# Are All the "Discovered" Patterns Interesting?

- Data mining may generate thousands of patterns: Not all of them are interesting

  - Suggested approach: Human-centered, query-based, focused mining

- **Interestingness measures**

  - A pattern is interesting if it is easily understood by humans, valid on new or test data with some degree of certainty, potentially useful, novel, or validates some hypothesis that a user seeks to confirm

- **Objective vs. subjective interestingness measures**

  - Objective: based on statistics and structures of patterns, e.g., support, confidence, etc.

  - Subjective: based on user's belief in the data, e.g., unexpectedness, novelty, actionability, etc.

# Find All and Only Interesting Patterns?

- Find all the interesting patterns: Completeness

    - Can a data mining system find all the interesting patterns? Do we need to find all of the interesting patterns?

    - Heuristic vs. exhaustive search

    - Association vs. classification vs. clustering

- Search for only interesting patterns: An optimization problem

    - Can a data mining system find only the interesting patterns?

    - Approaches

        - First general all the patterns and then filter out the uninteresting ones

        - Generate only the interesting patterns—mining query optimization

# Other Pattern Mining Issues

- Precise patterns vs. approximate patterns

  - Association and correlation mining: possible find sets of precise patterns

    - But approximate patterns can be more compact and sufficient

    - How to find high quality approximate patterns??

  - Gene sequence mining: approximate patterns are inherent

    - How to derive efficient approximate pattern mining algorithms??

- Constrained vs. non-constrained patterns

  - Why constraint-based mining?

  - What are the possible kinds of constraints? How to push constraints into the mining process?

# Pattern Interestingness Measure

- Simplicity

    e.g., (association) rule length, (decision) tree size

- Certainty

    e.g., confidence, P(A|B) = #(A and B)/ #(B), classification reliability or accuracy, certainty factor, rule strength, rule quality, discriminating weight, etc.

- Utility

    potential usefulness, e.g., support (association), noise threshold (description)

- Novelty

    not previously known, surprising (used to remove redundant rules, e.g., Illinois vs. Champaign rule implication support ratio)