



سمینار دوره کارشناسی ارشد

مهندسی کامپیوتر - نرم افزار

عنوان سمینار:

کاوش مجموعه داده های سودمند

دانشجو:

پروین تقوی

استاد راهنما:

دکتر نگین دانشپور

تیر ۹۹

صلى الله عليه وسلم

تقدیر و تشکر

با سپاس فراوان از راهنمایی‌ها و زحمات استاد محترم خانم دکتر دانشپور که مرا در این مسیر یاری نمودند و خانم دکتر ترابی که با مشورت خود مرا راهنمایی کردند.

چکیده

امروزه با وجود حجم بالای داده، یافتن رابطه‌ی با معنی بین داده‌ها و استخراج دانش مفید امری ضروری و کارآمد برای تصمیم‌گیری است. در اغلب روش‌های گذشته که به استخراج این روابط با معنی و یا داده‌های مفید پرداخته‌اند، به سود آیتم‌ها توجهی نشده‌است. حال آن که داشتن اطلاعات درباره داده‌های سودمند می‌تواند در تصمیم‌گیری‌ها کارآمدتر باشد. در این تحقیق سعی در بررسی روش‌هایی داریم که با توجه به میزان سودمندی هر آیتم، مجموعه‌ای از داده‌هایی را که نسبت به یک آستانه سودمندی از پیش تعیین شده سود بیشتری دارند کشف کنند.

کلیدواژه‌ها

داده‌کاوی، مجموعه داده‌ی سودمند، مجموعه داده‌ی تکراری، آستانه‌ی سودمندی

فهرست عناوین

فصل ۱: مقدمه و شرح مسئله	۱
۱-۱- مقدمه	۲
۱-۲- روش‌های پایه‌ی مسئله‌ی استخراج مجموعه‌داده‌های سودمند	۳
۱-۲-۱- روش کشف مجموعه‌داده‌های پرتکرار	۳
۱-۲-۲- روش کشف قواعد وابستگی	۴
۱-۳- تاریخچه و تعریف موضوع	۴
۱-۴- جمع‌بندی	۵
فصل ۲: مجموعه‌داده‌های پرسود؛ شرح راهکارها	۷
۲-۱- مقدمه	۸
۲-۲- تعریف اولیه مسئله	۹
۲-۳- روش حل دوفاز	۱۱
۲-۴- روش حل تک‌فاز	۱۲
۲-۴-۱- روش‌های مبتنی بر لیست سودمندی	۱۲
۲-۴-۲- روش‌های مبتنی بر پایگاه داده‌ی واکنشی شده	۱۵
۲-۵- بسط مسئله‌ی استخراج مجموعه‌داده‌های پرسود	۱۸

۱-۵-۲- مسئله k مجموعه داده‌ی پرسود	۱۸
۲-۵-۲- مسئله مجموعه داده‌های پرسود موجود در قفسه	۲۰
۳-۵-۲- مسئله مجموعه داده‌های پرسود با سود مثبت و منفی	۲۰
۴-۵-۲- مسئله مجموعه داده‌های متناوب پرسود	۲۵
۶-۲- جمع‌بندی	۲۸
فصل ۳: نتیجه‌گیری و جمع‌بندی مطالب	۲۹
۱-۳- مقدمه	۳۰
۲-۳- علت تفاوت روش‌ها	۳۰
۳-۳- مقایسه روش‌های تک‌فاز و دو‌فاز	۳۲
۱-۳-۳- ویژگی روش‌های دو‌فاز	۳۳
۲-۳-۳- ویژگی روش‌های تک‌فاز	۳۳
۴-۳- جمع‌بندی	۳۴
۵-۳- نتیجه‌گیری	۳۶
فهرست منابع و مراجع	۳۷

فصل اول

مقدمه و شرح مسئله

1-1- مقدمه

از سال ۱۹۵۱ به بعد که رایانه، در تحلیل و ذخیره‌سازی داده‌ها به کار رفت، حجم اطلاعات ذخیره شده در آن پس از حدود ۱۰ سال دو برابر شد و هم‌زمان با پیشرفت فناوری اطلاعات، حجم داده‌ها در پایگاه داده‌ها هر دو سال یک‌بار، دو برابر شد [۱] و همچنان با سرعت بیش‌تری نسبت به گذشته حجم اطلاعات ذخیره شده بیش‌تر و بیش‌تر می‌شود. باوجود منابع اطلاعاتی مانند شبکه جهانی وب، سیستم‌های یکپارچه اطلاعاتی، سیستم‌های یکپارچه بانکی و تجارت الکترونیک لحظه‌به‌لحظه به حجم داده‌ها در پایگاه داده‌ها اضافه شده و باعث به وجود آمدن انبارهای عظیمی از داده‌ها شده‌است، به‌طوری‌که ضرورت کشف و استخراج سریع و دقیق دانش از این پایگاه‌های داده را بیش‌ازپیش نمایان کرده است. شدت رقابت‌ها در عرصه‌های علمی، اجتماعی، اقتصادی، سیاسی و نظامی نیز اهمیت سرعت یا زمان دسترسی به اطلاعات را دوچندان کرده است. بنابراین نیاز به طراحی سیستم‌هایی که قادر به اکتشاف سریع اطلاعات مورد علاقه کاربران با تأکید بر حداقل مداخله انسانی باشند از یک‌سو و روی آوردن به روش‌های تحلیل متناسب با حجم داده‌های حجیم از سوی دیگر، به‌خوبی احساس می‌شود. در حال حاضر، داده‌کاوی^۳ مهم‌ترین فناوری برای بهره‌وری مؤثر، صحیح و سریع از داده‌های حجیم است و اهمیت آن رو به فزونی است. داده‌کاوی فرآیندی پیچیده جهت شناسایی الگوها و مدل‌های صحیح، جدید و به‌صورت بالقوه مفید، در حجم وسیعی از داده می‌باشد، به طریقی که این الگوها و مدل‌ها برای انسان‌ها قابل‌درک باشند. داده‌ها اغلب حجیم می‌باشند و به تنهایی قابل استفاده نیستند، اما دانش نهفته در داده‌ها قابل استفاده می‌باشد. بنابراین بهره‌گیری از قدرت فرآیند داده‌کاوی جهت شناسایی الگوها و مدل‌ها و نیز ارتباط عناصر مختلف در پایگاه داده جهت کشف دانش نهفته در داده‌ها و نهایتاً تبدیل

¹ Data

² Database

³ Datamining

داده به اطلاعات، روزبه‌روز ضروری‌تر می‌شود. در داده‌کاوی معمولاً به کشف الگوهای مفید از میان داده‌ها اشاره می‌شود. منظور از الگوی مفید مدلی در داده‌ها است که ارتباط میان یک زیرمجموعه از داده‌ها را توصیف می‌کند و معتبر، ساده، قابل‌فهم و جدید است. روش‌های مختلفی برای کشف این الگوها ارائه شده است که در ادامه به شرح مختصری از دو روشی که پایه‌ی مسئله‌ی استخراج مجموعه‌داده‌های پرسود هستند، می‌پردازیم.

1-2- روش‌های پایه‌ی مسئله‌ی استخراج مجموعه‌داده‌های سودمند

مزیت این روش‌ها که مبتنی بر الگویابی هستند بر چندین رویکرد داده‌کاوی دیگر این است که کشف الگوها نوعی یادگیری بدون نظارت است که نیازی به داده‌های برچسب‌دار ندارد. برای حل مسئله‌ی استخراج مجموعه‌داده‌های سودمند^۱ از بسط روش کشف مجموعه‌داده‌های پرتکرار استفاده شده است و از روش کشف قواعد وابستگی نیز ایده گرفته شده است.

۱-۲-۱- روش کشف مجموعه‌داده‌های پرتکرار

یکی از مسائل مشهور در زمینه‌ی داده‌کاوی مسئله‌ی استخراج مجموعه‌داده‌های پرتکرار (FIM) بوده است. هدف آن پیدا کردن گروهی از آیتم‌هاست که به صورت تکراری در پایگاه داده وجود دارند. از ویژگی‌های این مسئله بسته‌شدن رو به پایین است و بیان می‌کند که ابرمجموعه‌های یک مجموعه‌داده‌ی غیرتکراری، غیرتکراری هستند و زیرمجموعه‌های یک مجموعه‌داده‌ی مکرر، پرتکرار هستند [۲]. اگر فرض کنیم پایگاه داده‌ی ما شامل تراکنش‌های خرید باشد در این صورت یکی از مشکلات این روش این است که به تعداد اقلام خریداری‌شده در هر تراکنش توجه نمی‌شود. در واقع

¹ High Utility Itemset Mining

² Frequent Itemsets Mining

تنها وجود و یا عدم وجود اقلام مهم است. مسئله‌ی دیگر نیز این است که در این روش تمام آیتم‌ها از ارزش یکسانی برخوردار هستند.

۲-۱-۲- روش کشف قواعد وابستگی

کاوش قواعد وابستگی^۱ آن دسته از مسائل داده‌کاوی را شامل می‌شود که در آن به دنبال استخراج و تعریف قواعد و الگوهایی هستیم که توصیف دقیق‌تری را از فضای حاکم بر داده‌ها ارائه می‌دهند [۳]. کاربردهای گسترده‌ی این روش‌ها در هوش تجاری، شبکه‌های اجتماعی و مجازی، تجارت الکترونیک، صنعت بانکداری، وب‌کاوی و همچنین کاربرد آن در پیش‌بینی باعث می‌شود که این مسئله اهمیت زیادی داشته‌باشد. برای مثال با استفاده از این روش و با در نظر گرفتن یک پایگاه داده‌ی تراکنشی، رابطه‌ای را کشف می‌کنیم که در آن با خرید پنیر به احتمال ۷۰ درصد گردو و نان هم به فروش می‌رسد. با وجود تمام این‌ها در این روش نیز تمام اقلام از ارزش یکسانی برخوردار هستند، در صورتی که در واقعیت هر جنس و یا رویدادی از ارزش خاص خود برخوردار است.

3-1- تاریخچه و تعریف موضوع

مسئله‌ی استخراج مجموعه‌داده‌های پرسود زمانی مورد توجه قرار گرفت که یافتن داده‌های پرتکرار برای کاربردهای واقعی خیلی کارآمد نبودند و سود یکی از فاکتورهای مهم برای نتایج بهتر و موثرتر بود. اصطلاح مجموعه‌داده‌های پرسود در ابتدا در سال ۲۰۰۳ مطرح شد و در سال ۲۰۰۵ تعریف اولیه‌ی آن ارائه شد [۲]. در این مسئله یک پایگاه داده که شامل آیتم‌ها و تعداد آن‌ها است، یک جدول که شامل ارزش هر آیتم است و آستانه سودمندی که توسط کاربر تعیین می‌شود به عنوان ورودی به الگوریتم داده می‌شود. در خروجی اگر سود مجموعه‌داده‌ای از این آستانه بیشتر بود به عنوان یک داده‌ی پرسود استخراج می‌شود. در واقع راه حل‌های این مسئله به طور کلی به دو دسته تقسیم می‌شوند که عبارتند

¹ Association Rule Mining

از روش‌هایی که در دوفاز اجرا می‌شوند و روش‌های تک‌فاز [۲]. اولین الگوریتمی که برای این مسئله ارائه شد الگوریتم [۴] two phase نام داشت که در دوفاز اجرا می‌شد. در روش دوفاز یک معیار¹ TWU تعریف می‌شود. در فاز اول تمام فضای جست‌وجو بررسی می‌شود و آن دسته از مجموعه داده‌هایی که مقدار سود آن‌ها از TWU شان بیشتر بود به عنوان مجموعه داده‌های کاندید در نظر گرفته می‌شوند. در فاز دوم سود دقیق اعضای مجموعه‌ی کاندید را با اسکن پایگاه داده محاسبه می‌کنند.

در روش تک فاز که در یک مرحله انجام می‌شود ابتدا با روش‌های هرس مناسب فضای جست‌وجو را کاهش می‌دهند و سپس سود آیت‌هایی که باقی مانده است را با یک بار اسکن پایگاه داده محاسبه می‌کنند و با مقدار آستانه سودمندی مقایسه می‌کنند [۲].

فاکتورهایی که باعث می‌شود روش‌های ارائه‌شده با هم متفاوت باشند عبارتند از:

- تعداد مراحل اجرای الگوریتم
- نوع ساختمان داده‌ی به کاررفته
- روش جست‌وجو
- نمایش پایگاه داده
- معیارهای هرس فضای جست‌وجو

1-4- جمع‌بندی

در این فصل به چرایی ضرورت الگویابی در داده‌های حجیم پرداختیم و همچنین پیشینه و ویژگی‌های کلی مسئله‌ی HUIM را بیان کردیم.

¹ Transaction Weighted Utility

در فصل بعدی به شرح راهکارهای ارائه شده برای مسئله‌ی استخراج مجموعه داده‌های پرسود پرداخته می‌شود و فاکتورهای بالا به طور کامل توضیح داده می‌شوند.

فصل دوم

مجموعه داده‌های پرسود؛ شرح راهکارها

1-2- مقدمه

در علم داده‌کاوی، کاوش الگوهایی با کارایی و سود بالا در حال ظهور است، که شامل کشف الگوهایی است که دارای اهمیت بالایی در پایگاه‌های اطلاعاتی هستند. کارایی یک الگو را می‌توان از نظر معیارهای عینی مختلف مانند سود، فرکانس و وزن اندازه گرفت. در میان انواع مختلف الگوهای کاربردی که می‌توانند در پایگاه‌های اطلاعاتی کشف شوند، داده‌هایی با سود بالا بیشتر مورد مطالعه قرار می‌گیرند. یک مجموعه داده با سود بالا مجموعه‌ای از مقادیر است که در پایگاه داده ظاهر می‌شود، سود آن توسط یک تابع منفعت اندازه‌گیری می‌شود و برای کاربر دارای اهمیت بالایی است. استخراج مجموعه داده‌های سودمند مشکل استخراج داده‌های مکرر را با در نظر گرفتن مقادیر و ارزش آیتم‌ها حل می‌کند.

مسئله‌ی استخراج مجموعه داده‌های سودمند کاربردهای مختلفی در زمینه‌ی تجارت الکترونیک، پزشکی، وب‌کاوی و بررسی جریان کلیک دارد. برای مثال در زمینه پزشکی می‌توان بررسی نمود که کدام روش درمان و یا داروها برای یک بیماری خاص مفیدتر است. همچنین با بررسی جریان کلیک کاربر در صفحات وب می‌توان دریافت که به چه موضوعاتی علاقه‌ی بیشتری دارد و آن‌ها را در ورود بعدی به کاربر معرفی نمود و یا محصولات مورد علاقه‌ی وی را پیدا کرد که سودآور نیز هستند. یک کاربرد معروف استخراج مجموعه داده‌های سودمند کشف همه مجموعه موارد خریداری شده توسط مشتریان است که سود بالایی را به دست می‌آورد. در این فصل به توضیح کامل راهکارهای ارائه شده برای حل مسئله‌ی کشف مجموعه داده‌های پرسود پرداخته می‌شود. در بخش تعریف مسئله، نحوه‌ی کلی محاسبه‌ی سود اقلام، کران‌های بالا برای کاهش فضای جست‌وجو در هر دو روش تک‌فاز و دوفاز توضیح داده خواهد شد. در ادامه به صورت مبسوط به بررسی روش‌های تک‌فاز و دوفاز می‌پردازیم و در بخش پایانی بسط‌های مختلفی از مسئله‌ی پایه را توضیح خواهیم داد.

2-2- تعریف اولیه مسئله

تعریف‌هایی که در این بخش به آن پرداخته می‌شود تعاریف کلی هستند که در همه‌ی کارهای انجام شده یکسان است بنابراین برای فهم بیشتر و یکپارچگی موضوع از تعاریف موجود در مقاله‌ی [۵] استفاده شده‌است.

فرض کنید I مجموعه‌ای شامل تمام اقلام موجود در تراکنش‌ها باشد. و $D = \{T_1, T_2, \dots, T_n\}$ نیز پایگاه داده‌ای شامل تمام تراکنش‌ها است. هر آیت $i \in I$ نیز به یک مقدار مثبت $p(i)$ که سود خارجی حاصل از فروش یک قلم از آن است، مرتبط می‌شود و در هر تراکنش مانند T_c تعدادی از هر آیت مانند i وجود دارد که این مقدار مثبت را با $q(i, T_c)$ نشان می‌دهند. توجه داشته باشید که در مسئله‌ی استخراج مجموعه‌داده‌های پرسود پایه، سود هر آیت مثبت در نظر گرفته می‌شود.

پایگاه داده‌ی تراکنشی مورد استفاده در این نوشته در جدول (۱-۲) و سود خارجی هر آیت در جدول (۲-۲) نشان داده شده‌است.

جدول ۱-۲: پایگاه داده‌ی تراکنشی

TID	Transaction
T_1	$(a, 1)(c, 1)(d, 1)$
T_2	$(a, 2)(c, 6)(e, 2)(g, 5)$
T_3	$(a, 1)(b, 2)(c, 1)(d, 6)(e, 1)(f, 5)$
T_4	$(b, 4)(c, 3)(d, 3)(e, 1)$
T_5	$(b, 2)(c, 2)(e, 1)(g, 2)$

جدول ۲-۲: سود خارجی اقلام

Item	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
Profit	5	2	1	2	3	1	1

تعریف ۱: (سود یک آیتم و سود یک مجموعه داده). سود آیتمی مانند i در یک تراکنشی مانند T_c به روش زیر به دست می آید.

$$U(i, T_c) = q(i, T_c) \times p(i) \quad (۱-۲)$$

سود مجموعه داده ای مانند $X \subseteq I$ که در تراکنش T_c حضور دارد، به صورت زیر به دست می آید. که در این جا $g(X)$ مجموعه ای از تراکنش هاست که X در آن ها حضور دارد.

$$u(X, T_c) = \sum_{i \in X \wedge T_c \in g(X)} U(i, T_c) \quad (۲-۲)$$

تعریف ۲: (سود کلی یک مجموعه داده در کل پایگاه داده). سود کلی یک مجموعه داده از جمع سود آن در هر تراکنشی که حضور دارد به وجود می آید.

$$U(X) = \sum_{X \subseteq T_c \wedge T_c \in D} u(X, T_c) \quad (۳-۲)$$

تعریف ۳: (استخراج مجموعه داده های پرسود) زمانی مجموعه داده ای مانند X را پرسود نامیم که سود آن از مقدار آستانه ی سودمندی تعریف شده توسط کاربر ($minutil$) کمتر نباشد.

تعریف ۴: سود کلی یک تراکنش مانند T_c از جمع سود حاصل از هر آیتم موجود در آن به دست می آید.

^۱Minimum utility

$$TU(T_c) = \sum_{x \in T_c} U(x, T_c) \quad (4-2)$$

تعاریفی که در بالا آورده شده است برای تمام روش‌های تک‌فار و دوفاز مورد استفاده قرار می‌گیرد.

2-3- روش حل دوفاز

همان‌طور که قبلاً هم گفته شد الگوریتم‌های دوفاز از راه حل‌های اولیه‌ی ارائه شده در این زمینه بوده‌اند. عواملی که باعث می‌شوند روش‌های دوفاز با هم متفاوت باشند بستگی به این دارد که از چه روش جست‌وجویی استفاده شود و این که روش ارائه شده بسطی از چه الگوریتمی باشد. به طور کلی الگوریتم‌های دوفاز یا بسطی از الگوریتم Apriori هستند و یا بسطی از روش FP-Growth می‌باشند [۲]. هر دوی این الگوریتم‌ها روش‌هایی برای حل مسئله‌ی FIM هستند. در ابتدا از بسط روش Apriori استفاده شد که در این روش از جست‌وجوی سطحی برای کاوش فضای جست‌وجو بهره برده شد. مشکل این روش این بود که در جست‌وجوی سطحی حافظه‌ی زیادی مورد نیاز بود. بنابراین از بسط روش مبتنی بر رشد استفاده کردند. در این روش از جست‌وجوی عمقی استفاده شد که میزان مصرف حافظه در آن کمتر بوده است. بنا به تعریفی که در فصل گذشته برای روش‌های دوفاز بیان شد، روش‌های محاسبه‌ی سودی که برای اولین بار در این دسته از روش‌ها مورد استفاده قرار گرفت در زیر آورده شده است.

تعریف ۴: سودمندی وزنی تراکنش برای یک مجموعه داده مانند X از جمع سود تمام تراکنش‌هایی که X در آن‌ها حضور داشته است به دست می‌آید.

$$TWU(X) = \sum_{T_c \in g(X)} TU(T_c) \quad (5-2)$$

در جدول (۳-۲) مقدار TWU محاسبه شده برای تک‌آیتم‌ها آورده شده است.

جدول ۲-۳: مقدار TWU برای تک آیتم‌ها

<i>Item</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
<i>TWU</i>	65	41	96	58	88	30	38

ویژگی ۱: اگر مقدار TWU برای یک مجموعه داده X از مقدار minutil کمتر باشد، آنگاه X و ابرمجموعه‌های آن جز موارد پرسود نخواهند بود.

در فصل قبل گفته شد که مقدار TWU معیاری است که بتوان با استفاده از آن و ویژگی ۱ اعضای مجموعه‌ی کاندید را در روش‌های دوفاز مشخص کرد. سپس با اسکن پایگاه داده و بر اساس رابطه‌ی ۲-۳ مقدار سود هر عضو کاندید محاسبه می‌شود. به دلیل این که روش‌های دوفاز جز راه‌حل‌های اولیه برای حل این مسئله هستند، قدیمی بوده و از توضیح و بسط آن اجتناب می‌کنیم.

2-4- روش حل تک‌فاز

به طور معمول در روش‌های تک‌فاز از دو ساختمان داده استفاده می‌شود. یکی از ساختمان‌های داده لیست سودمندی است و دیگری پایگاه داده‌ی واکشی شده می‌باشد. الگوریتم HUI-Miner که در [۶] آمده‌است، اولین روشی است که به صورت تک فاز ارائه شد و کاستی‌های روش دوفاز را پوشش داد. در این روش از لیست سودمندی استفاده شده‌است. الگوریتم EFIM که در [۵] آمده‌است نیز از روش‌های پایه در زمینه تک‌فاز بودن است که در آن از پایگاه داده‌ی واکشی شده استفاده شده‌است. در ادامه به توضیح لیست سودمندی و پایگاه داده‌ی واکشی شده می‌پردازیم.

۲-۴-۱- روش‌های مبتنی بر لیست سودمندی

در مقالات [۷]، [۸]، و [۹] که از این روش استفاده کرده‌اند سعی شده‌است که با ارائه‌ی راه حلی هزینه ادغام لیست سودمندی را کاهش دهند. برای مثال در [۷] با ارائه‌ی روشی هر آیتم در هر تراکنش

با آیت‌های قبل خود در ارتباط است و مقایسه‌ی دو لیست سودمندی را به‌طور هوشمندانه‌ای انجام می‌دهد تا در مصرف زمان و حافظه صرفه‌جویی شود. یا در [۸] و [۹] با ارائه‌ی روشی سعی در کمتر کردن حافظه‌ی مصرفی توسط لیست سودمندی می‌شود. در مقالات دیگر مانند [۱۰]، [۱۱]، [۱۲]، [۱۳]، [۱۷] و [۱۵] هر یک با استفاده از این روش سعی در حل بسطی از مسئله‌ی پایه HUIM دارند. تعریف‌هایی که در این بخش به آن پرداخته می‌شود تعاریف کلی هستند که در همه‌ی کارهای مبتنی بر لیست سودمندی یکسان است. بنابراین برای فهم بیشتر و یکپارچگی موضوع از تعاریف موجود در مقاله‌ی [۶] استفاده شده‌است.

در این روش هر مجموعه داده با یک لیست سودمندی مرتبط است. لیست سودمندی برای تک آیت‌ها لیست سودمندی اولیه است که می‌تواند توسط دو اسکن پایگاه داده ساخته شوند. در اسکن اول مقدار TWU برای اقلام محاسبه می‌شود. اگر مقدار آن از آستانه سودمندی کمتر باشد دیگر در روند استخراج در نظر گرفته نمی‌شود. فرض کنیم \succ ترتیبی برای آیت‌ها باشد. معیاری که برای مرتب کردن آیت‌ها استفاده می‌کنیم بر اساس معیار TWU آن‌هاست که کمتر از آستانه سودمندی نبوند و به صورت صعودی مرتب می‌شوند. حال بنا بر این توضیحات پایگاه داده بر اساس ترتیب جدید مرتب می‌شود. این مراحل مرتب‌سازی طی اسکن دوم پایگاه داده انجام می‌شود. برای پایگاه داده‌های بزرگ، به منظور کاهش بیشتر فضای جستجو و افزایش کارایی HUIM به حدود بالای دقیق‌تری نسبت به TWU نیاز است. محاسبه حدود بالایی جدید متکی بر سود باقی مانده است که به صورت زیر تعریف می‌شود. فرض کنیم \succ ترتیبی برای آیت‌ها باشد. معیاری که برای مرتب کردن آیت‌ها استفاده می‌کنیم بر اساس معیار TWU آن‌هاست و به صورت صعودی مرتب می‌شوند. ترتیب آیت‌های مرتبط با پایگاه داده‌ی جدول ۱ به صورت $c \succ e \succ a \succ d \succ b \succ g \succ f$ خواهد بود.

لیست سودمندی مجموعه داده‌ی X یک سه تایی شامل (tid, iutil, rutil) خواهد بود که tid معرف شماره تراکنشی است که X در آن حضور دارد، iutil سود حاصل از مجموعه داده‌ی X موجود در

تراکنش مورد نظر است و ru_{til} سود باقی مانده‌ی حاصل از آیتم‌هایی است که با توجه به ترتیب تعریف شده بعد از X می‌آیند. تعریف سود باقی مانده در زیر آمده است.

تعریف ۵: (سود باقی مانده برای یک مجموعه داده). فرض کنید X یک مجموعه داده در تراکنش T_c است. سود باقی مانده برای آن برابر با مجموع سود آیتم‌هایی است که بعد از تمام آیتم‌های موجود در X بر اساس ترتیب \succ می‌آید.

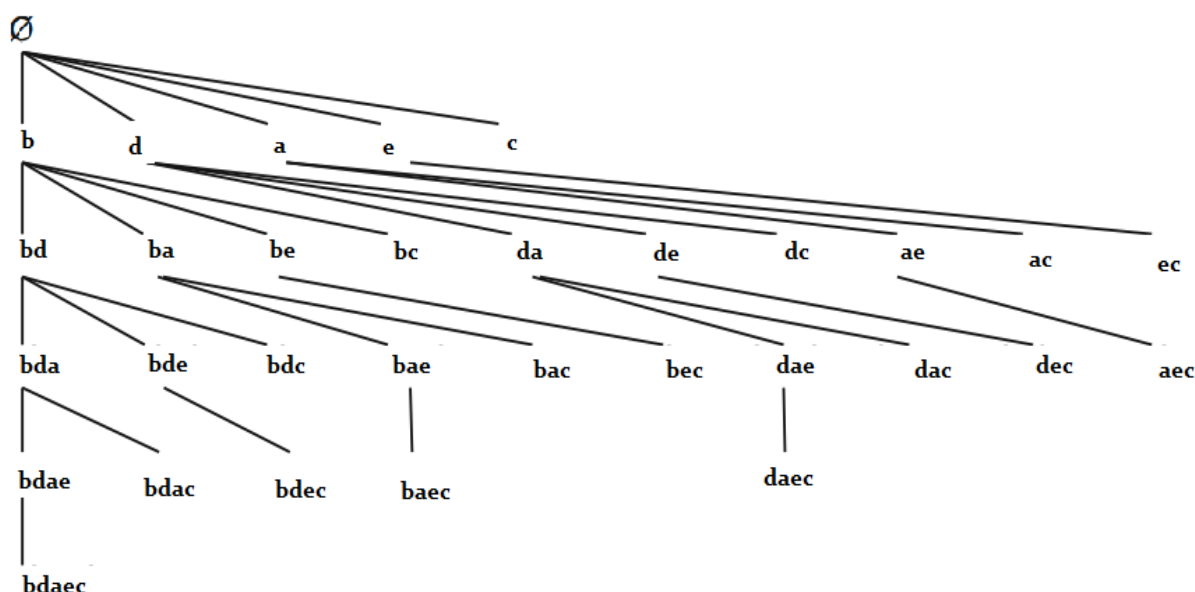
$$ru(X, T_c) = \sum_{i \in T_c \wedge i \succ x \forall x \in X} u(i, T_c) \quad (۶-۲)$$

برای ساخت لیست سودمندی مجموعه داده‌های دوتایی از روی لیست سودمندی تک آیتم‌ها، مقدار tid لیست سودمندی دو آیتمی که قرار است با هم مجموعه داده‌ی دوتایی را بسازند، مقایسه می‌شوند. این مقایسه به دلیل مرتب بودن شماره تراکنش‌ها به صورت صعودی، به روش دو طرفه انجام می‌شود. سپس به ازای تراکنش‌های هم شماره مقدار $iutil$ از جمع سود آن‌ها و مقدار ru_{til} از لیست سودمندی آیتمی که اولویت بیشتری داشته حاصل می‌شود.

برای ساخت لیست سودمندی مجموعه داده‌هایی با طول بیش از دو، مانند روش قبلی ابتدا تراکنش‌هایی که شماره یکسانی در هر دو لیست دارند مشخص می‌شوند. سپس برای محاسبه سود جدید، سود هر دو لیست پایه با هم جمع می‌شود و مقدار سود $k-2$ آیتم قبلی از این عدد کم می‌شود. محاسبه سود باقی مانده نیز مانند قبل است.

روند اجرای الگوریتم‌های مبتنی بر سودمندی به این صورت است که برای جست و جوی مجموعه داده‌های پرسود از درخت شمارشی و با استفاده از جست و جوی اول عمق عمل می‌کنند. اقلام موجود در درخت به ترتیب صعودی TWU مرتب می‌شوند. بدین ترتیب اگر جمع $iutil$ و ru_{til} یک مجموعه داده از مقدار آستانه سودمندی کمتر باشد آنگاه می‌توان خودش و زیردرخت‌اش را حذف نمود. اگر فرض کنیم مقدار آستانه سودمندی تعریف شده ۴۰ باشد آنگاه دو آیتم f و g حذف می‌شوند زیرا مقدار TWU

آن‌ها کمتر از *minutil* است. با این فرض در شکل ۱-۲ درخت شمارشی مرتبط با پایگاه داده‌ی جدول (۱-۲) آمده است.



شکل ۱-۲: درخت شمارشی برای نمایش فضای جست‌وجو

روند الگوریتم کاوش بدین صورت است که لیست سودمندی مجموعه داده‌ی مورد بررسی و تمام بسط‌های یک‌تایی آن به عنوان ورودی دریافت می‌شوند. اگر جمع مقدار *iutil* کمتر از آستانه سودمندی نبود آن را به عنوان مجموعه داده‌ی پرسود شناسایی می‌کند. در ادامه برای بسط بیشتر مجموعه داده شرط لازم این است که جمع *iutil* و *rutil* آن کمتر از آستانه سودمندی نباشد. اگر شرط برقرار بود در این صورت از بسط یک‌تایی، بسط‌های دو تایی ساخته می‌شوند و بررسی ادامه پیدا می‌کند تا تمام مجموعه داده‌های پرسود کشف شوند.

۲-۴-۲- روش‌های مبتنی بر پایگاه داده‌ی واكشی‌شده

در این روش‌ها که مقالات [۵]، [۱۶] و [۱۷] از آن بهره برده‌اند، برای عمل جست‌وجو از روش جست‌وجوی عمق اول استفاده می‌شود. تفاوت این روش‌ها در این است که برای مثال در [۱۷] به حل بسطی از مسئله‌ی پایه که مربوط به یافتن k مجموعه داده‌ی پرسود است، پرداخته شده است، و در [۱۶]

سعی شده است تا با ارائه‌ی یک الگوریتم موازی روش ارائه شده در [۵] را بهبود بخشد. تعریف‌هایی که در این بخش به آن پرداخته می‌شود تعاریف کلی هستند که در همه‌ی کارهای مبتنی بر پایگاه داده‌ی واکنشی شده یکسان است. بنابراین برای فهم بیشتر و یکپارچگی موضوع از تعاریف موجود در مقاله‌ی [۵] استفاده شده است. برای بسط هر مجموعه داده، با یک آیتم به صورت بازگشتی عمل می‌شود. برای این منظور مجموعه‌ی $E(X) = \{z \mid z \in I \wedge z \succ x, x \in X\}$ تعریف شده است.

تعریف ۶: (مجموعه‌ی بسط برای یک مجموعه داده). تمام آیتم‌هایی که در بسط یک مجموعه داده مانند X می‌توانند مورد استفاده قرار گیرند در یک مجموعه $E(X) = \{z \mid z \in I \wedge z \succ x, x \in X\}$ تعریف می‌شود.

برای کاهش هزینه اسکن پایگاه داده، کاهش اندازه پایگاه داده بهتر است. برای این منظور هنگامی که یک مجموعه α در طول جست‌وجوی اولیه در نظر گرفته می‌شود، همه آیتم‌هایی که در مجموعه‌ی بسط α نیستند را می‌توان در هنگام اسکن پایگاه داده برای محاسبه سودمندی مجموعه داده‌ها در زیر درخت α نادیده گرفت.

تعریف ۷: واکنشی تراکنش T برای مجموعه داده‌ای مانند α به صورت $\alpha - T = \{i \mid i \in T \wedge i \in E(\alpha)\}$ تعریف می‌شود. حال واکنشی پایگاه داده برای X به صورت $\alpha - D = \{\alpha - T \mid T \in D \wedge \alpha - T \neq \emptyset\}$ تعریف می‌شود.

این روش معمولاً هزینه اسکن پایگاه داده را تا حد زیادی کاهش می‌دهد زیرا با بزرگ‌تر شدن تعداد آیتم‌هایی که مورد بررسی قرار می‌گیرند، تراکنش‌ها کوچک‌تر می‌شوند. اما روش پیاده‌سازی آن مهم است. یک روش ساده و ناکارآمد ایجاد کپی‌های فیزیکی از تراکنش‌ها برای هر واکنشی است. برای این منظور ابتدا اقلام در هر تراکنش بر اساس ترتیبشان که قبلاً ذکر شد، مرتب می‌شوند. سپس هر تراکنش واکنشی شده با یک اشاره‌گر آفست بر روی پایگاه داده‌ی اصلی مورد اشاره قرار می‌گیرد. به عبارتی دیگر برای هر مجموعه داده به تعداد تراکنش‌هایی که در آن حضور دارد اشاره‌گر خواهیم داشت و به اقلامی

که در آن تراکنش اولویتشان بیشتر است اشاره می‌شود. در این صورت پیچیدگی محاسبه‌ی پایگاه‌داده‌ی واکنشی شده $O(nw)$ خواهد بود که n نشان‌دهنده‌ی تعداد تراکنش‌ها و w نشان‌دهنده‌ی طول متوسط هر تراکنش است. با این حال، با بررسی مجموعه‌داده‌های بزرگ‌تر، اندازه پایگاه‌های داده پیش‌بینی‌شده کاهش می‌یابد.

روش دیگر برای کاهش هزینه اسکن پایگاه‌داده بر این اساس است که پایگاه‌های داده تراکنشی اغلب دارای تراکنش‌های یکسان هستند. این روش شامل شناسایی این تراکنش‌ها و سپس جایگزین کردن آن‌ها با تراکنش‌های تکی، و در عین حال جمع سود داخلی آن‌ها است. ادغام تراکنش‌های یکسان اندازه پایگاه‌داده را کاهش می‌دهد. اما این کاهش کوچک است اگر پایگاه‌داده شامل تعداد کمی تراکنش یکسان باشد. برای کاهش بیشتر پایگاه‌داده، تراکنش‌ها را در پایگاه‌های داده واکنشی شده ادغام می‌کنیم. این امر به طور کلی به کاهش بسیار بالاتری دست می‌یابد زیرا تراکنش‌های پیش‌بینی‌شده کوچک‌تر از تراکنش‌های اصلی هستند، و بنابراین احتمال یکسان بودن آن‌ها بیشتر است. با این حال، مشکل اصلی پیاده‌سازی کارآمد آن است. روش ساده برای شناسایی تراکنش‌های یکسان، مقایسه تمام تراکنش‌ها با یکدیگر است. برای پیاده‌سازی کارآمد، این روش تراکنش‌ها را بر اساس یکسری تعاریف مرتب می‌کند. مثلاً برای تراکنش‌هایی که یکسان هستند، آن تراکنشی از اولویت بیشتری برخوردار است که شماره بزرگ‌تری دارد؛ یا اگر اشتراک دو تراکنش، تراکنشی شود که تعداد ارقام آن کمتر است؛ تراکنش بزرگ‌تر از اولویت بیشتری برخوردار است. به این ترتیب تراکنش‌ها در پایگاه داده تنها یک بار مرتب می‌شوند. در این حالت می‌شود که از زمان محاسبه‌ی این کار صرف نظر کرد. با استفاده از ویژگی فوق، تمام تراکنش‌های یکسان در یک پایگاه داده واکنشی‌شده را می‌توان تنها با مقایسه هر تراکنش با تراکنش بعدی در پایگاه‌داده شناسایی کرد. هر مقایسه بین دو تراکنش می‌تواند در زمان خطی با استفاده از یک مقایسه دو طرفه انجام شود. بنابراین، هزینه کلی ادغام تمام تراکنش‌ها در یک پایگاه داده واکنشی‌شده برابر است با $O(nw)$.

2-5- بسط مسئله‌ی استخراج مجموعه داده‌های پرسود

توضیحاتی که تاکنون داده شد مربوط به مسئله‌ی پایه‌ی HUIM بوده است که در آن مقدار سود اقلام تنها مثبت است و فرض می‌شود پایگاه داده حاوی آیتم‌هایی است که همیشه وجود دارند. اما در زندگی واقعی همیشه شرایط مسئله به این سادگی نیست. در ادامه به بسط‌های متفاوتی از مسئله‌ی HUIM خواهیم پرداخت.

۱-۵-۲- مسئله k مجموعه داده‌ی پرسود

تعیین مقدار آستانه سودمندی مناسب یک امر بسیار دقیق، حساس و زمان‌بر است به این معنی که اگر مقدار تعیین شده کوچک باشد تعداد بسیار زیادی مجموعه داده به عنوان اقلام پرسود استخراج می‌شوند. اگر این مقدار بزرگ باشد شاید هیچ مجموعه داده‌ی پرسودی استخراج نشود. بنابراین تعیین این مقدار توسط کاربر نیاز به تخصص لازم دارد. پس بر آن شدند که به جای تعیین مقدار آستانه سودمندی، کاربر تنها تعداد مجموعه داده‌های پرسودی که مد نظرش است را تعیین کند. برای حل این مسئله یک آستانه سودمندی داخلی تعریف می‌شود که آن را در ابتدا با صفر یا یک مقداردهی می‌کنند. سپس با استفاده از راهبردهای مناسب، این مقدار را افزایش می‌دهند تا k مجموعه داده‌ی مورد نظر استخراج شود. مشخص است که این مسئله از مسئله‌ی پایه پیچیده‌تر و زمان‌برتر است. مقالاتی که در آن‌ها به حل مسئله‌ی k مجموعه داده‌ی پرسود پرداخته شده است، [۱۵]، [۱۷] و [۱۸] هستند.

راهبردهای مختلفی برای افزایش مقدار آستانه سودمندی داخلی وجود دارد. در ادامه به برخی از این راهبردها اشاره خواهیم کرد. برای مثال از چهار تعریف زیر به عنوان راهبردی برای افزایش آستانه سودمندی در مقاله‌ی [۱۷] استفاده شده است.

تعریف ۸: راهبرد^۱ RIU مقدار سود واقعی تک آیتم‌ها را محاسبه می‌کند و مقدار آستانه سودمندی داخلی را از صفر یا یک به k امین عدد محاسبه شده که بیشترین سود را داراست، تغییر می‌دهد.

در بسیاری از راهبردهای معرفی شده از ساختار ماتریس مثلثی استفاده می‌شود. این ساختارها بسته به این که چه مقداری را در خود ذخیره می‌کنند، متفاوتند. برای مثال راهبردهای CUD^۳ و COV^۴ که در [۱۷] استفاده شده‌است، RIRU^۴ و RIRU2^۵ که در [۱۵] استفاده شده‌اند از ساختار ماتریس مثلثی استفاده می‌کنند.

تعریف ۹: راهبرد CUD از یک ساختار ماتریس مثلثی به نام CUDM^۶ استفاده می‌کند که در آن مقدار سود واقعی جفت آیتم‌هایی که از ترکیب هر آیتم با اقلام بعد از خودش به وجود می‌آید، ذخیره می‌شود. مقدار آستانه سودمندی داخلی به k امین عدد محاسبه شده که بیشترین مقدار را در این ماتریس داراست، تغییر می‌کند.

تعریف ۱۰: با فرض دو آیتم x و y اگر $g(X) \subseteq g(Y)$ باشد، آنگاه y پوششی برای x خواهد بود.

تعریف ۱۱: در راهبرد COV مقدار سود ترکیب هر آیتمی مانند i با اعضای مجموعه‌ی پوشش آن محاسبه می‌شود. مقدار آستانه سودمندی داخلی به k امین عدد محاسبه شده که بیشترین مقدار را داراست، تغییر می‌کند.

^۱ Rael Itemset Utility

^۲ Co-occurrence Utility Descending order

^۳ Coverage

^۴ Real 1-Itemset Relative Utility

^۵ Real 2-Itemset Relative Utility

^۶ Co-occurrence Utility Descending order Matrix

۲-۵-۲- مسئله مجموعه داده‌های پرسود موجود در قفسه

این مسئله از آن جایی مورد توجه قرار گرفت که در واقعیت اقلام همیشه موجود نیستند. برای مثال برخی اقلام تنها در تابستان هستند و یا برخی از اقلام تنها در زمان عید نوروز موجود هستند. در نتیجه پرداختن به این مسئله معیار عادلانه‌تری برای یافتن آیتم‌های پرسود است. برای حل این مسئله، دوره‌ی حضور هر آیتم در پایگاه داده، از قبل مشخص شده است. این دوره‌ها را با عدد مشخص می‌کنند. تمام تعاریفی که در مطالب قبلی برای محاسبه‌ی سود بیان شد، برای هر دوره زمانی مورد بررسی قرار می‌گیرد و مجموعه داده‌های سودمند بنا بر دوره‌ی زمانی که موجود هستند کشف می‌شوند. در [۱۵] به این مسئله نیز پرداخته شده است.

۲-۵-۳- مسئله مجموعه داده‌های پرسود با سود مثبت و منفی

در واقعیت ممکن است در پایگاه داده‌ی تراکنشی اقلام با سود منفی نیز وجود داشته باشند. این بدان معنی است که این اقلام زیر قیمت اصلی به فروش می‌رسند تا فروش اقلام دیگر بیشتر شود. در این حالت مسئله‌ی HIUM باید با در نظر گرفتن هر دو نوع سود اجرا شود. در مقالات [۱۰]، [۱۱] و [۱۵] به بررسی این مسئله پرداخته شده است.

در ادامه روشی را معرفی می‌کنیم که هر سه مسئله‌ی بالا را مورد بررسی قرار داده است. الگوریتم KOSHU [۱۵] به حل مسئله‌ی استخراج مجموعه داده‌ی پرسود می‌پردازد در حالی که هر دو نوع سود مثبت و منفی اقلام را در نظر می‌گیرد، سود اقلام را با توجه به دوره زمانی که موجود هستند محاسبه می‌کند و برای این مسئله به جای آستانه سودمندی از تعداد K آیتم پرسود که توسط کاربر تعیین می‌شود بهره می‌برد. در ادامه به تعاریف مورد نیاز برای حل این مسئله می‌پردازیم.

❖ تعاریف مربوط به اقلام پرسود دوره‌ای

فرض کنید PE مجموعه‌ای است که شامل دوره‌هایی است که اقلام در آن‌ها به فروش می‌رسند و هر کدام با یک عدد نمایش داده می‌شوند. به هر تراکنش نیز یکی از اعداد مجموعه‌ی PE اختصاص داده می‌شود، $pt(T_d) \in PE$ ، که نشان می‌دهد هر تراکنش در چه دوره‌ی زمانی‌ای رخ داده‌است.

تعریف ۱۲: دوره‌ی یک مجموعه داده مانند X، مجموعه‌ای از دوره‌هایی است که X در آن‌ها به فروش رسیده‌اند و به صورت $pi(X) = \{pt(T_d) | T_d \in D \wedge X \subseteq T_d\}$ تعریف می‌شود.

تعریف ۱۳: سود یک مجموعه داده مانند X در یک دوره‌ی زمانی خاص مانند h از جمع سود آن مجموعه داده در تراکنش‌هایی به دست می‌آید که در دوره زمانی h رخ دادند.

$$u(X, h) = \sum_{T_d \in D \wedge h = pt(T_d)} u(X, T_d) \quad (۷-۲)$$

سود یک مجموعه داده مانند X در پایگاه داده از رابطه (۴) به دست می‌آید.

$$u(X) = \sum_{h \in pi(X)} u(X, h) \quad (۸-۲)$$

تعریف ۱۴: سود یک تراکنش مانند T_d در پایگاه داده از جمع سود آیتم‌های موجود در آن و به روش $TU(T_d) = \sum_{i \in T_d} u(i, T_d)$ به دست می‌آید.

سود کلی یک دوره‌ی زمانی برای مجموعه داده‌ی X از جمع سود تراکنش‌هایی به دست می‌آید که دوره‌ی زمانی آن‌ها برابر با دوره‌های زمانی است که X در آن‌ها حضور داشته‌است.

$$to(X) = \sum_{h \in pi(X) \wedge T_d \in D \wedge h = pt(T_d)} TU(T_d) \quad (۹-۲)$$

تعریف ۱۵: سود نسبی یک مجموعه داده مانند X برابر با $ru(X) = u(X) / |to(X)|$ است.

سودمندی نسبی یک مجموعه داده مانند X نشان می‌دهد که چگونه سود و زیان تولید شده توسط X با سود و زیان کلی تولید شده در طول دوره‌های زمانی که X فروخته شد، مقایسه می‌شود.

مجموعه داده‌ی X یک HOU است اگر منفعت نسبی آن ($ru(X)$) کم‌تر از حداقل آستانه سودمندی مشخص شده توسط کاربر نباشد.

همان‌طور که قبلاً گفته شد معیار TWU معرفی شده برای پایگاه داده‌هایی که اقلام را تنها با سود مثبت در نظر می‌گیرند مناسب است. بنابراین برای پایگاه داده‌هایی با سود مثبت و منفی تعریف جدیدی ارائه شده است.

تعریف ۱۶: سود تراکنشی وزنی برای مجموعه داده‌ای مانند X در دوره زمانی h از جمع سود تراکنش‌هایی بدست می‌آید که دوره زمانی آن‌ها h است و X در آن‌ها حضور دارد.

$$TWU(X, h) = \sum_{X \subseteq T_d \wedge T_d \in D \wedge h = pt(T_d)} TU(T_d) \quad (۱۰-۲)$$

تعریف ۱۷: سود یک دوره‌ی زمانی مانند h از جمع سود تراکنش‌هایی که دوره‌ی زمانی آن‌ها برابر با h است به دست می‌آید.

$$pto(h) = \sum_{T_d \in D \wedge h = pt(T_d)} TU(T_d) \quad (۱۱-۲)$$

سود نسبی برای مجموعه داده‌ی X در دوره زمانی h برابر با $ru(X, h) = u(X) / |pto(h)|$ است.

ویژگی ۲: با داشتن مجموعه داده X و دوره زمانی h همواره $TWU(X, h) \geq u(X, h)$ است.

ویژگی ۳: اگر $X \subset Y$ آنگاه $TWU(X, h) \geq TWU(Y, h)$ خواهد بود.

¹ High On-Shelf Utility

ویژگی ۴: همواره $TWU(X, h) / pto(h) \geq ru(X, h)$ است. یعنی عبارت سمت چپ نامساوی یک

حد بالا برای سود نسبی مجموعه داده‌ی X در دوره زمانی h است.

ویژگی ۵: با داشتن مجموعه داده‌ی X ، اگر برای هیچ دوره‌ی زمانی عبارت

$TWU(X, h) / pto(h) \geq \min util$ صادق نباشد آنگاه X یک HOU نخواهد بود. در غیر این صورت

ممکن است HOU باشد یا نباشد.

• تعاریف مربوط به سود مثبت و منفی اقلام

تعریف ۱۸: سود تراکنش بازتعریف شده (RTU) باز تعریفی برای اقلام منفی است که از جمع سود

آیتم‌هایی با سود مثبت به دست می‌آید.

$$RTU(T_d) = \sum_{x \in T_d \wedge p(x) > 0} u(x, T_d) \quad (۱۲-۲)$$

سود وزنی تراکنش بازتعریف شده (RTWU)، باز تعریفی از TWU است و به صورت زیر به دست

می‌آید.

$$RTWU(X, h) = \sum_{T_d \in D \wedge X \subseteq T_d \wedge pt(T_d) = h} RTU(T_d) \quad (۱۳-۲)$$

فرض کنید $up(X)$ مجموعه‌ی آیتم‌هایی از X باشد که سود مثبت دارند و $un(X)$ مجموعه‌ای از

آیتم‌هایی با سود منفی باشند که در X وجود دارند.

ویژگی ۶: با داشتن مجموعه داده‌ی X ، خواهیم داشت که $u(X, h) \leq u(up(X), h)$.

¹ Redefined Transaction Utility

² Redefined Transaction Weighted Utility

ویژگی ۷: فرض کنیم که X یک مجموعه داده و z یک آیتم با سود منفی باشد که عضو X نیست، در این صورت $u(up(X \cup \{z\}), h) \leq u(up(X), h)$ خواهد بود.

ویژگی ۸: فرض کنید Y بسطی منفی از X باشد، آنگاه $u(up(Y), h) \leq u(up(X), h)$ خواهد بود.

ویژگی ۹: اگر بر اساس ترتیب \succ تنها اقلام منفی بتوانند با X ترکیب شوند و هیچ دوره‌ای وجود نداشته باشد که $u(up(X), h) / pto(h) \geq \min util$ ، آنگاه X و هر بسطی از آن HOU نخواهد بود.

❖ راهبردهای افزایش آستانه سودمندی

راهبردهای معرفی شده مربوط به الگوریتم KOSHU هستند که در [۱۸] آمده است.

۱- راهبرد افزایش آستانه سودمندی با استفاده از مقدار سودمندی نسبی تک‌آیتم‌ها

اولین روش، راهبرد افزایش آستانه سودمندی با استفاده از مقدار سودمندی نسبی تک‌آیتم‌ها (RIRU) نام دارد. در این روش که بعد از اولین اسکن پایگاه داده اجرا می‌شود، مقدار سود نسبی تک‌آیتم‌ها حساب می‌شود. سپس بر اساس این راهبرد، مقدار آستانه سودمندی که برابر صفر بود به مقدار k امین سود نسبی محاسبه شده تغییر می‌کند.

۲- راهبرد افزایش آستانه سودمندی با استفاده از مقدار سودمندی نسبی آیتم‌های دوتایی

دومین روش، راهبرد افزایش آستانه سودمندی با استفاده از مقدار سودمندی نسبی تک‌آیتم‌ها (RIRU2) نام دارد. در این روش که بعد از اسکن دوم انجام می‌شود، سود نسبی جفت آیتم‌ها بر اساس تعریف ۵ حساب می‌شود. سپس مقدار آستانه سودمندی از مقداری که با روش RIRU تعیین شده بود به مقدار k امین سود نسبی بالای محاسبه شده در این مرحله افزایش می‌یابد.

¹ Real 1-Itemset Relative Utility

² Real 2-Itemset Relative Utility

۴-۵-۲- مسئله مجموعه داده های متناوب پرسود

در مسئله ی HUIM پایه مجموعه داده های پرسودی که استخراج می شوند به پایگاه داده ای تعلق دارند که برای مدت طولانی اطلاعات در آن جمع آوری شده است. برای مثال برای یک فروشگاه اینترنتی مثل دیجی کالا، جواهرات می توانند به عنوان خروجی مسئله ی HUIM پایه استخراج شوند. این درحالی است که خرید جواهرات به ندرت اتفاق می افتد. اقلام دیگری که به طور متناوب خریداری می شوند می توانند به عنوان اقلام پرسود در دوره تناوب خودشان محسوب شوند. برای مثال از بین اقلامی که هفتگی خریداری می شوند، لبنیات می تواند جز اقلام پرسود این دوره باشد. این مسئله می تواند رفتار مشتریان را بهتر تحلیل کند. برای حل این مسئله نیاز است علاوه بر مقدار آستانه سودمندی، مقدار کمینه و بیشینه ی دوره تناوب و همچنین مقدار کمینه و بیشینه ی میانگین تناوب توسط کاربر تعیین شود. با این شرایط برای هرس فضای جست و جو باید این مقادیر نیز مورد توجه قرار گیرند. مقالاتی که به این موضوع پرداخته اند، [۱۳] و [۱۹] هستند. در ادامه، تعاریف مربوط به این مسئله آمده است. برای یکپارچگی در تعریف مسئله تعاریف ارائه شده در این بخش مربوط به [۱۳] است.

❖ الگوهای متناوب تکراری

تعریف ۱۹: (تناوب یک مجموعه داده). فرض کنید X یک مجموعه داده در پایگاه داده ای است که شامل n تراکنش است. $g(X) = \{T_{g_1}, T_{g_2}, \dots, T_{g_k}\}$ مجموعه ی تراکنش هایی است که X در آن ها حضور دارد و رابطه ی $1 \leq g_1 < g_2 < \dots < g_k \leq n$ برقرار است. مقدار تناوب برای دو تراکنش متوالی با شماره شناسایی x و y در رابطه با X به صورت $pe(T_x, T_y) = (y - x)$ به دست می آید که برابر با تعداد تراکنش هایی است که بین این دو تراکنش وجود دارد. تناوب یک مجموعه داده مانند X ، لیستی از مقدار تناوب برای هر تراکنشی است که در مجموعه ی $g(X)$ وجود دارند. و به صورت $ps(X) = \cup_{1 \leq z \leq k+1} (g_z - g_{z-1})$ به دست می آید.

تعریف ۲۰: (الگوهای متناوب تکراری). بیشینه‌ی تناوب یک مجموعه داده‌ای مانند X به صورت $\max_{per}(X) = \max(ps(X))$ تعریف می‌شود. یک مجموعه داده زمانی PFP محسوب می‌شود که $\max_{per}(X) < \max Per$ و $|g(X)| \geq \min sup$ باشد. مقادیر $\max Per$ و $\min sup$ توسط کاربر تعریف می‌شوند.

یکی از مهم‌ترین محدودیت‌ها در الگوریتم‌های کاوش الگوهای تکراری متناوب این است که تنها به تعداد دفعات تکرار توجه می‌کنند و سود حاصل از آن‌ها (PFP) مهم نیست.

❖ الگوهای پرسود متناوب

در این بخش معیارهایی که برای پیدا کردن مجموعه داده‌های پرسود متناوب استفاده می‌شود، معرفی خواهد شد.

معیار بیشینه‌ی تناوب در بسیاری از الگوریتم‌های PFP استفاده شده است. با این معیار اگر مجموعه داده‌ای، مجموعه تناوب‌اش شامل یک تناوب واحد باشد که از مقدار $\max Per$ بیشتر باشد کل آن مجموعه داده نادیده گرفته می‌شود. برای رفع این مشکل از معیار میانگین تناوب استفاده می‌شود.

تعریف ۲۱: (میانگین تناوب برای یک مجموعه داده). این معیار از رابطه (۱۰) به دست می‌آید.

$$avg_{per}(X) = \sum_{g_i \in ps(X)} g_i / |ps(X)| \quad (۱۴-۲)$$

ویژگی ۱۰: (رابطه‌ی بین میانگین تناوب و مقدار پشتیبانی). فرض کنید X مجموعه داده‌ای در پایگاه

داده‌ی D باشد. راه دیگری برای محاسبه میانگین تناوب به صورت (۱۱) است.

$$avg_{per}(X) = |D| / (|g(X)| + 1) \quad (۱۵-۲)$$

¹ Periodic Frequent Pattern

این روش به این دلیل مهم است که اندازه D تنها یکبار محاسبه می‌شود و میانگین تناوب هر مجموعه داده تنها با محاسبه $|g(X)+1|$ به دست می‌آید. علاوه بر این، این قیاس مهم است چون نشان می‌دهد که یک رابطه بین پشتیبانی مورد استفاده در FIM و تناوب متوسط یک الگو وجود دارد. اگر چه میانگین تناوب مفید است، نباید به عنوان تنها معیار برای ارزیابی تناوب‌پذیری یک الگو استفاده شود زیرا در نظر نمی‌گیرد که آیا یک مجموعه دارای دوره‌های متفاوت با طول زیاد است یا خیر. بنابراین این مجموعه داده نباید یک مجموعه‌ی متناوب باشد. برای اجتناب از یافتن الگوهایی با دوره‌های متفاوت، راه حل ترکیب اندازه تناوب متوسط با دیگر معیارها را درپیش می‌گیرند. یکی از معیارهایی که استفاده می‌شود معیار کمینه‌ی تناوب است که به صورت $\min_{\text{per}}(X) = \min(\text{ps}(X))$ تعریف می‌شود. این معیار برای جلوگیری از مجموعه داده‌هایی که تناوب کوتاهی دارند به کار گرفته می‌شود. اما تناوب‌هایی که برابر با ۰ و ۱ باشند مثل آن‌هایی که شامل تراکنش اول و آخر باشند را در نظر نمی‌گیرد. بنابراین راه حل این است که تناوب‌های اول و آخر از لیست تناوب هر مجموعه داده را در نظر نگیریم. اگر با این کار لیست تناوب خالی باشد آن را ∞ در نظر می‌گیریم. منطق استفاده از این مقیاس در ترکیب با تناوب متوسط این است که می‌تواند از کشف الگوهای دوره‌ای که برای دوره‌های طولانی رخ نمی‌دهند جلوگیری کند. دلیل استفاده از این سه معیار این است که از نظر محاسبه و مصرف حافظه کم هزینه هستند.

تعریف ۲۲: (مجموعه داده‌های پرسود متناوب). با فرض این که مقادیر مثبت \min_{Avg} , \min_{util} , \max_{Avg} , \min_{Per} و \max_{Per} را کاربر تعیین کند، مجموعه داده X پرسود است اگر

$$\max_{\text{per}}(X) \leq \max_{\text{Per}} \text{ و } \min_{\text{per}}(X) \geq \min_{\text{Per}} \text{ و } \min_{\text{Avg}} \leq \text{avgper}(X) \leq \max_{\text{Avg}}$$

در نهایت $u(X) \geq \min_{\text{util}}$ باشد.

با توجه به تعاریف ارائه شده در بالا می‌توان مجموعه داده‌های پرسود متناوب را استخراج نمود.

6-2- جمع‌بندی

در این فصل به شرح مفصل مسئله‌ی HUIM و بسط‌های آن پرداختیم و نحوه‌ی محاسبه‌ی سود و یافتن الگوی مناسب را بیان کردیم.

در نهایت لازم به ذکر است که برای مسئله‌ی HUIM پایه بسط‌های مختلف دیگری ارائه شده‌است که هر یک برای حل چالشی مناسب هستند. برای مثال مسئله‌ی HUIM را می‌توان با در نظر گرفتن داده‌های درجریان مورد بررسی قرار داد. مسئله‌ی مجموعه‌داده‌ی پرسود بسته مجموعه‌داده‌ی را استخراج می‌کند که پرسود است و از روی آن می‌توان مجموعه‌داده‌های پرسود دیگر را استخراج کرد. این روش باعث می‌شود تا زمان استخراج مجموعه‌داده‌های پرسود کمتر شود زیرا تعداد الگوهای بسته کمتر است و یا حتی می‌توان HUIM را با توجه به استراتژی‌های تخفیف مختلف حل کرد.

فصل سوم

نتیجه‌گیری و جمع‌بندی مطالب

3-1- مقدمه

حل مسئله‌ی HUIM نسبت به مسئله‌ی FIM سخت‌تر و چالش‌برانگیزتر است. زیرا در مسئله‌ی FIM از ویژگی غیریکنواخت بودن پشتیبانی استفاده می‌کنند و فضای جست‌وجو را به‌طور موثر هرس می‌کنند. اما معیار سود در HUIM نه یکنواخت است و نه غیریکنواخت. به عبارتی سود مجموعه‌داده‌ای می‌تواند کمتر، مساوی و یا بیشتر از ابرمجموعه‌اش باشد در حالی که مقدار پشتیبانی یک مجموعه‌داده همواره از زیرمجموعه‌اش کمتر و یا با آن برابر است. همین امر سبب می‌شود تا برای هرس فضای جست‌وجو نیازمند به راهبردهای مختلفی باشیم تا کاهش به صورت موثری انجام شود. بنابراین پیچیدگی در حل این مسئله سبب می‌شود تا زمان اجرا و مصرف حافظه از مسئله‌ی FIM بیشتر باشد. روش‌های ارائه‌شده برای انواع مختلف مسئله‌ی HUIM از راهبردهای مختلف هرس، ساختمان داده‌های متفاوت و روش‌های مختلفی برای جست‌وجوی فضای حالت استفاده می‌کنند که همه‌ی این‌ها سبب می‌شود تا میزان مصرف حافظه و زمان اجرای این روش‌ها با هم متفاوت باشد.

3-2- علت تفاوت روش‌ها

همان‌طور که قبلاً اشاره شد روش‌ها ممکن است به دلیل نوع نمایش پایگاه داده، تعداد مراحل انجام کار، روش جست‌وجو و راهبردهای هرس مورد استفاده در زمان اجرا و میزان حافظه‌ی مصرفی با هم متفاوت باشند. دلایل تفاوت روش‌ها در شکل (۳-۱) نشان داده می‌شود.



شکل ۳-۱: معیار تفاوت روش‌های حل

اگر از پایگاه داده‌ی تراکنشی و یا پایگاه داده‌ی واکنشی‌شده استفاده کنیم، نمایش پایگاه داده به صورت افقی خواهد بود ولی اگر از لیست سودمندی استفاده کنیم از پایگاه داده‌ی عمودی استفاده کرده‌ایم. استفاده از روش جست‌وجوی عمقی نیز از روش جست‌وجوی سطحی بهتر است. زیرا تنها مجموعه داده‌ای را بسط می‌دهد که در پایگاه داده موجود است. بنابراین میزان حافظه کمتری استفاده می‌کند. اما در جست‌وجوی سطحی مجموعه داده‌ها را بدون در نظر گرفتن این که در پایگاه داده موجود است یا نه بسط می‌دهد. این امر سبب می‌شود تا فضای حالت بزرگتری داشته باشیم که عملاً بسیاری از داده‌های مورد بررسی قرار نخواهند گرفت. در زیر الگوریتم‌های پایه‌ی مسئله‌ی HUIM مورد بررسی قرار گرفته‌اند.

جدول ۳-۱: دسته‌بندی روش‌های پایه

الگوریتم	سال	روش جست‌وجو	تعداد مراحل	نمایش پایگاه داده	الگوریتم پایه
Tow phase [۴]	۲۰۰۵	سطحی	دو فاز	افقی	Apriori
PB [۲۰]	۲۰۱۴	سطحی	دو فاز	افقی	Apriori
IHUP [۲۱]	۲۰۰۹	عمق اول	دو فاز	افقی (درخت پیشوندی)	FPGrowth
UPGrowth [۲۲]	۲۰۱۳	عمق اول	دو فاز	افقی (درخت پیشوندی)	FPGrowth
HUPGrowth [۲۳]	۲۰۱۱	عمق اول	دو فاز	افقی (درخت پیشوندی)	FPGrowth
HUIMiner [۶]	۲۰۱۲	عمق اول	تک فاز	عمودی (لیست سودمندی)	Eclat
FHM [۲۴]	۲۰۱۴	عمق اول	تک فاز	عمودی (لیست سودمندی)	Eclat
HUIMiner* [۷]	۲۰۱۸	عمق اول	تک فاز	عمودی (لیست سودمندی*)	Eclat
ULBMiner [۸]	۲۰۱۷	عمق اول	تک فاز	عمودی (لیست سودمندی)	Eclat
EFIM [۵]	۲۰۱۵	عمق اول	تک فاز	افقی (پایگاه داده واکنشی شده)	LCM

3-3- مقایسه روش‌های تک‌فاز و دوفاز

از آنجایی که روش‌های دوفاز از راه حل‌های اولیه به‌شمار می‌روند در نتیجه کاستی‌هایی نیز داشته-

اند. در زیر ویژگی‌های این دو روش را برخواهیم شمرد.

۱-۳-۳- ویژگی روش‌های دوفاز

❖ مزایا

- سادگی در فهم و پیاده‌سازی

❖ معایب

- مصرف زیاد حافظه به دلیل زیاد بودن تعداد کاندیدهای تولیدشده
- بالا بودن تعداد اسکن پایگاه داده برای محاسبه سود هر مجموعه داده
- استفاده از راهبردهای نه‌چندان دقیق برای هرس فضای جست‌وجو
- بالا بودن زمان اجرا به دلیل اسکن‌های پی‌درپی پایگاه داده

۲-۳-۳- ویژگی روش‌های تک‌فاز

❖ مزایا

- استفاده از ساختمان داده‌های جدید مانند لیست سودمندی
- کاهش دفعات اسکن پایگاه داده به دلیل استفاده از لیست سودمندی و یا پایگاه داده‌ی واکنشی‌شده
- کاهش میزان مصرف حافظه نسبت به روش‌های دوفاز به دلیل عدم تولید مجموعه‌ی کاندید
- کارآمدتر بودن از نظر زمان اجرا نسبت به روش‌های دوفاز به دلیل عدم اسکن پی‌درپی پایگاه داده
- کران‌های بالای دقیق‌تر

❖ معایب

- بالا بودن هزینه ادغام در روش‌های مبتنی بر لیست سودمندی
- پیچیدگی در پیاده‌سازی

3-4- جمع‌بندی

در این بخش قصد داریم تا راه‌کارهای مورد استفاده را براساس روش‌های مورد استفاده‌ی آنها، دسته‌بندی نماییم. در جدول (۲-۳) به‌طور خلاصه راه‌کارهای ارائه‌شده به همراه نوع مسئله‌ی تحت پوشش و تحلیل مختصری از آن آورده شده‌است.

جدول ۳-۲: خلاصه‌ای از راه کارهای ارائه شده

الگوریتم	سال انتشار	سود در نظر گرفته شده	نوع الگوریتم	ساختار داده‌ی مورد استفاده	راهنماهای هرس	تحلیل
HMiner_closed [۹]	۲۰۱۹	مثبت	تک فاز	لیست سودمندی فشرده	C-prune و LA-prune	کاهش زمان و حافظه با استفاده از فشرده سازی لیست و همچنین پرداختن به مسئله ی مجموعه داده های پرسود بسته برای کاهش تعداد HUIها
MHAUIPNU [۱۰]	۲۰۱۹	مثبت و منفی	تک فاز	لیست سودمندی	Tubpn، بر اساس آیت منفی و بر اساس بسط منفی	ارائه ی راهنماهای هرس دقیق برای کاهش فضای جست و جو
ULB-Miner [۸]	۲۰۱۷	مثبت	تک فاز	لیست سودمندی بافر شده	EUCP	بهبود مصرف حافظه و زمان با استفاده از لیست بافر شده و بهبود عملکرد الگوریتم های قبلی با استفاده از این روش
ABSA [۱۱]	۲۰۱۹	مثبت و منفی	تک فاز	لیست سودمندی	ساختار EUCS	شاخه ای از متن کاوی که با استفاده از تحلیل نظر کاربران ویژگی های برتر یک محصول را استخراج می کند. هر ویژگی از قبل ارزش گذاری می شود.
FUP-HAUIDM [۱۲]	۲۰۱۸	مثبت	تک فاز	لیست سودمندی AU-List	auub	محاسبه ی HAU در پایگاه داده ی پویا با در نظر گرفتن حذف تراکنش ها. به دلیل پویا بودن نتوانست کران بالای دقیق تری برای کاهش فضا استفاده کند.
PHM [۱۳]	۲۰۱۶	مثبت	تک فاز	لیست سودمندی	avgPer، maxPer و EUCP	تحلیل رفتار مشتری بر اساس دوره های خرید او، بهره گیری از یک روش کارآمد مانند FHM
PEFIM [۱۶]	۲۰۱۹	مثبت	تک فاز	آرایه سودمندی	su و Lc	واکشی پایگاه داده، ادغام تراکنش ها و کم کردن هزینه ادغام و حافظه در لیست سودمندی و موازی سازی برای کاهش زمان اجرا

MHUI [۱۴]	۲۰۱۷	مثبت	نقض	لیست سودمندی	U-M, TWU-M LA-M و EUCS-M	استفاده از چند آستانه سودمندی برای کاوش مجموعه داده‌های سودمند
KOSHU [۱۵]	۲۰۱۶	مثبت و منفی	نقض	لیست سودمندی	CE2P, EMRP, PUP	پرداختن به دو مسئله‌ی دوره‌ای بودن اقلام و k مجموعه‌ی برتر
TKEH [۱۷]	۲۰۱۸	مثبت	نقض	آرایه سودمندی و صف اولویت	sup, Tm	پرداختن به مسئله k مجموعه‌ی برتر و استفاده مفید از روش‌های قبلی و ارائه کران- های بالای دقیق
SPHUI _{TP} [۱۹]	۲۰۱۷	مثبت	نقض	لیست	MP و TWU	پرداختن به مسئله اقلام دوره‌ی زمانی کوتاه اما پرهزینه در مصرف زمان و حافظه
EFIM [۵]	۲۰۱۶	مثبت	نقض	آرایه سودمندی	su و Lc	واکشی پایگاه داده، ادغام تراکنش‌ها و کم کردن هزینه ادغام و حافظه در لیست سودمندی
HUI-Miner [۶]	۲۰۱۲	مثبت	نقض	لیست سودمندی	TWU و Ru	بسیار سریع بودن در مقابل الگوریتم‌های دوفاز
HUI-Miner* [۷]	۲۰۱۹	مثبت	نقض	لیست سودمندی	TWU و Ru	کم کردن زمان مقایسه تراکنش‌های مشترک هنگام ادغام دو لیست و موثر برای داده‌های پراکنده
LUIM [۲۵]	۲۰۱۹	مثبت	نقض	آرایه	NUL و TWU	پرداختن به مسئله‌ی مجموعه داده‌های کم- سود- بالا بودن مصرف حافظه و زمان

3-5- نتیجه‌گیری

به‌طور کلی روش‌هایی که تک‌فاز هستند از نظر سرعت و حافظه بهتر از روش‌های دوفاز عمل می‌کنند. اما درهمین راستا نیز اگر روشی از لیست سودمندی استفاده کند هزینه ادغام لیست‌ها خودیک چالش است که هنوز جای کار دارد. برای بهبود این کاستی از پایگاه‌های داده‌ی واکشی‌شده، لیست‌های سودمندی بافرشده و یا لیست سودمندی فشرده‌شده استفاده شده‌است. برای کاهش فضای حالت نیز استفاده از راهبردهای دقیق‌تر پیشنهاد می‌شود که از محاسبه سود دقیق اقلام بهره‌مند شوند.

فهرست منابع و مراجع

- [۱] تقوی، پروین؛ فاطمه، اکبرپور. ۱۳۹۵. طراحی و پیاده سازی سیستم تحلیل و بررسی میزان کتابخوانی در ایران. پایان نامه کارشناسی، دانشگاه تربیت دبیر شهید رجایی.
- [2] Fournier-Viger, Philippe, et al. "A survey of high utility itemset mining." *High-Utility Pattern Mining*. Springer, Cham, 2019. 1-45.
- [۳] مدیر آموزش فرادرس. ۱۳۹۲. فیلم آموزشی جامع کاوش قواعد وابستگی. متلب سایت.
<https://matlabsite.com/933/mvrdm9206ij-association-rule-mining-in-data-mining.html> (دسترسی در ۱۳۹۹/۳/۱۵)
- [4] Liu, Ying, Wei-keng Liao, and Alok Choudhary. "A two-phase algorithm for fast discovery of high utility itemsets." *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, Berlin, Heidelberg, 2005.
- [5] Zida, S., Fournier-Viger, P., Lin, J.C.-W., Wu, C.W., Tseng, V.S.: EFIM: a highly efficient algorithm for high-utility itemset mining. In: *Proceedings of the 14th Mexican International Conference Artificial Intelligence*, pp. 530–546. Springer (2015)
- [6] Liu, Mengchi, and Junfeng Qu. "Mining high utility itemsets without candidate generation." *Proceedings of the 21st ACM international conference on Information and knowledge management*. 2012.
- [7] Qu, Jun-Feng, Mengchi Liu, and Philippe Fournier-Viger. "Efficient algorithms for high utility itemset mining without candidate generation." *High-Utility Pattern Mining*. Springer, Cham, 2019. 131-160.
- [8] Duong, Quang-Huy, et al. "Efficient high utility itemset mining using buffered utility-lists." *Applied Intelligence* 48.7 (2018): 1859-1877.
- [9] Nguyen, Loan TT, et al. "An efficient method for mining high utility closed itemsets." *Information Sciences* 495 (2019): 78-99.

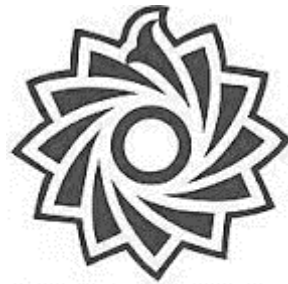
- [10] Yildirim, Irfan, and Mete Celik. "Mining High-Average Utility Itemsets with Positive and Negative External Utilities." *New Generation Computing* (2019): 1-34.
- [11] Demir, Seyfullah, et al. "Extracting Potentially High Profit Product Feature Groups by Using High Utility Pattern Mining and Aspect Based Sentiment Analysis." *High-Utility Pattern Mining*. Springer, Cham, 2019. 233-260.
- [12] Lin, Jerry Chun-Wei, et al. "Maintenance algorithm for high average-utility itemsets with transaction deletion." *Applied Intelligence* 48.10 (2018): 3691-3706.
- [13] Fournier-Viger, Philippe, et al. "PHM: mining periodic high-utility itemsets." *Industrial conference on data mining*. Springer, Cham, 2016.
- [14] Krishnamoorthy, Srikumar. "Efficient mining of high utility itemsets with multiple minimum utility thresholds." *Engineering Applications of Artificial Intelligence* 69 (2018): 112-126.
- [15] Dam, Thu-Lan, et al. "An efficient algorithm for mining top-k on-shelf high utility itemsets." *Knowledge and Information Systems* 52.3 (2017): 621-655.
- [16] Nguyen, Trinh DD, Loan TT Nguyen, and Bay Vo. "A Parallel Algorithm for Mining High Utility Itemsets." *International Conference on Information Systems Architecture and Technology*. Springer, Cham, 2018.
- [17] Singh, Kuldeep, et al. "TKEH: an efficient algorithm for mining top-k high utility itemsets." *Applied Intelligence* 49.3 (2019): 1078-1097.
- [18] Krishnamoorthy, Srikumar. "Mining top-k high utility itemsets with effective threshold raising strategies." *Expert Systems with Applications* 117 (2019): 148-165.
- [19] Lin, Jerry Chun-Wei, et al. "A two-phase approach to mine short-period high-utility itemsets in transactional databases." *Advanced Engineering Informatics* 33 (2017): 29-43.
- [20] Lan, Guo-Cheng, Tzung-Pei Hong, and Vincent S. Tseng. "An efficient projection-based indexing approach for mining high utility itemsets." *Knowledge and information systems* 38.1 (2014): 85-107.

- [21] Ahmed, Chowdhury Farhan, et al. "Efficient tree structures for high utility pattern mining in incremental databases." *IEEE Transactions on Knowledge and Data Engineering* 21.12 (2009): 1708-1721.
- [22] Tseng, Vincent S., et al. "Efficient algorithms for mining high utility itemsets from transactional databases." *IEEE transactions on knowledge and data engineering* 25.8 (2012): 1772-1786.
- [23] Lin, Chun-Wei, Tzung-Pei Hong, and Wen-Hsiang Lu. "An effective tree structure for mining high utility itemsets." *Expert Systems with Applications* 38.6 (2011): 7419-7424.
- [24] Fournier-Viger, Philippe, et al. "FHM: Faster high-utility itemset mining using estimated utility co-occurrence pruning." *International symposium on methodologies for intelligent systems*. Springer, Cham, 2014.
- [25] Alhusaini, Naji, et al. "LUIM: New Low-Utility Itemset Mining Framework." *IEEE Access* 7 (2019): 100535-100551.
- [26] Krishnamoorthy, Srikumar. "A comparative study of top-k high utility itemset mining methods." *High-Utility Pattern Mining*. Springer, Cham, 2019. 47-74.

Abstract

With the explosive growth of data and information, the need for strategies and tools to convert data into useful knowledge is increased. In the most of the past literatures, although they find meaningful information, but they do not care about the utility of the items. So having information about the high utility itemsets can be useful to make right decisions. To satisfy this need, in this literature we investigate some solutions to mine the itemsets that have utilities more than a utility threshold.

Keywords: Data mining, high utility itemsets, frequent itemset, minimum utility



Shahid Rajaee Teacher
Training University

Shahid Rajaee Teacher Training University
Faculty of Computer Engineering
Department of Software

Title:
High Utility Itemset Mining

By:
Parvin Taghavi

Supervisor:
Dr. Negin Daneshpour

June 2020