



تحلیل نمونه های متفاوت در داده کاوی

سوسن حسینی

دانشجوی کارشناسی ارشد، نرم افزار کامپیوتر، دانشگاه آزاد اسلامی مشهد

Sosan.hosseini@gmail.com

چکیده

داده کاوی وسیله ای برای کشف داده ها از منابع اطلاعاتی محسوب می گردد و به بررسی نمونه های مختلف و متشابه در پایگاه داده ها می پردازد. هر نمونه با چندین ویژگی مشخص شده است به طوری که انواع مختلفی از مقادیر برای هر ویژگی وجود دارد، که می توان با تحلیلی صحیح بر روی نمونه ها به دانش مفیدی و جدیدی دست یافت.

کلید واژه- داده کاوی، کشف داده، نمونه های متفاوت، پایگاه داده، دانش مفید.

1- مقدمه

داده کاوی (Data Mining)، یکی از مهمترین روش هایی است که به وسیله آن الگوهای مفید نهان در داده ها کشف شده، اطلاعاتی را در اختیار کاربران و تحلیل گران قرار می دهد تا بر اساس آنها تصمیمات مهم و حیاتی اتخاذ شوند. اطلاعات و دانش بدست آمده از این طریق، جهت کاربردهای گوناگونی از قبیل تحلیل بازار، تشخیص کلاهبرداری، حفظ مشتری، کنترل تولید و... را شامل می شود. [1] به کل فرایند بکارگیری متدولوژی مبتنی بر کامپیوتر از جمله روش های جدید برای دریافت دانش و اطلاعات ازداده ها داده کاوی گفته می شود. درحقیقت، داده کاوی جستجوی لازم برای یافتن اطلاعات ارزشمند و غیربدهی از میان حجم زیاد داده ها می باشد. فعالیتهای داده کاوی را می توان در دو گروه زیر طبقه بندی کرد:

الف- داده کاوی پیشگویانه، مدلی از سیستم ارائه می دهد که توسط مجموعه داده های مشخصی توصیف می شود.

ب- داده کاوی توصیفی، اطلاعات جدید و غیربدهی را بر اساس مجموعه داده های موجود ارائه می دهد.

داده کاوی بخشی از یک فرایند بزرگتر موسوم به کشف دانش در پایگاه داده (Knowledge Discovery in Databases) می باشد اما در اکثر موارد از واژه داده کاوی بجای KDD استفاده می شود [2]. ادامه این مقاله اینگونه سازماندهی شده است دربخش دوم مروری برنمونه های در داده کاوی و دربخش سوم به تحلیل نمونه های متفاوت و درانتها به نتیجه گیری پرداخته خواهد شد.

۲- بررسی نمونه

هر نمونه با چندین ویژگی مشخص می شود، طوری که انواع مختلفی از مقادیر برای هر ویژگی وجود دارد [3]. دو نوع رایج آن categorical و numeric می باشد، مقادیر numeric شامل متغیرهای real و یا متغیرهای integer می باشد، نوع numeric دو ویژگی، مهم در رابطه دارد:

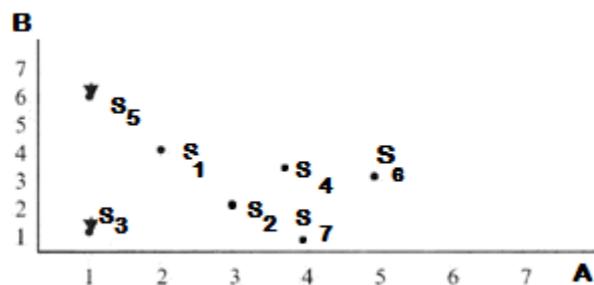
۱- منظم. فاصله در مقابل categorical، هیچ یک از این دو رابطه را ندارند. راه دیگر [4] برای نمونه بندی متغیرها، مبتنی بر مقادیرشان می باشد، به طوری که آیا متغیرهای پیوسته اندیاگسسته متغیرهای پیوسته به متغیرهای metric (کمی) نیز معروف اند. آنها هم با مقیاس فاصله ای و هم با مقیاس نسبت اندازه گیری می

شوند. تفاوت این روش با مقیاس در چگونگی تعریف نقطه صفی در مقیاس می باشد، نقطه صفر در مقیاس فاصله به طور قرار دادی مشخص شده است و چیزی نیست که بتوان آن را با اندازه گیری مشخص کرد. برای مثال، مقیاس دما، که درجه صفر، رابطه نسبت برای متغیرهایی که با استفاده از مقیاس فاصله اندازه گیری شده اند درست کار نمی کند. برای مثال ۸۰ درجه فاینرتهات بر دو برابر بودن گرمای ۴۰ درجه دلالت نمی کند. در مقابل مقیاس نسبی دارای نقطه صفر مطلق می باشد و به تبع آن، رابطه نسبی برای این نوع اندازه گیری به خوبی کار می کند، کمیت های مانند ارتفاع، قد و حقوق با استفاده از این مقیاس اندازه گیری می شود. متغیرهای پیوسته به صورت مجموعه داده بزرگ با مقادیری چون **real, Integer** نمایش داده می شوند. متغیرهای گسسته (متغیرهای کیفی) یکی از دو نوع مقیاس های کیفی ترتیبی یا اسمی اندازه گیری می شوند یا مقادیرشان مشخص می شود. یک مقیاس اسمی یک مقیاس بدون ترتیب می باشد که نشانه ها، کاراکترها و اعداد مختلف حالت های (مقادیر) مختلف متغیر را نشان می دهد، یک مثال از متغیرهای اسمی (غیر ترتیبی) ، یک کاربرد با مقادیر صنعتی، تجاری و خانگی می باشد. این مقادیر می تواند به صورتی الفبایی **C,B,A** و یا اعداد به صورت ۱ و ۲ و ۳ کد شوند، اما آنها دارای هیچگونه ترتیبی نمی باشند. مثال دیگر، **Zip-code** می باشد. در هر دو مثال، اعداد استفاده شده برای یک متغیر، نسبت به هم هیچگونه نظم و رابطه خاصی ندارند. یک مقیاس ترتیبی شامل، نظم و گسستگی می باشد، مانند **Ranking** یک متغیر ترتیبی یک متغیر **categorical** می باشد که دارای رابطه نظم می باشد، اما فاقد رابطه فاصله می باشد. مثالی از این نوع متغیر، رتبه دانش آموزان داخل یک کلاس و مدل های طلا و نقره و برنز در رقابت های ورزشی می باشد. برای مقیاس منظم [5]، خطی بودن نیاز ضروری نمی باشد. برای مثال، فاصله بین رتبه های ۴ و ۵ دانش آموزان، نیاز ندارد که فاصله ۱۵ و ۱۴ یکسان می باشد. عموماً متغیرهای منظم یک متغیر عددی را مجموعه کوچکی از مقادیر **overlap** متناظر با مقادیر یک متغیر منظم می کند. یک کلاس ویژه از متغیرهای گسسته، متغیرهای دوره ای یا تناوبی می باشد. متغیرهای تناوبی دارای ویژگی رابطه فاصله ای باشند اما رابطه نظم در آنها وجود ندارد. برای مثال، روزهای هفته، دوشنبه و سه شنبه به عنوان مقادیر این متغیر می باشند که فاصله بین آنها یک می باشد، اما دوشنبه می تواند بعد و یا قبل از جمع قرار گیرد. در نهایت، یک بعد دیگر از طبقه بندی داده ها مبنی رفتارشان نسبت به زمان می باشد. بعضی از داده ها با توجه به زمان تغییر نمی کنند و داده هایی استاتیک می باشند. در طرف دیگر بعضی از داده ها با توجه به زمان تغییر می کنند و داده های **Data mining** برای داده های استاتیک مناسب می باشند و بعضی از پیش پردازش ها و رسیدگی های ویژه ای برای داده های دینامیک لازم می باشد. اکثر مشکلات **Data mining** ناشی از وجود مقدار زیادی از مثال ها با انواع مختلفی از ویژگی ها می باشد.

3-تحلیل نمونه متفاوت:

بسیاری از مواقع، در مجموعه داده های بزرگ، نمونه هایی وجود دارد که از رفتار کلی مدل داده ها تبعیت نمی کنند. چنین نمونه هایی که تا حد زیادی متفاوت اند، مابقی مجموعه داده ها ناهماهنگ هستند، نمونه های متفاوت نامیده می باشند، که می توانند توسط خطای ارزیابی ایجاد گردند یا نتیجه تنوع فطری داده ها بوده اند. به عبارت دیگر، اگر در پایگاه داده ها، تعداد بچه های یک فرد ۲۵ باشد این داده غیر معمول است و باید چک شود. این اشتباه می تواند یک خطای چاپی یا می تواند صحیح بوده و نشان دهنده تنوع واقع برای خصیصه داده شده باشد. بسیاری از الگوریتم های داده کاوی در تلاش اند، تا تأثیر نمونه های متفاوت بر نمونه نهایی را کاهش داده یا آنها را در، مراحل پیش از پردازش ریشه کن سازند. تحلیل گر داده کاوی باید در ریشه کنی اتوماتیک نمونه های متفاوت بسیار مراقب باشد، زیرا اگر داده ها درست باشند می تواند منجر به فقدان اطلاعات مهم مخفی، گردد. برخی از کاربردهای داده کاوی بر روی شناسایی نمونه متفاوت متمرکز و این امر نتیجه فوری یک تحلیل داده می باشد. به عنوان مثال در حین شناسایی تعاملات کارت اعتباری جعلی در یک بانک نمونه های متفاوت

نمونه هایی معمولی هستند که ممکن است فعالیت جعلی را نشان دهند و کل فرآیند داده کاوی شوند نمونه های متفاوت چندان مفید نیستند و بیشتر نتیجه خطاهایی در جمع آوری داده ها و سپس مشخصه ای از یک مجموعه داده ها هستند. شناسایی نمونه متفاوت و رفع بالقوه آن از یک مجموعه داده ها می تواند به عنوان فرآیندی از انتخاب k در خارج از نمونه هایی توصیف شود که با توجه به داده های باقی مانده به طور قابل توجه ای ناهماهنگ هستند. مشکل تشریح نمونه های متفاوت، بی اهمیت نمی باشد، خصوصاً در نمونه های چند بعدی، شیوه های تجسم داده ها که شناسایی نمونه متفاوت برای ابعاد یک تا سه مفید هستند، در داده های چند بعدی ضعیف تر می باشند، نمونه ای از یک تجسم نمونه های دو بعدی و شناسایی بصری نمونه های متفاوت در نمودار (1) داده می شود.



نمودار (1): تجسم مجموعه داده های دو بعدی برای شناسایی نمونه های متفاوت [6].

ساده ترین رویکرد جهت شناسایی نمونه متفاوت برای نمونه ها، بر اساس آمار می باشد. با فرض اینکه تنوع مقادیر داده شود یافتن پارامترهای اساسی آماری از قبیل مقدار میانگین و واریانس ضروری است. بر اساس این مقادیر و عدد مورد انتظار (یا پیش بینی شده) در مورد نمونه های متفاوت، اثبات مقدار ورودی به عنوان عملکردی از واریانس، امکان پذیر است. تمام نمونه های خارج از مقدار ورودی، کاندیداهای نمونه های متفاوت می باشند. مشکل اساسی موجود در این شیوه ساده، یک برداشت پیشوا در مورد توزیع داده ها است. در بسیاری از نمونه های دنیای واقعی، توزیع داده ها ممکن است شناخته شده نباشد. اگر مجموعه ی معین از داده نشان دهنده مشخصه سن یا، بیست مقدار متفاوت باشد و مقدار ورودی توزیع نرمال داده، به صورت (انحراف استاندارد $+2$ * میانگین = ورودی) انتخاب شود، تمام داده هایی که خارج از دامنه $\{ ۱۳۱/۲ \}$ و $\{ -۵۴/۱ \}$ می باشند، نمونه های متفاوت بالقوه خواهند بود. (دانش اضافی از مشخصه های خصیصه سن همیشه بیشتر از صفر است). ممکن است دامنه به صورت $\{ ۱۳۱/۲ \}$ و $\{ ۰ \}$ بیشتر کاهش دهد. بر اساس ضابطه معین در نمونه ما سه مقدار وجود دارد که نمونه های متفاوت می باشند: ۱۵۶، ۱۳۹، ۶۷. شناسایی نمونه متفاوت بر اساس فاصله، شیوه ثانویه می باشد که برخی از محدودیت های تحمیلی از سوی رویکرد آماری را از بین می برد. مهم ترین تفاوت این است که این شیوه برای نمونه های چند بعدی قابل کاربردی باشد در حالی که توصیف کنندگان تنها یک بعد واحد یا چندین بعد را، به صورت جداگانه ارزیابی می نماید. دشواری، محاسباتی این شیوه، ارزیابی مقادیر فاصله بین، نمونه ها، در یک مجموعه n بعدی از داده ها می باشد. بنابراین یک نمونه s_i در یک مجموعه داده های s ، یک نمونه متفاوت می باشد. اگر حداقل یک بخش p ، از نمونه های s در فاصله ای بیشتر از d قرار داشته باشد و ملاک شناسایی نمونه متفاوت بر حسب دو پارامتر p, d می باشد که قبل از استفاده از دانش مربوط به داده ها داده شود یا در حین تکرار رویکرد آزمون و خطا، انتخاب گویاترین نمونه متفاوت، تغییر داده شوند. برای نشان دادن رویکرد، می توان مجموعه ای از نمونه های دو بعدی را تحلیل کرد. نتیجه نمونه های متفاوت، مقادیر ورودی $d \geq 3, p \geq 4$ می باشد. مسافت های اقلیدوسی $d = \sqrt{[(X1 - X2)^2 + (y1 - y2)^2]}$ ، برای مجموعه s در جدول ۱ دیده می شود.

	S_1	S_2	S_3	S_4	S_5	S_6	S_7
S_1		2.236	3.162	2.236	2.236	3.162	2.828
S_2			2.236	1.414	4.472	2.236	1.000
S_3				3.605	5.000	4.472	3.162
S_4					4.242	1.000	1.000
S_5						5.000	5.000
S_6							1.414

جدول ۱: مسافت های اقلیدسی برای $p \geq 4$ و $d \geq 3$ [7]

بر اساس جدول (۱) می توان با مسافت ورودی داده شده ($d=3$) برای هر نمونه، یک مقدار، را برای پارامتر p ، بالای مقدار ورودی است: ($p=4$). همین نتایج را می توان از طریق بازدید بصری یک مجموعه داده بدست آورد که در جدول ۲ ارائه شده است. البته مجموعه داده تعیین شده بسیار کوچک است و یک نمودار گرافیکی دو بعدی امکان پذیر و مفید است. برای تحلیل داده های n بعدی در دنیای واقعی، فرآیند تجسم بسیار دشوارتر است و رویکردهای تحلیل در شناسایی نمونه متفاوت غالباً کاربردی تر و قابل اطمینان تر هستند.

Sample	p
S_1	2
S_2	1
S_3	5
S_4	2
S_5	5
S_6	3

جدول ۲: نتایج مسافت های اقلیدسی [8]

تکنیک های بر اساس انحراف سومین طبقه از شیوه های شناسایی نمونه متفاوت هستند. این تکنیک ها را شبیه سازی می نمایند که بشر می تواند در آن، نمونه های غیر معمول را، از مجموعه نمونه های مشابه دیگر جدا سازی نماید. این شیوه ها خصوصیات اساسی مجموعه ستونه را تعریف می نماید و تمام نمونه هایی که از این خصوصیات منحرف می گردند، نمونه های متفاوت می باشند. تکنیک غیر متوالی یک رویکرد امکان پذیر است که بر اساس یک عملکرد عدم شباهت می باشد. یک عملکرد امکان پذیر، عدم شباهت برای مجموعه ای معین از نمونه های n عبارت است از انحراف کلی مجموعه نمونه می باشد و با تعریف کردن، کوچکترین زیر مجموعه نمونه ها است، که رفع کردن آن منجر به بزرگترین کاهش عملکرد عدم شباهت، مجموعه باقیمانده می گردد. وظیفه کلی برای پیدا کردن نمونه های متفاوت با، این شیوه می تواند بسیار پیچیده باشد (رشد سریع ترکیبی گزینه های متفاوت مجموعه نمونه های متفاوت بالقوه - مجموعه بسیار معروف استثنای می گردد) و می تواند به صورت تئوری به عنوان یک مشکل سخت NP تعریف شود. اگر به یک پاسخ کمتر از حد مطلوب رضایت بدهیم، سختی الگوریتم می تواند تا سطح تک بعدی کاهش یابد.

۴- نتیجه گیری

شناخت ویژگی ها، در تحلیل نمونه متفاوت بسیار مفید می باشد بر اساس ضابطه در، نمونه های متفاوت آزمایش شده (۱۳۹، ۱۵۶، ۶۷ -) می توان نتیجه گیری کرد، هر سه تا دارای خطاهای چاپی می باشند (داده های وارد شده با اعداد اضافی یا یک علامت « - » اضافی). با استفاده از یک رویکرد متوالی با استفاده کردن از شیوه مشتاق، الگوریتم با انتخاب نمونه ای در هر مرحله که سبب بیشترین کاهش در واریانس کلی می گردد، نمونه به نمونه اندازه به صورت متوالی کاهش می دهد.

مراجع

- [1] K.R.Venugopal, K.G.srinivasa and L.M.patnaik, soft computing for data mining applications, pp. 85-87, 2009.
- [2] IH, Frank "An instruction to data mining", pp.55-59, 2007.
- [3] Brunk.c, Kelly.j, "an integrated system for data mining", proc.3, int.1.conf.on knowledge discovery and data minig, pp.135-138, 1997.
- [4] Javad azimi, Reza Davoodi, Morteza analoui, "fast Convergence Clustering ensemble, In Conference of Data mining and Data warehouse", pp.14-16, 2006.
- [5] W. OCK .Day, "consensus methods as tools for data analysis.in classification and related method for data analysis", Elsiver science publishers B.V, pp.317-320, 1988.
- [6] Brunk.c, Kelly.j, "an integrated system for data mining", proc.3, int.1.conf.on knowledge discovery and data minig, pp.135-138, 1999.
- [7] David Hand, Heikki Mannila, "Principles of data mining", pp.246-248, 2001.
- [8] J.Han, M.Kamber, "Date Mining: Goncepts and Techinqueus", pp.54-56, 2000.