

Â¶HkhTwh#ÁI¿°¼¶pA#½sÄ»#R⁻IpTeH#»#nI¶A#½»q]

آمار توصیفی و احتمالات

جزوه درسی

WWW.SHOP.IRANESTEKHDAM.IR

دانلود نمونه سوالات استخدامی

آمار توصیفی

شامل مجموعه‌ای از روش‌ها برای برنامه‌ریزی آزمایش‌ها، بدست آوردن اطلاعات و سپس سازماندهی، خلاصه‌سازی، ارائه و تجزیه و تحلیل اطلاعات و در نهایت نتیجه‌گیری بر مبنای این اطلاعات می‌باشد. آمار در لغت به معنای شمارش است. آمار توصیفی، مجموعه روشهایی است که به خلاصه کردن، طبقه بندی، توصیف و تفسیر داده‌ها می‌پردازد. هدف آمار توصیفی، توصیف واقعیات موجود است که تنها خطای موجود در آن خطای اندازه‌گیری است، مشخص کردن قد دانش آموزان یک کلاس، بررسی روند رشد سرانه ملی در سه سال گذشته نمونه‌ای از مفاهیم مورد بررسی آمار توصیفی هستند. نمادهای مهم در آمار عبارتند از:

شاخص مورد مطالعه	نمونه (آماره)	جامعه (پارامتر)
میانگین	\bar{X}	μ (مو)
نسبت	ρ	\hat{P} (پی هد)
همبستگی	r_{xy}	ρ_{xy} (رو)
واریانس	S^2	δ^2 (سیگما به توان دو)
انحراف معیار	S	δ (سیگما)
تعداد	N	N

X = نمره یا نمره ها	X_c = حد میانی نمرات
L = کرانه طبقات	F_i = فراوانی مطلق
F_c = فراوانی تجمعی	x = انحراف نمره از میانگین
X^2 = مجذور (مربع نمره) انحراف نمره از میانگین	SS = مجموع مجذورات انحراف نمره از میانگین
MS = میانگین مجموع مجذور انحراف نمره از میانگین	\bar{i} = فاصله طبقاتی

حد واقعی اعداد

حدود واقعی: حدود واقعی نمرات بصورت کم کردن ۰/۵ نمره در اعداد صحیح و در اعداد اعشاری، نیم واحد یعنی ۰/۰۵ - ۰/۰۵ و مانند آن کسر می‌شود. یعنی: کرانه عدد ۲۵ ← ۲۴/۵ تا ۲۵/۵ و کرانه عدد ۲۵/۵ ← ۲۵/۴۵ تا ۲۵/۵۵ می‌باشد. مفهوم حدود واقعی مخصوصاً زمانی مفید است که اعداد گروه بندی یا طبقه بندی شوند. مثال: پس از اجرای یک آزمون ریاضی مشاهده می‌شود که ۱۰ نفر نمره ۱۲ گرفته‌اند. این بدان معنی نیست که همه توانایی یکسان دارند، بلکه دقیق نبودن وسیله اندازه‌گیری ممکن است موجب این امر شده باشد. به این خاطر نیاز به حدود واقعی می‌باشد؛ یعنی ۱۲/۵ - ۱۱/۵.

توزیع فراوانی

عبارتست از سازمان دادن اندازه‌ها یا مشاهدات به صورت طبقات همراه با فراوانی هر طبقه. توزیع فراوانی، داده‌ها را بصورت خلاصه و مرتب، به نحوی که تفسیر آنها آسان شود، نمایش می‌دهند.

مراحل ساخت جدول توزیع فراوانی

۱- مرتب کردن اعداد از کوچک به بزرگ یا برعکس.

۲- مشخص کردن تعداد دفعاتی که هر عدد تکرار شده است (تعداد فراوانی). زمانی که همه اعداد تک تک در جدول آورده شوند، جدول توزیع فراوانی منفرد یا طبقه بندی نشده گفته می شود. اما زمانی که نمره ها یا اعداد دارای دامنه گسترده ای هستند و تنظیم اعداد بصورت توزیع فراوانی طبقه بندی نشده وقتگیر و طاقت فرسا است، اعداد را طبقه بندی می کنیم و از جدول توزیع فراوانی طبقه بندی شده استفاده می کنیم. (در حال حاضر با وجود رایانه طبقه بندی اعداد غیر منطقی است) زیرا: در اثر طبقه بندی کردن اطلاعات برخی از اطلاعات از بین می رود. ستون داده ها (طبقات) را در جدول فراوانی با X نشان می دهند. فراوانی مطلق (f) برابر است با مقدار دفعات تکرار هر داده در هر طبقه.

مثال: در توزیع فراوانی درس آمار یک کلاس، نمرات به شرح زیر می باشد جدول فراوانی مربوط به توزیع را فراهم کنید؟

X	F
۱۵	۱
۱۲	۳
۱۱	۵
۱۰	۴

۱۱-۱۲-۱۱-۱۰-۱۲-۱۰-۱۲-۱۱-۱۰-۱۱-۱۲-۱۰

نکته: با توجه به جدول فوق، عدد ۵ در ستون f بیانگر این است که عدد ۱۱ پنج بار تکرار شده است. اگر داده های

ستون فراوانی (F) را با هم جمع کنیم تعداد کل داده ها بدست می آید. $N = \sum F$

یعنی در مثال فوق $N = ۱۳$

توزیع فراوانی طبقه بندی شده

زمانی که تعداد اعداد یک توزیع و همچنین فاصله بین آنها خیلی زیاد باشد، از توزیع فراوانی طبقه بندی شده استفاده می شود.

نکته: زمانی که تفاضل بین بزرگترین و کوچکترین نمره یا عدد مساوی یا بزرگتر از ۲۰ باشد از توزیع فراوانی طبقه بندی شده استفاده می شود.

طبقات بایستی ناسازگار باشند. یعنی یک عدد معین فقط در یک طبقه قرار داده شود. نمره های فردی در این نوع توزیع ها هویت خود را از دست می دهند.

نحوه ساختن توزیع فراوانی طبقه بندی شده

برای ساختن توزیع فراوانی طبقه بندی شده دو روش وجود دارد.

الف) روش یکم دارای مراحل زیر است:

$$R = X_H - X_L + ۱$$

۱- تعیین دامنه تغییرات

۲- تقسیم R به ترتیب بر ۱-۲-۳-۴ تا حاصل از ۲۰ کمتر شود. عدد بدست آمده تعداد طبقات و مقسوم علیه فاصله طبقات می باشد

۳- نوشتن طبقات: اولین عدد باید از کوچکترین عدد توزیع کوچکتر باشد و مضربی از فاصله طبقات نیز باشد.

۴- نوشتن فراوانی طبقات

ب) روش دوم دارای مراحل زیر است:

۱- تعیین دامنه تغییرات

$$K = 1 + \sqrt[3]{Log N}$$

۲- تعیین تعداد طبقات با استفاده از قانون استرژ

$$i = \frac{R}{K}$$

۳- تعیین اندازه یا حجم هر طبقه (فاصله طبقات)

۴- نوشتن طبقات

۵- نوشتن فراوانی طبقات

نکته: تعداد طبقات اختیاری است و معمولاً بین ۲۰-۱۰ است. و اگر تعداد طبقات بزرگتر از ۲۰ باشد، تهیه و تنظیم جدول نیاز به وقت و کار بیشتر دارد. اگر تعداد طبقات کوچکتر از ۱۰ باشد اندازه طبقات بزرگ می شود و اطلاعات بیشتری از دست می رود

نمایندگی طبقات (نقاط وسط طبقات): نماینده طبقات یا نقاط میانی را با X' نمایش می دهند و از طریق فرمول زیر به دست می آید:

$$X' = \frac{\text{حد بالای طبقه} + \text{حد پایین طبقه}}{2}$$

توزیع فراوانی تراکمی: اگر پژوهشگری علاقمند به دانستن تعداد افراد یا نمره هایی باشد که در پایین نمره یا عدد خاصی وجود دارند، نیاز به توزیع فراوانی تراکمی دارد. فراوانی تراکمی با (cf) نشان داده می شود که از جمع کردن فراوانی های ساده هر طبقه با طبقه بزرگتر به دست می آید.

نکته ۱: فراوانی تراکمی کوچکترین طبقه همیشه برابر با فراوانی ساده یا مطلق آن طبقه است.

نکته ۲: فراوانی تراکمی بزرگترین طبقه همیشه برابر با مجموع داده ها (ΣF) یا N می باشد.

$$\text{محاسبه ی فراوانی نسبی} = \frac{\text{فراوانی مطلق}}{\text{جمع کل فراوانی}}$$

درصد فراوانی مطلق و تراکمی

$$P = \frac{F}{N} \times 100 \quad P = \frac{\text{فراوانی مطلق هر طبقه}}{\text{تعداد کل فراوانی ها}} \times 100$$

$$CF\% = \frac{CF}{N} \times 100 \quad Cf\% = \frac{\text{فراوانی تراکمی هر طبقه}}{\text{تعداد کل فراوانی ها}} \times 100$$

نمودارها

نمودار دایره ای: نموداری است که با داده های اسمی و کیفی بکار می رود. برای محاسبه درجه و زاویه مرکزی متعلق به یک گروه از فرمول زیر استفاده می کنیم:

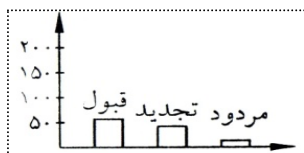
$$\text{درجه} = \frac{n}{N} \times 360$$

مثال: اگر داخل یک گروه ۱۲۰ نفری ۶۰ نفر قبول، ۴۰ نفر تجدید و ۲۰ نفر مردود شده باشند، درجه مربوط به تجدیدی‌ها چند است؟



$$\text{درجه تجدیدی} = \frac{40}{120} \times 360 = 120$$

نمودار ستونی (میله‌ای): نموداری است که با داده‌های اسمی بکار می‌رود که در محور عمودی، فراوانی و در محور افقی طبقات قرار می‌گیرند. در نمودار میله‌ای فاصله بین اعداد و نقاط یکسان و ثابت است. مثل تعداد فرزندان. ولی اگر برای متغیرهای طبقه‌ای مانند قبولی-مردودی، رنگ چشم و اعتقادات مذهبی نمودار رسم شود، هیچ ضرورتی ندارد که فاصله‌ها در محور X ثابت می‌ماند. برای این متغیرها از نمودار دایره‌ای و مانند آن استفاده می‌کنند (نمودار پای). مثال:

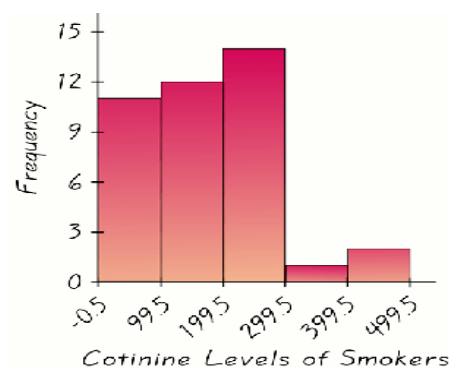


نمودار هیستوگرام: با داده‌های فاصله‌ای و نسبتی به کار می‌رود. نموداری شبیه نمودار ستونی است، ولی در آن ستونها به هم چسبیده است و در محور افقی (X) کرانه (حدود واقعی) طبقات و در محور عمودی (Y) فراوانی مطلق قرار می‌گیرد.

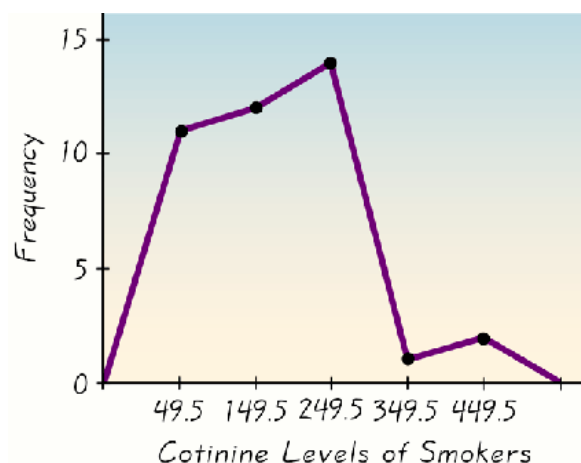
نکته: این نمودار، اغلب در مواردی استفاده می‌شود که بخواهیم فقط یک نمودار واحد را نشان دهیم، درحالی‌که در موقعیت‌های آزمایشی و شبه آزمایشی که مایلیم نمره‌های آزمودنی‌های دو گروه مجزا را با هم مقایسه کنیم ناگزیر هستیم از پلیگون استفاده کنیم.

مثال:

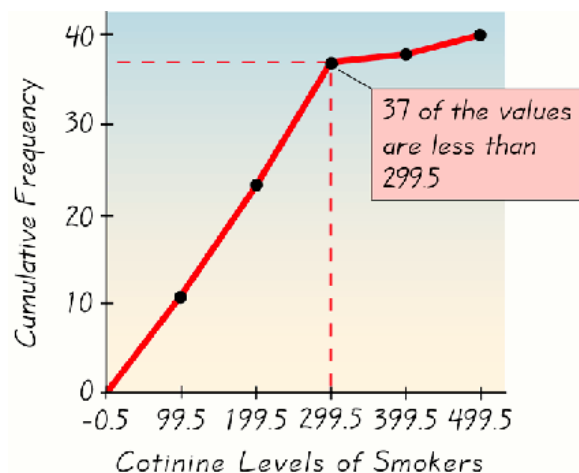
X	F_i
۴۰۰-۴۹۹	۲
۳۰۰-۳۹۹	۱
۲۰۰-۲۹۹	۱۴
۱۰۰-۱۹۹	۱۲
۰-۹۹	۱۱



نمودار چندضلعی (پلیگون): با داده‌های فاصله‌ای و نسبتی به کار می‌رود. نموداری شبیه نمودار هیستوگرام است که در محور X چند نمره میانی و در محور Y فراوانی مطلق قرار می‌گیرد.



نکته: برای مقایسه چند گروه در یک محور مختصات اگر تعداد آزمودنی‌ها در گروه‌ها برابر باشد، می‌توان از نمودار ستونی و یا هیستوگرام استفاده کرد، ولی اگر تعداد افراد گروه‌ها برابر نباشد باید از این توزیع استفاده کرد. نمودار چندضلعی تراکمی (اجایو): مانند نمودار چند ضلعی است که با داده‌های نسبتی و فاصله‌ای به‌کار می‌رود. در محور X کرانه طبقات و در محور Y فراوانی تجمعی قرار می‌گیرد. اگر به جای فراوانی تراکمی درصد فراوانی تراکمی را به صورت هندسی نمایش دهیم حاصل اجایو خوانده می‌شود. این نمودار وقتی مفید است که پژوهشگر علاقمند باشد وضعیت یک نمره یا یک فرد را نسبت به بقیه نمره‌ها یا افراد مشخص کند. برای ترسیم این نمودار در روی محور عرضی (Y) فراوانی تراکمی و در روی محور طولی (X) حدود واقعی طبقات قرار می‌گیرد. مثل نمودار چند ضلعی است.



شاخص‌های گرایش مرکزی نما (مد)

بی‌ثبات‌ترین شاخص گرایش مرکزی است که با داده‌های اسمی بکار می‌رود و عددی است که دارای بیشترین فراوانی می‌باشد. مثلاً: ۶-۵-۹-۶-۶-۶-۶-۶-۶-۶-۵-۶-۷-۵-۶-۷-۱. در اینجا نما عدد ۶ است. برخی موارد توزیع، دونمایی یا چندنمایی می‌شود مثل: ۱-۲-۳-۲-۴-۵-۵-۶-۷.

در اینجا ۲ و ۵ نما هستند که نما در اصل $\frac{2+5}{2} = 3.5$ است.

X	F _i
۵	۵
④	۹
۳	۳
۲	۲
۱	۶

در جدول فراوانی، نما عبارتست از عددی که F_i آن بزرگتر است.

در اینجا عدد ۴ نما است، زیرا عدد چهار، ۹ بار تکرار شده است.

نکته ۱: برای محاسبه نما باید ابتدا اعداد را از کوچک به بزرگ مرتب کنیم.

نکته ۲: در موقعیت هایی که دو عدد مجاور X فراوانی یکسان داشته باشند که بزرگتر از فراوانی سایر ارزش های X باشند، نما را می توان به طور قراردادی به عنوان میانگین دو ارزش مجاور X در نظر گرفت.

$$۱-۲-۳-۳-۴-۴-۵ \text{ میانه} = \frac{۳+۴}{۲} = ۳/۵$$

نکته ۳: در جایی که ارزش های غیر مجاور فراوانی های بزرگتر از فراوانی های طبقه های مجاور داشته باشند، هرکدام از این ارزش ها را می توان به عنوان نما در نظر گرفت. در چنین شرایطی توزیع دونمایی نامیده می شود.

۱۱-۱۱-۱۲-۱۲-۱۲-۱۳-۱۳-۱۳-۱۳-۱۴-۱۴-۱۴-۱۵-۱۵-۱۵-۱۵-۱۶-۱۶-۱۷-۱۷-۱۸

در اینجا عدد ۱۳ پنج بار تکرار شده به طوری که این فراوانی بزرگتر از فراوانی ارزش های مجاور است. همچنین عدد ۱۵ چهار بار تکرار شده و این فراوانی نیز بزرگتر از فراوانی ارزش های مجاور است. این مجموعه از مشاهدات دارای دونما است (فرگوسن، تاکانه، ۱۳۸۰).

محاسبه نما در داده ای طبقه بندی شده

از معادله زیر محاسبه می شود، که در آن:

L = حد واقعی پایین طبقه ای که دارای بیشترین فراوانی است

i = طول یا فاصله طبقات

d_1 = تفاوت فراوانی طبقه نمادار از فراوانی طبقه قبل

d_2 = تفاوت فراوانی طبقه نمادار از طبقه بعد.

مثال:

X	F
۵۴-۵۶	۱
۵۷-۵۹	۳
۶۰-۶۲	۶
۶۳-۶۵	۸
۶۶-۶۸	۲

طبقه نما ←

$$MO = L + i \left(\frac{d_1}{d_1 + d_2} \right)$$

$$d_1 = ۸ - ۶ = ۲$$

$$d_2 = ۸ - ۲ = ۶$$

$$MO = ۶۲/۵ + ۳ \left(\frac{۲}{۲+۶} \right) = ۶۳/۲۵$$

نکته: چنانچه توزیع نرمال باشد، نما از فرمول زیر محاسبه می شود: $MO = ۳Md - ۲\bar{X}$

میانۀ

میانۀ جایگاهی در توزیع نمره هاست و توزیع نمره ها را به دو قسمت مساوی تقسیم می کند؛ یعنی جایی است که دقیقاً ۵۰ درصد نمره ها بالای آن و ۵۰ درصد نمره ها زیر آن قرار می گیرند. میانۀ از نما باثبات تر و از میانگین بی ثبات تر است و با داده های رتبه ای بکار می رود. زیرا ما ابتدا نمره ها را از کوچک به بزرگ مرتب می کنیم.

طریقه محاسبه میانۀ نمرات خام در اعداد گسسته

ابتدا نمره ها را از کوچک به بزرگ مرتب می کنیم و سپس اگر تعداد اعداد فرد باشد، میانۀ، عدد وسط است.

مثال: ۵-۶-۱-۳-۴-۹-۱۲

مثال: ۱-۳-۴-۵-۶-۹-۱۲ در اینجا عدد ۵ میانۀ است

مثال: در توزیع اعداد ۱۲-۸-۱۷-۳۰-۳۰-۲۱-۵-۴-۳۱ میانۀ برابر است با ۱۷.

طریقه محاسبه میانۀ اعداد خام در تعداد زوج

۱- اعداد را از کوچک به بزرگ مرتب می کنیم.

۲- دو عدد وسط را با هم جمع و تقسیم بر دو می کنیم.

۸+۵ برابر ۱۳ و ۱۳ تقسیم بر ۲ میانۀ محاسبه می شود که برابر ۶/۵ می باشد.

نکته: هنگامی که نمره یا عددی که توزیع را به دو قسمت تقسیم می کند تکراری است، میانۀ از طریق محاسبه بدست می آید.

الف) ابتدا حد پایین عدد تکراری که میانۀ یکی از آنها است را می نویسیم.

ب) کسری را در نظر می گیریم که منخرج آن تعداد اعداد تکراری و صورت آن نشان دهنده تعداد اعداد تکراری است که در سمت چپ خط رسم کننده میانۀ قرار می گیرند.

ج) حاصل مراحل الف) و ب) را با هم جمع و میانۀ را بدست می آوریم. (میانۀ روی عدد ۴ می افتد، بنابراین نیمی از آن بعلاوه ی یک مورد ۴ قبل از آن جمعاً ۱/۵ تا از ۴ ها به حد پایین اضافه می شود) مانند:

$$۳-۴-۴-۵-۶ \quad \text{میانۀ} = ۳/۵ + \frac{۱/۵}{۲} = ۴/۲۵$$

$$۱-۲-۳-۳-۳-۴-۵-۶ \quad \text{میانۀ} = ۳/۵ + \frac{۲}{۳} = ۳/۱۷$$

محاسبه میانۀ در جدول اعداد طبقه بندی شده

۱- ابتدا نمرات را از کوچک به بزرگ مرتب می کنیم.

۲- سپس توزیع فراوانی و جدول را تشکیل و با فرمول نمرات طبقه بندی شده میانۀ را محاسبه می کنیم.

مثال: ۳-۵-۲-۴-۸-۴-۴-۴-۴-۲

X	F _i	F _c
۸	۱	۱۰
۵	۱	۹
۴	۵	۸
۳	۱	۳
۲	۲	۲
	ΣF=۱۰	

۸

$$L = \text{حد پایین طبقه میانه دار} = L + \frac{\frac{N}{2} - F_c}{F_i} \cdot i$$

طبقه ی میانه دار = مجموع فراوانی تقسیم بر دو = $\frac{N}{2}$

F_c = فراوانی تجمعی ماقبل سطری که در آن میانه واقع شده است.

F_i = فراوانی مطلق طبقه میانه دار

i = فاصله طبقاتی

$$F_c=3, \quad F_i=5, \quad i=1$$

مثال:

$$md = L + \frac{\frac{N}{2} - F_c}{F_i} \times i \Rightarrow md = 3/5 + \frac{5-3}{5} \times 1 \quad md = 3/9$$

نکته: در صورتی که داده در مقیاس فاصله ای یا نسبی باشند، بهترین شاخص، گرایش مرکزی میانگین است. ولی اگر در توزیعی که نمره ای در کرانه (نمره خیلی بزرگ یا خیلی کوچک) باشد (توزیع دارای کجی باشد) میانه شاخص مناسبتری است. به عنوان مثال در توزیع ۵-۶-۷-۹-۱۵-۳۰۰ میانه مناسبتر است.

ویژگی‌ها

۱- نسبت به اعداد بزرگ یا کوچک حساس نیست. بنابراین بهترین شاخص است که تمرکز اعداد را در وسط توزیع نشان می‌دهد. به عنوان مثال در توزیع های زیر که دارای میانگین های متفاوتی هستند، میانه برابر و مساوی ۲۰ می‌باشد.

$$25-24-20-7-5 \quad \text{و} \quad 63-52-20-15-10$$

۲- مورد استفاده میانه زمانی است که مقیاس اندازه گیری رتبه ای باشد، هرچند که می‌تواند برای داده‌هایی با مقیاس فاصله‌ای و نسبی هم استفاده شود.

۳- مجموع قدرمطلق انحرافهای نمره ها از میانه کوچکتر یا مساوی مجموع قدرمطلق انحرافهای نمره‌ها از هر عدد دیگری است (بدون در نظر گرفتن علامت).

$\sum X - Md \leq \sum X - C $					
نمره ها	قدر مطلق انحرافات				
	میانه (۶)	۴	۵	۷	۹
۴	۲	۰	۱	۳	۵
۵	۱	۱	۰	۲	۴
۶	۰	۲	۱	۱	۳
۷	۱	۳	۲	۰	۲
۹	۳	۵	۴	۲	۰
	۷	۱۱	۸	۸	۱۴

میانگین

میانگین، باثبات‌ترین شاخص گرایش مرکزی است که با داده‌های فاصله‌ای و نسبی بکار می‌رود و مرکز ثقل داده‌هاست.

طریقه محاسبه میانگین اعداد خام

مثال: ۱۱-۹-۱۰

$$\bar{X} = \frac{\sum X}{N} \Rightarrow \text{میانگین} = \frac{\text{مجموع نمره ها}}{\text{تعداد نمره ها}}$$

$$\bar{X} = \frac{10+9+11}{3} = 10 \quad \bar{X} = 10$$

طریقه محاسبه اعداد طبقه بندی شده با $i=1$

$$\bar{X} = \frac{\sum F.X}{n}$$

مثال:

X	Fi	F.X
۵	۳	۱۵
۴	۸	۳۲
۳	۵	۱۵
۲	۱۲	۲۴
$\sum F_i = n = 28$		$\sum F.X = 86$

$$\bar{X} = \frac{86}{28} = 3.07$$

طریقه محاسبه میانگین اعداد طبقه بندی شده با $i \neq 1$

$$\bar{X} = \frac{\sum F.X_c}{n}$$

مثال:

X	Fi	Xc	F.X
۱۷-۱۹	۵	۱۸	۹۰
۱۴-۱۶	۶	۱۵	۹۰
۱۱-۱۳	۱۱	۱۲	۱۳۲
۸-۱۰	۳	۹	۲۷
۵-۷	۹	۶	۵۴
$\sum F_i = 34$		$\sum F.X_c = 393$	


$$\bar{X} = \frac{393}{34} = 11.56 \quad \bar{X} = 11.56$$

حد میانی طبقه

X_c = برای پیدا کردن حد میانی نمرات، دو نمره طبقه را با هم جمع و تقسیم بر ۲ می‌نمائیم. $(7+5=12 \div 2 = 6)$

نکته

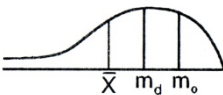
۱- اگر میانگین و میانه و نما با هم برابر باشد، توزیع متقارن است و کجی آن صفر است.

$$\bar{X} = md = m_o$$


۲- اگر میانگین بزرگتر از میانه و میانه بزرگتر از نما باشد کجی مثبت است، یعنی اکثر نمرات پایین بوده و امتحان سخت است.

$$\bar{X} > md > m_o$$

۳- اگر میانگین کوچکتر از میانه و میانه کوچکتر از نما باشد، کجی منفی است و اکثر نمرات بالاست و امتحان آسان بوده است.

$$\bar{X} < md < m_o$$


۴- کشیدگی^۱ ارتفاع (Y) منحنی را نشان می دهد. به برآمدگی یا خوابیدگی منحنی در مقایسه با منحنی طبیعی کشیدگی گویند.

۵- در صورتی که یک نمره ثابت (C) در کلیه نمرات یک توزیع ضرب، تقسیم، جمع و یا تفریق شود، در هر چهار حالت میانگین تغییر می کند، یعنی، میانگین قبلی در آن عدد ثابت ضرب، تقسیم، جمع و یا تفریق می شود.

$$\bar{X}_1 = \bar{X} - C \quad \bar{X}_1 = \bar{X} + C \quad \bar{X}_1 = \bar{X} : C \quad \bar{X}_1 = \bar{X} . C$$

۶- مجموع انحراف نمره ها از میانگین همیشه صفر است. $\sum (X - \bar{X}) = 0$ درواقع شاخصی از مکان مرکزی به معنای حداقل مجزورات است (فرگوسن، تاکانه، ۱۳۸۰: ۸۱).

۷- مجموع مجذور انحرافات از میانگین همیشه می نیمم است (کمتر از هر عدد دیگری است).

$$\sum (x - \mu)^2 < \sum (x - a)^2$$

رابطه شاخص های گرایش به مرکز

$$Mo = 3md - 2\bar{X} \quad md = \frac{mo + 2\bar{X}}{3} \quad \bar{X} = \frac{3md - mo}{2}$$

میانگین میانگین ها

$$\bar{X}_T = \frac{\sum \bar{X}_i}{N}$$

در صورتی که حجم نمونه ها برابر باشد یعنی $n_1 = n_2 = \dots = n_n$

مثال:

$$\bar{X}_1 = 5 \quad \bar{X}_2 = 10 \quad \bar{X}_3 = 15 \quad n_1 = 10 \quad n_2 = 10 \quad n_3 = 10$$

^۱. Kurtosis

خواهیم داشت:

$$\bar{X}_T = \frac{\bar{X}_1 + \bar{X}_2 + \bar{X}_3}{N} = \frac{5 + 10 + 15}{3} = 10 \quad \bar{X}_T = 10$$

در صورتی که حجم نمونه ها برابر نباشد به طریق زیر محاسبه می شود:

$$\bar{X}_T = \frac{\sum \bar{X}_i . n_i}{n_1 + n_2 + n_3} \Rightarrow \frac{\bar{X}_1 . n_1 + \bar{X}_2 . n_2 + \bar{X}_3 . n_3}{n_1 + n_2 + n_3}$$

مثال:

$$\bar{X}_1 = 5 \quad \bar{X}_2 = 10 \quad \bar{X}_3 = 15 \quad n_1 = 15 \quad n_2 = 10 \quad n_3 = 20$$

$$\bar{X}_T = \frac{(5 \times 15) + (10 \times 10) + (15 \times 20)}{15 + 10 + 20} = 10.55 \quad \bar{X}_T = 10.55$$

میانگین هارمونیک (همساز)

این نوع میانگین برای مواردی بکار می رود که مقیاس ترکیبی باشد مانند متر در ثانیه و کیلومتر بر ساعت. این میانگین در عینک سازی و مطالعه شبکه های برقی به کار می رود.

$$HM = \frac{1}{\frac{1}{N} \left[\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n} \right]}$$

مثال: میانگین سرعت های ۵ و ۶ و ۷ و ۲ کیلومتر در ساعت ۴ ماشین چند است؟

$$HM = \frac{1}{\frac{1}{4} \times \left[\frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{2} \right]} = 3.96$$

میانگین هندسی

نوعی میانگین است که با G نشان داده می شود و معمولاً هرگاه x_i ها از درصدها یا نسبت ها تشکیل شده باشند، استفاده می شود. در کارهای اقتصادی یا جمعیت شناسی به کار می رود.

$$G = \sqrt[n]{X_1 \cdot X_2 \cdot X_3 \cdot (X_n)}$$

مثال: میزان سود شرکت مهرداد در ۵ سال گذشته برحسب درصد به ترتیب ۲، ۳، ۴، ۴، ۳ شده است. کدام یک از گزینه های زیر وضع سود آوری شرکت را بهتر نشان می دهد.

$$G = \sqrt[5]{3 \times 4 \times 4 \times 3 \times 2} = 3.01$$

تفسیر: اعدادی مانند ۳ و ۶ را می توان به عنوان فاصله های بین نقاط نشان داد. حاصل ضرب این دو عدد را نیز می توان به عنوان یک مساحت نشان داد. بنابراین مستطیلی که عرض آن ۸ و طول آن ۱۵ سانتی متر است، مساحتی برابر ۱۲۰ سانتی متر مربع دارد. ریشه دوم ۱۲۰ یا ۱۰/۹۵ بعدی از یک مربع است که مساحت آن با مساحت مستطیل ۸ ضرب در ۱۵ برابر است. به همین ترتیب حاصل ضرب سه عدد، مبین یک حجم است و ریشه سوم حاصل ضرب آنها بعدی از یک مکعب با همان حجم است (فرگوسن، تاکانه، ۱۳۸۰: ۸۱).

رابطه بین سه میانگین فوق به صورت زیر است:

$$\bar{X} > G > HM$$

میانگین هارمونیک > میانگین هندسی > میانگین وزنی

نکته: اگر بین داده ها عدد صفر یا منفی وجود داشته باشد نمی توان از این روش استفاده کرد.

شاخص های پراکندگی

فرض کنید گفته شود درجه حرارت دما در مکانی ۶۵ درجه فارنهایت است. اگر صرفاً با نگاه کردن به میانگین نتیجه بگیریم که این مکان برای زندگی مناسب است، خطا کرده ایم. زیرا این عدد مقدار متوسط ماه های زمستان و تابستان است. ما نیاز داریم تغییرپذیری درجات حرارت را بدانیم. یعنی درجه حرارت روزانه چقدر از متوسط ۶۵ درجه فارنهایت متفاوت است. از این رو باید شاخص های پراکندگی که نشان می دهد داده ها از میانگین چقدر دور و یا به آن نزدیک هستند را محاسبه کنیم. شاخص های پراکندگی، میزان پراکنده بودن نمرات حول و حوش مرکز داده ها را نشان می دهد که به ترتیب عبارتند از: دامنه تغییرات- انحراف متوسط، انحراف چارکی، واریانس و انحراف معیار.

دامنه تغییرات (R)

دامنه تغییرات یک شاخص پراکندگی و درواقع بی ثبات ترین شاخص و حساس ترین شاخص پراکندگی است که با داده های فاصله ای بایستی بکار رود و تفاضل بین کوچکترین و بزرگترین عدد در توزیع است (بدون در نظر گرفتن حدود واقعی اعداد)

$$R = \max - \min$$

$$R = \max - \min + 1$$

مثال: در توزیع نمرات ۳-۲۰-۱۴-۱۵-۹-۶-۵ دامنه تغییرات را محاسبه کنید:

$$R = 20 - 3 + 1 = 18$$

$$R = 18$$

نکته

- ۱- دامنه تغییرات، برای نمونه های بزرگ شاخص بی ثباتی است.
- ۲- واریانس نمونه برداری دامنه تغییرات برای نمونه های کوچک، بزرگتر از واریانس نمونه برداری انحراف معیار نیست. ولی سریعاً با زیاد شدن N افزایش می یابد.
- ۳- دامنه تغییرات به جز در موارد خاص مستقل از حجم نمونه نیست.
- ۴- دامنه تغییرات برای نمونه های کوچک مناسب است.

انحراف چارکی

بیشتر در مواردی کاربرد دارد که نمره ها دارای مقیاس رتبه ای هستند و یا نمره ای در کرانه باشد. پراکندگی را در اطراف مرکز توزیع نشان می دهد. در داده های پرت به جای انحراف استاندارد استفاده می شود.

مثال: برای نمرات ۳۲-۶۰-۱۵-۹-۶-۵-۳ انحراف چارکی مناسب

$$Q_1 \quad md \quad Q_3$$

است.

$$SIRQ = \frac{Q_3 - Q_1}{2} = \frac{\text{چارک اول} - \text{چارک سوم}}{2}$$

$$Q = \frac{32 - 5}{2} = 13.5$$

برای محاسبه چارک متوسط اعداد طبقه بندی شده نیاز به محاسبه چارک اول و سوم داریم. شرایط استفاده از چارک متوسط (انحراف چارکی) مانند میانه است. بهترین مورد استفاده از انحراف چارکی هنگامی است که چولگی شدید است، زیرا تحت تأثیر نمرات بالا و پایین قرار نمی گیرد. این شاخص از دامنه تغییرات کوچکتر است.

محاسبه چارک ها در اعداد طبقه بندی شده

۱- اعداد را از کوچک به بزرگ مرتب کنید.

۲- میانه را حساب کنید. Q_2

سؤال: میانه اعداد سمت چپ و راست را محاسبه کنید. Q_1 و Q_3 مثال:

$$6 - 8 - 9 - 11 - 12 - 13 - 14 - 17$$

$$Q_1 = 8.5 \quad Q_2 = 11.5 \quad Q_3 = 13.5$$

نکته ۱

چارک یکم: یعنی نقطه ای که یک چهارم افراد زیر آن و سه چهارم بالای آن قرار دارند.

دهک یکم: نمره های یک دهم افراد زیر آن و ۹/۰ افراد از آن بزرگتر است.

صدک یکم: نمره های یک صدم افراد زیر آن و ۹۹/۰ افراد بالای آن قرار دارد.

نکته ۲: پس از محاسبه نمره چارک اول و سوم انحراف چارکی بدست می آید در صورتی که:

$$Q_3 - Q_2 < Q_2 - Q_1 \quad \text{کجی منفی} \quad Q_3 - Q_2 > Q_2 - Q_1 \quad \text{کجی مثبت} \quad Q_3 - Q_2 = Q_2 - Q_1 \quad \text{مقارن}$$

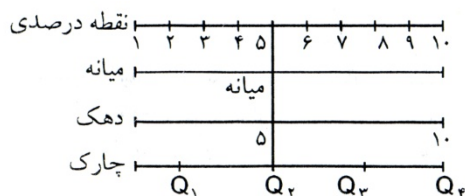
نقطه درصدی P_n

نقاط درصدی نمره را بر مقياس صد نشان می دهد. یعنی نقطه ۵۰ درصدی برابر میانه است و نقطه ۲۵ درصدی برابر چارک اول است و نقطه ۷۵ درصدی برابر چارک سوم است.

$$P_{.75} = Q_3 \quad P_{.50} = D_1 \quad P_{.25} = D_2 \quad P_{.50} = md = Q_2 \quad P_{.25} = Q_1$$

دهک هم مثل نقطه درصدی، موقعیت نمره را در مقياس ۱۰ تایی نشان می دهد و چارک هم موقعیت نمره را در

مقياس ۲۵ تایی در چهار موقعیت نشان می دهد، یعنی:



فرمول محاسبه میانه و دیگر نقاط

$$Q_r = L + \left(\frac{\frac{rN}{2} - F_c}{F_i} \right) i \quad Q_1 = L + \left(\frac{\frac{1N}{2} - F_c}{F_i} \right) i \quad md = L + \left(\frac{\frac{N}{2} - F_c}{F_i} \right) i$$

$$P_{i,} = L + \left(\frac{\frac{i \cdot N}{100} - F_c}{F_i} \right) i \quad D_1 = L + \left(\frac{\frac{1N}{2} - F_c}{F_i} \right) i \quad D_1 = L + \left(\frac{\frac{1N}{2} - F_c}{F_i} \right) i$$

مثال: چارک سوم، چارک اول انحراف چارکی دهک دوم، صدک شصت و هشتم را در داده های زیر محاسبه کنید.

X	F_i	F_c
۲۵-۲۹	۱۷	۱۰۰
۲۰-۲۴	۲۹	۸۳
۱۵-۱۹	۲۱	۵۴
۱۰-۱۴	۱۸	۳۳
۵-۹	۱۵	۱۵
	$\sum F_i = 100$	

$$Q_r = L + \frac{\frac{rN}{2} - F_c}{F_i} \cdot i \quad N = \frac{3 \times 100}{2} = 150$$

$$L = 19/5 \quad F_c = 54$$

$$F_i = 29 \quad i = 5$$

$$Q_r = 19/5 + \left(\frac{75 - 54}{29} \right) \times 5 = 23/12$$

$$Q_1 = L + \left(\frac{\frac{N}{2} - F_c}{F_i} \right) i \quad F_c = 15 \quad L = 9/5$$

$$Q_1 = 9/5 + \left(\frac{75 - 15}{18} \right) \times 5 = 12/28$$

$$Q = \frac{Q_r - Q_1}{2} = \frac{23/12 - 12/28}{2} = 5/42 \quad (\text{انحراف چارکی})$$

$$D_r = L + \left(\frac{\frac{rN}{2} - F_c}{F_i} \right) i \quad L = 9/5 \quad F_c = 15$$

$$D_r = 9/5 + \left(\frac{75 - 15}{18} \right) \times 5 \Rightarrow D_r = 10/89$$

$$P_{N(0.78)} = L + \left(\frac{\frac{78N}{100} - F_c}{F_i} \right) i \quad F_c = 54 \quad L = 19/5$$

$$I = 5 \quad F_i = 29 \quad \frac{0.78 \times 100}{100} = 78$$

$$P_{.78} = 19/5 + \left(\frac{78 - 54}{29} \right) \times 5 \Rightarrow P_{.78} = 21/91$$

رتبه درصدی

رتبه درصدی مثل نمره های استاندارد است که موقعیت نسبی فرد را در داخل گروه نشان می دهد یعنی با داشتن رتبه درصدی فرد می توانیم بگوییم که او از چند درصد گروه بهتر و یا بدتر عمل کرده است.

$$P_R = \left(\frac{F_c + \frac{F_i}{2}}{n} \right) \times 100$$

نکته: رتبه درصدی موقعیت فرد را در گروه و نقطه درصدی موقعیت نمره فرد را در داخل نمره ها (X)، نشان می دهد. بطور مثال کسی که با نمره ۱۸، رتبه درصدی ۸۵ کسب کند، یعنی نمره ۱۸ (نقطه درصدی) از ۸۵ درصد افراد (رتبه درصدی) بهتر و از ۱۵ درصد افراد بدتر عمل کرده است.

مثال: اگر بخواهیم رتبه درصدی عدد ۲۲ را در جدول فوق مشخص کنیم، خواهیم داشت.

$$P_R = \left(\frac{F_c + \frac{F_i}{2}}{n} \right) \times 100 \quad P_R = \left(\frac{54 + \frac{29}{2}}{100} \right) \times 100 = 67.5 \quad P_R = 68/5$$

نکته ۱: در توزیع های نامتقارن اغلب از میانه به عنوان شاخص مرکزی و انحراف چارکی به عنوان شاخص پراکندگی استفاده می شود.

نکته ۲: اگر انحراف چارکی اندازه ها برابر صفر باشد، ۵۰ درصد اندازه هایی که در وسط قرار گرفته اند، با هم برابرند. در نتیجه چارک های یکم و دوم و سوم با هم برابرند و برعکس.

پرسش: کدام پراکندگی برای توزیع فراوانی زیر مناسب تر است؟

حدود طبقات	۲۰-۵۰	۲۰-۴۰	۴۰-۶۰	۶۰ و بیشتر
فراوانی	۲۵	۳۵	۳۰	۱۰

الف) انحراف چارکی ب) انحراف معیار ج) ضریب تغییرات د) انحراف متوسط
پاسخ: گزینه (الف) صحیح است.

انحراف متوسط (MD)

انحراف متوسط، یک شاخص پراکندگی است که به آن میانگین قدر مطلق انحراف نمره از میانگین گفته می شود.

$$MD = \frac{\sum |X - \bar{X}|}{N}$$

مثال: برای اعداد ۱-۲-۳-۴-۵ خواهیم داشت:

X	X - \bar{X}	X - \bar{X}	
۵	۲	۲	
۴	۱	۱	
۳	۰	۰	
۲	-۱	۱	
۱	-۲	۲	
		$\sum X - \bar{X} = 6$	

$$md = \frac{6}{5} = 1.2$$

نکات مهم

- ۱- با این شاخص عملیات جبری را نمی‌توان انجام داد.
- ۲- در انحراف متوسط علائم اعداد و در انحراف چارکی کلیه اعداد مورد مطالعه قرار نمی‌گیرند.
- ۳- تأثیر انحرافات بزرگ را در شرایطی که تعداد زیادی انحرافات کوچک در برابر تعداد کمی انحراف بزرگ باشد، نشان نمی‌دهد (مهمترین کاستی).
- ۴- انحراف متوسط اعداد ثابت صفر است.

واریانس (s^2)

یک شاخص پراکندگی است که میانگین مجموع مجذورات انحراف نمره از میانگین می‌باشد؛ یعنی:

$$S^2 = \frac{\sum (X - \bar{X})^2}{n} = \frac{\sum (X)^2}{n} - \frac{ss}{n} = \frac{\sum X^2 - \frac{(\sum X)^2}{n}}{n}$$

$$S^2 = \frac{\sum (X - \bar{X})^2}{n}$$

محاسبه واریانس از راه انحراف نمره از میانگین

مثال: برای اعداد ۱-۲-۳-۴-۵

X	$X - \bar{X}$	X^2
۵	۲	۴
۴	۱	۱
۳	۰	۰
۲	-۱	۱
۱	-۲	۴
$\sum X - \bar{X} = ۰$		$\sum X^2 = ۱۰$

$$S^2 = \frac{۱۰}{۵} = ۲$$

$$S^2 = \frac{\sum X^2 - \frac{(\sum X)^2}{n}}{n}$$

محاسبه واریانس از راه نمرات خام

نکته: هنگامی که واریانس را در نمونه محاسبه می‌کنیم در مخرج کسر $n-1$ که به آن درجه آزادی (df) گفته می‌شود استفاده می‌کنیم.

درجه آزادی: تعداد مشاهداتی است که می‌تواند آزادانه تغییر کند. بنابراین درجه آزادی واریانس $df = n-1$ می‌باشد.

X	X^2
۵	۲۵
۴	۱۶
۳	۹
۲	۴
۱	۱
$\sum X = ۱۵$	$\sum X^2 = ۵۵$

نکته مهم: واریانس ۵ عدد متوالی در نمونه برابر ۲/۵ و

واریانس ۵ عدد متوالی در جامعه برابر ۲ می‌باشد.

$$S^2 = \frac{۵۵ - \frac{(۱۵)^2}{۵}}{۵} = ۲$$

انحراف معیار (انحراف استاندارد) S

جذر واریانس را انحراف معیار می گویند، یعنی در مثال فوق:

$$S = \sqrt{2} = 1/41$$

نکته: توان دوم انحراف معیار، واریانس نام دارد. انحراف معیار، باثبات ترین شاخص پراکندگی است. هرچقدر انحراف استاندارد بیشتر باشد، پراکندگی بیشتر است. شاید اساسی ترین فایده انحراف استاندارد این باشد که با استفاده از آن می توان مشخص کرد چه نسبتی از نمره ها در فاصله های مختلف نسبت به میانگین قرار گرفته است.

نکته مهم: زمانی که توزیع دارای کجی زیاد باشد از انحراف استاندارد باید با احتیاط استفاده شود. زمانی از S استفاده می شود که میانگین به عنوان شاخص مرکزی مورد استفاده قرار می گیرد. کلیه شاخص های پراکندگی (دامنه تغییرات، انحراف متوسط، انحراف معیار) با مقیاس حداقل فاصله ای بکار می روند.

ثبات پراکندگی به ترتیب

کمتر

بیشتر

انحراف معیار ← واریانس ← انحراف متوسط ← چارک متوسط ← دامنه تغییرات

$$R \leftarrow Q \leftarrow MD \leftarrow S^2 \leftarrow S$$

$$S > MD > Q$$

نکته مهم

تأثیر چهار عمل اصلی در واریانس

در شرایط جمع و تفریق یک عدد ثابت در نمرات یک توزیع واریانس تغییر نمی کند.

$$S_1^2 = S^2 \quad \text{تفریق} \quad S_1^2 = S^2 \quad \text{جمع}$$

ولی در شرایط ضرب و تقسیم، واریانس قدیم در توان دوم مجذور عدد ثابت ضرب یا تقسیم می شود.

$$S_1^2 = S^2 : C^2 \quad \text{تقسیم} \quad S_1^2 = S^2 \cdot C^2 \quad \text{ضرب}$$

مثال: اگر واریانس یک توزیع ۲ باشد و عدد ۳ در کلیه نمرات توزیع ضرب شود واریانس جدید برابر خواهد بود با:

$$S_1^2 = 2 \times 3^2 = 18$$

اگر واریانس ۵ و عدد ثابت ۶ با کلیه نمرات جمع شود، واریانس جدید برابر است با: $S_1^2 = 5$ (تغییری نمی کند).

تأثیر چهار عمل اصلی در انحراف معیار

در شرایط جمع و تفریق مثل واریانس است؛ یعنی، انحراف معیار تغییری نمی کند، ولی در شرایط ضرب و تقسیم

مثل میانگین است؛ یعنی، در همان عدد ثابت ضرب یا تقسیم می شود.

$$S_1 = S \quad \text{تفریق} \quad S_1 = S \quad \text{جمع}$$

$$S_1 = S : C \quad \text{تقسیم} \quad S_1 = S \cdot C \quad \text{ضرب}$$

مثال: اگر انحراف معیار یک توزیع ۶ و یک نمره ثابت ۲ در کلیه نمرات توزیع ضرب شود، انحراف معیار جدید

$$S_1 = 6 \times 2 = 12$$

برابر است با:

اگر انحراف معیار یک توزیع ۶ و نمره ثابت ۳ با کلیه نمرات جمع شود انحراف معیار جدید برابر است با:

$$S = 6 \quad \text{(تغییری نمی کند).}$$

نکات مهم

- ۱- میانه و نما در جمع و تفریق و ضرب و تقسیم یک عدد ثابت تابع میانگین بوده و مثل هم تغییر می کنند.
- ۲- در جمع و تفریق، واریانس - انحراف معیار - انحراف متوسط - دامنه تغییرات و انحراف چارکی مثل هم بوده و تغییر نمی کنند.
- ۳- در شریط ضرب و تقسیم، انحراف معیار - انحراف متوسط - انحراف چارکی و دامنه تغییرات در همان عدد ضرب یا تقسیم می شوند.
- ۴- واریانس در مجذور آن عدد ضرب یا تقسیم می شود.
- ۵- واریانس اعداد ثابت صفر است.
- ۶- اگر چند جامعه با هم ترکیب شوند، میانگین و واریانس جامعه کل ترکیبی، بزرگتر از میانگین و واریانس جوامع تشکیل دهنده خواهد بود. مگر آن که میانگین جوامع برابر باشد که در آن صورت واریانس آنها هم برابر است.
- ۷- واریانس تفاوت ها در نمونه های همبسته برابر است با:

$$S_{y_1}^2 - S_{y_2}^2 = S_{y_1}^2 + S_{y_2}^2 - 2\rho S_{y_1} S_{y_2}$$

- ۸- واریانس تفاوت ها در نمونه های مستقل برابر است با:

$$S_{x_1-x_2}^2 = S_{x_1}^2 + S_{x_2}^2 = \frac{S_1^2}{N_1} + \frac{S_2^2}{N_2}$$

- ۹- واریانس مجموع برابر است با:

$$S_{(X+Y)}^2 = S_X^2 + S_Y^2 \pm 2Cov(X, Y) = (S_X^2 + S_Y^2) \pm 2r S_X S_Y$$

تصحیح شپرد

- فرمول شپرد برای تصحیح انحراف معیار زمانی کاربرد دارد که فاصله طبقاتی بزرگ و تعداد طبقات کمتر از ۱۲ باشد.

$$S_c = \sqrt{S^2 - \frac{i^2}{12}}$$

مثال: اگر فاصله طبقاتی ۱۰ و انحراف معیار توزیع ۵ باشد، انحراف معیار تصحیح شده برابر خواهد بود با:

$$S_c = \sqrt{25 - \frac{10^2}{12}} = 4.97$$

ضریب پراکندگی (V)

- همان ضریب نسبی واریانس است که بر اساس آن پراکندگی ویژگی یک گروه را با پراکندگی ویژگی دیگر همان گروه مقایسه می کنند.

$$V = \frac{S}{\bar{X}} \times 100 = \frac{\text{انحراف معیار}}{\text{میانگین}} \times 100$$

مثال: کارخانه ای دو نوع لاستیک اتومبیل تولید می کند. برای نوع الف میانگین عمر ۱۰۰۰۰ کیلومتر، با انحراف استاندارد ۲۰۰۰ کیلومتر، و برای نوع ب میانگین عمر ۱۱۰۰۰ کیلومتر با انحراف استاندارد ۱۰۰۰ کیلومتر می باشد. کدام نوع لاستیک بهتر است؟

$$\text{الف) } V = \frac{2000}{10000} \times 100 = 20 \quad \text{ب) } V = \frac{1000}{11000} \times 100 = 9$$

پاسخ: نوع (ب) بهتر است، زیرا هم میانگین عمر آن بیشتر است و هم ضریب تغییر آن کوچکتر.

موارد استفاده

زمانی که دو یا چند جامعه در مقایسه با هم دارای مشاهدات ناهمگون از نظر واحد اندازه گیری باشند، مانند یک جامعه برحسب متر و یک جامعه برحسب اینچ، و یا چند جامعه دارای میانگین های متفاوتی باشند، استفاده می شود. گاهی نیز مقیاس صفت مورد اندازه گیری در دو جامعه یکسان است ولی بزرگی مشاهدات آنها به طور قابل ملاحظه ای تفاوت دارد. مانند مقایسه پراکندگی سود و زیان در صنایع دستی با صنایع سنگین.

نکته ۱: ضریب تغییرات اعداد ثابت برابر صفر است.

پرسش: برای تعیین آنکه در ۳۰ روز گذشته به نسبت، قیمت دلار از ثبات بیشتری برخوردار بوده است یا یورو، استفاده از کدام شاخص آماری مناسب تر است؟

الف) انحراف متوسط ب) ضریب پراکندگی ج) ضریب چولگی د) واریانس

پاسخ: گزینه (ب) صحیح است.

گشتاورهای پیرامون میانگین

$$m_1 = \frac{\sum (X - \bar{X})}{N} = 0 \quad \text{گشتاور اول:}$$

$$m_2 = \frac{\sum (X - \bar{X})^2}{n} = S^2 \quad \text{گشتاور دوم:}$$

$$m_3 = \frac{\sum (X - \bar{X})^3}{n} \quad \text{گشتاور سوم:}$$

$$m_4 = \frac{\sum (X - \bar{X})^4}{n} \quad \text{گشتاور چهارم:}$$

گشتاور اول همیشه صفر است. گشتاور دوم همان واریانس است. گشتاور سوم برای محاسبه چولگی (کجی) بکار

$$SK = \frac{m_3}{m_2 \sqrt{m_2}} \quad \text{می رود، یعنی:}$$

$$Kg = \frac{m_4}{(m_2)^2} - 3 \quad \text{گشتاور چهارم برای محاسبه کشیدگی بکار می رود، یعنی:}$$

نکته ۲: کشیدگی معیاری است بدون واحد که ارتفاع را نشان می دهد و رابطه معکوس با پراکندگی دارد.

اگر $k=0$ باشد، کشیدگی هم اندازه و هم ارتفاع توزیع نرمال است، اگر $k > 0$ باشد، از نرمال بزرگتر و پراکندگی کمتر و اگر $k < 0$ باشد، از نرمال کوتاهتر و پراکندگی بیشتر است.

$$g_1 = \frac{\bar{X} - m}{S}$$

فرمول کجی پیرسون

نکته

- ۱- اگر $|sk| \leq 0.1$ باشد، تقریباً کجی وجود ندارد و جامعه نرمال است.
- ۲- اگر $0.1 \leq |sk| \leq 0.5$ باشد، چولگی موجود اندک ولی غیر قابل اغماض است. درحقیقت جامعه از نظر تقارن اندکی با توزیع نرمال متفاوت است.
- ۳- اگر $|sk| > 0.5$ باشد، چولگی زیاد و غیر قابل اغماض است. به عبارت دیگر جامعه از نظر قرینگی دارای تفاوت فاحشی با توزیع نرمال است.

نمره های استاندارد

نمره های استاندارد موقعیت فرد را در گروه معین می کنند. با داده های مقیاس حداقل فاصله ای کاربرد دارند. در تبدیل نمرات خام به نمره استاندارد از فرمول مقابل استفاده می کنیم:

$$Z = \frac{X - \bar{X}}{S}$$

مثال: در امتحانی که میانگین آن ۳۸ و انحراف معیار آن ۳ باشد، کسی که نمره ۴۴ گرفته است دارای $Z = 2$ می باشد.

$$Z = \frac{44 - 38}{3} = \frac{6}{3} = 2$$

در تبدیل نمره Z به دیگر نمرات استاندارد از فرمول مقابل استفاده می کنیم:

$$X = \bar{X} - Z \frac{s}{\bar{s}}$$

$$t = 50 + 2 \times 10 = 70$$

مثال: کسی که نمره او $Z = 2$ باشد، نمره t او برابر خواهد بود با:

انحراف استاندارد	میانگین	نمره های استاندارد
۱	۰	Z
۱۰	۵۰	T
۱۰۰	۵۰۰	CEEB
۲	۵	نه گانه
۱۵	۱۰۰	هوش وکسلر
۱۶	۱۰۰	هوش بینه
۲۰	۱۰۰	AGCT ارتش آمریکا

سپس در تبدیل نمرات به نمرات استاندارد دیگر به روش زیر عمل می کنیم:

$$CEEB = 500 + 100Z$$

$$t = 50 + 10Z$$

$$100 + 15Z = \text{هوش وکسلر}$$

$$AGCT = 100 + 20Z$$

$$100 + 16Z = \text{هوش بینه}$$

$$5 + 2Z = \text{نه گانه}$$

مثال:

۱- در آزمون CEEB فردی ۸۰۰ گرفته است. نمره Z او چقدر است؟

$$Z = \frac{X - \bar{X}}{S} = \frac{800 - 500}{100} = 3$$

۲- فردی که در یک آزمون نمره نه گانه او ۳ شده است، نمره Z او چقدر است؟

$$\text{stanine} = 5 + 2Z \rightarrow Z = \frac{3 - 5}{2} = -1$$

۳- نمره IQ فردی ۱۴۰ شده است. نمره Z او چقدر است؟

$$Z = \frac{140 - 100}{15} = 2.67$$

۴- نمره t فردی ۳۰ شده است. نمره Z او چقدر است؟

$$Z = \frac{30 - 50}{10} = -2$$

$$t = 50 + 10Z$$

۵- نمره Z فردی ۲ شده است. نمره CEEB او چقدر است؟

$$X = \bar{X} + Z.S \Rightarrow 500 + (2 \times 100) = 700$$

۶- نمره Z فردی ۱/۵ می باشد. نمره نه گانه او چقدر است؟

$$X = 5 + (1/5 \times 2) = 8$$

۷- نمره Z فردی ۲- است. نمره AGCT او چقدر است؟

$$X = 100 + (-2) \times 20 = 60$$

۸- نمره Z فردی ۱/۵- شده است. نمره t او چقدر است؟

$$t = 50 + (-1/5) \times 10 = 35$$

۹- نمره CEEB فردی ۳۵۰ شده است. نمره هوشبهر او چقدر است؟

$$Z = \frac{X - \bar{X}}{S} \Rightarrow \frac{350 - 500}{100} = -1.5$$

$$X = \bar{X} + ZS \Rightarrow 100 + (-1.5) \times 15 = 77.5$$

۱۰- نمره نه گانه فردی ۸ شده است، نمره t وی چقدر است؟

$$Z = \frac{8 - 5}{2} = 1.5$$

$$t = 50 + 2/5 \times 10 = 75$$

۱۱- نمره IQ فردی ۴۵ شده است. نمره t او چقدر است؟

$$Z = \frac{45 - 100}{15} = -3.67$$

$$t = 50 + (-3.67) \times 10 = 13$$

منحنی طبیعی و سطوح زیر منحنی

منحنی طبیعی (زنگوله ای یا گوس) یک توزیع طبیعی است.

۱- شکل آن به میانگین و انحراف استاندارد بستگی دارد.

۲- میانگین آن صفر و انحراف استاندارد آن یک می باشد.

۳- در عمل یک منحنی طبیعی داریم و بالاترین ارتفاع در میانگین است.

۴- در منحنی طبیعی میانگین و میانه و نما با هم برابر هستند.

به طوریکه می بینیم مشخص می شود که $34/13$ درصد افراد بین $+1$ تا \bar{X}

قرار می گیرند. همینطور $13/59$ درصد افراد بین $+2$ تا $+1$ قرار می گیرند. و

$2/14$ درصد افراد بین $+3$ تا $+2$ و $0/14$ درصد افراد بین $+\infty$ تا $+3$ قرار

می گیرند. این درصدها برای طرف چپ منحنی به دلیل تقارن هم صادق



است. اگر بخواهیم رتبه درصدی افراد را بر اساس منحنی نرمال محاسبه کنیم به روش زیر عمل می کنیم:

اگر داده (Z) یک عدد باشد توجه می کنیم مثبت است یا منفی، در صورت مثبت بودن با رتبه درصدی اولیه 50

جمع می کنیم و اگر منفی بود از رتبه درصدی 50 کم می شود.

رتبه درصدی $Z=1$ چقدر است؟ $50 + 34/13 \Rightarrow 84/13$ درصد

رتبه درصدی $Z=-1$ چقدر است؟ $50 - 34/13 = 15/87$

در صورتیکه سطوح بین دو نمره Z خواسته شود دیگر 50 کاربردی ندارد و سطوح زیر منحنی محاسبه می شود.

مثال: سطوح زیر منحنی بین ± 1 چقدر است؟ $34/13 + 34/13 = 68/26$ درصد

سطوح دیگر عبارتند از:

\bar{X} تا $+3 \Rightarrow 49/86$ درصد

\bar{X} تا $+2 \Rightarrow 47/72$ درصد

\bar{X} تا $-3 \Rightarrow 49/86$ درصد

\bar{X} تا $-2 \Rightarrow 47/72$ درصد

سطوح بین مثبت و منفی

-2 تا $+2 \Rightarrow 95/44$ درصد -3 تا $+3 \Rightarrow 99/77$ درصد $+1$ تا $-2 \Rightarrow 81/85$ درصد

در صورتیکه بصورت اعشاری باشد با روش زیر عمل می کنیم:

رتبه درصدی $Z=1/5$ چقدر است؟

$$50 + 34/13 + 13/59 \times \frac{5}{10} = 90/93$$

رتبه درصدی $Z=-2/5$ چقدر است؟

$$50 - \left[34/13 + 13/59 + 2/14 \times \frac{25}{100} \right] = 1/75$$

رتبه درصدی $Z=1/2$ چقدر است؟

$$50 + 34/13 + 13/59 \times \frac{2}{100} = 84/40$$

سطوح زیر منحنی بین $Z = 1/75$ تا $Z = -1/5$ چقدر است؟

$$P_R = -1/5 \Rightarrow \left[34/13 + 13/59 \times \frac{5}{10} \right] = 40/93$$

$$P_R = 1/75 \Rightarrow \left[34/13 + 13/59 \times \frac{75}{10} \right] = 44/32$$

$$40/93 + 44/32 = 85/25$$

نکات مهم

- ۱- استفاده از روش بالا حدود رتبه درصدی را نشان می دهد. برای محاسبه دقیق باید از جداول مربوط که سطوح بین Z های مختلف را ارائه کرده استفاده کرد. (پیوست ۱)
- ۲- انتقال از نمره های خام به نمره تراز Z شکل نمره را تغییر نمی دهد.
- ۳- اگر توزیع نمره های خام دارای کجی باشد، توزیع Z نیز دارای کجی است.
- ۴- برخلاف رتبه های درصدی اختلاف در نمره های Z اختلاف در نمره های خام را نشان می دهد.
- ۵- نسبت اختلاف نمره ها در توزیع اصلی یا نمره های خام، مساوی نسبت اختلاف بین نمره های Z آنها است. بنابراین فاصله بین اندازه نمره های اصلی در تبدیل به نمره های Z تغییر نمی کند.
- ۶- در نمره های نه گانه، ۹ فاصله یا واحد وجود دارد که هر فاصله یا واحد آن، مساوی نصف واحد انحراف معیار است. واحد میانی (پنجمین نمره نه گانه) نقطه میانی توزیع را نشان می دهد که دربرگیرنده ۲۰ درصد موارد است. چهارمین، سومین، دومین، و نخستین نمره نه گانه، از مرکز توزیع به سمت کرانه انتهایی پایین توزیع، به ترتیب ۱۷، ۱۲، ۷، ۴ درصد از موارد را دربرمی گیرد. (و برعکس) نمره ۹ گانه مانند رتبه درصدی نقطه ای بر روی مقیاس نیست، بلکه دربرگیرنده محدوده وسیعی از توزیع است. به همین دلیل نسبت به اختلاف های قابل توجه در طول مقیاس حساس نیست، اما نسبت به اختلاف های موجود در کرانه های انتهایی توزیع، به صورت گمراه کننده ای حساس است. (آیزاک، ۱۳۸۴: ۱۱۲).
- ۷- نمره های کلاسی مانند نمره های هوشی دارای توزیع بهنجار هستند. و بیشتر برای دوره ابتدایی که مطالب از پیوستگی بیشتری برخوردارند مناسب است. (آیزاک، ۱۳۸۴: ۱۱۵).

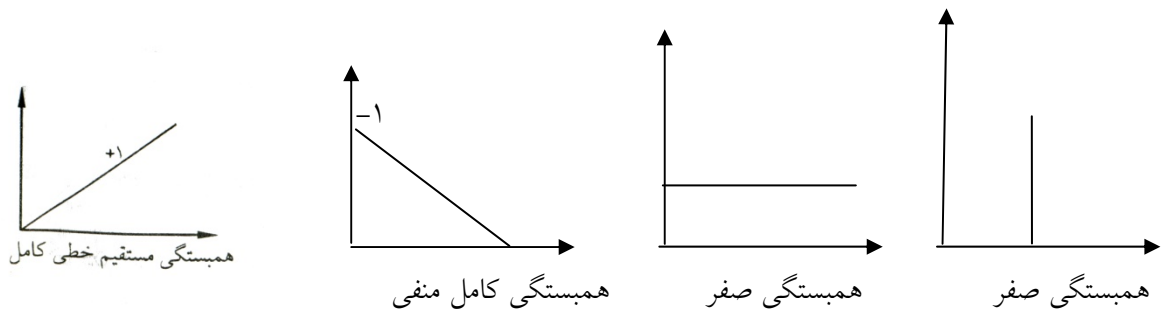
همبستگی

همبستگی بررسی رابطه بین دو متغیر با چند متغیر است، یعنی به دنبال این امر هستیم که آیا افزایش یا کاهش یک متغیر تأثیری روی متغیر دیگر دارد. فنون توصیفی در اینگونه مطالعات (که در آنها فقط ارتباط بین متغیرها، بدون آنکه هیچ یک از آنها دستکاری یا کنترل شود، مورد بررسی قرار می گیرد) عبارتست از میزان همخوانی (مانند ضریب فای یا ضرایب مشابه دیگر) برای مواردی که دو متغیر اندازه گیری شده طبقه ای باشد. ضریب رتبه ای اسپیرمن برای مواقعی که دو متغیر مورد اندازه گیری در مقیاس ترتیبی بیان می شوند، و سرانجام ضریب همبستگی پیرسون برای مواقعی که دو متغیر مورد بحث دارای مقیاس فاصله ای باشند. هرکدام از اینها نشان دهنده قدرت ارتباط بین دو متغیر است (هومن، ۱۳۸۲: ۱۲۷).

اگر هدف بررسی یک متغیر مثل X روی متغیر دیگر مثل Y باشد، بررسی همبستگی تک متغیری است. اصطلاحاً به پژوهش‌های تک متغیری موقعیت‌های مصنوعی گفته می‌شود. اگر هدف بررسی چند متغیر مثل X_1, X_2, \dots, X_n بر روی یک متغیر دیگر یعنی Y باشد، بررسی همبستگی چند متغیری است. اگر هدف بررسی چند متغیر X_1, X_2, \dots, X_n بر روی چند Y_1, Y_2, \dots, Y_n باشد، همبستگی کانونی است. در این مطالعات متغیر مستقل را متغیر پیش‌بینی‌کننده و متغیر وابسته را متغیر پیش‌بینی‌شونده گویند.

جهت همبستگی و شدت همبستگی

جهت همبستگی: اگر افزایش یک متغیر با افزایش متغیر دیگر همراه شود یا کاهش یک متغیر با کاهش متغیر دیگر همراه شود همبستگی مستقیم است. ولی اگر، افزایش یک متغیر با کاهش متغیر دیگر یا کاهش یک متغیر با افزایش متغیر دیگر همراه باشد، همبستگی منفی است. هنگامی که شیب پراکندگی (خط رگرسیون) افقی یا عمودی باشد، همبستگی صفر است.

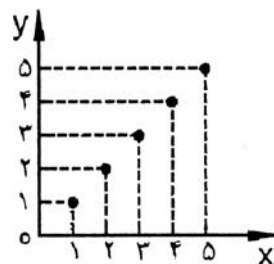


شدت همبستگی: دامنه همبستگی بین $+1$ و -1 یعنی همبستگی -1 معکوس و کامل، همبستگی $+1$ مستقیم و کامل است. مقدار همبستگی صفر تا 1 می‌باشد (شیولسون، ۱۳۸۰: ۱۷۶). هر داده‌ای بین این دو حد ناقص است. یعنی همبستگی $0/99$ ناقص است. هرچه همبستگی به قدر مطلق 1 نزدیکتر باشد شدت آن بیشتر است. مثال: همبستگی $-0/90$ از همبستگی $+0/70$ شدت بیشتری دارد.

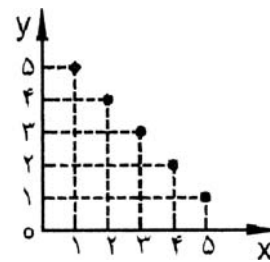
نکته: اگر ضریب همبستگی از $0/70$ بالاتر باشد می‌گوئیم رابطه خطی است. اگر پایین‌تر باشد رابطه غیرخطی است و در صورت غیر خطی بودن رابطه بین دو متغیر از ضریب همبستگی ایتا η استفاده کنیم.

$$\eta = 1 - \frac{SS_{reg}}{SS_{res}}$$

نام	X	Y
حسن	۵	۵
علی	۴	۴
احمد	۳	۳
مهدی	۲	۲
محمد	۱	۱



X	Y
۵	۱
۴	۲
۳	۳
۲	۴
۱	۵



نکته ۱: برای محاسبه همبستگی ابتدا باید نمودار پراکندگی را رسم کرد. اگر نقاط به خط فرضی نزدیک باشند، همبستگی زیاد است.

نکته ۲: چهار عمل اصلی $(-+ \div \times)$ یک نمره در کلیه نمرات یک توزیع هیچ تأثیری روی ضریب همبستگی ندارد. مثلاً اگر همبستگی بین Y, X (قد و وزن) $0/5$ باشد همبستگی بین $2X$ و Y همان $0/5$ خواهد بود.

مثال:

X	Y	X	Y
۱۵۰	۱۱۲	۱۵۰	۱۴۵
۱۶۵	۱۰۸	۱۶۰	۱۴۰
۱۷۵	۱۰۳	۱۷۰	۱۷۳
۱۸۰	۱۰۱	۱۶۰	۱۴۲
۱۹۲	۹۹	۱۵۰	۱۴۱
$r_{xy} = -1$		$r_{xy} = -1$	

انواع ضریب همبستگی

الف) در حالتی که دو متغیر اسمی و یا یکی اسمی و دیگری رتبه ای باشند: ضرایب عبارتند از کرامر- توافق C - لاندا- تاو گودمن و کراسکال.

- همبستگی کرامر: (جدول از 2×2 بیشتر باشد و حداقل یکی از متغیرها چند ارزشی است). رابطه مقارن نیست. مانند بررسی این فرض: بین جنس دانشجویان و گرایشات سیاسی آنها رابطه وجود دارد. جنس اسمی و گرایشات

سیاسی هم متغیر اسمی. (بی طرف، تندرو، محافظه کار). معادله آن عبارتست از: $\chi^2 = \frac{n(L-1)}{n(L-1)}$ که n تعداد حجم نمونه و L تعداد ردیف یا ستون، هرکدام که تعداد کمتری دارند.

سیاسی جنسیت	گرایش	بیطرف	تندرو	محافظه کار	جمع
خانم ها	۳۲	۲۲	۵	۵۹	
آقایان	۵	۳۰	۱۰	۴۵	
جمع	۳۷	۵۲	۱۵	۱۰۴	

$$v = \sqrt{\frac{21.0982}{104(2-1)}} = 4.5$$

نکته: هنگامی که تعداد کمتر سطر یا ستون برابر ۲ باشد، فرمول بالا به علت این که مقدار L برابر با یک می شود، به صورت زیر که به همبستگی فی معروف است، تبدیل می شود.

$$\varphi = \sqrt{\frac{\chi^2}{n}}$$

- **همبستگی c توافقی:** (جدول از ۲×۲ بیشتر باشد و حداقل یکی از متغیرها چند ارزشی است). معادله آن عبارتست از:

$$c = \sqrt{\frac{\chi^2}{\chi^2 + n}}$$

- **ضریب همبستگی لاندا:** این ضریب تفسیر روشن تری از همبستگی می دهد. شما تا چه اندازه می توانید از روی جنس افراد، گرایشات سیاسی آنها را پیش بینی کنید و یا برعکس. این ضریب امکان چنین پیش بینی را می دهد. معادله آن عبارتست از: $\gamma = \frac{e_1 - e_2}{e_1}$. در این رابطه e_1 اشتباه گروه بندی در موقعیت اول و e_2 اشتباه گروه بندی در موقعیت دوم می باشد.

(ب) **ضریب همبستگی در حالتی که هر دو متغیر دارای مقیاس رتبه ای باشند:** عبارتند از گاما- تاو کندال b- تاو کندال c- ضریب d سامرز.

نکته: همه ی این ضرایب گزینه هستند. بدین معنا که با تغییر متغیر مستقل و وابسته در مقدار ضریب تغییری ایجاد نمی شود.

- **ضریب گاما (G):** رابطه ی بین دو متغیر دو مقوله ای ترتیبی را به دست می دهد. مانند رابطه ی سطح تحصیلات مادران با نگرش آنان نسبت به تحصیلات دختران.

- **ضریب لامبدا (A):** رابطه ی بین دو متغیر نامتقارن است و رابطه ی ماقبل و مابعد را در یک توالی از رفتارها مطرح می کند. به ما می گوید متغیر a تا چه حد متغیر b را پیش بینی می کند. عکس آن ممکن نیست (سرمد و دیگران، ۱۳۸۰: ۲۲۳). مقدار آن بین صفر و یک است.

- **ضریب کپا (k):** اگر چند متغیر اسمی و رابطه متقارن باشد از این ضریب استفاده می شود. مثلاً اگر k داور N شیئی را به m مقوله اسمی تبدیل کنند، آماره کپا میزان توافق داورها در مورد اینگونه طبقه بندی را نشان می دهد. مقدار آن بین صفر و یک است.

- **ضریب هماهنگی (w) یا توافق کندال (u):** در پژوهش هایی که در آنها بیش از دو مجموعه رتبه وجود دارد و مایلیم بدانیم که بین رتبه هایی که توسط m داور به n فرد داده شده تا چه حد توافق وجود دارد، از این دو شاخص استفاده می شود. وقتی داده ها به صورت زوج های همتا و نه به صورت رتبه جمع آوری شده باشد، روش ضریب توافق کندال (u) مناسب تر است (سرمد و دیگران، ۱۳۸۰: ۲۲۴). مقدار w بین صفر و یک است. مقدار u اگر k زوج باشد $(\frac{-1}{k} - 1)$ و اگر k فرد باشد $(\frac{-1}{k})$ ، حداکثر برابر یک است.

$$\omega = \frac{SS_r}{\frac{1}{12} m^2 (n^2 - n)}$$

m = تعداد داوران

n = تعداد افراد رتبه بندی شده

σ^2 = واریانس مجموع رتبه ها

آزمون تقریبی معنادار بودن χ^2 نیز با استفاده از مشخصه آماری زیر بدست می آید:

$$X_{ab}^2 = m(n-1)\omega$$

نکته: هرگاه پژوهشگر بخواهد اثر متغیر سوم را ثابت نگهدارد و میزان رابطه‌ی دو متغیر دیگر را مشخص کند، از ضریب همبستگی رتبه ای تفکیکی کنдал باید استفاده کند (سرمد و دیگران، ۱۳۸۰: ۲۲۴)

همبستگی گشتاوری پیرسون

زمانی استفاده می شود که:

- ۱- هر دو متغیر توزیع نرمال داشته باشند.
- ۲- پراکندگی نمرات در هر دو متغیر یکسان باشد.
- ۳- توزیع خطی باشد.
- ۴- مقیاس متغیرها حداقل فاصله ای باشد (یا نسبی). مهمترین مفروضه این است که هر دو متغیر دارای حداقل مقیاس فاصله ای باشند. مثل: قد و وزن هوش و پیشرفت تحصیلی.

فرمول مهم

$$r_{xy} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{S_x \cdot S_y} = \frac{\text{cov}}{S_x \cdot S_y} = \frac{Zx \cdot Zy}{n}$$

ضریب همبستگی رتبه‌ای اسپیرمن

هنگامی که فقط دو دسته رتبه وجود داشته باشد (تعداد رتبه‌ها کم)، از ضریب اسپیرمن و یا تاو کنдал می توان استفاده کرد. ضریب اسپیرمن از زوج های همتا استفاده می کند. زمانی کاربرد دارد که داده‌ها رتبه‌ای باشد. یکی از متغیرها یا هر دو متغیر و یا اینکه داده‌ها فاصله ای، ولی مفروضه های آمار پارامتریک رعایت نشده باشد. مثل: ترتیب تولد و رتبه در کلاس. وقتی رابطه غیرخطی باشد، داده های فاصله‌ای باید به رتبه تبدیل شوند. (شیولسون، ۱۳۸۰: ۱۹۰). این ضریب برای رابطه های کاملاً غیر خطی مناسب نیست (شیولسون، ۱۳۸۰: ۱۹۶).

$$r_s = 1 - \frac{\sum D^2}{n(n^2 - 1)} \quad D = R_x - R_y$$

X خلاقیت	Y پیشرفت تحصیلی	R_x	R_y	D	D^2
۱۲۰	۱۸	۲	۳	-۱	۱
۱۱۲	۱۹	۴	۲	۲	۴
۱۱۸	۱۲	۳	۵	-۲	۴
۱۲۵	۱۷	۱	۴	-۳	۹
۱۰۱	۱۹/۵	۵	۱	۴	۱۶
					$\sum D^2 = ۳۴$

$$r_s = 1 - \frac{6 \times 34}{5(5^2 - 1)} \Rightarrow 1 - \frac{204}{120} = 1 - 1/7 \Rightarrow 0/7$$

نکته: زمانی که اعداد تکرار می‌شوند، برای تعیین رتبه اعداد مشابه باید بشرح زیر عمل کرد:

X	Y	R _x	R _y	D	D ²
۱۸	۲۹	۵	۶	-۱	۱
۱۹	۴۲	۳	۲	۱	۱
۱۷	۳۵	۷/۵	۴	۳/۵	۱۲/۲۵
۱۸	۳۵	۵	۴	۱	۱
۲۵	۲۸	۱	۷	-۶	۳۶
۱۲	۲۱	۹	۸	۱	۱
۲۳	۳۵	۲	۴	-۲	۴
۱۱	۶۶	۱۰	۱	۹	۸۱
۱۸	۱۱	۵	۱۰	-۵	۲۵
۱۷	۱۴	۷/۵	۹	-۱/۵	۲/۲۵
					$\sum D^2 = ۱۶۴/۵$

$$\text{رتبه عدد تکراری} \quad \frac{۴ + ۵ - ۶}{۳} = ۵ \quad \frac{۷ + ۸}{۲} = ۷/۵ \quad \frac{۳ + ۴ + ۵}{۳} = ۴$$

$$r_s = 1 - \frac{6 \times \sum D^2}{n(n^2 - 1)} = 1 - \frac{6 \times ۱۶۴/۵}{۱۰(۱۰^2 - 1)} = 1 - \frac{۹۸۷}{۹۹۰} = ۰/۰۰۳$$

نکته: اگر $n \leq ۳۰$ از جدول استفاده کنید. اگر $n \geq ۳۰$ از فرمول زیر مقدار بحرانی را بیابید. اگر بین مقادیر مثبت و منفی بود همبستگی وجود ندارد.

$$r_s = \frac{\pm z}{\sqrt{n-1}}$$

در مثال بالا اگر فرض کنیم افراد نمونه ۶۵ نفر هستند. بنابراین داریم:

$$r_s = \frac{\pm ۰/۰۰۳}{\sqrt{۶۵-1}} = \pm ۰/۲۴۵$$

است فرض صفر $\pm ۰/۲۴۵$ باشد، چون خارج از دامنه‌ی $\pm ۰/۲۴۵$ است فرض صفر رد می‌شود. یعنی بین رتبه‌ها همبستگی وجود دارد.

ضریب همبستگی دورشته‌ای نقطه‌ای^۱

زمانی مورد استفاده قرار می‌گیرد که یک متغیر پیوسته و دیگری دوارزشی واقعی باشد یا یک متغیر پیوسته و دیگری دوارزشی و فرض نرمال بودن توزیع رعایت نشده باشد. مثل: همبستگی بین نمره کل آزمون و یک سؤال یا جنسیت و نمره در یک آزمون که می‌تواند از یک کمتر باشد تابع p و q است. در حالت بیشینه مقدار آن از ۱ کمتر است. ($r = ۰/۷۹۷۸$)

^۱. Combination

۱- ضریب همبستگی دو رشته‌ای: یک متغیر پیوسته و دیگری دو ارزشی ساختگی است. می‌تواند از یک بیشتر باشد. توزیع نرمال و رابطه خطی است. منظور از ساختگی یعنی ما مثلاً نمرات ریاضی ۱۵-۹-۸ را به صورت مردود - قبول طبقه بندی کنیم (۰ و ۱).

X	X'
۱۵	۱
۱۵	۱
۹	۰
۸	۰

۲- ضریب همبستگی فی: زمانی کاربرد دارد که هر دو متغیر دوارزشی واقعی باشند یا هر دو متغیر دوارزشی و فرض نرمال بودن رعایت نشده باشد و رابطه متقارن باشد. یعنی جای متغیرها عوض شود نتیجه تغییر نمی‌کند. مانند رابطه جنسیت و وضع سواد (بی‌سواد- باسواد) و یا همبستگی بین دو سؤال (در کتاب خانم سرمد برای این مورد تراکوریک پیشنهاد شده است، ص ۲۲۲) که البته بیشتر فی بکار می‌رود.

۳- ضریب همبستگی تراکوریک: زمانی کاربرد دارد که یک یا هر دو متغیر دو ارزشی ساختگی باشند. یا هر دو متغیر دو ارزشی و فرض نرمال بودن رعایت شده باشد، مثل همبستگی بین پاسخ گروه قوی و ضعیف به سؤال ۱.

۴- ضریب همبستگی کانونی: هنگامی که بخواهیم رابطه بین دو یا چند متغیر را با دو یا چند متغیر دیگر بررسی کنیم، از این ضریب همبستگی استفاده می‌کنیم.

۵- همبستگی تفکیکی (پاره‌ای): ضریب همبستگی دو متغیری است. برای داده‌های ناپیوسته که اثر متغیر سوم و بعدی را می‌خواهیم روی هر دو متغیر حذف کنیم، بکار می‌رود.

$$r_{XY.Z}$$

مرتبه اول

$$r_{XY.ZF}$$

مرتبه دوم

$$r_{XY.Zfg}$$

مرتبه سوم

رابطه پیشرفت تحصیلی و خلاقیت با حذف اثر هوش بر روی هر دو متغیر (تفکیکی مرتبه اول)

۶- ضریب نیمه تفکیکی: اثر متغیر سوم و بعدی را روی یکی از متغیرها حذف می‌کنیم:

$$r_{X(Y.Z)}$$

مرتبه اول

$$r_{X(Y.Zf)}$$

مرتبه دوم

نکته: تفسیر ضریب همبستگی نباید برحسب درصد و نسبت باشد. مثلاً $r_{xy} = 0.70$ هفتاد درصد رابطه بین متغیرها

را تبیین نمی‌کند و $r_{xy} = 0.90$ دقیقاً دو برابر $r_{xy} = 0.45$ نیست.

خلاصه همبستگی ها

روش	نماد	متغیر ۱	متغیر ۲	توضیح
گشتاوری پیرسون	R	پیوسته	پیوسته	بائبات ترین روش است.
تفاوت رتبه ها (rho)	ρ	رتبه ای	رتبه ای	بیشتر به جای همبستگی گشتاوری به کار برده می شود، وقتی که تعداد موارد کمتر از ۳۰ است.
تای کندال	τ	رتبه ای	رتبه ای	زمانی که تعداد موارد ۱۰ کمتر است، به جای rho به کار برده می شود.
دو رشته ای	Γ_{bis}	دو ارزشی ساختگی	پیوسته	مقدار این همبستگی گاهی اوقات از ۱ بیشتر می شود و خطای استاندارد آن بزرگتر از ۱ است. از این روش در تجزیه و تحلیل سؤال استفاده می شود.
دو رشته ای گسترده	Γ_{wbis}	دو ارزشی ساختگی گسترده	پیوسته	زمانی که کار برده می شود که علاقمند به مطالعه افرادی هستید که در کرانه های انتهایی متغیر دو ارزشی قرار دارند.
دو رشته ای نقطه ای	Γ_{pbis}	دو ارزشی واقعی	پیوسته	ضریب همبستگی حاصل در این روش کوچکتر از ۱ و بسیار کوچکتر از Γ_{pbis} است.
چهارخانه ای (تتراکوریک)	Γ_t	دو ارزشی ساختگی	دو ارزشی ساختگی	زمانی که کار برده می شود که هر دو متغیر را بتوان در نقطه ای بحرانی معینی به دو قسمت تقسیم کرد.
ضریب فی	ϕ	دو ارزشی واقعی	دو ارزشی واقعی	برای محاسبه همبستگی درونی بین سؤال ها در آزمون های چندگزینه ای یا دوگزینه ای به کار برده می شود.
ضریب توافقی	c	دو یا چند طبقه ای	دو یا چند طبقه ای	در برخی شرایط معین، با Γ_t قابل مقایسه است. ارتباط نزدیکی با خی دو دارد.
شاخص رابطه نامتقارن سامرز	d	هر دو متغیر دو مقوله ای یا رتبه ای		اینکه کدام متغیر مستقل یا وابسته نامگذاری می شود در نوع محاسبه ای شاخص تأثیر دارد.

نکته

۱- متغیر پیوسته: به متغیری گفته می شود که پیوستار زیربنایی آن تمایل به توزیع طبیعی یا بهنجار دارد. نمونه‌هایی از این متغیر عبارتند از: وزن، توانایی یا پیشرفت تحصیلی که به وسیله آزمون های استاندارد اندازه‌گیری می‌شوند.

۲- متغیر دوارزشی ساختگی: وقتی حاصل می شود که یک متغیر پیوسته بر اساس یک معیار قراردادی که غالباً نزدیک مرکز داده‌هاست، به دو گروه تقسیم شود. به عنوان مثال، تقسیم بندی دانش آموزان به موفق- ناموفق، بالای متوسط، پایین متوسط، قبول- رد، و صمیمی- بی تفاوت در یک مقیاس نگرش سنج.

۳- متغیر دوارزشی واقعی: نقطه برش نسبتاً روشنی دارد (البته، نه ضرورتاً مطلق) به طوری که سرانجام داده های این متغیر به دو گروه تقسیم می‌شوند. نمونه‌هایی از این متغیر عبارتند از: مذکر- مؤنث، مرده- زنده، معلم- غیر معلم، مردود- نامردود، و سیگاری- غیرسیگاری. متغیرهای دیگری وجود دارند که برای محاسبه همبستگی، می‌توان آنها را دو ارزشی واقعی تلقی کرد. نظیر کوررنگی- غیرکوررنگی، الکلی- غیر الکلی و پاسخ های صحیح- غلط یک سؤال معین در یک آزمون. توزیع‌های زیربنایی دو ارزشی واقعی اگر دارای تفاوت های مطلق نباشند، دونمایی و یا به طور نسبی ناپیوسته هستند.

کوواریانس (واریانس مشترک)

کوواریانس تغییرپذیری مشترک بین دو متغیر را گویند (واریانس یک متغیر). مقدار و جهت رابطه بین دو متغیر را اندازه می گیرد. ولی دو محدودیت دارد. مقدار کوواریانس به مقدار تغییرپذیری نمره های X و Y بستگی دارد. اگر پراکندگی زیاد کوواریانس نیز زیاد باشد، در نتیجه نمی توان کوواریانس اندازه های مختلف را با هم مقایسه کرد. تفاوت در مقدار کوواریانس ممکن است در اثر رابطه بین متغیرها، تفاوت در انحراف معیارها و یا تفاوت در رابطه و انحراف معیارها باشد. شاخصی که از تفاوت در اندازه مقیاس ها تأثیر نپذیرد همبستگی است. این عمل تصحیح از رابطه الف بدست می آید (شیولسون، ۱۳۸۰: ۱۷۵).

$$\text{cov} = \frac{\sum(xy)}{n-1} = \frac{\sum(x-\bar{x})(y-\bar{y})}{n-1}$$

فرمول محاسبه

پرسش: اگر کوواریانس X و Y صفر باشد، کدام عبارت درباره X و Y صحیح است؟

الف) رابطه ای وجود ندارد. ب) رابطه ای غیر خطی وجود دارد.

ج) دو متغیر مستقل هستند. د) یا رابطه غیر خطی با استقلال وجود دارد.

پاسخ: گزینه (د) صحیح است. اگر کوواریانس صفر شد، نمی توان حکم به استقلال دو متغیر داد، چون ممکن است بین آنها رابطه‌ای غیرخطی وجود داشته باشد که بوسیله کوواریانس مشخص نمی شود. چون کوواریانس فقط نوع رابطه خطی دو متغیر را تعیین می کند (آذر، ۱۳۸۰: ۱۳۲).

مثال ۱: اگر cov برابر ۵۵ و $S_y^2 = ۶۴$ باشد، همبستگی بین Y و X چقدر است؟

$$r_{xy} = \frac{\text{cov}}{S_x \cdot S_y} = \frac{۵۵}{۸ \times ۹} = \frac{۵۵}{۷۲} = +۰/۷۶$$

مثال ۲: اگر $\sum (X - \bar{X}) = 5$ و $\sum (Y - \bar{Y}) = 8$ باشد و انحراف معیار X برابر ۱۰ و انحراف معیار Y برابر ۱۶ باشد، ضریب همبستگی چقدر است؟

$$r_{xy} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{S_x \cdot S_y} = \frac{5 \times 8}{10 \times 16} = \frac{40}{160} = +0.25$$

مثال ۳: اگر $Z_x = 2$ و $Z_y = 3/5$ باشد و حجم نمونه برابر ۴۰ نفر، ضریب همبستگی چقدر است؟

$$r_{xy} = \frac{Z_x \cdot Z_y}{n} = \frac{2 \times 3/5}{40} = \frac{7}{40} = +0.175$$

مثال ۴: چنانچه میانگین وزن دانش آموزان ($n=10$) برابر ۵۰ و میانگین قد همین گروه ۱۱۰cm باشد و فردی قد ۱۱۵cm و وزن ۵۸kg داشته باشد، میزان همبستگی بین قد و وزن او چقدر است؟ $S_1=5$ وزن و $S_2=3$ قد

$$\text{وزن } Z_x = \frac{X - \bar{X}}{S} = \frac{58 - 50}{5} = 1.6 \quad \text{قد } Z_y = \frac{Y - \bar{Y}}{S} = \frac{115 - 110}{3} = 1.67$$

$$r_{xy} = \frac{Z_x \cdot Z_y}{n} = \frac{1.6 \times 1.67}{10} = +0.27$$

$$V = r_{xy}^2 \times 100$$

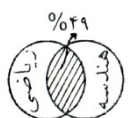
ضریب تعیین (درصد واریانس مشترک)

ضریب تعیین نشان دهنده میزان تأثیری است که متغیر X (مستقل) در متغیر Y (وابسته) ایجاد می کند یا به عبارتی با محاسبه این ضریب می توان تعیین کرد چند درصد از کل واریانس X ناشی از واریانس Y است.

$$r_{xy}^2 = \frac{\text{cov}^2}{S_x^2 \cdot S_y^2} = \frac{\text{واریانس مشترک}}{\text{کل واریانس}}$$

نکته: ضریب تعیین هیچ وقت منفی نخواهد شد، زیرا برای محاسبه آن ضریب همبستگی مجذور می شود.

سؤال: اگر همبستگی بین نمره ریاضی و نمره هندسه ۰/۷ باشد، درصد واریانس را بدست آورده و با شکل نشان دهید.



۵۱٪ عوامل دیگر

$$V = (0.7)^2 \times 100 = 49$$

عواملی که بر ضریب همبستگی تأثیر می گذارند:

۱- اساس رابطه از جامعه ای به جامعه دیگر فرق می کند. مثلاً در افراد بشر در سنین ۱۶-۱۰ سالگی بین سن تقویمی و توانایی فیزیولوژیکی همبستگی بالایی وجود دارد. ولی بین این دو متغیر در سنین ۲۶-۲۰ سالگی همبستگی وجود ندارد.

۲- پراکندگی متغیرها در جوامع مختلف متفاوت است. بدین معنی که هرچه تجانس بیشتر باشد (واریانس کمتر باشد)، همبستگی کمتر است. به عنوان مثال، در یک پژوهش اگر همه باهوش باشند، دامنه محدود و مقدار I کاهش می یابد.

۳- همبستگی بین دو متغیر تحت تأثیر همبستگی آنها با متغیر سوم قرار دارد. به عنوان مثال، همبستگی بین فیزیک و ریاضی ممکن است به دلیل همبستگی این متغیرها با هوش باشد.

۴- استفاده از گروه های انتهایی ضریب همبستگی را زیاد می کند. نمره های انتهایی در حجم های کوچک بر همبستگی تأثیر زیاد می گذارد.

۵- اگر n افزایش یابد، احتمال معناداری r بیشتر است.

۶- وقتی متغیرها نامرتبط باشند، r کاهش می یابد.

۷- رابطه غیر خطی امکان دارد مقدار ضریب همبستگی پیرسون را به صفر نزدیک کند.

۸- محدودیت در دامنه تغییرات ضریب همبستگی را کاهش می دهد. برای وجود محدودیت در دامنه تغییر باید واریانس ها یا انحراف معیارهای متغیرهایی که ضریب همبستگی آنها محاسبه می گردد مورد بررسی قرار گیرند. کوچک بودن واریانس ها می تواند نشانه محدودیت در دامنه تغییر باشد (شیولسون، ۱۳۸۰: ۱۹۷).

۹- اگر واریانس کوچک باشد، تفسیر ضریب همبستگی باید با احتیاط لازم صورت گیرد.

نکته ۱: آزمون معناداری برای همبستگی در دو متغیر مستقل آزمون t فریدلی است. و اگر گروه ها مستقل است آزمون Z فیشر و اگر ضریب همبستگی را در دو گروه وابسته مقایسه کنیم t استیودنت است. (فرگوسن، تاکانه، ۱۳۸۰).

نکته ۲: درجه آزادی همبستگی برابر است با: $d.f = n - 2$

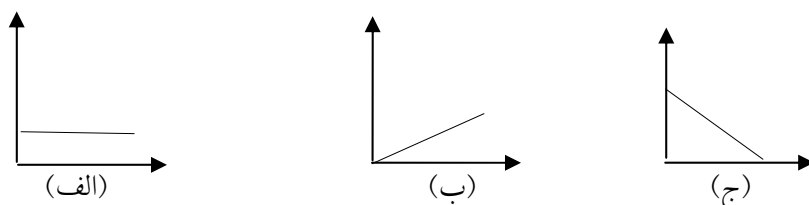
تحلیل همبستگی از نوع پیش بینی: رگرسیون و پیش بینی

زمانی که بین دو متغیر همبستگی وجود داشته باشد، می توان از طریق رگرسیون مقدار یک متغیر (Y) را از روی یک متغیر دیگر (X) پیش بینی یا برآورد کرد، و هرچه همبستگی بین متغیرها بالاتر باشد، به همان اندازه پیش بینی دقیقتر است. (معدل نمره دروس دانشگاهی را از روی نمره آزمون استعداد پیش بینی می کنیم) این پیش بینی از روی خط رگرسیون بدست می آید. این خط بهترین برازش از میان مجموعه نقاط نمودار پراکندگی است. راه دیگر که به داوری ذهنی بستگی ندارد، اصل کمترین مجذورات است که مجذور انحراف ها حول خط رگرسیون را کمینه می سازد. می توان معادله این خط را بدست آورد و سپس نمودار آن را رسم کرد. خط رگرسیون درحقیقت یک میانگین متحرک یا خط کمترین مجذورات است. رگرسیون درحقیقت معادله خط است. $Y=a+bx$

نکته: رابطه بین متغیر پیش بینی شونده (ملاک) (Y) و متغیر پیش بینی کننده (X) تابع علامت و شدت ضریب همبستگی است. اگر $r_{xy} = 1$ باشد، رگرسیون صفر و اگر $r_{xy} = 0$ باشد، همبستگی کامل است.

پیش بینی Y از روی X

شیوه شهودی: که از طریق ترسیم نمودار ممکن است. تعبیر و تفسیر آن تابع فرد مشاهده گر است.



شکل الف) چون خط افقی است، ارتفاع شیب آن صفر است. در این حالت پیش بینی دقیق امکان پذیر است.

این شکل دو عامل مهم در پیش بینی را دارد:

۱- میانگین متغیر Y یعنی \bar{Y} و شیب خط فرضی در نمودار پراکندگی. ۲- این نمودار نشان می دهد اگر شیب خط (ضریب زاویه) صفر باشد، برای هر مقدار از X نمره Y برابر است با میانگین نمره های Y (شیولسون، ۱۳۸۰: ۲۱۹).

شکل ب) بین X و Y رابطه مثبت و کامل وجود دارد. اگر X افزایش یابد مقدار Y به همان میزان افزایش می یابد. شکل ج) بین X و Y رابطه منفی و کامل وجود دارد. اگر X افزایش یابد مقدار Y به همان میزان کاهش می یابد.

$$\text{شیب و ضریب زاویه} = \frac{\text{تغییر در } Y}{\text{تغییر در } X} = \frac{Y_r - Y_1}{X_r - X_1}$$

برای پیش بینی بطور کل باید اطلاعات زیر را داشته باشیم:

الف) میانگین نمره های متغیر Y

ب) شیب خط فرضی که نمودار پراکندگی را توصیف می کند.

ج) موقعیت نسبی فرد در متغیر X (نمره مفروض). این مقدار باید به میانگین نمره های Y اضافه یا از آن کسر گردد.

$$Y = \bar{Y} + \frac{Y_r - Y_1}{X_r - X_1} (X - \bar{X})$$

پیش بینی با استفاده از معادله رگرسیون خطی

$$Y = a + bx$$

a = عرض از مبدأ (مقدار Y وقتی X صفر است).

$$b = \frac{Y_r - Y_1}{X_r - X_1} \quad B = \text{شیب خط (مقدار متغیر در } Y \text{ به ازای یک واحد تغییر در } X)$$

از این معادله هنگامی استفاده می شود که همه نقاط یک توزیع مشترک روی یک خط مستقیم قرار می گیرند.

پیش بینی نمره های استاندارد (Z)

مقدماتی ترین روشی که در استفاده از ضریب همبستگی پیرسون برای پیش بینی به کار برده می شود نمره های استاندارد است.

$$Z_Y = (Z_X)(r_{XY})$$

مهم: اگر همبستگی کامل (۱) باشد، نمره استاندارد پیش بینی نیز کامل، در غیر این صورت غیرکامل می باشد. به

عنوان مثال اگر نمره استاندارد Z برابر ۲ و همبستگی کامل (۱) باشد، نمره پیش بینی استاندارد نیز ۲ می باشد.

مهم: هنگامی که همبستگی بین دو متغیر کم باشد، نمره استاندارد می کنیم، نزدیک به میانگین خواهد بود.

نکته ۱: زمانی که همبستگی بین دو متغیر پایین باشد نمرات پیش بینی شده نزدیک به میانگین نمره پیش بینی شونده هستند تا نمره واقعی. به این پدیده رگرسیون می گویند.

نکته ۲: میزان همبستگی بین دو متغیر حدود یا مقدار اتفاق رگرسیون را تعیین می کند.

نکته ۳: پدیده رگرسیون اولین بار بوسیله گالتن مورد استفاده قرار گرفت. بر اساس مطالعات گالتن، فرزندان والدین بلندقد، بلندقد هستند اما نه به اندازه والدین خود. به همین ترتیب فرزندان والدین کوتاه قد، کوتاه قد هستند اما نه به کوتاهی والدین خود.

نکته ۴: خط رگرسیون، خطی است که خطاهای پیش بینی را به حداقل می رساند (خط حداقل مجذورها). یعنی اینکه مجموع مجذور فاصله Y ها از خط رگرسیون کوچکتر از فاصله هر خط دیگری تا محور Y ها می باشد. خط رگرسیون را خط برازنده نیز می نامند.

نکته ۵: اختلاف بین نمره واقعی (Y) و نمره پیش بینی شده (y') را خطای پیش بینی e می گویند.

$$e = y - y'$$

نکات بسیار مهم

- ۱- رگرسیون زمانی اتفاق می افتد که همبستگی کامل نباشد. یعنی نمرات از گروه های بالا و پایین جامعه انتخاب شده باشند (کرانه بالا و پایین).
- ۲- اگر همبستگی صفر باشد، رگرسیون کامل و اگر یک باشد رگرسیون صفر است.
- ۳- بین شدت همبستگی و رگرسیون همبستگی معکوس وجود دارد.
- ۴- انحراف نمره از خط رگرسیون از انحراف نمره از هر خط دیگری کوچکتر است.
- ۵- چنانچه پراکندگی نقاط در اطراف خط رگرسیون به شکل بیضی باشد، همبستگی بین دو متغیر در حد متوسط و مقدار آن نزدیک به ۰/۵ است (آیزاک، ۱۳۸۱: ۱۷۲).

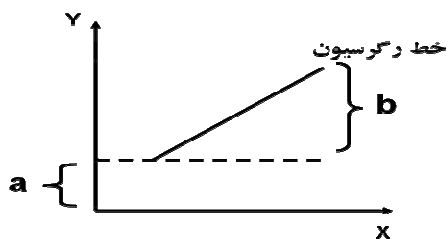
معادله خط رگرسیون

برای پیش بینی یک متغیر از روی متغیر دیگر از معادله خط رگرسیون استفاده می شود: $Y' = a + bx$ که در آن:

B = شیب خط $Y' =$ متغیری که می خواهیم پیش بینی کنیم

a = محل تلاقی خط رگرسیون با محور عرض ها یا عرض از مبدأ X = متغیری که مقدار آن را در اختیار داریم

نکته: بعضی مواقع مقدار a و b در اختیار است. در این صورت با جایگزینی مقادیر، محاسبه خیلی ساده می باشد.



مثال:

$$X = ۰$$

$$Y = ۴ + ۲(۰) = ۴$$

$$X = ۱$$

$$Y = ۴ + ۲(۱) = ۶$$

$$X = ۴$$

$$Y = ۴ + ۲(۴) = ۱۲$$

$$X = ۸$$

$$Y = ۴ + ۲(۸) = ۲۰$$

روش محاسبه ضرایب a و b :

$$b_{YX} = \frac{Cov_{XY}}{S_X^2} \quad b_{YX} = r_{XY} \frac{S_Y}{S_X} \quad b_{XY} = r_{XY} \frac{S_X}{S_Y}$$

شیب خط رگرسیون b_{yx} و Byx	انحراف استاندارد $S_y = y$
ضریب همبستگی r_{xy}	انحراف استاندارد $S_x = x$

محاسبه a و b از روی داده های خام

$$b_{yx} = \frac{N \sum XY - (\sum X)(\sum Y)}{N \sum X^2 - (\sum X)^2}$$

$$a_{yx} = \frac{\sum Y - b_{yx} \sum X}{N}$$

مثال: برای داده های زیر معادله خط رگرسیون را به دست آورید؟

X	Y	X ²	Y ²	XY
6	8	36	64	48
9	9	81	81	81
5	7	25	49	35
20	24	142	194	164

$$b_{yx} = \frac{3(164) - (20)(24)}{3(142) - (20)^2} = 0.46$$

$$a_y = \frac{24 - (0.46)(20)}{3} = 4.93$$

$$y = 4.93 + 0.46(x)$$

خطای استاندارد برآورد

اکثر مواقع بین نمرات پیش بینی شده و نمرات مشاهده تفاوت وجود دارد. به این حالت، اختلاف خطای استاندارد برآورد می گویند.

نکته: هرچه همبستگی بین متغیرها بیشتر باشد، خطای پیش بینی کمتر خواهد بود. این خطا به دو روش محاسبه می شود.

$$S_{yx} = S_y \sqrt{1 - r^2_{xy}} \quad \text{خطای استاندارد پیش بینی} \quad S_{yx} = \sqrt{\frac{\sum e^2}{N}} \quad \text{یا} \quad S_{yx} = \sqrt{\frac{\sum (Y - \bar{Y})^2}{N - 2}}$$

$$r_{xy} = \text{ضریب همبستگی} \quad e = \text{خطای پیش بینی}$$

نکته: $N - 2$ را درجه آزادی خطای استاندارد پیش بینی گویند.

رابطه پایایی و خطا

$$S_E = S \sqrt{1 - r_{xx}}$$

همانطور که ملاحظه می شود r_{xx} ضریب اعتبار (پایایی) است.

نکته ۱: اگر همبستگی کامل باشد خطا صفر است. یعنی همبستگی با خطا رابطه عکس دارند.

نکته ۲: انحراف معیار نمونه و خطای پیش بینی رابطه مستقیم دارند (به فرمول ها توجه کنید).

نکته ۳: بین خطای اندازه گیری و ضریب پایایی رابطه عکس وجود دارد.

انواع رگرسیون

تک متغیری: از روی یک متغیر، متغیر دیگر را پیش بینی می کنیم. از روی هوش، خلاقیت فرد را پیش بینی می کنیم.

چندگانه: از روی چند متغیر، یک متغیر پیش بینی می شود. از روی اعتماد به نفس و اضطراب و خلاقیت، پیشرفت

تحصیلی پیش بینی می شود.

چند متغیری: چند متغیر مستقل، چند متغیر وابسته را پیش‌بینی می‌کنند.

- نکته ۱:** وقتی همبستگی چند متغیر X را با چند متغیر Y محاسبه می‌کنیم، همبستگی قانونی یا قانونی گفته می‌شود.
- نکته ۲:** مقدار همبستگی چندگانه معمولاً از همبستگی بین هریک از متغیرهای پیش‌بین و ملاک بیشتر است. (کازی، ۱۳۸۶: ۲۷۸).

راهنمای تحلیلی در انتخاب متغیرهای مستقل برای تشکیل معادله‌ی رگرسیون

- ۱- **راهبرد همزمان:** همزمان همه‌ی متغیرهای مستقل با هم وارد تحلیل می‌شوند.
- ۲- **رگرسیون گام به گام:** در این روش نحوه‌ی ورود متغیرها به تحلیل در اختیار محقق نیست. از روش گام به گام برای پیش‌بینی متغیر وابسته استفاده می‌شود. از این روش برای تبیین استفاده نمی‌شود.
- ۳- **رگرسیون سلسله مراتبی:** ترتیب ورود متغیرها با توجه به چارچوب نظری و توسط محقق مشخص شده است. چنانچه پژوهشگر بخواهد در پیش‌بینی پیشرفت تحصیلی از متغیرهای خلاقیت و هوش استفاده کند، می‌تواند با استفاده از این روش و با توجه به این که هوش بر خلاقیت تقدم دارد، عمل کند.
- نکته ۱:** برای پی‌بردن به معنادار بودن آماره‌های محاسبه شده از آزمون‌های آماری زیر استفاده می‌شود: آزمون معنادار بودن R^2 و ΔR^2 که از توزیع F پیروی می‌کند و آزمون معناداری ضرایب رگرسیون که از توزیع T پیروی می‌کند.
- نکته ۲:** در رابطه با ضرایب استاندارد نشده b و استاندارد شده β باید توجه داشت، هرگاه بخواهید اثر متغیرهای متفاوت را در یک جامعه‌ی واحد مقایسه کنید از β و هرگاه خواسته باشید اثر متغیرهای خاصی را در جوامع متفاوت مقایسه کنید از b استفاده نمایید (سرمد و همکاران، ۱۳۸۰: ۲۳۲).

روش‌های دیگر تحلیل رگرسیون

- ۱- **آزمون χ^2 نسبت درست نمایی (آزمون G^2)**
این آزمون مانند χ^2 معمولی برای داده‌های 2×2 به کار می‌رود. اما به دلایل زیر بر آزمون χ^2 معمولی برتری دارد:
- کمتر تحت تأثیر حجم نمونه قرار می‌گیرد.
- قابل افراز به مؤلفه‌هایی است که مانند مجموع مجزورات در تحلیل واریانس جمع پذیرند.
- در مدل لگاریم خطی که در آن با جداول چند بعدی سروکار داریم مورد استفاده قرار می‌گیرد.

۲- تحلیل ممیز^۱

در مواردی استفاده می‌شود که متغیر وابسته اسمی و متغیرهای مستقل کمی باشند. در تحلیل ممیز باید گروه‌ها ناهمپوش و عضویت آزمودنی‌ها در آنها مشخص گردد.

برای مثال فرض می‌کنیم دانشگاهی می‌خواهد موفقیت یا شکست متقاضیان ورود به دوره‌ی دکترای رشته‌ی خاصی را پیش‌بینی کند. این دانشگاه با داشتن ویژگی‌های داوطلبان سال قبل مانند پیشرفت تحصیلی، توصیه‌نامه‌ها و نمرات امتحانات ورودی که به عنوان متغیر مستقل در نظر گرفته می‌شود، با تحلیل ممیز می‌تواند به پیش‌بینی دست بزند.

^۱. Discriminant

۳- رگرسیون لجستیک

در مدل تحلیل ممیز همه‌ی متغیرهای پیش بین با مقیاس کمی (فاصله ای یا نسبتی) اندازه گیری شده و متغیر ملاک (وابسته) مقوله ای است. در صورتی که متغیرهای پیش بینی هم در مقیاس کمی و هم در مقیاس مقوله‌ای اندازه‌گیری شده باشند، از تحلیل ممیز نمی‌توان استفاده کرد. در این موارد از رگرسیون لجستیک استفاده می‌شود که متغیر، متغیر وابسته مقوله ای و دو سطحی است. این دو مقوله معمولاً به عضویت یا عدم عضویت در یک گروه یا بلی و خیر اشاره دارد. در معادله‌ی رگرسیون معمولی از تعدادی متغیرهای پیش بین با ضرایبی (وزن‌ها) برای پیش‌بینی متغیر وابسته استفاده می‌شود. در رگرسیون لجستیک آنچه که پیش بینی می‌شود یک احتمال است که ارزش آن بین ۰ و ۱ تغییر می‌کند. در این رگرسیون از مفهوم بخت استفاده می‌شود که عبارتست از نسبت احتمال وقوع یک پدیده بر احتمال عدم وقوع آن.

نکته: واژه‌ی کلیدی در رگرسیون لجستیک سازه ای به نام لوجیت^۱ است که لگاریتم طبیعی بخت می‌باشد. مثال: فرض کنید پژوهشگری می‌خواهد کمک‌کردن یا نکردن افراد را به دیگری برحسب احساس همدردی، احساس توانایی و قومیت فرد، پیش بینی کند. در این مثال احساس همدردی و توانایی در مقیاس فاصله ای و قومیت اسمی می‌باشد.

۴- مدل لگاریتم خطی

در مواردی که تمام متغیرهای مدل در مقیاس اسمی اندازه گیری شده و داده‌ها به صورت فراوانی باشد، روش معمولی برای پی بردن به رابطه‌ی بین دو متغیر مقوله ای استفاده از آزمون χ^2 است. هرگاه تعداد متغیرهای مقوله‌ای بیش از دو باشد، تفسیر جداول χ^2 مشکل یا غیر ممکن می‌شود. در چنین مواردی می‌توان از مدل لگاریتم خطی استفاده کرد. که اساس آن شبیه روش رگرسیون چندگانه است. با این تفاوت که در این روش متغیر مستقل، متغیرهای مقوله‌ای و اثرمتقابل آنها، و متغیر وابسته لگاریتم طبیعی فراوانی‌های خانه‌ای است. مثال: اثر جنس، قومیت، سطح درآمد و اثرهای متقابل آنها بر فراوانی‌های یک جدول توافقی.

نکته: همانند تحلیل واریانس اثرهای متقابل نامگذاری می‌شوند. به عنوان مثال در سه عاملی اثر متقابل (مرتبه سوم) نامیده می‌شود. مدل‌هایی که تمامی اثرات ممکن را دربرداشته باشد، مدل اشباع شده نامیده می‌شود. که مثال بالا این چنین است. چون ساده نیست مطلوب ترین روش محسوب نمی‌شود. تحلیل مدل لگاریتم خطی به دو صورت سلسله مراتبی و غیر سلسله مراتبی است.

نکته: تفاوت عمده مدل لگاریتم خطی با مدل رگرسیون لجستیک که مدل لوجیت نیز نامیده می‌شود، این است که در مدل لوجیت یک متغیر مقوله ای دو سطحی به عنوان متغیر وابسته وجود دارد. درحالی‌که در لگاریتم خطی متغیر وابسته، فراوانی‌های خانه‌ها است (سرمد و دیگران، ۱۳۸۰: ۲۶۵).

^۱ . Logit

تحلیل پروبیت^۱ و لوجیت

برای پاسخگویی به پرسش هایی از قبیل: چند ساعت مطالعه در هفته برای بدست آوردن نمره ای معادل ۸۰۰ در یک آزمون مانند تافل لازم است؟ که در چنین پرسش هایی متغیر تأثیرگذار که در یک مقیاس پیوسته اندازه گیری می شود (ساعات مطالعه) و یک متغیر برون داد (وابسته) که مقوله ای است و موفقیت یا شکست را اندازه گیری می کند، از این تحلیل ها استفاده می شود. این شیوه ها برای تبدیل داده های خام به داده های استاندارد بکار می روند. مقدار هر دوی آنها بین ۰/۲ و ۰/۸ تغییر کرده و مانند Z عمل می کند. برای تبدیل به پروبیت، نمره ی Z را با ۵ جمع می کنیم. برای تبدیل به لوجیت از فرمول زیر استفاده می شود.

$$\text{نمره لوجیت} = \frac{\ln \frac{p}{(p-1)}}{2} + 5$$

Ln = لگاریتم طبیعی و p = نسبت پیامدهای موفقیت آمیز است.

مثال: اگر ۱۵۰ دانشجو ۳۰ ساعت در هفته مطالعه ی زبان انگلیسی کرده باشند و ۱۱۱ نفر از آنها توانسته باشند

$$\text{نمره ای بالاتر از } 1200 \text{ به دست آورند، نسبت موفقیت برابر است با: } \frac{111}{150} = 0.74$$

نمره ی Z متناظر با آن (۰/۷۴ درصد) در توزیع نرمال برابر ۰/۶۴ است. از این رو پروبیت این پیامد برابر است با:

$$0.64 + 5 = 5.64 \text{ لوجیت این پیامد نیز عبارتست از:}$$

$$\text{لوجیت} = \text{Ln} \frac{0.74(1-0.74)}{2} + 5 = 5.52$$

تحلیل عاملی

برای بررسی روایی سازه از روش آماری تحلیل عاملی استفاده می شود که با تشکیل ماتریس همبستگی سازه ی مورد نظر بررسی می شود. برای انجام تحلیل عاملی ابتدا باید اطلاعات موجود در ماتریس معنادار باشد. برای این منظور از آزمون χ^2 بارتلت استفاده می شود. معنادار بودن χ^2 و آزمون بارتلت^۲ حداقل شرط لازم برای انجام محاسبات تحلیل عاملی است. در آزمون بارتلت فرض صفر این است که متغیرها فقط با خودشان همبستگی دارند. رد فرض صفر به این معنا است که ماتریس حاوی اطلاعات معناداری است. این آزمون را آزمون کرویت^۳ گویند.

احتمالات و توزیع دو جمله ای

احتمال یک نسبت است که صورت آن تعداد دفعاتی است که پیشامد رخ داده و مخرجش کل کوشش ها است.

$$P = \frac{f}{N}$$

مثال: همیشه تعداد تولد پسرها بیشتر از دخترها است. در یک گروه معمولی ۲۰۵ نوزاد به دنیا آمده است. از بین این نوزادان ۱۰۵ نفر پسر هستند. اگر یک نوزاد به طور اتفاقی از بین این گروه انتخاب شود، احتمال اینکه این نوزاد پسر نباشد چقدر است؟

^۱. Probit

^۲. Bartlett

^۳. Sphericity

چون از ۲۰۵ نوزاد ۱۰۵ نفر پسر هستند پس ۱۰۰ نوزاد دخترند. در نتیجه:

$$P_c = \frac{100}{205} = 0/488$$

نکته: احتمال یک پیشامد غیر ممکن صفر است. احتمال رخ دادن یک پیشامد مسلم یک است.

فضای نمونه: از تمام پیشامدهای ساده تشکیل می شود. یعنی فضای نمونه از تمامی پیشامدهایی تشکیل می شود که دیگر نمی توان آنها را به مؤلفه های ساده تر تبدیل کرد.

متمم رویداد A: شامل تمام پیشامدهایی می شود که در آنها پیشامد A اتفاق نمی افتد.

نکته: وقتی مقدار احتمال را بیان می کنیم باید یا باید عدد دقیق اعشاری را بیان کنیم یا اینکه عدد را تا ۳ رقم معنادار

گرد کنیم. اگر عدد کسر ساده مانند $\frac{2}{3}$ است لازم نیست گرد شود.

پیشامد مرکب^۱: به پیشامدی گفته می شود که دارای ۲ یا بیشتر پیشامد ساده باشد.

$$S = P^n$$

فرمول فضای نمونه کل حالتها می ممکن:

محاسبه ی احتمال

در احتمال محاسبات بر اساس دو قانون جمع و ضرب انجام می شود. قانون جمع زمانی به کار می رود که یکی از دو حادثه اتفاق بیفتد، در صورتیکه قانون ضرب در شرایطی به کار برده می شود که هر دو حادثه با هم اتفاق بیفتند.

دو رویداد ناسازگار: به حوادثی گفته می شود که وقوع یکی از آنها مانع وقوع دیگری باشد. برای مثال در پرتاب سکه، نشستن یک روی سکه مانع نشستن طرف دیگر می شود. در این صورت احتمال اینکه یکی از دو رویداد اتفاق بیفتد مساوی است با مجموع احتمال هریک از آنها.

$$P(A \text{ or } B) = P_A + P_B \Rightarrow P(A \cap B) = 0$$

دو رویداد سازگار: هنگامی که دو رویداد سازگارند امکان وقوع آنها به صورت همزمان وجود دارد.

$$P(A \text{ or } B) = P_A + P_B - P(A \cap B)$$

مثال: اگر از یک دسته ورق دو کارت بکشیم، احتمال آمدن ۱۰ تا دل چقدر است؟

$$P(A \text{ یا } B) = \frac{4}{52} + \frac{13}{52} - \frac{1}{52} = \frac{16}{52}$$

قاعده ی ضرب احتمال ها

رویدادهای مستقل: در ضرب، هر دو رویداد با هم اتفاق می افتد. به هم این و هم آن معرف است. در این صورت دو حالت وجود دارد. یا رویدادها مستقلند یا وابسته:

- **ضرب رویدادهای مستقل:** به رویدادهایی مستقل گویند که وقوع یا عدم وقوع هریک در احتمال وقوع یا عدم وقوع حادثه ی دیگر تأثیر نداشته باشد. برای مثال تولد فرزندان و یا احتمال آمدن عدد ۵ در انداختن دو تاس سبز و آبی.

$$P(A \times B) = P(A) \times P(B) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

^۱. Compound even

- ضرب رویدادهای وابسته: اگر دو رویداد وابسته باشند، یعنی احتمال وقوع یکی از آنها به احتمال وقوع دیگری وابسته باشد، احتمال آن که هر دو اتفاق بیفتند عبارتست از:

$$P(A \cap B) = P_A \times P_{B/A} = P_A \times \frac{P(A \cap B)}{P(A)} = P\left(\frac{B}{A}\right) = \frac{P(A \cap B)}{P(A)}$$

مثال: در یک کیسه ۱۸ مهره وجود دارد که ۶ مهره مشکی و ۱۲ مهره آن سفید است. احتمال اینکه اگر دو مهره خارج کنیم هر دو سفید باشد، چقدر است؟ (بدون جایگزینی)

$$P_A = \frac{12}{18} \Rightarrow P_B = \frac{11}{17} \rightarrow P = \frac{12}{18} \times \frac{11}{17} = \frac{132}{306}$$

فاکتوریل: $n!$

نشان‌دهنده حاصلضرب اعداد صحیح مثبت است، به طوری که هر عدد در عدد صحیح قبل از خود ضرب می‌شود. این عمل تا رسیدن به عدد یک ادامه پیدا می‌کند.

$$4! = 4 \times 3 \times 2 \times 1 = 24$$

مثال: به چند روش می‌توان ۴ نوع کتاب را کنار هم چید؟

نکته: $0! = 1$

فرمول ترکیب^۱: زمانی به کار می‌رود که ترتیب چیدن افراد یا اشیاء برای ما مهم نیست. ترکیب abc با cba فرقی نمی‌کند.

$$C_r^n = \frac{n!}{r!(n-r)!}$$

نکته: باید ایتهم‌ها متفاوت باشند. اگر مشابه باشند از این قاعده استفاده نمی‌شود.

ما باید ۲ مورد را از n آیتهم انتخاب کنیم (بدون جایگزینی)

مثال: ۶ نفر می‌خواهند دو به دو کنار هم بنشینند. اینکه چه افرادی کنار هم بنشینند، مهم نیست. تعداد حالت‌های ممکن برای نشستن را محاسبه کنید.

$$C_2^6 = \frac{6!}{2!(4!)} = \frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{2 \times 1 \times (4 \times 3 \times 2 \times 1)} = 5$$

فرمول ترتیب: زمانی بکار می‌رود که ترتیب قرار گرفتن موارد برایمان مهم است.

$$A_r^n = \frac{n!}{(n-r)!}$$

مثال: ۶ نفر می‌خواهند دو به دو کنار هم بنشینند. این که چه افرادی کنار هم بنشینند، مهم است. تعداد حالت‌های ممکن برای نشستن را محاسبه کنید.

$$A_2^6 = \frac{6!}{(6-2)!} = \frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{4 \times 3 \times 2 \times 1} = 30$$

بسط دو جمله ای نیوتن (پاسکال)

$$[p+q]^n = p^n + c.p^{n-1}.q + c.p^{n-2}.q^{n+1} + \dots q^n$$

^۱. Combination

شرایط استفاده

- ۱- آزمایش‌ها فقط نشان دو نتیجه پیروزی و شکست باشد.
- ۲- آزمایش‌ها مستقل از هم باشند.
- ۳- تعداد دفعات کوشش مشخص باشد.

هنگام مواجه با مجموعه‌ای از پرسش‌های مشابه زیر از بسط دو جمله‌ای که راحت‌تر است استفاده کنید:

مثال: در امتحان ۳ سؤالی، ۴ جوابی، احتمال اینکه دانش آموزی بصورت تصادفی به یک سؤال پاسخ دهد، چقدر است؟

$$\frac{1}{4}, q = \frac{3}{4} \rightarrow (p+q)^3 = p^3 + 3p^2q + 3p^1q^2 + q^3 = 1$$

- احتمال اینکه به همه‌ی سؤال‌ها از روی شانس پاسخ درست دهد برابر است با:

$$p^3 = \left(\frac{1}{4}\right)^3 = \frac{1}{64}$$

- احتمال اینکه به دو سؤال پاسخ درست دهد برابر است با:

$$3p^2q = 3\left(\frac{1}{4}\right)^2\left(\frac{3}{4}\right) = \frac{9}{64}$$

- احتمال اینکه به یک سؤال پاسخ درست دهد برابر است با:

$$3p^1q^2 = 3\left(\frac{1}{4}\right)\left(\frac{3}{4}\right)^2 = \frac{27}{64}$$

- احتمال اینکه اصلاً پاسخ درست ندهد برابر است با:

$$q^3 = \left(\frac{3}{4}\right)^3 = \frac{27}{64}$$

- احتمال اینکه حداقل به یک سؤال پاسخ درست دهد چقدر است؟

$$p^2 + 3p^2q + 3p^1q^2 = \left(\frac{1}{64}\right) + \left(\frac{9}{64}\right) + \left(\frac{27}{64}\right) = \frac{37}{64}$$

$$\bar{X} = P.n$$

فرمول میانگین در توزیع دو جمله‌ای

n: تعداد دفعات تکرار آزمایش

P: احتمال پیروزی

$$\bar{X} = \frac{1}{2} \times 10 = 5$$

مثال: اگر سکه‌ای ده بار پرتاب شود، میانگین آن چقدر است؟

$$S_f = \sqrt{N.p.q}$$

فرمول خطای استاندارد برآورد در توزیع دو جمله‌ای:

مثال: اگر در امتحان ۱۰ سؤالی شرکت کرده باشیم و امتحان ۴ گزینه‌ای باشد، خطای استاندارد برآورد آنرا محاسبه

$$S_p = \sqrt{10 \times \frac{1}{4} \times \frac{3}{4}} = \sqrt{\frac{30}{16}} = 1.37$$

نمایید؟

انواع تحلیل داده

تحلیل داده‌های کیفی: مستلزم سه فعالیت است:

- ۱- تلخیص داده‌ها: که منظور انتخاب، تمرکز، تنظیم و تبدیل داده‌ها است. رمزگذاری داده‌ها و اینکه پژوهشگر مشخص کند چه قسمتی از داده‌ها مفید و چه قسمتی نمی‌تواند مورد استفاده قرار گیرد در این بخش قرار دارد.

- ۲- **عرضه داده ها**^۱: منظور ظاهر ساختن مجموعه‌ای از داده های سازمان یافته است. مانند مطالب روزنامه، رایانه. شکل رایج عرضه داده های کیفی در گذشته استفاده از متن های داستان گونه بوده است.
- ۳- **نتیجه گیری / تأیید**^۲: پی بردن به معنای هر رویداد، الگوهای وقوع آنها، تبیین آنها. با نتیجه گیری نمی توان فعالیت پژوهشی را پایان یافته تلقی کرد. زیرا این فعالیت باید مورد تأیید واقع شود. یعنی باید از جنبه موجه بودن، استحکام و قابلیت تأیید مورد بازبینی قرار گیرد. این فرایند همان تعیین اعتبار نتایج است.
- تحلیل داده های کمی**: چنانچه داده ها کمی و پیوسته باشند، دارای توزیع بهنجار بوده و آزمودنی ها با نمونه گیری تصادفی انتخاب شده باشند، برای تحلیل داده ها می توان از آزمون های پارامتری استفاده کرد. این آزمون ها دارای مفروضه های نرمال بودن توزیع نمرات، یکسانی واریانس در گروه ها و داشتن مقیاس فاصله ای و نسبتی و پیوسته بودن داده ها است. مانند آزمون t. اگر این مفروضه ها وجود نداشته باشند، باید از آزمون های غیر پارامتری که با متغیرهای اسمی و رتبه ای سروکار دارد استفاده کرد. در مواردی که مقیاس فاصله ای است، ولی حجم نمونه کوچک است نیز از آزمون های ناپارامتری استفاده می شود.
- نکته**: در تحلیل های پارامتری می توان پارامترهای جامعه را برآورد کرد، درحالی که در تحلیل های ناپارامتری فقط آزمون فرض صورت می گیرد (سرمد و دیگران، ۱۳۸۰: ۲۱۳).

^۱ . Data display

^۲ . Conclusion drawing/ verification