



مدرس‌ان شریف

فصل اول

«آمار توصیفی»

درسنامه (۱): مفاهیم اولیه



معرفی علم آمار: علم آمار به سه شاخه‌ی آمار توصیفی، آمار استنباطی (پارامتری) و آمار ناپارامتری تقسیم می‌شود.

آمار توصیفی (Descriptive Statistics)

این شاخه از علم آمار با کمک ابزارهایی مانند جداول فراوانی، نمودارهای آماری و محاسبه‌ی شاخص‌های آماری به توصیف داده‌ها می‌پردازد. در آمار توصیفی همه‌ی داده‌ها به طریق سرشماری جمع‌آوری می‌شوند و هیچ‌گونه استنباطی بر روی آن‌ها انجام نمی‌شود و فقط به توصیف آن‌ها پرداخته می‌شود.

آمار استنباطی (پارامتری) (Inferential Statistics)

در آمار استنباطی به کمک نمونه‌گیری به تحلیل و استنتاج در مورد پارامترهای کلی جامعه‌ی آماری پرداخته می‌شود. توجه کنید که در آمار استنباطی توزیع یا نحوه‌ی پخش داده‌های آماری مشخص است.

آمار ناپارامتری (آزاد از توزیع) (Non – Parametric Statistics)

آمار ناپارامتری دقیقاً مانند آمار استنباطی عمل می‌کند، با این تفاوت که در آمار ناپارامتری اغلب صفات کیفی هستند و توزیع یا نحوه‌ی پخش داده‌های آماری مشخص نیست.

کلمه مثال: مرکز آمار ایران هر چند سال یک بار سرشماری انجام می‌دهد. اطلاعات جمع‌آوری شده با استفاده از کدام شاخه علم آمار مورد بررسی قرار می‌گیرند؟

- (۱) آمار توصیفی (۲) آمار استنباطی (۳) آمار ناپارامتری (۴) هیچکدام

پاسخ: گزینه «۱» در آمار استنباطی و آمار ناپارامتری نمونه‌گیری انجام می‌گیرد در حالی که در سرشماری، مرکز آمار ایران همه‌ی واحدهای جامعه را مورد بررسی قرار می‌دهد.

جامعه‌ی آماری یا جمعیت (Population)

مجموعه‌ی تمام افراد یا اشیایی که مطالعات آماری در مورد یک یا چند صفت آن‌ها در یک مکان و زمان معین انجام می‌گیرد، جامعه‌ی آماری یا جمعیت گفته می‌شود. تعداد اعضای جامعه، حجم جامعه نامیده می‌شود که آن را با N نشان می‌دهند.

نمونه (Sample): در بررسی‌های آماری به دلیل هزینه‌ی زیاد، کمبود وقت و در بعضی مواقع غیرممکن بودن انجام کار، زیرمجموعه‌ای از جامعه با قاعده و ضابطه‌ی خاصی انتخاب می‌شود که به آن نمونه‌ی آماری گفته می‌شود. تعداد اعضای نمونه، حجم نمونه نامیده می‌شود و آن را با n نشان می‌دهند.

صفت مشخصه (Attribute): صفتی است که بین همه‌ی عناصر جامعه‌ی آماری مشترک بوده و جداکننده‌ی جامعه‌ی آماری از سایر جوامع است. به‌طور کلی صفات، خود به دو دسته‌ی کمی و کیفی تقسیم می‌شوند.

صفات کمی: صفات کمی صفاتی هستند که می‌توانند به صورت عددی بیان شوند. تعداد دانشجویان یک دانشگاه، درآمد یک خانواده و ... دارای صفت کمی هستند.

صفات کیفی: صفات کیفی صفاتی هستند که نمی‌توانند به صورت عددی بیان شوند. گروه خون، رنگ چشم انسان‌ها و ... دارای صفت کیفی هستند.



کدام یک از گزینه‌های زیر تعریف صفت مشخصه است؟

(مدیریت و حسابداری - سراسری ۹۳)

- (۲) صفت مشترک بین کلیه افراد جامعه است.
 (۴) عنصر مشترک جوامع آماری مختلف است.

- (۱) از فردی به فرد دیگر تغییر می‌کند.
 (۳) متمایزکننده عناصر جامعه از یکدیگر است.

پاسخ: گزینه «۲» صفت مشترک بین اعضاء جامعه است.

انواع مقیاس‌های اندازه‌گیری صفات (مقیاس‌های استیونز)

الف- کیفی (گروهی)

این مقیاس، خود به دو مقیاس اسمی و رتبه‌ای تقسیم می‌شود. ماهیت عددی ندارند، قابل اندازه‌گیری نیستند و واحد ندارند. وضعیت تأهل، رنگ چشم، جنسیت و مهارت، مثال‌هایی از مقیاس کیفی هستند.

۱- مقیاس اسمی (Nominal Scale): زمانی که صفتی دارای حالت‌های مختلف باشد، مفهوم بهتر (بدتر) و بزرگتر (کوچکتر) بودن را نداشته باشد، همچنین صفت، ماهیت عددی نداشته باشد و بتوان آن را در طبقات جداگانه‌ای قرار داد، مقیاس این صفت اسمی است. به‌طور مثال جنسیت افراد را در نظر بگیرد. دو حالت مرد و زن دارد و مرد و زن بودن بر یکدیگر ارجحیت ندارند. همچنین ماهیت آن عددی نیست.

۲- مقیاس رتبه‌ای (Rank Scale): زمانی که صفتی مفهوم بهتر (بدتر) و بزرگتر (کوچکتر) بودن را داشته باشد و نتوان آن را اندازه‌گیری کرد، از این مقیاس استفاده می‌شود. مثلاً طبقه‌بندی مردم یک کشور به سه طبقه‌ی پر درآمد، با درآمد متوسط و کم درآمد تقسیم شود.

ب - مقیاس کمی (عددی)

وسیله یا روشی خاص برای اندازه‌گیری یک صفت است که ماهیت و حاصل آن بصورت عددی است، مانند تعداد فرزندان، طول قد افراد، درجه‌ی حرارت این مقیاس خود به سه دسته‌ی شمارشی، فاصله‌ای و نسبی (نسبتی) تقسیم می‌شود.

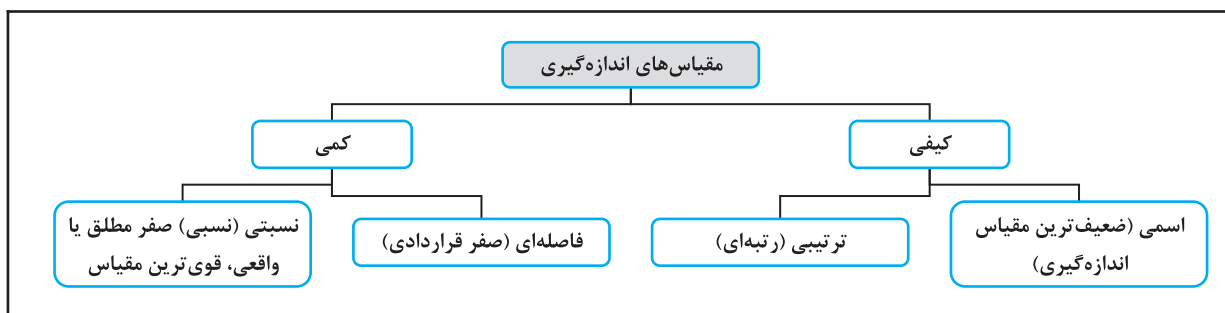
۱- مقیاس فاصله‌ای (Interval Scale): زمانی که صفتی دارای ماهیت عددی باشد، اختلاف بین حالتها را بیان کند و مفهوم بهتر (بدتر) و بزرگتر (کوچکتر) بودن را داشته باشد مقیاس آن فاصله‌ای است. اعداد روی دماسنج و یا نمره‌های مربوط به هوش افراد دارای مقیاس فاصله‌ای هستند.

۲- مقیاس نسبی (نسبتی) (Ratio Scale): این مقیاس دارای صفر مطلق است و چهار عمل اصلی روی این مقیاس انجام می‌گیرد. طول قد افراد، درآمد و هزینه‌ها همه مثال‌هایی از مقیاس نسبی هستند. در مقیاس نسبی نسبت حفظ می‌شود، مثلاً اگر برای دو جسم و ویژگی وزن با X_1 و X_2 اندازه‌گیری شود باید $\frac{X_1}{X_2}$ ثابت بماند و به واحد اندازه‌گیری بستگی نداشته باشد.

کدام مثال ۳: در اندازه‌گیری داده‌های سطح تحصیلات، دمای اجسام و طول قد دانشجویان از چه مقیاس‌هایی استفاده شده است؟

- (۱) ترتیبی - فاصله‌ای - نسبی (۲) فاصله‌ای - فاصله‌ای - نسبی (۳) ترتیبی - ترتیبی - اسمی (۴) اسمی - فاصله‌ای - فاصله‌ای

پاسخ: گزینه «۱» در داده‌های سطح تحصیلات ترتیب رعایت می‌شود - دمای اجسام دارای مقیاس فاصله‌ای و طول قد دانشجویان دارای مقیاس نسبی است.



کدام یک از مقیاس‌ها دارای صفر قرار دادی است؟

(مدیریت و حسابداری - سراسری ۹۴ - طراحی و برنامه‌ریزی شهری - سراسری ۸۹)

- (۱) نسبی (۲) فاصله‌ای (۳) اسمی (۴) رتبه‌ای

پاسخ: گزینه «۲» صفر قراردادی به صورت قرارداد صفر در نظر گرفته شده است و مفهوم آن عدد صفر واقعی نیست.

داده‌های آماری (Statistics data)

در بررسی‌های آماری باید صفت مورد بررسی به صورت اعداد و ارقام نمایش داده شود. اگر صفت مورد بررسی کمی باشد، این عمل به سادگی امکان‌پذیر است ولی اگر داده‌ها کیفی باشند، باید طبق ضوابط خاصی با عدد و رقم نمایش داده شوند. به‌طور کلی داده‌های آماری به دو دسته‌ی گسسته (طبقه‌بندی‌نشده) و پیوسته (طبقه‌بندی شده) تقسیم می‌شوند.

داده‌های گسسته (Discrete data): داده‌هایی هستند که بین هر دو مقدار متوالی a مورد نظر از آن‌ها هیچ عددی نمی‌تواند قرار گیرد. تعداد دانشجویان یک دانشکده، تعداد فرزندان یک خانواده و ... به صورت گسسته بیان می‌شوند.

داده‌های پیوسته (Continous data): داده‌هایی هستند که بین هر دو مقدار مورد نظر از آن‌ها بی‌شمار عدد می‌تواند قرار گیرد. طول قد و یا وزن افراد، طول عمر انسان‌ها و ... همه مثالهایی از داده‌های پیوسته هستند. توجه کنید که ماهیت داده‌های پیوسته به صورت اعشاری است، اگرچه در پارهای از مواقع ممکن است آن‌ها را به دلیل رند کردن یا به طور تصادفی به صورت گسسته نمایش دهند، اما این نمایش مهم نیست و ماهیت اعداد برای ما مهم است. به طور مثال اگر وزن فردی ۷۴ کیلوگرم نمایش داده شود این عدد واقعی نیست چرا که وزن، ماهیتی پیوسته دارد و باید اعشاری باشد لذا ما به ماهیت عدد توجه می‌کنیم نه به نوع نمایش آن.

کج مثال ۵: کدام یک از موارد زیر داده‌های گسسته می‌باشد؟

(۲) وزن افراد یک شرکت

(۱) قد افراد یک کلاس

(۴) تعداد افرادی که به یک فروشگاه وارد می‌شوند.

(۳) زمان رسیدن اتوبوس به یک ایستگاه

پاسخ: گزینه «۴» بنا به تعریف داده‌های گسسته داریم: داده‌هایی که بین هر دو مقدار متوالی a مورد نظر از آن‌ها هیچ عددی نمی‌تواند قرار گیرد بنابراین گزینه‌ی (۴) جواب است.

مراحل یک پژوهش علمی در آمار

در تحقیقات علمی به روش آماری چند مرحله‌ی مهم باید به‌طور کامل انجام شود:

(الف) مشخص کردن هدف: در این مرحله با کمک روش‌های تحقیق، تلاش برای افزایش اطلاعات از موضوع و بررسی اطلاعات پایدارتر برای رسیدن به هدف انجام می‌گیرد.

(ب) جمع‌آوری داده‌ها: در هر پژوهش آماری، تهیه‌ی داده‌های واقعی با توجه به هدفی که از پژوهش داریم، اهمیت اساسی دارد. فرآیند جمع‌آوری داده‌ها ممکن است به روش‌های مختلفی انجام گیرد.

(ج) تجزیه و تحلیل داده‌ها: داده‌هایی که به روش‌های مناسب گردآوری شده، منبع اساسی برای کسب اطلاعات جدید درباره‌ی پدیده مورد مطالعه هستند در این مرحله داده‌ها را بررسی کرده و معلومات جدیدی به‌دست می‌آوریم که تعیین‌کننده‌ی نقاط قوت و ضعف آن‌ها است.

(د) بیان یافته‌ها: در این مرحله، تحلیل پایانی بر روی داده‌های آماری که در مرحله‌ی اول تحقیق مشخص شده‌اند، انجام می‌شود و همچنین در این مرحله پاسخ به سؤالات اولیه امکان‌پذیر است.

کج مثال ۶: اولین مرحله در تحقیق علمی کدام است؟

(مدیریت بازرگانی - آزاد ۸۸)

(۴) جمع‌آوری داده‌ها

(۳) هدف‌های پژوهش

(۲) تحلیل یافته‌ها

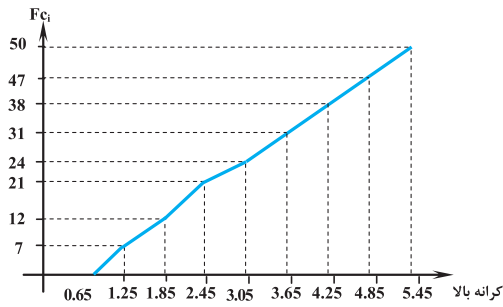
(۱) فرضیه‌سازی

پاسخ: گزینه «۳» با توجه به ترتیب مراحل گفته شده در بالا هدف‌های پژوهش اولین مرحله است.

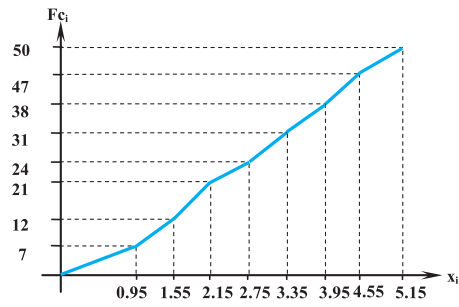
مطالعه توصیفی داده‌ها (آمار توصیفی)

مفهوم و کاربرد نمادها \sum

مجموعه‌ی داده‌های آماری را معمولاً با نماد X_1, X_2, \dots, X_N نشان می‌دهند. نماد \sum برای جمع کردن داده‌های آماری بکار می‌رود. علامت \sum برای عدم استفاده مکرر از علامت "+" بکار می‌رود. در نماد \sum یک کران پایین، یک کران بالا و یک متغیر (چند متغیر) در جلوی سیگما نوشته می‌شود. جمله‌ای که بعد از نماد \sum نوشته می‌شود، نشان دهنده‌ی مقادیری است که باید با هم جمع شوند و کران‌های پایین و بالای \sum نشان‌دهنده‌ی ابتدا و انتهای جمع هستند.



«نمودار فراوانی تجمعی»



«نمودار پلی‌گان فراوانی تجمعی»

۲. نمودارهای وصفی

این نوع از نمودارها برای داده‌های کیفی بکار برده می‌شوند. مهم‌ترین این نمودارها عبارتند از:

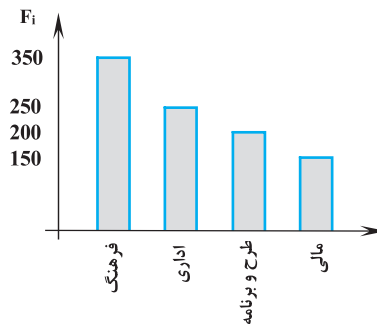
الف) نمودار ستونی (میله‌ای) (Bar-chart)

در این نوع نمودار، معمولاً محور X نشان‌دهنده کیفیت مشاهدات و محور عمودی آن نشان‌دهنده فراوانی مطلق یا نسبی هر گروه است.

کلمه کلیدی مثال ۱۳۲: تعداد کارکنان یک سازمان بر حسب معاونت به شرح زیر است، نمودار ستونی آن را رسم کنید.

فرهنگ	طرح و برنامه	مالی	اداری	معاونت
۳۵۰	۲۰۰	۱۵۰	۲۵۰	تعداد کارکنان

پاسخ: می‌توانیم به جای مستطیل‌ها از خطوط استفاده کنیم ولی در اینجا مستطیل رسم می‌کنیم که نشان می‌دهد فراوانی‌ها بر روی موضوعات پخش شده است.



ب) نمودار دایره‌ای (کلوچهای) (Pie chart)

برای رسم این نمودار، ابتدا دایره‌ای رسم کرده و سپس قطاع‌هایی بر حسب درصد یا درجه از دایره جدا می‌کنیم. سهم این قطاع‌ها از دایره به صورت زیر محاسبه می‌گردد:

$$\text{سهم بر حسب درصد} = \frac{F_i}{n} \times 100 \quad , \quad \text{سهم بر حسب درجه} = \frac{F_i}{n} \times 360$$

(علوم اقتصادی - سراسری ۹۷)

کلمه کلیدی مثال ۱۳۳: کدام نمودار، برای نشان دادن وضعیت توزیع فراوانی داده‌های کیفی، مناسب است؟

- (۱) هیستوگرام (۲) پراکنش (۳) میله‌ای (۴) چند بر فراوانی انباشته

پاسخ: گزینه «۳» به بررسی گزینه‌ها می‌پردازیم. نمودار هیستوگرام برای داده‌های فاصله‌ای و کمی می‌باشد، نمودار پراکنش برای داده‌های کمی و دومتغیره می‌باشد و چند بر فراوانی انباشته از طریق نمودار هیستوگرام به وجود می‌آید که اگر بخواهیم دو جامعه را از نظر داده‌های کمی با هم مقایسه نماییم نمودار مناسبی می‌باشد ولی برای داده‌های کیفی، نمودارهای دایره‌ای، میله‌ای و غیره مناسب می‌باشد.

کلمه کلیدی مثال ۱۳۴: در یک نمودار دایره‌ای، قطاعی به اندازه 36° جدا شده است. اگر مجموع فراوانی‌های مطلق برابر با ۱۰۰ باشد، فراوانی مطلق طبقه‌ی

مورد نظر کدام است؟

۴۰ (۴)

۳۰ (۳)

۲۰ (۲)

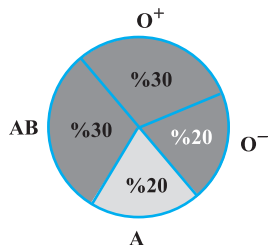
۱۰ (۱)

پاسخ: گزینه «۱» با توجه به رابطه‌ی $\text{سهم بر حسب درصد} = \frac{F_i}{n} \times 360$ داریم: $36 = \frac{F_i}{100} \times 360 \Rightarrow 3600 = 360 F_i \Rightarrow F_i = 10$

مثال ۱۳۵: تعدادی بیمار برحسب گروه خونی به صورت زیر مرتب شده‌اند. نمودار دایره‌ای را رسم کنید.

گروه‌های خونی	A	AB	O ⁺	O ⁻
F _i	۱۰۰۰	۱۵۰۰	۱۵۰۰	۱۰۰۰

پاسخ: با توجه به رابطه‌های گفته شده در بالا خواهیم داشت:



$$A \text{ سهم گروه خونی } = \frac{F_i}{N} \times 100 = \frac{1000}{5000} \times 100 = 20\%$$

$$AB \text{ سهم گروه خونی } = \frac{F_i}{N} \times 100 = \frac{1500}{5000} \times 100 = 30\%$$

$$O^+ \text{ سهم گروه خونی } = \frac{F_i}{N} \times 100 = \frac{1500}{5000} \times 100 = 30\%$$

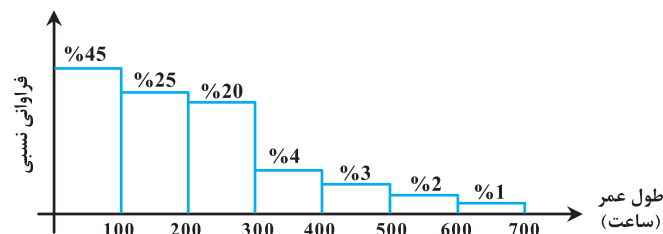
$$O^- \text{ سهم گروه خونی } = \frac{F_i}{N} \times 100 = \frac{1000}{5000} \times 100 = 20\%$$

ج) نمودار پارتو (Pareto chart)

در نمودار پارتو، فراوانی نسبی هر موضوع روی محور عمودی و نوع آن روی محور افقی آورده می‌شود. نمودار پارتو همیشه به ترتیب نزولی فراوانی‌ها رسم می‌شود. پر وقوع‌ترین موضوع در سمت چپ قرار می‌گیرد. توجه کنید که این نمودار سه محور دارد که محور سوم آن نشان دهنده‌ی فراوانی‌های نسبی تجمعی است. این نمودار برای تحلیل موجودی انبار کالا، نواقص سیستمها، توزیع درآمد و توزیع کارمندان سازمانها، کاربرد فراوانی دارد.

مثال ۱۳۶: توزیع طول عمر ۱۰۰۰ تولید نوعی قطعه ترانزیستوری که در یک کارخانه تولید شده است به شکل زیر است. اگر بدانیم ۲۰٪ از قطعه‌ها

از درجه کیفیت کمتری برخوردار می‌باشند و دارای طول عمر پایین هستند پس این تعداد کمتر از چه مدت عمر می‌کنند؟ (محیط زیست - سراسری ۸۷)



۲۰ (۱)

۴۴/۴ (۲)

۵۰/۵ (۳)

۷۰ (۴)

پاسخ: گزینه «۲» با کمی دقت از روی نمودار، جدول فراوانی را به صورت زیر می‌کشیم:

طول عمر (ساعت)	۱۰۰-۲۰۰	۲۰۰-۳۰۰	۳۰۰-۴۰۰	۴۰۰-۵۰۰	۵۰۰-۶۰۰	۶۰۰-۷۰۰
(فراوانی نسبی) f _i	۰/۴۵	۰/۲۵	۰/۲۰	۰/۰۴	۰/۰۳	۰/۰۲

در سؤال از ما صدک ۲۰ام طول عمرها خواسته شده است یعنی نقطه‌ای که ۲۰٪ از قطعات کمتر یا مساوی آن نقطه عمر می‌کنند برای بدست آوردن صدک ۲۰ام به صورت زیر عمل می‌کنیم. توجه کنید اولین طبقه‌ای که فراوانی تجمعی نسبی آن بزرگ‌تر یا مساوی ۰/۲ است طبقه‌ی اول است بنابراین صدک بیستم در طبقه‌ی اول قرار دارد. از رابطه چندک‌ها در داده‌های پیوسته خواهیم داشت:

$$P_{20} = L + \frac{0/2 - fC_{i-1}}{f_i} \times I = 0 + \frac{0/2 - 0}{0/45} \times 100 = \frac{20}{0/45} = 44/4$$

تحلیل اکتشافی داده‌ها (Exploratory data Analysis)

شناسایی انفرادی داده‌ها را تحلیل اکتشافی داده‌ها می‌گویند. نمودارهایی در مراحل اولیه‌ی داده‌ها به‌عنوان نمودارهای تحلیل اکتشافی رسم می‌شوند در اینجا به دو نمونه از آن‌ها اشاره می‌کنیم.

۱. نمودار شاخه و برگ (ساقه و برگ)، (Stem and leaf display)

فرض کنید X_1, X_2, \dots, X_n مشاهدات باشند و هر مشاهده دست کم دورقمی باشد. برای تهیه‌ی نمودار شاخه و برگ، ارقام را به دو بخش تقسیم می‌کنیم: شاخه شامل یک یا چند رقم اولیه و برگ آن‌ها شامل رقم‌های باقیمانده است.

۰, ۱۷, ۱۳, ۱۳, ۱۸, ۱۰, ۱۹, ۵, ۷, ۴, ۱, ۲۱, ۲۵, ۲۴, ۲۹, ۳۵

شاخه	برگ
۰	۵, ۴, ۷, ۰, ۱
۱	۷, ۳, ۳, ۸, ۹, ۰
۲	۱, ۵, ۴, ۹
۳	۵

مثال ۱۳۷: برای داده‌های روبرو، نمودار شاخه و برگ را رسم کنید.

پاسخ: رقم‌های یکان برگ و رقم‌های دهگان به عنوان شاخه هستند. بنابراین شکل به صورت روبرو خواهد شد:

مثال ۱۳۸: داده‌های آماری به صورت نمودار شاخه و برگ زیر داده شده است. با حذف ۲۵ درصدی پایین و ۲۵ درصدی بالا، تفاضل میانه از میانگین

(مدیریت و حسابداری - سراسری ۹۳)

پیراسته کدام است؟

شاخه	برگ					
۷	۱	۱	۲	۴	۷	۹
۸	۰	۰	۳	۵	۶	۷
۹	۲	۴	۵	۵	۷	۸

۰/۱ (۴)

۰/۳ (۳)

۰/۲ (۲)

۰/۴ (۱)

پاسخ: گزینه «۳» داده‌ها را از نمودار بیرون می‌کشیم، اکنون می‌توانیم برای راحتی از تک‌تک داده‌ها ۸۰ واحد کم کنیم تا مقیاس کوچک شود.

۷۱, ۷۱, ۷۲, ۷۴, ۷۷, ۷۹, ۸۰, ۸۰, ۸۳, ۸۵, ۸۶, ۸۷, ۹۲, ۹۴, ۹۵, ۹۵, ۹۷, ۹۸

$x_i - 80: -9, -9, -8, -6, -3, -1, 0, 3, 5, 6, 7, 12, 14, 15, 15, 17, 18$

تعداد داده‌ها ۱۸ تا می‌باشد، بنابراین $n = 18$ زوج است. پس میانه $Me = \frac{x_9 + x_{10}}{2}$ است که برابر است با: $Me = \frac{3 + 5}{2} = 4$. اکنون $\frac{1}{4}$ داده‌ها از بالا و

$$[np] = [18 \times \frac{1}{4}] = [4/5] = 4$$

پایین حذف می‌کنیم (میانگین‌گیری پیراسته)

$x_i: -3, -1, 0, 0, 3, 5, 6, 7, 12, 14$

پس ۴ تا از بالا و پایین حذف کرده میانگین‌گیری می‌کنیم:

$$\mu = \frac{\sum_{i=1}^N x_i}{N} = \frac{-3 + -1 + 0 + 0 + 3 + 5 + 6 + 7 + 12 + 14}{10} = \frac{43}{10} = 4.3 \Rightarrow \mu - Me = 4.3 - 4 = 0.3$$

۲. نمودار جعبه‌ای (Box Plot)

مراحل رسم این نمودار به شرح زیر است:

- داده‌های خام را به صورت غیر نزولی مرتب می‌کنیم، سپس یک خط افقی را چنان مدرج می‌کنیم که بتوان همه‌ی داده‌ها را روی آن نشان داد.
- میانه و چارکهای اول و سوم داده‌ها را محاسبه می‌کنیم.
- بالای خط مدرج شده مستطیلی رسم می‌کنیم که طول آن برابر با $Q_3 - Q_1$ بوده و از نقطه Q_1 شروع و به Q_3 ختم شود، عرض مستطیل به اندازه‌ی معقول در نظر گرفته می‌شود. این مستطیل را جعبه یا باکس می‌نامیم.
- اندازه میانه را به صورت خطی به موازات عرض مستطیل رسم نموده و مستطیل را به وسیله یک خط منقطع به موازات خط مدرج شده به دو قسمت تقسیم می‌کنیم (شکل روبرو را مشاهده کنید).



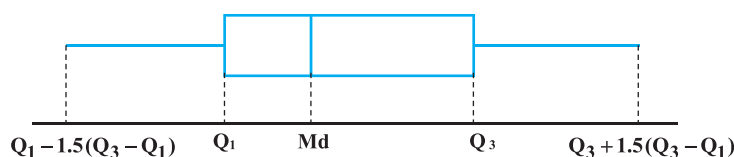
۵- مرزهای داخلی و خارجی داده‌ها به صورت زیر تعیین می‌شوند:

$$\text{مرز داخلی پایین} = Q_1 - 1/5(Q_3 - Q_1)$$

$$\text{مرز داخلی بالا} = Q_3 + 1/5(Q_3 - Q_1)$$

$$\text{مرز خارجی پایین} = Q_1 - 3/5(Q_3 - Q_1)$$

$$\text{مرز خارجی بالا} = Q_3 + 3/5(Q_3 - Q_1)$$



۶- هر داده‌ای که خارج از مرزهای داخلی قرار گرفته باشد را یک داده‌ی پرت می‌نامند و چنانچه بین مرزهای داخلی و خارجی قرار گیرد، آن را داده‌ی پرت ضعیف نامیده و با علامت ۰ نشان می‌دهیم و چنانچه خارج از مرزهای خارجی قرار گیرد آن را داده‌ی پرت قوی نامیده و با علامت • نشان می‌دهیم.

با استفاده از نمودار جعبه‌ای می‌توان اطلاعات زیر را در مورد داده‌ها کسب نمود:

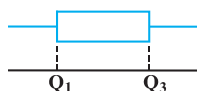
- ۱- اگر میانه نزدیک وسط مستطیل (جعبه) باشد، توزیع داده‌ها تقریباً متقارن است.
- ۲- اگر میانه در طرف چپ وسط مستطیل باشد، توزیع چوله به راست و اگر میانه در طرف راست وسط مستطیل قرار گیرد، توزیع چوله به چپ است.
- ۳- در مقایسه‌ی نمودار جعبه‌ای دو مجموعه از داده‌ها می‌توان پراکندگی آن‌ها را با توجه به طول مستطیل‌های نمودار با یکدیگر مقایسه نمود. مستطیلی که طول بزرگتری دارد، داده‌های آن دارای پراکندگی بیشتری است.

کلمه مثال ۱۳۹: داده‌های آماری به صورت شاخه و برگ زیر، نشان داده شده است. در نمودار جعبه‌ای این داده‌ها، واریانس داده‌های داخل جعبه، تقریباً کدام است؟

(مدیریت و حسابداری - سراسری ۹۴)

شاخه	برگ							
۴	۱	۱	۲	۴	۶	۸		۱۵/۱ (۱)
۵	۰	۲	۲	۳	۵	۶		۱۴/۹ (۲)
۶	۰	۰	۱	۴	۵	۷	۹	۱۵/۳ (۳)
								۱۶/۲ (۴)

پاسخ: گزینه «۳» نمودار جعبه‌ای به فرم زیر است:



در واقع سؤال واریانس بین چارک اول و سوم را خواسته است. (داخل جعبه)

از روی نمودار شاخه و برگ داده شده ابتدا چارک اول و سوم را به دست می‌آوریم. برای چارک اول، میانه کل داده‌ها را به دست آورده و مجدداً میانه را در ۵۰٪ اول داده‌ها مشخص می‌کنیم و در ۵۰٪ دوم باز هم میانه را بدست می‌آوریم که همان چارک سوم است.

۴۱, ۴۱, ۴۲, ۴۴, ۴۶, ۴۸, ۵۰, ۵۲, ۵۲, ۵۳, ۵۵, ۵۶, ۶۰, ۶۰, ۶۱, ۶۴, ۶۵, ۶۷, ۶۹

\downarrow Q_1 \downarrow Me \downarrow Q_3

اکنون داده‌ها بین چارک اول و سوم عبارتند از:

۴۸, ۵۰, ۵۲, ۵۲, ۵۳, ۵۵, ۵۶, ۶۰, ۶۰

می‌دانیم که اضافه یا کم کردن تأثیری در واریانس ندارد. بنابراین می‌توانیم با کم کردن عددی دلخواه مقیاس داده‌ها را کوچک کنیم. از هر کدام از داده‌ها ۵۳ واحد کم می‌کنیم تا داده‌های جدید با مقیاس کوچک‌تر و بهتر حاصل شود.

$$-5, -3, -1, -1, 0, 2, 3, 7, 7 \Rightarrow \sigma^2 = \frac{\sum x_i^2}{N} - \mu^2$$

$$\sigma_x^2 = \frac{25 + 9 + 1 + 1 + 4 + 9 + 49 + 49}{9} - \left(\frac{9}{9}\right)^2 = 16/33 - 1 = 15/33$$

(مدیریت - دکتری ۹۵)

کلمه مثال ۱۴۰: کدام یک از نمودارهای ذیل برای نمایش متغیرهای کیفی کاربرد دارد؟

- (۱) شاخه و برگ (۲) هیستوگرام (۳) جعبه‌ای (۴) پارتو

پاسخ: گزینه «۴» برای رسم داده‌های کمی (فاصله‌ای و نسبی) از نمودارهای ۱- هیستوگرام ۲- چندصنفی ۳- شاخ و برگ ۴- جعبه‌ای استفاده می‌شود و برای رسم داده‌های کیفی (اسمی و رتبه‌ای) از نمودارهای ۱- ستونی (میله‌ای) و ۲- دایره‌ای استفاده می‌شود.

توجه: داوطلبانی که بعد از مطالعه کل کتاب نیاز به تست بیشتری برای مرور و تمرین دارند، می‌توانند با مراجعه به

سایت www.modaresanesharif.ac.ir بانک تست‌های مربوط به همه فصول کتاب را دانلود نمایند.

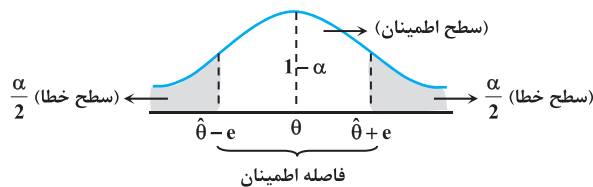
درسنامه (۲): برآورد فاصله‌ای



در بخش قبل میانگین نمونه (\bar{X}) یک برآوردگر خوب برای میانگین جامعه (μ) و همچنین واریانس نمونه (S^2) و نسبت نمونه (\bar{P}) برآوردگرهای خوبی برای واریانس و نسبت کل جامعه بودند. اما میانگین، واریانس و نسبت نمونه از نمونه‌ای به نمونه‌ای دیگر تغییر می‌کند، لذا نمی‌توان دقیقاً ادعا کرد که μ, σ^2, p دقیقاً برابر با \bar{X}, S^2, \bar{P} است. بنابراین برای آنکه برای تمام نمونه‌ها بتوانیم برآورد دقیقی از پارامترهای جامعه داشته باشیم، یک برآورد فاصله‌ای به نام فاصله‌ی اطمینان تعریف می‌کنیم. که این فاصله به صورت $(\hat{\theta} - e, \hat{\theta} + e)$ است و به e خطای برآورد گفته می‌شود. در برآورد فاصله‌ای به این علت که یک فاصله با درجه‌ی خطا یا اطمینان برای پارامتر به دست می‌آوریم درجه‌ی صحت برآورد مشخص می‌شود.

برآورد فاصله‌ای (فاصله اطمینان) یک فاصله است که بیان می‌دارد که با درصد اطمینان $100(1-\alpha)$ پارامتر جامعه (θ) در فاصله‌ای به صورت $(\hat{\theta} - e, \hat{\theta} + e)$ قرار می‌گیرد، مثلاً اگر برای میانگین جامعه یک فاصله‌ی اطمینان ۹۵٪ درصد به صورت $(-2, 4)$ بدست آید، این به مفهوم آن است که با اطمینان ۹۵ درصد میانگین جامعه در فاصله $(-2, 4)$ قرار دارد. توجه کنید که درصد اطمینان برای ساختن فواصل اطمینان معمولاً ۹۰٪ و ۹۵٪ و ۹۹٪ در نظر گرفته می‌شود که به آن‌ها سطوح اطمینان گفته می‌شود. در اینصورت سطح خطاها به صورت ۱۰٪ و ۵٪ و ۱٪ خواهد بود. رابطه‌ی احتمالی برای فاصله‌ی اطمینان به صورت زیر است:

$$P(\hat{\theta} - e < \theta < \hat{\theta} + e) = 1 - \alpha$$



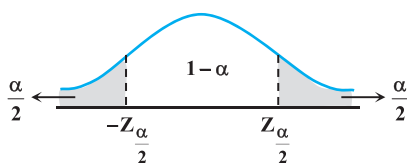
به شکل مقابل توجه کنید:

این شکل بیان می‌دارد که با احتمال $(1-\alpha)$ اطمینان داریم که فاصله‌ی $(\hat{\theta} - e, \hat{\theta} + e)$ پارامتر θ را در بر می‌گیرد و با احتمال $(\frac{\alpha}{2} + \frac{\alpha}{2} = \alpha)$ احتمال می‌دهیم که پارامتر (θ) در فاصله‌ی $(\hat{\theta} - e, \hat{\theta} + e)$ قرار نگیرد و همچنین می‌خواهیم به اندازه‌ی مقدار e ، برآورد $\hat{\theta}$ یعنی θ از پارامتر θ تفاوت داشته باشد. بنابراین $1-\alpha$ سطح اطمینان و α سطح خطا و خطای برآورد $e = |\hat{\theta} - \theta|$ است. اکنون در اینجا به معرفی فواصل اطمینان برای پارامترهای جامعه μ, σ^2, p در شرایط مختلف می‌پردازیم.

مثال ۴۳: هر چقدر سطح خطای α کاهش یابد کدام گزینه صحیح است؟ (بقیه عوامل ثابت)

- (۱) سطح اطمینان $1-\alpha$ افزایش می‌یابد.
- (۲) خطای برآورد افزایش می‌یابد.
- (۳) ضریب اطمینان افزایش می‌یابد.
- (۴) هر سه گزینه

پاسخ: گزینه «۴» به شکل مقابل دقت کنید:



هر چقدر سطوح هاشور خورده کمتر شود، سطح هاشور نخورده افزایش می‌یابد.

و مقادیر $Z_{\frac{\alpha}{2}}$ افزایش می‌یابد و خطای برآورد (e) نیز افزایش می‌یابد. به مثال‌های عددی زیر توجه کنید:

$$\begin{cases} \alpha = 0.05 \Rightarrow 1 - \alpha = 1 - 0.05 = 0.95 \\ \frac{\alpha}{2} = 0.025 \Rightarrow Z_{0.025} = 1.96 \end{cases} \Rightarrow e = Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \quad (\text{به مقادیر } Z_{\frac{\alpha}{2}} \text{ یا } Z_{\alpha} \text{ ضریب اطمینان گفته می‌شود}).$$

$$\begin{cases} \alpha = 0.01 \Rightarrow 1 - \alpha = 1 - 0.01 = 0.99 \\ \frac{\alpha}{2} = 0.005 \Rightarrow Z_{0.005} = 2.58 \end{cases} \Rightarrow e = Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$

پس اگر $\alpha = 0.05$ به $\alpha = 0.01$ کاهش یابد، هر سه گزینه‌ی دیگر افزایش می‌یابند.

مثال ۴۴: فرض کنید X دارای تابع چگالی $f_x(x) = \frac{1}{\theta} e^{-\frac{x}{\theta}}$ ، $x > 0$ است. مقدار c چقدر باشد تا فاصله‌ی $(0, cX)$ یک فاصله اطمینان ۹۰٪ برای θ در نظر گرفته شود؟

(علوم اقتصادی - دکتری ۹۴)

(۴) $-\ln(0.9)$

(۳) $-\frac{1}{\ln(0.9)}$

(۲) $\ln\left(\frac{10}{9}\right)$

(۱) $\ln(9)$

پاسخ: گزینه «۳» مفهوم فاصله اطمینان عبارت است از: $P(a < \theta < b) = 1 - \alpha \xrightarrow{\text{طبق گفته مسئله}} P(0 < \theta < CX) = 0/90$

برای محاسبه احتمال در یک بازه روی آن بازه از تابع چگالی احتمال انتگرال گیری می‌کنیم:

$$P\left(x > \frac{\theta}{C}\right) = 0/90 \Rightarrow \int_{\frac{\theta}{C}}^{\infty} \frac{1}{\theta} e^{-\frac{x}{\theta}} dx = 0/90 \Rightarrow -e^{-\frac{x}{\theta}} \Big|_{\frac{\theta}{C}}^{\infty} = 0/90 \Rightarrow e^{-\frac{1}{C}} = 0/90$$

$$\Rightarrow -\frac{1}{C} = \ln 0/90 \Rightarrow C = \frac{-1}{\ln 0/90}$$

مثال ۴۵: فرض کنید X طول عمر یک لامپ باشد که دارای توزیع یکنواخت روی فاصله $(0, \theta)$ می‌باشد و θ پارامتر نامعلوم است. مقدار c چقدر

باشد تا فاصله تصادفی $\left(\frac{c}{2X}, \frac{c}{X}\right)$ مقدار $\frac{1}{\theta}$ را با احتمال $0/90$ شامل شود؟ (علوم اقتصادی - سراسری ۹۶)

۱ (۴)

۱/۸ (۳)

۰/۱ (۲)

۰/۹ (۱)

پاسخ: گزینه «۳» طبق تعریف فاصله اطمینان $P(a < X < b) = 1 - \alpha$ می‌باشد. رابطه داده شده را به این فرم تبدیل می‌کنیم:

$$P\left(\frac{c}{2X} < \frac{1}{\theta} < \frac{c}{X}\right) = 0/90 \Rightarrow P\left(\frac{c\theta}{2} < X < c\theta\right) = 0/90$$

متغیر تصادفی X دارای توزیع یکنواخت بر روی فاصله $(0, \theta)$ است، بنابراین تابع چگالی آن $0 < X < \theta$ و $f(x) = \frac{1}{\theta}$ می‌باشد. برای محاسبه احتمال روی بازه خواسته شده از تابع چگالی انتگرال گیری می‌کنیم:

$$\int_{\frac{c\theta}{2}}^{c\theta} \frac{1}{\theta} dx = 0/90 \Rightarrow \frac{x}{\theta} \Big|_{\frac{c\theta}{2}}^{c\theta} = 0/90 \Rightarrow \left(\frac{c\theta}{\theta} - \frac{c\theta}{2\theta}\right) = 0/90 \Rightarrow c - \frac{c}{2} = 0/90 \Rightarrow \frac{1}{2}c = 0/90 \Rightarrow c = 1/8$$

توجه: در مسائل آمار استنباطی معمولاً ۳ سطح خطا و ۳ سطح اطمینان بیشتر مورد نظر است. جدول زیر این سطوح را مشخص کرده و ضرایب اطمینان آن‌ها را با توجه به جدول نرمال استاندارد مشخص کرده است. به دانشجویان توصیه می‌شود این جدول را به‌خاطر بسپارند:

ضریب اطمینان	سطح خطا	سطح اطمینان
$Z_{\alpha} = 1/28 ; Z_{\frac{\alpha}{2}} = 1/645$	$\alpha = 0/1 ; \frac{\alpha}{2} = \frac{0/1}{2} = 0/05$	$(1 - \alpha) = 0/90$
افزایش $Z_{\alpha} = 1/645 ; Z_{\frac{\alpha}{2}} = 1/96$	کاهش $\alpha = 0/05 ; \frac{\alpha}{2} = \frac{0/05}{2} = 0/025$	افزایش $(1 - \alpha) = 0/95$
$Z_{\alpha} = 2/32 ; Z_{\frac{\alpha}{2}} = 2/58$	$\alpha = 0/01 ; \frac{\alpha}{2} = \frac{0/01}{2} = 0/005$	$(1 - \alpha) = 0/99$

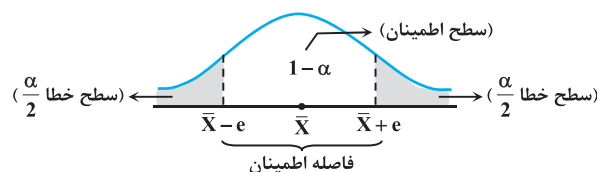
نکته ۱: هرگاه مسأله مقدار α را مشخص نکرد پیش‌فرض بر آن است که $\alpha = 0/05$ است.

فاصله اطمینان یا برآورد فاصله‌ای برای میانگین جامعه μ

یک فاصله اطمینان برای میانگین، به حالت‌های زیر بوجود می‌آید:

۱. جامعه نرمال و واریانس جامعه σ^2 معلوم

در جامعه‌ای نرمال با میانگین μ و واریانس معلوم σ^2 ، نمونه‌ای تصادفی به حجم n اختیار می‌کنیم. اگر میانگین این نمونه n تایی \bar{X} باشد یک فاصله اطمینان $(1 - \alpha)100\%$ درصد برای میانگین جامعه به‌صورت زیر است:



$Z_{\frac{\alpha}{2}}$: از روی جدول نرمال استاندارد بدست می‌آید.

$$\bar{x} - \frac{\sigma}{\sqrt{n}} Z_{\frac{\alpha}{2}} < \mu < \bar{x} + \frac{\sigma}{\sqrt{n}} Z_{\frac{\alpha}{2}}$$

مثال ۴۶: از جامعه آماری در نمونه‌ای به حجم ۴۰۰ مقدار میانگین ۵۲ و انحراف معیار ۱۶ محاسبه شده است. با احتمال ۹۵ درصد میانگین این جامعه در کدام بازه است؟ (مدیریت و حسابداری - سراسری ۹۲)

(۱) $(50/432, 53/568)$ (۲) $(50/324, 53/676)$ (۳) $(50/716, 53/284)$ (۴) $(51/024, 52/976)$

پاسخ: گزینه «۱» از رابطه فاصله اطمینان برای میانگین استفاده می‌کنیم:

$$\bar{x} - \frac{\sigma}{\sqrt{n}} Z_{\frac{\alpha}{2}} < \mu < \bar{x} + \frac{\sigma}{\sqrt{n}} Z_{\frac{\alpha}{2}}$$

در این مسئله $n = 400$, $\bar{X} = 52$ و $\sigma = 16$ و $\frac{\alpha}{2} = 0.025 \rightarrow \alpha = 0.05 \rightarrow 1 - \alpha = 0.95$ مقادیر را در فرمول بالا جایگذاری می‌کنیم:

$$52 - \frac{16}{\sqrt{400}} \times 2 < \mu < 52 + \frac{16}{\sqrt{400}} \times 2 \Rightarrow 52 - \frac{32}{20} < \mu < 52 + \frac{32}{20} \Rightarrow 50/432 < \mu < 53/568$$

مثال ۴۷: از توزیع نرمال یک نمونه به حجم ۱۶ اختیار شده است. می‌دانیم واریانس جامعه ۶۴ می‌باشد. اگر میانگین نمونه ۱۰ باشد، با اطمینان ۹۵٪

(محیط زیست - سراسری ۸۵)

میانگین جامعه در چه فاصله‌ای قرار دارد؟ $\int_{-4}^{-2} f(y) dy \approx 0.25$

(۱) $[0, 6]$ (۲) $[6, 8]$ (۳) $[6, 14]$ (۴) $[8, 14]$

پاسخ: گزینه «۳» توزیع جامعه نرمال و واریانس جامعه معلوم است. بنابراین مقدار نمونه (n) هر چه باشد توزیع \bar{X} نرمال خواهد بود و فاصله

اطمینان μ بنابر مورد (الف) عبارت است از:

$$\mu: \bar{x} \pm Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \Rightarrow \mu: 10 \pm 2 \frac{\sqrt{64}}{\sqrt{16}} = 10 \pm 2 \times \frac{8}{4} = 10 \pm 4 = (6, 14)$$

$$\bar{x} = 10, \sigma^2 = 64, n = 16, \alpha = 0.05 \rightarrow Z_{\frac{\alpha}{2}} = Z_{0.025} \approx 2$$

مثال ۴۸: در سطح معنی‌داری ۵ درصد، برای m مشاهده از جامعه‌ای نرمال با واریانس ۴، حدود اطمینان میانگین جامعه (۵ و ۳) به دست آمده است. اگر اندازه نمونه را ۴ برابر کنیم، عرض فاصله اطمینان جدید با همان سطح اطمینان قبلی برابر خواهد بود با:

(حسابداری و مدیریت - دکتری ۹۳)

(۱) ۳ (۲) ۲ (۳) ۱ (۴) ۴

پاسخ: گزینه «۳» در یک توزیع نرمال با واریانس معلوم فاصله اطمینان به صورت:

$$\left(\bar{x} - \frac{\sigma}{\sqrt{m}} Z_{\frac{\alpha}{2}}, \bar{x} + \frac{\sigma}{\sqrt{m}} Z_{\frac{\alpha}{2}} \right)$$

می‌باشد که عرض این فاصله اطمینان به صورت:

$$\left(\bar{x} + \frac{\sigma}{\sqrt{m}} Z_{\frac{\alpha}{2}} \right) - \left(\bar{x} - \frac{\sigma}{\sqrt{m}} Z_{\frac{\alpha}{2}} \right) = \frac{2\sigma}{\sqrt{m}} Z_{\frac{\alpha}{2}}$$

می‌باشد اکنون m تبدیل به $4m$ شود عرض به صورت روبرو تبدیل می‌شود:

$$\frac{2\sigma}{\sqrt{4m}} Z_{\frac{\alpha}{2}} \Rightarrow \frac{\cancel{2}\sigma}{\cancel{2}\sqrt{m}} Z_{\frac{\alpha}{2}}$$

بنابراین عرض فاصله اطمینان جدید نصف خواهد شد. اما توجه کنید که عرض فاصله اطمینان قبلی $2 = 5 - 3$ بوده است بنابراین عرض جدید نصف ۲ یعنی ۱ است.

تذکره: در فاصله‌ی اطمینان بالا به مقدار $e = Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$ خطای نمونه‌گیری یا خطای حدی یا خطای برآورد گفته می‌شود. همچنین به راحتی

می‌توانیم با محاسبه‌ی تفاضل حد بالا و پایین فاصله‌ی اطمینان، طول فاصله را به دست آوریم.

$$\text{بنابراین طول فاصله اطمینان برابر با } 2e = \left(\bar{x} + Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \right) - \left(\bar{x} - Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \right) = 2 Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} = 2e$$

است.



مثال ۵۳: تابع مصرف از نمونه‌ای با حجم $N=66$ به صورت زیر برآورد گردیده است:

اگر ضریب همبستگی بین درآمد قابل تصرف (Y_{dt}) و مصرف (C_t) برابر $0/6$ برآورد شده باشد، مقدار آماره آزمون معنی دار بودن میل نهایی به مصرف کدام یک از گزینه‌های زیر است؟

- (۱) $0/6 -$ (۲) 14 (۳) $1/4$ (۴) $6 +$

پاسخ: گزینه «۴» مقدار آماره آزمون برای معنی دار بودن رگرسیون عبارت است از: $t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$ بنابراین:

$$t = \frac{0/6\sqrt{66-2}}{\sqrt{1-(0/6)^2}} = 6$$

مثال ۵۴: در صورتی که در یک رگرسیون با دو متغیر مستقل که دارای ۲۳ مورد است ($N=23$)، همبستگی متغیرهای مستقل با متغیر وابسته ($0/8$) باشد، مقدار عددی F برابر است با:

- (۱) $1/782$ (۲) $8/78$ (۳) $12/78$ (۴) $17/8$

پاسخ: گزینه «۴» آماره آزمون F با استفاده از ضریب تعیین عبارت است از:

$$F_{k-1, n-k} = \frac{\frac{R^2}{k-1}}{\frac{1-R^2}{n-k}}$$

R^2 : ضریب تعیین متغیرهای مستقل با متغیر وابسته است. k : تعداد متغیرهای مستقل و وابسته.

که در این جا: $k=3$ است، دو متغیر مستقل و یک متغیر وابسته، مقادیر را جایگذاری می‌کنیم:

$$F_{2,20} = \frac{0/64}{\frac{3-1}{1-0/64}} = \frac{0/32}{0/518} = 17/78 \approx 17/8, \quad R^2 = (0/8)^2 = 0/64$$

۳- آزمون صفر بودن زیر مجموعه‌ای از ضرایب رگرسیون (آزمون F جزئی):

فرض کنید R_k^2 ضریب همبستگی چندگانه باشد وقتی همه k متغیر در مدل رگرسیونی باشند و R_m^2 ضریب همبستگی چندگانه باشد وقتی m متغیر مشخص داخل مدل باشد پس فرضیه صفر بیان می‌کند که $k-m$ متغیر تعیین شده، دارای ضریب رگرسیونی صفر است.

$$F = \frac{\frac{R_k^2 - R_m^2}{k-m}}{\frac{1-R_k^2}{n-k-1}} \quad F > F_{\alpha, k-k, n-k-1} \quad (\text{ناحیه بحرانی})$$

اگر هر یک از آزمون‌های t مربوط به ضرایب انفرادی رگرسیون معنی‌دار باشد (H_0 رد شود) آنگاه F برای آزمون معنی‌داری مدل نیز معنی‌دار است (H_0 رد می‌شود).

شکل ماتریسی رگرسیون خطی چندگانه:

مدل رگرسیون $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + e$ را می‌توان به شکل ماتریسی $y = X\beta + e$ نوشت که y بردار $n \times 1$ از مشاهدات پاسخ، X یک ماتریس $n \times p$ از متغیرهای مستقل، β یک بردار $p \times 1$ از ضرایب رگرسیونی و e بردار $n \times 1$ از خطاهای تصادفی است.

برآورد کمترین مربعات بردار β : برآوردگر حداقل مربعات بردار β به شرط وجود ماتریس XX' به صورت مقابل می‌باشد:

$$\hat{\beta} = (X'X)^{-1} X'y$$

نکته ۷: می‌توان دید که $\hat{y} = X\hat{\beta}$ یعنی $\hat{y} = X(X'X)^{-1} X'y$ که به صورت $\hat{y} = Hy$ نمایش داده می‌شود که H ماتریس برازش نامیده می‌شود.

$$e = y - \hat{y} = y - X\hat{\beta} = y - Hy = (I - H)y$$

خواص برآوردگرهای حداقل مربعات:

$$1. \hat{\beta} \text{ برآورد نااریب بردار } \beta \text{ است یعنی } E(\hat{\beta}) = \beta \quad 2. \text{ کوواریانس بردار } \hat{\beta} \text{ یک ماتریس متقارن می‌باشد.} \quad 3. \text{Cov}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$$

معیارهای انتخاب مدل رگرسیونی

پس از برآورد کردن ضرایب رگرسیونی، معیارهایی وجود دارد که باید تست شوند؛ چرا که آنها فرضیات مدل‌های رگرسیونی را آزمایش می‌کنند و نشان می‌دهند که مدلی که به دست آمده است مدل مناسبی است و یا خطاهایی در آن‌ها وجود دارد.

۱- ضریب تعیین چندگانه: یکی از معیارهای ارزیابی مدل رگرسیون ضریب تعیین چندگانه است که به صورت زیر تعریف می‌شود:

$$R^2 = \frac{SS_{\text{Reg}}}{SS_y} = \frac{SS_y - SSE}{SS_y} = 1 - \frac{SSE}{SS_y}$$

که در اینجا SS_{Reg} مجموع مربعات رگرسیونی و SSE مجموع مربعات خطا و SS_y مجموع مربعات انحراف داده‌ها از میانگین نام دارند. R^2 به عنوان معیاری برای کاهش در تغییرپذیری Y معرفی می‌شود و $0 \leq R^2 \leq 1$ ولی R^2 بزرگ لزوماً یک رگرسیون خوب را نتیجه نمی‌دهد، چرا که با اضافه کردن یک متغیر مستقل به مدل، R^2 افزایش می‌یابد. لذا از آماره‌ای به نام R^2 اصلاح شده (تصحیح شده) به صورت $\bar{R}_k^2 = 1 - \frac{n-1}{n-k-1}(1-R^2)$ استفاده می‌کنند. در یک مدل رگرسیونی مقدار $R^2 = 95\%$ شده است این به مفهوم آن است که 95% از تغییرپذیری Y به وسیله متغیرهای مستقل که در مدل وجود دارند، توضیح داده شده‌اند.

کلمه مثال ۵۵: اگر ضریب همبستگی بین دو متغیر $0/6$ و دو متغیر دیگر $0/3$ باشد، می‌توان گفت همبستگی دو متغیر اول: «چند برابر قوی‌تر» از دو متغیر دوم است؟

(مدیریت و حسابداری - دکتری ۹۲)

(۴) سه

(۳) نه

(۲) چهار

(۱) دو

پاسخ: گزینه «۲» یکی از معیارهای مهم برای مقایسه دو متغیر ضریب تعیین بین آنها است و نباید براساس ضریب همبستگی آنها را مقایسه کرد:

$$r_1 = 0/6 \Rightarrow r_1^2 = 0/36 \Rightarrow \text{همبستگی دو متغیر اول ۴ برابر همبستگی دو متغیر دوم است}$$

$$r_2 = 0/3 \Rightarrow r_2^2 = 0/9$$

۲- خود همبستگی (آزمون دوربین - واتسون): یکی از فرض‌های مهم در رگرسیون آن است که خطاهای مدل e_t ها متغیرهای تصادفی ناهمبسته باشند. در بسیاری از کاربردهای تحلیل رگرسیون با داده‌هایی سروکار داریم که ممکن است این فرض برای آن‌ها برقرار نباشد. برای آزمون کردن این فرضیه از آزمونی به نام آزمون دوربین - واتسون استفاده می‌کنیم. دوربین - واتسون یک آزمون است که فرضیه‌ها، آماره‌ی آزمون و ناحیه‌ی بحرانی آن به صورت زیر است:

$$D = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2} \quad (\text{آماره‌ی آزمون})$$

$$\begin{cases} H_0: \rho = 0 & (\text{همبستگی بین خطا وجود ندارد}) \\ H_1: \rho > 0 & (\text{همبستگی بین خطاها وجود دارد}) \end{cases}$$

مقدار آماره‌ی آزمون D ، دارای حداقل مقدار صفر و حداکثر مقدار 4 است. اگر فرضیه‌ی H_0 درست باشد مقدار مورد انتظار آماره‌ی آزمون برابر با 2 خواهد بود. ناحیه‌ی بحرانی این آزمون به صورت زیر است:

اگر $D > d_u$ باشد فرض صفر $H_0: \rho = 0$ رد نمی‌شود. اگر $D < d_L$ باشد فرض صفر $H_0: \rho = 0$ رد می‌شود.

اگر $d_L < D < d_u$ باشد، آزمون بی‌نتیجه خواهد بود.

به d_u و d_L مقادیر بحرانی این آزمون گفته می‌شود که از جدولی به نام جدول دوربین واتسون مشخص می‌شوند.

توجه: وجود خود همبستگی در خطاها آثار سوء متعددی دارد که در زیر به آن‌ها اشاره خواهیم کرد:

(۱) ضرایب رگرسیونی برآورد شده اگر چه هنوز نااریب هستند، اما دیگر برآوردهایی با کمترین واریانس نیستند، بنابراین با تورم واریانس ضرایب، روبرو خواهیم شد.

(۲) فواصل اطمینان برای ضرایب کوچکتر، از آن چه در واقع باید باشند بزرگ‌تر خواهند شد و همچنین $\hat{\sigma}^2 = S_e^2$ به طور کلی کم برآورد خواهد شد.

(۳) به دست آوردن فواصل اطمینان و آزمون فرضیه‌ها بر اساس توزیع‌های t و F دیگر مناسب نیست.

مثال ۵۶: در یک خروجی کامپیوتری مقدار آماری آزمون دوربین = واتسون برابر با $2/08$ به دست آمده است اگر مقادیر بحرانی به ترتیب $0/95$ و $1/54$ به دست آمده باشد کدام گزینه صحیح است؟

- (۱) فرضیه همبستگی خطاها برابر با صفر رد می‌شود.
 (۲) فرضیه همبستگی خطاها برابر با صفر رد نمی‌شود.
 (۳) آزمون بی‌نتیجه است.
 (۴) نمی‌توان قضاوت کرد.

پاسخ: گزینه «۲» از آنجائیکه مقدار آماری آزمون $D = 2/08$ از مقدار بحرانی $1/54$ بزرگتر است، لذا فرض $H_0: \rho = 0$ رد نمی‌شود.

مثال ۵۷: برای قضاوت در مورد همبستگی خطاها در یک مدل رگرسیون چند متغیره مقادیر $\sum_{t=2}^{20} (e_t - e_{t-1})^2 = 3$ و $\sum_{t=1}^{20} e_t^2 = 2$ به دست آمده‌اند.

مقدار آماری آزمون مناسب کدام است؟

- (۱) $1/08$ (۲) $1/5$ (۳) $1/01$ (۴) 1

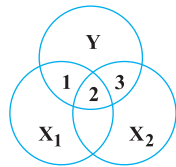
پاسخ: گزینه «۲» طبق تعریف آماری آزمون دوربین - واتسون خواهیم داشت:

$$D = \frac{\sum_{t=2}^{20} (e_t - e_{t-1})^2}{\sum_{t=1}^{20} e_t^2} = \frac{3}{2} = 1/5$$

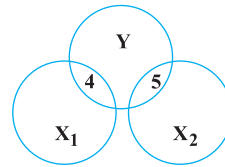
توجه: در حالتی که فرضیه‌ها به صورت مقابل باشند $\begin{cases} H_0: \rho = 0 \\ H_1: \rho < 0 \end{cases}$ آماری آزمون مناسب به صورت $D^* = 4 - D$ می‌باشد و قضاوت و نواحی بحرانی

به صورت قبل است.

۳- هم خطی: در اغلب مسائل رگرسیون متغیرهای مستقل با یکدیگر همبسته‌اند در مواردی که این همبستگی متقابل خیلی شدید باشد، گوئیم هم خطی وجود دارد که تأثیر بسیار بدی بر برآورد پارامترهای ضرایب رگرسیونی خواهد داشت و باعث تورم واریانس‌های برآوردگرها خواهد شد و روش حداقل مربعات خطا، برآوردهای بسیار ضعیفی از هر یک از پارامترهای مدل به دست خواهد داد. به شکل‌ها توجه کنید.



وجود هم خطی بین X_2, X_1



عدم وجود هم خطی بین X_2, X_1

در شکل بالا ناحیه‌ی (۲) تغییرات Y را نشان می‌دهد که توسط هر دو متغیر X_1 و X_2 توضیح داده شده است که به دلیل وجود هم خطی باعث افزایش تغییرات واقعی Y است و به عبارت دیگر ناحیه (۲) دو بار در محاسبات واریانس Y محاسبه می‌شود و واریانس Y بیش از مقدار واقعی نشان داده خواهد شد.

مثال ۵۸: برای به دست آوردن معادله رگرسیون چند متغیره (چندگانه) خطی حاصل از 10 متغیر مستقل، کدام مورد درست است؟

(مدیریت - دکتری ۹۴)

- (۱) بهترین ترکیب از متغیرهای مستقل، یک ترکیب حداقل سه تایی است.
 (۲) بهترین روش تهیه معادله رگرسیون، روش کل‌گرا (Enter) است.
 (۳) بهترین ترکیب از متغیرهای مستقل، یک ترکیب ده تایی است.
 (۴) بین متغیرهای مستقل هم خطی وجود نداشته باشد.

پاسخ: گزینه «۴» یکی از موضوعات بسیار مهم در رگرسیون، هم خطی چندگانه است که باید بین متغیرهای مستقل همبستگی وجود داشته باشد؛ در مواردی که بین این متغیرها همبستگی متقابل خیلی شدید باشد، می‌گوییم هم خطی وجود دارد که تأثیر بسیار بدی بر برآورد پارامترهای ضرایب رگرسیونی خواهد داشت و باعث تورم واریانس این برآوردها خواهد شد.

درسنامه (۲): آنتروپی (Entropy)



کلمه آنتروپی به مفهوم بی‌نظمی است که اولین بار توسط ژول مطرح شد. این مفهوم در آمار توسط شانون در اواسط قرن بیستم معرفی شد. در آمار، آنتروپی به عنوان یک معیار عدم قطعیت به کار برده می‌شود. در مسائل مدیریتی واژه‌ی آنتروپی حاکی از تمایل سیستم‌ها به بی‌نظمی است. سیستم‌های بسته به مرور زمان از بین می‌روند و اصطلاحاً آنتروپی آنها مثبت است ولی سیستم‌های باز دارای آنتروپی منفی هستند.

❖ **تعریف:** اگر X یک متغیر تصادفی دلخواه با تابع چگالی احتمال $f(x)$ باشد تابع $H(X) = E(-\log f(X))$ را آنتروپی X یا توزیع X گویند. (لگاریتم غالباً در پایه‌ی ۲ یا e منظور می‌شود)

لگاریتم بر مبنای ۲ مهم‌ترین معیار می‌باشد که وقتی بر مبنای ۲ است، آنتروپی براساس بیت می‌آید. یعنی وقتی $H(X) = ۲ \text{ bit}$ شود، می‌توانیم با پرسیدن ۲ سؤال که جواب آنها بله یا خیر است به جواب اصلی دست پیدا کنیم.

اگر X یک متغیر تصادفی گسسته با توزیع احتمال $P_i = f(x_i) = P(X = x_i) \quad i = 1, 2, \dots, n$ باشد، آنگاه آنتروپی X عبارت است از:

$$H(X) = E(-\log f(X)) = -\sum f(x) \log f(x) = -\sum_{i=1}^n P_i \log P_i$$

به سادگی مشاهده می‌شود که تابع آنتروپی وابسته به احتمال‌های P_1, P_2, \dots, P_n است. لذا هر چقدر آنتروپی بیشتر باشد، بی‌نظمی یا عدم قطعیت بیشتر است.

📌 **مثال ۵۹:** کدام یک از عبارات‌های زیر «درجه نامعینی» یک آزمایش α را که k نتیجه ممکن $A_1, A_2, A_3, \dots, A_k$ هم احتمال دارد. آنتروپی (Entropy) آزمایش α می‌نامند: (مدیریت صنعتی - آزاد ۸۸)

$$H(\alpha) = \underbrace{k \log k + k \log k + \dots + k \log k}_k \quad (۲)$$

$$H(\alpha) = \underbrace{-k \log k - k \log k - \dots - k \log k}_k \quad (۱)$$

$$H(\alpha) = \underbrace{-\frac{1}{k} \log \frac{1}{k} - \frac{1}{k} \log \frac{1}{k} - \dots - \frac{1}{k} \log \frac{1}{k}}_k \quad (۴)$$

$$H(\alpha) = \underbrace{\frac{1}{k} \log \frac{1}{k} + \frac{1}{k} \log \frac{1}{k} + \dots + \frac{1}{k} \log \frac{1}{k}}_k \quad (۳)$$

☑ **پاسخ:** گزینه «۴» طبق تعریف آنتروپی که در بالا گفته شد می‌توانیم نتیجه بگیریم که همه A_i ها هم شانس هستند یعنی $P(A_i) = \frac{1}{k}$ بنابراین:

$$i = 1, 2, \dots, k \quad H(\alpha) = -\frac{1}{k} \log \frac{1}{k} - \frac{1}{k} \log \frac{1}{k} - \dots - \frac{1}{k} \log \frac{1}{k}$$

(مدیریت صنعتی - آزاد ۸۸)

📌 **مثال ۶۰:** نتیجه آزمایش α توسط جدول زیر بیان می‌شود. آنتروپی این آزمایش کدام است؟

x	A_1	A_2	A_3
$P(A_i)$	$\frac{2}{5}$	$\frac{2}{5}$	$\frac{1}{5}$

= ۱

$$H(\alpha) = 0.4791 \quad (۲) \quad H(\alpha) = 0.4581 \quad (۱)$$

$$H(\alpha) = 0.4381 \quad (۴) \quad H(\alpha) = 0.4497 \quad (۳)$$

$$H(\alpha) = -\sum_{i=1}^n P_i \log P_i = -P_1 \log P_1 - P_2 \log P_2 - P_3 \log P_3$$

☑ **پاسخ:** گزینه «۱» با توجه به رابطه آنتروپی خواهیم داشت:

$$= -0.4 \log 0.4 - 0.4 \log 0.4 - 0.2 \log 0.2 = (-0.4)(-0.398) - (0.4)(-0.398) - (0.2)(-0.699) = 0.4581$$

📌 **مثال ۶۱:** دو متغیر تصادفی گسسته X و Y با توزیع‌های زیر مفروض است، آنتروپی‌ها X و Y را محاسبه کنید.

$$q_i = P(Y = x_i) = \begin{cases} \frac{1}{64} & i = 1, 2, 3 \\ \frac{61}{64} & i = 4 \end{cases} \quad P_i = P(X = x_i) = \frac{1}{4} \quad i = 1, 2, 3, 4$$

☑ **پاسخ:** از رابطه گفته شده بالا استفاده می‌کنیم، توجه کنید که متغیرهای تصادفی X و Y گسسته‌اند.

$$H(X) = -\sum_{i=1}^4 P_i \log P_i = -\frac{1}{4} \sum_{i=1}^3 \log \frac{1}{4} = -\frac{1}{4} \times \sum_{i=1}^3 (\log_2 1 - \log_2 4) = -\frac{1}{4} \times \sum_{i=1}^3 -\log_2 4 = 2$$

$$H(Y) = -\sum_{i=1}^4 q_i \log q_i = -\frac{1}{64} \sum_{i=1}^3 \log \frac{1}{64} - \frac{61}{64} \log \frac{61}{64} = 0.215$$