 0$ .

Ačkoliv zpětné implikace v obou případech neplatí, přesto nám korelace umožňuje nějakým způsobem zachytit jistou míru kauzální závislosti dvou veličin.

**Poznámka:** Uvědomme si, že existuje několik stupňů "nezávislosti" veličin:

$$X \text{ a } Y \text{ jsou nezávislé} \xrightarrow{\text{(pokud cov ex.)}} \text{cov}(X, Y) = 0 \xrightarrow{\text{(pokud } X, Y \text{ nejsou konst.)}} X \text{ a } Y \text{ jsou lineár. nezáv.}$$

(tj.  $X - E(X)$  a  $Y - E(Y)$  jsou kolmé)

Konstantní veličina  $X$  spolu s jakoukoliv jinou veličinou  $Y$  vždy tvoří vzájemně nezávislé veličiny  $X$  a  $Y$  (tento případ je ale celkem nezajímavý).

**Definice:** Náhodný vektor  $(X, Y)$  má *diskrétní rozdělení*  $\Leftrightarrow$  existuje  $A \subseteq \mathbb{R}^2$ , která je konečná nebo spočetná a taková, že  $P((X, Y) \in A) = 1$ . (Tedy vektor má nejvýše spočetně mnoho “zajímavých” hodnot.)  
V tomto případě pak pro *sdrúženou distribuční funkci* máme

$$F_{(X,Y)}(a, b) = P(X \leq a, Y \leq b) = \sum_{\substack{u \leq a \\ t \leq b}} P(X = u, Y = t).$$

**Věta:** Nechtě náhodný vektor  $(X, Y)$  má *diskrétní rozdělení*. Pak:

$X$  a  $Y$  jsou nezávislé  $\Leftrightarrow P(X = i, Y = j) = P(X = i) \cdot P(Y = j)$  pro všechna  $i, j \in \mathbb{R}$ .

(Srovnejte to s obecně NEdiskrétním případem, kdy je nezávislost popsána jen nerovnostmi:

$$P(X \leq i, Y \leq j) = P(X \leq i) \cdot P(Y \leq j) \text{ pro všechna } i, j \in \mathbb{R}.)$$

**Příklad 9.1** Sdrúžené pravděpodobnosti náhodných veličin  $X$  a  $Y$  jsou dány následující tabulkou:

	$X = 0$	$X = 1$	$X = 2$
$Y = 0$	1/4	1/8	0
$Y = 1$	1/4	1/4	1/8

- (a) Jaká jsou jejich marginální rozdělení? Určete pravděpodobnost  $P(X + Y > 1)$ .
- (b) Jsou veličiny  $X$  a  $Y$  nezávislé? Zdůvodněte. Pokud  $X$  a  $Y$  jsou závislé, popište (jednoznačně určené) rozdělení náhodného vektoru  $(X', Y')$  se stejnými marginálními rozděleními jako  $(X, Y)$ , jehož složky jsou nezávislé.
- (c) Určete střední hodnotu  $E(X \cdot Y)$ . Určete varianční a korelační matici.

### Řešení:

Na začátku bychom si měli pro pořádek ještě ověřit, že součet všech pravděpodobností v tabulce je = 1 (pokud by byl např. < 1, pak nemáme úplnou informaci o rozdělení a nemůžeme dál pokračovat).

- (a) U veličin  $X$  a  $Y$  předpokládáme obory hodnot určené tabulkou. Pak máme

$$X + Y > 1 \Leftrightarrow (X, Y) \in \{(0, 2), (1, 1), (1, 2)\}$$

a tedy

$$P(X + Y > 1) = P((X, Y) \in \{(0, 2), (1, 1), (1, 2)\}) = \frac{1}{4} + \frac{1}{8} + 0 = \frac{3}{8}.$$

Marginální (česky: okrajová) rozdělení náhodného vektoru  $(X, Y)$  jsou rozdělení jeho jednotlivých složek, tedy veličin  $X$  a  $Y$ . Vektor  $(X, Y)$  má diskrétní rozdělení a obě veličiny  $X$  a  $Y$  budou proto mít také diskrétní rozdělení a pro jejich rozdělení platí:

$$P(X = i) = P(X = i, Y \in \mathbb{R}) = \sum_{j \in \mathbb{R}} P(X = i, Y = j)$$

$$P(Y = j) = P(X \in \mathbb{R}, Y = j) = \sum_{i \in \mathbb{R}} P(X = i, Y = j)$$

Hodnoty pravděpodobností získáme tedy sečtením v řádcích (pro  $X$ ) a sloupcích (pro  $Y$ ) naší tabulky:

	$X = 0$	$X = 1$	$X = 2$	$P(Y = j)$
$Y = 0$	1/4	1/8	0	3/8
$Y = 1$	1/4	1/4	1/8	5/8
$P(X = i)$	1/2	3/8	1/8	

Tedy

$$P(X = i) = \begin{cases} 1/2, & i = 0 \\ 3/8, & i = 1 \\ 1/8, & i = 2 \\ 0, & \text{jinak} \end{cases} \quad \text{a} \quad P(Y = j) = \begin{cases} 3/8, & j = 0 \\ 5/8, & j = 1 \\ 0, & \text{jinak.} \end{cases}$$

(b) Protože nyní je např.

$$P(X = 2, Y = 0) = 0 \neq \frac{1}{8} \cdot \frac{3}{8} = P(X = 2) \cdot P(Y = 0).$$

jsou  $X$  a  $Y$  **závislé**.

Nechť  $(X', Y')$  je nyní náhodný vektor s **nezávislými** složkami, které mají stejná marginální rozdělení jako má vektor  $(X, Y)$  tedy

$$P(X' = i) = P(X = i) \quad \text{a} \quad P(Y' = j) = P(Y = j) \quad \text{pro všechna } i, j \in \mathbb{R}.$$

Pro sdružené pravděpodobnosti vektoru  $(X', Y')$  pak tedy platí, že

$$P(X' = i, Y' = j) = P(X' = i) \cdot P(Y' = j) = P(X = i) \cdot P(Y = j)$$

a můžeme je tak popsat následující tabulkou:

	$Y' = 0$	$Y' = 1$	$Y' = 2$	$P(X' = i)$
$X' = 0$	$\frac{1}{2} \cdot \frac{3}{8} = \frac{3}{16}$	$\frac{3}{8} \cdot \frac{3}{8} = \frac{9}{64}$	$\frac{1}{8} \cdot \frac{3}{8} = \frac{3}{64}$	3/8
$X' = 1$	$\frac{1}{2} \cdot \frac{5}{8} = \frac{5}{16}$	$\frac{3}{8} \cdot \frac{5}{8} = \frac{15}{64}$	$\frac{1}{8} \cdot \frac{5}{8} = \frac{5}{64}$	5/8
$P(Y' = j)$	1/2	3/8	1/8	

(c) Pro diskrétní vektor  $(X, Y)$  a borelovskou (tj. téměř každou, např. spojitou) funkci  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$  platí podobná věc jako jsme měli už u náhodné veličiny a sice:

$$E(h(X, Y)) = \sum_{(i,j) \in \mathbb{R}^2} h(i, j) \cdot P(X = i, Y = j).$$

Odsud tedy snadno spočítáme střední hodnotu  $E(XY)$ :

$$\begin{aligned} E(X \cdot Y) &= \sum_{(i,j) \in \mathbb{R}^2} i \cdot j \cdot P(X = i, Y = j) = \\ &= 0 \cdot 0 \cdot \frac{1}{4} + 0 \cdot 1 \cdot \frac{1}{8} + 0 \cdot 2 \cdot 0 + 1 \cdot 0 \cdot \frac{1}{4} + 1 \cdot 1 \cdot \frac{1}{4} + 1 \cdot 2 \cdot \frac{1}{8} = \frac{1}{2} \end{aligned}$$

Kovarianci pak vypočteme ze vztahu

$$\text{cov}(X, Y) = E(XY) - E(X) \cdot E(Y)$$

K tomu potřebujeme znát také

$$E(X) = 0 \cdot \frac{1}{2} + 1 \cdot \frac{3}{8} + 2 \cdot \frac{1}{8} = \frac{5}{8}$$

$$E(Y) = 0 \cdot \frac{3}{8} + 1 \cdot \frac{5}{8} = \frac{5}{8}$$

Takže

$$\text{cov}(X, Y) = E(XY) - E(X) \cdot E(Y) = \frac{1}{2} - \frac{5}{8} \cdot \frac{5}{8} = \frac{7}{64}.$$

Pro korelaci

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)}\sqrt{\text{var}(Y)}}$$

potřebujeme ještě znát rozptyly, takže si je dopočítáme:

$$E(X^2) = 0^2 \cdot \frac{1}{2} + 1^2 \cdot \frac{3}{8} + 2^2 \cdot \frac{1}{8} = \frac{7}{8},$$

$$E(Y^2) = 0^2 \cdot \frac{3}{8} + 1^2 \cdot \frac{5}{8} = \frac{5}{8},$$

$$\text{var}(X) = E(X^2) - (E(X))^2 = \frac{7}{8} - \left(\frac{5}{8}\right)^2 = \frac{31}{64},$$

$$\text{var}(Y) = E(Y^2) - (E(Y))^2 = \frac{5}{8} - \left(\frac{5}{8}\right)^2 = \frac{15}{64}.$$

Všimněme si ještě, že  $Y \sim \text{Alt}\left(\frac{5}{8}\right)$ , takže rozptyl jsme mohli spočítat jako  $\text{var}(Y) = \frac{5}{8} \cdot \left(1 - \frac{5}{8}\right) = \frac{15}{64}$ .

Korelace tedy je

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)} \cdot \sqrt{\text{var}(Y)}} = \frac{\frac{7}{64}}{\sqrt{\frac{31}{64}} \cdot \sqrt{\frac{15}{64}}} = \frac{7}{\sqrt{465}} \doteq 0.32462,$$

Varianční matice je tudíž

$$\text{Var}(X, Y) = \begin{pmatrix} \text{cov}(X, X) & \text{cov}(X, Y) \\ \text{cov}(Y, X) & \text{cov}(Y, Y) \end{pmatrix} = \begin{pmatrix} \text{var}(X) & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \text{var}(Y) \end{pmatrix} = \begin{pmatrix} \frac{31}{64} & \frac{7}{64} \\ \frac{7}{64} & \frac{15}{64} \end{pmatrix}$$

a korelační matice je

$$\text{Corr}(X, Y) = \begin{pmatrix} \text{corr}(X, X) & \text{corr}(X, Y) \\ \text{corr}(Y, X) & \text{corr}(Y, Y) \end{pmatrix} = \begin{pmatrix} 1 & \text{corr}(X, Y) \\ \text{corr}(X, Y) & 1 \end{pmatrix} = \begin{pmatrix} 1 & \frac{7}{\sqrt{465}} \\ \frac{7}{\sqrt{465}} & 1 \end{pmatrix}.$$

Pro zajímavost si ještě můžeme zjistit úhel  $\alpha \in \langle 0, \pi \rangle$  mezi našimi náhodnými veličinami  $X - EX$  a  $Y - EY$ :

$$\alpha = \arccos(\text{corr}(X, Y)) \doteq \arccos(0.32462) \doteq 71.06^\circ.$$

(Důvodem pro závislost  $X$  a  $Y$  je také fakt, že  $\text{cov}(X, Y) \neq 0$ . Zdůrazněme ale, že pokud je kovariance nulová, nemůžeme (pouze na základě její znalosti) o nezávislosti obecně nic říct!)

**Definice:** Náhodný vektor  $(X, Y)$  má spojité rozdělení se *sduženou hustotou pravděpodobnosti*  $f_{X,Y} : \mathbb{R}^2 \rightarrow \langle 0, +\infty \rangle \Leftrightarrow f_{X,Y}$  je integrabilní funkce a pro každou “rozumnou” množinu  $A \subseteq \mathbb{R}^2$  (tj. takovou, která se dá získat z intervalu v  $\mathbb{R}^2$  pomocí sjednocování, průniku a doplňku) platí, že

$$P((X, Y) \in A) = \iint_A f_{X,Y}(x, y) \, dx dy .$$

To nastává právě když

$$F_{X,Y}(a, b) = \int_{-\infty}^a \int_{-\infty}^b f_{X,Y}(x, y) \, dx dy$$

pro každé  $a, b \in \mathbb{R}$ .

Sdužená hustota  $f_{X,Y}$  opět (jako u veličin) NENÍ zdaleka určena jednoznačně, co se týče její funkční hodnoty, ale pouze hodnotami integrálů z této funkce (např. její změnou v konečně mnoha bodech nebo na nějaké hladké křivce se nezmění příslušné integrály, takže i změněná funkce bude také hustotou). Přesněji, dvě nezáporné funkce  $f_{X,Y}$  a  $g_{X,Y}$  (s integrálem rovným jedné) jsou hustotami pro tutéž sduženou distribuční funkci  $F_{X,Y}$  právě když se rovnají *skoro všude* a zapisuje se to jako

$$f_{X,Y} = g_{X,Y} \quad (\text{s.v.}) .$$

(tj. mohou se lišit jen na takové množině  $A \subseteq \mathbb{R}^2$ , že  $\iint_A 1 \, dx dy = 0$ , tj. pokud  $A$  má nulový plošný obsah).

**Příklad 9.2** Sdužená hustota náhodných veličin  $X$  a  $Y$  je

$$f_{(X,Y)}(x, y) = \begin{cases} \frac{1}{2}e^{-x-\frac{y}{2}}, & x > 0, y > 0, \\ 0, & \text{jinak.} \end{cases}$$

- Jaká jsou jejich marginální rozdělení?
- Jsou veličiny  $X$  a  $Y$  nezávislé? Zdůvodněte.
- Jak vypadá jejich korelační matice?
- Určete pravděpodobnost  $P(X > Y)$ .

**Řešení:**

(a) Marginální hustoty (tj. hustoty jednotlivých veličin  $X$  a  $Y$ ) jsou

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) \, dy = \begin{cases} \int_0^{\infty} \frac{1}{2}e^{-x-\frac{y}{2}} \, dy = e^{-x} \cdot [-e^{-\frac{y}{2}}]_0^{\infty} = e^{-x} & \text{pro } x > 0, \\ \int_{-\infty}^{\infty} 0 \, dy = 0 & \text{pro } x \leq 0. \end{cases}$$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) \, dx = \begin{cases} \int_0^{\infty} \frac{1}{2}e^{-x-\frac{y}{2}} \, dx = \frac{1}{2}e^{-\frac{y}{2}} \cdot [-e^{-x}]_0^{\infty} = \frac{1}{2}e^{-\frac{y}{2}} & \text{pro } y > 0, \\ \int_{-\infty}^{\infty} 0 \, dx = 0 & \text{pro } y \leq 0. \end{cases}$$

Vidíme tedy, že obě rozdělení jsou exponenciální, konkrétně  $X \sim \text{Exp}(1)$  a  $Y \sim \text{Exp}(\frac{1}{2})$ .

(b) Složky  $X$  a  $Y$  jsou nezávislé právě tehdy, když

$$f_{X,Y}(x, y) = f_X(x) \cdot f_Y(y) \quad \text{pro skoro všechna } (x, y) \in \mathbb{R}^2,$$

což znamená, že množina bodů, kde uvedená rovnost neplatí má nulový plošný obsah.

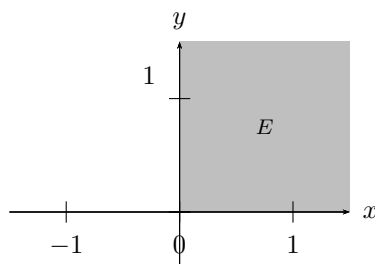
(Podmínce “skoro všude” se nelze vyhnout z toho důvodu, že hustoty nejsou jednoznačně definovány svými hodnotami, ale svými integrály.)

Jak je hned vidět, v našem případě je rovnost splněna dokonce všude, takže  $X$  a  $Y$  JSOU nezávislé.

(c) Z nezávislosti  $X, Y$  plyne okamžitě  $\text{cov}(X, Y) = 0$ , tedy také  $\text{corr}(X, Y) = 0$  a korelační matice je tak

$$\text{Corr}(X, Y) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

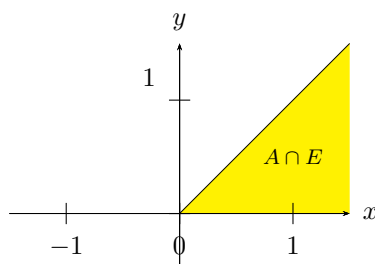
(d) Sdružená hustota  $f_{(X,Y)}$  je nenulová na množině  $E := (0, +\infty)^2$  (vyznačena šedě):



Jev “ $X > Y$ ” je popsán jako  $(X, Y) \in \underbrace{\{(x, y) \in \mathbb{R}^2 \mid x > y\}}_A$ . Pravděpodobnost tohoto jevu tak

dostaneme zintegrováním hustoty přes množinu  $A$ , neboli

$$P(X > Y) = P((X, Y) \in A) = \iint_A f_{X,Y}(x, y) dx dy = \iint_{A \cap E} f_{X,Y}(x, y) dx dy =$$



$$\begin{aligned} &= \{A \cap E : 0 < y < x\} = \int_0^\infty \int_0^x \frac{1}{2} e^{-x-\frac{y}{2}} dy dx = \int_0^\infty e^{-x} \left[ -e^{-\frac{y}{2}} \right]_{y=0}^{y=x} dx = \\ &= \int_0^\infty e^{-x} (1 - e^{-\frac{x}{2}}) dy = \int_0^\infty e^{-x} - e^{-\frac{3}{2}x} dx = \left[ -e^{-x} + \frac{2}{3} e^{-\frac{3}{2}x} \right]_{x=0}^{x=\infty} = 1 - \frac{2}{3} = \frac{1}{3}. \end{aligned}$$

**Poznámky ke kovarianci:**

Kovariance  $\text{cov}(\cdot, \cdot)$  má tyto vlastnosti ( $X, Y, Z$  jsou veličiny,  $a, b, c, d \in \mathbb{R}$ ):

- je lineární v každé složce zvlášť (tj. je bilineární), tedy:

$$\text{cov}(aX + bY, Z) = a \cdot \text{cov}(X, Z) + b \cdot \text{cov}(Y, Z)$$

$$\text{cov}(Z, aX + bY) = a \cdot \text{cov}(Z, X) + b \cdot \text{cov}(Z, Y)$$

- symetrická, tj.  $\text{cov}(X, Y) = \text{cov}(Y, X)$
- pozitivně semi-definitní, tj.  $\text{cov}(X, X) \geq 0$ , kde navíc platí, že:  
 $\text{cov}(X, X) = 0 \Rightarrow \exists \alpha \in \mathbb{R}$ , že  $P(X = \alpha) = 1$  (neboli:  $X$  odpovídá konstantní veličině)
- $\text{cov}(X + c, Y + d) = \text{cov}(X, Y)$ ,
- $\text{cov}(X, X) = \text{var}(X) = (\sigma_X)^2$ .

Pro rozptyl  $\text{var}(\cdot)$  díky tomu máme:

- $\text{var}(aX + c) = \text{var}(aX) = a^2 \cdot \text{var}(X)$ ,
- $\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2 \cdot \text{cov}(X, Y)$ .

**Příklad 9.3** Náhodný vektor  $(X, Y)$  má kovarianční matici.

$$\begin{pmatrix} 3 & -2 \\ -2 & 4 \end{pmatrix}$$

(a) Rozhodněte, zda jsou náhodné veličiny  $X$  a  $Y$  závislé či nezávislé.

(b) Určete  $\text{corr}(X, Y)$ ,  $\text{var}(X + Y)$  a  $\text{cov}(X + Y, -2Y + 1)$ .

**Řešení:**

(a) Kovarianční matice pro  $(X, Y)$  představuje tyto parametry:

$$\begin{pmatrix} \text{cov}(X, X) & \text{cov}(X, Y) \\ \text{cov}(Y, X) & \text{cov}(Y, Y) \end{pmatrix} = \begin{pmatrix} \text{var}(X) & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \text{var}(Y) \end{pmatrix} = \begin{pmatrix} 3 & -2 \\ -2 & 4 \end{pmatrix}$$

Máme tedy  $\text{cov}(X, Y) = -2 \neq 0$  a proto veličiny  $X$  a  $Y$  musí být závislé. (Jestliže by kovariance byla nulová, nemohli bychom udělat o nezávislosti žádný závěr.)

(b) Korelaci určíme jako

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \cdot \text{var}(Y)}} = \frac{-2}{\sqrt{3 \cdot 4}} = -\frac{1}{\sqrt{3}}.$$

a dále máme

$$\text{var}(X + Y) = \text{var}(X) + 2\text{cov}(X, Y) + \text{var}(Y) = 3 - 2 \cdot 2 + 4 = 3$$

a

$$\begin{aligned} \text{cov}(X + Y, -2Y + 1) &= \text{cov}(X + Y, -2Y) = \text{cov}(X, -2Y) + \text{cov}(Y, -2Y) = \\ &= -2\text{cov}(X, Y) - 2\text{cov}(Y, Y) = (-2) \cdot (-2) + (-2) \cdot 4 = -4. \end{aligned}$$

Připomeňme si, co říká **Centrální limitní věta (CLV)**:

Nechť  $X_i$ , pro  $i = 1, 2, \dots$  je posloupnost nezávislých náhodných veličin, které mají stejná rozdělení se střední hodnotou  $\mu$  a (konečným) rozptylem  $\sigma^2$ . Pak pro veličiny

$$Z_n = \sum_{i=1}^n X_i$$

platí, že

$$\lim_{n \rightarrow \infty} P(\text{norm}(Z_n) \leq t) = \Phi(t) \quad \text{pro každé } t \in \mathbb{R}.$$

Neboli: pro velká  $n$  má veličina  $\text{norm}(Z_n)$  přibližně normální rozdělení  $N(0, 1)$ .

Centrální limitní větu můžeme formulovat (namísto pro  $Z_n$ ) také pro tzv. výběrový průměr, tj. veličiny

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

protože pro ně platí  $\text{norm}(\bar{X}_n) = \text{norm}(Z_n)$ .

**Poznámka:** V rámci Centrální limitní věty (níže) se vyskytuje posloupnost nezávislých náhodných veličin, která nejčastěji vzniká následujícím způsobem:

Mějme náhodnou veličinu  $X : \Omega \rightarrow \mathbb{R}$  na pravděpodobnostním prostoru  $\Omega$  (např. pro házení mincí je  $\Omega = \{\text{rub}, \text{líc}\}$  a veličina třeba  $X(\text{líc}) = 1$  a  $X(\text{rub}) = 0$  s rozdělením  $\text{Alt}(p)$ ). Jestliže nyní budeme opakovat (nekonečně) nezávislých pokusů, pak jejich výsledky tvoří posloupnost  $\tilde{\omega} = (\omega_1, \omega_2, \dots)$ , kde  $\omega_i \in \Omega$  pro  $i \in \mathbb{N}$ . Množina všech takovýchto možných posloupností je tedy  $\Omega^{\mathbb{N}}$  (tj. spočetná kartézská mocnina množiny  $\Omega$ ).

Na této množině  $\Omega^{\mathbb{N}}$  lze opět vybudovat pravděpodobnostní prostor tj.  $\sigma$ -algebru  $\tilde{\mathcal{A}}$  na  $\Omega^{\mathbb{N}}$  (která se bude skládat ze spočetných sjednocení množin typu  $\bigtimes_{i=1}^{\infty} A_i = A_1 \times A_2 \times \dots$ , kde  $A_i \subseteq \Omega$  je jev pro každé  $i$ ) a pravděpodobnost bude dána jako  $\tilde{P}\left(\bigtimes_{i=1}^{\infty} A_i\right) = \prod_{i=1}^{\infty} P(A_i)$ .

Výsledek při  $i$ -tém pokusu nyní bude veličina  $X_i : \Omega^{\mathbb{N}} \rightarrow \mathbb{R}$ , definovaná prostě jako  $X_i(\tilde{\omega}) = \omega_i$  pro  $\tilde{\omega} = (\omega_1, \omega_2, \dots)$ . Takovéto veličiny pak budou nezávislé a budou mít rozdělení stejné jako veličina  $X$ .

**Rychlost konvergence v CLV:** Pokud pro veličiny  $X_i$  v CLV navíc ještě je  $\varrho := E(|X_i - \mu|^3) < \infty$ , pak platí Berry–Esseenův odhad chyby (pro všechna  $t \in \mathbb{R}$  a  $n \in \mathbb{N}$ ):

$$\left| F_{\text{norm}(Z_n)}(t) - \Phi(t) \right| < 0.4748 \cdot \frac{\varrho}{\sigma^3 \sqrt{n}}$$

**Odhad chyby v CLV pro alternativní rozdělení:** Pokud mají veličiny  $X_i$  alternativní rozdělení s parametrem  $p$ , tj.  $P(X_i = 1) = p$ , pak

$$\mu = E(X_i) = p, \quad \sigma = \sqrt{D(X_i)} = \sqrt{p(1-p)}, \quad \varrho = E(|X_i - \mu|^3) = p(1-p)(p^2 + (1-p)^2) = \sigma^2(p^2 + (1-p)^2)$$

čímž pro binomické rozdělení  $Z_n = \sum_{i=1}^n X_i \sim \text{Bi}(n, p)$  dostáváme odhad

$$\left| F_{\text{norm}(Z_n)}(t) - \Phi(t) \right| < 0.4748 \cdot \frac{p^2 + (1-p)^2}{\sqrt{np(1-p)}} = 0.4748 \cdot \frac{p^2 + (1-p)^2}{\sqrt{D(Z_n)}}.$$

Aproximace CLV se obvykle používá pro  $D(Z_n) \geq 9$ . Pak je odhad chyby nejvýše:  $\left| F_{\text{norm}(Z_n)}(t) - \Phi(t) \right| < \frac{0.4748}{\sqrt{9}} = 0.159$ .

**Příklad 9.4** *Pravděpodobnost toho, že se za dobu  $T$  porouchá přístroj je  $p = 0.2$ . S jakou pravděpodobností se za dobu  $T$  ze 100 (nezávisle pracujících) přístrojů porouchá*



- (a) alespoň 20,  
 (b) méně než 28,  
 (c) 14 až 26 přístrojů?

**Řešení:**

Pro  $i = 1, \dots, n$  (kde  $n = 100$ ) si zavedeme veličiny

$$X_i = \begin{cases} 1 & , \text{ } i\text{-tý přístroj se porouchá,} \\ 0 & , \text{ } i\text{-tý přístroj bude v pořádku.} \end{cases}$$

Veličiny  $X_i$  budou nezávislé s alternativním rozdělením  $\text{Alt}(p) = \text{Alt}(0.2)$ , protože  $P(X_i = 1) = 0.2$ . Počet porouchaných přístrojů je tedy veličina

$$Z = \sum_{i=1}^{100} X_i$$

kteřá má tudíž binomické rozdělení  $\text{Bi}(n, p) = \text{Bi}(100, 0.2)$ . To sice umíme přesně popsat, ale vyčíslování součtu mnoha velmi malých členů by vedlo ke značným numerickým chybám (a bez softwaru by ani nebylo možné). Proto použijeme CLV, která velmi dobře aproximuje hledané pravděpodobnosti.

Pro  $Z = \sum_{i=1}^n X_i$  tedy máme

$$\begin{aligned} E(Z) &= n \cdot E(X_1) = n \cdot p = 100 \cdot 0.2 = 20 \\ \text{var}(Z) &= n \cdot \text{var}(X_1) = n \cdot p \cdot (1 - p) = 100 \cdot 0.2 \cdot 0.8 = 16 \\ &\Rightarrow \sqrt{\text{var}(Z)} = \sqrt{16} = 4 \end{aligned}$$

Podle CLV můžeme předpokládat, že veličina  $\text{norm}(Z) = \frac{Z - E(Z)}{\sqrt{\text{var}(Z)}} = \frac{Z - 20}{\sqrt{16}}$  má přibližně normované normální rozdělení  $N(0, 1)$ .

To také můžeme chápat tak, že veličina  $Z$  má *přibližně* normální rozdělení

$$N(E(Z), \text{var}(Z)) = N(20, 4^2)$$

tedy že

$$F_Z(t) \doteq \Phi\left(\frac{t - 20}{\sqrt{16}}\right) \text{ pro } t \in \mathbb{R} .$$

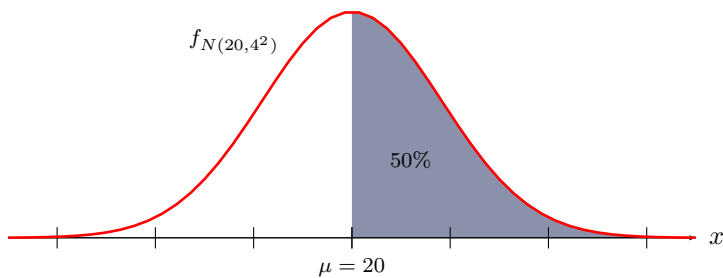
Pak tedy máme:

(a)

$$\begin{aligned} P(Z \geq 20) &= P\left(\frac{Z - 20}{\sqrt{16}} \geq \frac{20 - 20}{\sqrt{16}}\right) = P(\text{norm}(Z) \geq 0) = \\ &= 1 - P(\text{norm}(Z) < 0) \stackrel{\text{(CLV)}}{\doteq} 1 - \Phi(0) = 1 - 0.5 = \mathbf{0.5} . \end{aligned}$$

(Pro srovnání: skutečná hodnota pro binomické rozdělení je **0.5398**.)

Když budeme uvažovat (spojité) rozdělení  $N(20, 4^2)$ , které ALE POUZE aproximuje *původní diskrétní* binomické rozdělení veličiny  $Z$ , můžeme si představit hledanou pravděpodobnost (z HLEDISKA VÝPOČTU) takto:



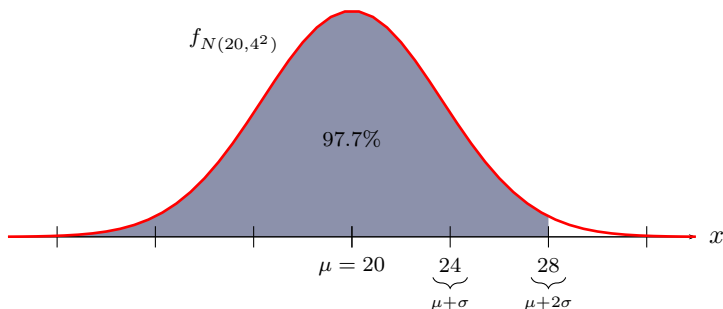
(b)

$$P(Z < 28) = P\left(\frac{Z - 20}{\sqrt{16}} < \frac{28 - 20}{\sqrt{16}}\right) = P(\text{norm}(Z) < 2) \stackrel{(CLV)}{=} \\ \stackrel{(CLV)}{=} \Phi(2) \doteq \mathbf{0.977} .$$

(Pro srovnání: skutečná hodnota pro binomické rozdělení je **0.9658**.)

Když opět budeme uvažovat (spojité) rozdělení  $N(20, 4^2)$ , které POUZE aproximuje *původní diskrétní* binomické rozdělení veličiny  $Z$  a vezmeme do úvahy pravidlo tří sigma, můžeme uvažovat (z HLEDISKA VÝPOČTU) takto:

Protože  $28 = \mu + 2\sigma$ , tak hodnota  $P(Z < 28) \doteq P(N(20, 4^2) < 28)$  nám musí vyjít větší než 95% (viz náčrt).



(c)

$$P(14 \leq Z \leq 26) = P\left(\frac{14 - 20}{\sqrt{16}} \leq \frac{Z - 20}{\sqrt{16}} \leq \frac{26 - 20}{\sqrt{16}}\right) = P(-1.5 \leq \text{norm}(Z) \leq 1.5) = \\ = P(\text{norm}(Z) \leq 1.5) - P(\text{norm}(Z) < -1.5) \stackrel{(CLV)}{=} \Phi(1.5) - \Phi(-1.5) = \\ = 2 \cdot \Phi(1.5) - 1 \doteq 2 \cdot 0.933 - 1 = \mathbf{0.866} .$$

(Pro srovnání: skutečná hodnota pro binomické rozdělení je **0.8973**.)

Když opět budeme uvažovat (spojité) rozdělení  $N(20, 4^2)$ , které POUZE aproximuje *původní diskrétní* binomické rozdělení veličiny  $Z$  a vezmeme do úvahy pravidlo tří sigma, můžeme uvažovat (z HLEDISKA VÝPOČTU) takto:

Protože  $\mu - 2\sigma = 12 < 14 < 16 = \mu - \sigma$  a  $\mu + \sigma = 24 < 16 < 28 = \mu + 2\sigma$ , tak hodnota  $P(14 \leq Z \leq 26) \doteq P(14 \leq N(20, 4^2) \leq 26)$  nám musí vyjít mezi 68% a 95% (viz náčrt).

