

MACHINE LEARNING: why should you care?

Filip Wójcik

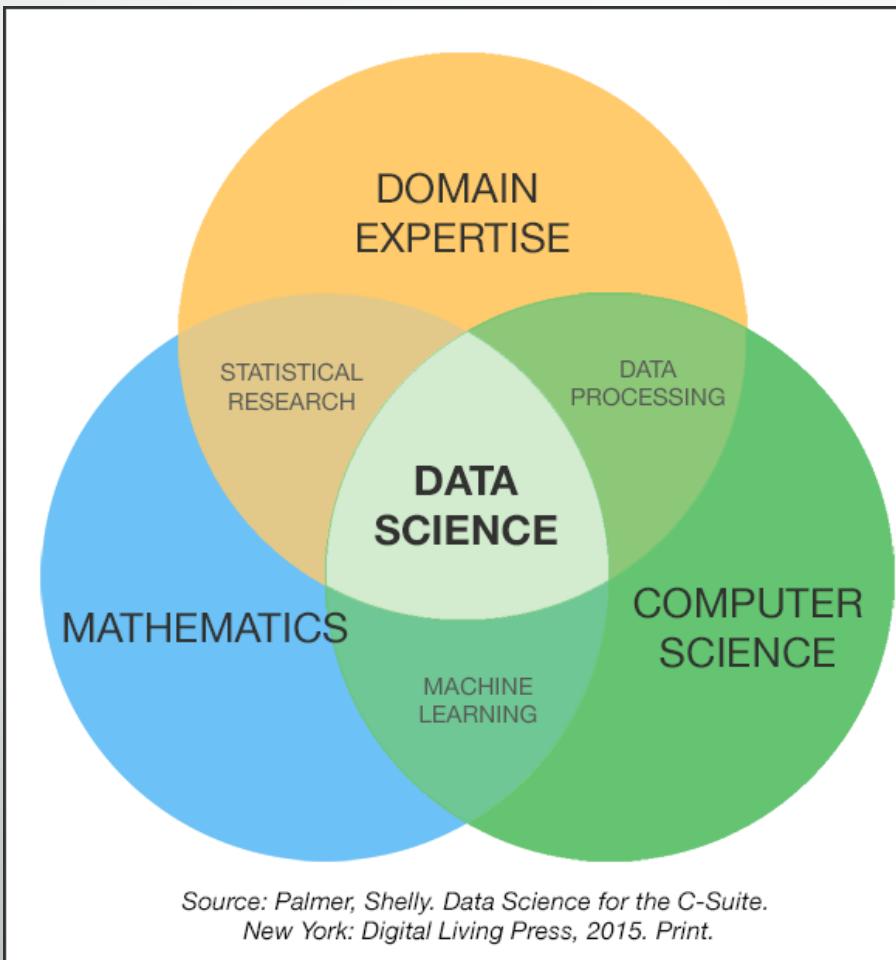
Credit Suisse

Fx Sales Analytics VICE PRESIDENT

University lecturer

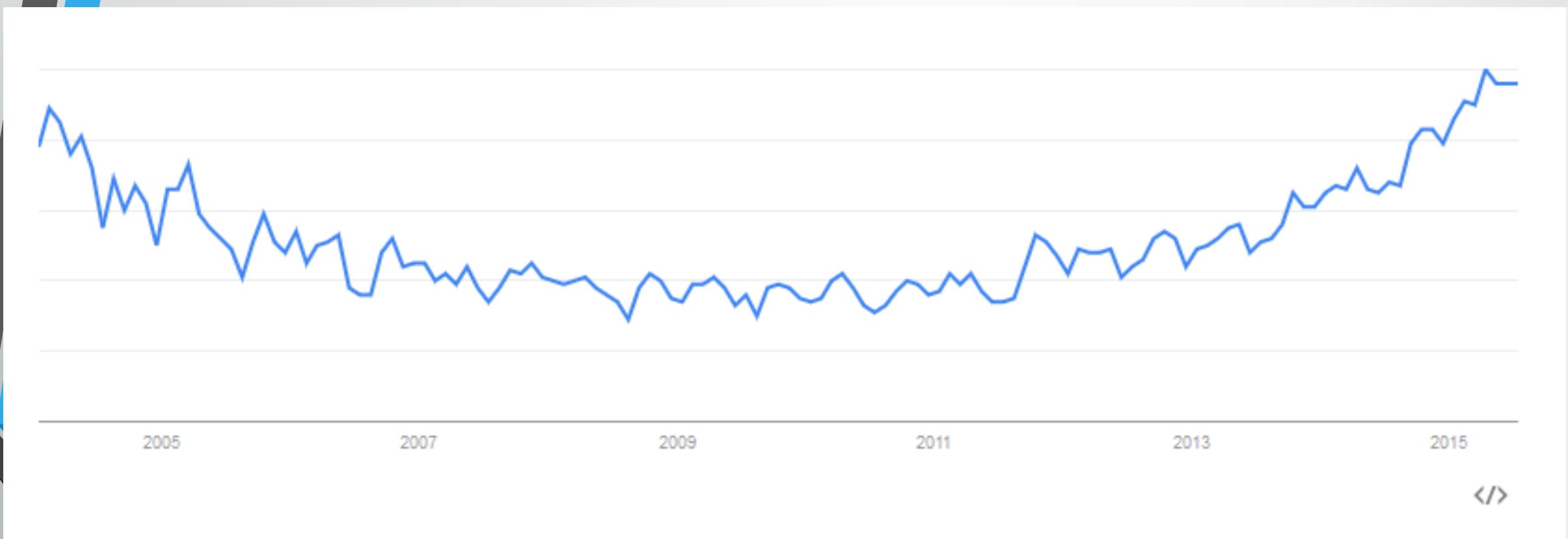
[filip.wojcik@outlook.com](mailto:filiw.wojcik@outlook.com)

WHAT IS MACHINE LEARNING?



WHAT IS MACHINE LEARNING?

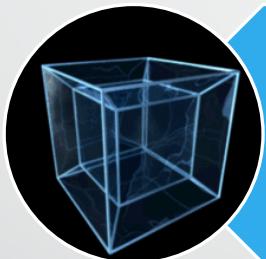
Google trends analysis for phrase „Machine learning”



WHAT IS MACHINE LEARNING?



More and more data
produced by corporations



Need to process more
complex data



Automation of data
analysis processes

WHAT IS MACHINE LEARNING?

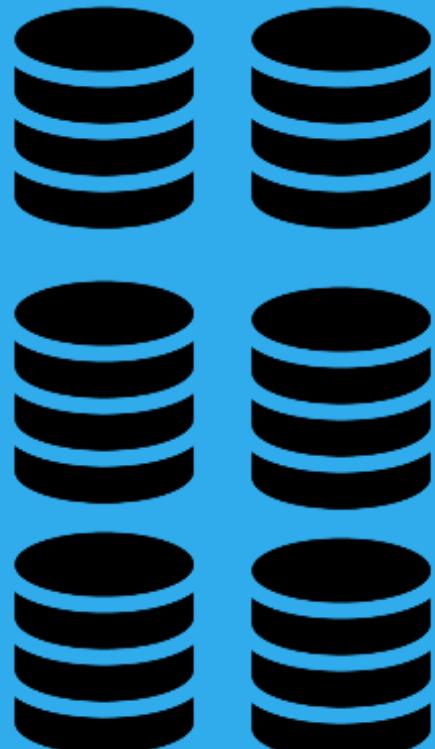
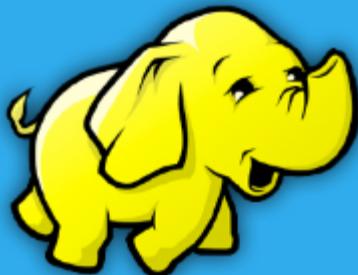
Big data

Machine learning

- Large volumes of data storage & processing
- Highly parallelized algorithms
- Sophisticated architecture
- Hardware-related (clusters, nodes, server machines)

- Smart data processing methods
- Domain-agnostic
- Technology-agnostic
- Hardware-agnostic
- Predictions and modelling
- Strongly related to statistics

WHAT IS MACHINE LEARNING?



Hadoop cluster
(data storage)

← Apply to



Machine learning algorithms
(data processing)

MACHINE LEARNING USE CASES

- Customer preferences discovery
- Automated expert systems construction
- Assigning new data to groups

- Financial trends discovery
- Statistical analysis
- Prediction of numerical values/outcomes

SUPERVISED

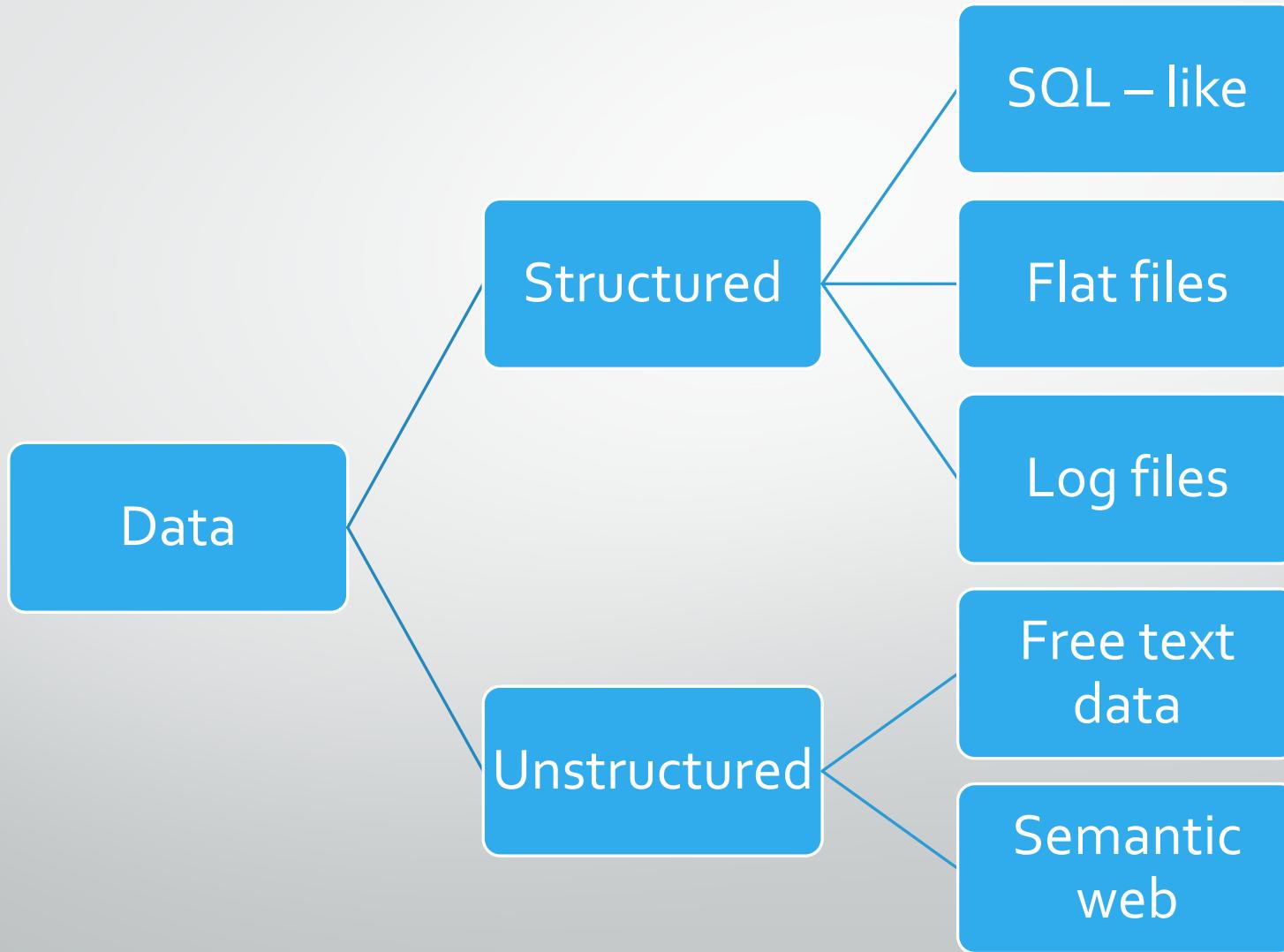
CL.

Reducing amount of data!!!

- Market basket analysis
- Discovering preferences
- Explaining data

- Detecting irrelevant features/columns
 - Detecting highly correlated features/columns
 - Detecting noise
- Customers grouping
 - Discovering similarities
 - Features importance recognition

DATA FOR MACHINE LEARNIG



DATA FRAME – BASIC DATA STRUCTURE OVERVIEW

Features/attributes



Discrete features

Boolean feature

Numerical feature

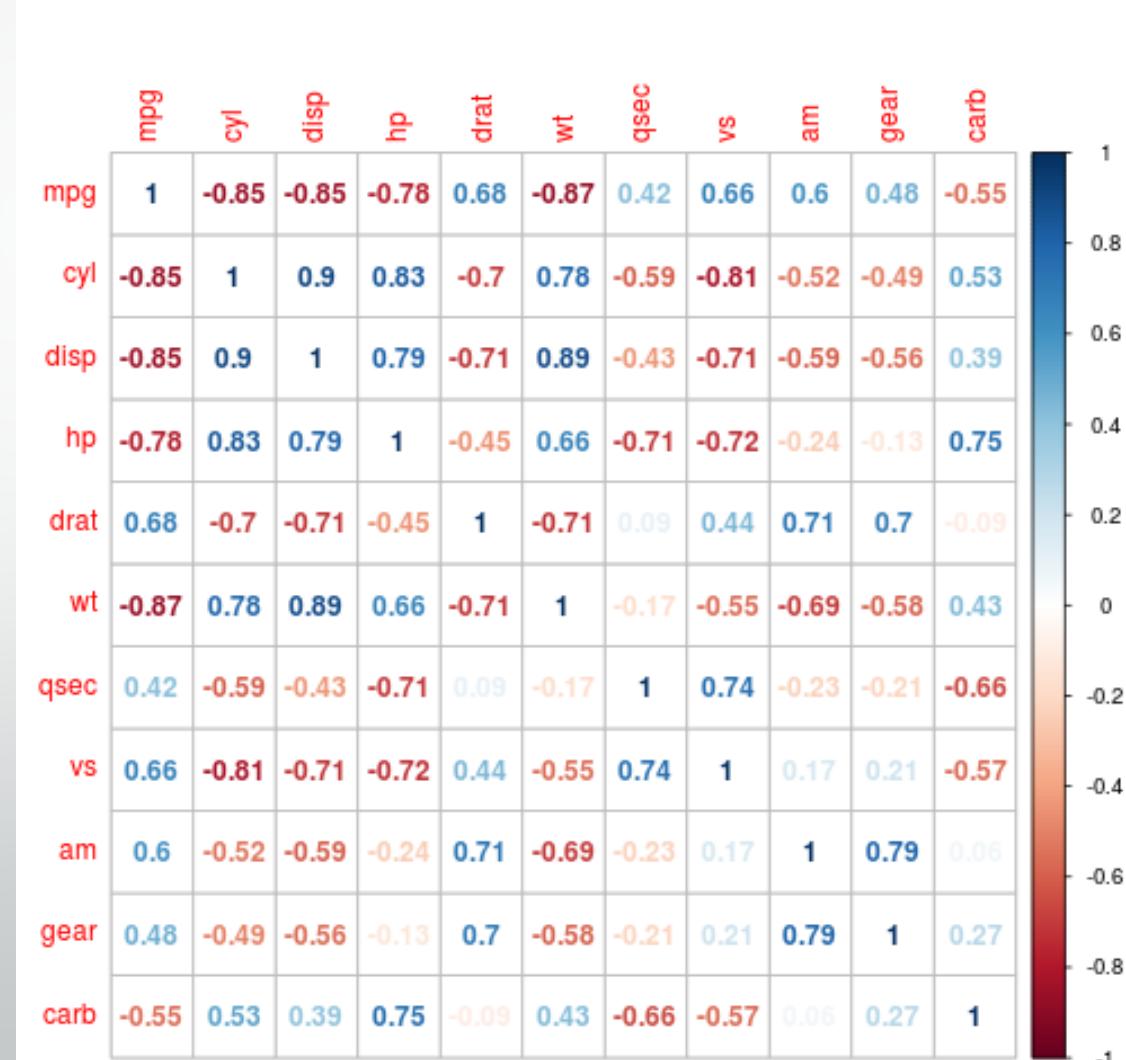
Company	Financial instruments	Status	Revenue
Company X	Equities	Open	0.6
Company Y	Corporate Bonds	Open	0.03
Company Z	Structure hybrid	Closed	0.02



Instances

DATA FRAME – BASIC DATA STRUCTURE

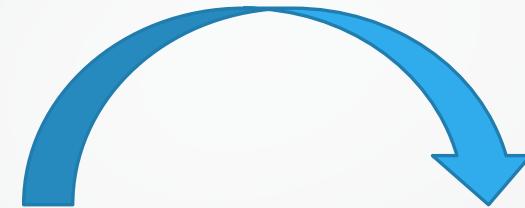
FEATURES INVESTIGATION



DATA FRAME – BASIC DATA STRUCTURE

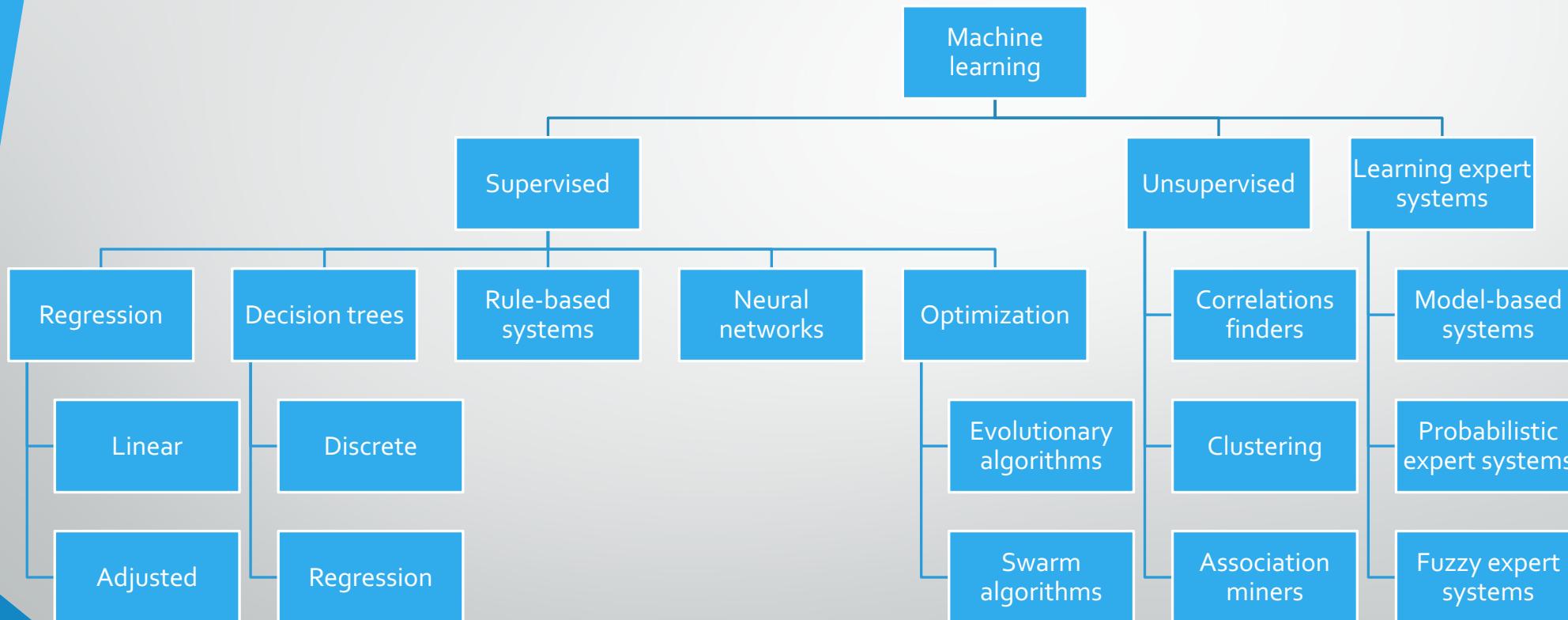
FEATURES ENCODING

Company	Financial instruments	Status	Revenue
Company X	Equities	Open	0.6
Company Y	Corporate Bonds	Open	0.03
Company Z	Structure hybrid	Closed	0.02



Company	Financial instruments	Status	Revenue
001	001	1	0.6
010	010	1	0.03
100	100	0	0.02

ALGORITHMS OVERVIEW



SUPERVISED LEARNING

SUPERVISED LEARNING

Two data sets

- Training– known „answers”, given to algorithm
- Test– known „answers”, not given to algorithm

“Teacher/oracle”

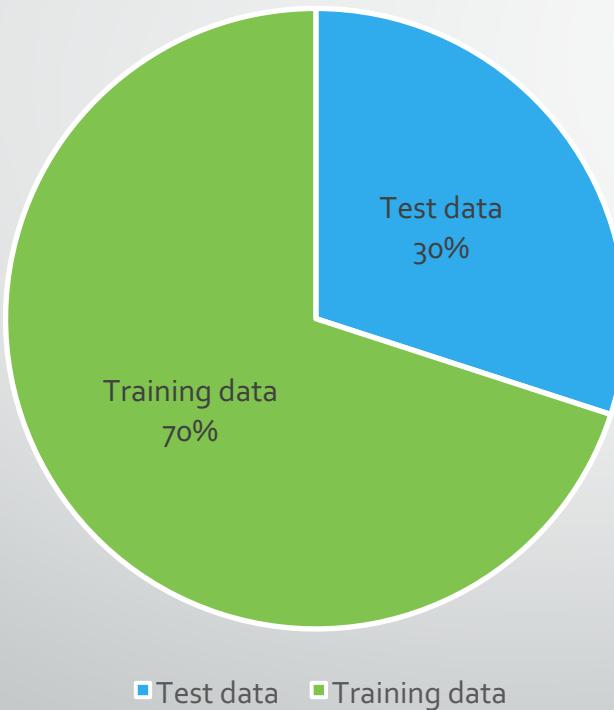
- Objective rating function
- Checks the algorithm progress

Learning based
on the
experience

- Application of teachers/oracle suggestions to improve score
- Avoiding overfitting

SUPERVISED LEARNING

Data partitioning

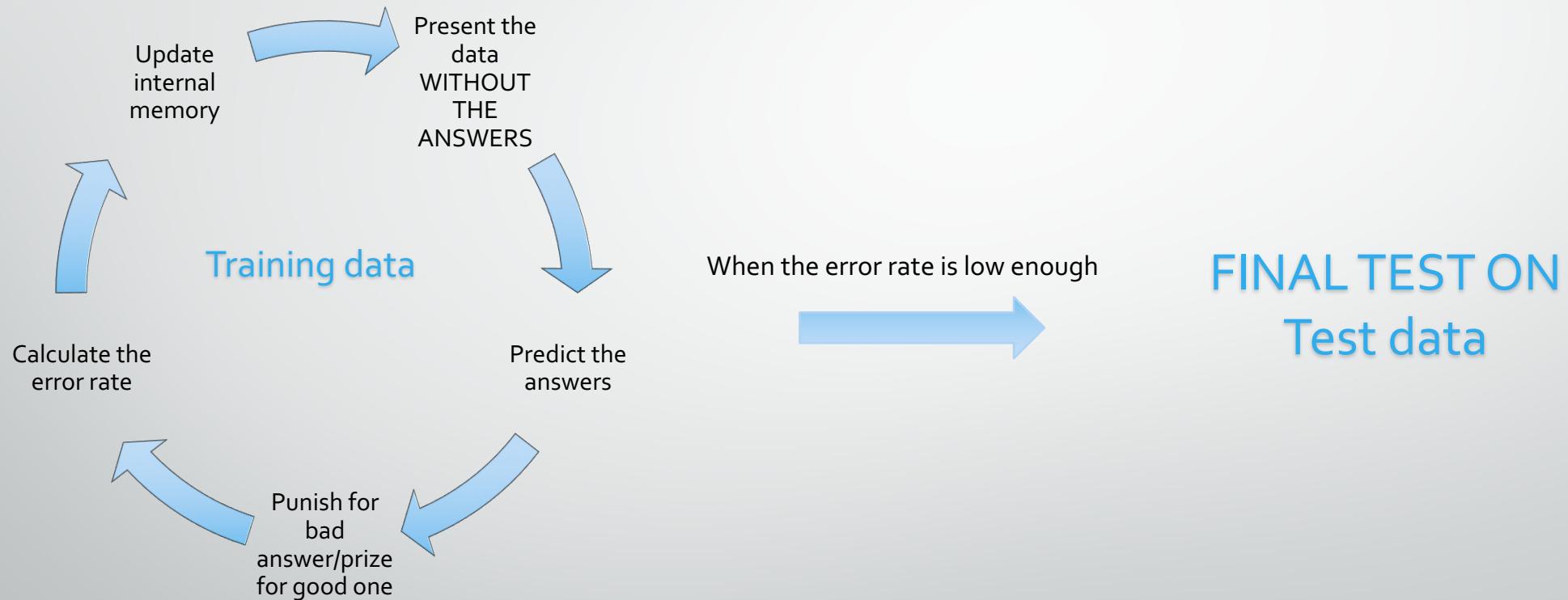


Sometimes the amount of data with known „answers” is limited

Data division helps in better controlling the learning process

Improving the effectiveness of data usage

SUPERVISED LEARNING



SUPERVISED LEARNING

Decision trees

General approach

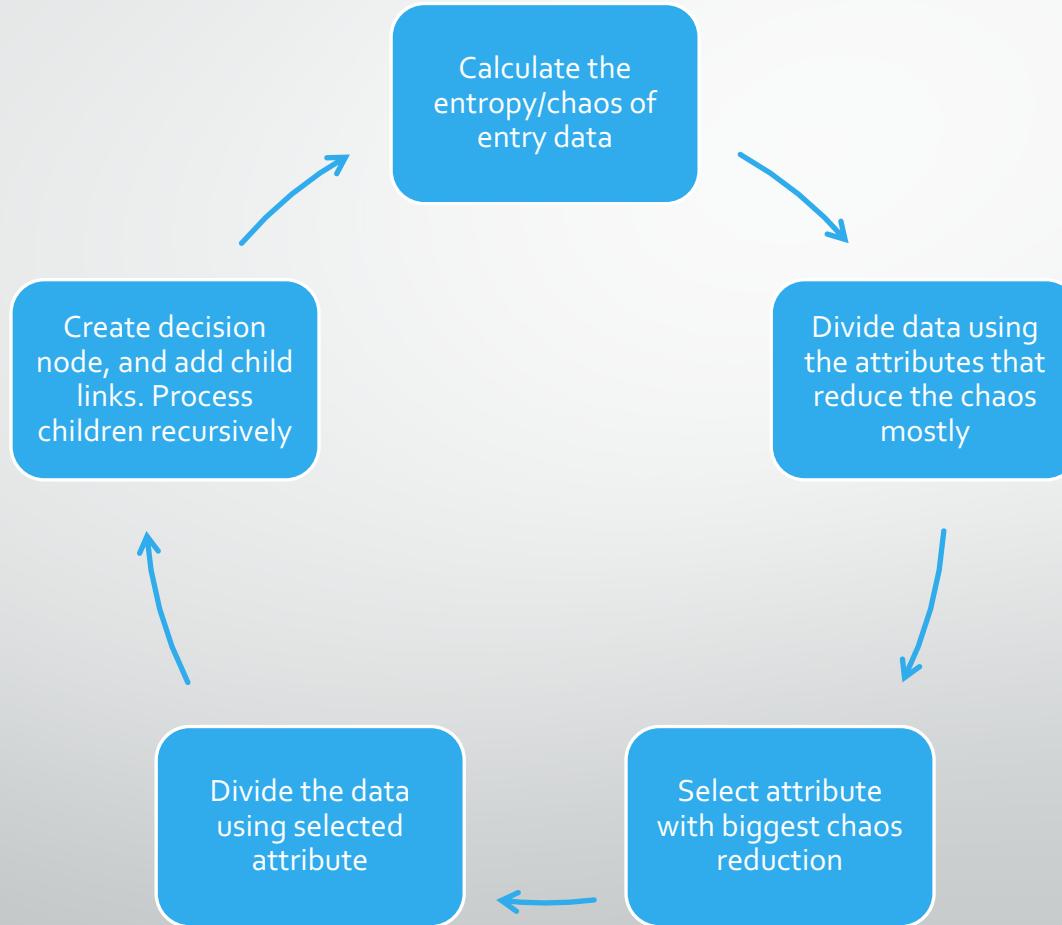
- Uses structured data
- Recursive top-down approach: divide and conquer, based on the best-promising attributes
- Can use numerical and discrete data as well

Pros

- Very flexible
- Easy to implement
- Easy to interpret by humans
- Can be translated to easy-to-read rules and included in reports/documentations

SUPERVISED LEARNING

Decision trees



SUPERVISED LEARNING

Decision trees

client	hotel	addons	money_spent	offer
business	Hilton	trip	40000	deluxe
business	Hilton	full board	38000	deluxe
business	Hilton	trip	40000	deluxe
middle class	Meta	none	800	basic
middle class	Meta	meal	900	basic
manager	Meta	spa	1500	premium

Value	Count	%
Deluxe	3	0.5
Basic	2	0.333
Premium	1	0.16666



client	hotel	addons	money_spent	offer
business	Hilton	trip	40000	deluxe
business	Hilton	full board	38000	deluxe
business	Hilton	trip	40000	deluxe
middle class	Meta	none	800	basic
middle class	Meta	meal	900	basic
manager	Meta	spa	1500	premium



Split on: „client”
Value to split: „business”

Client == business ?

True

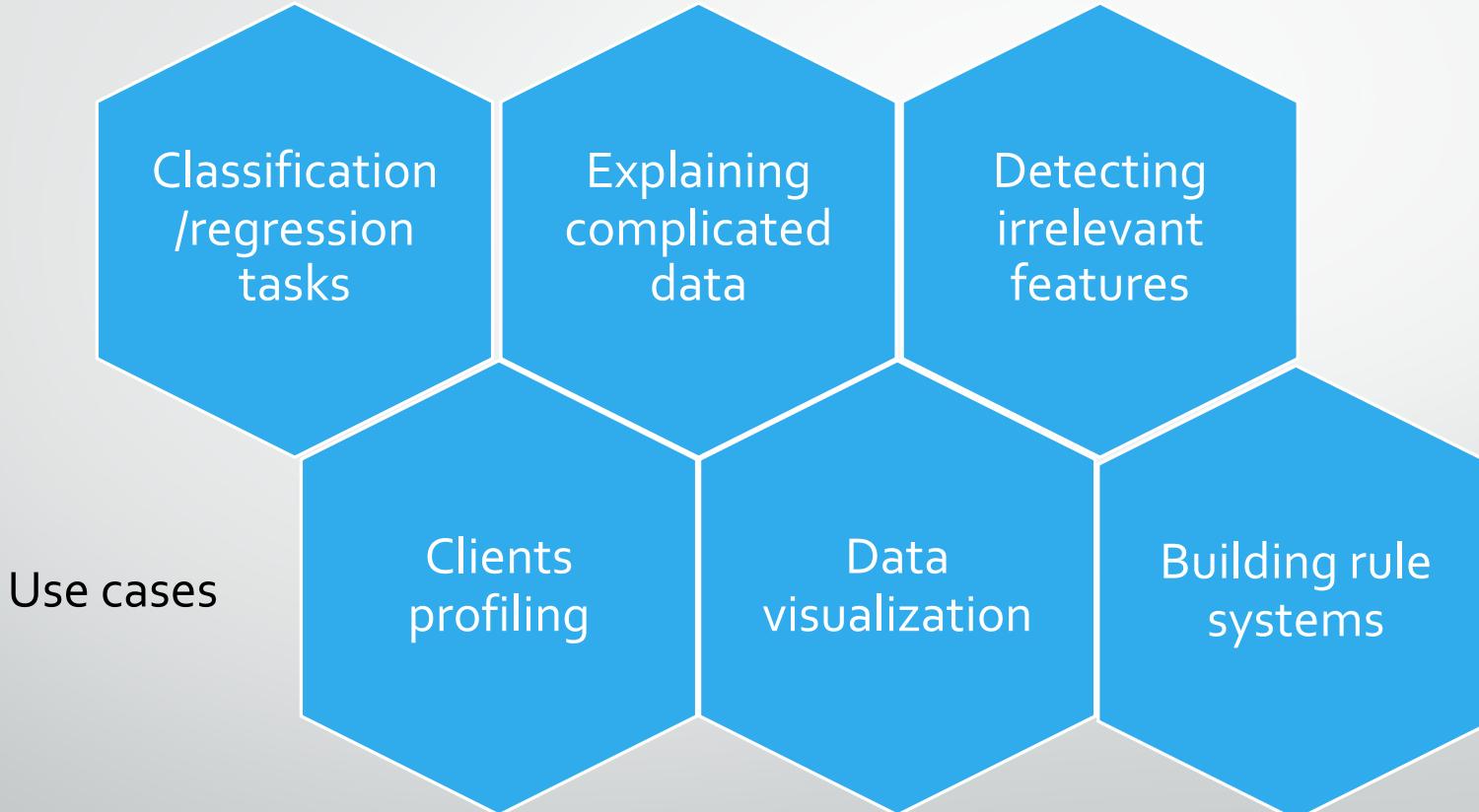
False

hotel	addons	money_spent	offer
Hilton	trip	40000	deluxe
Hilton	full board	38000	deluxe
Hilton	trip	40000	deluxe

hotel	addons	money_spent	offer
Meta	none	800	basic
Meta	meal	900	basic
Meta	spa	1500	premium

SUPERVISED LEARNING

Decision trees



UNSUPERVISED LEARNING

UNSUPERVISED LEARNING

One data set

- Single set of data
- No „good answers” provided (in most cases)

No
teacher/oracle

- No option to evaluate prediction against „correct answers”
- Algorithm evaluation based on similarity measures/chaos measures/etc.

Algorithm
operates on
data on its own

- Algorithm explores the possible data partitioning
- Algorithm maintains its internal error measures

UNSUPERVISED LEARNING

Association rules learning

General approach

- Ordered data
- Searching for coincidences/correlations in data

Features

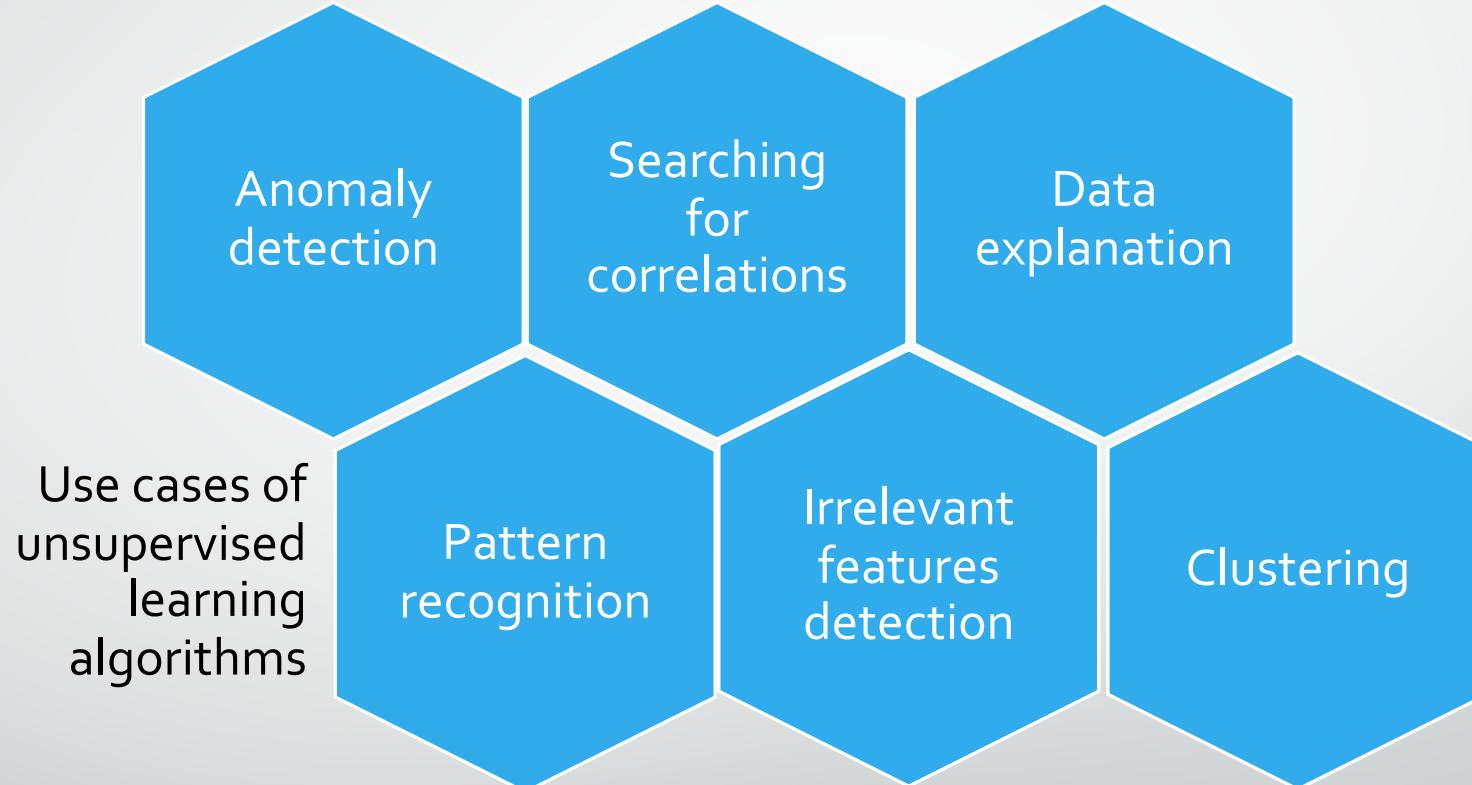
- Easy to implement
- Flexible
- Easy to interpret by humans
- Can significantly reduce the amount of irrelevant features

UNSUPERVISED LEARNING

Association rules learning

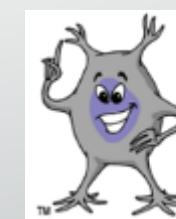
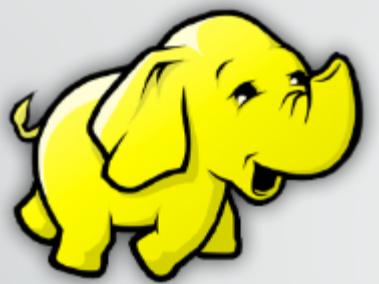
Transaction no.	Products
1.	1. Soya milk 2. Salad
2.	1. Salad 2. Walnuts 3. Wine 4. Bread
3.	1. Soya milk 2. Walnuts 3. Wine 4. Juice
4.	1. Salad 2. Soya milk 3. Walnuts 4. Wine
5.	1. Salad 2. Soya milk 3. Walnuts 4. Juice

UNSUPERVISED LEARNING



TOOLS OF THE TRADE

MACHINE LEARNING TOOLS



Encog