# CSE4334/5334   Data Mining

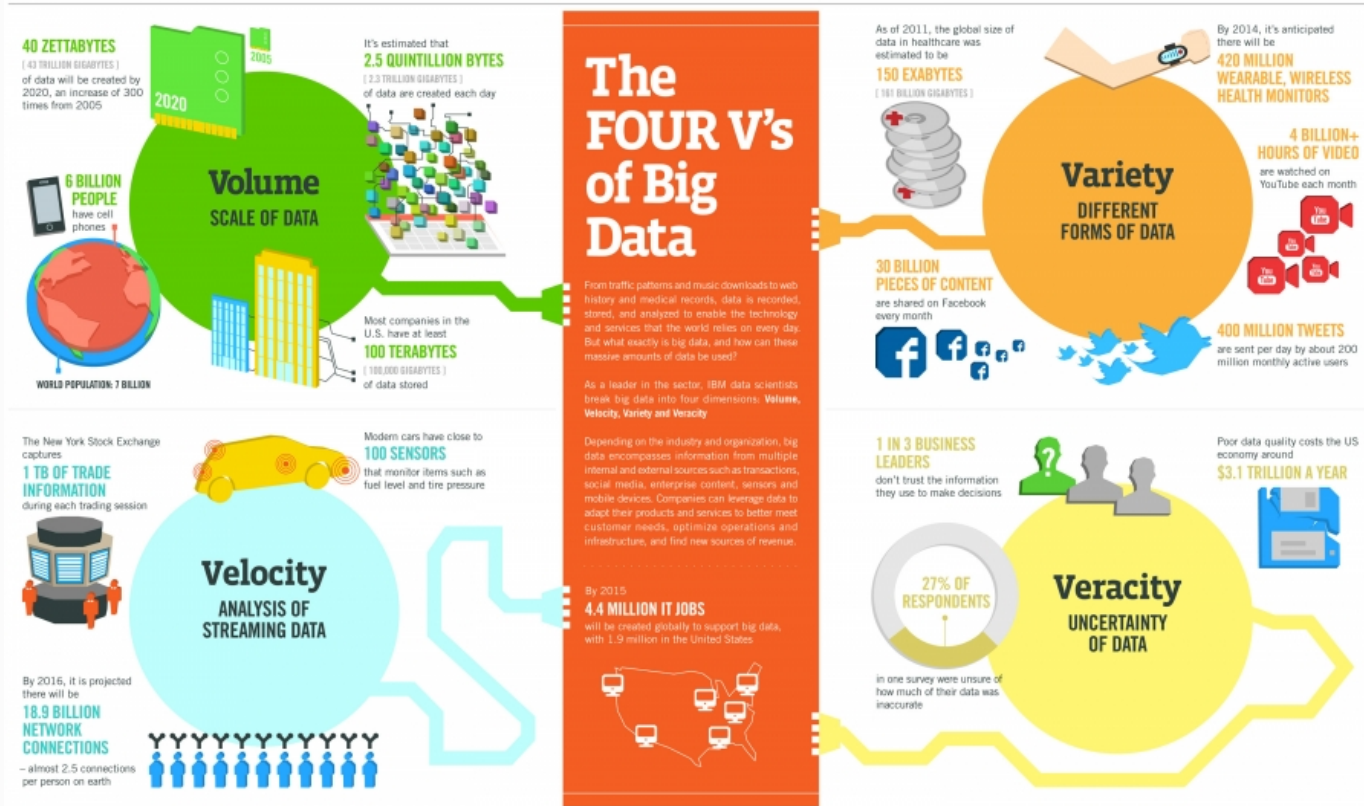# Overview of Data Mining

Chengkai Li

Department of Computer Science and Engineering

University of Texas at Arlington

Fall 2020      (Slides partly courtesy of Pang-Ning Tan, Michael Steinbach and Vipin Kumar, and Jiawei Han, Micheline Kamber and Jian Pei)

# Big Data



http://www.ibmbigdatahub.com/infographic/four-vs-big-data

# Big Data

## The 4 Vs

- Volume
- Variety
- Velocity
- Veracity

# Volume: How much data is out there?

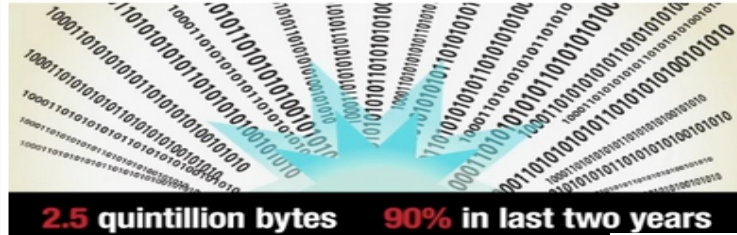## Every Day We Create 2.5 Quintillion Bytes of Data

*IBM study of 1,734 chief marketing officers from 64 countries*

This is a Press Release edited by StorageNewsletter.com on 2011.10.21

BOOKMARK

http://www.sciencedaily.com/releases/2013/05/130522085217.htm

A new IBM Corp.'s study of more than 1,700 chief marketing officers from 64 countries and 19 industries reveals that the majority of the world's top marketing executives recognize a critical and permanent shift occurring in the way they engage with their customers, but question whether their marketing organizations are prepared to manage the change.

**2.5 quintillion bytes    90% in last two years**

## Big Data, for better or worse: 90% of world's data generated over last two years

Date: May 22, 2013

Source: SINTEF

Summary: A full 90 percent of all the data in the world has been generated over the last two years. Internet-based companies are awash with data that can be grouped and utilized. Is this a good thing?

**Share This**

> ✉ Email to a friend
> f Facebook
> 🐦 Twitter
> in LinkedIn
> 8+ Google+
> 🖨 Print this page

http://www.storagenewsletter.com/rubriques/market-reportsresearch/ibm-cmo-study/

# Volume: How much data is out there?

**In total, 2.7 Zettabytes of data exists in our digital universe.**
("A terabyte is equal to 1,024 gigabytes. A petabyte is equal to 1,024 terabytes.
An exabyte is equal to 1,024 petabytes. A zettabyte is equal to 1,024 exabytes.")

**Online Activity**
- Every minute:
  - 149, 513 emails are sent
  - 3.3 million Facebook posts are created
  - 3.8 million Google searches are performed
  - 65,972 Instagram photos are uploaded.
  - 448,800 tweets are constructed
  - 500 hours of YouTube videos are uploaded

https://www.nodegraph.se/big-data-facts/

# Variety: Types of Data

Structured data
- o (relational) database tables
- o CSV/TSV files

Semi-structured data
- o XML, JSON, RDF

Unstructured data
- o text data (documents, Web pages, short texts, e.g., social media)

Multimedia data
- o images, videos, audios

Other types of data
- o matrices, graphs, sequences, time-series, spatio-temporal

# Velocity: Streaming Data

- ❖ Stock trades
- ❖ Highway sensors
- ❖ Weather data
- ❖ Social media
- ❖ Telephone calls
- ❖ Video streaming

http://mashable.com/2012/06/22/data-created-every-minute/

7

# Veracity: uncertain and imprecise data

- ❖ Biases
- ❖ Data Lineage
- ❖ Bugs
- ❖ Noise
- ❖ Abnormalities
- ❖ Information security
- ❖ Unreliable sources
- ❖ Falsification
- ❖ Uncertainty
- ❖ Out of date
- ❖ Human error



https://simplicable.com/new/data-veracity

# Data in Every Application Area

o   Business: e-commerce, transactions (retailers, banking, credit cards), ratings, reviews, stock trading, …

o   Web, social media (YouTube, Flickr, …), and social networks (Facebook, Twitter, …)

o   News

o   Science: bioinformatics, scientific experiments, environment, climate, astronomy

o   Logs and measurements

o   Personal information: emails, calendars, digital photos, videos

o   Transportation

o   Telecommunication

o   Education

o   Entertainment (film, music, gaming, …)

o   Sports

o   Health care

o   Crime, security

# What is Data Mining?

Data mining (knowledge discovery from data)
- o Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
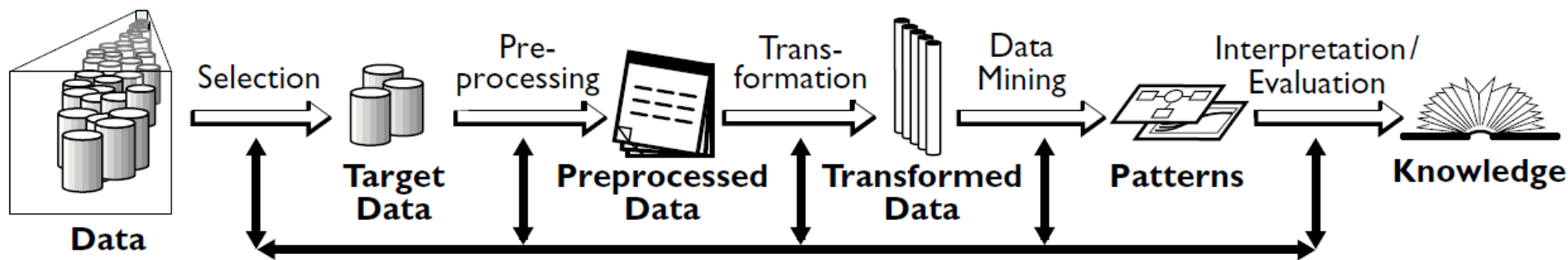
What is not Data Mining?
- o Retrieve data instead of knowledge or pattern
- o Not interesting (trivial, explicit, known, useless)
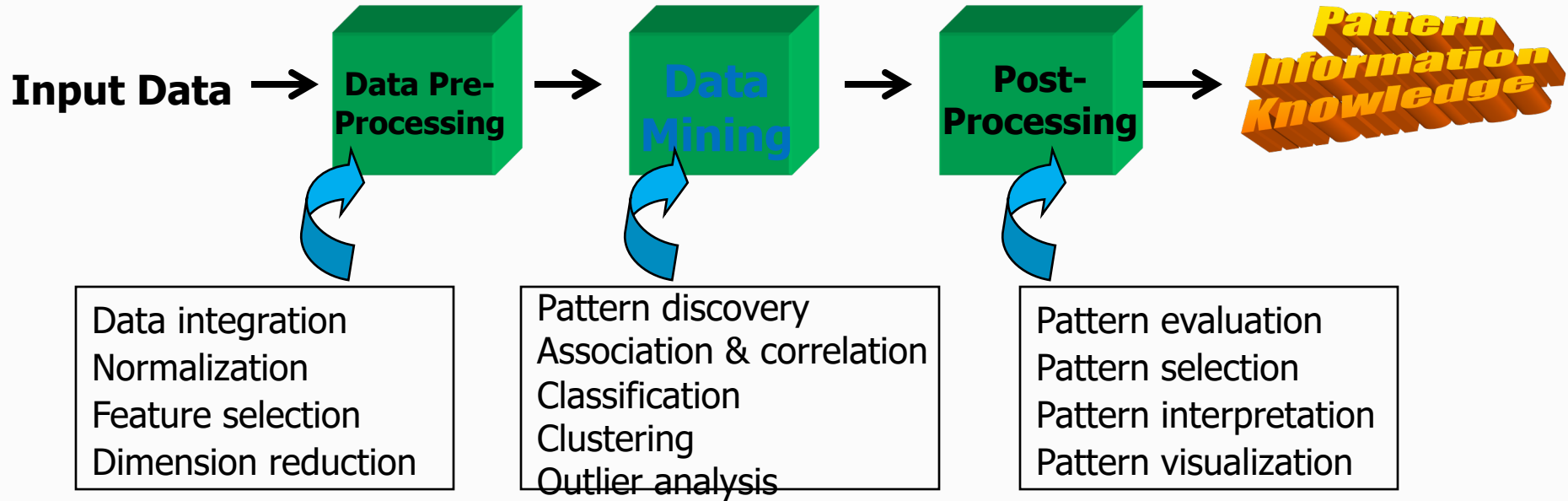
# Knowledge Discovery (KDD) Process

❖ Data mining plays an essential role in the knowledge discovery process
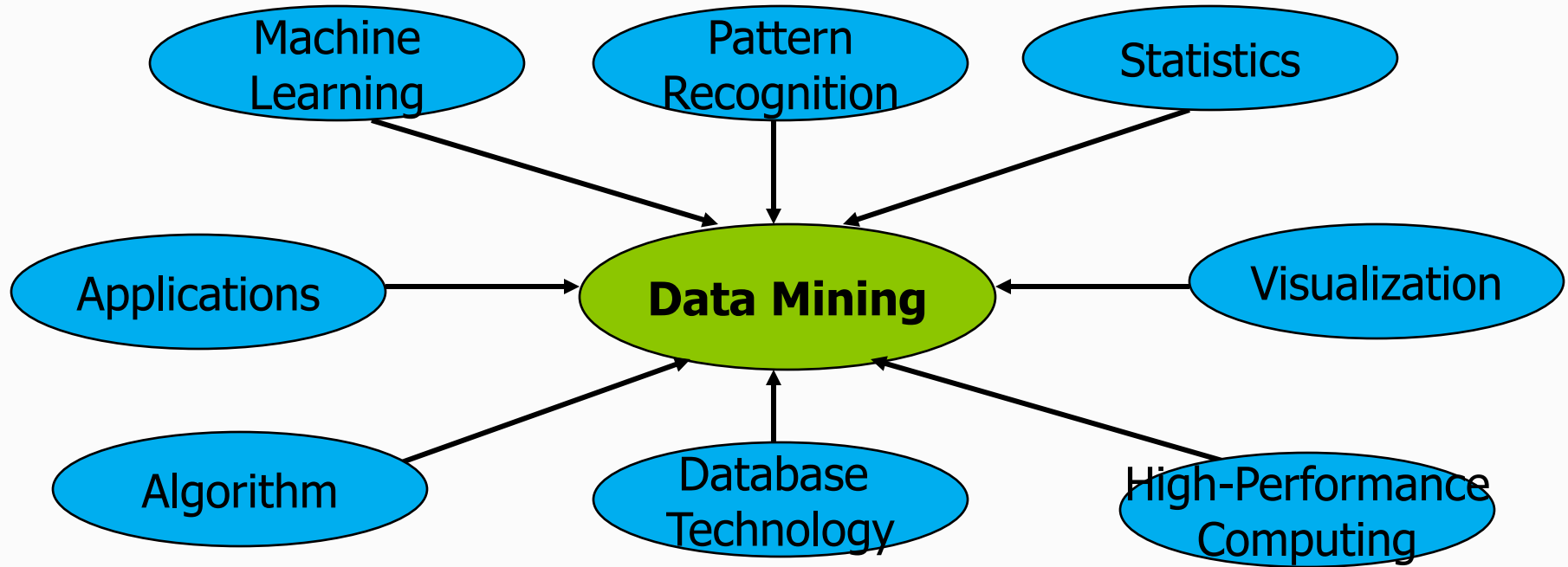


http://cacm.acm.org/magazines/1996/11/8517-the-kdd-process-for-extracting-useful-knowledge-from-volumes-of-data/abstract

# KDD Process: A Typical View from ML and Statistics

**Input Data** → Data Pre-Processing → Data Mining → Post-Processing → *Pattern Information Knowledge*

| Data integration<br>Normalization<br>Feature selection<br>Dimension reduction | Pattern discovery<br>Association & correlation<br>Classification<br>Clustering<br>Outlier analysis | Pattern evaluation<br>Pattern selection<br>Pattern interpretation<br>Pattern visualization |
|---|---|---|

This is a view from typical machine learning and statistics communities

# Data Mining: Confluence of Multiple Disciplines

# Data Mining Software

**Free, open-source**

- RapidMiner
- Weka: Data mining tool in java
- SCaVis: scientific computation and visualization, Java
- Orange: Python suite
- Scikit-learn: Python machine learning lbirary
- NumPy/SciPy/Ipython/ mlpy (python modules for scientific computing, scientific library, interactive computing, machine learning)
- R: statistical computing and graphic
- RattleGUI: data mining GUI using R
- Octave: numerical analysis
- Shogun: machine learning toolkit in C++

**Text Mining Tools**

- NLTK (NLP Toolkit): NLP suite for Python
- SenticNet API: sentiment analysis
- Stanford NLP software
- UIMA

**Large-Scale Data Processing, Machine Learning**

- Apache Mahout
- GraphLab
- MapReduce/Hadoop
- Spark
- Pregel/Giraph

**Commercial Products**

- Matlab
- Oracle Data Mining
- SAS
- IBM SPSS
- Microsoft SQL Server Analysis Services
- HP Vertica