

# ANOVA & spol.

---

JAN ŠEREK & STANDA JEŽEK

PSYB2520 STATISTICKÁ ANALÝZA  
DAT II

# Program dnešní přednášky

---

jednofaktorová (one-way) ANOVA

faktoriální (two...-way) ANOVA

ANCOVA (ANOVA s kovariáty)

MANOVA (ANOVA s více závislými)

ANOVA pro opakovaná měření

# ANOVA (analysis of variance)

---

Liší se 2 skupiny v průměru nějaké proměnné? → **t-test**

- $H_0: \mu_1 = \mu_2$

Liší se 3 (a více) skupiny v průměru nějaké proměnné? → **ANOVA**

- $H_0: \mu_1 = \mu_2 = \mu_3 \dots$

- „Liší se děti z úplných rodin, neúplných rodin a náhradní péče ve své partnerské spokojenosti?“
- „Liší se průměrná tepová frekvence participantů, kteří byli vystavení podnětu A, podnětu B a žádnému podnětu (kontrolní skupina)?“

# ANOVA (**a**nalysis of **v**ariance)

---

Liší se 2 skupiny v průměru nějaké proměnné? → **t-test**

- $H_0: \mu_1 = \mu_2$

Liší se 3 (a více) skupiny v průměru nějaké proměnné? → **ANOVA**

- $H_0: \mu_1 = \mu_2 = \mu_3 \dots$

- „Liší se děti z úplných rodin, neúplných rodin a náhradní péče ve své partnerské spokojenosti?“
- „Liší se průměrná tepová frekvence participantů, kteří byli vystavení podnětu A, podnětu B a žádnému podnětu (kontrolní skupina)?“

1 nezávislá kategorická → 1 závislá intervalová

# ANOVA (analysis of variance)

---

Liší se 2 skupiny v průměru nějaké proměnné? → **t-test**

- $H_0: \mu_1 = \mu_2$

Liší se 3 (a více) skupiny v průměru nějaké proměnné? → **ANOVA**

- $H_0: \mu_1 = \mu_2 = \mu_3 \dots$

v jazyku ANOVY se tato nezávislá kategorická proměnná nazývá **faktor**, který má určité **úrovně**

in, neúplných rodin a náhradní péče  
nosti?"

frekvence participantů, kteří byli  
nětu B a žádnému podnětu

1 nezávislá kategorická → 1 závislá intervalová

# ANOVA – 2 základní kroky

---

## **KROK 1:** Jsou průměry skupin stejné?

- Testujeme  $H_0: \mu_1 = \mu_2 = \mu_3 \dots$
- Testová statistika je  $F$  ptáme se po  $p$

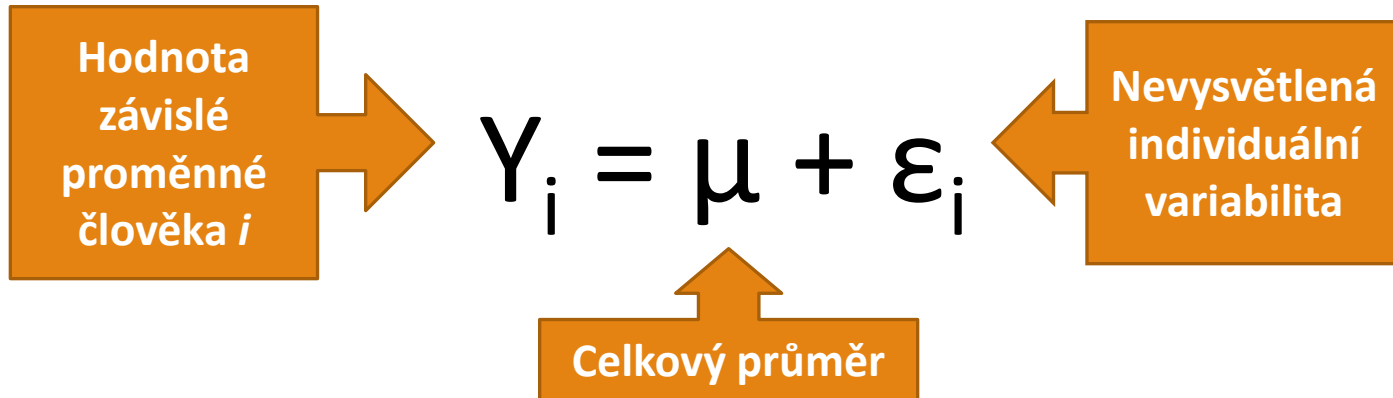
$P(\text{diff}M|H_0)$ ,  
kde  $\text{diff}M$  zastupuje  
„nejméně taková  
míra rozdílnosti  
průměrů, jakou  
pozorujeme“

Pokud ANO → konec

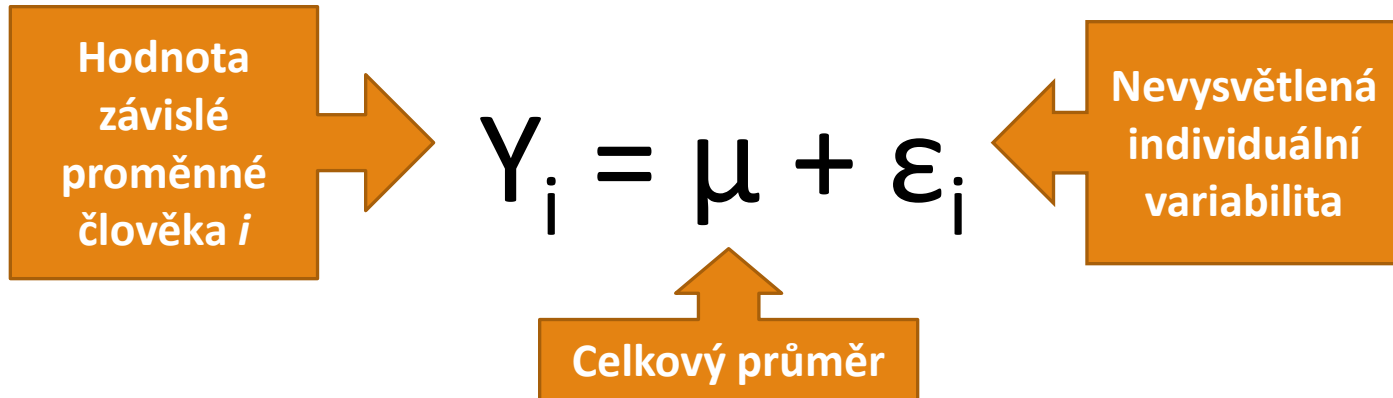
Pokud NE → **KROK 2:** Které skupiny konkrétně se mezi sebou svými průměry liší?

- máme o tom hypotézy → plánované kontrasty
- nemáme o tom hypotézy → post-hoc testy

# ANOVA jako regrese



# ANOVA jako regrese

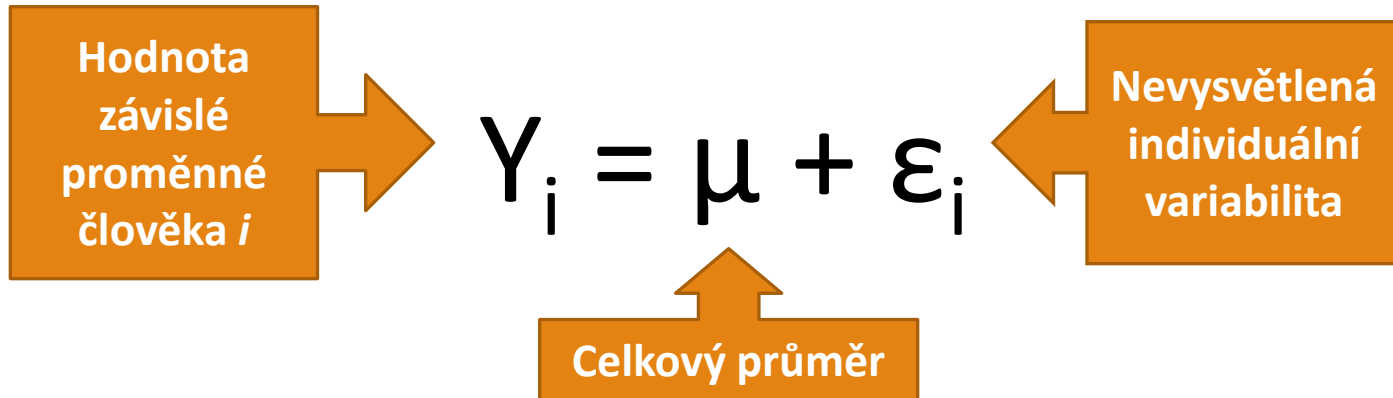


$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$$

Vliv toho, že je člověk členem skupiny  $j$



# ANOVA jako regrese



$$Y_i = \mu + \alpha_j + \varepsilon_i$$

## Podstata ANOVY

Jak dobře je závislá proměnná vysvětlena modelem, který předpokládá odlišnost skupin ( $\alpha \neq 0$ )? Nepostačí nám stejně dobře model, který předpokládá, že se skupiny neliší?

# ANOVA jako regrese

---

Souvisí socioekonomický status rodiny s tím, jak často dítě používá internet?

Nezávislá kategorická proměnná (*faktor*):

**socioekonomický status**

3 hodnoty (*úrovně*): nízký, střední, vysoký

Závislá intervalová proměnná: **frekvence používání internetu**

Liší se děti z rodin s nízkým, středním a vysokým SES v tom, jak často používají internet?

# ANOVA jako regrese

$$INET_i = \mu + \varepsilon_i$$

Celkový průměr

$$INET_i = \mu + \text{SES}_j + \varepsilon_i$$

Vliv toho, že je člověk  
členem skupiny  $j$

# ANOVA jako regrese

---

$$[\text{inet}]_i = [\text{průměrný inet}] + \varepsilon_i$$

# ANOVA jako regrese

---

$$[\text{inet}]_i = [\text{průměrný inet}] + \varepsilon_i$$

$$[\text{inet}]_i = [\text{průměrný inet}] + \mathbf{b}[\text{ses}] + \varepsilon_i$$

# ANOVA jako regrese

---

$$[\text{inet}]_i = [\text{průměrný inet}] + \varepsilon_i$$

$$[\text{inet}]_i = [\text{průměrný inet}] + \mathbf{b}[\mathbf{ses}] + \varepsilon_i$$

Každá kategorická proměnná o  $k$  hodnotách (úrovních) může být vyjádřena souborem  $k-1$  binárních dummy proměnných.

3 typy SES  $\rightarrow$  2 binární proměnné **vys** a **str**

**vys** = 1 & **str** = 0  $\rightarrow$  vysoký SES

**vys** = 0 & **str** = 1  $\rightarrow$  střední SES

**vys** = 0 & **str** = 0  $\rightarrow$  nízký SES

# ANOVA jako regrese

---

$$[\text{inet}]_i = [\text{průměrný inet}] + \varepsilon_i$$

$$[\text{inet}]_i = [\text{průměrný inet}] + \mathbf{b}[\mathbf{ses}] + \varepsilon_i$$

$$[\text{inet}]_i = \mathbf{b}_0 + \mathbf{b}_1[\mathbf{vys}] + \mathbf{b}_2[\mathbf{str}] + \varepsilon_i$$

Každá kategorická proměnná o  $k$  hodnotách (úrovních) může být vyjádřena souborem  $k-1$  binárních dummy proměnných.

3 typy SES  $\rightarrow$  2 binární proměnné **vys** a **str**

**vys** = 1 & **str** = 0  $\rightarrow$  vysoký SES

**vys** = 0 & **str** = 1  $\rightarrow$  střední SES

**vys** = 0 & **str** = 0  $\rightarrow$  nízký SES

# ANOVA i

Průměrná  
frekvence dětí z  
rodin s nízkým  
SES

O kolik se liší  
průměrná  
frekvence dětí z  
rodin s vysokým  
SES

O kolik se liší  
průměrná  
frekvence dětí z  
rodin se  
středním SES

[inet]<sub>i</sub> =  $b_0 + b_1[vys] + b_2[str] + \epsilon_i$

$$[inet]_i = b_0 + b_1[vys] + b_2[str] + \epsilon_i$$

Každá kategorická proměnná o  $k$  hodnotách (úrovních) může být vyjádřena souborem  $k-1$  binárních dummy proměnných.

3 typy SES → 2 binární proměnné **vys** a **str**

**vys** = 1 & **str** = 0 → vysoký SES

**vys** = 0 & **str** = 1 → střední SES

**vys** = 0 & **str** = 0 → nízký SES



# ANOVA

O kolik se liší

O kolik se liší

Jestliže  $b_1 = 0$  a  $b_2 = 0$ , znamená to, že SES nemá žádný vliv a všechny skupiny mají stejnou průměrnou frekvenci.

Potom by nám stačil základní model predikující frekvenci pouze z celkové průměrné frekvence a nevysvětlitelné individuální variability.

Vysvětlí nám model předpokládající nenulové  $b_1$  a/nebo  $b_2$  něco navíc?

# ANOVA jako regrese

Lze SES nakódovat tak, aby  $b_0$  byl celkový průměr?

$$[\text{inet}]_i = [\text{průměrný inet}] + \mathbf{b}[\text{ses}] + \varepsilon_i$$

$$[\text{inet}]_i = \mathbf{b}_0 + \mathbf{b}_1[\mathbf{vys}] + \mathbf{b}_2[\mathbf{str}] + \varepsilon_i$$

efektové kódování SES  $\rightarrow$  **vys** a **str**

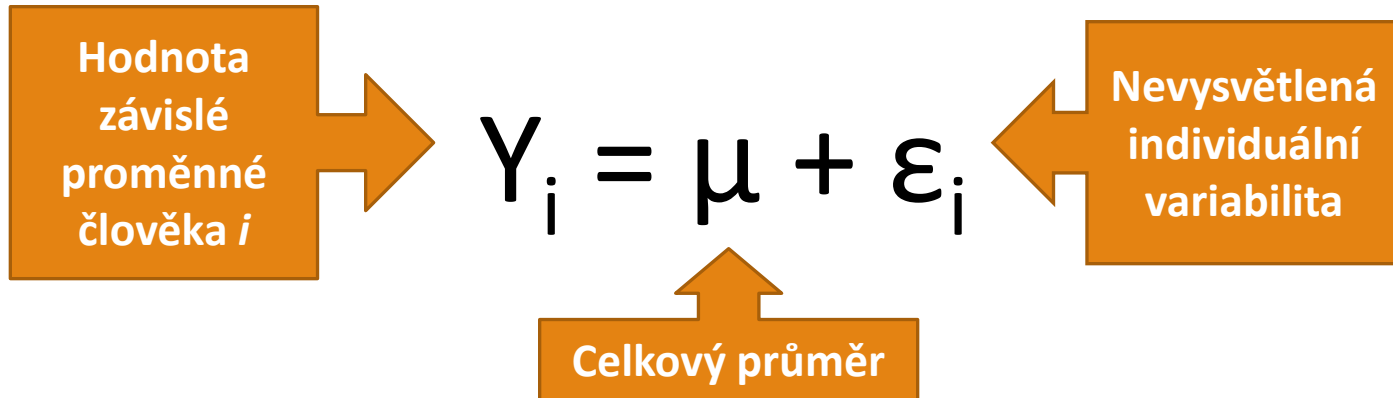
**vys** = 1 & **str** = 0 pro vysoký SES

**vys** = 0 & **str** = 1 pro střední SES

**vys** = -1 & **str** = -1 pro nízký SES

nízký SES stále referenční, ale  $b$  vyjadřují rozdíl skupinového průměru proti celkovému

# ANOVA jako regrese



$$Y_i = \mu + \alpha_j + \varepsilon_i$$

## Podstata ANOVY

Jak dobře je závislá proměnná vysvětlena modelem, který předpokládá odlišnost skupin ( $\alpha \neq 0$ )? Nepostačí nám stejně dobře model, který předpokládá, že se skupiny neliší?

# ANOVA jako ANOVA

---

V rámci lineární regrese umíme modelovat vliv kategorické nezávislé pomoci dummy proměnných a víme, že regresní koeficienty  $b$  udávají rozdíly průměrů indikovaných skupin oproti referenční skupině.

Dokázali bychom použít i efektové kódování místo indikátorového a testovat tak rozdíly průměrů indikovaných skupin oproti celkovému průměru.

Jak ale srovnává průměry ANOVA?

Jiným způsobem, který je ale ve výsledku ekvivalentní regresi.

# ANOVA je Analýza rozptylu

---

Je tedy založena na analýze (~dělení) rozptylu závislé proměnné.

- $s^2 = \frac{\sum_{i=1}^N (X_i - M)^2}{N-1} = \frac{SS}{N-1}$

Rozptyl reprezentuje variabilitu, tedy *rozdíly mezi lidmi*. Reprezentuje je součtem kvadratických odchylek od průměru – **suma čtverců - SS.\***

- Pro rozptyl dělíme SS konstantou (N-1), ale tím, co zachycuje kvantitu těch rozdílů, je SS.

**Analýzou** miníme úvahu o tom, co mohlo tyto rozdíly způsobit – třeba členství ve skupině.

Když počítáme rozptyl pro všechny dohromady, aniž bychom vzali do úvahy, do jaké skupiny patří, počítáme **celkovou sumu čtverců –  $SS_T$ ,  $SS_{Total}$**

$$SS_T = \sum_i (\text{hodnota člověka } i - \text{celkový průměr})^2$$

$$MS_T = SS_T / (N-1) = SS_T / df_T$$

**Rozptyl se v analýze rozptylu nejmenuje rozptyl, ale MEAN SQUARE ;-)\*\***

**df: počet kousků informace – 1**

# Suma čtverců meziskupinová (modelová) $SS_B, SS_M$

Měřítko toho, jak moc se průměry skupin liší.

Jak vysoká je  $SS$  jenom díky rozdílům skupinových průměrů?

= Jaká by byla  $SS$ , kdyby měli všichni členové skupin hodnotu právě rovnou průměru skupiny? Tj. odchylky členů od průměru skupiny ignorujeme.

= Kolik variability ZP lze připsat odlišnostem mezi průměry skupin (modelu)?

$$SS_M = \sum_j \text{velikost skupiny } j * (\text{průměr skupiny } j - \text{celkový průměr})^2$$

$$MS_M = SS_M / df_M \quad \text{kde } df_M = (\text{počet skupin} - 1)^*$$

# Skoro bychom mohli být hotoví...

---

$SS_M / SS_T$  vyjadřuje, jaký podíl z celkového rozptylu proměnné je vysvětlen tím, že různí lidé pocházejí z různých skupin, které se liší svým průměrem

- Je to ekvivalent  $R^2$  Jmenuje se  $\eta^2$  éta na druhou, eta squared
- „Přeložili“ jsme rozdíly mezi průměry do vysvětleného rozptylu.
- $\eta^2$  je stejně jako  $R^2$  nadhodnocené, a tak si ukážeme jeho korigovanou verzi  $\omega^2$

Stále ještě nevíme, jak moc může být  $SS_M$  nadhodnocené jen díky výběrové chybě, tak abychom mohli otestovat  $H_0$ .

# Suma čtverců vnitroskupinová (reziduální)

$SS_W$ ,  $SS_R$

Měřítko toho, jak moc se lidé liší mezi sebou uvnitř skupin.

Jak velkou část  $SS_T$  tvoří odchylky jednotlivce od průměru jeho skupiny?  $SS_W$  jako WITHIN-GROUP  $SS_R$  jako RESIDUAL\*

Lze také říci: Jaká by byla  $SS$ , kdyby měly všechny skupiny stejný průměr?

Lze interpretovat jako vážený průměr rozptylů uvnitř skupin.

$$SS_R = \sum_{ij} (\text{hodnota člověka } i \text{ ze skupiny } j - \text{průměr skupiny } j)^2$$

$$MS_R = SS_R / df_R \quad df_R = (\text{celkový počet lidí} - \text{počet skupin})$$



---

Dílčí sumy čtverců se nasčítají do celkové

$$SS_T = SS_M + SS_R$$

Stupně volnosti se se také nasčítají –  
reprezentují „kousky využitých informací“

$$df_T = df_M + df_R$$

Pro střední čtverce to neplatí, protože ve jmenovateli jsou pokaždé jiné stupně volnosti. Neplatí tedy  $MST = MSM + MSR$ !

Ale pokud **platí**  $H_0$ , pak  $MS_M = MS_R$

# ANOVA – statistika F

$$F = MS_M / MS_R$$

Platí-li  $H_0$  očekáváme  $F$  kolem 1.

Čím vyšší  $F$ , tím více záleží na rozdělení lidí do jednotlivých skupin, **tj. tím více se skupiny od sebe liší v závislé proměnné**

$F$  je výběrová statistika, která má *Fisherovo* rozložení, definované dvojicí stupňů volnosti ( $df_M, df_R$ )

Můžeme testovat, zda ji hodnota  $F$  v našem výzkumu překračuje, **tj. testovat statistickou významnost nalezených rozdílů mezi skupinami**

# ANOVA – předpoklady $F$ -testu

---

**nezávislost pozorování** ( $\rightarrow$  ANOVA pro opakovaná měření)

**normalita rozložení** (v rámci každé skupiny)

- narušení nevadí, pokud jsou skupiny stejně velké + mají velikost alespoň okolo 30
- neparametrická alternativa – Kruskal-Wallisův test

**homogenita rozptylů** (skupiny mají stejné rozptyly)

- Levenův test – chceme, aby byl nesignifikantní
- $s^2_{\max} / s^2_{\min} < 3$
- narušení by nemělo vadit, pokud jsou skupiny stejně velké
- při narušení lze použít **Welchovo F**

# Co jsme zatím získali oproti regresi s kategorickým prediktorem?

---

Technicky vzato, nic moc, protože výsledný model ANOVA je stejný – lineární model **predikující** každému průměr jeho skupiny.

Zatímco v regresi byl trik s dummy proměnnými jakýsi *hack*, který nám umožnil zařadit kategorické proměnné, v ANOVA jsme přímo vyšli ze srovnávání průměrů.

Uvědomili jsme si tím, jak jsou **rozdíly skupinových průměrů** přeloženy do **vysvětleného rozptylu**.

To znamená, že dokážeme uvažovat o kvantifikaci **vlivu** kategorické proměnné na nějakou metrickou ZP bez ohledu na počet kategorií.

Ale F-test je pouze test  $H_0$  – po něm se chceme interpretačně vrátit k rozdílům mezi skupinami.

# ANOVA – SPSS

---

Analyze → Compare Means → One-Way ANOVA

	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	$SS_M$	$df_M$	$MS_M$		
Within Groups	$SS_R$	$df_R$	$MS_R$		
Total	$SS_T$	$df_M + df_R$			

DCtimeuse Estimated minutes online each day

SES	M	SD	N
1 High	102,9	63,3	6 274
2 Medium	107,7	65,2	7 989
3 Low	96,1	64,2	3 555
<i>Celkem</i>	<i>103,7</i>	<i>64,5</i>	<i>17 818</i>
Nevážený průměr	102,2		

ANOVA

DCtimeuse Estimated minutes online each day					
	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	338 736,712	2	169 368,356	40,941	0,000
Within Groups	73698975,20	17 815	4 136,906		
Total	74037711,92	17 817			

# ANOVA

---

Máme hypotézy o konkrétních rozdílech mezi skupinami.

**H1:** Děti z rodin s nízkým SES používají internet méně často než ostatní děti.

**H2:** Děti z rodin se středním SES používají internet méně často než děti z rodin s vysokým SES.

# ANOVA – plánované kontrasty

---

Umožňují porovnat jednotlivé skupiny v jednom kroku bez nutnosti korigovat hladinu významnosti (bez snížení síly testu)

Jen když máme dopředu hypotézy

Kontrastů lze provést tolik, kolik je **počet skupin – 1**

Každý kontrast srovnává **2 průměry**

- průměr skupiny nebo průměr více skupin dohromady
- např. **NÍZ vs. STŘ+VYS** nebo **STŘ vs. VYS**

**ortogonální (nezávislé) kontrasty**

- skupina použitá v jednom srovnání není použitá v dalším

**neortogonální kontrasty**



# ANOVA – plánované kontrasty

Zkoumáme, zda daný kontrast (rozdíl mezi dvěma průměry) signifikantně přispívá k variabilitě vysvětlené modelem ( $SS_M$ )

Abychom to zjistili, jakoby překódujeme hodnoty dummy proměnných, aby odhadnuté parametry ( $b_1, b_2$  atd.) odrážely požadované kontrasty

$$[\text{inet}]_i = b_0 + b_1[\text{vys}] + b_2[\text{str}] + \varepsilon_i$$

$$[\text{inet}]_i = b_0 + b_1[\text{kontrast1}] + b_2[\text{kontrast2}] + \varepsilon_i$$

Kategorie	Kontrast 1 NÍZ vs. STŘ+VYS	Kontrast 2 STŘ vs. VYS
Vysoký SES	1/2	-1
Střední SES	1/2	1
Nízký SES	-1	0

# ANOVA – plánované kontrasty

Zkoumáme, zda daný kontrast (rozdíl mezi dvěma skupinami) významně přispívá k variabilitě vysvětlené modelem ( $SS_M$ )

Abychom to zjistili, jakoby překódujeme hodnoty odhadnuté parametry ( $b_1, b_2$  atd.) odrážely požadovaný kontrast, aby

Součet pro každý kontrast musí být 0

$$[inet]_i = b_0 + b_1[kontrast1] + b_2[kontrast2] + \varepsilon_i$$

Srovnávané skupiny musí mít odlišná znaménka

Skupina, kterou nechceme zahrnout  $\rightarrow 0$

Kategorie	Kontrast 1 NÍZ vs. STŘ+VYS	Kontrast 2 STŘ vs. VYS
Vysoká	1/2	-1
Střední	1/2	1
Nízká	-1	0

Skupiny brané dohromady musí mít stejné číslo

# SPSS prezentuje kontrasty jako t-testy

Contrast Tests

		Contrast	Value of Contrast	Std. Error	t	df	Sig. (2-tailed)
DCtimeuse Estimated minutes online each day	Assume equal variances	1	9,1842	1,20747	7,606	17 815	0,000
		2	-4,8649	1,08499	-4,484	17 815	0,000
	Does not assume equal variances	1	9,1842	1,20513	7,621	5 481,033	0,000
		2	-4,8649	1,08150	-4,498	13 646,398	0,000

Pokud součet kontrastových koeficientů se stejným znaménkem = 1, pak hodnota kontrastu (je vlastně jako  $b$ ) vypovídá o velikosti rozdílu průměrů srovnávaných skupin (či sloučených skupin)

# ANOVA – post-hoc testy

---

Používáme, pokud nemáme dopředu jasné hypotézy

Srovnávají vše se vším – každou skupinu s každou (ale neumí slučovat skupiny jako kontrasty)

Mají v sobě mechanismy zohledňující zvýšené riziko chyby I. typu

Z principu jsou oboustranné

Je jich mnoho – liší se v několika parametrech:

- konzervativní (ch. II. typu!) / liberální (ch. I. typu!)
- ne/vhodné pro rozdílně velké skupiny
- ne/vhodné pro rozdílné skupinové rozptyly

# ANOVA – post-hoc testy

---

Doporučení podle A. Fielda:

- stejně velké skupiny a skupinové rozptyly (ideální situace): **REGWQ** nebo **Tukey**
- pokud si chceme být jistí, že P chyby I. typu nepřekročí zvolenou hladinu: **Bonferroni**
- pokud jsou velikosti skupin trochu/hodně rozdílné: **Gabriel/Hochberg GT2**
- pokud pochybujeme o shodnosti skupinových rozptylů: **Games-Howell**

Multiple Comparisons

Dependent Variable: DCtimeuse Estimated minutes online each day

	(I) DPSESHH3 Socio-economic status - Household	(J) DPSESHH3 Socio-economic status - Household	Mean Difference (I-J)	Std. Error	Sig.
Tukey HSD	1 High	2 Medium	-4,86487*	1,08499	0,000
		3 Low	6,75176*	1,35021	0,000
	2 Medium	1 High	4,86487*	1,08499	0,000
		3 Low	11,61663*	1,29673	0,000
	3 Low	1 High	-6,75176*	1,35021	0,000
		2 Medium	-11,61663*	1,29673	0,000

\*. The mean difference is significant at the 0,050 level.

DCtimeuse Estimated minutes online each day

	DPSESHH3 Socio-economic status - Household	N	Subset for alpha = 0,050		
			1	2	3
Ryan-Einot-Gabriel-Welsch Range	3 Low	3 555	96,1097		
	1 High	6 274		102,8615	
	2 Medium	7 989			107,7263
	Sig.		1,000	1,000	1,000
Tukey HSD <sup>a,b</sup>	3 Low	3 555	96,1097		
	1 High	6 274		102,8615	
	2 Medium	7 989			107,7263
	Sig.		1,000	1,000	1,000

Means for groups in homogeneous subsets are displayed.

a. Uses Harmonic Mean Sample Size = 5301,721.

b. The group sizes are unequal. The harmonic mean of the group sizes is used. Type I error levels are not guaranteed.

# One-way ANOVA – reportování

---

$$F(df_M, df_R) = \dots, p = \dots, \eta^2 \text{ nebo } \omega^2 = \dots$$

Vždy uvést **deskriptivy pro každou skupinu** – alespoň velikost, průměr, směrodatnou odch.

Vždy dopočítat **velikost účinku** (interpretujeme jako  $R^2$  v lineární regresii)

$$\eta^2 = SS_M / SS_T$$

$$\omega^2 = [SS_M - (df_M)MS_R] / [SS_T + MS_R] \text{ (jako Adj. } R^2)$$

$df_M$  a  $df_R$  musejí být uváděny v tomto pořadí

U kontrastů uvádíme:  $t(df) = \dots, p = \dots, d$  nebo  $r = \dots$   $r = \sqrt{[t^2 / (t^2 + df)]}$

Neuvádíme Anova Table! Vše je v textu.

# One-way ANOVA - shrnutí

---

- Výsledkově shodná s lineární regresí (lineární model)
- Specifikace modelu optimalizovaná pro kategorické prediktory – faktory – tedy pro porovnávání průměrů
- Zdůrazňuje myšlenku dělení rozptylu závislé proměnné na části, které lze připsat různým zdrojům rozptylu (faktorů, náhodné chybě...).



„V modelu je pouze jeden faktor. Člověk je však ve skutečnosti obvykle členem více typů skupin najednou, což může mít vliv!“

„Provedeme více ANOV pro různé faktory (skupiny).“

„Tím se však vrátí známý problém s nárůstem rizika chyby I. typu. Navíc přijdeme o možnost posoudit vliv všech faktorů najednou v jednom modelu.“

„Můžeme přidat přímo do modelu další nezávislé kategorické proměnné – a spočítat tzv. **faktoriální ANOVU.**“

# Faktoriální ANOVA

---

ANOVA s více kategorickými nezávislými (faktory)

uplatnění v **experimentálních** designech, kde pracujeme s několika druhy experimentální manipulace nebo kde chceme zohlednit kromě experimentální manipulace i další proměnné (např. pohlaví)

uplatnění v **neexperimentálních** designech, kde chceme posoudit vliv více kategorických prediktorů najednou

# Typy faktorů (platí i pro one-way)

---

## Fixed factors

- Všechny úrovně faktoru, o které nám jde, jsou v našem výzkumu zahrnuty
- Obvykle máme hypotézy o rozdílech mezi konkrétními skupinami.
- „Liší se užívání internetu mezi třemi typy SES?“

## Random factors

- úrovně faktoru, zahrnuté v našem výzkumu, představují pouze náhodný vzorek z větší populace.
- Obvykle nás nezajímají rozdíly mezi konkrétními skupinami.
- Do F-testu je zahrnuta tato přidaná míra nejistoty → nižší síla testu
- „Liší se užívání internetu mezi zeměmi?“
- „Liší se užívání internetu podle školy, kterou adolescent navštěvuje?“

# One-way ANOVA

---

$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$$

# Faktoriální ANOVA

$$Y_{ijk} = \mu + \alpha_j + \beta_k + \gamma_{j \times k} + \varepsilon_{ijk}$$



# One-way ANOVA

---

$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$$

# Faktoriální ANOVA

$$Y_{ijk} = \mu + \alpha_j + \beta_k + \gamma_{j \times k} + \varepsilon_{ijk}$$

Vliv toho, že  
je člověk  
členem  
skupiny  $j$   
**MAIN EFFECT**

Vliv toho, že  
je člověk  
členem  
skupiny  $k$   
**MAIN EFFECT**

Vliv toho, že je člověk  
členem kombinace  
skupin  $j$  a  $k$

# One-way ANOVA

---

$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$$

# Faktoriální ANOVA

$$Y_{ijk} = \mu + \alpha_j + \beta_k + \gamma_{j \times k} + \varepsilon_{ijk}$$



# One-way ANOVA

---

$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$$

# Faktoriální ANOVA

$$Y_{ijk} = \mu + \alpha_j + \beta_k + \gamma_{j \times k} + \varepsilon_{ijk}$$



# Interakce (moderace)

---

- V různých úrovních jednoho faktoru se rozdíly mezi průměry úrovní druhého faktoru liší (rozdíl rozdílů).
- S měnící se úrovní jedné nezávislé se mění vliv druhé nezávislé na závislou proměnnou
- Nezávislá proměnná nemusí mít **hlavní efekt** (main effect) na závislou proměnnou, ale může ji ovlivňovat tím, že ovlivňuje vliv druhé nezávislé
- Při interpretaci interakcí je užitečné znázornění grafem.
- Jde o totéž, co jsme měli u regrese!
- V ANOVě bude interakce zahrnuta automaticky (lze změnit)



# Interakce (moderace)

---

**dva faktory** (případ faktoriální ANOVY)

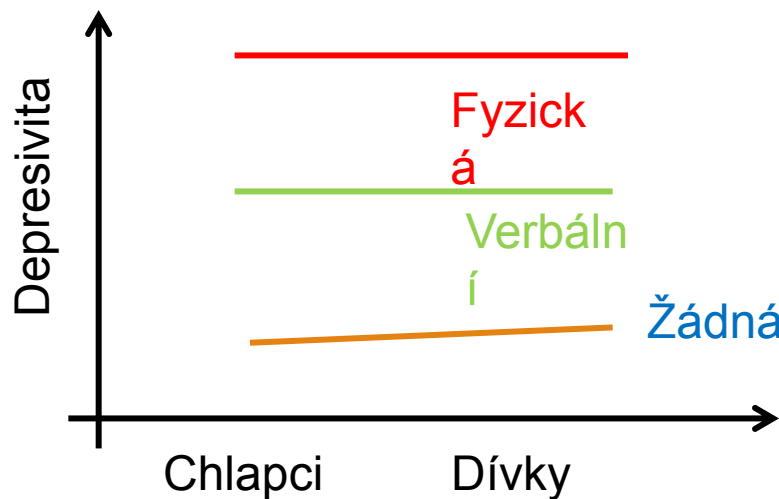
- Zážitek s různými typy školní šikany má jiný vliv na průměrnou depresivitu u dívek a u chlapců.

# interakce (moderace)

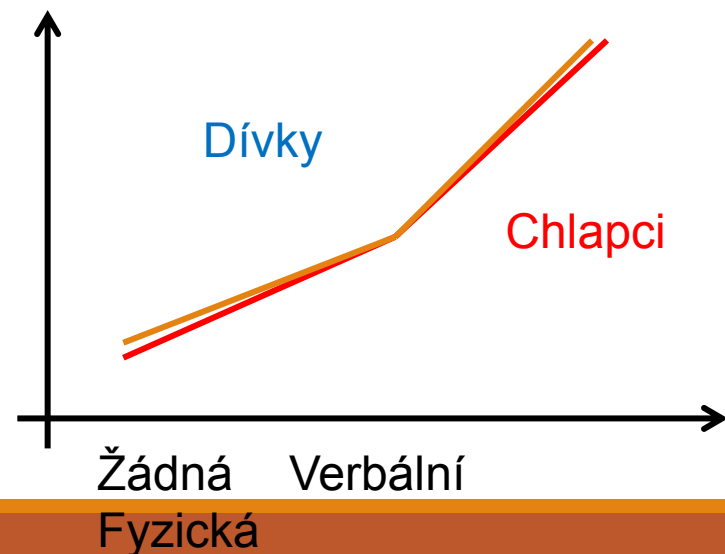
**dva faktory** (případ faktoriální ANOVY)

- Zážitek s různými typy školní šikany má jiný vliv na průměrnou depresivitu u dívek a u chlapců.

žádná interakce



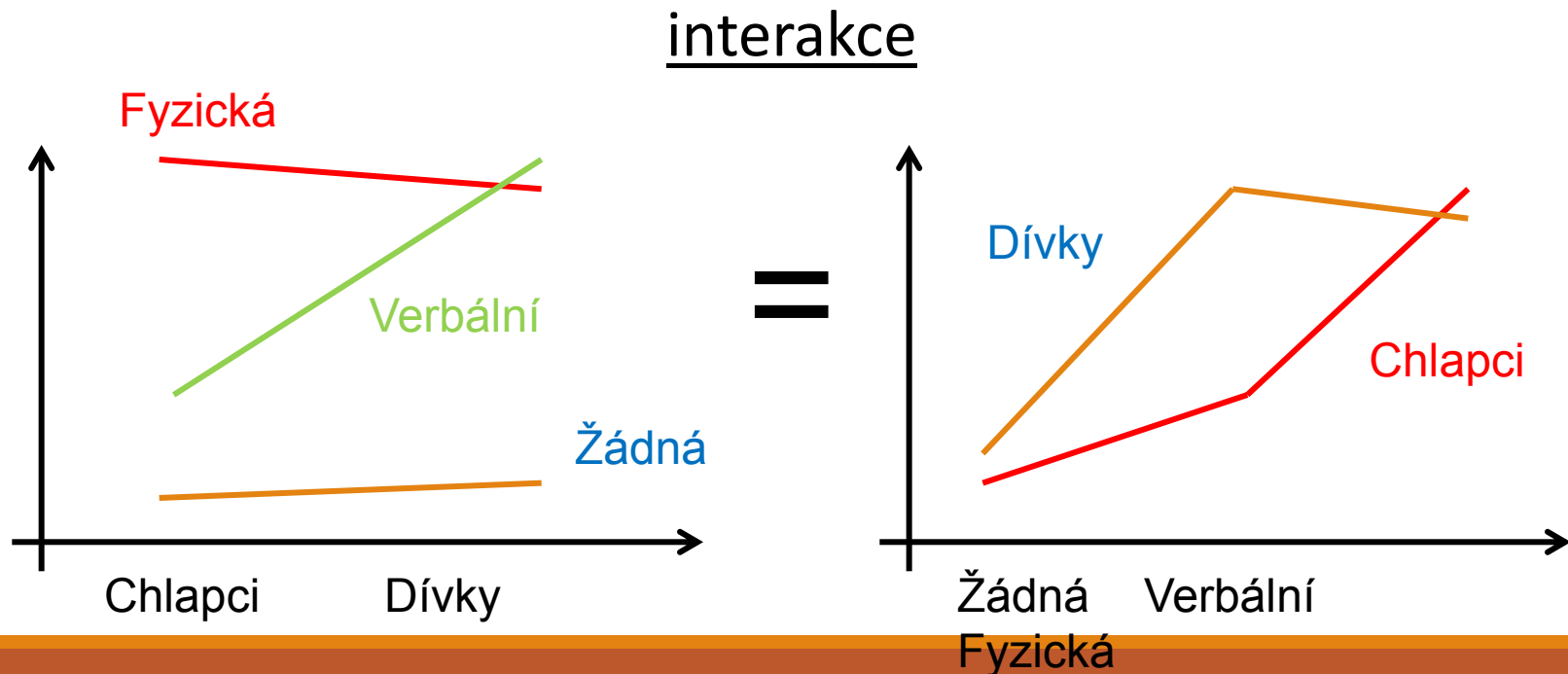
=



# interakce (moderace)

**dva faktory** (případ faktoriální ANOVY)

- Zážitek s různými typy školní šikany má jiný vliv na průměrnou depresivitu u dívek a u chlapců.



# interakce (moderace)

---

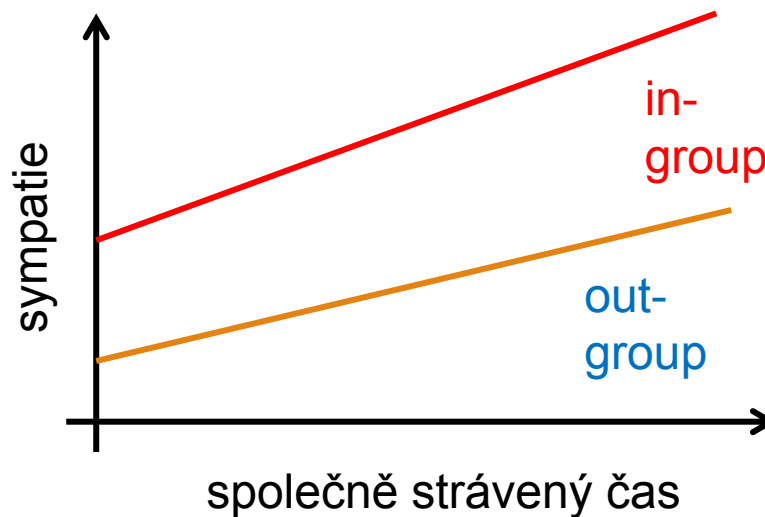
## kategorická a intervalová proměnná

- Společně strávený čas posiluje naše sympatie pouze k členům in-group, nikoli out-group.

# Interakce (moderace)

## kategorická a intervalová proměnná

- Společně strávený čas posiluje naše sympatie pouze k členům in-group, nikoli out-group.

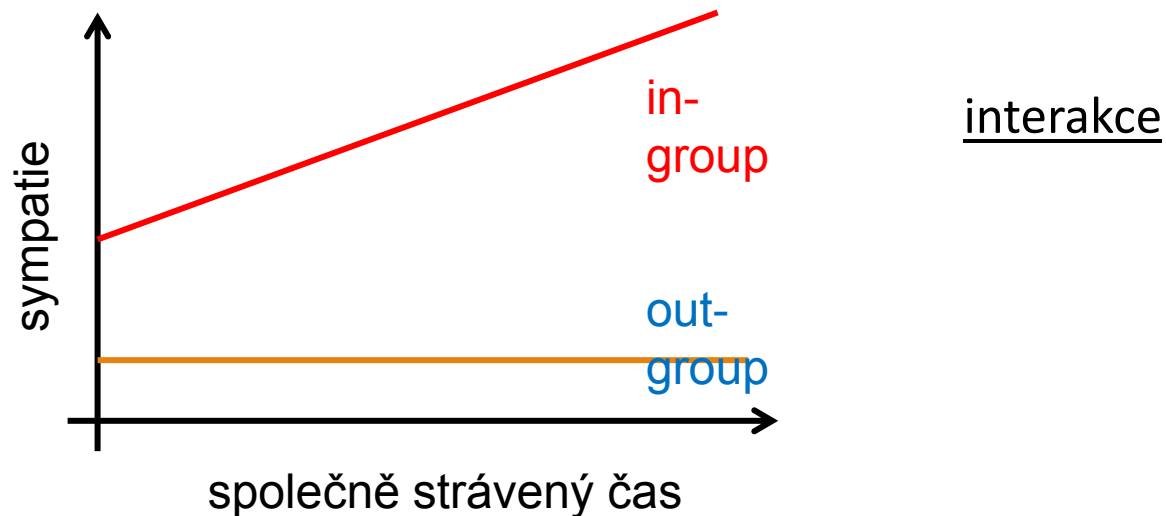


žádná interakce

# Interakce (moderace)

## kategorická a intervalová proměnná

- Společně strávený čas posiluje naše sympatie pouze k členům in-group, nikoli out-group.



# Interakce (moderace)

---

**dva faktory** (případ faktoriální ANOVY)

- Zážitek s různými typy školní šikany má jiný vliv na depresivitu u dívek a u chlapců.

**kategorická a intervalová proměnná**

- Společně strávený čas posiluje naše sympatie pouze k členům in-group, nikoli out-group.

**dvě intervalové proměnné**

- S rostoucím příjmem se oslabuje vztah mezi spokojeností v práci a celkovou životní spokojeností.

# Faktoriální ANOVA

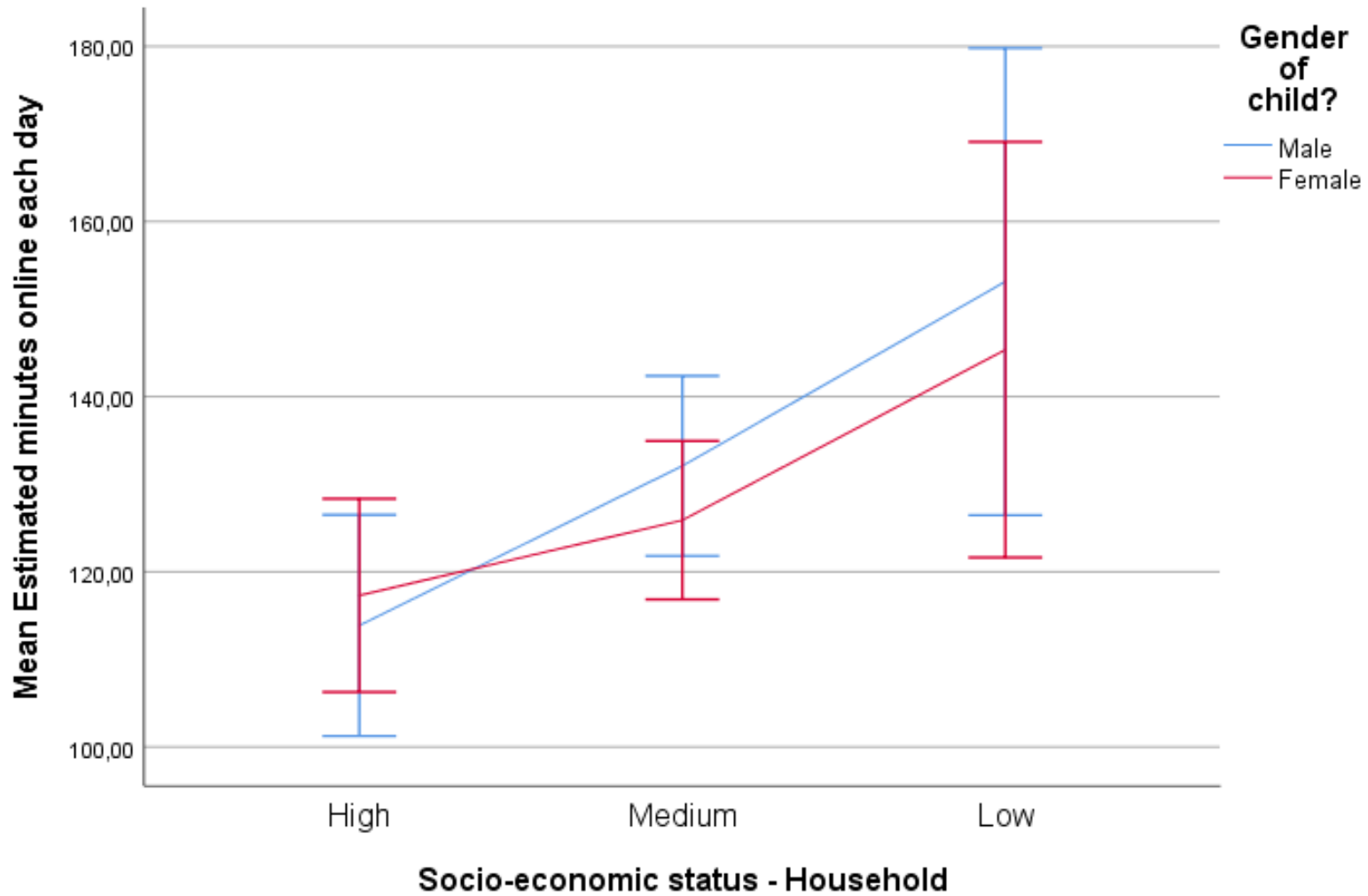
**SES:** Souvisí SES s frekvencí používání internetu?

**pohlaví:** Souvisí pohlaví s frekvencí používání internetu?

**interakce:** Má SES jinou souvislost s používáním internetu u chlapců než u dívek?

	Nízký SES	Střední SES	Vysoký SES
Chlapci	153	132	114
Dívky	145	126	117





Error bars: 95% CI

# Faktoriální ANOVA - předpoklady

Vše, co v případě one-way ANOVY

Pro každou kombinaci faktorů by  
měl být zastoupený dostatečný  
počet případů.

Lze posoudit na základě  
jednoduché kontingenční tabulky.

Počet případů	Nízký SES	Střední SES	Vysoký SES
Kluci	26	202	114
Holky	32	205	130

# Faktoriální ANOVA v SPSS

Analyze → Generalized Linear Model → Univariate...

celková vysvětlená variabilita ( $SS_M$ ) je rozsekána zvlášť pro jednotlivé faktory

Source	Type X Sum of Squares	df	Mean Square	F	Sig.
Corrected Model	$SS_M$	$df_M$	$MS_M$		
intercept					
Faktor1	$SS_{\text{Faktor1}}$	$df_{\text{Faktor1}}$	$MS_{\text{Faktor1}}$		
Faktor2	$SS_{\text{Faktor2}}$	$df_{\text{Faktor2}}$	$MS_{\text{Faktor2}}$		
Faktor1*Faktor2	$SS_{\text{interakce F1*F2}}$	$df_{\text{Int. F1*F2}}$	$MS_{\text{Int. F1*F2}}$		
Error	$SS_R$	$df_R$	$MS_R$		
Total					
Corrected Total	$SS_T$	$df_M+df_R$			

Každý faktor a interakce má vlastní statistiku F, proto lze posoudit, zda je signifikantním prediktorem závislé proměnné

Dependent Variable: DCtimeuse Estimated minutes online each day

Source	Type III Sum of Squares	df	Mean Square	F	Sig.
Corrected Model	65230,568 <sup>a</sup>	5	13046,114	2,802	0,016
Corrected Total	3338064,519	708			
DPSESHH3	60982,524	2	30491,262	6,549	0,002
QP201b	1155,893	1	1155,893	0,248	0,618
DPSESHH3 * QP201b	3890,106	2	1945,053	0,418	0,659
Error	3272833,951	703	4655,525		
Corrected Total	3338064,519	708			

a. R Squared = 0,020 (Adjusted R Squared = 0,013)

# Faktoriální ANOVA – reportování

---

Uvádíme zvlášť, jaký efekt měl každý faktor (main effect) nebo interakce faktorů:

$F(df_{Faktor}, df_R) = \dots, p = \dots, \text{parciální } \eta^2 \dots$

$$\text{parciální } \eta^2 = SS_{Faktor} / (SS_{Faktor} + SS_R)$$

$$*\text{parciální } \omega^2 : \omega_p^2 = \frac{df_{effect} \times (MS_{effect} - MS_{error})}{df_{effect} \times MS_{effect} + (N - df_{effect}) \times MS_{error}}$$

$$\hat{\omega}_p^2 = (F - 1) / (F + (df_{Error} + 1) / df_{Effect})$$

+ případné kontrasty a post-hoc testy jako u ANOVY

---



V některých situacích má smysl předpokládat, že je závislá proměnná ovlivňována nejen faktory, ale i intervalovými nezávislými proměnnými. Potřebujeme tedy model, který bude **kombinovat kategorické a intervalové nezávislé proměnné**.

---

Proč zavádět intervalové nezávislé do ANOVY:

snížíme množství nevysvětlené variability v modelu

kontrolujeme, zda není vliv faktorů zkreslen nějakou související intervalovou proměnnou

→ přesnější posouzení vlivu faktorů

**Příklad:** Používání internetu může souviset s věkem člověka. Pokud budeme tuto proměnnou kontrolovat, získáme představu o vlivu SES na frekvenci používání internetu, který je „očistěný“ od možného vlivu věku.

# ANCOVA (**a**nalysis of **cov**ariance)

ANOVA s jednou či více nezávislými intervalovými proměnnými (tzv. kovariáty)

---

zavádět jen kovariáty, pro které existují **dobré důvody** (nenacpat tam vše, co jsme měřili)

**dobře zvolené kovariáty** → zvýšení síly testu

- kovariát odebere část nevysvětlené variability ( $SS_R$ ) závislé proměnné, čímž se lépe projeví případný vliv faktorů

**špatně zvolené kovariáty** → snížení síly testu

- za každý přidaný kovariát ztrácíme jeden stupeň volnosti

uplatnění v **experimentálních designech**, kde chceme **statisticky kontrolovat** nežádoucí rozdíly mezi skupinami

uplatnění v **neexperimentálních designech**, kde chceme statisticky kontrolovat intervalové prediktory a posoudit tak nezkraslený vliv kategorických prediktorů



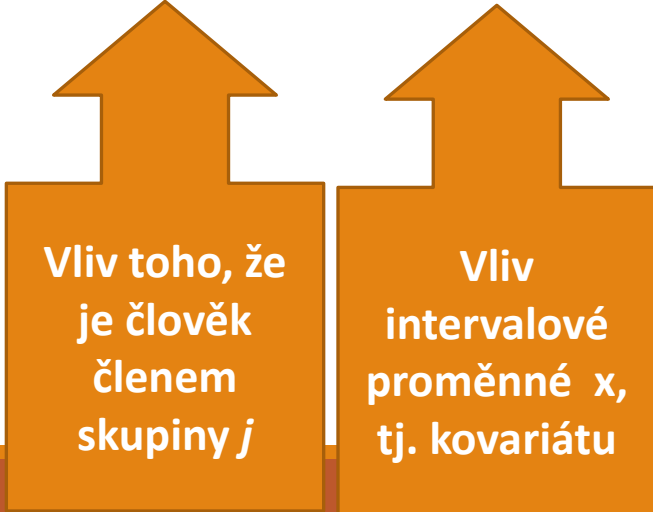
# One-way ANOVA

---

$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$$

## ANCOVA

$$Y_{ijk} = \mu + \alpha_j + \beta X_{ij} + \varepsilon_{ijk}$$



Vliv toho, že  
je člověk  
členem  
skupiny  $j$

Vliv  
intervalové  
proměnné  $x$ ,  
tj. kovariátu

# ANCOVA - předpoklady

---

Předpoklady ANOVY + předpoklady lineární regrese

Kovariát a faktor musí být nezávislé

- pokud nejsou, je obtížné interpretovat výsledky

Kovariát musí mít ve všech skupinách stejně silný vliv na závislou proměnnou (stejný regr. koef.)

- lze testovat zavedením interakce mezi faktorem a kovariátem do modelu (chceme, aby vyšla nesignifikantní)

# ANCOVA v SPSS

Analyze → Generalized Linear Model → Univariate...

celková vysvětlená variabilita ( $SS_M$ ) je rozsekána zvlášť pro kovariát(y) a faktor(y)

Source	Type X Sum of Squares	df	Mean Square	F	Sig.
Corrected Model	$SS_M$	$df_M$	$MS_M$		
intercept					
Kovariát1	$SS_{Kovariát1}$	$df_{Kovariát1}$	$MS_{Kovariát1}$		
Faktor1	$SS_{Faktor1}$	$df_{Faktor1}$	$MS_{Faktor1}$		
Error	$SS_R$	$df_R$	$MS_R$		
Total					
Corrected Total	$SS_T$	$df_M + df_R$			

můžeme si nechat zobrazit tzv. „**marginal means**“ (= jaké by byly skupinové průměry, kdyby se úroveň kovariátu nelišila napříč skupinami)

# ANCOVA – reportování

---

Uvádíme, jaký efekt měl každý kovariát:

$$F(df_{\text{Kovariát}}, df_R) = \dots, p = \dots, r = \dots$$

pro jednotlivé kovariáty vždy  $df_{\text{Kovariát}} = 1$

$$r = \text{odmocnina}[t^2 / (t^2 + df)]$$

A uvádíme, jaký efekt měl každý faktor:

$$F(df_{\text{Faktor}}, df_R) = \dots, p = \dots, \text{parciální } \eta^2 = \dots$$

$$\text{parciální } \eta^2 = SS_{\text{Faktor}} / (SS_{\text{Faktor}} + SS_R) \text{ lépe } \omega_p^2$$

+ případné kontrasty a post-hoc testy jako u ANOVY

# MANOVA (multivariační ANOVA)

---

ANOVA s více **závislými** intervalovými proměnnými

- posuzujeme vliv nezávislých proměnných na lineární kombinaci závislých proměnných
- pracujeme s multivariační obdobou F
- bereme v úvahu nejen (ne)vysvětlený rozptyl, ale i (ne)vysvětlenou kovarianci mezi závislými proměnnými

## výhody oproti sérii více ANOV

- kontrolujeme nárůst rizika chyby I. typu
- lze odhalit vztah ke kombinaci závislých proměnných

## nevýhody

- obtížná interpretace výsledků
- málokdy přinese nové informace oproti ANOVĚ
- vyžaduje splnění dalších předpokladů, které nelze jednoduše otestovat v SPSS (multivariační normalita)

---



# ANALÝZA ROZPTYLU PRO OPAKOVANÁ MĚŘENÍ

PSY252

Statistická analýza dat v psychologii II

---

# Opakovaná měření

---

Vnitrosubjektové a long designy

Sledujeme vývoj nějaké proměnné v čase

Vystavujeme jedince několika experimentálním podmínkám a hledáme rozdíl ve změně

Hledáme rozdíly v určitém znaku mezi příbuznými jedinci

**Výhoda:** větší síla, potřeba menšího vzorku

**Nevýhoda:** složitější statistika



ID	EDA klid	EDA stres1	EDA stres2
101A	1	2	3
102A	4	5	6
...			
199A	5	3	5

ID	Stres	EDA
101A	Klid	1
101A	Stres1	2
101A	Stres2	3
102A	Klid	4
102A	Stres1	5
102A	Stres2	6
...		
199A	Klid	5
199A	Stres1	3
199A	Stres2	5

---

Při opakovaných měřeních je porušen předpoklad ANOVA či lineární regrese o nezávislosti pozorování

funguje podobně jako faktoriální ANOVA

Nový předpoklad – sféricita (compound symmetry) –  
Mauchlyho test

- Splněna pokud **rozptyly** jednotlivých opakovaných měření jsou **stejné** a **kovariance** mezi jednotlivými opakovanými měřeními jsou **stejné**
- V longitudinálních designech obvykle problém – měření, která jsou si blízká v čase, obvykle více korelují
- Při nesplnění – korekce (G-G, H-F) či MANOVA

Méně spolehlivé post-hoc testy

# Dělení variability

---

Variabilita mezi subjekty – různí lidé mají různou průměrnou hodnotu závislé

Variabilita mezi měřeními (treatments) – rozdílnost průměrů měření

Chybový rozptyl – náhodná variabilita kolem hodnoty závislé predikované osobou a pořadím měření (treatmentem)

(Variabilita způsobená rozdílným efektem treatments na různé jedince)

# Příklad

---

EDA – elektrodermální aktivita (=pocení se)

3 úrovně stresu – klid, nekonfliktní Stroop, konfliktní Stroop – v tomto pořadí

„Soulad“ EDA na pravé a levé dlani

- Koeficient laterality (-30;30) (levopotivý – pravopotivý)
- PTI – synchronizace křivek pocení (0; 25)

Psychopatologie – BDI, SAS, TSC40

# Velikost účinku

---

U kontrastů můžeme počítat Cohenovo  $d$

- problematická je smysluplná volba SD, kterou bychom rozdíl průměrů standardizovali
- SD baseline měření
- střední SD napříč měřeními

Nebo můžeme spočítat  $r$  (F: s. 567)

$$r = \sqrt{\frac{F(1, df_R)}{F(1, df_R) + df_R}}$$

- velikost efektu „očistěnou“ o korelaci mezi měřeními nadhodnocenou
- vhodné pro usuzování na sílu testu

$\omega^2$  pro celý faktor F: s. 566

- nápověda  $SS_{\text{TOTAL}} = s^2 * (N-1)$

# Kontrasty a post-hoc testy

---

Kontrasty pro vnitrosubjektový faktor jako u faktoriální anovy.

- Transformation matrix v Options pro kontrolu

Post-hoc testy pro vnitrosubjektový a mezisubjektový faktor na jiných místech.

- Field: Vezměte na vědomí dopad odchylek od sféricity na platnost post-hoc testů

# Rozšíření

---

Faktoriální vnitrosubjektová/repeated Anova – více než 1 vnitrosubjektový faktor

Mixed ANOVA - kombinace vnitrosubjektových a mezisubjektových faktorů (tj. repeated+normální ANOVA)

# Nevýhody Repeated ANOVA

---

Požadavek sféricity - při výrazném nesplnění, či jiném očekávání je na místě hledat jiné modely

Neumí se vypořádat s chybějícími hodnotami

Flexibilní řešení obou problémů nabízí multi-level lineární modely (v SPSS Analyze ->Mixed models)