

۱-۶ پارامترهای پراکندگی

شاخص های مرکزی تنها یک منطقه را به عنوان محل تمرکز داده ها معرفی می کنند حال آنکه ممکن است دو دسته داده با پراکندگی های متفاوت دارای میانگین برابری باشند به عبارتی نیاز به شاخصی برای نمایش میزان پراکندگی داده ها خواهیم داشت که در ادامه به معرفی این پارامترها می پردازیم.

۱- دامنه داده‌ها: حاصل تفاضل بزرگترین داده از کوچکترین داده را R یا دامنه داده‌ها گویند.

$$R = X_N - X_1$$

که در آن داده‌ها به فرم $X_{(1)} < X_{(2)} < \dots < X_N$ با فرض اینکه مرتب شده‌اند.

این شاخص معیار خوبی برای محاسبه میزان پراکندگی داده‌ها نیست. زیرا در محاسبه تنها کوچکترین و بزرگترین داده وارد می‌شود.

مثال ۴: نمرات ۱۰ دانش آموز در دو کلاس متفاوت در درس ریاضی به قرار زیر است:

A کلاس: ۰, 4, 4, 12, 12, 12, 14, 16, 16, 20

B کلاس: 8, 8, 9, 9, 12, 12, 12, 13, 13, 14

با محاسبه مقادیر \bar{X} و m و M دیده می‌شود که برای هر دو کلاس داریم:

$$\bar{X} = 11, \quad M = 12, \quad m = 12$$

اما با توجه به مقادیر تک تک نمرات واضح است که میزان پراکندگی نمرات در دو کلاس کاملاً متفاوت است و برای این منظور نیاز به شاخص‌های پراکندگی می‌باشد تا این مطلب را بتوان با مقایسه آنها نشان داد.

برای دو کلاس مقدار دامنه را محاسبه می‌کنیم:

$$\text{A کلاس: } R = 20 - 0 = 20$$

$$\text{B کلاس: } R = 14 - 8 = 6$$

۲- میانگین انحرافات: فاصله داده X_i از میانگین را انحراف از میانگین داده X_i گویند که به صورت $|X_i - \bar{X}|$ محاسبه می‌شود. اگر این مقدار را

$$D = \frac{1}{N} \sum_{i=1}^K f_i |X_i - \bar{X}|$$

برای تمامی داده‌ها محاسبه کنیم و از نتیجه میانگین بگیریم میانگین انحرافات بدست خواهد آمد که عبارتست از:

از آنجا که میانگین انحرافات به تمام داده‌ها وابسته است معیار مناسبی برای سنجش پراکندگی داده‌ها محسوب می‌شود اما بدلیل وجود قدر مطلق در فرمول، محاسبه آن مشکل است و نمی‌توان آنرا ساده نمود بنابر این از واریانس و انحراف استاندارد استفاده می‌کنیم.

۳- واریانس و انحراف استاندارد: میانگین مجذور انحرافات را واریانس می‌نامیم و با نماد σ^2 یا S_b^2 نمایش می‌دهیم. که عبارتست از:

$$S_b^2 = \frac{1}{N} \sum_{i=1}^K f_i (X_i - \bar{X})^2$$

در مبحث استنباط آماری واریانس را از مجموع مجذور انحرافات داده ها تقسیم بر $N-1$ بدست می‌آورند و آنرا با S^2 نمایش می‌دهند.

$$S^2 = \frac{1}{N-1} \sum_{i=1}^K f_i (X_i - \bar{X})^2$$

در مباحث این درس هر جا صحبت از واریانس می‌کنیم منظور S^2 می‌باشد.

اگر از واریانس جذر بگیریم یعنی $S = \sqrt{S^2}$ در این صورت S را انحراف استاندارد می‌نامیم که معیار مناسبی برای سنجش پراکندگی می‌باشد.

همچنین S^2 را می‌توان از فرمول زیر محاسبه نمود:

$$S^2 = \frac{1}{N-1} \left[\sum_{i=1}^K f_i X_i^2 - \frac{1}{n} \left(\sum_{i=1}^K f_i X_i \right)^2 \right] = \frac{1}{n-1} \left[\sum_{i=1}^K f_i X_i^2 - n \bar{X}^2 \right]$$

اثبات:

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N f_i (x_i - \bar{x})^2 = \frac{1}{N-1} \sum_{i=1}^n f_i (x_i^2 - 2\bar{x}_i \bar{x} + \bar{x}^2)$$

$$= \frac{1}{N-1} \left(\sum_{i=1}^n f_i x_i^2 \right) - 2n\bar{x}^2 + n\bar{x}^2 = \frac{1}{N-1} \left(\sum_{i=1}^n f_i x_i^2 \right) - n\bar{x}^2$$

مثال ۵: مقدار واریانس را برای مثال ۲ محاسبه کنید:

$$\bar{X} = 14/9 \quad N=40$$

$$S^2 = \frac{1}{40-1} \left[(50 \times 64 + 8 \times 121 + 9 \times 196 + 9 \times (17)^2 + 6 \times 40 + 3 \times (23)^2 - 40 \times (14/9)^2 \right]$$

۴- ضریب تغییرات: در محاسبه میزان پراکندگی داده‌ها همواره با داده‌هایی سروکار داریم که با مقیاس‌های مختلفی اندازه‌گیری شده‌اند بنابراین برای مقایسه میزان پراکندگی داده‌های بدست آمده از دو جامعه آماری که با مقیاس‌های مختلفی اندازه‌گیری شده‌اند استفاده از واریانس مناسب نمی‌باشد

زیرا واریانس به مقیاس اندازه‌گیری وابسته می‌باشد بنابراین این از مقیاس مناسبتری به نام ضریب تغییرات استفاده می‌کنیم که از رابطه $CV = \frac{S}{\bar{X}}$ بدست می‌آید و معمولاً با ضریب آن در عدد صد بر حسب درصد بیان می‌شود.

مثال ۶: یک کارخانه تولید لاستیک دو نوع محصول A و B تولید می‌کند. لاستیک نوع A دارای میانگین طول عمر ۲۰۰۰۰ کیلو متر و انحراف استاندارد ۲۰۰۰ کیلو متر می‌باشد و لاستیک نوع B دارای میانگین طول عمر ۱۸۰۰۰ کیلو متر و انحراف استاندارد ۲۰۰ کیلو متر می‌باشد، کدام نوع لاستیک برای خرید مناسب‌تر می‌باشد؟

$$X_A = 20000 \quad \bar{X}_B = 18000 \quad \Rightarrow \quad CV_A = \frac{2000}{20000} = 0.1$$

$$S_A = 2000 \quad S_B = 200 \quad \Rightarrow \quad CV_B = \frac{200}{18000} = 0.01$$

$$CV_A = 0.1 \times 100 = 10\%$$

$$CV_B = 0.01 \times 100 = 1\%$$

همانطور که ملاحظه می‌کنید میانگین طول عمر لاستیک دوم از لاستیک اول کمتر است ولی با توجه به اینکه ضریب تغییرات لاستیک دوم کمتر از لاستیک اول است خرید لاستیک دوم به صرفه‌تر می‌باشد.

۷-۱ تغییر مقیاس و مبدأ

داده‌ها را با واحدهای متفاوتی می‌توان از جامعه آماری جمع آوری نمود به عنوان مثال فرض کنید داده‌های مربوط به وزن ۴۰ نفر از دانشجویان یک کلاس را با واحد کیلوگرم جمع آوری کرده باشید و بخواهید مقادیر میانگین و واریانس را بر حسب پاوند بدست بیاورید برای این منظور نیازی به محاسبه مجدد میانگین و واریانس نمی‌باشد بلکه کافیست از روش تغییر مقیاس و مبدأ استفاده کنید.

۱-۷-۱ تغییر مقیاس

اگر تمامی داده‌ها در عدد a ضرب شوند در این صورت داریم:

$$x_1, x_2, \dots, x_n \rightarrow ax_1, ax_2, \dots, ax_n$$

$$\Rightarrow \bar{X} = \frac{1}{N} \sum_{i=1}^n x_i \rightarrow \bar{X}_a = \frac{1}{N} \sum_{i=1}^n ax_i = a \left(\frac{1}{N} \sum_{i=1}^n x_i \right) = a\bar{X} \Rightarrow \bar{X}_a = a\bar{X}$$

به همین ترتیب برای محاسبه واریانس بدست می‌آید:

تمرین ۱:

$$S_a^2 = a^2 S^2 \rightarrow S_a = |a| S$$

$$S_a = \frac{1}{n-1} \sum_1^n (a x_i - a \bar{x})^2 = a^2 \frac{1}{n-1} \sum_1^n (x_i - \bar{x})^2 = a^2 S^2$$

$$\Rightarrow S_a = \sqrt{a^2 S^2} = |a| S$$

تمرین ۲:

$$X \rightarrow X + b$$

$$\bar{X}_b = \bar{X} + b$$

$$\bar{X}_b = \frac{1}{n} \sum_1^n (x_i + b) = \frac{1}{n} \left[\sum_1^n x_i + \sum_1^n b \right] = \frac{1}{n} \sum_1^n x_i + \frac{b}{n} \sum_1^n (1) = \bar{X} + b \frac{n}{n} = \bar{X} + b$$

$$S_b^2 = S^2$$

$$S_b^2 = \frac{1}{n-1} \sum_1^n (x_i + b - \bar{X} - b)^2 = \frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2 = S^2$$

تمرین ۳ -

$$X \rightarrow Y = aX + b$$

$$\bar{Y} = a\bar{X} + b, \quad S_Y^2 = a^2 S^2$$

$$\text{اثبات: } \bar{Y} = \frac{1}{n} \sum_1^n (ax_i + b) = \frac{a}{n} \sum_1^n x_i + \frac{bn}{n} = a\bar{X} + b$$

$$S_Y^2 = \frac{1}{n-1} \sum_1^n (ax_i + b - a\bar{X} - b)^2 = \frac{a^2}{n-1} \sum_1^n (x_i - \bar{X})^2 = a^2 S^2$$

استاندارد سازی

مثال: نمره علی از امتحان فیزیک و ریاضی به ترتیب برابر ۴۰ و ۶۰ شده است اگر میانگین نمرات امتحان فیزیک و ریاضی به ترتیب برابر ۲۰ و ۵۰ باشد و انحراف معیار امتحان فیزیک و ریاضی به ترتیب برابر ۱ و ۲ باشد علی کدام درس را بهتر امتحان داده است.
حل: برای اینکه بتوان نمرات دو درس را با یکدیگر مقایسه نمود می‌بایستی ابتدا نمرات را استاندارد سازی نمود و سپس آنها را با یکدیگر مقایسه نمود.

$$\text{نمرات استاندارد شده علی در درس ریاضی} = \frac{۶۰ - ۵۰}{۲} = ۵$$

$$\text{نمرات استاندارد شده علی در درس فیزیک} = \frac{۴۰ - ۲۰}{۱} = ۲۰$$

با وجود اینکه نمره علی در درس فیزیک کمتر از ریاضی می‌باشد اما با استاندارد نمودن نمره دو درس مشاهده می‌کنیم که نمره وی در درس فیزیک بالاتر از درس ریاضی می‌باشد به عبارتی علی درس فیزیک را بهتر از درس ریاضی امتحان داده است.

۱-۷-۲ تغییر مبدأ

در صورتی که به تمام داده‌ها مقدار b را اضافه یا کم کنیم می‌توان نشان داد که مقادیر \bar{X} و S^2 جدید از روابط زیر محاسبه می‌شوند:

$$\bar{X}_b = \bar{X} + b$$

$$S_b^2 = S^2 \quad \text{تغییر مبدأ روی واریانس بی تأثیر است}$$

با اعمال همزمان تغییر مبدأ و مقیاس خواهیم داشت:

$$\bar{Y} = a\bar{X} + b$$

$$S_a^2 = a^2 S^2$$

مطالب فوق برای میانه و مد نیز صادق می‌باشند و داریم:

$$m^1 = am + b$$

$$M^1 = a m + b$$

۱۳-۱ استاندارد سازی

از یک جامعه آماری n نمونه X_1, X_2, \dots, X_n بصورت تصادفی انتخاب می‌کنیم بطوریکه میانگین و واریانس نمونه‌ها بترتیب \bar{X} و S_X^2 می‌باشد. با توجه به تغییر مبدأ و مقیاس مقدار هر نمونه را از میانگین نمونه‌ها کم می‌کنیم و حاصل را بر S_X تقسیم می‌کنیم بنابر این داده‌های

$$y_1 = \frac{X_1 - \bar{X}}{S_X}, \quad y_2 = \frac{X_2 - \bar{X}}{S_X}, \quad y_n = \frac{X_n - \bar{X}}{S_X} \quad \text{را خواهیم داشت.}$$

$$\bar{y} = \frac{\bar{X} - \bar{X}}{S_X} = 0; \quad S_y^2 = \left(\frac{1}{S_X^2}\right) S_X^2 = 1 \quad \text{با محاسبه مقایر میانگین و واریانس داده‌های جدید داریم:}$$

همانطور که ملاحظه می‌کنید داده‌های جدید دارای میانگین صفر و واریانس ۱ می‌باشند که به آنها داده‌های استاندارد شده می‌گوییم. همینطور اگر

$$X \text{ داده‌های } X_1 \text{ تا } X_n \text{ را با متغیر تصادفی } X \text{ نمایش دهیم در این صورت } Y = \frac{X - \bar{X}}{S_X} \text{ فرم استاندارد شده یا صورت معیاری متغیر تصادفی } X$$

می‌باشد.

مسائل فصل اول :

۱- دو جامعه با اندازه‌های میانگین \bar{X}_1, \bar{X}_2 و انحراف معیار S_1, S_2 را با یکدیگر ادغام می‌کنیم ثابت کنید میانگین و انحراف معیار جدید از روابط زیر بدست می‌آید:

$$\bar{X} = \frac{N_1 \bar{X}_1 + N_2 \bar{X}_2}{N_1 + N_2}$$

$$S^2 = \frac{N_1 S_1^2 + N_2 S_2^2}{N_1 + N_2} + \frac{N_1 N_2}{(N_1 + N_2)^2} (\bar{X}_1 - \bar{X}_2)^2$$

۲- میانگین و واریانس ۲ داده به ترتیب ۱۵ و ۵ می‌باشد. اگر به جای عدد ۲۵ اشتباهاً عدد ۱۵ را در محاسبات اعمال کرده باشیم میانگین و واریانس جدید را بدست بیاورید.

۳- نشان دهید تغییر مقیاس داده بر روی مقدار ضریب تغییرات $CV = \frac{S}{X}$ بی‌اثر می‌باشد، آیا این مطلب در مورد تغییر مبدأ نیز صادق است؟

۴- ثابت کنید برای میانگین حسابی، هندسی و همساز رابطه زیر برقرار است.

$$\bar{X}_H \leq \bar{X}_G \leq \bar{X}$$

۵- اگر میانگین را از داده‌های یک جامعه آماری کم کنیم و نتیجه را بر انحراف معیار تقسیم کنیم (یعنی $y_i = \frac{x_i - \bar{X}}{S}$) نشان دهید میانگین و انحراف معیار جدید به ترتیب صفر و یک می‌باشد.

۶- برای بدست آوردن یک معیار پراکنندگی جدید داده‌ها را دو به دو با یکدیگر مقایسه می‌کنیم و میانگین n^2 داده جدید $(X_i - X_j)$ را با S_{ij}^2 نمایش می‌دهیم:

$$S_{ij}^2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (X_i - X_j)^2$$

نشان دهید $S_{ij}^2 = 2S^2$ (راهنمایی: داخل پرانتز مقدار \bar{X} را اضافه و کم کنید)

۷- نشان دهید میانگین حسابی و واریانس نخستین n عدد طبیعی به ترتیب $\frac{n+1}{2}$ و $\frac{n^2-1}{12}$ می‌باشد.

۸- جدول زیر را برای داده‌ها و فراوانی آنها در نظر بگیرید:

x	۰	۱	۲	...	n
f	$\binom{n}{0}$	$\binom{n}{1}$	$\binom{n}{2}$...	$\binom{n}{n}$

نشان دهید میانگین و واریانس این داده‌ها به ترتیب عبارتست از:

$$\bar{X} = \frac{n}{2}, \quad S^2 = \frac{n}{4}$$

۹- اگر متحرکی مسافت X_1 را با سرعت V_1 و ... و مسافت X_n را با سرعت V_n طی کنید ثابت کنید سرعت متوسط این متحرک با استفاده از رابطه میانگین همساز یا هارمونیک بدست می‌آید که در آن V_i معادل مقادیر داده‌ها و X_i معادل فراوانی آنهاست.

۱۰- عدد QP که $0 < P < 1$ را چندک P م داده‌ها تعریف می‌کنیم هر گاه فراوانی تجمعی نسبی (I_i) آن بزرگتر یا مساوی با عدد P باشد. به عبارت دیگر هر گاه XP ۱۰۰٪ از داده‌ها قبل از آن قرار گیرند. به عنوان مثال $Q_{.5}$ که به آن چارک دوم می‌گوییم همان میانه می‌باشد چرا که ۵۰٪ داده‌ها قبل از آن قرار دارند. در حالت کلی $Q_{.25}$ و $Q_{.5}$ و $Q_{.75}$ را به ترتیب با Q_1, Q_2, Q_3 نمایش می‌دهیم و به آنها چارکهای اول و دوم و سوم می‌گوییم.

الف) اگر برای داده‌های گسسته X_i داشته باشیم $X_1 < X_2 < \dots < X_i < \dots < X_N$ انگاه نشان دهید $Q_P = (1-\omega) X_r + \omega X_{r+1}$ که در آن $\omega = (n+1)P - r$, $r = (n+1)P$.

ب) برای داده‌های پیوسته رده‌ای که فراوانی تجمعی آن بزرگتر یا مساوی با عدد P باشد رده QP می‌نامیم نشان دهید چندک P م برای داده‌های پیوسته از رابطه زیر بدست می‌آید:

$$Q_P = L_p + \frac{(np - g_p)}{f_p}$$

که در آن :

L_p : کران پایین رده QP

g_p : فراوانی تجمعی رده قبل از رده QP

f_p : فراوانی رده QP

ω : طول رده QP

۱۱- جدول زیر تعداد کتب فروخته شده توسط کتابفروشی را در طول ۳۰ روز نمایش می‌دهد

۱۵	۱۰	۷	۲۰	۱۱	۱۳	۱۸	۶	۵	۴
۱۱	۱۹	۱۲	۱۶	۹	۱۰	۲۱	۱۳	۸	۱۴
۱۷	۲۰	۱۰	۱۲	۱۶	۱۳	۱۱	۱۲	۷	۱۱

برای داده‌های فوق :

الف: یک جدول فراوانی تشکیل دهید و نمودار میله‌ای داده‌ها را رسم کنید.

ب: میانگین داده‌ها \bar{X} ، مد M و میانه m را بدست آورید.

ج: دامنه داده‌ها R ، میانگین انحرافات، واریانس S^2 و ضریب تغییر را بدست آورید.

د: اگر به بزرگترین داده مقدار X $[X \geq 0]$ واحد اضافه کنیم کدامیک از مقادیر مد یا میانه بدون تغییر باقی می‌مانند.

ه: چارک اول و سوم را بدست آورید و دامنه چارکها را محاسبه کنید. (دامنه چارکها : $Q_3 - Q_1$)

و: مقدار دهک چهارم را محاسبه کنید. (دهک چهارم همان $Q_{0.4}$ می‌باشد)

۱۲- در یک شهر میزان درجه حرارت در طول ۳۰ روز به قرار زیر است:

۵	۷	۸/۳	۱۰	۱۱	۱۲/۵	۱۳/۸	۱۳	۱۲	۱۳/۱
۱۴	۱۵/۲	۱۵/۶	۱۶	۱۵/۴	۱۵	۱۶/۵	۱۷	۱۹	۱۹/۷
۲۰/۶	۲۱	۲۱/۳	۲۰/۵	۲۲	۲۲/۸	۲۱/۷	۲۳	۲۴/۱	۲۵

الف: برای داده‌های فوق یک جدول فراوانی تشکیل دهید و هستیوگرام و چند بر فراوانی را رسم کنید

ب: مقادیر میانگین، میانه و مد را محاسبه کنید.

ج: مقادیر واریانس و ضریب تغییر را محاسبه کنید.

د: چند درصد داده‌ها در فاصله $(\bar{X} - S, \bar{X} + S)$ و چند درصد داده‌ها در فاصله $(\bar{X} - 2S, \bar{X} + 2S)$ قرار دارند؟

ه: چند درصد داده‌ها بین چارک اول و سوم قرار دارند؟