

Lesson A-1

Introduction to Data Mining

Topics

- Introduction to Data Mining/Big Data
- Data mining project development process

POPULAR SCIENCE

THE
FUTURE
NOW

THE CONTROL CENTERS

Using Data to Feed the World,
Solve Cold Cases, Battle Malware,
Predict Our Fate p.52

OFFICER ALGORITHM

Can a Crime Be Prevented
Before It Begins? p.38

NEW WAYS OF SEEING

A Gallery of
Extraordinary
Infographics p.69

SPECIAL ISSUE

DATA IS POWER

HOW INFORMATION
IS DRIVING
THE FUTURE

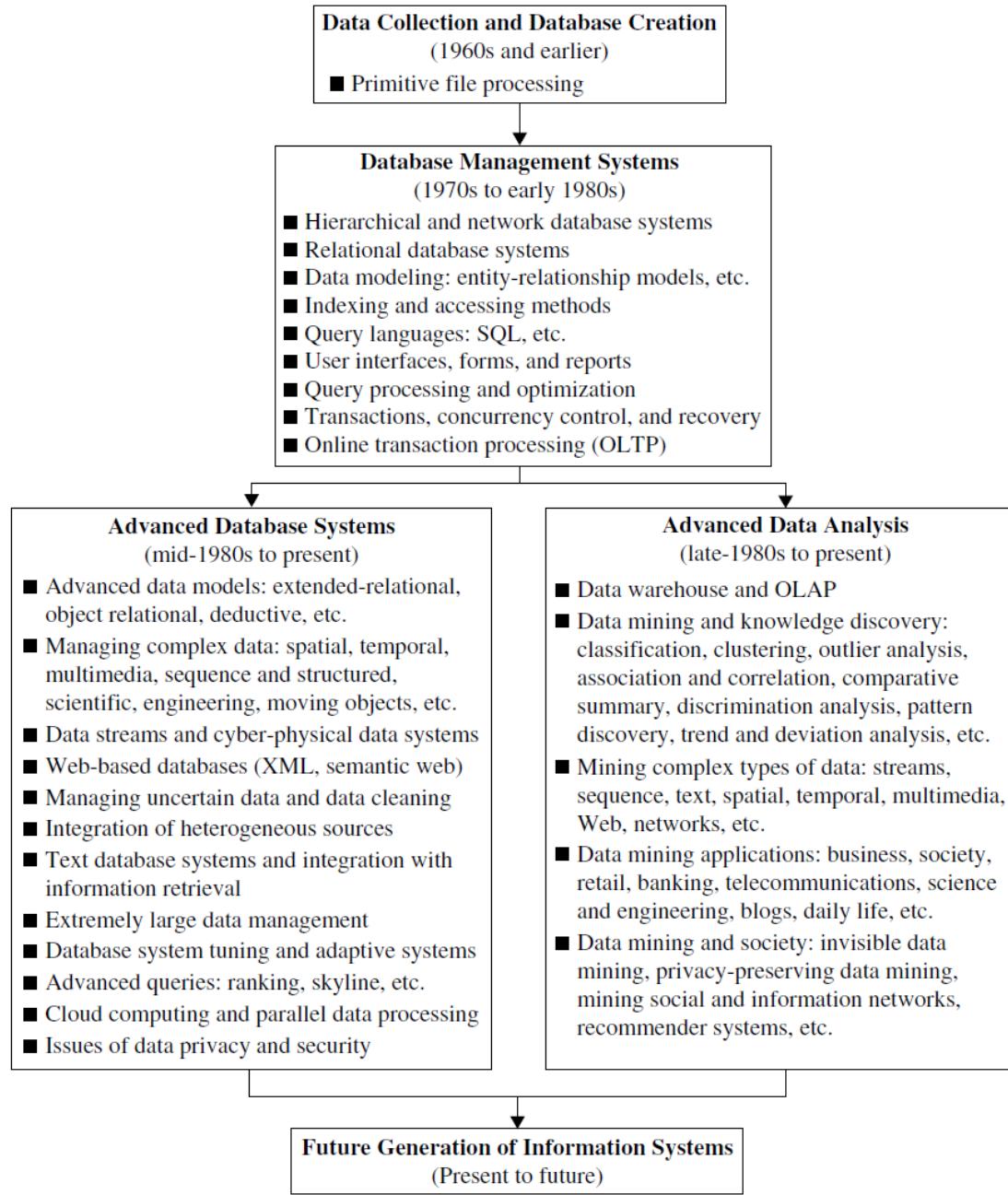
PLUS

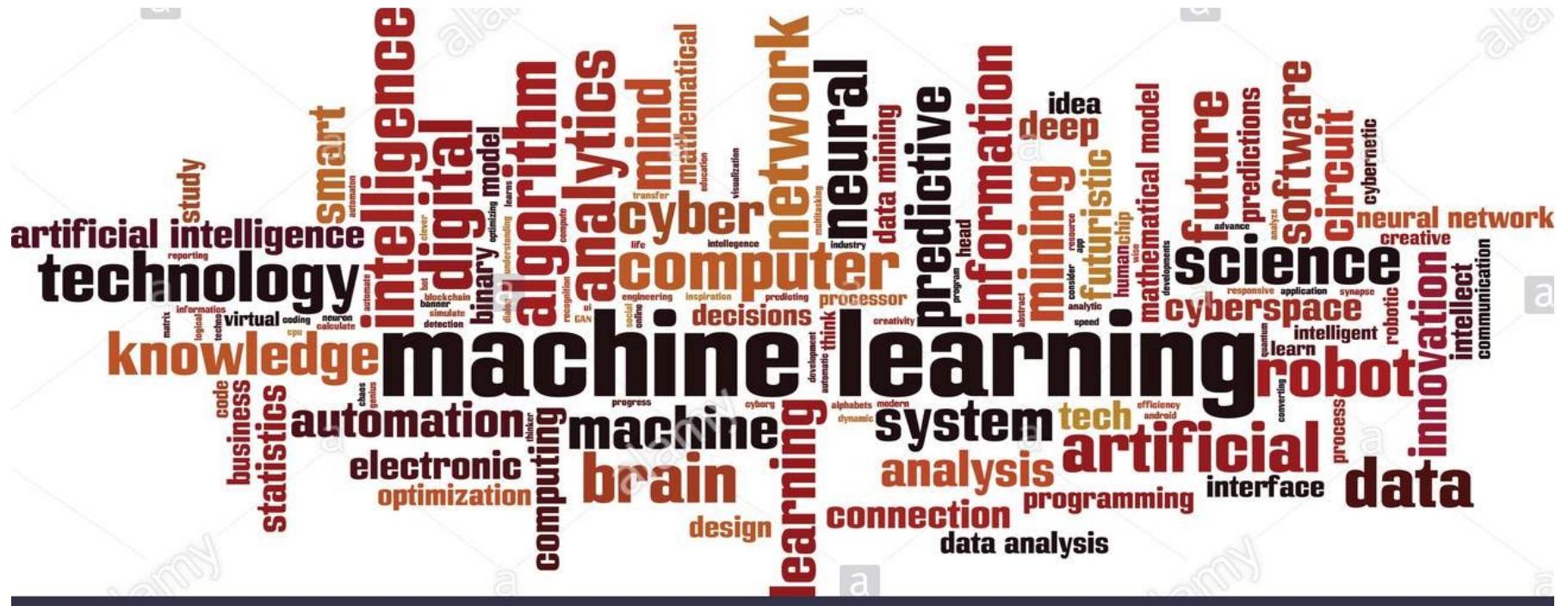
Juan Enriquez
Reprograms Life
p.31

James Gleick
Unsplits the Bit
p.58

AND
Lawrence
Weschler
Questions the
Cloud
p.76







 alamy stock photo

WJF2DY
www.alamy.com

DATA SCIENCE

DETECTION
SOCIAL MEDIA
SERVICES
MULTIMEDIA
NETWORK
PROJECTS
PREDICTIVE
PROGRAM
ANALYTICS
DATA SCIENCE
INFORMATION
SOLUTIONS
MACHINE LEARNING
WEB SERVICES
MATHS PATTERN
ENGINEERING
PLANNING
MEDIA
STATISTICS
TARGET
BIG DATA
CONTENT
PROCESSING
CONSUMER
ORGANIZATION
PLANNING
EVENTS
PROGRAMMING
SOFTWARE
MODELS
E-MARKETING
COMMUNICATION
BRANDING
CONSUMER DEMAND MARKETS
WEB MARKETING
DATA MINING
VISION
WEB DEV
STRATEGY
MOBILE
INFORMATION
DIGITAL
SERVICE
PRICING
CODE
SEGMENTATION
SOCIAL NETWORKS
RESEARCH
PROBABILITY
COMPUTING
KDD
WORLDWIDE
VISUALIZATION
DATA
SOCIAL NETWORK

ARTIFICIAL INTELLIGENCE

MACHINES
ADVERTISING
SOCIAL MEDIA
SERVICES
BIG DATA
CONTENT
SMART
CONSUMER
ORGANIZATION
SEO PLANNING
PROGRAMMING
PROJECTS
WEB
SMART
PROCESS
SOFTWARE
COMPUTING
REAL-TIME
DYNAMIC
WEB SERVICES
CLOUD
EVERYTHING
AUTOMATION
SMART
REPORT
MULTIMEDIA
NETWORK
IDENTIFY
ANALYSIS
DETECT
BRANDS
SOLUTIONS
APPs
AUTOMATION
INFORMATION
COMMUNICATION
SERVICES
KDD
SOFTWARE
PREDICTIVE
ANALYTICS
STATISTICS
OPERATION
PLANNING
AUTOMATE
DECISION
ENGINEERING
RESEARCH
ROBOT
COM
DATA
STRATEGY
WORLDWIDE
AUTOMATION
ORGANIZATION
PERFORMANCE
SOLUTIONS
COMPUTER
PROCESSING
#87886806

ALSO
SET
USE
CONTINUES
LARGER
TENS
COMPLEX
ANALYTICS
USED
CREATED
NOW
CAPACITY
HUNDREDS
RECORDS
NETWORKS
DATABASES
SEARCH
CONNECTOMICS
ORGANIZATIONS
RELATIONAL
SOCIAL
INDEXING
CITATION
BIOLOGICAL
PROCESSING
UBIQUITOUS
WORLD'S
DESKTOP
CURRENTLY
SOLID
TYPES
GARTNER
DIFFICULTY
BIG
DATA
EVERY
EXAMPLES
TOOLS
DISK
TARGET
APPLIED
SHARED
SENSOR
DEFINITION
CURRENT
MOVING
MAY
WITHIN
THOUGHT
RECONSIDER
PRACTITIONERS
ZETTABYTES
INTERNET
DESCRIBING
RADION-FREQUENCY
MANAGEMENT
DISTRIBUTED
SETS
GENOMICS
TERABYTES
CASE
COMPLEXITY
NEEDED
TIME
SOFTWARE
PERFORMANCE
RESEARCH
LOGS
COST
QUALITIES
SIZE
ABILITY
INCLUDE
TOLERABLE
SYSTEMS
FIRMS
ELAPSED
PETABYTES
STORAGE
PARALLEL
SAN
MASSIVELY
GROW
SIGNAL
DEFINING
RELATED
COMPUTING
MANAGE
ARCHIVES
BIOGEOCHEMICAL
TECHNOLOGIES
ONE
LARGE
PROCESS
STORE
PRESENTATIONS
COMBAT

What is Data Mining?

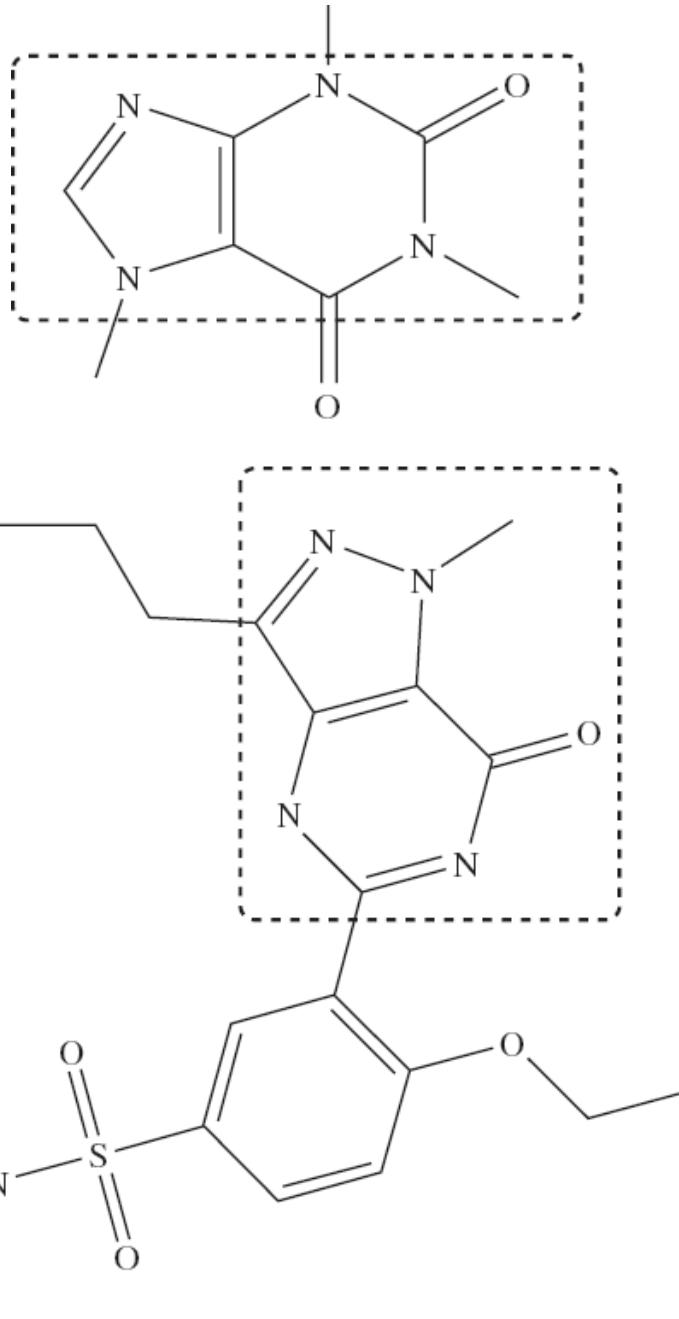
- *Data mining is the process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems.*
- *Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data.*

Market Basket Analysis

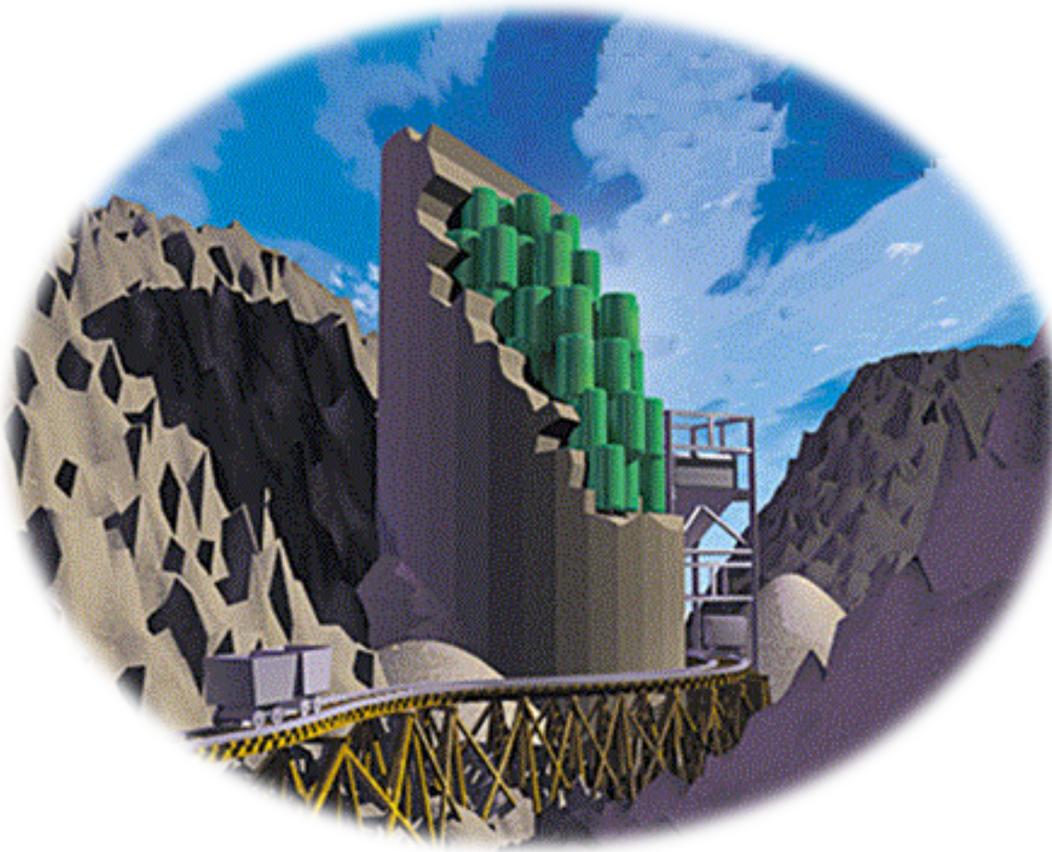


Chemistry Informatics

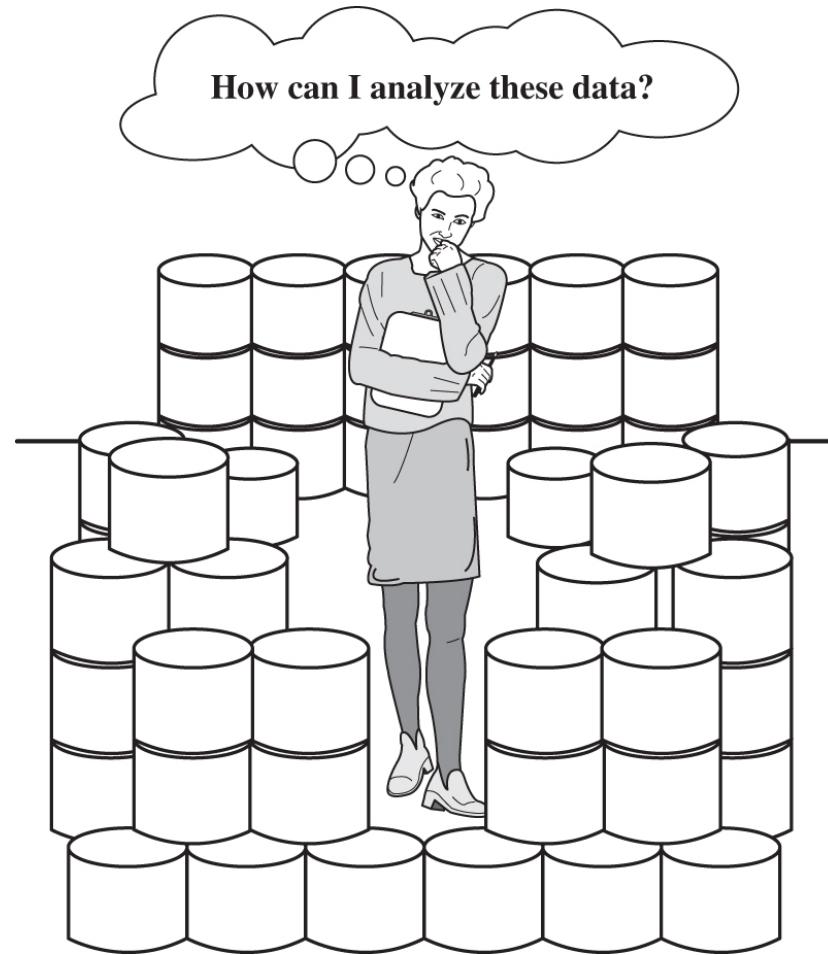
A Chemical database.

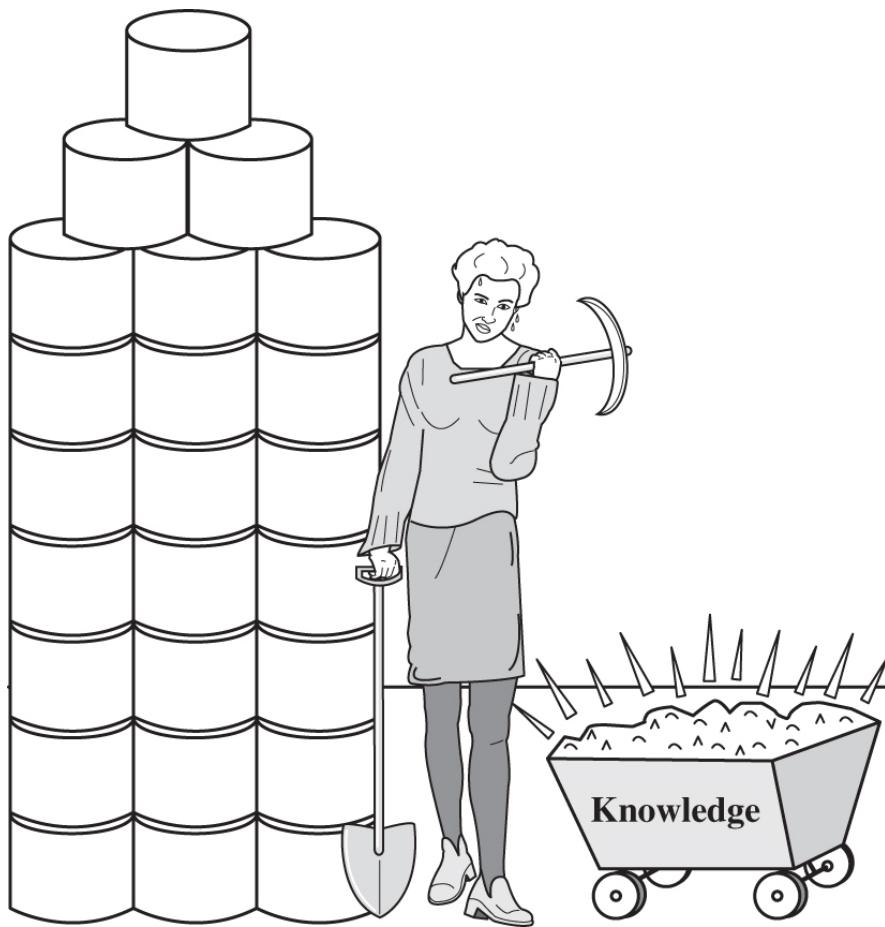


What is Data Mining?



Source: Cover page of *Advanced in Knowledge Discovery and Data Mining*,
edited by U. Fayyad, G. Piatesky-Shapiro, P. Smyth and R. Uthurusamy, MIT Press





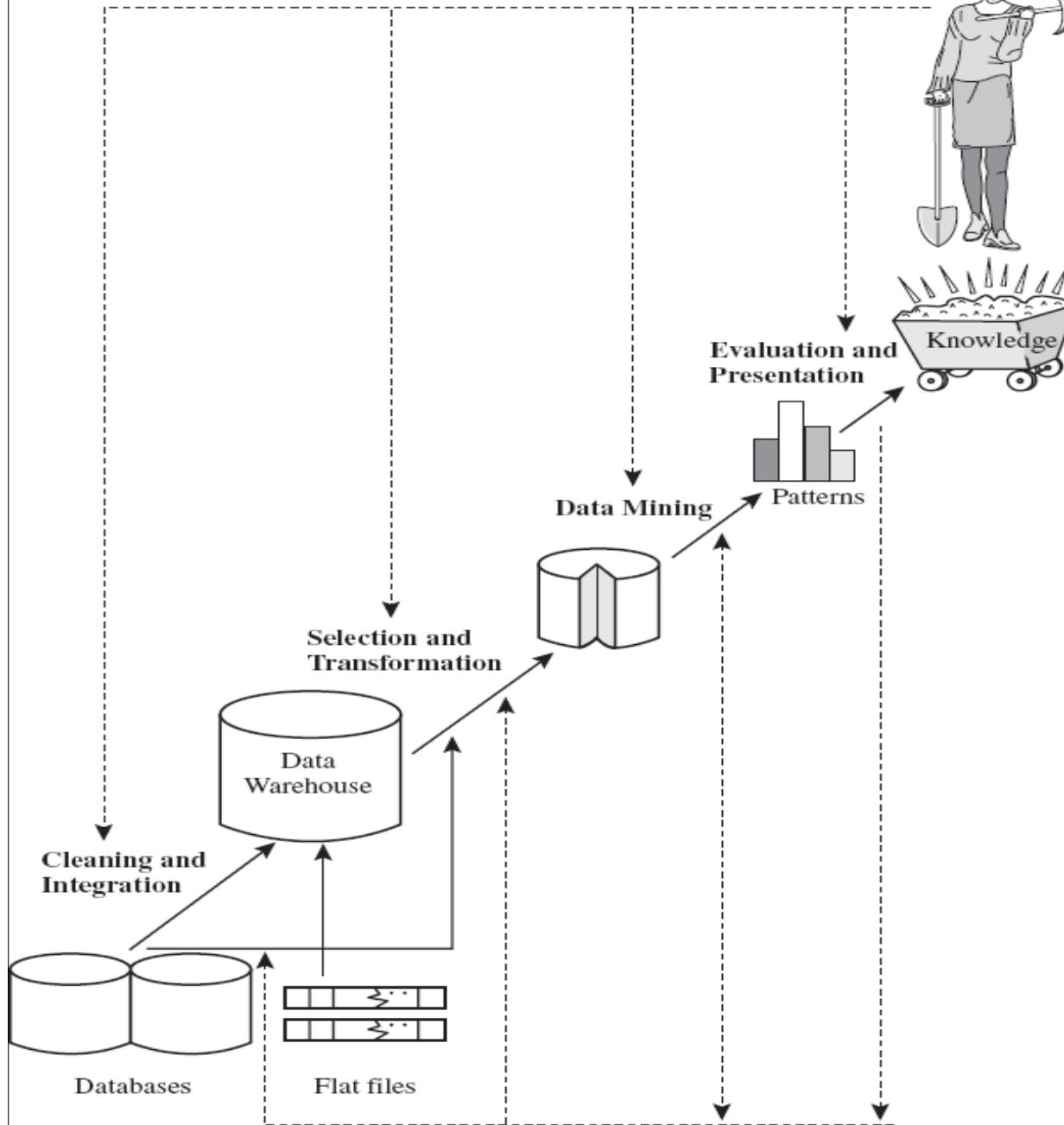
What Is Data Mining?

- 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”?
 - Simple search and query processing
 - (Deductive) expert systems

discovery

Process of knowledge discovery

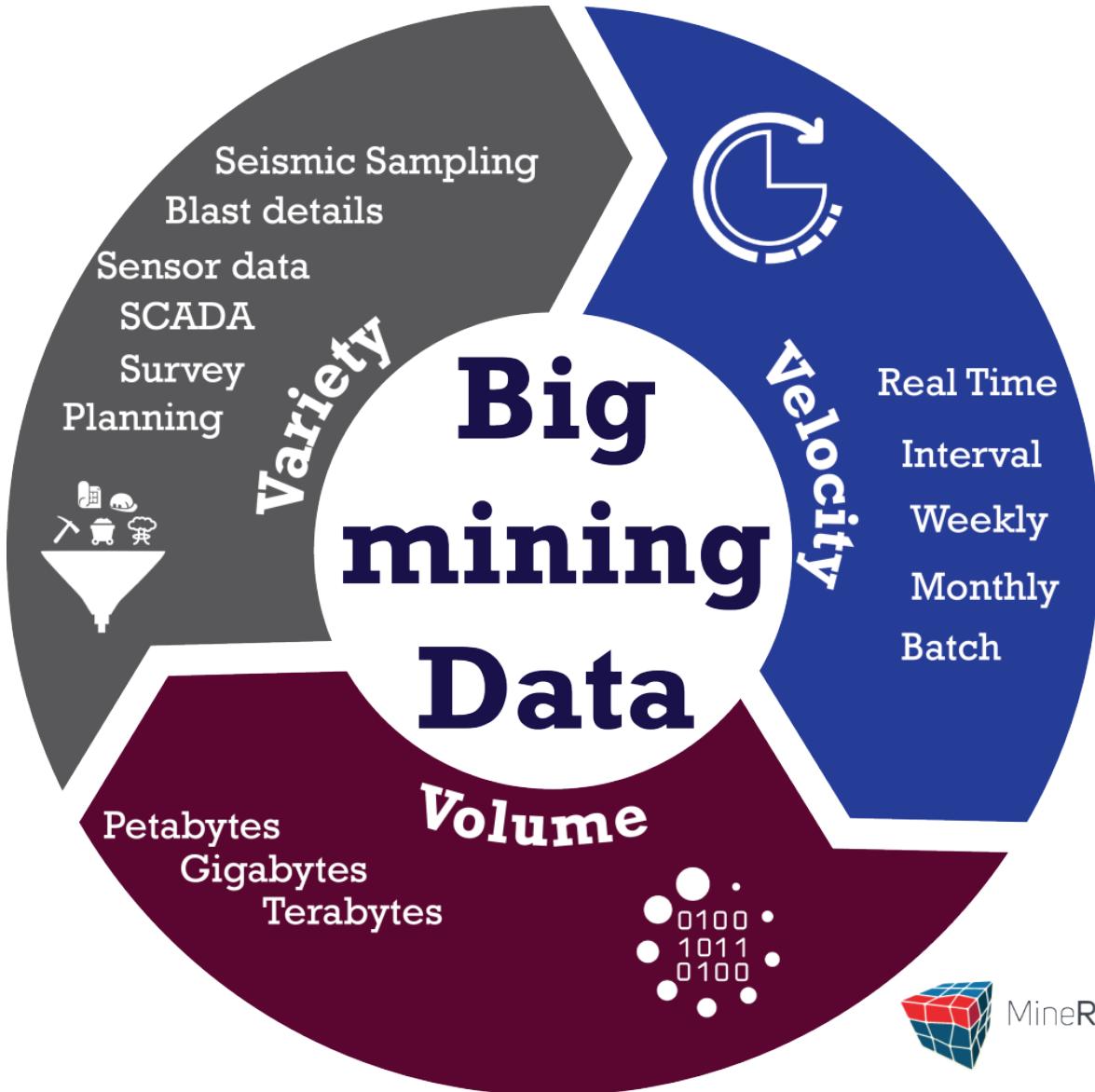
Data Mining is a

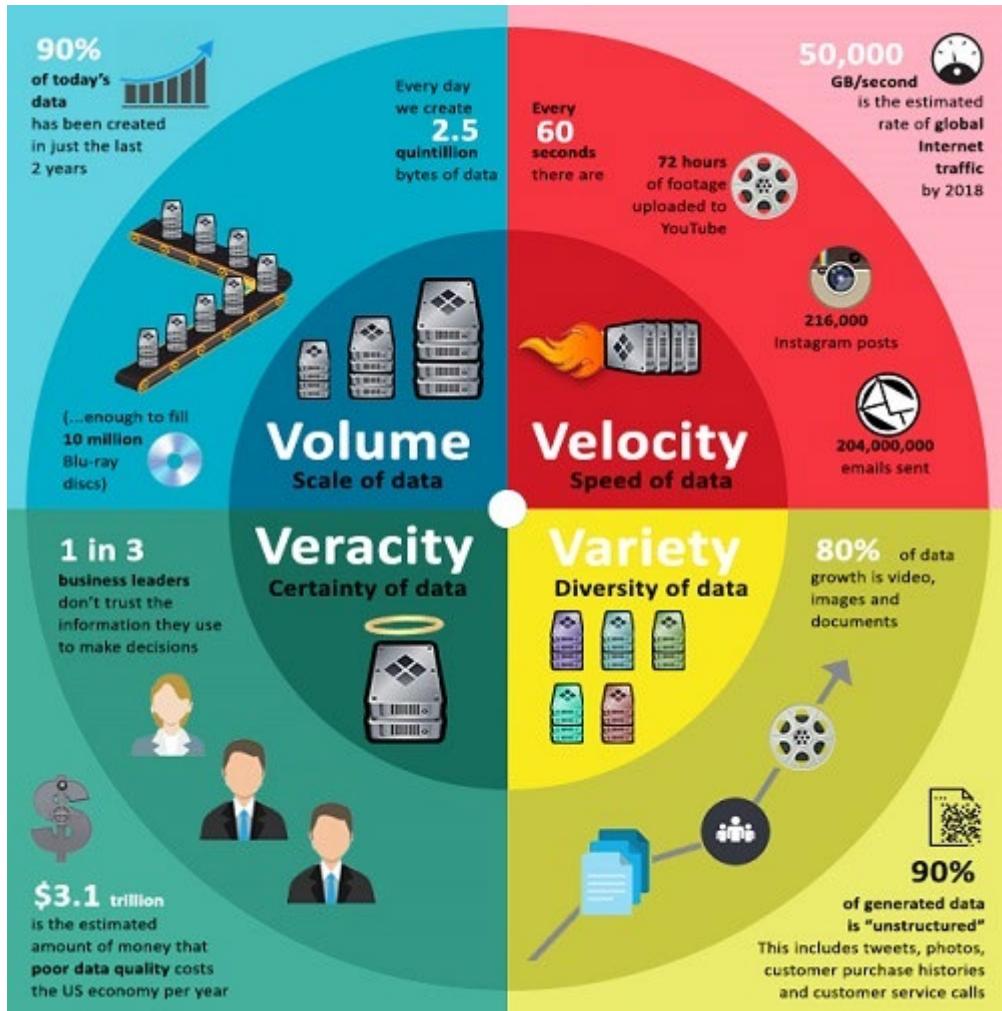


Data mining as a step in the process of knowledge discovery



<https://youtu.be/zDAYZU4A3w0>





The Five V's of Big Data



Scale of Data

This refers to the sheer volume of data being generated every second.



40 Zettabytes
of data will be created by 2020 and
increase of 300 times from 2005

6 Billion People
have cell phones

Most companies in the U.S. have at least
100 Terabytes
of data stored.



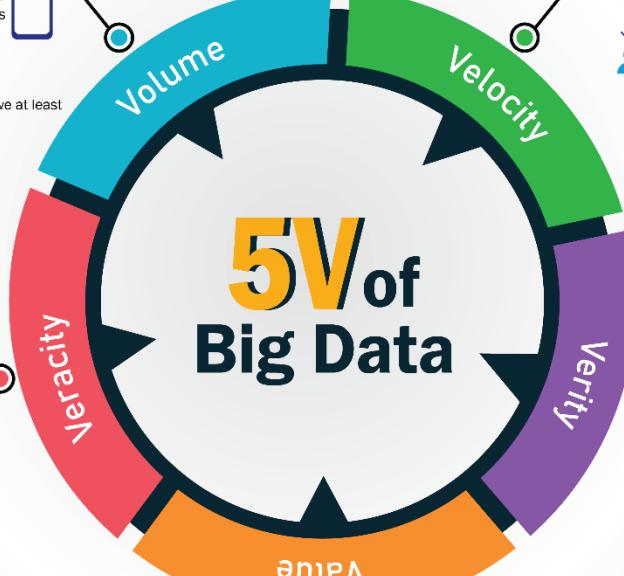
Uncertainty Of Data

1 in 3 Business leaders
don't trust the information they use to make decisions



This refers to the discrepancies found in the data.

Poor data quality costs the US economy around
\$ 3.1 Trillion a year



Analysis of Streaming Data

The New York Stock Exchange
capture **1 TB of Trade Information**

Denotes the speed at which data is emanating and changes are occurring between the diverse data sets.



By 2016 it is projected there will be **18.9 Billion** network connections



Modern cars have close to **100 Sensors**



4 Billion+
hours of video are watched on You Tube each month



30 Billion
pieces of content are shared on facebook every month



400 Million
tweets are sent per day by about 200 million monthly active users



As more and more data is being digitized.



Value Of Data

Having access to big data is all well and good but that's only useful if we can turn it into a value.

Exhibit 1

Big data can generate significant financial value across sectors



US health care

- \$300 billion value per year
- ~0.7 percent annual productivity growth



Europe public sector administration

- €250 billion value per year
- ~0.5 percent annual productivity growth



Global personal location data

- \$100 billion+ revenue for service providers
- Up to \$700 billion value to end users



US retail

- 60+% increase in net margin possible
- 0.5–1.0 percent annual productivity growth



Manufacturing

- Up to 50 percent decrease in product development, assembly costs
- Up to 7 percent reduction in working capital

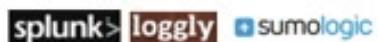
SOURCE: McKinsey Global Institute analysis

Big Data Landscape

Vertical Apps



Log Data Apps



Ad/Media Apps



Business Intelligence



Analytics and Visualization



Data As A Service



Analytics Infrastructure



Operational Infrastructure



Infrastructure As A Service



Structured Databases

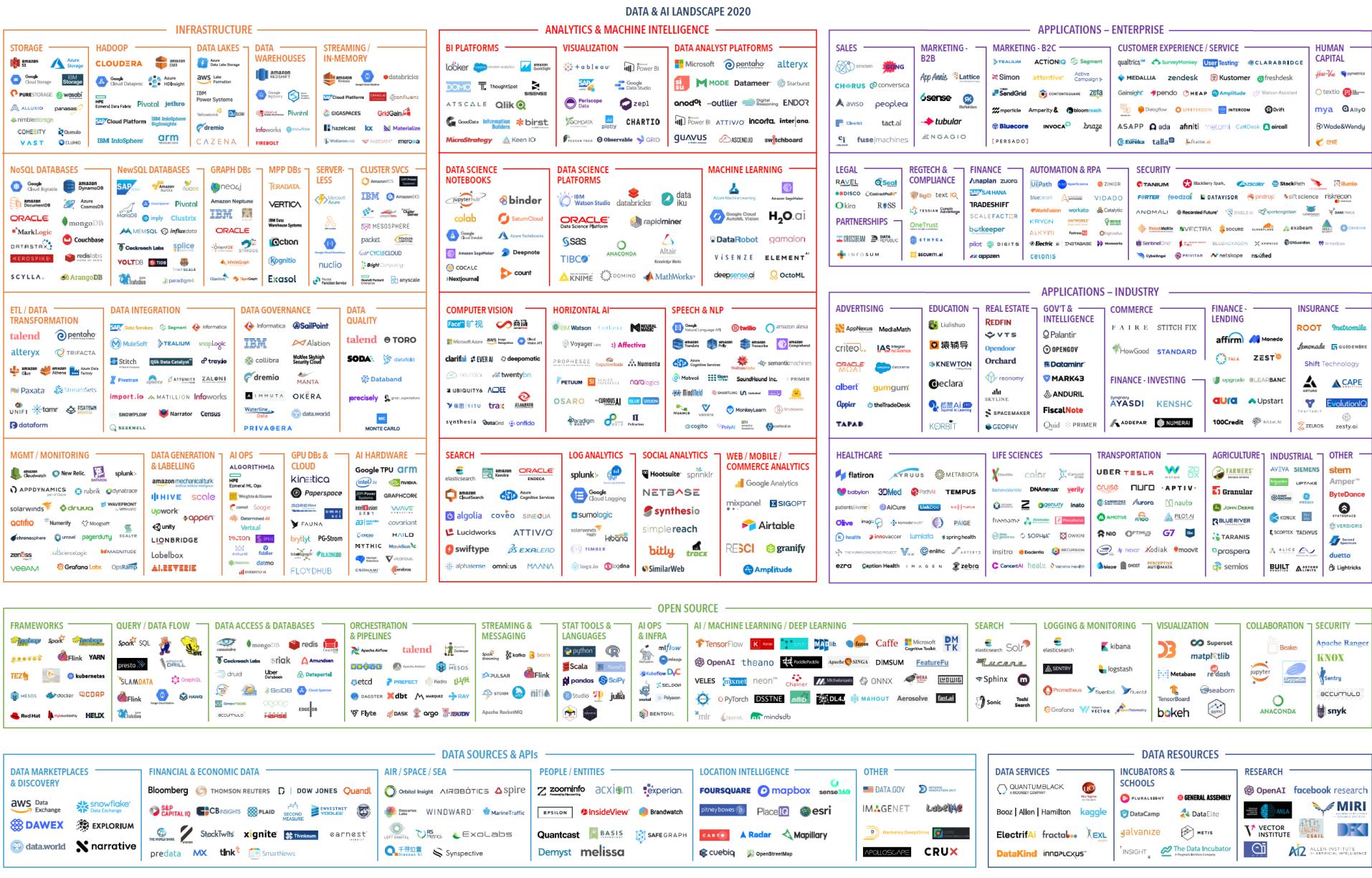


Technologies

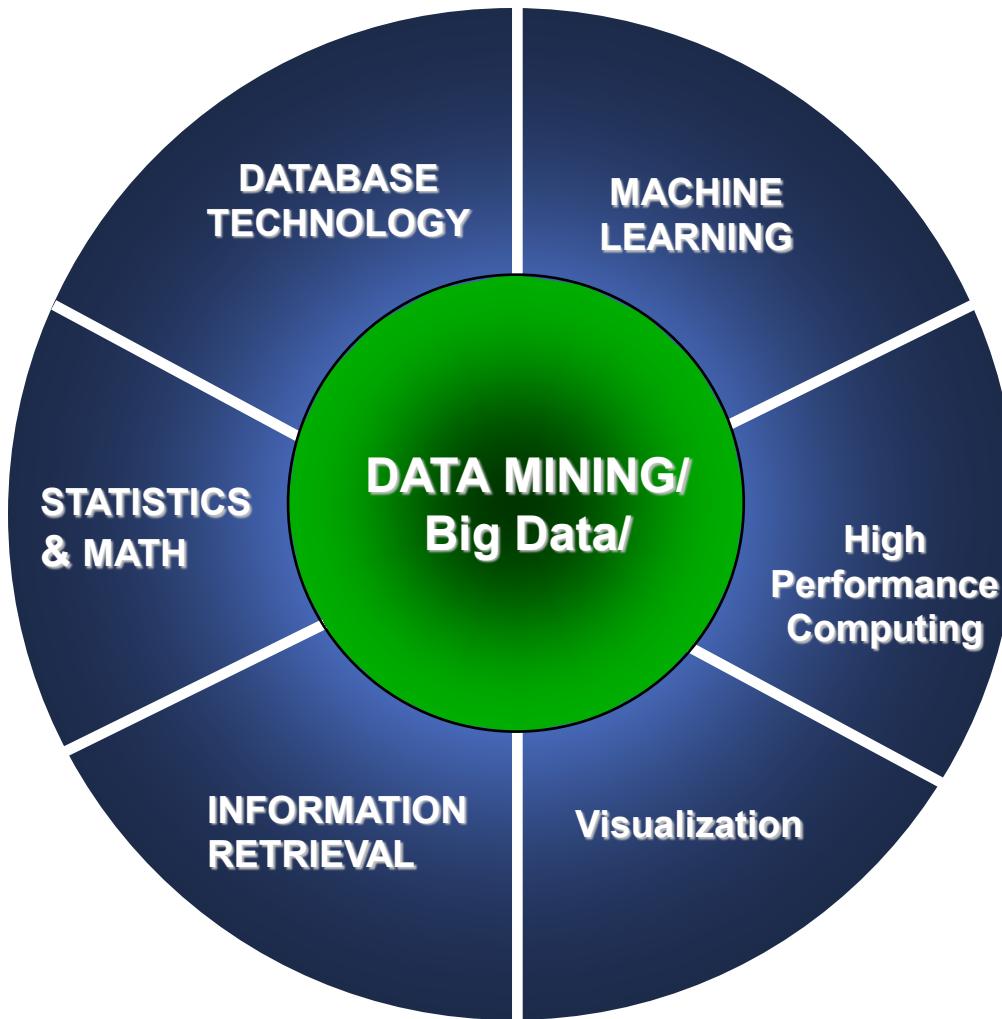


BIG DATA & AI LANDSCAPE 2018





Data Mining/Big Data/Data Science is an Interdisciplinary and Multidisciplinary Field



Why Data Mining/Big Data is Important?

- *McKinsey predicts that data--driven technologies will bring an additional \$300 billion of value to the U.S. health care sector alone, by 2020, 1.5 million more “data--savvy managers” will be needed to capitalize on the potential of data*

– May 1, 2011

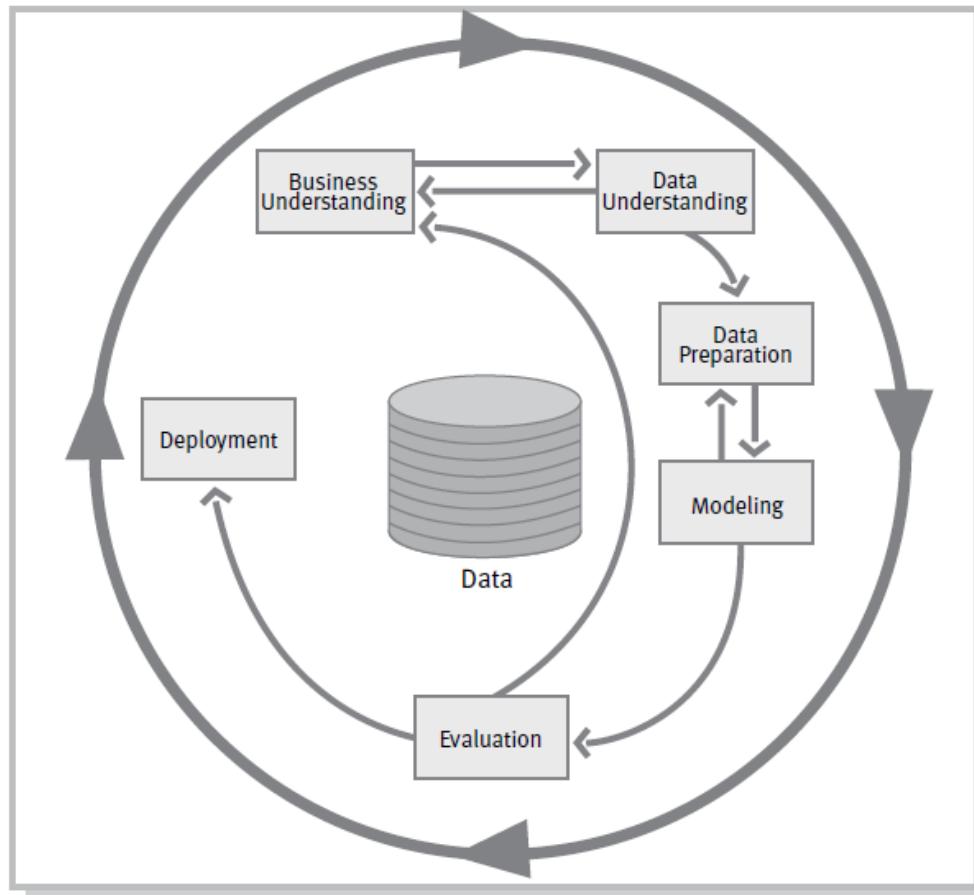
Multi-Dimensional View of Data Mining/Big Data

- **Data to be mined**
 - Structured data, semi-structured data, unstructured data, relational data, data warehouse, transactional data, stream, spatiotemporal, time-series, sequence, text and web, multi-media data, graphs data, social network, etc.
- **Data preprocessing and feature engineering**
 - Cleaning, integration, reduction, transformation, generalization, feature engineering, etc.
- **Knowledge to be mined (Data mining functions and models)**
 - generalization, association, classification, clustering, etc.
- **Techniques utilized**
 - data warehouse (OLAP), frequent pattern mining, machine learning, statistics, pattern recognition, visualization, high-performance computing, parallel and distributed computing, cloud computing, etc.
- **Applications adapted**
 - Market basket analysis, healthcare, manufacturing engineering, education, CRM, fraud detection, customer segmentation, intrusion detection, financial banking, criminal investigation, bio Informatics, direct marketing, classifying stars, medical diagnosis, computer vision, drug discovery, speech recognition, handwriting recognition, credit scoring, human genetic clustering, medical imaging, market segmentation, product positioning, new product development, social network analysis, image segmentation, recommender systems, anomaly detection, crime analysis, climatology, market research, customer trend analysis, market movement tracking, customer buying patterns, analyze tax records, decision making, on line operations and report, etc.

CRISP-DM

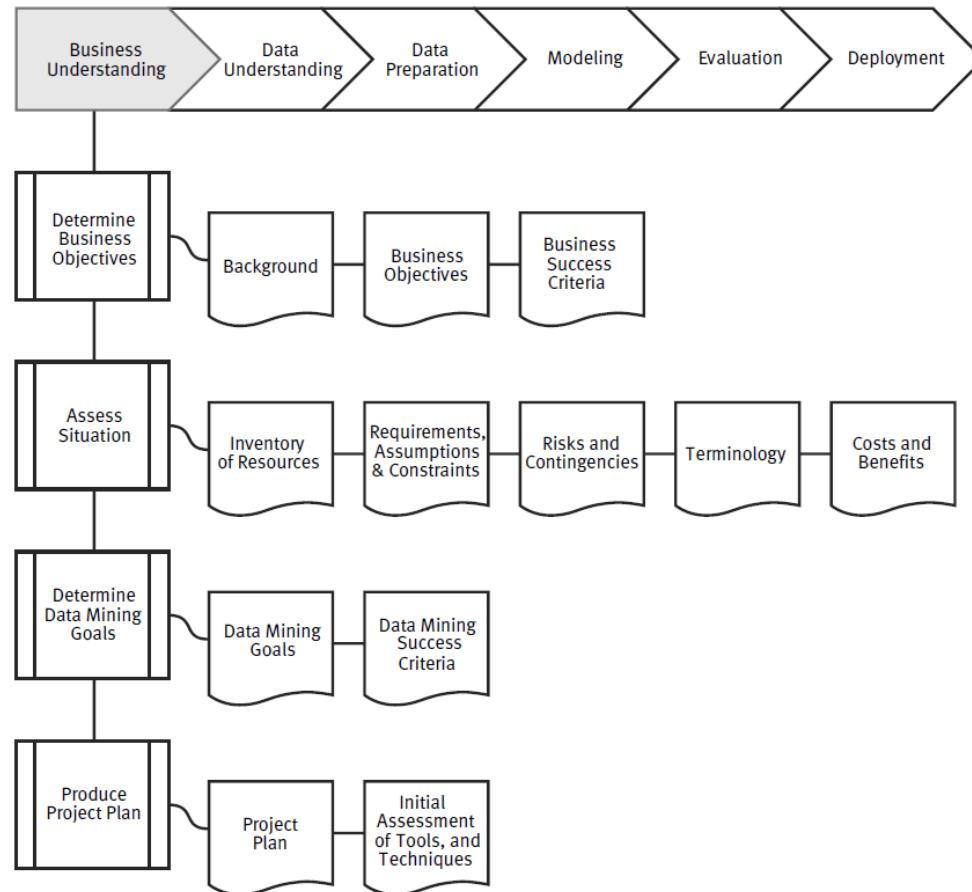
- CRoss-Industry Standard Process for Data Mining
 - is described in terms of a hierarchical process model, consisting of sets of tasks described at four levels of abstraction (from general to specific): phase, generic task, specialized task, and process instance

The life cycle of a data mining project

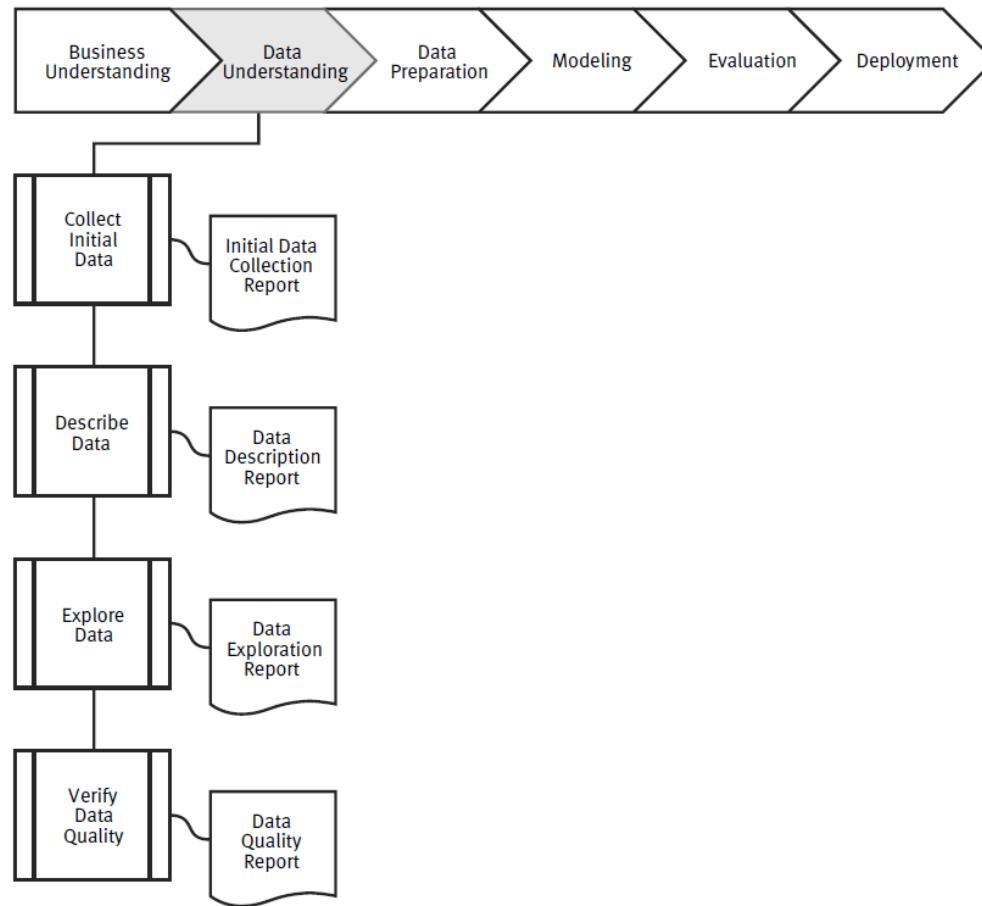


Business Understanding	Data Understanding	Data Preparation	Modeling	Evaluation	Deployment
<p>Determine Business Objectives <i>Background</i> <i>Business Objectives</i> <i>Business Success Criteria</i></p> <p>Assess Situation <i>Inventory of Resources Requirements, Assumptions, and Constraints</i> <i>Risks and Contingencies</i> <i>Terminology</i> <i>Costs and Benefits</i></p> <p>Determine Data Mining Goals <i>Data Mining Goals</i> <i>Data Mining Success Criteria</i></p> <p>Produce Project Plan <i>Project Plan</i> <i>Initial Assessment of Tools and Techniques</i></p>	<p>Collect Initial Data <i>Initial Data Collection Report</i></p> <p>Describe Data <i>Data Description Report</i></p> <p>Explore Data <i>Data Exploration Report</i></p> <p>Verify Data Quality <i>Data Quality Report</i></p>	<p>Select Data <i>Rationale for Inclusion/Exclusion</i></p> <p>Clean Data <i>Data Cleaning Report</i></p> <p>Construct Data <i>Derived Attributes</i> <i>Generated Records</i></p> <p>Integrate Data <i>Merged Data</i></p> <p>Format Data <i>Reformatted Data</i></p> <p><i>Dataset</i> <i>Dataset Description</i></p>	<p>Select Modeling Techniques <i>Modeling Technique</i> <i>Modeling Assumptions</i></p> <p>Generate Test Design <i>Test Design</i></p> <p>Build Model <i>Parameter Settings</i> <i>Models</i> <i>Model Descriptions</i></p> <p>Assess Model <i>Model Assessment</i> <i>Revised Parameter Settings</i></p>	<p>Evaluate Results <i>Assessment of Data Mining Results w.r.t. Business Success Criteria</i> <i>Approved Models</i></p> <p>Review Process <i>Review of Process</i></p> <p>Determine Next Steps <i>List of Possible Actions</i> <i>Decision</i></p>	<p>Plan Deployment <i>Deployment Plan</i></p> <p>Plan Monitoring and Maintenance <i>Monitoring and Maintenance Plan</i></p> <p>Produce Final Report <i>Final Report</i> <i>Final Presentation</i></p> <p>Review Project <i>Experience Documentation</i></p>

Business understanding



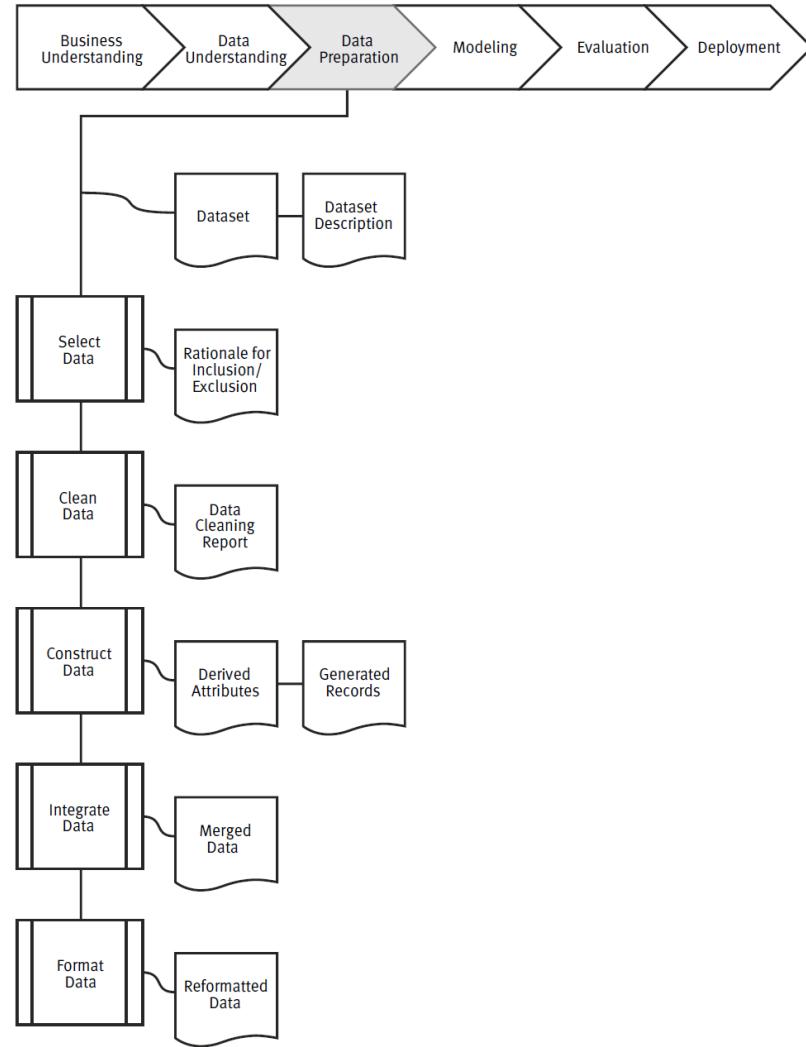
Data understanding



Structured
Semi-structured
Unstructured
Categorical
 Nominal
 Binary
 Ordinal
Numeric
 Interval
 Ratio
Discrete
Continuous

Descriptive statistics
 Count, mean, mode, standard deviation, quantiles, etc.
Class distribution
Feature skewness
Correlations between features
 Chi-square
 Correlation
Data visualization

Data preparation



Data preprocessing	
Data cleaning	
	Missing values Use the most probable value to fill in the missing value (and five other methods)
	Noisy data Binning; regression; clustering
Data integration	
	Entity ID problem Metadata
	Redundancy Correlation analysis (Correlation coefficient, chi-square test)
Data transformation	
	Smoothing Data cleaning
	Aggregation Data reduction
	Generalization Data reduction
	Normalization Min-max; z-score; decimal scaling
	Attribute construction New attributes are constructed and added from the given set of attributes to help mining process
Data reduction	
	Data cube aggregation Data cube store multidimensional aggregated information
	Attribute selection Stepwise forward selection; stepwise backward selection; combination; decision tree induction
	Dimensionality reduction Discrete wavelet transforms (DWT); Principle components analysis (PCA)
	Numerosity reduction Parametric data reduction; clustering; data cube aggregation;
	Data discretization Binning; histogram analysis; entropy-based discretization; Interval merging by chi-square analysis; cluster analysis; intuitive partitioning
	Concept hierarchy Concept hierarchy generation

Feature Engineering

Curse of dimensionality

Feature selection

Filter methods

Wrapper methods

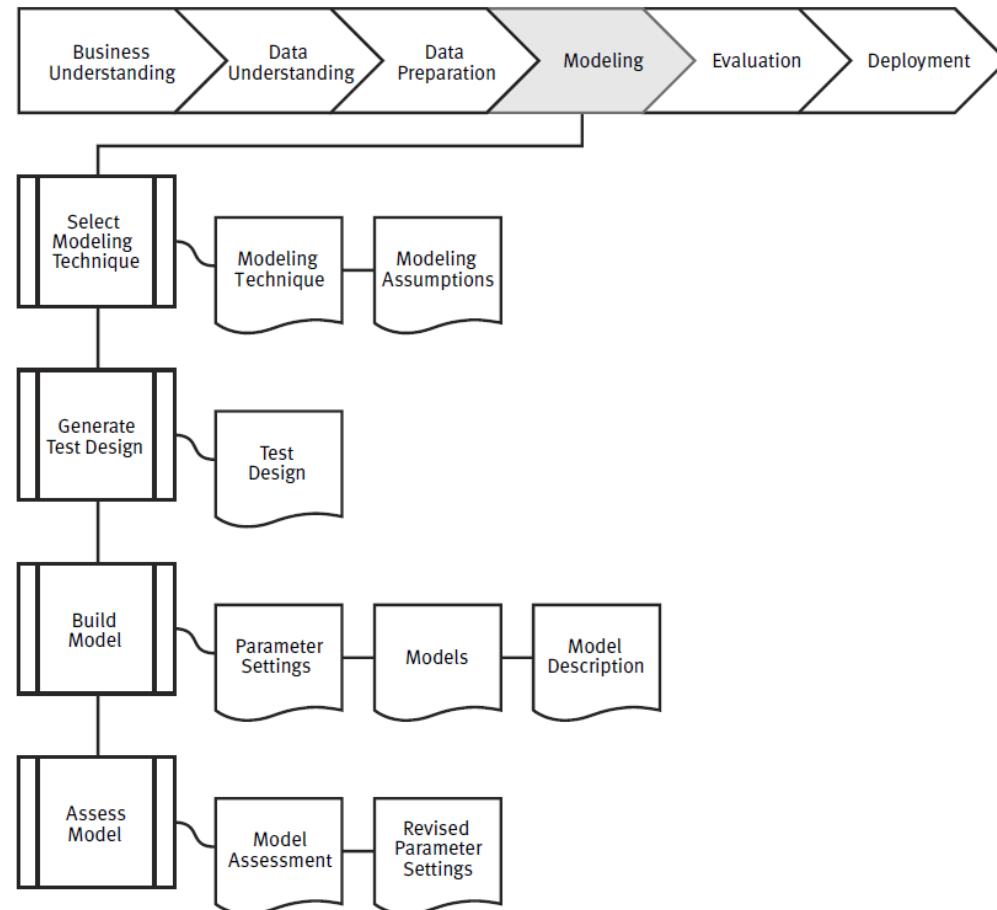
Embedded methods

Hybrid methods

Feature construction

Feature creation

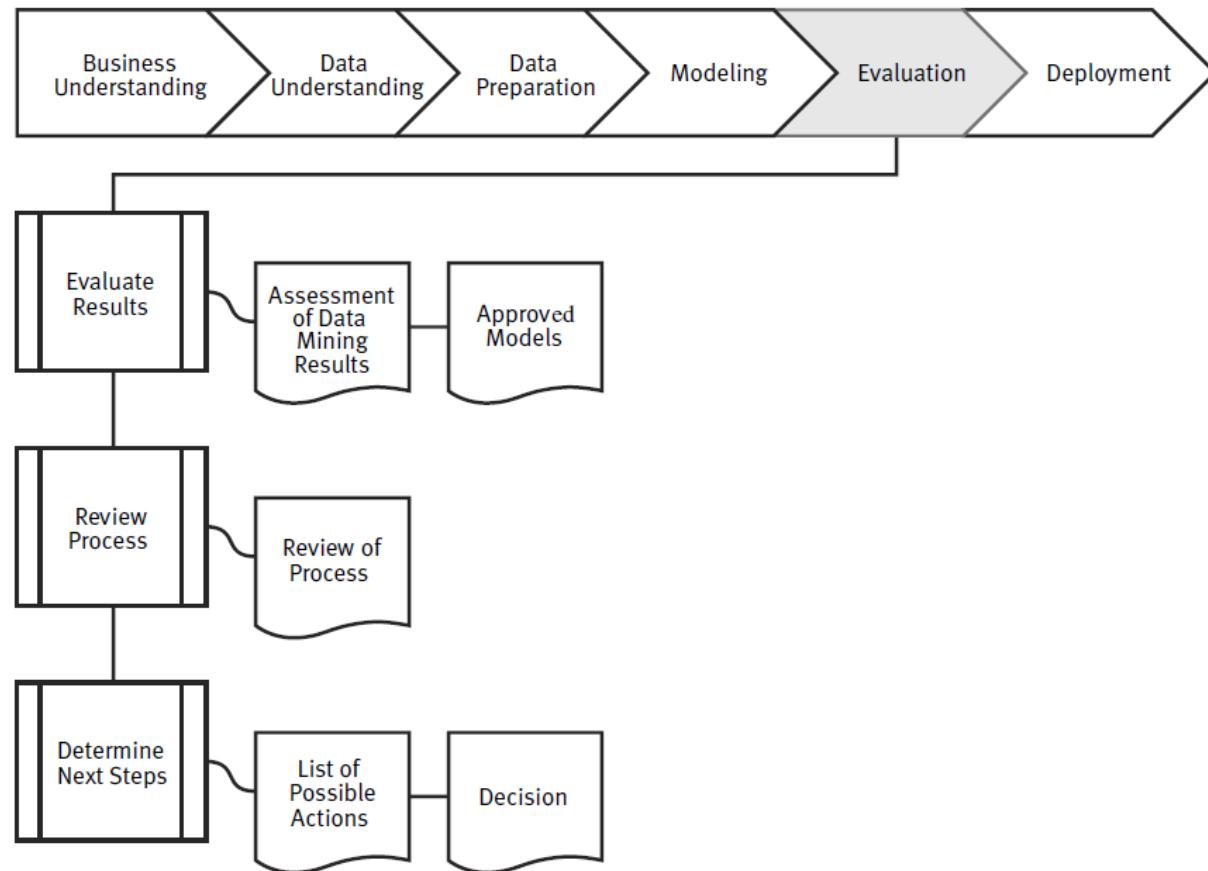
Modeling



Decision tree
Bayesian classifier
Regression
 Simple linear regression
 Multiple linear regression
 Polynomial regression
 Logistic regression
k-Nearest Neighbor classifier (kNN)
Perception
Gradient Descent
Neural Network
 Feed-Forward Neural Network (FNN)
 Backpropagation
 Recurrent Neural Network (RNN)
 Convolutional Neural Network (CNN)
 Probabilistic Neural Network (PNN)
Support Vector Machine (SVM)
Binary Classification
Multiclass Classification
Rule-Based Classification
Bayesian Belief Networks
Genetic Algorithms
Rough Set Approach
Fuzzy Set Approaches
Etc.

Multidimensional Data Modeling
 Data Cube: A Multidimensional Data Model
 Stars, Snowflakes, and Fact Constellations
 Dimensions: The Role of Concept Hierarchies
 OLAP Operations
Mining Frequent Patterns, Associations, and Correlations
 Apriori Algorithm
 Frequent Itemset
 Generating Association Rules from Frequent Itemsets
Clustering
 Partitional methods
 k-means
 k-medoids
 Hierarchical methods
 Agglomerative hierarchical clustering method
 Divisive hierarchical clustering method
 Density-Based Methods
 DBSCAN: Density-Based Clustering Based on Connected Regions with High Density
 OPTICS: Ordering Points to Identify the Clustering Structure
 DENCLUE: Clustering Based on Density Distribution Functions
 Grid-Based Methods
 STING: STatistical INformation Grid
 CLIQUE: An Apriori-like Subspace Clustering Method

Evaluation



confusion matrix

accuracy, recognition rate

error rate, misclassification rate

sensitivity, true positive rate, recall

specificity, true negative rate, precision

$F1$ -score, harmonic mean of precision and recall

F -beta , where β is a non-negative real number

k-fold cross-validation

bootstrap

.632 bootstrap

model selection using statistical tests of significance

comparing classifiers based on cost-benefit and ROC curves

Ensemble Methods

Random Forests

Bagging

Boosting

AdaBoost

Deployment

