# Mining frequent Patterns, Associations and co-relations

## Basics :-

**Frequent Pattern** :- a Pattern (set of items, subsequences, substructures) that occur frequently in a dataset. [An intrinsic & important property of dataset]

**motivation** → finding inherent regularities in data
- what products are often purchased together
- what are subsequent purchases after buying a pc
- what kinds of DNA are sensitive to this new drug?
- can we automatically classify web documents

**Applications**
- Basket-data analysis
- cross-marketing
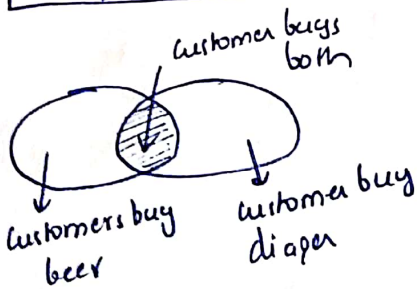- Catalogue design
- DNA sequence analysis

Frequent Patterns lay foundation for many essential data mining tasks JJ
- classification: discriminative, frequent Pattern analysis
- clustering: Frequent Pattern-based clustering
- Pattern analysis & semantic data compression & other broad applications

| TID | Items Bought |
|-----|--------------|
| 10 | Beer, nuts, diaper |
| 20 | Beer, coffee, diaper |
| 30 | Beer, Diaper eggs |
| 40 | Nuts, eggs, Milk |
| 50 | Nuts, coffee, Diaper, eggs, Milk |

**Itemset** :- A set of one or more items
eg: k-itemset $X = \{x_1, x_2, x_3 -- x_k\}$


customer buys both
customers buy beer    customer buy diaper

An itemset x is frequent if x's support is no less than minimum Support Threshold JJ

**(absolute) Support / Support count of X** : Frequency of occurrence of an itemset x

**(relative) support** :- s, is the fraction of transactions that contain X (Probability of transaction Contains X) (P(AUB))

- If a subsequence appears frequently in dataset then it is a frequential sequence pattern
- If a substructure appears frequently in dataset then it is a frequential structure pattern :)
  Structure → lattice, subgraphs, subtrees, Sub-lattices.

**Confidence** :- c, conditional Probability that a transaction having x also has y (P(B|A) = P(AnB)/P(B))

**Association rules** : finding all rules x →y with min support & confidence

**Market-Basket Analysis** (ARM is used)
: customer analysis, buying habits shopping trends & above all concepts come here

## Long Patterns problem

- A long pattern contains combinatorial number of sub patterns

eg: $(a_1 -- a_{100})$

$$\downarrow$$

$$\binom{100}{1} + \binom{100}{2} + \cdots$$

$$= 2^{100} - 1$$

$$= 1.27 \times 10^{30} \text{ sub patterns !}$$

: Solution :
mine Closed patterns & max patterns instead

### worst case in generating itemsets

: $M^N$

- M : distinct items
- N : Max length.

why? → mine Information, extract knowledge

### Downward closure property of frequent patterns

ᄄ Any subset of a frequent itemset must be frequent ⟫

## Closed Pattern

: An itemset X is closed, if X is frequent and there exist no supper-pattern Y such that

$$Y \supset X$$

and with same support as X

## Max Pattern

: An itemset X is Max, if X is frequent and there exist no super-pattern Y such that

$$Y \supset X$$

ᄄ closed Pattern is a lossless compression of frequent Patterns i.e
reducing number of patterns & rules ⟫

### Mining Association Rules : Process of finding frequent Patterns or associations within dataset or from set of data sets.

### How it is done / scalable Frequent itemset Mining Methods

- Apriori (candidate generation & test)
- FP growth (Frequent Pattern growth)
- Vertical data format (ECLAT)

## Apriori Algorithm :- [R. Agarwal, R. Srikant in 1994] [iterative]

- Apriori Pruning Principle :- ᄄ If there is any itemset which is infrequent, its superset should not be generated / tested ! ⟫

- Method :-
  1. Initially scan DB for getting frequent 1-itemset
  2. Generate length (k+1) candidate itemsets from length k frequent itemsets
  3. Test the candidates against DB
  4. Terminate when no frequent or candidate set can be generated.

ᄄ Association Rule Mining is of two types single level & multi-level ⟫

| TID | Itemset |
|-----|---------|
| T1 | F, A, D, B |
| T2 | D, A, C, E, B |
| T3 | C, A, B, E |
| T4 | B, A, D |

## An example of Apriori Algorithm

- Support-count : 60%.
- min-confidence : 80%.

**Sol:-** Finding the 1-itemset

| Items | count | SupportCount |
|-------|-------|--------------|
| A | 4 | $4/4 \times 100 = 100\%$ | ✓ |
| B | 4 | $4/4 \times 100 = 100\%$ | ✓ |
| D | 3 | $3/4 \times 100 = 75\%$ | ✓ |
| C | 2 | $2/4 \times 100 = 50\%$ |
| E | 2 | $2/4 \times 100 = 50\%$ |
| F | 1 | $1/4 \times 100 = 25\%$ |

Finding the 2-itemset

| Items | Count | Support Count % |
|-------|-------|-----------------|
| A, D | 3 | $3/4 \times 100 = 75\%$ | ✓ |
| A, B | 4 | $4/4 \times 100 = 100\%$ | ✓ |

finding the 3-itemset

| Items | Count | Support Count % |
|-------|-------|-----------------|
| A, D, B | 3 | $3/4 \times 100\% = 75\%$ | ✓ |

Association rule :-

$$A, D \rightarrow B \qquad 3/3 \times 100\% = 100\% \checkmark$$
$$B \rightarrow A, D \qquad 3/4 \times 100\% = 75\%$$
$$A, B \rightarrow D \qquad 3/4 \times 100\% = 75\%$$
$$D \rightarrow A, B \qquad 3/3 \times 100\% = 100\% \checkmark$$
$$B, D \rightarrow A \qquad \text{INVALID}$$
$$A \rightarrow B, D \qquad \text{INVALID}$$

## The Apriori Algorithm pseudocode

$C_k$ : candidate itemset of size k

$L_k$ : frequent itemset of size k

$L_1$ : {frequent itemsets};

```
for (k=1; Lk!=0; k++) do began
    Ck+1 = Candidates generated from Lk;
    for each transaction t in database do
        Increment the count of all candidates in Ck+1
        that are contained in t

    Lk+1 = Candidates in Ck+1 with min support

end

return Uk Lk
```

How to generate candidate keys

Step1: Self-joining Lk

Step2: Pruning

eg: L3 = {abc, abd, acd, ace, bcd}

self joining L3*L3
• abcd from abc & abd
• acde from acd & ace

→ Pruning:
    acde is removed, because
    ade is not in L3

→ C4 = {abcd}

Why counting supports of a candidate a Problem

↳ Total number of candidates can be very huge
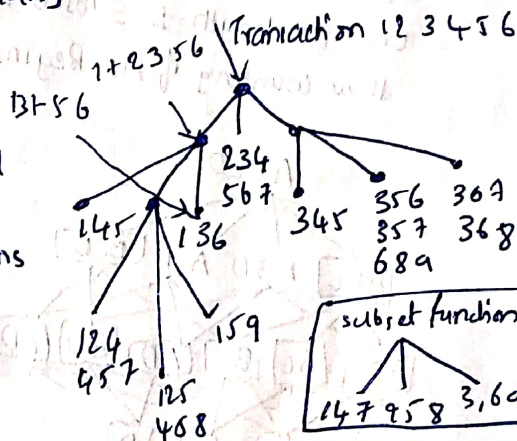→ One transactions may contain many candidates

★ Solution /Hashing itemset count
    • Candidate itemsets are stored in hash tree
    • leaf node of hash tree contains a list of itemsets & count
    • Interior node contains a hash table
    • subset function, finds all candidates contained in a transaction



1+2356   Transaction 1 2 3 4 5 6
BF 56

234
567
145   136   345   356 367
                  357 368
                  689

124
457   159
125
468

subset function

147 958  3,6 a

Anti monotonicity
    • If a set cannot pass a test, all of its
      supersets will fail same test as well

Improving the efficiency of Apriori Algorithm

1. Partition : Scan database only twice
    - Any itemset that is potentially frequent in DB, must be frequent
      in atleast one of partitions of DB
    Scan1 : Partition database & find local frequent patterns
    Scan2 : consolidate global frequent patterns

2. DHP : Reduce Number of candidates :-
    - A k-itemset whose corresponding hashing bucket count is
      below the threshold cannot be frequent.

An effective
hash-based
Algorithm for
mining association
rules

| Count | itemlets |
|-------|----------|
| 35 | {ab, ad, ae} |
| 88 | {bd, be, de} |
| ⋮ | ⋮ |
| 102 | {yz, ys, wt} |

SIGMOD'95

- Frequent 1-itemset : a, b, d,
- ab is not a candidate - 2
  itemset If sum of count (
  ae) is below support-
  heshold.

3. Sampling for frequent Patterns
   - select a sample of original database, mining frequent Patterns
     within Sample using APriori
   - Scan database once to verify frequent itemsets found in

VLDB96
↓
sampling
large databases
for Association
rules

Sample, only borders of closure of frequent Patterns au
checked

eg: check abcd, instead of ab, ac, -- etc

- Scan database again to find missing frequent Patterns

4. DIC: Reduce number of scans [dynamic itemset counting]
   - once both A & D are determined frequent, the counting
     of AD begins
   - once all length-2 Subsets of BCD are determined frequent
     then counting of BCD Begins



SIGMOD 97

itemset lattice

[Pattern growth approach]

- Bottlenecks of Apriori Approach:
  • Breadth first (level-wise) search
  • candidate generation & test
    +often generates a huge number of candidates

- FP Growth Approach
  • Depth first Search
  • Avoid explicit candidate generation

ccc Growing long
Patterns from short
ones using local
frequent items only

eg: If d is a local
frequent Pattern, abc→abcd is frequent
Pattern.

• Compress a large database into a compact, frequent-pattern tree (fp-tree) structure.
    —avoids costly database scans.
    —highly condensed, but complete.

## Construction of FP-Tree

1. First create the root of tree, labelled with Null
2. scan database D a second time.
    The items in each transaction are processed in L order (Sorted according to descending support count) and a branch is created for each transaction.
3. when considered the branch to be added for a transaction, the count of each node among a common prefix is incremented by '1'.
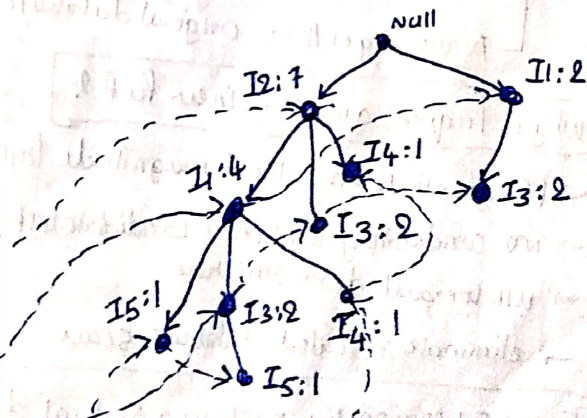
## Steps to create FP Tree

1. Scan DB once, find frequent -1 itemset (single item pattren)
2. Order frequent items in frequency descending order
3. scan DB again, construct FP-Tree

eg:- T100 : I1, I2, I5
     T200 : I2, I4 .
     T300 : I2, I3
     T400 : I1, I2, I4
     T500 : I1, I3
     T600 : I2, I3
     T700 : I3, I1
     T800 : I1, I2, I3, I5
     T900 : I1, I2, I3

| Item | Supp count |
|------|-----------|
| I1 | 6 |
| I2 | 7 |
| I3 | 6 |
| I4 | 2 |
| I5 | 2 |

In descending order:

| Item | Support count | Node-link |
|------|--------------|-----------|
| I2 | 7 | |
| I1 | 6 | |
| I3 | 6 | |
| I4 | 2 | |
| I5 | 2 | |

# Mining Frequent Patterns using FP-Tree

- divide & conquer
  - recursively grow frequent pattern path using FP Tree
- Method
  - For each item, construct its conditional Pattern-Base and then its conditional FP-Tree
  - Repeat the process on each newly created conditional FP-Tree
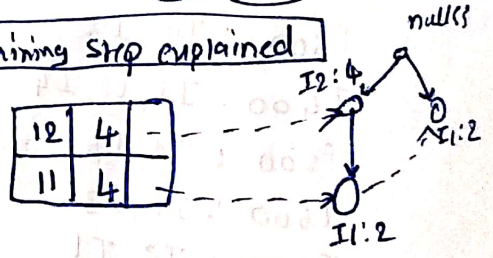  - until resulting FP-tree is empty(or) it contains only one path

| Item | Conditional Pattern Base | Conditional FP Tree | Frequent Patterns |
|------|--------------------------|---------------------|-------------------|
| $I_5$ | $\{(I_2 I_1 :1), (I_2, I_1, I_3 :1)\}$ | $\{I_2:2, I_1:2\}$ | $(I_2, I_5 :2)(I_1, I_5:2)$ $(I_2, I_1, I_5 :2)$ |
| $I_4$ | $\{(I_2, I_1 :1)(I_2 :1)\}$ | $\{I_2:2\}$ | $\{I_2 I_4:2\}$ |
| $I_3$ | $\{(I_2, I_1 :2), (I_2:2)(I_1 :2)\}$ | $\{I_2:4, I_1:2\}\{I_1:2\}$ | $(I_2 I_3 :4),$ $(I_1, I_3: 2)$ $(I_2 I_1 I_3:2)$ |
| $I_1$ | $\{(I_2 :4)\}$ | $\{I_2:4\}$ | $(I_2 I_1:4)$ |

Association rules can be drawn out same as apriori

## Benifits of Fp-Tree structure

→ completeness
  - Preserve complete information for frequent pattern mining

→ compactness
  - Reduce irrelevant information - infrequent items are gone
  - Items in frequency descending order
  - Never larger than original database

### mining step explained

| 12 | 4 | - |
| 11 | 4 | - |



null$\{\}$

$I_2:4$

$I_1:2$

$I_1':2$

## Why is frequent growth pattern fast ?

→ FP-growth is an order of magnitude faster than Apriori & also tree projection
→ No candidate generation & candidate test
→ use compact data structure
→ eliminate repeated database scan

## ECLAT : Mining By exploring vertical data format   [mining closed patterns]

- vertical format : $t(AB) = \{T_{11}, T_{25} - \}$   (Trans-id)
- Deriving frequent patterns Based on vertical intersections
  - $t(x) = t(y)$ : x & y alway happen together
  - $t(x) \subset t(y)$ : transaction having x alway have Y

[which patterns are interesting?] — [Pattern evaluation methods]

• Most Association rule mining algorithms employ a support-confidence framework.

→• many of rules generated are still not interesting to user.
• The above statement is true when mining for long patterns on at low support Thresholds.

℃℃ Strong Rules are not necessarily Interesting ℧℧

✓
based on subjective (or)
objective perspectives.

• The Interesting new measures are objective statistics.

Then which Strong association rules are interesting ?

• There are several Co-relation measures help us to choose Association Rules.

[Lift] : The occurrence of itemset A is independent of occurrence of itemset B if $P(A \cup B) = P(A) / P(B)$;

Otherwise itemsets A & B are dependent & correlated as events

∴ $Lift(A,B) = \dfrac{P(A \cup B)}{P(A)(P(B))}$    $\left( on \quad P(B/A) / P(B) \right)$

• If $Lift(A,B) < 1$, then occurrence of A is negatively correlated with occurrence of B (occurrence of one likely would be absence of other)

• If $Lift(A,B) > 1$, then A & B are positively correlated.
                    (one occurrence implies another occurrence)

• If $Lift(A,B) = 1$ then A & B are idependent & there is no correlation Between them.

eg:- Calculating the chi-square value $(x^2)$ for the given data

|  | game | game | $\Sigma row$ |
|---|---|---|---|
| video | 4000 | 3500 | 7500 |
| video | 2000 | 500 | 2500 |
| $\Sigma col$ | 6000 | 4000 | 10000 |

→ P(game) = 60%.     P(game, video) / P(game) * P(video) = 6%
P(video) = 75%.

$< 1$ ∴ negative Correlation

customer Purchasing both
two independent purchases.

$$x^2 = \Sigma \frac{(observed - expected)^2}{expected}$$

$$= \frac{(4000 - \$500)^2}{4500} + \frac{(3500 - 3000)^2}{3000} + \frac{(2000 - 1500)^2}{1500}$$

$$+ \frac{(500 - 100)^2}{100} = 555.6 \quad x^2 (>1)$$

☞ Lift & $x^2$ are not null invariant ☜

## A comparision of Pattern evaluation measures

→ If two itemsets are given then

· All confidence : all_conf(A,B) = $\frac{sup(A \cup B)}{max[sup(A), sup(B)]}$

$= min[P(A/B), P(B/A)]$

__MAX confidence__ : max_conf(A,B) = $max[P(A/B), P(B/A)]$

__KULCZYNSKI__ : kulc(A,B) = $\frac{1}{2}(P(A/B) + P(B/A))$

Cosine measure : $cosine(A, B) = \frac{P(A \cup B)}{\sqrt{P(A) \times P(B)}} = \frac{sup(A \cup B)}{\sqrt{sup(A) \times sup(B)}}$

$= \sqrt{P(A/B) \times P(B/A)}$

| Multiple-level Association Rules | Imbalance Yatio |

· It is not always easy to find the strong association
· So we must use multi levels of abstraction

## Kinds of patterns & rules

- **Basic Patterns**
  - Frequent patterns
  - association rule
  - closed/max pattern
  - generation

- **Multilevel & multidimensional patterns**
  - Multilevel (uniform/varied/reduced support)
  - multidimensional Pattern (high dimensional pattern)
  - continuous data & discretization (based on statistical)

- **Extended Patterns**
  - approximate pattern
  - uncertain pattern
  - compressed pattern
  - rare pattern/negative pattern
  - high dimensional & colossal pattern

## Mining Methods

- **Basic Mining Methods**
  - Candidate generation (Apriori, Partition, Sampling)
  - Pattern growth (FPGrowth, Hmine, FPmax, closet +)
  - vertical format (Eclat, CHARM)

- **Mining Interesting Patterns**
  - interestingness (subjective vs objective)
  - constraint based mining
  - correlation rules
  - exception rules

- **Distributed, parallel & Incremental**
  - distributed/parallel mining
  - incremental mining
  - stream pattern

## Extensions & Applications

- **Extended data types**
  - sequential & time series patterns
  - structural (tree, lattice)
  - spatial
  - temporal (periodic, involuntary)
  - image video etc
  - network pattern

- **Applications**
  - pattern-based classification
  - pattern-based clustering
  - pattern-based semantic annotation
  - Collaborative filtering
  - Privacy preserving

### Pattern mining in Multi-level & multi dimensional

- Sometimes we also want intresting or rare pattrens (occurs rarely but of critical importance) & negative pattrens (Pattrens with negative correlation between them).

**℃℃** Based on the abstraction levels involved in a Pattern. Pattrens or association rules may have items that are residing at high, low, multiple abstraction levels. JJ

eg: $buys(x, \text{"computer"}) \rightarrow buys(X, \text{"Printer"})$ ↓ high level abstraction

$buys(x, \text{"laptop"}) \rightarrow buys(X, \text{"laser Printer"})$ ↕ low level abstration

→ These are multi-level association rules.

**℃℃℃** If the items or attributes in an association rule or pattren reference only one dimension then it is a single dimenjimal association rule/Pattren JJ

**℃℃** Otherwise multi-dimensional association rule/Pattren JJ

eg: $age(X, \text{"20..29"}) \wedge income(X, \text{"52k..58k"}) \rightarrow buys(x, \text{"ipad"})$
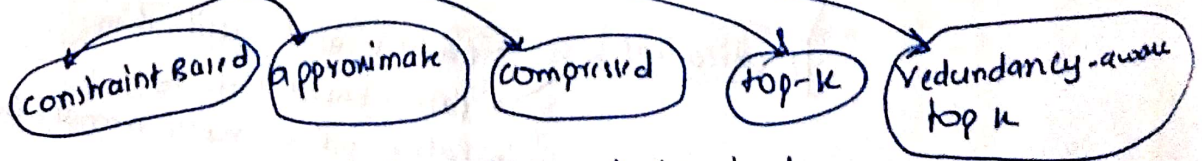
<u>Boolean association rule</u> : If a rule involve the associations between Presence or absence of items.
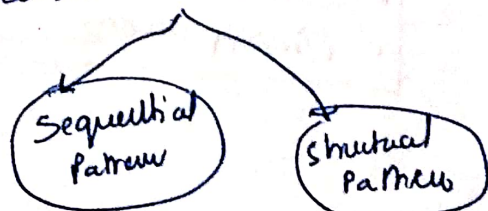
<u>quantitative association rule</u> : If a rule describes association between quantitative attributes or terms.

→ Based on constraints or criteria used to select patterns → patterns or rules discovered are

```
           ┌──────────────┬──────────────┬────────────┬───────┐
           ↓              ↓              ↓            ↓       ↓
```
( constraint Based ) ( approximate ) ( compressed ) ( top-k ) ( redundancy-aware top k )

→ Based on kinds of data & features to be mined

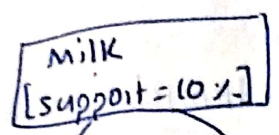( sequential Pattren ) ( structural Pattren )

# Mining Multiple level Association rules

- Items often form hierarchies
- Flexible support settings { Items at lower level are expected to have low support }
- eg :-

level1
min-sup=5%.

**Milk** [support = 10%.]

level 1
min sup. = 10%.

It is easy to find generic items to find interesting pattern

level 2
min-sup=5% **2% Milk** [sup = 6%.]

**Skim Milk** [sup = 4%.]

level 2
min-sup = 3%.

(uniform support)

(reduced support)

---

## Flexible support & Redundancy filtering

- **Flexible support** : some items are rare but more valuable
  └→ less frequent
  
  → use non-uniform, group based min support
  eg : (diamond, watch, camera) : 0.05%.
  (milk, bread) → 5%.

- **Redundancy filtering** : some rules may be redundant due to "ancestor" relationships between items.

  eg :-   milk → bread [s = 8%., L = 70%.] (ancestor)
          2% milk → bread [s = 2%. L = 72%.]

  ⟦⟦ A rule is redundant if it's support is close to the expected value, based on ancestor rule. ⟧⟧
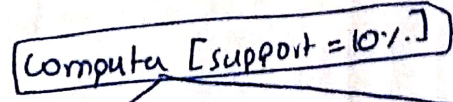
---

## uniform support vs Reduced support

- **uniform support** : same minimum support for all levels
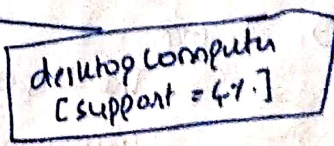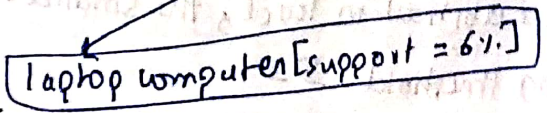
  because if (sup-threshold)
              ↓
          too high ⇒ Miss low level association
          too low ⇒ generate two many high level associations

level-1    eg :-
minsup =5%.    **computer [support = 10%.]**

level 2
minsup=5%.   **laptop computer [support = 6%.]**

**desktop computer [support = 4%.]**

- **Reduced support** : reduced minimum support at lower levels

  - **4 Search strategies**
  - → level-by-level independent (full breadth search)

    searched full
  - → level-cross filtering by k-itemset

    If node is frequent, its children are examined
    Otherwise the descendent nodes are pruned
    from search
  - → level-cross filtering by single item :-

    If {computer, printer} is frequent then other
    nodes are examined
  - → Controlled level-cross filtering by
    Single item :- (same as above except)

    A threshold called level passage threshold
    can be setup for passing down relatively
    frequent items to lower levels

---

**Mining Multi dimensional Association**

- Single dimensional rules :-
    buys(x, "milk") → buys(x, "bread")
- Multidimensional rules (2 ≥ predicates)
  - → Inter-dimension assoc. rules (no repeated predicates)
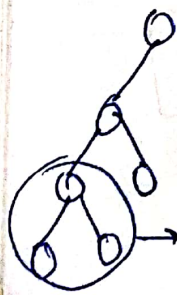  - eg:- age(x, "19-25") ∧ occupation(x, "student") → buys(x, "ok")
  - → hybrid-dimension assoc rules (repeated predicates)
    age(x, "19-25") ∧ buys(x, "Popcorn") → buys(x, "coke").

- Categorical attributes : finite no. of possible values & no ordering
    [data cube]

- quantitative Attributes : Numeric, implicit ordering among
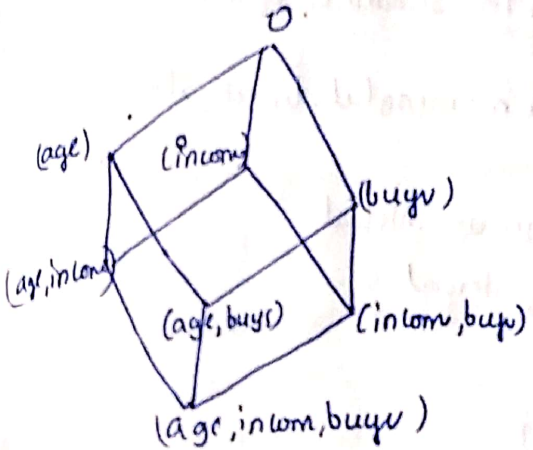    values — discretization, clustering, gradient approaches



  The deeper the abstraction level, the smaller the
  corresponding threshold

→ for each group different support, threshold _____

## mining quantitative associations

- Techniques can be categorized by how numerical are attributes are treated

1. **Static discretization**, based on predefined concept hierarchies (data cubes method)

2. **Dynamic discretization**, based on data distribution (discreted into bins dynamically)

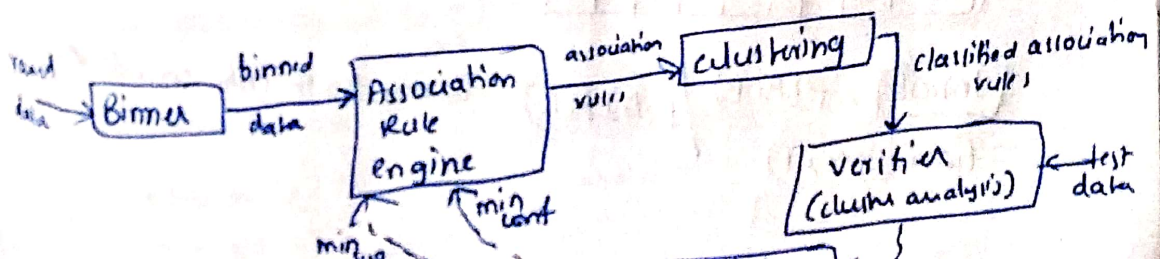3. **Distance-Based association** (clustering)

### static discretization



- Mining from data cubes can be much faster
- data cubes are well suited for mining
- Numerical values are replaced by ranges.

### dynamic discretization

for $(Sex = female) \Rightarrow$ wage : Mean = $7/hr (overall mean = $9)

- LHS : subset of population
- RHS : An extra ordinary behaviour of this subset
- This rule is accepted only if a statistical test (z-test) confirms the interference with high confidence.

→ Numeric values are dynamically discretized.

Such that confidence on compactness of rules is maximized

### ARCS [Association Rule clustering system]

## clustering the association rule

- Age $(x, 34) \wedge Income(x, "31k..40k") \rightarrow buys(x, "Iphone")$
  Age $(x, 35) \wedge Income(x, "31k..49k") \rightarrow buys(x, "Iphone")$
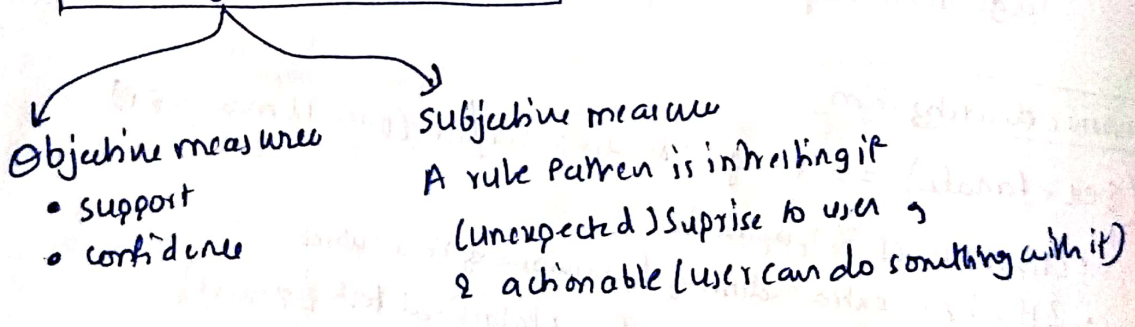  Age $(x, 34) \wedge Income(x, "31k..50k") \rightarrow buys(x, "Iphone")$

  $\downarrow$

  They can be clustered & replaced by

  Age $(x, 34..35) \wedge income(x, "31k..50k") \Rightarrow buys(x, "Iphone")$

## Mining distance base association rules

- Binning methods donot Capture the semantics of interval data
- distance-based partitioning more meaningful diurchization considering
  - density, number of points in an interval
  - "closeness" of points in an interval

## Intresting new measurements

Objective measures
- support
- confidence

Subjective measure
A rule pattern is interesting if
(unexpected) Suprise to user &
& actionable (user can do something with it)

## Correlation (lift)

$$corr(A,B) = \frac{P(A \cup B)}{P(A)P(B)} = \frac{P(B|A)}{P(B)}$$

If $corr < 1$   -ve corr
If $corr > 1$   +ve corr
If $corr = 1$   independent

eg:-

| X | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| Y | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Z | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

Intrest : $P(A \cap B) / P(A) P(B)$ : Taking Both into Consideration
(correlation) lift)

**Constraint Based mining** | mining by constraints

- **Knowledge Based mining constraint**
  - specific type of knowledge (classification)
- **Data constraint**
  - task-relevant data
- **Mining new constraint**
  - specifying threshold to support & confidence
- **Rule Based mining constraint**
  - mining by specific forms of rules.

→ **Dimensional constraint**

improves efficiency of mining process.

Such constraints may be expressed as metarules [Syntactic form of rules]

as Max or min # of predicates that can occur in rule antecedent, consequent or relationship among attributes or in aggregates.

↳ template of metarule.

$$P1 \wedge P2 \wedge P3 \wedge -- P_i \Rightarrow Q1 \wedge Q2 \wedge Q3 - Aq?$$

eg: Data mining can search for rules that match with metarule

age $(x, "30..39") \wedge$ income $(x, "41k..60k") \rightarrow$ buys $(x, slw")$

**Monotonicity:** If a set S satisfies a constraint then any superset of S satisfies a constraint.

eg: Sum $(1.price) >= 100$

**Anti monotonicity:** If a set S satisfies Violates a constraint then any superset of S also Violates constraint.

eg:1 sum $(s.price) \leq V \rightarrow$ is antimonotone.

eg:2 Avg $(1.price) <= 100$ is not antimonotone.

**Succint constraint:** only sets that satisfy constraint are enumerated

**Convertible constraint:** neither monotonic nor Anti-monotonic But convertible

**Relationships among Category of constraints**

- Succinctness
- Anti monotonicity
- monotonicity
- Convertible constraints