

Challenge: There Are Too Many Frequent Patterns!

- □ A long pattern contains a combinatorial number of sub-patterns
- □ How many frequent itemsets does the following TDB₁ contain?
 - \Box TDB_{1:} T₁: {a₁, ..., a₅₀}; T₂: {a₁, ..., a₁₀₀}
 - Assuming (absolute) minsup = 1
 - Let's have a try

```
1-itemsets: (a<sub>1</sub>): 2, (a<sub>2</sub>): 2, ..., (a<sub>50</sub>): 2, (a<sub>51</sub>): 1, ..., (a<sub>100</sub>): 1, 2-itemsets: (a<sub>1</sub>, a<sub>2</sub>): 2, ..., (a<sub>1</sub>, a<sub>50</sub>): 2, (a<sub>1</sub>, a<sub>51</sub>): 1 ..., ..., (a<sub>99</sub>, a<sub>100</sub>): 1, ..., ..., ...
```

99-itemsets: {a₁, a₂, ..., a₉₉}: 1, ..., {a₂, a₃, ..., a₁₀₀}: 1

100-itemset: {a₁, a₂, ..., a₁₀₀}: 1

□ In total: $\binom{100}{1} + \binom{100}{2} + \dots + \binom{1}{1} \binom{0}{0} = 2^{100} - 1$ sub-patterns!

A too huge set for any computer to compute or store!

Expressing Patterns in Compressed Form: Closed Patterns

- How to handle such a challenge?
- □ Solution 1: **Closed patterns**: A pattern (itemset) X is **closed** if X is *frequent*, and there exists *no super-pattern* Y ⊃ X, *with the same* support as X
 - □ Let Transaction DB TDB₁: T_1 : {a₁, ..., a₅₀}; T_2 : {a₁, ..., a₁₀₀}
 - □ Suppose minsup = 1. How many closed patterns does TDB₁ contain?
 - □ Two: P_1 : "(a_1 , ..., a_{50}): 2"; P_2 : "(a_1 , ..., a_{100}): 1"
- Closed pattern is a lossless compression of frequent patterns
 - Reduces the # of patterns but does not lose the support information!
 - □ You will still be able to say: " $(a_2, ..., a_{40})$: 2", " (a_5, a_{51}) : 1"

Expressing Patterns in Compressed Form: Max-Patterns

- □ Solution 2: Max-patterns: A pattern X is a max-pattern if X is frequent and there exists no frequent super-pattern Y ⊃ X
- □ Difference from close-patterns?
 - Do not care the real support of the sub-patterns of a max-pattern
 - □ Let Transaction DB TDB₁: T_1 : {a₁, ..., a₅₀}; T_2 : {a₁, ..., a₁₀₀}
 - □ Suppose minsup = 1. How many max-patterns does TDB₁ contain?
 - □ One: P: "(a₁, ..., a₁₀₀): 1"
- Max-pattern is a lossy compression!
 - \square We only know (a₁, ..., a₄₀) is frequent
 - But we do not know the real support of $(a_1, ..., a_{40})$, ..., any more!
- ☐ Thus in many applications, mining close-patterns is more desirable than mining max-patterns

Recommended Readings

- □ R. Agrawal, T. Imielinski, and A. Swami, "Mining association rules between sets of items in large databases", in Proc. of SIGMOD'93
- □ R. J. Bayardo, "Efficiently mining long patterns from databases", in Proc. of SIGMOD'98
- □ N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal, "Discovering frequent closed itemsets for association rules", in Proc. of ICDT'99
- □ J. Han, H. Cheng, D. Xin, and X. Yan, "Frequent Pattern Mining: Current Status and Future Directions", Data Mining and Knowledge Discovery, 15(1): 55-86, 2007