

فصل ۴

تمرین ۱

(الف)

مرحله ۱:

تعداد تکرار هر آیتم در تراکنش‌ها را نسبت به تعداد کل تراکنش‌ها محاسبه می‌کنیم.

ردیف	آیتم	پشتیبان (Support)	درصد	
۱	A	۱-۵	۲۰	
۲	C	۲-۵	۴۰	
۳	D	۱-۵	۲۰	
۴	E	۴-۵	۸۰	✓
۵	I	۱-۵	۲۰	
۶	K	۵-۵	۱۰۰	✓
۷	M	۳-۵	۶۰	✓
۸	N	۲-۵	۴۰	
۹	O	۲-۵	۶۰	✓
۱۰	U	۱-۵	۲۰	
۱۱	Y	۳-۵	۶۰	✓

مرحله ۲:

از مرحله ۱ اقلامی که پشتیبان ۶۰ درصد و بیشتر از ۶۰ درصد دارند را انتخاب می‌کنیم. مجموعه اقلام ۲ تایی از آن‌ها می‌سازیم و مجدد پشتیبان را برای این مجموعه جدید محاسبه می‌کنیم.

ردیف	آیتم	پشتیبان (Support)	درصد	
۱	E, k	۴۵	۸۰	✓
۲	E, M	۳۵	۴۰	
۳	E, O	۳۵	۶۰	✓
۴	E, Y	۳۵	۴۰	
۵	K, M	۳۵	۶۰	✓
۶	K, O	۳۵	۶۰	✓
۷	K, Y	۳۵	۶۰	✓
۸	M, O	۱۵	۲۰	
۹	M, Y	۲۵	۴۰	
۱۰	O, Y	۲۵	۴۰	

مرحله ۳:

با اقلام ۲ تایی از مرحله ۲ که پشتیبان برابر یا بیشتر از ۶۰ درصد دارند و اقلامی که از مرحله ۱ قبلاً انتخاب شده‌اند، مجموعه اقلام ۳ تایی می‌سازیم و پشتیبان را برای این مجموعه اقلام ۳ تایی محاسبه می‌کنیم.

ردیف	آیتم	پشتیبان (Support)	درصد	
۱	E, K, M	۲۵	۴۰	
۲	E, K, O	۳۵	۶۰	✓
۳	E, K, Y	۳۵	۴۰	
۴	E, M, O	۱۵	۲۰	
۵	E, O, Y	۲۵	۴۰	
۶	K, M, O	۱۵	۲۰	
۷	K, M, Y	۲۵	۴۰	
۸	K, O, Y	۲۵	۴۰	

مرحله ۴:

با اقلام ۳ تایی از مرحله ۳ که پشتیبان برابر یا بیشتر از ۶۰ درصد دارند و اقلامی که از مرحله ۱ قبلا انتخاب شده‌اند، مجموعه اقلام ۴ تایی می‌سازیم و پشتیبان را برای این مجموعه اقلام ۴ تایی محاسبه می‌کنیم.

ردیف	آیتم	پشتیبان (Support)	درصد
۱	E, K, M, O	$\frac{1}{5}$	۲۰
۲	E, K, O, Y	$\frac{2}{5}$	۴۰

در این مرحله چون هیچ مجموعه اقلام ۴ تایی پشتیبان برابر یا بیشتر از ۶۰ درصد ندارد، بنابراین الگوریتم متوقف می‌شود.

نتیجه نهایی:

{{E}, {K}, {M}, {O}, {Y}, {E, K}, {E, O}, {K, M}, {K, O}, {K, Y}, {E, K, O}}

(ب)

مرحله ۱:

برای هر آیتم، پشتیبان (support) را محاسبه می‌کنیم:

ردیف	آیتم	پشتیبان (support)	درصد
۱	M	$\frac{3}{5}$	۶۰
۲	O	$\frac{3}{5}$	۶۰
۳	N	$\frac{2}{5}$	۴۰
۴	K	$\frac{5}{5}$	۱۰۰
۵	E	$\frac{4}{5}$	۸۰
۶	Y	$\frac{3}{5}$	۶۰
۷	D	$\frac{1}{5}$	۲۰
۸	A	$\frac{1}{5}$	۲۰
۹	U	$\frac{1}{5}$	۲۰
۱۰	C	$\frac{2}{5}$	۴۰
۱۱	I	$\frac{1}{5}$	۲۰

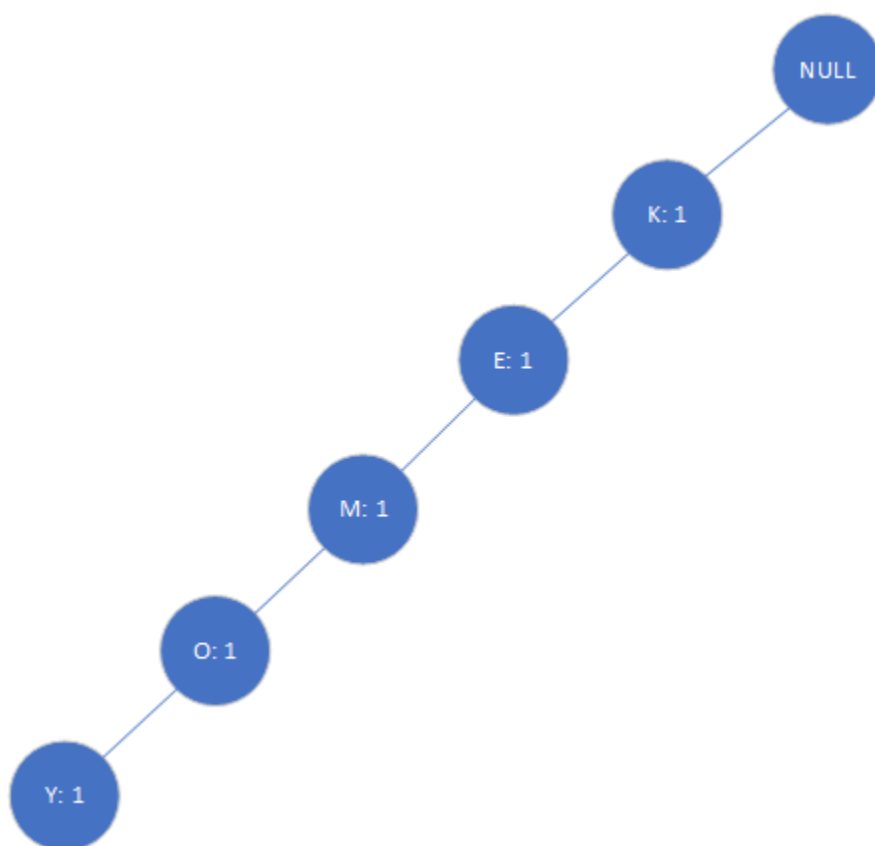
مرحله ۲:

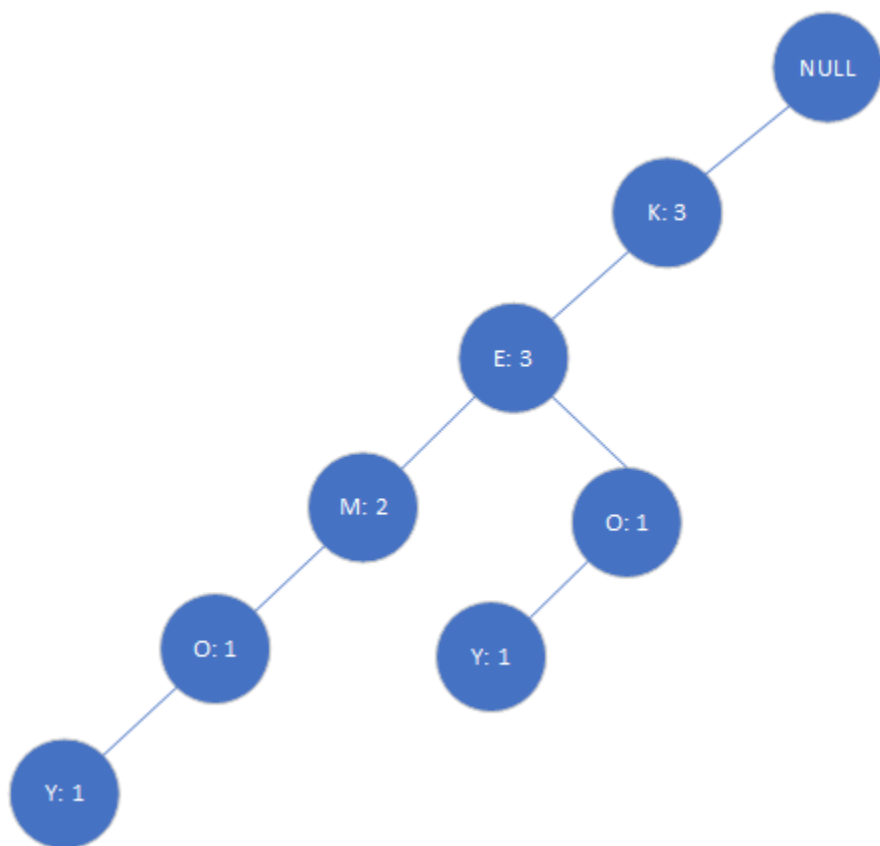
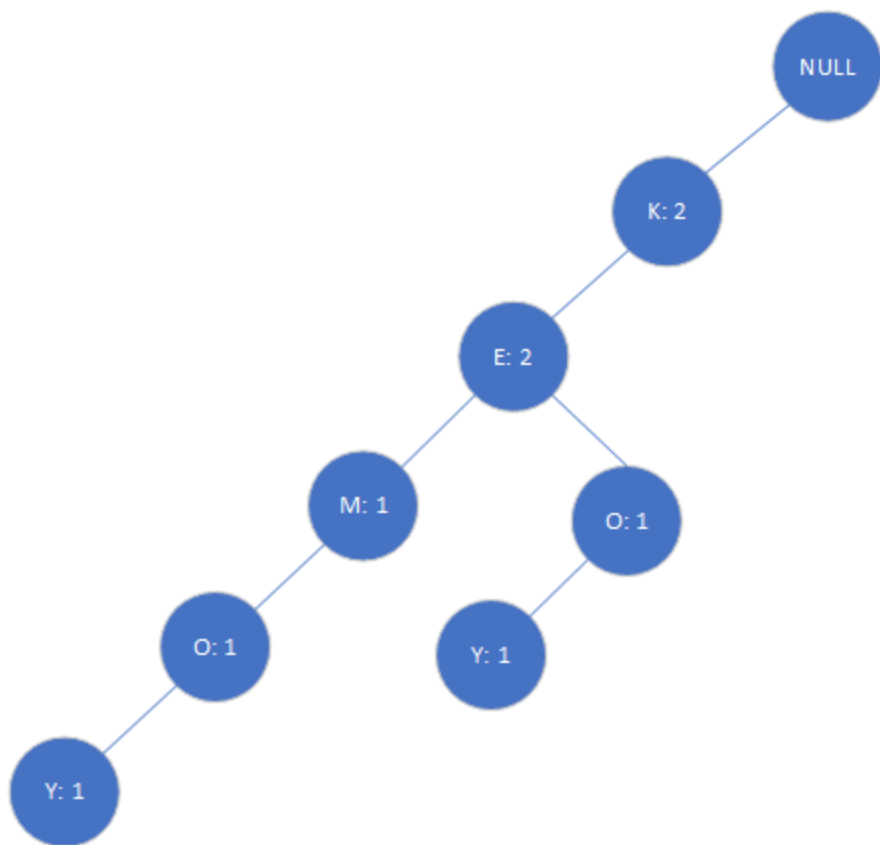
در هر تراکنش آیتم‌هایی که مقدار پشتیبان زیر ۶۰ درصد دارند را حذف می‌کنیم و آیتم‌های باقی مانده را به ترتیب مقدار پشتیبان محاسبه شده در مرحله قبل به صورت نزولی مرتب می‌کنیم:

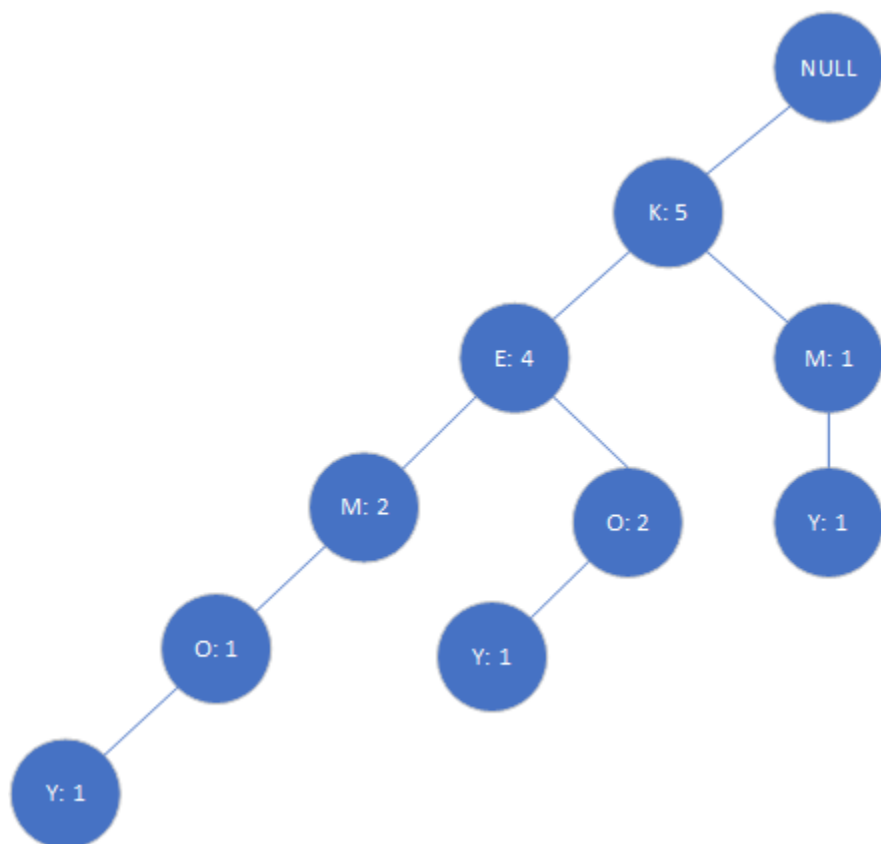
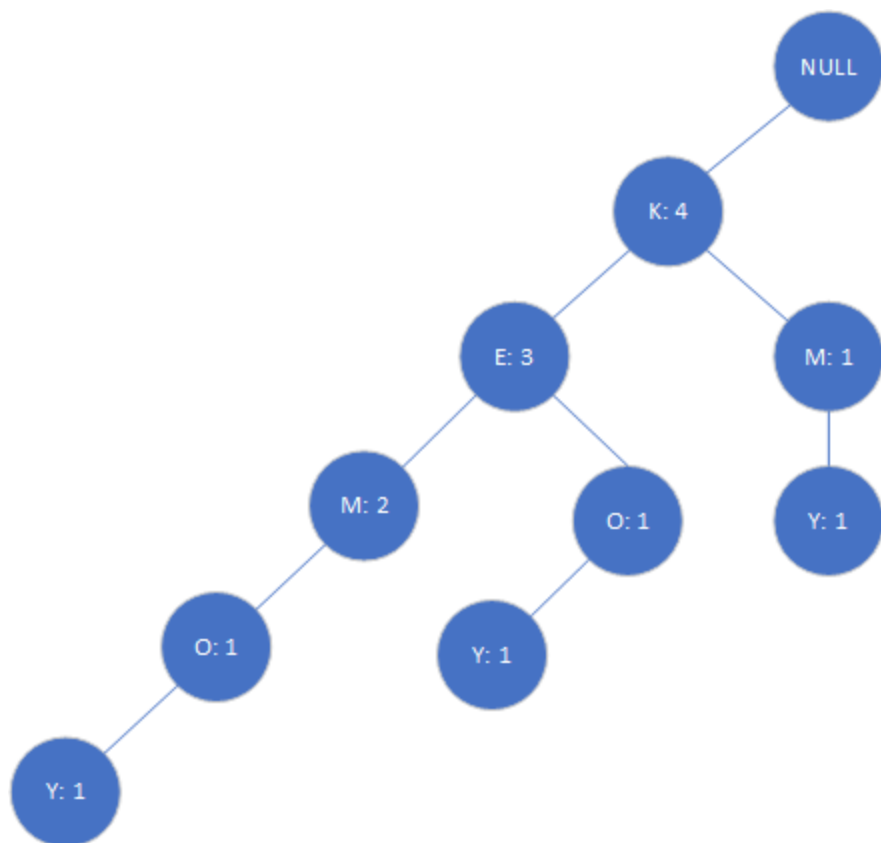
#	TID	Items_bought
1	T100	{K, E, M, O, Y}
2	T200	{K, E, O, Y}
3	T300	{K, E, M}
4	T400	{K, M, Y}
5	T500	{K, E, O}

مرحله ۳:

درخت FP-tree را ایجاد می‌کنیم؛ به این صورت که یک گره با مقدار null ایجاد می‌کنیم. سپس برای هر تراکنش یک شاخه به گره null اضافه می‌کنیم:







مرحله ۴:

پایگاه الگوهای شرطی را ایجاد می‌کنیم.

برای هر آیت Conditional Pattern Base را ایجاد می‌کنیم:

#	Item	Conditional Pattern Base (CPB)
1	Y	{K, E, M, O: 1}, {K, E, O: 1}, {K, M: 1}
2	O	{K, E, M: 1}, {K, E: 2}
3	M	{K, E: 2}, {K: 1}
4	E	{K: 4}

در پایگاه داده جدید ساخته شده مجدد برای هر تراکنش، بازای هر آیت در ستون CPB، پشتیبان را محاسبه کرده و آیت‌هایی که پشتیبان کمتر از ۶۰ درصد دارند را حذف می‌کنیم:

#	Item	Conditional Pattern Base (CPB)	Conditional FP-tree
1	Y	{K, E, M, O: 1}, {K, E, O: 1}, {K, M: 1}	K: 3
2	O	{K, E, M: 1}, {K, E: 2}	K: 3, E: 3
3	M	{K, E: 2}, {K: 1}	K: 3
4	E	{K: 4}	K: 4

در آخر ترکیب اقلام ستون Conditional FP-tree را با ستون Item تولید می‌کنیم:

#	Item	CPB	Conditional FP-tree	Frequent Pattern Generated
1	Y	{K, E, M, O: 1}, {K, E, O: 1}, {K, M: 1}	K: 3	{K, Y}
2	O	{K, E, M: 1}, {K, E: 2}	K: 3, E: 3	{K, Y}, {E, Y}, {K, E, Y}
3	M	{K, E: 2}, {K: 1}	K: 3	{K, M}
4	E	{K: 4}	K: 4	{K, E}

جواب نهایی شامل اقلام محاسبه شده در مرحله اول و Frequent Pattern Generated هستند:

{K}, {E}, {M}, {O}, {Y}, {E, K}, {E, O}, {K, M}, {K, O}, {K, Y}, {E, K, O}

(پ)

رابطه:

$$\forall x \in transaction, buys(X, item_1) \wedge buys(X, item_2) \Rightarrow buys(X, item_3)$$

باید سه آیتی را انتخاب کنیم که:

$$item_1, item_2 \rightarrow item_3$$

برای اینکه در رابطه بالا حداقل مقدار پشتیبان (minsup) را داشته باشیم، باید زیرمجموعه‌هایی ۳تایی از مجموعه ۱۱ آیتم را بیابیم که حداقل مقدار پشتیبان را داشته باشند. تنها مجموعه {E, K, O} حداقل مقدار پشتیبان (۶۰ درصد) را دارد. بنابراین مقدار minconf زیرمجموعه‌های دو عضوی از مجموعه {E, K, O} را محاسبه می‌کنیم:

$$conf(E, K \Rightarrow O) = \frac{3}{4} = 75\%$$

$$conf(K, O \Rightarrow E) = \frac{3}{3} = 100\%$$

$$conf(E, O \Rightarrow K) = \frac{3}{3} = 100\%$$

تمرین ۲

(الف)

ابتدا تراکنش‌ها را براساس کالا بازنویسی می‌کنیم:

TID	Items_bought
T100	{Crab, Milk, Cheese, Bread}
T200	{Cheese, Milk, Apple, Pie, Bread}
T300	{Apple, Milk, Bread, Pie}
T400	{Bread, Milk, Cheese}

سپس الگوهای مکرر را براساس حداقل مقدار پشتیبان (minsup) یعنی ۶۰ درصد پیدا می‌کنیم (برای مثال از روش Apriori):

{{Milk: 100 %}, {Cheese: 75 %}, {Bread: 100 %}, {Milk, Cheese: 75 %}, {Milk, Bread: 100 %}, {Cheese, Bread: 75 %}, {Milk, Cheese, Bread: 75 %}}

مشاهده می‌شود که مجموعه اقلام ۳تایی بزرگ‌ترین مجموعه اقلام است؛ بنابراین برای حالت‌های مختلف مقدار پشتیبان و مقدار اطمینان را براساس قانون کلی زیر محاسبه می‌کنیم:

$$\forall x \in transaction, buys(X, item_1) \wedge buys(X, item_2) \Rightarrow buys(X, item_3)$$

#	Rule	Support	Confidence
1	Milk, Cheese -> Bread	75 %	100 %
2	Milk, Bread -> Cheese	75 %	75 %
3	Cheese, Bread -> Milk	75 %	100 %

بنابراین دو قانون زیر حداقل مقدار پشتیبان (minsup) ۶۰ درصد و حداقل مقدار اطمینان (minconf) ۸۰ درصد را دارند:

Milk, Cheese -> Bread

Cheese, Bread -> Milk

(ب)

ابتدا مقدار پشتیبان را برای هر برند-کالا محاسبه می‌کنیم:

{Best-Bread: 25%}, {Best-Cheese: 25%}, {Dairyland-Cheese: 50%}, {Dairyland-Milk: 50%}, {Goldenfarm-Apple: 25%}, {King's-Crab: 25%}, {Sunset-Milk: 50%}, {Tasty-Pie: 50%}, {Westcoast-Apple: 25%}, {Wonder-Bread: 75%}

مشاهده می‌شود که فقط Wonder-Bread بیشتر از حداقل مقدار پشتیبان است.