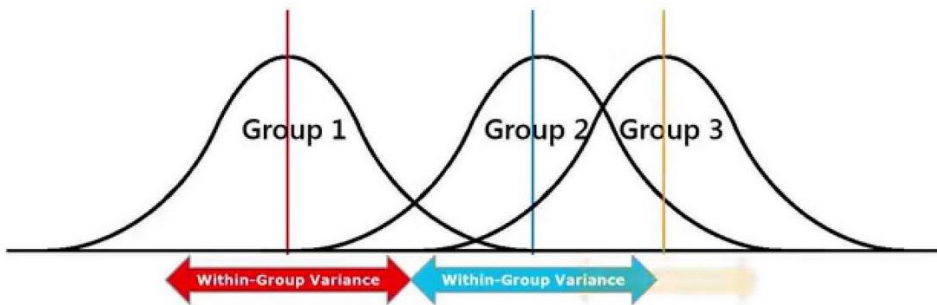


ضروریات آمار برای علم داده

تحلیل واریانس

تحلیل واریانس



- ❑ واریانس کل داده‌ها براساس یک یا چند متغیر عامل به دو یا چند بخش تفکیک شده است.
- ❑ براساس آزمون‌های مربوط به واریانس می‌توان همگون یا ناهمگون بودن گروه‌ها را آزمود.

شرایط :

- نمونه‌های حاصل از هر گروه یا جامعه، تصادفی و از توزیع نرمال باشند.
- واریانس‌های گروه‌ها برابر ولی نامشخص هستند. (آزمون فلینجر)
- گروه‌ها از یکدیگر مستقل هستند.

تحليل واريانس

$$\begin{cases} H_0 : \mu_1 = \mu_2 = \dots = \mu_k \\ H_1 : \sim H_0 \end{cases}$$

$$\begin{cases} H_0 : \text{میانگین همه گروه‌ها یکسان است} \\ H_1 : \text{حداقل میانگین یکی از گروه‌ها با بقیه متفاوت است} \end{cases}$$

$$F = \frac{MS_{between}}{MS_{within}} \sim F_{a-1, a(n-1)}$$

تحلیل واریانس یک طرفه

تنها یک عامل در گروه‌های مختلف بررسی می‌شود.

مثال:

یک مدیر ارشد فروشگاه‌های زنجیره‌ای ادعا نموده است که خلاقیت‌های متفاوت مدیران، اثرهای متفاوتی بر متوسط میزان فروش می‌گذارد. وی برای تعیین صحت و سقم ادعای خود، ۵ تا از فروشگاه‌هایی که به جز مدیر آن دارای شرایط یکسانی بوده‌اند انتخاب و فروش (به میلیون تومان) هفت روز هفته‌ای را که به تصادف انتخاب کرده است ثبت نموده است. چنانچه داده‌ها از توزیع نرمال پیروی کنند، آن‌ها را در سطح معنی داری ۵ درصد تحلیل نمایید.

مدیر منابع انسانی کارخانه‌ای ادعا کرده است که کارایی سبک اجرای قوانین کار در کارخانه‌ها به طور متوسط با یکدیگر متفاوت است. وی برای اثبات ادعای خود نمونه‌ای از سه نوع سبک اجرای مختلف از یک قانون کار را مشخص و کارایی آن‌ها را در موارد متعددی که در کارخانه‌های مختلف اجرا شده است به وسیله تخصیص نمره به آن‌ها اندازه‌گیری نموده است. چنانچه نمره‌ها از توزیع نرمال پیروی کنند، داده‌ها را در سطح معنی داری ۵ درصد تحلیل کنید.

میزان فروش روزانه‌ی یک هفته									
میانگین ($\bar{y}_{i.}$)	جمع ($y_{i.}$)	۷	۶	۵	۴	۳	۲	۱	
$\bar{y}_1 = 5/43$	$y_1 = 38$	۵	۵	۴	۳	۷	۹	۵	۱
$\bar{y}_2 = 2/186$	$y_2 = 20$	۳	۲	۳	۴	۳	۲	۳	۲
$\bar{y}_3 = 7/57$	$y_3 = 53$	۷	۷	۹	۷	۶	۸	۹	۳
$\bar{y}_4 = 4/71$	$y_4 = 33$	۴	۵	۴	۶	۵	۵	۴	۴
$\bar{y}_5 = 6/29$	$y_5 = 44$	۶	۷	۵	۷	۶	۶	۷	۵
جمع کل									
$y_{..} = 188$									
میانگین کل									
$\bar{y}_{..} = 5/37$									

خلاقیات
مدیران

منبع تغییر پذیری SV	مجموع توان های دوم SS	درجه ی آزادی DF	میانگین توان های دوم MS	مقدار آماره ی آزمون F_0
تیمار (A)	$SSA = \sum_{i=1}^{a=3} \frac{y_{i.}^2}{n_i} - \frac{y_{..}^2}{N}$ $= \left(\frac{17^2}{3} + \frac{18^2}{2} + \frac{14^2}{4} \right) - \frac{49^2}{9}$ $= 40/55$	$a-1 = 3-1$ $= 2$	$MSA = \frac{SSA}{a-1}$ $= \frac{40/55}{2} = 20/275$	$F_0 = \frac{MSA}{MSE}$ $= \frac{20/275}{2/61} = 7/77$
خطا (E)	$SSE = SST - SSA$ $= 56/23 - 40/55 = 15/68$	$N - a = 9 - 3$ $= 6$	$MSE = \frac{SSE}{N - a}$ $= \frac{15/68}{6} = 2/61$	
کل (T)	$SST = \sum_{i=1}^{a=3} \sum_{j=1}^{n_i} y_{ij}^2 - \frac{y_{..}^2}{N}$ $= (1^2 + 5^2 + \dots + 5^2) - \frac{49^2}{9}$ $= 323 - 266/77 = 56/23$	$N - 1 = 9 - 1$ $= 8$		

y	تکرار				
	۴	۳	۲	۱	
$y_1 = 17$		۴	۵	۸	۱
$y_2 = 18$			۸	۱۰	۲
$y_3 = 14$	۵	۲	۳	۴	۳
$y_{..} = 49$					

سبک اجرای
کار

منبع تغییر پذیری (SV)	مجموع توان‌های دوم (SS)	درجه‌ی آزادی (DF)	میانگین توان‌های دوم (MS)	مقدار آماره‌ی آزمون F_0
--------------------------	----------------------------	----------------------	------------------------------	------------------------------

$$F_0 = \frac{MSA}{MSE}$$

$$= \frac{21/76}{1/37}$$

$$= 15/88$$

$$MSA = \frac{SSA}{a-1}$$

$$= \frac{87/03}{4}$$

$$= 21/76$$

$$a-1$$

$$= 5-1$$

$$= 4$$

$$SSA = \frac{1}{n} \sum_{i=1}^{a=5} y_i^2 - \frac{y_{..}^2}{an}$$

$$= \frac{1}{7} \left(38^2 + 20^2 + 53^2 + 33^2 + 44^2 \right) - \frac{188^2}{5 \times 7}$$

$$= 87/03$$

خلاقیت مدیران
(A)

$$MSE = \frac{SSE}{a(n-1)}$$

$$= \frac{41/14}{30}$$

$$= 1/37$$

$$a(n-1)$$

$$= 5(7-1)$$

$$= 30$$

$$SSE = SST - SSA$$

$$= 128/17 - 87/03$$

$$= 41/14$$

خطا (E)

$$an-1$$

$$= 5 \times 7 - 1$$

$$= 34$$

$$SST = \sum_{i=1}^{a=5} \sum_{j=1}^{n=7} y_{ij}^2 - \frac{y_{..}^2}{an}$$

$$= (5^2 + 9^2 + \dots + 6^2) - \frac{188^2}{5 \times 7}$$

$$= 128/17$$

کل (T)

تحلیل واریانس دو طرفه

اثر دو عامل و همچنین اثر متقابل آن‌ها در گروه‌های مختلف آزمون می‌شود.

مثال

فرض کنید با یک مجموعه داده با ۴۸ سطر و ۳ ستون (متغیر) مواجه هستیم که مربوط به زمان اثر سم به یک نوع خوک با توجه به نوع سم و نحوه درمان خوک‌ها است. می‌خواهیم بدانیم که آیا نوع سم در کاهش طول عمر خوک‌ها موثر است یا خیر. اگر شیوه درمان هم به عامل نوع سم اضافه کنیم یعنی می‌خواهیم تحت نوع سم، شیوه درمان و اثر متقابل شیوه درمان و نوع سم، میانگین طول عمر خوک‌ها را تحت تاثیر قرار می‌دهد.

آزمون تعقیبی (Post Hock)

تحلیل واریانس نشان می‌دهد که آیا نمونه‌ها متعلق به جامعه هستند یا خیر. در صورتی که فرض صفر رد شود، معلوم نیست که کدام یک از نمونه‌ها در کدام جامعه قرار دارند. به عبارت دیگر، معنی‌دار شدن نسبت F به ما نمی‌گوید که اختلاف بین کدام جفت از میانگین‌ها معنی‌دار است. بلکه با آماره F ، تنها می‌توانیم پی ببریم که اختلاف بین میانگین گروه‌ها معنی‌دار است.

برابری واریانس‌ها

آزمون لون

$$\begin{cases} H_0 : \sigma_1^2 = \sigma_2^2 \\ H_1 : \sigma_1^2 \neq \sigma_2^2 \end{cases}$$

$$W = \frac{(N - k)}{(k - 1)} \cdot \frac{\sum_{i=1}^k N_i (Z_{i\cdot} - Z_{\cdot\cdot})^2}{\sum_{i=1}^k \sum_{j=1}^{N_i} (Z_{ij} - Z_{i\cdot})^2} \quad W > F(\alpha; k - 1, N - k), \text{ Reject } H_0$$

ناپارامتری

آزمون کروسکال والیس، آزمونی ناپارامتری برای تحلیل واریانس یک طرفه است
فریدمن، آزمونی ناپارامتری برای تحلیل واریانس دو طرفه است