# *Data Mining*

## George Karypis
Department of Computer Science
Digital Technology Center
University of Minnesota, Minneapolis, USA.
http://www.cs.umn.edu/~karypis
karypis@cs.umn.edu

# *Overview*

- Data-mining
  - What is it and why do we care?
- Commercial & Scientific Applications
  - What can it do for me?
- Ongoing Research Activities
  - What is George doing?
- From Research to Technology Transfer
  - How can we help you?

# *Problem*

> **We have lots of data!**

> **We have little information!**

> **We have no knowledge!**

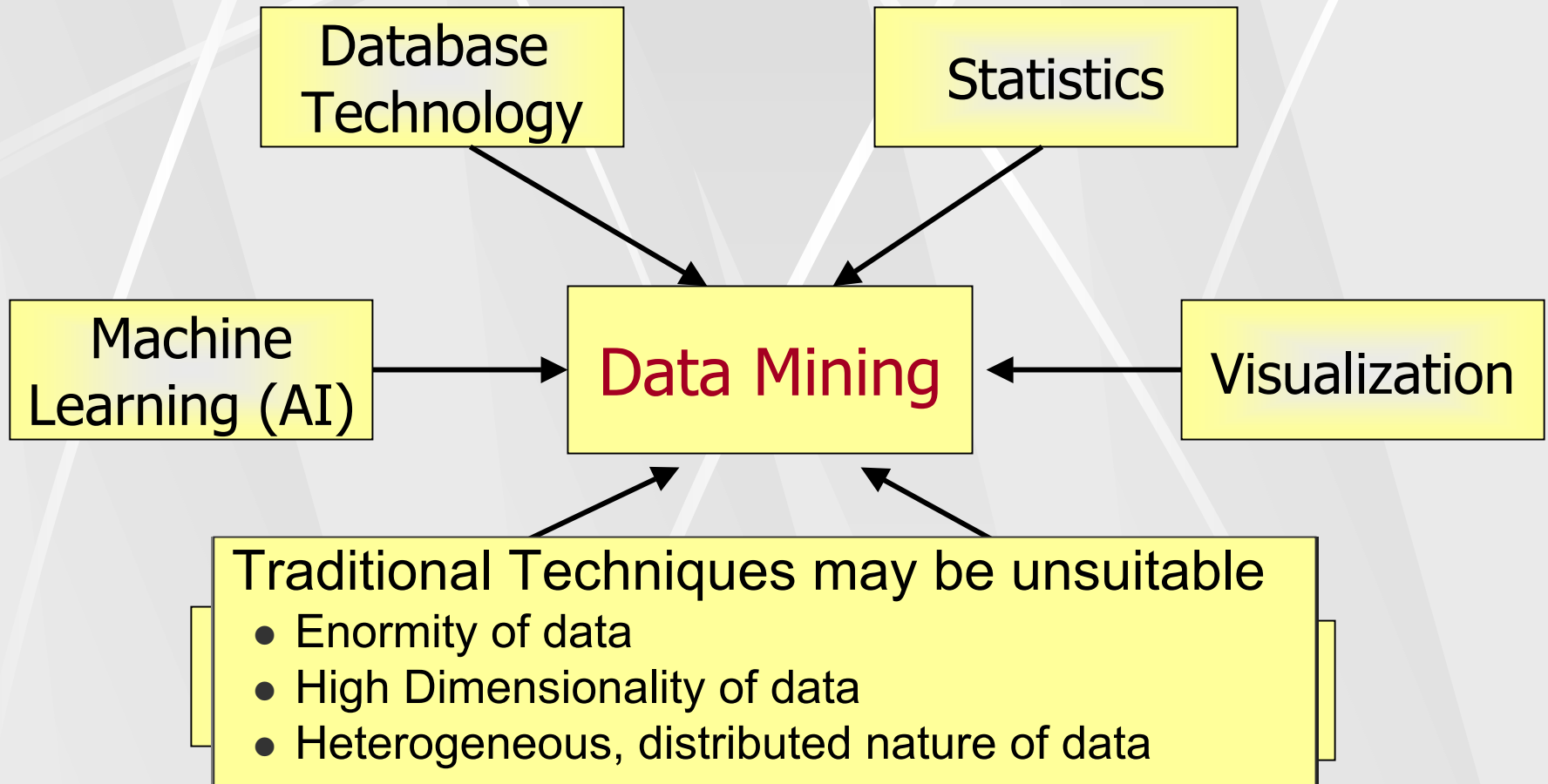# *What is Data Mining?*

Many Definitions…

- A short one…

  Search for Valuable Information in Large Volumes of Data.

- A long one…

  Exploration & Analysis, by Automatic or Semi-Automatic Means, of Large Quantities of Data in order to Discover Meaningful Patterns & Rules.

**DIGITAL TECHNOLGY CENTER**

# Origins of Data Mining



**Database Technology**

**Statistics**

**Machine Learning (AI)**

**Data Mining**

**Visualization**

**Traditional Techniques may be unsuitable**
- Enormity of data
- High Dimensionality of data
- Heterogeneous, distributed nature of data

# A Brief History of Data Mining Activities

- <u>1989 IJCAI Workshop on Knowledge Discovery in Databases</u>
  - Knowledge Discovery in Databases (G. Piatetsky-Shapiro and W. Frawley, 1991)
- <u>1991-1994 Workshops on Knowledge Discovery in Databases</u>
  - Advances in Knowledge Discovery and Data Mining (U. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, 1996)
- <u>1995-1998 International Conferences on Knowledge Discovery in Databases and Data Mining (KDD'95-98)</u>
  - Journal of Data Mining and Knowledge Discovery (1997)
- <u>1998 ACM SIGKDD, SIGKDD'1999-2003 conferences, and SIGKDD Explorations</u>
- <u>More conferences on data mining</u>
  - PAKDD, PKDD, SIAM-Data Mining, (IEEE) ICDM, DaWaK, SPIE-DM, etc.

**DIGITAL TECHNOLGY CENTER**

# *Why Mine Data?  Why Now?*

- Lots of data is being produced, collected, and warehoused.

- Computing has become affordable.

- Competitive pressures are strong
  - Provide better, customized services for an *edge*.
  - Information is becoming *product* in its own right.

- Data mining has become an integral part of modern CRM.

# *Data Mining Tasks*

- ## Predictive Tasks
  - Use some variables to predict unknown or future values of other variables.


- ## Descriptive Tasks
  - Find human-interpretable patterns that describe the data.

From [Fayyad, et.al.] Advances in Knowledge Discovery and Data Mining, 1996

# *What Problems Can Data Mining Solve?*

- Better/Effective/Personalized/Real-time Marketing
  - Identify the subset of customers that is most likely to buy products from a particular catalog.
  - Create a product catalog that will lead to the highest profit.
  - Identify which of the users that browse my website are most likely to purchase something.

- Fraud/Anomaly Detection
  - Identify a fraudulent credit card transaction given the past transactions of a particular customer.

- Customer Attrition/Churn
  - Identify my customers that are most likely to be lost to a competitor.
  - Identify any actions that I can take in order to retain them.

- Effective Information Filtering/Compression/Navigation
  - Find today's news articles that I will like to read.
  - Help me find what I'm looking on the web.
  - Please organize my hard-drive (mail folders, bookmarks, contacts, etc).

# *Ongoing Research Activities*

- Customer segmentation
- Marketing campaigns
- Recommender systems
- Document categorization & clustering
- Meta-search engines
- Analysis of web-browsing behavior
- Discovery of complex patterns
- Finding patterns in relational data
- Mining scientific and biological databases

# *Marketing Campaigns*

- Problem:
  - Identify the set of customers that are most likely to buy a set of products.
- Solution:
  - Developed prediction algorithms based on past purchasing and demographic information.

- Problem
  - Create a set of product catalogs and for each catalog identify the subset of customers that it should be mailed to.
- Solution:
  - Developed algorithms to cluster the customers based on predicted future purchases and create the catalogs by analyzing the characteristics of each cluster.

# Recommender Systems

- Problem: Information Overload!



Recommender systems help us identify worthwhile stuff!

- Filter articles that we will like.
- Predict how much we will like a particular book or movie.
- Recommend the top-$N$ products that we will most likely buy.

# Recommender Systems (cont)

- Developed *scalable* item-based collaborative filtering-based approaches for rating prediction and top-*N* recommendations.

- Collaborative Filtering
  - Key insight—No person is an island!
    - Each individual is a member of a group(s) &
    - The group's collective knowledge can help filter the information!

# Document Clustering

- Problem:
  - As the amount of textual information increases, there is a need to automatically organize them into meaningful groups and hierarchies, and provide effective summaries.
    - How do you navigate through the 1000+ results that comes back from a Google query?

- Solution:
  - Document Clustering!

# Document Clustering (cont)

Developed scalable document clustering algorithms and effective cluster summarization approaches.
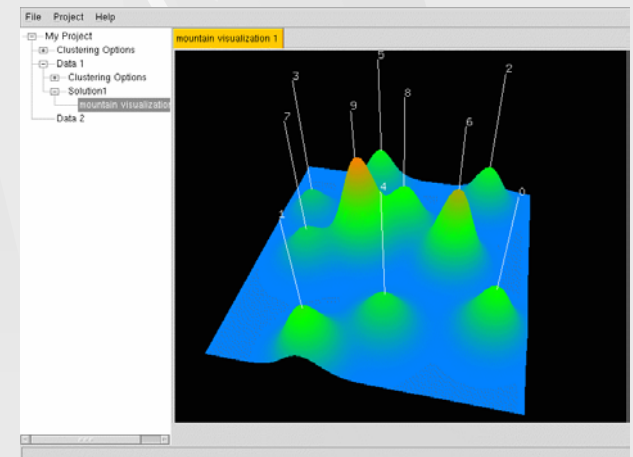
*Search to Browse…*

# Meta-Search Engines

# *From Research to Technology Transfer*

- We put a considerable effort to provide industrial-strength implementations of the various algorithms that we are developing.

- Currently available software tools:
  - CLUTO for clustering
    - http://www.cs.umn.edu/~cluto
  - PAFI for frequent pattern discovery
    - http://www.cs.umn.edu/~pafi
  - SUGGEST for top-*N* recommendation
    - http://www.cs.umn.edu/~karypis/suggest

- These tools are used extensively by various academic, government, and commercial entities.

# CLUTO

- Clustering Algorithms
  - High-performance & High-quality partitional clustering
  - High-quality agglomerative clustering
  - High-quality graph-partitioning-based clustering
  - Hybrid partitional & agglomerative algorithms for building trees for very large datasets.
- Cluster Analysis Tools
  - Cluster signature identification
  - Cluster organization identification
- Visualization Tools
  - Hierarchical Trees
  - High-dimensional datasets
  - Cluster relations
- Interfaces
  - Stand-alone programs
  - Library with a fully published API
- Available on Windows, Sun, and Linux

http://www.cs.umn.edu/~cluto

# *Thank you*

## George Karypis
Department of Computer Science
Digital Technology Center
University of Minnesota, Minneapolis, USA.
http://www.cs.umn.edu/~karypis
karypis@cs.umn.edu