

فصل هشتم: آمار و مدل سازی

درسنامه

فصل اول: اندازه گیری و مدل سازی

علم آمار علم جمع آوری اطلاعات و تجزیه و تحلیل آن هاست. اطلاعات به دو دسته‌ی «کمی» و «کیفی یا توصیفی» تقسیم بندی می‌شوند. اطلاعات کمی به آن دسته از اطلاعات می‌گوییم که به صورت عدد و رقم مطرح می‌شوند و قابل تفسیر نیستند. اطلاعات کیفی یا توصیفی به آن دسته از اطلاعات می‌گوییم که قابل تفسیر باشند و ممکن است شنونده‌های مختلف برداشت‌های مختلف از آن داشته باشند. اولین قدم در رسیدن به اطلاعات عددی «اندازه‌گیری» است.

تعریف: (مدل سازی): بیان یک مسأله به زبان ریاضی را «مدل سازی» آن مسأله می‌گوییم.

تفاضل مقدار واقعی و مقدار اندازه‌گیری شده را «خطای اندازه‌گیری» می‌گوییم. معمولاً خطا را با E نمایش می‌دهند، و محاسبه آن به صورت «مقدار تقریبی - مقدار واقعی = E » است. توجه کنید که مقدار خطا (E)، می‌تواند مثبت یا منفی باشد و باید از واحد اندازه‌گیری کم‌تر باشد.

مثال: اگر طول ضلع مربعی ۴ سانتی‌متر اندازه‌گیری شده باشد آن را به صورت $a = 4 + E$ نشان می‌دهند که در آن مقدار خطا در محدوده‌ی $|E| < 1$ سانتی‌متر است.

مثال: فرض کنید وزن فردی $83/5$ کیلو گزارش شده است. در این صورت وزن او را با مدل $W = 83/5 + E$ نمایش می‌دهیم که قدر مطلق E کم‌تر از $0/5$ کیلوگرم یا 500 گرم است.

تذکر: در مسائل مدل سازی از جملاتی که شامل توان دوم خطا یا بالاتر باشند (مانند E^2, E^3, \dots) صرف نظر می‌کنیم.

مثال: اگر شعاع دایره‌ای به صورت $R = 2 + E$ باشد، مدل مساحت این دایره (S) به صورت زیر است:

$$S = \pi R^2 = \pi(2 + E)^2 = \pi(4 + 4E + E^2)$$

طبق تذکر بالا می‌توان از E^2 صرف نظر کرد، پس مدل مساحت این دایره (به صورت تقریبی) برابر است با:

$$S = \pi(4 + 4E)$$

تست ۱: میزان خطای یک ترازوی دیجیتالی ۳- گرم است. اگر جرم جسمی با این ترازو ۱۷۷ گرم نشان داده شود، جرم واقعی این جسم کدام است؟

۱) ۱۸۰ گرم ۲) ۱۷۴ گرم ۳) ۱۷۸/۵ گرم ۴) ۱۷۵/۵ گرم

پاسخ: می‌دانیم «مقدار تقریبی - مقدار واقعی = مقدار خطا». بنابراین می‌توان نوشت:

$$177 - \text{مقدار واقعی} = -3 \Rightarrow \text{مقدار واقعی} = 174$$

بنابراین گزینه‌ی (۲) درست است.

فصل دوم: جامعه و نمونه

تعریف: (جامعه آماری): مجموعه‌ای از اشیا یا افراد که درباره‌ی اعضای آن می‌خواهیم موضوع و یا موضوعاتی را مطالعه کنیم، «جامعه‌ی آماری» می‌گوییم. تعداد اعضای جامعه آماری را «اندازه‌ی جامعه» یا «حجم جامعه» می‌گوییم.

مثال: فرض کنید بخواهیم نمره‌ی ریاضی ترم اول دانش‌آموزان سال دوم دبیرستان کل کشور را در سال ۱۳۹۴ بررسی کنیم. در این بررسی آماری، جامعه‌ی آماری کل دانش‌آموزان سال دوم دبیرستان در سال ۱۳۹۴ هستند و موضوع مورد مطالعه، نمره‌ی ریاضی ترم اول آن هاست.

تعریف: (سرشماری): اگر در یک بررسی آماری، تمام افراد جامعه‌ی آماری را مورد مطالعه قرار دهیم، می‌گوییم «سرشماری» کرده‌ایم. معمولاً در سرشماری با مشکلاتی مثل، در دسترس نبودن تمام اعضای جامعه، وقت‌گیر بودن، گران تمام شدن، از بین رفتن جامعه در برخی از مطالعات و ... مواجه هستیم.

(نمونه): زیرمجموعه‌ی یک جامعه‌ی آماری را «نمونه» می‌گوییم. تعداد اعضای نمونه را «اندازه‌ی نمونه» یا «حجم نمونه» می‌گوییم.

تذکر عمل «نمونه‌گیری» مهم‌ترین بخش آمار است. یک نمونه گروه کوچکی از جامعه‌ی آماری است که به نحوی انتخاب می‌شود که نمایانگر خصوصیات جامعه باشد.

نتایج حاصل از بررسی یا اندازه‌گیری نمونه را «داده» می‌گوییم. داده‌ها را به چهار روش جمع‌آوری می‌کنند که عبارت است از:

- ۱- داده‌های از پیش تهیه شده
- ۲- جمع‌آوری داده‌ها از طریق پرسش
- ۳- جمع‌آوری داده‌ها از طریق آزمایش
- ۴- جمع‌آوری داده‌ها از طریق مشاهده

فرض کنید بخواهیم عددی را از مجموعه‌ای از اعداد به طور تصادفی انتخاب کنیم. انتخاب تصادفی به معنای آن است که ذهنیت ما هیچ تأثیری در انتخاب عدد تصادفی نداشته باشد. برای انتخاب عددی تصادفی در بین تعدادی از اعداد کافی است عدد تصادفی تولید شده توسط ماشین حساب را در اندازه‌ی جامعه ضرب کنیم، سپس قسمت اعشاری عدد به دست آمده را حذف کرده و یک واحد به آن اضافه کنیم.

تذکر اعداد تصادفی که توسط ماشین حساب تولید می‌شوند همواره غیر منفی و کوچک‌تر از ۱ هستند.

مثال: فرض کنید عدد تصادفی تولید شده توسط ماشین حساب 0.273 باشد و اندازه‌ی جامعه‌ای که می‌خواهیم از آن نمونه را انتخاب کنیم 25 باشد، در این صورت خواهیم داشت: $0.273 \times 25 = 6.825$

اکنون قسمت اعشاری عدد را حذف کرده و یک واحد به آن اضافه می‌کنیم که به عدد 69 می‌رسیم.

تست ۲: می‌خواهیم از بین اعداد طبیعی 41 تا 90 یک عدد را به تصادف انتخاب کنیم. اگر عدد تصادفی تولید شده توسط ماشین حساب برابر 0.23 باشد، عدد انتخاب شده کدام است؟

51 (۱) 53 (۲) 52 (۳) 54 (۴)

پاسخ: ابتدا توجه کنید که تعداد اعداد طبیعی 41 تا 90 برابر است با:

اکنون می‌توان نوشت $11/5 = 23/5 \times 0.5$ ، پس باید دوازدهمین عدد را با شروع از 41 انتخاب کنیم. بنابراین باید عدد 52 را انتخاب کنیم.

بنابراین گزینه‌ی (۳) درست است.

فصل سوم: متغیرهای تصادفی

موضوع و مشخصه‌ای از اشیاء جامعه که مورد بررسی قرار می‌گیرد، «متغیر تصادفی» نامیده می‌شود. انواع متغیرهای تصادفی به صورت زیر است:

متغیرهایی هستند که قابل «اندازه‌گیری» یا «شمارش» هستند و به دو دسته تقسیم می‌شوند:

الف) متغیرهای کمی گسسته: این متغیرها «قابل شمارش» هستند مانند تعداد نامه‌های یک صندوق، تعداد بیماران مراجعه کننده به یک پزشک در طول روز و ...

ب) متغیرهای کمی پیوسته: این متغیرها «قابل اندازه‌گیری» هستند مانند وزن نامه‌های موجود در یک صندوق، زمانی که یک بیمار در اتاق انتظار مطب، منتظر است.

این متغیرها قابل اندازه‌گیری و شمارش نیستند و به دو دسته تقسیم می‌شوند:

الف) متغیر کیفی ترتیبی: در این متغیرها نوعی ترتیب طبیعی وجود دارد مانند مراحل تحصیلی که به صورت دبستان، راهنمایی، دبیرستان و دانشگاه است. یا مراحل زندگی که به صورت نوزادی، کودکی، نوجوانی، جوانی، میان‌سالی و کهنسالی است.

ب) متغیر کیفی اسمی: در این نوع متغیر کیفی، هیچ‌گونه ترتیبی وجود ندارد مانند گروه خونی، RH خون، رنگ اشیا.

تست ۳: میزان «آلودگی هوا» کدام نوع متغیر است؟

(۱) کمی - گسسته (۲) کمی - پیوسته (۳) کیفی - ترتیبی (۴) کیفی - اسمی

پاسخ: میزان آلودگی هوا را به کمک دستگاههایی اندازه گیری می کنند، پس متغیری کمی و پیوسته است. بنابراین گزینه ی (۲) درست است.

فصل چهارم: دسته بندی داده ها و جدول فراوانی

معمولاً داده ها در قالب یک جدول نوشته می شوند که به آن «جدول فراوانی» می گوئیم.

فراوانی مطلق داده ی X_i برابر تعداد دفعاتی است که آن داده تکرار شده است. فراوانی مطلق داده ی X_i را با نماد f_i نشان می دهیم. «فراوانی مطلق» را گاهی به اختصار «فراوانی» می گویند.

تذکر مجموع تمام فراوانی های مطلق برابر تعداد کل داده هاست.

مثال: فرض کنید نمرات یک کلاس ۱۰ نفره به صورت ۸، ۸، ۹، ۱۲، ۱۴، ۱۴، ۱۶، ۱۶، ۱۶، ۱۹ است. در این صورت جدول فراوانی این

داده ها به صورت مقابل است:

X_i	۸	۹	۱۲	۱۴	۱۶	۱۹
f_i	۲	۱	۱	۲	۳	۱

به عنوان نمونه، طبق جدول بالا، فراوانی مطلق داده ی $X_i=16$ برابر ۳ است.

برای بررسی تعدادی داده، گاهی بهتر است آن ها را دسته بندی کنیم. دسته بندی به این صورت انجام می گیرد که داده هایی که نزدیک به هم هستند در یک دسته قرار می گیرند. واضح است که هر چقدر تعداد دسته ها بیش تر باشد، نتایج حاصل از آن دسته بندی دقیق تر خواهد بود. قبل از بیان روش دسته بندی، به تعریف «دامنه ی تغییرات» می پردازیم.

اختلاف بین بزرگ ترین و کوچک ترین داده را «دامنه ی تغییرات» می گوئیم. دامنه ی تغییرات را معمولاً با R نمایش می دهند.

مثال: در داده های ۱۰، ۱۱، ۱۱، ۱۲، ۱۳، ۱۴، ۱۷، ۱۷، ۱۷، ۱۸، ۱۸، ۱۹، ۱۹ دامنه ی تغییرات برابر است با: $R=19-10=9$

تذکر اگر به همه ی داده ها عدد ثابت k را اضافه کنیم یا از همه ی داده ها عدد ثابت k را کم کنیم دامنه ی تغییرات داده ها تغییری نمی کند. به همین ترتیب، اگر همه ی داده ها را در عدد k ضرب کنیم دامنه ی تغییرات در $|k|$ ضرب می شود.

مثال: اگر دامنه ی تغییرات داده های X_1, X_2, \dots, X_p برابر ۱۷ باشد، دامنه ی تغییرات داده های $X_1+3, X_2+3, \dots, X_p+3$ نیز برابر ۱۷ خواهد بود.

مثال: اگر دامنه ی تغییرات X_1, X_2, \dots, X_p برابر ۱۴ باشد، دامنه ی تغییرات داده های $-2X_1, -2X_2, \dots, -2X_p$ برابر است با:

$$|-2| \times 14 = 28$$

تست ۴: اگر دامنه ی تغییرات داده های X_1, X_2, \dots, X_p برابر ۱۹ باشد، دامنه ی تغییرات داده های $3X_1+7, 3X_2+7, \dots, 3X_p+7$ کدام است؟

(۱) ۲۶ (۲) ۵۷ (۳) ۶۴ (۴) ۱۹

پاسخ: همان طور که در مطلب بالا ملاحظه کردید اضافه یا کم کردن عددی ثابت به همه ی داده ها تأثیری در دامنه ی تغییرات ندارد. بنابراین می توان نتیجه گرفت دامنه ی تغییرات داده های $3X_1+7, 3X_2+7, \dots, 3X_p+7$ برابر است با $3 \times 19 = 57$.

بنابراین گزینه ی (۲) درست است.

برای دسته بندی داده ها، ابتدا دامنه ی تغییرات را محاسبه می کنیم. اگر دامنه ی تغییرات را بر تعداد دسته ها تقسیم کنیم، طول دسته ها محاسبه می شود. به جز دسته ی آخر، بقیه ی دسته ها را به صورت $[a_i, b_i]$ در نظر می گیریم که در آن ها به a_i ، کران پایین و به b_i کران بالای دسته می گوئیم. توجه کنید که دسته ی آخر را به صورت $[a_i, b_i]$ در نظر می گیریم.

تعریف: (مرکز دسته): به میانگین کران بالا و پایین هر دسته، «نشان دسته» یا «مرکز دسته» می‌گوییم.

مثال: فرض کنید نمرات یک کلاس ۱۵ نفره به صورت ۸، ۱۲، ۱۲، ۱۲، ۱۴، ۱۴، ۱۶، ۱۶، ۱۷، ۱۷، ۱۷، ۱۸، ۱۹، ۲۰ است. اگر بخواهیم

این داده‌ها را به ۴ دسته تقسیم کنیم، با توجه به این‌که، دامنه‌ی تغییرات برابر $R = 20 - 8 = 12$ است، طول دسته‌ها برابر $\frac{12}{4} = 3$

است. پس دسته‌ها به صورت زیر خواهند بود:

$[8, 11), [11, 14), [14, 17), [17, 20]$

اکنون جدول فراوانی این داده‌های دسته‌بندی شده به صورت زیر است:

دسته‌ها	مرکز دسته‌ها	فراوانی هلالی
$[8, 11)$	۹/۵	۱
$[11, 14)$	۱۲/۵	۳
$[14, 17)$	۱۵/۵	۴
$[17, 20]$	۱۸/۵	۷

نکته: ۱- اختلاف مرکز دو دسته‌ی متوالی برابر است با طول دسته

۲- اختلاف کران پایین و کران بالای هر دسته برابر است با طول دسته

۳- فاصله‌ی کران پایین (یا کران بالا) هر دسته تا مرکز آن دسته برابر است با نصف طول آن دسته

مثال: در یک آمارگیری، تعداد دسته‌ها ۵ و دامنه‌ی تغییرات ۵۰ است. اگر کم‌ترین داده ۲ باشد، کران پایین دسته‌ی سوم و مرکز دسته‌ی

چهارم را محاسبه می‌کنیم. برای این منظور می‌دانیم طول دسته‌ها برابر است با $\frac{50}{5} = 10$. اکنون با توجه به این‌که طول دسته‌ها برابر

۱۰ و کران پایین دسته‌ی اول برابر ۲ است، این دسته‌ها به صورت زیر هستند:

...، $[32, 42)$: دسته‌ی چهارم ، $[22, 32)$: دسته‌ی سوم ، $[12, 22)$: دسته‌ی دوم ، $[2, 12)$: دسته‌ی اول

پس کران پایین دسته‌ی سوم برابر ۲۲ و مرکز دسته‌ی چهارم برابر است با: $\frac{32+42}{2} = 37$

تست ۵: یک سری داده‌ی آماری را در دسته‌هایی با طول‌های مساوی، دسته‌بندی کرده‌ایم. اگر مرکز دسته‌ی دوم $15/8$ و کران بالای دسته‌ی پنجم

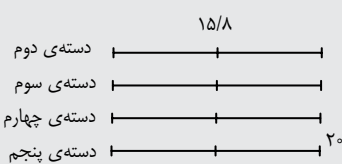
۲۰ باشد، کران بالای دسته‌ی هشتم کدام است؟

۲۴/۶ (۴)

۲۲/۶ (۳)

۲۳/۶ (۲)

۲۶/۳ (۱)



پاسخ: با توجه به شکل مقابل، اگر طول دسته‌ها را x در نظر بگیریم فاصله‌ی مرکز دسته‌ی دوم تا

کران بالای دسته‌ی پنجم برابر $3/5x$ است پس می‌توان نوشت:

$$3/5x = 20 - 15/8 \Rightarrow x = \frac{6}{5} = 1/2$$

$$20 + 3x = 20 + 3 \cdot 1/2 = 23/2$$

اکنون می‌توان نتیجه گرفت کران بالای دسته‌ی هشتم برابر است با:

بنابراین گزینه‌ی (۲) درست است.

اگر فراوانی هر دسته را به تعداد کل داده‌ها تقسیم کنیم، عدد به‌دست آمده را «فراوانی نسبی» آن دسته می‌گویند.

تذکر مجموع فراوانی‌های نسبی همه‌ی دسته‌ها برابر یک است.

تست ۶: در جدول توزیع فراوانی روبه‌رو حاصل $x+y+z$ کدام است؟

دسته‌ها	فراوانی	فراوانی نسبی
۰-۵	۳۰	۰/۲
۵-۱۰	x	۰/۱
۱۰-۱۵	y	z

(۱) ۱۲۰/۷

(۲) ۱۲۲/۷

(۳) ۱۲۵/۷

(۴) ۱۲۴/۳

پاسخ: با توجه به این که مجموع فراوانی‌های نسبی برابر ۱ است می‌توان فهمید:

$$0/2 + 0/1 + z = 1 \Rightarrow z = 0/7$$

اکنون با توجه به این که فراوانی‌های مطلق و نسبی دسته‌ی اول به ترتیب ۳۰ و ۰/۲ است می‌توان تعداد کل داده‌ها را محاسبه کرد. اگر

تعداد کل داده‌ها را n فرض کنیم، داریم:

$$0/2 = \frac{30}{n} \Rightarrow n = 150$$

از آن‌جا که فراوانی نسبی دسته‌ی دوم و تعداد کل داده‌ها را می‌دانیم می‌توانیم مقدار x و y را محاسبه کنیم:

$$0/1 = \frac{x}{150} \Rightarrow x = 150, \quad y = 150 - (30 + 15) = 105$$

پس حاصل $x+y+z$ برابر ۱۲۰/۷ است.

بنابراین گزینه‌ی (۱) درست است.

اگر فراوانی نسبی هر دسته را در عدد ۱۰۰ ضرب کنیم، عدد به‌دست آمده را «درصد فراوانی نسبی» آن دسته می‌گوییم.

تذکر مجموع درصدهای فراوانی نسبی همه‌ی دسته‌ها برابر ۱۰۰ است.

تست ۷: در یک جدول توزیع فراوانی، فراوانی مطلق داده‌ها در دسته‌های اول تا پنجم دارای روندی صعودی است. اگر فراوانی مطلق دسته‌ی سوم برابر ۸ باشد و درصد فراوانی نسبی دسته‌های دوم و چهارم به ترتیب ۲۵ و ۴۰ باشد، آن‌گاه فراوانی کل کدام می‌تواند باشد؟

(۱) ۳۴ (۲) ۳۶ (۳) ۳۰ (۴) ۱۸

پاسخ: اگر فراوانی کل داده‌ها را n در نظر بگیریم، درصد فراوانی نسبی دسته‌ی سوم برابر $\frac{8}{n} \times 100$ می‌باشد.

از طرفی، طبق صورت مسأله، فراوانی‌های دسته‌ها روندی صعودی دارند پس می‌توان فهمید درصد فراوانی‌های نسبی هم روندی صعودی دارند پس می‌توان نوشت:

$$25 \leq \frac{800}{n} \leq 40 \Rightarrow 20 \leq n \leq 32$$

در بین گزینه‌ها، فقط عدد ۳۰ در این محدوده قرار دارد.

بنابراین گزینه‌ی (۳) درست است.

مجموع فراوانی‌های مطلق دسته‌های اول، دوم، ... تا \sum را «فراوانی تجمعی دسته‌ی \sum » می‌گوییم.

تذکر فراوانی تجمعی دسته‌ی اول با فراوانی مطلق آن برابر است. همچنین فراوانی تجمعی دسته‌ی آخر برابر است با تعداد کل داده‌ها.

اگر فراوانی تجمعی هر دسته را بر تعداد کل داده‌ها تقسیم کنیم، به آن «فراوانی تجمعی نسبی» آن دسته می‌گوییم و اگر این عدد را در ۱۰۰ ضرب کنیم، به آن «درصد فراوانی تجمعی نسبی» آن دسته می‌گوییم.

تذکر در جدول فراوانی، فراوانی تجمعی نسبی و درصد فراوانی تجمعی نسبی دسته‌ی آخر به ترتیب برابر ۱ و ۱۰۰ است.

مثال: برای کلاس ۱۵ نفره با نمرات ۸، ۱۲، ۱۲، ۱۲، ۱۴، ۱۴، ۱۶، ۱۶، ۱۷، ۱۷، ۱۷، ۱۷، ۱۸، ۱۹ و ۲۰ جدول فراوانی را در مبحث روش دسته‌بندی نوشتیم، می‌توانیم جدول را به صورت زیر بنویسیم.

درصد فراوانی نسبی	فراوانی تجمعی نسبی	فراوانی تجمعی	درصد فراوانی نسبی	فراوانی نسبی	فراوانی مطلق	مرکز دسته	دسته‌ها
$\frac{1}{15} \times 100$	$\frac{1}{15}$	۱	$\frac{1}{15}$	$\frac{1}{15}$	۱	۹/۵	[۸, ۱۱)
$\frac{3}{15} \times 100$	$\frac{4}{15}$	۴	$\frac{3}{15}$	$\frac{3}{15}$	۳	۱۲/۵	[۱۱, ۱۴)
$\frac{4}{15} \times 100$	$\frac{8}{15}$	۸	$\frac{4}{15}$	$\frac{4}{15}$	۴	۱۵/۵	[۱۴, ۱۷)
$\frac{7}{15} \times 100$	۱	۱۵	$\frac{7}{15}$	$\frac{7}{15}$	۷	۱۸/۵	[۱۷, ۲۰]
مجموع = ۱۰۰			مجموع = ۱		مجموع = ۱۵		

تست ۸: اگر در یک جدول فراوانی، فراوانی تجمعی دسته‌های یازدهم و دوازدهم با هم برابر باشند، درصد فراوانی نسبی دسته‌ی دوازدهم کدام است؟
 (۱) ۱٪ (۲) ۱۰٪ (۳) صفر (۴) ۱۰۰٪

پاسخ: می‌دانیم «فراوانی مطلق دسته‌ی دوازدهم + فراوانی تجمعی دسته‌ی یازدهم = فراوانی تجمعی دسته‌ی دوازدهم» اکنون با توجه به این که فراوانی تجمعی دسته‌های یازدهم و دوازدهم با هم برابر هستند می‌توان فهمید فراوانی مطلق دسته‌ی دوازدهم برابر صفر است. پس فراوانی نسبی (و درصد فراوانی نسبی) نیز در این دسته برابر صفر است. بنابراین گزینه‌ی (۳) درست است.

فصل پنجم: نمودارها و تحلیل داده‌ها

نمودارها وسیله‌ای سودمند برای به تصویر کشیدن ویژگی‌های یک جامعه‌ی آماری هستند. در ادامه انواع نمودارها را معرفی می‌کنیم:

نمودار میله‌ای از سه قسمت اصلی تشکیل شده است.

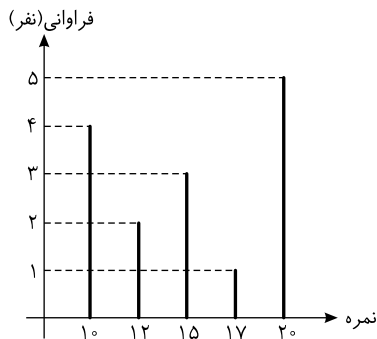
الف) عنوان: زیر هر نمودار باید موضوع مورد مطالعه به‌طور خلاصه نوشته شود.

ب) برچسب محورها: هر یک از محورها برچسبی دارند که مشخص‌کننده‌ی متغیری است که آن محور نشان می‌دهد.

پ) مقیاس: مقیاس هر محور باید مشخص باشد.

مثال: اگر نمرات امتحان نهایی درس ریاضی یک کلاس به صورت ۱۰، ۱۰، ۱۰، ۱۰، ۱۰، ۱۰، ۱۲، ۱۲، ۱۵، ۱۵، ۱۵، ۱۷، ۲۰، ۲۰، ۲۰، ۲۰ باشد.

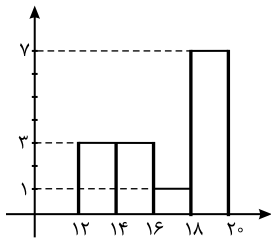
نمودار میله‌ای این نمرات به صورت زیر است:



نمودار نمرات امتحان نهایی درس ریاضی کلاس (الف)

در مثال بالا داده‌ها دسته‌بندی نشده بودند. در رسم نمودار میله‌ای برای داده‌های دسته‌بندی شده مرکز دسته‌ها را روی محور افقی و فراوانی هر دسته را روی محور عمودی (محور y ها) در نظر می‌گیریم. در این حالت واضح است که فاصله‌ی هر دو میله‌ی متوالی برابر طول دسته‌ها است. نمودار میله‌ای بیش‌تر برای متغیرهای کمی گسسته و متغیرهای کیفی مناسب است.

نمودار مستطیلی برای متغیرهای کمی پیوسته مناسب است. این نمودار، نمایشی از داده‌های دسته‌بندی شده است که در آن سطح مستطیل‌ها متناسب با فراوانی دسته‌ها است. در این نمودار هر دسته مستطیلی دارد که اگر طول همه‌ی دسته‌ها یکسان باشد یکی از اضلاع آن طول دسته و ضلع دیگرش فراوانی آن دسته است.



مثال: اگر نمرات یک کلاس ۱۴ نفره به صورت $۱۸, ۱۷, ۱۵, ۱۴/۵, ۱۴, ۱۳, ۱۲/۷, ۱۲, ۱۱, ۱۰, ۹, ۸, ۷, ۶, ۵, ۴, ۳, ۲, ۱$ باشد و بخواهیم آن‌ها را به ۴ دسته تقسیم کنیم، با توجه به این که دامنه‌ی تغییرات برابر $R=20-12=8$ است، طول دسته‌ها برابر ۲ خواهد بود و نمودار مستطیلی آن به صورت مقابل است:

تذکر: با توجه به نمودار مستطیلی می‌توان فراوانی نسبی هر دسته را از طریق تقسیم «مساحت مربوط به آن دسته» بر «مجموع مساحت

همه‌ی مستطیل‌ها» محاسبه کرد. به عنوان مثال در نمودار مستطیلی مثال قبل، فراوانی نسبی دسته‌ی (۱۲, ۱۴) برابر است با: $\frac{6}{28}$

تذکر: مساحت زیر یک نمودار مستطیلی وقتی طول همه‌ی دسته‌ها یکسان باشند برابر است با «تعداد کل داده‌ها \times طول دسته».

تست ۹: مساحت زیر نمودار مستطیلی n داده که در دسته‌هایی با طول‌های مساوی دسته‌بندی شده‌اند ۵۴۰ است. اگر تعداد کل داده‌ها $n=۹۰$ باشد، طول هر دسته کدام است؟

۴ (۱) ۶ (۲) ۹ (۳) ۱۰ (۴)

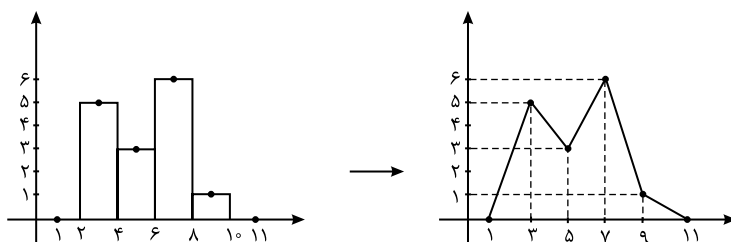
پاسخ: همان‌طور که گفته شد مساحت زیر یک نمودار مستطیلی در چنین شرایطی برابر است با «طول دسته \times تعداد کل داده‌ها». پس می‌توان نوشت: (طول دسته‌ها را x در نظر می‌گیریم)

$$۵۴۰ = ۹۰ \cdot x \Rightarrow x = ۶$$

بنابراین گزینه‌ی (۲) درست است.

نمودار چندبر فراوانی تغییرات یک متغیر پیوسته (مانند وزن، قد، دما و ...) را بهتر از نمودار مستطیلی نمایش می‌دهد. برای رسم نمودار «چندبر فراوانی»، نقاطی از صفحه را در نظر می‌گیریم که طول آن‌ها مرکز دسته‌ها و عرض هر یک از آن‌ها فراوانی مربوط به آن دسته باشد. سپس آن نقاط را به ترتیب به یک‌دیگر متصل می‌کنیم. در انتها برای آن‌ها منحنی «چندبر فراوانی» به شکل بسته ایجاد شود و سطح زیر چندبر با سطح زیر نمودار مستطیلی یکسان شود، دو دسته با فراوانی صفر در ابتدا و انتهای نمودار در نظر می‌گیریم.

تذکر: نمودار چندبر فراوانی را می‌توان با داشتن فراوانی نسبی نیز رسم کرد که آن‌را «چندبر فراوانی نسبی» می‌گوییم.



مثال: در شکل زیر از روی نمودار مستطیلی (نمودار

سمت چپ) و به کمک نقاط مشخص شده،

نمودار «چندبر فراوانی» (سمت راست) رسم

شده است. همان‌طور که ملاحظه می‌کنید برای

رسم «چندبر فراوانی» از نقاط

$(1,0), (3,5), (5,3), (7,6), (9,1), (11,0)$

استفاده شده است. یادآوری می‌کنیم که

مساحت زیر نمودار «مستطیلی» با مساحت

زیر نمودار «چندبر فراوانی» برابر است.

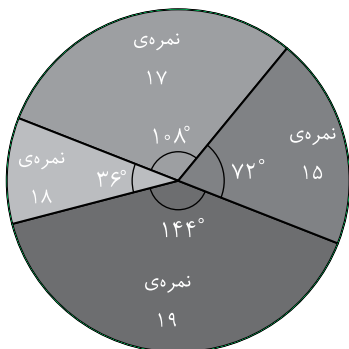
تست ۱۰: هرگاه دسته‌ی اول یک جدول فراوانی به صورت $(۹, ۱۵)$ و فراوانی آن ۵ باشد، طول‌های نقاط اول و دوم نمودار چندبر فراوانی تشکیل شده به ترتیب چه اعدادی هستند؟

(۱) ۱۲, ۶ (۲) صفر, ۵ (۳) ۶, ۵ (۴) صفر, ۱۲

پاسخ: توجه کنید که در نمودار چندبر فراوانی مرکز هر دسته را بر روی محور x قرار می‌دهیم و قبل از دسته‌ی اول یک دسته با فراوانی صفر به نمودار اضافه می‌کنیم. توجه کنید که فاصله‌ی بین نقاط متوالی برابر طول دسته‌ها است. پس می‌توان فهمید نقطه‌ی دوم چندبر فراوانی برابر مرکز دسته‌ی $(۹, ۱۵)$ است که برابر با $\frac{۹+۱۵}{۲} = ۱۲$ و طول نقطه‌ی اول برابر است با $۱۲ - ۶ = ۶$. بنابراین گزینه‌ی (۱) درست است.

یکی دیگر از نمودارهایی که می‌تواند اطلاعات موجود در داده‌ها را به سرعت در معرض دید قرار دهد «نمودار دایره‌ای» است. برای رسم این نمودار چنین عمل می‌کنیم:

اگر متغیر تصادفی مورد مطالعه‌ی ما دارای k حالت باشد که فراوانی‌های آن‌ها به ترتیب f_1, f_2, \dots, f_k باشد و $f_1 + f_2 + \dots + f_k = n$ ، ابتدا دایره‌ای به شعاع دلخواه رسم می‌کنیم. سپس این دایره را به کمک زاویه‌های مرکزی به k قسمت تقسیم می‌کنیم به طوری که زاویه‌ی مرکزی دسته‌ی i ام (با فراوانی f_i)، برابر است با:

$$\frac{f_i}{n} \times ۳۶۰$$


مثال: اگر نمرات یک کلاس ۱۰ نفره به صورت ۱۹, ۱۹, ۱۹, ۱۸, ۱۷, ۱۷, ۱۷, ۱۵, ۱۵ باشد، نمودار دایره‌ای به صورت زیر است:

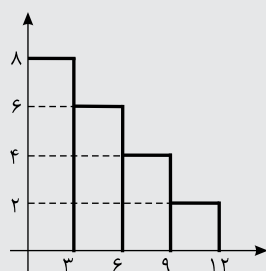
$$۱۵ \text{ نمره‌ی } ۱۵ = \frac{۲}{۱۰} \times ۳۶۰ = ۷۲^\circ$$

$$۱۷ \text{ نمره‌ی } ۱۷ = \frac{۳}{۱۰} \times ۳۶۰ = ۱۰۸^\circ$$

$$۱۸ \text{ نمره‌ی } ۱۸ = \frac{۱}{۱۰} \times ۳۶۰ = ۳۶^\circ$$

$$۱۹ \text{ نمره‌ی } ۱۹ = \frac{۴}{۱۰} \times ۳۶۰ = ۱۴۴^\circ$$

تست ۱۱: در نمودار مستطیلی مقابل، زاویه‌ی مرکزی مربوط به دسته‌ای با مرکز $۷/۵$ در نمودار دایره‌ای کدام است؟



(۱) ۶۰°

(۲) ۷۲°

(۳) ۸۰°

(۴) ۷۵°

پاسخ: در دسته‌ی $(۶, ۹)$ ، مرکز دسته برابر $۷/۵$ است. بنابراین برای تعیین زاویه‌ی مرکزی این دسته در نمودار دایره‌ای باید فراوانی نسبی این دسته را در ۳۶۰ ضرب کنیم:

$$\text{زاویه‌ی مرکزی مربوط به دسته‌ای با مرکز } ۷/۵ = \frac{۴}{۲+۴+۶+۸} \times ۳۶۰ = \frac{۴}{۲۰} \times ۳۶۰ = ۷۲^\circ$$

بنابراین گزینه‌ی (۲) درست است.

به داده‌های آماری مقابل دقت کنید:

۴۰, ۴۰, ۴۰, ۴۰, ۴۱, ۴۱, ۴۲, ۴۳, ۵۰, ۵۲, ۵۴, ۵۶, ۵۶, ۵۶, ۵۷

این داده‌ها را می‌توان به صورت نمودار زیر که «نمودار ساقه و برگ» نامیده می‌شود نمایش داد:

ساقه	برگ							
۴	۰	۰	۰	۰	۱	۱	۲	۳
۵	۰	۲	۴	۶	۶	۶	۷	

کلید نمودار: ۴۱=۴

در نمودار بالا، ستون سمت چپ را که بیانگر قسمت مشترک تعدادی عدد است «ساقه» و اعداد روبه‌روی این ستون را «برگ» می‌نامند.

مثال: داده‌های اعشاری زیر را با نمودار ساقه و برگ نمایش می‌دهیم:

۴۱, ۴۱/۱, ۴۱/۲, ۴۱/۲, ۴۲, ۴۲, ۴۲/۳, ۴۲/۳, ۴۲/۳, ۴۲/۵, ۴۳, ۴۳, ۴۳, ۴۳/۹

ساقه	برگ							
۴۱	۰	۱	۲	۲				
۴۲	۰	۰	۳	۳	۳	۵		
۴۳	۰	۰	۰	۹				

کلید نمودار: ۴۱/۱=۴۱

همان‌طور که ملاحظه می‌کنید برای نمایش اعداد اعشاری که قسمت صحیح آن‌ها شبیه هم است می‌توان از نمودار ساقه و برگ استفاده کرد. نمودار جعبه‌ای را بعد از معرفی میانه توضیح خواهیم داد.

فصل ششم: شاخص‌های مرکزی

شاخص‌های مرکزی به شاخص‌هایی گفته می‌شود که محل تمرکز داده‌ها را معرفی می‌کند. شاخص‌های مرکزی عبارت هستند از: مُد (نَما)، میانه و میانگین که به بررسی هر یک می‌پردازیم.

مُد داده‌ای است که بیش‌ترین فراوانی را دارد.

مثال: در داده‌های ۲۰, ۲۰, ۱۹, ۱۸, ۱۸, ۱۸, ۱۷, ۱۷, ۱۷, ۱۵ عدد ۱۸ مُد است، چون بیش‌ترین فراوانی را دارد.

تذکر: مُد ممکن است منحصر به فرد نباشد به عنوان مثال داده‌های ۲۰, ۱۹, ۱۸, ۱۷, ۱۶, ۱۶, ۱۵ دارای دو مُد ۱۶ و ۱۹ هستند به این نوع جوامع، ۲ مُدی می‌گویند. جامعه ممکن است «چند مُدی» هم باشد. مُد در این قبیل جامعه‌ها شاخص معتبری نیست.

تست ۱۲: در نمودار ساقه و برگ مقابل مُد کدام است؟

ساقه	برگ								۲۴ (۲)	۱۴ (۱)	
۱	۰	۱	۲	۳	۳					۲۱ (۴)	۱۰ (۳)
۲	۱	۲	۳	۴	۴	۴					

پاسخ: با توجه به نمودار، داده‌ی ۲۴ دارای بیش‌ترین فراوانی است، پس مقدار مُد برابر ۲۴ است. بنابراین گزینه‌ی (۲) درست است.

پس از مرتب کردن داده‌ها، داده‌ای که وسط همه‌ی داده‌ها قرار می‌گیرد را «میانه» می‌گویند. (مشابه تعریف میانه‌ی مثلث در هندسه!!) اگر تعداد داده‌ها زوج باشد، میانگین دو داده‌ی وسط را میانه در نظر می‌گیریم.

مثال: میانه‌ی داده‌های ۱۲, ۲۰, ۱۹, ۱۸, ۱۴, ۱۵, ۱۷، پس از مرتب کردن داده‌ها، برابر ۱۷ است.

مثال: میانه‌ی داده‌های ۲۰, ۱۷, ۱۶, ۱۵, ۱۴, ۱۲، برابر است با: $\frac{۱۵+۱۶}{۲} = ۱۵/۵$

تذکر: همان‌طور که ملاحظه می‌کنید میانه، لزوماً عضوی از داده‌ها نیست.

تذکر: میانه و مد هیچ‌کدام نسبت به اندازه‌ی داده‌ها حساسیتی نشان نمی‌دهند. به عنوان مثال در هر دو جامعه‌ی (۰, ۲, ۳) و (۰, ۲, ۷۰۰) میانه عدد ۲ است در حالی که اندازه‌ی داده‌ها در این دو جامعه خیلی متفاوت است. همین مطلب در مورد مُد نیز برقرار است.

تست ۱۳: هر یک از داده‌های زیر را با چه عددی جمع کنیم تا میانه‌ی داده‌ها برابر ۱۵ باشد؟

۸, ۱۰, ۱۱, ۱۷, ۱۸, ۱۹

۱ (۴)

۲/۵ (۳)

۱/۵ (۲)

۲ (۱)

پاسخ: فرض کنید که به همه‌ی داده‌ها x واحد اضافه کرده‌ایم. در این صورت میانه‌ی داده‌ها به صورت زیر خواهد بود:

$$\text{میانه} = \frac{(11+x) + (17+x)}{2} \Rightarrow 15 = \frac{28+2x}{2} \Rightarrow 30 = 28+2x \Rightarrow x=1$$

بنابراین گزینه‌ی (۴) درست است.

نمودارهایی که تاکنون شناخته‌ایم (میله‌ای، مستطیلی، چندبر، دایره‌ای و ساقه‌وبرگ) همگی، برای مقایسه‌ی داده‌ها بسیار مفیدند ولی هیچ کدام به سؤالاتی از قبیل «آیا داده‌ها به هم نزدیک هستند؟» و «داده‌ها بیش تر در اطراف کدام داده متمرکزند؟ میانگین، بیش ترین داده یا کم ترین داده؟» پاسخ نمی‌دهند. برای پاسخ به این گونه سؤالات نمودار جعبه‌ای کاربرد دارد که به روش رسم آن می‌پردازیم.

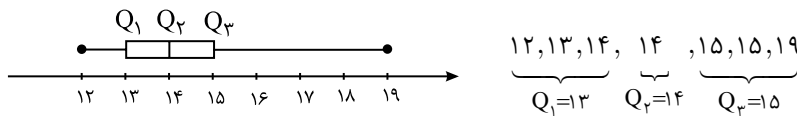
چارک اول: به میانه‌ی نیمه‌ی اول داده‌ها (داده‌هایی که قبل از میانه قرار دارند)، چارک اول می‌گوییم و آن را با Q_1 نمایش می‌دهیم.

چارک دوم: به میانه‌ی همه‌ی داده‌ها چارک دوم می‌گوییم و آن را با Q_2 نمایش می‌دهیم. پس چارک دوم همان میانه است.

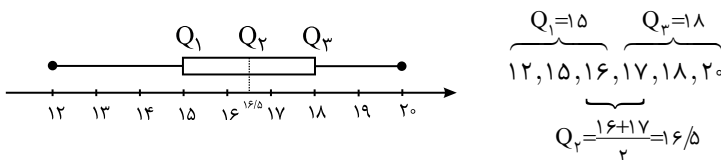
چارک سوم: به میانه‌ی نیمه‌ی دوم داده‌ها (داده‌هایی که بعد از میانه قرار دارند)، چارک سوم می‌گوییم و آن را با Q_3 نمایش می‌دهیم.

تعریف: نمودار جعبه‌ای نموداری تصویری است که داده‌ها را بر اساس پنج مقدار کوچک ترین داده، چارک اول، چارک دوم (میانه)، چارک سوم و بزرگ ترین داده نمایش می‌دهد.

مثال: نمودار جعبه‌ای داده‌های مرتب شده‌ی ۱۲, ۱۳, ۱۴, ۱۴, ۱۵, ۱۵, ۱۹ به صورت زیر تعیین می‌شود:



مثال: نمودار جعبه‌ای داده‌های مرتب شده‌ی ۱۲, ۱۵, ۱۶, ۱۷, ۱۸, ۲۰ به صورت زیر تعیین می‌شود:



تست ۱۴: چه تعدادی از داده‌های یک مجموعه‌ی ۱۲ عضوی بین چارک اول و سوم قرار دارند؟

۳ (۴)

۶ (۳)

۷ (۲)

۵ (۱)

پاسخ: اگر داده‌های مورد نظر را به صورت مرتب a_1, a_2, \dots, a_{12} در نظر بگیریم (توجه کنید که داده‌ها متمایز هستند). می‌توان فهمید میانه‌ی

نیمه‌ی اول داده‌ها برابر $Q_1 = \frac{a_3 + a_4}{2}$ و میانه‌ی نیمه‌ی دوم داده‌ها برابر $Q_3 = \frac{a_9 + a_{10}}{2}$ است. پس داده‌های a_4, a_5, \dots, a_9 بین

چارک اول و چارک سوم هستند که تعداد آن‌ها ۶ تا است. بنابراین گزینه‌ی (۳) درست است.

میانگین به معنای معدل داده‌ها است و شاخص خوبی برای نشان دادن مرکزیت داده‌ها است. میانگین داده‌های x_1, x_2, \dots, x_n را با \bar{x} نمایش

می‌دهیم و به صورت زیر تعریف می‌شود:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

مثال: میانگین داده‌های ۱۲، ۱۳، ۱۴، ۱۵، ۱۵، ۱۵، ۱۶، ۳۶ برابر است با: $\bar{x} = \frac{12+13+14+15+15+15+16+36}{8} = \frac{136}{8} = 17$

میانگین وزن دار (وزنی)

فرض کنید داده‌های x_1, x_2, \dots, x_n به ترتیب دارای ضرایب‌های f_1, f_2, \dots, f_n باشند، این داده‌ها را برای سادگی می‌توانیم در جدول زیر خلاصه کنیم:

داده‌ها	x_1	x_2	...	x_n
ضریب (وزن)	f_1	f_2	...	f_n

در این صورت میانگین داده‌های بالا با احتساب ضرایب مربوط، به صورت زیر محاسبه می‌شود:

$$\text{میانگین وزنی} = \frac{f_1 \cdot x_1 + f_2 \cdot x_2 + \dots + f_n \cdot x_n}{f_1 + f_2 + \dots + f_n} = \frac{\sum_{i=1}^n f_i \cdot x_i}{\sum_{i=1}^n f_i}$$

مثال: نمرات دانش‌آموزی به صورت مقابل است. میانگین وزن دار نمرات او چقدر است؟

نمره	۱۹	۱۸	۱۵	۲۰
ضریب	۳	۲	۱	۲

$$\text{میانگین وزنی} = \frac{3 \times 19 + 2 \times 18 + 1 \times 15 + 2 \times 20}{3 + 2 + 1 + 2} = \frac{148}{8} = 18.5$$

تذکر میانگین وزن دار را نیز با نماد \bar{x} نشان می‌دهند و اگر بخواهند تمایزی بین میانگین معمولی و میانگین وزن دار قائل شوند، میانگین وزن دار را با نماد \bar{x}_w نشان می‌دهند.

تذکر با توجه به فرمول میانگین در مواردی که با داده‌های دسته‌بندی شده مواجه هستیم، به جای x_i ‌ها از مرکز هر دسته استفاده می‌کنیم.

مثال: به جدول مقابل توجه کنید، اکنون میانگین داده‌های این جدول را به دست می‌آوریم:

دسته	$[0, 5)$	$[5, 10)$	$[10, 15)$	$[15, 20]$
فراوانی	۲	۱	۷	۶

می‌دانیم مرکز دسته‌ها به ترتیب، $2.5, 7.5, 12.5$ و 17.5 است، پس میانگین وزنی این داده‌ها برابر است با:

$$\bar{x} = \frac{2 \times 2.5 + 1 \times 7.5 + 7 \times 12.5 + 6 \times 17.5}{2 + 1 + 7 + 6} = \frac{205}{16} = 12.8$$

درس	ادبیات	معارف	زبان	اختصاصی
درصد	۳	۹۰	۸۱	۷۰
ضریب	۴	۲	۳	۸

تست ۱۵: در جدول مقابل درصد نمرات داوطلبی نمایش داده شده است. اگر حداقل

میانگین برای پذیرش ۷۵ باشد، حداقل نمره‌ی ادبیات وی برای پذیرش کدام است؟

$$71 \quad (1)$$

$$72 \quad (2)$$

$$73 \quad (3)$$

$$74 \quad (4)$$

پاسخ: ابتدا به کمک فرمول میانگین وزن دار، میانگین این فرد را محاسبه می‌کنیم: (درصد درس ادبیات را a در نظر می‌گیریم).

$$\bar{x} = \frac{4a + 2(90) + 3(81) + 8(70)}{4 + 2 + 3 + 8} = \frac{4a + 983}{17}$$

اکنون می‌توان نوشت:

$$\frac{4a + 983}{17} \geq 75 \Rightarrow 4a + 983 \geq 1275 \Rightarrow a \geq 73$$

پس حداقل باید درس ادبیات را ۷۳ درصد بزند. بنابراین گزینه‌ی (۳) درست است.

ویژگی‌های میانگین

۱- اگر \bar{x} میانگین داده‌های x_1, x_2, \dots, x_n باشد در این صورت میانگین داده‌های $ax_1 \pm b, ax_2 \pm b, \dots, ax_n \pm b$ را با نماد $a\bar{x} \pm b$ نشان می‌دهیم و داریم: $a\bar{x} \pm b = \overline{ax \pm b}$

تساوی بالا به معنای آن است که اگر همه‌ی داده‌ها را با یک عدد ثابت جمع کنیم، میانگین اولیه نیز با همان عدد جمع می‌شود. همچنین اگر همه‌ی داده‌ها را در یک عدد ثابت ضرب کنیم، میانگین اولیه نیز در همان عدد ضرب می‌شود.

۲- اگر \bar{x} میانگین داده‌های x_1, x_2, \dots, x_n باشد، آنگاه مجموع اختلاف داده‌ها از میانگین برابر صفر است. یعنی:

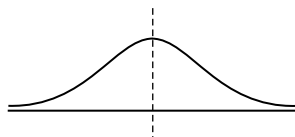
$$(x_1 - \bar{x}) + (x_2 - \bar{x}) + \dots + (x_n - \bar{x}) = \sum_{i=1}^n (x_i - \bar{x}) = 0$$

مثال: میانگین داده‌های ۱، ۴، ۱۷، ۱۸، ۲۰ برابر ۱۲ است. مجموع اختلاف داده‌ها از میانگین برابر است با:

$$(1-12) + (4-12) + (17-12) + (18-12) + (20-12) = -11 - 8 + 5 + 6 + 8 = 0$$

مقایسه‌ی میانگین و میانه

اگر میانه از میانگین کوچک‌تر باشد، به معنای آن است که اکثریت داده‌ها نسبت به میانگین کوچک‌تر هستند. اگر میانه از میانگین بزرگ‌تر باشد، به معنای آن است که اکثریت داده‌ها از میانگین بزرگ‌تر هستند.



میانه، مد، میانگین

در منحنی نرمال (شکل مقابل) مقادیر میانگین، میانه و مد با هم برابر هستند.

روش سریع محاسبه‌ی میانگین

فرض کنید می‌خواهیم میانگین اعداد ۱۵، ۱۵، ۱۵، ۱۶، ۱۷، ۱۷، ۱۸، ۱۹، ۲۰، ۲۰، ۲۰، ۲۰ را محاسبه کنیم. ابتدا، مقداری حدسی برای میانگین در نظر می‌گیریم، ما در این‌جا مقدار ۱۸ را در نظر می‌گیریم. اکنون کافی است $x_i - \bar{x}$ (تفاضل داده با میانگین) هر یک از داده‌ها را محاسبه کرده و میانگین آن‌ها را محاسبه کنیم و آن را به میانگین حدس زده شده اضافه کنیم.

$$\frac{4(20-18) + (19-18) + (18-18) + 2(17-18) + (16-18) + 3(15-18)}{12} = \frac{-4}{12} = -\frac{1}{3}$$

پس میانگین واقعی برابر است با: $18 - \frac{1}{3} = \frac{53}{3}$

نکته: اگر داده‌ها تعدادی از جملات متوالی از یک دنباله‌ی حسابی باشند، میانگین داده‌ها برابر جمله‌ی وسط (یا میانگین دو جمله‌ی وسط) است.

مثال: میانگین داده‌های ۱۳۱، ۱۲۹، ۱۲۷، ۱۲۵، ۱۲۳ برابر ۱۲۷ است.

مثال: اگر داده‌های $a, b, 19, 8, y$ و جملات متوالی یک دنباله‌ی حسابی باشند میانگین آن‌ها برابر $13/5 = \frac{8+19}{2}$ است.

فصل هفتم: شاخص‌های پراکندگی

فرض کنید نمرات دو کلاس A و B به صورت مقابل باشند:

$$\begin{cases} \text{نمرات کلاس A: } 10, 10, 10, 20, 20, 20 \\ \text{نمرات کلاس B: } 14, 14, 14, 16, 16, 16 \end{cases}$$

به راحتی می‌توان فهمید که میانگین نمرات هر دو کلاس برابر ۱۵ است. ولی این دو کلاس از جهتی با هم تفاوت دارند و آن «پراکندگی نمرات» است. توجه کنید پراکندگی نمرات در کلاس A خیلی بیش‌تر از کلاس B است. به عبارت دیگر در کلاس B نمرات به هم نزدیک‌تر هستند. شاخص‌هایی که میزان پراکندگی داده‌ها را در یک جامعه‌ی آماری مشخص می‌کنند، «شاخص‌های پراکندگی» نام دارند. این شاخص‌ها عبارت‌اند از دامنه‌ی تغییرات، واریانس، انحراف معیار و ضریب تغییرات که به بررسی آن‌ها می‌پردازیم.

همان‌طور که قبلاً گفتیم اختلاف بین بزرگ‌ترین و کوچک‌ترین داده را «دامنه‌ی تغییرات» می‌گوییم. به عنوان نمونه، دامنه‌ی تغییرات نمرات کلاس A، برابر ۱۰ و دامنه‌ی تغییرات نمرات کلاس B، برابر ۲ است. همان‌طور که ملاحظه می‌کنید، دامنه‌ی تغییرات کلاس A بیش‌تر از دامنه‌ی تغییرات کلاس B است.

واریانس داده‌های آماری x_1, x_2, \dots, x_n را با نماد σ^2 نمایش می‌دهند و به یکی از دو صورت زیر محاسبه می‌شود: (σ از حروف کوچک یونانی است و سیگما خوانده می‌شود. حرف بزرگ آن Σ است.)

$$\left\{ \begin{aligned} \sigma^2 &= \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n} = \frac{\sum (x_i - \bar{x})^2}{n} \\ \sigma^2 &= \frac{x_1^2 + x_2^2 + \dots + x_n^2}{n} - (\bar{x})^2 = \frac{\sum x_i^2}{n} - (\bar{x})^2 \end{aligned} \right.$$

مثال: برای محاسبه‌ی واریانس داده‌های ۱، ۲، ۳، ۴، ۵، با هر دو فرمول ابتدا میانگین را محاسبه می‌کنیم:

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

اکنون طبق فرمول اول واریانس می‌توان نوشت:

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5} = \frac{4+1+0+1+4}{5} = 2$$

حال از فرمول دوم واریانس این داده‌ها را محاسبه می‌کنیم:

$$\sigma^2 = \frac{1^2 + 2^2 + 3^2 + 4^2 + 5^2}{5} - (3)^2 = \frac{55}{5} - 9 = 2$$

تست ۱۶: اگر تفاضل داده‌ها از میانگین داده‌ها به صورت ۹، ۲، ۰، ۴- و a باشد، واریانس داده‌ها کدام است؟

۳۵ (۴)

۳۰ (۳)

۲۵ (۲)

۲۰ (۱)

پاسخ: می‌دانیم مجموع انحراف از میانگین داده‌ها همواره برابر صفر است، پس می‌توان نوشت:

$$a - 4 + 0 + 2 + 9 = 0 \Rightarrow a = -7$$

اکنون می‌توان واریانس را به صورت زیر محاسبه کرد:

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n} = \frac{(-7)^2 + (-4)^2 + 0^2 + 2^2 + 9^2}{5} = \frac{150}{5} = 30$$

بنابراین گزینه‌ی (۳) درست است.

تذکر همان‌طور که ملاحظه کردید فرقی نمی‌کند از کدام فرمول واریانس استفاده شود. فقط در مواقعی که مجموع مربعات داده‌ها به نحوی در مسأله مطرح باشد، بهتر است از فرمول دوم واریانس استفاده کنیم.

مثال: اگر مجموع ۲۰ داده‌ی آماری برابر ۱۴۰ و مجموع مربعات آن‌ها برابر ۱۹۰۰ باشد، برای محاسبه‌ی واریانس داده‌ها می‌توان نوشت:

$$\bar{x} = \frac{140}{20} = 7$$

اکنون می‌توان واریانس را از فرمول دوم محاسبه کرد: (توجه کنید که $x_1^2 + x_2^2 + \dots + x_n^2 = 1900$)

$$\sigma^2 = \frac{1900}{20} - 7^2 = 95 - 49 = 46$$

تذکر اگر داده‌ها دسته‌بندی شده باشند و نشان دسته‌های آن‌ها به ترتیب x_1, x_2, \dots, x_n باشد، واریانس آن‌ها را به صورت مقابل

محاسبه می‌کنیم:

$$\sigma^2 = \frac{\sum f_i (x_i - \bar{x})^2}{\sum f_i}$$

تست ۱۷: در جدول فراوانی تجمعی داده‌های دسته‌بندی شدهی مقابل، واریانس چقدر است؟

حدود دسته	۰-۲	۲-۴	۴-۶	۶-۸
فراوانی تجمعی	۱	۳	۱۲	۱۶

۱) ۱/۷۵ (۲)

۳) ۲/۲۵ (۴)

پاسخ: جدول فراوانی مطلق داده‌های مورد نظر به صورت مقابل است:

حدود دسته	۰-۲	۲-۴	۴-۶	۶-۸
فراوانی مطلق	۱	۲	۹	۴

بنابراین میانگین و واریانس این داده‌ها به صورت زیر به دست می‌آید:

$$\bar{x} = \frac{1 \times 1 + 2 \times 2 + 9 \times 5 + 4 \times 7}{1 + 2 + 9 + 4} = \frac{40}{16} = 2.5, \quad \sigma^2 = \frac{1(1-2.5)^2 + 2(2-2.5)^2 + 9(5-2.5)^2 + 4(7-2.5)^2}{16} = 2.5$$

بنابراین گزینه‌ی (۴) درست است.

نکاتی درباره‌ی واریانس

۱- اگر داده‌ها برابر باشند، واریانس آن‌ها صفر است و برعکس، یعنی اگر واریانس صفر باشد، در این صورت همه‌ی داده‌ها با هم برابرند.

مثال: واریانس داده‌های ۱۳۹۳، ۱۳۹۳، ۱۳۹۳، ۱۳۹۳ برابر صفر است.

مثال: اگر واریانس داده‌های ۷، ۱-ا، ۲-ب برابر صفر باشد می‌توان دریافت همه‌ی داده‌ها با هم برابرند. پس می‌توان نوشت:

$$a - 1 = 7 \Rightarrow a = 8, \quad b - 2 = 7 \Rightarrow b = 9$$

۲- واحد واریانس مجذور واحد داده‌هاست. به عنوان مثال اگر داده‌ها برحسب متر باشند، واریانس برحسب مترمربع است. (این موضوع، یکی از معایب استفاده از واریانس است.)

۳- اگر داده‌های x_1, x_2, \dots, x_n را با عدد ثابت b جمع (یا تفریق) کنیم، واریانس داده‌ها تغییری نمی‌کند. ولی اگر این داده‌ها را در عدد حقیقی a ضرب کنیم، واریانس اولیه در a^2 ضرب می‌شود. این مطلب را به صورت زیر می‌توان نوشت:

$$\sigma_{ax \pm b}^2 = a^2 \sigma_x^2$$

مثال: اگر واریانس داده‌های x_1, x_2, \dots, x_n برابر ۸ باشد، واریانس داده‌های $x_1 + 3, x_2 + 3, \dots, x_n + 3$ نیز برابر ۸ است.

مثال: اگر واریانس داده‌های x_1, x_2, \dots, x_n برابر ۵ باشد، واریانس داده‌های $2x_1, 2x_2, \dots, 2x_n$ برابر است با: $4 \times 5 = 20$

مثال: اگر واریانس داده‌های x_1, x_2, \dots, x_n برابر ۷ باشد، واریانس داده‌های $3x_1 + 13, 3x_2 + 13, \dots, 3x_n + 13$ برابر است با: $9 \times 7 = 63$

تست ۱۸: واریانس ۲۵ داده‌ی آماری برابر ۴ است. اگر ۵ تا از داده‌ها را که با میانگین برابر هستند از بین داده‌ها حذف کنیم، واریانس ۲۰ داده‌ی باقی‌مانده کدام است؟

۱) ۴ (۲) ۵ (۳) ۹ (۴) ۸

پاسخ: از آن‌جا که ۵ تا از ۲۵ داده‌ی اولیه برابر با میانگین است، می‌توان داده‌های اولیه را به صورت زیر در نظر گرفت:

$$\bar{x}, \bar{x}, \bar{x}, \bar{x}, \bar{x}, x_1, x_2, \dots, x_p$$

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n} \Rightarrow 4 = \frac{\sum (x_i - \bar{x})^2}{25} \Rightarrow \sum (x_i - \bar{x})^2 = 100$$

حال می‌توان نوشت:

با حذف ۵ داده‌ای که با میانگین برابر هستند، تغییری در میانگین و حاصل عبارت $\sum (x_i - \bar{x})^2$ رخ نمی‌دهد. پس واریانس ۲۰ داده‌ی

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{20} = \frac{100}{20} = 5$$

باقی‌مانده برابر است با:

بنابراین گزینه‌ی (۲) درست است.

همان طور که ملاحظه کردید واحد واریانس مجذور واحد داده‌هاست. برای رفع این مشکل، معمولاً در بررسی‌های آماری از جذر واریانس استفاده می‌کنند. جذر واریانس را «انحراف معیار» می‌گویند که آن را با σ نمایش می‌دهند و به کمک یکی از دو فرمول زیر محاسبه می‌شود:

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}} = \sqrt{\frac{\sum x_i^2}{n} - (\bar{x})^2}$$

تذکر: همان طور که از فرمول نیز مشخص است، واحد انحراف معیار با واحد داده‌ها یکسان است.

مثال: واریانس داده‌های ۱، ۲، ۳، ۴، ۵ برابر ۲ است، پس انحراف معیار آن‌ها برابر $\sqrt{2}$ است.

ویژگی‌های انحراف معیار

۱- اگر همه‌ی داده‌ها با هم برابر باشند، انحراف معیار برابر صفر است و برعکس.

۲- اگر به داده‌ها عدد ثابت b را اضافه (یا کم) کنیم، انحراف معیار آن‌ها تغییری نمی‌کند. ولی اگر همه‌ی داده‌ها را در عدد a ضرب کنیم، انحراف معیار آن‌ها در $|a|$ ضرب می‌شود. به طور خلاصه می‌توان نوشت:

$$\sigma_{ax \pm b} = |a| \sigma_x$$

مثال: اگر انحراف معیار داده‌های x_1, x_2, \dots, x_n برابر ۴ باشد، انحراف معیار داده‌های $-2x_1 + 3, -2x_2 + 3, \dots, -2x_n + 3$ برابر است با:

$$|-2| \times 4 = 8$$

تست ۱۹: اگر انحراف معیار داده‌های مجموعه‌ی $A = \{2, 4, 6, 8, 10, 12, 14\}$ برابر a باشد، آنگاه انحراف معیار داده‌های مجموعه‌ی

$$B = \{5, 6, 7, 8, 9, 10, 11\}$$

$$B \text{ کدام است؟}$$

$$4a \quad (۴)$$

$$\frac{a}{4} \quad (۳)$$

$$\frac{a}{2} \quad (۲)$$

$$2a \quad (۱)$$

پاسخ: توجه کنید که اگر داده‌های مجموعه‌ی A را نصف کنیم و به هر کدام از آن‌ها ۴ واحد اضافه کنیم داده‌های مجموعه‌ی B به دست می‌آیند. بنابراین اگر داده‌های مجموعه‌ی A را با x_i نمایش دهیم، داده‌های مجموعه‌ی B به صورت $\frac{1}{2}x_i + 4$ هستند پس انحراف

معیار داده‌های مجموعه‌ی B برابر است با $\frac{1}{2}a$. بنابراین گزینه‌ی (۲) درست است.

نکته: اگر داده‌های x_1, x_2, \dots, x_n دنباله‌ای حسابی با قدر نسبت d باشند، واریانس و انحراف معیار آن‌ها را می‌توان از فرمول زیر محاسبه کرد:

$$\sigma^2 = \frac{(n^2 - 1)}{12} d^2, \quad \sigma = |d| \sqrt{\frac{n^2 - 1}{12}}$$

مثال: مقادیر واریانس و انحراف معیار داده‌های ۱، ۴، ۷، ۱۰، ۱۳ را محاسبه می‌کنیم. توجه کنید که این داده‌ها دنباله‌ای حسابی با قدر نسبت $d=3$ هستند و تعداد آن‌ها $n=5$ است. پس می‌توان نوشت:

$$\sigma^2 = \frac{(25-1)}{12} \times 9 = 18, \quad \sigma = \sqrt{18} = 3\sqrt{2}$$

همان طور که ملاحظه کردید انحراف معیار با داده‌ها «هم‌واحد» است. به عنوان مثال انحراف معیار قد برحسب «متر» و انحراف معیار وزن برحسب «کیلوگرم» است. برای آن که بتوانیم پراکندگی‌های دو نوع داده، با واحدهای متفاوت را با هم مقایسه کنیم از شاخصی «بدون واحد» که به آن «ضریب تغییرات» گفته می‌شود، استفاده می‌کنیم. ضریب تغییرات را با CV نمایش می‌دهند و به صورت زیر تعریف می‌شود:

$$CV = \frac{\sigma}{\bar{x}}$$

پس ضریب تغییرات از تقسیم انحراف معیار بر میانگین محاسبه می‌شود.

مثال: انحراف معیار و میانگین داده‌های ۱, ۲, ۳, ۴, ۵ به ترتیب برابر $\sqrt{2}$ و ۳ است. پس ضریب تغییرات این داده‌ها برابر است با:

$$CV = \frac{\sigma}{\bar{x}} = \frac{\sqrt{2}}{3}$$

تذکر از آنجایی که ضریب تغییرات، معیاری برای میزان پراکندگی است، باید مثبت باشد. لذا ضریب تغییرات را فقط برای داده‌های مثبت تعریف می‌کنیم.

تست ۲۰: اگر x_1, x_2, \dots, x_n داده‌هایی با میانگین بزرگ‌تر از ۳ باشند، آن‌گاه ضریب تغییرات کدام دسته از داده‌های زیر بیش‌تر است؟

$$(۱) \quad 3x_1 + 2, 3x_2 + 2, \dots, 3x_n + 2$$

$$(۲) \quad 3x_1 + 1, 3x_2 + 1, \dots, 3x_n + 1$$

$$(۳) \quad x_1 + 4, x_2 + 4, \dots, x_n + 4$$

$$(۴) \quad 3x_1 - 1, 3x_2 - 1, \dots, 3x_n - 1$$

پاسخ: اگر میانگین و انحراف معیار داده‌های x_1, x_2, \dots, x_n را به ترتیب با \bar{x} و σ نمایش دهیم، می‌توانیم گزینه‌ها را به صورت زیر بررسی کنیم.

گزینه‌ی (۱): انحراف معیار داده‌های $3x_1 + 2, 3x_2 + 2, \dots, 3x_n + 2$ برابر است با: $CV_1 = \frac{3\sigma}{3\bar{x} + 2}$

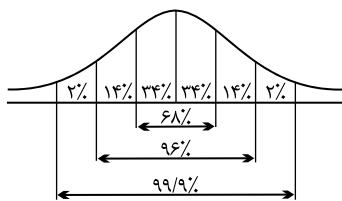
گزینه‌ی (۲): انحراف معیار داده‌های $3x_1 + 1, 3x_2 + 1, \dots, 3x_n + 1$ برابر است با: $CV_2 = \frac{3\sigma}{3\bar{x} + 1}$

گزینه‌ی (۳): انحراف معیار داده‌های $x_1 + 4, x_2 + 4, \dots, x_n + 4$ برابر است با: $CV_3 = \frac{\sigma}{\bar{x} + 4} = \frac{3\sigma}{3\bar{x} + 12}$

گزینه‌ی (۴): انحراف معیار داده‌های $3x_1 - 1, 3x_2 - 1, \dots, 3x_n - 1$ برابر است با: $CV_4 = \frac{3\sigma}{3\bar{x} - 1}$

اکنون با مقایسه‌ی مقادیر CV_1, CV_2, CV_3, CV_4 مقدار CV_4 از همه‌ی مقادیر بزرگ‌تر است (زیرا مخرج عبارت CV_4 از بقیه‌ی عبارت‌ها کوچک‌تر است!) بنابراین گزینه‌ی (۴) درست است.

پراکندگی در منحنی نرمال



منحنی نرمال مقابل را در نظر بگیرید. تقریباً ۶۸٪ داده‌ها در بازه‌ای به مرکز میانگین و شعاع انحراف معیار هستند.

به همین ترتیب تقریباً ۹۶٪ از داده‌ها در بازه‌ای به مرکز میانگین و شعاع دو برابر انحراف معیار هستند.

و سرانجام تقریباً ۱۰۰٪ از داده‌ها در بازه‌ای به مرکز میانگین و شعاع سه برابر انحراف معیار قرار دارند.