



دانشگاه تهران  
دانشکده‌ی مهندسی کامپیوتر

تمرین سوم

داده کاوی

خانم دکتر شاکری

محمد امین عرب خراسانی

۸۱۰۱۰۲۲۰۵

بهار ۱۴۰۳

## بخش تشریحی

### ۱ سوال اول

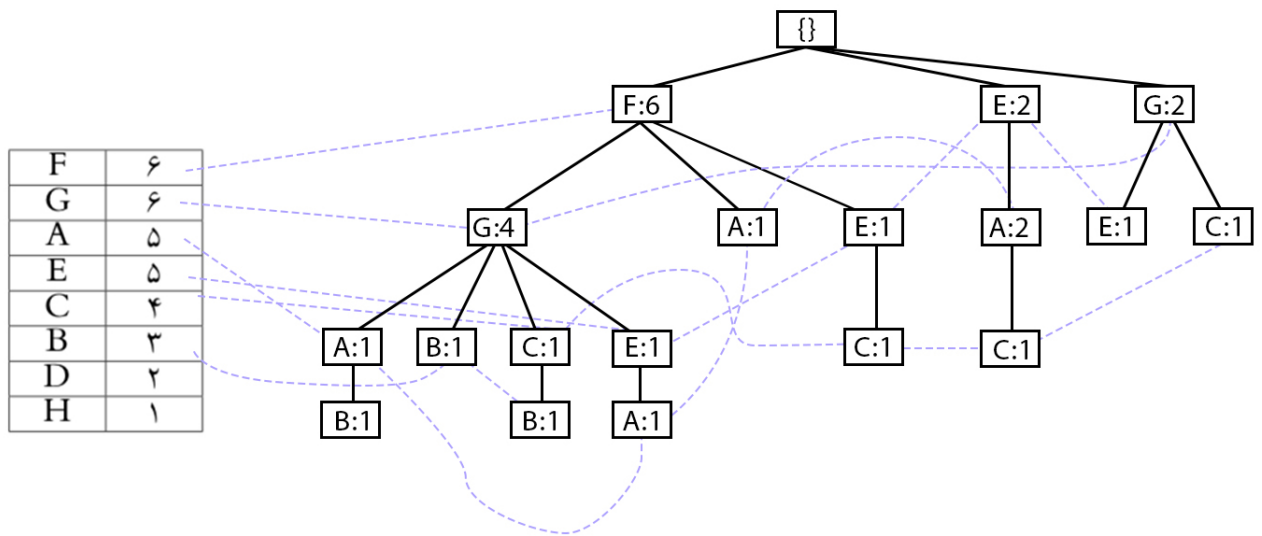
الف) برای رسم FP-tree ابتدا تعداد دفعات تکرار هر item-bought در محاسبه شده و به ترتیب تعداد دفعات تکرار از بیشترین به کمترین دفعات تکرار مرتب می‌شود. در نتیجه F-list به فرم زیر خواهد بود.

F	۶
G	۶
A	۵
E	۵
C	۴
B	۳
D	۲
H	۱

همانطور که از جدول بالا مشخص است، D و H مقدار min-support را ارضا نمی‌کنند بنابراین در FP-tree لحاظ نمی‌شوند. در ادامه transaction ها بر اساس تعداد تکرار مرتب می‌شوند تا FP-tree رسم شود. نتیجه‌ی این مرتب‌سازی در جدول زیر آورده شده است.

Item-bought	rearranged
G, B, A, F	F, G, A, B
H, A, E	E, A
F, B, G	F, G, B
A, C, D, E	E, A, C
C, F, G, B	F, G, C, B
A, F, D	F, A
E, F, G, A	F, G, E, A
C, F, E	F, E, C
G, E	G, E
H, C, G	G, C

در نهایت FP-tree به شکل زیر خواهد بود.



ب) برای به دست آوردن B's conditional تمامی شاخه‌هایی که به B ختم می‌شوند از FP-tree به دست می‌آیند. برای B's conditional داریم:

B's conditional database: FGA:۱، FG:۱، FGC:۱

ج) برای پیدا کردن closed pattern و max-pattern به صورت sequential هر itemset محاسبه می‌شود.

۱-D itemset

item	count
F	۶
G	۶
A	۵
E	۵
C	۴
B	۳
D	۲
H	۱

با توجه به min-support و جدول بالا، H و D کمتر از ۳ بار تکرار شده‌اند بنابراین در ادامه‌ی itemset ها لحاظ نمی‌شوند. بقیه‌ی item-bought ها این شرط را ارضا می‌کنند.

۲-D itemset

item	count
F, G	۴
F, A	۳
F, E	۲
F, C	۲
F, B	۳
G, A	۲
G, E	۲
G, C	۲
G, B	۳
A, E	۳
A, C	۱
A, B	۱
E, C	۲
E, B	۰
C, B	۱

با توجه به شرط minimum-support سطرهایی از جدول بالا که صورتی هستند حذف می‌شوند.

$$F(۶) \Rightarrow FG(۴), FA(۳), FB(۳), FE(۲), FC(۲)$$

با توجه به آن که F(۶) از تمام زیرمجموعه‌ها بزرگتر است بنابراین F، closed می‌باشد. از آنجایی که زیرمجموعه‌ای وجود دارد که شرط minimum-support را ارضا می‌کند بنابراین F، maximal نمی‌باشد.

$$G(۶) \Rightarrow FG(۴), GA(۲), GE(۲), GC(۲), GB(۳)$$

با توجه به آن که G(۶) از تمام زیرمجموعه‌ها بزرگتر است بنابراین G، closed می‌باشد. از آنجایی که زیرمجموعه‌ای وجود دارد که شرط minimum-support را ارضا می‌کند بنابراین G، maximal نمی‌باشد.

$$A(۵) \Rightarrow FA(۳), GA(۲), AC(۱), AB(۱)$$

با توجه به آن که A(۵) از تمام زیرمجموعه‌ها بزرگتر است بنابراین A، closed می‌باشد.

از آنجایی که زیرمجموعه‌ای وجود دارد که شرط minimum-support را ارضا می‌کند بنابراین A، maximal نمی‌باشد.

$$E(5) \Rightarrow FE(2), GE(2), AE(3), EC(2), EB(0)$$

با توجه به آن که E(5) از تمام زیرمجموعه‌ها بزرگتر است بنابراین E، closed می‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود دارد که شرط minimum-support را ارضا می‌کند بنابراین E، maximal نمی‌باشد.

$$C(4) \Rightarrow FC(2), GC(2), AC(1), EC(2), CB(1)$$

با توجه به آن که C(4) از تمام زیرمجموعه‌ها بزرگتر است بنابراین C، closed می‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود ندارد که شرط minimum-support را ارضا کند بنابراین C، maximal می‌باشد.

$$B(3) \Rightarrow FB(3), GB(3), AB(1), EB(0), CB(1)$$

با توجه به آن که B(3) از تمام زیرمجموعه‌ها بزرگتر نیست بنابراین B، closed نمی‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود دارد که شرط minimum-support را ارضا می‌کند بنابراین B، maximal نمی‌باشد.

۳-D itemset

item	count
F, G, A	۲
F, G, B	۳
F, A, B	۱
F, E, A	۱

از آنجایی که ۱۰ تا از زیرمجموعه‌های itemset قبلی frequent نیستند فقط برای زیرمجموعه‌های frequent بررسی می‌شوند.

$$FG(4) \Rightarrow FGA(2), FGB(3)$$

با توجه به آن که FG(4) از تمام زیرمجموعه‌ها بزرگتر است بنابراین FG، closed می‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود دارد که شرط minimum-support را ارضا می‌کند بنابراین FG، maximal نمی‌باشد.

$$FA(3) \Rightarrow FGA(2), FAB(1), FEA(1)$$

با توجه به آن که  $FA(3)$  از تمام زیرمجموعه‌ها بزرگتر است بنابراین  $FA$  closed می‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود ندارد که شرط  $minimum\text{-}support$  را ارضا کند بنابراین  $FA$  maximal می‌باشد.

$$FB(3) \Rightarrow FGB(3), FAB(1)$$

با توجه به آن که  $FB(3)$  از تمام زیرمجموعه‌ها بزرگتر نیست بنابراین  $FB$  closed نمی‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود دارد که شرط  $minimum\text{-}support$  را ارضا می‌کند بنابراین  $FB$  maximal نمی‌باشد.

$$GB(3) \Rightarrow FGB(3)$$

با توجه به آن که  $GB(3)$  از تمام زیرمجموعه‌ها بزرگتر نیست بنابراین  $GB$  closed نمی‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود دارد که شرط  $minimum\text{-}support$  را ارضا می‌کند بنابراین  $GB$  maximal نمی‌باشد.

$$AE(3) \Rightarrow FEA(1)$$

با توجه به آن که  $AE(3)$  از تمام زیرمجموعه‌ها بزرگتر است بنابراین  $AE$  closed می‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود ندارد که شرط  $minimum\text{-}support$  را ارضا کند بنابراین  $AE$  maximal می‌باشد.

۴-D itemset

item	count
F, G, B, A	۱
F, G, C, B	۱

$$FGB(3) \Rightarrow FGAB(1), FGCB(1)$$

با توجه به آن که  $FGB(3)$  از تمام زیرمجموعه‌ها بزرگتر است بنابراین  $FGB$  closed می‌باشد.  
از آنجایی که زیرمجموعه‌ای وجود ندارد که شرط  $minimum\text{-}support$  را ارضا کند بنابراین  $FGB$  maximal می‌باشد.

در نهایت جدول زیر  $closed\ pattern$  و  $max\text{-}pattern$  آورده شده است.

$closed\ pattern$ : F, G, A, E, C, FG, FA, AE, FGB

$max\text{-}pattern$ : C, FA, AE, FGB

(د) با توجه به itemset های حاصل در بخش قبل، و فرمول زیر جدول مربوط به association rule تکمیل می شود.

$$\{X\} \rightarrow \{Y\} : \frac{\text{support}(\{X, Y\})}{\text{support}(\{X\})}$$

$\{F\} \rightarrow \{G\} : 0.67$	$\{G\} \rightarrow \{F\} : 0.67$	$\{A\} \rightarrow \{F\} : 0.6$	$\{A\} \rightarrow \{E\} : 0.6$
$\{E\} \rightarrow \{A\} : 0.6$	$\{B\} \rightarrow \{F\} : 1$	$\{B\} \rightarrow \{G\} : 1$	$\{F, G\} \rightarrow \{B\} : 0.75$
$\{B\} \rightarrow \{F, G\} : 1$	$\{F, B\} \rightarrow \{G\} : 1$	$\{G\} \rightarrow \{F, B\} : 1$	$\{B, G\} \rightarrow \{F\} : 1$
$\{F\} \rightarrow \{B, G\} : 1$			

همانطور که از نتایج مشخص است association rules قوی با فرض minimum-confidence ارائه شده حاصل می شود.

(ه) ابتدا یک جدول تشکیل داده می شود که شامل support count های هر آیت می باشد. ترتیب این جدول بر اساس تعداد دفعات تکرار است.

### ۱-D itemset

itemset	support-count
A	۴
D	۴
F	۴
G	۴
B	۲
C	۲
E	۱
H	۱

با توجه به minimum-support داده شده، item هایی که frequent نمی باشند حذف می شوند.

### ۲-D itemset

itemset	support-count
F, G	۴
A, D	۳
A, F	۳
A, G	۳
D, F	۳
D, G	۳

در این مرحله تمام زیر مجموعه ها حفظ می شوند.

### ۳-D itemset

itemset	support-count
D, F, G	۳
A, F, G	۳
A, D, F	۲
A, D, G	۲

برای تشکیل مراحل فوق برای مجموعه های چهارتایی فقط ۲ آیتمست سه تایی موجود است و بنابراین یک مجموعه ۴ تایی از ادغام آن ها می توان ساخت که به صورت A, D, F, G است، که از زیرمجموعه های همین مجموعه چهارتایی، موردی مثل A, D, G در لیست پیشین وجود ندارد یا Frequent نبوده و تاییدی است بر اینکه مراحل در همین جا متوقف شده و هیچ subset دیگری نمی توان ساخت.



و) همانطور که در مراحل قبل اشاره شد با توجه به مراحل الگوریتم، بزرگترین itemset ممکن مربوط به A، D، F و A، D، G می باشد. زیرا در A، D، F، G که خود مکرر نمی باشد آیتمست های دیگری نظیر A، D، F وجود دارد که frequent نمی باشد. اما در A، D، F و A، D، G خود frequent نمی باشند و آیتمست دیگری غیر از تهی ندارند.

ز) با توجه به meta rule ی که در صورت سوال اشاره شده است، فقط آن آیتمست هایی لحاظ می شوند که به فرم زیر باشند.

$$\{X, Y\} \rightarrow \{Z\}$$

در نتیجه برای association rules قوی با توجه به minimum-confidence عنوان شده داریم:

$$\{F, G\} \rightarrow \{A\} : \frac{\text{support}(\{A, F, G\})}{\text{support}(\{F, G\})} = \frac{3}{4} \times 100 = 75\%$$

$$\{A, F\} \rightarrow \{G\} : \frac{\text{support}(\{A, F, G\})}{\text{support}(\{A, F\})} = \frac{3}{3} \times 100 = 100\%$$

$$\{A, G\} \rightarrow \{F\} : \frac{\text{support}(\{A, F, G\})}{\text{support}(\{A, G\})} = \frac{3}{3} \times 100 = 100\%$$

$$\{F, G\} \rightarrow \{D\} : \frac{\text{support}(\{D, F, G\})}{\text{support}(\{F, G\})} = \frac{3}{4} \times 100 = 75\%$$

$$\{D, F\} \rightarrow \{G\} : \frac{\text{support}(\{D, F, G\})}{\text{support}(\{D, F\})} = \frac{3}{3} \times 100 = 100\%$$

$$\{D, G\} \rightarrow \{F\} : \frac{\text{support}(\{D, F, G\})}{\text{support}(\{D, G\})} = \frac{3}{3} \times 100 = 100\%$$

## ۲ سوال دوم

الف) جمله ی اول: درست است. از آنجایی که این association قوی بین s و l برقرار است به این معناست که شرط min-confidence ارضا شده است.

$$\{s\} \rightarrow \{l\} : \frac{\text{support}(\{s, l\})}{\text{support}(\{s\})} \geq \text{min-confidence}$$

برای آیت‌ست a نیز داریم:

$$\{a\} \rightarrow \{l\} : \frac{\text{support}(\{a, l\})}{\text{support}(\{a\})}$$

از آنجایی که a زیرمجموعه‌ی s است بنابراین هر ترنزاکشنی از a در s وجود دارد پس داریم:

$$\text{support}(a) \leq \text{support}(s) \rightarrow \text{support}(a, l) \leq \text{support}(s, l)$$

$$\frac{\text{support}(\{a, l\})}{\text{support}(\{a\})} \leq \frac{\text{support}(\{s, l\})}{\text{support}(\{s\})}$$

جمله‌ی دوم: الف) ۱. نادرست است.

برای آیت‌ست a داریم:

$$\{a\} \rightarrow \{l\} : \frac{\text{support}(\{a, l\})}{\text{support}(\{a\})}$$

فرض می‌کنیم a و s برابر مجموعه‌ی زیر است.

a: {bread, egg, butter}

s: {bread, egg}

$$\{bread, egg, butter, l\} \rightarrow \{l\} : \frac{\text{support}(\{a, l\})}{\text{support}(\{a\})}$$

ب) برای مقایسه‌ی confidence این دو و اثبات عبارت موجود در سوال ابتدا مقادیر confi-

dence محاسبه می‌شود. بنابراین داریم:

$$\text{conf}(a \rightarrow (c - a)) = \frac{sc - sa}{sa}$$

$$\text{conf}(b \rightarrow (c - b)) = \frac{sc - sb}{sb}$$

از آن جایی که  $a \subseteq b$  و  $b \subseteq c$  است بنابراین هر  $sa \leq sb$  است زیرا همه‌ی زیرمجموعه‌های ممکن برای  $a$  در آیت‌ست  $b$  می‌باشد. با توجه به این نکته confidence ها با یکدیگر مقایسه می‌شوند.

$$\frac{sc - sb}{sb} \leq \frac{sc - sa}{sa}$$

$$\frac{sc}{sb} - 1 \leq \frac{sc}{sa} - 1$$

$$\frac{sc}{sb} \leq \frac{sc}{sa}$$

همانطور که مشخص است در نهایت عبارت موجود در سوال اثبات می‌شود.

### ۳ سوال سوم

الف) استخراج الگوهای متداول در sequential pattern mining به معنای کشف دنباله‌هایی است که به طور قابل ملاحظه در داده‌ها ظاهر می‌شوند. این الگوها نشان‌دهنده رخداد‌های مکرر و الگوهایی هستند که در داده‌ها به صورت پنهانی وجود دارند. در sequential pattern mining زمان ثبت یک ترنزاکشن ثبت نمی‌شود. این مورد در شرایطی که کشف الگوهای متوالی حیاتی است به کار می‌رود. از طرفی Association rule هدف اصلی شناسایی روابط یا ارتباطات جالب بین موارد مختلف در مجموعه داده است، و اغلب در داده‌های تراکشی می‌بینیم. همچنین این مورد در درک ارتباطات بین موارد مهم موثر می‌باشد.

ب) الگوریتم PrefixSpan از نوع الگوریتم‌های DFS است. این الگوریتم به صورت بازگشتی

در گراف جستجوی دنباله حرکت می‌کند. به عبارت دیگر، ابتدا به عمق دنباله می‌رود و سپس به طور بازگشتی بازگشته و به حالت‌های دیگری از گراف حرکت می‌کند. این رویکرد برای کاهش زمان و حافظه مصرفی مفید است زیرا فقط دنباله‌هایی که قابلیت گسترش دارند را بررسی می‌کند و دنباله‌های غیرقابل گسترش را از لیست حذف می‌کند، بدون این که تمام فضای جستجو را بررسی کند. این ویژگی‌ها به الگوریتم PrefixSpan کمک می‌کند تا به صورت کارآمد الگوهای دنباله‌ای متداول را در مجموعه داده‌های بزرگی کشف کند. در قسمت ج یک مثال از این الگوریتم حل می‌شود.

ج) برای حل این قسمت با توجه به الگوریتم ذکر شده و با توجه به مقدار min-support مسئله حل می‌شود.

ID	Sequence
۱S	a , b , c , b , b , c , d
۲S	d , c , b , a , b , c
۳S	c , b , b , c , d

برای هر مرحله min-support چک می‌شود. نتایج به شرح زیر است.

<a>	۲
<b>	۳
<c>	۳
<d>	۳

b	c	(d
cbbcd	bbcd	(cbabc
abc	babc	)
bcd	bbcd	)

در ادامه ترکیبات مختلفی از هر جدول در نظر گرفته می‌شود تا جایی که min-support ارضا نشود

## ۴ سوال چهارم

برای حل این سوال از روابط داده شده در صورت سوال استفاده می‌شود.

$$\chi^2 = \frac{(400-240)^2}{240} + \frac{(200-360)^2}{360} + \frac{(400-560)^2}{560} + \frac{(1000-840)^2}{840} = 253.8$$

با توجه به عدد حاصل، همبستگی بالای بین این دو مولفه شناسایی می‌شود.

$$lift(carrot, onion) = \frac{\frac{400}{2000}}{\frac{600}{2000} \cdot \frac{800}{2000}} = 1/667$$

با توجه به آن که مقدار lift بیشتر از یک به دست آمد بنابراین این دو با یک دیگر همبستگی مثبت دارن و این بدان معناست که با خرید هویج احتمال خرید پیاز زیاد است.

$$Kulczynski(carrot, onion) = \frac{400}{2000} \left( \frac{1}{800} + \frac{1}{600} \right) = 0/58$$

عدد به دست آمدع نشان دهنده‌ی وابستگی این دو به یکدیگر می‌باشد.

$$IR(carrot, onion) = \frac{|800-600|}{800+600-400} = 0/2$$

با توجه به نزدیکی این عدد به صفر می‌توان بالانس بودن این دو مولفه را نتیجه گرفت.

## ۵ سوال پنجم

الف) anti-monotone است. زیرا اگر یک عضو در مجموعه‌ی S محدودیتی را نقض کند، هر عضو جایگزین s' آن نیز محدودیت را نقض خواهد کرد. به عبارت دیگر، هرگاه یک مورد از مجموعه S قانونی را نقض کند، هر عضو جایگزین آن نیز قانونی را نقض می‌کند. علاوه بر این، اگر مجموع قیمت‌های مورد s برابر یا بیشتر از یک حداقل تعیین شده باشد، مجموع قیمت‌های مورد s' هم برابر یا بیشتر از این حد خواهد بود. در نتیجه، اگر تراکم قیمت‌های s بالاتر از حداقل تعیین شده باشد، تراکم قیمت‌های s' نیز بالاتر از آن خواهد بود.

ب) monotone است. به عبارت دیگر، اگر مجموعه s این ویژگی را از خود نشان دهد، یعنی اگر ترتیب مقادیر درون آن افزایشی یا کاهشی باشد، در این صورت برای هر یک از s' ها نیز این خاصیت حفظ خواهد شد. به این معنا که مجموع قیمت‌های s' همواره برابر یا بیشتر از حداقل مقدار تعیین شده خواهد بود.

ج) monotone است. اگر حداقل یک شیر در مجموعه s وجود داشته باشد، این نشان می‌دهد

که در هر مجموعه  $s'$  نیز حداقل یک شیر وجود دارد و شرایط محدودیت همچنان برقرار است. به عبارت دیگر، اگر در یک مجموعه شامل حداقل یک شیر باشد، این ویژگی نسبت به هر مجموعه مشابه دیگری هم حفظ می‌شود و محدودیت‌ها همچنان رعایت می‌شود.

(د) succinct است. اگر تنها داده‌هایی که محدوده‌ی قیمت آن‌ها بین مقادیر  $u$  و  $v$  است را در نظر بگیریم، این محدودیت اعمال می‌شود.