

# Analýza a klasifikace dat – přednáška 7 – doplnění

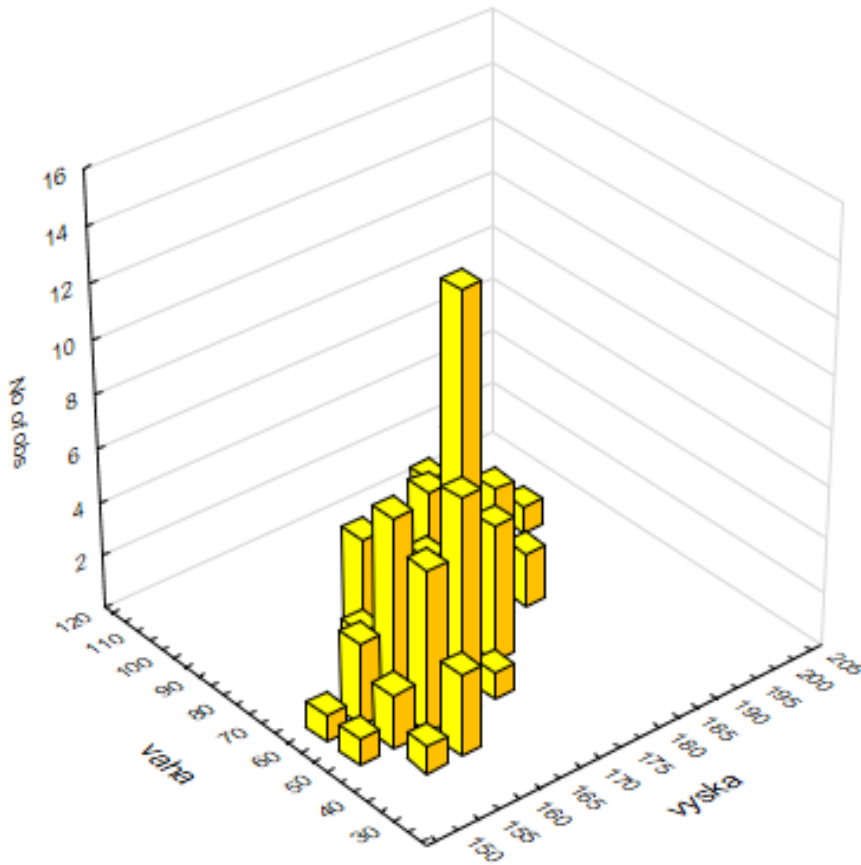


RNDr. Eva Janoušová

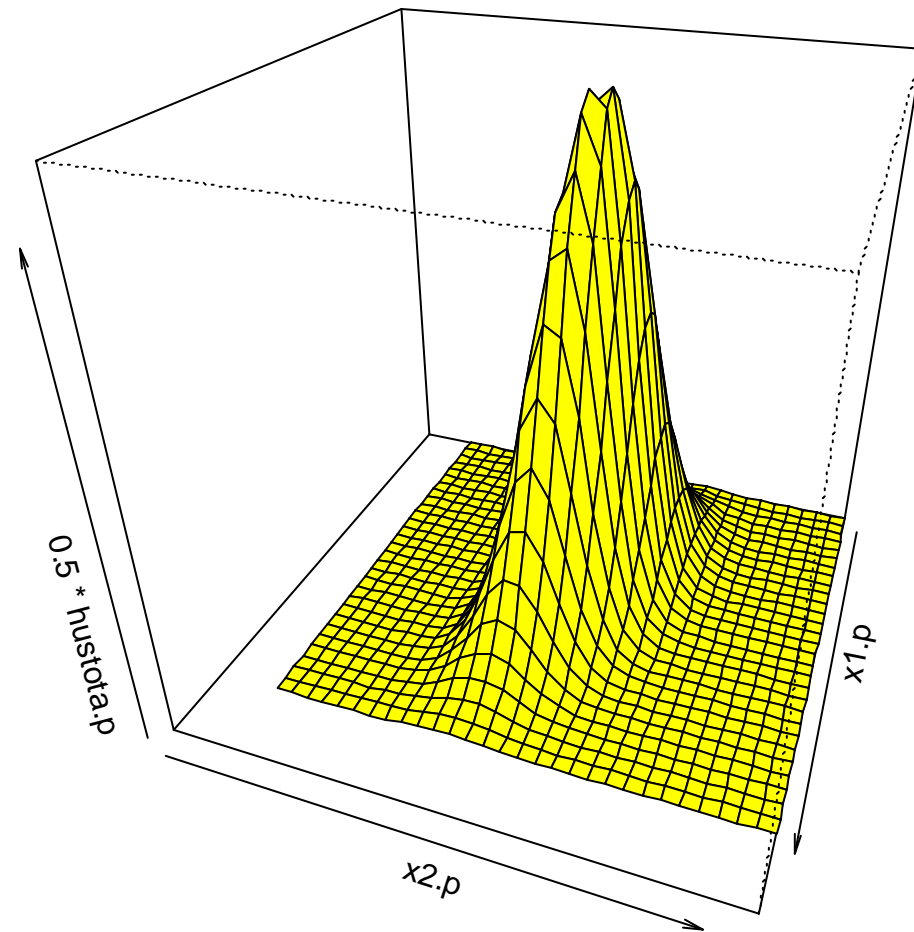
Podzim 2014

# Motivace

Dvourozměrný histogram



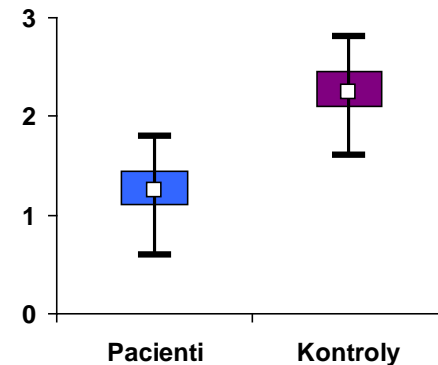
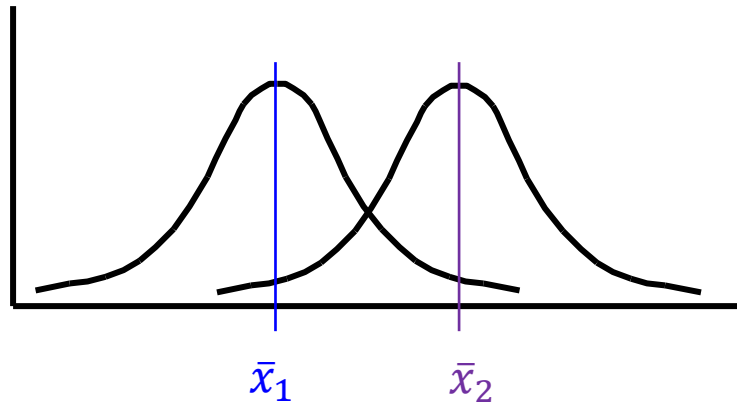
Hustota dvourozměrného normálního rozdělení



# Vícerozměrný t-test

# Jednorozměrný dvouvýběrový t-test

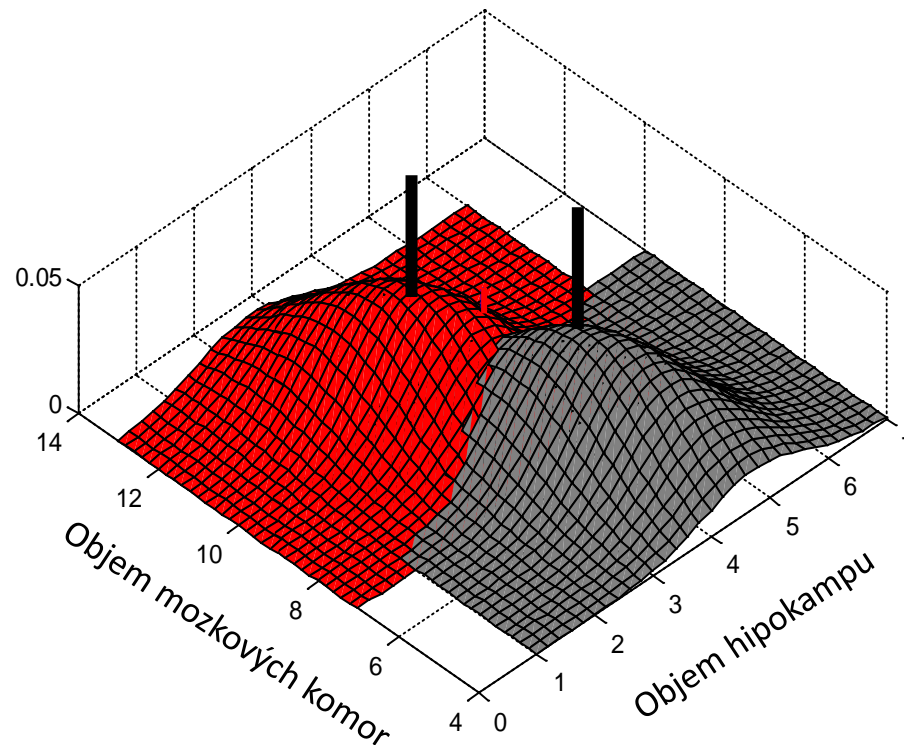
- Srovnáváme dvě skupiny dat, které jsou na sobě nezávislé – mezi objekty neexistuje vazba.
- Příklady: srovnání objemu hipokampu u mužů a u žen, srovnání kognitivního výkonu podle dvou kategorií věku,...



- Předpoklad: **normalita dat v OBOU skupinách, shodnost (homogenita) rozptylů** v obou skupinách
- Testová statistika: 
$$t = \frac{\bar{x}_1 - \bar{x}_2 - c}{s_* \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$
, kde  $s_*$  je vážená směrodatná odchylka,  $c$  je konstanta, o kterou se rozdíl průměrů má lišit (většinou rovna 0)

# Vícerozměrný t-test

- Srovnáváme dvě skupiny dat, které jsou na sobě nezávislé – mezi objekty neexistuje vazba.
- Na rozdíl od jednorozměrného dvouvýběrového t-testu jsou dvě skupiny dat popsány více proměnnými.



# Vícerozměrný t-test

## Jednorozměrný dvouvýběrový t-test:

- testová statistika:  $t = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{s \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}$ , kde  $t \sim T(n_x + n_y - 2)$  ← Studentovo rozdělení
- $s$  je vážená směrodatná odchylka  $s^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1)s_y^2}{(n_x - 1) + (n_y - 1)}$
- $(\mu_x - \mu_y) = c$  je konstanta, o kterou se rozdíl průměrů má lišit (většinou  $c = 0$ )
- nulová hypotéza zamítnuta, pokud  $|t| > t_{crit}$

## Je ekvivalentní testu:

- $t^2 = \left( \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{s \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} \right)^2 = (\bar{z} - \mu_z) \left[ s^2 \left( \frac{1}{n_x} + \frac{1}{n_y} \right) \right]^{-1} (\bar{z} - \mu_z)$ , kde  $t^2 \sim F(1, n_x + n_y - 2)$  ← F rozdělení  
 $\bar{z} = \bar{x} - \bar{y}$  a  $\mu_z = \mu_x - \mu_y$

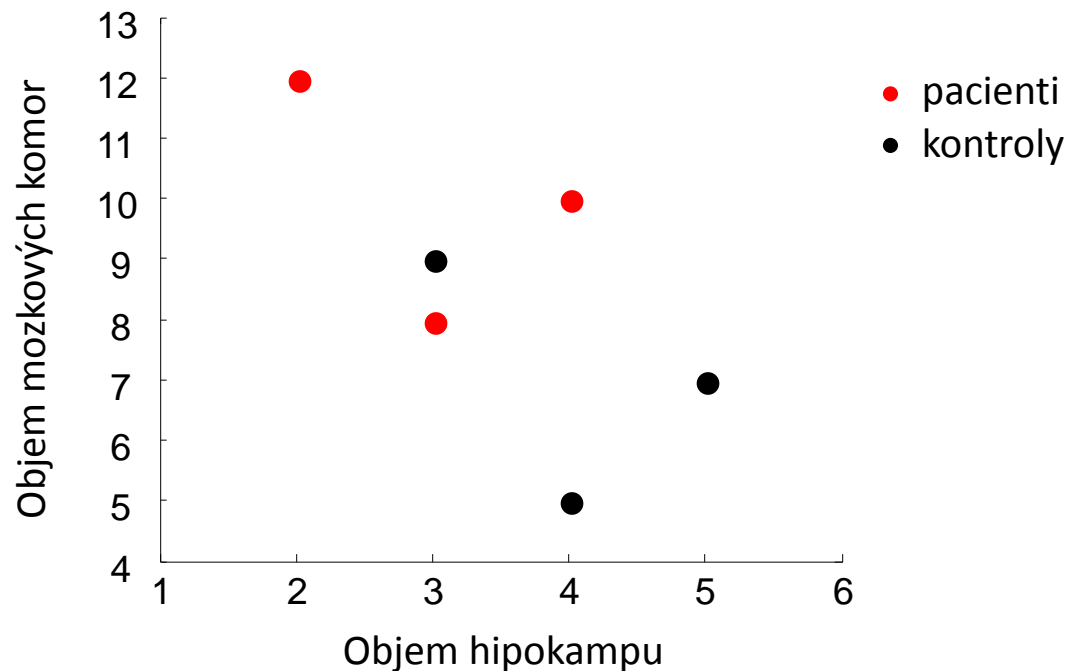
## Vícerozměrný t-test:

- dvouvýběrová Hotellingova  $T^2$  testová statistika:  $T^2 = (\bar{X} - \bar{Y})^T \left[ S \left( \frac{1}{n_x} + \frac{1}{n_y} \right) \right]^{-1} (\bar{X} - \bar{Y})$
- kde  $S$  je vážená kovarianční matice  $S = \frac{(n_x - 1)S_X + (n_y - 1)S_Y}{(n_x - 1) + (n_y - 1)}$
- $T^2 \sim \chi^2(k)$ ; pro malé  $n_x$  a  $n_y$  je lepší použít:  $F = \frac{n-k}{k(n-1)} T^2 \sim F(k, n-k)$ , kde  $n = n_x + n_y - 1$  ← F rozdělení
- nulová hypotéza zamítnuta, pokud  $F > F_{crit}$

# Úkol 1

- Zjistěte, zda se liší skupina pacientů se schizofrenií od zdravých subjektů na základě parametrů popisujících objem mozkových struktur subjektů.

$$\mathbf{X}_D = \begin{bmatrix} 2 & 12 \\ 4 & 10 \\ 3 & 8 \end{bmatrix}, \mathbf{X}_H = \begin{bmatrix} 5 & 7 \\ 3 & 9 \\ 4 & 5 \end{bmatrix}$$



# Úkol 1

- Zjistěte, zda se liší skupina pacientů se schizofrenií od zdravých subjektů na základě parametrů popisujících objem mozkových struktur subjektů.

$$\mathbf{X}_D = \begin{bmatrix} 2 & 12 \\ 4 & 10 \\ 3 & 8 \end{bmatrix}, \mathbf{X}_H = \begin{bmatrix} 5 & 7 \\ 3 & 9 \\ 4 & 5 \end{bmatrix}$$



# Úkol 1 - řešení

Vícerozměrné průměry:

$$\bar{\mathbf{x}}_D = \left[ \frac{1}{n_D} \sum_{i=1}^{n_D} x_{i1} \quad \frac{1}{n_D} \sum_{i=1}^{n_D} x_{i2} \right] = [3 \quad 10]$$

$$\bar{\mathbf{x}}_H = \left[ \frac{1}{n_H} \sum_{i=1}^{n_H} x_{i1} \quad \frac{1}{n_H} \sum_{i=1}^{n_H} x_{i2} \right] = [4 \quad 7]$$

Výběrové kovarianční matice:

$$\mathbf{S}_D = \begin{bmatrix} s_{11}^D & s_{12}^D \\ s_{21}^D & s_{22}^D \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 4 \end{bmatrix}$$

$$\mathbf{S}_H = \begin{bmatrix} s_{11}^H & s_{12}^H \\ s_{21}^H & s_{22}^H \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 4 \end{bmatrix}$$

Vážená kovarianční matice:

$$\mathbf{S} = \begin{bmatrix} 1 & -1 \\ -1 & 4 \end{bmatrix}$$

# Úkol 1 - řešení

Vícerozměrné průměry:

$$\bar{\mathbf{x}}_D = \left[ \frac{1}{n_D} \sum_{i=1}^{n_D} x_{i1} \quad \frac{1}{n_D} \sum_{i=1}^{n_D} x_{i2} \right] = [3 \quad 10]$$

$$\bar{\mathbf{x}}_H = \left[ \frac{1}{n_H} \sum_{i=1}^{n_H} x_{i1} \quad \frac{1}{n_H} \sum_{i=1}^{n_H} x_{i2} \right] = [4 \quad 7]$$

Výběrové kovarianční matice:

$$\mathbf{S}_D = \begin{bmatrix} s_{11}^D & s_{12}^D \\ s_{21}^D & s_{22}^D \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 4 \end{bmatrix}$$

$$\mathbf{S}_H = \begin{bmatrix} s_{11}^H & s_{12}^H \\ s_{21}^H & s_{22}^H \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 4 \end{bmatrix}$$

Vážená kovarianční matice:

$$\mathbf{S} = \begin{bmatrix} 1 & -1 \\ -1 & 4 \end{bmatrix}$$

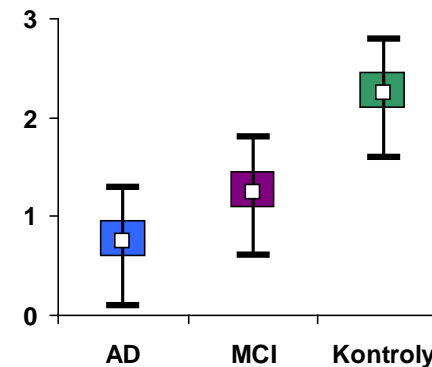
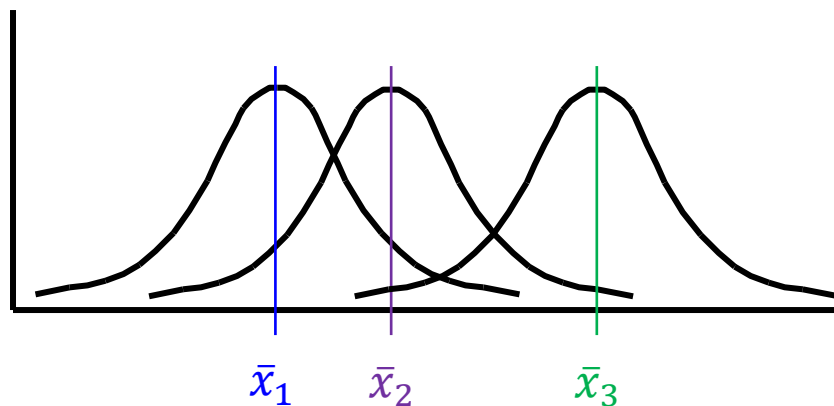
Vícerozměrný t-test:

n	5
k	2
T <sup>2</sup>	3,5
F	1,31
df1	2
df2	3
α	0,05
F-crit	9,55
p-hodnota	0,389

# Vícerozměrná analýza rozptylu

# Analýza rozptylu (ANOVA) jednoduchého třídění

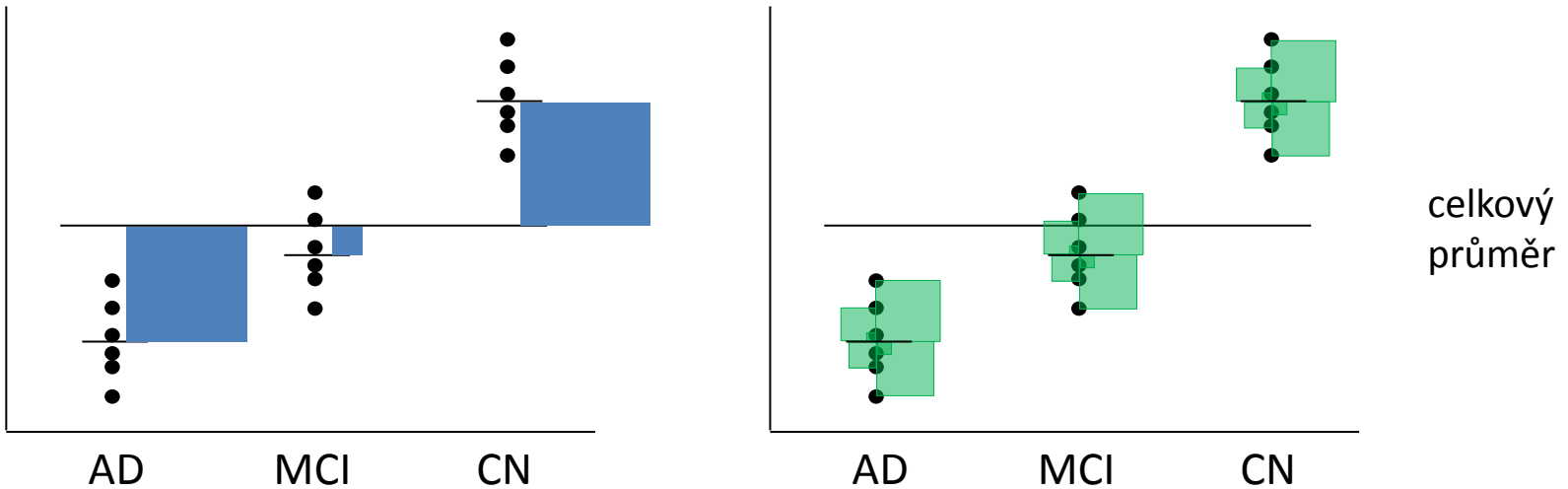
- Srovnáváme tři a více skupin dat, které jsou na sobě nezávislé (mezi objekty neexistuje vazba).
- Příklady: srovnání objemu hipokampu u pacientů s AD, pacientů s MCI a kontrol; srovnání kognitivního výkonu podle čtyř kategorií věku.



- Předpoklady: **normalita dat ve VŠECH skupinách, shodnost (homogenita) rozptylů VŠECH srovnávaných skupin**, nezávislost jednotlivých pozorování.
- Testová statistika: 
$$F = \frac{S_A / df_A}{S_e / df_e}$$

# Analýza rozptylu (ANOVA) – princip

- Srovnání variability (rozptylu) mezi výběry s variabilitou uvnitř výběrů.



- Tabulka analýzy rozptylu jednoduchého třídění (One-Way ANOVA):

Variabilita	Součet čtverců	Počet stupňů volnosti	Průměrný čtverec	F statistika	p-hodnota
Mezi skupinami	$S_A$	$df_A = k - 1$	$MS_A = S_A / df_A$	$F = \frac{S_A / df_A}{S_e / df_e}$	$p$
Uvnitř skupin (reziduální var.)	$S_e$	$df_e = n - k$	$MS_e = S_e / df_e$		
Celkem	$S_T$	$df_T = n - 1$			

# Analýza rozptylu jako lineární model

- Analýza rozptylu pro jednu vysvětlující proměnnou (jednoduché třídění) lze zapsat jako lineární model:

$$Y_{ij} = \mu_i + e_{ij} = \mu + \alpha_i + e_{ij}$$

Populační průměr       $\alpha_i$        $e_{ij}$

Reziduum  
 $i$ -tý efekt faktoru A

- Nulovou hypotézu pak lze vyjádřit jako:  $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_k$
- Rozšířením tohoto zápisu můžeme definovat další modely ANOVA:** více faktorů, hodnocení interakcí, opakovaná měření na jednom subjektu.

# Analýza rozptylu dvojného třídění

- Uvažujeme dvě vysvětlující proměnné zároveň.
- Zápis modelu:

$$Y_{ij} = \mu + \alpha_i + \beta_j + e_{ij}$$

Diagrammatic explanation of the model components:
 

- $\mu$ : Populační průměr (Population mean)
- $\alpha_i$ :  $i$ -tý efekt faktoru A
- $\beta_j$ :  $j$ -tý efekt faktoru B
- $e_{ij}$ : Reziduum (Residual)

- Nulové hypotézy pak máme dvě:  $H_{01} : \alpha_1 = \alpha_2 = \dots = \alpha_k$  ,  $H_{02} : \beta_1 = \beta_2 = \dots = \beta_r$

Variabilita	Součet čtverců	Počet stupňů volnosti	Průměrný čtverec	F statistika	p-hodnota
Faktor A	$S_A$	$df_A = k - 1$	$MS_A = S_A / df_A$	$F_A$	$p$
Faktor B	$S_B$	$df_B = r - 1$	$MS_B = S_B / df_B$	$F_B$	$p$
Rezidua	$S_e$	$df_e = (k - 1)(r - 1)$	$MS_e = S_e / df_e$		
Celkem	$S_T$	$df_T = n - 1 = kr - 1$			

# Analýza rozptylu dvojného třídění s interakcí

- Uvažujeme dvě vysvětlující proměnné a zároveň i jejich společné působení.

- Zápis modelu:

$$Y_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ij}$$

Diagrammatic explanation of the model components:
 

- $\mu$ : Populační průměr (Population mean)
- $\alpha_i$ :  $i$ -tý efekt faktoru A
- $\beta_j$ :  $j$ -tý efekt faktoru B
- $\gamma_{ij}$ : Interakce (Interaction)
- $e_{ij}$ : Reziduum (Residual)

- Nulové hypotézy pak máme tři:

$$H_{01} : \gamma_{11} = \gamma_{12} = \dots = \gamma_{kr} \quad H_{02} : \alpha_1 = \alpha_2 = \dots = \alpha_k \quad H_{03} : \beta_1 = \beta_2 = \dots = \beta_r$$

Variabilita	Součet čtverců	Počet stupňů volnosti	Průměrný čtverec	F statistika	$\rho$ -hodnota
Faktor A	$S_A$	$df_A = k - 1$	$MS_A = S_A / df_A$	$F_A$	$\rho$
Faktor B	$S_B$	$df_B = r - 1$	$MS_B = S_B / df_B$	$F_B$	$\rho$
Interakce AxB	$S_{AB}$	$df_{AB} = (k - 1)(r - 1)$	$MS_{AB} = S_{AB} / df_{AB}$	$F_{AB}$	$\rho$
Rezidua	$S_e$	$df_e = n - kr$	$MS_e = S_e / df_e$		
Celkem	$S_T$	$df_T = n - 1$			



# Úkol 2

Zjistěte, zda má vliv pohlaví a typ léku na počet uzdravených pacientů s leukémií.

Pohlaví	Typ léku	Počet uzdravených pacientů
M	placebo	1
M	lék 1	1
M	lék 2	6
Z	placebo	3
Z	lék 1	4
Z	lék 2	9

# Úkol 2 - řešení

Zjistěte, zda má vliv pohlaví a typ léku na počet uzdravených pacientů s leukémií.

Překódování:

Pohlaví	Typ léku	Počet uzdravených pacientů
1	1	1
1	2	1
1	3	6
2	1	3
2	2	4
2	3	9

Legenda:

Pohlaví: 1=M  
2=Z

Typ léku: 1=placebo  
2=lék 1  
3=lék 2

# Úkol 2 - řešení

Pohlaví	Typ léku	Počet uzdrav. pacientů
1	1	1
1	2	1
1	3	6
2	1	3
2	2	4
2	3	9

# Úkol 2 - řešení

Pohlaví	Typ léku	Počet uzdrav. pacientů	
1	1	$X_{1..} = 8$ $M_{1..} = 8/3$	1
1	2		1
1	3		6
2	1	$X_{2..} = 16$ $M_{2..} = 16/3$	3
2	2		4
2	3		9

$$a = 2; \quad b = 3; \quad c = 1; \quad n = 6;$$

$$X_{.1.} = 4; \quad M_{.1.} = 4/2 = 2$$

$$X_{.2.} = 5; \quad M_{.2.} = 5/2 = 2,5$$

$$X_{.3.} = 15; \quad M_{.3.} = 15/2 = 7,5$$

$$X_{...} = 24; \quad M_{...} = 24/6 = 4$$

**Součet čtverců pro faktor A (pohlaví):**

počet stupňů volnosti:  $f_A = a - 1 = 1$

$$S_A = bc \sum_{i=1}^a (M_{i..} - M_{...})^2 = 3 \cdot ((8/3 - 4)^2 + (16/3 - 4)^2) = 32/3 = 10,67$$

**Součet čtverců pro faktor B (typ léku):**

počet stupňů volnosti:  $f_B = b - 1 = 2$

$$S_B = ac \sum_{j=1}^b (M_{.j.} - M_{...})^2 = 2 \cdot ((2 - 4)^2 + (2,5 - 4)^2 + (7,5 - 4)^2) = 37$$

**Celkový součet čtverců :**

počet stupňů volnosti:  $f_T = n - 1 = 5$

$$S_T = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (X_{ijk} - M_{...})^2 = (1 - 4)^2 + (1 - 4)^2 + \dots + (9 - 4)^2 = 48$$

**Reziduální součet čtverců :**

počet stupňů volnosti:  $f_E = n - a - b + 1 = 2$

$$S_E = S_T - S_A - S_B = 0,33$$

# Úkol 2 - řešení

Tabulka analýzy rozptylu dvojného třídění:

Zdroj variability	Součet čtverců	Stupně volnosti	Podíl S/f	$F = \frac{S/f}{S_E/f_E}$
Faktor A (pohlaví)	$S_A = 10,67$	$f_A = 1$	10,67	63,99
Faktor B (typ léku)	$S_B = 37$	$f_B = 2$	18,5	110,98
Reziduální	$S_E = 0,33$	$f_E = 2$	0,16	-
Celkový	$S_T = 48$	$f_T = 5$	-	-

Srovnání s kvantily:

$F_A = 63,99 > F_{0,95}(1,2) = 18,1 \rightarrow$  pohlaví má vliv na počet uzdravených pacientů

$F_B = 110,98 > F_{0,95}(2,2) = 19 \rightarrow$  typ léku má vliv na počet uzdravených pacientů

# Úkol 2 – řešení v softwaru STATISTICA

Zjistěte, zda má vliv pohlaví a typ léku na počet uzdravených pacientů s leukémií.

Pohlaví	Typ léku	Počet uzdrav. pacientů
M	placebo	1
M	lék 1	1
M	lék 2	6
Z	placebo	3
Z	lék 1	4
Z	lék 2	9

**V softwaru STATISTICA:** Statistics – ANOVA – Main effects ANOVA – Quick specs dialog – OK – Variables – Dependent variable list: X, Categorical predictors (factors): A, B – OK – All effects.

*Post hoc testy:* More results – Post hoc – zvolit Effect – Tukey HSD (nebo Scheffé)

*Levenův test:* More results – Assumptions – zvolit proměnnou – Levene's test (ANOVA)

*Vykreslení krabicových grafů podle obou proměnných:* Graphs – 2D Graphs – Box Plots... – zvolit spojitou proměnnou jako Dependent variable, zvolit jednu kategoriální proměnnou jako Grouping variable – na listu Categorized u X-Categories zatrhnout On a Layout změnit na Overlaid – OK

*Pokud bychom uvažovali model s interakcemi, zvolíme Factorial ANOVA (namísto Main effects A.)*

# Úkol 2 – řešení v softwaru SPSS

Zjistěte, zda má vliv pohlaví a typ léku na počet uzdravených pacientů s leukémií.

Pohlaví	Typ léku	Počet uzdrav. pacientů
M	placebo	1
M	lék 1	1
M	lék 2	6
Z	placebo	3
Z	lék 1	4
Z	lék 2	9

**V softwaru SPSS:** Analyze – General Linear Model – Univariate – Dependent Variable: spojitá proměnná, Fixed Factor(s): kategoriální proměnné →

- Model – zatrhneme Custom – vybereme Typ:Main effects – do Model přetáhneme A, B (*pokud bychom chtěli model s interakcemi necháme zatržené Full factorial*) – odškrtneme Include intercept in model – Continue
- Post Hoc – Post hoc Tests for: zvolit kategoriální proměnnou – zatrhneme Tukey's-b – Continue
- Plots: zvolit proměnné do Horizontal Axis a Separte Lines – Add – Continue
- Options... – Homogeneity tests – Continue

*Vykreslení krabicových grafů podle obou proměnných:* Graphs – Legacy Dialogs – Boxplot... – Clustered – Define – zvolit Variable Category Axis a Define Clusters by - OK

# Úkol 2 – řešení v softwaru R

Zjistěte, zda má vliv pohlaví a typ léku na počet uzdravených pacientů s leukémií.

## V softwaru R:

```
data <- data.frame(pohl=c(1,1,1,2,2,2),lek=c(1,2,3,1,2,3),pocet=c(1,1,6,3,4,9))
data
```

```
model_bez_interakce <- aov(data$pocet ~ (as.factor(data$poohl)+as.factor(data$lek)))
summary(model_bez_interakce)
TukeyHSD(model_bez_interakce) # post-hoc test
```

```
# 2. způsob: anova(lm(data$pocet ~ (as.factor(data$poohl)+as.factor(data$lek))))
```

```
model_s_interakci <- aov(data$pocet ~ (as.factor(data$poohl)*as.factor(data$lek)))
summary(model_s_interakci)
```

```
boxplot(data$pocet ~(as.factor(data$poohl)*as.factor(data$lek)))
```

```
library("car") # instalace balíku car pomocí: install.packages("car")
leveneTest(data$pocet ~ (as.factor(data$poohl)*as.factor(data$lek)),center=mean)
```



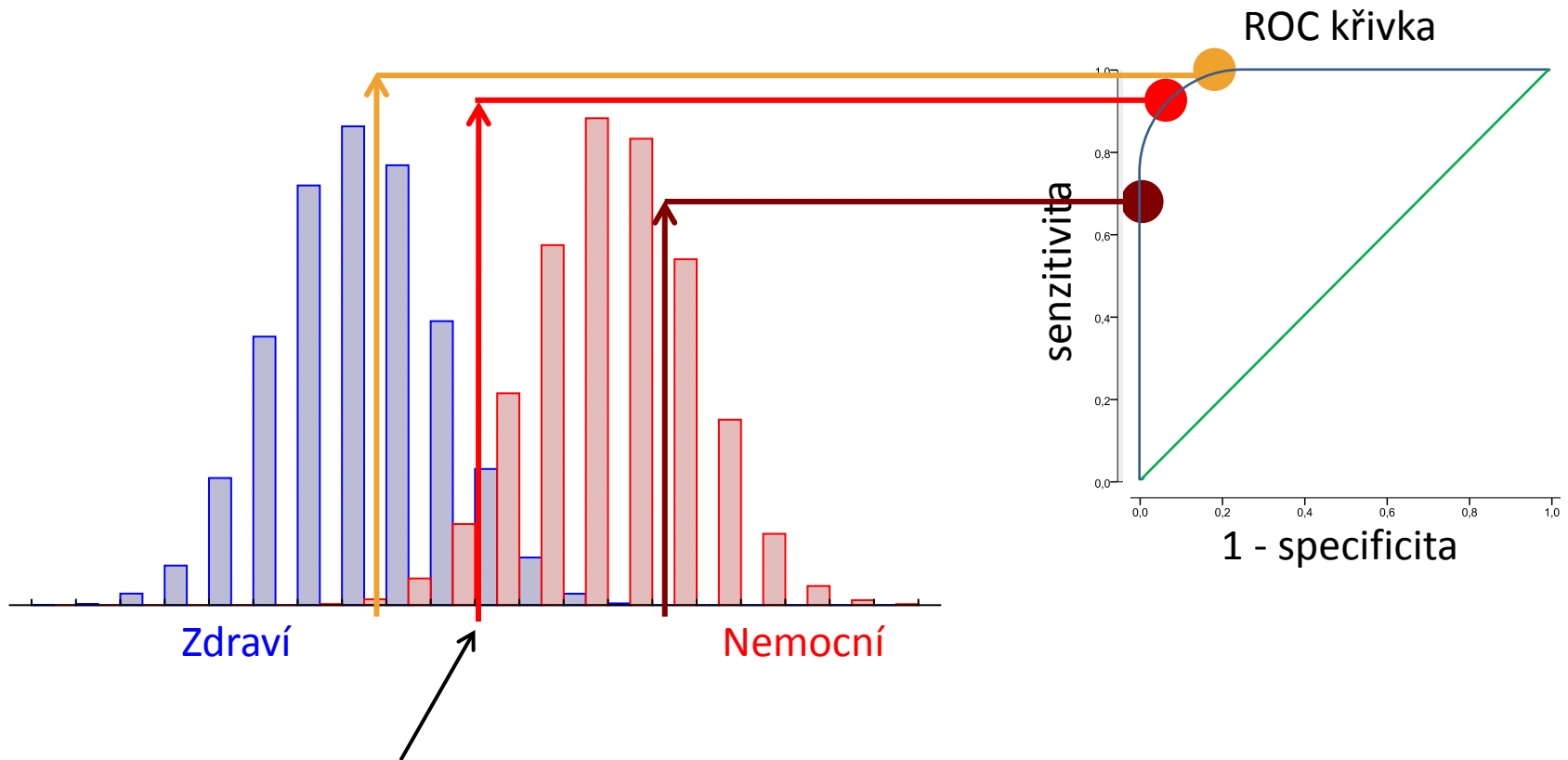
# Hledání diagnostického cut-off pomocí ROC křivek

# ROC analýza – motivace

- Dříve probrané ukazatele diagnostické síly testů (senzitivita, specificita apod.) **nelze použít u diagnostických testů, jejichž výstupem je spojitá (kvantitativní) proměnná** (např. koncentrace analytu v krevním séru, systolický krevní tlak).
- Na základě předchozích výzkumů známe dělicí body, které odlišují normální a patologické hodnoty spojitě proměnné, pomocí nichž můžeme spojitou proměnnou binarizovat – tzn. vytvoření dvou kategorií „pozitivní“ / „negativní“ (např. „pod normou“ / „v normě“).
- Pokud dělicí body nejsou známy předem, můžeme se je snažit nalézt pomocí **ROC („Receiver Operating Characteristic“) křivky**.
- **Cíle ROC analýzy:**
  1. Určit, zda je spojitá proměnná vhodná pro diagnostické odlišování zdravých a nemocných jedinců.
  2. Nalezení dělicího bodu („cut-off point“) na škále hodnot spojitě proměnné, který nejlépe odlišuje zdravé a nemocné jedince.

# ROC analýza

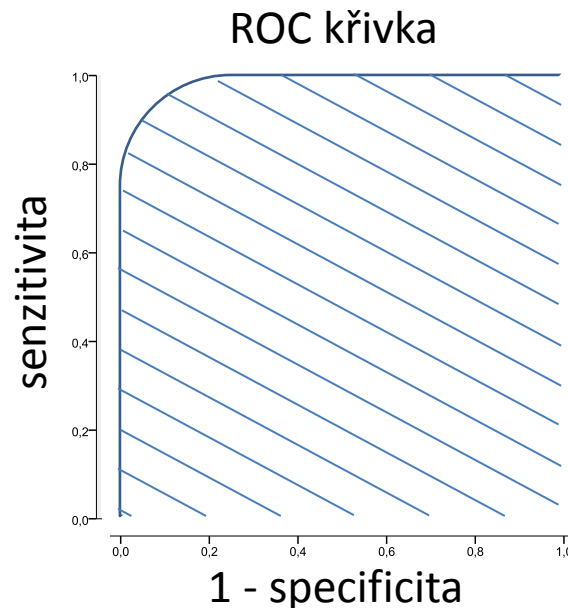
- Princip: Jakákoli hodnota spojité proměnné nějak rozlišuje zdravé a nemocné jedince, tzn. je spojena s nějakou senzitivitou a specificitou.



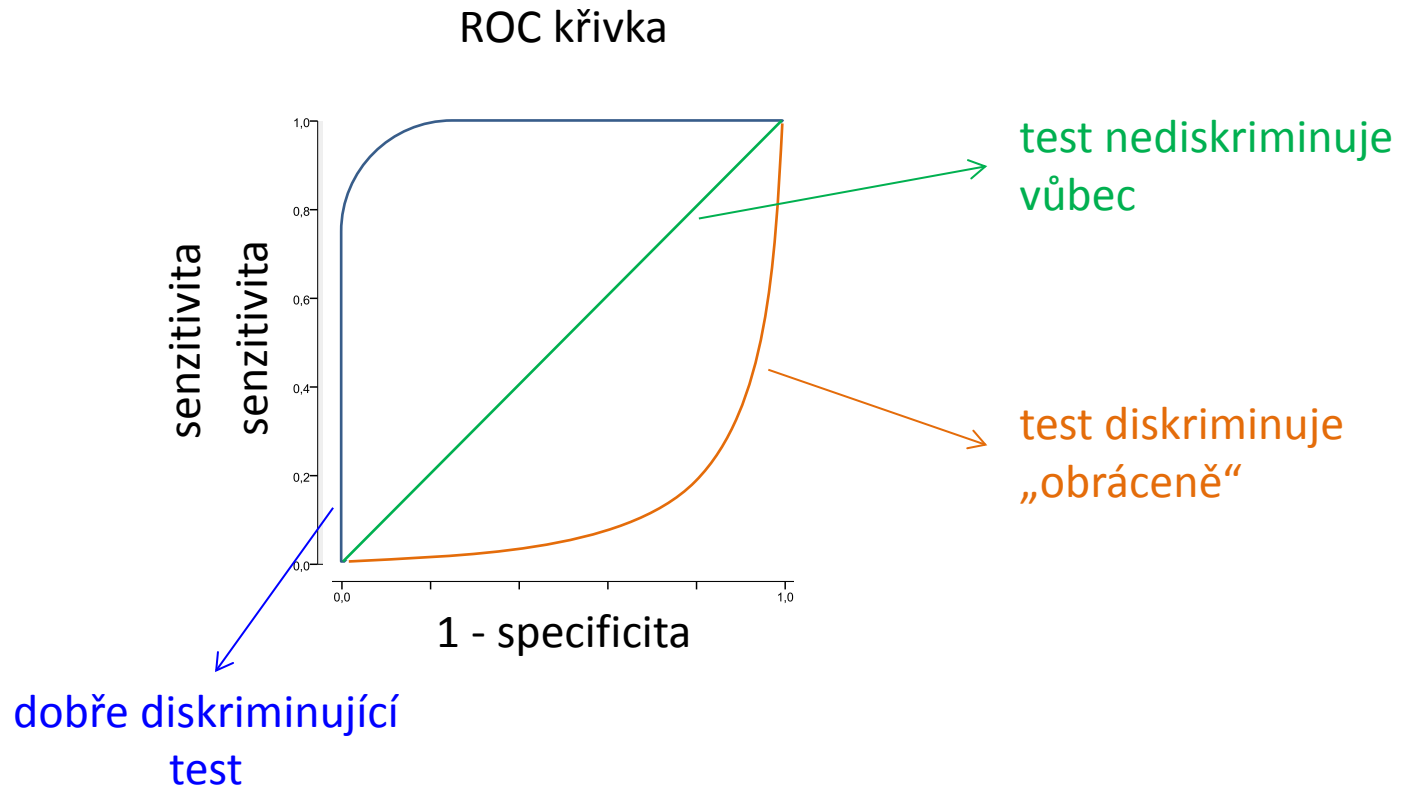
Nejlepší dělicí bod („cut-off“) – nejvyšší senzitivita a specificita pro odlišení skupin – tzn. maximální součet hodnot senzitivity a specificity.

# ROC analýza – plocha pod ROC křivkou

- Plocha pod ROC křivkou = „Area Under the Curve“ (AUC).
- Nabývá hodnot od 0 do 1.
- Slouží k vyjádření diagnostické síly (efektivity) testu.
- Čím větší hodnota AUC, tím lepší diagnostický test je (hodnota AUC nad 0,75 většinou poukazuje na uspokojivou diskriminační schopnost testu).

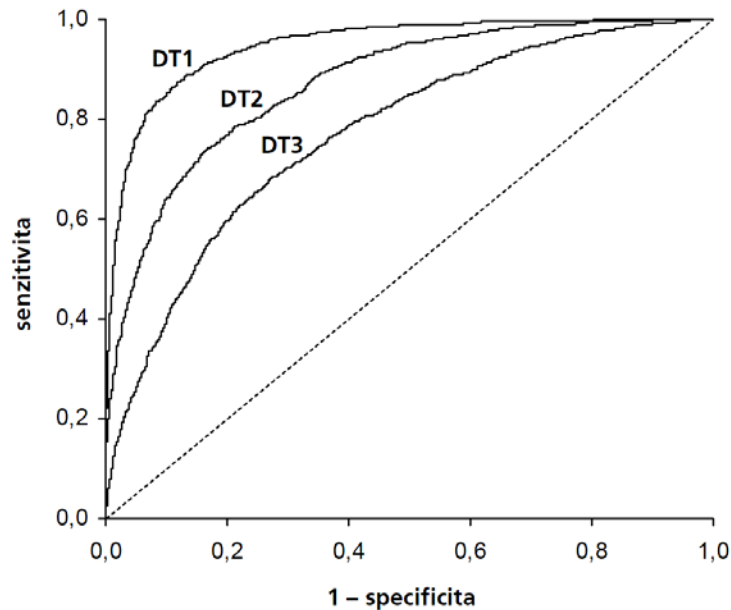


# ROC analýza – srovnání diagnostické síly různých testů



# ROC analýza – srovnání diagnostické síly různých testů

- Lze srovnat i velmi rozdílné testy (např. testy založené na různých proměnných).



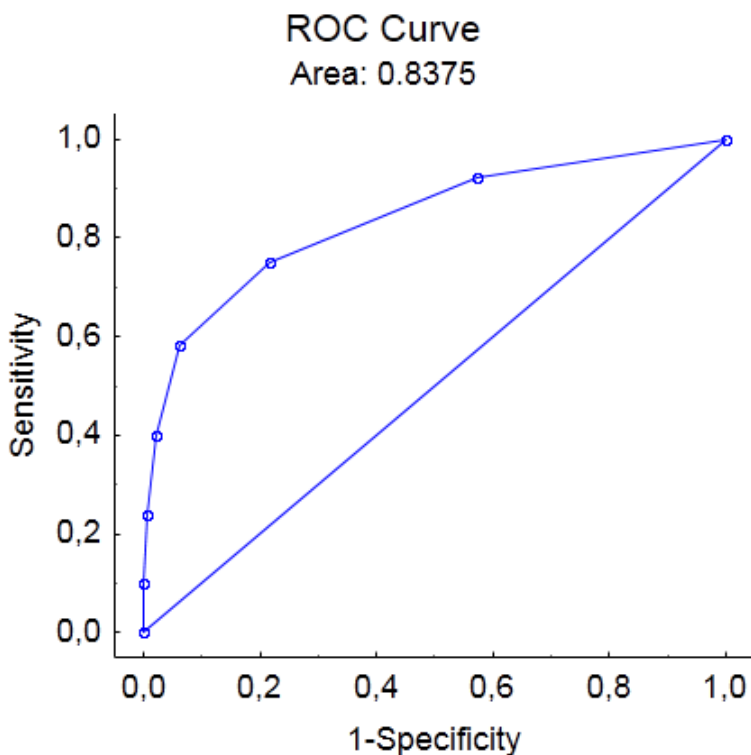
Diagnostický test	AUC
DT1	0,949
DT2	0,872
DT3	0,770

→ nejlepší

→ nejhorší

# ROC analýza

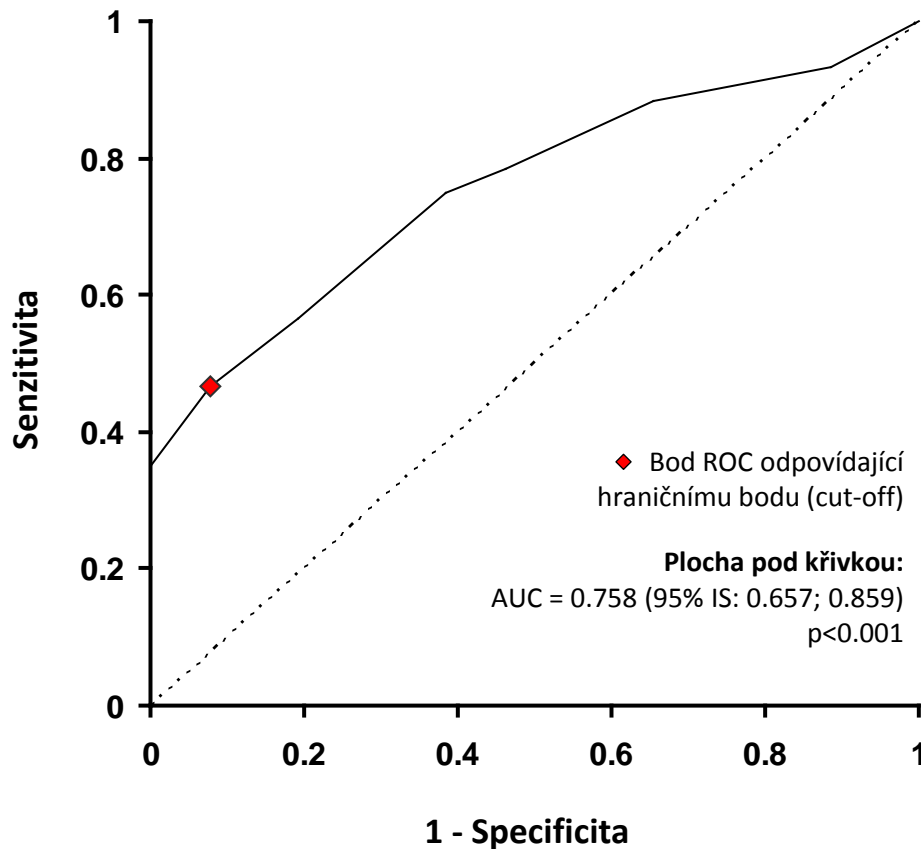
**Příklad:** Zjistěte, zda je MMSE skóre vhodné na diagnostiku mírné kognitivní poruchy (MCI). Najděte dělicí bod (cut-off), který nejlépe odlišuje pacienty s MCI od kontrolních subjektů.



MMSE skóre	Sensitivity	1-Specificity	Specificity	Sensitivity + Specificity
-23	0,002	0,000	1,000	1,002
-24	0,101	0,000	1,000	1,101
-25	0,239	0,004	0,996	1,235
-26	0,399	0,022	0,978	1,377
-27	0,581	0,061	0,939	1,520
<b>-28</b>	<b>0,749</b>	<b>0,217</b>	<b>0,783</b>	<b>1,531</b>
-29	0,924	0,574	0,426	1,350
-30	1,000	1,000	0,000	1,000

# Hledání cut-off – doplnění

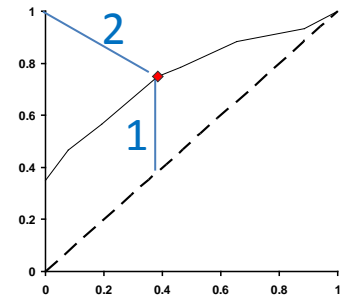
Příklad:



Sens	Spec	Sens+Spec
1.000	0.000	1.000
0.933	0.115	1.049
0.883	0.346	1.229
0.783	0.538	1.322
0.750	0.615	1.365
0.567	0.808	1.374
<b>0.467</b>	<b>0.923</b>	<b>1.390</b>
0.350	1.000	1.350
0.217	1.000	1.217
0.150	1.000	1.150
0.050	1.000	1.050
0.033	1.000	1.033
0.000	1.000	1.000



# Hledání cut-off – kritéria



Kritérium	Vzoreček	Reference
<b>1. Youdenova J statistika</b> <sup>1</sup> – maximalizace vzdálenosti od diagonály	$\max(se + sp)$	<ul style="list-style-type: none"> <li>• W. J. Youden (1950) “Index for rating diagnostic tests”. Cancer, 3, 32–35.</li> <li>• R-kový balík pROC</li> <li>• <a href="http://www.medicalbiostatistics.com/roccurve.pdf">http://www.medicalbiostatistics.com/roccurve.pdf</a></li> </ul>
<b>2. Nejbližší bod levému hornímu rohu grafu</b>	$\min((1 - se)^2 + (1 - sp)^2)$	<ul style="list-style-type: none"> <li>• R-kový balík pROC</li> <li>• <a href="http://www.medicalbiostatistics.com/roccurve.pdf">http://www.medicalbiostatistics.com/roccurve.pdf</a></li> </ul>
<b>3. Maximalizace součinu senzitivity a specificity</b>	$\max(se * sp)$	<ul style="list-style-type: none"> <li>• R-kový balík OptimalCutpoints</li> <li>• dr. Budíková používá maximalizaci geometrického průměru sens a spec</li> </ul>

<sup>1</sup> Youdenova J statistika je definována jako:  $J = se + sp - 1$ ; při hledání maxima lze ale člen (-1) zanedbat

# Hledání cut-off – vážená kritéria (dle R balíku pROC)

Kritérium	Vzoreček
<b>Youdenova J statistika</b> <sup>1</sup> – maximalizace vzdálenosti od diagonály	$\max(se + r * sp)$
Nejbližší bod levému hornímu rohu grafu	$\min((1 - se)^2 + r * (1 - sp)^2)$

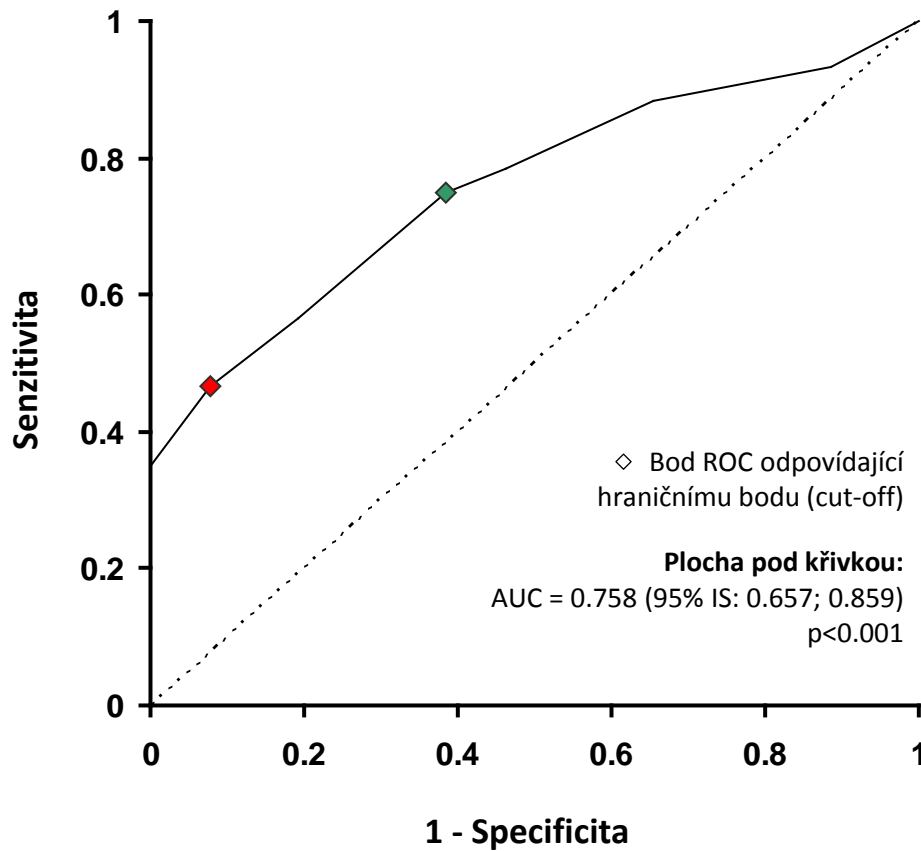
kde: 
$$r = \frac{1 - prevalence}{cost * prevalence}$$

$$prevalence = \frac{n_{cases}}{n_{cases} + n_{controls}}$$

*cost* – penalizace falešně negativních výsledků

defaultně: *prevalence* = 0,5 a *cost* = 1

# Příklad - pokračování



Sens	Spec	Sens+ Spec	closest. topleft	Sens* Spec
1.000	0.000	1.000	1.000	0.000
0.933	0.115	1.049	0.787	0.108
0.883	0.346	1.229	0.441	0.306
0.783	0.538	1.322	0.260	0.422
0.750	0.615	1.365	0.210	0.462
0.567	0.808	1.374	0.225	0.458
0.467	0.923	1.390	0.290	0.431
0.350	1.000	1.350	0.423	0.350
0.217	1.000	1.217	0.614	0.217
0.150	1.000	1.150	0.723	0.150
0.050	1.000	1.050	0.903	0.050
0.033	1.000	1.033	0.934	0.033
0.000	1.000	1.000	1.000	0.000