



Samueli
Computer Science



CS145 Discussion: Week 9

Frequent Pattern Mining, Sequential Pattern Mining & DTW

Junheng Hao

Friday, 12/04/2020

- **Part 1: Frequent Pattern Mining and Association Rules**
 - Apriori
 - FP-Growth
 - Association Rules & Pattern Evaluation
- **Part 2: Sequential Pattern Mining**
 - GSP
 - PrefixSpan
- **Part 3: Time Series**
 - DTW

Homework #5	Due on 11:59 PM, Dec 4 (Friday, Week 9)
Homework #6	Will be released on Dec 7 (Monday Week 10), due on 11:59 PM, Dec 14 (Monday, Final week)
Project: Kaggle Deadline	Due on 11:59 PM, Dec 6 (Sunday, Week 9)
Project: Final Report	Due on 11:59 PM, Dec 18 (Friday, Final week)
Final Exam (100 minutes via CCLE)	Morning Session: Around 8:00-9:40 AM, Dec 16 (Wednesday, Final Week) Evening Session: Around 6:00-7:40 PM, Dec 16 (Wednesday, Final Week)

- Frequent Pattern Mining and Association Rules
 - Apriori
 - FP-Growth
 - Association Rules & Pattern Evaluation
- Sequential Pattern Mining
 - GSP
 - PrefixSpan
- Time Series
 - DTW

- Given a transactional database, two itemsets X and Y , and an association rule $X \rightarrow Y$
 - Absolute/relative support of X : absolute/relative frequency of X
 - A frequent itemset X (pattern X): support of X is no less than a minsup threshold
 - Support of $X \rightarrow Y$: probability that a transaction contains $X \cup Y$
 - Confidence of $X \rightarrow Y$: conditional probability that a transaction with X also contains Y



- An itemset X is **closed** if X is *frequent* and there exists *no super-pattern* $Y \supset X$, **with the same support** as X (proposed by Pasquier, et al. @ ICDT'99)
- An itemset X is a **max-pattern** if X is frequent and there exists no frequent super-pattern $Y \supset X$ (proposed by Bayardo @ SIGMOD'98)
- Closed pattern is a lossless compression of freq. patterns
 - Reducing the # of patterns and rules

- **Apriori pruning principle**: If there is any itemset which is infrequent, its superset should not be generated/tested! (Agrawal & Srikant @VLDB'94, Mannila, et al. @ KDD' 94)
- Method:
 - **Initially**, scan DB once to get frequent 1-itemset
 - Generate length k **candidate itemsets** from length k-1 frequent itemsets
 - **Test** the candidates against DB
 - **Terminate** when no frequent or candidate set can be generated

Apriori: Example 1

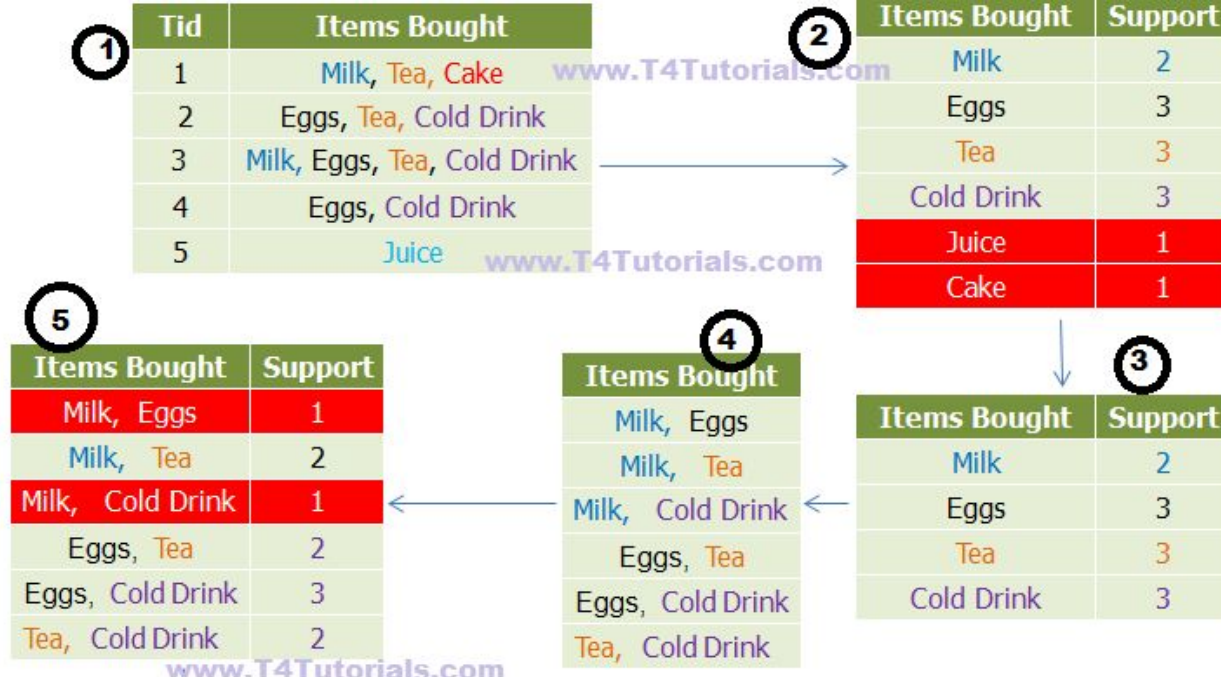
①

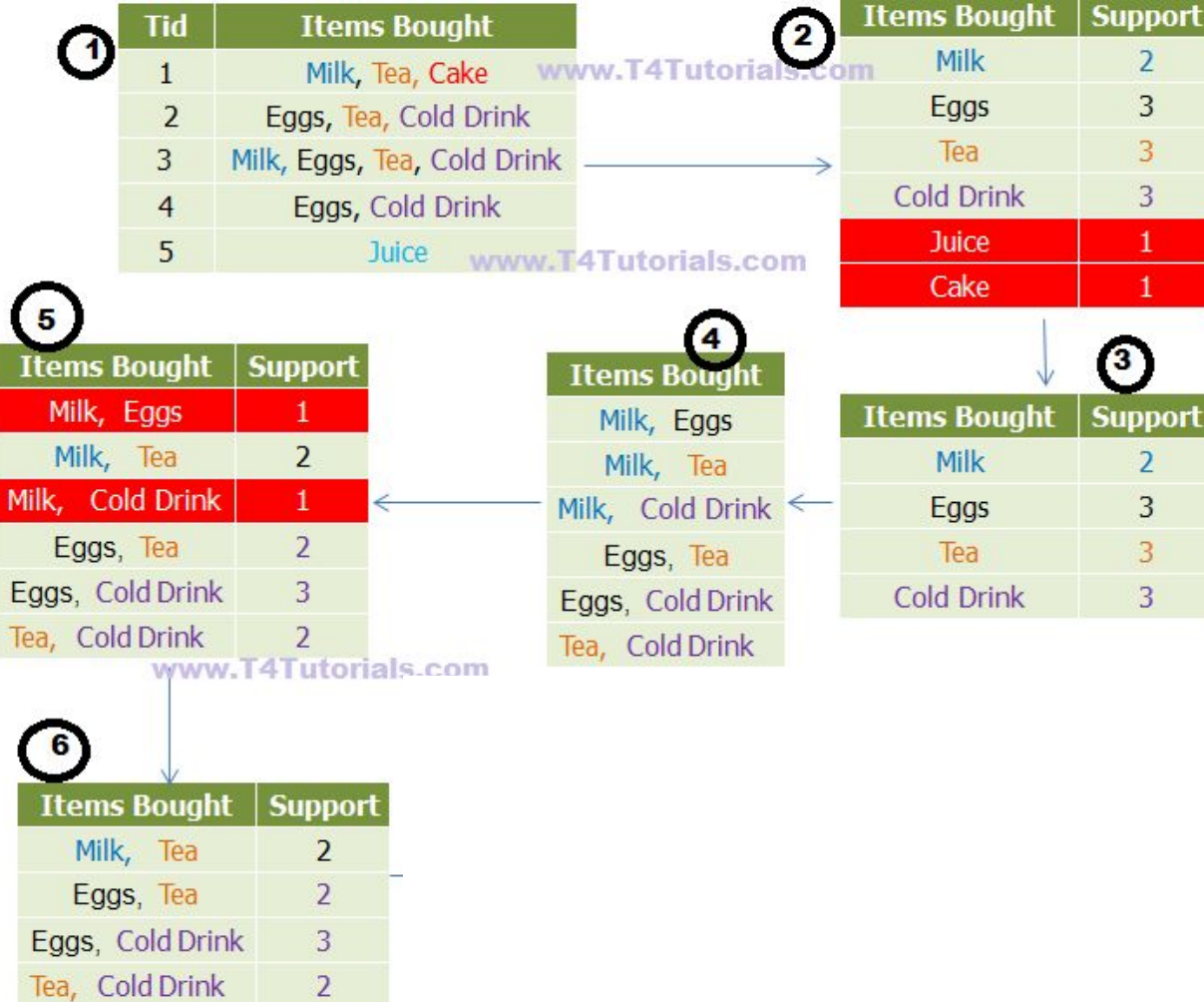
Tid	Items Bought
1	Milk, Tea, Cake
2	Eggs, Tea, Cold Drink
3	Milk, Eggs, Tea, Cold Drink
4	Eggs, Cold Drink
5	Juice

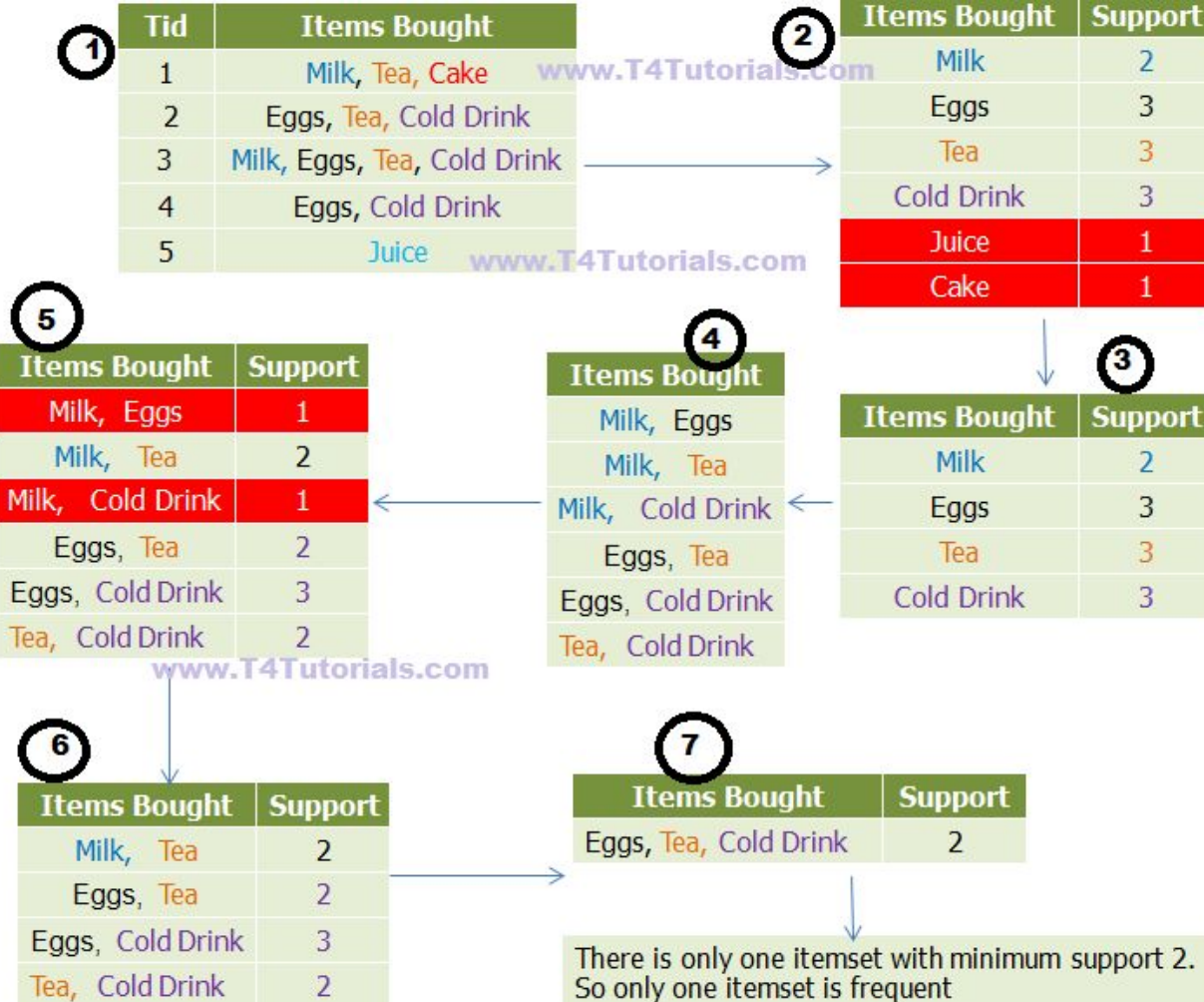
①	Tid	Items Bought	②	Items Bought	Support
	1	Milk, Tea, Cake		Milk	2
	2	Eggs, Tea, Cold Drink		Eggs	3
	3	Milk, Eggs, Tea, Cold Drink		Tea	3
	4	Eggs, Cold Drink		Cold Drink	3
	5	Juice		Juice	1
				Cake	1











Apriori: Example 2

1

Tid	Items Bought
1	Milk, Tea, Cake
2	Eggs, Tea, Cold Drink
3	Milk, Eggs, Tea, Cold Drink
4	Eggs, Cold Drink
5	Juice

1

Tid	Items Bought
1	Milk, Tea, Cake
2	Eggs, Tea, Cold Drink
3	Milk, Eggs, Tea, Cold Drink
4	Eggs, Cold Drink
5	Juice

2

Items Bought	Support
Milk	2
Eggs	3
Tea	3
Cold Drink	3
Juice	1
Cake	1

1

Tid	Items Bought
1	Milk, Tea, Cake
2	Eggs, Tea, Cold Drink
3	Milk, Eggs, Tea, Cold Drink
4	Eggs, Cold Drink
5	Juice

www.T4Tutorials.com

2

Items Bought	Support
Milk	2
Eggs	3
Tea	3
Cold Drink	3
Juice	1
Cake	1

www.T4Tutorials.com

3

Items Bought	Support
Eggs	3
Tea	3
Cold Drink	3

1

Tid	Items Bought
1	Milk, Tea, Cake
2	Eggs, Tea, Cold Drink
3	Milk, Eggs, Tea, Cold Drink
4	Eggs, Cold Drink
5	Juice

www.T4Tutorials.com

2

Items Bought	Support
Milk	2
Eggs	3
Tea	3
Cold Drink	3
Juice	1
Cake	1

4

Items Bought
Eggs, Tea
Eggs, Cold Drink
Tea, Cold Drink

als.com

3

Items Bought	Support
Eggs	3
Tea	3
Cold Drink	3

1

Tid	Items Bought
1	Milk, Tea, Cake
2	Eggs, Tea, Cold Drink
3	Milk, Eggs, Tea, Cold Drink
4	Eggs, Cold Drink
5	Juice

www.T4Tutorials.com

2

Items Bought	Support
Milk	2
Eggs	3
Tea	3
Cold Drink	3
Juice	1
Cake	1

5

Items Bought	Support
Eggs, Tea	2
Eggs, Cold Drink	3
Tea, Cold Drink	2

www.T4Tutorials.com

4

Items Bought	Support
Eggs, Tea	2
Eggs, Cold Drink	3
Tea, Cold Drink	2

3

Items Bought	Support
Eggs	3
Tea	3
Cold Drink	3

1

Tid	Items Bought
1	Milk, Tea, Cake
2	Eggs, Tea, Cold Drink
3	Milk, Eggs, Tea, Cold Drink
4	Eggs, Cold Drink
5	Juice

www.T4Tutorials.com

2

Items Bought	Support
Milk	2
Eggs	3
Tea	3
Cold Drink	3
Juice	1
Cake	1

5

Items Bought	Support
Eggs, Tea	2
Eggs, Cold Drink	3
Tea, Cold Drink	2

www.T4Tutorials.com

4

Items Bought	Support
Eggs, Tea	2
Eggs, Cold Drink	3
Tea, Cold Drink	2

3

Items Bought	Support
Eggs	3
Tea	3
Cold Drink	3

6

Items Bought	Support
Eggs, Cold Drink	3

- Grow long patterns from short ones using local frequent items only
 - “abc” is a frequent pattern
 - Get all transactions having “abc”, i.e., project DB on abc: $DB|abc$
 - “d” is a local frequent item in $DB|abc \rightarrow abcd$ is a frequent pattern

FP-Growth



<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o, w}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

min_support = 3

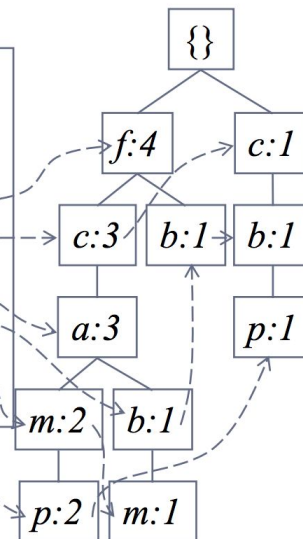
1. **Scan** DB once, find frequent 1-itemset (single item pattern)
2. **Sort** frequent items in frequency descending order, f-list
3. **Scan** DB again, construct FP-tree

Header Table

Item frequency head

<i>f</i>	4
<i>c</i>	4
<i>a</i>	3
<i>b</i>	3
<i>m</i>	3
<i>p</i>	3

F-list = f-c-a-b-m-p



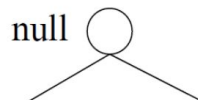
FP-Growth: Example

TID	Items
1	{a,b}
2	{b,c,d}
3	{a,c,d,e}
4	{a,d,e}
5	{a,b,c}
6	{a,b,c,d}
7	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c,e}

**Transaction
Database**

min support = 2

F-list = a-b-c-d-e



Header table

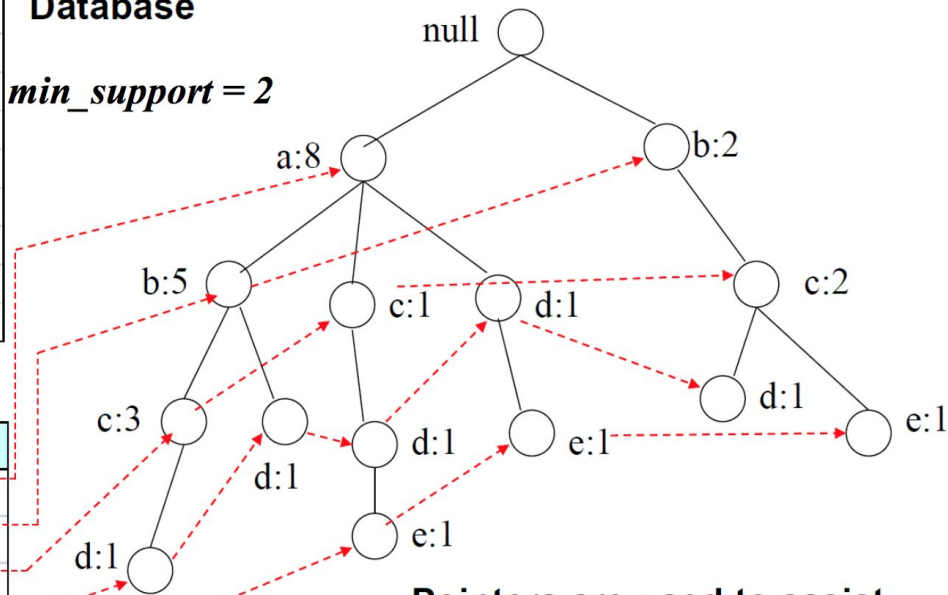
Item	Pointer
a	
b	
c	
d	
e	

TID	Items
1	{a,b}
2	{b,c,d}
3	{a,c,d,e}
4	{a,d,e}
5	{a,b,c}
6	{a,b,c,d}
7	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c,e}

Transaction Database

min_support = 2

F-list = a-b-c-d-e



Header table

Item	Pointer
a	
b	
c	
d	
e	

Pointers are used to assist frequent itemset generation

TID	Items
1	{a,b}
2	{b,c,d}
3	{a,c,d,e}
4	{a,d,e}
5	{a,b,c}
6	{a,b,c,d}
7	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c,e}

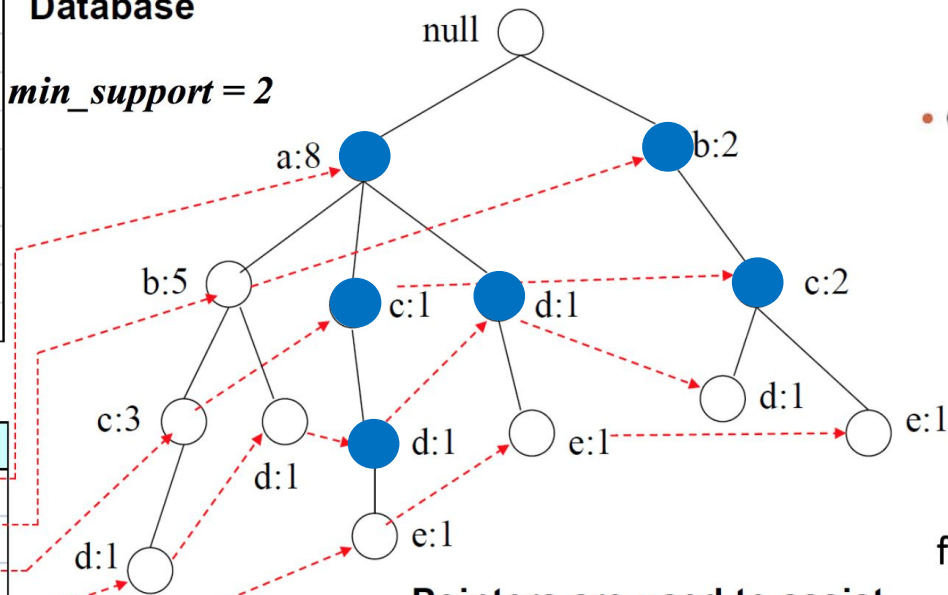
Transaction Database

min_support = 2

Header table

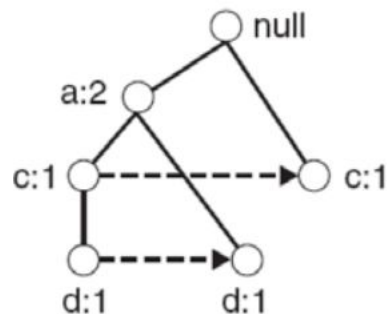
Item	Pointer
a	
b	
c	
d	
e	

F-list = a-b-c-d-e



Pointers are used to assist frequent itemset generation

- Conditional pattern base for e:
{acd:1; ad:1; bc:1}
- Conditional FP-tree for e:



frequent patterns with e are: {ade:2, de:2, ce:2, ae:2, e:3}

- (Relative) support

$$\text{Support} \{ \text{🍎} \} = \frac{4}{8}$$

- Confidence

$$\text{Confidence} \{ \text{🍎} \rightarrow \text{🍺} \} = \frac{\text{Support} \{ \text{🍎}, \text{🍺} \}}{\text{Support} \{ \text{🍎} \}}$$

- Strong association rules

- Satisfying minimum support and minimum confidence
- Recall: $\text{Confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{support}(A \cup B)}{\text{support}(A)}$

Transaction 1	🍎 🍺 🍲 🍗
Transaction 2	🍎 🍺 🍲
Transaction 3	🍎 🍺
Transaction 4	🍎 🍏
Transaction 5	🍼 🍺 🍲 🍗
Transaction 6	🍼 🍺 🍲
Transaction 7	🍼 🍺
Transaction 8	🍼 🍏

- Not all strong association rules are interesting

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

- Shall we target people who play basketball for cereal ads? $play\ basketball \Rightarrow eat\ cereal$ [40%, 66.7%]
- Hint: What is the overall probability of people who eat cereal?
 - $3750/5000 = 75\% > 66.7\%$!
- Confidence measure of a rule could be misleading

- Lift

$$\text{Lift} \{ \text{🍎} \rightarrow \text{🍺} \} = \frac{\text{Support} \{ \text{🍎}, \text{🍺} \}}{\text{Support} \{ \text{🍎} \} \times \text{Support} \{ \text{🍺} \}}$$

$$\text{lift} = \frac{P(A \cup B)}{P(A)P(B)}$$

1: independent

>1: positively correlated

<1: negatively correlated

Transaction 1	🍎 🍺 🍲 🍗
Transaction 2	🍎 🍺 🍲
Transaction 3	🍎 🍺
Transaction 4	🍎 🍏
Transaction 5	🍼 🍺 🍲 🍗
Transaction 6	🍼 🍺 🍲
Transaction 7	🍼 🍺
Transaction 8	🍼 🍏

$$\text{Support} \{\text{🍎}\} = \frac{4}{8}$$

$$\text{Confidence} \{\text{🍎} \rightarrow \text{🍺}\} =$$

$$\frac{\text{Support} \{\text{🍎, 🍺}\}}{\text{Support} \{\text{🍎}\}}$$

$$\text{Lift} \{\text{🍎} \rightarrow \text{🍺}\} =$$

$$\frac{\text{Support} \{\text{🍎, 🍺}\}}{\text{Support} \{\text{🍎}\} \times \text{Support} \{\text{🍺}\}}$$



Rule	Support	Confidence	Lift
$A \Rightarrow D$			
$C \Rightarrow A$			
$A \Rightarrow C$			
$B \& C \Rightarrow D$			

$$\text{Support} \{\text{🍎}\} = \frac{4}{8}$$

$$\text{Confidence} \{\text{🍎} \rightarrow \text{🍺}\} =$$

$$\frac{\text{Support} \{\text{🍎, 🍺}\}}{\text{Support} \{\text{🍎}\}}$$

$$\text{Lift} \{\text{🍎} \rightarrow \text{🍺}\} =$$

$$\frac{\text{Support} \{\text{🍎, 🍺}\}}{\text{Support} \{\text{🍎}\} \times \text{Support} \{\text{🍺}\}}$$



Rule	Support	Confidence	Lift
$A \Rightarrow D$	2/5	2/3	10/9
$C \Rightarrow A$	2/5	2/4	5/6
$A \Rightarrow C$	2/5	2/3	5/6
$B \& C \Rightarrow D$	1/5	1/3	5/9

Chi-Square Test



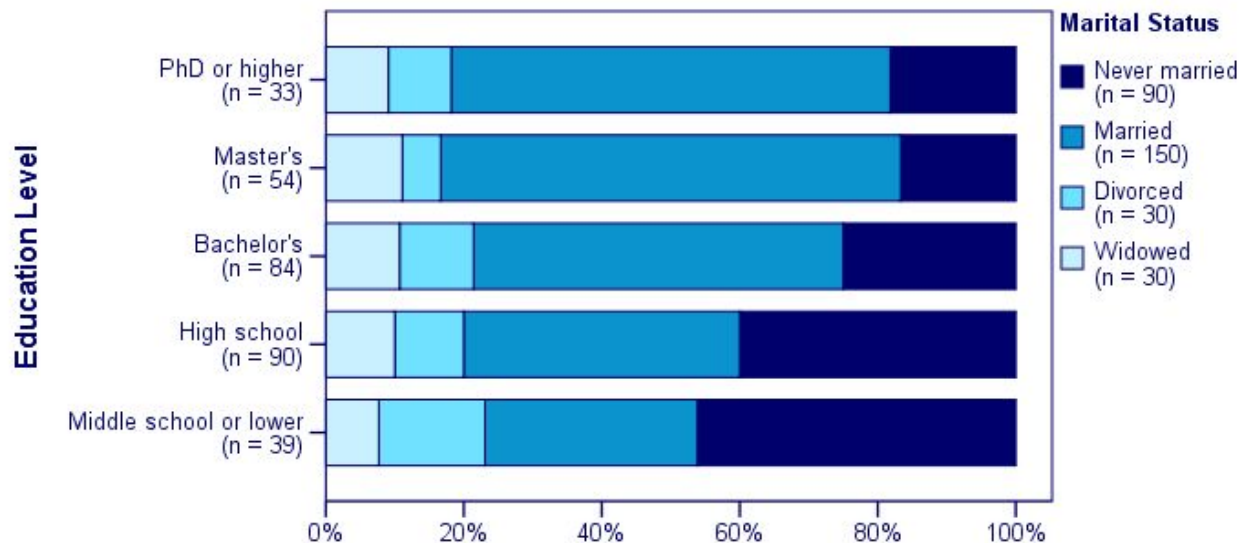
- Question: Are education level and marital status related?

Marital Status by Education | n = 300

	Middle school or lower	High school	Bachelor's	Master's	PhD or higher	Total
Never married	18	36	21	9	6	90
Married	12	36	45	36	21	150
Divorced	6	9	9	3	3	30
Widowed	3	9	9	6	3	30
Total	39	90	84	54	33	300

- Marital status is related to education level.

Marital Status by Education Level | N = 300

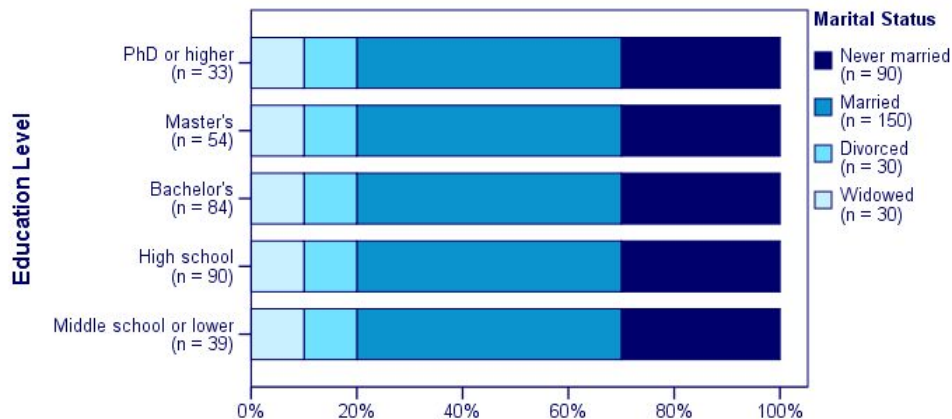


Chi-Square Test



- The null hypothesis for a chi-square independence test is that
 - two categorical variables are independent in some population.

Marital Status by Education Level | N = 300



- Statistical independence means that
 - the frequency distribution of a variable is the same for all levels of some other variable.

Chi-Square Test



- Expected frequencies are
 - the frequencies we expect in our sample if the null hypothesis holds.

Expected Frequencies for Perfectly Independent Variables

	Middle school or lower	High school	Bachelor's	Master's	PhD or higher	Total
Never married	11.7	27.0	25.2	16.2	9.9	90.0
Married	19.5	45.0	42.0	27.0	16.5	150.0
Divorced	3.9	9.0	8.4	5.4	3.3	30.0
Widowed	3.9	9.0	8.4	5.4	3.3	30.0
Total	39.0	90.0	84.0	54.0	33.0	300.0

- **Null hypothesis (independent) → expected frequencies**

$$P(\text{middle, never}) = P(\text{middle})P(\text{never}) = (39/300) * (90/300)$$

$$\text{Expected \# of (middle, never)} = 300 * P(\text{middle, never}) = 39 * 90 / 300 = 11.7$$

Expected Frequencies for Perfectly Independent Variables

	Middle school or lower	High school	Bachelor's	Master's	PhD or higher	Total
Never married	11.7	27.0	25.2	16.2	9.9	90.0
Married	19.5	45.0	42.0	27.0	16.5	150.0
Divorced	3.9	9.0	8.4	5.4	3.3	30.0
Widowed	3.9	9.0	8.4	5.4	3.3	30.0
Total	39.0	90.0	84.0	54.0	33.0	300.0

Chi-Square Test



- Real data → observed frequencies:

Marital Status by Education | n = 300

	Middle school or lower	High school	Bachelor's	Master's	PhD or higher	Total
Never married	18	36	21	9	6	90
Married	12	36	45	36	21	150
Divorced	6	9	9	3	3	30
Widowed	3	9	9	6	3	30
Total	39	90	84	54	33	300

Chi-Square Test



- Add up the differences for each of the 5*4=20 cells
 - $\rightarrow \chi^2$

$$\bullet X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

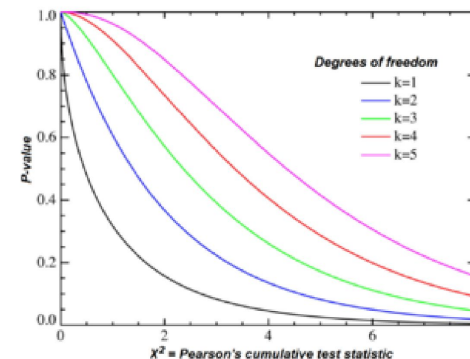
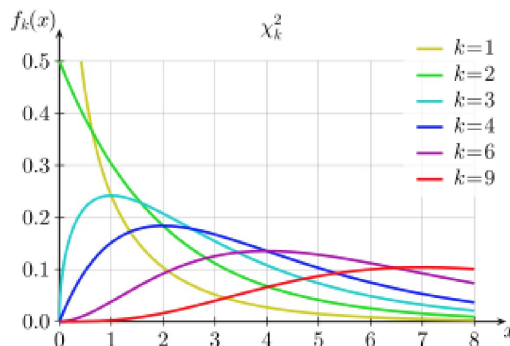
$$\chi^2 = \frac{(18 - 11.7)^2}{11.7} + \frac{(36 - 27)^2}{27} + \dots + \frac{(6 - 5.4)^2}{5.4} = 23.57$$

Chi-Square Test



- Is $\chi^2=23.57$ a large value?
 - If yes, reject the null hypothesis \rightarrow A and B are dependent
 - But how to tell if it is a large value?
 - χ^2 • Follows Chi-squared distribution with degree of freedom as $(r - 1) \times (c - 1)$

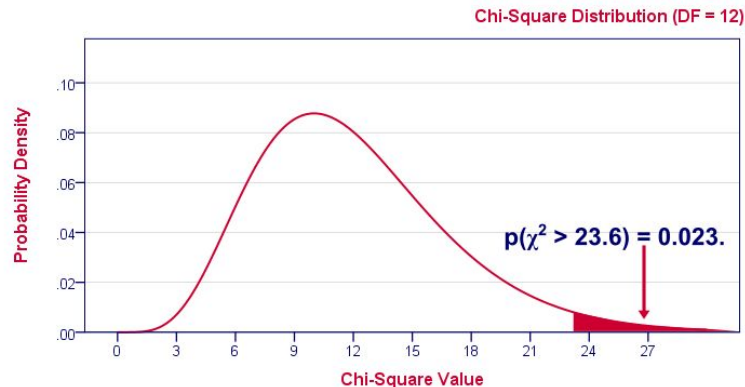
Pearson
established
it in 1900.
[See more.](#)



Chi-Square Test



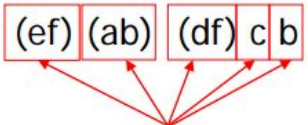
- In this example, $df = (5 - 1) \cdot (4 - 1) = 12$.
 - How to interpret $P(\chi^2 > 23.57) = 0.023$?
 - The probability of _____ under _____ hypothesis is very small, 2.3%.
 - A small p-value basically means that the data are unlikely under the null hypothesis. The convention is to reject the null hypothesis if $p < 0.05$.
 - Should we reject the null hypothesis in this case? Yes!
 - ***“An association between education and marital status was observed, $\chi^2(12) = 23.57, p = 0.023$.”***



- Frequent Pattern Mining and Association Rules
 - Apriori
 - FP-Growth
 - Association Rules & Pattern Evaluation
- Sequential Pattern Mining
 - GSP
 - PrefixSpan
- Time Series
 - DTW

- Given a set of sequences, find the complete set of *frequent* subsequences

A sequence: $\langle (ef) (ab) (df) c b \rangle$



A sequence database

SID	sequence
10	$\langle a(\underline{abc})(\underline{ac})d(cf) \rangle$
20	$\langle (ad)c(bc)(ae) \rangle$
30	$\langle (ef)(\underline{ab})(df)\underline{cb} \rangle$
40	$\langle eg(af)cbc \rangle$

An element may contain a set of items. Items within an element are unordered and we list them alphabetically.

$\langle a(bc)dc \rangle$ is a subsequence of $\langle \underline{a}(\underline{abc})(ac)\underline{d}(\underline{cf}) \rangle$

Given support threshold $min_sup = 2$, $\langle (ab)c \rangle$ is a sequential pattern

The Apriori Property of Sequential Patterns

- A basic property: Apriori (Agrawal & Srikant'94)
 - If a sequence S is not frequent
 - Then none of the super-sequences of S is frequent
 - E.g, $\langle hb \rangle$ is infrequent \rightarrow so do $\langle hab \rangle$ and $\langle (ah)b \rangle$

Seq. ID	Sequence
10	$\langle (bd)cb(ac) \rangle$
20	$\langle (bf)(ce)b(fg) \rangle$
30	$\langle (ah)(bf)abf \rangle$
40	$\langle (be)(ce)d \rangle$
50	$\langle a(bd)bcb(ade) \rangle$

Given support threshold
 $min_sup = 2$

- s_1 and s_2 can join, if dropping first item in s_1 is the same as dropping the last item in s_2
- Examples:
 - $\langle (12)3 \rangle \text{ join } \langle (2)34 \rangle = \langle (12)34 \rangle$
 - $\langle (12)3 \rangle \text{ join } \langle (2)(34) \rangle = \langle (12)(34) \rangle$

GSP: Example

- Initial candidates: all singleton sequences
 - $\langle a \rangle, \langle b \rangle, \langle c \rangle, \langle d \rangle, \langle e \rangle, \langle f \rangle, \langle g \rangle, \langle h \rangle$
- Scan database once, count support for candidates

min_sup = 2

Seq. ID	Sequence
1	$\langle (cd)(abc)(abf)(acdf) \rangle$
2	$\langle (abf)(e) \rangle$
3	$\langle (abf) \rangle$
4	$\langle (dgh)(bf)(agh) \rangle$

Cand	Sup
$\langle a \rangle$	
$\langle b \rangle$	
$\langle c \rangle$	
$\langle d \rangle$	
$\langle e \rangle$	
$\langle f \rangle$	
$\langle g \rangle$	
$\langle h \rangle$	

- Initial candidates: all singleton sequences
 - $\langle a \rangle, \langle b \rangle, \langle c \rangle, \langle d \rangle, \langle e \rangle, \langle f \rangle, \langle g \rangle, \langle h \rangle$
- Scan database once, count support for candidates

$\text{min_sup} = 2$

Seq. ID	Sequence
1	$\langle (cd)(abc)(abf)(acdf) \rangle$
2	$\langle (abf)(e) \rangle$
3	$\langle (abf) \rangle$
4	$\langle (dgh)(bf)(agh) \rangle$

Cand	Sup
$\langle a \rangle$	4
$\langle b \rangle$	4
$\langle c \rangle$	1
$\langle d \rangle$	2
$\langle e \rangle$	1
$\langle f \rangle$	4
$\langle g \rangle$	1
$\langle h \rangle$	1

min_sup = 2

Seq. ID	Sequence
1	<(cd)(abc)(abf)(acdf)>
2	<(abf)(e)>
3	<(abf)>
4	<(dgh)(bf)(agh)>

Cand	Sup
<a>	4
	4
<d>	2
<f>	4

Length 2 Candidates generated by join

Length 2 Frequent Sequences

min_sup = 2

Seq. ID	Sequence
1	<(cd)(abc)(abf)(acdf)>
2	<(abf)(e)>
3	<(abf)>
4	<(dgh)(bf)(agh)>

Cand	Sup
<a>	4
	4
<d>	2
<f>	4

Length 2 Candidates generated by join

<aa> <ab> <ad> <af> <ba> <bb> <bd> <bf>
 <da> <db> <dd> <df> <fa> <fb> <fd> <ff>
 <(ab)> <(ad)> <(af)> <(bd)> <(bf)> <(df)>

Length 2 Frequent Sequences

min_sup = 2

Seq. ID	Sequence
1	<(cd)(abc)(abf)(acdf)>
2	<(abf)(e)>
3	<(abf)>
4	<(dgh)(bf)(agh)>

Cand	Sup
<a>	4
	4
<d>	2
<f>	4

Length 2 Candidates generated by join

<aa> <ab> <ad> <af> <ba> <bb> <bd> <bf>
 <da> <db> <dd> <df> <fa> <fb> <fd> <ff>
 <(ab)> <(ad)> <(af)> <(bd)> <(bf)> <(df)>

Length 2 Frequent Sequences

<ba> <da> <db> <df> <fa>
 <(ab)> <(af)> <(bf)>

Length 2 Frequent Sequences

<ba> <da> <db> <df> <fa>
<(ab)> <(af)> <(bf)>

Length 3 Candidates generated by join

Seq. ID	Sequence
1	<(cd)(abc)(abf)(acdf)>
2	<(abf)(e)>
3	<(abf)>
4	<(dgh)(bf)(agh)>

min_sup = 2

Length 3 Frequent Sequences

Length 4 Candidates generated by join

Length 2 Frequent Sequences

<ba> <da> <db> <df> <fa>
<(ab)> <(af)> <(bf)>

Length 3 Candidates generated by join

<ba> and <(ab)> - <b(ab)> {1}
<ba> and <(af)> - <b(af)> {1}
<da> and <(ab)> - <d(ab)> {1}
<da> and <(af)> - <d(af)> {1}
<db> and <(bf)> - <d(bf)> {1, 4}
<db> and <ba> - <dba> {1, 4}
<df> and <fa> - <dfa> {1, 4}
<fa> and <(ab)> - <f(ab)> -
<fa> and <(af)> - <f(af)> {1}
<(ab)> and <(bf)> - <(abf)> {1,2,3}
<(ab)> and <ba> - <(ab)a> {1}
<(af)> and <fa> - <(af)a> {1}
<(bf)> and <fa> - <(bf)a> {1, 4}

Seq. ID

Sequence

1

<(cd)(abc)(abf)(acdf)>

2

<(abf)(e)>

3

<(abf)>

4

<(dgh)(bf)(agh)>

min_sup = 2

Length 3 Frequent Sequences

Length 4 Candidates generated by join

Length 2 Frequent Sequences

<ba> <da> <db> <df> <fa>
<(ab)> <(af)> <(bf)>

Length 3 Candidates generated by join

<ba> and <(ab)> - <b(ab)> {1}
<ba> and <(af)> - <b(af)> {1}
<da> and <(ab)> - <d(ab)> {1}
<da> and <(af)> - <d(af)> {1}
<db> and <(bf)> - <d(bf)> {1, 4}
<db> and <ba> - <dba> {1, 4}
<df> and <fa> - <dfa> {1, 4}
<fa> and <(ab)> - <f(ab)> -
<fa> and <(af)> - <f(af)> {1}
<(ab)> and <(bf)> - <(abf)> {1,2,3}
<(ab)> and <ba> - <(ab)a> {1}
<(af)> and <fa> - <(af)a> {1}
<(bf)> and <fa> - <(bf)a> {1, 4}

Seq. ID	Sequence
1	<(cd)(abc)(abf)(acdf)>
2	<(abf)(e)>
3	<(abf)>
4	<(dgh)(bf)(agh)>

min_sup = 2

Length 3 Frequent Sequences

<dba> <dfa> <(abf)> <(bf)a> <d(bf)>

Length 4 Candidates generated by join

Length 2 Frequent Sequences

<ba> <da> <db> <df> <fa>
<(ab)> <(af)> <(bf)>

Length 3 Candidates generated by join

<ba> and <(ab)> - <b(ab)> {1}
<ba> and <(af)> - <b(af)> {1}
<da> and <(ab)> - <d(ab)> {1}
<da> and <(af)> - <d(af)> {1}
<db> and <(bf)> - <d(bf)> {1, 4}
<db> and <ba> - <dba> {1, 4}
<df> and <fa> - <dfa> {1, 4}
<fa> and <(ab)> - <f(ab)> -
<fa> and <(af)> - <f(af)> {1}
<(ab)> and <(bf)> - <(abf)> {1,2,3}
<(ab)> and <ba> - <(ab)a> {1}
<(af)> and <fa> - <(af)a> {1}
<(bf)> and <fa> - <(bf)a> {1, 4}

Seq. ID	Sequence
1	<(cd)(abc)(abf)(acdf)>
2	<(abf)(e)>
3	<(abf)>
4	<(dgh)(bf)(agh)>

min_sup = 2

Length 3 Frequent Sequences

<dba> <dfa> <(abf)> <(bf)a> <d(bf)>

Length 4 Candidates generated by join

<d(bf)> and <(bf)a> - <d(bf)a> {1, 4}
<(abf)> and <(bf)a> - <(abf)a> {1}

PrefixSpan



Assume a pre-specified order on items, e.g., alphabetical order

- $\langle a \rangle$, $\langle aa \rangle$, $\langle a(ab) \rangle$ and $\langle a(abc) \rangle$ are prefixes of sequence $\langle a(abc)(ac)d(cf) \rangle$
 - Note $\langle a(ac) \rangle$ is not a prefix of $\langle a(abc)(ac)d(cf) \rangle$
- Given sequence $\langle a(abc)(ac)d(cf) \rangle$

Prefix	<u>Suffix</u>
$\langle a \rangle$	$\langle (abc)(ac)d(cf) \rangle$
$\langle aa \rangle$	$\langle (_bc)(ac)d(cf) \rangle$
$\langle a(ab) \rangle$	$\langle (_c)(ac)d(cf) \rangle$

- $(_bc)$ means: the last element in the prefix together with (bc) form one element

- Given a sequence, α , let α' be subsequence of α
 - α' is called a projection of α w.r.t. **prefix** β , if and only if
 - α' has prefix β , and
 - α' is the **maximum** subsequence of α with prefix β
- **Example:**
 - $\langle \text{ad}(\text{cf}) \rangle$ is a projection of $\langle \text{a}(\text{abc})(\text{ac})\text{d}(\text{cf}) \rangle$ w.r.t. prefix $\langle \text{ad} \rangle$

SID	sequence
10	$\langle \text{a}(\text{abc})(\text{ac})\text{d}(\text{cf}) \rangle$
20	$\langle (\text{ad})\text{c}(\text{bc})(\text{ae}) \rangle$
30	$\langle (\text{ef})(\text{ab})(\text{df})\text{cb} \rangle$
40	$\langle \text{eg}(\text{af})\text{cbc} \rangle$

PrefixSpan: Example

min_sup = 2

- 1. Find length-1 sequential patterns:

id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

<a>		<c>	<d>	<e>	<f>	<g>

PrefixSpan



min_sup = 2

- 1. Find length-1 sequential patterns:

id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

<a>		<c>	<d>	<e>	<f>	<g>
4	4	4	3	3	3	1

<a><c><d><e><f>

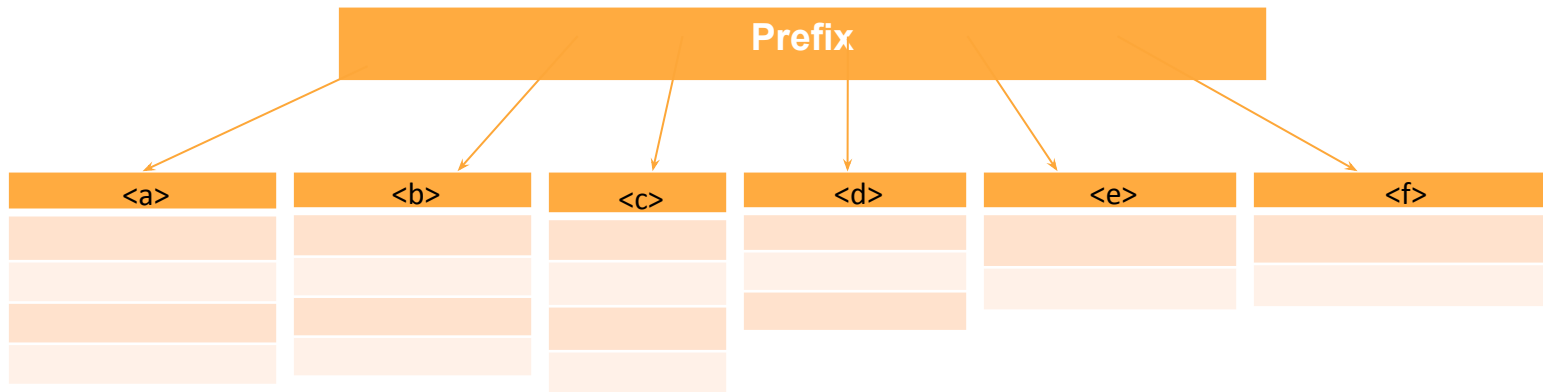
PrefixSpan



- 2. Divide search space

id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

min_sup = 2



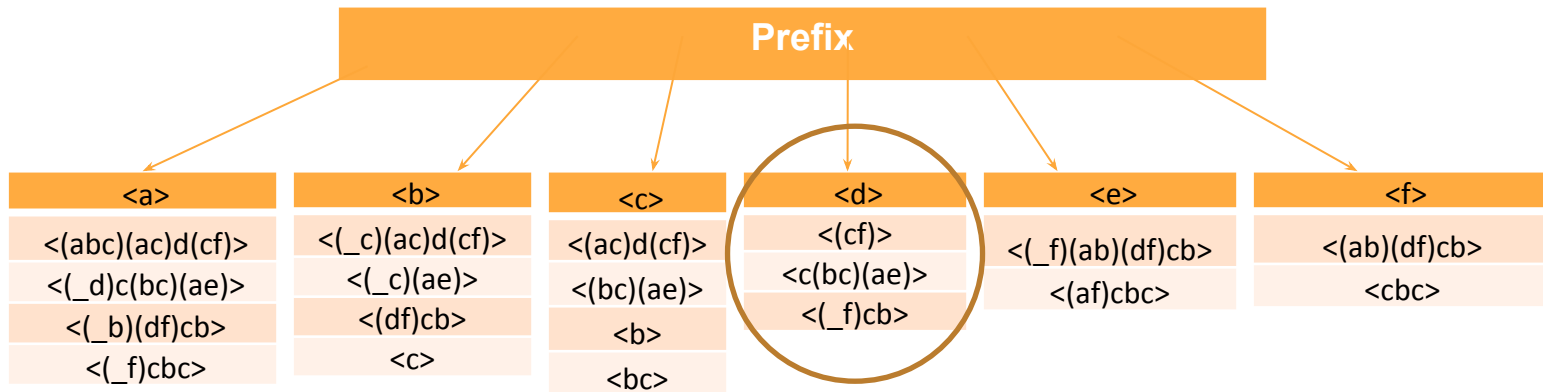
PrefixSpan



- 2. Divide search space

id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

min_sup = 2



PrefixSpan



- 3. Find subsets of sequential patterns

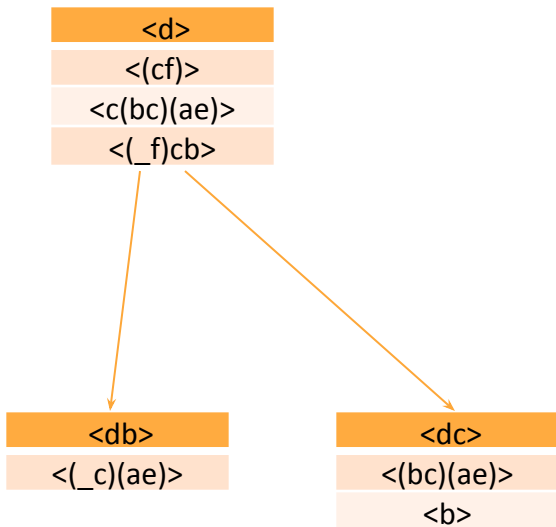
min_sup = 2

<d>
<(cf)>
<c(bc)(ae)>
<(_f)cb>

<a>		<c>	<d>	<e>	<f>	<_f>

min_sup = 2

- 3. Find subsets of sequential patterns

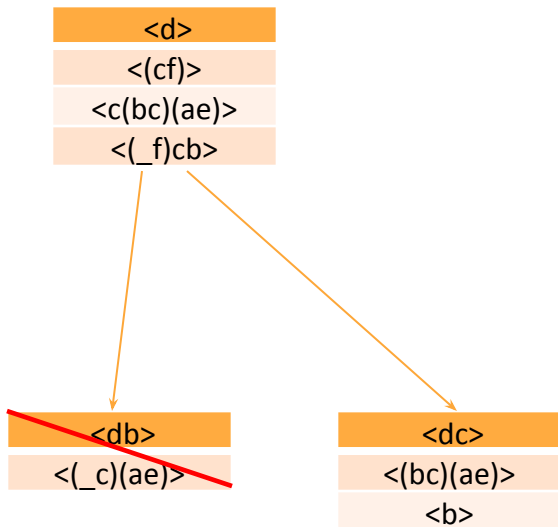


<a>		<c>	<d>	<e>	<f>	< f>
1	2	3	0	1	1	1

<db> <dc>

min_sup = 2

- 3. Find subsets of sequential patterns



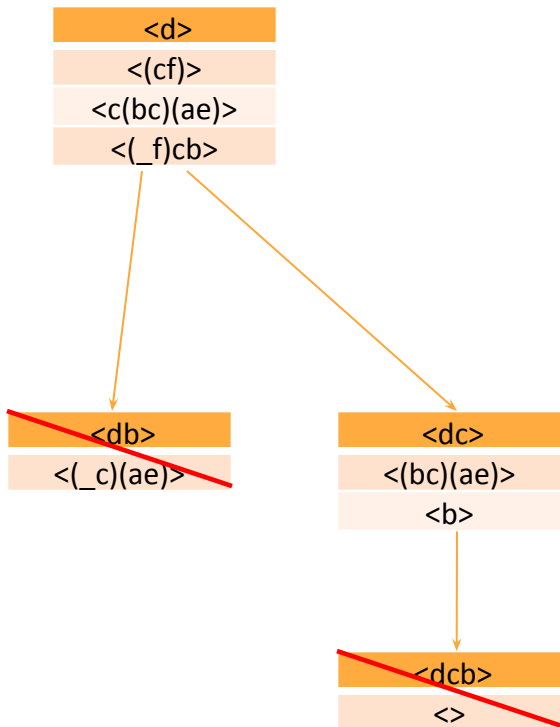
<a>		<c>	<d>	<e>	<f>	< f>
1	2	3	0	1	1	1

`<db>` `<dc>`

	<a>	<e>	<c>

min_sup = 2

- 3. Find subsets of sequential patterns



<a>		<c>	<d>	<e>	<f>	< f>
1	2	3	0	1	1	1

<db> <dc>

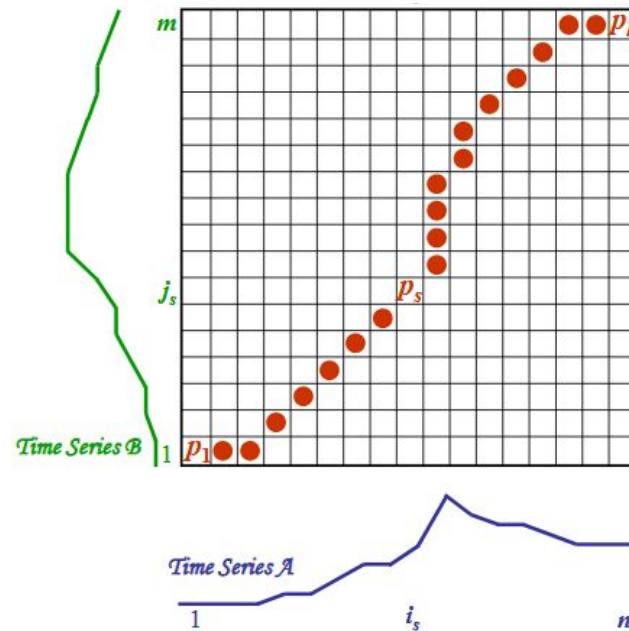
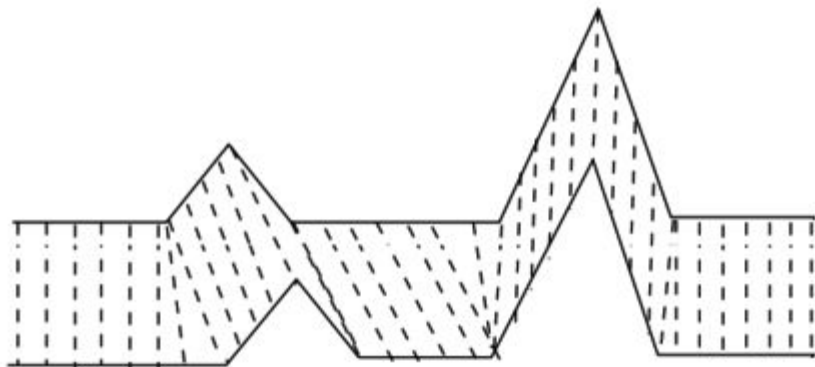
	<a>	<e>	<c>
2	1	1	1

<dcb>

- Frequent Pattern Mining and Association Rules
 - Apriori
 - FP-Growth
 - Association Rules & Pattern Evaluation
- Sequential Pattern Mining
 - GSP
 - PrefixSpan
- Time Series
 - DTW

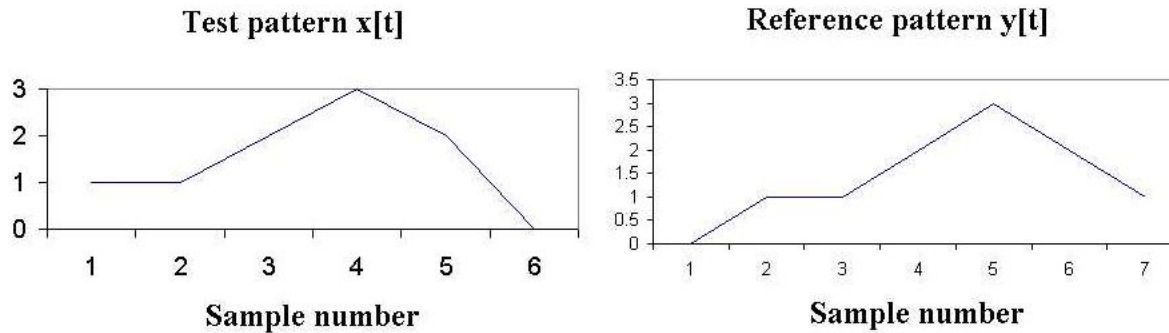
- **Given**
 - Two sequences (with possible different lengths):
 - $X = \{x_1, x_2, \dots, x_N\}$
 - $Y = \{y_1, y_2, \dots, y_M\}$
 - A local distance (cost) measure between x_n and y_m : $c(x_n, y_m)$
- **Goal:**
 - Find an alignment between X and Y, such that, the overall cost is minimized

DTW



DTW: Example

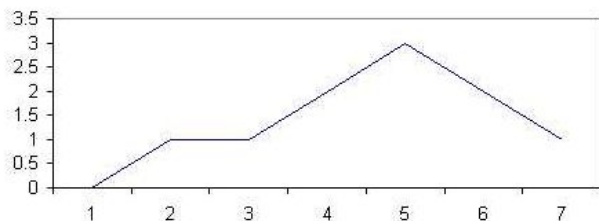
- How similar are these two peaked functions



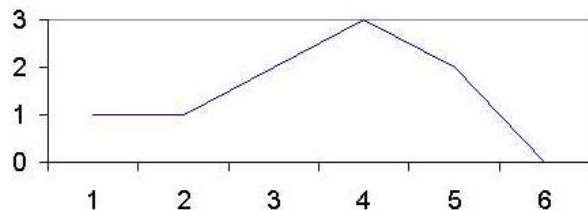
Distance function:

$$C(x, y) = (x - y)^2$$

Reference pattern $y[t]$



Test pattern $x[t]$

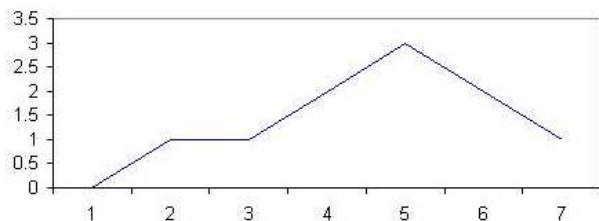


7						
6						
5						
4						
3						
2						
1						
	1	2	3	4	5	6

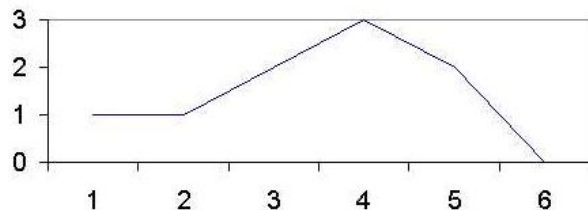
Distance function:

$$C(x, y) = (x - y)^2$$

Reference pattern $y[t]$

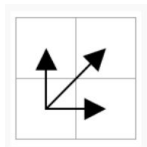


Test pattern $x[t]$



7	0	0	1	4	1	1
6	1	1	0	1	0	4
5	4	4	1	0	1	9
4	1	1	0	1	0	4
3	0	0	1	4	1	1
2	0	0	1	4	1	1
1	1	1	4	9	4	0
	1	2	3	4	5	6

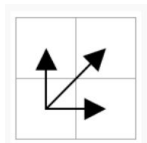
Update rule:



```
DTW[i, j] := cost + minimum(DTW[i-1, j ],
                             DTW[i , j-1],
                             DTW[i-1, j-1])
```

7	0	0	1	4	1	1
6	1	1	0	1	0	4
5	4	4	1	0	1	9
4	1	1	0	1	0	4
3	0	0	1	4	1	1
2	0	0	1	4	1	1
1	1	1	4	9	4	0
	1	2	3	4	5	6

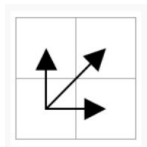
Update rule:



```
DTW[i, j] := cost + minimum(DTW[i-1, j ],
                             DTW[i , j-1],
                             DTW[i-1, j-1])
```

7	0 ->7	0 ->7	1 ->3	4 ->6	1 ->2	1 ->2
6	1 ->7	1 ->7	0 ->2	1 ->2	0 ->1	4 ->5
5	4 ->6	4 ->6	1 ->2	0 ->1	1 ->2	9 ->11
4	1 ->2	1 ->2	0 ->1	1 ->2	0 ->2	4 ->6
3	0 ->1	0 ->1	1 ->2	4 ->6	1 ->7	1 ->8
2	0 ->1	0 ->1	1 ->2	4 ->6	1 ->7	1 ->8
1	1 ->1	1 ->2	4 ->6	9 ->15	4 ->19	0 ->19
	1	2	3	4	5	6

Update rule:



```
DTW[i, j] := cost + minimum(DTW[i-1, j],
                             DTW[i, j-1],
                             DTW[i-1, j-1])
```

7	0 ->7	0 ->7	1 ->3	4 ->6	1 ->2	1 ->2
6	1 ->7	1 ->7	0 ->2	1 ->2	0 ->1	4 ->5
5	4 ->6	4 ->6	1 ->2	0 ->1	1 ->2	9 ->11
4	1 ->2	1 ->2	0 ->1	1 ->2	0 ->2	4 ->6
3	0 ->1	0 ->1	1 ->2	4 ->6	1 ->7	1 ->8
2	0 ->1	0 ->1	1 ->2	4 ->6	1 ->7	1 ->8
1	1 ->1	1 ->2	4 ->6	9 ->15	4 ->19	0 ->19
	1	2	3	4	5	6



Samueli
Computer Science



Thank you!

Q & A