

MA012 Statistika II

2. Dvoutfaktorová analýza rozptylu (Two-Way ANOVA)

Ondřej Pokora (pokora@math.muni.cz)

Ústav matematiky a statistiky, Přírodovědecká fakulta, Masarykova univerzita, Brno

(podzim 2015)



Motivační příklad

Příklad 1

Zkoumají se výnosy sena v tunách na hektar v závislosti na typu půdy a na způsobu hnojení. Každá kombinace byla realizována čtyřikrát, nezávisle na sobě.

[t/ha]	způsob hnojení (B):		
typ půdy (A)	bez hnojení	chlévká mrva	vápenaté hnojivo
normální	2,8; 3,2; 3,0; 3,0	3,7; 3,6; 3,9; 3,6	3,4; 3,8; 3,7; 3,6
kyselá	3,1; 2,7; 3,0; 2,9	3,4; 3,4; 3,0; 3,8	4,2; 4,0; 4,1; 3,9

Na hladině významnosti 0,05 testujte hypotézy

- Typ půdy nemá vliv na výnosy.
- Způsob hnojení nemá vliv na výnosy.
- Typ půdy a způsob hnojení jsou nezávislé, tj. neinteragují.

V některých reálných situacích se celkový soubor sledovaných náhodných veličin rozpadá na dílčí výběry takovým způsobem, že přihlížíme ke dvěma faktorům (třídícím znakům).

Podobně jako v jednofaktorové analýze rozptylu (One-Way ANOVA), se zajímáme o statistické posouzení toho, zda lze některým z faktorů, či oběma faktory, vysvětlit variabilitu pozorovaných hodnot.

Princip dvoufaktorové analýzy rozptylu (Two-Way ANOVA) je analogický jednofaktorové variantě. Konstruuje se řetězec submodelů a postupně se porovnávají rozptyly, které jednotlivým modelům odpovídají.

Předpoklady

- Uvažujeme dva faktory, A a B .
- Faktor A má $a \geq 2$ úrovní, faktor B má $b \geq 2$ úrovní.
- Pro každou kombinaci úrovní obou faktorů, tzn. pro $(A = i, B = j)$, máme n_{ij} výsledků $(Y_{ij1}, \dots, Y_{ijn_{ij}})$, které tvoří náhodný výběr z rozložení $N(\mu_{ij}, \sigma^2)$, $i = 1, \dots, a, j = 1, \dots, b$.
- V označení Y_{ijk} první index označuje skupinu podle úrovně faktoru A , druhý index označuje skupinu podle úrovně faktoru B , třetí index značí pořadí měření v dané skupině.
- Jednotlivé náhodné výběry jsou stochasticky nezávislé, tedy

$$Y_{ijk} = \mu_{ij} + \varepsilon_{ijk},$$

kde ε_{ijk} jsou stochasticky nezávislé náhodné veličiny s rozložením $N(0, \sigma^2)$, kde $i = 1, \dots, a, j = 1, \dots, b$ a $k = 1, \dots, n_{ij}$. Rozptyl náhodných chyb σ^2 přitom není známý.

Dvojné třídění

		faktor B		
		1	...	b
faktor A	1	$(Y_{111}, \dots, Y_{11n_{11}})$...	$(Y_{1b1}, \dots, Y_{1bn_{1b}})$
	\vdots	\vdots	$(Y_{ij1}, \dots, Y_{ijn_{ij}})$	\vdots
	a	$(Y_{a11}, \dots, Y_{a1n_{a1}})$...	$(Y_{ab1}, \dots, Y_{abn_{ab}})$

Příklad 1

[t/ha]	způsob hnojení (B):		
typ půdy (A)	Bez hnojení	chlévká Mrva	Vápenaté hnojivo
Normální	2,8; 3,2; 3,0; 3,0	3,7; 3,6; 3,9; 3,6	3,4; 3,8; 3,7; 3,6
Kyselá	3,1; 2,7; 3,0; 2,9	3,4; 3,4; 3,0; 3,8	4,2; 4,0; 4,1; 3,9

Součty a rozsahy výběrů

součty ve skupinách: $Y_{ij.} = \sum_{k=1}^{n_{ij}} Y_{ijk}$

součty v řádcích: $Y_{i..} = \sum_{j=1}^b \sum_{k=1}^{n_{ij}} Y_{ijk}$, součty ve sloupcích: $Y_{.j.} = \sum_{i=1}^a \sum_{k=1}^{n_{ij}} Y_{ijk}$

celkový součet: $Y_{...} = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} Y_{ijk}$

počet měření pro $A = i$: $n_{i.} = \sum_{j=1}^b n_{ij}$, počet měření pro $B = j$: $n_{.j} = \sum_{i=1}^a n_{ij}$

celkový rozsah souboru: $n = \sum_{i=1}^a \sum_{j=1}^b n_{ij} = \sum_{i=1}^a n_{i.} = \sum_{j=1}^b n_{.j}$

Průměry výběrů

průměry ve skupinách: $\bar{Y}_{ij\cdot} = \frac{1}{n_{ij}} \sum_{k=1}^{n_{ij}} Y_{ijk}$

v řádcích: $\bar{Y}_{i..} = \frac{1}{n_{i\cdot}} \sum_{j=1}^b \sum_{k=1}^{n_{ij}} Y_{ijk}$ ve sloupcích: $\bar{Y}_{\cdot j\cdot} = \frac{1}{n_{\cdot j}} \sum_{i=1}^a \sum_{k=1}^{n_{ij}} Y_{ijk}$

celkový průměr: $\bar{Y}_{\dots} = \frac{1}{n} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} Y_{ijk}$

Počty a průměry ve skupinách

		faktor B					v řádcích
		1	...	j	...	b	
faktor A	1	n_{11}	...	n_{1j}	...	n_{1b}	$n_{1.}$
	⋮	⋮		⋮		⋮	⋮
	i	n_{i1}	...	n_{ij}	...	n_{ib}	$n_{i.}$
	⋮	⋮		⋮		⋮	⋮
	a	n_{a1}	...	n_{aj}	...	n_{ab}	$n_{a.}$
ve sloupcích		$n_{.1}$...	$n_{.j}$...	$n_{.b}$	n

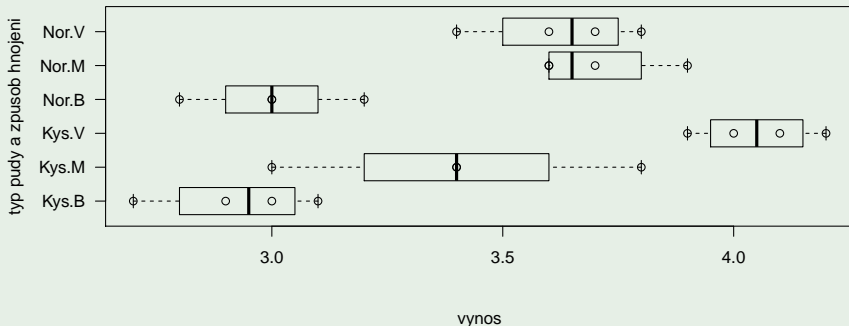
		faktor B					v řádcích
		1	...	j	...	b	
faktor A	1	\bar{Y}_{11}	...	\bar{Y}_{1j}	...	\bar{Y}_{1b}	$\bar{Y}_{1.}$
	⋮	⋮		⋮		⋮	⋮
	i	\bar{Y}_{i1}	...	\bar{Y}_{ij}	...	\bar{Y}_{ib}	$\bar{Y}_{i.}$
	⋮	⋮		⋮		⋮	⋮
	a	\bar{Y}_{a1}	...	\bar{Y}_{aj}	...	\bar{Y}_{ab}	$\bar{Y}_{a.}$
ve sloupcích		$\bar{Y}_{.1}$...	$\bar{Y}_{.j}$...	$\bar{Y}_{.b}$	$\bar{Y}_{...}$

Průměry a boxploty

Příklad 1

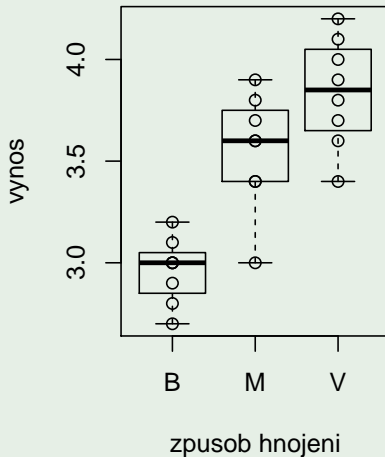
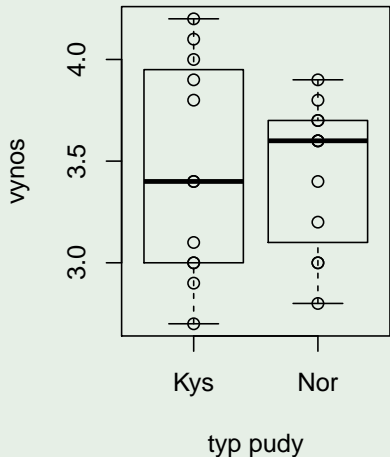
\bar{Y}_{ij}	způsob hnojení (B):			průměry
typ půdy (A)	Bez hnojení	chlévká M rva	Vápenaté hnojivo	$\bar{Y}_{i.}$
Normální	3,00	3,70	3,63	3,44
Kyselá	2,93	3,40	4,05	3,46
průměry $\bar{Y}_{.j}$	2,96	3,55	3,84	$\bar{Y}_{..} = 3,45$

$$n_{ij} = p = 4, \quad n_{i.} = 12, \quad n_{.j} = 8, \quad n = 24$$



Boxploty pro řádky a sloupce

Příklad 1



Hypotézy

Stanovíme hypotézy, které na hladině významnosti α chceme testovat.

H_{A0} : všechny střední hodnoty v řádcích jsou stejné, tzn. faktor A nemá vliv
oproti alternativní hypotéze

H_{A1} : některá(é) dvojice středních hodnot v řádcích se liší, tzn. faktor A má vliv

H_{B0} : všechny střední hodnoty ve sloupcích jsou stejné, tzn. faktor B nemá vliv
oproti alternativní hypotéze

H_{B1} : některá(é) dvojice středních hodnot ve sloupcích se liší, tzn. faktor B má vliv

Základní model M

Definice 1 (model M)

Náhodné veličiny Y_{ijk} se řídí modelem M , pokud

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}$$

pro $i = 1, \dots, a$, $j = 1, \dots, b$ a $k = 1, \dots, n_{ij}$, přičemž ε_{ijk} jsou stochasticky nezávislé náhodné veličiny s rozložením $\varepsilon_{ijk} \sim N(0, \sigma^2)$.

Interpretace neznámých parametrů:

- μ je společná část střední hodnoty sledované veličiny
- $\alpha_1, \dots, \alpha_a$ jsou efekty faktoru A , odchylky od μ způsobené vlivem A
- β_1, \dots, β_b jsou efekty faktoru B , odchylky od μ způsobené vlivem B
- σ^2 je rozptyl náhodných chyb

M jako lineární regresní model

$$\underbrace{\begin{pmatrix} Y_{111} \\ \vdots \\ Y_{abn_{ab}} \end{pmatrix}}_Y = \underbrace{\begin{pmatrix} \mathbf{1}_{n_{11}} & \mathbf{1}_{n_{11}} & \mathbf{0}_{n_{11}} & \cdots & \mathbf{0}_{n_{11}} & \mathbf{1}_{n_{11}} & \mathbf{0}_{n_{11}} & \cdots & \mathbf{0}_{n_{11}} \\ \mathbf{1}_{n_{12}} & \mathbf{1}_{n_{12}} & \mathbf{0}_{n_{12}} & \cdots & \mathbf{0}_{n_{12}} & \mathbf{0}_{n_{12}} & \mathbf{1}_{n_{12}} & \cdots & \mathbf{0}_{n_{12}} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_{n_{1b}} & \mathbf{1}_{n_{1b}} & \mathbf{0}_{n_{1b}} & \cdots & \mathbf{0}_{n_{1b}} & \mathbf{0}_{n_{1b}} & \mathbf{0}_{n_{1b}} & \cdots & \mathbf{1}_{n_{1b}} \\ \hline \mathbf{1}_{n_{21}} & \mathbf{0}_{n_{21}} & \mathbf{1}_{n_{21}} & \cdots & \mathbf{0}_{n_{21}} & \mathbf{1}_{n_{21}} & \mathbf{0}_{n_{21}} & \cdots & \mathbf{0}_{n_{21}} \\ \mathbf{1}_{n_{22}} & \mathbf{0}_{n_{22}} & \mathbf{1}_{n_{22}} & \cdots & \mathbf{0}_{n_{22}} & \mathbf{0}_{n_{22}} & \mathbf{1}_{n_{22}} & \cdots & \mathbf{0}_{n_{22}} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_{n_{2b}} & \mathbf{0}_{n_{2b}} & \mathbf{1}_{n_{2b}} & \cdots & \mathbf{0}_{n_{2b}} & \mathbf{0}_{n_{2b}} & \mathbf{0}_{n_{2b}} & \cdots & \mathbf{1}_{n_{2b}} \\ \hline \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \hline \mathbf{1}_{n_{a1}} & \mathbf{0}_{n_{a1}} & \mathbf{0}_{n_{a1}} & \cdots & \mathbf{1}_{n_{a1}} & \mathbf{1}_{n_{a1}} & \mathbf{0}_{n_{a1}} & \cdots & \mathbf{0}_{n_{a1}} \\ \mathbf{1}_{n_{a2}} & \mathbf{0}_{n_{a2}} & \mathbf{0}_{n_{a2}} & \cdots & \mathbf{1}_{n_{a2}} & \mathbf{0}_{n_{a2}} & \mathbf{1}_{n_{a2}} & \cdots & \mathbf{0}_{n_{a2}} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_{n_{ab}} & \mathbf{0}_{n_{ab}} & \mathbf{0}_{n_{ab}} & \cdots & \mathbf{1}_{n_{ab}} & \mathbf{0}_{n_{ab}} & \mathbf{0}_{n_{ab}} & \cdots & \mathbf{1}_{n_{ab}} \end{pmatrix}}_X \underbrace{\begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_a \\ \beta_1 \\ \vdots \\ \beta_b \end{pmatrix}}_\theta$$

$$Y = X\theta$$

Model M

Jaké rozměry a jakou hodnotu má matice plánu X ?

Model M

Jaké rozměry a jakou hodnotu má matice plánu \mathbf{X} ?

$$\mathbf{X} : n \times (1 + a + b), \quad h(\mathbf{X}) = a + b - 1$$

Matice plánu tedy není plně hodnosti.

Odhad vektoru parametrů $\boldsymbol{\theta}$ se proto počítá pomocí pseudoinverzní matice

$$\boldsymbol{\theta} = (\mathbf{X}'\mathbf{X})^{-} \mathbf{X}'\mathbf{Y},$$

anebo se přidávají tzv. reparametrizační rovnice (proč dvě a co vyjadřují?)

$$\sum_{i=1}^a \alpha_i = 0, \quad \sum_{j=1}^b \beta_j = 0.$$

Odhady parametrů v modelu M

Věta 2 (odhady parametrů modelu M)

Odhady parametrů modelu

$$M : Y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}$$

metodou nejmenších čtverců jsou rovny

$$\mu = \bar{Y}_{...},$$

$$\alpha_i = \bar{Y}_{i..} - \bar{Y}_{...},$$

$$\beta_j = \bar{Y}_{.j.} - \bar{Y}_{...}.$$

Odhad střední hodnoty měření ve skupině ($A = i, B = j$) modelu M je rovný

$$\hat{\mu}_{ijk} = \bar{Y}_{i..} + \bar{Y}_{.j.} - \bar{Y}_{...}.$$

Tvrzení se dokazuje přímým výpočtem maticové rovnice pro odhad parametrů v lineárním regresním modelu.

Porovnejte s odvozením v přednášce k jednofaktorové analýze rozptylu.

Model M_B

Při zkoumání vlivu faktoru B testujeme v modelu M hypotézu

$$H_{B0} : \beta_1 = \dots = \beta_b = 0,$$

$$H_{B1} : \exists i \in \{1, \dots, b\} : \beta_i \neq 0.$$

Definice 3 (model M_B)

Náhodné veličiny Y_{ijk} se řídí modelem M_B , pokud

$$Y_{ijk} = \mu + \alpha_i + \varepsilon_{ijk}$$

pro $i = 1, \dots, a$, $j = 1, \dots, b$ a $k = 1, \dots, n_{ij}$, přičemž ε_{ijk} jsou stochasticky nezávislé náhodné veličiny s rozložením $\varepsilon_{ijk} \sim N(0, \sigma^2)$.

Interpretace neznámých parametrů:

- μ je společná část střední hodnoty sledované veličiny
- $\alpha_1, \dots, \alpha_a$ jsou efekty faktoru A , odchylky od μ způsobené vlivem A
- σ^2 je rozptyl náhodných chyb

M_B jako lineární regresní model

$$\underbrace{\begin{pmatrix} Y_{111} \\ \vdots \\ Y_{abn_{ab}} \end{pmatrix}}_Y = X_B \underbrace{\begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_a \end{pmatrix}}_{\theta_B} \quad \text{tj. } Y = X_B \theta_B$$

Jak vypadá matice plánu X_B odpovídající modelu M_B ?
Jaké má rozměry a hodnoty?

M_B jako lineární regresní model

$$\underbrace{\begin{pmatrix} Y_{111} \\ \vdots \\ Y_{abn_{ab}} \end{pmatrix}}_{\mathbf{Y}} = \mathbf{X}_B \underbrace{\begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_a \end{pmatrix}}_{\boldsymbol{\theta}_B} \quad \text{tj. } \mathbf{Y} = \mathbf{X}_B \boldsymbol{\theta}_B$$

Jak vypadá matice plánu \mathbf{X}_B odpovídající modelu M_B ?
Jaké má rozměry a hodnoty?

\mathbf{X}_B získáme vynecháním posledních b sloupců z matice \mathbf{X} .

$$\mathbf{X}_B : n \times (1 + a), \quad h(\mathbf{X}_B) = a$$

Matice plánu tedy opět není plné hodnosti.

Odhad vektoru parametrů $\boldsymbol{\theta}_B$ se proto počítá pomocí pseudoinverzní matice

$$\boldsymbol{\theta}_B = (\mathbf{X}_B' \mathbf{X}_B)^- \mathbf{X}_B' \mathbf{Y},$$

anebo se přidá reparametrizační rovnice $\sum_{i=1}^a \alpha_i = 0$.

Odhady parametrů v modelu M_B

Řešením maticové rovnice pro lineární regresní model M_B lze spočítat:

Věta 4 (odhady parametrů modelu M_B)

Odhady parametrů modelu

$$M_B: Y_{ijk} = \mu + \alpha_i + \varepsilon_{ijk}$$

metodou nejmenších čtverců jsou rovny

$$\mu = \bar{Y}_{\dots},$$

$$\alpha_i = \bar{Y}_{i..} - \bar{Y}_{\dots}.$$

Odhad střední hodnoty měření ve skupině ($A = i, B = j$) modelu M_B je rovný

$$\hat{\mu}_{B,ijk} = \bar{Y}_{i..}$$

Model M_A

Při zkoumání vlivu faktoru A testujeme v modelu M hypotézu

$$H_{A0} : \alpha_1 = \dots = \alpha_a = 0,$$

$$H_{A1} : \exists i \in \{1, \dots, a\} : \alpha_i \neq 0.$$

Definice 5 (model M_A)

Náhodné veličiny Y_{ijk} se řídí modelem M_A , pokud

$$Y_{ijk} = \mu + \beta_j + \varepsilon_{ijk}$$

pro $i = 1, \dots, a$, $j = 1, \dots, b$ a $k = 1, \dots, n_{ij}$, přičemž ε_{ijk} jsou stochasticky nezávislé náhodné veličiny s rozložením $\varepsilon_{ijk} \sim N(0, \sigma^2)$.

Interpretace neznámých parametrů:

- μ je společná část střední hodnoty sledované veličiny
- β_1, \dots, β_b jsou efekty faktoru B , odchylky od μ způsobené vlivem B
- σ^2 je rozptyl náhodných chyb

M_A jako lineární regresní model, odhady parametrů

Jak vypadá matice plánu X_A pro model M_A ? Jaké má rozměry a hodnoty?

M_A jako lineární regresní model, odhady parametrů

Jak vypadá matice plánu X_A pro model M_A ? Jaké má rozměry a hodnota?

X_A získáme vynecháním 2. až $(a+1)$ -ního sloupce z matice X .

$$X_A : n \times (1 + b), \quad h(X_A) = b$$

Věta 6 (odhady parametrů modelu M_A)

Odhady parametrů modelu

$$M_A : Y_{ijk} = \mu + \beta_j + \varepsilon_{ijk}$$

metodou nejmenších čtverců jsou rovny

$$\mu = \bar{Y}_{\dots},$$

$$\beta_j = \bar{Y}_{.j.} - \bar{Y}_{\dots}.$$

Odhad střední hodnoty měření ve skupině ($A = i, B = j$) modelu M_A je rovný

$$\hat{\mu}_{A,ijk} = \bar{Y}_{.j.}$$

Minimální model M_0

Při zkoumání vlivu obou faktorů vyjdeme testujeme v modelu M_B hypotézu

$$H_{A0} : \alpha_1 = \dots = \alpha_a = 0,$$

$$H_{A1} : \exists i \in \{1, \dots, a\} : \alpha_i \neq 0.$$

Definice 7 (model M_0)

Náhodné veličiny Y_{ijk} se řídí modelem M_0 , pokud

$$Y_{ijk} = \mu + \varepsilon_{ijk}$$

pro $i = 1, \dots, a$, $j = 1, \dots, b$ a $k = 1, \dots, n_{ij}$, přičemž ε_{ijk} jsou stochasticky nezávislé náhodné veličiny s rozložením $\varepsilon_{ijk} \sim N(0, \sigma^2)$.

Interpretace neznámých parametrů:

- μ je střední hodnota sledované veličiny bez ohledu na kategorizaci dle faktorů
- σ^2 je rozptyl náhodných chyb

M_0 jako lineární regresní model, odhad μ

$$Y = X_0 \mu$$

Jak vypadá matice plánu X_0 pro model M_0 ? Jaké má rozměry a hodnoty?

M_0 jako lineární regresní model, odhad μ

$$Y = X_0 \mu$$

Jak vypadá matice plánu X_0 pro model M_0 ? Jaké má rozměry a hodnost?

$$X_0 = \mathbf{1}_n, \quad X_A : n \times 1, \quad h(X_A) = 1$$

Matice je v tomto modelu plné hodnosti a odhad parametru μ počítáme klasicky:

Věta 8 (odhady parametrů modelu M_0)

V modelu

$$M_0 : Y_{ijk} = \mu + \varepsilon_{ijk}$$

je odhad parametru μ metodou nejmenších čtverců roven

$$\hat{\mu}_{A,ijk} = (X_0' X_0)^{-1} X_0' Y = \bar{Y} \dots$$

Shrnutí: modely a submodely

$$M : Y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}$$

$$M_B : Y_{ijk} = \mu + \alpha_i + \varepsilon_{ijk}$$

$$M_A : Y_{ijk} = \mu + \beta_j + \varepsilon_{ijk}$$

$$M_0 : Y_{ijk} = \mu + \varepsilon_{ijk}$$

Jaké řetězce submodelů lze vytvořit?

Shrnutí: modely a submodely

$$M : Y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}$$

$$M_B : Y_{ijk} = \mu + \alpha_i + \varepsilon_{ijk}$$

$$M_A : Y_{ijk} = \mu + \beta_j + \varepsilon_{ijk}$$

$$M_0 : Y_{ijk} = \mu + \varepsilon_{ijk}$$

Jaké řetězce submodelů lze vytvořit?

$$M \xrightarrow{H_{B0}} M_B \xrightarrow{H_{A0}} M_0, \quad M \xrightarrow{H_{A0}} M_A \xrightarrow{H_{B0}} M_0.$$

Pro postupné testování hypotéz v analýze rozptylu se volí vždy jeden z řetězců submodelů, v praxi obvykle $M \rightarrow M_B \rightarrow M_0$.

V praxi volí většina statistických softwarů (mj. i R) odlišnou stavbu modelů s

$$\alpha_1 = 0, \quad \beta_1 = 0,$$

kdy se neuvažují reparametrizační rovnice. Tato odlišná parametrizace nemá vliv na vlastní výpočet, pouze parametry modelů mají odlišnou interpretaci.

První úroveň faktorů A a B jsou tedy stanoveny jako referenční. Parametr μ zde interpretujeme jako střední hodnotu kategorie ($A = 1, B = 1$), a efekty $\alpha_2, \dots, \alpha_a$, resp. β_2, \dots, β_b , vyjadřují odchylky vlivem faktoru A , resp. B , od této střední hodnoty.

Definice 9 (vyvážené třídění)

O vyváženém třídění hovoříme, pokud je počet ve všech kategoriích stejný,

$$n_{ij} = p, \quad i = 1, \dots, a, j = 1, \dots, b.$$

Při vyváženém třídění dostáváme $n = abp$, $n_{i\cdot} = bp$, $n_{\cdot j} = ap$.

Součty čtverců

Skupinový součet čtverců S_B charakterizuje variabilitu mezi jednotlivými náhodnými výběry ve skupinách faktoru B . Jde tedy o součet čtverců rozdílů odhadů mezi modely M a M_B :

$$S_B = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} \left(\hat{\mu}_{ijk} - \hat{\mu}_{B,ijk} \right)^2 = \sum_{i=1}^a \sum_{j=1}^b n_{ij} (Y_{.j.} - \bar{Y}_{...})^2 = ap \sum_{j=1}^b Y_{.j.}^2 - n \bar{Y}_{...}^2$$

Podobně S_A charakterizuje variabilitu mezi jednotlivými náhodnými výběry ve skupinách faktoru A . Musíme však dodržet pořadí submodelů v řetězci, počítáme proto součet čtverců rozdílů odhadů mezi modely M_B a M_0 :

$$S_A = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} \left(\hat{\mu}_{B,ijk} - \hat{\mu}_{0,ijk} \right)^2 = \sum_{i=1}^a \sum_{j=1}^b n_{ij} (Y_{i..} - \bar{Y}_{...})^2 = bp \sum_{i=1}^a Y_{i..}^2 - n \bar{Y}_{...}^2$$

Celkový součet čtverců charakterizuje variabilitu jednotlivých pozorování kolem celkového průměru, bez ohledu na faktory. Jde tedy o součet čtverců rozdílů odhadů mezi pozorovanými veličinami a modelem M_0 :

$$S_T = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} \left(Y_{ijk} - \hat{\mu}_{0,ijk} \right)^2 = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} \left(Y_{ijk} - \bar{Y}_{...} \right)^2 = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} Y_{ijk}^2 - n \bar{Y}_{...}^2$$

Součty čtverců

Reziduální součet čtverců charakterizuje varibilitu nevysvětlenou modelem M :

$$S_e = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} (Y_{ijk} - \hat{\mu}_{ijk})^2 = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} (Y_{ijk} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...})^2$$

Věta 10

Platí

$$S_e = S_T - S_A - S_B$$

Analýza rozptylu dvojného třídění je založena na porovnání podílů rozptylů S_B/S_e a S_A/S_e vážených odpovídajícími stupni volnosti.

ANOVA tabulka dvojného třídění bez interakcí

Předcházející pojmy se shrnují v tzv. **tabulce analýzy rozptylu dvojného třídění**:

<i>Zdroj variability</i>	<i>Součet čtverců</i> SS	<i>Stupně volnosti</i> df	<i>Podíl</i> $MS = \frac{SS}{df}$	$F = \frac{MS}{s^2}$	<i>p-hodnota</i>
<i>řádky (A)</i>	S_A	$df_A = a - 1$	$MS_A = \frac{S_A}{df_A}$	$F_A = \frac{MS_A}{MS_e}$	$P(F \geq F_A)$
<i>sloupce (B)</i>	S_B	$df_B = b - 1$	$MS_B = \frac{S_B}{df_B}$	$F_B = \frac{MS_B}{MS_e}$	$P(F \geq F_B)$
<i>reziduální</i>	S_e	$df_e = n - a - b + 1$	$MS_e = \frac{S_e}{df_e}$	–	–
<i>celkový</i>	S_T	$df_T = n - 1$	–	–	–

Věta 11

Rozdíl mezi modely M a M_B ověřujeme pomocí testové statistiky

$$F_B = \frac{S_B/df_B}{S_e/df_e},$$

kteřá má za platnosti H_{B0} rozdělení pravděpodobnosti

$$F_B \sim F(df_B, df_e) = F(b - 1, n - a - b + 1).$$

Rozdíl mezi modely M_B a M_0 ověřujeme pomocí testové statistiky

$$F_A = \frac{S_A/df_A}{S_e/df_e},$$

kteřá má za platnosti H_{A0} rozdělení pravděpodobnosti

$$F_A \sim F(df_A, df_e) = F(a - 1, n - a - b + 1).$$

Věta 12

Porovnáváme modely M a M_B . Pokud

$$F_B = \frac{S_B/df_B}{S_e/df_e} \geq F_{1-\alpha}(b-1, n-a-b+1),$$

zamítneme na hladině významnosti α nulovou hypotézu H_{B0} , tzn. statisticky prokážeme vliv faktoru B (sloupce).

Pokračujeme porovnáním modelů M_B a M_0 . Pokud

$$F_A = \frac{S_A/df_A}{S_e/df_e} \geq F_{1-\alpha}(a-1, n-a-b+1),$$

zamítneme navíc na hladině významnosti α nulovou hypotézu H_{A0} , tzn. statisticky prokážeme také vliv faktoru A (řádky).

Postup testování ve dvojném třídění

1. $M \rightarrow M_B$

Pomocí F_B testujeme rozdíly mezi sloupci (faktor B) a přitom přihlížíme k eventuálním řádkovým efektům.

2. $M_B \rightarrow M_0$

Pomocí F_A testujeme rozdíly mezi řádky (faktor A), **nebereme** však v úvahu případný vliv sloupcových efektů.

Postup testování ve dvojném třídění

1. $M \longrightarrow M_B$

Pomocí F_B testujeme rozdíly mezi sloupci (faktor B) a přitom přihlížíme k eventuálním řádkovým efektům.

2. $M_B \longrightarrow M_0$

Pomocí F_A testujeme rozdíly mezi řádky (faktor A), **nebereme** však v úvahu případný vliv sloupcových efektů.

jiný postup: $M \longrightarrow M_A \longrightarrow M_0$

Analýzu rozptylu můžeme ale provést také v řetězci $M \longrightarrow M_A \longrightarrow M_0$, čímž dokážeme testovat rozdíly mezi řádky (faktor A) při přihlédnutí k eventuálním sloupcovým efektům.

V případě vyváženého designu vyjdou analýzy obou řetězců stejně. Pro nevyvážený design mohou analýzy obou řetězců vyjít odlišné.

Mnohonásobné porovnávání

Zjistíme-li významný rozdíl mezi řádky (faktor A), zjišťujeme dále, které řádky (úrovně faktorů A) se od sebe signifikantně liší. K ověření máme opět Scheffého metodu a Tukeyovu metodu, z nichž v praxi vybíráme tu, která je citlivější.

Využíváme přitom následující tvrzení:

Věta 13

Za platnosti M jsou veličiny $(\bar{Y}_{1..}, \dots, \bar{Y}_{a..})$ stochasticky nezávislé a platí

$$\bar{Y}_{i..} \sim N \left(\mu + \alpha_i + \frac{1}{b} \sum_{j=1}^b \beta_j, \frac{\sigma^2}{bp} \right), \quad i = 1, \dots, a.$$

Zjistíme-li významný rozdíl mezi sloupci (faktor B), analogicky zjišťujeme, které sloupce (úrovně faktorů B) se od sebe signifikantně liší.

Porovnání řádků (faktor A)

Porovnávání řádků ($A = u$) a ($A = v$) odpovídá porovnávání efektů α_u a α_v .

Věta 14 (Scheffého metoda)

Hypotézu o rovnosti $\alpha_u = \alpha_v$ zamítáme, pokud

$$|\bar{Y}_{u..} - \bar{Y}_{v..}| > \sqrt{\frac{2(a-1)S_e}{bp(n-a-b+1)} F_{1-\alpha}(a-1, n-a-b+1)}.$$

Věta 15 (Tukeyova metoda)

Hypotézu o rovnosti $\alpha_u = \alpha_v$ zamítáme, pokud

$$|\bar{Y}_{u..} - \bar{Y}_{v..}| > \sqrt{\frac{S_e}{bp(n-a-b+1)} q_{1-\alpha}(a, n-a-b+1)}.$$

Porovnání sloupců (faktor B)

Porovnávání řádků ($B = u$) a ($B = v$) odpovídá porovnávání efektů β_u a β_v .

Věta 16 (Scheffého metoda)

Hypotézu o rovnosti $\beta_u = \beta_v$ zamítáme, pokud

$$|\bar{Y}_{..u.} - \bar{Y}_{..v.}| > \sqrt{\frac{2(b-1)S_e}{ap(n-a-b+1)}} F_{1-\alpha}(b-1, n-a-b+1).$$

Věta 17 (Tukeyova metoda)

Hypotézu o rovnosti $\beta_u = \beta_v$ zamítáme, pokud

$$|\bar{Y}_{..u.} - \bar{Y}_{..v.}| > \sqrt{\frac{S_e}{ap(n-a-b+1)}} q_{1-\alpha}(b, n-a-b+1).$$

Dvojné třídění s interakcemi (vyvážené)

U dvojného třídění se často stává, že se řádkové a sloupcové efekty jen prostě nesčítají, jak to předpokládá náš model M .

V takových situacích uvažujeme následující komplexnější model.

Definice 18 (model M_*)

Náhodné veličiny Y_{ijk} se řídí modelem M_* , pokud

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \lambda_{ij} + \varepsilon_{ijk}$$

pro $i = 1, \dots, a$, $j = 1, \dots, b$ a $k = 1, \dots, p$, přičemž ε_{ijk} jsou stochasticky nezávislé náhodné veličiny s rozložením $\varepsilon_{ijk} \sim N(0, \sigma^2)$.

Interpretace neznámých parametrů:

- μ je společná část střední hodnoty sledované veličiny
- $\alpha_1, \dots, \alpha_a$ jsou efekty faktoru A ,
- β_1, \dots, β_b jsou efekty faktoru B ,
- λ_{ij} ($i = 1, \dots, a$, $j = 1, \dots, b$) jsou tzv. **interakce**,
- σ^2 je rozptyl náhodných chyb

Model M_*

Zapišeme M_* jako lineární regresní model:

$$Y = X_* \underbrace{(\mu, \alpha_1, \dots, \alpha_a, \beta_1, \dots, \beta_b, \lambda_{11}, \dots, \lambda_{ab})'}_{\text{vektor parametrů } \theta_*}$$

Jaké rozměry a jakou hodnotu má matice plánu X_* ?

Model M_*

Zapišeme M_* jako lineární regresní model:

$$Y = X_* \underbrace{(\mu, \alpha_1, \dots, \alpha_a, \beta_1, \dots, \beta_b, \lambda_{11}, \dots, \lambda_{ab})'}_{\text{vektor parametrů } \theta_*}$$

Jaké rozměry a jakou hodnotu má matice plánu X_* ?

$$X_* : n \times (1 + a + b + ab), \quad h(X_*) = ab$$

Matice plánu není plně hodnosti.

Model se obvykle řeší přidáním reparametrizačních rovnic

$$\sum_{i=1}^a \alpha_i = 0, \quad \sum_{j=1}^b \beta_j = 0, \quad \sum_{i=1}^a \lambda_{ij} = 0, \quad \sum_{j=1}^b \lambda_{ij} = 0.$$

Odhady parametrů v modelu M_*

Věta 19 (odhady parametrů modelu M_*)

Odhady parametrů modelu

$$M_* : \mu + \alpha_i + \beta_j + \lambda_{ij} + \varepsilon_{ijk}$$

metodou nejmenších čtverců jsou rovny

$$\mu = \bar{Y}_{\dots},$$

$$\alpha_i = \bar{Y}_{i..} - \bar{Y}_{\dots},$$

$$\beta_j = \bar{Y}_{.j.} - \bar{Y}_{\dots},$$

$$\lambda_{ij} = \bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{\dots}.$$

Odhad střední hodnoty měření ve skupině ($A = i, B = j$) modelu M_ je rovný*

$$\hat{\mu}_{*ijk} = \bar{Y}_{ij.}$$

Porovnejte s odhady pro model M .

Věta 20

Při vyváženém dvojném třídění s interakcemi platí:

$$S_B = ap \sum_{j=1}^b Y_{.j.}^2 - n \bar{Y}_{...}^2$$

$$S_A = bp \sum_{i=1}^a Y_{i..}^2 - n \bar{Y}_{...}^2$$

$$S_T = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} Y_{ijk}^2 - n \bar{Y}_{...}^2$$

$$S_e = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} Y_{ijk}^2 - p \sum_{i=1}^a \sum_{j=1}^b Y_{ij.}^2$$

$$S_{AB} = S_T - S_A - S_B - S_e.$$

Porovnejte se součty čtverců pro model M.

ANOVA tabulka dvojného třídění s interakcemi

Tabulka analýzy rozptylu dvojného třídění s interakcemi:

Zdroj variability	Součet čtverců SS	Stupně volnosti df	Podíl $MS = \frac{SS}{df}$	$F = \frac{MS}{s^2}$	p -hodnota
řádky (A)	S_A	$df_A = a - 1$	$MS_A = \frac{S_A}{df_A}$	$F_A = \frac{MS_A}{MS_e}$	$P(F \geq F_A)$
sloupce (B)	S_B	$df_B = b - 1$	$MS_B = \frac{S_B}{df_B}$	$F_B = \frac{MS_B}{MS_e}$	$P(F \geq F_B)$
interakce	S_{AB}	$df_{AB} = (a - 1)(b - 1)$	$MS_{AB} = \frac{S_{AB}}{df_{AB}}$	$F_{AB} = \frac{MS_{AB}}{MS_e}$	$P(F \geq F_{AB})$
reziduální	S_e	$df_e = n - ab$	$MS_e = \frac{S_e}{df_e}$	–	–
celkový	S_T	$df_T = n - 1$	–	–	–

Testování v modelu s interakcemi

Testování ve dvojném třídění s interakcemi probíhá obvykle v řetězci submodelů

$$M_* \longrightarrow M \longrightarrow M_B \longrightarrow M_0.$$

Při vyváženém třídění lze řetězec submodelů libovolně měnit bez změny výsledku.

Věta 21

Porovnáváme modely M_ a M . Pokud*

$$F_{AB} = \frac{S_{AB}/df_{AB}}{S_e/df_e} \geq F_{1-\alpha}((a-1)(b-1), n-ab),$$

zamítneme na hladině významnosti α nulovou hypotézu

$$H_{*0} : \lambda_{ij} = 0, \quad i = 1, \dots, a, \quad j = 1, \dots, b,$$

tztn. statisticky prokážeme vliv interakcí řádků a sloupců.

Porovnávání dvojic řádků a dvojic sloupců se dále provádí podle Vět 14–17, v nichž se upraví počet stupňů volnosti df_e tak, že $(n - a - b + 1)$ nahradíme $(n - ab)$.

- k čemu se používá jednoduché a dvojné třídění (s interakcemi)
- testy ověření podmínek – normalita dat, homogenita rozptylů
- lineární regresní modely, význam parametrů, řetězec submodelů
- hypotézy a jejich souvislost s testováním submodelů
- testovací statistiky v analýze rozptylu založené na součtech čtverců
- interpretace výsledků v ANOVA tabulce
- výpočet efektů a interakcí a jejich interpretace
- výpočet středních hodnot ve skupinách
- metody mnohonásobného porovnávání

ANOVA v R

Funkce pro analýzu rozptylu:

```
model <- aov (formule , data)
```

Zápis *formule*:

```
# jednoduché tridení
```

```
Y ~ faktorA
```

```
# dvojnásobné tridení bez interakcí
```

```
Y ~ faktorA + faktorB
```

```
# dvojnásobné tridení s interakcemi
```

```
Y ~ faktorA + faktorB + faktorA:faktorB
```

```
# dvojnásobné tridení s interakcemi
```

```
Y ~ faktorA * faktorB
```

Porovnání regresních modelů:

```
model1 <- lm (formule1 , data)
```

```
model2 <- lm (formule2 , data)
```

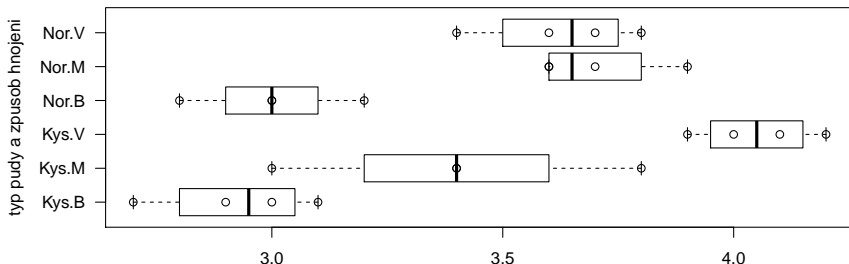
```
model <- anova (model1 , model2)
```


Příklad

Příklad 1

Zkoumají se výnosy sena v tunách na hektar v závislosti na typu půdy a na způsobu hnojení. Každá kombinace byla realizována čtyřikrát, nezávisle na sobě.

[t/ha]	způsob hnojení (B):		
typ půdy (A)	Bez hnojení	chlévká M rva	Vápenaté hnojivo
Normální	2,8; 3,2; 3,0; 3,0	3,7; 3,6; 3,9; 3,6	3,4; 3,8; 3,7; 3,6
Kyselá	3,1; 2,7; 3,0; 2,9	3,4; 3,4; 3,0; 3,8	4,2; 4,0; 4,1; 3,9



ANOVA tabulka dvojného třídění

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
puda	1	0.002	0.0017	0.027	0.871	
hnojeni	2	3.182	1.5912	25.752	2.93e-06	**
Residuals	20	1.236	0.0618			

Tables of effects

```

puda
puda
      Kys      Nor
0.008333 -0.008333
hnojeni
hnojeni
      B      M      V
-0.4875  0.1000  0.3875
  
```

Tables of means

```

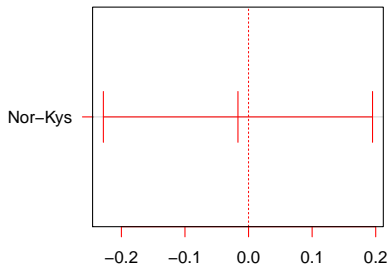
Grand mean
      3.45
puda
puda
      Kys      Nor
3.458  3.442
hnojeni
hnojeni
      B      M      V
2.963  3.550  3.838
  
```

(Intercept)	pudaNor	hnojeniM	hnojeniV
2.97083333	-0.01666667	0.58750000	0.87500000

Scheffé (typ půdy)

	trt	means	M
1	Kys	3.458333	a
2	Nor	3.441667	a

95% family-wise confidence level

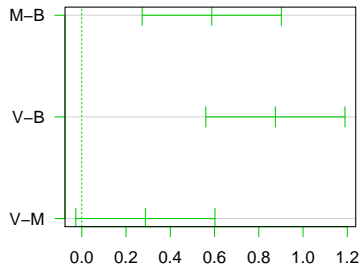


Differences in mean levels of puda

Scheffé (způsob hnojení)

	trt	means	M
1	V	3.8375	a
2	M	3.5500	a
3	B	2.9625	b

95% family-wise confidence level



Differences in mean levels of hnojeni

ANOVA tabulka dvojného třídění s interakcemi

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
puda	1	0.002	0.0017	0.044	0.83658	
hnojeni	2	3.182	1.5912	41.814	1.72e-07	***
puda:hnojeni	2	0.551	0.2754	7.237	0.00494	**
Residuals	18	0.685	0.0381			

parametry (nahore efekty, dole interakce)

(Intercept)	pudaNor	hnojeniM	hnojeniV
2.925	0.075	0.475	1.125
pudaNor:hnojeniM	pudaNor:hnojeniV		
0.225	-0.500		

Tables of effects

puda

puda

Kys

Nor

0.008333 -0.008333

hnojeni

hnojeni

B

M

V

-0.4875 0.1000 0.3875

puda:hnojeni

hnojeni

puda

B

M

V

Kys -0.04583 -0.15833 0.20417

Nor 0.04583 0.15833 -0.20417

Tables of means

Grand mean

3.45

puda

puda

Kys

Nor

3.458 3.442

hnojeni

hnojeni

B

M

V

2.963 3.550 3.838

puda:hnojeni

hnojeni

puda

B

M

V

Kys 2.925 3.400 4.050

Nor 3.000 3.700 3.625

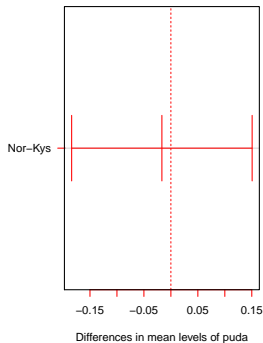
Scheffé (typ půdy)

	trt	means	M
1	Kys	3.458333	a
2	Nor	3.441667	a

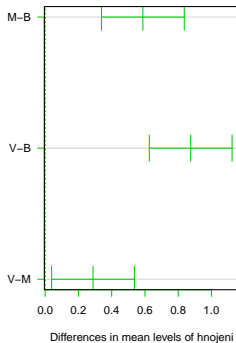
Scheffé (způsob hnojení)

	trt	means	M
1	V	3.8375	a
2	M	3.5500	b
3	B	2.9625	c

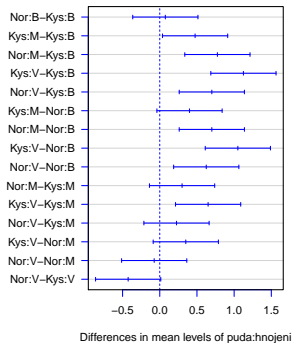
95% family-wise confidence level



95% family-wise confidence level



95% family-wise confidence level



ANOVA tabulka jednoduchého třídění

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
group	5	3.735	0.7470	19.63	1.02e-06 ***
Residuals	18	0.685	0.0381		

parametry

(Intercept)

2.925

groupKys.M	groupKys.V	groupNor.B	groupNor.M	groupNor.V
0.475	1.125	0.075	0.775	0.700

Tables of means

Grand mean

3.45

group

group

Kys.B	Kys.M	Kys.V	Nor.B	Nor.M	Nor.V
-------	-------	-------	-------	-------	-------

2.925	3.400	4.050	3.000	3.700	3.625
-------	-------	-------	-------	-------	-------

Tables of effects

group

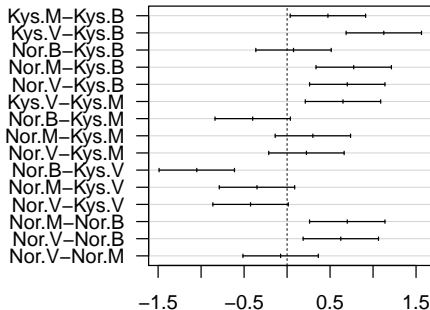
group

Kys.B	Kys.M	Kys.V	Nor.B	Nor.M	Nor.V
-0.525	-0.050	0.600	-0.450	0.250	0.175

95% family-wise confidence level

Scheffé (přůtahy)

	trt	means	M
1	Kys.V	4.050	a
2	Nor.M	3.700	ab
3	Nor.V	3.625	ab
4	Kys.M	3.400	bc
5	Nor.B	3.000	c
6	Kys.B	2.925	c



Differences in mean levels of group