

## کاوش الگوهای متناوب از درخت FP با کمک ماتریس دوبعدی

بابک توسلی<sup>۱</sup>؛ مهرداد جلالی<sup>۲</sup>؛ مجید وفایی جهان<sup>۳</sup>

### چکیده

کشف الگوهای پنهان و ارزشمند از درون حجم وسیعی از داده‌های خام، اخیراً توجه بسیاری از محققان را به خود جلب کرده‌است. اولین گام برای استخراج روابط بین ارقام موجود در داده‌ها، مشخص نمودن تعداد رخداد‌های همزمان ارقام با یکدیگر است. الگوریتم‌های متعددی برای نیل به این هدف تاکنون ارائه شده‌اند اما می‌توان از نظر نحوه عملکرد آن‌ها را در دو دسته کلی قرار داد. در دسته نخست کاندیدهایی از مجموعه قلم‌داده تولید می‌شوند و سپس با پویش تراکنش‌ها، تکرار هر کدام به دست می‌آید و آنهایی که تکرارشان بیشتر از حد آستانه پشتیبانی کمینه تعیین شده باشد به عنوان الگوی متناوب معرفی می‌شوند. مشکل اصلی این روش هزینه بالای تولید کاندیدها و پویش‌های متعدد تراکنش‌هاست. در دسته دوم الگوریتم‌های رشد الگو قرار می‌گیرند که برای استخراج الگوهای متناوب نیازی به تولید کاندید ندارند. کارایی این الگوریتم‌ها در پایگاه‌های داده بزرگ، جایی که الگوریتم‌های دسته اول به خاطر تعداد زیاد کاندیدها بسیار ناکارآمد عمل می‌کنند، بسیار محسوس است.

در الگوریتم‌های رشد الگو، از ساختار داده‌ای به نام درخت الگوی متناوب استفاده می‌شود. با بررسی روش رشد الگوی متناوب مشخص گردیده‌است که حدود ۸۰٪ از زمان اجرای الگوریتم صرف پیمایش درخت الگوی متناوب می‌شود. از آنجاییکه پیمایش لیست پیوندی فرآیندی زمانبر است لذا در روش پیشنهادی از ساختار داده آرایه برای نگهداری درخت الگوی متناوب استفاده شده‌است تا سرعت پیمایش افزایش یابد. همچنین برای عملیات کاوش نیز دیگر نیازی به ساختن درخت‌های شرطی نیست و بر خلاف الگوریتم‌های فعلی که در هر بار پیمایش شاخه‌های درخت، الگوی جاری به طول  $n$  را تا  $n+1$  گسترش می‌دهند، در روش پیشنهادی با کمک ماتریس دوبعدی، الگوی متناوب جاری به طول  $n$  تا حداکثر  $n+2$  گسترش می‌یابد. نتایج تجربی نشان می‌دهد که تغییر ساختار نگهداری درخت و روش کاوش پیشنهادی در پایگاه‌های تراکنش متوسط تا بزرگ باعث بهبود در سرعت اجرای الگوریتم به خصوص در آستانه‌های کوچکتر می‌شود.

### کلمات کلیدی

داده کاوی، مجموعه قلم‌داده، رشد الگو، درخت الگوی متناوب، ماتریس دوبعدی

# کنفرانس داده کاوی ایران

<sup>۱</sup> دانشجوی کارشناسی ارشد کامپیوتر- نرم افزار ، دانشگاه آزاد اسلامی واحد مشهد، babak.tavassoli@gmail.com

<sup>۲</sup> استادیار گروه کامپیوتر دانشکده فنی مهندسی دانشگاه آزاد اسلامی واحد مشهد، mehrjalali@gmail.com

<sup>۳</sup> استادیار گروه کامپیوتر دانشکده فنی مهندسی دانشگاه آزاد اسلامی واحد مشهد، vafaeija@yahoo.com

# Mining Frequent Patterns from FP-Tree

## By Using 2D Matrix

Babak Tavassoli; Mehrdad Jalali; Majid Vafaei Jahan

### Abstract

Discovering hidden pattern from large volume of raw data has recently attracted the attention of many researchers. The first step to extract relations among data items is specifying the number of simultaneous occurrence of items together. Several algorithms have been proposed to achieve this goal so far but they can be placed in two major categories based on their functionality. In the first category candidates generated from data items and their number of occurrence can be calculated after the first scan of database and those with frequency equal or greater than the minimum support are marked as frequent pattern. The main bottleneck of this method is high cost of generating candidates and multiple scans of database. In the second category there are pattern growth algorithms that for mining frequent patterns are not required to generate candidates. Performance of these algorithms in large databases where the first category algorithms are very inefficient because of large number of candidates is very noticeable. Pattern growth algorithms use a data structure called frequent pattern tree or FP-Tree. From numerous experiments, we found out that about 80 percent of the CPU time was used for traversing FP-trees. Since linked list traversal is a time consuming process in the proposed algorithm we have used array structure for storing transactions in order to increase the speed. Constructing conditional fp-trees is no longer needed in the proposed method and unlike current pattern-growth algorithms which extend a frequent pattern with length  $n$  to  $n+1$  in each step of mining, in the proposed method each frequent pattern with length  $n$  can be extended to maximum  $n+2$ . Experimental results show that changing the structure for constructing tree and new method of mining frequent patterns can improve speed of the algorithm for medium to large data bases especially in the smaller thresholds.

### Keywords:

Data Mining; Itemset; Pattern Growth; Frequent Pattern Tree; 2D Matrix

### ۱. مقدمه

الگوهای متناوب<sup>۱</sup> مجموعه قلم‌داده‌ها<sup>۲</sup>، توالی‌ها<sup>۳</sup> و یا ساختارهایی<sup>۴</sup> هستند که در یک مجموعه داده به تعداد دفعاتی بیشتر یا مساوی با آنچه توسط کاربر به عنوان حد آستانه تعیین شده است، تکرار شده باشند. به عنوان مثال مجموعه‌ای از چند قلم‌داده (مانند شیر و نان) که معمولاً با هم در تراکنش‌های موجود در پایگاه تراکنش یک فروشگاه به کار می‌روند، یک مجموعه قلم‌داده متناوب<sup>۵</sup> هستند. در یک توالی علاوه بر قلم‌داده‌ها، ترتیب زمانی به کار رفتن و یا روی دادن آنها در یک تراکنش مهم است. به عنوان مثال خرید یک کامپیوتر و پس از آن میز کامپیوتر و سپس سی‌دی‌های بازی، یک توالی در پایگاه تراکنش فروشگاه محسوب می‌شود. یک ساختار می‌تواند اشکال مختلف ساختارها از قبیل گراف، درخت و شبکه را در برگیرد. اگر ساختاری در یک پایگاه گرافیکی بیش از حد مشخصی تکرار شود آن را الگوی ساختاری متناوب<sup>۶</sup> می‌نامند. یافتن الگوی متناوب، نقشی کلیدی در کاوش قوانین وابستگی، همبستگی‌ها و بسیاری از روابط جالب دیگر در بین داده‌ها ایفا می‌نماید. علاوه بر آن در اندیس نمودن داده‌ها، طبقه بندی، خوشه بندی و سایر عملیات داده کاوی نیز کارایی دارد.

### ۱.۱ قوانین وابستگی

با در نظر گرفتن یک مجموعه از قلم‌داده‌ها و یک پایگاه تراکنش متشکل از این مجموعه، با کمک قوانین وابستگی می‌توان وابستگی میان قلم‌داده‌ها را پیدانمود. این وابستگی‌ها بیان می‌کنند که حضور هر قلم‌داده و یا مجموعه قلم‌داده‌ها در یک تراکنش موجب حضور برخی دیگر از قلم‌داده‌ها و یا مجموعه قلم‌داده‌ها می‌شود. مثال متداول در رابطه با کشف قوانین وابستگی، «تحلیل سبد خرید» است [۱]. در این فرآیند با توجه به اقلام مختلفی که مشتریان در سبد خریدشان قرار می‌دهند، عادات و رفتار خرید آن‌ها مورد تحلیل قرار می‌گیرد. هدف در این فرآیند پیدا کردن خودکار قوانینی مثل «۶۰٪ افرادی که نان خریداری می‌کنند شیر هم می‌خرند» می‌باشد که البته برای قابل قبول بودن قوانین معیارهایی نیز وجود

دارد. به طور کلی هدف اصلی در کاوش قوانین وابستگی پیدا کردن تمامی قوانینی است که آستانه پشتیبان کمینه<sup>۷</sup> و ضریب اطمینان کمینه<sup>۸</sup> مدنظر کاربر را رعایت کرده باشند. مسئله کشف قوانین وابستگی را می‌توان به دو زیر مساله تقسیم کرد [۲]:

۱. پیدا کردن تمام مجموعه قلم داده‌هایی که حداقل آستانه پشتیبان را داشته باشند. به مجموعه قلم داده‌هایی که حداقل آستانه پشتیبان را داشته باشند مجموعه قلم داده‌های متناوب گفته می‌شود و مجموعه قلم داده‌هایی که این حداقل آستانه را نداشته باشند، مجموعه قلم داده‌های غیرمتناوب نامیده می‌شوند.  
۲. استفاده از مجموعه قلم داده‌های متناوب جهت تولید قوانین وابستگی.

از آنجایی که زیر مساله اول؛ یعنی کاوش مجموعه قلم داده‌های متناوب مساله زمانبرتری است؛ لذا در سال‌های اخیر بررسی‌ها و پژوهش‌های بیشتری در این زمینه صورت گرفته است.

## ۲،۱ کاوش قوانین وابستگی<sup>۹</sup>

کاوش قوانین وابستگی یکی از مهمترین تکنیک‌های داده‌کاوی است که هدف اصلی آن استخراج ارتباطات پنهان و وابستگی‌های جالب بین قلم داده‌های موجود در یک پایگاه‌تراکنش است. عملیات یافتن قوانین وابستگی را کاوش قوانین وابستگی گویند. این عملیات شامل دو مرحله اساسی است: مرحله اول کاوش مجموعه قلم داده‌های متناوب و مرحله دوم تولید قوانین وابستگی با استفاده از مجموعه قلم داده‌های متناوب. بار اصلی این عملیات بر دوش مرحله اول یعنی پیدا کردن مجموعه قلم داده‌های متناوب است. کاوش مجموعه قلم داده‌های متناوب نخستین بار در سال ۱۹۹۳ توسط آگروال<sup>۱۰</sup> برای تحلیل سبد خرید و در قالب کاوش قوانین وابستگی ارائه شد [۳]. نتیجه چنین تحلیل‌هایی نشان می‌داد مشتریانی که کالای خاصی را می‌خرند به همراه این کالا به خرید چه نوع کالاهای دیگری رغبت نشان می‌دهند. بر اساس این نتایج صاحبان فروشگاه‌ها می‌توانند نحوه چینش کالاها در قفسه‌ها را بر اساس نیاز مشتریان تنظیم نمایند.

## ۳،۱ روش های یافتن مجموعه قلم داده‌های متناوب

تاکنون روش‌های متعددی برای یافتن مجموعه قلم داده‌های متناوب پیشنهاد شده‌اند که به دو دسته کلی طبقه بندی میشوند؛ روش‌های مبتنی بر تولید کاندید<sup>۱۱</sup> و روش‌های مبتنی بر رشد الگو<sup>۱۲</sup>. برای نخستین بار آگروال در سال ۱۹۹۴ خاصیت جالبی به نام Apriori را در مجموعه آیتم‌های متناوب مشاهده نمود [۴]. بر اساس این خاصیت، «یک مجموعه  $k$ - قلم داده‌ای فقط وقتی متناوب است که همه زیر مجموعه قلم داده‌های آن متناوب باشند». نتیجه اصلی این خاصیت، «متناوب نبودن ابرالگوهای یک الگوی نامتناوب» است که به الگوریتم کاوش امکان کنار گذاشتن مجموعه‌های  $k$ - قلم داده‌ای نامتناوب در کاوش مجموعه‌های  $(k+1)$ - قلم داده‌ای را خواهد داد. بر این اساس برای یافتن مجموعه کامل همه مجموعه قلم داده‌های متناوب موجود در یک پایگاه تراکنش، باید ابتدا با یکبار پویش<sup>۱۳</sup> پایگاه تراکنش، همه مجموعه‌های  $1$ - قلم داده‌ای متناوب (قلم داده‌های متناوب) موجود در آن پیدا شده و سپس تمامی مجموعه‌های متناوب  $2$ - قلم داده‌ای قابل ساخت از این قلم داده‌های متناوب به عنوان کاندیدای متناوب بودن تولید شوند. هر یک از این مجموعه‌های  $2$ - قلم داده‌ای، چون تنها از قلم داده‌های متناوب تشکیل شده‌اند (همه زیرمجموعه قلم داده‌های آنها متناوب هستند)، به صورت بالقوه امکان متناوب بودن در پایگاه تراکنش را دارند. برای اطمینان از این تناوب، لازم است بار دیگر پایگاه تراکنش را پویش نموده و تناوب یا عدم تناوب هر یک از کاندیداهای  $2$ - قلم داده‌ای مشخص شود. پس از تعیین مجموعه‌های متناوب  $2$ - قلم داده‌ای، به صورت مشابه از آنها برای ساخت کاندیداهای  $3$ - قلم داده‌ای و بررسی تناوب یا عدم تناوب آنها با پویش مجدد پایگاه تراکنش استفاده می‌شود. روند استفاده از مجموعه‌های متناوب  $k$ - قلم داده‌ای در تولید کاندیداهای  $(k+1)$ - قلم داده‌ای و بررسی تناوب یا عدم تناوب آنها با یک پویش کامل پایگاه تراکنش، تا استخراج همه مجموعه قلم داده‌های متناوب موجود در پایگاه تراکنش ادامه می‌یابد. اگرچه در بسیاری از مسائل و نمونه‌ها، روش‌های مبتنی بر تولید کاندید با استفاده از اصل «عدم تناوب ابرالگوهای یک نامتناوب»، سبب کاهش قابل ملاحظه‌ای در اندازه مجموعه کاندیدهای تولیدشده برای حل مساله می‌گردند اما همواره دو اشکال اصلی بر این روش‌ها وارد است. نخستین اشکال تولید تعداد فوق العاده زیادی کاندیدا در قالب مجموعه‌های کاندید است که با افزایش تعداد قلم داده‌ها بسیار محسوس میشود. دومین مشکل اساسی پویش مکرر و چندین باره پایگاه تراکنش برای مطابقت کاندیداهای تولیدشده با تراکنش‌ها به منظور تشخیص تناوب یا عدم تناوب آنهاست.

هان<sup>۱۴</sup> در سال ۲۰۰۰ روشی به نام رشد الگوی متناوب<sup>۱۵</sup> را ابداع نمود که مجموعه کامل همه مجموعه قلم داده‌های متناوب را بدون تولید کاندیداها کاوش می‌نمود. این روش مبتنی بر تقسیم وحل بوده و با یک بار پویش پایگاه تراکنش، همه قلم داده‌های متناوب موجود در پایگاه تراکنش را یافته و بر اساس تعداد تکرارشان در پایگاه تراکنش (میزان تناوب)، به صورت نزولی مرتب می‌کند. بر مبنای لیست مرتب قلم داده‌های متناوب، پایگاه تراکنش در قالب یک درخت الگوی متناوب<sup>۱۶</sup> فشرده‌سازی می‌شود که این درخت شامل همه اطلاعات موجود در پایگاه تراکنش است و می‌تواند در فرآیند کاوش مورد استفاده قرارگیرد. کاوش درخت الگوی متناوب با شروع از هر قلم داده متناوب (به عنوان الگوی پسوندی اولیه)

آغاز می‌گردد و پایگاه الگوهای شرطی آن قلم‌داده‌ها ساخته می‌شود. پایگاه الگوهای شرطی یک الگو شامل مجموعه پیشوندی مسیریابی از درخت الگوی متناوب است که الگوی مورد نظریه عنوان پسوند در آن مسیرها به کار رفته باشد. پس از این مرحله، درخت الگوی متناوب شرطی الگوی مزبور ساخته می‌شود و عملیات کاوش به صورت بازگشتی بر روی این درخت اعمال می‌گردد. رشد الگوها با اتصال الگوهای پسوندی به الگوهای متناوب کاوش شده از درخت‌های الگوی متناوب شرطی، ساخته می‌شوند. الگوریتم رشد الگوی متناوب مساله یافتن الگوهای متناوب بزرگ را به مساله جستجوی الگوهای متناوب کوچکتر به صورت بازگشتی تبدیل می‌نماید و سپس پسوندها را به الگوهای یافته شده متصل می‌کند. این روش با انتخاب درست قلم‌داده‌ها می‌تواند کمترین قلم‌داده‌ها را به عنوان پسوند در نظر بگیرد. مطالعات انجام شده در ارتباط با کارایی این الگوریتم نشان داده‌است که روش مزبور، زمان جستجو را به شکل قابل ملاحظه‌ای کاهش می‌دهد [۵].

## ۴٫۱ کارهای مرتبط

تاکنون پیشنهادهای زیادی برای بهبود روش رشد الگوهای متناوب ارائه شده‌است. در روش رشد الگوی متناوب یکبار در ابتدای پویش پایگاه تراکنش درخت الگوی متناوب ساخته شده و در ضمن انجام فرآیند کاوش، به دفعات، درخت‌های شرطی به منظور رشد الگو تولید و پیمایش می‌شوند. یکی از راه کارهای سریعتر کردن روش رشد الگوی متناوب به نام رشد الگوی متناوب استار<sup>۱۷</sup> در سال ۲۰۰۵ توسط گراهنه و ژو<sup>۱۸</sup> ارائه شده است [۶]. در این روش نشان داده شده است که بیش از ۸۰٪ زمان پردازنده صرف ساخت و پیمایش درخت‌های الگوی متناوب می‌شود. تولید هر درخت شرطی جدید نیازمند دو بار پویش درخت جاری است. بار اول تعداد تکرارهای هر قلم‌داده به دست می‌آید و در بار دوم تراکنش‌های دربردارنده الگوی جاری از شاخه‌های درخت خوانده شده و پس از اصلاح و حذف قلم‌داده‌های کم‌تکرار، در درخت شرطی جدید قرار می‌گیرند. با بکارگیری ساختار داده ماتریس می‌توان با یکبار پیمایش درخت جاری، درخت شرطی جدید را ساخت. در این روش در هنگام تولید درخت شرطی جدید دیگر نیازی به پیمایش درخت جاری برای یافتن تعداد تکرار قلم‌داده‌ها نیست؛ بلکه با مراجعه به سطر مربوطه می‌توان جدول سرآیند<sup>۱۹</sup> درخت شرطی جدید را بدست آورد و تنها با یکبار پیمایش درخت جاری به درخت شرطی رسید. بررسی نتایج و نمودارهای مربوط به روش استار نشان می‌دهد که زمان اجرای الگوریتم کاهش چشمگیری پیدا می‌کند. با این وجود الگوریتم همچنان برای نگهداری تراکنش‌ها از ساختار لیست پیوندی استفاده میکند که پیمایش آن زمانبر بوده و نیز برای تولید الگوهای متناوب در هر مرحله از عملیات کاوش باید درخت‌های شرطی را تولید و پیمایش نماید که با افزایش عمق درخت اصلی باعث کاهش کارایی الگوریتم میشود.

در سال ۲۰۱۱ کی-چانگ<sup>۲۰</sup> و همکاران پیشنهاد دیگری برای بهبود روش رشد الگوی متناوب ارائه کردند [۷]. بر اساس تحقیقات آن‌ها عملیات درج تراکنش‌ها به منظور ساختن درخت الگوی متناوب بخشی از زمان اجرای الگوریتم را به خود اختصاص می‌دهد. اگرچه این زمان زیاد نیست اما مشکلی که باعث این اتلاف وقت می‌شود در مراحل کاوش درخت، خود را نشان می‌دهد. مشکل بدین صورت است که نگهداری فرزندان هر گره در یک لیست باعث می‌شود تا هنگام پیدا کردن شاخه‌های درخت، زمان محسوسی به خاطر جستجو در این لیست از دست داده شود. کی-چانگ برای حل این مشکل ساختار جدیدی به نام جدول آدرس<sup>۲۱</sup> را پیشنهاد کرده‌است. این ساختار آدرس هر گره فرزند را در یک جدول درهم نگهداری می‌کند و در هنگام مراجعه به هر گره به راحتی می‌توان فرزندان را پیدا کرد. برای حل مشکل ساختن متوالی درخت‌های شرطی، گره کی-چانگ روش جدیدی برای کاوش درخت پیشنهاد کرده‌اند. این روش بیان می‌کند که برای قلم‌داده‌های پرتکرارتر بهتر است درخت‌های شرطی بر روی درخت اصلی تشکیل شوند. با این کار علاوه بر صرفه‌جویی در حافظه، سرعت عملیات کاوش نیز بهبود می‌یابد. پارامتر سطح درخت<sup>۲۲</sup>، قلم‌داده‌هایی را مشخص می‌کند که درخت‌های شرطی آنها باید بر روی درخت اصلی ساخته شود. با توجه به آزمایشات فراوانی که بر روی پایگاه‌های تراکنش انجام گرفته است، مقدار ۵۰٪ به عنوان مقدار پایه‌ای برای این پارامتر ارائه شده‌است؛ یعنی اینکه اگر برای قلم‌داده‌های نیمه بالایی جدول سرآیند درخت‌های شرطی را بر روی درخت اصلی و برای قلم‌داده‌های نیمه پایینی جدول سرآیند همانند روش رشد الگوی متناوب عمل کنیم به بهترین زمان‌های اجرا دست می‌یابیم. اگرچه در این روش میزان نیاز به تولید درخت‌های شرطی تا ۵۰٪ کاهش می‌یابد اما مشکل درخت‌های شرطی همچنان در آن وجود دارد. از طرف دیگر تمامی پایگاه‌های تراکنش با مقدار ۵۰٪ برای پارامتر سطح درخت سازگار نیستند و لازم است تا برای مقادری این پارامتر از روش‌های بهتری استفاده شود. مشکل مشترک در هر دو روش ذکر شده این است که همچنان برای عملیات کاوش به درخت‌های شرطی نیاز دارند؛ هرچند که در روش کی-چانگ این نیاز بسیار کم‌تر است. در الگوریتم پیشنهادی ما سعی بر این بوده‌است تا با تغییر ساختار نگهداری درخت از لیست پیوندی به ترکیبی از جداول درهم و آرایه زمان پیمایش شاخه‌ها را کاهش داده و با تغییر روش کاوش از پیمایش درخت دربردارنده تراکنش‌ها و تولید درختان شرطی به یک روش کاملاً مبتنی بر ماتریس دوبعدی، ضمن افزایش سرعت عملیات کاوش، تعداد بیشتری الگوی متناوب در هر مرحله کاوش تولید نماییم.



## ۲. روش پیشنهادی

این بخش به بررسی قسمت‌های مختلف الگوریتم پیشنهادی بر روی یک پایگاه تراکنش نمونه می‌پردازد. در شکل‌های (۱) و (۲) کلیات الگوریتم پیشنهادی آورده شده است.

### ۱.۲ فرضیات

در کاوش قوانین وابستگی رخ دادن مجموعه قلم‌داده‌ها مورد توجه است و تعداد دفعات تکرار هر قلم‌داده تأثیری در ارتباط بین آن‌ها ندارد. با توجه به این فرض پایگاه تراکنش مانند  $D$  را بر روی مجموعه قلم‌داده  $L = \{۱, ۲, ۳, ۴, ۵, ۶, ۷, ۸, ۹, ۱۰, ۱۱, ۱۲, ۱۳, ۱۴, ۱۵, ۱۶, ۱۸\}$  که در جدول (۱) آمده است را در نظر بگیرید. این پایگاه‌داده شامل تراکنش‌هایی است که هر کدام از آن‌ها شامل چند قلم‌داده می‌باشند.

جدول (۱) پایگاه تراکنش  $D$

تراکنش	مجموعه قلم‌داده‌ها
۱	۱۶, ۳, ۴, ۶, ۹, ۷, ۱۳, ۱
۲	۶, ۲, ۱۳, ۳, ۱, ۱۲, ۱۵
۳	۲, ۶, ۱۵, ۱۰, ۸
۴	۳, ۱۱, ۲, ۱۸, ۱۶
۵	۱۴, ۳, ۵, ۱۶, ۱۲, ۱۳, ۱۶

همانطور که مشاهده می‌شود برای راحتی از اعداد به عنوان نماد قلم‌داده‌ها استفاده شده است. با یکبار پیمایش پایگاه تراکنش تعداد تکرار هر قلم‌داده مشخص می‌شود. در گام بعدی و با در نظر داشتن حد آستانه پشتیبان کمینه<sup>۲۳</sup>، قلم‌داده‌های کم تکرار هر تراکنش حذف شده و سایر قلم‌داده‌ها بر اساس میزان تکرار به طور نزولی مرتب می‌شوند. در الگوریتم پیشنهادی پایگاه تراکنش دو بار پویش می‌شود. در پویش نخست قلم‌داده‌های متناوب<sup>۲۴</sup> شناسایی می‌شوند و در پویش دوم با حذف قلم‌داده‌های کم تکرار، با در نظر گرفتن حد آستانه پشتیبان کمینه، تراکنش‌ها اصلاح می‌شوند و در شاخه‌های درخت قرار می‌گیرند. حد آستانه پشتیبان کمینه برابر با سه،  $\delta = 3$ ، فرض شده است. در جدول (۲) قلم‌داده‌های متناوب بر اساس میزان تکرار به صورت نزولی نشان داده شده‌اند. در صورت تکرار مساوی برای چند قلم‌داده، ارزش عددی آن‌ها به صورت نزولی، مد نظر قرار گرفته است.

جدول (۲) قلم‌داده‌های متناوب

قلم داده	۳	۶	۱	۲	۱۳	۱۶
تکرار	۴	۴	۳	۳	۳	۳

### ۲,۲ فاز اول: پیش پردازش

همانطور که در شکل (۱) مشاهده می‌شود پایگاه تراکنش‌ها ابتدا یکبار مورد پویش قرار می‌گیرد و با توجه به حد آستانه قلم‌داده‌های متناوب پیدامی‌شوند. سپس در پویش بعدی قلم‌داده‌های غیرمتناوب از تراکنش‌ها حذف شده و قلم‌داده‌های متناوب در تراکنش‌ها بر اساس میزان تکرارشان به طور نزولی قرار می‌گیرند. در اینجا مرحله پیش پردازش به اتمام رسیده و الگوریتم وارد فاز کاوش می‌شود. در جدول (۳) پایگاه تراکنش بعد از اتمام مرحله پیش پردازش نشان داده شده است.

جدول (۳) پایگاه تراکنش  $D$  بعد از اصلاح تراکنش‌ها

تراکنش	مجموعه قلم‌داده‌ها
۱	۳, ۶, ۱, ۱۳, ۱۶
۲	۳, ۱, ۱۶
۳	۳, ۶, ۱, ۲, ۱۳
۴	۶, ۲
۵	۳, ۶, ۱, ۱۳, ۱۶

برای مواردی که تعداد تکرار قلم‌داده‌ها یکسان است، اولویت با عدد کمتر است. برای مثال هر دو قلم‌داده ۳ و ۶ به تعداد یکسان تکرار شده‌اند لذا اگر هر دو در یک تراکنش با هم بیایند، ابتدا قلم‌داده ۳ و سپس ۶ قرار می‌گیرد.

### ۱,۲,۲ حد آستانه

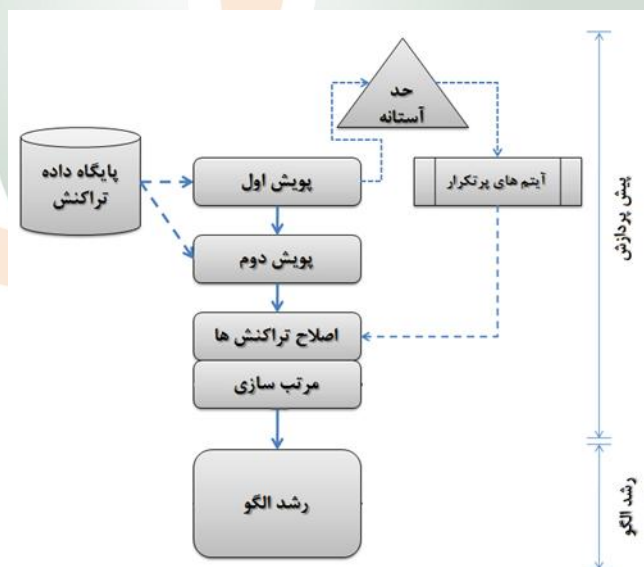
حد آستانه یک مقدار عددی است که گاه به صورت مطلق و گاه به صورت درصد توسط کاربر تعیین می‌شود. این عدد تعداد تکرار لازم برای متناوب بودن یک قلم‌داده را مشخص می‌کند و با علامت  $\delta$  نمایش داده می‌شود.

### ۳.۲ فاز دوم: رشد الگو

این مرحله از الگوریتم به سه بخش کلی تقسیم می‌شود که در ادامه مورد بررسی قرار می‌گیرند

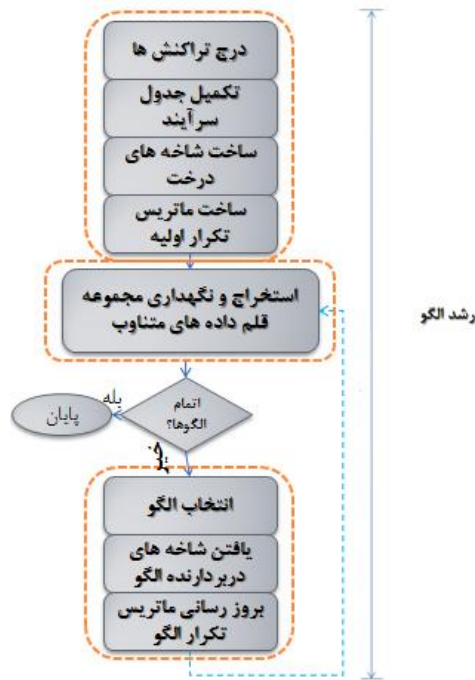
#### ۱,۳,۲ آماده سازی

دو هدف اصلی این بخش به ترتیب عبارتند از خواندن تراکنش‌های اصلاح‌شده و قراردادن آن‌ها در شاخه‌های درخت و تشکیل ماتریس تکرار اولیه. هنگام خواندن هر تراکنش با مراجعه به جدول راهنما وجود شاخه‌ای با قلم‌داده آغازین مشابه به قلم‌داده آغازین تراکنش جاری مورد بررسی قرار می‌گیرد. در صورتیکه شاخه‌ای وجود داشته باشد تراکنش جاری در آن قرار گرفته و در صورت عدم وجود، شاخه جدیدی ایجاد شده و تراکنش به‌طور کامل در آن درج می‌شود.



شکل (۱) شماتیک الگوریتم پیشنهادی - پیش پردازش

بعد از اتمام درج با توجه به قلم‌داده‌های بکار رفته در تراکنش، جدول راهنما بروز می‌شود و موقعیت قلم‌داده‌ها در بخش موقعیت جدول قرار می‌گیرند. با اتمام خواندن هر تراکنش ماتریس تکرار اولیه بروز می‌شود. تعداد سطرها و ستون‌های این ماتریس به تعداد قلم‌داده‌های متناوب درون پایگاه تراکنش است. اگر قلم‌داده‌های متناوب بر اساس میزان تکرارشان به صورت نزولی مرتب شده و به صورت  $\{F_1, F_2, \dots, F_n\}$  نشان داده شوند، برای  $n = 5$ ، ماتریس تکرار اولیه ساختاری شبیه جدول (۴) خواهد داشت.



شکل (۲) شماتیک الگوریتم پیشنهادی - کاوش

مقادیر درون خانه‌های ماتریس تکرار که به صورت  $V_{ij}$  نشان داده شده‌است، بیانگر میزان تکرار قلم‌داده‌های  $F_i$  و  $F_j$  با یکدیگر است. جدول (۵) ساختار شاخه‌های درخت را برای شاخه‌ای با  $n=5$  قلم‌داده متناوب نشان می‌دهد. الگوریتم برای دسترسی و پیمایش شاخه‌های درخت از ساختاری به نام جدول راهنما استفاده می‌کند. در حقیقت این جدول بهبود یافته جدول سرآیند در روش رشد الگوی متناوب است که با کمک آن دیگر نیازی به استفاده از اشاره‌گرهای درونی برای یافتن موقعیت قلم‌داده‌ها در درخت نیست. جدول ۶ نمونه‌ای از این ساختار داده را نشان می‌دهد.

جدول (۴) ساختار ماتریس تکرار

$F_1$	$V_{11}$	-	-	-	-
$F_2$	$V_{12}$	$V_{22}$	-	-	-
$F_3$	$V_{13}$	$V_{23}$	$V_{33}$	-	-
$F_4$	$V_{14}$	$V_{24}$	$V_{34}$	$V_{44}$	-
$F_5$	$V_{15}$	$V_{25}$	$V_{35}$	$V_{45}$	$V_{55}$
	$F_1$	$F_2$	$F_3$	$F_4$	$F_5$

جدول (۵) ساختار شاخه درخت

قلم داده	$F_1$	$F_2$	...	$F_n$
تکرار	$V_1$	$V_2$	...	$V_n$

جدول (۶) جدول راهنما

موقعیت‌ها شروع قلم داده

$F_1$	$S_1$	$\{R_{1i}, \dots, R_{1k}\}$
$F_2$	$S_2$	$\{R_{2i}, \dots, R_{2p}\}$
...	...	...
$F_n$	$S_n$	$\{R_{ni}, \dots, R_{nq}\}$

ستون سمت چپ قلم‌داده‌های متناوب را نشان می‌دهد، ستون میانی ( $S_i$ ) شاخه‌ای را نشان می‌دهد که با  $F_i$  شروع شده است و ستون سمت راست ( $R_i$ ) دنباله‌ای از اعداد را شامل می‌شود که بیانگر شاخه‌هایی از درخت است که قلم‌داده  $F_i$  در آنها وجود دارد. برای روشن‌تر شدن مطلب تراکنش نخست از پایگاه تراکنش  $D$  را در نظر بگیرید. این تراکنش با قلم‌داده ۳ آغاز می‌شود. از آنجاییکه در ابتدا جدول راهنما خالی است لذا با مراجعه به آن شاخه‌ای برای قلم‌داده ۳ پیدا نمی‌شود و در نتیجه شاخه‌ای جدید ایجاد شده و قلم‌داده‌های این تراکنش در آن قرار می‌گیرند. نتیجه این عملیات در شکل (۳) نشان داده شده است.

قلم داده	موقعیت‌ها	شروع	قلم داده	3	6	1	13	16
3	1	1	تکرار	1	1	1	1	1
6	-	1						
1	-	1						
2	-	-						
13	-	1						
16	-	1						

شاخه‌های درخت

جدول راهنما

$$T_1 = \{3, 6, 1, 13, 16\}$$

شکل (۳) وضعیت شاخه‌های درخت بعد از ورود اولین تراکنش

در ادامه تراکنش دوم از پایگاه تراکنش خوانده می‌شود. این تراکنش نیز با قلم‌داده ۳ آغاز می‌شود. این بار با مراجعه به جدول راهنما، شاخه شماره یک به عنوان شاخه مقصد بازگردانده می‌شود. مقدار عددی خانه تکرار برای این قلم‌داده در شاخه شماره یک، یک واحد افزایش می‌یابد. قلم‌داده دوم در این تراکنش قلم‌داده ۲ می‌باشد و این درحالیست که قلم‌داده بعدی در شاخه شماره یک، قلم‌داده ۶ است؛ یعنی اینکه برای درج ادامه تراکنش باید شاخه دیگری انتخاب شود. در جدول راهنما در حال حاضر شاخه‌ای برای قلم‌داده ۲ ثبت نشده است و لذا برای ادامه درج باید شاخه جدیدی ساخته شود. نتیجه این عملیات در شکل (۴) نشان داده شده است.



موقعیت ها	شروع	قلم داده
1,2	1	3
1	-	6
1	-	1
2	-	2
1	-	13
1,2	-	16

جدول راهنما

قلم داده	3	6	1	13	16
تکرار	2	1	1	1	1

قلم داده	3	2	16
تکرار	۱	1	1

شاخه های درخت

$$T_2 = \{3, 2, 16\}$$

#### شکل (۴) وضعیت شاخه‌های درخت بعد از ورود دومین تراکنش

با ادامه این روند تمامی تراکنش‌ها در شاخه‌های مناسب قرار می‌گیرند. شکل (۵) نتیجه نهایی این عملیات را نشان می‌دهد. بعد از درج شدن هر تراکنش، ماتریس تکرار نیز بروز می‌شود. تراکنش نخست شامل قلم‌داده‌های ۳، ۶، ۱، ۱۳ و ۱۶ است. ستون مربوط به قلم‌داده ۳ را در نظر بگیرید. از آنجاییکه این قلم‌داده با قلم‌داده‌های ۶، ۱، ۱۳ و ۱۶ تکرار شده‌است لذا سطرهای مربوط به این قلم‌داده‌ها برای ستون قلم‌داده ۳ مقدار یک می‌گیرند. در صورتیکه این خانه‌ها از قبل مقدار داشته‌باشند آن مقدار یک واحد افزایش می‌یابد.

موقعیت ها	شروع	قلم داده
1,2,3	1	3
1,3,4	4	6
1,3	-	1
2,3,4	-	2
1,3	-	13
1,2	-	16

جدول راهنما

قلم داده	3	6	1	13	16
تکرار	4	3	3	2	2

قلم داده	3	2	16
تکرار	۱	1	1

قلم داده	3	6	1	2	13
تکرار	۱	۱	۱	1	1

قلم داده	6	2
تکرار	1	1

شاخه های درخت

#### شکل (۵) وضعیت نهایی بعد از ورود آخرین تراکنش

جدول (۷) وضعیت نهایی این ماتریس را بعد از ورود تمامی تراکنش‌ها نشان می‌دهد. قطر اصلی ماتریس در این مرحله تعداد تکرار قلم‌داده‌های متناوب پایگاه تراکنش را نشان می‌دهد. از آنجاییکه قلم‌داده‌ها در تراکنش‌ها بر اساس میزان تکرار مرتب شده‌اند لذا ماتریس تکرار متقارن می‌شود و از نیمه بالایی آن استفاده نمی‌شود؛ به همین علت مقدار صفر برای خانه‌های بالای قطر اصلی در نظر گرفته شده‌است.

# کنفرانس داده کاوی ایران

جدول (۷) ماتریس تکرار بعد از ورود تمام تراکنش‌ها

3	4	-	-	-	-	-
6	3	4	-	-	-	-
1	3	3	3	-	-	-
2	2	2	1	3	-	-
13	3	3	3	1	3	-
16	3	2	2	1	2	3
	3	6	1	2	13	16

### ۲,۳,۲ استخراج الگو

بعد از اتمام خواندن تراکنش‌ها و تکمیل ماتریس تکرار اولیه و یا بعد از مرحله بازسازی، الگوریتم وارد بخش استخراج الگو می‌شود. با مراجعه به ماتریس تکرار و در نظر گرفتن حد آستانه، مکان‌هایی از ماتریس را که مقدارشان بیشتر از حد آستانه است برای استخراج الگو در نظر گرفته می‌شوند. ماتریس تکرار جدول (۷) را در نظر بگیرید، اگر  $V_{ij} \geq \delta$  باشد و الگوی متناوب جاری  $P_k$  باشد، الگوهای جدید به صورت زیر تولید می‌شوند:

$$P_{k+1} = P_k \wedge F_i \quad (1)$$

$$P_{k+2} = P_k \wedge F_j \quad (2)$$

$$P_{k+3} = P_k \wedge F_i \wedge F_j \quad (3)$$

در فرمول‌های (۱)، (۲) و (۳) نماد  $\wedge$  بیانگر عملگر اتصال است و به صورت زیر تعریف می‌شود:

$$F_i \wedge F_j = \begin{cases} F_i F_j & i \neq j \\ F_i & i = j \end{cases} \quad (5)$$

باید توجه داشت که در اولین ورود به بخش کاوش، هیچ الگویی از قبل وجود ندارد و لذا داریم:  $P_k = \phi$

برای مثال با بررسی ماتریس تکرار در این مرحله تمامی الگوهای متناوب به طول ۲ قابل استخراج هستند. نتیجه این عملیات در شکل ۶ نشان داده شده است.

3	4	0	0	0	0	0
6	3	4	0	0	0	0
1	3	3	3	0	0	0
2	2	2	1	3	0	0
13	3	3	3	1	3	0
16	3	2	2	1	2	3
	3	6	1	2	13	16

{3,6}  
{3,1}  
{3,13}  
{3,16}  
{6,1}  
{6,13}  
{1,13}

الگوهای متناوب به طول دو

ماتریس تکرار

شکل (۶) استخراج الگوهای متناوب به طول دو

## ۳,۳,۲ بازسازی

در این مرحله یک الگوی متناوب از الگوهای استخراج شده در مرحله استخراج الگو انتخاب شده و رشد می‌یابد. فرآیند رشد بدین شکل است که ابتدا از جدول راهنما موقعیت‌های مربوط به قلم‌داده‌های درون الگوی انتخاب شده به‌دست آمده و سپس با اشتراک‌گیری از این موقعیت‌ها شاخه‌های در بردارنده الگوی انتخابی مشخص می‌شود. حال با مراجعه به این شاخه‌ها و خواندن قلم‌داده‌هایی که شامل می‌شوند الگوریتم مقادیر درون خانه‌های ماتریس تکرار را بروز می‌کند. در اینجا مرحله بازسازی به اتمام می‌رسد و الگوریتم دوباره وارد بخش استخراج الگو می‌شود؛ با این تفاوت که این بار الگوهای داری پیشوند الگوی انتخابی استخراج می‌شوند. در شکل (۷) و (۸) نتیجه عملیات بازسازی برای الگوی {۳,۶} نشان داده شده‌است.

1

قلم داده	3	6	1	13	16
تکرار	4	3	3	2	2

قلم داده	3	2	16
تکرار	۱	1	1

3

قلم داده	3	6	1	2	13
تکرار	۱	۱	۱	1	1

قلم داده	6	2
تکرار	1	1

شماره های درخت

شکل (۷) پیداکردن شاخه‌ها برای الگوی {۳,۶}

با بررسی ماتریس بازسازی شده سه الگوی جدید استخراج می‌شود؛ این الگوهای در سمت راست شکل ۳-۸ نشان داده شده‌اند. با ادامه این روند برای سایر الگوها تمامی مجموعه قلم‌داده‌های پایگاه تراکنش به‌دست می‌آیند.

3	0	0	0	0	0	0
6	0	0	0	0	0	0
1	0	0	3	0	0	0
2	0	0	1	1	0	0
13	0	0	3	1	3	0
16	0	0	2	1	2	2
	3	6	1	2	13	16

ماتریس تکرار

الگوهای متناوب حاصل از رشد الگوی ۳ و ۶

{3,6,1}  
{3,6,13}  
{3,6,1,13}

شکل (۸) ماتریس تکرار بازسازی شده برای الگوی {۳,۶}

## ۳. نتایج تجربی

در این بخش الگوریتم پیشنهادی به همراه چند الگوریتم دیگر در شرایط یکسان پیاده سازی و مورد آزمایش قرار گرفته‌اند تا به طور عملی کارایی هر یک از آنها مشخص شود. سعی بر این بوده است تا از کدهای استاندارد و موجود برای الگوریتم‌های ذکر شده استفاده شود. پیاده‌سازی الگوریتم‌ها با زبان سی‌شارپ نسخه ۴ و در محیط ویژوال استودیو ۲۰۱۰ انجام شده‌است. در جدول (۸) مشخصات سیستمی که آزمایشات روی آن انجام شده، آورده شده است.

#### جدول (۸) مشخصات سیستم مورد استفاده در آزمایشات

Processor	Intel Core 2 Quad 6600
Ram	2 GB
OS	MS Windows 7 SP1

#### ۱,۳ الگوریتم‌های مورد آزمایش

الگوریتم پیشنهادی به همراه الگوریتم‌های رشدالگوی متناوب، FP-Growth-Star و IFP بر روی دیتاست‌های مختلف اجرا شده‌اند تا تاثیر پیشنهادات داده شده در عمل مورد بررسی قرار بگیرد.

#### ۲,۳ دیتاست‌های استفاده شده در آزمایش

در جداول (۹) و (۱۰) اطلاعات کلیدی در مورد دیتاست‌های استاندارد که در آزمایش مورد استفاده قرار گرفته‌اند، به طور خلاصه آورده شده است.

جدول (۹) دیتاست T.I.D.K

T10I4D100K	نام دیتاست
100,000	تعداد تراکنش‌ها
870	تعداد آیتم‌ها
10.1	میانگین طول تراکنش

#### ۳,۳ ارزیابی کارایی

تمامی الگوریتم‌هایی که در بخش ۱,۳ معرفی شدند توسط زبان سی‌شارپ پیاده‌سازی کامپایل شده‌اند و از ساختارهای داده استاندارد ۴,۰ Framework استفاده شده‌است. بستر امن<sup>۲۵</sup> انجام شده و در هیچ یک از آنها از تکنیک‌های موازی سازی استفاده نشده است. زمان‌های اندازه گیری شده از اولین نمونه‌سازی اشیاء مربوط به الگوریتم‌های مختلف تا زمان اتمام پروسه کاوش<sup>۲۶</sup> هر شیئی در نظر گرفته شده است. کلاس‌های مربوط به کار با پایگاه‌تراکنش‌ها در تمامی روش‌ها یکسان پیاده سازی شده است و تا جای ممکن برای عملیات یکسان، از کدها و ساختارهای داده‌ای مشابه استفاده شده است تا نتایج ارزیابی‌ها دقیق باشند.

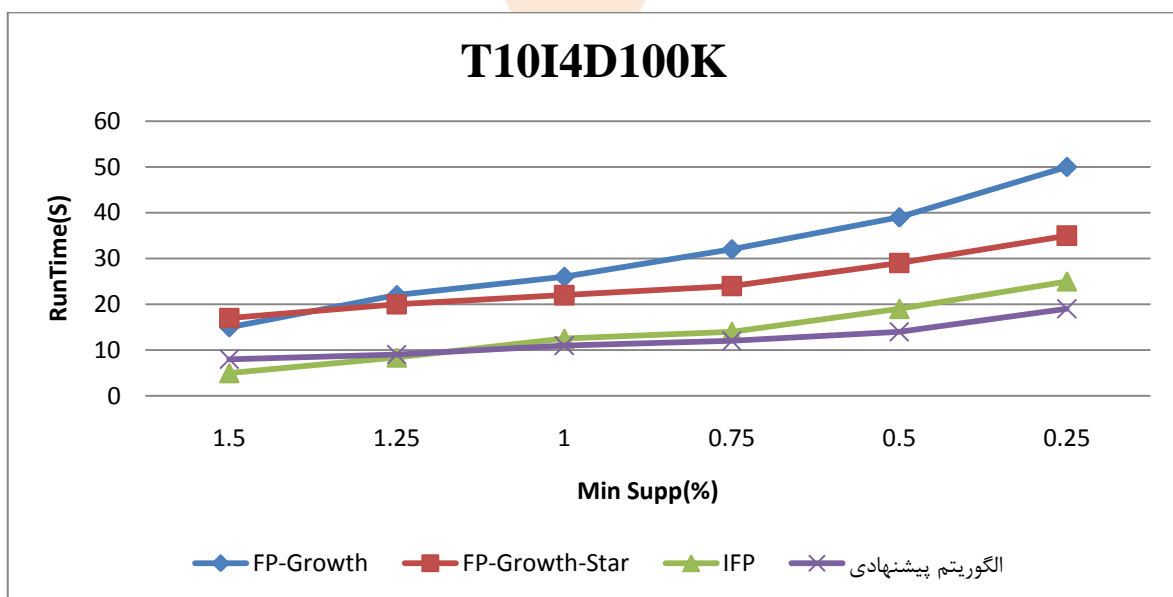
جدول (۱۰) دیتاست T.I.D.K

T40I10D100K	نام دیتاست
100,000	تعداد تراکنش‌ها
942	تعداد آیتم‌ها
39.6	میانگین طول تراکنش

#### ۱,۳,۳ نتایج اجرا بر روی دیتاست T10I4D100K

نتایج اجرای الگوریتم‌های FP-Growth، FP-Growth-Star، IFP و الگوریتم پیشنهادی بر روی این دیتاست در شکل (۹) آورده شده است. در این نمودار محور افقی بیانگر حد آستانه پشتیبانی کمینه و محور عمودی بیانگر مدت زمان اجرای الگوریتم بر حسب ثانیه است. الگوریتم استار با استفاده از تکنیک آرایه برای ساختن درخت‌های شرطی یک بار کمتر نسبت به الگوریتم الگوی متناوب مجبور به پیمایش درخت اصلی است و لذا بهتر از آن عمل می‌کند. هرچه حد آستانه کمتر باشد، درخت عریض‌تر و عمیق‌تر می‌شود در نتیجه الگوریتم استار کارایی خود را بهتر نشان

می‌دهد. با افزایش حد آستانه تعداد قلم‌داده‌های متناوب کمتر شده و در نتیجه آن، الگوهای متناوب نیز کمتر می‌شود. با کاهش تعداد الگوها، درخت کوچکتر و کوتاه‌تر می‌شود و در نتیجه الگوریتم رشد الگوی متناوب سریعتر آن را پیمایش می‌کند و چون درخت‌های شرطی نیز کمتر و کوتاه‌تر می‌شوند لذا سرعت اجرای الگوریتم بیشتر می‌شود. با افزایش حد آستانه و کوچک‌تر شدن درخت، سربار تولید ماتریس تکرار و پیمایش آن باعث کند شدن الگوریتم استار نسبت به الگوریتم الگوی متناوب می‌شود به گونه‌ای که در آستانه ۱,۵٪ الگوریتم استار بدتر عمل می‌کند. الگوریتم IFP به علت حذف ۵۰٪ درخت‌های شرطی در همه آستانه‌ها بهتر از الگوریتم استار و الگوی متناوب عمل می‌کند. الگوریتم پیشنهادی در آستانه‌های کمتر که درخت عریض‌تر و عمیق‌تر است کارایی خود را نشان می‌دهد؛ چرا که با هر بار پیمایش ماتریس، اکثر خانه‌ها مقداری بیشتر از حد آستانه دارند و در نتیجه زمانی که صرف پیمایش ماتریس می‌شود با استخراج الگو همراه است. بهبود روش کاوش نیز باعث می‌شود تا در هر بار پیمایش ماتریس، الگوهای بیشتری استخراج شوند. با افزایش حد آستانه و کوچک‌تر شدن درخت، تاثیر سربار تشکیل ماتریس دوبعدی همانند الگوریتم استار مشخص‌تر می‌شود. در حد آستانه ۱٪ هر دو روش پیشنهادی و IFP بسیار نزدیک به هم می‌شوند اما با گذشت از این حد آستانه، الگوریتم IFP عملکرد بهتری از خود نشان می‌دهد.

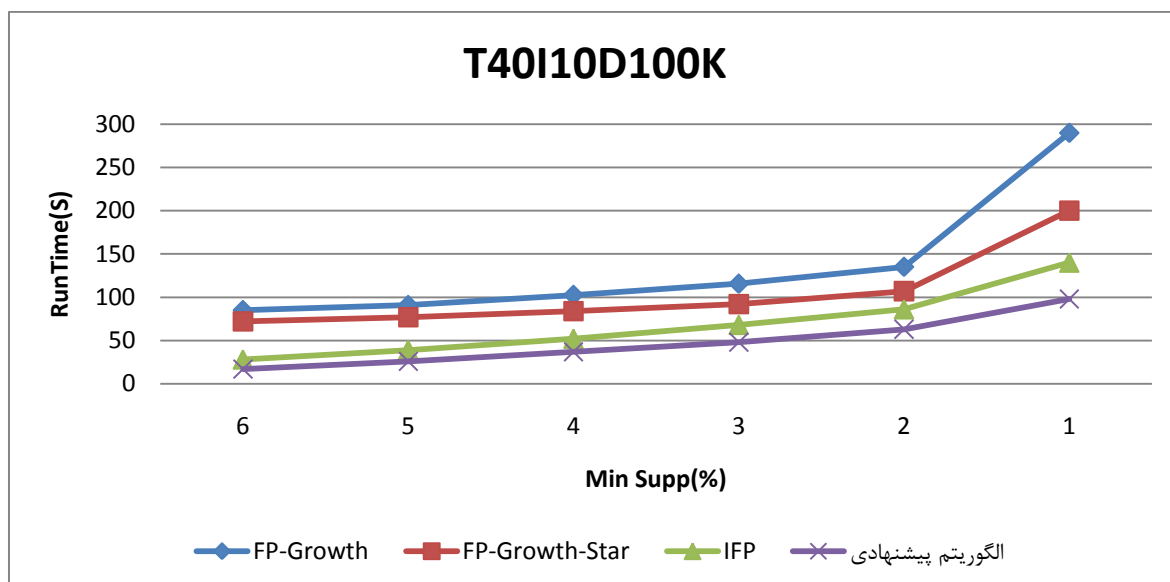


شکل (۹) نتایج اجرا بر روی دیتاست T<sub>10</sub>I<sub>4</sub>D<sub>100</sub>K

### ۲,۳,۳ نتایج اجرا بر روی دیتاست T<sub>4</sub>·I<sub>4</sub>D<sub>100</sub>K

نتایج اجرای الگوریتم‌های IFP, FP-Growth-Star, FP-Growth و الگوریتم پیشنهادی بر روی این دیتاست در شکل (۱۰) آورده شده است. در این دیتاست به علت افزایش متوسط طول تراکنش‌ها عمق درخت بیشتر می‌شود. از طرف دیگر افزایش طول تراکنش‌ها باعث شده تا قلم‌داده‌های بیشتری در یک تراکنش قرار گیرند و لذا تعداد تکرار هر قلم‌داده بیشتر شده است. با در نظر گرفتن این نکته کاملاً مشهود است که با کمتر کردن حد آستانه پشتیبان کمینه تعداد قلم‌داده‌های متناوب زیاد شده و در نتیجه تعداد الگوهای متناوب نیز زیاد می‌شوند. این امر تاثیرش را بر روی هر الگوریتم با افزایش زمان اجرا نشان می‌دهد. به همین خاطر است که حد آستانه پشتیبان کمینه در محدوده ۱٪ تا ۶٪ انتخاب شده است. با افزایش قلم‌داده‌های متناوب تعداد شاخه‌های درخت نیز افزایش می‌یابد و درخت عریض‌تر می‌شود. در این شرایط الگوریتم‌هایی که درخت را بیشتر پیمایش می‌کنند کندتر می‌شوند. برای آستانه‌های بزرگتر درخت کوچک‌تر می‌شود و در نتیجه روش پیشنهادی به خاطر پیمایش ماتریس دچار سربار می‌شود؛ چرا که اکثر خانه‌های ماتریس از حد آستانه کوچک‌تر درخت انبوه‌تر شده و لذا روش پیشنهادی با هر بار پیمایش ماتریس، با توجه به تغییر روش کاوش، می‌تواند تعداد بیشتری الگوی متناوب تولید کند و به همین خاطر بهتر از سایر الگوریتم‌ها عمل می‌کند.





شکل (۱۰) نتایج اجرا بر روی دیتاست T40I10D100K

#### ۴. نتیجه گیری و کارهای آینده

مهم‌ترین بخش در کشف قوانین وابستگی یافتن الگوهای متناوب است. از سال ۱۹۹۴ که آگروال خاصیت Apriori را در مجموعه قلم‌داده‌های متناوب ارائه کرد تاکنون روش‌های متعددی برای یافتن الگوهای متناوب ارائه شده‌اند. تمامی این روش‌ها در دو دسته کلی قرار می‌گیرند. دسته اول آنهایی هستند که بر اساس تولید کاندیدها و شمارش آن‌ها عمل می‌کنند و دسته دوم قلم‌داده‌های متناوب را رشد می‌دهند تا الگوهای جدید تولید شوند. یکی از معروف‌ترین روش‌های مبتنی بر رشد الگو، الگوریتم رشد الگوی متناوب است که پایگاه‌تراکنش‌ها را در یک ساختار درختی به نام درخت الگوی متناوب به صورت فشرده نگهداری می‌کند و از این ساختار برای کاوش الگوهای متناوب بهره می‌برد. مهم‌ترین مشکل این روش هزینه زیاد ساختن درخت‌های شرطی است که به طور بازگشتی ساخته و پیمایش می‌شوند. الگوریتم‌های استار و IFP با ارائه راه‌کارهایی سعی در بهبود زمان اجرای الگوریتم رشد الگوی متناوب داشته‌اند. الگوریتم استار با استفاده از یک ساختار ماتریسی تعداد دفعات پیمایش درخت را کاهش داده و باعث سرعت بخشیدن به الگوریتم می‌شود. در الگوریتم IFP با تغییر ساختار گره و تغییر روش کاوش و حذف حداکثر نیمی از درختان شرطی باعث بهبود الگوریتم می‌شود. بر اساس مطالعات انجام شده در این تحقیق، یکی از روش‌های ممکن برای بهبود بخشیدن به الگوریتم رشد الگوی متناوب تغییر ساختار داده مورد استفاده در پیاده‌سازی درخت است. نگهداری درخت در لیست‌های پیوندی اتلاف زمان زیادی را در هنگام پیمایش به الگوریتم تحمیل می‌کند، در صورتیکه با تغییر ساختار و نگهداری درخت در ساختار داده آرایه این اتلاف زمان به حداقل می‌رسد. از طرف دیگر بجای ساختن متوالی و بازگشتی درختان شرطی می‌توان از ساختار ماتریس برای کاوش و رشد الگوهای جدید استفاده کرد. این کار نه تنها از سربار ساختن درختان شرطی به شکل محسوسی می‌کاهد، بلکه این امکان را می‌دهد تا در هر بار اجرای عملیات کاوش، الگوی جاری تا حداکثر دو مرحله رشد داده شود. همانطور که نتایج تجربی نشان می‌دهد الگوریتم پیشنهادی بر روی دیتاست T10I4D100K در اکثر آستانه‌ها و بر روی دیتاست T40I10D100K در تمامی آستانه‌ها بهتر عمل می‌کند. اگرچه با استفاده از ساختار ماتریس دوبعدی در الگوریتم پیشنهادی عملیات کاوش سریعتر شده است اما حالت‌هایی وجود دارد که ماتریس به صورت ماتریس خلوت در می‌آید. در این حالت‌ها بیشتر خانه‌های ماتریس دارای مقدار عددی ۰ هستند و لذا پیمایش آن‌ها نه تنها باعث تولید الگوی متناوب جدید نمی‌شود بلکه سرباری را در زمان ایجاد و زمان پیمایش بر الگوریتم تحمیل می‌نماید. طراحی ساختارهای داده‌ای که از پیمایش خانه‌های دارای مقدار ۰ در ماتریس دو بعدی جلوگیری می‌کند می‌تواند به عنوان یکی از زمینه‌های کاری آینده مورد بررسی قرار بگیرد.

- [۱] R. Agrawal, T. Imielinski, and A. Swami, “Mining Association rules Between Sets of Items in Large Databases,” Proceedings of ACM SIGMOD, May, ۱۹۹۳.
- [۲] K. Cios, W. Pedrycz, R. Swiniarski, and LA. Kurgan, *Data Mining: A Knowledge Discovery Approach*, Springer, ۲۰۰۷.
- [۳] R. Agrawal, T. Imielinski, and A. Swami, “Mining association rules between sets of terms in large databases,” In Proc. of the ACM- SIGMOD Int. Conf. on management of data (SIGMOD’۹۳), Washington, DC, pp. ۲۰۷–۲۱۶, ۱۹۹۳.
- [۴] R. Agrawal, and R. Srikant, “Fast algorithms for mining association rules,” In Proc. Int. Conf. Very Large Data Bases (VLDB’۹۴), pages ۴۸۷–۴۹۹, Santiago, Chile ,Sept. ۱۹۹۴.
- [۵] J. Han, J. Pei, and Y. Yin, “Mining frequent patterns without candidate generation,” In Proc. of the ACM-SIGMOD Int. Conf. on management of data (SIGMOD’۰۰), Dallas, TX, pp. ۱–۱۲, ۲۰۰۰.
- [۶] G. Grahne, and J. Zhu, “Fast Algorithms for Frequent Itemset Mining Using FP-Trees,” IEEE Trans. On Knowledge and Data Engineering, vol. ۱۷, no. ۱۰, October ۲۰۰۵.
- [۷] K. C. Lin, I. E. Liao, and Z. S. Chen, “An Improved Frequent Pattern Growth Method for Mining Association Rules,” Expert Systems with Applications, vol. ۳۸, no. ۵, pp. ۵۱۵۴–۵۱۶۱, ۲۰۱۱

زیر نویس ها

- 
- <sup>۱</sup> Frequent Patterns  
<sup>۲</sup> Itemsets  
<sup>۳</sup> Sequences  
<sup>۴</sup> Structures  
<sup>۵</sup> Frequent Itemset  
<sup>۶</sup> Frequent Structural Pattern  
<sup>۷</sup> Minimum Support  
<sup>۸</sup> Minimum Confidence  
<sup>۹</sup> Association Rule Mining  
<sup>۱۰</sup> Agrawal  
<sup>۱۱</sup> Candid Generation  
<sup>۱۲</sup> Pattern Growth  
<sup>۱۳</sup> Scan  
<sup>۱۴</sup> Han  
<sup>۱۵</sup> Frequent Pattern Growth  
<sup>۱۶</sup> FP-Tree  
<sup>۱۷</sup> FP-Growth\*  
<sup>۱۸</sup> Grahne & Zhu  
<sup>۱۹</sup> Header Table  
<sup>۲۰</sup> Ke-Chung  
<sup>۲۱</sup> Address Table  
<sup>۲۲</sup> Tree Level  
<sup>۲۳</sup> Minimum Support Threshold  
<sup>۲۴</sup> Frequent Items  
<sup>۲۵</sup> Safe Mode  
<sup>۲۶</sup> Mining Procedure

کنفرانس داده کاوی ایران