

PSY117/454

Statistická analýza dat v psychologii

Přednáška 12

Analýza rozptylu

Srovnávání více než dvou průměrů

If your experiment needs statistics, you ought to have done a better experiment.

Ernest Rutherford

Omezení t -testu (i jeho n Par alternativ)

t -test umožňuje srovnání pouze dvou průměrů

- Více skupin (j) \gg mnoho porovnání: $j(j-1)/2$

Více srovnání způsobuje strmý růst pravděpodobnosti chyby I. typu

- např. při $\alpha=0,05$ a 20 testech $p=0,64$ (1 nebo více chyb)
 - aplikace binomického rozložení
- Platí to pro jakékoli statistické testy (zejm. korelace)

Je *nevhodné* provádět velké množství testů na jedněch datech (cca >5)

- Zneužití se označuje jako rybaření v datech – capitalizing on chance
- Lze kompenzovat korekcí hladiny α (Bonferroniho korekce), avšak za cenu značného snížení síly testu ($1-\beta$).
 - Místo α testujeme na hladině $\alpha' = \alpha/N$, kde N je počet prováděných testů.

Řešení = Analýza rozptylu (ANOVA)

Testuje na více skupinách jen jednu hypotézu:

- Je někde mezi skupinovými průměry někde rozdíl?
 - Je mezi Pražáky, Brňáky a Ostraváky rozdíl v průměrné lakotě?
 - $H_0: \mu_{\text{Pražáci}} = \mu_{\text{Brňáci}} = \mu_{\text{Ostraváci}}$
- Je-li odpověď „**ano**“ ($p < \alpha$), pak se můžeme podívat na jednotlivé rozdíly detailněji (post-hoc testy)
- Je-li odpověď „**ne**“ ($p > \alpha$), pak bychom neměli (rybaření)

Terminologická vložka - ANOVA

- ANOVA = ANalysis Of Variance = analýza rozptylu
 - i přes svůj název jde o srovnávání **průměrů**
 - ANOVA zjišťuje vztah mezi **kategoriální nezávislou a intervalovou závislou**.
 - kategoriální nezávislá = **faktor** (factor, „-way“)
 - hodnoty kategoriální nez. = **úrovně** (level, treatment)
 - Zjištěný rozdíl = efekt, účinek (effect)
-

Princip ANOVY 1.

□ rozptyl = MS = mean square = $SS/df = (\sum(x-m))/(n-1)$

□ MS_{within} : variabilita uvnitř skupin ($MS_{e, error}$)

□ $MS_{within} = SS_{within}/n - j$

□ $SS_{within} = \sum_j \sum_i (x_i - m_j)^2$

□ $MS_{between}$: s^2 spočítaný ze skupinových průměrů, variabilita uvnitř skupiny je ignorována (též $MS_{A, B, treatment}$)

□ $MS_{between} = SS_{between}/j - 1$

□ $SS_{between} = \sum_j (n_j(m_j - m))^2$

Platí-li H_0 , jaký čekáme vztah mezi $MS_{between}$ a MS_{within} ?

	sk1	sk2	sk3	Celk.		sk1	sk2	sk3	Celkem
čl1	2	4	6		čl1	0	6	4	
čl2	2	4	6		čl2	4	2	8	
čl3	2	4	6		čl3	0	6	4	
čl4	2	4	6		čl4	4	2	8	
čl5	2	4	6		čl5	2	4	6	
m	2	4	6	4	m	2	4	6	4
s²	0	0	0	2,9	s²	4,0	4,0	4,0	6,3
			MS_{bg}	20				MS_{bg}	20
			MS_w	0				MS_w	4

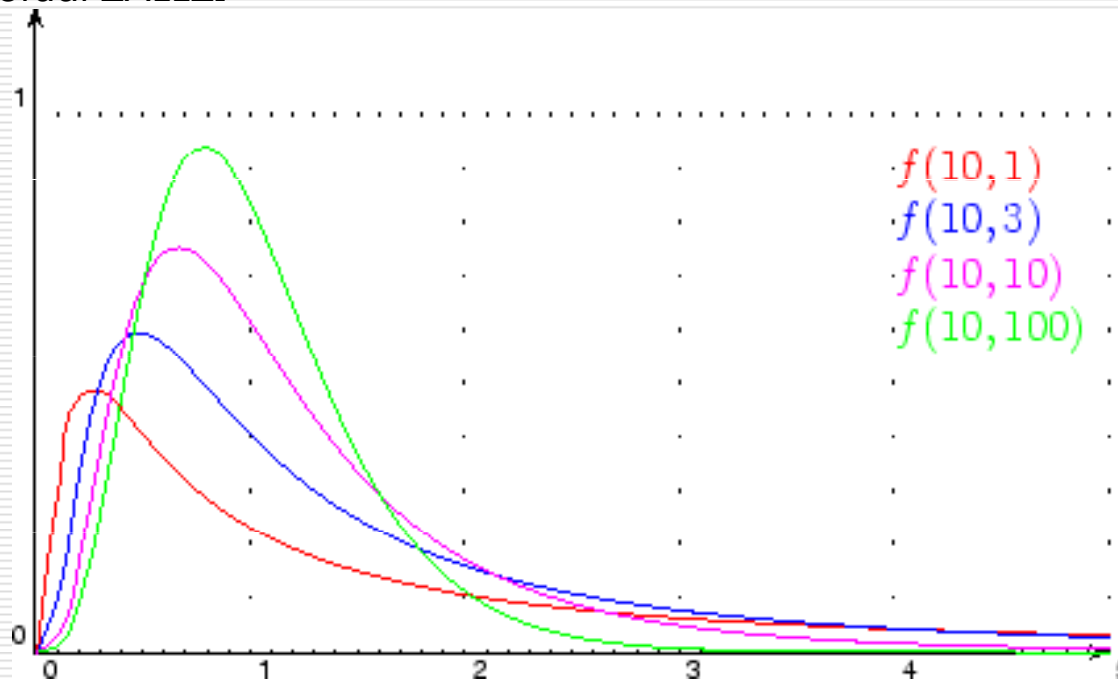
	sk1	sk2	sk3	Celkem		F	5
čl1	1	4	2				
čl2	3	5	5			$_{0,95}F(2,12)$	3,8853
čl3	5	1	3			p	0,0263
čl4	4	2	1				
čl5	2	3	4				
m	3	3	3	3			
s²	2,5	2,5	2,5	2,1			
			MS_{bg}	0			
			MS_w	2,5			

Princip ANOVY – F -test

- Čím jsou si průměry podobnější, tím je rozptyl mezi skupinami nižší (Platí-li H_0 , MS_{between} se blíží s^2)
 - Čím nižší je rozptyl uvnitř skupin (MS_{within} se blíží 0), tím průkaznější se průměry mezi skupinami zdají být.
 - Důležitý je **poměr těchto dvou odhadů rozptylu:**
$$F = \frac{MS_{\text{between}}}{MS_{\text{within}}}$$
 - Čím vyšší je F -poměr, tím průkaznější jsou rozdíly mezi skupinovými průměry (rozsah je 0 až ∞)
 - F -poměr má při platnosti H_0 jako výběrová statistika **F – rozložení** s $(df1, df2)$, které má průměr přibližně 1 (přesně $df2/(df2-2)$)
-

Fisherovo-Snedecorovo F -rozložení

- Podobně jako t -rozložení, je F -rozložení vlastně rodina mnoha rozložení mírně se lišící svým tvarem $(F(1; v) = t(v)^2)$
- Tato rozložení se liší tentokrát dvěma parametry – stupni volnosti
 - $v_1 = \text{počet skupin} - 1$: stupně volnosti čitatele - MS_{between}
 - $v_2 = \text{počet lidí} - \text{počet skupin}$: stupně volnosti jmenovatele - MS_{within}
 - na pořadí ZÁLEŽÍ



Princip ANOVY – dělení rozptylu.

- Dělení variability (rozptylu) podle zdrojů **jako u lineární regrese**

$$X_{ij} = \mu + \alpha_j + e_{ij}$$

$$Y_i = a + b_1X_1 + b_2X_2 + \dots + b_{j-1}X_{j-1} + e_i$$

- X_{ij} = skóre jedince (i -tý jedinec v j -té skupině)
- μ = průměr populace
- α = vliv příslušnosti ke skupině (vliv úrovně faktoru)
- e_{ij} = chyba (vše, s čím nepočítáme, individuální prom.)

$$X_{ij} - m = (m - m_j) + (X_{ij} - m_j)$$

odchylka od celkového průměru = odchylka od skupinového průměru +
odchylka skupinového průměru od celkového průměru

- ... odchylky umocněné na druhou = cesta k rozptylu

$$SS_{\text{Total}} = SS_{\text{Between (A, treatment)}} + SS_{\text{Within(Error)}}$$

$$MS_{\text{Total}}; MS_{\text{Error}}; MS_A$$

Velikost účinku (efektu)

- Podobně jako u regrese chceme vědět, jaká část rozptylu závislé je vysvětlená nezávislou
 - Ekvivalentem R^2 je u anovy η^2 (eta)
 - $\eta^2 = SS_{\text{Between}} / SS_{\text{Total}}$
 - Poněkud přesnější je $\omega^2 = (SS_{\text{Between}} - df_{\text{Between}} \cdot MS_{\text{Within}}) / (SS_{\text{Total}} + MS_{\text{Within}})$
 - Pro konkrétní rozdíl průměrů $d_{\text{Coh}} = m_1 - m_2 / \sqrt{MS_{\text{Within}}}$
 - Velikost účinku je vždy třeba uvádět
-

Předpoklady použití ANOVY

- normální rozložení uvnitř skupin
 - při $n_j > 30$ a $n_1 = n_2 = \dots = n_j$ je ANOVA robustní
- stejné rozptyly uvnitř skupin:
homoskedascita
 - do $s_{\max}/s_{\min} < 3$ je ANOVA robustní, zvláště při $n_1 = n_2 = \dots = n_j$
- nezávislost všech pozorování
 - při opakovaných měřeních je třeba použít ANOVU pro opakovaná měření

viz Hendl 343

Post-hoc testy (simultánní porovnávání)

- Po (a pouze po) prokázání „nějakých“ rozdílů mezi průměry obvykle chceme vědět, mezi kterými skupinami konkrétně rozdíly jsou: **post-hoc testy**
 - Srovnáváme každou skupinu s každou způsobem, který nezpůsobí nárůst α .
 - Je-li důležité udržet α pod kontrolou, pak je správnou volbou **Scheffe**ho test – volba pro *rybaření*
 - Pokud to není tak kritické a máte-li pár *kvazi*-hypotéz na mysli, pak je volbou **Student-Neuman-Keuls (S-N-K)**
 - Extrémně „dajný“ a nepříliš vhodný pro více než 3 skupiny je **LSD** a proto se nedoporučuje.
-

Další varianty a rozšíření ANOVA

- ANOVA pro opakovaná měření (jako párový t -test)
- ANOVA s 2 a více faktory (faktoriální ANOVA)
- MANOVA – s více závislými proměnnými

To vše v SPSS skryto pod GLM – general linear model

- Pořadovou (neparametrickou) alternativou ANOVY jsou
 - **Kruskal-Wallis H:** $H_0: Md_1 = Md_2 = \dots = Md_j$ $H_1: Md_1 \neq Md_2 \neq \dots \neq Md_j$
 - **Jonckheere-Terpstra Test:** $H_1: Md_1 \leq Md_2 \leq \dots \leq Md_j$