



COMP9318: Data Warehousing and Data Mining

— L1: Introduction —

Chapter 1. Introduction

- Motivation: Why data mining?
- What is data mining?
- Data Mining: On what kind of data?
- Data mining functionality
- Are all the patterns interesting?
- Classification of data mining systems
- Major issues in data mining

Necessity Is the Mother of Invention



- Data explosion problem
 - Automated data collection tools and mature database technology lead to tremendous amounts of data accumulated and/or to be analyzed in databases, data warehouses, and other information repositories
- We are drowning in data, but starving for knowledge!

Who could be expected to digest millions of records, each having tens or hundreds of fields?

- Solution: Data warehousing and data mining
 - Data warehousing and on-line analytical processing
 - Mining interesting knowledge (rules, regularities, patterns, constraints) from data in large databases

Evolution of Database Technology

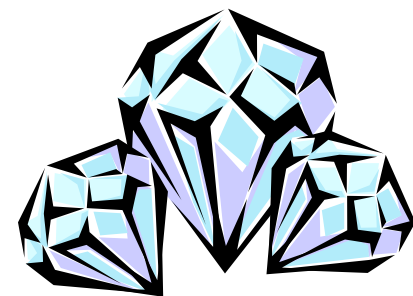


- 1960s:
 - Data collection, database creation, IMS and network DBMS
- 1970s:
 - Relational data model, relational DBMS implementation
- 1980s:
 - RDBMS, advanced data models (extended-relational, OO, deductive, etc.)
 - Application-oriented DBMS (spatial, scientific, engineering, etc.)
- 1990s:
 - Data mining, data warehousing, multimedia databases, and Web databases
- 2000s
 - Stream data management and mining
 - Data mining with a variety of applications
 - Web technology and global information systems

What Is Data Mining?



- Data mining (knowledge discovery from data)
 - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from **huge amount** of data
 - Data mining: a misnomer?
- Alternative names
 - Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
- Watch out: Is everything “data mining”?
 - (Deductive) query processing.
 - Expert systems or small ML/statistical programs



Why Data Mining?—Potential Applications

- Data analysis and decision support
 - Market analysis and management
 - Target marketing, customer relationship management (CRM), market basket analysis, cross selling, market segmentation
 - Risk analysis and management
 - Forecasting, customer retention, improved underwriting, quality control, competitive analysis
 - Fraud detection and detection of unusual patterns (outliers)
- Other Applications
 - Text mining (news group, email, documents) and Web mining
 - Stream data mining
 - DNA and bio-data analysis

Market Analysis and Management



- Where does the data come from?
 - Credit card transactions, loyalty cards, discount coupons, customer complaint calls, plus (public) lifestyle studies
- Target marketing
 - Find clusters of “model” customers who share the same characteristics: interest, income level, spending habits, etc.
 - Determine customer purchasing patterns over time
- Cross-market analysis
 - Associations/co-relations between product sales, & prediction based on such association
- Customer profiling
 - What types of customers buy what products (clustering or classification)
- Customer requirement analysis
 - identifying the best products for different customers
 - predict what factors will attract new customers
- Provision of summary information
 - multidimensional summary reports
 - statistical summary information (data central tendency and variation)

Corporate Analysis & Risk Management

- Finance planning and asset evaluation
 - cash flow analysis and prediction
 - contingent claim analysis to evaluate assets
 - cross-sectional and time series analysis (financial-ratio, trend analysis, etc.)
- Resource planning
 - summarize and compare the resources and spending
- Competition
 - monitor competitors and market directions
 - group customers into classes and a class-based pricing procedure
 - set pricing strategy in a highly competitive market

Fraud Detection & Mining Unusual Patterns



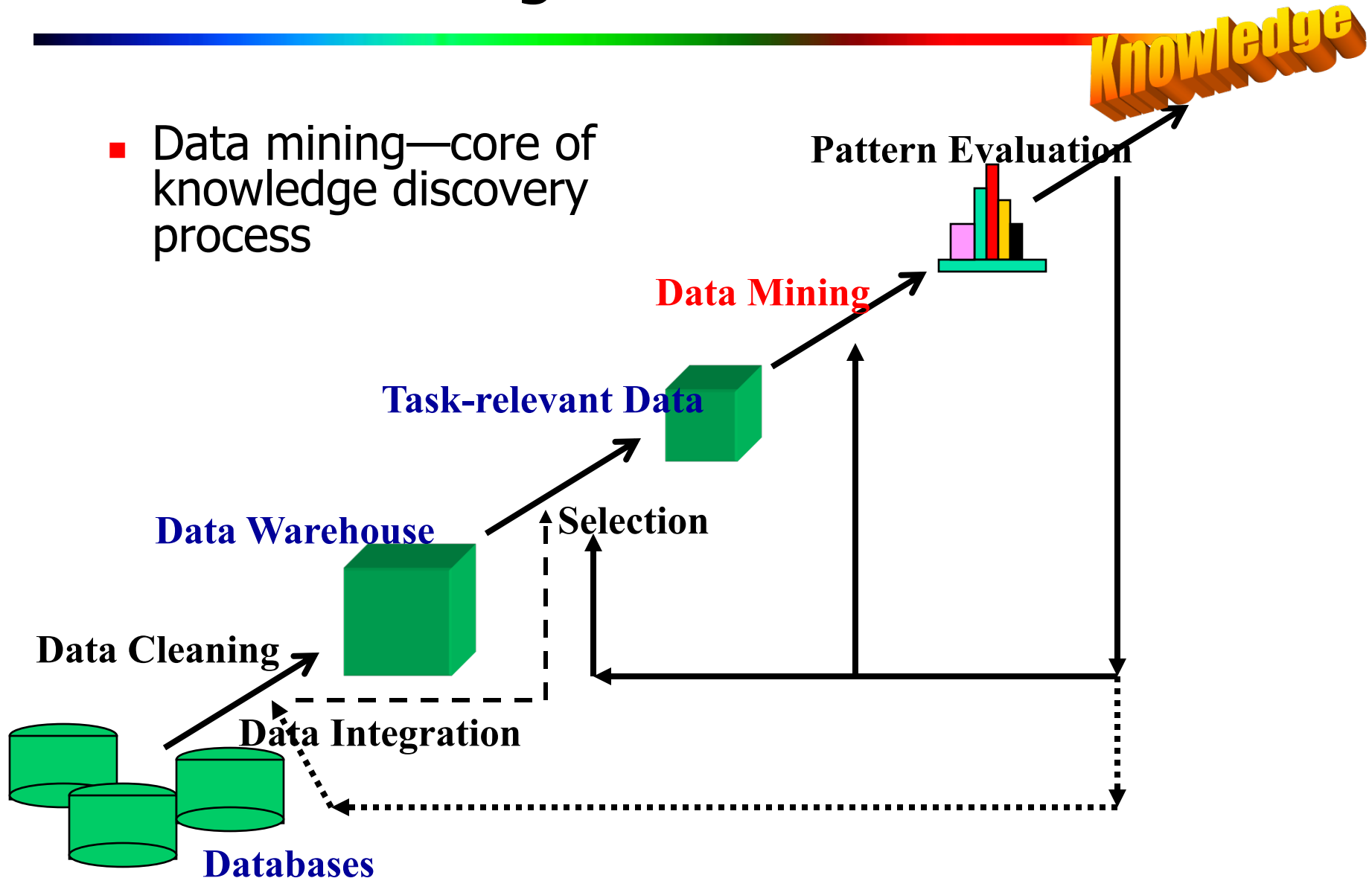
- Approaches: Clustering & model construction for frauds, outlier analysis
- Applications: Health care, retail, credit card service, telecomm.
 - Auto insurance: ring of collisions
 - Money laundering: suspicious monetary transactions
 - Medical insurance
 - Professional patients, ring of doctors, and ring of references
 - Unnecessary or correlated screening tests
 - Telecommunications: phone-call fraud
 - Phone call model: destination of the call, duration, time of day or week. Analyze patterns that deviate from an expected norm
 - Retail industry
 - Analysts estimate that 38% of retail shrink is due to dishonest employees
 - Anti-terrorism

Other Applications

- Sports
 - IBM Advanced Scout analyzed NBA game statistics (shots blocked, assists, and fouls) to gain competitive advantage for New York Knicks and Miami Heat
- Astronomy
 - JPL and the Palomar Observatory discovered 22 quasars with the help of data mining
- Internet Web Surf-Aid
 - IBM Surf-Aid applies data mining algorithms to Web access logs for market-related pages to discover customer preference and behavior pages, analyzing effectiveness of Web marketing, improving Web site organization, etc.

Data Mining: A KDD Process

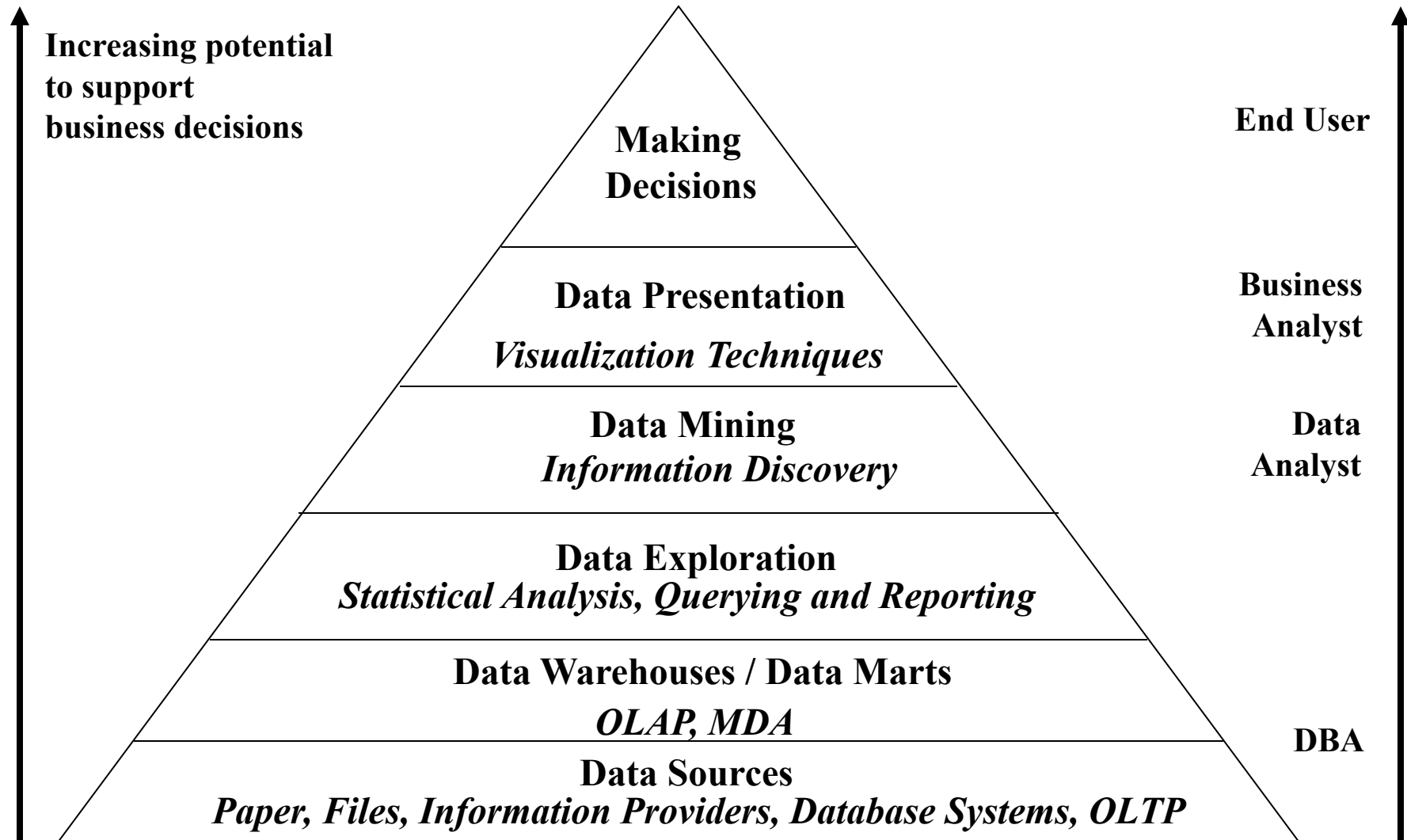
- Data mining—core of knowledge discovery process



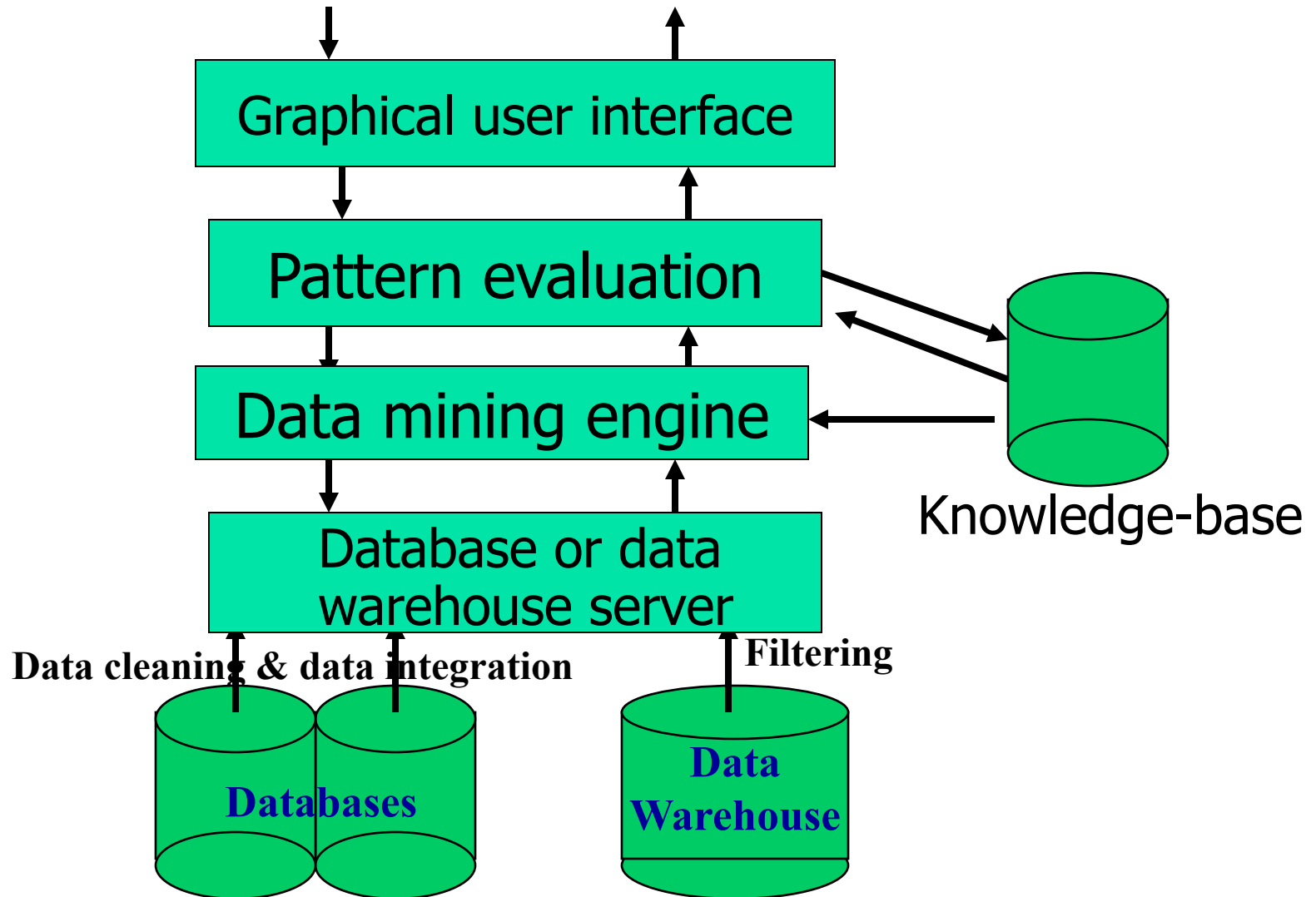
Steps of a KDD Process

- Learning the application domain
 - relevant prior knowledge and goals of application
- Creating a target data set: data selection
- Data cleaning and preprocessing: (*may take 60% of effort!*)
- Data reduction and transformation
 - Find useful features, dimensionality/variable reduction, invariant representation.
- Choosing functions of data mining
 - summarization, classification, regression, association, clustering.
- Choosing the mining algorithm(s)
- Data mining: search for patterns of interest
- Pattern evaluation and knowledge presentation
 - visualization, transformation, removing redundant patterns, etc.
- Use of discovered knowledge

Data Mining and *Business Intelligence*



Architecture: Typical Data Mining System



Data Mining: On What Kinds of Data?

- Relational database
- Data warehouse
- Transactional database
- Advanced database and information repository
 - Object-relational database
 - Spatial and temporal data
 - Time-series data
 - Stream data
 - Multimedia database
 - Heterogeneous and legacy database
 - Text databases & WWW

Data Mining Functionalities

- Concept description: Characterization and discrimination
 - Generalize, summarize, and contrast data characteristics, e.g., dry vs. wet regions
- Association (correlation and causality)
 - Diaper → Beer [0.5%, 75%]
- Classification and Prediction
 - Construct models (functions) that describe and distinguish classes or concepts for future prediction
 - E.g., classify countries based on climate, or classify cars based on gas mileage
 - Presentation: decision-tree, classification rule, neural network
 - Predict some unknown or missing numerical values

Data Mining Functionalities (2)

- Cluster analysis
 - Class label is unknown: Group data to form new classes, e.g., cluster houses to find distribution patterns
 - Maximizing intra-class similarity & minimizing interclass similarity
- Outlier analysis
 - Outlier: a data object that does not comply with the general behavior of the data
 - Noise or exception? No! useful in fraud detection, rare events analysis
- Trend and evolution analysis
 - Trend and deviation: regression analysis
 - Sequential pattern mining, periodicity analysis
 - Similarity-based analysis
- Other pattern-directed or statistical analyses

Are All the “Discovered” Patterns Interesting?



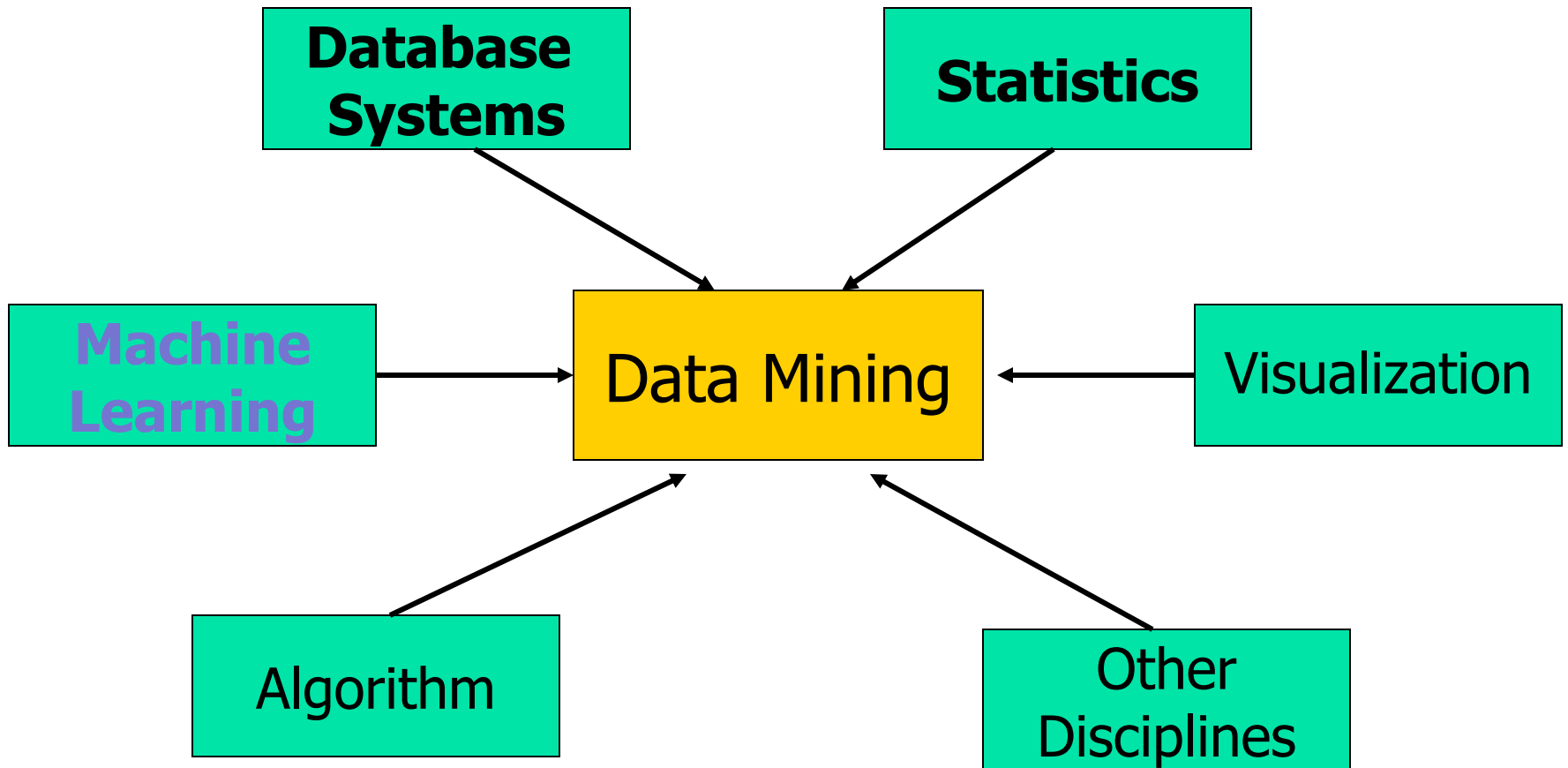
- Data mining may generate thousands of patterns: Not all of them are interesting
 - Suggested approach: Human-centered, query-based, focused mining
- **Interestingness measures**
 - A pattern is **interesting** if it is easily understood by humans, valid on new or test data with some degree of **certainty**, 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.

Can We Find All and Only Interesting Patterns?



- Find all the interesting patterns: [Completeness](#)
 - Can a data mining system find all 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 generate all the patterns and then filter out the uninteresting ones.
 - Generate only the interesting patterns—mining query optimization

Data Mining: Confluence of Multiple Disciplines



Data Mining: Classification Schemes

- General functionality
 - Descriptive data mining
 - Predictive data mining
- Different views, different classifications
 - Kinds of data to be mined
 - Kinds of knowledge to be discovered
 - Kinds of techniques utilized
 - Kinds of applications adapted

Multi-Dimensional View of Data Mining

■ **Data to be mined**

- Relational, data warehouse, transactional, stream, object-oriented/relational, active, spatial, time-series, text, multi-media, heterogeneous, legacy, WWW

■ **Knowledge to be mined**

- Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
- Multiple/integrated functions and mining at multiple levels

■ **Techniques utilized**

- Database-oriented, data warehouse (OLAP), machine learning, statistics, visualization, etc.

■ **Applications adapted**

- Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, Web mining, etc.

Major Issues in Data Mining

■ Mining methodology

- Mining different kinds of knowledge from diverse data types, e.g., bio, stream, Web
- Performance: efficiency, effectiveness, and scalability
- Pattern evaluation: the interestingness problem
- Incorporation of background knowledge
- Handling noise and incomplete data
- Parallel, distributed and incremental mining methods
- Integration of the discovered knowledge with existing one: knowledge fusion

■ User interaction

- Data mining query languages and ad-hoc mining
- Expression and visualization of data mining results
- Interactive mining of knowledge at multiple levels of abstraction

■ Applications and social impacts

- Domain-specific data mining & invisible data mining
- Protection of data security, integrity, and privacy

Summary

- Data mining: discovering interesting patterns from large amounts of data
- A natural evolution of database technology, in great demand, with wide applications
- A KDD process includes data cleaning, data integration, data selection, transformation, data mining, pattern evaluation, and knowledge presentation
- Mining can be performed in a variety of information repositories
- Data mining functionalities: characterization, discrimination, association, classification, clustering, outlier and trend analysis, etc.
- Data mining systems and architectures
- Major issues in data mining

A Brief History of Data Mining Society



- 1989 IJCAI Workshop on Knowledge Discovery in Databases (Piatetsky-Shapiro)
 - Knowledge Discovery in Databases (G. Piatetsky-Shapiro and W. Frawley, 1991)
- 1991-1994 Workshops on Knowledge Discovery in Databases
 - Advances in Knowledge Discovery and Data Mining (U. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, 1996)
- 1995-1998 International Conferences on Knowledge Discovery in Databases and Data Mining (KDD'95-98)
 - Journal of Data Mining and Knowledge Discovery (1997)
- 1998 ACM SIGKDD, SIGKDD'1999-2001 conferences, and SIGKDD Explorations
- More conferences on data mining
 - PAKDD (1997), PKDD (1997), SIAM-Data Mining (2001), (IEEE) ICDM (2001), etc.

1. [DBLP](#)
2. [Google](#)
3. [Citeseer](#)
4. [DL@lib](#)

Where to Find References?

■ Data mining and KDD

- Conferences: ACM-SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, PAKDD, etc.
- Journal: Data Mining and Knowledge Discovery, KDD Explorations

■ Database systems

- Conferences: ACM-SIGMOD, ACM-PODS, VLDB, IEEE-ICDE, EDBT, ICDT, DASFAA
- Journals: ACM-TODS, IEEE-TKDE, JIIS, J. ACM, VLDBJ, etc.

■ AI & Machine Learning

- Conferences: Machine learning (ML), AAAI, IJCAI, COLT (Learning Theory), etc.
- Journals: Machine Learning, Artificial Intelligence, etc.

■ Statistics

- Conferences: Joint Stat. Meeting, etc.
- Journals: Annals of statistics, etc.

■ Visualization

- Conference proceedings: CHI, ACM-SIGGraph, etc.
- Journals: IEEE Trans. visualization and computer graphics, etc.

Recommended Reference Books



- I. H. Witten and E. Frank, Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations, Morgan Kaufmann, 2001
- **C. C. Aggarwal, Data Mining: The Textbook, Springer, 2015□□**
- **J. Leskovec, A. Rajaraman, and J. Ullman, Mining of Massive Datasets (v2.1), Cambridge University Press, 2014.**
- Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin, Learning From Data. AMLBook, 2012.
- **J. Han and M. Kamber. Data Mining: Concepts and Techniques. Morgan Kaufmann, 2001**
- **D. J. Hand, H. Mannila, and P. Smyth, Principles of Data Mining, MIT Press, 2001**
- **T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Springer-Verlag, 2001**
- **T. M. Mitchell, Machine Learning, McGraw Hill, 1997**
- **P-N. Tan, M. Steinbach, and V. Kumar, Introduction to Data Mining,. Addison-Wesley, 2005**
- S. M. Weiss and N. Indurkha, Predictive Data Mining, Morgan Kaufmann, 1998

Jai's Project (COMP9318, 2016s2)

■ Problem

- <http://kentandlime.com.au/>, a startup company helping male customers to stay in fashion but out of the shops.
- Status-quo:
 - Ask questions, and stylists makes a list of recommended items, and send them to customers
 - If happy, customers pay for the product.
- Recommendation is the key!

■ Challenges

- Dirty data
- Not an easy/typical recommendation system settings
- Customer feedbacks
- Real-time recommendations

<http://www.news.com.au/lifestyle/fashion/fashion-trends/fashions-most-unlikely-trend-would-you-buy-clothes-chosen-for-you/news-story/8634b5f06f608b9f2fd7c27758f9c10a>

Solutions - Highlight

- Use domain-knowledge and quick evaluations to guide the whole process
- Data preprocessing
 - Data source: CRM (profile) + NoSQL DB (transactions)
 - Missing data: e.g., due to schema changes
 - Data normalization: $A's\ XL = B's\ L$
 - Data noise: k-means / binning
 - Data selection: remove sparse columns/rows
- Feature engineering
 - weight-to-height ratio

Solutions – Highlight /2

- Product class clustering and prediction
- Collaborative filtering with smoothing and weighting
- Content-based recommendation (solve the cold start problem)
- Incorporate customer feedbacks
- Association rule mining:
 - LSShirts_1, Shorts_2 → Socks_3
- Ensemble of the above

- Plus many engineering efforts

Results

- Test set:
 - Classification rate: 74%, on par with humans
- Deployed to production on 18-24 Nov 2016:
 - Customers rejecting on average 2.36 items out of a basket of 10-12 items → (76.4%, 80.3%)
 - Latency: 2.3s
- Future work identified
 - e.g., seasonality
- Check Jai's presentation slides for more details.