

سیزده راه نگرش به

ضریب همبستگی*

جوزف لی راجرز، و. ال. نایسواندر*

ترجمه علی عمیدی

است.

ایده اصلی همبستگی اساساً قبل از ۱۸۸۵ عنوان شده بود [۱۳]. پیرسن در سال ۱۹۲۵ ارائه رویه نرمال دو متغیر همبسته (در ۱۸۲۳) را به گاوس نسبت داد. اما گاوس به همبستگی، به عنوان یک مفهوم متمایز، توجه خاصی نداشت و در معادله‌های توزیع‌هایش، همبستگی را به عنوان یکی از پارامترها تعبیر کرد. پیرسن در یک مقاله تاریخی قبلی که در ۱۸۹۵ انتشار یافت، ارائه توزیع نرمال دو متغیره (در ۱۸۴۶) را به اوگوست براوه، ستاره‌شناس فرانسوی، نسبت داد (سک. [۳۵]). براوه در واقع، به پارامتری از توزیع نرمال دو متغیره، عنوان "همبستگی" را اطلاق کرد، اما نظیر گاوس اهمیت همبستگی را به عنوان معیار پیوند بین متغیرها تشخیص نداد. [قبل از ۱۹۲۵، پیرسن امتیازی را که به براوه داده بود پس گرفت. اما والکر [۲۶] و سیل [۲۴] تاریخچه‌ای را که پیرسن گزارش داده و در پروردادن آن کمک کرده بود مرور کردند و ادعای براوه را در تقدم تاریخی تأیید نمودند.] چارلز داروین، دایی زاده گالتن، در ۱۸۶۸ با اشاره به اینکه "تمام اجزاء یک نظام به میزانی معین به هم مربوط یا همبسته‌اند" مفهوم همبستگی را به کار برد. سپس، در ۱۸۷۷، گالتن در یک سخنرانی مربوط به بستگی مشخصه‌های جسمانی نسلهایی از والدین و فرزندان، برای اولین بار به اصطلاح "Reversion" اشاره کرد. "قانون Reversion" اولین توصیف رسمی از مفهومی است که گالتن بعدها "رگرسیون" را به جای آن باب کرد.

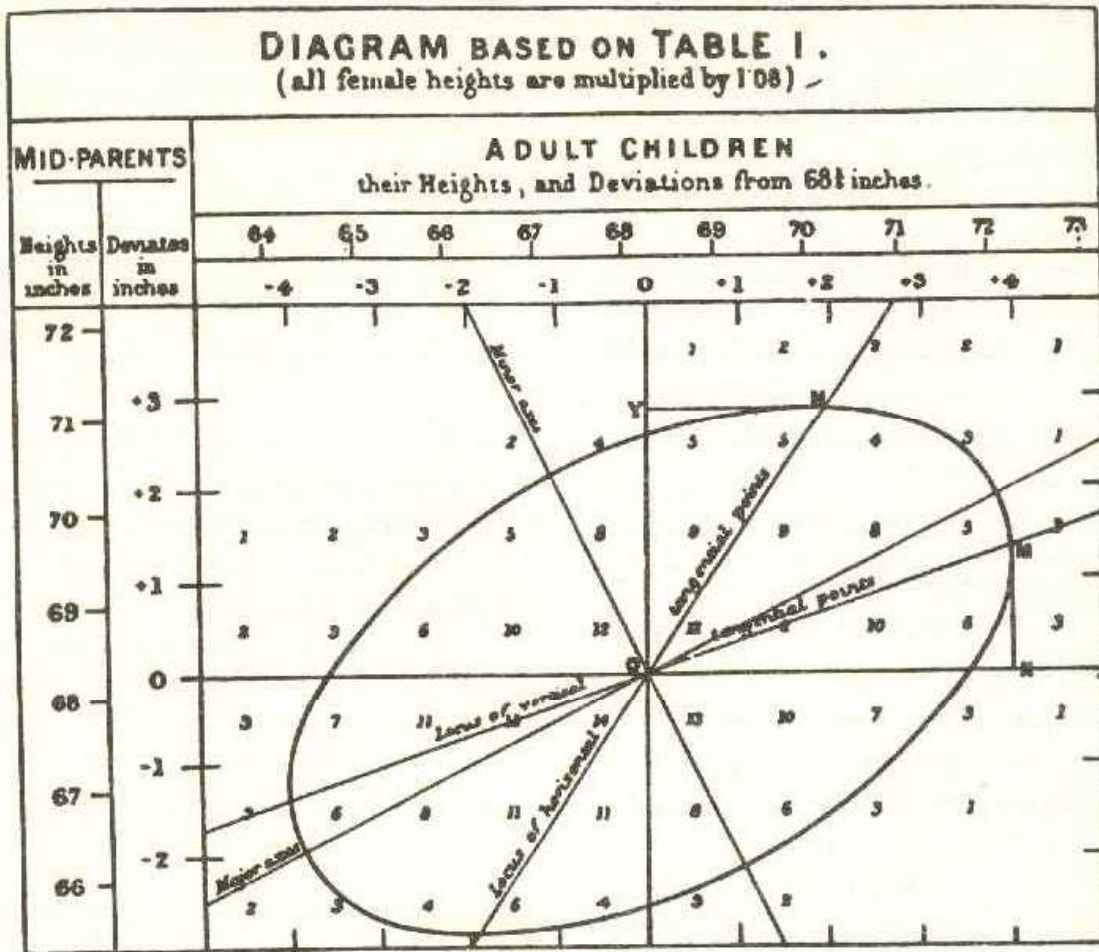
در طی این دوره، پیشرفت مهم فلسفه نیز سبب افزایش بار معنایی مفاهیم همبستگی و رگرسیون شد. در ۱۸۴۳، جان استوارت میل فیلسوف بریتانیایی برای اولین بار "پنج قانون تحقیق تجربی" خود را ارائه داد. بین اینها روش تغییرات ملازم هم را گنجانیده بود: "هر وقت تغییر پدیده‌ای به روشی خاص، به نحوی موجب تغییر

در ۱۸۸۵، سرفرانسیس گالتن برای اولین بار اصطلاح "رگرسیون" را تعریف، و نظریه همبستگی دو متغیره را کامل کرد. یک دهه بعد، کارل پیرسن، شاخص r پیرسن را که هنوز هم برای اندازه‌گیری همبستگی به کار می‌رود، ارائه داد. مقاله ما برای بزرگداشت صدمین سال اولین بحث گالتن از رگرسیون و همبستگی نوشته شده است. مقاله را با تاریخچه کوتاهی آغاز می‌کنیم و سپس ۱۳ فرمول مختلف را می‌آوریم که تعریفهای محاسباتی و مفهومی متفاوتی از آنند. هسرفرمول راهی برای اندیشیدن درباره این شاخص، از دیدگاه جبری، هندسی، و مثلثاتی ارائه می‌کند. نشان می‌دهیم که r پیرسن (یا تابعهای ساده r) را می‌توان به صورتهای مختلف، نظیر نوعی خاص از میانگین، نوعی خاص از واریانس، نسبت دو میانگین، نسبت دو واریانس، شیب یک خط، کسینوس یک زاویه، و مماس بزرگ بیضی در نظر گرفت، و همچنین نشان می‌دهیم که ممکن است از دیدگاههای جالب دیگری نیز به آن نگریست.

مقدمه

ما این روزها در اواسط "دهه بعد از صد سالگی" همبستگی و رگرسیون هستیم. دستاوردهای تجربی و نظری که موجب شد رگرسیون و همبستگی به صورت مطالب آماری تعریف شوند در ۱۸۸۵ به وسیله سرفرانسیس گالتن ارائه شدند. سپس، کارل پیرسن در ۱۸۹۵، r پیرسن را انتشار داد. این مقاله که درباره ضریب همبستگی پیرسن گفته می‌کند، هم زیر بنا و هم چند استنباط از آن را مورد بحث قرار می‌دهد که برای مدرسین آمار مفید است.

ما مقاله را با تاریخچه‌ای کوتاه از گسترش همبستگی و رگرسیون آغاز و به دنبال آن، راههای تعبیر ضریب همبستگی را با تفصیل بیشتر مرور می‌کنیم. این مرور نشان می‌دهد که همبستگی به صورت شاخصی کلی در آمده است که استنباطهای گوناگونی از آن می‌شود. مع هذا، گذشت زمان بر این شاخص صد ساله چندان تأثیری نداشته



شکل ۱. اولین نمودار پراکنش دو متغیره (از گالتن ۱۸۸۵)

یکی است]. گالتن، با همیاری هامیلتن دیکسن، ریاضیدانی از کمبریج، توانست فرمولی نظری برای توزیع نرمال دو متغیره به دست آورد. این فرمول، از نظر ریاضی، به موضوعی که گاوس و براوه نیم قرن قبل روی آن کار کرده بودند رسمیت داد. پیرسن [۴۵، ص ۳۷] اظهار داشت که "در ۱۸۸۵، گالتن نظریه همبستگی دو متغیره را کامل کرد."

در سالهای بعد از ۱۸۸۵، چند پیشامد دیگر بر اهمیت ریاضی اثر ۱۸۸۵ گالتن افزودند. در ۱۸۸۸، گالتن متذکر شد که r ، دقت "همبستگی" را اندازه می گیرد و (هر چند ایده همبستگی منفی هنوز به ذهن او نرسیده بود) اظهار کرد که r نمی تواند از ۱ بیشتر باشد. هفت سال بعد، پیرسن فرمول ریاضی را که هنوز متداولترین فرمول برای اندازه گیری همبستگی است، یعنی ضریب همبستگی گشتاور حاصلضربی پیرسن را ارائه داد. از دیدگاه تاریخی، به نظر مناسبتر می رسد که نام این شاخص معروف، r گالتن-پیرسن باشد. تحولهای مهم در حکایت همبستگی و رگرسیون در جدول ۱ خلاصه شده اند.

اینک، پس از یک قرن، دانشمندان معاصر اغلب ضریب همبستگی را مطلبی بدیهی و مسلم می دانند. به این مطالب توجه نمی شود که قبل از گالتن و پیرسن، تنها وسیله استقرار بستگی بین متغیرها، یافتن ارتباط علّی بود. حتی راهی برای بحث درباره پیوند بین متغیرهایی که فاقد بستگی علت و معلولی بودند وجود نداشت، چه رسد به

پدیده دیگری شود، پدیده اول علت یا معلول پدیده دوم است، یا از طریق واقعیتی علّی بدان مربوط است. میل برای معتبر بودن استنباط علّی سه شرط پیشنهاد کرد [۶]. اولاً، علت باید از نظر زمانی مقدم بر معلول باشد. ثانیاً، علت و معلول باید به هم مربوط باشند. ثالثاً، توضیحات به ظاهر موجه دیگر باید رد شوند. بنا بر این تفکیک پذیری همبستگی و علّیت، و توصیف اولی به صورت شرطی لازم و نه کافی برای دومی، در نظام جافاناده فلسفه و نظام کم سابقه زیست سنجی تقریباً به طور همزمان مشخص داده شد.

تا ۱۸۸۵، زمینه برای ارائه چند اثر مهم فراهم آمد. در طول آن سال، گالتن رئیس بخش مردم شناسی انجمن بریتانیا بود. وی در خطابه ای که به مناسبت انتصابش ایراد نمود، اولین بار به رگرسیون عنوان تمیم "قانون Reversion" را اطلاق کرد. کمی بعد در همان سال [۹]، خطابه مزبور را به همراه اولین نمودار پراکنش دو متغیره که همبستگی را نشان می داد منتشر کرد (شکل ۱).

وی در این نمودار، فراوانی ترکیبهای بلندی قد فرزندان و بلندی قد والدین را به نمایش گذاشت. وقتی وی نتایج را هموار کرد و خطوطی بر نقاط با فراوانی برابر گذراند، دریافت که "خطوط مابعد درایه های هم مقدار، یک سری از بیضیهای هم مرکز و مشابه تشکیل می دهند." این اولین نمایش تجربی خمهای تک چگالی از توزیع نرمال دو متغیره بود [فراوانی روی هر یک از این خمها

جدول ۱. پیشامدهای برجسته تاریخچه همبستگی و رگرسیون

تاریخ	شخص	پیشامد
۱۸۲۳	کارل فریدریش گاوس، ریاضیدان آلمانی	رویة نرمال N متنفر تصادفی همبسته را ارائه داد.
۱۸۴۳	جان استوارت میل، فیلسوف بریتانیایی	پنج قانون استقرار، از جمله تغییرات ملازم را مطرح کرد.
۱۸۴۶	اوگوست براده، افسر نیروی دریایی فرانسه و ستاره‌شناس	بسا اشاره به "یک همبستگی"، روی توزیع نرمال دو متغیره کار کرد.
۱۸۶۸	چارلز داروین دایمی زاده گالتن، طبیعیدان بریتانیایی	تمام اجزاء یک نظام... بهم مربوط یا همبسته اند.
۱۸۷۷	سرفرانسیس گالتن، بریتانیایی، اولین متخصص زیست‌سنجی	برای اولین بار از "Reversion"، سلف رگرسیون، بحث کرد.
۱۸۸۵	سرفرانسیس گالتن	برای اولین بار به "رگرسیون" اشاره کرد. نمودار براکتش دو متغیره با خمهای تک‌چکالی نرمال دو متغیره، اولین نمودار همبستگی، را انتشار داد. نظریة همبستگی نرمال دو متغیره را کامل کرد [۲۰].
۱۸۸۸	سرفرانسیس گالتن	r را به‌طور مفهومی تعریف کرد، و کران بالایی آن را مشخص نمود.
۱۸۹۵	کارل پیرسن، آماردان بریتانیایی	ضریب همبستگی گشتاور حاصلضرب بی (گالتن-پیرسن) را تعریف کرد.
۱۹۲۰	کارل پیرسن	"یادداشت‌هایی بر تاریخچه همبستگی" را نوشت.
۱۹۸۵		صد سالگی رگرسیون و همبستگی

اندازه‌گیری آنها، امروزه، ضریب همبستگی و معادله رگرسیون همتای آن‌ها در بسیاری از زمینه‌ها برای آزمایش‌های مبتنی بر مشاهدات، یک ابزار آماری اصلی است. کارول [۳، ص ۳۴۷] در خطایه‌ای که به مناسبت انتصابش به ریاست انجمن روانسنجی ایراد کرده، ضریب همبستگی را "یکی از متداولترین ابزارهایی که در روانسنجی به کار می‌رود... و شاید یکی از متداولترین ابزارهایی که بد به کار گرفته شده است" خواند. در تجزیه عاملی، در مدل‌های ژنتیکی رفتاری، در مدل‌های معادله‌های ساختاری (مثلاً LISREL)، و در دیگر روش‌های وابسته، ضریب همبستگی به عنوان واحد اساسی داده‌ها به کار می‌رود.

بحث ما حول ضریب همبستگی گشتاور حاصلضرب پیرسن دور می‌زند. ۳ پیرسن اولین معیار رسمی اندازه‌گیری همبستگی بود، و هنوز هم معیاری برای بستگی است که کاربرد فراوان دارد. بدون شک، بسیاری از شاخص‌های همبستگی "رقیب"، در واقع حالت‌های خاصی از فرمول پیرسن‌اند. پی اسپیرمن، همبستگی نقطه‌ای مربوط به دوسری از داده‌ها، و ضریب فی، مثال‌هایی هستند که هر یک با به کار بردن ۳ پیرسن در مورد انواع خاصی از داده‌ها قابل محاسبه است (مثلاً، ر.ک. [۱۵]).

آنها را معرفی خواهیم کرد. به پیروی از پیرسن، تأکید ما بر ضریب همبستگی به عنوان یک شاخص محاسباتی است که برای اندازه‌گیری پیوند دو متغیره به کار می‌رود. از نظر آماری، درک عمیق‌تری از همبستگی مستلزم توجه به مدل نمونه‌گیری است که فرض می‌شود زیربنای مشاهدات باشد (مثلاً، ر.ک. [۳] و [۱۴])، و نیز مستلزم درک تعمیم همبستگی به همبستگی چندگانه و جزئی است، اما در اینجا تأکید ما بر همبستگی جنبه بنیادینتری دارد. اولاً، توجه اصلی را به وضعیت‌های دو متغیره محدود می‌کنیم. ثانیاً، اکثر تعبیرهای ما آزاد-توزیع‌اند، زیرا محاسبه همبستگی نمونه‌ای نیاز به هیچ فرضی درباره‌ی جامعه ندارد (ر.ک. [۱۶]). برای بررسی مسائل زیادی که به کاربرد استنباطی r (مثلاً، محدود کردن، و کوچک کردن آن) مربوط‌اند خواننده را به روش‌های دیگری ارجاع می‌دهیم (مثلاً، ر.ک. [۱۳]). برای استنباط این اصل‌ترین معیار بستگی دو متغیره سیزده راه مختلف معرفی می‌کنیم. ادعا نمی‌کنیم که این مقاله تمام تعبیرهای ممکن ضریب همبستگی را شامل می‌شود. محققاً تعبیرهای دیگری هم وجود دارند، و یقیناً تعبیرهای جدیدی نیز پیشنهاد خواهند شد.

۱. همبستگی به صورت تابعی از اندازه‌های خام و میانگین‌ها

مطلب را با کمی چاشنی آموزشی عرضه می‌کنیم. در نظر اول، معیار همبستگی، ساده و سرراست است. اما، اختلاف‌های جزئی شگفت‌انگیزی در مفهوم ضریب همبستگی وجود دارند که برخی از

مطلب را با کمی چاشنی آموزشی عرضه می‌کنیم. در نظر اول، معیار همبستگی، ساده و سرراست است. اما، اختلاف‌های جزئی شگفت‌انگیزی در مفهوم ضریب همبستگی وجود دارند که برخی از

می‌کنند. در اینجا همبستگی به صورت تابعی از شیب یکی از دو خط رگرسیون و انحراف معیارهای دو متغیر بیان می‌شود. نسبت انحراف معیارها دارای این اثر است که واحد شیب رگرسیون را به واحد همبستگی تبدیل می‌کند. بنابراین همبستگی، یک شیب استاندارد شده است.

تعبیری مشابه، همبستگی را به عنوان شیب خط رگرسیون استاندارد شده مطرح می‌کند. وقتی دو متغیر خام را استاندارد می‌کنیم، انحراف معیارها برابر واحد می‌شوند و شیب خط رگرسیون به صورت همبستگی درمی‌آید. در این حالت، عرض از مبدأ خط صفر است، و خط رگرسیون به آسانی به صورت

$$\hat{z}_Y = r z_X \quad (۲.۳)$$

بیان می‌شود. از این تعبیر آشکار است که همبستگی، برای پیشگویی واحدهای متغیر استاندارد شده Z_Y ، واحدهای متغیر استاندارد شده Z_X را از نو مقیاس بندی می‌کند. توجه کنید که شیب رگرسیون z_Y روی z_X ، خط رگرسیون را ملزم می‌کند که در ناحیه سایه خورده شکل ۲، بین نیمسازهای محورهای مختصات بیفتد. همبستگیهای مثبت ایجاب می‌کنند که خط از ربعهای اول و سوم بگذرد؛ همبستگیهای منفی ایجاب می‌کنند که خط از ربعهای دوم و چهارم عبور کند. زاویه خط رگرسیون z_Y روی z_X با محور Z_X ، برابر زاویه خط رگرسیون z_Y روی z_X با محور X است، و این خط همان طور که نشان داده ایم در ناحیه سایه نخورده شکل ۲ می‌افتد.

۴. همبستگی به صورت میانگین هندسی دوشیب رگرسیون

همبستگی را ممکن است به صورت تابعی همزمان از دوشیب خطهای رگرسیون استاندارد شده، $b_{Y.X}$ و $b_{X.Y}$ بیان کرد. در واقع این تابع، میانگین هندسی است، و اولین تعبیر از چند تعبیر r است که آن را به عنوان نوعی خاص از میانگین بیان می‌کند:

$$r = \pm \sqrt{b_{Y.X} b_{X.Y}} \quad (۱.۴)$$

این رابطه را می‌توان از معادله‌های (۱.۳) نتیجه گرفت: جمله‌های دوم و سوم برابریها را در هم ضرب می‌کنیم تا r^2 به دست آید، آنگاه انحراف معیارها را حذف می‌کنیم، و از دو طرف جذر می‌گیریم.

تعمیمی از این تعبیر وجود دارد که متضمن رگرسیون چندمتغیره است. وقتی $B_{Y.X}$ و $B_{X.Y}$ ، ماتریسهای ضرایب رگرسیون مربوط به دو مجموعه از متغیرها، معلوم اند، جذرها و ویژه مقادیر حاصل ضرب این ماتریسها، همبستگیهای کانونیک این دو مجموعه از متغیرها هستند. وقتی یک تک متغیر X و یک تک متغیر Y وجود داشته باشند، این مقادیر به ضریب همبستگی ساده تبدیل می‌شوند.

۵. همبستگی به صورت جذر نسبت دو واریانس

گاه از همبستگی، به این دلیل که مقدارش تعبیر روشنی ندارد، انتقاد می‌شود. این انتقاد با مربع کردن همبستگی و با تعبیر زیر تخفیف می‌یابد. شاخص مربع شده را اغلب ضریب تعیین می‌خوانند. مقدار این ضریب را ممکن است به عنوان کسری از واریانس یکی از متغیرها تعبیر کرد ([۱۸] را، برای آگاهی از بحثی که مربوط به چند

پیرسن در ۱۸۹۵، فرمولی ریاضی برای این معیار مهم ارائه داد:

$$r = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{[\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2]^{1/2}} \quad (۱.۱)$$

این فرمول، یا صورت‌های ساده جبری دیگر آن، فرمولهای متداولی در کتابهای درسی آمار مقدماتی هستند. در صورت این کسر، اندازه‌های خام هر متغیر، به وسیله کم کردن میانگین متغیر از آن، "مرکزی" شده‌اند و مجموع حاصلضربهای مقاطع متغیرهای مرکزی شده به دست آمده است. مخرج کسر، مقیاس متغیرها را برای اینکه واحدهای یکسان داشته باشند تبدیل می‌کند. بنا بر این معادله (۱.۱)، r را به صورت مجموع حاصلضربهای مقاطع دو متغیر مرکزی شده و استاندارد شده توصیف می‌کند. با استفاده از نابرابری کوشی-شوارتس، می‌توان نشان داد که قدر مطلق صورت از مخرج بیشتر نیست (مثلاً، [۱۴، ص ۸۷] را ببینید). بنا بر این برای r ، حدود ± 1 به دست می‌آیند. برای منظورهای محاسباتی، چندین تبدیل ساده جبری این فرمول را می‌توان به کار برد.

۴. همبستگی به صورت کوواریانس استاندارد شده

کوواریانس، شبیه همبستگی، معیار پیوند خطی بین متغیرهاست. کوواریانس روی مجموع حاصلضربهای مقاطع متغیرهای مرکزی شده، که مقیاس آنها تبدیل نیافته است، تعریف می‌شود. هر چند در کتابهای درسی مقدماتی کوواریانس را اغلب نسبیده می‌گیرند ولی واریانس (که از آن گفتگو می‌شود)، واقعاً حالتی خاص از کوواریانس است. بدین معنا که واریانس، کوواریانس یک متغیر با خود آن متغیر است. کوواریانس دو متغیر، در جامعه‌های نامتناهی، دارای کرانهای مشخصی نیست، و در نمونه کرانهای نامعین (و تعبیری نامناسب) دارد. بنابراین، کوواریانس اغلب معیار توصیفی مفیدی برای پیوند نیست، زیرا مقدار آن به مقیاسهای اندازه گیری X و Y بستگی دارد. ضریب همبستگی، از تغییر مقیاس کوواریانس به دست می‌آید:

$$r = \frac{s_{XY}}{s_X \cdot s_Y} \quad (۱.۲)$$

که در آن s_{XY} ، کوواریانس نمونه‌ای است و s_X و s_Y انحراف معیارهای نمونه‌ای هستند. وقتی کوواریانس به دو انحراف معیار تقسیم شود، برد کوواریانس به بازه $(-1, +1)$ محدود می‌شود. لذا، تعبیر همبستگی به عنوان معیار بستگی معمولاً ساده‌تر از تعبیر کوواریانس به عنوان معیار بستگی است (و به وسیله آن همبستگیهای مختلف با سهولت بیشتر مقایسه می‌شوند).

۳. همبستگی به صورت شیب استاندارد شده خط رگرسیون

بستگی همبستگی و رگرسیون را می‌توان به صورتی ساده‌تر یا رابطه

$$r = b_{Y.X} \left(\frac{s_X}{s_Y} \right) = b_{X.Y} \left(\frac{s_Y}{s_X} \right) \quad (۱.۳)$$

نمایش داد که در آن، $b_{Y.X}$ و $b_{X.Y}$ شیبهای خطهای رگرسیون اند که به ترتیب Y را از روی X یا X را از روی Y پیشگویی

مقتارن وجود دارد. تعبیر بعدی، که باز تعبیری مثلثاتی است، اساساً ارزش مفهومی بیشتری دارد.

۸. همبستگی به صورت تابعی از زاویه بین دو بردار متغیر

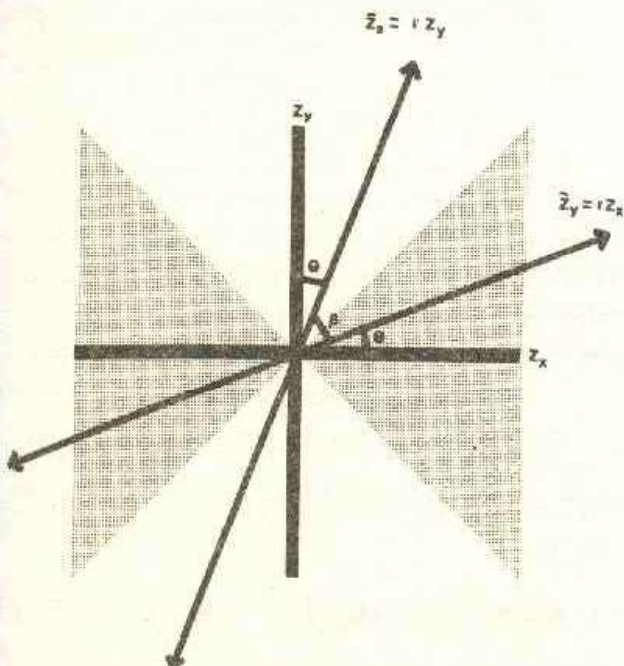
مدل هندسی استاندارد برای نمایش تصویری بستگی بین متغیرها، نمودار پراکنش است. در این مدل، مشاهدات را به صورت نقطه‌هایی از فضای نمایش می‌دهیم که با محورهای متغیرها تعریف می‌شود. صورت "وارونه" این فضا را که معمولاً به "فضای شخصی" موسوم است، می‌توان با فرض اینکه هر محور مشاهده‌ای را نمایش دهد، تعریف کرد. این فضا دو نقطه را در بردارده یکی برای هر متغیر - که نقاط انتهایی بردارهای واقع در این فضای (بالقوه) کلان بعدی را تعریف می‌کند. گرچه چند بعدی بودن این فضا مانع تجسم آن است، اما دو بردار متغیر، یک زیر فضای دو بعدی را تعریف می‌کنند که به آسانی قابل تجسم است.

اگر بردارهای متغیر مبتنی بر متغیرهای مرکزی شده باشند، آنگاه همبستگی با α که زاویه بین بردارهای متغیر است بستگی سرراست دارد [۳۱]:

$$r = \cos(\alpha) \quad (1.8)$$

وقتی زاویه α است، بردارها روی یک خط می‌افتند و در نتیجه $\cos(\alpha) = \pm 1$. وقتی زاویه 90° است، بردارها برهم عمودند و $\cos(\alpha) = 0$. (راجرز، نایواندر، توتاکر [۳۲] بستگی بین بردارهای متغیر متعامد و ناهمبسته را در فضای شخصی نشان داده‌اند.)

برای تجسم همبستگی، مشاهده یک زاویه خیلی ساده‌تر از مشاهده چگونگی گرد آمدن نقاط حول خط رگرسیون است. به عقیده ما، این تعبیر در مقایسه با سایر تعبیرها آسانترین راه "ملاحظه" اندازه



شکل ۳. تصویر هندسی همبستگی دو متغیر برای متغیرهای استاندارد شده.

تعبیر از ضریب تعیین است، ببینید). مجموع مربعات کل (SS) برای Y را ممکن است به دو مجموع مربعات حاصل از رگرسیون (رگرسیون SS) و مجموع مربعات ناشی از خطا (خطا SS) افراز کرد. نسبتی از کل تغییرات Y که از تغییرات X ناشی می‌شود، نسبت رگرسیون SS به کل SS است، و r ، جذر این نسبت است:

$$r = \sqrt{\frac{\sum(Y_i - \hat{Y}_i)^2}{\sum(Y_i - \bar{Y})^2}} = \sqrt{\frac{SS_{\text{رگرسیون}}}{SS_{\text{کل}}}}$$

هم ارز آن، اگر صورت و مخرج این معادله را به $(N-1)^{1/2}$ تقسیم کنیم، r مساوی با جذر نسبت واریانسها (یا نسبت انحراف معیارها)ی متغیرهای پیش‌بینی شده و مشاهده شده می‌شود:

$$r = \sqrt{\frac{s_Y^2}{s_Y^2}} = \frac{s_Y^2}{s_Y^2} \quad (1.5)$$

(توجه کنید که s_Y^2 بر آوردی اریب از σ_Y^2 است، در حالی که s_Y^2 بر آوردی نااریب است). این تعبیر، تعبیری است که انگیزه استنباط اولیه پیوسن از شاخص همبستگی بوده است. ([۱۵]، ص ۴) را ببینید). همبستگی به عنوان نسبت دو واریانس را می‌توان با تعبیر دیگری از همبستگی به صورت نسبت دو میانگین (منسوب به گالان) مقایسه کرد. ما این تعبیر را در بخش ۱۳ ارائه خواهیم داد.

۹. همبستگی به صورت میانگین حاصلضرب متقاطع متغیرهای استاندارد شده

راه دیگر تعبیر همبستگی به صورت میانگین (بخش ۴ را ببینید)، بیان آن به صورت متوسط حاصلضرب متقاطع متغیرهای استاندارد شده است:

$$r = \frac{\sum z_x z_y}{N} \quad (1.6)$$

معادله (۱.۶) را می‌توان مستقیماً از تقسیم صورت و مخرج معادله (۱.۱) بر حاصلضرب انحراف معیارهای دو نمونه به دست آورد. چون میانگین یک توزیع، اولین گشتاور آن است، این فرمول، به کار بردن مفهوم "گشتاور حاصلضرب" به جای ضریب همبستگی را توجیه می‌کند. دو توصیف بعدی، متضمن تعبیرهای مثلثاتی از همبستگی هستند.

۷. همبستگی به صورت تابعی از زاویه بین دو خط رگرسیون استاندارد شده

همان‌طور که در بخش ۳ اشاره شد، دو خط رگرسیون استاندارد شده نسبت به هر نیمساز متقارن‌اند. فرض کنیم زاویه بین دو خط β باشد (شکل ۲ را ببینید). در این صورت

$$r = \sec(\beta) \pm \tan(\beta) \quad (1.7)$$

برهان ساده‌ای از این بستگی را در اختیار داریم. معادله (۱.۷) به‌طور شهودی واضح نیست و برای هدفهای محاسباتی یا مفهومی هم به اندازه دیگرها مفید نیست. مقدار r در این رابطه نشان می‌دهد که بین همبستگی و ناصله زاویه‌ای، دو خط رگرسیون یک بستگی

۹۰. همبستگی بر آورد شده از روی قاعده بادکنکی

این تعبیر منسوب به شاتیبون [۴] است. او پیشنهاد کرد که پیرامون نمودار پراکنش يك بستگی دو متغیره، "بادکنکی" رسم کنیم. این بادکنک در واقع يك بیضی تقریبی است که از روی آن دو اندازه h و H به دست می آیند (شکل ۳ را ببینید). h ، قطر قائم بیضی است که از مرکز توزیع واقع بر محور X می گذرد؛ H ، برد تغییرات عرضی بیضی روی محور Y است. شاتیبون نشان داد که همبستگی را می توان به طور تقریبی به صورت

$$r = \sqrt{1 - \left(\frac{h}{H}\right)^2} \quad (1.10)$$

محاسبه کرد. وی با فرض توزیع نرمال دو متغیره و توزیع یکنواخت دو متغیره، درباره کارایی این شیوه محاسباتی تقریبی و سهل الوصول، توجیهی نظری ارائه داد. او همچنین مثالهایی چند معرفی کرد که در آنها این تکنیک به خوبی قابل استفاده است. پیشنهاد اغوا کننده ای که وی عرضه کرد این است که "قاعده بادکنکی" را می توان برای بنای تقریبی يك بستگی دو متغیره با همبستگی معینی به کار برد. بدین طریق که بیضی می کشیم که r مطلوب را تولید کند و آنگاه سراسر بیضی را به طور یکنواخت از نقطه ها پر می کنیم. توماس [۲۵] يك "نوموگراف جیبی" ساخت که نواری $5" \times 3"$ است که می توان از آن برای "مشاهده" يك بستگی دو متغیره استفاده کرد و همبستگی مبتنی بر قاعده بادکنکی را بر آورد نمود.

۹۱. همبستگی در ارتباط با بیضیهای دو متغیره تک چگالی

دو نویسنده مختلف، تعبیرهایی از r را که در ارتباط با بیضیهای دو متغیره تک چگالی اند پیشنهاد کرده اند. توجه کنید که این بیضیها صورتهای رسمیت "بادکنک" بخش ۹۰ اند و ساختارهای هندسی هستند که گالتن در داده های تجریش مشاهده کرده است (شکل ۱ را ببینید). شاتیبون [۵] ردهای از توزیعهای دو متغیره (شامل نرمال، یکنواخت، و آمیزه های از یکنواخت) را ارائه داده است که دارای خمهای تک چگالی بیضی شکل هستند. با معلوم بودن همبستگی جامعه ای، برای هر ثابت مثبت يك بیضی وجود دارد. بادکنکی که پیرامون يك نمودار پراکنش کشیده می شود به ازای ثابت مثبت بزرگی یکی از این بیضیها را تقریب می کند. اگر متغیرها استاندارد شده باشند، آنگاه مرکز این بیضیها در مبدأ است. قطر اطول برای $\rho > 0$ بر نیمساز ربع اول وسوم و برای $\rho < 0$ بر نیمساز دوربع دیگر می افتد.

مارکس [۱۴]، با محاسبه ای ساده نشان داد که شیب خط مماس در $z_Y = 0$ ، برابر با همبستگی است. شکل ۳ این خط مماس را نشان می دهد که شیب آن مساوی با r است. وقتی همبستگی ۰ است، بیضی يك دایره است و مماس دارای شیب ۰ است. وقتی همبستگی ۱ است، بیضی به خطی مستقیم میل می کند که همان نیمساز (با شیب ۱) است. توجه کنید که چون تمام بیضیهای تک چگالی موازی اند، تعبیر ارائه شده به انتخاب بیضی بستگی ندارد. همچنین شایان توجه است که شیب خط مماس در $z_Y = 0$ همان شیب خط رگرسیون

همبستگی است؛ زیرا می توان مستقیماً اندازه زاویه بین دو بردار را مشاهده کرد. اما، این فضای "وارونه"، که اجازه می دهد r را به صورت کسینوس يك زاویه نمایش دهیم به عنوان يك ابزار تغییر کننده، نسبتاً از نظر دور مانده است. چند تعبیر مربوط به تجزیه عاملی، نمایشهای هندسی تحلیل رگرسیون چندگانه در پیرا و اسمیت [۷] صص (۲۰۱-۲۰۳)، و هک و ساندر [۹۱، ص ۵۲]، از جمله استثنایهایی هستند که باید بر شمرد. فیشر نیز برای تفهیم بینشهای آماری استادانه اش کرا را از این فضا استفاده کرده است ([۹] را ببینید).

۹۲. همبستگی به صورت يك واریانس تجدید مقیاس شده از تفاضل بین اندازه های استاندارد شده

$Z_Y - Z_X$ را به عنوان تفاضل بین متغیرهای استاندارد شده X و Y برای هر مشاهده تعریف می کنیم. در این صورت

$$r = 1 - \frac{s^2(Z_Y - Z_X)}{2} \quad (1.9)$$

برای نشان دادن این رابطه می توانیم با واریانس تفاضلی

$$s^2_{Z_Y - Z_X} = s^2_X + s^2_Y - 2r s_X s_Y$$

شروع کنیم. چون وقتی متغیرها را استاندارد می کنیم، انحراف میانه ها و واریانسها برابر واحد می شوند، می توانیم به آسانی رابطه بالا را نسبت به r حل کنیم و معادله (۱.۹) را به دست آوریم.

توجه به این نکته جالب است که در این معادله، چون همبستگی به بازه از -1 تا 1 محدود است، واریانس این اندازه تفاضلی به بازه از 0 تا 4 محدود می شود. بنا بر این واریانس تفاضلی اندازه های استاندارد شده هرگز از 4 تجاوز نمی کند. وقتی همبستگی برابر -1 است، واریانس به کران بالایی می رسد. می توانیم r را به صورت واریانس مجموعی از متغیرهای استاندارد شده نیز تعریف کنیم:

$$r = \frac{s^2(Z_Y + Z_X)}{2} - 1 \quad (2.9)$$

در اینجا، واریانس مجموع نیز از 0 تا 4 تغییر می کند، و موقعی که همبستگی برابر $+1$ است به کران بالای خود می رسد. مقدار این نهمین تعبیر است که نشان می دهد همبستگی، تبدیل خطی نوع خاصی از واریانس است. لذا، با معلوم بودن همبستگی، می توانیم مستقیماً واریانس مجموع یا واریانس تفاضل متغیرهای استاندارد شده را تعریف کنیم، و برعکس.

تمام نه تعبیر قبلی ضریب همبستگی ماهیتاً جبری و مثلثاتی هستند. تا اینجا درباره ماهیت توزیعهای تک متغیره یا دو متغیره X و Y هیچ فرضی نشد. در تعبیرهای پایانی، نرمال بودن دو متغیره را فرض خواهیم کرد. توجه خود را همچنان به صورتهای مفهومی و محاسباتی r معطوف می داریم، اما چند تعبیر آخری را بر این فرض مشترک درباره توزیع جامعه، استوار می کنیم.

همبستگی را می‌توان به عنوان معیاری برای قوت یک اثر تیماری، در مقابل معنادار بودن یک اثر به کار برد. در این وضعیت آزمون معنادار بودن r ، همان آزمون t معمولی را به دست می‌دهد. بنابراین، r به وضوح می‌تواند به همان خوبی که معیاری از پیوند را در وضعینهای مشاهده‌ای فراهم می‌کند به عنوان یک آماره آزمون نیز در یک آزمایش طرح شده به کار رود.

در حالت تحلیل واریانس با گروههای بیشتری یا عملهای چندگانه، تعمیم این بستگی، ضریبهای همبستگی چندگانه مربوط به اثرهای اصلی و اثرهای متقابل در آزمایشهای پیچیده تر را تعریف می‌کند. مثلاً، در حالت تحلیل واریانس یکطرفه با k گروه و مجموع کل N آزمودنی، مربع همبستگی چندگانه بین متغیر وابسته و متغیرهای ماتریس طرح، از طریق فرمول زیر به آماره F مربوط می‌شود. [۹۳]، ص:

$$R^2 = \frac{F(k-1)}{[F(k-1) + (N-k)]}$$

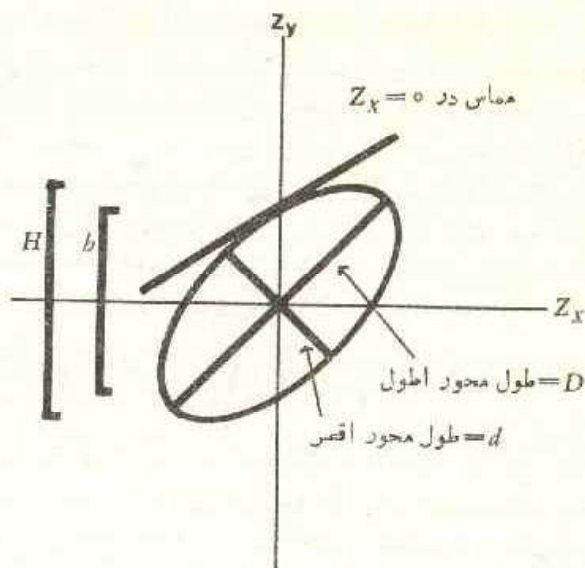
۱۳. همبستگی به صورت نسبت دو میانگین

این سومین تعبیر همبستگی است که در ارتباط با میانگینهاست (بخشهای ۴ و ۶ را ببینید). این تعبیر، پایانی مناسب بر مقاله ماست، زیرا اولین بار گالتن آن را پیشنهاد کرده است. بدعلاوه، قدیمترین استنباط گالتن از همبستگی و محاسبه آن، برای تعبیر مبتنی بود. نایواندر و پرایس [۱۷] صورت پیشرفته تری از این تعبیر را ارائه داده‌اند.

اما گالتن، طبیعی بود که توجه خود را بر همبستگی به صورت نسبت میانگینها متمرکز کند؛ زیرا وی به سؤالی از این قبیل علاقه مند بود که چگونه باید متوسط قد پدرهایی را که به طور غیر معمول بلندقدند با متوسط قد پسرانشان مقایسه کرد. در بحث زیر به جای نماد نمونه‌ای، از نماد جامعه‌ای استفاده می‌شود، زیرا این تنها در حد (حجم نمونه‌ای افزایشی) است که عبارت نسبت میانگینها همان مقادیر r پیرس را به دست می‌دهد.

وضعیتی مشابه مورد توجه گالتن بود در نظر می‌گیریم. فرض کنید X متغیری باشد که بهره هوشی مادر را نشان دهد، و Y متغیری باشد که بهره هوشی بزرگترین فرزند این مادر را نمایش دهد. بدعلاوه فرض کنید که $\mu(X)$ و $\mu(Y)$ برابر σ و انحراف معیارهای $\sigma(X)$ و $\sigma(Y)$ برابر یک باشند. حال مقداری بد دلخواه بزرگ از X (مثل X_c) را بر می‌گزینیم، و میانگین بهره هوشی مادرهایی را که بهره هوشی آنها از X بزرگتر است محاسبه می‌کنیم. این میانگین را با $\mu(X|X > X_c)$ نشان می‌دهیم. این، متوسط بهره هوشی مادرانی است که بهره هوشی آنها بزرگتر از X_c است. سپس متوسط اندازه‌های Y ، یعنی بهره هوشی بزرگترین فرزند این مادرهای خاص را محاسبه می‌کنیم. این میانگین را با $\mu(Y|X > X_c)$ نمایش می‌دهیم. یعنی، متوسط بهره هوشی بزرگترین فرزند مادرانی که بهره هوشی آنها بزرگتر از X_c است. در این صورت، می‌توان نشان داد که

$$r = \frac{\mu(Y|X > X_c) - \mu_Y}{\mu(X|X > X_c) - \mu_X} = \frac{\mu(Y|X > X_c)}{\mu(X|X > X_c)} \quad (1.13)$$



شکل ۳. همبستگی در ارتباط با تابعی از بیضیهای تک چکالی

استاندارد شده است (بخش ۳ را ببینید).

شیلینگ [۲۳] نیز برای رسیدن به بستگی مشابهی، از این چارچوب استفاده کرده است. فرض کنید متغیرها استاندارد شده باشند به قسمی که، مثل قبل، مرکز بیضیها در مبدأ باشد. اگر D ، درازای قطر اطول یکی از بیضیهای تک چکالی و d ، درازای قطر اقصر آن باشد، آنگاه

$$r = \frac{(D^2 - d^2)}{(D^2 + d^2)} \quad (1.11)$$

این محورها نیز در شکل ۳ رسم شده‌اند، و تعبیر، مثل قبل، به انتخاب بیضی بستگی ندارد.

۱۴. همبستگی به صورت تابعی از آماره آزمونی از آزمایشهای طرح شده

تعبیرهای قبلی r بر متغیرهای کمی مبتنی بودند. دوازدهمین نمایش همبستگی، بستگی آن را با آماره آزمونی از آزمایشهای طرح شده نشان می‌دهد که در آنها یکی از متغیرها (متغیر مستقل)، متغیری رسته‌ای است. این تعبیر، ساختگی بودن نمایش آزمایشها از همبستگی را در بحث طرح آزمایشها ثابت می‌کند. در واقع، فیشر (۱۹۵۲) در اصل، تحلیل واریانس را بر حسب ضریب همبستگی درون-رده‌ای ارائه داد (ر. ک. [۱]).

فرض می‌کنیم آزمایش طرح شده‌ای با دوشروط تیماری داشته باشیم. مدل آماری استاندارد برای آزمون تفاوت بین شرایط، آزمون t مربوط به دو نمونه مستقل است. اگر X به صورت یک متغیر دو حالتی معرف عضویت گروه، تعریف شود (0 ، اگر گروه ۱، ۱، اگر گروه ۲)، آنگاه همبستگی بین X و متغیر وابسته Y عبارت است از

$$r = \frac{t}{\sqrt{t^2 + n - 2}} \quad (1.12)$$

که در آن، n تعداد کل مشاهدات دو گروه تیماری است. این ضریب

