

داده کاوی

میثم مدنی

مباحثی در کاوش
الگو

مردم کتاب

کاوش الگو در فضای چند
بعدي

کاوش الگوی تکواری مقید

کاوش الگوهای بسیار بزرگ و
چند بعدي

کاوش الگوهای تخمینی و
فشرده

دوره درسی ترم ۲ سال ۱۳۹۱

داده کاوی



میثم مدنی

دانشگاه صنعتی شریف

فصل هفتم،

<http://www.madani.pro>

Email: me@madani.pro

چشم انداز

داده کاوی

میشم مدنی

مباحثی در کاوش
الگو

مروارثی

کاوش الگو در فضای چند
بعدي

کاوش الگوی تکواری مفید

کاوش الگوهای بسیار بزرگ و
چند بعدي

کاوش الگوهای تخمینی و
فشرده

۱ مباحثی در کاوش الگو

داده کاوی

میشم مدنی

مباحثی در کاوش
الگو

مرور کلی

کاوش الگو در فضای چند
بعدي

کاوش الگوی تکراری مفید

کاوش الگوهای بسیار بزرگ و
چند بعدي

کاوش الگوهای تخمینی و
فشرده

مرور کلی

کاوش قواعد شراکت چند سطحی

- ۱ قاعده یک جد قاعده دوم است اگر قاعده دوم با جایگزینی مواردش با اجدادشان به قاعده یک تبدیل شود.
- ۲ به این صورت یک سلسله مراتب روی قواعد شراکت ایجاد خواهد شد.
- ۳ یک قاعده زائد است اگر با توجه به اجدادش پشتیبان و اطمینانش به مقدار مورد انتظار نزدیک باشد.
- ۴ مثلا ما در laptop ۴ برند داشتیم و انتظار می رود پشتیبان آن یک چهارم شود.

$$\text{buys}(X, \text{"laptop computer"}) \Rightarrow \text{buys}(X, \text{"HP printer"})$$

$$[\text{support} = 8\%, \text{confidence} = 70\%]$$

$$\text{buys}(X, \text{"Dell laptop computer"}) \Rightarrow \text{buys}(X, \text{"HP printer"})$$

$$[\text{support} = 2\%, \text{confidence} = 72\%]$$

کاوش قواعد شراکت چند بعدی

داده کاوی

میشم مدنی

مباحثی در کاوش
الگو

روز کنی

کاوش الگو در فضای چند
بعدی

کاوش الگوی تکراری مفید

کاوش الگوهای بسیار بزرگ و
چند بعدیکاوش الگوهای تخمینی و
فشرده

۱ قواعدی که دو یا بیشتر بعد دارند را **قاعده شراکت چند بعدی** می گویند.

۲ قواعد شراکتی چند بعدی که پیش بینی تکراری دارند را **هیبرید** می گویند.

$$age(X, "20 \dots 29") \wedge occupation(X, "student") \Rightarrow buys(X, "laptop").$$

$$age(X, "20 \dots 29") \wedge buys(X, "laptop") \Rightarrow buys(X, "HP printer").$$

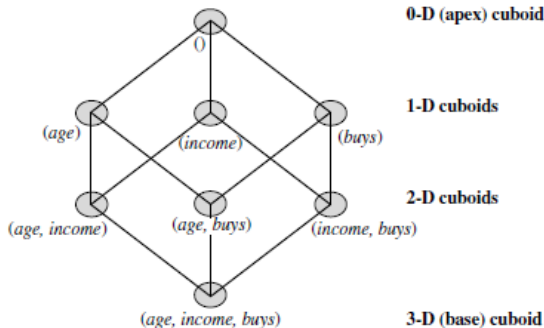
روش های کشف قواعد شراکت را می توان به دو دسته تقسیم کرد.

۱ کمی

۲ اسمی

کاوش قواعد شراکت کمی

۱ بر مبنای مکعب داده



۲ بر مبنای خوشه (بالا به پایین، پایین به بالا)

۳ روش آماری (تشخیص استثناها)

$$\text{sex} = \text{female} \Rightarrow \text{meanwage} = \$7.90/\text{hr} \text{ (overall_mean_wage} = \$9.02/\text{hr)}.$$

$$\text{population_subset} \Rightarrow \text{mean_of_values_for_the_subset},$$

کاوش قواعد شراکت کمی

- ۱ یک الگوی کمیاب الگویی است که تکراری نباشد!
- ۲ فرض کنید دنبال الگویی می گردیم که شامل حداقل یک مورد با ارزشی بیش از ۵۰۰ دلار داشته باشد.
- ۳ کاوش کارا از این الگو کمی مشکل است.
- ۴ باید چاره ای اندیشید

تعریف: اگر مجموعه های تکراری X, Y به ندرت با هم رخ دهند

$\sup(X \cap Y) < \sup(X) \times \sup(Y)$ آنگاه این دو مجموعه را همبسته منفی می گویند. و $X \cap Y$ را الگوی همبسته منفی می گویند.

اگر $\sup(X \cap Y) \ll \sup(X) \times \sup(Y)$ دو مجموعه به طور قویا منفی همبسته اند.

فرض کنید یک فروشگاه ابزار خیاطی داریم

از هر بسته A, B ۱۰۰ بسته فروخته باشیم. و تنها یک فروش از هر دو داشته باشیم.

$$\sup(X \cap Y) = 0.005 \ll \sup(X) \times \sup(Y) = 0.25$$

فرض کنید فروش 10^6 فاکتور داریم.

$$\sup(X \cap Y) = 10^{-6} \gg \sup(X) \times \sup(Y) = 10^{-8}!!!$$

کاوش قواعد شراکت کمی

مشکل تراکنش های خالی را این گونه م میتوان حل کرد NullTransaction
تعریف:

مجموعه های تکراری X, Y را به طور منفی همبسته می گویند اگر
 $(P(X|Y) + P(Y|X))/2 < \epsilon$ که ϵ آستانه الگوی منفی است.

فرض کنید یک فروشگاه ابزار خیاطی داریم

از هر بسته A, B ۱۰۰ بسته فروخته باشیم. و تنها یک فروش از هر دو داشته باشیم.

$$\epsilon = 0.02, \text{minsup} = 0.0001$$

$$(P(X|Y) + P(Y|X))/2 = (0.01 + 0.01)/2 = 0.01 < 0.02$$

فرض کنید فروش 10^6 فاکتور داریم.

$$(P(X|Y) + P(Y|X))/2 = (0.01 + 0.01)/2 = 0.01 < 0.02$$

داده کاوی

میشم مدنی

مباحثی در کاوش
الگو

مروز کئی

کاوش الگو در فضای چند
بعدي

کاوش الگوی تکراری مقید

کاوش الگوهای بسیار بزرگ و
چند بعدي

کاوش الگوهای تخمینی و
فشرده

کاوش الگوی تکراری مقید

- فرآیند داده کاوی ممکن است هزاران قاعده از داده ها کشف کند که بسیاری از آنها بی ربط یا بی استفاده برای کاربر است.
- کاربران اغلب دید خوبی از الگوهای مناسب و شکل الگوی مورد نظر را دارند.
- همچنین شرایطی ممکن است برای قواعد مورد نظر وجود داشته باشد.
- یک روش مناسب تعیین شرایط و خصوصیات به عنوان قید است که فضای جستجو را کاهش می دهد.
- به این فرآیند کاوش مقید می گویند **Constraint – Based Mininn** قیود می تواند شامل موارد زیر باشند

۱ **قیدهای تعیین نوع دانش شراکت، همبستگی، خوشه بندی و طبقه بندی**

۲ **قیود داده تعیین نوع داده ها**

۳ **قید بعد، سطح**

۴ **قید تمایل تعیین آستانه هایی مانند پشتیبان، اطمینان و همبستگی**

۵ **قید قاعده تعیین نوع، شکل یا شرایط قاعده ای است که می خواهد کاوش شود**

این ها به عنوان ابر قاعده توصیف خواهند شد.

کاوش با ابرقاعده قواعد شراکت

- ابرقاعده این امکان را می دهد تا شکل های از قواعد را که برای ما جلب هستند پیدا کنیم
- فرض کنید که دسترسی به مشخصات خریداران داریم
- می خواهیم الگوهایی را کشف کنیم که خصیصه خریدار را به کالاهای خریداری شده وی مرتبط کند.
- به طور تخصصی تر می خواهیم بدانیم کدام جفت خصیصه ها منجر به خرید آفیس می شود؟

$$P_1(X, Y) \wedge P_2(X, W) \Rightarrow \text{buys}(X, \text{"office software"});$$

$$\text{age}(X, \text{"30..39"}) \wedge \text{income}(X, \text{"41K..60K"}) \Rightarrow \text{buys}(X, \text{"office software"}).$$

$$P_1 \wedge P_2 \wedge \dots \wedge P_l \Rightarrow Q_1 \wedge Q_2 \wedge \dots \wedge Q_r,)$$

کاوش با ابرقاعده قواعد شراکت

قید قاعده ارتباط بین متغیرها، مقادیر اولیه متغیرها و قیود در توابع تجمعی در قواعد اکتشافی است.

مثال: فرض کنید که دسترسی به مشخصات خریداران داریم

- *item(item_ID, item_name, description, category, price)*
- *sales(transaction_ID, day, month, year, store_ID, city)*
- *trans_item(item_ID, transaction_ID)*

■ یک حالت این است که در مورد الگوهای ارزان جستجو کنیم (با مجموع قیمت

$$\text{sum}(I.\text{price}) < 10 \text{ (کمتر از ده دلار)}$$

■ الگوهای گران قیمت $\text{sum}(I.\text{price}) > 50$

هرس فضای الگوها با قیود هرس کننده

۵ نوع قید داریم

- ۱ **پادیکنوا** $\text{sum}(I.\text{price}) < 10$ را در نظر بگیرید k امین تکرار را در الگوریتم استقرایی در نظر بگیرید اگر مجموعه ای این قید را ارضا نکند دیگر نیازی نیست ادامه دهیم. این شاخه را هرس می کنیم. در واقع اگر یک مجموعه این قید را برآورده نکند مجموعه های شامل آن هم آن قید را ارضا نمی کنند.
 - ۲ **یکنوا** مشابه قسمت قبل اما برعکس! $\text{sum}(I.\text{price}) > 50$ در واقع اگر یک مجموعه این قید را برآورده نکند مجموعه های مشمول در آن هم آن قید را ارضا نمی کنند.
 - ۳ **فشرده** می توانیم مستقیماً با قید تمام فضای مربوط به آن را مشخص کنیم. $\min(I.\text{price}) > 50$
 - ۴ **تبدیل پذیر** هیچ یک از موارد قبل نیست اما اگر آیتیم ها به نحو مناسبی مرتب شوند یکی از موارد قبل خواهد شد. مثلاً $\text{avg}(I.\text{price}) > 50$ هیچ یک از موارد زیر نیست اما اگر از کوچک به بزرگ وارد کنیم پادیکنوا خواهد شد.
 - ۵ **تبدیل ناپذیر** هیچ کدام!
- خبر خوب اینکه اکثر قیود از چهار نوع اول هستند.

مباحثی در کاوش الگو

مروارذ

کاوش الگو در فضای چند

بعدي

کاوش الگوی تکراری مقید

کاوش الگوهای بسیار بزرگ و

چند بعدي

کاوش الگوهای تخمینی و

فشرده

<i>Constraint</i>	<i>Antimonotonic</i>	<i>Monotonic</i>	<i>Succinct</i>
$v \in S$	no	yes	yes
$S \supseteq V$	no	yes	yes
$S \subseteq V$	yes	no	yes
$\min(S) \leq v$	no	yes	yes
$\min(S) \geq v$	yes	no	yes
$\max(S) \leq v$	yes	no	yes
$\max(S) \geq v$	no	yes	yes
$\text{count}(S) \leq v$	yes	no	weakly
$\text{count}(S) \geq v$	no	yes	weakly
$\text{sum}(S) \leq v (\forall a \in S, a \geq 0)$	yes	no	no
$\text{sum}(S) \geq v (\forall a \in S, a \geq 0)$	no	yes	no
$\text{range}(S) \leq v$	yes	no	no
$\text{range}(S) \geq v$	no	yes	no
$\text{avg}(S) \theta v, \theta \in \{\leq, \geq\}$	convertible	convertible	no
$\text{support}(S) \geq \xi$	yes	no	no
$\text{support}(S) \leq \xi$	no	yes	no
$\text{all_confidence}(S) \geq \xi$	yes	no	no
$\text{all_confidence}(S) \leq \xi$	no	yes	no

داده کاوی

میشم مدنی

مباحثی در کاوش
الگو

مروز کئی

کاوش الگو در فضای چند
بعدی

کاوش الگوی تکراری مفید

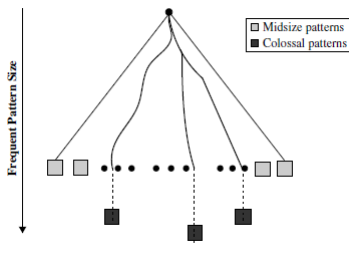
کاوش الگوهای بسیار بزرگ و
چند بعدی

کاوش الگوهای تخمینی و
فشرده

کاوش الگوهای بسیار بزرگ و چند بعدی

کاوش الگوهای بسیار بزرگ با آمیختگی الگو

فرض کنید داده های زیر را داریم



row/col	1	2	3	4	...	38	39
1	2	3	4	5	...	39	40
2	1	3	4	5	...	39	40
3	1	2	4	5	...	39	40
4	1	2	3	5	...	39	40
5	1	2	3	4	...	39	40
...
39	1	2	3	4	...	38	40
40	1	2	3	4	...	38	39
41	41	42	43	44	...	78	79
42	41	42	43	44	...	78	79
...
60	41	42	43	44	...	78	79

در اغلب روش های جستجو از رشد تعداد آیتم ها استفاده می کنیم.
این روش ها نمی توانند از پس تولید حجم عظیم الگوهای با اندازه میانه بر بیایند.
به روش جدیدی نیازی داریم.

۱ فاصله دو الگوی تکراری

$$Pat_Dist(P_1, P_2) = 1 - \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|}.$$

۲ مثلاً

$$T(P_1) = \{t_1, t_2, t_3, t_4, t_5\}$$

$$T(P_2) = \{t_1, t_2, t_3, t_4, t_6\},$$

$$Pat_Dist(P_1, P_2) = 1 - \frac{4}{6} = \frac{1}{3}.$$

۳ یک الگو P با الگوی دیگر P' ، λ - پوشیده می شود اگر $O(P) \subseteq O(P')$ و $Dist(P, P') \leq \lambda$ که منظور از $O(P)$ مجموعه موارد P است

۴ تخمین

$$\delta \geq Pat_Dist(P, P_r) = 1 - \frac{|T(P_r)|}{|T(P)|} \geq 1 - \frac{k}{min_sup}.$$

$$\delta \geq Pat_Dist(P, P_r) = 1 - \frac{|T(P_r)|}{|T(P)|} \geq 1 - \frac{k}{min_sup}.$$

k is support P_r .

داده کاوی

میشم مدنی

مباحثی در کاوش الگو

مردم کل

کاوش الگو در فضای چند بعدی

کاوش الگوی تکراری مفید

کاوش الگوهای بسیار بزرگ و چند بعدی

کاوش الگوهای تخمینی و فشرده