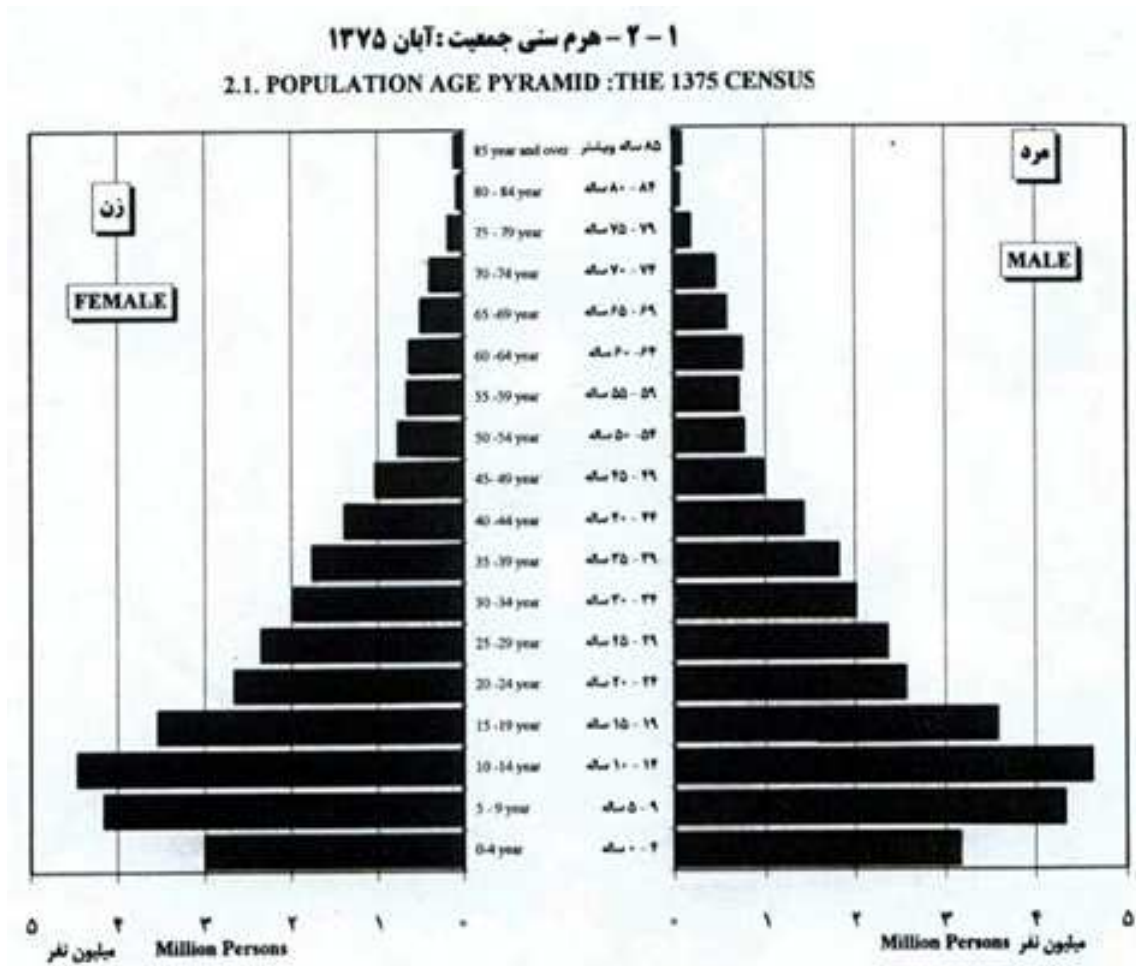




آمار و احتمالات شاہپور نصرتی



اردیبهشت ۸۶

فصل ۱

آمار

۱.۱ مقدمه

۱.۱.۱ تاریخچه

بشر از زمانی که حکومت کردن را آموخت، لازم دید تا اطلاعاتی را از جمعیت و منابع موجود تحت سلطه خود را برای عملکرد بهتر در دست داشته باشد. این گونه آگاهی از میزان کارگزار و مقدار منابع در تواریخ عهد عتیق، مصریان قدیم، بابلیان و رومیان در هزاره قبل از میلاد دیده می شود. با مرور زمان اهمیت جمع آوری فهرستهای اطلاعاتی بیشتر شده و تهیه فهرست هایی از نیروی انسانی و املاک متعلقه مهم تلقی شده تا جایی که در قرون وسطی این امر را زیر نظر کلیسا انجام می دادند. در سال ۱۳۰۰ در انگلستان بدستور ویلیام فاتح^۱ اولین کتاب اطلاعات از جمعیت و منابع اقتصادی نگاشته شد. بین سالهای ۱۶۶۱ الی ۱۶۶۹ میلادی کاپیتان جان گرانت آماری از زاد و ولد، مرگ و میر تهیه کرد و با ارائه آن در کتاب خود، نشان داد که در انگلیس عدد زنان بیشتر از مردان است.

برای اولین بار کلمه آمار^۲ از کلمه یونانی سیاستمدار^۳ توسط گاتفرید آچنوال^۴ (۱۷۱۹-۱۷۷۲) بکار رفت و در سال ۱۷۸۷ کلمه آمار در انگلستان توسط زیمرمان^۵ بکار گرفته شد و سپس سینکلیر^۶ در کتابش به نام شمارش آماری اسکاتلند^۷ آنرا عمومیت داد. در اواسط قرن هفدهم مجله آماری در فرانسه ایجاد شد و اولین سرشماری در این کشور در ۱۷۹۱ م. انجام گرفت، سپس طی بیست سال بعدی سرشماری های دیگری در کشورهای مختلف انجام گرفت و لزوم استفاده از آمار و تحلیل های عددی را نشان داد. تا آنجا که فرانسیس بیکن (۱۵۶۱-۱۶۲۶) اظهار داشت که تنها با مطالعه داده های پدیده های تجربی و طبیعی و تحلیل آنها می توان دانش بشری را سامان داد.

استفاده از آمار بعنوان شاخه ای از ریاضیات کاربردی، مختص یک شاخه بخصوص از علوم نبوده و امروزه می توان کاربرد آنرا در بسیاری از علوم مانند اقتصاد، روانشناسی، جامعه شناسی، علوم زیستی، فیزیک، شیمی، پزشکی، کشاورزی و ... یافت.

^۱ William the Conqueror

^۲ Statistics

^۳ Statista

^۴ Gottfried Achenwal

^۵ E.A.W. Zimmermann

^۶ John Sinclair

^۷ Statistical Account of Scotland

با استفاده از داده های آماری و بهره گیری از علم احتمالات می توان نتایج تجربی را تحلیل نمود و قاعده ای کلی را از آنها استنتاج کرد. شناخت آماری مسائل اکنون برای مدیران و صاحبان صنایع و موسسات از ضروریات بشمار رفته و به آنها کمک می کند که مشکلات را بشناسند و راه مقابله را مشکل را بیابند. آمار بمعنای یک علم تصمیم گیرنده^۸ زیربنای توسعه و پیشرفت بوده و در برنامه ریزی ها و تصمیمات کلان آینده یک کشور نقش عمده ای ایفا می کند.

از پیشگامان آمار لامبرت کتله^۹ (۱۸۷۴-۱۷۹۴) منجم و آمارگر بلژیکی که اولین بار درباره میانگین و انحراف و خطاهای آماری در زمینه های اجتماعی به بحث پرداخت. کتله قانون اعداد بزرگ را که ژاکوب برنولی (۱۷۰۵-۱۶۵۴) و سیمون پواسون (۱۸۴۰-۱۷۸۱) به آن دست یافته بودند مورد استفاده قرار داد. با مرور زمان آماردانان و ریاضیدانان به وسعت این علم افزودند و راه حل های آماری و احتمال را در آن وسعت بخشیدند. اکنون محاسبات داده ها توسط رایانه با سرعت و دقت بیشتری انجام می شود و کاوش ماشینی داده ها^{۱۰} جای محاسبات دستی را گرفته است. نرم افزارهای آماری اکنون در تحلیل داده ها نقش بسیار مثبتی را عهده دار بوده و تصمیمات نهائی را بهینه می کنند.

۲.۱.۱ تعریف مساله

علم آمار عبارتست از جمع آوری اطلاعات و پردازش آنها و سپس نتیجه گیری و استنباط یک یا چند صفت است. آمار در واقع ابزاری علمی است که اگر آنرا روی مجموعه ای از اطلاعات بکار ببریم آگاهی ما را نسبت به آنها افزایش می دهد. تحقق این امر، در ابتدا منوط به ایجاد ساختاری ریاضی بر روی مجموعه داده ها با شاخصی معین می باشد، بنابراین مساله را با زبان ریاضی بیان می کنیم تا بتوانیم از طرفندهای ریاضی برای حل مساله بهره بگیریم. در این حالت مساله را مدلسازی ریاضی نموده و از این طریق استنباط روی شاخص مورد نظر را انجام می دهیم. مدلسازی بهینه، مفاهیم ریاضی بکار برده شده را ساده تر و ابتدائی تر و نتیجه کار را به شخصیت اعضاء مجموعه مورد نظر نزدیکتر می سازد. پس از اتخاذ مدلسازی، روشهای آماری را برای تجزیه و تحلیل داده ها بکار می بریم.

مجموعه ای از انسانها یا اشیاء که مورد مطالعه آماری قرار می گیرند و حداقل در یک صفت مشترک می باشند را جامعه آماری نامیم. سنجش این صفت مشترک اعم از اینکه صفت کمی یا کیفی باشد، بر روی اعضاء جامعه اعمال می شود و بستگی به تک تک افراد آن دارد. پس جامعه آماری مجموعه ای از افراد یا اشیاء است که درباره اعضاء آن می خواهیم موضوع و یا موضوعاتی را مطالعه کنیم. جامعه ممکن است محدود (مثلاً افراد یک کلاس) یا نامحدود (مثلاً جمعیت کشور) باشد. اگر تک تک افراد یک جامعه را مورد مطالعه قرار دهیم، در اینصورت سرشماری کرده ایم و از آنجائیکه سرشماری بدلائل مختلف همیشه امکانپذیر نیست (بخصوص در جوامع نامحدود) یا قادر به سرشماری کل جامعه نیستیم و یا در مواردی ضرورتی به سرشماری نیست. در اینصورت بجای استفاده از کل اعضاء جامعه، تعداد محدودی از اعضاء را در نظر گرفته و شاخص مورد نظر را روی این تعداد دلخواه می سنجیم. این تعداد دلخواه از اعضاء جامعه را که بصورت تصادفی انتخاب می شوند، نمونه تصادفی گوئیم و نمونه زیرمجموعه ای از جامعه آماری محسوب می شود.

^۸ برخی نیز آمار را شاخه ای از نظریه تصمیم ها *Decision Theory* بشمار می آورند.

^۹ Lambert A.J. Quetelet

^{۱۰} Data Mining

۳.۱.۱ نمونه

کار آمار نمونه گیری از تعدادی از افراد جامعه و پس از تصمیم گیری، عمومیت دادن آن به کل جامعه است. اما نمونه گیری از افراد یک جامعه بدون خطا نخواهد بود، جمع آوری داده های مفید، صحیح و سالم و تجزیه و تحلیل دقیق آنها مسلماً نتایج مفید و درخور اعتباری را در بر خواهد داشت. بنابراین می بایست خطاهای مورد نظر در مسیر دستیابی به عمومیت شاخص را در نظر گرفت و تحقق این موضوع به علم احتمالات وابسته است. برای برآورد دقیق یک شاخص یا پارامتر روی جامعه بایستی از نمونه تصادفی استفاده کنیم. طبعاً شاخص مورد نظر روی اعضاء جامعه (نمونه) متفاوت بوده و صحت برآورد شاخص روی نمونه به تصادفی بودن اعضاء نمونه کاملاً وابسته است. اطلاعات حاصله اطلاعات خام محسوب شده و با روشهای آماری نمونه را بررسی می کنیم و نتیجه گیری می نمائیم. اکنون جامعه را کنار گذاشته و به شناخت دقیق یک نمونه می پردازیم. تعداد اعضاء یک جامعه را با N و تعداد اعضاء نمونه را با n نشان خواهیم داد. n را اندازه نمونه نیز گویند.

۴.۱.۱ متغیر یک نمونه

گفتیم که اعضاء جامعه دارای صفتی مشترک هستند، شاخصی را که می بایست روی یک نمونه اندازه بگیریم را متغیر نامیم. این شاخص می تواند سن یا قد افراد یا دوام محصول یک کارخانه یا ... باشد. متغیر تصادفی به دو دسته متغیر کمی و متغیر کیفی تقسیم بندی می شود. متغیرهای کمی برخلاف متغیرهای کیفی، متغیرهایی قابل اندازه گیری اند. در اینصورت اگر بخواهیم متغیر کیفی را از قبیل شادی و اندوه و میزان هوش و ... بسنجیم می بایست آن را تبدیل به متغیر کمی کنیم. متغیر کمی ممکن است پیوسته باشد و متغیر پیوسته یک متغیر کمی است که اندازه های ممکن یک فاصله را اختیار کند، عبارتی دیگر اگر دو مقدار را بتواند اختیار کند، هر مقدار بین آنها را نیز می تواند اختیار نماید. به متغیر کمی که پیوسته نباشد متغیر گسسته گوئیم. متغیرهای کیفی که در آنها نوعی ترتیب طبیعی وجود دارد متغیرهای کیفی ترتیبی می گوئیم و متغیر کیفی که ترتیبی نباشد متغیر کیفی اسمی می گوئیم.

به اعضاء یک نمونه داده آماری یا بطور خلاصه داده می گوئیم و آن را با x_i نشان می دهیم. در واقع داده های آماری بجای متغیر یک نمونه می نشینند. فرض کنید می خواهیم میزان سن دانشجویان یک دانشگاه را مورد سنجش قرار دهیم. در اینجا دانشجویان این دانشگاه جامعه آماری خواهند بود و شاخص مورد سنجش سن دانشجویان است که متغیر نمونه محسوب خواهد شد و داده های حاصل از نمونه برداری تصادفی نیز بجای x_i قرار می گیرند.

۵.۱.۱ جدول فراوانی داده ها

بعد از نمونه گیری و بدست آوردن داده ها، لازم است تا میانگین و انحراف معیار داده ها - که بعداً راجع به آنها صحبت خواهیم کرد - را بدست آوریم. داده ها را پس از جمع آوری و مرتب نمودن، در قالب یک جدول می ریزیم که به آن جدول آماری یا جدول فراوانی داده ها گوئیم. در ستون اول این جدول x_i (داده های نمونه) را جای داده و در ستون دوم فراوانی مطلق آنها را یادداشت می کنیم. فراوانی مطلق یک داده x_i برابر با تعداد دفعاتی است که آن داده تکرار شده است و آنرا با f_i نشان می دهیم. اگر f_i فراوانی مطلق داده i ام و تعداد داده های نمونه n باشد، مقدار کسر $\frac{f_i}{n}$ را \bar{f}_i فراوانی نسبی دسته i -ام می

گوئیم و آنرا در ستون \bar{f}_i جای می دهیم. فراوانی تجمعی F_i هر داده، برابر است با مجموع فراوانی آن داده و داده‌های قبل از آن. فراوانی های نسبی و تجمعی را می توان در انتهای جدول نیز جای داد.

x_i	f_i	F_i	\bar{f}_i	$f_i x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$f_i(x_i - \bar{x})^2$...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

میانگین داده‌های یک نمونه، همان متوسط داده‌هاست که با \bar{x} نشان می دهیم و برای n داده x_1 و x_2 و ... و x_n برابرست با

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

فرمول میانگین برای n داده x_1 و x_2 و ... و x_n با فراوانی های f_1 و f_2 و ... و f_n بصورت $\bar{x} = \frac{\sum f_i x_i}{n}$ خلاصه می شود. واریانس یک نمونه از n داده x_1 و x_2 و ... و x_n با فراوانی های f_1 و f_2 و ... و f_n عبارتست از:

$$S_x^2 = \frac{\sum f_i(x_i - \bar{x})^2}{n - 1}$$

که آنرا واریانس نمونه ای نیز گوئیم. برای بدست آوردن این مقدار ستون $x_i - \bar{x}$ و سپس ستون های $(x_i - \bar{x})^2$ و $f_i(x_i - \bar{x})^2$ را به جدول اضافه می کنیم تا کار محاسبه آسانتر شود. بعد از بدست آوردن واریانس، مقدار انحراف معیار را که جذر آن خواهد بود، بدست می آوریم پس $S_x = \sqrt{S_x^2}$. به مثال زیر دقت کنید:

مثال ۱.۱.۱ در نمونه گیری از دانشجویان یک دانشگاه، نمرات ۲۰ نفر از آنها در درس آمار بصورت زیر بوده است:

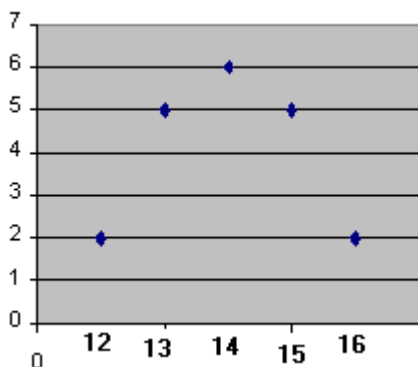
۱۲، ۱۶، ۱۵، ۱۴، ۱۳، ۱۴، ۱۵، ۱۳، ۱۲، ۱۳، ۱۴، ۱۴، ۱۵، ۱۴، ۱۳، ۱۵، ۱۶، ۱۳، ۱۴، ۱۵

مطابق آنچه گفته شد جدول فراوانی داده ها چنین می نویسیم:

x_i	f_i	F_i	\bar{f}_i	$f_i x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$f_i(x_i - \bar{x})^2$
۱۲	۲	۲	۰/۱۰	۲۴	-۲	۴	۸
۱۳	۵	۷	۰/۲۵	۶۵	-۱	۱	۵
۱۴	۶	۱۳	۰/۳۰	۸۴	۰	۰	۰
۱۵	۵	۱۸	۰/۲۵	۷۵	-۱	۱	۵
۱۶	۲	۲۰	۰/۱۰	۳۲	-۲	۴	۸
جمع	۲۰		۱	۲۸۰			۲۶

از ستون پنجم می نویسیم $\bar{x} = \frac{\sum f_i x_i}{n} = \frac{280}{20} = 14$ و از ستون آخر می توان واریانس را بدست آورد:

$$S_x^2 = \frac{\sum f_i(x_i - \bar{x})^2}{n - 1} = \frac{26}{19} = 1/37 \Rightarrow S_x = 1/17$$



دقت کنید که سطر آخر جمع فراوانی های مطلق می بایست برابر تعداد داده ها و مجموع فراوانی های نسبی می بایست ۱ شود. همچنین آخرین مقدار فراوانی تجمعی n خواهد بود. نمودار نقطه‌ای این داده ها بصورت مقابل است.

۲.۱ تحلیل داده‌ها

پس از مرتب سازی داده ها و جدول فراوانی، برای تصمیم گیری نهائی لازم است پارامترها و معیارهای مختلفی را روی آنها پیاده نمود تا تصمیم گیری بر اساس این معیارها انجام شود. تحلیل اطلاعات با دو نوع شاخص انجام می شود: شاخص های عددی و شاخص های هندسی. شاخص عددی مقدار است که از داده ها حاصل شده و خصوصیات کمی جامعه بر اساس آنها تعیین می گردد. انواع شاخص های عددی عبارتند از شاخص های مرکزی، شاخص های پراکندگی، چولگی و کشیدگی. شاخص های هندسی نیز همان نمودارها هستند.

۱.۲.۱ شاخص های مرکزی

شاخصی که تمرکز صفت مورد نظر را در بین داده ها مشخص می کند شاخص مرکزی نامیده می شود. با توجه به داده های یک نمونه کاملاً تصادفی، دیده می شود که قرارگیری تعدادی از داده ها تمرکز و ویژگی خاصی نسبت به بقیه ایجاد نموده است. بنابراین شاخص مرکزی، محل تمرکز اکثر داده ها است. مهمترین شاخص های مرکزی عبارتند از میانگین، میانه و مد. مد یا نما (M) داده ای است که دارای بیشترین فراوانی است. در مثال ۱.۱.۱ داده ۱۴ دارای بیشترین فراوانی بوده و بنابراین $M = 14$. در مواردی که داده ها دارای دو فراوانی بیشین هستند - وضعیت داده های دو مدی - اگر داده ها در کنار هم واقع شوند میانگین آنها را بعنوان مد در نظر می گیریم و در حالتیکه داده ها در کنار هم نباشند تمام آنها مد محسوب می شوند. اگر داده ها همه فراوانی برابر داشته باشند مد نداریم. منظور از میانه (m) عبارتست از داده وسط. بنابراین برای بدست آوردن میانه، پس از مرتب کردن داده ها، مقداری را که تعداد داده های بعد از آن با تعداد داده های قبل از آن برابر است را میانه می نامیم. اگر تعداد داده ها فرد باشد داده وسطی میانه خواهد بود $m = x_{\frac{n+1}{2}}$ و در صورت زوج بودن تعداد داده ها، میانه عبارتست از میانگین دو داده وسط $m = \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1})$. در مثال ۱.۱.۱ اگر داده ها را مرتب کنیم:

۱۲، ۱۲، ۱۳، ۱۳، ۱۳، ۱۳، ۱۳، ۱۳، ۱۴، ۱۴، ۱۴، ۱۴، ۱۴، ۱۴، ۱۵، ۱۵، ۱۵، ۱۵، ۱۵، ۱۶، ۱۶
چون تعداد داده ها زوج است بنابراین $m = 14$.

مثال ۱.۲.۱ شاخص های مرکزی را برای داده های زیر بیابید ۵، ۳، ۲، ۴، ۳، ۴، ۳، ۲، ۳، ۴، ۶، ۵

حل. داده ها را مرتب می کنیم ۲، ۲، ۳، ۳، ۳، ۳، ۴، ۴، ۵، ۵، ۶

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{\sum_{i=1}^{10} x_i}{10} = \frac{2+2+3+3+3+3+4+4+5+5+6}{10} = 3.7$$

داده ۳ دارای بیشترین فراوانی است بنابراین $M = 3$ و چون تعداد داده ها زوج هستند پس $m = \frac{3+4}{2} = 3.5$

تمرین ۱.۲.۱ کلاسی. در یک نمونه ۲۰ نفره دانش آموزان یک مدرسه نمونه سن آنها چنین بدست آمده است:

۷ و ۸ و ۱۰ و ۸ و ۹ و ۶ و ۷ و ۷ و ۶ و ۷ و ۷ و ۸ و ۷ و ۸ و ۹ و ۸ و ۹ و ۷ و ۸ و ۶ و ۷ و ۸ و ۷

پس از رسم جدول فراوانی داده ها، میانگین و واریانس و شاخص های مرکزی را بیابید.

۲.۲.۱ شاخص های پراکندگی

با اینکه میانگین شاخص مهمی در تصمیم گیری است، ولی بالا بودن میانگین در یک نمونه به معنی بالا بودن کل داده ها نیست مثلاً بالا بودن میانگین نمرات یک کلاس ممکن است تنها بخاطر چند نمره ۲۰ باشد و سایر نمرات نزدیک ۱۰ باشند و بنابراین داده ها می بایست مثلاً با میانه نیز مقایسه شوند. معیاری که می تواند تفاوت داده ها و میزان این تفاوت و بخصوص دوری یا نزدیکی آنها را از میانگین برای ما بیان کند شاخص پراکندگی نامیم. مهمترین شاخص های پراکندگی عبارتند از واریانس، انحراف معیار و ضریب تغییرات. بعنوان بهترین شاخص، می توان به واریانس اشاره کرد که به معنی تفاوت و تغییر است و عبارتست از میانگین مجذور انحرافات از میانگین. اگر داده ها نزدیک \bar{x} باشند، واریانس کوچک می شود، بنابراین واریانس معیار خوبی برای سنجش پراکندگی داده ها از \bar{x} است. اما چون واریانس واحد داده ها را مربع می کند، از جذر آن که انحراف معیار یا انحراف استاندارد نامیده می شود استفاده می کنیم و با s_x نشان می دهیم. واحد انحراف معیار همان واحد x_i است و پراکندگی نسبی داده ها را مشخص می کند. در واقع انحراف معیار، انحراف از میانگین \bar{x} است و در حالت خاص اگر داده ها همه برابر باشند واریانس آنها صفر است که دارای پراکندگی نیستند.

مثال ۲.۲.۱ از دانشجویان دو کلاس آمار نمره های زیر بدست آمده است:

$$A: 14 \text{ و } 18 \text{ و } 14 \text{ و } 15 \text{ و } 13 \text{ و } 17 \text{ و } 14$$

$$B: 11 \text{ و } 15 \text{ و } 20 \text{ و } 18 \text{ و } 14 \text{ و } 17 \text{ و } 10 \text{ و } 19 \text{ و } 11$$

با اینکه میانگین نمره ها در این هر دو برابر ۱۵ است ولی پراکندگی کلاس B بیشتر است و $S_A = 1/82$ و $S_B = 3/74$. معیار دیگری برای میزان پراکندگی ضریب تغییرات V_x است که عبارتست از خارج قسمت انحراف معیار بر میانگین $V_x = \frac{S_x}{\bar{x}}$. عبارتی دیگر ضریب تغییرات عبارتست از میزان پراکندگی بازای یک واحد از میانگین. این ضریب مستقل از واحد است.

مثال ۳.۲.۱ کارخانه ای دو نوع لامپ ۱۰۰ واتی تولید می کند. نوع A دارای میانگین عمر ۹۰۰۰ ساعت و انحراف استاندارد ۱۸۰۰ ساعت و نوع B دارای میانگین عمر ۱۰۰۰۰ ساعت و انحراف استاندارد ۱۲۰۰ ساعت است. کدام نوع لامپ بهتر است؟

حل. ضریب تغییرات را برای هر دو نوع حساب می کنیم:

$$V_A = \frac{1800}{9000} = 0/20 = \%20 \quad V_B = \frac{1200}{10000} = 0/12 = \%12$$

نوع B بهتر است زیرا هم میانگین بیشتری دارد و هم دارای ضریب تغییر کمتری است.

تمرین ۲.۲.۱ کلاسی. ضریب تغییرات دو کلاس مذکور در مثال ۲.۲.۱ را یافته، این دو کلاس را با هم مقایسه کنید.

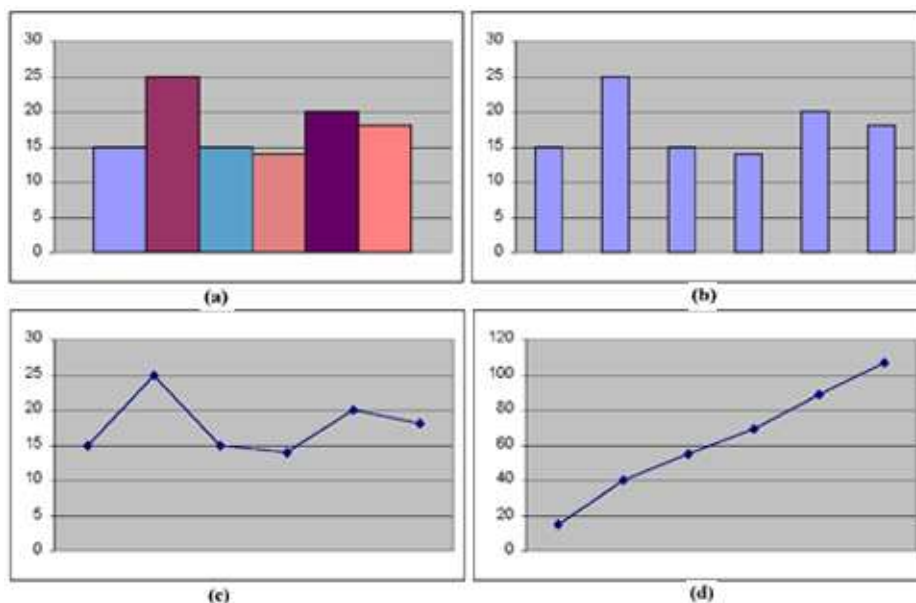
۳.۲.۱ نمودارها

نمودار نمایش هندسی داده ها و روشی برای تحلیل اطلاعات است. نمودارها شاخص های هندسی و کمکی برای ایجاد تصویر ذهنی از داده ها بوده و برای تصویرسازی از یک نمونه، نمودار تا حدی میزان جایگیری و پراکندگی داده ها را مشخص می کند و بدین ترتیب شاخصی برای تصمیم گیری محسوب می گردد. انواع مختلفی از نمودارها وجود دارد که مهمترین آنها را بیان می

کنیم. البته تنوع نمودار برای ارائه شکل قرارگیری داده‌ها زیاد است (تمرین ۲۱) که تنها چند نوع آن دارای کاربرد بیشتری است. از انواع مهم نمودار می‌توان به موارد ذیل اشاره نمود:

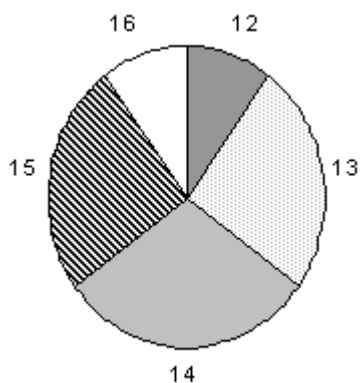
الف) نمودار بافتی یا هیستوگرام یا نمودار ستونی که بیشتر برای نشان دادن داده‌های پیوسته و داده‌های دسته بندی شده بکار می‌رود و از انواع دیگر ساده تر است، در این نوع نمودار، دسته‌ها عرض مستطیل‌ها بوده و ارتفاع مستطیل‌ها را متناسب با فراوانی دسته‌ها انتخاب می‌کنیم (تصویر a).

ب) نمودار میله‌ای مانند نمودار بافتی ولی با ستونهای مجزاست و عرض میله‌ها نیز برابر است (تصویر b).



ج) نمودار چندبرفراوانی یا نمودار خطی (چندضلعی) که در آن نقاط داده‌ها را بهم متصل می‌کنیم و این نوع نمودار علمی تر بوده و می‌توان آنرا با نمودار نرمال که بعداً خواهیم آورد مقایسه نمود، نقاط نمودار مانند نقاط صفحه مختصات ترسیم می‌شوند که در آنها x_i طول نقاط و f_i عرض نقاط را تشکیل خواهند داد. در این نوع نمودار، اگر وسط بالای مستطیل‌ها را بصورت نقطه‌ای در نظر بگیریم و آنها را بهم وصل کنیم، در این صورت آنرا نمودار چندبرفراوانی نامیم (تصویر c).

د) نمودار فراوانی تجمعی که بجای فراوانی مطلق در ستون عرضها از فراوانی تجمعی استفاده می‌کنیم، این نمودار صعودی بوده در این حالت اگر متغیر پیوسته باشد توزیعی تراکمی یا اوجیو بوجود می‌آید و این نمودار در مواردی مختلفی مانند پزشکی کاربرد دارد (تصویر d).



ه) چندبرفراوانی نسبی در این نوع بجای فراوانی مطلق، نمودار فراوانی نسبی

را رسم کنیم، چون اطلاعات با کل نمونه مقایسه می‌شود این نمودار رسم شده اطلاعات جالبی را بیان خواهد داشت.

و) نمودار دایره‌ای (کلوچه‌ای) که برای تجسم داده‌ها و مقایسه آنهاست. در این نوع تقسیم بندی که در خلاف جهت عقربه‌های ساعت انجام می‌شود، ترسیم فراوانی نسبی داده‌ها بر روی یک دایره خواهد بود. مقدار زاویه در نظر گرفته شده برابر

$f \times 360 = 2\pi f$ (درجه) خواهد بود. نمودار دایره‌ای برای داده‌های مثال ۱.۱.۱ نشان داده شده است.

۴.۲.۱ دسته‌بندی داده‌ها

اگر تنوع داده‌ها زیاد بوده و فراوانی هر کدام از آنها کم باشد، برای جلوگیری از حجم زیاد محاسبات، داده‌ها را دسته‌بندی می‌کنیم. در این روش ابتدا دامنه تغییرات R را بدست آورده و سپس تعداد دسته‌ها و طول هر دسته را معین می‌کنیم و بعد جدول را مطابق قبل رسم می‌کنیم. دامنه تغییرات R یک متغیر عبارتست از طول بازه‌ای که متغیر در آن تغییر می‌کند و آن عبارتست از تفاضل بیشترین داده و کمترین داده، می‌نویسیم $R = x_{max} - x_{min} + 1$. تعیین تعداد دسته‌ها T قانون مشخصی ندارد و نظرات متفاوتی درباره آن ارائه شده، ولی غالباً ما آنرا بین ۵ تا ۲۰ دسته انتخاب می‌کنیم. یک فرمول از استارچیس^{۱۱} بصورت $T = 1 + 3/322 \log n$ پیشنهاد شده است. بهرطریق پس از انتخاب T ، طول دسته‌ها از رابطه زیر بدست می‌آید:

$$c = \frac{R}{T} = \text{طول دسته}$$

مجموع طول دسته‌ها نباید از دامنه‌ی تغییرات کمتر باشد. بهترین حالت انتخاب هر دسته بصورت نیم بازه $[a_i, b_i)$ است که a_i را کران پائین دسته و b_i را کران بالای دسته می‌نامیم. نماینده دسته نیز عبارتست از $x_i = \frac{a_i + b_i}{2}$ و مسلماً تفاضل دو کران پائین متوالی یا دو کران بالای متوالی طول دسته c خواهد بود. نکته‌ای که در انتخاب دسته‌ها باید دقت شود این است که کران بالا و پائین هر دسته می‌بایست بگونه‌ای انتخاب شود که هر داده تنها در یک دسته شمارش شود. بفرض اگر داده‌ای در دو دسته قرار گرفت می‌بایست از تمامی کرانها مقداری (مثلاً ۵/۰) کم شود تا داده تنها در یک دسته قرار گیرد، یا اینکه دسته بعدی بفاصله کمی از کران بالای دسته قبلی شروع شود $b_i \leq a_{i+1}$ بدین ترتیب بین طبقات تداخل داده ایجاد نخواهد شد. مثال زیر را ببینید.

مثال ۴.۲.۱ در نمونه گیری از سن دانشجویان یک دانشگاه، سن ۵۰ نفر از آنها بصورت زیر بوده است:

۲۱	۱۸	۲۲	۱۹	۲۱	۳۰	۲۰	۱۸	۳۵	۳۰
۲۰	۲۵	۱۸	۱۷	۲۵	۲۷	۲۲	۲۲	۲۰	۲۱
۲۳	۲۰	۲۱	۱۸	۱۹	۲۱	۲۰	۱۹	۳۳	۲۹
۲۰	۱۸	۳۵	۲۲	۲۰	۳۰	۲۵	۲۴	۲۲	۱۷
۱۹	۱۸	۲۰	۱۹	۲۵	۲۶	۲۰	۲۱	۱۸	۲۰

مطابق این نمونه ۵۰ تائی داریم $R = x_{max} - x_{min} + 1 = 35 - 17 + 1 = 19$ اگر تعداد دسته را ۵ انتخاب کنیم پس

$$c = \frac{19}{5} \sim 4 = \text{طول دسته}$$

با انتخاب دسته‌ها بصورت ۲۱ - ۲۵ و ۱۷ - ۲۱ و ... انتخاب کنیم، داده ۲۱ متعلق به دو دسته خواهد شد برای رفع این دودستگی، دسته‌ها را با اختلاف ۵/۰ انتخاب می‌کنیم و جدول فراوانی داده‌ها بصورت زیر درمی‌آید:

دسته‌ها	x_i دسته‌نماینده	f_i	F_i	\bar{f}_i	$f_i x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$f_i(x_i - \bar{x})^2$
۱۶/۵ - ۲۰/۵	۱۸/۵	۲۴	۲۴	۰/۴۸	۴۴۴	-۳/۹۲	۱۵/۳۷	۳۶۸/۷۹
۲۰/۵ - ۲۴/۵	۲۲/۵	۱۳	۳۷	۰/۲۶	۲۹۳	-۰/۰۸	۰/۰۱	۰/۰۸۳۲
۲۴/۵ - ۲۸/۵	۲۶/۵	۶	۴۳	۰/۱۲	۱۵۹	۴/۰۸	۱۶/۶۵	۹۹/۸۷۸
۲۸/۵ - ۳۲/۵	۳۰/۵	۴	۴۷	۰/۰۸	۱۲۲	۸/۰۸	۶۵/۲۹	۲۶۱/۱۵
۳۲/۵ - ۳۶/۵	۳۴/۵	۳	۵۰	۰/۰۶	۱۰۴	۱۲/۰۸	۱۴۵/۹۳	۴۳۷/۷۸
جمع		۵۰		۱	۱۱۲۱			۱۱۶۷/۷

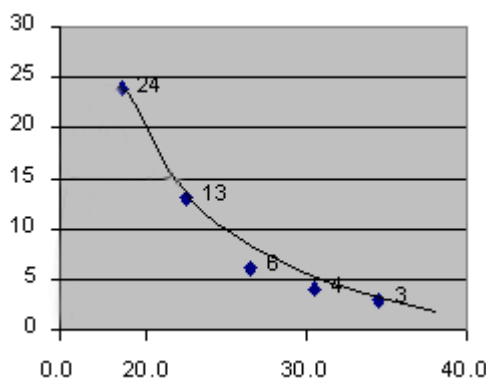
^{۱۱}Sturges، این فرمول از $n = 2^{T-1}$ نتیجه شده است.

از جدول بالا می توان مقادیر زیر را محاسبه نمود:

$$\bar{x} = \frac{\sum x_i f_i}{n} = \frac{1121}{50} = 22.42, \quad S_x^2 = \frac{\sum f_i (x_i - \bar{x})^2}{n-1} = \frac{11677}{49} = 23.83 \Rightarrow S_x = 4.88$$

$$V_x = \frac{S_x}{\bar{x}} = \frac{4.88}{22.42} = 0.218 = 21.8\%$$

نمودار داده‌ها بصورت مقابل است.



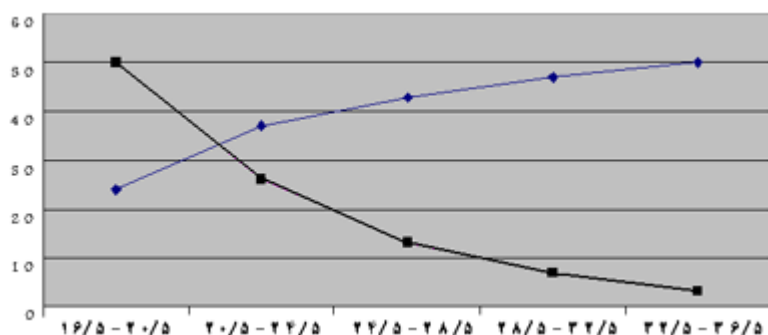
میانه در داده های دسته بندی شده:

برای بدست آوردن میانه با استفاده از جدول فراوانی داده ها اگر n تعداد داده ها و a_i حد پائین دسته ای که فراوانی تجمعی آن بلافاصله از $\frac{n}{2}$ بیشتر یا مساوی باشد و f_i فراوانی مطلق دسته i - ام و F_{i-1} فراوانی تجمعی دسته ماقبل دسته i - ام باشد سپس میانه عبارتست از ^{۱۲}

$$m = a_i + \frac{\frac{n}{2} - F_{i-1}}{f_i} c$$

که c فاصله طبقات است ($c = a_i - a_{i-1}$). در مثال قبل چون $\frac{n}{2} = 25$ پس فراوانی تجمعی $F_2 = 37$ از آن بیشتر یا مساوی است و $i = 2$ و $c = 20/5 - 16/5 = 4$ طبق فرمول بالا

$$m = 20/5 + \frac{25 - 24}{13} 4 = 20/8 \in [20/5, 24/5]$$



اگر داده ها را مرتب کنید و از روش اصلی داده^{۱۲} وسط را بیابید میانه واقعی برابر $m = 21$ خواهد بود. برای بدست آوردن میانه با استفاده از نمودار کفایت فراوانی تجمعی صعودی و نزولی را در یک صفحه مختصات رسم کنیم، تقاطع آنها میانه خواهد بود.

مد در داده های دسته بندی شده: برای بدست آوردن مد در داده های دسته بندی از فرمول زیر استفاده می کنیم:

$$M = a_i + \frac{d_1}{d_1 + d_2} c$$

با یافتن دسته ای که دارای بیشترین فراوانی است (دسته یا رده^{۱۲} نمائی)، اگر a_i کران پائین دسته و c فاصله طبقات باشد، d_1 تفاوت فراوانی دسته با دسته قبلی و d_2 تفاوت فراوانی دسته با دسته بعدی است. برای مثال بالا

$$M = a_i + \frac{d_1}{d_1 + d_2} c = 16/5 + \frac{24}{24 + 11} 4 = 19/24$$

و اگر داده ها را مرتب کنیم $M = 20$ خواهد بود (البته این مقدار را تنها برای مقایسه و درستی مطلب بدست آورده ایم).

^{۱۲} کندال *kendall* آماردان انگلیسی می گوید در حالت n زوج بجای $\frac{n}{2}$ از $\frac{n+1}{2}$ استفاده کنید.

۵.۲.۱ گشتاورها

برای n داده x_1 و x_2 و \dots و x_n با فراوانی های f_1 و f_2 و \dots و f_n ، مقدار m_r گشتاور مرکزی مرتبه r و m'_r گشتاور مرتبه r و چنین تعریف می شوند:

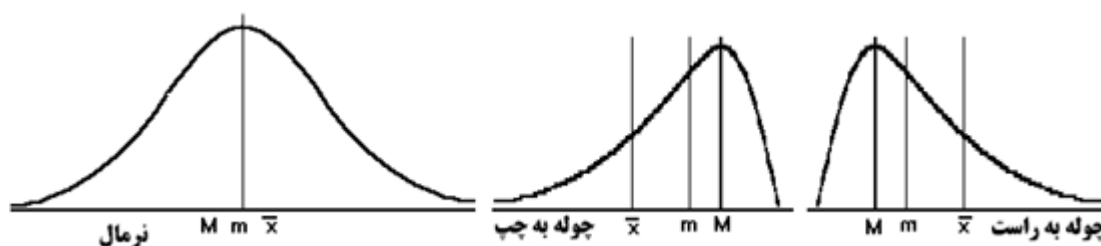
$$m_r = \frac{1}{n} \sum_{i=1}^n f_i (x_i - \bar{x})^r \quad , \quad m'_r = \frac{1}{n} \sum_{i=1}^n f_i x_i^r$$

می توان دید که $m'_0 = m_0 = 1$ و $m'_1 = \bar{x}$ و همچنین $m_2 \simeq S_x^2$.

تمرین ۳.۲.۱ کلاسی. با استفاده از جدول فراوانی مثال ۴.۲.۱، دو ستون جدید بعنوان ستون های $f_i x_i^2$ و $f_i (x_i - \bar{x})^3$ اضافه کرده و با محاسبه مجموع این ستونها مقادیر گشتاورهای m_3 و m'_3 را بدست آورید.

۶.۲.۱ چولگی و کشیدگی

در یک جامعه نرمال، با بررسی داده ها و رسم شاخص هندسی دیده می شود که منحنی جامعه بصورت منحنی زیر است. این منحنی به منحنی زنگوله ای یا منحنی نرمال مشهور بوده و شکل جامعه کاملاً نرمال را نشان می دهد. در این توزیع متقارن، مقدار مد، میانه و میانگین برابر بوده و روی محور تقارن نمودار واقع می شوند. از آنجاکه عملاً هیچ نمونه کاملاً نرمالی یافت نمی شود، با رسم نمودار داده های یک نمونه و مقایسه آن با نمودار نرمال، میزان انحراف در داده ها را می توان دریافت. بنابراین پس از ترسیم نمودار چندبفراوانی، معمولاً شکل آن شبیه منحنی نرمال، چاوله، J -شکل (مانند مثال ۴.۲.۱) یا در برخی موارد U -شکل خواهد بود. اگر معیار ما نرمال بودن جامعه آماری از لحاظ صفت مورد نظر باشد، پس می توانیم منحنی داده ها را با منحنی نرمال سنجیده و مقدار انحراف داده ها را بدست آوریم. این مقدار انحراف که مقدار دوری داده ها را از توزیع متقارن نشان می دهد تحت عنوان چولگی شناخته می شود.



دو نمودار فوق که با کمی مچاله شدن داده ها به نمودار نرمال شبیهند، تحت عنوان منحنی چوله به راست و چوله به چپ مطرحند. وقتی در این منحنی ها دامنه به سمت چپ کشیده تر باشد آنرا چوله به چپ و اگر دامنه به سمت راست کشیده تر باشد آنرا چوله به راست گوئیم. پس چوله بر راست یعنی اکثر داده ها در طرف چپ جمع شده اند. میزان چولگی یک نمودار را توسط مقیاسهای چولگی اول پیرسن b_1 و چولگی دوم پیرسن b_2 و ضریب گشتاوری چولگی g مطابق فرمول های زیر بدست می آوریم.

$$b_1 = \frac{\bar{x} - M}{S_x} \quad , \quad b_2 = 3 \frac{\bar{x} - m}{S_x} \quad , \quad g = \frac{m_3}{S_x^3}$$

برای نمودار نرمال $m_2 = 0$ و در نتیجه g صفر است. در نمودارهای دیگر اگر این سه معیار چولگی مثبت باشند نمودار چوله به راست و اگر منفی باشند نمودار چوله به چپ خواهد بود. از آنجائیکه مد در برخی مسائل بدقت محاسبه نمی شود (مثل جداول چندمدی)، بجای ضریب b_1 ، می توان از ضریب b_2 استفاده کرد که از میانه بهره می برد. این دو کمابیش با هم برابرند و از نظر تجربی دیده شده که $b_1, b_2 < 1$ است.

تمرین ۴.۲.۱ در مثال ۴.۲.۱ که نمودار چوله به راست است، ضرایب چولگی را محاسبه کنید.

همچنین ممکن است داده های مختلف از لحاظ میانگین و پراکندگی و حتی چاولگی مساوی باشند ولی از لحاظ شکل ظاهری متفاوت باشند مثلاً بسمت بالا کشیده تر باشد، این مقدار کشیدگی (برجستگی یا پخی) با یک ضریب بصورت $\beta_2 = \frac{m_2}{S_x^2}$ بیان می شود که آنرا ضریب کشیدگی پیرسن نامند. مقدار کشیدگی برای منحنی نرمال محاسبه شده و برابر $\beta_2 = 3$ است. بنابراین اگر $\beta_2 > 3$ توزیع داده ها از منحنی نرمال بلندتر و اگر $\beta_2 < 3$ از منحنی نرمال کوتاهتر و پخ تر است. فیشر برای ضریب کشیدگی مقدار $\gamma = \beta_2 - 3$ را پیشنهاد می کند پس اگر $\gamma > 0$ منحنی داده ها از منحنی نرمال بلندتر بطرف بالا کشیده است و اگر $\gamma < 0$ از منحنی نرمال کوتاهتر و بطرف کناره ها متمایل تر است.

مثال ۵.۲.۱ در یک نمونه گیری از درختان یک باغ 30° هکتاری، نمونه قطر ۲۵ درخت برحسب سانتیمتر بصورت زیر بدست آمده است:

۳۰	۲۵	۲۳	۲۰	۳۹
۱۸	۳۵	۳۴	۲۵	۳۲
۲۷	۲۳	۳۱	۳۷	۲۶
۲۵	۴۲	۳۷	۲۴	۳۰
۲۸	۲۲	۳۵	۳۱	۲۹

با رسم جدول توزیع فراوانی، شاخص های مرکزی و پراکندگی را یافته و گشتاورهای مرتبه اول، دوم و سوم را بدست آورید. با رسم نمودار چندبرفراوانی و میله ای، چولگی داده ها را مشخص نموده و ضرایب چولگی را بنویسید.

حل. در این نمونه ۲۵ تائی داریم $R = x_{max} - x_{min} + 1 = 42 - 18 + 1 = 25$ اگر تعداد دسته را ۵ انتخاب کنیم پس $T = \frac{25}{5} = 5$ و طول دسته را برابر ۵ می گیریم. مانند آنچه ذکر شد برای اینکه هر داده در یک دسته قرار گیرد، داده ها را از ۱۷/۵ آغاز می کنیم و جدول فراوانی بصورت زیر خواهد بود:

دسته ها	x_i	f_i	F_i	\bar{f}_i	$f_i x_i$	$x_i - \bar{x}$	$f_i(x_i - \bar{x})^2$	x_i^2	$f_i x_i^2$	$f_i x_i^3$	$f_i(x_i - \bar{x})^3$
۱۷/۵ - ۲۲/۵	۲۰	۳	۳	۰/۱۲	۶۰	-۹	۲۴۳	۴۰۰	۱۲۰۰	۲۴۰۰۰	-۲۱۸۷
۲۲/۵ - ۲۷/۵	۲۵	۸	۱۱	۰/۳۲	۲۰۰	-۴	۱۲۸	۶۲۵	۵۰۰۰	۱۲۵۰۰۰	-۵۱۲
۲۷/۵ - ۳۲/۵	۳۰	۷	۱۸	۰/۲۸	۲۱۰	۱	۷	۹۰۰	۶۳۰۰	۱۸۹۰۰۰	۷
۳۲/۵ - ۳۷/۵	۳۵	۵	۲۳	۰/۲۰	۱۷۵	۶	۱۸۰	۱۲۲۵	۶۱۲۵	۲۱۴۳۷۵	۱۰۸۰
۳۷/۵ - ۴۲/۵	۴۰	۲	۲۵	۰/۰۸	۸۰	۱۱	۲۴۲	۱۶۰۰	۳۲۰۰	۱۲۸۰۰۰	۲۶۶۲
		۲۵		۱	۷۲۵		۸۰۰	۴۷۵۰	۲۱۸۲۵	۶۸۰۳۷۵	۱۰۵۰

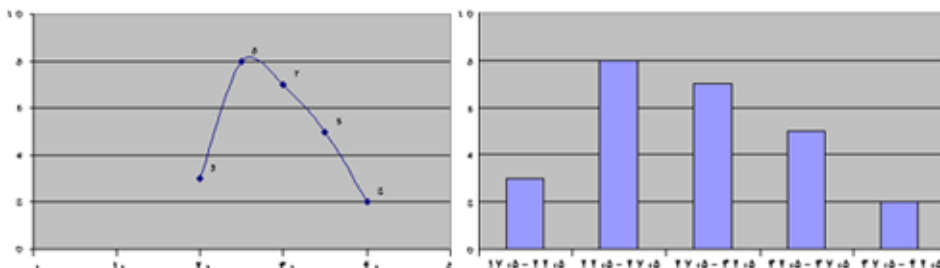
$$\bar{x} = 29, \quad S_x^2 = 33/34, \quad S_x = 5/77, \quad V_x = 0/2 = \%20$$

$$m = a_i + \frac{\frac{n}{f_i} - F_{i-1}}{f_i} c = 27/5 + \frac{12/5 - 11}{7} 5 = 28/57, \quad M = a_i + \frac{d_1}{d_1 + d_2} c = 22/5 + \frac{5}{5+1} 5 = 26/67$$

$$m_1 = 0, \quad m_2 = 32, \quad m_3 = 42, \quad m'_1 = 29, \quad m'_2 = 873, \quad m'_3 = 27215$$

$$b_1 = 0/4, \quad b_2 = 0/22, \quad g = 0/22, \quad \gamma = \frac{m_2}{S_x^2} - 3 = \frac{2200}{1024} - 3 = -0/75$$

نمودار داده‌ها بصورت ذیل بوده و از نمودار چندبرفراوانی دیده می‌شود که نمودار داده‌ها چوله به راست بوده و مقادیر b_1 و b_2 و g این را بخوبی تایید می‌کنند. مقدار ضریب کشیدگی فیشر $\gamma = -0.75$ مشخص می‌کند که نمودار داده‌ها پخ است.



۷.۲.۱ تغییر داده‌ها

اگر n داده x_1 و x_2 و \dots و x_n داده‌های یک نمونه مفروض باشند و n داده y_1 و y_2 و \dots و y_n چنان باشند که $y_i = ax_i + b$ در آنصورت با داشتن \bar{x} و S_x می‌توان \bar{y} و S_y را بصورت زیر بدست آورد

$$\bar{y} = a\bar{x} + b, \quad S_y^2 = a^2 \cdot S_x^2, \quad S_y = |a|S_x$$

برای مثال داده‌های ۴۳ و ۴۱ و ۴۲ و ۴۴ و ۴۰ را در نظر گرفته و آنها را متغیر x_i می‌نامیم. فرض کنید که $y_i = x_i - 40$ پس y_i —ها بترتیب عبارتند از ۳ و ۱ و ۲ و ۴ و ۰ و از آنجاکه $\bar{y} = 2$ و $\bar{y} = \bar{x} - 40$ نتیجه می‌گیریم که $\bar{x} = 42$ است. بعلاوه

$$S_y^2 = \frac{(3-2)^2 + (1-2)^2 + (2-2)^2 + (4-2)^2 + (0-2)^2}{4} = \frac{10}{4} = 2.5 \implies S_y = |1|S_x \implies S_x = 2.5$$

در برخی موارد ساده تر است برای بدست آوردن میانگین داده‌های دسته بندی شده از روش کدگذاری استفاده کنیم. اگر x_i نماینده دسته‌های موجود در جدول فراوانی بوده و c فاصله دو دسته و A داده دسته وسطی باشد. با تغییر داده‌ها بصورت بالا، فرض می‌گیریم که $y_i = \frac{x_i - A}{c}$ با اضافه کردن ستون y_i به جدول و یافتن \bar{y} از فرمول $\bar{x} = c\bar{y} + A$ مقدار \bar{x} را پیدا می‌کنیم.

۸.۲.۱ تصحیح شپارد برای واریانس

در دسته بندی داده‌ها بعلت داخل شدن خطا در دسته بندی داده‌ها، مقدار $\frac{c^2}{12}$ را از واریانس S_x^2 کم می‌کنیم تا جبران خطا شود این مقدار به تصحیح شپارد معروف است. بنابراین واریانس تصحیح شده \hat{S}_x^2 از رابطه زیر بدست می‌آید:

$$\hat{S}_x^2 = S_x^2 - \frac{c^2}{12}$$

۹.۲.۱ میانگین اصلاح شده

در نمونه‌هایی که یک یا چند داده پرت دارند چون این تعداد کم بهرحال روی میانگین تاثیر می‌گذارند در این موارد بهتر است بجای میانگین از میانه استفاده کنیم، زیرا میانه تحت تاثیر داده پرت قرار نمی‌گیرد. روش دیگر حذف داده‌های پرت از فهرست داده‌هاست. در این روش پس از مرتب کردن داده‌ها، به تعداد داده‌های پرت از ابتدا و انتهای فهرست حذف نموده و میانگین

مابقی داده‌ها را بعنوان میانگین کل داده‌ها بکار می‌گیریم. برای مثال در داده‌های ۱ و ۲ و ۳ و ۴ و ۵ و ۴۵ مقدار داده ۴۵ مقداری پرت بوده و میانگین ۱۰ معیار تمرکز مفیدی نیست. در این حالت می‌توان میانه $m = 3/5$ را بجای میانگین بکار برد و یا با استفاده از روش دوم میانگین ۲ و ۳ و ۴ و ۵ را که برابر $3/5$ است، برای نمونه انتخاب نمود. اگر تعداد داده‌های پرت نمونه مرتب شده برابر r باشد، میانگین اصلاح شده مرتبه r عبارتست از

$$\hat{x} = \frac{1}{n - 2r} \sum_{i=r+1}^{n-r} x_i$$

۱۰.۲.۱ چندک‌ها

دیدیم که میانه داده‌ها را به دو دسته تقسیم می‌کرد و در حالیکه داده‌ها دارای فراوانی بودند، میانه فراوانی را دو قسمت کرده و داده «وسط» را انتخاب می‌نمود. اگر در یک نمونه، فراوانی داده‌ها را به چهار دسته تقسیم کنیم، داده‌های چارکی Q_j بدست می‌آیند. بطور کلی ما تنها دارای سه چارک Q_1 و Q_2 و Q_3 هستیم، که چارک Q_2 همان میانه خواهد بود. در داده‌های فاقد دسته بندی - پس از مرتب کردن داده‌ها - برای محاسبه چارک Q_j عبارت $\frac{j}{4}(n+1)$ را محاسبه کرده و مقدار صحیح آن را با k و اعشار آنرا با ω نشان می‌دهیم. سپس فرمول $Q_j = x_k + \omega(x_{k+1} - x_k)$ مقدار چارک را بدست می‌دهد.

مثال ۶.۲.۱ چارک‌ها را در داده‌های ۱۴ و ۱۷ و ۲۴ و ۱۸ و ۲۱ بدست آورید.

با مرتب کردن داده‌ها داریم: ۲۴ و ۲۱ و ۱۸ و ۱۷ و ۱۴ از آنجا که $n = 5$

$$\text{چارک اول } \frac{1}{4}(5+1) = \frac{3}{4} = 0.75 \Rightarrow k=1, \omega=0.25 \Rightarrow Q_1 = x_1 + 0.25(x_2 - x_1) = 15/5$$

$$\text{چارک دوم } \frac{2}{4}(5+1) = \frac{6}{4} = 1.5 \Rightarrow k=2, \omega=0 \Rightarrow Q_2 = x_2 + 0(x_3 - x_2) = 18$$

$$\text{چارک سوم } \frac{3}{4}(5+1) = \frac{9}{4} = 2.25 \Rightarrow k=3, \omega=0.25 \Rightarrow Q_3 = x_3 + 0.25(x_4 - x_3) = 22/5$$

تمرین ۵.۲.۱ کلاسی. چارک‌ها را در داده‌های زیر محاسبه کنید.

۵, ۷, ۹, ۱۰, ۱۱, ۱۵

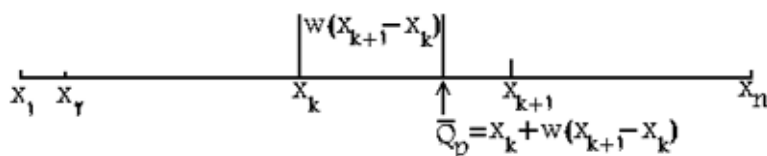
در داده‌های دسته بندی شده مانند فرمول میانه، دسته i که فراوانی جمع آن بیشتر یا مساوی $\frac{j}{4}n$ باشد را یافته و از فرمول

$$Q_j = a_i + \frac{j \frac{n}{4} - F_{i-1}}{f_i} c, \quad j = 1, 2, 3$$

برای محاسبه چارک‌ها استفاده می‌کنیم. مانند چارک، اگر دامنه به ۱۰ بخش تقسیم شود، هر قسمت را دهک گفته و با D_1, D_2, \dots, D_9 نشان می‌دهیم، بعلاوه دهک پنجم همان میانه است. اگر دامنه ۱۰۰ قسمت شود آنرا صدک گوئیم و با P_1, P_2, \dots, P_{99} نشان داده و صدک پنجاهم همان میانه خواهد بود.

عدد \bar{Q}_p را با $0 < p < 1$ حقیقی، چندک مرتبه p نامیم، هرگاه تقریباً $100p\%$ داده‌ها کوچکتر از آن باشند. مثلاً صدک دوم ($p = \frac{2}{100}$) را با $\bar{Q}_{\frac{2}{100}}$ نشان داده و مقداری است که دو درصد داده‌ها کمتر از آن قرار می‌گیرند. به همین ترتیب دهک هفتم ($p = \frac{7}{10}$) را با $\bar{Q}_{\frac{7}{10}}$ و چارک سوم ($p = \frac{3}{4}$) را با $\bar{Q}_{\frac{3}{4}}$ نشان می‌دهیم، روشن است که $m = \bar{Q}_{\frac{1}{4}} = Q_2 = D_5 = P_{50}$.

در داده های فاقد دسته بندی - پس از مرتب کردن داده ها- برای محاسبه چندک \bar{Q}_p عبارت $p(n+1)$ را محاسبه کرده، مقدار صحیح آن را با k و اعشار آنرا با w نشان می دهیم. سپس $\bar{Q}_p = x_k + w(x_{k+1} - x_k)$ مقدار چندک را بدست می دهد.



برای n داده دسته بندی شده، مقادیر چندکی چنین بدست می آیند که ابتدا دسته i -ام که فراوانی تجمعی آن بیشتر یا مساوی pn باشد را یافته و فرمول زیر را بکار می گیریم:

$$\bar{Q}_p = a_i + \frac{pn - F_{i-1}}{f_i} c$$

تمرین ۶.۲.۱ منزل.

(۱) سرشماری چیست؟ تفاوت آن را با نمونه گیری بیان کنید.

(۲) ویژگیهای یک نمونه برداری سالم و کارآمد چیست؟ مثال بزنید.

(۳) نماد سیگما \sum به معنای جمع چند عدد است، بدینصورت

$$\sum_{i=m}^n x_i = x_m + x_{m+1} + x_{m+2} + \dots + x_n$$

مثلاً

$$\sum_{i=1}^4 \frac{i}{i+5} = \frac{1}{1+5} + \frac{2}{2+5} + \frac{3}{3+5} + \frac{4}{4+5}$$

با استفاده از این تعریف مقادیر زیر را بدست آورید:

$$(a) \sum_{i=1}^8 \frac{2i+6}{5}, \quad (b) \sum_{i=1}^{10} 7i-3, \quad (c) \sum_{i=1}^4 \frac{i}{4}, \quad (d) \sum_{i=1}^4 (2i-6), \quad (e) \sum_{i=1}^{10} \frac{i+2}{3}, \quad (f) \sum_{i=1}^5 i^2+i+1$$

(۴) نماد سیگمای تعریف شده در تمرین قبل دارای خواصی بصورت زیر است که از آنها بهره می گیریم:

$$\sum_{i=1}^n C = nC, \quad \sum_{i=1}^n ax_i = a \sum_{i=1}^n x_i, \quad \sum_{i=1}^n (x_i + y_i) = \sum_{i=1}^n x_i + \sum_{i=1}^n y_i$$

اگر $\sum_{i=1}^5 x_i = 4$ و $\sum_{i=1}^5 x_i^2 = 13$ باشد مطلوبست مقادیر

$$(a) \sum_{i=1}^5 (2x_i + 3), \quad (b) \sum_{i=1}^5 x_i(x_i + 2), \quad (c) \sum_{i=1}^5 (x_i + 3)^2, \quad (d) \sum_{i=1}^5 \frac{x_i^2 - 1}{x_i - 1}$$

(۵) در یک نمونه گیری از ۵۰ خانواده ایرانی تعداد فرزندان این خانواده ها بصورت زیر بدست آمده است:

۲ و ۱ و ۳ و ۶ و ۱ و ۲ و ۵ و ۱ و ۰ و ۱ و ۲ و ۴ و ۱ و ۵ و ۲ و ۱ و ۰ و ۱ و ۳ و ۲ و ۱ و ۰ و ۱ و ۲ و ۱ و ۰ و ۲ و ۳ و ۳ و ۰ و ۰ و ۲ و ۱ و ۲ و ۳ و ۱ و ۲ و ۱ و ۰ و ۲ و ۴ و ۵ و ۱ و ۳ و ۲ و ۱ و ۲ و ۱ و ۰ و ۲ و ۳ و ۴

با رسم جدول توزیع فراوانی داده ها، شاخص های مرکزی و پراکندگی را یافته و گشتاورهای مرتبه اول، دوم و سوم و ضریب کشیدگی را بدست آورید. با رسم نمودار میله ای و چند بر فراوانی، چولگی داده ها را مشخص نموده و ضرایب چولگی را بنویسید.

(۶) بصورت تجربی مشخص شده که بین میانگین، میانه و مد رابطه تقریبی زیر برقرار است $\bar{x} - M \simeq 3(\bar{x} - m)$ صحت این قاعده را برای مثال های ۱.۲.۱ و ۴.۲.۱ و تمرین ۵ بررسی نمائید. (دقت کنید که این رابطه برای بدست آوردن یک پارامتر از دو پارامتر دیگر نبوده و تنها روشی برای کنترل اطلاعات و صحت ارقام بدست آمده محسوب می شود).

(۷) یک کارخانه کنسروسازی، قوطی هائی با وزن ۲۰۰ گرم عرضه می کند. در یک نمونه گیری از ۴۰ قوطی عرضه شده به بازار وزن آنها برحسب گرم بصورت زیر اندازه گیری شده است:

۱۹۸ و ۱۸۱ و ۱۹۲ و ۱۹۹ و ۲۰۳ و ۲۰۵ و ۱۹۴ و ۱۸۰ و ۲۰۵ و ۲۰۷ و ۱۹۹ و ۲۰۰ و ۲۰۷ و ۲۱۸ و ۱۹۰ و ۱۸۸ و ۱۹۱ و ۲۱۰ و ۱۹۹ و ۲۰۱ و ۲۰۷ و ۲۱۴ و ۲۰۶ و ۲۰۴ و ۱۹۹ و ۱۹۶ و ۱۹۴ و ۲۰۱ و ۱۸۴ و ۲۰۲ و ۱۹۸ و ۲۰۰ و ۲۰۲ و ۱۸۲ و ۱۹۵ و ۱۹۴ و ۱۸۲ و ۲۰۴ و ۱۹۹ و ۲۰۳

با رسم جدول توزیع فراوانی داده ها، شاخص های مرکزی و پراکندگی را یافته و گشتاورهای مرتبه اول، دوم و سوم را بدست آورید. با رسم هیستوگرام و نمودار داده ها، چولگی داده ها را مشخص نموده و ضرایب چولگی را بنویسید. نمودار چندبرفراوانی صعودی و نزولی را رسم و میانه را با استفاده از آن بیابید.

(۸) از نظر هندسی و نمودارها، چندک \bar{Q}_p چه معنی دارد؟ چارکها، دهک سوم D_3 و دهک هفتم D_7 تمرین ۷ را حساب کنید.

(۹) برای داده های ۶ و ۲ و ۱۲ و ۹ و ۴ و ۸ و ۵ مقادیر چندکی $\bar{Q}_{\frac{1}{3}}$ ، $\bar{Q}_{\frac{2}{3}}$ ، $\bar{Q}_{\frac{1}{4}}$ ، $\bar{Q}_{\frac{3}{4}}$ ، P_{10} ، D_2 ، Q_1 را حساب کنید.

(۱۰) اگر میانگین داده های x_1, x_2, \dots, x_n برابر با \bar{x} باشد ثابت کنید $\sum_1^n (x_i - \bar{x}) = 0$

(۱۱) برای داده های مثبت x_1, x_2, \dots, x_n میانگین هندسی G و میانگین توافقی H را چنین تعریف می کنیم:

$$G = \sqrt[n]{x_1 x_2 \dots x_n} \quad , \quad H = \frac{n}{\sum \frac{1}{x_i}}$$

رابطه ای که بین \bar{x} و G و H برقرار است بصورت $H \leq G \leq \bar{x}$ این رابطه را برای داده های مثال ۱.۲.۱ بررسی نمائید.

(۱۲) گشتاورهای مرتبه اول تا چهارم را برای داده های ۱ و ۳ و ۵ و ۷ و ۹ حساب کنید.

(۱۳) گاهی دامنه تغییرات R را بعنوان شاخص پراکندگی محسوب می کنند، بنظر شما آیا می توان از این شاخص نتیجه ای آماری گرفت؟ میانبرد که با $R' = \frac{1}{4}(x_{max} + x_{min})$ تعریف می شود چطور؟

(۱۴) میانه داده های ۵ و ۶ و ۷ و ۸ با فراوانی های بترتیب برابر ۴ و ۵ و ۳ و ۵ چند است. مد آنها چند است؟

(۱۵) ثابت کنید که در فرمول واریانس اگر بجای \bar{x} هر عدد دیگری قرار دهیم، مقدار بدست آمده از واریانس بیشتر خواهد شد. عبارتی مقدار واریانس کوچکترین عددی است که با مقدار میانگین از فرمول حاصل می شود.

(۱۶) یکی از مقیاس های سنجش چولگی در نمودارها با استفاده از چارکها و بررسی علامت $Q_1 + Q_3 - 2Q_2$ است. اگر این مقدار مثبت باشد منحنی چوله به راست و اگر منفی چوله به چپ است و در حالت نرمال صفر خواهد بود. این مقدار را برای مثالهای ۴.۲.۱ و ۵.۲.۱ بررسی و چولگی آنها را بسنجید.

(۱۷) برای تعیین میانگین اعداد وزن دار مانند نمرات دارای ضریب، از فرمول میانگین وزن دار بصورت زیر استفاده می کنیم:

$$\bar{x}_w = \text{میانگین وزنی} = \frac{w_1x_1 + w_2x_2 + \dots + w_nx_n}{w_1 + w_2 + \dots + w_n}$$

ضرایب فردی در کنکور در درس ریاضی، فیزیک، شیمی و ادبیات بترتیب ۴۰، ۲۸، ۳۲ و ۵۰ درصد است اگر این دروس بترتیب دارای ضرایب ۴، ۳، ۳ و ۴ باشند رتبه فرد را محاسبه کنید.

(۱۸) ثابت کنید اگر به تمام داده ها مقدار ثابت a اضافه شود، گشتاورهای مرکزی، معیارهای چولگی و برجستگی داده ها تغییری نخواهند کرد.

(۱۹) کالفورد^{۱۳} آمارشناس معروف نشان داد که اگر یک توزیع متقارن یا خفیفاً چوله باشد و اندازه نمونه نیز بزرگ باشد انحراف معیار یک ششم دامنه تغییرات است یعنی $R = 6S_x$. درستی این رابطه را برای مثال های ۱.۲.۱، ۴.۲.۱ و ۵.۲.۱ بررسی نمایید. (این رابطه صرفاً برای بررسی صحت عملیات است).

(۲۰) برای استاندارد کردن داده های x_i از $z_i = \frac{x_i - \bar{x}}{S_x}$ استفاده کرده و z_i را متغیر استاندارد شده گوئیم. ثابت کنید $\bar{z} = 0$ و $S_z = 1$.

(۲۱) انواع نمودار در اکسل. در برنامه اکسل^{۱۴} کامپیوتر، انواع متنوع و مختلفی از نمودارها معرفی شده اند. برای آشنائی بیشتر با این نمودارها، شما بایستی با در نظر گرفتن داده های مثال های ۱.۲.۱ و ۵.۲.۱ نمایش مختلف داده ها را بصورت نمودار رسم نمایید: ستونی، میله ای، خطی، کلوچه ای، XY -پراکندگی، مسطح، دوناتی، راداری، سطحی، حبابی، انباشته، استوانه ای، مخروطی و هرمی.