# PRESENTATION
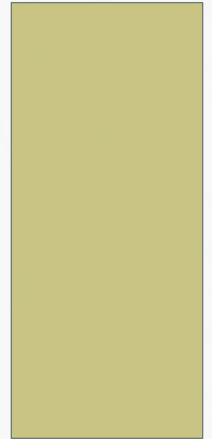
## BIG DATA INTELLIGENCE

### METHODS AND TECHNOLOGIES

## A.K.A. MACHINE LEARNING I

RICARDO ALER MUR (aler@inf.uc3m.es). 2.2B29

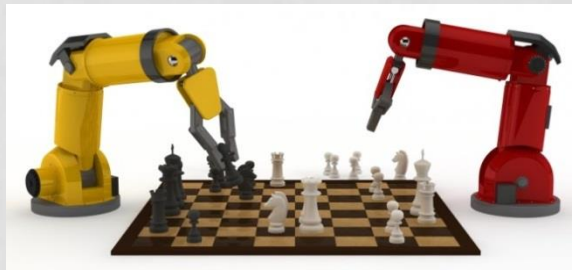**MASTER IN BIG DATA ANALYTICS**

# GOALS

1. To introduce the basics of **Machine Learning**: training, testing, models, hyper-parameter tuning, etc.
2. To show how to do Machine Learning in a **Big Data** Context
3. To apply the above goals with currently used **tools**
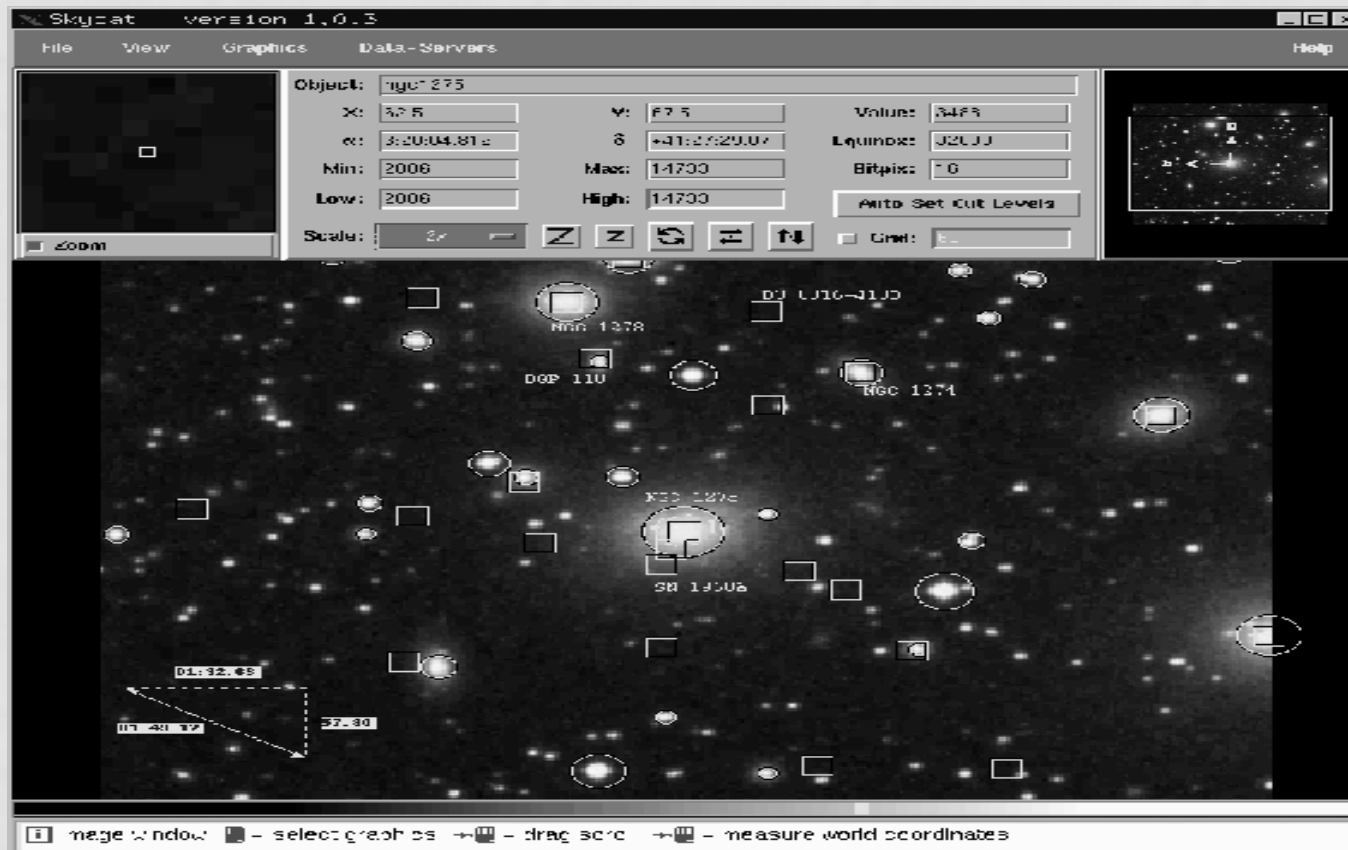
# MACHINE LEARNING

- In general, it's a subfield of **Artificial Intelligence** that tries to make computers and machines learn
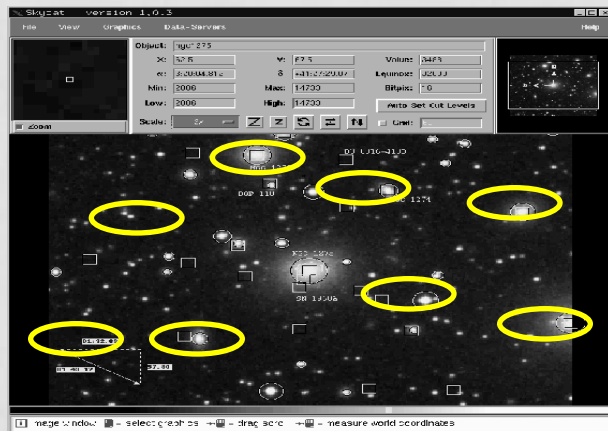


- In practice, it tries to create models from data and thus is closely related to statistics. This is the point of view we will follows in this course
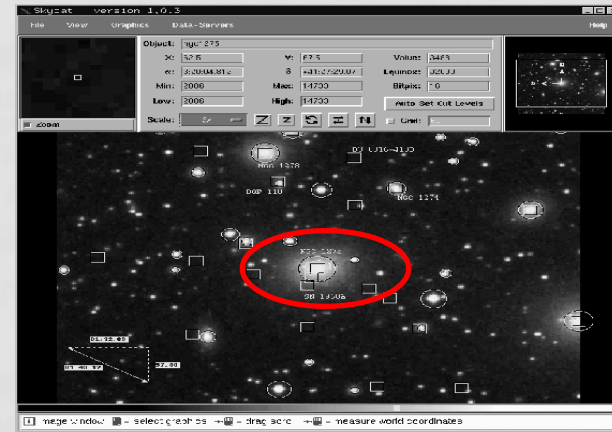
# WHAT IS MACHINE LEARNING

- Example: Skycat: AUTOMATIC CLASSIFICATION OF OBJECTS IN THE SKY

?

**Training data (labeled pictures of sky objects: galaxies, stars, nebulae, …)**

ML Algorithm

**Model**

Spiral galaxy

Pictures in the catalog have been labeled by a human expert (astronomer)

# APPLICATIONS

- Finances and banking
  - Credit card fraud detection
  - Credit default prediction
- Market analysis:
  - Market basket analysis
  - Market segmentation
- Insurance:
  - Expensive clients
- Education:
  - Prediction of school dropouts
- Industry:
  - Electric (energy) load forecasting
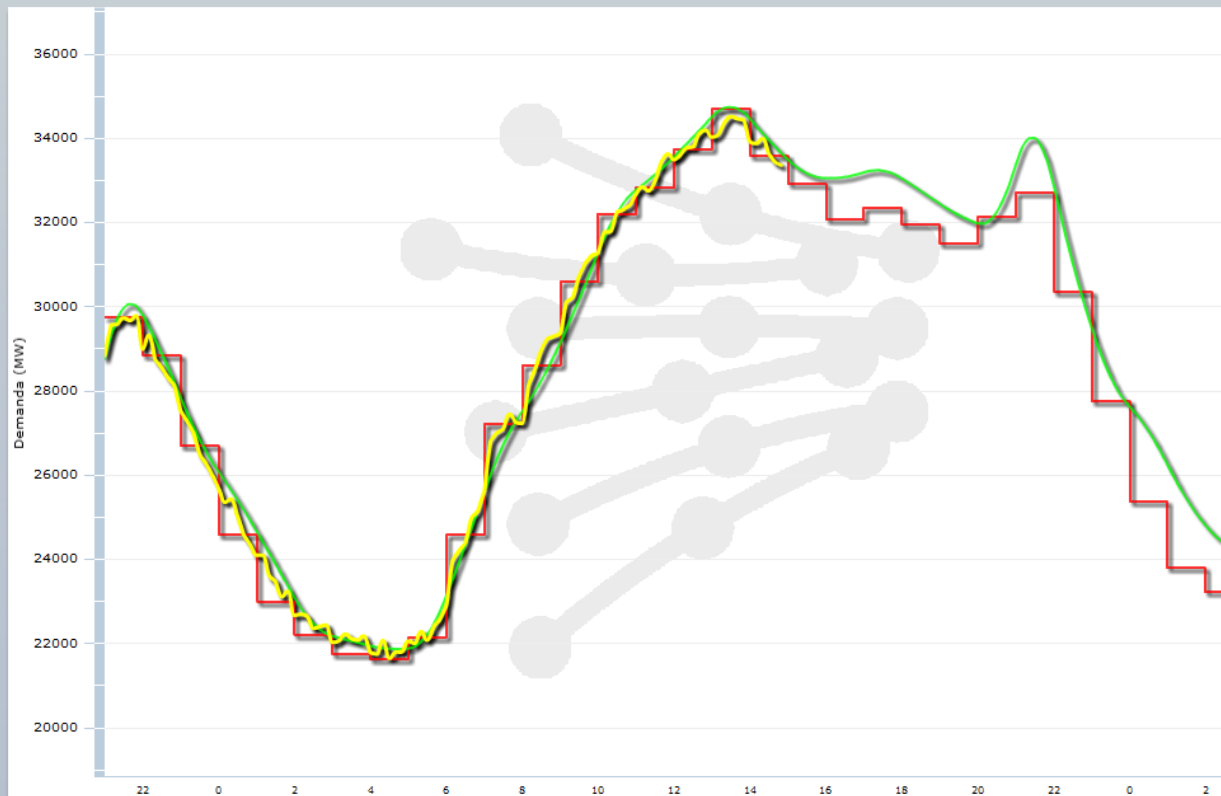  - Solar / wind energy forecasting

https://demanda.ree.es/demanda.html

# ELECTRIC LOAD FORECASTING

# APPLICATIONS II

- Medicine:
  - Illness diagnosis
- Science:
  - Illness prediction from DNA analysis
  - Prediction if a new substance causes cancer
  - SKYCAT
- Internet:
  - Spam detection (SpamAssassin)
  - Web: book recommendation (amazon.com)

- Longitud: 1574 páginas (estimación) ☑
- ¿No tienes un Kindle? Consigue un Kindle aquí.

**Descubre cómo ahorrar hasta un 90% en un título diferente cada día**
Inscríbete en la Newsletter Kindle Flash y recibe directamente en tu bandeja de entrada la oferta del día Kindle Flash para no perderte ni un título en promoción. Más información

minotauro
Edición **Kindle**

⊕ **ZOOM**
Ver imagen ampliada (con el zoom)
Compartir mis imágenes de cliente

Añadir a la Lista de deseos
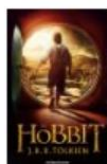
**Pruébalo gratis**
Lee el principio de este eBook gratis

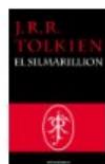Enviar fragmento ya

**Enviar a mi Kindle o a otro dispositivo**

Cómo obtener los fragmentos
Disponible en Windows

## Los clientes que compraron este producto también compraron

Página 1 de 2

HOBBIT
J.R.R.TOLKIEN

**El Hobbit**
J. R. R. Tolkien
★★★★☆ (16)
Versión Kindle
EUR 6,64

J.R.R.TOLKIEN
EL SILMARILLION

**El Silmarillion**
J. R. R. Tolkien
★★★★★ (2)
Versión Kindle
EUR 6,64

MISAL ROMANO

**Misal Romano Completo**
Jose Quintana Velasco
★★★☆☆ (2)
Versión Kindle
EUR 4,24

J.R.R.TOLKIEN
EL LIBRO DE LOS CUENTOS PERDIDOS I

**El Libro de los Cuentos Perdidos, 1. Historia de la Tierra Media, I ...**
J. R. R. Tolkien
Versión Kindle
EUR 6,64

Lee libros en tu ordenador o en otros dispositivos móviles gracias a nuestras Aplicaciones de lectura Kindle **GRATUITAS.**

Compartir ✉ ⬛ ⬛ ⬛

## Descripción del producto

### Descripción del producto

Concebida en un primer momento como una continuación de El Hobbit, acabó por convertirse en una historia independiente por derecho propio de mucho más alcance y extensión. En 1999 la trilogía de El Señor de los Anillos fue elegida como «Libro del Milenio» por los participantes de una encuesta de Amazon.com. En la adormecida e idílica Comarca, un joven hobbit recibe un encargo: custodiar el Anillo Único y emprender el viaje para su destrucción en las Grietas del Destino. Consciente de la importancia de su misión, Frodo abandona la Comarca e inicia el camino hacia Mordor con la compañía inesperada de Sam, Pippin y Merry. Pero sólo con la ayuda de Aragorn conseguirán vencer a los Jinetes Negros y alcanzar el refugio de la Casa de Elrond en Rivendel.

# SYLLABUS

1. Overview and introduction to Machine Learning: tasks and models.
2. Predictive models:
   - Decision trees, regression trees
   - K Nearest Neighbour (KNN)
   - Machine Learning pipeline: training, => ML algorithm => model => test / evaluation . Preprocessing, hyperparameter tuning, …
3. Ensemble methods: bagging, boosting, stacking
4. Preprocessing: selection of attributes and methods of dimensionality reduction
5. Machine learning software for Big Data:
   1. Python: scikit-learn, numpy
   2. Mapreduce
   3. Spark: pyspark, MLLIB
6. Other topics:
   1. Online learning
   2. Metaheuristics: genetic algorithms, genetic programming, …

# SYLLABUS

1. Introduction to Machine Learning: tasks, algorithms & models

2. Methods for training classification and regression models:

   - Nearest Neighbour (KNN)
   - Decision / regression trees & rules

3. Methodology and the Machine Learning pipeline

4. Methods for attributes (feature selection, transformation)

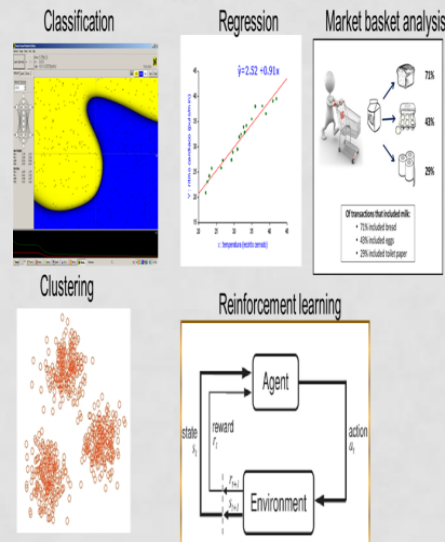5. Methods based on ensembles of models: Bagging, boosting, stacking

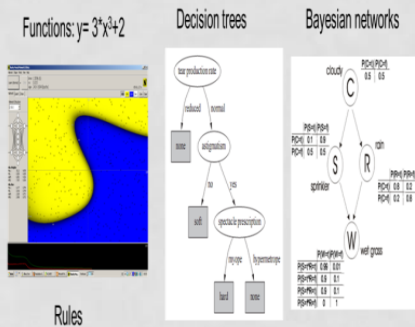6. Large Scale Machine Learning. Big Data

   - Map-reduce & Spark

7. Software tools:

   - Python + scikit-learn & Pyspark + MLIB



TASKS

Classification   Regression   Market basket analysis

Clustering   Reinforcement learning

MODELS

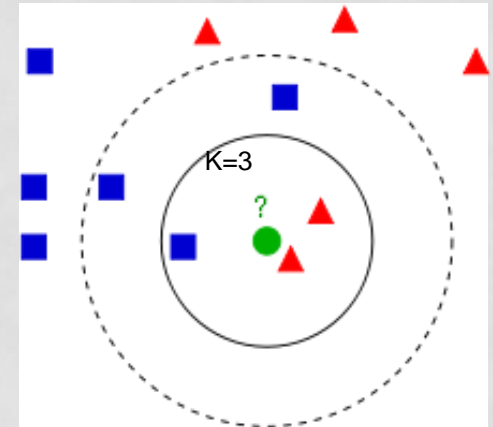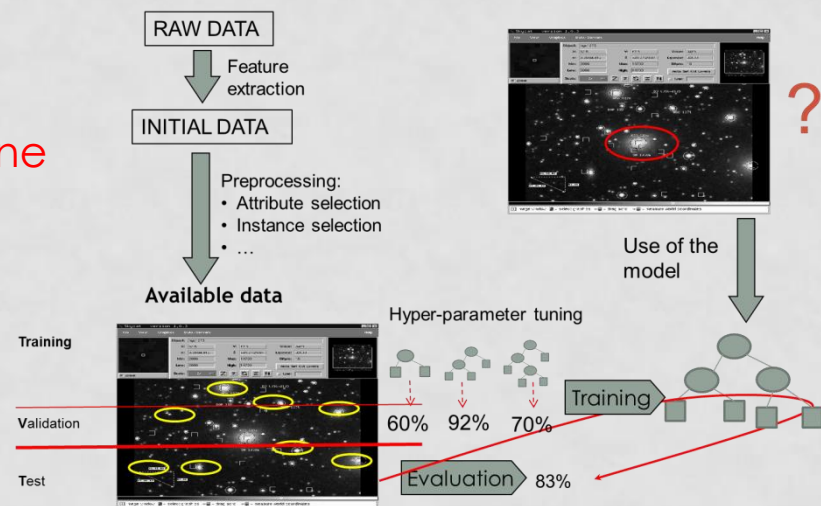Functions: y= 3*x³+2   Decision trees   Bayesian networks

Rules

If humidity = normal and windy = false then play = yes

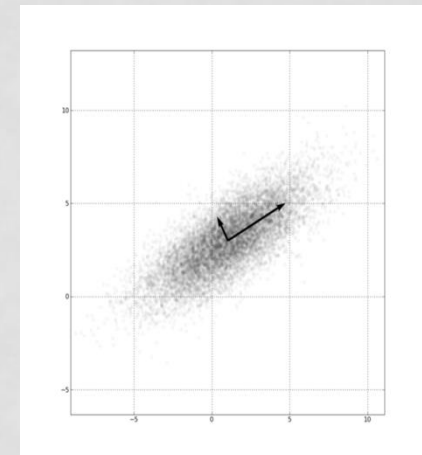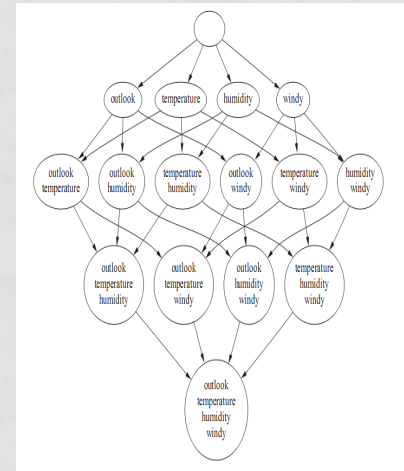And many more: neural networks, nearest neighbor, ...

# SYLLABUS

1. Introduction to Machine Learning: tasks, algorithms & models

2. Methods for training classification and regression models:

- Nearest Neighbour (KNN)
- Decision / regression trees & rules

3. Methodology and the Machine Learning pipeline

4. Methods for attributes (feature selection, transformation)

5. Methods based on ensembles of models: Bagging, boosting, stacking

6. Large Scale Machine Learning. Big Data

- Map-reduce & Spark

7. Software tools:
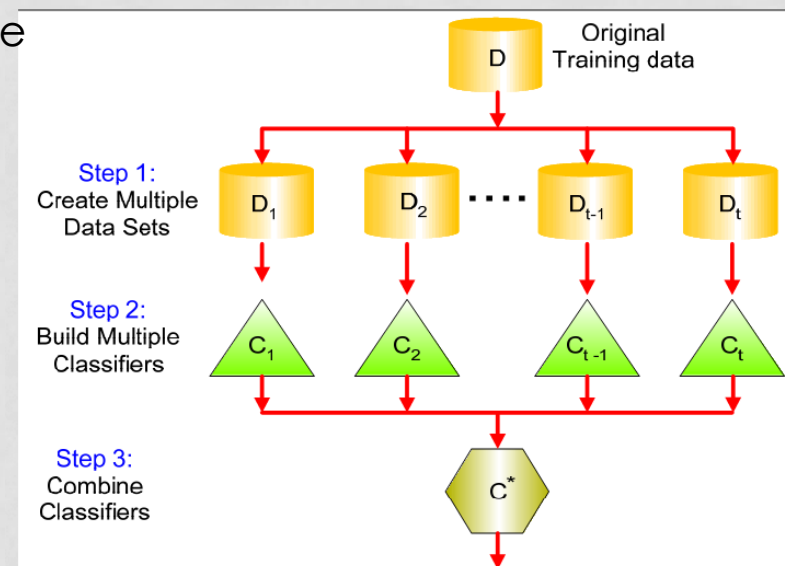
- Python + scikit-learn & Pyspark + MLIB

# SYLLABUS

1. Introduction to Machine Learning: tasks, algorithms & models

2. Methods for training classification and regression models:

- Nearest Neighbour (KNN)
- Decision / regression trees & rules

3. Methodology and the Machine Learning pipeline

4. Methods for attributes (feature selection, transformation)

5. Methods based on ensembles of models: Bagging, boosting, stacking

6. Large Scale Machine Learning. Big Data

- Map-reduce & Spark

7. Software tools:

- Python + scikit-learn & Pyspark + MLIB

# SYLLABUS

1. Introduction to Machine Learning: tasks, algorithms & models

2. Methods for training classification and regression models:

   - Nearest Neighbour (KNN)
   - Decision / regression trees & rules

3. Methodology and the Machine Learning pipeline

4. Methods for attributes (feature selection, transformation)

5. Methods based on ensembles of models: Bagging, boosting, stacking

6. Large Scale Machine Learning. Big Data

   - Map-reduce & Spark

7. Software tools:

   - Python + scikit-learn & Pyspark + MLIB

# SYLLABUS

1. Introduction to Machine Learning: tasks, algorithms & models

2. Methods for training classification and regression models:

- Nearest Neighbour (KNN)
- Decision / regression trees & rules

3. Methodology and the Machine Learning pipeline

4. Methods for attributes (feature selection, transformation)

5. Methods based on ensembles of models: Bagging, boosting, stacking

6. Large Scale Machine Learning. Big Data

- Map-reduce & Spark

7. Software tools:

- Python + scikit-learn & Pyspark + MLIB

# SYLLABUS

1. Introduction to Machine Learning: tasks, algorithms & models

2. Methods for training classification and regression models:

- Nearest Neighbour (KNN)
- Decision / regression trees & rules

3. Methodology and the Machine Learning pipeline

4. Methods for attributes (feature selection, transformation)

5. Methods based on ensembles of models: Bagging, boosting, stacking

6. Large Scale Machine Learning. Big Data

- Map-reduce & Spark

7. Software tools:

- Python + scikit-learn & Pyspark + MLIB

# SYLLABUS:

## 7. SOFWARE TOOLS

**FOR MACHINE LEARNING BASICS ANACONDA (python + scikit-learn)**



**PYSPARK AND MLIB**



**IPYTHON NOTEBOOKS**

# GRADING

- A=30% FINAL EXAM
- B = 70% ASSIGNMENTS
  - A1: Programming
  - A2: Scikit-learn
  - A3: Pyspark / MLLIB

- Pass if A+B>=50%

- What can be done?

Classification



Regression



Market basket analysis



Clustering



Reinforcement learning

# MODELS

- What models can be obtained?

Linear

Non linear

# MODELS

- What models can be obtained?

Functions: y= $3*x^3+2$

Decision trees

Bayesian networks



Rules

If humidity = normal and windy = false then play = yes

And many more: neural networks, nearest neighbor, …

# Decision trees and regression trees

# Ensembles of classifiers

# ATTRIBUTE SELECTIÓN AND TRANSFORMATION



Attribute selection



Principal Component Analysis and Random Projections

# BIG DATA / MAP-REDUCE, SPARK (MLLIB)

# TASKS / MODELS / ALGORITHMS

- <u>What can be done? Tasks:</u>
  - Supervised ML: classification, regression, …
  - Unsupervised ML:  clustering, association, …
  - Semi-supervised ML
  - Reinforcement learning
- <u>What kind of models can be learned?</u>
  - Attribute-value:
    - Trees
    - Nearest neighbor
    - Functions: neural networks, support vector machines, …
    - Bayesian networks
    - Ensembles (bagging, boosting, stacking, …)
  - Relational
- <u>How can models be learned? Algorithms:</u>
  - Linear models: linear regression, simple perceptron, naive bayes, SVM with linear kernel, …
  - Neural networks: backpropagation, rprop, …
  - Decision trees:  ID3, C4.5, C5.0, …
  - Nearest neighbour: IB1,  …

**Training data (labeled pictures of sky objects: galaxies, stars, nebulae, …)**

Pictures in the catalog have been labeled by a human expert (astronomer)

ML Algorithm

**Model**

- Trees
- Nearest neighbor
- Functions: neural networks, support vector machines, …
- Bayesian networks
- Ensembles (bagging, boosting, stacking, …)

Spiral galaxy

# TASKS

- **Inductive learning** (from instances)
  - **Supervised learning:**
    - **Classification:**
    - Regression
  - Semi-supervised learning
  - Unsupervised learning:
    - Clustering
    - Association
  - Reinforcement learning

# CLASSIFICATION TASK. AN EXAMPLE:

- Bank credit approval:
  - An Internet bank owns a large data base with information about clients whose credits were approved or rejected
  - The banks requires a model to determine if a new customer will repay the loan or not
  - Instances (client records in the database):
    - Input attributes : credit  time-length (years), amount, overdue accounts?, own house?
    - Class: yes/no
  - Rule-based model:
    - **IF** (overdue accounts > 0) **THEN** repay loan = no
    - **IF** (overdue accounts = 0) **AND** ((salary  > 2500) **OR** (years > 10)) **THEN** repay loan = yes

# CLASSIFICATION TASK. AN EXAMPLE:

**test set**

| Years | Amount | Salary | Own house? | Overdue accounts? | Repay loan |
|-------|--------|--------|------------|-------------------|------------|
| 10 | 50000 | 3000 | Yes | 0 | ?? |

**T = training set (instances)**

| Years | Amount | Salary | Own house? | Overdue accounts? | Repay loan |
|-------|--------|--------|------------|-------------------|------------|
| 15 | 60000 | 1900 | Yes | 2 | No |
| 2 | 30000 | 3500 | Yes | 0 | Yes |
| 9 | 9000 | 1700 | Yes | 1 | No |
| 15 | 18000 | 3000 | No | 0 | Yes |
| 10 | 24000 | 2100 | No | 0 | No |
| … | … | … | … | … | … |

x (or input attributes)

y (class, or output attribute)

Algorithm

**Model**
**IF** OA >0 **THEN** NO

**IF** OA==0 **AND** S>2500 **THEN** Yes

Repay loan = yes

# IMPORTANT: MODELS

- In the previous slide, the model built from training data is a set of rules:

**IF** OA >0 **THEN** NO **ELSEIF** OA==0 **AND** S>2500 **THEN** Yes

- But there are many more that can be learned:

Functions: y= 3*x³+2

Decision trees

Bayesian networks



And many more: neural networks, nearest neighbor, support vector machines (SVMs).

# TASKS

- **Inductive learning** (from instances)
  - **Supervised learning:**
    - Classification
    - **Regression**
  - Semi-supervised learning
  - Unsupervised learning:
    - Clustering
    - Association
  - Reinforcement learning

# REGRESSION

- If the class is continuous, it is a **regression** problem
- Models are typically mathematical functions *y=g(x)*
  - Linear: y = ax+b
  - Non linear: y = a*$x^2$+bx+c / y = log(sin(x))

# REGRESSION EXAMPLE

- **A wind power forecasting problem: predicting hourly power generation at 7 wind farms**





Wind
(u, v)
ws
wd

Some input variables:
- ws: wind speed
- wd: wind direction
- (u,v): wind direction vector

Model to estimate electricity production from ws, wd, u, v?
wp = f(ws, wd, u, v, …)

# REGRESSION EXAMPLE

DATA

| | date | hors | u | v | ws | wd | dateB | wp1 |
|---|---|---|---|---|---|---|---|---|
| 1 | 2009-07-01 01:00:00 | 1 | 2.34 | -0.79 | 2.47 | 108.68 | 2009-07-01 | 0.085 |
| 2 | 2009-07-01 02:00:00 | 2 | 2.18 | -0.99 | 2.40 | 114.31 | 2009-07-01 | 0.020 |
| 3 | 2009-07-01 03:00:00 | 3 | 2.20 | -1.21 | 2.51 | 118.71 | 2009-07-01 | 0.060 |
| 4 | 2009-07-01 04:00:00 | 4 | 2.35 | -1.40 | 2.73 | 120.86 | 2009-07-01 | 0.045 |
| 5 | 2009-07-01 05:00:00 | 5 | 2.53 | -1.47 | 2.93 | 120.13 | 2009-07-01 | 0.035 |
| 6 | 2009-07-01 06:00:00 | 6 | 2.66 | -1.29 | 2.96 | 115.79 | 2009-07-01 | 0.005 |

Some input variables:
- ws: wind speed
- wd: wind direction
- (u,v): wind direction vector

# REGRESSION EXAMPLE

DATA

| | date hors | u | v | ws | wd | dateB | wp1 |
|---|---|---|---|---|---|---|---|
| 1 | 2009-07-01 01:00:00 | 1 2.34 | -0.79 | 2.47 | 108.68 | 2009-07-01 | 0.085 |
| 2 | 2009-07-01 02:00:00 | 2 2.18 | -0.99 | 2.40 | 114.31 | 2009-07-01 | 0.020 |
| 3 | 2009-07-01 03:00:00 | 3 2.20 | -1.21 | 2.51 | 118.71 | 2009-07-01 | 0.060 |
| 4 | 2009-07-01 04:00:00 | 4 2.35 | -1.40 | 2.73 | 120.86 | 2009-07-01 | 0.045 |
| 5 | 2009-07-01 05:00:00 | 5 2.53 | -1.47 | 2.93 | 120.13 | 2009-07-01 | 0.035 |
| 6 | 2009-07-01 06:00:00 | 6 2.66 | -1.29 | 2.96 | 115.79 | 2009-07-01 | 0.005 |

Linear model:

$wp = f(ws, wd, u, v)$
$wp = a_1*ws + a_2*wd + a_3*u + a_4*v + b$

Obviously, a non-linear model could do better

# TASKS

- **Inductive learning**(from instances)
  - Supervised learning:
    - Classification
    - Regression
  - **Semi-supervised learning**
  - Unsupervised learning:
    - Clustering
    - Association
  - Reinforcement learning

# SEMISUPERVISED LEARNING

- When both labelled and unlabelled instances are available
- Why: labelling instances may be costly (ex: to perform a biopsy to determine if a person has cancer)

# TASKS

- **Inductive learning** (from instances)
  - Supervised learning:
    - Classification
    - Regression
  - Semi-supervised learning
  - **Unsupervised learning:**
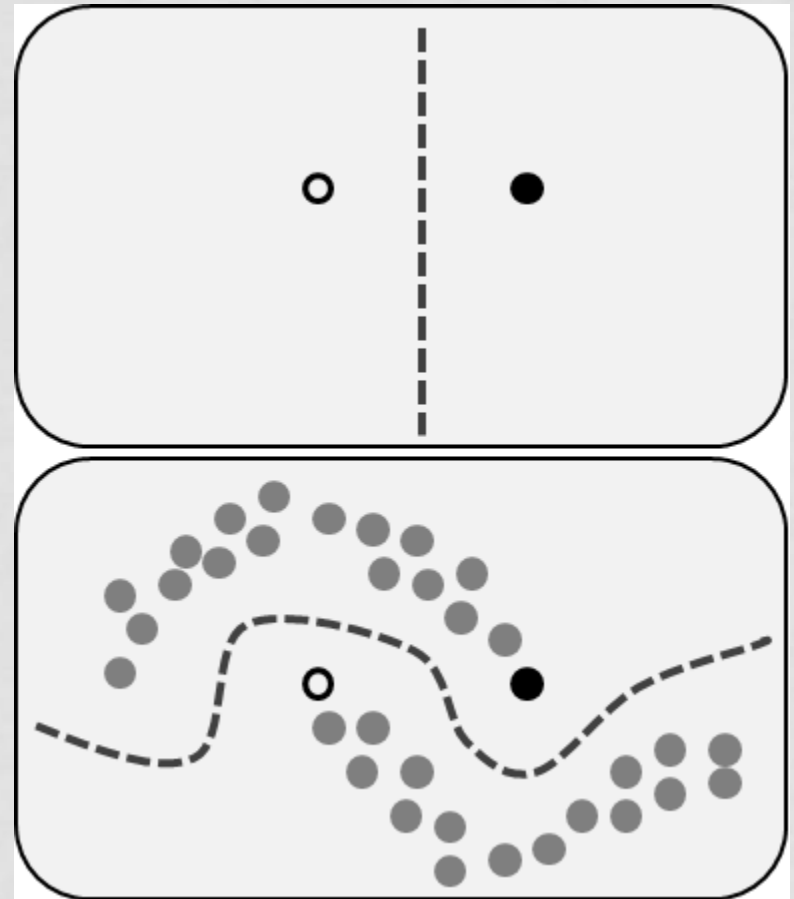    - **Clustering**
    - Association
  - Reinforcement learning

# UNSUPERVISED LEARNING (NO LABELS): CLUSTERING

- To determine natural clusterings in instance space, based on the input attributes (no labels)

X2: Sentence Average length

X1: Word average length

Example:each point is a different book. 2 groups:

* Long words and sentences (philosophy?)

* Short words and sentences (best-sellers?)

# CLUSTER REPRESENTATION

- Most commonly: centroids (ex: k-means algorithm)

K-MEANS: http://www.youtube.com/watch?v=74rv4snLl70

# CLUSTERING

- Clustering is not so well defined as classification: clustering based on neighbourhood or connectivity?

# CLUSTERING EXAMPLE

- Human resources department would like to cluster employees in order to understand the different types of employee and treat them accordingly (fire problematic workers? ☺ ).

# CLUSTERING EXAMPLE. TRAINING DATA

| Id | Salary | Married | Car | Offsp ring | Own-house | Syndicate | Sick leave | Years working | Sex |
|----|--------|---------|-----|-----------|-----------|-----------|-----------|---------------|-----|
| 1 | 1000 | Yes | No | 0 | No | No | 7 | 15 | M |
| 2 | 2000 | No | Yes | 1 | No | Yes | 3 | 3 | F |
| 3 | 1500 | Yes | Yes | 2 | Yes | Yes | 5 | 10 | M |
| 4 | 3000 | Yes | Yes | 1 | No | No | 15 | 7 | F |
| 5 | 1000 | Yes | Yes | 0 | Yes | Yes | 1 | 6 | F |
| ... | ... | ... | ... | ... | | ... | ... | ... | ... |

# MODEL (CLUSTERS)

| | GROUP 1 | GROUP 2 | GROUP 3 |
|---|---|---|---|
| Salary | 1535 | 1428 | 1233 |
| Married (No/Yes) | 77%/22% | 98%/2% | 0%/100% |
| Car | 82%/18% | 1%/99% | 5%/95% |
| Offspring | 0.05 | 0.3 | 2.3 |
| Own-house | 99%/1% | 75%/25% | 17%/83% |
| Syndicated | 80%/20% | 0%/100% | 67%/33% |
| Sick leave | 8.3 | 2.3 | 5.1 |
| Years working | 8.7 | 8 | 8.1 |
| Sex (M/W) | 61%/39% | 25%/75% | 83%/17% |

# MODEL (CLUSTERS)

- Cluster 1: No offspring and rented house. Low level of syndication. Lots of sick leaves
- Cluster 2: No offspring and own-car. High syndication level. Few sick leaves. Tipically women living in rented houses
- Cluster 3: Married men with children and own-car and own-houses. Low syndication level

# TASKS

- **Inductive learning**(from instances)
  - Supervised learning:
    - Classification
    - Regression
  - Semi-supervised learning
  - **Unsupervised learning:**
    - Clustering
    - **Association**
  - Reinforcement learning

# MARKET BASKET ANALYSIS (**ASSOCIATION**)

- A supermarket needs to know customer behavior.
  - Ex: if customer buys X then s/he also buys Y
- Service might be improved (putting together products bought together, etc.)

# TRAINING DATA (CUSTOMER BASKETS)

| Id | Eggs | Oil | Napies | Wine | Milk | Butter | Salmon | Lettuce | ... |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Yes | No | No | Yes | No | Yes | Yes | Yes | ... |
| 2 | No | Yes | No | No | Yes | No | No | Yes | ... |
| 3 | No | No | Yes | No | Yes | No | No | No | ... |
| 4 | No | Yes | Yes | No | Yes | No | No | No | ... |
| 5 | Yes | Yes | No | No | No | Yes | No | Yes | ... |
| 6 | Yes | No | No | Yes | Yes | Yes | Yes | No | ... |
| 7 | No | No | No | No | No | No | No | No | ... |
| 8 | Yes | Yes | Yes | Yes | Yes | Yes | Yes | No | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

# MODEL

- Rules **IF** $At_1=a$ **AND** $At_2=b$ y ... **THEN** $At_n=c$
  - **IF** nappies=Yes **THEN** milk=Yes
  - **IF** butter = Yes **AND** salmon = Yes **THEN** wine = Yes

- Also: **IF** $At_1=a$ AND $At_2=b$ **THEN** $At_n=c$, **$At_4=D$**

Service might be improved (putting together nappies and milk, etc.)

# ASSOCIATION

# TASKS

- **Inductive learning** (from instances)
  - Supervised learning:
    - Classification
    - Regression
  - Semi-supervised learning
  - Unsupervised learning:
    - Clustering
    - Association
  - **Reinforcement learning**

# TASK: REINFORCEMENT LEARNING

- The goal of learning is a "policy" π so that the agent (mouse) knows what to do at each situation (in the case of the mouse, a situation is a particular location within the maze). Robotics.

- Actions:
  - forward
  - turn left
  - turn right

# TASKS

- **Inductive learning**(from instances)
  - Attribute-value models
    - Supervised learning:
    - Semi-supervised learning
    - Unsupervised learning:
    - Reinforcement learning
  - **Relational learning**

# Relational Learning

- For instance, learn the concept of "being a daughter"
- **IF** X is female **AND** Y is the mother of X **THEN** X is a daugther of Y
- Compare this rule with:

**IF** Overdue Accounts ==0 **AND** Salary >2500 **THEN** Repay loan = Yes

- Relational rules use variables (X, Y) and relations

# Relational Learning: ILP (inductive logic programming

| Training examples | Background knowledge | |
|---|---|---|
| daughter(mary, ann). ⊕ | parent(ann, mary). | female(ann). |
| daughter(eve, tom). ⊕ | parent(ann, tom). | female(mary). |
| daughter(tom, ann). ⊖ | parent(tom, eve). | female(eve). |
| daughter(eve, ann). ⊖ | parent(tom, ian). | |

Learned Knowlege:

$$daughter(X, Y) \leftarrow female(X), mother(Y, X).$$
$$daughter(X, Y) \leftarrow female(X), father(Y, X).$$

# Bibliography

**SCIKIT-LEARN**

- Learning scikit-learn: Machine Learning in Python:
    - http://m.proquest.safaribooksonline.com/hd/catalog?isbn=9781783281930
- Mastering Machine Learning with scikit-learn:
    - http://m.proquest.safaribooksonline.com/hd/catalog?isbn=9781783988365
- scikit-learn Cookbook
    - http://m.proquest.safaribooksonline.com/hd/catalog?isbn=9781783989485

**SPARK**

- Holden Karau, Andy Konwinski, Patrick Wendell, Matei Zaharia, Learning Spark, O'Reilly Media, 2015. ISBN: 978-1-449-35862-4
- Sandy Ryza, Uri Laserson, Sean Owen, Josh Wills, Advanced Analytics with Spark, O'Reilly, 2015. ISBN: 978-1-491-91276-8