- Data Mining
  - : KDD (Knowledge Discovery from Data)
  - : Data Intelligence
  - : Data Management

- Association Metrics   _Relative support_   Transaction이 많아?
  - ① Support (S)   ← Transaction
    $S = \sigma(\text{itemset}) / |T|$
  - ② Confidence (C) ← itemset X의 조건부 (비중)
    $C = \sigma(X \text{ and } Y) / \sigma(X)$

  Support Count
  $\Downarrow$

  _Finding frequent itemset_ 이 중요하다
  $\text{support} \geq \text{minsup (threshold)}$

Association Rule   ┌ Apriori, FP-growth
  ① Frequent Itemset generation
  ② Rule generation
     $(X \to Y$ 구조로! $)$
     ┌─ * Same frequent itemset에서 발생하는 rule
         & support는 같고   (binary partition)
           Confidence 다르다

# Association Rule Mining — ① Frequent Itemset Mining

**b) Apriori approach**

- **anti-monotone** property of support

  $\forall X, Y : (X \subseteq Y) \quad S(X) \geq S(Y)$
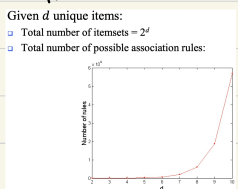
  ⊟ **Subset의** support가 < minsupport면

  superset의 support도 < minsup

Found to be Infrequent

Pruned supersets

**a) Brute-force approach** ←2탄

: 각 Itemset의 candidate itemset의

support 다 계산하고

Given $d$ unique items:
- Total number of itemsets = $2^d$
- Total number of possible association rules:

- **Algorithm**

k-1 size의 frequent itemset서
k size의 candidate itemset 만들어

```
1)  L_1 = {large 1-itemsets};
2)  for ( k = 2; L_{k-1} ≠ ∅; k++ ) do begin
3)      C_k = apriori-gen(L_{k-1});  // New candidates
4)      forall transactions t ∈ D do begin
5)          C_t = subset(C_k, t);  // Candidates contained in t
6)          forall candidates c ∈ C_t do
7)              c.count++;
8)      end
9)      L_k = {c ∈ C_k | c.count ≥ minsup}
10) end
11) Answer = ∪_k L_k;
```

## Apriori-gen function

- **Self-join**

  insert into $C_k$
  select $p.item_1, p.item_2, ..., p.item_{k-1}, q.item_{k-1}$
  from $L_{k-1}\ p,\ L_{k-1}\ q$
  where $p.item_1 = q.item_1, ..., p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1};$

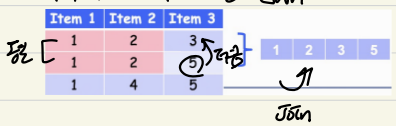- **Pruning**

  forall itemsets $c \in C_k$ do
      forall $(k-1)$-subsets $s$ of $c$ do
          if $(s \notin L_{k-1})$ then
              delete $c$ from $C_k$;

- Generate Ck+1 from Lk **(self-join)**
  - (k-1) items들은 **동일** 두 row에 대하서
  - 각 row는 Increasing order
  - later가 더 크도록 Join

| Item 1 | Item 2 | Item 3 |
|--------|--------|--------|
| 1 | 2 | 3 |
| 1 | 2 | 5 |
| 1 | 4 | 5 |

동일 [

병합 ⑤

1  2  3  5

Join

- only frequent candidate itemsets ! **(pruning)**
  - anti-monotone 이용
    Subset이 infrequent하면 superset 다 prune

⤷ Apriori needs multiple DB scan ⊕ generate candidate / test support

C) **FP Growth**

Frequent - Pattern mining
without Candidate generation

· **HeuristIc Property**

  P (frequent Itemset)

  S (set of transaction with P)

  X (Item)

  (if) X = frequent Itemset in S

  (then) {x} U P must be frequent Itemset

· FP-tree에서 Frequent Pattern get?

① **Conditional Pattern bases**

  Pattern-base | α

  : Item α의 Prefix sets



  Pattern-base[m]
  ⇒ <f,c,a> : support = 2
  ⇒ <f,c,a,b> : support = 1

② **Conditional FP-trees**

  FP-tree | α

  : Pattern-base | α 에서 minsup 만족하는 애들만 ㉿

③ **Find Frequent Pattern**

  : 갱 조합? { minsup 넘는애들로! )

| Item | Conditional Pattern Base | Conditional FP-tree | Frequent Patterns Generated |
|------|--------------------------|---------------------|------------------------------|
| I4 | {I2,I1:3},{I2,I3:1} | {I2:2, I3:2} | {I2,I4:2},{I3,I4:2},{I2,I3,I4:2} |
| I3 | {I2,I1:3},{I2:1} | {I2:4, I1:3} | {I2,I3:4}, {I1:I3:3}, {I2,I1,I3:3} |
| I1 | {I2:4} | {I2:4} | {I2,I1:4} |

---

* **FP-tree**

  : frequent pattern data 정보담고있음

- How to make FP-tree

① **Make F-list**

  · DB scan하나 1-size Item을 frequent한걸 나열,
    (desc order)

  오직 (count만necessary)  DB scan(1)

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

② **order Items in Itemset (F-list순으로나열)**

| TID | Items | Ordered items |
|-----|-------|---------------|
| 1 | f, a, c, d, g, i, m, p | f, c, a, m, p |
| 2 | a, b, c, f, l, m, o | f, c, a, b, m |
| 3 | b, f, h, j, o, w | f, b |
| 4 | b, c, k, s, p | c, b, p |
| 5 | a, f, c, e, l, p, m, n | f, c, a, m, p |

③ **row-wise 하나씩 F-list에서랑 Tree 만들기**



DB scan(2)

■ The final FP-tree

Frequent 1-items in
frequency descending order:
f,c,a,b,m,p



Min support = 3

**Implication in Property.**

❑ Process of mining frequent patterns can be viewed as first mining frequent 1-itemsets and then progressively growing each such itemset by mining its conditional pattern base, which can in turn be done similarly

❑ We successfully transform a frequent k-itemset mining problem into a sequence of $k$ frequent 1-itemset mining problems via a set of conditional pattern bases

### Why Is FP-Growth the Winner?

■ Divide-and-conquer
  ❑ Decomposing both the mining task and database according to the frequent patterns obtained so far
  ❑ Leading to focused search of smaller databases
■ Other factors
  ❑ **No candidate generation, no candidate test**
  ❑ Compressed database: FP-tree structure
  ❑ No repeated scan of the entire database
  ❑ Cheap operations: counting local frequent items and building sub FP-trees, but no pattern search and matching
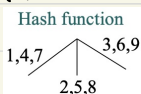
# ＊ reduce # of Comparison

- Support Count = Transaction Scan
  줄이기 ①

## Store Candidates in HashTable

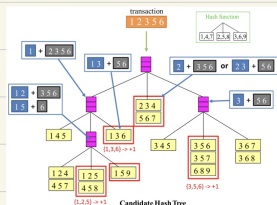- key : Candidate Itemset
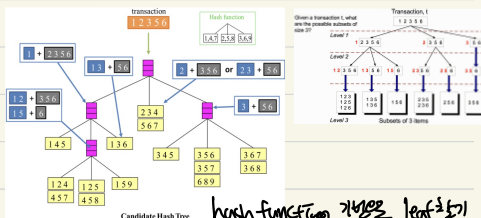- value : Support Count

① generate hash tree

Hash function
1,4,7     3,6,9
2,5,8

- hash function
  : hash-tree에서 특정 key의 location
- Max leaf size

② update hash table



**Candidate Hash Tree**

hash function 기준으로 leaf 찾기



**Candidate Hash Tree**

일치하는 candidate에 update

# ＊ reduce complexity

- Minsup 낮추면
  : 줄이면 frequent itemset 개수 증가수감소

- Dimensionality ( # of Items)
  : 아늘러면서 computation / I/O 늘어남수감소

- Size of DB
  : Apriori는 multiple pass스캔,
    run time of algorithm may increase
    with # of transactions

# Association Rule Mining ② - Rule Generation

- Given frequent Itemset 에서

  **Candidate rule** 만들기

  (응용, 전체집합에서! $2^d-2$ 개)

If {A,B,C,D} is a frequent itemset, candidate rules:

| | | | |
|---|---|---|---|
| ABC→D, | ABD→C, | ACD→B, | BCD→A, |
| A→BCD, | B→ACD, | C→ABD, | D→ABC |
| AB→CD, | AC→BD, | AD→BC, | BC→AD, |
| BD→AC, | CD→AB, | | |

- Who to remove?

  Confidence가 minconfidence가 작은애들 @

  ＊ Confidence는 anti-monotone 특징지녀,

e.g., L = {A,B,C,D} : $c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$

$$\frac{\sigma(ABCD)}{\sigma(ABC)} \geq \frac{\sigma(ABCD)}{\sigma(AB)} \geq \frac{\sigma(ABCD)}{\sigma(A)}$$

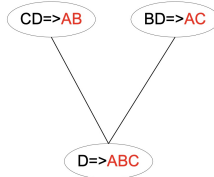Confidence is anti-monotone w.r.t. the number of items on the RHS of the rule

이거는 특이 @

1DW키 superset이 infrequent = 나머지다 infrequent

　　　　　　subset이 frequent = 나머지다 frequent

■ A candidate rule is generated by merging two rules that share the same prefix in the rule consequent

■ Example:
　□ join(CD→AB, BD→AC) would produce the candidate rule D→ABC
　□ **D→ABC is pruned if its subset BD→AC does not have high confidence**

CD=>AB　　BD=>AC

D=>ABC

Superset이 infrequent해도 @

anti -monotone



| | |
|---|---|
| Low Confidence Rule | ABCD=>{} |
| | BCD=>A　ACD=>B　ABD=>C　ABC=>D |
| | CD=>AB　BD=>AC　BC=>AD　AD=>BC　AC=>BD　AB=>CD |
| | D=>ABC　C=>ABD　B=>ACD　A=>BCD |
| Pruned Rules | |