# Data and Data Mining

CSE 4334 / 5334 Data Mining
Spring 2019

## Won Hwa Kim

(Slides courtesy of Pang-Ning Tan, Michael Steinbach and Vipin Kumar)
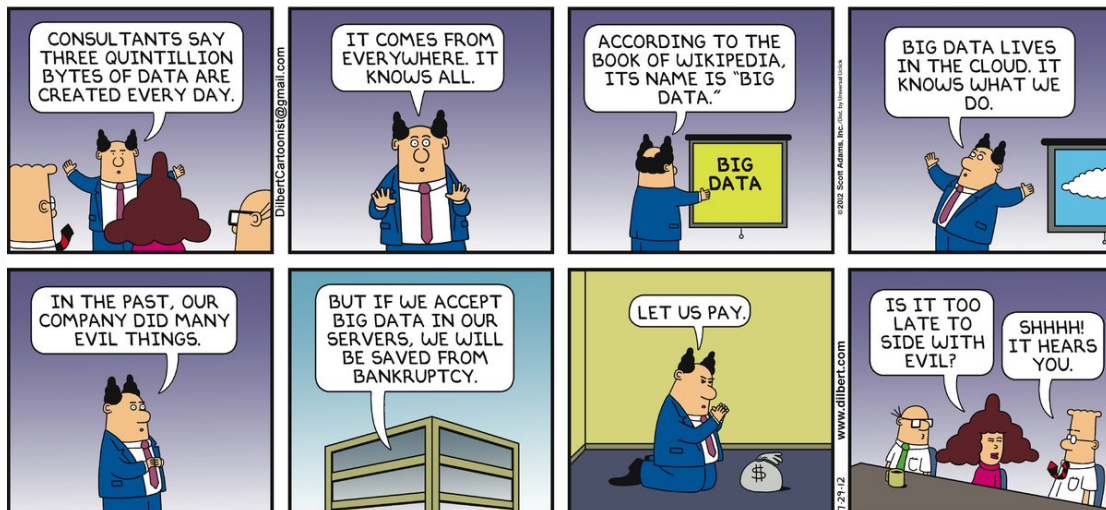
---

## Big Data



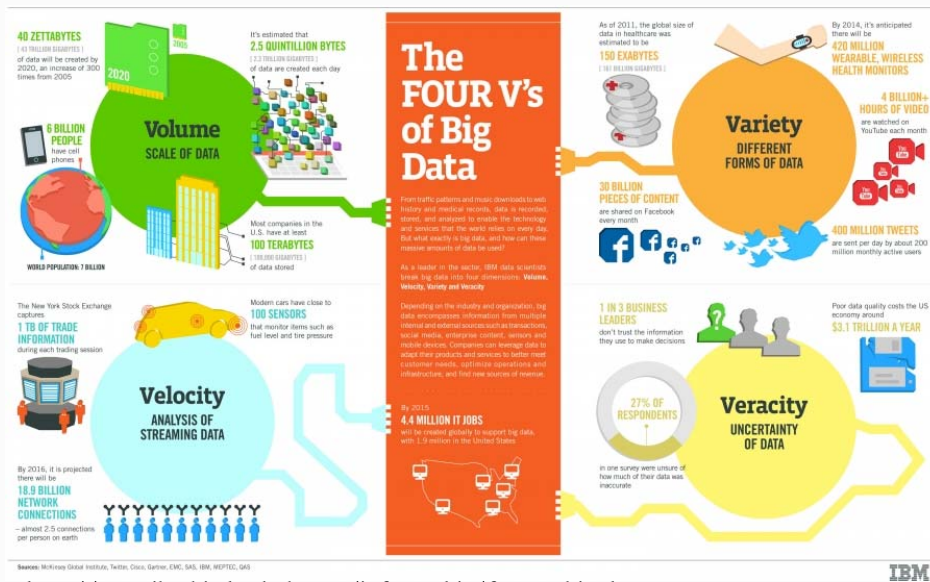http://dilbert.com/strip/2012-07-29

2

# Big Data



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

3

# Big Data

## The 4 Vs

- o Volume
- o Variety
- o Velocity
- o Veracity

4

# Volume: How much data is out there?
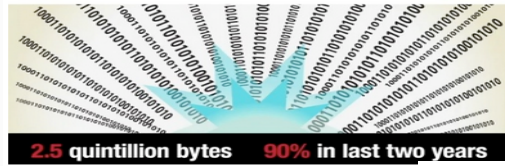
**Every Day We Create 2.5 Quintillion Bytes of Data**
*IBM study of 1,734 chief marketing officers from 64 countries*

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

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
> Google+
> Print this page

5

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

# 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

6

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# 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**

Copyright ©2007-2017 The University of

# Veracity: uncertain and imprecise data

❖ Quality and origin of data
❖ Consistent? Complete? Integrity?
❖ Untrusted and Uncleaned
❖ Fake stories

❖ Lots of cost to justify the data…

**8**

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Datasets

- ❖ Amazon Public Data Sets
- ❖ Data.gov
- ❖ Linked Open Data, Knowledge Bases, Encyclopedia
- ❖ Yahoo! Webscope
- ❖ Stanford Large Network Dataset Collection
- ❖ UCI Machine Learning Repository
- ❖ UCR Time Series Classification/Clustering
- ❖ Time Series Data Library　　http://robjhyndman.com/TSDL/
- ❖ KDnuggets Dataset List　　http://www.kdnuggets.com/datasets/index.html
- ❖ KDD Cup Datasets　　http://www.sigkdd.org/kddcup/index.php

**9**

# Amazon Public Data Sets

http://aws.amazon.com/public-data-sets/

o   NASA NEX: A collection of Earth science data sets maintained by NASA, including climate change projections and satellite images of the Earth's surface

o   Common Crawl Corpus: A corpus of web crawl data composed of over 5 billion web pages

o   1000 Genomes Project: A detailed map of human genetic variation

o   Google Books Ngrams: A data set containing Google Books n-gram corpuses

o   US Census Data: US demographic data from 1980, 1990, and 2000 US Censuses

o   Freebase Data Dump: A data dump of all the current facts and assertions in the Freebase system, an open database covering millions of topics

**10**

# Data.gov

http://www.data.gov/ (137,608 datasets)

o   Consumer Complaint Database
o   U.S. International Trade in Goods and Services: Monthly report that provides national trade data including imports, exports, and balance of payments for goods and services.
o   DTV Reception Maps
o   Food Access Research Atlas — presents a spatial overview of food access indicators for low-income and other census tracts using different measures of supermarket...
o   U.S. Hourly Precipitation Data
o   Great Chile Earthquake of May 22, 1960
o   Consumer Expenditure Survey
o   Farmers Markets Geographic Data: longitude and latitude, state, address, name, and zip code of Farmers Markets in the United States
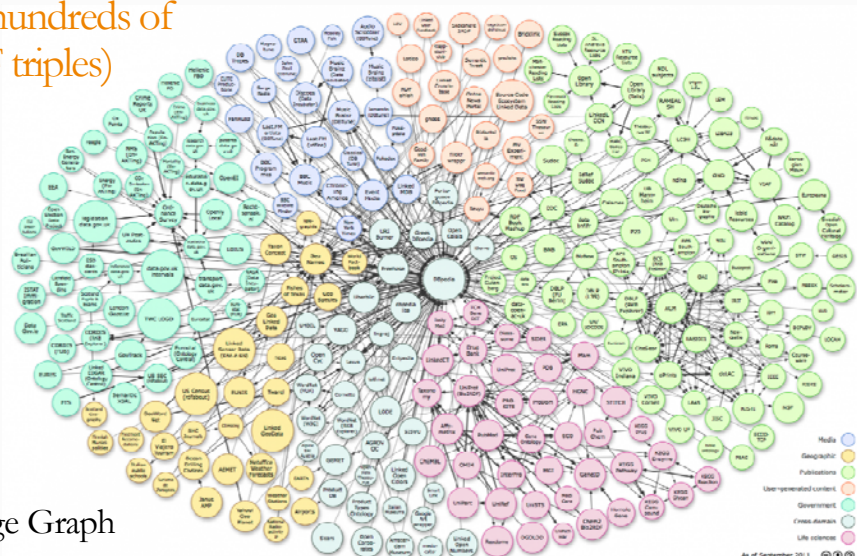o   Crimes - 2001 to present (City of Chicago)

11

# Linked Data, Knowledge Bases, Encyclopedia

http://linkeddata.org/ (hundreds of datasets, billions of RDF triples)

IMDB
DBLP
PubMed
Wikipedia, DBpedia
YAGO
Freebase/Google Knowledge Graph

12

# Stanford Large Network Dataset Collection

http://snap.stanford.edu/data/
o    Social networks : online social networks, edges represent interactions between people
o    Communication networks : email communication networks with edges representing communication
o    Citation networks : nodes represent papers, edges represent citations
o    Collaboration networks : nodes represent scientists, edges represent collaborations (co-authoring a paper)
o    Web graphs : nodes represent webpages and edges are hyperlinks
o    Amazon networks : nodes represent products and edges link commonly co-purchased products
o    Internet networks : nodes represent computers and edges communication
o    Road networks : nodes represent intersections and edges roads connecting the intersections

**13**

# 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
**14**    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
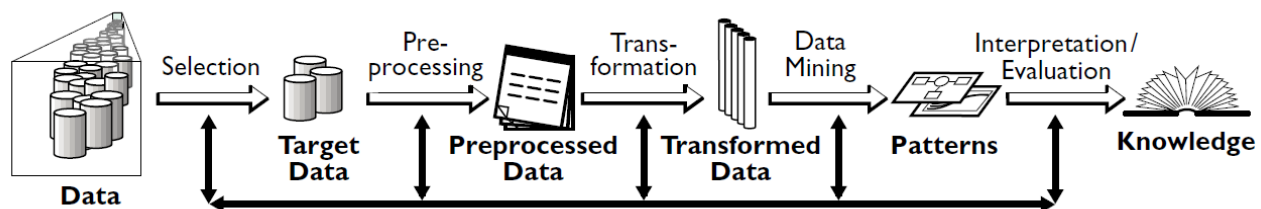  o Not interesting (trivial, explicit, known, useless)

**15**

# Knowledge Discovery (KDD) Process

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



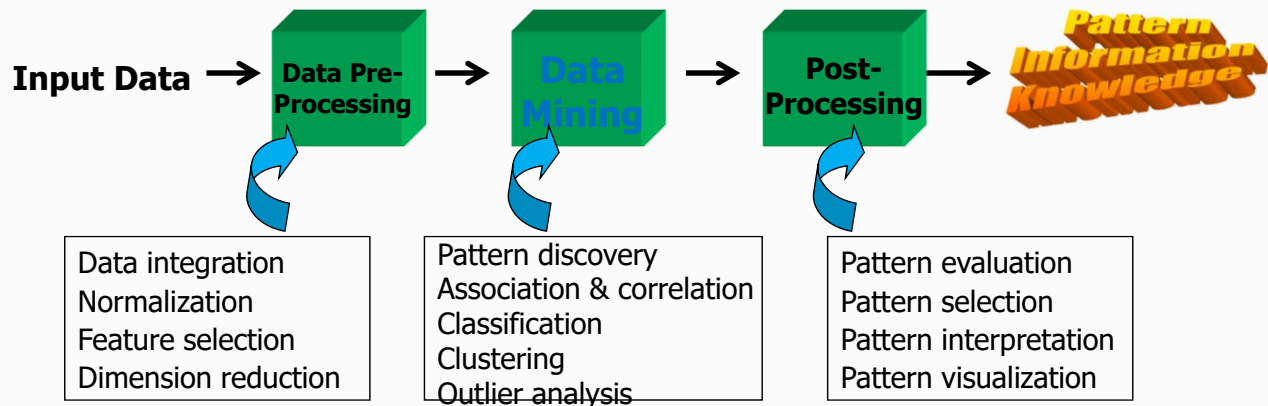http://cacm.acm.org/magazines/1996/11/8517-the-kdd-process-for-extracting-useful-knowledge-from-volumes-of-data/abstract

**16**

## 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*

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Data Mining: Confluence of Multiple Disciplines

Machine Learning · Pattern Recognition · Statistics

Applications → **Data Mining** ← Visualization

Algorithm · Database Technology · High-Performance Computing

18

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# 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

**19**