

Machine Learning

João Gama

jgama@fep.up.pt



LIAAD-INESC Porto, University of Porto, Portugal

December 2018

1 Introduction

2 Applications

3 Introduction to Data Mining

4 Overview of Data Mining Problems

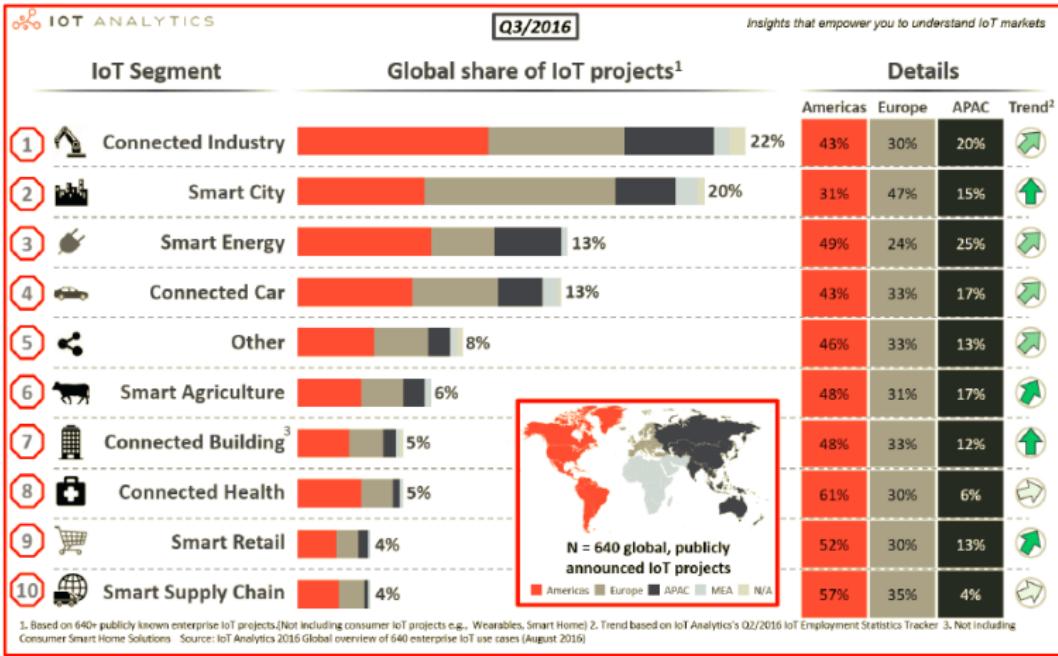
5 Resources

Industry 4.0

We have machines that collect, process, and send information to other machines



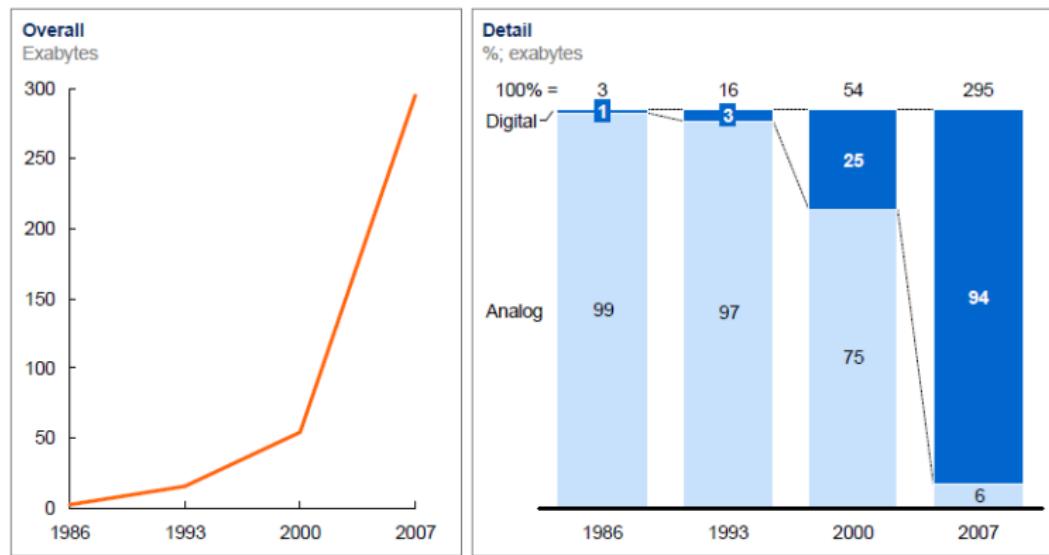
Internet of Things



The Big Bang of digital data ...

Data storage has grown significantly, shifting markedly from analog to digital after 2000

Global installed, optimally compressed, storage

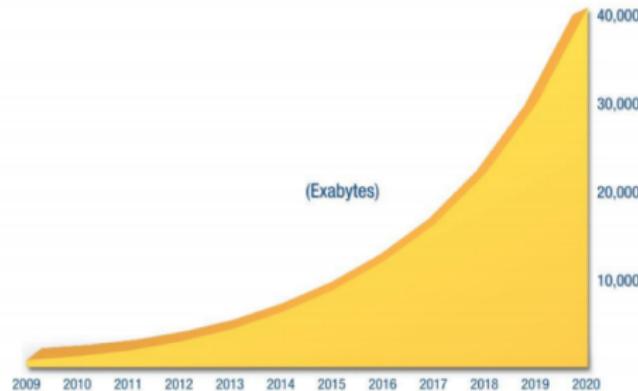


NOTE: Numbers may not sum due to rounding.

SOURCE: Hilbert and López, "The world's technological capacity to store, communicate, and compute information," Science, 2011

The Growth of Digital Data...

The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020



Source: IDC's Digital Universe Study, sponsored by EMC, December 2012

Memory unit	Size	Binary size
kilobyte (kB/KB)	10^3	2^{10}
megabyte (MB)	10^6	2^{20}
gigabyte (GB)	10^9	2^{30}
terabyte (TB)	10^{12}	2^{40}
petabyte (PB)	10^{15}	2^{50}
exabyte (EB)	10^{18}	2^{60}
zettabyte (ZB)	10^{21}	2^{70}
yottabyte (YB)	10^{24}	2^{80}

Tools: time ago ...

Tools seemed quite powerful



Tools



Problems

Tools: nowadays ...

Last few years



Main Goal: Understanding Data

A brief history of big data, the Noam Chomsky way



Text Size - +

Published: Saturday, 23 Nov 2013 | 7:00 AM ET

By: Eric Rosenbaum | CNBC.com

[Facebook](#) Recommend 47

[Twitter](#) 104

[Google+](#) 1 6

[LinkedIn](#) 15

[Share](#)



ChinaFotoPress | Getty Images

Noam Chomsky

The latest news from the fast-evolving world of the [Data Economy](#).

For those familiar with Noam Chomsky, the pioneering linguist whose theory of recursion seeks to find the universal in all human languages, you probably also know that Chomsky often has not-so-nice things to say about the U.S. government, and has also made a career of finding the universal

Big data is a step forward, but our problems are not lack of access to data, but understanding them. Big data is very useful if I want to find out something without going to the library, but I have to understand it, and that's the problem.

Motivation



• Consultant
• Speaker
• Trends
• Scenarios
• Planning

Global Future
Strategies for a Global Age

Home	What's New	Global Future Reports™	Book Reviews	Bibliographies	Contact Us
------	------------	------------------------	--------------	----------------	------------

Global Future Report™

January 15th, 2001

10 Emerging Technologies That Will Change the World

© Dr. Terry J. van der Werff, CMC

MIT's *Technology Review* has identified 10 emerging areas of technology that will soon have a profound impact on the economy and how we live and work.

Regular readers of *Global Future Report™* know I am a sucker for lists of things that matter. I even write lists of my own, e.g. my "[Ten Tips for Harnessing the Future](#)" or the four forces converging to alter global telecommunications in "[Calling the Future](#)."

To launch the New Millennium the January/February issue of *Technology Review*, MIT's magazine of innovation, focuses on "The Technology Review Ten" - "10 emerging areas of technology that will soon have a profound impact on the economy and how we live and work." For each, one innovator's work is highlighted.

Drum roll, please! The ten emerging technologies that will change the world are:

- **Brain-Machine Interfaces** - In essence, researchers try both to understand how the brain works and to use this knowledge to implant electrodes in specific parts of the brain to permit control of computers, robotic arms, or other artificial devices designed to restore lost sensory and motor functions.
- **Flexible Transistors** - Silicon does not bend readily, so a new class of hybrid materials are being developed that marry the speed of inorganic compounds with the flexibility of organic polymers. They have the advantage of being able to be dissolved and printed onto paper or plastic as if they were ink particles.
- **Data Mining** - Ever get an e-mail from amazon.com suggesting a book that relates to an earlier one you ordered from them? You have been the subject of data mining, which is nothing more than the extraction of meaningful information and patterns from huge data sets.
- **Digital Rights Management** - Think Napster! The Internet permits the sharing of digital content far and wide at little cost. But originators of the content - articles, data, graphics, songs - may lose control of their intellectual property. Digital rights management combines encryption with payment software to

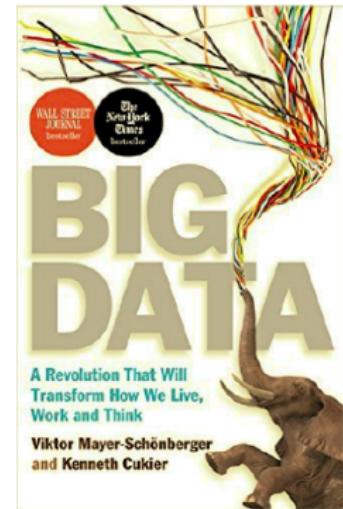
Motivation

McKinsey Global Institute



June 2011

Big data: The next frontier
for innovation, competition,
and productivity



Definitions

Extração de Conhecimento de Dados;
Gama, Carvalho, Faceli, Lorena, Oliveira; 3 edition,
Silabo, 2017



Definitions

- *Self-constructing or self-modifying representations of what is being experienced for possible future use* Michalski, 1990
- *Analysis of observational data to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful for the data owner* Hand, Mannila, Smyth, 2001
- *Obter representações em compreensão a partir de representações em extensão.*

Scientific Areas

Scientific Areas

- Statistics
 - Inference
- Computer Science
 - Artificial Intelligence
 - Machine Learning
- Data Bases
 - Multidimensional Data Bases

Outline

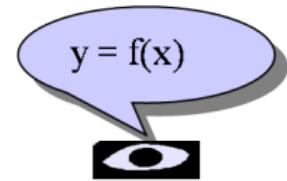
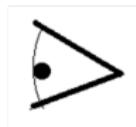
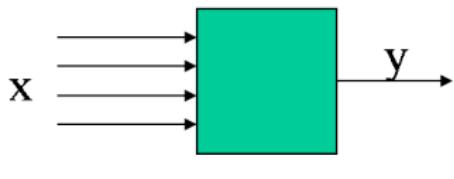
1 Introduction

2 Applications

3 Introduction to Data Mining

4 Overview of Data Mining Problems

5 Resources



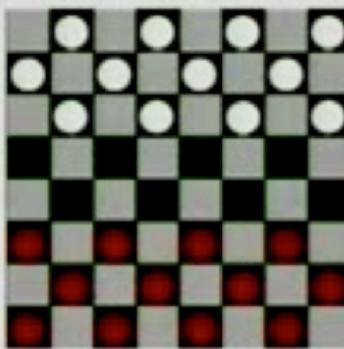
Self-Programming machines



Lecture 1 - The Motivation & Applications of Machine Learning (1)

Machine Learning definition

- Arthur Samuel (1959). Machine Learning:
Field of study that gives computers the ability
to learn without being explicitly programmed.



STANFORD
UNIVERSITY

00:34:19

Bank Check

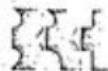
secretariat Prive President J. B. ARISTIDE
siecle National
Port-au-Prince, Haiti
N° 226-2012

PAU-Pac 13/07/02

94
PP 10-1100
0000 4
USD 600,000
DOLLARS

ayez ce chèque
l'ordre de

Banque Populaire Haïtienne
six cent mille et 000



Banque de la République d' Haiti
Au Capital Social de 50.000.000 GOES
Port-au-Prince, Tel: 299-1000, Fax: 1145

EMO:

N° 94 101100000000

1110110301

000600000000

Postal Codes



#123456789

Postal Codes

x	3	8	6	9	6	4
y	‘3’	‘8’	‘6’	‘9’	‘6’	‘4’
	↓	↓	↓	↓	↓	↓

Postal Codes

0000000000000000
1111111111111111
2222222222222222
3333333333333333
4444444444444444
5555555555555555
6666666666666666
7777777777777777
8888888888888888
9999999999999999

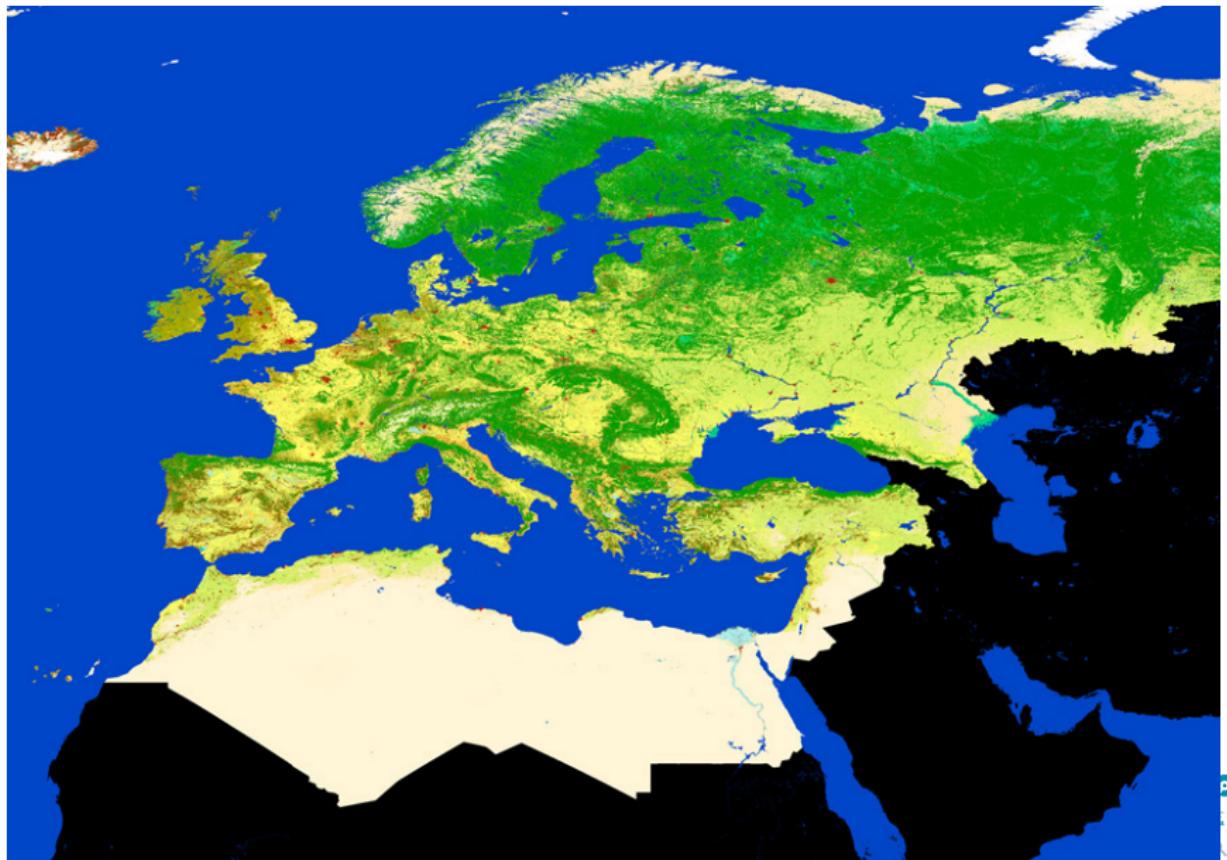
Satélite traça mapa ambiental do continente europeu em tempo recorde - Publico 25.11.2010

Segundo a Agência Espacial Europeia (ESA), o mapa foi realizado em nove meses, a partir dos dados recolhidos pelo satélite Envisat, graças ao seu espectógrafo MERIS (Medium Resolution Imaging Spectrometer). Este mapa, resultado do projecto GlobCorine, é *o primeiro do género a ser produzido em tão pouco tempo, nove meses em comparação com vários anos*, salienta.

O objectivo foi desenvolver um serviço automático de produção de cartografia de ocupação e utilização do solo à escala europeia, incluindo os países da bacia mediterrânica e a Rússia. Agora, o mapa relativo a 2009, deverá ser actualizado todos os anos, algo considerado essencial para as agências de Ambiente dos vários países. *Satisfazer as necessidades ambientais da Europa exige informação sobre o uso do solo consistente e actualizada*, explica a ESA.

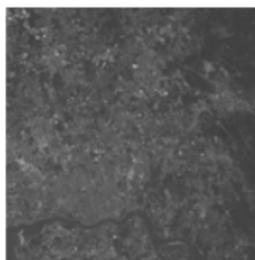
Finalmente podemos ter informação relevante e actual que suporte os processos de tomada de decisão, comentou Chris Steenmans, responsável pelo projecto na ESA.

Predicting Land Usage

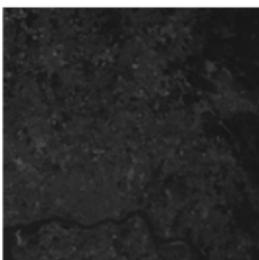


Porto

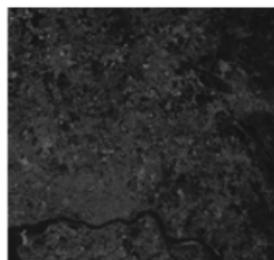
LandSat - Porto



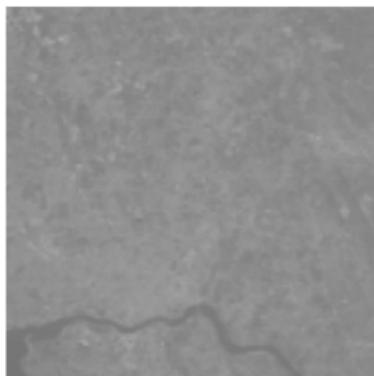
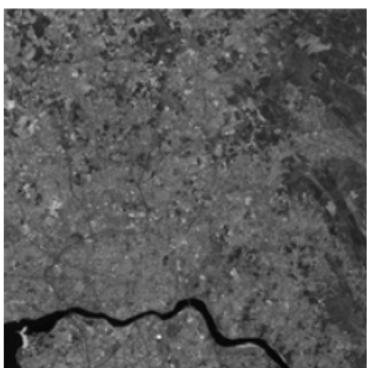
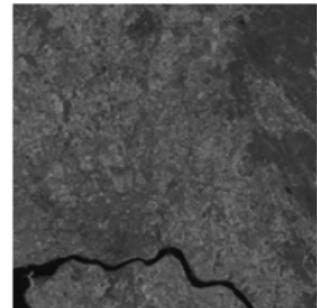
Verde-azul



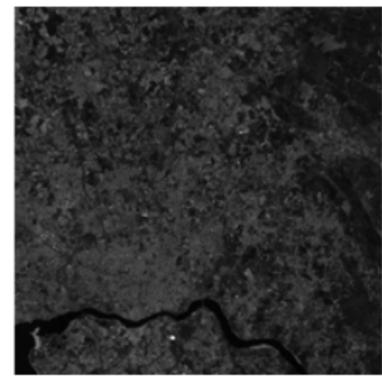
vermelho



vermelho



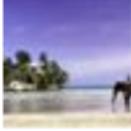
(Infravermelhos)



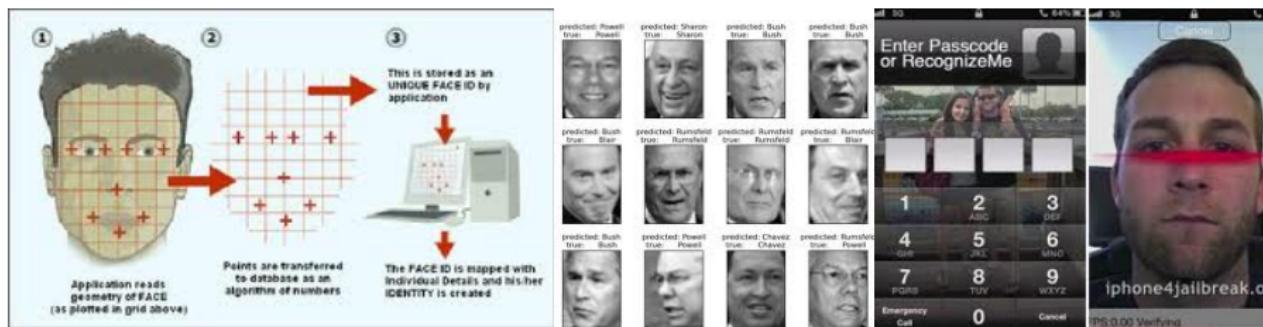
Pixel:30mx30m

FACEBOOK: The Missing Link

Friend Requests [Find Friends - Settings](#)

	Val Verbić	Confirm	Delete Request
	Yedan Cevokop	Confirm	Delete Request
	9 mutual friends	Confirm	Delete Request
		Confirm	Delete Request
	Gystal Edmund 4 mutual friends	Confirm	Delete Request

Face Recognition



Google (2)

The screenshot shows the Google Translate interface. At the top, there is a search bar with the word "Tradutor". Below it, the language selection is set to "De: Português" and "Para: Inglês". The main area contains two text boxes. The left box, labeled "Português", contains the text "Bom dia,". The right box, labeled "Inglês", contains the translated text "Good morning,". There are tabs for "Inglês", "Português", and "Francês" at the top of each box. A blue "Traduzir" button is located between the two boxes. The entire interface is contained within a light gray box.

Google (3)



Google Photos

From Wikipedia, the free encyclopedia

Google Photos is a photo sharing and storage service developed by Google. It was announced in May 2015 and spun out from Google+, the company's social network.

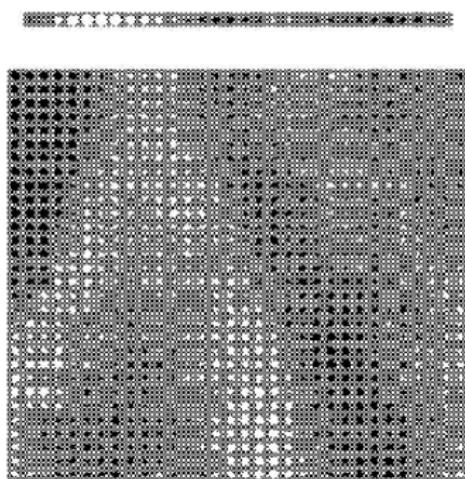
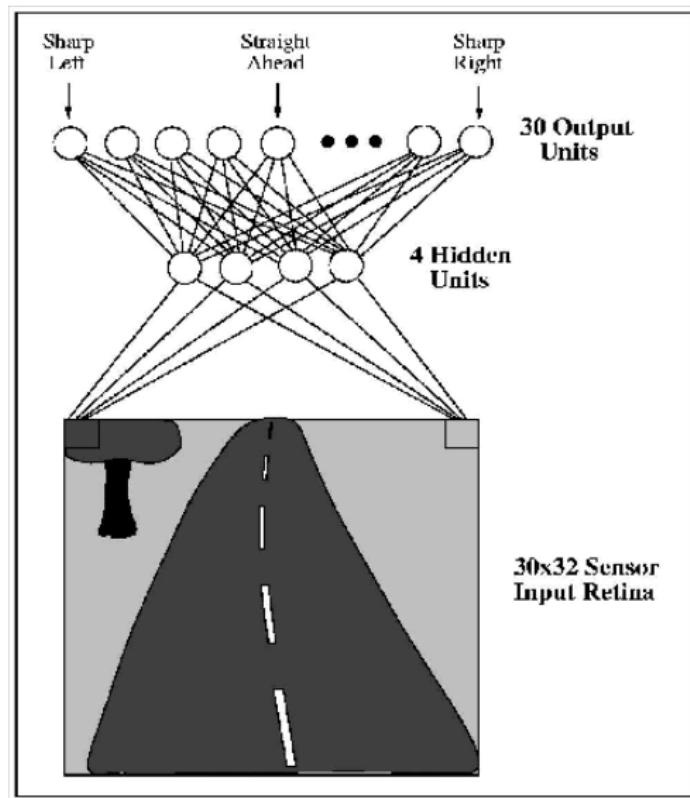
Google Photos gives users free, unlimited storage for photos up to 16 megapixels and videos up to 1080p resolution. The service automatically analyzes photos, identifying various visual features and subjects. Users can search for anything in photos, with the service returning results from three major categories: People, Places, and Things. Google Photos recognizes faces, grouping similar ones together; geographic landmarks (such as the Eiffel Tower); and subject matter, including birthdays, buildings, animals, food, and more. Google implements different forms of machine learning into the Photos service, particularly its recognition of photo contents, as well as enabling features that can automatically generate albums, animate similar photos into quick videos, surface past memories at significant times, and improve the quality of photos and videos. In May 2017, Google announced several updates to Google Photos, including reminders for suggested sharing of photos, shared photo libraries between two users, and physical albums, with Photos automatically suggesting collections based on face, location, trip, or other distinction.

Google Photos received critical acclaim after its decoupling from Google+ in 2015. Reviewers liked the updated Photos service for its recognition technology, search, apps, and loading times. Nevertheless, privacy concerns were raised, including Google's motivation for building the service, as well as its relationship to governments and possible laws requiring Google to hand over a user's entire photo history. Google Photos has seen strong user adoption. It reached 100 million users after five months, 200 million after one year, and 500 million as of May 2017, with Google announcing that over 1.2 billion photos are uploaded to the service every day, with the grand total of all uploaded content measuring over 13.7 petabytes of storage.

Autonomous Vehicles



Autonomous Vehicles



Recommender Web Sites

Amazon.com: Recommended For You - Microsoft Internet Explorer

Ficheiro Editar Ver Favoritos Ferramentas Ajuda

Retroceder Avançar Recurso Procurar Favoritos Multimédia E-mail Imprimir Imprimir para PDF Imprimir para Word Imprimir para Excel Imprimir para Word 2007 Imprimir para Excel 2007 Imprimir para Word 2010 Imprimir para Excel 2010 Imprimir para Word 2013 Imprimir para Excel 2013 Imprimir para Word 2016 Imprimir para Excel 2016

Enderijo: https://www.amazon.com/exec/obidos/tg/stores/recs/instant-recs/-/recs/104-7366371-5740701 Ir para Hiperligações

Google amazon Search Web 9 blocked AutoFill Options amazon

amazon.com | VIEW CART | WISH LIST | YOUR ACCOUNT | HELP

Shop in Pampers Health & Personal Care (Beta-What is this?)

WELCOME JOÃO'S STORE BOOKS APPAREL & ACCESSORIES ELECTRONICS TOYS & GAMES MUSIC COMPUTER & VIDEO GAMES SEE MORE STORES

RECOMMENDATIONS WIZARD IMPROVE YOUR RECOMMENDATIONS FRIENDS & FAVORITES LEARN MORE

João's Gold Box

Recommended for João Gama (If you're not João Gama, [click here.](#))

BROWSE RECOMMENDED

Recommendations

All Stores

- Baby
- Books
- DVD
- Electronics
- Outdoor Living
- Tools & Hardware
- Kitchen & Housewares
- Magazine Subscriptions
- Music
- Computers
- Camera & Photo
- Software
- Toys & Games
- Video
- Computer & Video Games

Your recommendations are based on [1 items you own](#) and more. [More results](#)

view: All | [New Releases](#) | [Coming Soon](#) | [Bargains](#)

1. **Machine Learning**
by Tom M. Mitchell
Average Customer Review: ★★★★☆
Publication Date: March 1, 1997
Our Price: \$143.45 Used & new from \$49.00

[Add to cart](#) [Add to Wish List](#)

[See related items](#) [Why was I recommended this?](#)
Rate this item I own it Not interested

2. **Artificial Intelligence: A Modern Approach (2nd Edition)**
by Stuart J. Russell, Peter Norvig
Average Customer Review: ★★★★☆
Publication Date: December 20, 2002
Our Price: \$78.32 Used & new from \$39.00

[Add to cart](#) [Add to Wish List](#)

[See related items](#) [Why was I recommended this?](#)
Rate this item I own it Not interested

3. **Neural Networks for Pattern Recognition**

[Iniciar](#) [CARF](#) Microsoft PowerPoint ... Amazon.com: Recom...

PT 20:49

LIAAD INSCRIÇÃO

Jeopardy!



Diagnosis

- Sensors:

- gsr_low_average
- heat_flux_high_average
- near_body_temp_average
- pedometer
- skin_temp_average
- longitudinal_accelerometer_SAD
- longitudinal_accelerometer_average
- transverse_accelerometer_SAD
- transverse_accelerometer_average



The SenseWear armband, shown in the figure below, is a sleek, wireless and accurate wearable body monitor that enables continuous physiological monitoring outside the laboratory.

Help Systems

ERIC HORVITZ, a researcher at Microsoft and a guru in the field of Bayesian statistics, feels bad about the paperclip, but he hopes his latest creation will make up for it. The paperclip in question, as even casual users of Microsoft's Office software will be aware, is a cheery character who pops up on the screen to offer advice on writing a letter or formatting a spreadsheet. That was the idea, anyway. But many people regard the paperclip as annoyingly over-enthusiastic, since it appears without warning and gets in the way.

Mobile Manager evaluates incoming e-mails on a user's PC and decides which are important enough to forward to a pager, mobile phone or other e-mail address. Its Bayesian innards give it an almost telepathic ability to distinguish junk mail from genuinely important messages.

The *Clip* is an interface for a Bayesian Network:



Outline

1 Introduction

2 Applications

3 Introduction to Data Mining

4 Overview of Data Mining Problems

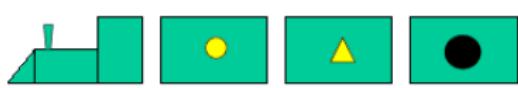
5 Resources

The Michalsky trains: data

□ Problema

- O que caracteriza os comboios que vão para o Ocidente (\leftarrow) relativamente aos que vão para o Oriente (\rightarrow)?

□ Dados



The Michalsky trains: patterns

□ Padrões observados (entre muitos)

- Os que vão para Ocidente têm um vagão com um pequeno triângulo amarelo, os que vão para Oriente não têm.
- Todos os comboios que têm um triângulo amarelo, têm também um círculo



The Iris Problem



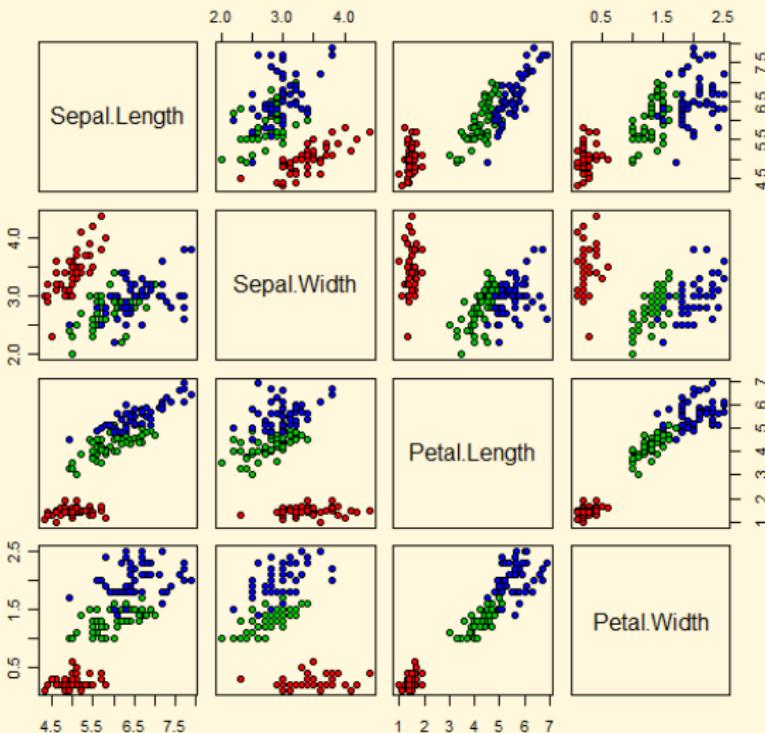
The Iris problem: data

```
> data(iris)
> str(iris)
'data.frame': 150 obs. of 5 variables:
 $ Sepal.Length: num 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
 $ Sepal.Width : num 3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
 $ Petal.Length: num 1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
 $ Petal.Width : num 0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
 $ Species      : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1
> |
```

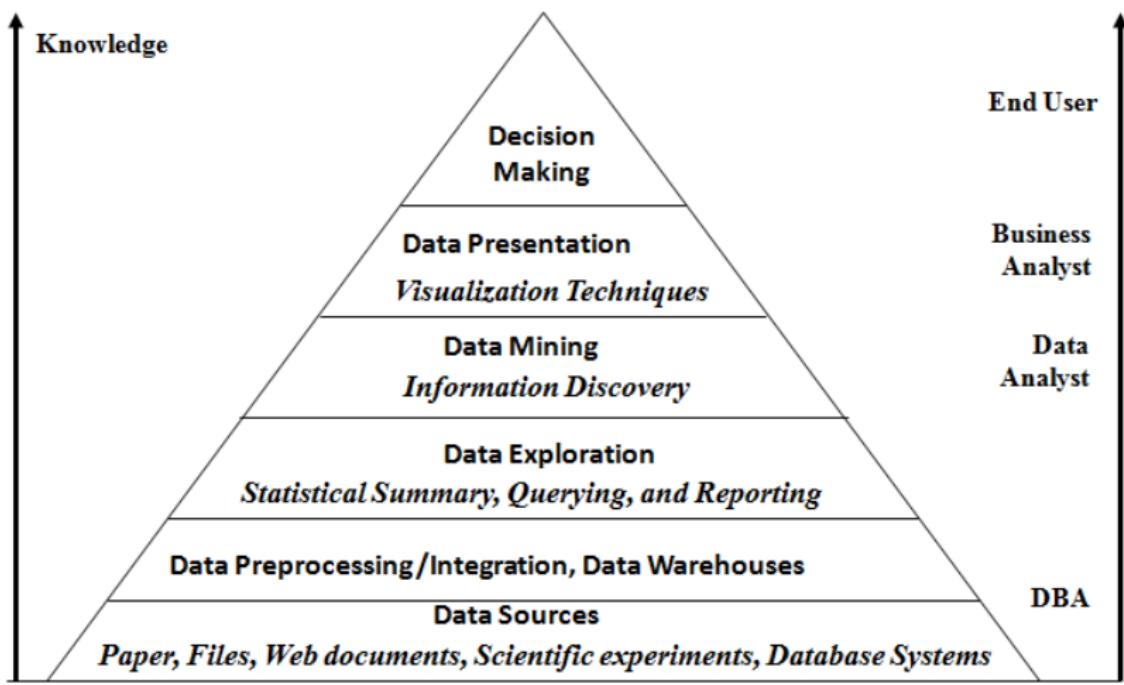
Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
5.8	2.7	5.1	1.9	virginica
7.1	3.0	5.9	2.1	virginica
6.3	2.9	5.6	1.8	virginica

Where are the patterns?

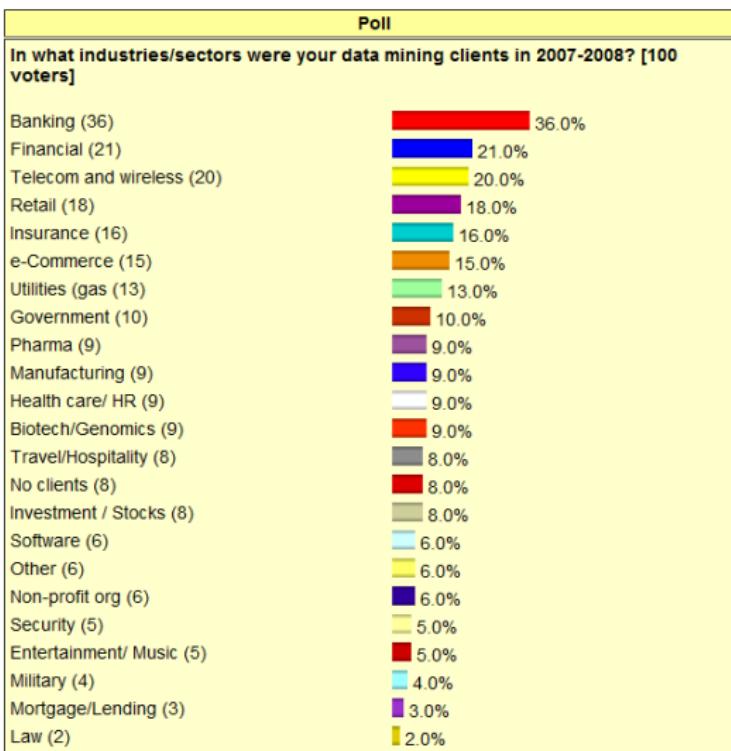
Edgar Anderson's Iris Data



Data and Knowledge

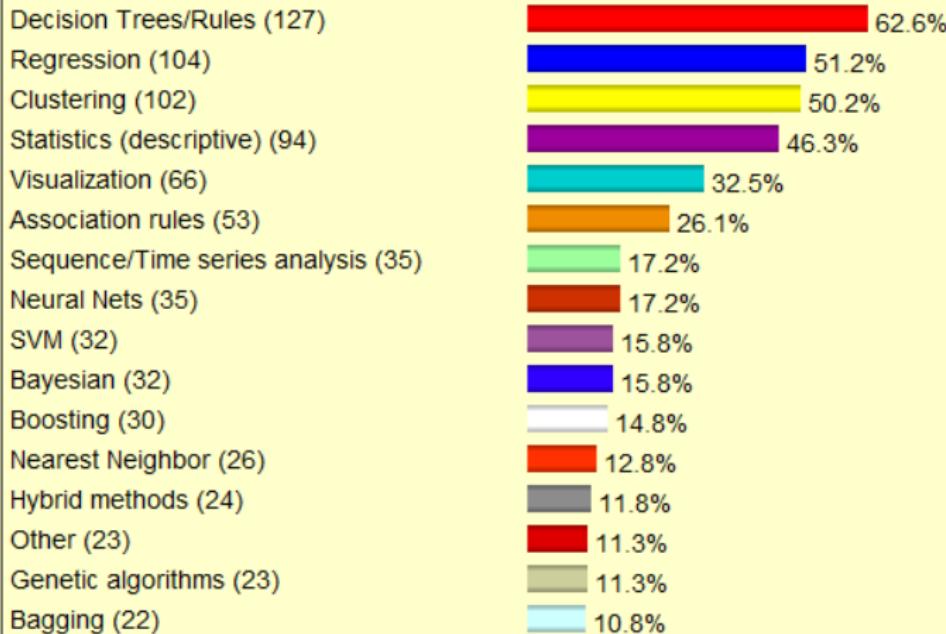


Where is DM used?



The most used techniques

Data mining/analytic methods you used frequently in the past 12 months: [203 voters]

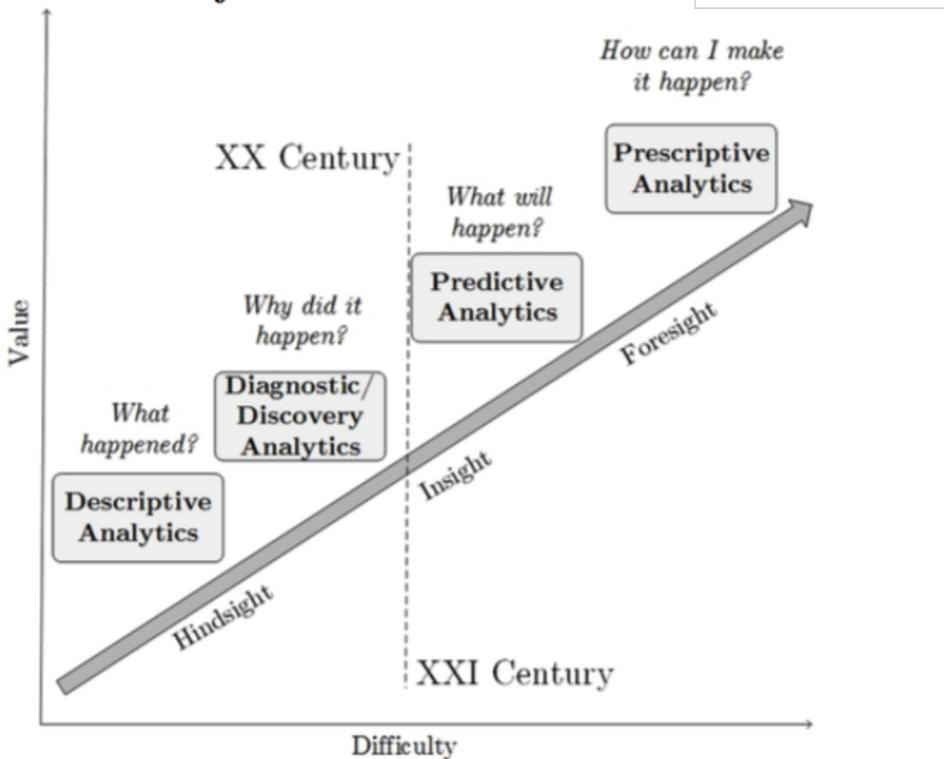


Outline

- 1 Introduction
- 2 Applications
- 3 Introduction to Data Mining
- 4 Overview of Data Mining Problems
- 5 Resources

The Evolution of Computational Data Analysis - II

Gartner Analytics Value Escalator.



Examples of What People are Doing with Data Mining

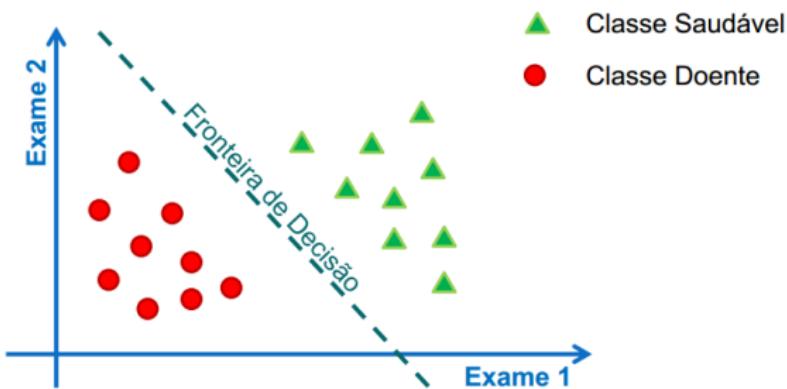
- Fraud/Non-Compliance Anomaly detection
 - Isolate the factors that lead to fraud, waste and abuse
 - Target auditing and investigative efforts more effectively
- Credit/Risk Scoring
- Intrusion detection
- Parts failure prediction
- Recruiting/Attracting customers
- Maximizing profitability (cross selling, identifying profitable customers)
- Service Delivery and Customer Retention
- Build profiles of customers likely to use which services
- Web Mining
- Social Networks analysis

Types of Problems

- Supervised
 - **Classifying** people or things into groups by recognizing patterns
 - **Regression** relationship between a dependent continuous variable and one or more independent variables.
 - **Forecasting** what may happen in the future
- Unsupervised
 - **Clustering** people or things into groups based on their attributes
 - **Associating** what events are likely to occur together
 - **Sequencing** what events are likely to lead to later events
 - **Link Analysis** who is connected with

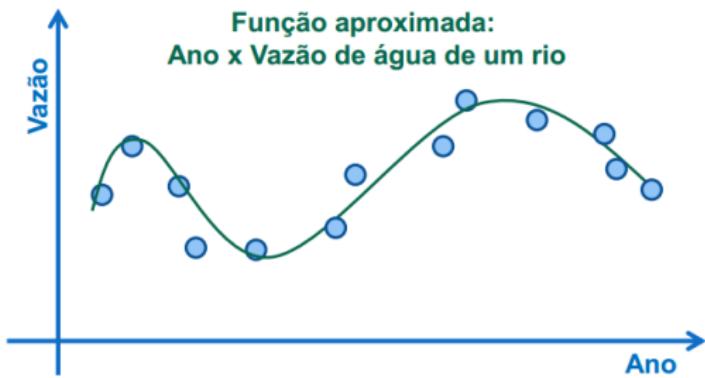
Classification

Classifying people or things into groups by recognizing patterns



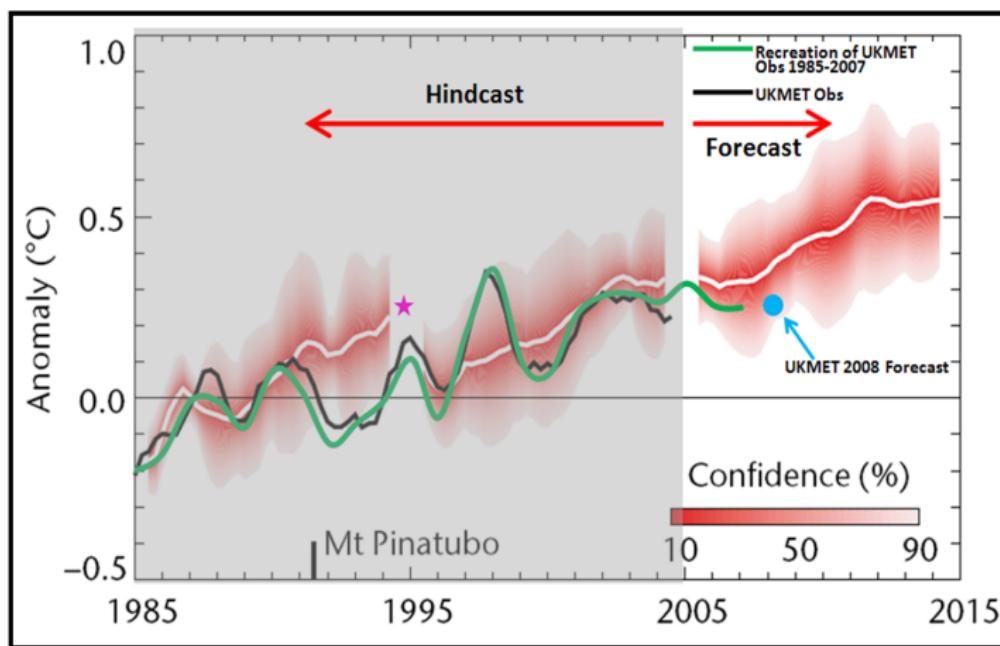
Regression: Function Approximation

relationship between a dependent continuous variable and one or more independent variables.



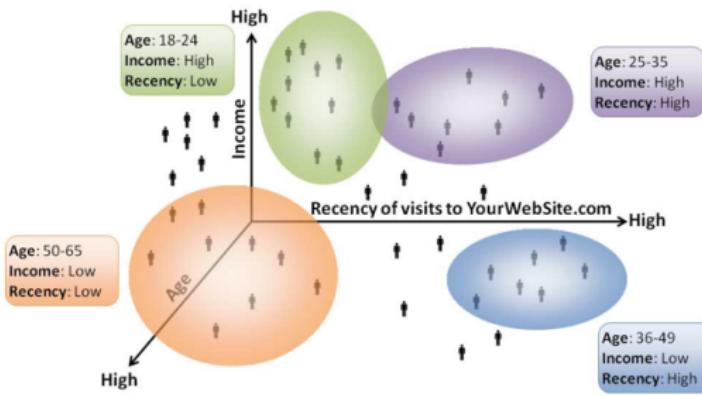
Forecast

Predicting the Future:



Cluster Analysis

Clustering people or things into groups based on their attributes



Association Analysis

What does it goes with?



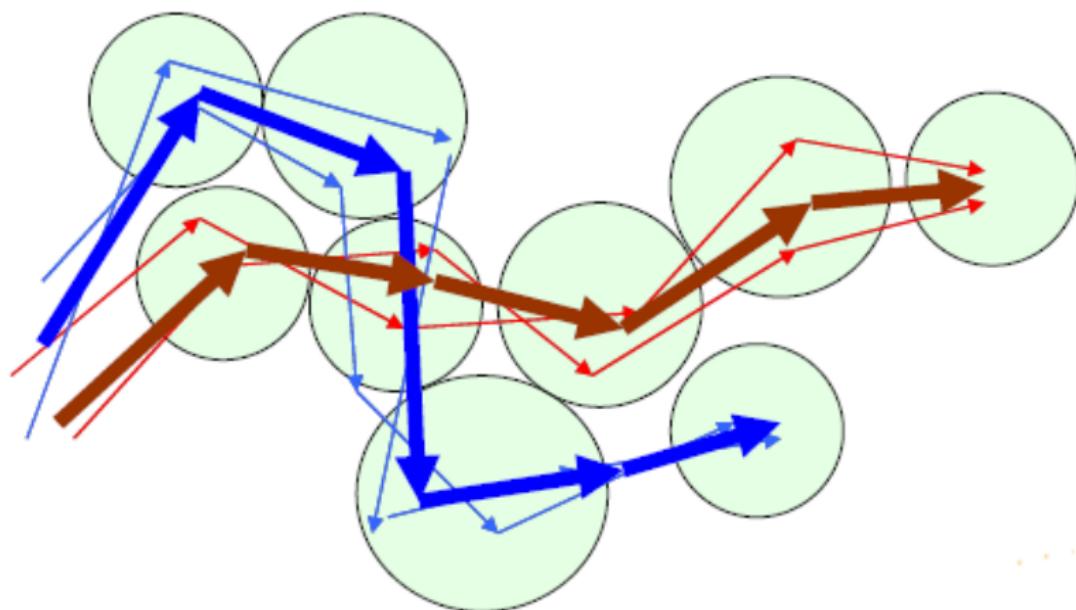
[3, 75%]



[3, 100%]

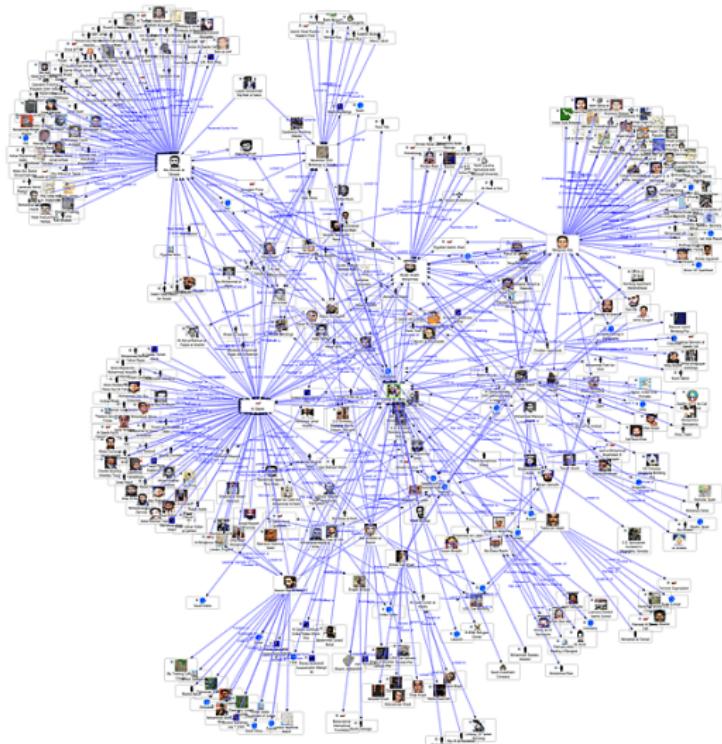
Sequence Analysis

Record trajectories from GPS devices. Which ones are most frequent?



Social Network Analysis

People communicate with people. Who are the pivots?



Outline

1 Introduction

2 Applications

3 Introduction to Data Mining

4 Overview of Data Mining Problems

5 Resources

Recursos

- Video Lectures:
http://videolectures.net/Top/Computer_Science/Machine_Learning
- KDD nuggets: <http://www.kdnuggets.com>
- Data Sets at UCI: <http://archive.ics.uci.edu/ml/>
- [http://www.sciencemag.org/feature/
data/compsci/machine_learning.dtl](http://www.sciencemag.org/feature/data/compsci/machine_learning.dtl)
- <http://www.microstrategy.pt/data-mining/>
- <http://www.kdubiq.org/>
- <http://www.acm.org/sigs/sigkdd/explorations/>
- [http://home.earthlink.net/~dwaha/research/machine-
learning.html](http://home.earthlink.net/~dwaha/research/machine-learning.html)