

قواعد وابستگی (هم باشی)

Association Rules

(Part 1)

M.M. Pedram
pedram@khu.ac.ir
 Kharazmi University
 (Fall 2007)

1

تحلیل سبد بازار (Market Basket Analysis)

یکی از استراتژیهای داده کاوی

هدف، یافتن ارتباطات جالب بین محصولات خرد فروشی (retail products) است، و تجزیه و تحلیل اقلامی که متحمل است با هم خریداری شوند؛ و با کاهش هزینه ها، سود افزایش یابد.

کاربردهای نتایج تحلیل توسط خرد فروشان
 در ارتقا طراحی،

- ▶ مرتب سازی قفسه ها (shelf)،
- ▶ مواردی که باید در کاتالوگ ذکر شود،
- ▶ استراتژی های micro-marketing و cross-marketing

- ✓ **Cross-marketing**: suggesting related products or services to a customer who is considering buying something.
 - If you're buying a book on Amazon.com, for example, you may be shown a list of books similar to the one you've chosen or books purchased by other customers that bought the same book you did.
 - A search on a company's Web site for bed linens might also bring up listings of matching draperies. The most ubiquitous example of cross-sell is likely the oft-spoken fast food phrase: "Would you like fries with that?"
- ✓ **Micro-marketing**: The study of activities of an organization, i.e. The activities a firm practices in order to react controllably to external forces, e.g., setting objectives and selecting target markets.

قواعد همباشی (Association Rules)

- ❖ کشف وابستگیهای جالب بین attribute هایی که در پایگاه داده وجود دارند،
- ❖ مشهورترین روش برای تحلیل سبد بازار،
- ❖ در نظر گرفتن تمام ترکیبات گروه محصولاتی که به طور بالقوه جالبند،
- ❖ قابلیت تولید هزاران قاعده همباشی با تعداد اندکی attribute.

▶ 3 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

3

تجزیه و تحلیل وابستگی (Affinity Analysis)

- ❖ فرآیند کلی تعیین اینکه چه چیزهایی با هم هستند و جلو می روند.
- ❖ این تجزیه و تحلیل در تجزیه و تحلیل سبد بازار به کار می رود.
- ❖ خروجی تجزیه و تحلیل سبد بازار مجموعه‌ای از همباشی‌های مربوط به رفتار خرید مشتری است.
- ❖ **قواعد همباشی (association rules):** اگر همباشها توسط قواعد توصیف شوند به آنها قواعد همباشی گویند.
▶ کاربرد این قواعد در تعیین استراتژیهای مناسب برای خرید محصولات است.

▶ 4 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

4

تفاوت قواعد همباشی با قواعد عادی کلاس‌بندی

1. در قواعد همباشی، هر **attribute** در یک قاعده می‌تواند به عنوان یک پیش شرط (یا مقدم) استفاده شود و در یک قاعده دیگر در بخش "نتیجه" (تالی) ظاهر شود در حالی که در قواعد عادی اینگونه نیست.
2. در قواعد عادی، در بخش "نتیجه" (تالی) فقط از یک **attribute** استفاده می‌شود، اما تولید کننده‌های قواعد همباشی این امکان را می‌دهند که در نتیجه‌ی (تالی) قواعد چندین **attribute** را استفاده کرد.

▶ 5 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

5

مثالی از قواعد هم باشی

هدف آن است که تعیین شود در خرید اقلام زیر توسط مشتری، چه ارتباطاتی وجود دارد:

- Milk ▪
- Cheese ▪
- Bread ▪
- Eggs ▪

قواعد همباش زیر می‌تواند مطرح باشد:

1. اگر مشتری Milk بخرد، آنگاه Bread هم خواهد خرید.
2. اگر مشتری Bread بخرد، آنگاه Milk هم خواهد خرید.
3. اگر مشتری Milk و Eggs بخرد، آنگاه Cheese و Bread هم خواهد خرید.
4. اگر مشتری Milk و Cheese بخرد، آنگاه Bread هم خواهد خرید.

توضیح: برای قاعده اول (و هر قاعده دیگر) سوالی به صورت زیر می‌تواند مطرح شود:
"چقدر متحمل است که خرید Milk منجر به خرید Bread شود؟"

▶ 6 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

6

”دامنهٔ پشتیبانی قاعده“ و ”میزان اعتماد به قاعده“

درجهٔ پشتیبانی قاعده (support for a rule)

بیانگر درصدی از تمام معاملات (خریدها) است که شامل attribute‌های آن قاعده است.

فرض کنید در مثال قبل، در مجموع 14000 خرید صورت گرفته است و در 5000 مورد خرید Milk و Bread با هم بوده است لذا دامنه کاربرد قاعده (وهمین طور) که این دو صفت را شامل می‌شود، به صورت زیر است:

$$\begin{aligned} Supp(A \Rightarrow B) &= Supp(B \Rightarrow A) = Supp(A, B) \\ &= \frac{|A \cap B|}{|X|} = P(A \cap B) = \frac{5000}{14000} = 0.357 \equiv 35.7\% \end{aligned}$$

درجهٔ اعتماد به قاعده (rule confidence value)

میزان اعتمادیه هر قاعده، یعنی احتمال ظاهر شدن attribute تالی به شرط وجود attribute مقدم، بیان می‌گردد. مثلاً برای قاعده اول اگر 10000 خرید مشتریها شامل Milk باشد و از آنها 5000 خرید شامل خرید Bread هم باشد در آن صورت میزان اعتماد به آن، به صورت زیر است:

$$\begin{aligned} Conf(A \Rightarrow B) &= P(B | A) = \frac{|A \cap B|}{|A|} = \frac{Supp(A \rightarrow B)}{supp(A)}; \quad supp(A) = \frac{|A|}{|X|} \\ P(\text{Bread} | \text{Milk}) &= \frac{5000}{10000} = 0.5 \equiv 50\% \end{aligned}$$

▶ 7 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

7

قاعده قوی (Strong Rule): قاعده‌ای است که دامنه پشتیبانی و میزان اعتماد آن بالا باشد.

▶ هدف استخراج قواعد هم‌باشی، استخراج قواعد قوی است.

خاصیت Apriori: که اگر مجموعه‌ای از صفت-مقدارها، پر تکرار (frequent) باشد، هر زیر مجموعه‌ای از آن نیز پر تکرار است.

اگر $A = \{x_1, x_2, x_3, x_4\}$ و $AB = \{x_1, x_2, x_3\}$ پر تکرار باشد، لذا $Supp(AB)$ آنها باید بزرگ است، بنابراین نسبت آنها به عنوان $Conf$ تعریف می‌شود، و مطلوب است به یک نزدیک باشد.

▶ 8 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

8

مثال

Transaction	Items
t_1	Bread,Jelly,PeanutButter
t_2	Bread,PeanutButter
t_3	Bread,Milk,PeanutButter
t_4	Beer,Bread
t_5	Beer,Milk

$A \Rightarrow B$	Supp	Conf
$\text{Bread} \Rightarrow \text{PeanutButter}$	60%	75%
$\text{PeanutButter} \Rightarrow \text{Bread}$	60%	100%
$\text{Beer} \Rightarrow \text{Bread}$	20%	50%
$\text{PeanutButter} \Rightarrow \text{Jelly}$	20%	33.3%
$\text{Jelly} \Rightarrow \text{PeanutButter}$	20%	100%
$\text{Jelly} \Rightarrow \text{Milk}$	0%	0%

استخراج قواعد هم باش (Mining Association Rules)

❖ از آنجا که در مسائل عملی، تعداد ترکیبیهای ممکن برای تالی قواعد هم باشی می تواند بسیار زیاد باشد(از دید ریاضی)، لذا روش‌های عادی تولید قواعد برای تولید این قواعد استفاده نمی شود!

❖ الگوریتم **Apriori**، برای استخراج قواعد هم باشی استفاده می شود.

الگوریتم Apriori

❖ توسط Agrawal و همکارانش در ۱۹۹۳ مطرح شد.

R.Agrawal,T. Imielinski and A. Swami, "Mining association rules between sets of items in massive databases", Proc.ACm SIGMOD Conf. on Management of Data, Washington DC, USA, 1993, pp. 207-216.

❖ این الگوریتم Item Sets را تولید می کند.

▶ (مجموعه اقلام): ترکیباتی از صفت-مقدارها که شرط شمول (coverage) Item Sets برآورده کنند.

▶ شرط شمول (coverage): دامنه پشتیبانی (Support) توسط یک مجموعه اقلام (پشتیبانی کمینه یا min_supp)

▶ در این الگوریتم ترکیب‌های مختلف صفت-مقدار تولید می‌شود و آنها باید شرط شمول را برآورده نمی‌سازند، کنار گذارده می‌شوند. لذا این الگوریتم در زمانی معقول قواعد همباشی را تولید می کند.

▶ خاصیت Aperiori استفاده می‌شود.

▶ 11 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

11

گام‌های الگوریتم Apriori

۱. تولید مجموعه های اقلام (Item Sets): تولید ترکیباتی که شرط شمول را برآورده می سازند. (Item Sets برای support)

۲. تولید قواعد همبash به کمک مجموعه اقلام تولیدشده در گام ۱: در این مرحله نه تنها قواعدی که میزان اعتماد بالایی دارند، مورد توجه قرار می گیرند بلکه یافتا زیر نیز برای آنها برقرار باشد:

$$\text{Conf}(A \Rightarrow B) - \text{Supp}(B) > d$$

or

$$\text{Supp}(A \rightarrow B) - \text{Supp}(A) \cdot \text{Supp}(B) > k$$

که d و k ثابت های مناسبی هستند.

▶ 12 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

12

مثال

* مثال مربوط به پذیرش آگهی بیمه عمر را در نظر بگیرید:

فرض: یک مجموعه داده از رفтар مشتریان در ارتباط با خرید چند Item در اختیار داریم.

هدف: بر اساس رفtar خرید مشتریها قواعدی بیان شود که نحوه ای ارتباط اقلام خریداری شده با هم مشخص گردد.

● A Subset of the Credit Card Promotion Database

Magazine Promotion	Watch Promotion	Life Insurance Promotion	Credit Card Insurance	Sex
Yes	No	No	No	Male
Yes	Yes	Yes	No	Female
No	No	No	No	Male
Yes	Yes	Yes	Yes	Male
Yes	No	Yes	No	Female
No	No	No	No	Female
Yes	No	Yes	Yes	Male
No	Yes	No	No	Male
Yes	No	No	No	Male
Yes	Yes	Yes	No	Female

▶ 13 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

13

گام 1 : تولید Item Sets

.a ابتدا Single Item Sets (یعنی انواع ترکیب یک صفت-مقدار) را مدنظر قرار می‌دهیم. اولین عبارت است از Magazine Promotion attribute که دارای دو مقدار Yes, No است. بر اساس این صفت دو Item Set می‌توان تصور کرد که تعداد تکرار هر یک به قرار ذیل است:

Single Item Sets	No. of Item	accept
Magazine pro = Yes	7	✓
Magazine pro = No	3	✗

چون مورد دوم کم است (یعنی شرط شمول حداقل 40٪ را نمی‌پوشاند)، لذا نادیده گرفته می‌شود.

▶ 14 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

14

جدول Single Item Set به صورت زیر به دست خواهد آمد:

● Single-Item Sets

Single-Item Sets	Number of Items
Magazine Promotion = Yes	7
Watch Promotion = Yes	4
Watch Promotion = No	6
Life Insurance Promotion = Yes	5
Life Insurance Promotion = No	5
Credit Card Insurance = No	8
Sex = Male	6
Sex = Female	4

: Two Item Set برای به دست آوردن Single Item Set .b

از دو مورد اول جدول صفحه قبل شروع می کنیم:

Magazine promotion = Yes , Watch Promotion=Yes

dataset 3 مورد از 10 مورد

چون زیر 40٪ است، آنرا نمی پذیریم.

دو مورد بعدی را بررسی می کنیم:

Magazine Promotion = Yes , Watch Promotion = No

dataset 4 مورد از 10 مورد

شرط شمول (40٪) را براورده می کند، لذا پذیرفته می شود.

بنابراین:

ادامه روش فوق منجر به تولید 11 مورد **two Item Set** زیر می شود:

● **Two-Item Sets**

Two-Item Sets	Number of Items
Magazine Promotion = Yes & Watch Promotion = No	4
Magazine Promotion = Yes & Life Insurance Promotion = Yes	5
Magazine Promotion = Yes & Credit Card Insurance = No	5
Magazine Promotion = Yes & Sex = Male	4
Watch Promotion = No & Life Insurance Promotion = No	4
Watch Promotion = No & Credit Card Insurance = No	5
Watch Promotion = No & Sex = Male	4
Life Insurance Promotion = No & Credit Card Insurance = No	5
Life Insurance Promotion = No & Sex = Male	4
Credit Card Insurance = No & Sex = Male	4
Credit Card Insurance = No & Sex = Female	4

▶ 17 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

17

ترکیب **tree Item Set** برای به دست آوردن **two Item Set** .C از جدول **two Item Sets** هر دو درایه که یک صفت مشترک دارند را با هم در نظر می گیریم پس دو تای اول را در نظر می گیریم که ترکیب آنها یک **three Item Set** خواهد بود:

Magazine Promotion=Yes & Watch Promotion=No &
Life Insurance Promotion=Yes

data set از 1 مورد

آن را نمی پذیریم.

▶ 18 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

18

تنهایی که حداقل 4 مورد از موارد **data set** اولیه را پوشاند، از جدول **2-Item Set** اخذ شده باشد، مورد زیر است:

Watch Promotion =No & Life Insurance Promotion =No & Credit Card Insurance = No

data set 4

این روند را برای مسائلی که می توانند **5-Item Set**, **4-Item Set**, ... داشته باشند، ادامه می دهیم.

▶ 19 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

19

گام 2 : تولید قواعد هم باش

از جدول **2-Item Set** قوانینی که شرط حداقل اعتماد (مثلا 80%) را برآورد می کنند، استخراج می کنیم. a

مثلا اولین مورد جدول **2-Item Set** را در نظر می گیریم:

IF Magazine Promotion=Yes

THEN Watch Promotion =No

$$P(\text{Watch Pro.} | \text{Magazine Pro.}) = 3 \div 7 < 80\%$$

پس آنرا نادیده می گیریم!

البته عکس قاعده بالا هم بررسی می کنیم.

▶ 20 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

20

با بررسی سایر موارد می بینیم که دو قاعده زیر از روی 2-ItemSets قابل بیان هستند که حداقل ۸۰% confidence را برآورده می کنند.

IF Magazine Promotion =Yes

THEN Life Insurance Promotion =Yes

(5/7) => confidence

IF Life Insurance Promotion =Yes

THEN Magazine Promotion =Yes

(5/5) => confidence

از جدول 3-قواعدی که ترکیب سه تابی از صفت-مقدار را می توانند شامل شوند، به دست می آوریم و آنها یک شرط حداقل confidence را برآورد می سازند، نگاه می داریم. در این مثال فقط یک 3-Item Set داریم. از روی آن برخی قاعده ها که سه تای آنها در ذیل ذکر می شود به دست می آوریم.

IF Watch Promotion = No & Life Insurance= No

THEN Credit Card Ins. = No (4/4) => confidence

IF Watch Promotion = No

THEN Life Insurance= No & Credit Card Ins. = No (4/6)

IF Credit Card Ins. = No

THEN Watch Promotion = No & Life Insurance= No (4/8)

مثال

Transaction	Items
t_1	Blouse
t_2	Shoes,Skirt,TShirt
t_3	Jeans,TShirt
t_4	Jeans,Shoes,TShirt
t_5	Jeans,Shorts
t_6	Shoes,TShirt
t_7	Jeans,Skirt
t_8	Jeans,Shoes,Shorts,TShirt
t_9	Jeans
t_{10}	Jeans,Shoes,TShirt
t_{11}	TShirt
t_{12}	Blouse,Jeans,Shoes,Skirt,TShirt
t_{13}	Jeans,Shoes,Shorts,TShirt
t_{14}	Shoes,Skirt,TShirt
t_{15}	Jeans,TShirt
t_{16}	Skirt,TShirt
t_{17}	Blouse,Jeans,Skirt
t_{18}	Jeans,Shoes,Shorts,TShirt
t_{19}	Jeans
t_{20}	Jeans,Shoes,Shorts,TShirt

▶ 23 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

23

Scan	Candidates	Large Itemsets
1	{Blouse},{Jeans},{Shoes}, {Shorts},{Skirt},{TShirt}	{Jeans},{Shoes},{Shorts} {Skirt},{Tshirt}
2	{Jeans,Shoes},{Jeans,Shorts},{Jeans,Skirt}, {Jeans,TShirt},{Shoes,Shorts},{Shoes,Skirt}, {Shoes,TShirt},{Shorts,Skirt},{Shorts,TShirt}, {Skirt,TShirt}	{Jeans,Shoes},{Jeans,Shorts}, {Jeans,TShirt},{Shoes,Shorts}, {Shoes,TShirt},{Shorts,TShirt}, {Skirt,TShirt}
3	{Jeans,Shoes,Shorts},{Jeans,Shoes,TShirt}, {Jeans,Shorts,TShirt},{Jeans,Skirt,TShirt}, {Shoes,Shorts,TShirt},{Shoes,Skirt,TShirt}, {Shorts,Skirt,TShirt}	{Jeans,Shoes,Shorts}, {Jeans,Shoes,TShirt}, {Jeans,Shorts,TShirt}, {Shoes,Shorts,TShirt}
4	{Jeans,Shoes,Shorts,TShirt}	{Jeans,Shoes,Shorts,TShirt}
5	Ø	Ø

▶ 24 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

24

مزایا و معایب الگوریتم Apriori

مزایا

- ❖ اجرای آن راحت است.
- ❖ به طور موازی قابل اجراست.

معایب

- ❖ زمانبر است.
- ❖ اگر attribute‌ها، از نوع غیرباينري باشند، کارايی زمانی آن کمتر است.

ملاحظات کلی در مورد قواعد هم‌باشی

- ❖ در تعبیر قواعد هم‌باشی باید دقت کرد زیرا بسیاری از روابطی که این قوانین کشف می‌کنند، بدیهی (trivial) هستند.
- ❖ مثال: از ۱۰۰۰۰ خرید، ۷۰٪ شامل خرید شیر و نیز ۵۰٪ شامل نان می‌شوند. پس قاعدة زیر محتتماً بدست می‌آید:

IF customers purchase Milk **THEN** they also Purchase Bread.

فرض کنید **confidence** این قاعده هم‌باشی بالای ۴۰٪ باشد. اما این قاعده برای ما ارزش زیادی ندارد زیرا بیشتر مشتریها هر دو را می‌خرند. بنابراین، این قاعده اطلاع جدیدی به ما نمی‌دهد. (در مورد اینکه خرید نان با شیر را اگر آگهی کنیم به نفع ماست، زیرا این را می‌دانیم.)

دو نوع ارتباط که در قواعد هم باشی مورد توجه هستند

۱. به آن قواعد هم باشی علاقه مندیم که ارتفاق فروش یک محصول خاص را ناشی از هم باشی با یک یا چند محصول دیگر بیان دارد.
به این ترتیب، با این اطلاعات می توانیم فروش آنرا به عنوان نتیجه هم باشی بالا ببریم.

۲. به آن قواعد هم باشی علاقه مندیم که برای یک هم باشی خاص مقدار **confidence** کمتری نسبت به آنچه انتظار داشتیم را بیان کنند.
به این ترتیب می توانیم نتیجه بگیریم که محصولاتی که در آن قاعده هم باش به کار رفته اند با هم رقابت می کنند.

▶ 27 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

27

یک نکته

هر قاعده قوی، برای نمایش لزوماً مناسب نیست. بلکه قواعدی که در یافtar ذکر شده هم صدق کنند، مورد توجه خواهند بود.

مثال:

دبیرستانی با 5000 دانش آموز و یک فروشنده صبحانه را در نظر بگیرید.
۶۰٪ دانش آموزان بسکتبال بازی می کنند.
۷۵٪ دانش آموزان صبحانه می خورند.
۴۰٪ دانش آموزان بسکتبال بازی می کنند و صبحانه می خورند.
شرط شمول = ۴٪، شرط حداقل میزان اعتماد = ۶۰٪
قاعده روپرو تولید می شود: (play basketball) → (eat cereal)
چون شرط شمول را برا آورده می کند، و نیز $Conf = 2000/3000 = 0.66$

اما:

$$Supp(A,B) - Supp(A).Supp(B) = 0.4 - 0.6 * 0.75 = -0.05 < 0 \quad !!$$

▶ 28 - Association Rules

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2007

28