

## Kapitola 6.: Analýza rozptylu jednoduchého třídění

### 6.1. Motivace

Zajímáme se o problém, zda lze určitým faktorem (tj. nominální náhodnou veličinou A) vysvětlit variabilitu pozorovaných hodnot náhodné veličiny X, která je intervalového či poměrového typu. Např. zkoumáme, zda metoda výuky určitého předmětu (faktor A) ovlivňuje počet bodů dosažených studenty v závěrečném testu (náhodná veličina X).

Předpokládáme, že faktor A má  $r \geq 3$  úrovně a  $i$ -té úrovni odpovídá  $n_i$  výsledků  $X_{i1}, \dots, X_{in_i}$ , které tvoří náhodný výběr z rozložení  $N(\mu_i, \sigma^2)$ ,  $i = 1, \dots, r$  a jednotlivé náhodné výběry jsou stochasticky nezávislé, tedy  $X_{ij} = \mu_i + \varepsilon_{ij}$ , kde  $\varepsilon_{ij}$  jsou stochasticky nezávislé náhodné veličiny s rozložením  $N(0, \sigma^2)$ ,  $i = 1, \dots, r, j = 1, \dots, n_i$ .

Výsledky lze zapsat do tabulky

faktor A	výsledky
úroveň 1	$X_{11}, \dots, X_{1n_1}$
úroveň 2	$X_{21}, \dots, X_{2n_2}$
...	...
úroveň r	$X_{r1}, \dots, X_{rn_r}$

Na hladině významnosti  $\alpha$  testujeme nulovou hypotézu, která tvrdí, že všechny střední hodnoty jsou stejné proti alternativní hypotéze, která tvrdí, že aspoň jedna dvojice středních hodnot se liší. Jedná se tedy o zobecnění dvouvýběrového t-testu a na první pohled se zdá, že stačí utvořit  $\binom{r}{2}$  dvojic náhodných výběrů a na každou dvojici aplikovat dvouvýběrový t-test.

Tento postup však nelze použít, neboť nezaručuje splnění podmínky, že pravděpodobnost chyby 1. druhu je  $\alpha$ . Proto ve 30. letech 20. století vytvořil R. A. Fisher metodu ANOVA (analýza rozptylu, v popsané situaci konkrétně analýza rozptylu jednoduchého třídění), která uvedenou podmínku splňuje.

Pokud na hladině významnosti  $\alpha$  zamítneme nulovou hypotézu, zajímá nás, které dvojice středních hodnot se od sebe liší. K řešení tohoto problému slouží metody mnohonásobného porovnávání, např. Scheffého nebo Tukeyova metoda.

### 6.2. Označení

V analýze rozptylu jednoduchého třídění se používá následující označení.

$$n = \sum_{i=1}^r n_i \dots \text{celkový rozsah všech } r \text{ výběrů}$$

$$X_{i.} = \sum_{j=1}^{n_i} X_{ij} \dots \text{součet hodnot v } i\text{-tém výběru}$$

$$M_{i.} = \frac{1}{n_i} X_{i.} \dots \text{výběrový průměr v } i\text{-tém výběru}$$

$$X_{..} = \sum_{i=1}^r \sum_{j=1}^{n_i} X_{ij} \dots \text{součet hodnot všech výběrů}$$

$$M_{..} = \frac{1}{n} X_{..} \dots \text{celkový průměr všech } r \text{ výběrů}$$

### 6.3. Testování hypotézy o shodě středních hodnot

Náhodné veličiny  $X_{ij}$  se řídí modelem

$M_0: X_{ijk} = \mu + \alpha_i + \varepsilon_{ij}$  pro  $i = 1, \dots, r, j = 1, \dots, n_i$ , přičemž  $\varepsilon_{ij}$  jsou stochasticky nezávislé náhodné veličiny s rozložením  $N(0, \sigma^2)$ ,  $\mu$  je společná část střední hodnoty závisle proměnné veličiny,  $\alpha_i$  je efekt faktoru A na úrovni i.

Parametry  $\mu, \alpha_i$  neznáme. Požadujeme, aby platila tzv. reparametrizační rovnice:

$$\sum_{i=1}^r \alpha_i = 0.$$

Zavedeme součty čtverců

$S_T = \sum_{i=1}^r \sum_{j=1}^{n_i} (X_{ij} - M_{..})^2$  ... celkový součet čtverců (charakterizuje variabilitu jednotlivých pozorování kolem celkového průměru), má počet stupňů volnosti  $f_T = n - 1$ ,

$S_A = \sum_{i=1}^r n_i (M_{i..} - M_{..})^2$  ... skupinový součet čtverců (charakterizuje variabilitu mezi jednotlivými náhodnými výběry), má počet stupňů volnosti  $f_A = r - 1$ ,

$S_E = \sum_{i=1}^r \sum_{j=1}^{n_i} (X_{ij} - M_{i.})^2$  ... reziduální součet čtverců (charakterizuje variabilitu uvnitř jednotlivých výběrů), má počet stupňů volnosti  $f_E = n - r$ .

Lze dokázat, že  $S_T = S_A + S_E$ .

Sčítanec  $(M_{i.} - M_{..})$  představuje bodový odhad efektu  $\alpha_i$ .

Kdyby nezáleželo na faktoru A, platila by hypotéza  $\alpha_1 = \dots = \alpha_r = 0$  a dostali bychom model

$M_1: X_{ij} = \mu + \varepsilon_{ij}$ .

Rozdíl mezi modely  $M_0$  a  $M_1$  ověříme pomocí testové statistiky

$F_A = \frac{S_A / f_A}{S_E / f_E}$ , která se řídí rozložením  $F(r-1, n-r)$ , je-li model  $M_1$  správný. Hypotézu o nevý-

znamnosti faktoru A tedy zamítneme na hladině významnosti  $\alpha$ , když platí:

$$F_A \geq F_{1-\alpha}(r-1, n-r).$$

Výsledky výpočtů zapisujeme do tabulky analýzy rozptylu jednoduchého třídění.

Zdroj variability	součet čtverců	stupně volnosti	podíl	F
skupiny	$S_A$	$f_A = r-1$	$S_A/f_A$	$\frac{S_A/f_A}{S_E/f_E}$
reziduální	$S_E$	$f_E = n-r$	$S_E/f_E$	-
celkový	$S_T$	$f_T = n-1$	-	-

### 6.4. Testy shody rozptylů

Před provedením analýzy rozptylu je zapotřebí ověřit předpoklad o shodě rozptylů v daných r výběrech.

### 6.4.1. Levenův test

Položme  $Z_{ij} = |X_{ij} - M_i|$ . Označíme

$$M_{Zi} = \frac{1}{n_i} \sum_{j=1}^{n_i} Z_{ij}, M_Z = \frac{1}{n} \sum_{i=1}^r \sum_{j=1}^{n_i} Z_{ij}, S_{ZE} = \sum_{i=1}^r \sum_{j=1}^{n_i} (Z_{ij} - M_{Zi})^2, S_{ZA} = \sum_{i=1}^r n_i (M_{Zi} - M_Z)^2.$$

Platí-li hypotéza o shodě rozptylů, pak statistika  $F_Z = \frac{S_{ZA}/(r-1)}{S_{ZE}/(n-r)} \sim F(r-1, n-r)$ .  $H_0$  tedy zamítáme na hladině významnosti  $\alpha$ , když  $F_Z \geq F_{1-\alpha}(r-1, n-r)$ .

### 6.4.2. Bartlettův test

Platí-li hypotéza o shodě rozptylů, pak statistika

$$B = \frac{1}{C} \left[ (n-r) \ln S_*^2 - \sum_{i=1}^r (n_i - 1) \ln S_i^2 \right] \text{ má přibližně rozložení } \chi^2(r-1), \text{ kde}$$

$$C = 1 + \frac{1}{3(r-1)} \left( \sum_{i=1}^r \frac{1}{n_i - 1} - \frac{1}{n-r} \right), S_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (X_{ij} - M_i)^2 \text{ je výběrový rozptyl } i\text{-tého vý-}$$

běru,  $S_*^2 = \frac{\sum_{i=1}^r (n_i - 1) S_i^2}{n-r} = \frac{S_E}{n-r}$  je vážený průměr výběrových rozptylů.

$H_0$  zamítáme na přibližné hladině významnosti  $\alpha$ , když  $B \geq \chi^2_{1-\alpha}(r-1, n-r)$ . Bartlettův test lze použít, pokud rozsahy všech výběrů jsou aspoň 7.

## 6.5. Metody mnohonásobného porovnávání

Zamítneme-li na hladině významnosti  $\alpha$  hypotézu o shodě středních hodnot, chceme zjistit, které dvojice středních hodnot se liší na dané hladině významnosti  $\alpha$ .

### 6.5.1. Tukeyova metoda

Mají-li všechny výběry týž rozsah  $p$  (říkáme, že třídění je vyvážené), použijeme Tukeyovu metodu: rovnost středních hodnot  $\mu_k$  a  $\mu_l$  zamítneme na hladině významnosti  $\alpha$ , když

$$|M_k - M_l| \geq q_{1-\alpha}(r, n-r) \frac{S_*}{\sqrt{p}}, \text{ kde hodnoty } q_{1-\alpha}(r, n-r) \text{ jsou kvantily studentizovaného rozpětí}$$

a najdeme je ve statistických tabulkách.

### 6.5.2. Scheffého metoda

Nemají-li všechny výběry stejný rozsah, použijeme Scheffého metodu: rovnost středních hodnot  $\mu_k$  a  $\mu_l$  zamítneme na hladině významnosti  $\alpha$ , když

$$|M_k - M_l| \geq S_* \sqrt{(r-1) \left( \frac{1}{n_k} + \frac{1}{n_l} \right) F_{1-\alpha}(r-1, n-r)}.$$

Pozor, může nastat situace, kdy při zamítnutí  $H_0$  nenajdeme významný rozdíl u žádné dvojice středních hodnot. Pak je významně rozdílná některá složitější kombinace středních hodnot, tzv. kontrast.

## 6.6. Příklad

U čtyř odrůd brambor (označených symboly A, B, C, D) se zjišťovala celková hmotnost brambor vyrostlých vždy z jednoho trsu. Výsledky (v kg):

odrůda	hmotnost
A	0,9 0,8 0,6 0,9
B	1,3 1,0 1,3
C	1,3 1,5 1,6 1,1 1,5
D	1,1 1,2 1,0

Na hladině významnosti 0,05 testujte hypotézu, že střední hodnota hmotnosti trsu brambor nezávisí na odrůdě. Zamítnete-li nulovou hypotézu, zjistěte, které dvojice odrůd se liší na hladině významnosti 0,05.

### Řešení

Data považujeme za realizace čtyř nezávislých náhodných výběrů ze čtyř normálních rozložení se stejným rozptylem. Testujeme hypotézu, že všechny čtyři střední hodnoty jsou stejné.

$M_1 = 0,8$ ,  $M_2 = 1,2$ ,  $M_3 = 1,4$ ,  $M_4 = 1,1$ ,  $M_{\cdot} = 1,14$ ,  $S_E = 0,3$ ,  $S_A = 0,816$ ,  $S_T = 1,116$ ,  $F = 9,97$ ,  $F_{0,95}(3,11) = 3,59$ . Protože testová statistika se realizuje v kritickém oboru,  $H_0$  zamítáme na hladině významnosti 0,05.

Výsledky zapíšeme do tabulky ANOVA

Zdroj variability	Součet čtverců	Stupně volnosti	podíl	F
skupiny	$S_A = 0,816$	$f_A = 3$	$S_A/3 = 0,272$	$\frac{S_A/f_A}{S_E/f_E} = 9,97$
reziduální	$S_E = 0,3$	$f_E = 11$	$S_E/11 = 0,02727$	-
celkový	$S_T = 1,116$	$f_T = 14$	-	-

Nyní pomocí Scheffého metody zjistíme, které dvojice odrůd se liší na hladině významnosti 0,05.

Srovnávané odrůdy	Rozdíly $ M_{k_{\cdot}} - M_{l_{\cdot}} $	Pravá strana vzorce
A, B	0,4	0,41
A, C	0,67	0,36
A, D	0,3	0,41
B, C	0,2	0,40
B, D	0,1	0,44
C, D	0,3	0,40

Na hladině významnosti 0,05 se liší odrůdy A a C.

## 6.7. Význam předpokladů v analýze rozptylu

- Nezávislost jednotlivých náhodných výběrů – velmi důležitý předpoklad, musí být splněn, jinak dostaneme nesmyslné výsledky.
- Normalita – ANOVA není příliš citlivá na porušení normality, zvláště pokud mají všechny výběry rozsah nad 20 (důsledek centrální limitní věty). Při výraznějším porušení normality se doporučuje Kruskalův – Wallisův test.
- Shoda rozptylů – mírné porušení nevádí, při větším se doporučuje Kruskalův – Wallisův test. Test shody rozptylů má smysl provádět až po ověření předpokladu normality.

### Kontrolní otázky

- Jaký problém řeší analýza rozptylu jednoduchého třídění?
- Jak je definován celkový, skupinový a reziduální součet čtverců a co tyto součty čtverců charakterizují?
- Popište způsob testování hypotézy o shodě středních hodnot.
- Jak se testuje hypotéza o shodě rozptylů?
- Které metody mnohonásobného porovnání se používají v analýze rozptylu jednoduchého třídění?
- Pojednejte o významu předpokladů v analýze rozptylu jednoduchého třídění.

### Příklady

1. Jsou známy měsíční tržby (v tisících Kč) tří prodavačů za dobu půl roku.

- prodavač: 12 10 9 10 11 9
- prodavač: 10 12 11 12 14 13
- prodavač: 19 18 16 16 17 15

Na hladině významnosti 0,05 testujte hypotézu, že střední hodnoty tržeb všech tří prodavačů jsou stejné. Pokud zamítneme nulovou hypotézu, zjistěte, tržby kterých dvou prodavačů se liší na hladině významnosti 0,05.

Výsledek:

$M_1 = 10,17$ ,  $M_2 = 12$ ,  $M_3 = 16,83$ ,  $M_{..} = 13$ ,  $S_E = 27,7$ ,  $S_A = 142,3$ ,  $S_T = 170$ ,  $F = 38,58$ ,  $F_{0,975}(2,015) = 3,6823$ ,  $H_0$  tedy zamítáme na hladině významnosti 0,05.

Výsledky zapíšeme do tabulky ANOVA

Zdroj variability	Součet čtverců	Stupně volnosti	podíl	F
skupiny	$S_A = 142,3$	$f_A = 2$	$S_A/f_A = 71,17$	$\frac{S_A/f_A}{S_E/f_E} = 38,58$
reziduální	$S_E = 27,7$	$f_E = 15$	$S_E/f_E = 1,84$	-
celkový	$S_T = 170$	$f_T = 17$	-	-

Nyní pomocí Tukeyovy metody zjistíme, které dvojice prodavačů se liší na hladině významnosti 0,05.

Srovnávání prodavači	Rozdíly $ M_k - M_l $	Pravá strana vzorce
1, 2	1,83	2,03
1, 3	6,67*	2,03
2, 3	4,83*	2,03

Pravá strana:  $q_{1-\alpha}(r, n-r) \frac{S_*}{\sqrt{p}} = q_{0,95}(3,15) \frac{\sqrt{1,84}}{\sqrt{6}} = 4,83 \frac{\sqrt{1,84}}{\sqrt{6}} = 2,03$ , kde  $S_*^2 = \frac{S_E}{n-r} = 1,84$

Na hladině významnosti 0,05 se liší tržby prodavačů 1, 3 a 2, 3.

2. Je dáno pět nezávislých náhodných výběrů o rozsazích 5, 7, 6, 8, 5, přičemž  $i$ -tý výběr pochází z rozložení  $N(\mu_i, \sigma^2)$ ,  $i = 1, \dots, 5$ . Byl vypočten celkový součet čtverců  $S_T = 15$  a reziduální součet čtverců  $S_E = 3$ . Na hladině významnosti 0,05 testujte hypotézu o shodě středních hodnot.

Výsledek:

$$n = 5 + 7 + 6 + 8 + 5 = 31, r = 5, S_A = S_T - S_E = 15 - 3 = 12$$

$$F = \frac{S_A/(r-1)}{S_E/(n-r)} = \frac{12/4}{3/26} = 26, F_{0,95}(4,26) = 2,9752$$

Protože  $F \geq F_{0,95}(4,26)$ ,  $H_0$  zamítáme na hladině významnosti 0,05.

3. Je dána neúplná tabulka ANOVA. Místo otazníků doplňte chybějící čísla.

zdroj variability	součet čtverců	stupně volnosti	podíl	F
skupiny	?	2	?	?
reziduální	16,033	?	?	-
celkový	17,301	35	-	-

Výsledek:

zdroj variability	součet čtverců	stupně volnosti	podíl	F
skupiny	<b>1,268</b>	2	<b>0,634</b>	<b>1,304</b>
reziduální	16,033	<b>33</b>	<b>0,486</b>	-
celkový	17,301	35	-	-

4. (S) Studenti byli vyučováni předmětu za využití pěti pedagogických metod: tradiční způsob, programová výuka, audioteknika, audiovizuální technika a vizuální technika. Z každé skupiny byl vybrán náhodný vzorek studentů a všichni byli podrobni témuž písemnému testu. Na hladině významnosti 0,05 testujte hypotézu, že znalosti všech studentů jsou stejné a nezávisí na použité pedagogické metodě. V případě zamítnutí nulové hypotézy zjistěte, které výběry se liší na hladině významnosti 0,05.

metoda	počet bodů					
tradiční	76,2	48,3	85,1	63,7	91,6	87,2
programová	85,2	74,3	76,5	80,3	67,4	67,9 72,1 60,4
audio	67,3	60,1	55,4	72,3	40	
audiovizuální	75,8	81,6	90,3	78	67,8	57,6
vizuální	50,5	70,2	88,8	67,1	77,7	73,9

Výsledek:

Všech pět náhodných výběrů má rozložení blízké normálnímu rozložení. Levenův test shody rozptylů má testové kritérium 0,819, počet stupňů volnosti je 4 a 26, odpovídající  $p$ -hodnota je 0,5248, tedy na hladině významnosti 0,05 hypotézu o shodě rozptylů nezamítáme. Analýza rozptylu má testové kritérium 1,6236, počet stupňů volnosti je 4 a 26, odpovídající  $p$ -hodnota je 0,1983, tedy na hladině významnosti 0,05 hypotézu o shodě středních hodnot nezamítáme. Znamená to, že na hladině významnosti 0,05 se neprokázaly odlišnosti ve znalostech studentů.

5. (S) Pan Novák může cestovat z místa bydliště do místa pracoviště třemi různými způsoby: tramvají (způsob A), autobusem (způsob B) a metrem s následným přestupem na tramvaj (způsob C). Máme k dispozici jeho naměřené časy cestování do práce v době ranní špičky (včetně čekání na příslušný spoj) v minutách.

Způsob A: 32, 39, 42, 37, 34, 38

Způsob B: 30, 34, 28, 26, 32

Způsob C: 40, 37, 31, 39, 38, 33, 34

Pro všechny tři způsoby dopravy vypočtěte průměrné časy cestování. Na hladině významnosti 0,05 testujte hypotézu, že doba cestování do práce nezávisí na způsobu dopravy. V případě zamítnutí nulové hypotézy zjistěte, které způsoby dopravy do práce se od sebe liší na hladině významnosti 0,05.

Výsledek:

Průměrné časy cestování pro tři způsoby dopravy jsou 37 min, 30 min, 36 min.

Všechny tři náhodné výběry mají rozložení blízké normálnímu rozložení. Levenův test shody rozptylů má testové kritérium 0,1054, počet stupňů volnosti je 2 a 15, odpovídající p-hodnota je 0,9007, tedy na hladině významnosti 0,05 hypotézu o shodě rozptylů nezamítáme. Analýza rozptylu má testové kritérium 6,7151, počet stupňů volnosti je 2 a 15, odpovídající p-hodnota je 0,0083, tedy na hladině významnosti 0,05 hypotézu o shodě středních hodnot zamítáme.

Scheffého metoda mnohonásobného porovnávání prokázala na hladině významnosti 0,05 rozdíl mezi způsoby A a B a mezi způsoby Ba C.