

افقیهای جدید آمار: تحلیل داده‌های نمادین

محمد رضا مشکانی^۱

چکیده

بیشرفت شگفت‌انگیز در محاسبات آماری با رایانه‌های سریع، موجب تحولاتی ژرف در اکثر علوم و از آن جمله در آمار شده است. تا چند دهه پیش، تحلیل مجموعه داده‌هایی که شامل بردارها یا ماتریس‌های با ابعاد زیاد بودند، مشکلی عمده به شمار می‌آمد، ولی امروزه چنین محاسباتی امری عادی به شمار می‌روند، اما تحلیل چنین داده‌های حجیم با ابزارها و روشهای سابق تا حدودی به معنای نادیده انگاشتن اطلاعات جالب توجهی است که در آنها نهفته‌اند. بنابراین، از چندی قبل، ضرورت ابداع روشهایی جدید که از این پیشرفت‌های محاسباتی بهره‌برگیرند، احساس شده است. برای بهره‌گیری از این امکانات جدید، باید نگرش ما به مجموعه داده‌ها تغییر یابد و مفاهیمی جدید تعریف شوند تا بر اساس آنها روشهای جدید بنا نهاده شوند. با تعریف داده‌های نمادین که می‌توانند رده‌هایی از بردارها یا ماتریس‌هایی اندازه‌گیری شده از هر فردی آماری باشند، به تعمیم جالبی می‌توان دست یافت.

در این مقاله، ابتدا ضرورت استفاده از داده‌هایی نمادین و تحلیل آنها را مورد بحث قرار می‌دهیم. سپس با معرفی مفاهیم و اصطلاحات مورد استفاده در این رشته به اجمال و همراه با مثال روشهایی از تحلیل داده‌های نمادین را عرضه می‌کنیم.

واژه‌های کلیدی: داده‌های نمادین، متغیرهای نمادین، داده کاوی، رده‌بندی، خوشه‌بندی، متغیرهای بازه-مقدار، چند مقدار.

۱. مقدمه

نیستند. زیرا در اغلب موارد، داده‌ها یک «کل» را تشکیل می‌دهند. مثل سوابق حسابهای کارتهای اعتباری همه کسانی که از کارت استفاده می‌کنند. سؤال مربوط دیگر آن است که آیا داده‌های یک نقطه زمانی مشخص را می‌توان از همان جامعه‌ای دانست که در زمان دیگر گردآوری می‌شوند. مثلاً، آیا داده‌های کارتهای اعتباری این هفته، از همان جامعه داده‌های هفته «بعد»ند؟ یا اگر دو مجموعه داده‌ها در هم ادغام شوند، باز هم متعلق به همان جامعه اولیه‌اند؟ دلیل عمده دیگر که روشهای شناخته شده را نمی‌توان به کار برد، صرفاً اندازه مجموعه داده‌هاست. برای مثال، فرض کنید n مشاهده در مورد p متغیر از هر فرد

با اختراع رایانه‌ها، انواع پایگاههای داده‌ای بزرگ و بسیار بزرگ امری عادی شده است. چیزی که عادی نیست، نحوه تحلیل داده‌های موجود در این پایگاهها و یا چگونگی استخراج اطلاعات مفید از درون این پایگاههای گسترده است، اما روشن است که در مواردی هم که نظریه و روش‌شناسی موجود به نظر می‌رسد که قابل کاربست باشد، استفاده روزمره از این گونه روشهای آماری اغلب مناسب نیست. دلیل‌های زیادی وجود دارند. به طور کلی، یکی از دلایل عمده آن است که این گونه مجموعه‌های داده‌ها واقعاً نمونه‌ای از یک جامعه

^۱ دانشکده علوم ریاضی، گروه آمار، دانشگاه شهید بهشتی

یا پیش قلب (مثل «۷۲ و ۶۰»)، فشار خون سیستولیک («۱۳۰ و ۱۲۰») برای مثلاً $n=120$ بیمار (یا مثلاً برای $n=1/2 \times 10^6$ بیمار). یا، ممکن است $n=20$ دانشجو داشته باشیم که به وسیله بافتنگار یا توزیع نمره‌هایشان برای چند متغیر ریاضی، فیزیک، آمار و غیره مشخص شده باشند. ممکن است پرندگان بر حسب رنگشان مشخص شوند. مثلاً پرنده ۱ = (سیاه)، پرنده ۲ = (زرد و آبی)، پرنده ۳ = (نیمی زرد و نیمی قرمز) ... یعنی، متغیر «رنگ» تنها یک رنگ ممکن را برای هر پرنده اختیار نمی‌کند، بلکه فهرستی از رنگ‌ها را اختیار می‌کند، یا حتی فهرستی را با نسبت‌های متناظر از هر رنگ را برای هر پرنده اختیار می‌کند. از سوی دیگر، نقطه داده‌ای {سیاه} ممکن است گردابه‌ای از پرندگان را که همگی سیاه‌اند، نشان دهد و نقطه {قرمز (۰/۶) و زرد (۰/۴)} گردابه‌ای از پرندگانی خواهد بود که ۴۰٪ زرد و ۶۰٪ قرمزند یا گردابه‌ای که از آنها ۴۰٪ کلاً زرد و ۶۰٪ کلاً قرمزند. مثالهایی بی‌شمار را می‌توان ذکر کرد. در راستایی دیگر، ممکن است پرنده خاصی را مدنظر نداشته باشیم، بلکه به مفهوم پرنده‌ای سیاه، یا پرنده‌ای زرد و قرمز علاقه‌مند باشیم. به همین ترتیب، می‌توانیم یک شرکت مهندسی را با پایگاه علمی آنکه مشتمل بر تجارب کارمندان است، مشخص کنیم. این گونه تجارب بیشتر با مفاهیم توصیف می‌شوند تا با داده‌های استاندارد و بدین طریق، مثالهایی از داده‌های نمادین‌اند. در مورد مجموعه‌های کوچک از داده‌های نمادین سؤال این است که تحلیل چگونه صورت می‌گیرد. در مورد مجموعه‌های بزرگ، سؤال نخست این است که رویکرد انتخابی برای خلاصه‌سازی چیست تا لزوماً به مجموعه‌ای کوچکتر خلاصه شوند. برخی روشهای خلاصه‌سازی الزاماً شامل داده‌های نمادین و تحلیل نمادین در قالبی از انواع قالبها هستند، در حالی که در برخی دیگر لازم نیست که آن گونه عمل شود.

در این مقاله، سعی بر آن است که مفاهیم و روشهای گوناگون ابداع شده تحت عنوان تحلیل داده‌های نمادین یا سایر عنوان‌ها را مرور کنیم. در حقیقت، این روشها تاکنون محدود به ابداع روش شناسی‌هایی بوده‌اند که داده‌ها را در قالبهایی معقول و قابل اداره سازمان دهند، کاری تا حدودی مشابه با تشکیل جدول‌های توزیع فراوانی و بافتنگارها و دیگر آماره‌های خلاصه پایه‌ای در دوره قبل از سال ۱۹۰۰ است. همین اقدامات خلاصه‌سازی بود که بنیاد آمار استنباطی را در قرن بیستم میلادی پایه‌گذاری کرد. در زیر مروری کوتاه بر روشهای آماری نمادین موجود می‌آوریم. موضوعی که سریعاً آشکار می‌شود این است

آماري داریم. در اغلب موارد، ماتریس داده‌های X که از بعد $n \times p$ است، به صورت $(X'X)^{-1}$ که از بعد $p \times p$ است، در محاسبه وارد می‌شود. اگر n از بعد صدها هزار و p از بعد صد یا بیشتر باشد، در عین حال که از حیث نظری، چنین کاری ممکن است، از حیث عملی کاری بسیار سنگین خواهد بود. حتی وقتی توانایی رایانه‌ها (مثلاً برای وارون کردن ماتریس‌ها در مدت زمانی معقول) افزایش یابند، در نتیجه، مجموعه داده‌های بسیار بسیار بزرگتر تولید خواهند شد و مساله تقریباً به همین شکل باقی خواهد ماند. بنابراین، در حالی که روشهای سنتی در مورد مجموعه‌هایی کوچک از داده‌ها که در گذشته وجه غالب بوده‌اند، به خوبی کار کرده است. اکنون انتظار از آمارشناسان تحلیل داده‌ها آن است که شیوه‌هایی را ابداع کنند که در مورد مجموعه‌های بزرگ داده‌های جدید خوب کار کنند. این شیوه‌ها، باید چنان باشند که ما را از اطلاعات زیر ساختی (یا دانش) نهفته در داده‌ها آگاه سازند.

یک رویکرد آن است که مجموعه‌های بزرگ داده‌ها را به طریقی خلاصه کنیم که مجموعه داده‌های خلاصه اندازه‌ای کوچک داشته باشد تا بتوانیم از عهده‌اش برآیم. بدین ترتیب، در مثال کارت اعتباری به عوض صورت زیر عملیات بانکی برای هر فرد (یا هر کارت) در طی زمان، خلاصه‌ای از این عملیات برای هر کارت (یا برای هر واحد زمان مثل هفته، ماه) تهیه کنیم. یک چنین خلاصه‌ای می‌تواند گروه‌بندی معاملات بر حسب میزان پول خرج شده (مثلاً ۱۰۰-۱۰ هزار ریال، غیره)، یا بر حسب خرید نوع کالا (مثلاً غذا، لباس، بنزین، غیره)، یا بر حسب پول خرج شده و نوع کالا (مثل «غذا-۲۰-۱۵»، «لباس-۳۰-۱۰»...) یا هر نوع خلاصه‌سازی دیگر باشد.

در هر یک از این مثالها، داده‌ها دیگر مقداری تک مثل داده‌های سنتی از قبیل: ۱۶۴۰۰ ریال کل هزینه کارت اعتباری، یا ۵۲ مورد معامله، برای هر نفر در هر واحد زمان ندارند. این گونه داده‌ها، مثالهایی از داده‌های نمادین‌اند. بدیهی است که کاربست روشی ساده مثل گرفتن میانگین حساسی در مورد آنها چندان مناسب ندارد. برای چنین داده‌های نمادین، باید از روشهای تحلیل داده‌های نمادین استفاده کرد.

در حالی که خلاصه‌سازی مجموعه‌هایی بزرگ از داده‌ها می‌تواند مجموعه‌هایی کوچک از داده‌ها مشتمل بر داده‌های نمادین را تولید کند، داده‌های نمادین اصالتاً در مورد هر اندازه از مجموعه داده‌ها، خواه کوچک یا بزرگ متمایزند. مثلاً، نامعقول نیست که داده‌هایی داشته باشیم که مشتمل بر متغیرهایی‌اند که در دامنه‌ای ثبت شده‌اند، مانند نبض

ضروری برای تحلیل‌های آماری است، در بخش ۳ بررسی می‌شود. در بخش ۵، به اجمال روشهای موجود برای تحلیل داده‌های نمادین را ارائه می‌کنیم. آنچه که از این بحث‌ها آشکار می‌شود آن است که داده‌های نمادین به طور فزاینده‌ای شایع می‌شوند و بنابراین نیازی مبرم به ابداع روش‌شناسی‌های آماری برای این گونه داده‌ها وجود دارد. همچنین، ملاحظه خواهد شد که در حال حاضر، روشهایی محدود وجود دارند و حتی برای آنهایی که وجود دارند، نیاز به تاسیس مبانی آماری-ریاضی مطرح است.

۲. مجموعه داده‌های نمادین

یک مجموعه داده‌ای از همان ابتدا ممکن است به صورت مجموعه‌ای از داده‌های نمادین تشکیل شود. به شق دیگر، ممکن است از ابتدا به صورت داده‌های کلاسیک فراهم آید، اما بعداً به صورت داده‌های نمادین سازمان داده شود تا به قالبی بیشتر قابل اداره در آید، به ویژه هنگامی که در ابتدا دارای اندازه‌ای بسیار بزرگ باشد.

در این بخش، مثالهایی از هر دو نوع داده‌های کلاسیک و داده‌های نمادین را عرضه می‌کنیم.

همچنین، نوعی نمادگذاری را معرفی می‌کنیم که مجموعه‌های داده‌های نمادین آماده برای تحلیل را توصیف می‌کند. این فرآیند وضعیت‌هایی را در برمی‌گیرد که مثلاً دو یا چند مجموعه از داده‌ها در هم ادغام می‌شوند یا گاهی که بر ویژگی‌های مختلفی از داده‌ها باید تأکید شود.

فرض کنید مجموعه‌ای از داده‌ها را داریم که مشتمل بر سوابق پزشکی افراد در کشور است. فرض کنید برای هر فرد، سابقه‌ای از متغیرهای مکان جغرافیایی، مانند ناحیه (شمال، شمال شرق، جنوب، ...)، شهر (بابلسر، سبزوار، کرمان، ...)، شهر یا روستا (بلی، خیر) و غیره وجود دارد. متغیرهای جمعیت شناختی از قبیل جنس، وضع تأهل، سن، اطلاعات درباره پدر و مادر (هنوز زنده یا نه)، برادر و خواهر، تعداد فرزندان، کارفرما، عرضه کننده خدمات تندرستی و غیره نیز در کار خواهند بود. متغیرهای پایه‌ای پزشکی می‌توانند شامل وزن، تپش نبض، فشار خون و غیره باشند. دیگر متغیرهای تندرستی (که برای آنها فهرست متغیرهای ممکن بی‌پایان است) شامل بروز ناراحتی‌ها و بیماریهای معین خواهند بود. همچنین، برای شیوع یا تشخیص یک

که تاکنون روش‌شناسی آماری اندکی برای صورت‌های مختلف داده‌های نمادین ابداع شده‌اند. اما به تعبیری دیگر، تحلیل اکتشافی داده‌ها که از سوی توکی و همکارانش (توکی، ۱۹۹۷) ابداع شدند، نویدی است از آنچه هم اکنون در حال بسط و گسترش است.

تحلیل اکتشافی داده‌ها، داده‌کاوی، کشف دانش در پایگاههای داده‌ها، آماره‌ها، داده‌های نمادین و مانند آنها اصطلاحاتی هستند که امروزه به فراوانی به چشم می‌خورند. تحلیل داده‌های نمادین اندیشه‌های به کار رفته در تحلیل اکتشافی سنتی را به داده‌های کلی‌تر و پیچیده‌تر گسترش می‌دهد.

سایس (۱۹۹۸) سعی دارد که داده‌کاوی را به عنوان گامی که طی آن الگوهای موجود در داده‌ها به طور خودکار (مثلاً با استفاده از الگوریتمهای محاسباتی) کشف می‌شوند، معرفی کند. در حالی که کشف دانش نه تنها مرحله داده‌کاوی، بلکه گامهای پیش‌فرآوری (از قبیل پاک‌سازی داده‌ها) و گامهای پس‌فرآوری (نظیر تفسیر نتایج) را نیز در بر دارد. روشن است که همین مرحله پس‌فرآوری است که نقش سنتی آمارشناسان بوده است. الدر و بریگیون (۱۹۹۶) دورنمایی از کشف دانش در پایگاههای داده‌ها را عرضه می‌کنند. هند و همکاران (۲۰۰۰) داده‌کاوی را به عنوان «تحلیل ثانی پایگاههای داده‌ای بزرگ با هدف یافتن روابط غیر منتظره که مورد علاقه صاحبان پایگاهها هستند و برایشان ارزشمندند» تعریف می‌کنند.

غالباً اندازه پایگاه داده‌ای چنان بزرگ است که تحلیل‌های اکتشافی داده‌ها به روش کلاسیک برای آنها نارسا هستند. چون برخی مسائل موجود در داده‌کاوی و کشف دانش در پایگاههای داده‌ای به طور طبیعی به صورت‌های داده‌های نمادین می‌انجامند، تحلیل‌های داده‌های نمادین در اینجا نیز نقشی برای ایفا دارند. بدین ترتیب، همکاری تیمهای علمی بین رشته‌ای در هنگام کار با مجموعه‌های بزرگ داده‌ها (مثل متخصصان رایانه و آمارشناسان) به امری اساسی تبدیل می‌شود.

هدف این مقاله، آن است که مفاهیم داده‌های نمادین و شیوه‌های تحلیل آنها را که در حال حاضر در آثار مکتوب موجودند، مرور کند. بنابراین، داده‌های نمادین که گاهی به اصطلاح «اتم‌های دانش» خوانده می‌شوند، در بخش ۲ تعریف شده با داده‌های کلاسیک مقایسه می‌شوند. ساختن رده‌هایی از اشیاء، نمادین که وقتی اندازه داده‌ها برای تحلیل‌های کلاسیک بسیار بزرگ باشد، یا هنگامی که به عوض داده‌های استاندارد دانش به عنوان درونداد داده شود، از مقدمات

بر عکس، یک نقطه داده‌ای نمادین ممکن است یک ابر مکعب در فضای p بعدی با حاصلضرب دکارتی توزیعها باشد.

درایه‌های موجود در مجموعه داده‌های نمادین (که با ξ_{ij} نمایانده می‌شوند) مقید به مقادیر مشخص تک نیستند. بدین ترتیب، سن را می‌شود به صورتی که در یک بازه واقع باشد، ثبت کرد. مثل، [۰، ۱۰]، [۱۰، ۲۰] و ... این وضع وقتی می‌تواند رخ دهد که نقطه داده‌ای نمایشگر سن یک خانواده، یا گروهی از افراد باشد که سن آنها به طور دستجمعی در بازه‌ای قرار دارد (مانند زن و شوهری مسن [۶۰، ۷۰]) یا داده‌ها ممکن است مربوط به یک فرد تنها باشد که سن دقیقش معلوم نیست، جز آنکه می‌دانیم در یک بازه سنی قرار دارد، یا سن او در طی زمان همراه با اجرای آزمایشی که داده‌ها را ایجاد کرده است، تغییر کرده است، یا ترکیبی از این حالتها و انواع مختلف آنها باشد، که داده‌های با دامنه بازه‌ای را به وجود می‌آورند.

در راستایی متفاوت، ممکن است اندازه‌گیری مشخصه‌ای دقیقاً به صورت مقداری تک میسر نباشد، نتوان اندازه گرفت که نبض ۶۴ است، بلکه بتوان اندازه گرفت که مقدار آن (64 ± 1) یا به طور کلی اندازه متغیری $(x \pm \delta)$ است.

شخصی ممکن است $x \leq 2$ یا $x > 2$ خواهر و برادر (یا فرزند یا ...) داشته باشد. متغیر فشار خون ممکن است به صورت مقادیر [بالا، پایین] آن، مثلاً $\xi_{ij} = [78, 120]$ ثبت شود. این متغیرها، متغیرهای نمادین بازه-مقدارند.

نوعی متفاوت از متغیر می‌تواند متغیری مربوط به سرطان باشد که ممکن است قلمرو $\{...، کبد، سینه، استخوان، ریه\}$ را داشته باشد که فهرستی از همه سرطان‌های ممکن است و شخص معینی مثلاً دارای مقادیر خاص $\{کبد، ریه\} = \xi_{ij}$ است. در مثالی دیگر، فرض کنید متغیر Y نمایشگر نوع خودرویی باشد که خانواده‌ای داراست با $\{...، تویوتا، بنز، پژو، پیکان\} = Y_j$. خانواده خاص i ممکن است مقدار $\{بنز، پیکان\} = \xi_{ij}$ را دارا باشد. چنین متغیرهایی، متغیرهای چند مقدار نامیده می‌شود.

نوع سوم متغیر نمادین متغیر مدی است. متغیرهای مدی متغیرهای چند وضعی با فراوانی، احتمال یا وزن منسوب به هر یک از مقادیر خاص در داده‌ها هستند. یعنی، متغیر مدی Y نگاشتی است مانند

$$Y(i) = \{U(i), \pi_i\}, \quad i \in \Omega$$

بیماری معین معالجات و سایر متغیرهای مربوط به آن بیماری ثبت می‌شوند. برای مثال، مجموعه‌ای نوعی از این گونه داده‌ها مثل جدول ۱ خواهد بود.

گیریم p تعداد متغیرهای مربوط به فرد i ، به ازای $i \in \Omega = \{1, \dots, n\}$ باشد، روشن است که p و n می‌توانند عددهایی بزرگ یا حتی فوق‌العاده بزرگ باشند و گیریم Y_j به ازای $j = 1, \dots, p$ نشانگر متغیر زام باشد. گیریم $Y_j = x_{ij}$ مقدار خاصی باشد که Y_j برای فرد i در آرایش کلاسیک اختیار می‌کند و به عنوان ماتریس $n \times p$ شامل کل مجموعه داده‌ها بنویسیم $X = (x_{ij})$ گیریم قلمرو Y_j عبارت باشد از Y_j به طوری که $X = (Y_1, \dots, Y_p)$ مقادیری را در فضای $X = \prod_{j=1}^p Y_j$ اختیار کند.

چون بود یا نبود مقدارهای گمشده در حال حاضر اهمیتی ندارد، گیریم همه مقادیر وجود دارند؛ اگر چه در مورد مجموعه‌های بزرگ داده‌ها محتمل‌ترین وضع آن است که داده گمشده داشته باشیم. متغیرها می‌توانند کمی باشند، مثل سن با $Y_{age} = \{x \geq 0\} = Y_+$ به عنوان متغیر تصادفی پیوسته، یا به صورت $Y_{age} = \{0, 1, 2, \dots\} = \eta_0$ به عنوان متغیری تصادفی گسسته. متغیرها ممکن است رسته‌ای باشند، مثل شهر با $Y_{city} = \{...، اصفهان، تهران\}$ یا رمزیده باشند، به صورت $Y = \{1, 2, \dots\}$ متغیرهای بیماری ممکن است به صورت رسته‌هایی (رمزیده یا نارمزیده) از یک متغیر با قلمرو $\{...، کبد، سرطان، قلب\} = Y$ یا به طور محتمل‌تر به صورت متغیری نشانگر، مثل سرطان $Y = \{...، قلمرو، بلی، خیر\} = Y$ یا $\{0, 1\}$ یا با سطوح رمزیده دیگر که نشانگر مراحل بیماری است، ثبت شوند.

به همین ترتیب، برای ثبت انواع بسیار زیاد سرطان‌های ممکن، هر نوع سرطان ممکن است با متغیری مثل Y نمایش داده شود، یا با رسته‌ای از متغیر سرطان نمایش داده شود.

ماهیت دقیق توصیف متغیرها چندان مهم نیست. آنچه که اساسی است این است که در آرایش کلاسیک برای هر x_{ij} موجود در X ، دقیقاً تنها یک مقدار محقق ممکن وجود دارد. یعنی، مثلاً مشخصات یک فرد عبارتند از $Y_{age}=24$ ، یا کاشان Y_{city} ، بلی Y_{cancer} ، $Y_{pulse}=24$ و مانند آن. بدین ترتیب یک نقطه داده‌ای کلاسیک تک نقطه‌ای است در فضای p بعدی X .

B_j زیر مجموعه‌ای از خط حقیقی R است، یعنی $B_j \subseteq R$ ، اگر Y_j متغیری بازه‌ای باشد، $B_j = \{[\alpha, \beta], -\infty < \alpha, \beta < \infty\}$ ؛ اگر Y_j رشته‌ای (نامی، ترتیبی، زیر مجموعه‌ای از یک قلمرو متناهی Y_j) باشد، $B_j = \{B \mid B \subseteq \{\text{مثلاً فهرست همه سرطان‌ها}\}\}$ ؛ اگر Y_j متغیری مدی باشد، $B_j = M(Y_j)$ که در آن $M(Y)$ خانواده‌ای از همه اندازه‌های نامنفی روی Y است. در این صورت، داده نمادین برای مجموعه اشیاء E با ماتریس $X = (\xi_{ij})$ از مرتبه $N \times p$ نمایش داده می‌شود که در آن $z \in B_j, Y_j(u) = \xi_{ij}$ مقدار نمادین مشاهده شده برای متغیر $Y_j, z = 1, \dots, p$ در مورد شیئی $u \in E$ است. سطر x'_i از X توصیف نمادین شیئی u خوانده می‌شود. بدین ترتیب برای داده‌های جدول ۲ سطر اول،

$$x'_1 = ([۷۰ و ۸۰], \{\text{مرد}\}, \{\text{سرطان مغز}\}, \{\text{بوشهر}\}, [۱۲۰ و ۱۷۹] \text{ و } [۳۰ و ۲۰])$$

نمایشگر مردی است در سنین ۲۰ سالگی که سرطان مغز دارد، با فشار خون ۱۲۰/۷۹، وزنی بین ۷۰-۸۰ کیلوگرم و در بوشهر می‌زید.

شیئی u مربوط به این x'_i می‌تواند مرد خاصی باشد که طی دوره ده ساله‌ای تحت نظارت بوده که در این مدت وزنش بین ۷۰ و ۸۰ کیلوگرم تغییر کرده، یا u می‌توانست گردآیه‌ای از افراد باشد که نشان از ۲۰ تا ۳۰ سال متغیر بوده و دارای مشخصاتی بوده‌اند که با x'_i توصیف شده است. داده‌های x'_i در جدول ۲ ممکن است معرف همان فردی باشد که به عنوان فرد $i = 1$ چهار از جدول ۱ معرفی شده، اما تنها این امر معلوم است که او (با احتمال p) دارای سرطان سینه یا (با احتمال $1-p$) دارای سرطان ریه است، ولی معلوم نیست کدام یک. از سوی دیگر، داده‌های فوق می‌توانند معرف مجموعه زنان ۴۷ ساله از بابلسر باشند که از آنها به نسبت p دارای سرطان ریه و $(1-p)$ دارای سرطان سینه‌اند، یا می‌تواند نمایشگر افرادی باشد که دارای هر دو سرطان ریه و سینه‌اند و مانند آن (در مرحله‌ای از کار، اینکه متغیر (در اینجا نوع سرطان) از نوع رشته‌ای، فهرستی، مدی یا هرچه هست، باید به صراحت تعریف شود).

موضوع دیگر، به متغیرهای وابسته مربوط می‌شود که برای داده‌های نمادین به معنای وابستگی منطقی، وابستگی سلسله مراتبی، یا وابستگی تصادفی است. وابستگی منطقی همان گونه که از معنای واژه پیداست به صورت اگر-آنگاه است که در مثال فوق اگر $[۱۰ \leq \text{سن}]$ آنگاه $[۰ = \text{تعداد فرزندان}]$. وابستگی سلسله مراتبی وقتی رخ می‌دهد که

که در آن π_i اندازه‌ای نامنفی یا توزیعی روی قلمرو Y متشکل از مقادیر مشاهده‌ای ممکن و $U(i) \subseteq Y$ تکیه π_i گاه است. مثلاً اگر سه نفر از خواهران و برادران شخصی مرض قند داشته باشند و یک نفر نداشته باشد، متغیر توصیف‌گر استعداد مرض قند می‌تواند مقدار خاص $\{ \frac{1}{4} \text{ سالم}, \frac{3}{4} \text{ مرض قند} \} = \xi_{ij}$ را اختیار کند. به طور کلی‌تر، ξ_{ij} می‌تواند یک بافتنگار، توزیع تجربی، توزیع احتمال، یک مدل یا مانند آنها باشد. در واقع شرایتر (۱۹۸۴) نظر داد که «توزیعها عددهای آینده‌اند». در حالی که در این مثال، وزن‌های $(\frac{3}{4}, \frac{1}{4})$ ممکن است معرف فراوانی‌های نسبی باشند، انواع دیگر وزن‌ها از قبیل «گنجایش‌ها»، «باورمندی‌ها»، «بایستگی‌ها»، «امکانات» و غیره می‌توانند به کار گرفته شوند. در اینجا «گنجایش» را به تعبیر شوکه (۱۹۵۴) به عنوان احتمال اینکه دست کم یک فرد در رده مذکور دارای مقدار Y معین (مثلاً مرض قند) است، تعریف می‌کنیم و «باورمندی» به تعبیر شیفتر (۱۹۷۶) تعریف می‌شود به احتمال اینکه هر فرد موجود در رده دارای آن مشخصه باشد (دیدی، ۱۹۹۵).

پس به طور کلی، بر خلاف داده‌های کلاسیک که برای آنها هر نقطه داده‌ای مشتمل بر تک مقدار (رسته‌ای یا کمی) است، داده‌های نمادین می‌توانند تغییراتی درونی را در برداشته باشند و می‌توانند ساختار یافته باشند. وجود همین تغییرات درونی است که نیاز به فنون جدید برای تحلیل را لازم می‌نماید. این فنون به طور کلی با فنون مربوط به داده‌های کلاسیک تفاوت خواهد داشت.

از نظر نمادگذاری، مجموعه‌ای پایه‌ای از اشیاء را داریم که اعضاء یا هستی‌هایی هستند، $E = \{1, \dots, N\}$ مجموعه اشیاء خوانده می‌شود. این مجموعه اشیاء می‌تواند معرف کیهانی از افراد $E = \Omega$ (مثل بالا) باشد که در آن حالت $N=n$ یا اگر داشته باشیم $N \leq n$ ، هر یک مجموعه از اشیای مذکور زیر مجموعه‌ای از Ω است. همچنین، چنانکه در داده کاوی یا تحلیل‌های نمادین کمرأ رخ می‌دهد، اشیاء u در E رده‌های C_1, \dots, C_m از افراد موجود در Ω هستند، با $E = \{C_1, \dots, C_m\}$ و $N=m$ بدین ترتیب، مثلاً رده C_1 ممکن مشتمل بر همه افراد موجود در Ω باشد که سرطان داشته بودند. هر شیئی $u \in E$ با p متغیر نمادین $Y_j, j=1, \dots, p$ با قلمرو Y_j تعریف می‌شود که Y_j نگاهی از مجموعه اشیاء E به برد Y_j است که به نوع متغیر Y_j بستگی دارد. بدین ترتیب، اگر Y_j متغیری کمی کلاسیک باشد، قلمرو

تحلیل آماری به دست دهد. توجه کنید که این ساخت ممکن است از رده‌های حاصل از یک شیوه خوشه‌بندی فرق داشته باشد، اما لزوماً چنین نیست. توجه داشته باشید که این m رده انبوهیده ممکن است معرف m الگوی به دست آمده از یک شیوه داده کاوی باشد.

این توصیف ما را به مفهوم شیئی نمادین رهنمون می‌شود که در سلسله‌ای از مقالات دیدی و همکارانش بسط یافته است (مثلاً دیدی، ۱۹۸۷، ۱۹۸۹، ۱۹۹۰، باک و دیدی، ۲۰۰۰ و استفن و همکاران، ۲۰۰۰). این مطلب را ابتدا با چند مثال معرفی می‌کنیم و در پایان بخش تعریفی دقیق ارائه می‌شود.

چند مثال

فرض کنید به مفهوم «شمال خاوری» علاقه‌مندیم. بدین ترتیب، توصیفی مانند d داریم که معرف مقادیر خاص $\{ \dots \}$ و سایر شهرهای شمال خاوری، ... نیشابور، مشهد} در قلمرو Y_{city} است و رابطه‌ای مانند R (در اینجا \in) داریم که متغیرهای Y_{city} را با توصیف خاص مورد نظر ربط می‌دهد. این مطلب را به صورت مثلاً

$$a = \{ \text{سایر شهرهای ش-خ، ...، نیشابور، مشهد} \} \in Y_{city}$$

می‌نویسیم. در این صورت، هر فرد i در $\Omega = \{1, \dots, n\}$ یا یک فرد شمال خاوری است یا نیست. یعنی، a نگاشتی است از Ω بر $\{0, 1\}$ که برای فردی مانند i که در شمال خاوری زندگی می‌کند $a(i) = 1$ و اگر نه $a(i) = 0$ است. پس، اگر فرد i در مشهد زندگی کند (یعنی، $\text{مشهد} = Y_{city}(i)$)، آنگاه

$$a(i) = 1 = \{ \text{سایر شهرهای ش-خ، ...، نیشابور، مشهد} \} \in \text{مشهد}$$

مجموعه $i \in \Omega$ که برای آنها $a(i) = 1$ ، گستره a در Ω نامیده می‌شود. سه گانه $S = (a, R, d)$ یک شیئی نمادین است که در آن R رابطه‌ای بین توصیف $Y(i)$ از متغیر (بی صدای) Y و توصیف d و a نگاشتی از Ω به L است که L به R و d بستگی دارد (در مثال شمال خاوری، $L = \{0, 1\}$). توصیف d می‌تواند توصیف قصدی باشد. مثلاً همان طور که از نامش پیداست، قصد داریم که مجموعه افراد موجود در Ω را که در «شمال خاوری» زندگی می‌کنند، بیابیم. بدین ترتیب، مفهوم «شمال خاوری» تا حدی شبیه به مفهوم جامعه کلاسیک است و گستره موجود در Ω متناظر با نمونه افراد از شمال خاوری در مجموعه داده‌های واقعی است، اما به خاطر بیاورید که همان طور که در بخش ۲

برآمد یک متغیر (مثلاً معالجه سرطان Y_p) با $\{ \dots \}$ پرتو درمانی، شیمی درمانی $\{ Y_p \}$ به برآمد واقعی تحقق یافته متغیری دیگر (مثلاً سرطان دارد $Y_1 = \{ \text{بلبی، خیر} \}$) وابسته است. اگر Y_1 دارای مقدار $\{ \text{بلبی} \}$ باشد، آنگاه $\{ \text{مثلاً شیمی درمانی} \} = Y_p$. در حالی که اگر $\{ \text{خیر} \} = Y_1$ آنگاه روشن است که Y_p کاربرد ندارد [به منظور توضیح مطلب فرض می‌کنیم که معالجه شیمی درمانی به دلیلی دیگر صورت نگیرد]. در این موارد، متغیر کاربرد ندارد Z را با قلمرو $\{ \text{ک ن} \} = Z$ تعریف می‌کنیم. این گونه متغیرها را متغیرهای مادر (Y_1) -دختر (Y_p) نیز گویند.

۳. رده‌ها و طرز تشکیل آنها؛ اشیاء نمادین

مجموعه داده‌های نمادین ممکن است از آغاز کار پیشاپیش از حیث اندازه به قدر کافی کوچک باشند که به منظور انجام هر نوع تحلیل آماری نمادین به طور مستقیم، نیاز به پیش‌فرآوری نباشد. مثالی از این نوع داده‌ها جدول ۳ است که برای تحلیل مولفه‌های اصلی نمادین به کار رفته‌اند. اما به طور کلی تر و تقریباً به طرز ناپذیر قبل از آنکه هر نوع تحلیل نمادین از داده‌ها را بتوان انجام داد، به ویژه برای مجموعه‌های بزرگ داده‌ها، نیاز خواهد بود که درجات مختلفی از داده‌آمایی برای سازمان‌دهی اطلاعات در رده‌هایی که مناسب برای سؤال مورد نظر باشند، صورت گیرند.

در برخی موارد، اشیاء موجود در E (یا Ω) پیشاپیش در رده‌هایی جمع شده‌اند، حتی در این موارد هم سؤال‌هایی خاص ممکن است سازمان‌دهی مجدد در رده‌بندی متفاوتی را ضروری سازند، صرف نظر از اینکه مجموعه داده‌ها کوچک است یا بزرگ. مثلاً یک مجموعه از رده‌های C_1, \dots, C_m ممکن است افراد را که بر حسب m نوع مختلف از بیماری‌های مهم رسته‌بندی شده‌اند، نمایش دهد. در حالی که تحلیل دیگر ممکن است ضروری نماید که ساختار رده‌ها بر حسب شهر، جنس یا بر حسب جنس و سن و غیره باشد. مرحله ابتدایی دیگر وقتی است که داده‌ها در ابتدا برای پایگاه‌های داده‌های آمار کلاسیک یا علوم رایانه برای هر فرد $i \in \Omega = \{1, \dots, n\}$ با n فوق‌العاده بزرگ ثبت شده‌اند. همین طور است وقتی که پایگاه‌های داده‌ای نمادین بسیار بزرگ باشد. در این صورت، این مرحله از تحلیل داده‌های نمادین متناظر است با انبوهیدن n شیئی در m رده که m بسیار کوچکتر از n است و طوری طراحی شده که قالب‌های قابل اداره‌تری را پیش از هر

مجموعه جدیدی از داده‌ها را به یک «مشاهده» برای هر رده به دست بدهد. در این حالت اخیر، صرفنظر از اینکه مقادیر اولیه داده‌هایی کلاسیک یا نمادین باشند، مقادیر مجموعه داده‌ها به احتمال زیاد داده‌های نمادین خواهند بود. مثلاً حتی اگر هر فرد موجود در Ω به صورت دارای سرطان و یا بدون سرطان (بلی، خیر) $(Y_{cancer} = \text{بله}, \text{خیر})$ ، یعنی به صورت یک مقدار داده‌ای کلاسیک ثبت شده باشد، وقتی این متغیر با رده مربوط به شهر ارتباط داده شود، به صورت $\{0/9\}$ خیر، $\{0/1\}$ بلی، یعنی 0.1 دارای سرطان و 0.9 بدون سرطان در خواهند آمد. پس اکنون متغیر Y_{cancer} متغیری مدی است.

به همین نحو، رده‌ای که از گستره یک شیئی نمادین ساخته می‌شود، نوعاً به وسیله یک مجموعه داده‌های نمادین توصیف می‌شود. مثلاً، فرض کنید که علاقه ما متوجه است به «کسانی که در مشهد زندگی می‌کنند» یعنی، $[Y_{city} = \text{مشهد}] = a$ و فرض کنید که متغیر Y_{child} تعداد فرزندان هر فرد $i \in \Omega$ است که مقادیر ممکن $\{0, 1, 2, 3, \dots\}$ را دارد. فرض کنید مقدار داده‌ای برای هر i مقداری کلاسیک است (تبدیل آنها به یک داده نمادین از قبیل اینکه فرد i دارای ۱ یا ۲ فرزند است، یعنی $\{1, 2\} = \{i\}$ به آسانی نتیجه می‌شود). در این صورت، شیئی که معرف همه کسانی است که در مشهد زندگی می‌کنند، اکنون دارای متغیر نمادین Y_{child} خواهد بود با مقدار خاص،

$$Y_{child} = \{(0, f_0), (1, f_1), (2, f_2), (\geq 3, f_3)\}$$

که در آن $f_i, i = 0, 1, 2, \geq 3$ ، فراوانی نسبی افراد این رده است که دارای i فرزندند.

حالتی خاص از یک شیئی نمادین عبارت است از یک حکم. حکم‌ها، مخصوصاً هنگام انبوهیدن افراد درون رده‌ها از یک پایگاه داده‌ای آغازین (رابطه‌ای)، دارای اهمیت‌اند.

گیریم توصیف لازم مورد نظر از یک فرد یا از یک مفهوم w را با $Z = (Z_1, \dots, Z_p)$ بنمایانیم. در اینجا، Z_i می‌تواند هستی تک مقداره کلاسیک Z_i یا هستی نمادین Z_i باشد. یعنی در حالی که Z_i معرف مقدار داده‌ای کلاسیک تحقیق یافته‌ای است و Z_i معرف مقدار داده‌ای نمادین تحقیق یافته‌ای است، Z_i مقداری است که مشخصاً به دنبال آن هستیم مقدار تصریح شده است. بدین ترتیب، مثلاً فرض کنید که به شیئی نمادینی علاقه‌مندیم که معرف کسانی است که در شمال خاوری

ذکر شده Ω ممکن است پیشاپیش خود «جامعه» باشد یا ممکن است به معنای نمونه‌گیری آمار کلاسیک یک «نمونه» باشد.

اشیاء نمادین به یکی از سه طریق در محدوده عمل تحلیل‌های نمادین داده‌ها نقش ایفا می‌کند. نخست، یک شیئی نمادین ممکن است به کمک قصد آن (مثلاً توصیفش و راهی برای محاسبه گستره‌اش) نمایشگر مفهومی باشد و می‌توان به عنوان درون‌داد تحلیل داده‌های نمادین به کار گرفته شود. بدین ترتیب، مفهوم «شمال خاوری» را می‌توان با یک شیئی نمادین نمایش داد که قصد آن از طریقی توصیفی سرشت‌نما و راهی برای پیدا کردن گستره‌اش که مجموعه افرادی است که در شمال خاوری زندگی می‌کنند، تعریف می‌شود.

مجموعه‌ای از این گونه ناحیه‌ها و اشیای نمادین همراه با آنها می‌توانند درون‌داد یک تحلیل داده‌های نمادین را تشکیل دهند. دوم، می‌توان آن را به عنوان برون‌دادی از یک تحلیل داده‌های نمادین مورد استفاده قرار داد، مثل وقتی که تحلیل خوشه‌ای حاکی از آن است که شمال خاوری متعلق به خوشه‌ای خاص است که خود خوشه را به عنوان یک مفهوم می‌توان تلقی کرد و به وسیله یک شیئی نمادین نمایش داد. وضعیت سوم زمانی که فردی جدید (i') را داریم که دارای توصیف d' است و می‌خواهیم بدانیم که آیا این فرد (i') با شیئی نمادینی که توصیفش d است، تطبیق می‌کند یا خیر. یعنی d و d' را به کمک R مقایسه می‌کنیم تا $[d'Rd] \in L = \{0, 1\}$ را بدهد که در آن $[d'Rd] = 1$ به معنای آن است که ارتباطی بین d و d' وجود دارد. این فرد «جدید» ممکن است فردی «قدیم» اما با داده‌های روزآمد شده باشد یا ممکن است فردی جدید باشد که به پایگاه داده‌ها اضافه شده و امکان دارد که او با یکی از رده‌های قبلاً موجود اشیاء نمادین «همسان» باشد یا نباشد (مثلاً آیا باید برای این شخص پوشش بیمه‌ای خاص را فراهم ساخت؟).

در خصوص انبوهیدن داده‌هایمان در تعداد کمی از رده‌ها، اگر قرار باشد که افراد موجود در Ω را بر حسب شهر، یعنی بر حسب مقدار متغیر u_{city} بیان‌بوییم، رده‌های مربوطه C_{ii} با $u \in \{1, \dots, m\}$ متشکل از آن افرادی در Ω است که در گستره نگاشت متناظر a_{ii} قرار دارند. تحلیل آماری متعاقب می‌تواند یکی از دو راستای وسیع را در پیش گیرد. یا آنکه برای هر رده بر حسب اقتضا، به طور جداگانه داده‌های کلاسیک یا نمادین را برای افراد موجود در C_{ii} به عنوان نمونه‌ای مرکب از Ω_{ii} مشاهده تحلیل می‌کنیم، یا داده‌های هر رده را خلاصه می‌کنیم تا

رده افرادی را که دارای سرطان کبد، یا ریه یا هر دو هستند، توصیف می کرد.

همچنین، یک حکم می تواند صورت زیر را اختیار کند.

$$a = [Y_{age} \subseteq [20, 30]] \wedge [Y_{city} \in \{\text{مشهد}\}] \quad (7)$$

یعنی، این حکم کسانی را که در مشهد زندگی می کنند و از لحاظ سنی در سنین بیست سالگی هستند، توصیف می کند.

رابطه های R می تواند هر یک از صورت های $=, \neq, \in, \subseteq, \supseteq$ و غیره را اختیار کند یا می تواند رابطه ای تطبیقی (مانند مقایسه ای از توزیعهای احتمال) یا دنباله ای ساختار یافته و غیره باشند. این روابط پیوند بین متغیر نمادین Y و توصیف خاص مورد نظر Z را تشکیل می دهند. قلمرو شیئی نمادین را می توان به صورت زیر نوشت.

$$D = D_1 \times \dots \times D_p \subseteq X = \prod_{j=1}^p Y_j$$

که در آن $D_j \subseteq Y_j$ هر p گانه (D_1, \dots, D_p) از مجموعه ها یک نظام توصیف خوانده می شود و هر زیر مجموعه D یک مجموعه توصیف مشتمل بر بردارهای توصیف $Z = (Z_1, \dots, Z_p)$ است. در حالی که ترکیب عناصر $Z_j \in Y_j$ و مجموعه های $D_j \subseteq Y_j$ مثلاً مثل آنچه در (7) دیده می شود، صرفاً یک توصیف است. وقتی قیودی در مورد هر یک از متغیرها وجود داشته باشد (مانند مواردی که وابستگی های منطقی وجود دارند)، آنگاه فضای D در خود «سوراخی» دارد که متناظر با مقادیری است که با قیود تطبیق می کنند. کل همه توصیف های D عبارت است از فضای توصیف D . بنابراین، حکم را می توان چنین نوشت.

$$a = [Y \in D] \equiv \bigwedge_{k=1}^v [Y_{jk} R_{jk} Z_{jk}] = [YRZ]$$

که در آن $R = \bigwedge_{k=1}^v R_{jk}$ رابطه حاصلضربی خوانده می شود.

دقت کنید که اگر حکمی تلویحاً متغیری خاص مانند Y_w را دربر نداشته باشد، قلمرو مربوط به آن متغیر در Y_w بی تغییر باقی می ماند. مثلاً، اگر داشته باشیم $p = 3$ و شهر Y_1 و سن Y_2 و وزن Y_3 آنگاه حکم $a = [Y_1 \in \{\text{نیشابور، مشهد}\}]$ دارای قلمرو $Y_2 \times Y_3 \times Y_3$ است و در جستجوی همه آن کسانی است که (صرفنظر از سن و وزن) فقط در مشهد و نیشابورند.

زندگی می کنند. در این صورت، Z_{city} مجموعه شهرهای شمال خاوری است. این موضوع را با حکم،

$$a = [Y_{city} \in \{\text{نیشابور، مشهد، ...}\}] \quad (1)$$

که در آن a نگاشتی از Ω بر $\{0, 1\}$ است، به قسمی است که برای فرد یا شیئی w داریم $a(w) = 1$ اگر سایر شهرهای $Y_{city}(w) \in \{\text{مشهد، ...}\}$ درست باشد.

به طور کلی، هر حکمی صورت زیر را اختیار می کند

$$a = [Y_{j_1} R_{j_1} Z_{j_1}] \wedge [Y_{j_2} R_{j_2} Z_{j_2}] \wedge \dots \wedge [Y_{j_v} R_{j_v} Z_{j_v}] \quad (2)$$

به ازای $1 \leq j_1, \dots, j_v \leq p$ که در آن « \wedge » به معنی «و» ضرب منطقی است و R رابطه تصریح شده بین متغیرهای نمادین Y_j و مقدار توصیفی Z_j است. برای هر فرد $i \in \Omega$ وقتی که حکم برای آن فرد درست (یا غلط) باشد، داریم $a(i) = 1$ (یا 0). به طور دقیق تر بگوییم، هر حکم عبارت از ترکیب عطفی v پیشامد $[Y_K R_K Z_K]$ ، $k = 1, \dots, v$ است. مثلاً حکم های

$$a = [Y_{cancer} = \text{بلی}] \quad \text{و} \quad a = [Y_{age} \geq 60] \quad (3)$$

به ترتیب، نمایندگان افراد دارای سرطان و افراد 60 سال به بالا هستند. حکم

$$a = [Y_{cancer} = \text{بلی}] \wedge [Y_{age} \geq 60] \quad (4)$$

نمایندگان همه بیماران سرطانی است که 60 سال به بالا هستند، در حالی که حکم

$$a = [Y_{age} < 20] \wedge [Y_{abc} > 70] \quad (5)$$

در جستجوی همه افرادی است که زیر 20 سال و بالای 70 سال اند. در هر مورد، با مفهومی مانند «افراد 60 سال به بالا»، «افراد 60 سال به بالای دارای سرطان» و غیره سروکار داریم و افراد خاص حاضر در Ω را بر فضای $\{0, 1\}$ می نگارد.

اگر به عوض ثبت متغیر سرطان به صورت متغیر رسته ای $\{بلی، خیر\}$ به صورت متغیری مانند $\{\dots، سینه، کبد، ریه\}$ ثبت می کردیم، حکم

$$a = [Y_{cancer} \in \{\text{کبد، ریه}\}] \quad (6)$$

$$Ext\alpha(a) = Q\alpha = \{i \in a(i) \geq \alpha\}.$$

تعریف‌های رسمی

اکنون مفهوم‌های زیر را به طور رسمی تعریف می‌کنیم.

دیدنی و امیلیون (۱۹۹۶) و دیدنی و همکاران (۱۹۹۶) اشیاء نمادین مدی را بررسی می‌کنند و برخی از زیر ساخت‌های نظری آن را نیز عرضه می‌دارند. اکنون تعریف رسمی یک شیئی نمادین را به صورت زیر داریم.

تعریف: یک شیئی نمادین عبارت است از سه گانه $S = (a, R, d)$ که در آن a نگاشتی است $L : \Omega \rightarrow L$ که افراد $i \in \Omega$ را بر فضای L می‌نگارد و این فضا به رابطه R بین توصیف‌ها و توصیف d بستگی دارد. وقتی $L = \{0, 1\}$ ، یعنی، هنگامی که a نگاشتی دودویی است، S یک شیئی نمادین بولی است. هنگامی که $L = [0, 1]$ ، آنگاه S یک شیئی نمادین مدی است. یعنی، یک شیئی نمادین مدلی ریاضی از یک مفهوم است (ن. ک. به دیدنی، ۱۹۹۵). اگر داشته باشیم $[a, R, d] \in \{0, 1\}$ ، آنگاه S یک شیئی نمادین بولی و اگر داشته باشیم $[a, R, d] \in [0, 1]$ ، آنگاه S یک شیئی نمادین مدی است. مثلاً در (۷) رابطه $R = (\subseteq, \in)$ و توصیف (مشهد) و $d = \{[20, 30]\}$ را داریم. قصد آن است که «۲۰-۳۰ ساله‌هایی را که در مشهد زندگی می‌کنند» بیایم. گستره متشکل از همه افراد موجود در Ω که با این توصیف تطبیق دارند، یعنی آن‌هایی که برای آنها $a(i) = 1$.

در حالی که اذعان به ضرورت ابداع روشهایی برای تحلیل نمادین و ابزارهایی برای توصیف اشیاء نمادین ابداع نسبتاً جدید است، اندیشه بررسی واحدهای سطح بالاتر مثل مفاهیم در واقع امری باستانی است. ارغنون ارسطو در سده چهارم قبل از میلاد مسیح (ارسطو، IVBC، ۱۹۹۴) به روشنی افراد مرتبه نخست (از قبیل یک اسب یا یک انسان) را که معرف واحدهایی در جهان (جامعه آماری) اند از افراد مرتبه دوم (از قبیل اسب یا انسان) که به عنوان واحدهایی از رده‌ای از افراد معرفی می‌شوند، متمایز می‌سازد.

بعدها آرنو و نیکول (۱۶۶۲) مفهوم را به کمک پنداشته‌های قصد و گستره (که به معنای آنها در ۱۶۶۲ با معانی آنها در این مقاله تطبیق می‌کنند) به شرح زیر تعریف کردند.

اکنون در اندیشه‌های کیهانی دو چیز وجود دارند که تمایز کامل آنها اهمیت دارد؛ فراگیری و گسترش (به جای «قصد» و «گستره»). فراگیری یک اندیشه را صفاتی می‌نامیم که آن اندیشه در بردارد و بدون از بین بردن اندیشه نمی‌توان آنها را از آن گرفت. بدین ترتیب

تعریف: وقتی حکمی درباره شیئی خاص $i \in \Omega$ بیان می‌شود، اگر آن حکم درباره آن شیئی صادق باشد، مقدار راست ($a = 1$) را و اگر دروغ باشد ($a = 0$) را اختیار می‌کند. می‌نویسیم

$$a(i) = [Y(i) \in D] = \begin{cases} 1, & Y(i) \in D \\ 0, & o.w. \end{cases} \quad (8)$$

تابع $a(i)$ تابع راستی خوانده می‌شود و نمایشگر نگاشتی از Ω بر $\{0, 1\}$ است.

مجموعه همه مقادیر $i \in \Omega$ که برای آنها حکم صادق است، گسترش در Ω خوانده می‌شود که با $Ext(a)$ یا Q نمایانده می‌شود،

$$Ext(a) = Q = \{i \in \Omega \mid Y(i) \in D\} = \{i \in \Omega \mid a(i) = 1\} \quad (9)$$

نوعاً، هر رده با شیئی نمادینی که توصیفش می‌کند، شناسایی می‌شود. مثلاً حکم (۷) متناظر است با شیئی نمادین «۲۰-۳۰ ساله‌هایی که در مشهد زندگی می‌کنند». گسترش این حکم، $Ext(a)$ ، رده متشکل از همه اشیاء $i \in \Omega$ را که با این توصیف تطبیق دارند، تولید می‌کند.

یک رده ممکن است، مثل بالا، با جستجوی گسترش یک حکم a ساخته شود. به شق دیگر، ممکن است مطلوب آن باشد که همه آن افراد (دیگر) در Ω را که با توصیف $Y_j(u)$ مربوط به فردی معین $u \in \Omega$ تطبیق دارند، پیدا کنیم. بدین ترتیب در این حالت، حکم عبارت است

$$a_i = \bigwedge_{j=1}^p [Y_j = Y_j(u)] \quad \text{از}$$

که در آن اکنون $Z_j = Y_j(u)$ و دارای گسترش زیر است،

$$Ext(a_i) = \{i \in \Omega \mid Y_j(i) = Y_j(u), j = 1, \dots, p\}$$

به طور کلی‌تر، حکمی ممکن است با مقداری مدی (از قبیل یک احتمال) صادق باشد، یعنی $0 \leq a(i) \leq 1$ که نمایشگر درجه‌ای از تطبیق یک شیئی i با یک حکم a است.

در این حالت‌ها، نگاشت a بر بازه $[0, 1]$ است، یعنی، $[0, 1] \rightarrow \Omega \rightarrow a$ در این صورت، گسترش a دارای سطح α $0 \leq \alpha \leq 1$ است. که حال

در این صورت، بافتنگار استاندارد مبتنی بر پرندگان طوری است که فراوانی «بلی» دو برابر فراوانی «خیر» است. شکل ۱ الف را ببینید. برعکس، اگر افراد آماری را گونه‌ها در نظر بگیریم، چون دو گونه بی‌پرواز و یک گونه با پرواز وجود دارند، فراوانی «خیر» اکنون دو برابر فراوانی «بلی» است. شکل ۱ ب را ببینید.

دقت کنید که بافتنگار گونه‌ها صرفاً یک بافتنگار است. چیزی که این مثال نشان می‌دهد این است که نسبت به سطح افراد آماری (در اینجا پرندگان) و نسبت به مفهوم (در اینجا گونه‌ها) بافتنگارهای کاملاً متفاوتی را به دست می‌آوریم. در اینجا گونه‌ها واقعا در درون خود سطح دیگری از اطلاعات متناظر با نوع پرنده و فراوانی هر یک را دارا هستند. یعنی، «گونه‌ها» خود در بردارنده ساختاری در سطحی دیگر از متغیر گونه‌اند. به نحوی که در تولید بافتنگار شکل ۱ ب به کار رفته است. بافتنگاری نمادین (که در زیر تعریف خواهد شد) این ساختار را به حساب می‌آورد.

آماره‌های تک متغیری را (برای داده‌های $p \geq 1$ متغیر) در نظر بگیریم. نظایر آنها آماره‌های دو متغیری هستند که به خاطر نبود فضا و دقت فعلاً آنها را مورد بحث قرار نمی‌دهیم. در مورد متغیرهای با مقدار عدد صحیح و بازه-مقدار از رویکردی که برتران و گوپیل (۲۰۰۰) پیش گرفته‌اند، پیروی می‌کنیم. دو کارواو (۱۹۹۴-۱۹۹۵) و چوآکریا و همکاران (۱۹۹۸) روشهای متفاوت، ولی معادل را برای یافتن بافتنگار و احتمال‌های بازه‌ای برای متغیرهای بازه-مقدار به کار گرفته‌اند (رابطه (۲۶) را در زیر ببینید).

پیش از توصیف این کمیت‌ها، ابتدا نیاز داریم که مفهوم گسترش‌های واقعی را معرفی کنیم. از بخش ۳ به یاد بیاورید که توصیف نمادین یک شیئی $u \in E$ (یا به طور معادل $i \in \Omega$) با بردار توصیف

$$d_u = (\xi_{u1}, \dots, \xi_{up}), \quad u = 1, \dots, m$$

یا به طور کلی‌تر با $d \in (D_1, \dots, D_p)$ در فضای $D = \times_{j=1}^p D_j$ داده شد که در آن در هر حالت خاص تحقق Y_j می‌تواند یک x_j برای داده‌های کلاسیک یا یک ξ_j برای داده‌های نمادین باشد. توصیف‌های انفرادی، که با x نمایانده شده، همان توصیف‌هایی‌اند که برای آنها هر D_j یک مجموعه دارای تنها یک مقدار است. یعنی،

فراگ‌بری اندیشه یک مثلث، در حد ظاهر، شامل شکل، سه خط، سه زاویه، برابری مجموع این سه زاویه با دو زاویه قائمه و غیره است. گسترش این اندیشه را اشیایی می‌نامیم که این اندیشه درباره آنها مصداق دارد. این اشیاء زیرین‌های یک گزاره کیهانی نیز نامیده می‌شوند و آن گزاره زیرین آنهاست. بدین ترتیب، اندیشه یک مثلث به طور کلی به همه انواع مختلف مثلث گسترش می‌یابد.

سرانجام، کاربست محاسباتی ایجاد رده‌های مقتضی را می‌توان به وسیله احکام مورد استفاده در موتورهای کاوش اجرا کرد. انواع متعددی از این گونه الگوریتمها وجود دارند، مانند نرم‌افزار (SQL) standard query language یا C++ یا JAVA. نمونه‌هایی از به کارگیری این الگوریتمها در تعیین رده‌ها در مراجع ذکر شده‌اند.

۴. تحلیل داده‌های نمادین

با در دست داشتن یک مجموعه از داده‌های نمادین، گام بعد انجام تحلیل‌های آماری مقتضی است. مثل داده‌های کلاسیک در اینجا نیز امکانات بی‌پایان هستند. برخلاف داده‌های کلاسیک که یک قرن سعی و کوشش، کتابخانه‌ای از روش‌شناسی‌های تحلیلی یا آماری را فراهم کرده است، تحلیل‌های آماری برای داده‌های نمادین جدید بوده و از نظر تعداد بسیار اندک‌اند. در زیر برخی از آنها را مرور می‌کنیم.

۵. آماره‌های تک متغیره توصیفی

۱-۵- مقدمات

آماره‌های توصیفی پایه‌ای عبارتند از بافتنگارهای فراوانی و میانگین و واریانس نمونه. مشابه‌های این آماره‌ها را برای داده‌های نمادین مربوط به متغیرهای چند مقداره و بازه مقدار قاعده‌مند و متغیرهای مدی در نظر می‌گیریم. در حین اینکه این آماره‌ها را بسط می‌دهیم، مثال زیر را در نظر داشته باشیم که نشان می‌دهد باید بین سطوح در داده‌های نمادین و داده‌های کلاسیک تمایز قایل شد. این نکته به ویژه هنگام تشکیل بافتنگار این دو نوع داده آشکارتر می‌شود.

فرض کنید جزیره‌ای دوره افتاده دارای هزار پنگوئن و هزار شتر مرغ است که هر دو گونه‌هایی بی‌پرواز از پرندگان‌اند و چهار هزار کبوتر دارد که گونه‌ای پرنده با پرواز است. فرض کنید به متغیر «پرواز» علاقه‌مندیم که در اینجا دو مقدار ممکن «بلی» یا «خیر» را می‌گیرد.

نوعی از متغیرهای تصادفی گسسته باشند. فراوانی مشاهده شده ξ را به صورت

$$O_Z(\xi) = \sum_{u \in E} \pi_Z(\xi; u) \quad (12)$$

تعریف می‌کنیم که در آن مجموع یابی روی $\{1, \dots, m\}$ $u \in E$ است و

$$\pi_Z(x; u) = \frac{|\{x \in \text{vir}(d_u) \mid x_Z = \xi\}|}{|\text{vir}(d_u)|} \quad (13)$$

درصد بردارهای توصیف انفرادی در $\text{vir}(d_u)$ است به طوری که داشته باشیم $\xi = x_Z$ و $|A|$ تعداد توصیف‌های انفرادی موجود در فضای A است. در مجموع یابی (۱۲) هر u که برای آن $\text{vir}(d_u)$ تهی باشد، نادیده گرفته می‌شود. ملاحظه می‌کنیم که این فراوانی مشاهده شده یک عدد مثبت حقیقی است و مثل مورد داده‌های کلاسیک لزوماً عددی صحیح نیست.

در حالت کلاسیک، $|\text{vir}(d_u)| = 1$ و بنابراین حالتی خاص از (۱۳) است. به آسانی می‌توان نشان داد که

$$\sum_{\xi \in Z} O_Z(\xi) = m' \quad (14)$$

که $m' = (m - m_0)$ و m_0 تعداد u هایی است که برای آنها $|\text{vir}(d_u)| = 0$.

برای یک متغیر نمادین چند مقداره Z ، که مقادیر $\xi \in Z$ را اختیار می‌کند، توزیع فراوانی تجربی عبارت است از مجموعه زوج‌های $[\xi, O_Z(\xi)]$ به ازای $\xi \in Z$ و توزیع فراوانی نسبی یا بافتنگار فراوانی عبارت است از مجموعه

$$[\xi, (m')^{-1} O_Z(\xi)] \quad (15)$$

تعریف‌های زیر بلافاصله نتیجه می‌شود.

تابع توزیع تجربی Z به صورت زیر تعریف می‌شود.

$$F_Z(\xi) = (m')^{-1} \sum_{\xi_k \leq \xi} O_Z(\xi_k) \quad (16)$$

هنگامی که مقادیر ممکن ξ برای متغیر نمادین چند مقداره $Y_j = Z$ کمی باشند، می‌توانیم به شرح زیر میانگین، واریانس و میانه را تعریف کنیم.

میانگین نمونه نمادین برابر است با

$$\bar{Z} = (m')^{-1} \sum_{\xi_k} O_Z(\xi_k) \quad (17)$$

واریانس نمونه نمادین عبارت است از

$$x \equiv d = (\{x_1\}, \dots, \{x_p\}), x \in X = \prod_{j=1}^p Y_j$$

محاسبه بافتنگار فراوانی نمادین شامل شمارش تعداد توصیف‌های انفرادی است که با وابستگی‌های منطقی تلویحی موجود در داده‌ها تطبیق می‌کند. وابستگی منطقی را می‌توان با قاعده‌ای مانند v نمایش داد،

$$v = [x \in A] \Rightarrow [x \in B] \quad (10)$$

با ازای $A \subseteq D, B \subseteq D, x \in X$ و اینکه v نگاشتی از X بر $\{0, 1\}$ با $v(x) = 0(1)$ اگر x در قاعده صدق نکند (صدق کند). مثلاً فرض کنیم $x = (x_1, x_2) = (1, 0) = Y_1^{(d)}$ و تعداد فرزندان $Y_1^{(d)}$ برای $u^i \in \Omega$ باشد، فرض کنیم $A = \{x \leq 12\}$ و $B = \{0\}$. در این صورت، قاعده‌ای که تصریح کند فردی که سنش کمتر از ۱۲ سال است، دال بر این است که فرزندش نداشته است منطقی است، در حالی که فردی که سنش کمتر از ۱۲ سال است مستلزم آن باشد که دارای دو فرزند است، منطقی است نیست. نتیجه می‌شود که بردار توصیف انفرادی x در قاعده v صدق می‌کند اگر و تنها اگر $x \in A \cap B$ یا $x \notin A$. این فرمول‌بندی قاعده وابستگی منطقی برای مقاصد محاسبه آماره‌های توصیفی پایه‌ای کافی است. ورده و دو کارووالو (۱۹۹۸) انواع گوناگون قواعد مرتبط (از قبیل هم ارزی منطقی، استلزام منطقی، وابستگی‌های چند گانه، وابستگی‌های سلسله مراتبی و مانند آن را) مورد بحث قرار می‌دهند. پس تعریف رسمی زیر را داریم

تعریف: توصیف واقعی بردار توصیف d عبارت است از مجموعه همه بردارهای توصیف انفرادی x که در همه قواعد (وابستگی منطقی) v در X صدق می‌کنند. این تعریف را به صورت زیر می‌نویسیم،

$$\text{vir}(d) = \{x \in D, v(x) = 1, \forall v \in V_x\} \quad (11)$$

که در آن V_x مجموعه همه قواعد v است که روی x عمل می‌کنند.

۵-۲- متغیرهای چند مقداره - آماره‌های تک متغیره

فرض کنید می‌خواهیم توزیع فراوانی را برای متغیر نمادین چند مقداره خاص $Y_j \equiv Z$ که مقادیر خاص $\xi \in Z$ را اختیار می‌کند، بیابیم. این مقادیر می‌توانند مقادیر رسته‌ای (مثل انواع سرطان) یا هر

بدین ترتیب، برای توصیف اول d_1 داریم،

$$\{vir(d_1) = \{x \in \{1,0\} \times \{2\} : v(x) = 1\}\}$$

توصیف‌های انفرادی $\{x \in \{1,0\} \times \{2\}\}$ عبارتند از $(0,2)$ و $(1,2)$ که از آنها تنها یکی، $x = (1,2)$ در فضای C نیز هست. بنابراین، $vir(d_1) = \{(1,2)\}$.

در واقع، این عملیات عبارت است از به اصطلاح پالایش داده‌ها از نظر ریاضی از راه شناسایی فقط آن داده‌هایی که (با صدق کردن در قاعده وابستگی منطقی رابطه (20)) معنای منطقی دارند. بدین ترتیب، مقادیر داده‌ای $(Y_1, Y_2) = (0,2)$ که این اطلاع را ثبت می‌کنند که در عین حال که سرطان وجود نداشته دو معالجه مرتبط با سرطان صورت گرفته (تحت شرایط موجودی که با قاعده v مشخص می‌شوند) به عنوان مقادیر داده‌ای غلط شناسایی می‌شوند و بنابراین برای محاسبه میانگین در این تحلیل به کار نمی‌روند. در حالی که سعی در انجام چنین کاری را در اینجا نداریم، این شناسایی ممانع از آن نمی‌شود که شیوه‌هایی دیگر را که متعاقباً به کار خواهند رفت، تا آنچه را که این داده‌ها می‌توانستند باشند جانشین آنها کنند، در تحلیل بگنجانیم. در مورد مجموعه‌های کوچک داده‌ها، ممکن است میسر باشد که داده‌ها را به طور بصری (یا به نوعی دیگر) «تصحیح» کنیم. برای مجموعه‌های بسیار بزرگ داده‌ها، این امر همیشه ممکن نیست، از این رو یک قاعده وابستگی منطقی برای انجام عمل تصحیح به طور ریاضی/محاسباتی امری اساسی است.

به طور مشابه برای توصیف دوم d_2 می‌توانیم به دست آوریم.

$$vir(d_2) = \{x \in \{0,1\} \times \{0,1\} : v(x) = 1\}$$

در اینجا، بردارهای توصیف انفرادی $\{x \in \{0,1\} \times \{0,1\}\}$ عبارتند از $(0,0)$ ، $(0,1)$ ، $(1,0)$ و $(1,1)$ که از آنها $x = (0,0)$ و $x = (1,1)$ در C هستند. پس $vir(d_2) = \{(0,0), (1,0), (1,1)\}$. گسترش‌های واقعی برای همه d_u ، $u = 1, \dots, 9$ در جدول ۳ ارائه شده‌اند. توجه کنید $vir(d_0) = \emptyset$ مجموعه تهی است، زیرا مقدار داده‌ای d_0 نمی‌تواند با وجود قاعده v منطقاً درست باشد.

اکنون می‌توانیم توزیع فراوانی را بیابیم. فرض کنید ابتدا این توزیع را برای Y_1 پیدا کنیم. بنابه تعریف (12) ، فراوانی‌های مشاهده شده

$$O_{Y_1}(0) = \sum_{u \in E'} \frac{|\{x \in vir(d_u) | x_{Y_1} = 0\}|}{|vir(d_u)|}$$

$$S_Z^T = (m')^{-1} \sum_{\xi_k} O_Z(\xi_k) [\xi_k - \bar{Z}]^2 \quad (18)$$

و میانه نمادین آن مقدار از ξ است که به ازای آن

$$F_Z(\xi) \geq \frac{1}{2}, F_Z(\xi) \leq \frac{1}{2} \quad (19)$$

مثال

برای توضیح مطالب بالا، فرض کنید مجموعه‌ای بزرگ از داده‌ها (مشتمل بر بیمارانی که از سوی تامین اجتماعی، خدمات بهداشتی و درمانی دریافت می‌کنند) به طریقی انبوهیده شده‌اند که داده‌های جدول ۳ را داده‌اند. این داده‌ها برآمدهای مربوط به وجود سرطان Y_1 با $\{1=بلی، 0=خیر\}$ و Y_2 و تعداد معالجات مرتبط با سرطان Y_2 با $\{0,1,2,3\}$ را درباره $m = 9$ شیئی نمادین نشان می‌دهند.

بدین ترتیب، مثلاً، $d_1 = (\{0,1\}, \{2\})$ بردار توصیف برای شیئی مندرج در سطر ۱ جدول ۳ است. پس، برای افرادی که با این توصیف نمایش داده می‌شوند، مشاهده $Y_1 = \{0,1\}$ بیان می‌کند که یا بعضی افراد سرطان دارند و برخی ندارند یا آنکه تشخیص خیر/بلی دقیق را برای افراد رده‌بندی شده در اینجا نمی‌دانیم، در حالی که مشاهده $Y_2 = \{2\}$ به ما می‌گوید که همه افراد معرفی شده با d_1 دو مورد معالجه مرتبط با سرطان داشته‌اند. برعکس، $d_2 = (\{1\}, \{2,3\})$ معرف افرادی است که همه آنان دارای تشخیص سرطانی بوده‌اند ($Y_1 = 1$) و آنهایی که ۲ یا ۳ معالجه داشته‌اند، $Y_2 = \{2,3\}$. افزون بر آن فرض کنید یک وابستگی منطقی

$$v = y_1 \in \{0\} \Rightarrow y_2 \in \{0\} \quad (20)$$

وجود داشته باشد، یعنی اگر سرطانی تشخیص داده نشده است، نباید معالجات سرطانی در بین باشد. دقت کنید که

$$y_1 \in \{0\} \Rightarrow A = \{(0,0), (0,1), (0,2), (0,3)\}$$

و

$$y_2 \in \{0\} \Rightarrow B = \{(0,0), (1,0), (2,0), (3,0)\}$$

از (20) نتیجه می‌شود که یک قاعده توصیف انفرادی x که در این قاعده صدق می‌کند، عبارت است از $\{x \in A \cap B = \{(0,0)\}\}$ یا $x \in A$ یعنی، $x \in \{(1,0), (1,1), (1,2), (1,3)\}$. گیریم همه موارد ممکنه را که در این قاعده صدق می‌کنند با $x = \{(0,0), (1,0), \dots, (1,3)\}$ نشان دهیم. رابطه (20) را در هر مورد هر $u = 1, \dots, m$ ، d_u در داده‌ها به کار خواهیم بست تا گسترش‌های واقعی $vir(d_u)$ را بیابیم.

حساب کنیم که $w \geq 0$ و $\sum w_v = 1$. مثلاً اگر شیئی $u \in E$ رده‌های C_u متشکل از افرادی از مجموعه افراد $\Omega = \{1, \dots, n\}$ باشند، یعنی $w_u = |C_u|/|\Omega|$ یک وزن ممکن عبارت است از $u = 1, \dots, m$ ، C_u رده‌های نسبی رده‌های u است یا به طور کلی‌تر، اگر هر بردار توصیف انفرادی x را به صورت واحدی ابتدایی از شیئی u در نظر بگیریم، می‌توانیم برای $u \in E$ وزن‌های

$$w_u = \frac{|vir(d_u)|}{\sum_{v \in E} |vir(d_v)|} \quad (22)$$

را به کار ببریم.

مثلاً، اگر ۱۱‌های جدول ۳ رده‌های C_u با اندازه $|C_u|$ را نمایش دهند، با وجود $|\Omega| = 1000$ همانگونه که در جدول ۳ نشان داده شده اگر از وزن $w_u = |C_u|/|\Omega|$ استفاده کنیم می‌توانیم نشان دهیم که

$$O_{Y_1}(0) = \frac{1}{1000} \left[128 \times \frac{0}{1} + 75 \times \frac{1}{3} + 249 \times \frac{0}{1} + \dots + 12 \times \frac{0}{2} \right]$$

$$= 0.229$$

$$O_{Y_1}(1) = 0.771$$

و از آنجا که فراوانی نسبی Y_1 عبارت است از،

$$[0.229, 0.771]: Y_1 \text{ نسبی}$$

به همین ترتیب،

$$O_{Y_2}(2) = 0.2395, \quad O_{Y_2}(3) = 0.495$$

$$O_{Y_2}(0) = 0.254, \quad O_{Y_2}(1) = 0.97.$$

و بنابراین فراوانی نسبی Y_2 چنین است.

فراوانی نسبی Y_3 :

$$[0.254, 0.97, 0.2395, 0.495]$$

تابع توزیع تجربی Y_3 تبدیل می‌شود به

$$F_{Y_3}(\xi) = 0.254, \quad \xi < 1$$

$$0.351, \quad 1 < \xi \leq 2$$

$$0.5905, \quad 2 \leq \xi \leq 3$$

$$1, \quad \xi \geq 3$$

به طور مشابه، می‌توانیم نشان دهیم که میانگین نمونه نمادین Y_1 برابر

است با $\bar{Y}_1 = 0.771$ و از آن عبارت است از $\bar{Y}_2 = 1/8.45$ و

واریانس نمونه نمادین Y_1 برابر است با $S_1^2 = 0.176$ و از آن Y_2

برابر است با $S_2^2 = 1/4843$.

$$= \frac{0}{1} + \frac{1}{3} + \frac{0}{1} + \frac{0}{2} + \frac{1}{1} + \frac{0}{2} + \frac{0}{2} + \frac{0}{2} = \frac{4}{3}$$

و به همین ترتیب $O_{Y_1}(1) = \frac{20}{3}$ که در آنها $E' = E - (u=0)$ به

طوری که $|E'| = 8 = m'$ ، بنابراین، توزیع فراوانی نسبی برای Y_1

طبق رابطه (۱۵) عبارت است از

فراوانی نسبی Y_1 :

$$[(0, O_{Y_1}(0)/m'), (1, O_{Y_1}(1)/m')] = \left[\left(0, \frac{1}{6}\right), \left(1, \frac{5}{6}\right) \right]$$

به طور مشابه، فراوانی‌های مشاهده شده برای مقادیر ممکن Y_2 ، یعنی

$\xi = 0, 1, 2, 3$ را به ترتیب چنین داریم

$$O_{Y_2}(0) = \frac{|\{x \in vir(d_u) | x_{Y_2} = 0\}|}{|vir(d_u)|}$$

$$= \frac{0}{1} + \frac{2}{3} + \frac{0}{1} + \frac{0}{2} + \frac{1}{1} + \frac{0}{2} + \frac{0}{2} + \frac{0}{2} = \frac{5}{3}$$

$$O_{Y_2}(2) = 2/5 \quad \text{و} \quad O_{Y_2}(3) = 2/5$$

$$O_{Y_2}(1) = \frac{4}{3}$$

و از این رو، فراوانی نسبی Y_2 طبق رابطه (۱۵) عبارت است از،

$$[(0, \frac{5}{24}), (1, \frac{1}{6}), (2, \frac{5}{16}), (3, \frac{5}{16})]: Y_2 \text{ نسبی}$$

تابع توزیع تجربی برای Y_2 طبق رابطه (۱۶) چنین است،

$$F_{Y_2}(\xi) = \frac{5}{24}, \quad \xi < 1$$

$$\frac{3}{8}, \quad 1 < \xi \leq 2$$

$$\frac{11}{16}, \quad 2 \leq \xi \leq 3$$

$$1, \quad \xi \geq 3$$

از روی رابطه (۱۷) میانگین نمونه نمادین Y_1 و Y_2 به ترتیب برابر است

با $\bar{Y}_1 = \frac{5}{6}$ و $\bar{Y}_2 = \frac{83}{84}$ از روی رابطه (۱۸) واریانس نمونه نمادین

Y_1 و Y_2 به ترتیب $S_1^2 = 0.1739$ و $S_2^2 = 1/1915$ است و میانه

Y_2 برابر ۲ است.

سرانجام ملاحظه می‌کنیم که می‌توانیم فراوانی‌های موزون،

میانگین‌های موزون و واریانس‌های موزون را به جای (۱۲) با استفاده از

$$O_Z(\xi) = \sum_{u \in E} W_u \prod Z(\xi, u) \quad (21)$$

۳-۵- متغیرهای بازه مقدار - آماره‌های تک متغیره

خاطر نشان می‌سازند که توزیع حدی واقعی Z وقتی $m \rightarrow \infty$ فقط با توزیع دقیق $f(\xi)$ مذکور در (۲۵) تقریب زده می‌شود، زیرا این توزیع به درستی توزیع یکنواخت درون هر بازه بستگی دارد.

برای ساختن بافتنگار، گیریم $I = \left[\min_{u \in E} a_u, \max_{u \in E} b_u \right]$ بازه‌ای باشد که همه مقادیر مشاهده شده Z در X ایجاد می‌کنند و فرض کنید که I را به r زیر بازه $I_g = (\xi_{g-1}, \xi_g)$ ، $g = 1, \dots, r-1$ ، و $I_r = (\xi_{r-1}, \xi_r)$ افراز کنیم. آنگاه بافتنگار Z نمایش نموداری توزیع فراوانی $\{I_g, p_g\}$ ، $g = 1, \dots, r$ است که در آن

$$p_g = \frac{1}{m} \sum_{u \in E} \frac{\|Z(u) \cap I_g\|}{\|Z(u)\|} \quad (26)$$

یعنی، p_g احتمال آن است که یک بردار توصیف انفرادی دلخواه x در بازه I_g بیفتد. اگر بخواهیم بافتنگار را با ارتفاع f_g روی بازه I_g رسم کنیم، به طوری که «مساحت» برابر p_g باشد، آنگاه

$$p_g = (\xi_g - \xi_{g-1}) \times f_g \quad (27)$$

زیر ساخت‌های ریاضی موجود مربوط به بافتنگار در دیدی (۱۹۹۵) با استفاده از قانون قوی عددهای بزرگ و مفاهیم t -نرم‌ها و t -هم‌نرم‌ها که از سوی شوایتزر و اسکالر (۱۹۸۳) ابداع شده‌اند، بسط یافته‌اند.

میانگین نمونه نمادین

برای یک متغیر بازه-مقدار Z با

$$\bar{Z} = \frac{1}{\gamma m} \sum_{u \in E} (b_u + a_u) \quad (28)$$

داده می‌شود. برای تحقیق (۲۸) به خاطر بیاورید که میانگین تجربی \bar{Z} بر حسب تابع چگالی تجربی عبارت است از،

$$\bar{Z} = \int_{-\infty}^{\infty} \xi f(\xi) d\xi$$

پس از جایگذاری مقادیر از (۲۵) داریم،

$$\begin{aligned} \bar{Z} &= \frac{1}{m} \sum_{u \in E} \int_{-\infty}^{\infty} \frac{I_u(\xi)}{\|Z(u)\|} \xi d\xi \\ &= \frac{1}{m} \sum_{u \in E} \frac{1}{b_u - a_u} \int_{\xi \in Z(u)} \xi d\xi \\ &= \frac{1}{\gamma m} \sum_{u \in E} \frac{b_u^\gamma - a_u^\gamma}{b_u - a_u} = \frac{1}{\gamma m} \sum_{u \in E} (b_u + a_u) \end{aligned}$$

به همین ترتیب، می‌توانیم واریانس نمونه نمادین را به دست آوریم.

آماره‌های توصیفی متناظر برای متغیرهای بازه-مقدار به طرز مشابه به دست می‌آیند. برتران و گوپیل (۲۰۰۰) را ببینید. گیریم به متغیر خاص $Y_j \equiv Z$ علاقه‌مند باشیم و فرض کنیم که مقدار مشاهده شده مربوط به شیئی u عبارت است از بازه $Z(u) = [a_u, b_u]$ ، به ازای $u \in E = \{1, \dots, m\}$ بردارهای توصیف انفرادی $x \in \text{vir}(d_u)$ فرض می‌شوند که به طور یکنواخت بر روی بازه $Z(u)$ توزیع شده‌اند. بنابراین، نتیجه می‌شود که برای هر ξ ،

$$P\{x \leq \xi \mid x \in \text{vir}(d_u)\} = \begin{cases} 0 & , \xi < a_u \\ \frac{\xi - a_u}{b_u - a_u} & , a_u \leq \xi < b_u \\ 1 & , \xi \geq b_u \end{cases}$$

بردار توصیف انفرادی به طور یکجا در $U_{u \in E} \text{vir}(d_u)$ اختیار می‌کند. به علاوه فرض بر آن است که هر شیئی با احتمال $1/m$ و به طور همشانس مشاهده می‌شود. بنابراین، تابع توزیع تجربی $F_Z(\xi)$ ، تابع توزیع آمیزه‌ای از m توزیع یکنواخت $\{Z(u), u = \{1, \dots, m\}\}$ است. بنابراین بنا بر (۲۳)،

$$\begin{aligned} F_Z(\xi) &= \frac{1}{m} \sum_{u \in E} P\{x \leq \xi \mid x \in \text{vir}(d_u)\} \\ &= \frac{1}{m} \left\{ \sum_{\xi \in Z(u)} \frac{\xi - a_u}{b_u - a_u} + | \{u \mid \xi \geq b_u\} | \right\} \end{aligned}$$

از آنجا که از راه مشتق‌گیری نسبت به ξ ، تابع چگالی تجربی Z را به صورت زیر به دست می‌آوریم.

$$f(\xi) = \frac{1}{m} \sum_{u: \xi \in Z(u)} \left(\frac{1}{b_u - a_u} \right) \quad (24)$$

دقت کنید که مجموع‌یابی موجود در (۲۴) تنها روی آن اشیاء u است که برای آنها $\xi \in Z(u)$ می‌توانیم (۲۴) را به صورتی دیگر بنویسیم،

$$f(\xi) = \frac{1}{m} \sum_{u \in E} \frac{I_u(\xi)}{\|Z(u)\|} \quad , \xi \in \mathcal{R} \quad (25)$$

که در آن $I_u(\cdot)$ تابع نشانگر است که ξ در بازه $Z(u)$ هست یا نیست و $\|Z(u)\|$ طول آن بازه است. شباهت با (۱۵) و (۱۶) آشکار است. برتران و گوپیل (۲۰۰۰) با استفاده از قانون عددهای بزرگ،

همچنین بازه‌ها و احتمال‌های مربوط به متغیر بازه-مقدار Y_p را که معرف فشار خون سیستولیک است و Y_p را که نمایشگر فشار خون دیاستولیک است، برای همان ۱۰ بیمارستان نشان می‌دهد.

با استفاده از روابط (۲۸) و (۲۹) به ترتیب میانگین و واریانس نمونه نمادین را به دست می‌آوریم. بدین ترتیب، میانگین نبض $\bar{Y}_1 = 79/1$ با واریانس $S_1^2 = 162/29$ ؛ میانگین فشار خون سیستولیک $\bar{Y}_p = 131/3$ با واریانس $S_p^2 = 495/41$ و میانگین فشار خون دیاستولیک $\bar{Y}_p = 84/6$ با واریانس $S_p^2 = 182/44$.

نظیر مفاهیم و روشهای محاسباتی بالا را می‌توان در مورد متغیرهای مدی، چند مقدره، و بازه-مقدار نیز تعریف کرد. همچنین می‌توان این گونه مفاهیم و روشها را به آماره‌های توصیفی دو متغیری مربوط به متغیرهای چند مقدره و بازه-مقدار تعمیم داد.

راستایی دیگر برای تحویل مشاهدات به مجموعه‌هایی کوچکتر استفاده از روشهای تحویل داده‌ها نظیر تحلیل مؤلفه‌های اصلی و خوشه‌بندی نمادین است. کوشش‌هایی برای گسترش به داده‌های سه طرفه نیز به عمل آمده‌اند که برای اطلاع از آنها به بیلارد و دیدی (۲۰۰۱) مراجعه شود.

۶. نتیجه‌گیری

با وجود قالبهای داده‌ای و اندازه‌های مجموعه‌های داده‌ها، نیاز به ابداع روشهای آماری به منظور تحلیل آنها اهمیت روزافزون دارد. جهان آمار گنجینه‌ای از روش‌شناسی‌ها را در طی قرن بیستم میلادی ابداع کرده است. روشهایی که (در مقام قیاس) عمدتاً محدود به مجموعه‌های کوچک داده‌ها و قالبهای داده‌ای کلاسیک (اسکالری، برداری و ماتریسی)‌اند. این مقاله به اجمال فرمول‌بندی و ساخت اشیاء ورده‌های نمادین را مرور کردیم. به کوتاهی به روشهای تحلیل نمادین اشاره کرده، با مثالهایی طرز کار را بیان کردیم. نیاز به پژوهش و ابداع بیشتر به منظور انجام استنباطهای آماری درباره داده‌های نمادین آشکار بوده، افقی است که اخیراً پیش چشم آمارشناسان گشوده شده است.

(۲۹)

$$S^2 = \frac{1}{fm} \sum_{u \in E} (b_u + a_u)^2 - \frac{1}{fm^2} \sum_{u \in E} (b_u + a_u)^2$$

در مورد متغیرهای چند مقدره، اگر شینی u دارای ناسازگاری‌های درونی نسبت به یک قاعده منطقی باشد، یعنی اگر u چنان باشد که $|vir(d_u)| = 0$ ، آنگاه مجموع‌یابی (۲۸) و (۲۹) فقط روی u هایی است که برای آنها $|vir(d_u)| \neq 0$ یعنی بر روی $u \in E'$ است، و به جای m عدد m' گذاشته شود (که برابر با تعداد اشیاء u موجود در E' است).

در زیر قرار می‌گذاریم که m و E به آن u هایی اشاره کنند که برای آنها این قواعد صادق‌اند.

مثال: برای توضیح مطلب، داده‌های راجو (۱۹۹۷) را که در جدول ۴ آمده‌اند، در نظر بگیرید.

در آنجا نبض (Y_1)، فشار خون سیستولیک (Y_p) و دیاستولیک (Y_r) به صورت متغیرهای بازه-مقدار به ازای هر $u = 1, \dots, 10$ بیمار ثبت شده‌اند.

گیریم داده‌های نبض (Y_1) را برداریم. مجموعه کل داده‌ها بازه $I = [4, 112]$ را ایجاد می‌کند که در آن

$$\min_{u \in E} a_u = 44, \quad \max_{u \in E} b_u = 112$$

فرض کنید می‌خواهیم بافتنگاری را روی $r = 8$ بازه $I_8 = [11, 120], \dots, I_1 = [4, 50]$ بسازیم. با استفاده از رابطه (۲۶) می‌توانیم احتمال p_g را که بردار توصیف انفرادی دلخواه x در بازه $I_g, g = 1, \dots, 8$ بیفتد، حساب کنیم. مثلاً وقتی $g = 4$ احتمال آنکه یک x در بازه $I_4 = [70, 80]$ بیفتد، برابر است با

$$p_4 = \frac{1}{10} \left\{ 0 + \frac{2}{12} + \frac{10}{24} + \frac{10}{42} + \frac{2}{18} + \frac{10}{30} + \frac{8}{28} + \frac{4}{22} + 0 + 0 \right\} = 0.1611$$

از آنجا، بنابر رابطه (۲۷) می‌توانیم ارتفاع f_g از بافتنگار رسم شده را برای آن بازه برابر با

$$f_g = p_g / (\xi_g - \xi_{g-1})$$

یعنی،

$$f_4 = (0.1611) / 10 = 0.01611$$

حساب کنیم. خلاصه‌ای از مقادیر محاسبه شده P_g و f_g در جدول ۵ داده شده‌اند و نمودار بافتنگار در شکل ۲ رسم شده است. جدول ۵

جدول ۱

i	Y_1	Y_2	Y_3	...	Y_8	...	Y_{28}	Y_{29}	Y_{30}	شرح متغیرها
۱	بوشهر	M	۲۴		۰		N	N	۰	Y_1 : شهر محل سکونت
۲	بوشهر	M	۵۶		۲		N	N	۱	Y_2 : جنس (F,M)
۳	ایلام	F	۴۸		۲		y	N	۱	Y_3 : سن به سال
۴	ایلام	F	۴۷		۱		Y	۰	۰	Y_4 : نژاد (سفید W، سیاه B، دیگر O)
۵	آبادان	M	۷۹		۴		N	۰	۰	Y_5 : تأهل (متاهل S، متأهل M)
۶	لامرد	M	۱۲		۰		N	۰	۰	Y_6 : تعداد والدین زنده (۰، ۱، ...)
۷	لامرد	M	۶۷		۰		Y	۶	۰	Y_7 : تعداد برادران و خواهران (۰، ۱، ...)
۸	لامرد	F	۷۳		۴		N	۰	۲	Y_8 : تعداد فرزندان (۰، ۱، ...)
۹	شیراز	M	۲۹		۲		N	۰	۰	.
۱۰	کرمان	F	۴۴		۳		y	۰	۰	.
۰	Y_{28} : تشخیص سرطان (بلی y، خیر N)
۰	Y_{29} : تعداد دفعات معالجه سرطان پستان (۱، ۲، ... = N مصداق ندارد)
۰	Y_{30} : تعداد دفعات معالجه سرطان ریه (۱، ۲، ...)

جدول ۲- نمونه‌ای از داده‌های نمادین

U	سن	فشار خون	شهر	نوع سرطان	جنس	وزن
۱	[۲۰,۳۰]	(۷۹,۱۲۰)	بوشهر	{سرطان مغزی}	{زن}	[۷۰,۸۰]
۲	[۵۰,۶۰]	(۹۰,۱۳۰)	بوشهر	{کبد، ریه}	{مرد}	[۶۰,۷۵]
۳	[۴۵,۵۵]	(۸۰,۱۳۰)	ایلام	{پروستات}	{مرد}	[۶۰,۶۵]
۴	[۴۷,۴۷]	(۸۶,۱۲۱)	ایلام	{ریه (۱-p)، پستان p}	{زن}	[۵۵,۶۵]
.
.
.

جدول ۳- مجموعه‌ای از داده‌های نمادین

u	Y_1	Y_2	$vir(d_u)$	$ vir(d_u) $	$ C_u $
۱	{۰,۱}	{۲}	{(۱,۲)}	۱	۱۲۸
۲	{۰,۱}	{۰,۱}	{(۰,۰),(۱,۰),(۱,۱)}	۳	۷۵
۳	{۰,۱}	{۳}	{(۱,۳)}	۱	۲۴۹
۴	{۰,۱}	{۲,۳}	{(۱,۲),(۱,۳)}	۲	۱۱۳
۵	{۰}	{۱}	ϕ	۰	۲
۶	{۰}	{۰,۱}	{(۰,۰)}	۱	۲۰۴
۷	{۱}	{۲,۳}	{(۱,۲),(۱,۳)}	۲	۸۷
۸	{۱}	{۱,۲}	{(۱,۱),(۱,۲)}	۲	۲۳
۹	{۱}	{۰,۳}	{(۱,۱),(۱,۳)}	۲	۱۲۱

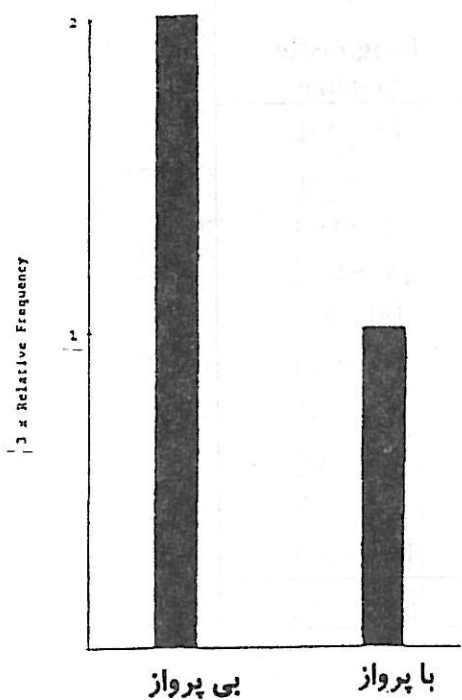
جدول ۴- سه متغیر نمادین

u	$:Y_1$ Pulse rate	$:Y_2$ Systolic Pressure	$:Y_3$ Diagnostic ressure
۱	[۴۴-۶۸]	[۹۰-۱۱۰]	[۵۰-۷۰]
۲	[۶۰-۷۲]	[۹۰-۱۳۰]	[۷۰-۹۰]
۳	[۵۶-۹۰]	[۱۴۰-۱۸۰]	[۹۰-۱۰۰]
۴	[۷۰-۱۱۲]	[۱۱۰-۱۴۲]	[۸۰-۱۰۸]
۵	[۵۴-۷۲]	[۹۰-۱۰۰]	[۵۰-۷۰]
۶	[۷۰-۱۰۰]	[۱۳۴-۱۴۲]	[۸۰-۱۱۰]
۷	[۷۲-۱۰۰]	[۱۳۰-۱۶۰]	[۷۶-۹۰]
۸	[۷۶-۹۸]	[۱۱۰-۱۹۰]	[۷۰-۱۱۰]
۹	[۸۶-۹۶]	[۱۲۸-۱۸۰]	[۹۰-۱۱۰]
۱۰	[۸۶-۱۰۰]	[۱۱۰-۱۵۰]	[۷۸-۱۰۰]
۱۱	[۶۳-۷۵]	[۶۰-۱۰۰]	[۱۴۰-۱۵۰]

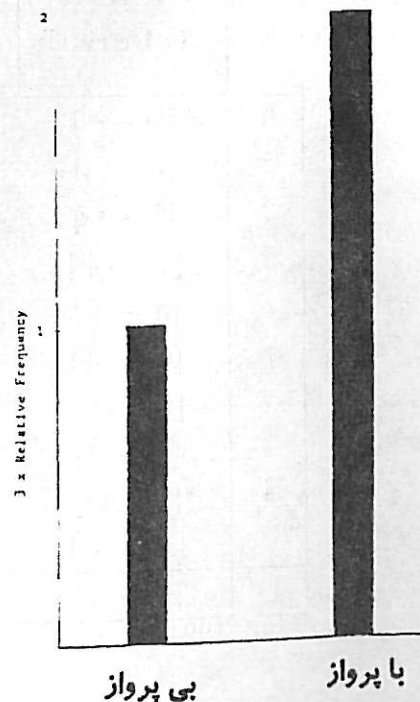
جدول ۵

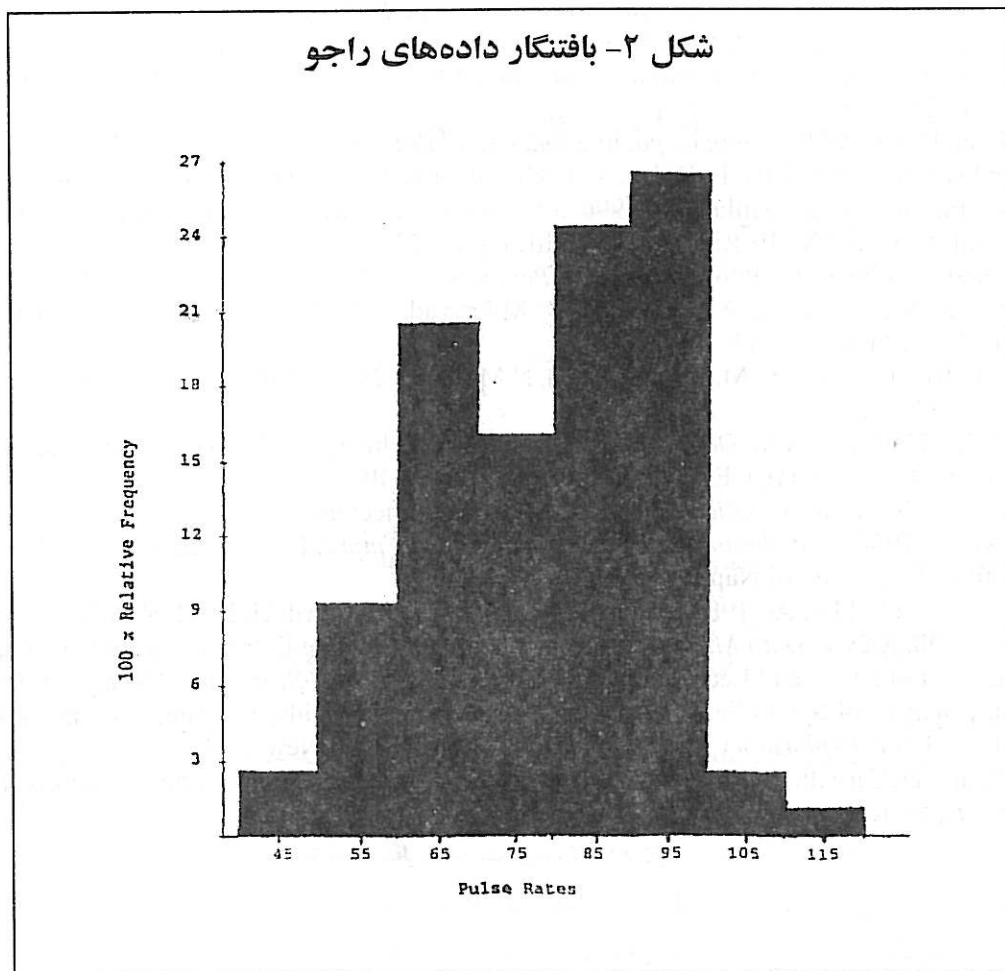
Y_1 : Pulse Rate		Y_2 : Systolic pressure		Y_3 : Diagnostic pressure	
I_g	P_g	I_g	P_g	I_g	P_g
[۴,۵۰)	۰/۰۲۵۰	[۹,۱۰۰)	۰/۱۷۵۰	[۵,۶۰)	۰/۱۰۰۰
[۵,۶۰)	۰/۰۸۶۸	[۱۰,۱۱۰)	۰/۰۷۵۰	[۶,۷۰)	۰/۱۰۰۰
[۶,۷۰)	۰/۲۰۱۶	[۱۱,۱۲۰)	۰/۰۹۳۸	[۷,۸۰)	۰/۱۱۲۷
[۷,۸۰)	۰/۱۶۱۱	[۱۲,۱۳۰)	۰/۰۹۳۸	[۸,۹۰)	۰/۲۶۰۹
[۸,۹۰)	۰/۲۳۶۳	[۱۳,۱۴۰)	۰/۱۸۱۸	[۹,۱۰۰)	۰/۲۸۹۵
[۹,۱۰۰)	۰/۲۶۰۶	[۱۴,۱۵۰)	۰/۱۵۰۹	[۱۰,۱۱۰)	۰/۱۳۶۹
[۱۰,۱۱۰)	۰/۰۲۳۸	[۱۵,۱۶۰)	۰/۰۹۴۶		
[۱۱,۱۲۰)	۰/۰۰۴۸	[۱۶,۱۷۰)	۰/۰۶۱۳		
		[۱۷,۱۸۰)	۰/۰۶۱۳		
		[۱۸,۱۹۰)	۰/۰۱۲۵		
Mean	$\bar{Y}_1 = 79/1$	$\bar{Y}_2 = 131/1$	$\bar{Y}_3 = 84/6$		
Variance	$S_1^2 = 162/29$	$S_2^2 = 162/29$	$S_3^2 = 182/44$		
Covariance	$S_{12} = 194/170$	$S_{23} = 257/920$	$S_{13} = 141/040$		
Correlation	$r(Y_1, Y_2) = 0/685$	$r(Y_2, Y_3) = 0/858$	$r(Y_1, Y_3) = 0/820$		

شکل ۱-ب بافتنگار پرندگان با پرواز/ بی پرواز



شکل ۱-الف بافتنگار گونه‌های با پرواز/ بی پرواز





مراجع

- [1] Aristotle (IVBC), 1994. *Des Categories De l'Interpretation*, Organ, Librarie Philosophique Journal Vrin.
- [2] Arnault, A. and Nicole, P., 1662. *La Logic ou l'Art Depensur*, Reprinted by Forman, Stuttgart (1965).
- [3] Bertrand, P., 1995. *Structural Properties of Pyramidal Clustering In: Partitioning Data Sets*, (eds. I. Corc, P. Hansen and B. Julesz), American Mathematical Society, 19, pp. 352-353.
- [4] Bertrand, P. and Goupil, F., 2000. *Descriptive Statistics for Symbolic Data*, In: Analysis of Symbolic Data (eds. H.H. Bock and E. Diday), Springer-Verlag, pp. 103-124.
- [5] Billard, L. and Diday, E., 2001. *From the Statistics of Data to the Statistics of knowledge: Symbolic Data Analysis*, Manuscript.
- [6] Bock, H.H. and Diday, E., 2000. *Symbolic Objects*, In: Analysis of Symbolic Data, (eds. H.H. Bock and E. Diday), Springer-Verlag, pp. 54-77.
- [7] Choquet, G., 1954. *Theory of Capacities*, Annals de l'Institute Fourier, 5, pp. 131-295.
- [8] Chouakria, A., 1998. *Extensions des methodes d'analyse factorielle a des donnee de type intervalle*, Ph.D. Thesis, University of Paris.
- [9] DeCarvalho, F.A.T., 1994. *Proximity Coefficients Boolean Symbolic Objects*, In: New Approaches in Classification and Data Analysis, (eds. E. Diday, Y. Lechevallier, M. Schader, P. Bertrand and B. Burtschy). Springer-Verlag, pp. 387-394.
- [10] DeCarvalho, F.A.T., 1995. *Histograms in Symbolic Data Analysis*, Annala of Operation Research. 55, pp.299-322.
- [11] Diday, E., 1987. *Introduction a l'Approach Symbolique en Analyse des Donnees*, Premere Jouneles Symbolique-Namerique, CEREMADE, Universite Paris-Dauphine, pp. 21-56.

- [12] Diday, E., (ed.), 1989. *Data Analysis, Learning Symbolic and Numerical Knowledge*, Nova Science, Antibes.
- [13] Diday, E., 1990. *Knowledge Representation and Symbolic Data Analysis*, In: *Knowledge Data and Computer Assisted Decisions*, (eds. M. Schader and W. Gaul). Springer-Verlag, pp. 17-34.
- [14] Diday, E., 1995. *Probabilistic, Possibilistic and Belief Objects for Knowledge Analysis*, *Annals of Operation Research*, 55, pp. 227-276.
- [15] Diday, E. and Emilion, R., 1996. *Capacities and Credibilities in Analysis of Probabilistic Objects*, In: *Ordinal and Symbolic Data Analysis*, (Eds. E. Diday, Y. Lechevallier and O. Opitz), Springer-Verlag, pp. 13-30.
- [16] Diday, E., Emilion, R. and Hillali, Y., 1996. *Symbolic Data Analysis of Probabilistic Objects by Capacities and Credibilities*, *Atti Della XXXVIII Riunione Scientifica*, pp. 5-22.
- [17] Elder, J. and Pregibon, D., 1996. *A Statistical Perspective on Knowledge Discover in Databases*, In: *Advances in Knowledge Discovery and Data Mining* (eds. U.M.Fayyad, G. Piatetsky-Shapiro, P. Symth and R. Uthurusamy), AAAI Press, 83-113.
- [18] Hand, D.J., Blunt, G., Kelly, M.G. and Adams, N.M., 2000. *Data Mining for Fun and Profit*, *Statistical Science*, 15, pp. 111-131.
- [19] Raju, S.R.K., 1997. *Symbolic Data Analysis in Cardiology*, In: *Symbolic Data Analysis and its Applications* (eds. E. Diaay and K.C. Gowda), CEREMADE, Paris, pp. 245-249.
- [20] Schafer, G., 1976. *A Mathematical Theory of Evidence*, Princeton.
- [21] Schweizer, B., 1984. *Distributions are the Numbers of the Future*, In: *Proceedings The Mathematics of Fuzzy Systems Meeting*. University of Naples, pp. 137-149.
- [22] Schweizer, B. and Sklar, A., 1983. *Probabilistic Metric Spaces*, North Holland, New York.
- [23] Siebs, A., 1998. *KESO: Data Mining and Statistics*, In: *Knowledge Extraction from Statistical Data*, pp. 1-13.
- [24] Stephan, V., Herbial, G. and Lechevallier, Y., 2000. *Generation of Symbolic Objects from Relational Databases*, In: *Analysis of Symbolic Data*, (Eds. H. -H. Bock and E. Diday). Springer-Verlag, pp. 78-105.
- [25] Tukey, J.W., 1977. *Exploratory Data Analysis*, Addison-Wesley, New York.
- [26] Verbe, R. and DeCarvallo, F.A.T., 1998 *Dependence Rules Influence on Factorial Representation of Boolean Symbolic Objects*, In: *Knowledge Extraction from Statistical Data*, pp. 14-23.