

Tour Version



# Data Mining

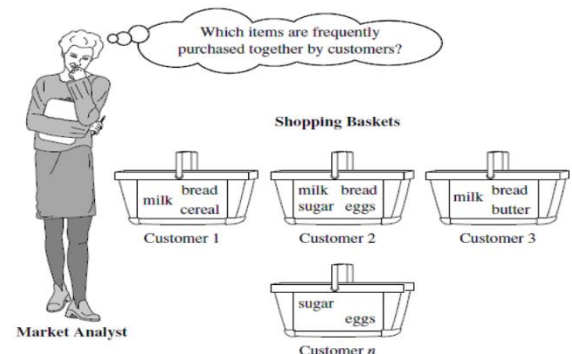
## Tour 3

# Mining Frequent Pattern

# Mining Frequent Pattern, Association & Correlation

## Frequent Pattern

- a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set
- A Big role for mining association and correlation
- بس ليه بدور على الأشكال دي في الداتا أيه الحافز؟؟
- Motivation: Finding inherent regularities in data
- لما تلاقى ان السؤال بتاعك شبه الأسئلة الجاية يبقى فعلا محتاج الموضوع ده



- What products were often purchased together?— Beer and diapers?!
- What are the subsequent purchases after buying a PC?
- What kinds of DNA are sensitive to this new drug?
- Can we automatically classify web documents?
- و بيخش في مجموعة من التطبيقات اللي في الحياة
- Applications: **Basket data analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, and DNA sequence analysis**
- Frequent Pattern are itemsets that appear frequently in a data set (e.g. Transaction record)
- Items that are frequently associated (e.g purchased) together can be represented as association rules. **Association Rule like =>**
- Computer → antivirus\_SW [Support = 2% , Confidence =60%]
- معنى كده ان شراء الكمبيوتر بيأثر على شراء الانتى فيروس باجمالى 2% من بياعة المنتجات و نسبة 60% من اللي اشتروا كمبيوترات اشتروا كمان انتى فيروس بس مش العكس؛؛ احتمال شرطي أفكر دكتور شريف
- **Support** and **Confidence** are measures of rule interestingness
- و يعتبروا هم اللي بقارن بينهم التريشولد بتوعى و بختبر بينهم مدى القرار على المنتجات و هو اللي بيعرفنى الباترن
- طبعا فيه تريشولد أكثر من دول بس المحاضرة بنتكلم عن 2
- 2% Support means 2% of Transactions Show that computers and antivirus\_SW are bought Together
- 60% Confidence means 60 % of customers who bought a computer also bought antivirus\_SW

- Basics About **Association Rules**:

- If frequency of itemset I satisfies min\_support count then I is a frequent itemset
- $\text{Support}(x) \geq \text{min\_support} \Rightarrow x$  is frequent item set
- If a rule satisfies min\_support and min\_confidence thresholds, it is said to be strong
  - problem of mining association rules reduced to mining frequent itemsets
- $\text{Support}(X) \geq \text{min\_support} \ \&\& \ \text{Conf}(X) \geq \text{min\_confidence} \Rightarrow$  strong Association
- Association rules mining becomes a two-step process:
  - Find all frequent itemsets that occur at least as frequently as a predetermined min\_support count
  - Generate strong association rules from the frequent itemsets that satisfy min\_support and min\_confidence

- Itemset  $X = \{x_1, \dots, x_k\}$       ex:  $X = \{A, B, C, D, E, F\}$
- Find all the rules  $X \rightarrow Y$  with minimum support and confidence
  - support, s, probability that a transaction contains  $X \cup Y$
  - confidence, c, conditional probability that a transaction having X also contains Y

$$\text{support } X \rightarrow Y = P(X \cup Y) = \frac{n(X \cup Y)}{N}$$

$$\text{confidence } (X \rightarrow Y) = P(Y|X) = \frac{n(X \cup Y)}{n(X)}$$

Ex: Let min\_Sup. = 50%, min\_conf. = 50%

Frequent Patterns:

$\{A:3, B:3, D:4, E:3, AD:3\}$

Association rules:

$A \rightarrow D$  (60%, 100%)

$D \rightarrow A$  (60%, 75%)

$$\text{conf } (A \rightarrow D) = \frac{3}{3} = 100 \%$$

$$\text{conf } (D \rightarrow A) = \frac{3}{4} = 75 \%$$

Transaction-id	Items bought
10	A, B, D
20	A, C, D
30	A, D, E
40	B, E, F
50	B, C, D, E, F

o9

- نبدأ نحل المسألة دي ب 2 الجوريزم و هختبر الحلين بالبایشون كمان الله المستعان
  - Apriori
  - FP-Growth

TID	List of items
T100	11, 12, 15
T200	12, 14
T300	12, 13
T400	11, 12, 14
T500	11, 13
T600	12, 13
T700	11, 13
T800	11, 12, 13, 15
T900	11, 12, 13
T1000	11, 12

## Apriori Algorithm

أول حاجة محتاجها ان نشوف! باهم

Frequent Itemsets?

و دي بدسترقاعته طريقته! min-Support threshold وبعد كده  
نشوف Association Rules اللي من خلالها مبيشوف

min-Confidence  
Threshold

و محتار كل Pattern تبع

واحد ابدأ ب items وعددهم

min-Support = 2

min-Confidence = 70%

C<sub>1</sub>

Item Set	# Support
{I <sub>1</sub> }	7 > 2 ✓
{I <sub>2</sub> }	8 ✓
{I <sub>3</sub> }	6 ✓
{I <sub>4</sub> }	2 ✓
{I <sub>5</sub> }	2 ✓

منفصلة  
ال min-Support  
⇒L<sub>1</sub> يبقي دول كالم items

Item Set	# Support
{I <sub>1</sub> }	7
{I <sub>2</sub> }	8
{I <sub>3</sub> }	6
{I <sub>4</sub> }	2
{I <sub>5</sub> }	2



من أجل بقى Compose الـ Set  
ونعد هم حصلوا كما مر

C<sub>2</sub>

Itemset	# Support	
{I <sub>1</sub> , I <sub>2</sub> }	5 > 2	✓
{I <sub>1</sub> , I <sub>3</sub> }	4	✓
{I <sub>1</sub> , I <sub>4</sub> }	1	✗
{I <sub>1</sub> , I <sub>5</sub> }	3	✓ min-Sup
{I <sub>2</sub> , I <sub>3</sub> }	4	✓ ⇒
{I <sub>2</sub> , I <sub>4</sub> }	2	✓
{I <sub>2</sub> , I <sub>5</sub> }	2	✓
{I <sub>3</sub> , I <sub>4</sub> }	0	✗
{I <sub>3</sub> , I <sub>5</sub> }	1	✗
{I <sub>4</sub> , I <sub>5</sub> }	0	✗

L<sub>2</sub>

Itemset	# Support
{I <sub>1</sub> , I <sub>2</sub> }	5
{I <sub>1</sub> , I <sub>3</sub> }	4
{I <sub>1</sub> , I <sub>5</sub> }	3
{I <sub>2</sub> , I <sub>3</sub> }	4
{I <sub>2</sub> , I <sub>4</sub> }	2
{I <sub>2</sub> , I <sub>5</sub> }	2

نعمل Compose ثاني مرة  
نأخذ بالنامت الـ Pruning  
لوفيه Subset ما يتحققش الـ min-Sup  
الـ Superset كانت متة محقق

C<sub>3</sub>

Itemset	# Support	
{I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> }	2 > 2	✓
{I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub> }	2	✓
<del>{I<sub>1</sub>, I<sub>2</sub>, I<sub>4</sub>}</del>		⇒

L<sub>3</sub>

Itemset	# Support
{I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> }	2
{I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub> }	2

وطبعا لو عملنا Compose مني متبقى حاجة لانه حصل  
Pruning الـ {I<sub>3</sub>, I<sub>5</sub>} هم Subset ما يتحققش

لوحالہ سوال عن مینہ م ال frequent itemset بیقی کل الی  
حققوا

frequent itemset

- { I<sub>1</sub> }
- { I<sub>2</sub> }
- { I<sub>3</sub> }
- { I<sub>4</sub> }
- { I<sub>5</sub> }
- { I<sub>1</sub>, I<sub>2</sub> }
- { I<sub>1</sub>, I<sub>5</sub> }
- { I<sub>1</sub>, I<sub>3</sub> }
- { I<sub>2</sub>, I<sub>3</sub> }
- { I<sub>2</sub>, I<sub>4</sub> }
- { I<sub>2</sub>, I<sub>5</sub> }
- { I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub> }
- { I<sub>1</sub>, I<sub>2</sub>, I<sub>5</sub> }

تمام ناضد مثلاً  
منهم عشان نحدد فينت  
Association Rules

ببدأ نطلع منه ال subset عشان نرشوف ال Pattern  
احتمال شرطی

Association Rules	Confidence
{ I <sub>1</sub> } → { I <sub>2</sub> , I <sub>5</sub> }	28% ✗
{ I <sub>2</sub> } → { I <sub>1</sub> , I <sub>5</sub> }	25% ✗
{ I <sub>5</sub> } → { I <sub>1</sub> , I <sub>2</sub> }	100% ✓
{ I <sub>1</sub> , I <sub>2</sub> } → { I <sub>5</sub> }	40% ✗
{ I <sub>1</sub> , I <sub>5</sub> } → { I <sub>2</sub> }	100% ✓
{ I <sub>2</sub> , I <sub>5</sub> } → { I <sub>1</sub> }	100% ✓

$$Conf = \frac{P(\{I_1, I_2, I_5\})}{P(I_1)}$$

$$= \frac{2}{7} = 28\%$$

دمكن

بيقی دول ال strong Association  
بيس احنا اخبرنا على frequent itemset واحدة بيقی

- أنا اسف خطي وحش جدا و حاسس بالموقف اللي انت فيه بس ربنا يستربس،، ده كله عشان ما بتقرأش المصادر يا عزيزي
- و بعد ما اتحلت منكم نحلها بايثون تمام و هنشوف ايه اللي فاضل من السترونج
- وده الليك بتاع الكود اقرا المصادر بقى



- [https://github.com/AhmedKhalil777/DataScience.Learning/blob/master/Extracting%20Frequent%20pattern/Frequent\\_pattern.ipynb](https://github.com/AhmedKhalil777/DataScience.Learning/blob/master/Extracting%20Frequent%20pattern/Frequent_pattern.ipynb)

	support	itemsets
0	0.7	(11)
1	0.8	(12)
2	0.6	(13)
3	0.2	(14)
4	0.2	(15)
5	0.5	(11, 12)
6	0.4	(13, 11)
7	0.2	(15, 11)
8	0.4	(13, 12)
9	0.2	(14, 12)
10	0.2	(15, 12)
11	0.2	(13, 11, 12)
12	0.2	(15, 11, 12)

	antecedents	confidence	consequents
0	(11)	0.714286	(12)
1	(15)	1.000000	(11)
2	(14)	1.000000	(12)
3	(15)	1.000000	(12)
4	(15, 11)	1.000000	(12)
5	(15, 12)	1.000000	(11)
6	(15)	1.000000	(11, 12)

## FP-Growth

- أولاً ما اتوقعش ان هي تيجي في امتحان بس للحظر بقول افهمها من المحاضرة
- اه هي تلخبط بس اللي فاهم ال
- Merge and Conquer
- هيفهمها كويس
- بس انا هعملها كود وحسب ال
- Time cost
- Big data بتاع الالجوريزمين و هتلاحظ فرق الوقت لو بتتعامل مع

```
[7] ▶ ML
%timeit -n 100 -r 10 apriori(df, min_support=0.6)
7.55 ms ± 1.31 ms per loop (mean ± std. dev. of 10 runs, 100 loops each)

[8] ▶ ML
%timeit -n 100 -r 10 fpgrowth(df, min_support=0.6)
3.36 ms ± 681 µs per loop (mean ± std. dev. of 10 runs, 100 loops each)

[26] ▶ ML 7.55/3.36
2.2470238095238098
```

- test Benchmarking بعد ال
- اتضح ان في داتا بسيطة زي المسألة ان سرعة الالجوريزم ده أسرع مرتين وربع و طبعاً هيفتخلف لو الداتا كبرت
- To avoid costly candidate generation
- Divide-and-conquer strategy:
  - Compress database representing frequent items into a frequent pattern tree (FPtree) – 2 passes over dataset
  - Divide compressed database (FP-tree) into conditional databases, then mine each for frequent itemsets – traverse through the FP-tree



دلوقة منقش على حل المسألة بطريقة الـ FP-Growth

نرى ما عملنا المرة اللي فاتت  
 حصلنا المرة دي بيبي صيرناهم  
 C<sub>1</sub>

ItemSet	# Support
{I <sub>1</sub> }	7
{I <sub>2</sub> }	8
{I <sub>3</sub> }	6
{I <sub>4</sub> }	2
{I <sub>5</sub> }	2

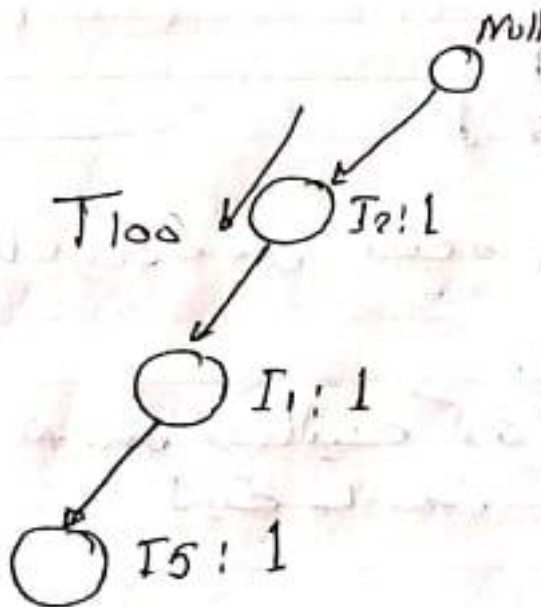
L<sub>1</sub>

ItemSet	# Support
{I <sub>2</sub> }	8
{I <sub>1</sub> }	7
{I <sub>3</sub> }	6
{I <sub>4</sub> }	2
{I <sub>5</sub> }	2

Transactional data example  
 N=10, min\_supp count=2

TID	List of items
T100	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>
T200	I <sub>2</sub> , I <sub>4</sub>
T300	I <sub>2</sub> , I <sub>3</sub>
T400	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>
T500	I <sub>1</sub> , I <sub>3</sub>
T600	I <sub>2</sub> , I <sub>3</sub>
T700	I <sub>1</sub> , I <sub>3</sub>
T800	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>
T900	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>
T1000	I <sub>1</sub> , I <sub>2</sub>

دلوقة هنعمل الـ FP-Tree وهنبداً بـ Null Set في الأول وبترتيب  
 العناصر كـ Node في I: Support



دلوقة هتتووع للعنود الأيسر

بتاع الـ Transactions ونكتب بالترتيب

اللي عملناه في الجدول اللي فوق

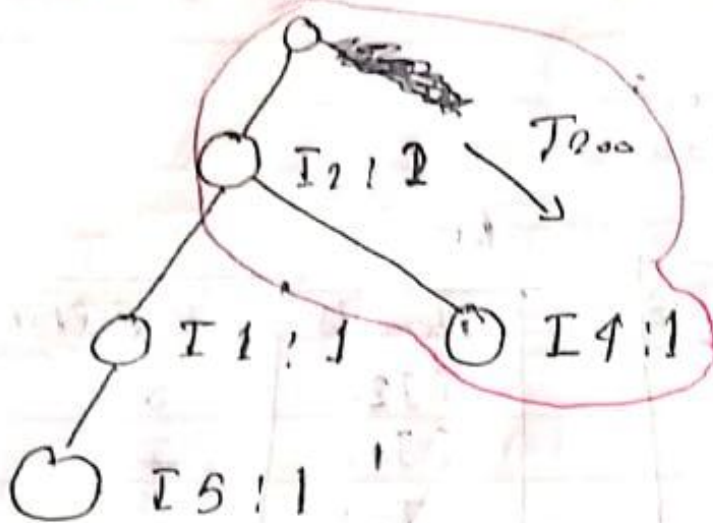
يعني اللي في الـ Tree دي T100

ديا بترتيب العنود L1

وحدة Support كل واحدة بـ 1 لأن دي زول Transactions

T200 Null

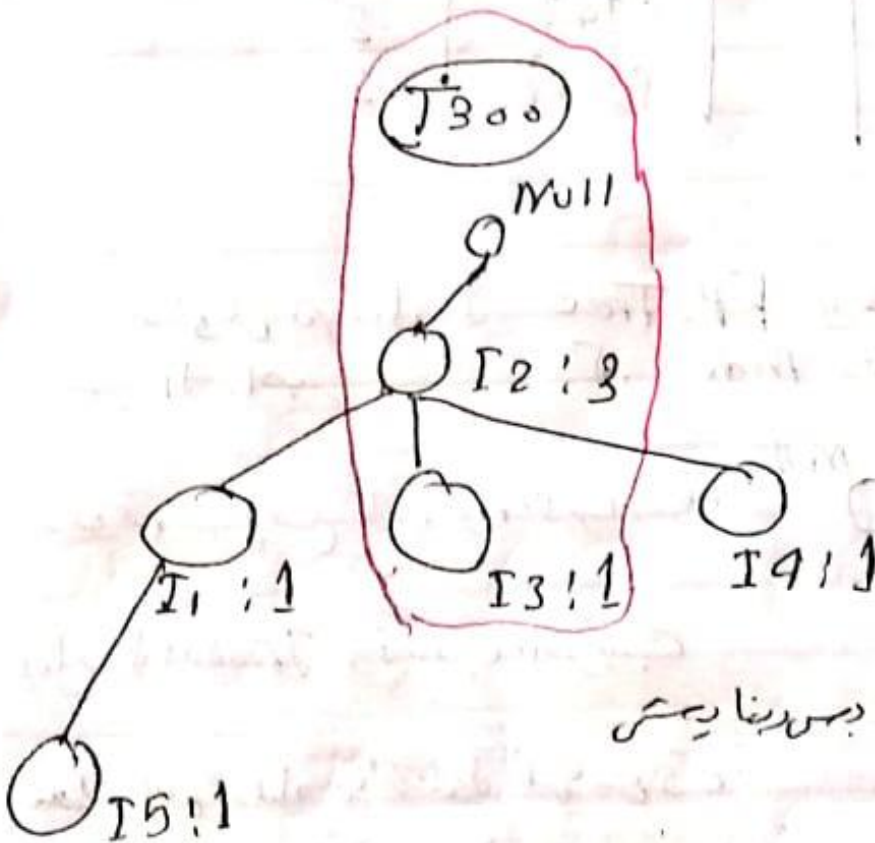
فردخل علی تاسیسات



دانشه شیت ولی I2  
مرتبیت لفظ

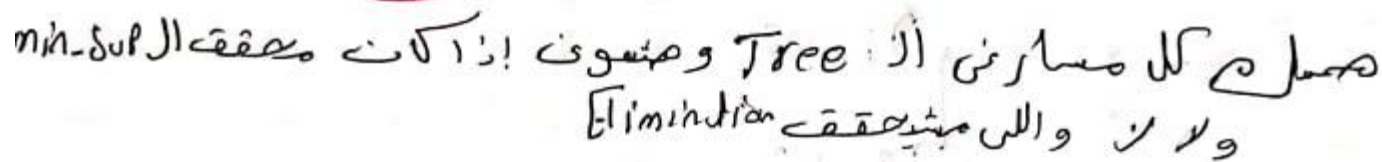
Transactional data example  
N=10, min\_supp count=2

TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2



انعارف است نهی ترهفک برهنا ریسر

مفصله شغالین که  
لغیت عارضی کی ال Transact



وہنا صیغہ ظہور مشککہ! ان ال I3 لہا ظہورت صی السبب! ان

وہمینے سے تاسی الجدول بیسی بطریقہ ال Condition 1 - fP

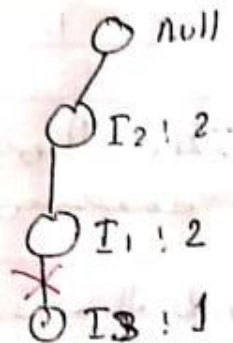


و سنستمر في ال Transaction ال فيها T5 مثلاً

اللى هم (T800 , T100)

T100	I <sub>1</sub> , I <sub>2</sub> , <u>I<sub>5</sub></u>
T800	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , <u>I<sub>5</sub></u>

منزل  
eliminate



تقاربه من حذف  
 $I_3 < \text{min-sup}$

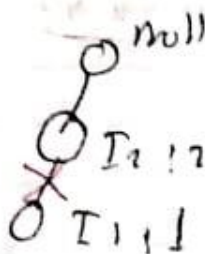
نتجى على تاني واحدة حذف ال Transaction ال فيها T5

I4

و براضى ال Transaction ال فيها T9 منزل فيهم زي T8

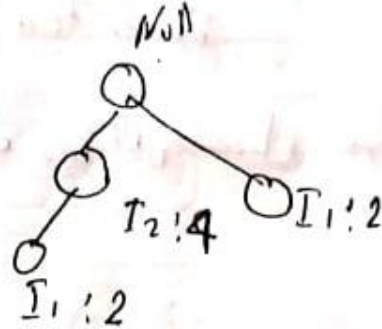
T200	I <sub>2</sub> , <u>I<sub>9</sub></u>
T400	I <sub>1</sub> , I <sub>2</sub> , <u>I<sub>9</sub></u>

eliminate





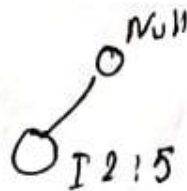
for  $I_3$



T300	I2, <u>I3</u>
T500	I1, <u>I3</u>
T600	I2, <u>I3</u>
T700	I1, <u>I3</u>
T800	I1, I2, <u>I3</u> , I5
T900	I1, I2, <u>I3</u>

هنا نعمل Cond-Tree التي تحتوي  
 $I_3$  يعني من غيرهم  
 $I_5$  ثاني

for  $I_1$



TID	List of items
T100	<del>I1</del> , I2, I5
T200	<del>I1</del> , I2, I4
T300	<del>I1</del> , I2, I3
T400	<del>I1</del> , I2, I4
T500	<del>I1</del> , I3
T600	<del>I1</del> , I2, I3
T700	<del>I1</del> , I3
T800	<del>I1</del> , I2, I3, I5
T900	<del>I1</del> , I2, I3
T1000	<del>I1</del> , I2

وبدرونا نعمل Elimination للبائع  
 (عنا ما عارفين)

صباحاً (Condition- $fp$ ) Trees هي أشجار نضيف منهم (frequent Itemset) ونحذف من جدول

Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	{{I2, I1: 1}, {I2, I1, I3: 1}}	<I2:2, I1:2>	{I2, I5: 2}, {I1, I5: 2}, {I2, I1, I5: 2}
I4	{{I2, I1: 1}, {I2: 1}}	<I2:2>	{I2, I4: 2}
I3	{{I2, I1: 2}, {I2: 2}, {I1: 2}}	<I2:4, I1:2>, <I1:2>	{I2, I3: 4}, {I1, I3: 4}, {I2, I1, I3: 2}
I1	{{I2: 5}}	<I2:5>	{I2, I1: 5}

تمام وناقص كل Non empty subset من D  
Generated patterns

وهو عمل زون الذي حصل هناك في Apriori

### Pattern Evaluation Method

- Not all association rules are interesting
  - Buys(X, "Computer games" → buys(X, "Videos") [40%, 66%]
  - P("videos") is 75% > 66%
  - The two items are
    - negatively associated means buying one decreases the likelihood of buying the other
  - We need to measure "real strength" of rule
  - Correlation analysis

- $A \rightarrow B$  [support , confidence , correlation]

○ یعنی نضیف قاعدة زیادة عشان نشوف اذا كان نشوف هو ایجابی او سلبی

1. **Lift** =  $\frac{P(A \cup B)}{P(A)P(B)}$

- A and B are independent if  $P(A \cup B) = P(A)P(B)$
  - Otherwise, **dependent and correlated** occurrence
  - If lift < 1, A is **Negatively correlated** with B
  - If lift > 1, A is **Positively correlated** with B ..... A's occurrence "lifts" the occurrence of B
-