



## Statistická rozdělení

Václav Adamec  
vadamec@mendelu.cz

## Úvod



- Náhodná proměnná: matematická veličina, jejíž hodnoty oscilují. Produkt náhodného procesu – lze charakterizovat funkcí
- Hodnoty proměnné v oboru přípustných hodnot
- Rozdělení definují funkční vztah mezi hodnotami náhodné proměnné a četnostmi jejich výskytu
- Spojité rozdělení: nekonečný počet možných hodnot, funkce pravděpodobnostní hustoty (p.d.f.,  $f(y)$ )
- Nespojitá rozdělení: konečný počet možných hodnot, pravděpodobnostní funkce (p.m.f.,  $p(y)$ )
- Kumulativní distribuční funkce (c.d.f.,  $F(y)$ )

## Typy proměnných



- Náhodné proměnné: numerické (kvantitativní)  
nominální (kvalitativní): barva, pohlaví
- Numerické: kardinální (měřitelné)  
ordinální (pořadové): třídy jakosti, stup. klasifikace
- Kardinální: spojité (kontinuální): hmotnost, masná užitkovost  
nespojité (diskrétní): počty mláďat, defektů

## Funkce popisující rozdělení



- Pravděpodobnostní funkce (p.m.f.,  $p(y)$ ) vyjadřuje pravděpodobnost výskytu diskrétní hodnoty  $y$  v oboru možných hodnot

$$p(y) = P(Y=y) = F(y_i) - F(y_{i-1}) \quad \sum_y p(y) = 1.0$$

- Funkce pravděpodobnostní hustoty (p.d.f.,  $f(y)$ ) spojité proměnné

$$f(y) = \frac{dF(y)}{dy} = F'(y)$$

- Kumulativní distribuční funkce (c.d.f.,  $F(y)$ ) vyjadřuje pravděpodob. výskytu hodnoty  $y$  menší nebo rovné  $Y$ .

$$F(y) = P(Y \leq y) = \sum_{y \leq Y} p(y) \quad \text{nespojité případy}$$

$$F(y) = P(Y \leq y) = \int_{-\infty}^y f(y) dy \quad \text{spojité případy}$$

## Střední hodnota (Expectation)



- Expectation (E) definujeme jako první obecný moment.

- Pro spojitou náhodnou proměnnou:

$$M_1' = \int_{-\infty}^{+\infty} y^1 f(y) dy \qquad E(Y) = \int_{-\infty}^{+\infty} yf(y) dy$$

- Pro nespojitou náhodnou proměnnou:

$$M_1' = \sum_y y^1 P(y) \qquad E(Y) = \sum_{i=1}^k y_i P(y_i)$$

## Variance



- Varianci (Var) definujeme jako druhý centrální moment

- Pro spojitou proměnnou:

$$Var(y) = M_2 = \int_{-\infty}^{+\infty} [y - M_1']^2 f(y) dy$$

- Pro nespojitou proměnnou:

$$Var(y) = M_2 = \sum_y [y - M_1']^2 P(y)$$

- Obecně platí:

$$Var(Y) = E(y - E(y))^2 = E(y^2) - (E(y))^2$$

## Příklad



- Je dáno rozdělení pravděpodobností diskrétní proměnné Y: Jaká je střední hodnota a variance ?

y	p(y)	F(y)
0	0,15	0,15
1	0,25	0,40
2	0,25	0,65
3	0,35	1,00

$$E(y) = 0 * 0,15 + 1 * 0,25 + 2 * 0,25 + 3 * 0,35 = 1,8$$

$$Var(y) = (0-1,8)^2 * 0,15 + (1-1,8)^2 * 0,25 + (2-1,8)^2 * 0,25 + (3-1,8)^2 * 0,35 = 0 + 0,16 + 0,01 + 0,504 = 1,16$$

## Pravidla pro expectation



$$E(c) = c$$

$$E(Y_i) = \mu_y$$

$$E(cY) = cE(Y) = c\mu_y$$

$$E(Y \pm c) = E(Y) \pm E(c) = \mu_y \pm c$$

$$E(Y \pm X) = E(Y) \pm E(X) = \mu_y \pm \mu_x$$

$$E(y_1 + y_2) = \mu_y + \mu_y = 2\mu_y$$

$$E(g(Y)) = \sum_{i=1}^k g(y_i) P(y_i)$$

Y diskrétní

$$E(g(Y)) = \int_{-\infty}^{+\infty} g(y) f(y) dy$$

Y spojitá

## Pravidla pro varianci



$$\text{Var}(c) = 0$$

$$\text{Cov}(Y, c) = 0$$

$$\text{Var}(y_i) = \sigma_y^2$$

$$\text{Var}(cY) = c^2 \text{Var}(Y) = c^2 \sigma_y^2$$

$$\text{Var}(Y \pm c) = \text{Var}(Y) + \text{Var}(c) \pm 2 \text{cov}(Y, c) = \sigma_y^2 + 0 \pm 0 = \sigma_y^2$$

$$\text{Var}(Y \pm X) = \text{Var}(Y) + \text{Var}(X) \pm 2 \text{cov}(Y, X) = \sigma_y^2 + \sigma_x^2 \pm 2\sigma_{yx}$$

Najděte

$$E(\bar{y}) = E\left(\frac{\sum_{i=1}^n y_i}{n}\right)$$

$$\text{Var}(\bar{y}) = \text{Var}\left(\frac{\sum_{i=1}^n y_i}{n}\right)$$

## Bernoulliho proměnná



- Binární nespojitá proměnná: pacient přežil ( $y = 1$ )  
pacient nepřežil ( $y = 0$ )

$$Y \sim \text{Bernoulli}(\pi)$$

kde  $\pi$  je parametr rozdělení (pravděpodobnost přežití,  $y = 1$ )

a  $1 - \pi$  je pravděpodobnost nepřežití  $y = 0$

Pravděpodobnostní f-ce:  $P(y) = \pi^y (1 - \pi)^{1-y}$

$$E(y) = \pi$$

$$\text{Var}(y) = \pi(1 - \pi)$$

Jaké jsou hodnoty  $F(y=0)$  a  $F(y=1)$  ?

## Binomické rozdělení



- Bernoulliho experiment opakovaný  $n$  - krát  
 $Y \sim \text{Bin}(n, \pi)$
- Bernoulliho opakování jsou vzájemně nezávislá
- Parametr  $\pi$  je stálý (Bernoulliho pokusy jsou identické, s vrácením)
- Pravděpodobnostní funkce:

$$P(y) = \binom{n}{y} \pi^y (1 - \pi)^{n-y}; y \geq 0; n > 0$$

$$\binom{n}{y} = \frac{n!}{y!(n-y)!}$$

$$E(y) = n\pi$$

$$\text{Var}(y) = n\pi(1 - \pi)$$

## Binomické rozdělení



$$\text{Bin}(n = 6, \pi = 0,5)$$

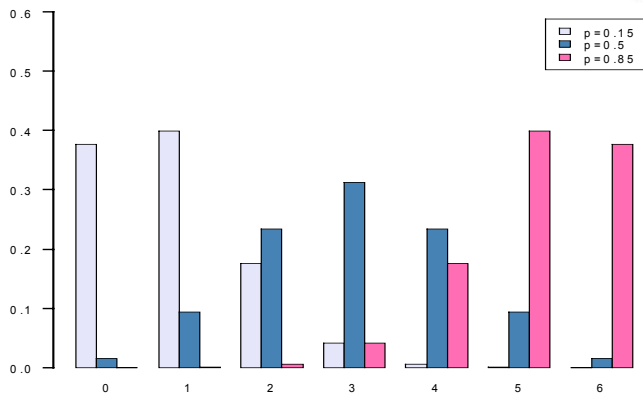
y	c	p(y)	F(y)
0	1	0,015625	0,015625
1	6	0,093750	0,109375
2	15	0,234375	0,343750
3	20	0,312500	0,656250
4	15	0,234375	0,890625
5	6	0,093750	0,984375
6	1	0,015625	1,000000

- Pravděpodobnosti nejvíce 2 synů, nejméně 2 synů ?
- Pravděpodobnost nejvíce 2 dcer ?  $E(y)$  ?  $\text{Var}(y)$  ?

## Binomická rozdělení pro různá $\pi$



Binomial p.m.f.



## Příklad: multinomické rozdělení



- Pravděpodobnost jedináčka  $\pi_1 = 0,6$ ; dvojčat  $\pi_2 = 0,3$ ; trojčat  $\pi_3 = 0,1$

- Jaká je pravděpodobnost, že ve vrzích 13 matek bude 7x jedináček, 4x dvojče a 2x trojče ?

$$P(y_1=7, y_2=4, y_3=2 | \pi) = \frac{13!}{(7!4!2!)} * 0,6^7 * 0,3^4 * 0,1^2 = 0,0584$$

- Jaká je pravděpodobnost, že ve vrzích 13 matek nebude ani jednou vrh s trojčaty ?

$$\text{Výsledek} = 0,2542$$

## Multinomické rozdělení



- Rozšíření binomického opakování na více ( $k > 2$ ) možných výstupů

$$Y_1, Y_2, \dots, Y_k \sim \text{Multinom}(n, \pi_1, \dots, \pi_k)$$

- Opakování jsou opět nezávislá

- Parametry  $\pi_1, \dots, \pi_k$  jsou stálé

- Pravděpodobnostní funkce:

$$P(y_1, y_2, \dots, y_k) = \frac{n!}{\prod_{i=1}^k y_i!} \cdot \prod_{i=1}^k \pi_i^{y_i}; \forall y_i \geq 0$$

$$E(y) = \sum_i y_i p_i$$

$$\text{Var}(y) = n \prod_i \pi_i$$

## Poissonovo rozdělení



- Proměnná: Počty bez přirozeného jmenovatele

$$Y \sim \text{Poisson}(\lambda)$$

- Binomické případy s  $n \rightarrow \infty$  a s malým  $\pi$
- Distribuční parametr  $\lambda = n\pi$  z Binomického rozdělení
- Parametr  $\lambda$  je stálý
- Pravděpodobnostní funkce:

$$P(y) = \frac{e^{-\lambda} \lambda^y}{y!}; y \geq 0$$

$$E(y) = \text{Var}(y) = \lambda$$

- Příklad: Na části chromozomu o dané délce se vyskytují rekombinace v průměru ( $=\lambda$ ) 1,05x za meiozi. Jaké jsou pravděpodobnosti výskytu  $y = 0, 1, 2, \dots, 9$  crossoverů na úseku?

## Poissonovo rozdělení



Poisson( $\lambda = 1,05$ )

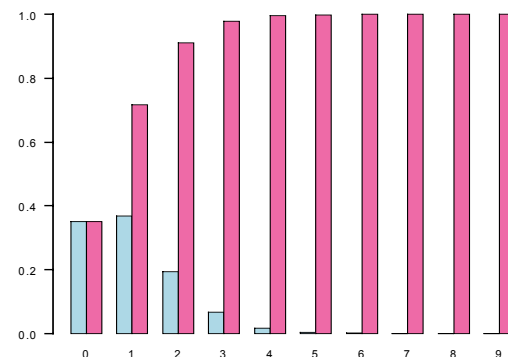
y	p(y)	F(y)
0	0,349938	0,34994
1	0,367435	0,71737
2	0,192903	0,91028
3	0,067516	0,97779
4	0,017723	0,99552
5	0,003722	0,99924
6	0,000651	0,99989
7	0,000098	0,99999
8	0,000013	1,00000
9	0,000001	1,00000

- Přesvědčete se, že  $E(y) = \text{Var}(y) = \lambda$

## Poissonovo rozdělení



Poisson p.m.f. a c.d.f.



## Gaussovo rozdělení



- Spojitá proměnná Y generovaná polyfaktoriální sumací
- Určujících faktorů je mnoho a jsou nezávislé
- Možné hodnoty Y v oboru reálných čísel od  $-\infty$  do  $+\infty$

$$Y \sim N(\mu, \sigma^2)$$

- Funkce pravděpodobnostní hustoty (p.d.f.):

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

- Hodnota f-ce pravděpodobnostní hustoty f(y) není pravděpodobnost !
- $P(y = Y) = 0 !$

$$E(y) = \mu_y$$

$$\text{Var}(y) = \sigma_y^2$$

## Gaussovo rozdělení



- Atributy:
  - Normálních rozdělení je nekonečně mnoho
  - Parametry  $\mu$  a  $\sigma^2$  definují každé normální rozdělení
  - Rozdělení je symetrické podle osy procházející průměrem
  - Lokační míry průměr, medián a modus jsou totožné
  - Plocha pod Gaussovou křivkou odpovídá  $P = 1,0$
  - Pravidlo 34 – 14 – 2 se týká pravděpodobnosti výskytu hodnot (%) mezi  $\mu$  a  $\sigma$ ,  $\sigma$  a  $2\sigma$ ,  $2\sigma$  a  $+\infty$

## Standardizované Gaussovo rozdělení



- Hodnoty z každého normálního rozdělení lze standardizovat
- Standardní normální proměnná  $z$ :

$$z = \frac{y - \mu_y}{\sigma_y} \sim N(\mu_z = 0, \sigma_z^2 = 1)$$

- P.d.f. Std. normálního rozdělení se značí  $\phi(z)$

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

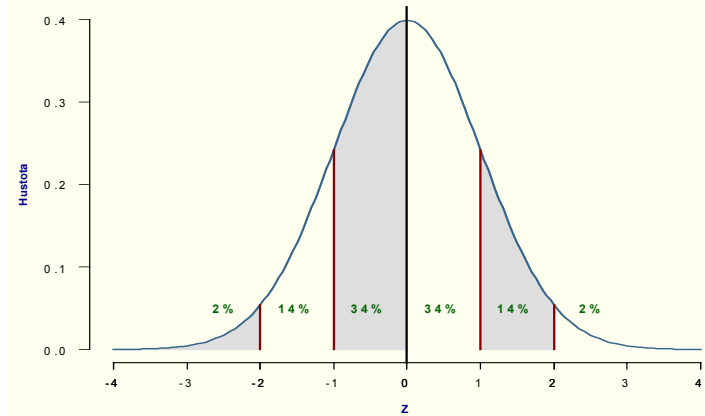
- C.d.f. Std. normálního rozdělení se značí  $\Phi(z)$

$$\Phi(z) = P(Z \leq z) = \int_{-\infty}^z f(z) dz$$

## Proměnná $Z \sim \text{IID}, N(0,1)$



Normovaná Gaussova křivka



## Kalkulace pravděpodobností



- Řešíme integrálem

$$P(z \leq Z) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^z e^{-\frac{(z-\mu)^2}{2\sigma^2}} dz$$

$$P(a \leq z \leq b) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_a^b e^{-\frac{(z-\mu)^2}{2\sigma^2}} dz$$

- Platí že:

$$\phi(z) = \phi(-z) \text{ (důsledek souměrnosti)}$$

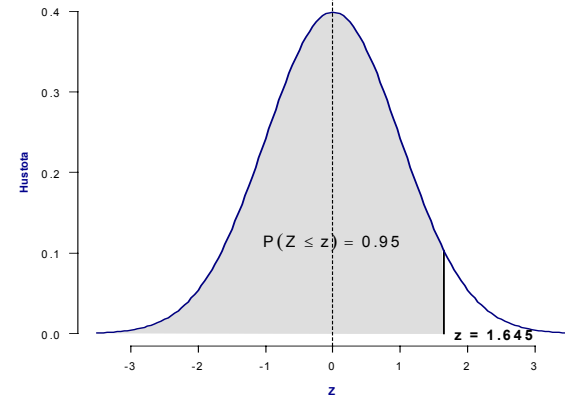
$$\Phi(-z) = 1 - \Phi(z)$$

$$z_{1-p} = -z_p \text{ (vyplývá z předchozího výrazu)}$$

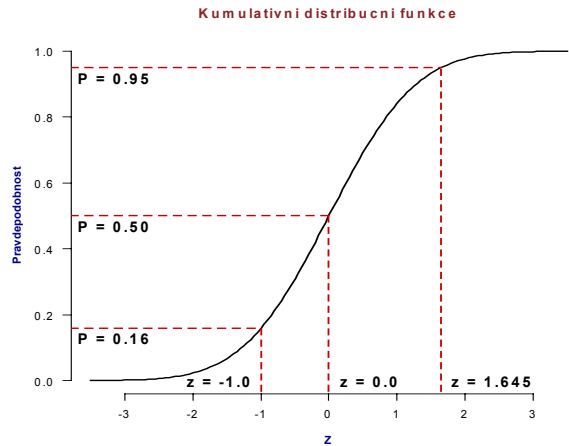
## Kalkulace pravděpodobností



Levostranná pravděpodobnost  $z = 1,645$



## Kumulativní distribuční funkce F(z)



## Pravděpodobnostní výrazy



- Princip:  
Kvantil z lze převést na levostrannou pravděpodobnost P a obráceně při využití souměrnosti rozdělení  $Z \sim N(0,1)$
- Kolik % dojnic se nachází v populaci s průměrem 4500 l a směrodatnou odchylkou 650 l mezi 3800 l až 5000 l ?  

$$z_1 = (3800 - 4500) / 650 = -1,07692$$

$$z_2 = (5000 - 4500) / 650 = 0,76923$$

$$1 - P(-1,07692) - (1 - P(0,76923))$$

$$1 - 0,140758 - (1 - 0,779122)$$

$$1 - 0,140758 - 0,220878 = 0,638364, \text{ tedy } 64 \%$$
- Jaká je pravděpodobnost výskytu dojnice s užitkovostí nad 6000 l ?  

$$z_1 = (6000 - 4500) / 650 = 2,30769$$

$$1 - P(2,30769) = 1 - 0,989492 = 0,0105082, \text{ tedy } 1,05 \%$$

## Pravděpodobnostní výrazy



- Jaká je pravděpodobnost výskytu dojnice s užitkovostí pod 3300 l ?  

$$z_1 = (3300 - 4500) / 650 = -1,84615$$

$$P(-1,84615) = 1 - P(1,84615) = 0,0324352, \text{ tedy } 3,2 \%$$
- 5 % nejlepších dojnic budou využity v ET. Stanovte selekční limit užitkovosti.  

$$z(0,95) = 1,64485$$

$$4500 + 1,64485 * 650 = 5569,15 \text{ l}, \text{ tedy } 5570 \text{ l}$$
- 15 % nejhorších dojnic nebudou zapojeny do reprodukce stáda. Stanovte limit užitkovosti pro vyřazení.  

$$z(0,15) = -1,03643$$

$$4500 - 1,03643 * 650 = 3826,32, \text{ tedy } 3830 \text{ l}$$

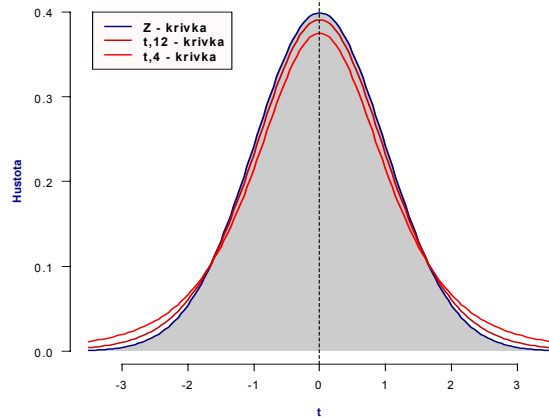
## Studentovo t - rozdělení



- Gossetovo t - rozdělení
- Spojité rozdělení derivátů výběrových veličin mající vztah výběrovému rozptylu s limitovanými stupni volnosti  $\nu$
- Možné hodnoty t v oboru reálných čísel od  $-\infty$  do  $+\infty$
- Rozdělení je unimodální a souměrné kolem nuly
- Platí, že 
$$t_{1-p;\nu} = -t_{p;\nu}$$
- Tvar p.d.f. definován parametrem  $\nu$  (stupně volnosti)
- Vztah k proměnné Z dán výrazem 
$$z_p = t_{p;\nu=\infty}$$
- V praktických případech, je-li přibližně  $\nu > 120$

## Studentovo t - rozdělení

Gaussova a Gossettova krivka



## Rozdělení Chí-kvadrát (Pearsonovo)



- Chí - kvadrát  $\chi^2$  je spojité rozdělení (p.d.f.) nezáporné veličiny
- Součet čtverců standardních normálních odchylek má Chí-kvadrát rozdělení s  $\nu = n - 1$  stupni volnosti

$$\sum_{i=1}^n z_i^2 = z_1^2 + z_2^2 + \dots + z_n^2 \sim \chi_{n-1}^2$$

$$\sum_{i=1}^n z_i^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sigma^2} = \frac{(n-1)s_{n-1}^2}{\sigma^2} \sim \chi_{\nu=n-1}^2$$

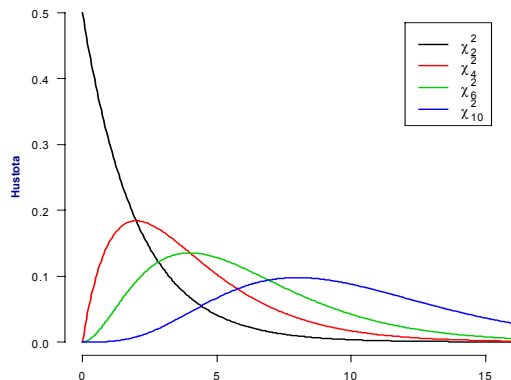
- Parametr rozdělení: stupně volnosti  $\nu$  dány počtem nezávislých odchylek od průměru
- Počet stupňů volnosti  $\nu$  určuje tvar křivky p.d.f.
- Užitečné při testování rozptylu a jeho derivátů (sumy čtverců)

$$E(\chi_{\nu}^2) = \nu$$

$$Var(\chi_{\nu}^2) = 2\nu$$

## Rozdělení Chí-kvadrát (Pearsonovo)

Chi-kvadrát densita a stupne volnosti



## Fisherovo - Snedecorovo rozdělení



- Je spojité rozdělení pro podíl dvou nezávislých nezáporných veličin (rozptylů, součtu čtverců)
- U každé z veličin se předpokládá Chí-kvadrát rozdělení  $\chi_{\nu_1}^2$  a  $\chi_{\nu_2}^2$
- Podíl těchto veličin má F rozdělení se stupni volnosti  $\nu_1$  (proměnná v čitateli) a  $\nu_2$  (proměnná ve jmenovateli)
- F-rozdělení je vždy asymetrické
- Platí, že:

$$F_{p;\nu_1;\nu_2} = \frac{1}{F_{1-p;\nu_2;\nu_1}}$$

$$F_{p;1;\nu_2} = t_{p/2;\nu_2}^2$$

$$F_{0.95;3;7} = 1 / F_{0.05;7;3} = 1 / 0.230053 = 4.34683$$

$$F_{0.95;1;4} = t_{0.975;4}^2 = 2.77645^2 = 7.70865$$



# F-rozdělení

F - rozdělení

