



INTRODUCTION TO DATA MINING

Eduardo Ogasawara
eogasawara@ieee.org
<https://eic.cefet-rj.br/~eogasawara>

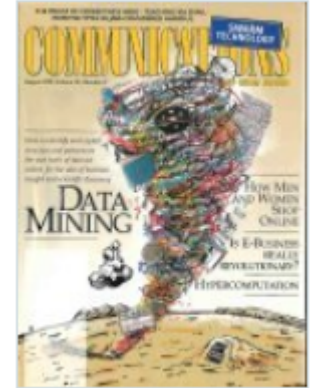
Why Data Mining?

- *Big Data scenario:*
 - *The explosive growth of data: from terabytes to petabytes*
 - *Data collection and data availability*
 - *Automated data collection tools, database systems, Web*
 - *Major sources of abundant and diverse data*
 - *Business: Web, e-commerce, transactions*
 - *Science: sensors, astronomy, bioinformatics, simulation*
 - *Society and everyone: news, photos, videos, open data, IoT*
- *We are drowning in data but starving for knowledge!*
- *"Need is the mother of invention"*
 - *Data mining - Automated analysis of massive data sets*



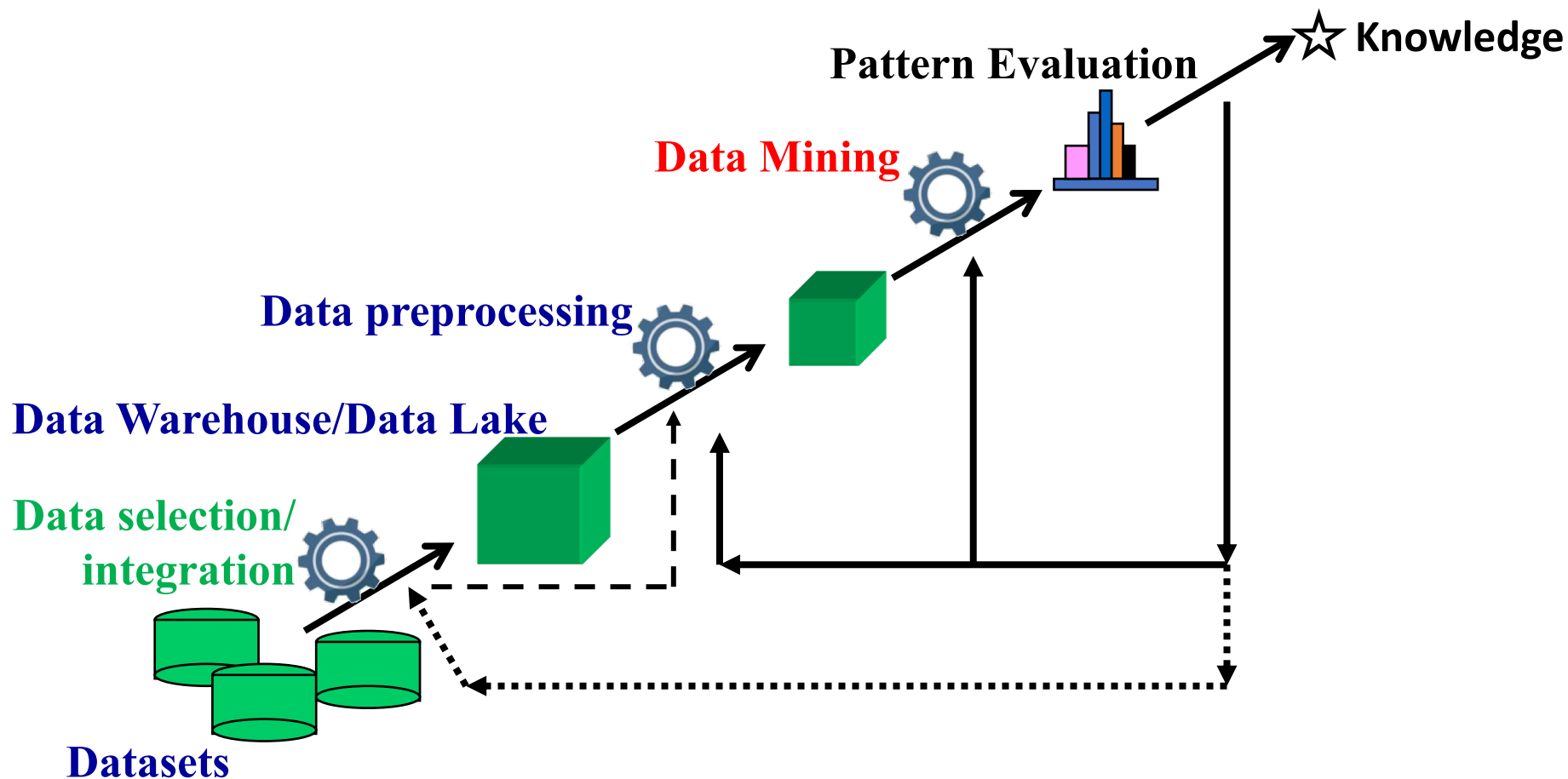
What is Data Mining?

- *Data mining (knowledge discovery from data)*
 - *Extraction of interesting (**non-trivial**, implicit, **previously unknown** and **potentially useful**) patterns or knowledge from a **massive** amount of data*
- *Alternative names*
 - *Knowledge Discovery in Databases (KDD)*
 - *Knowledge Extraction*
 - *Business Intelligence*
 - *Data Analysis*

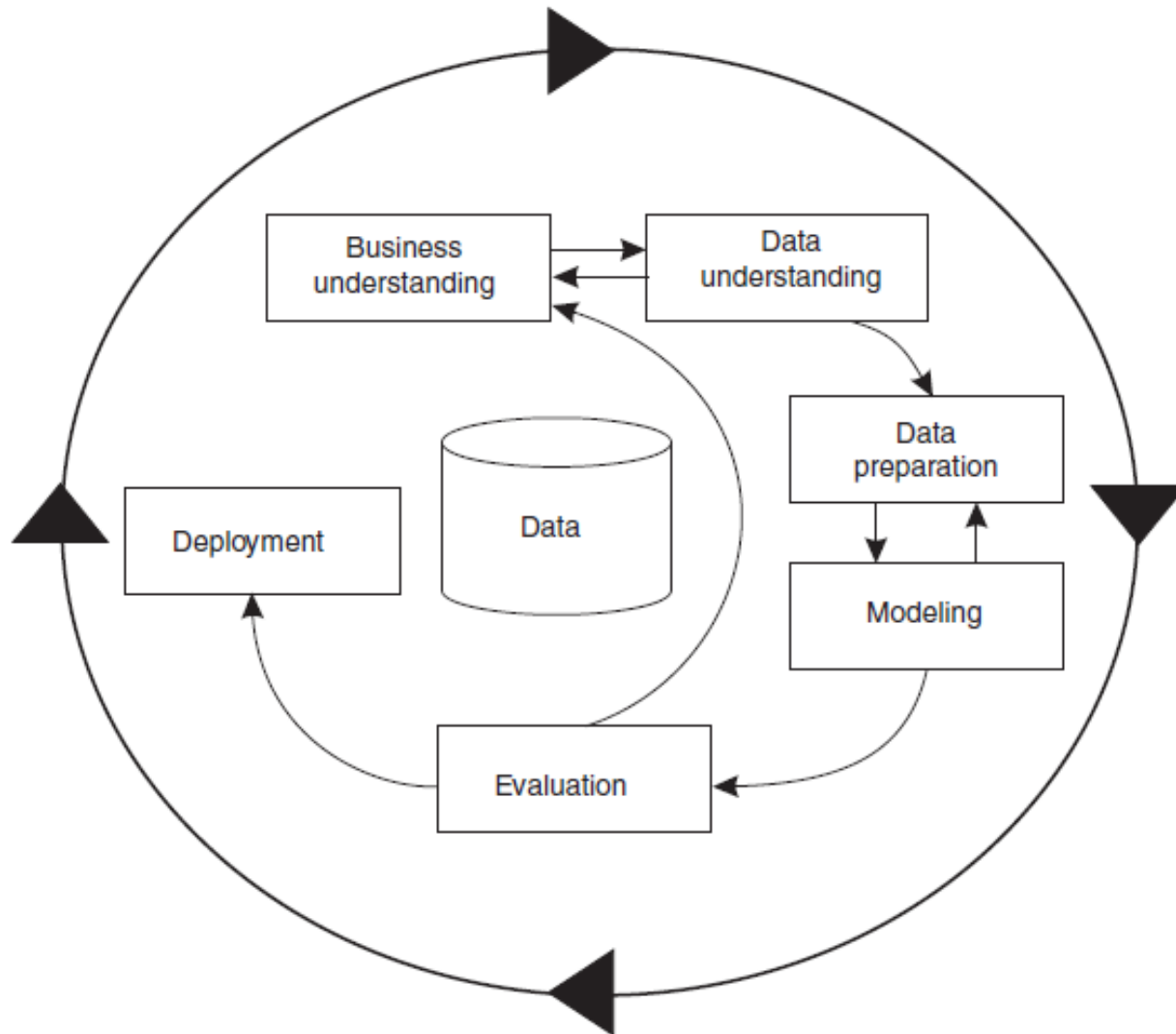


Knowledge discovery from data (KDD) process (Database perspective)

- This is a view from typical database systems
- Data mining plays an essential role in the KDD process








Cross-Industry Standard Process for Data Mining (CRISP-DM)



Knowledge discovery in databases (KDD) process vs. Cross-Industry Standard Process for Data Mining (CRISP-DM)

Methodology	Phases					
CRISP-DM	Business understanding	Data understanding	Data preparation	Modeling	Evaluation	Deployment
KDD-Detailed	Learning the application domain	Creating a target data set	Data cleaning and pre-processing	Choosing the function of DM	Interpretation	Using discovered knowledge
			Data reduction and projection	Choosing the DM algorithm		
			Data mining			
KDD-Outlined		Selection	Pre-processing	Modeling	Interpretation / Evaluation	
		Pre-processing	Transformation	Data mining		

Is everything "data mining"?

-  *Query processing*
-  *Information retrieval*
-  *Data Warehouse & OLAP queries*
-  *(Deductive) expert systems*
-  *Data exploration*

Multi-Dimensional View of Data Mining

Data models

Knowledge to be mined

Data mining functions

Methods/Techniques used

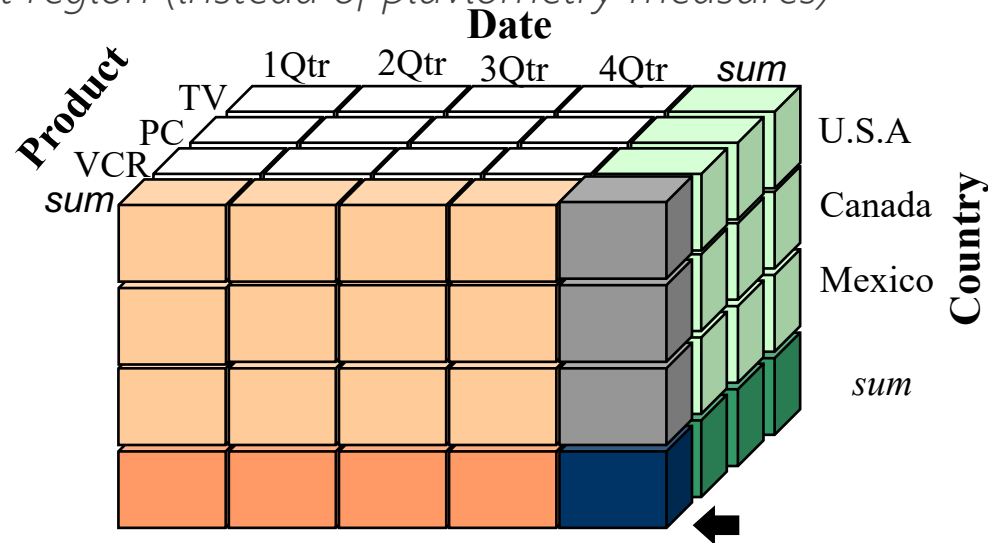
Applications

Data Mining: on what kinds of data models?

- *Database-oriented data sets and applications*
 - *A Relational database, data warehouse, transactional database*
 - *Object-relational databases, Heterogeneous databases, and legacy databases*
- *Advanced data sets and advanced applications*
 - *Data streams and sensor data*
 - *Time-series data, temporal data, sequence data (incl. bio-sequences)*
 - *Structure data, graphs, social networks, and information networks*
 - *Spatial data and spatiotemporal data*
 - *Multimedia database*
 - *Text databases (text mining)*
 - *The World-Wide Web*

Data Mining Function: (1) Generalization

- Information integration and data warehouse construction
 - Data cleaning, transformation, integration, and multidimensional data model
- Data cube technology
 - Scalable methods for computing (i.e., materializing) multidimensional aggregates
 - OLAP (online analytical processing)
- Multidimensional concept description: characterization and discrimination
 - Generalize, summarize, and contrast data characteristics,
 - E.g.: dry vs. wet region (instead of pluviometry measures)



Data Mining Function: (2) Association and Correlation Analysis

- Frequent patterns (or frequent item sets)
 - What items are frequently purchased together in your supermarket?
- Association, correlation vs. causality
 - Typical association rules
 - Diaper → Beer [0.5%, 75%] (support, confidence)
 - Are strongly associated items also strongly correlated?
- How to mine such patterns and identify such rules efficiently in large datasets?
- How to use (rank) such patterns?

Mining Association Rules between Sets of Items in Large Databases

Rakesh Agrawal, Tomasz Imieliński, Arun Swami
IBM Almaden Research Center
650 Harry Road, San Jose, CA 95120

Abstract

We present a new technique of mining transaction data. The problem is to find frequent item sets in a database. The algorithm presented here generates all frequent item sets in a database. The algorithm is efficient and runs in linear time.

1 Introduction

Consider a supermarket with a large collection of items. Typical business decisions that the management of the supermarket has to make include which to put on sale, how to stock shelves, how to place merchandise on shelves to make it attractive to people, etc. Analysis of past transaction data is a commonly used approach to make these decisions. The quality of such decisions can be improved by using data mining techniques.

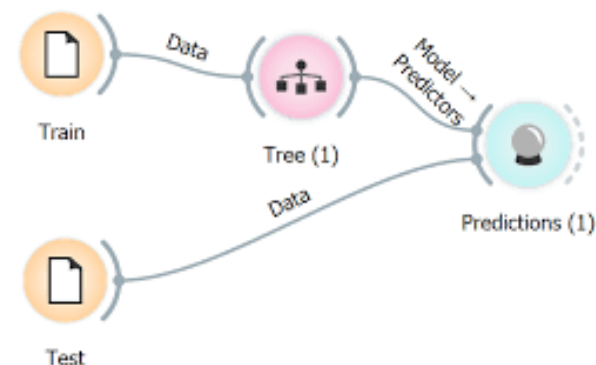
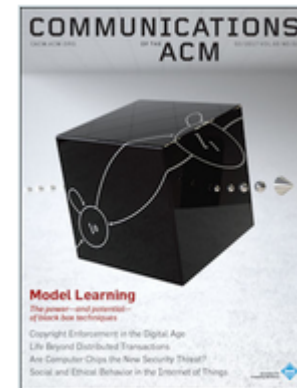
One of the most interesting and useful data mining techniques is association rule mining. An association rule is a statement that certain items are bought together. For example, "people who buy bread also buy butter" is an association rule. Association rules are used to analyze the buying behavior of customers and to make decisions about which items to stock, how to place merchandise on shelves, etc. Association rules are also used to analyze the buying behavior of customers and to make decisions about which items to stock, how to place merchandise on shelves, etc.

Association rules are used to analyze the buying behavior of customers and to make decisions about which items to stock, how to place merchandise on shelves, etc. Association rules are also used to analyze the buying behavior of customers and to make decisions about which items to stock, how to place merchandise on shelves, etc.

Association rules are used to analyze the buying behavior of customers and to make decisions about which items to stock, how to place merchandise on shelves, etc. Association rules are also used to analyze the buying behavior of customers and to make decisions about which items to stock, how to place merchandise on shelves, etc.

Data Mining Function: (3) Prediction

- *Classification and label prediction*
 - *Construct models based on some training examples*
 - *Data-driven models*
 - *Describe and distinguish classes or concepts for future prediction*
 - *E.g., classify countries based on (climate)*
- *Typical methods:*
 - *Decision trees*
 - *Naive Bayes classifier*
 - *Support vector machines*
 - *Neural networks & deep learning*
 - *Random Forest*
 - *Linear/Logistic regression*
- *Typical applications*
 - *Scientific*
 - *Industry & enterprises*
 - *Government & society*

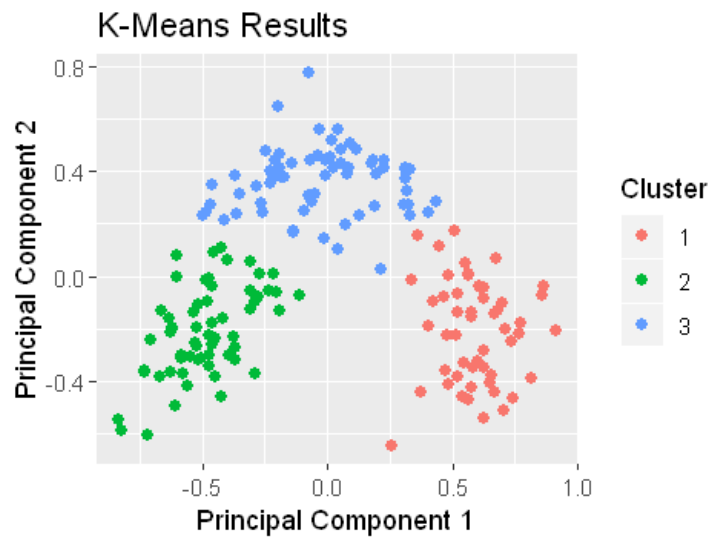


[1] F. Vaandrager, 2017, Model learning, Communications of the ACM, v. 60, n. 2, p. 86–95.

[2] <https://towardsdatascience.com/data-science-made-easy-data-modeling-and-prediction-using-orange-f451f17061fa>

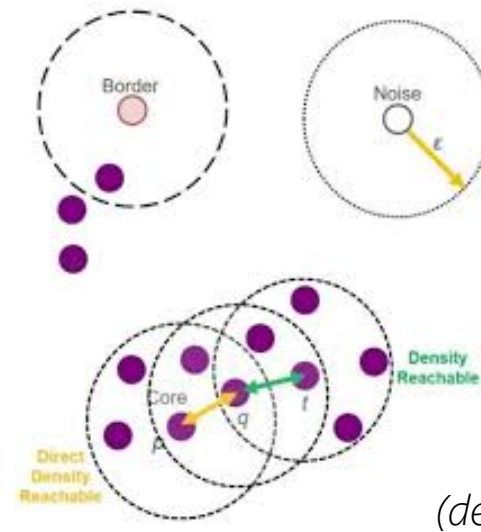
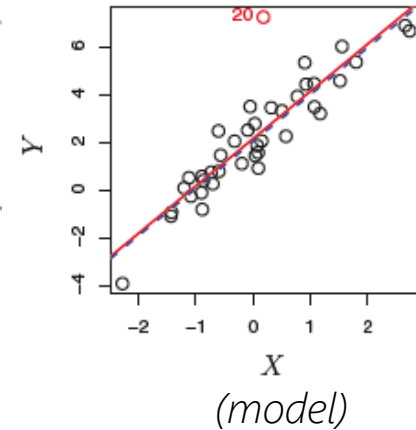
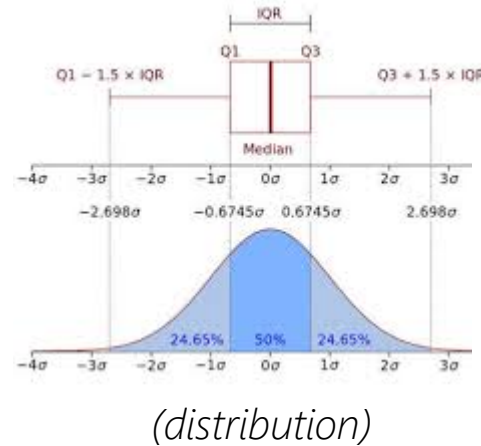
Data Mining Function: (4) Cluster Analysis

- *Unsupervised learning (i.e., Class label is unknown)*
- *Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns*
- *Principle: Maximizing intra-class similarity & minimizing interclass similarity*
- *Many methods and applications*



Data Mining Function: (5) Outlier Analysis

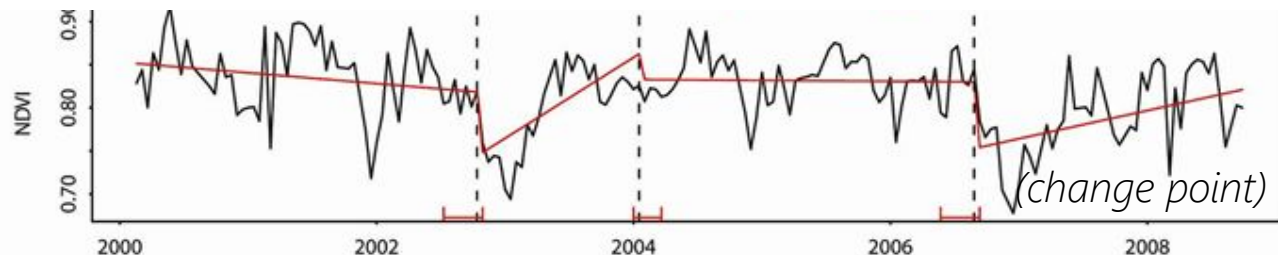
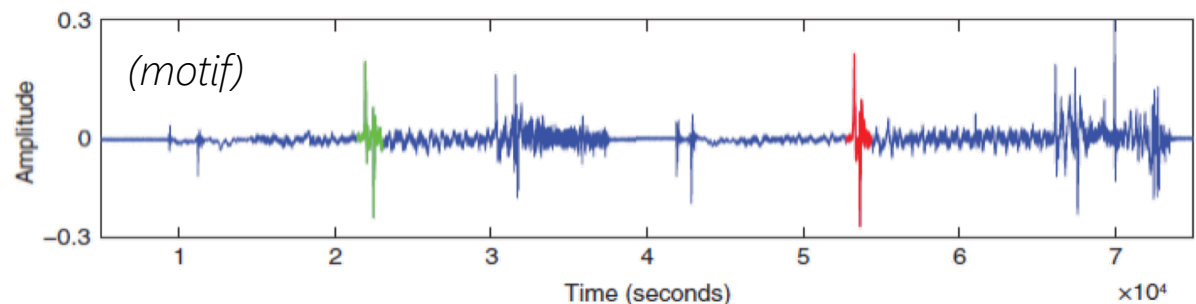
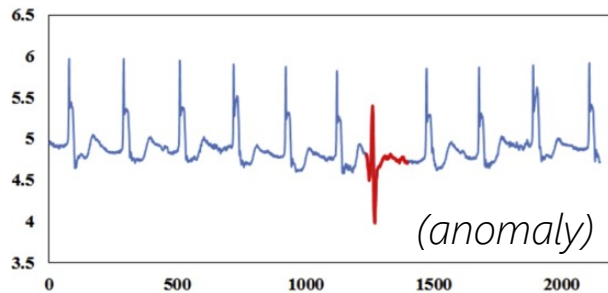
- *Outlier: A data object that does not comply with the general behavior of the data*
- *Noise or exception? — One person's garbage could be another person's treasure*
- *Methods: statistical, clustering or regression analysis*
- *Useful in fraud detection, rare events analysis*



Data Mining Function:

(6) Event detection

- *Anomalies (distribution, distance from a model, volatility)*
 - *A pattern or observation that do not conform to expected behavior*
 - *Build from another process*
- *Motifs*
 - *A pattern (unknown) that occurs a significant number of times in a dataset*
- *Change points / concept drifts*
 - *Points (or time intervals) that mark significant change in dataset behavior*
 - *They separate different states in the process that generates the dataset*



Data Mining Function: (7) Sequential Pattern

- *Sequential pattern mining*
 - *Detect and analyze frequent subsequences of events, items, or tokens occurring in an ordered metric space*



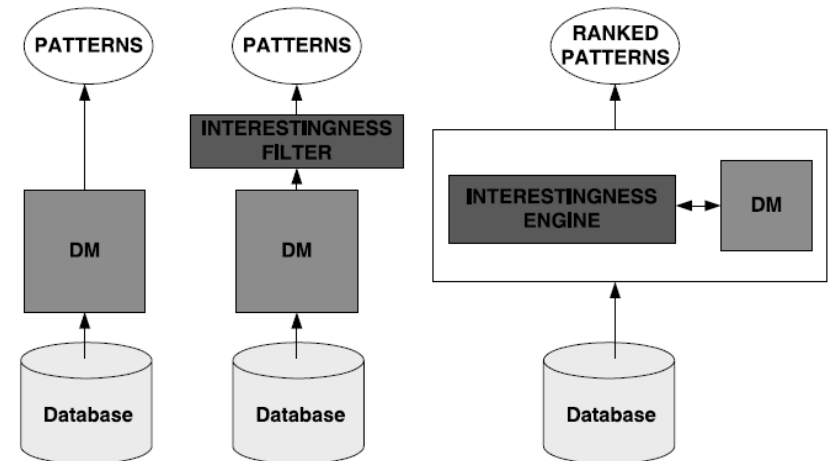
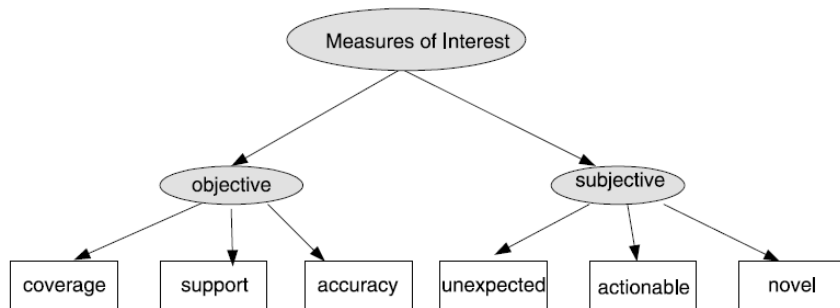
Data Mining Function: (8) Structure and Network Analysis

- *Graph mining*
 - *Finding frequent subgraphs (e.g., chemical compounds), trees (XML), substructures (web fragments)*
- *Information network analysis*
 - *Social networks: actors (objects, nodes) and relationships (edges)*
 - *e.g., Web community discovery*
- *Web mining*
 - *Opinion mining, topic detection, sentiment analysis*

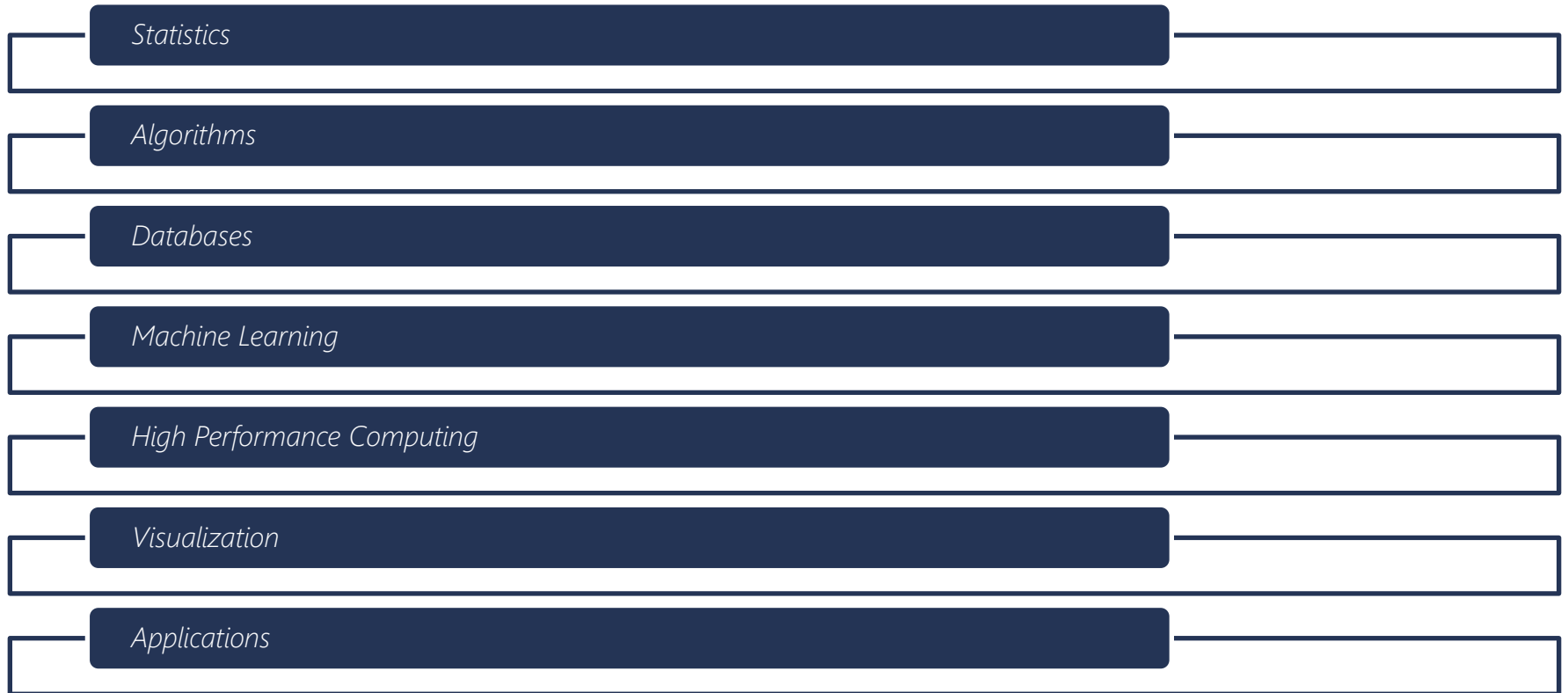


Evaluation of Knowledge

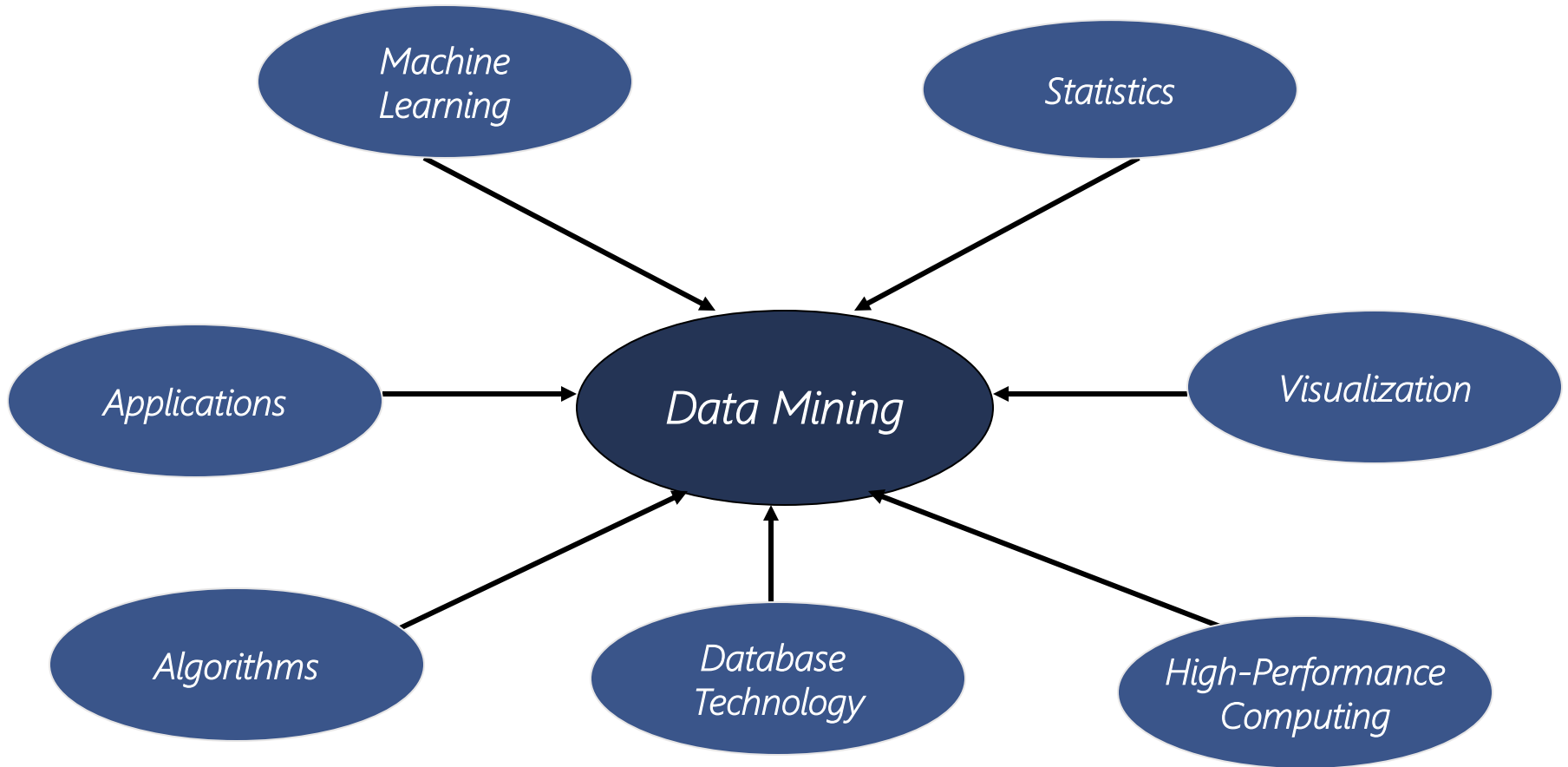
- Are all mined knowledge interesting?
 - One can mine the tremendous amount of "patterns"
 - Some may fit only certain dimension space (time, location)
 - Some may not be representative, may be transient
- Evaluation of mined knowledge → directly mine only interesting knowledge?



Data Mining: Confluence of Multiple Disciplines



Data Mining: Confluence of Multiple Disciplines



Why Confluence of Multiple Disciplines?

- *Tremendous amount of data*
 - *Algorithms must be scalable to handle big data*
- *High-dimensionality of data*
- *High complexity of data*
 - *Data streams, sensor data, spatial-temporal, text, multimedia*
- *New and sophisticated applications*

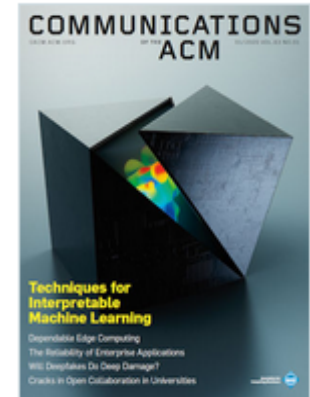
Availability of data

- *Do you have access to the data?*
- *Can you use the data?*
- *Can you publish your results?*
- *Is it big or small data?*
 - *Is enough to be considered data mining?*

Characteristic	Small data	Big data
Volume	Limited to large	Very large
Exhaustivity	Samples	Entire populations
Resolution and indexicality	Coarse and weak to tight and strong	Tight and strong
Relationality	Weak to strong	Strong
Velocity	Slow, freeze-framed	Fast
Variety	Limited to wide	Wide
Flexible and scalable	Low to middling	High

Major Issues in Data Mining

- *Mining Methodology*
 - *Mining knowledge in multi-dimensional space*
 - *Handling noise, uncertainty, and incompleteness of data*
 - *Pattern evaluation and constraint-guided pattern mining*
- *User Interaction*
 - *Interactive mining*
 - *Incorporation of background knowledge*
 - *Presentation and visualization of data mining results*
- *Efficiency and Scalability*
 - *Efficiency and scalability of data mining algorithms*
 - *Parallel, distributed, stream, and incremental mining methods*
- *Diversity of data types*
 - *Handling complex types of data*
 - *Mining dynamic, networked, and global data repositories*
- *Data mining and society*
 - *Social impacts of data mining*
 - *Privacy-preserving data mining*



Data Mining & Data Science

- Overview
 - *Science of data or the study of data*
- *Disciplinary view*
 - *Data Science is a new interdisciplinary field that synthesizes and builds on Statistics, Computer Science, Communication, Management and Sociology to study data and its environments (including domain and contextual aspects) in a way to transform data into knowledge and decisions*



[1] V. Dhar, 2013, *Data science and prediction*, *Communications of the ACM*, v. 56, n. 12, p. 64–73.

[2] L. Cao, 2017, *Data science: Challenges and directions*, *Communications of the ACM*, v. 60, n. 8, p. 59–68.

Data Mining & Data Analytics

- *Data Analytics*
 - *Theories, technologies, tools and processes that allow the understanding and discovery from data*
 - *Entire process of knowledge discovery: selection, pre-processing, transformation, data mining and interpretation*
- *Descriptive Analytics*
 - *Refers to the type of data analysis that normally uses statistics to describe the data used*
- *Predictive Analytics*
 - *Refers to the type of data analysis that makes predictions about unknown future events and reveals the reasons behind them, usually through advanced analysis*
- *Prescriptive Analytics*
 - *Refers to the type of data analysis that optimizes referrals and recommends actions for smart decision making*
- *Business Intelligence*
 - *Application of Data Analytics to support business decisions*



[1] C.-W. Tsai, C.-F. Lai, H.-C. Chao, e A.V. Vasilakos, 2015, *Big data analytics: a survey*, *Journal of Big Data*, v. 2, n. 1

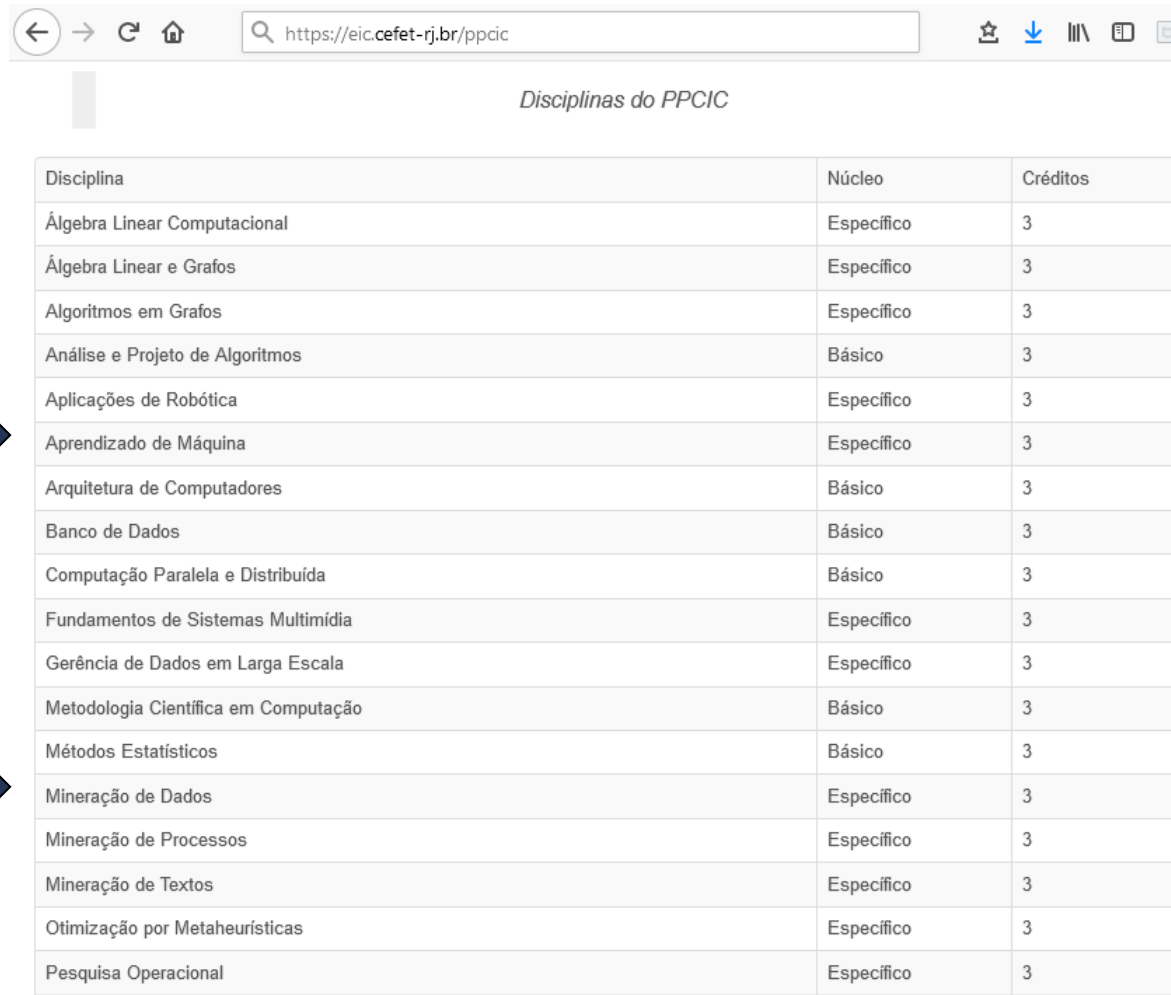
[2] L. Cao, 2017, *Data science: A comprehensive overview*, *ACM Computing Surveys*, v. 50, n. 3

Where to publish?

- *Data mining, KDD, Data Science*
 - *Conferences: SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, IEEE-DSAA*
 - *Journal: Data Mining and Knowledge Discovery, Statistical Analysis and Data Mining, ACM Transactions on Knowledge Discovery from Data*
- *Database systems*
 - *Conferences: SIGMOD, PODS, VLDB, IEEE-ICDE, EDBT, ICDT, SSDBM*
 - *Journals: IEEE-TKDE, VLDB J., Info. Sys.*
- *AI & Machine Learning*
 - *Conferences: AAI, ML, IJCNN, IJCAI, NeurIPS*
 - *Journals: Machine Learning, Artificial Intelligence, Knowledge and Information Systems*
- *Statistics*
 - *Journals: Journal of Applied Statistics, Annals of Data Science*
- *Applications*
 - *Journals*

PPCIC: Master Degree on Computer Science @ CEFET/RJ

Focus on Data Science



Disciplinas do PPCIC

Disciplina	Núcleo	Créditos
Álgebra Linear Computacional	Específico	3
Álgebra Linear e Grafos	Específico	3
Algoritmos em Grafos	Específico	3
Análise e Projeto de Algoritmos	Básico	3
Aplicações de Robótica	Específico	3
Aprendizado de Máquina	Específico	3
Arquitetura de Computadores	Básico	3
Banco de Dados	Básico	3
Computação Paralela e Distribuída	Básico	3
Fundamentos de Sistemas Multimídia	Específico	3
Gerência de Dados em Larga Escala	Específico	3
Metodologia Científica em Computação	Básico	3
Métodos Estatísticos	Básico	3
Mineração de Dados	Específico	3
Mineração de Processos	Específico	3
Mineração de Textos	Específico	3
Otimização por Metaheurísticas	Específico	3
Pesquisa Operacional	Específico	3

Trending Data Mining Languages

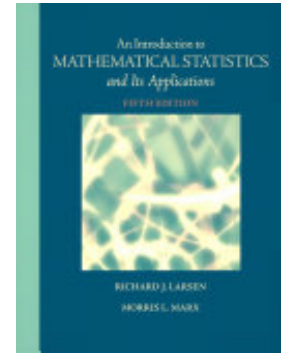
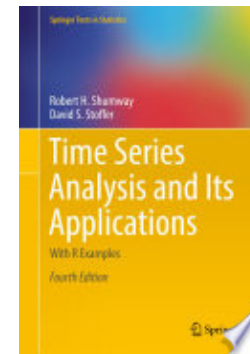
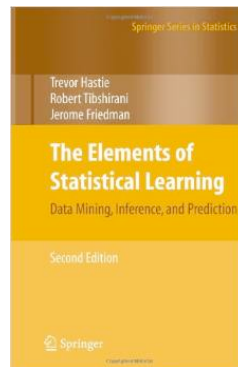
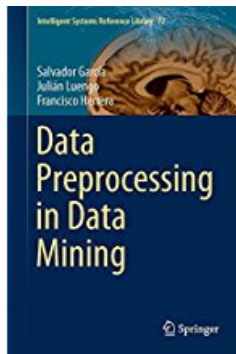
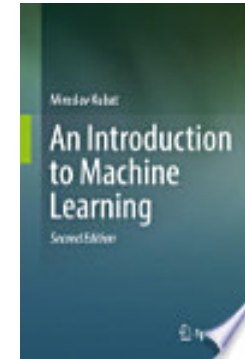
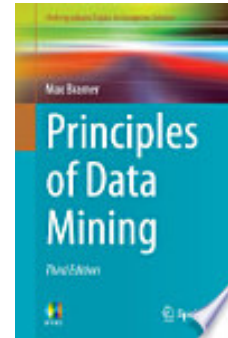
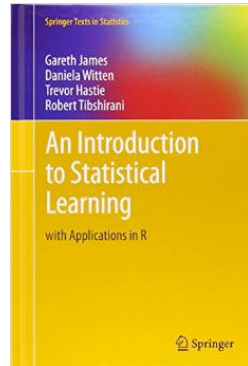
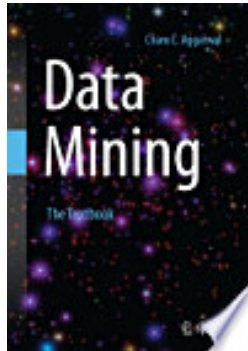
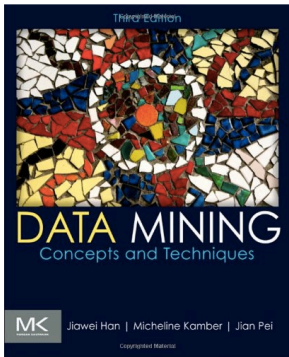
- *Python (Machine Learning Course)*
 - *Scikit learning*
- *R (Data Mining Course)*
 - *Myriad of packages*
- *Spark (Parallel and Distributed Computing)*
 - *Mlib*

Data Mining Tools

- *Rapid Miner (open source)*
 - <https://rapidminer.com>
- *Orange (open source)*
 - <https://orange.biolab.si>
- *Weka (open source)*
 - <https://www.cs.waikato.ac.nz/ml/weka>
- *Knime (open source)*
 - <https://www.knime.com>
- *Apache Mahout (open source)*
 - <http://mahout.apache.org>
- *Rattle (open source)*
 - <https://rattle.togaware.com>



Main References



Next Topics

- *Next classes*
 - *R programming language*
 - *Exploratory data analysis*
 - *Data preprocessing*
- *Slides and videos*
 - <https://eic.cefet-rj.br/~eogasawara/>
 - *Follow up YouTube channel*
 - <https://www.youtube.com/channel/UCAm1hAXWEqYJfXz4Ez zBhVg>

