# Association
# (Part II)

nanopoulos@ismll.de

---

## Outline

- Improving Apriori (FP-Growth, ECLAT)
- Questioning confidence measure
- Questioning support measure

2

# FP-growth Algorithm

Use a compressed representation of the database using an FP-tree

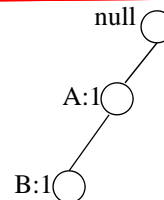Once an FP-tree has been constructed, it uses a recursive divide-and-conquer approach to mine the frequent itemsets
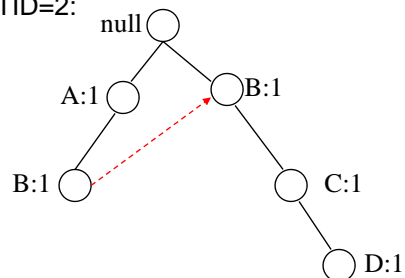
# FP-tree construction

**minSup = 2**

| 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 | {B,C} |
| 8 | {A,B,C} |
| 9 | {A,B,D} |
| 10 | {B,C,E} |

After reading TID=1:

null
A:1
B:1

After reading TID=2:

null
A:1   B:1
B:1   C:1
      D:1

# FP-tree construction

| 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 | {B,C} |
| 8 | {A,B,C} |
| 9 | {A,B,D} |
| 10 | {B,C,E} |

TID=2:

After reading TID=3:

# FP-Tree Construction

| 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 | {B,C} |
| 8 | {A,B,C} |
| 9 | {A,B,D} |
| 10 | {B,C,E} |

Transaction Database

Header table

| Item | Pointer |
|------|---------|
| A | |
| B | |
| C | |
| D | |
| E | |

Pointers are used to assist frequent itemset generation

# FP-growth

**E is frequent**

**Perhaps also frequent AE, ABE, etc.**

**Conditional pattern base and fptree for E:**

null ◯
A:7 ◯   ◯ B:3
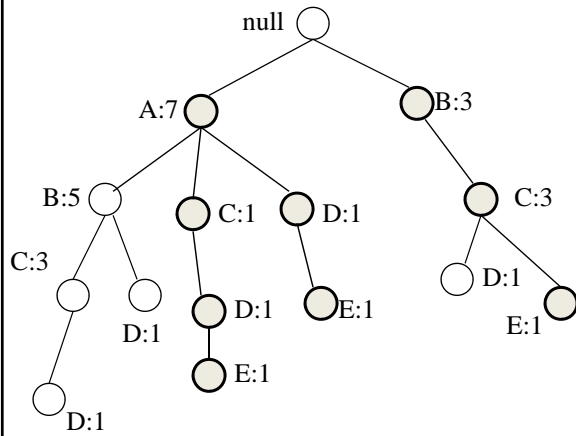B:5 ◯   ◯ C:1   ◯ D:1   ◯ C:3
C:3 ◯   ◯   ◯ D:1   ◯ E:1   ◯ D:1   ◯
D:1 ◯   ◯ E:1   E:1
   D:1

---

# FP-growth

**Conditional base and tree for E:**

null ◯
A:2 ◯   ◯ C:1
C:1 ◯   ◯ D:1
D:1 ◯

| 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} |

**Conditional Pattern base for E:**
   P = {(A:1,C:1,D:1),
       (A:1,D:1),
       (B:1,C:1)}

**Prune B**

**Build conditional FP-tree**

**Recursively apply FP-growth**

# FP-growth

**Conditional base and tree for D within conditional tree for E:**

null ○

A:2 ○

**Conditional pattern base for D within conditional base for E:**
    P = {(A:1,C:1),
          (A:1)}

**Prune C**

**Build conditional FP-tree**

**ADE and all its subsets are frequent**

---

# FP-growth

**Conditional tree for A within D within E:**

null ○

A:2 ○

**Count for A is 2: {A,D,E} is frequent itemset**

**Next step:**

**Construct conditional tree C within conditional tree E**

**Continue until exploring conditional tree for A (which has only node A)**

# Result

► Frequent itemsets found (ordered by suffix and order in which they are found):

| Suffix | Frequent Itemsets |
|--------|-------------------|
| e | {e}, {d,e}, {a,d,e}, {c,e},{a,e} |
| d | {d}, {c,d}, {b,c,d}, {a,c,d}, {b,d}, {a,b,d}, {a,d} |
| c | {c}, {b,c}, {a,b,c}, {a,c} |
| b | {b}, {a,b} |
| a | {a} |

11

# Benefits of the FP-tree Structure

- **Performance study shows**
  - FP-growth is an order of magnitude faster than Apriori, and is also faster than tree-projection
- **Reasoning**
  - No candidate generation, no candidate test
  - Use compact data structure
  - Eliminate repeated database scan
  - Basic operation is counting and FP-tree building

# ECLAT

For each item, store a list of transaction ids (tids)

Horizontal Data Layout

| TID | Items |
|-----|-------|
| 1 | A,B,E |
| 2 | B,C,D |
| 3 | C,E |
| 4 | A,C,D |
| 5 | A,B,C,D |
| 6 | A,E |
| 7 | A,B |
| 8 | A,B,C |
| 9 | A,C,D |
| 10 | B |

Vertical Data Layout

| A | B | C | D | E |
|---|---|---|---|---|
| 1 | 1 | 2 | 2 | 1 |
| 4 | 2 | 3 | 4 | 3 |
| 5 | 5 | 4 | 5 | 6 |
| 6 | 7 | 8 | 9 | |
| 7 | 8 | 9 | | |
| 8 | 10 | | | |
| 9 | | | | |

TID-list

---

# ECLAT

Determine support of any k-itemset by intersecting tid-lists of two of its (k-1) subsets.

| A |
|---|
| 1 |
| 4 |
| 5 |
| 6 |
| 7 |
| 8 |
| 9 |

$\wedge$

| B |
|---|
| 1 |
| 2 |
| 5 |
| 7 |
| 8 |
| 10 |

$\rightarrow$

| AB |
|----|
| 1 |
| 5 |
| 7 |
| 8 |

3 traversal approaches:

top-down, bottom-up and hybrid

Advantage: very fast support counting

Disadvantage: intermediate tid-lists may become too large for memory

# Pattern Evaluation

Association rule algorithms tend to produce too many rules

  many of them are uninteresting or redundant

  Redundant if {A,B,C} → {D} and {A,B} → {D}
      have same support & confidence

Interestingness measures can be used to prune/rank the derived patterns

In the original formulation of association rules, support & confidence are the only measures used

# Interestingness Measure

# Computing Interestingness Measure

Given a rule $X \rightarrow Y$, information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for $X \rightarrow Y$

|  | Y | $\overline{Y}$ |  |
|---|---|---|---|
| X | $f_{11}$ | $f_{10}$ | $f_{1+}$ |
| $\overline{X}$ | $f_{01}$ | $f_{00}$ | $f_{o+}$ |
|  | $f_{+1}$ | $f_{+0}$ | \|T\| |

$f_{11}$: support of X and Y
$f_{10}$: support of $\underline{X}$ and $\overline{Y}$
$f_{01}$: support of $\overline{X}$ and $\underline{Y}$
$f_{00}$: support of $\overline{X}$ and $\overline{Y}$

Used to define various measures

◆ support, confidence, lift, Gini, J-measure, etc.

# Drawback of Confidence

|  | Coffee | $\overline{\text{Coffee}}$ |  |
|---|---|---|---|
| Tea | 15 | 5 | 20 |
| $\overline{\text{Tea}}$ | 75 | 5 | 80 |
|  | 90 | 10 | 100 |

Association Rule: Tea $\rightarrow$ Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

$\Rightarrow$ Although confidence is high, rule is misleading

$\Rightarrow$ P(Coffee|$\overline{\text{Tea}}$) = 0.9375

# Statistical Independence

Population of 1000 students
  600 students know how to swim (S)
  700 students know how to bike (B)
  420 students know how to swim and bike (S,B)

P(S∧B) = 420/1000 = 0.42
P(S) × P(B) = 0.6 × 0.7 = 0.42

P(S∧B) = P(S) × P(B) => Statistical independence
P(S∧B) > P(S) × P(B) => Positively correlated
P(S∧B) < P(S) × P(B) => Negatively correlated

# Statistical-based Measures

Measures that take into account statistical dependence

$$Lift = \frac{P(Y \mid X)}{P(Y)}$$

$$Interest = \frac{P(X,Y)}{P(X)P(Y)}$$

$$PS = P(X,Y) - P(X)P(Y)$$

$$\phi - coefficient = \frac{P(X,Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

# Example: Lift/Interest

|  | Coffee | $\overline{\text{Coffee}}$ |  |
|---|---|---|---|
| Tea | 15 | 5 | 20 |
| $\overline{\text{Tea}}$ | 75 | 5 | 80 |
|  | 90 | 10 | 100 |

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

⇒ Lift = 0.75/0.9= 0.8333 (< 1, therefore is negatively associated)

# Drawback of Lift & Interest

|  | Y | $\overline{Y}$ |  |
|---|---|---|---|
| X | 10 | 0 | 10 |
| $\overline{X}$ | 0 | 90 | 90 |
|  | 10 | 90 | 100 |

|  | Y | $\overline{Y}$ |  |
|---|---|---|---|
| X | 90 | 0 | 90 |
| $\overline{X}$ | 0 | 10 | 10 |
|  | 90 | 10 | 100 |

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

Statistical independence:

If P(X,Y)=P(X)P(Y)  => Lift = 1

There are lots of measures proposed in the literature

Some measures are good for certain applications, but not for others

What criteria should we use to determine whether a measure is good or bad?

| # | Measure | Formula |
|---|---------|---------|
| 1 | $\phi$-coefficient | $\dfrac{P(A,B)-P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$ |
| 2 | Goodman-Kruskal's $(\lambda)$ | $\dfrac{\sum_j \max_k P(A_j,B_k)+\sum_k \max_j P(A_j,B_k)-\max_j P(A_j)-\max_k P(B_k)}{2-\max_j P(A_j)-\max_k P(B_k)}$ |
| 3 | Odds ratio $(\alpha)$ | $\dfrac{P(A,B)P(\overline{A},\overline{B})}{P(A,\overline{B})P(\overline{A},B)}$ |
| 4 | Yule's $Q$ | $\dfrac{P(A,B)P(\overline{AB})-P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{AB})+P(A,\overline{B})P(\overline{A},B)}=\dfrac{\alpha-1}{\alpha+1}$ |
| 5 | Yule's $Y$ | $\dfrac{\sqrt{P(A,B)P(\overline{AB})}-\sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{AB})}+\sqrt{P(A,\overline{B})P(\overline{A},B)}}=\dfrac{\sqrt{\alpha}-1}{\sqrt{\alpha}+1}$ |
| 6 | Kappa $(\kappa)$ | $\dfrac{P(A,B)+P(\overline{A},\overline{B})-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A)P(B)-P(\overline{A})P(\overline{B})}$ |
| 7 | Mutual Information $(M)$ | $\dfrac{\sum_i \sum_j P(A_i,B_j)\log\frac{P(A_i,B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i)\log P(A_i),-\sum_j P(B_j)\log P(B_j))}$ |
| 8 | J-Measure $(J)$ | $\max\Big(P(A,B)\log(\frac{P(B\mid A)}{P(B)})+P(A\overline{B})\log(\frac{P(\overline{B}\mid A)}{P(\overline{B})}),$ $P(A,B)\log(\frac{P(A\mid B)}{P(A)})+P(\overline{A}B)\log(\frac{P(\overline{A}\mid B)}{P(A)})\Big)$ |
| 9 | Gini index $(G)$ | $\max\Big(P(A)[P(B\mid A)^2+P(\overline{B}\mid A)^2]+P(\overline{A})[P(B\mid\overline{A})^2+P(\overline{B}\mid\overline{A})^2]$ $-P(B)^2-P(\overline{B})^2,$ $P(B)[P(A\mid B)^2+P(\overline{A}\mid B)^2]+P(\overline{B})[P(A\mid\overline{B})^2+P(\overline{A}\mid\overline{B})^2]$ $-P(A)^2-P(\overline{A})^2\Big)$ |
| 10 | Support $(s)$ | $P(A,B)$ |
| 11 | Confidence $(c)$ | $\max(P(B\mid A),P(A\mid B))$ |
| 12 | Laplace $(L)$ | $\max\Big(\frac{NP(A,B)+1}{NP(A)+2},\frac{NP(A,B)+1}{NP(B)+2}\Big)$ |
| 13 | Conviction $(V)$ | $\max\Big(\frac{P(A)P(\overline{B})}{P(A\overline{B})},\frac{P(B)P(\overline{A})}{P(B\overline{A})}\Big)$ |
| 14 | Interest $(I)$ | $\dfrac{P(A,B)}{P(A)P(B)}$ |
| 15 | cosine $(IS)$ | $\dfrac{P(A,B)}{\sqrt{P(A)P(B)}}$ |
| 16 | Piatetsky-Shapiro's $(PS)$ | $P(A,B)-P(A)P(B)$ |
| 17 | Certainty factor $(F)$ | $\max\Big(\frac{P(B\mid A)-P(B)}{1-P(B)},\frac{P(A\mid B)-P(A)}{1-P(A)}\Big)$ |
| 18 | Added Value $(AV)$ | $\max(P(B\mid A)-P(B),P(A\mid B)-P(A))$ |
| 19 | Collective strength $(S)$ | $\dfrac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})}\times\dfrac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}$ |
| 20 | Jaccard $(\zeta)$ | $\dfrac{P(A,B)}{P(A)+P(B)-P(A,B)}$ |
| 21 | Klosgen $(K)$ | $\sqrt{P(A,B)}\max(P(B\mid A)-P(B),P(A\mid B)-P(A))$ |

---

## Compact Representation of Frequent Itemsets

Some itemsets are redundant because they have identical support as their supersets

| TID | A1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 | A9 | A10 | B1 | B2 | B3 | B4 | B5 | B6 | B7 | B8 | B9 | B10 | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

Number of frequent itemsets $= 3\times\sum\limits_{k=1}^{10}\binom{10}{k}$
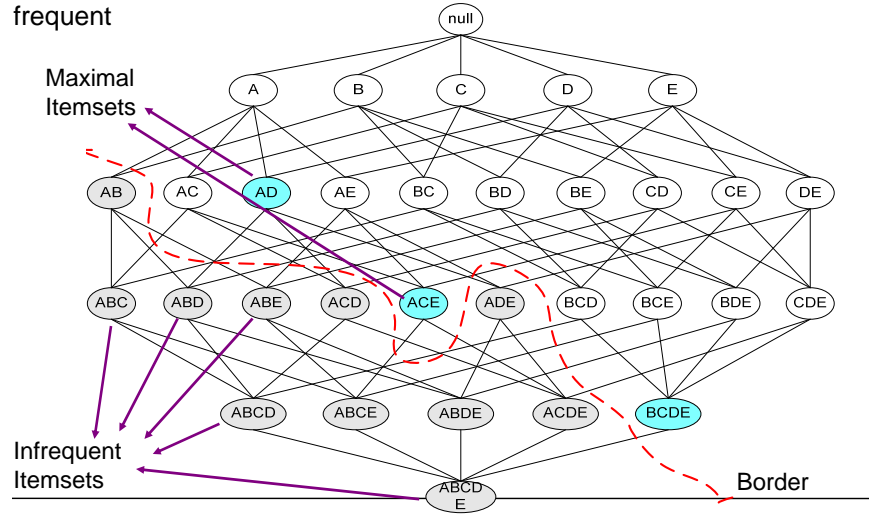
Need a compact representation

## Maximal Frequent Itemset

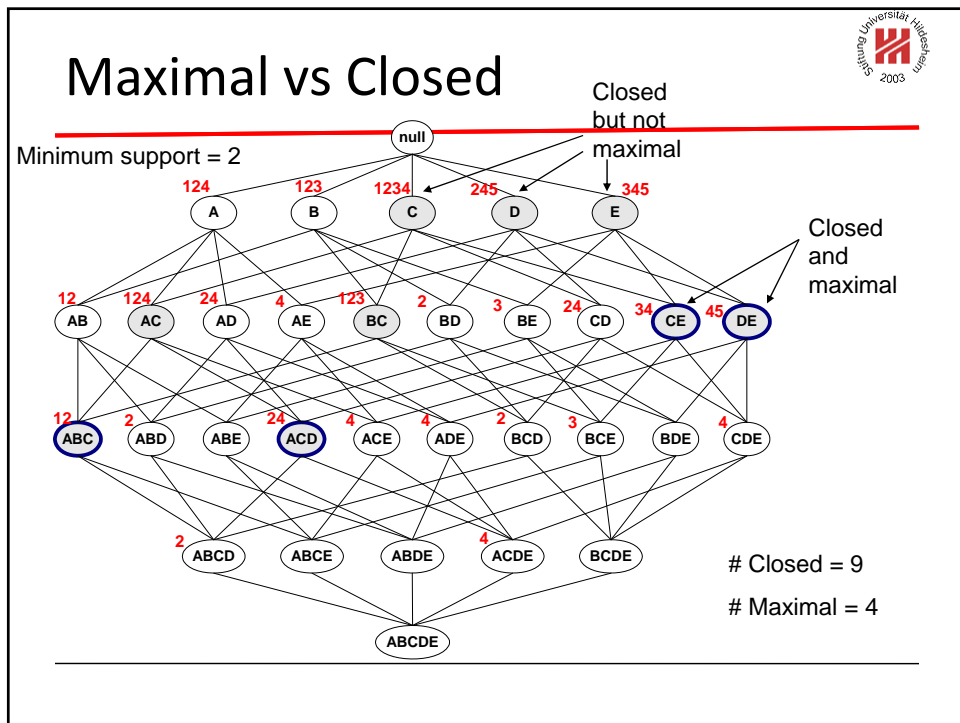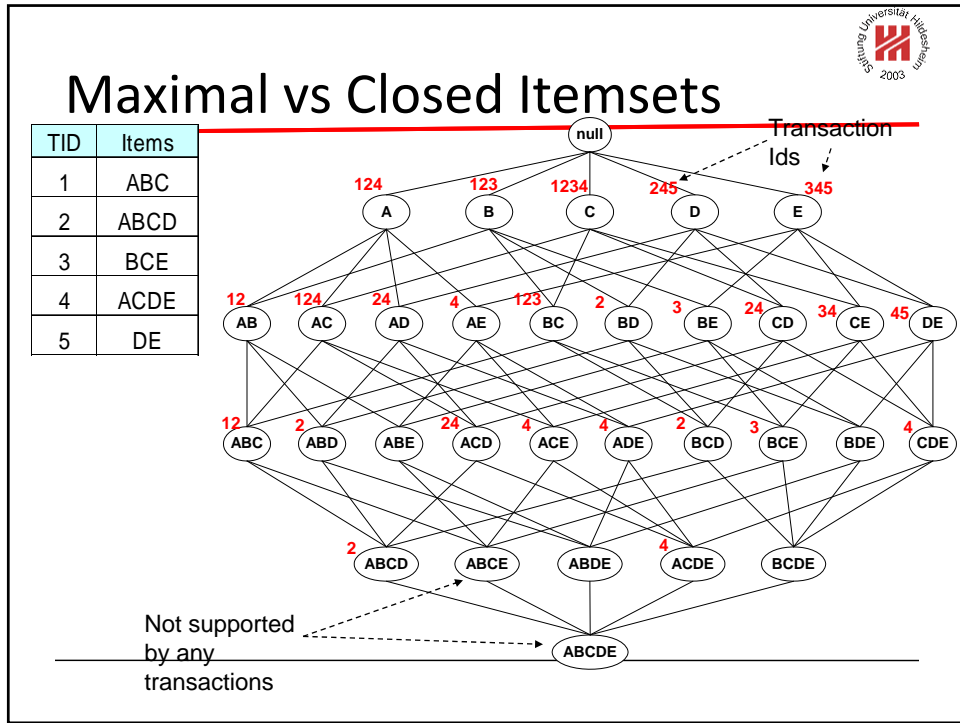An itemset is maximal frequent if none of its immediate supersets is frequent

Maximal Itemsets

Infrequent Itemsets

Border

## Closed Itemset

An itemset is closed if none of its immediate supersets has the same support as the itemset

| TID | Items |
|-----|-------|
| 1 | {A,B} |
| 2 | {B,C,D} |
| 3 | {A,B,C,D} |
| 4 | {A,B,D} |
| 5 | {A,B,C,D} |

| Itemset | Support |
|---------|---------|
| {A} | 4 |
| {B} | 5 |
| {C} | 3 |
| {D} | 4 |
| {A,B} | 4 |
| {A,C} | 2 |
| {A,D} | 3 |
| {B,C} | 3 |
| {B,D} | 4 |
| {C,D} | 3 |

| Itemset | Support |
|---------|---------|
| {A,B,C} | 2 |
| {A,B,D} | 3 |
| {A,C,D} | 2 |
| {B,C,D} | 3 |
| {A,B,C,D} | 2 |

# Maximal vs Closed Itemsets

| TID | Items |
|-----|-------|
| 1 | ABC |
| 2 | ABCD |
| 3 | BCE |
| 4 | ACDE |
| 5 | DE |



Transaction Ids

Not supported by any transactions

# Maximal vs Closed

Minimum support = 2

Closed but not maximal

Closed and maximal



# Closed = 9

# Maximal = 4

# Maximal vs Closed Itemsets

Frequent
Itemsets

Closed
Frequent
Itemsets

Frequent
Itemsets