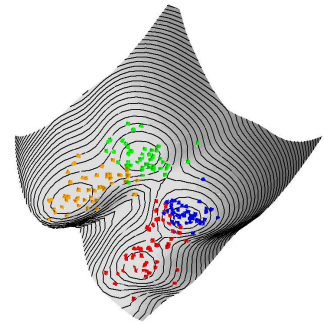# IIC2433 Minería de Datos

*Profesor: Karim Pichara*

*Departamento de Ciencia de la Computación*
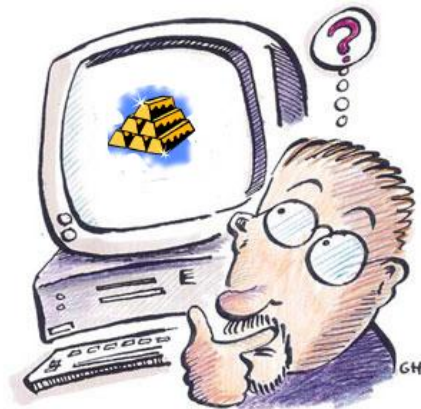
# Who is the professor?

# Data Mining

- Research area that belongs to Computer Science and Statistics (and may be Physics, and Math)

- It studies how to develop automatic tools to analyze data

- Solutions have to be:
  - fast enough (big data)
  - accurate enough (flexible & adaptive models)

Karim Pichara B.

Mining of gold from rocks ⟹ Gold Mining



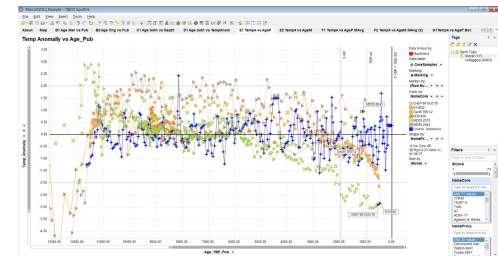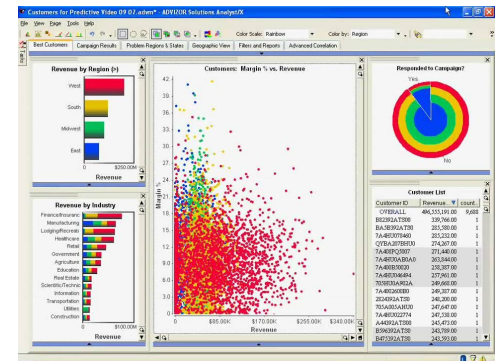Mining of Data from Information sources ⟹ Data Mining

# What name to use?

- Data Science:  If you want to find a job

- Machine Learning: If you want to hire people

- Artificial Intelligence: If you want to raise money

- Data Mining: Do not use it anymore

# Importance of Data Mining

- Current databases can not be analyzed manually
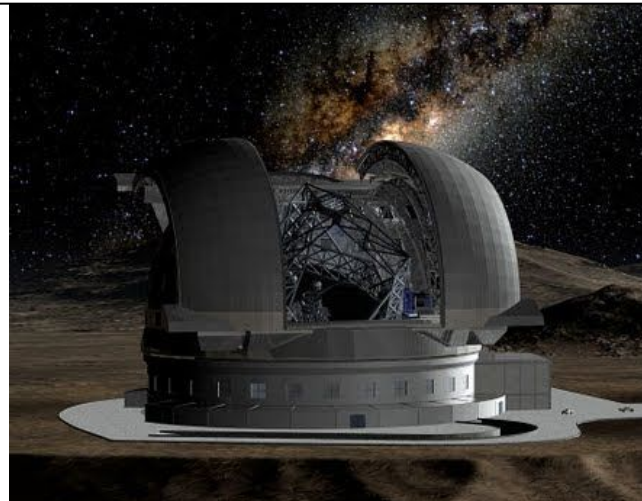- Automatic tools are needed to perform different tasks on data

# Data Everywhere!!



2) Please rate your satisfaction with each of the features of our restaurant.

| | Very satisfied | Somewhat satisfied | Neither satisfied nor dissatisfied | Somewhat dissatisfied | Very dissatisfied |
|---|---|---|---|---|---|
| Overall Food Quality | ○ | ● | ○ | ● | ○ |
| Your Server | ○ | ● | ○ | ● | ○ |
| Dining Area Cleanliness | ○ | ● | ○ | ● | ○ |
| Valet Parking | ○ | ● | ○ | ● | ○ |
| Decor, Music, Lighting | ○ | ● | ○ | ● | ○ |
| Restroom Cleanliness | ○ | ● | ○ | ● | ○ |

Additional comments:

DATA

# Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

| SUMMARY | SAVE | SHARE | COMMENT 5 | TEXT SIZE | PRINT | $8.95 BUY COPIES |

**W**hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and

DECISION MAKING

# Big Data: The Management Revolution

by Andrew McAfee and Erik Brynjolfsson
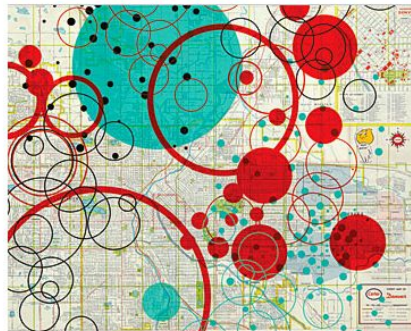
FROM THE OCTOBER 2012 ISSUE

| SUMMARY | SAVE | SHARE | COMMENT 0 | TEXT SIZE | PRINT | $8.95 BUY COPIES |



ARTWORK: TAMAR COHEN, HAPPY MOTORING, 2010, SILK SCREEN ON VINTAGE ROAD MAP, 26″ X 18″

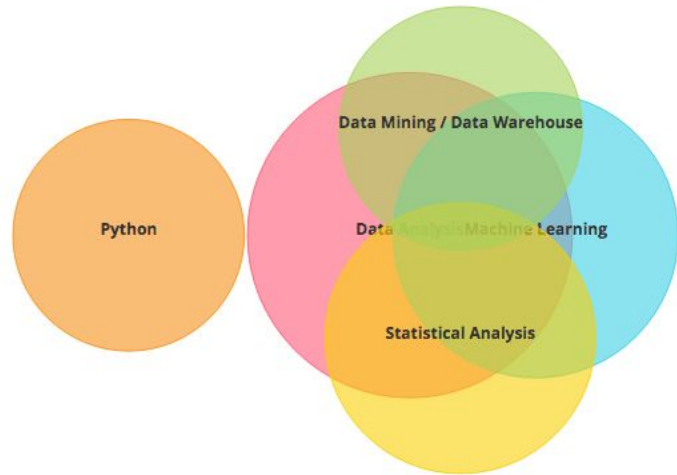**"Y**ou can't manage what you don't measure."

There's much wisdom in that saying, which has been attributed to both W. Edwards Deming and Peter Drucker, and it explains why the recent explosion of digital data is so important. Simply put, because of big data, managers can measure, and hence know, radically more about their businesses, and directly translate that knowledge into improved decision making and performance.

# Skills you will obtain from IIC2433



**Popular Skills for Data Scientist, IT**

*This chart shows the most popular skills for this job and what effect each skill has on pay.*
http://www.payscale.com/research/US/Job=Data_Scientist,_IT/Salary

**MODERN DATA SCIENTIST**

**MATH & STATISTICS**
- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants
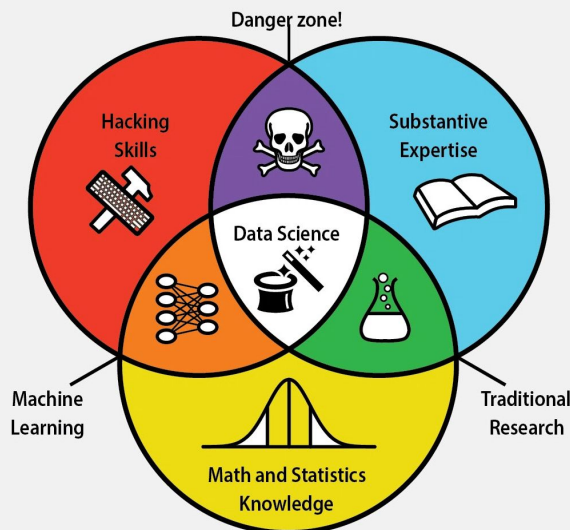
**PROGRAMMING & DATABASE**
- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing packages, e.g., R
- ☆ Databases: SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

**DOMAIN KNOWLEDGE & SOFT SKILLS**
- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative

**COMMUNICATION & VISUALIZATION**
- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

**DATA SCIENCE SKILLSET**

Data science, due to its interdisciplinary nature, requires an intersection of abilities: **hacking skills**, **math and statistics knowledge**, and **substantive expertise** in a field of science.

**Hacking skills** are necessary for working with massive amounts of electronic data that must be acquired, cleaned, and manipulated.

**Math and statistics knowledge** allows a data scientist to choose appropriate methods and tools in order to extract insight from data.

**Substantive expertise** in a scientific field is crucial for generating motivating questions and hypotheses and interpreting results.

**Traditional research** lies at the intersection of knowledge of math and statistics with substantive expertise in a scientific field.

**Machine learning** stems from combining hacking skills with math and statistics knowledge, but does not require scientific motivation.

**Danger zone!** Hacking skills combined with substantive scientific expertise without rigorous methods can beget incorrect analyses.

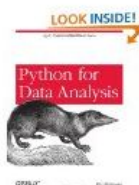http://berkeleysciencereview.com/article/first-rule-data-science/

# Recommender Systems

**amazon**.com

Your Amazon.com > **Recommended for You**
(If you're not Karim Pichara Baksai, click here.)

---

1.    LOOK INSIDE!

Python for
Data Analysis

O'REILLY

**Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython**
by Wes McKinney (October 29, 2012)
Average Customer Review: ★★★★☆ ☑ (58)
In Stock

**List Price:** $39.99
**Price: $25.24**
64 used & new from $16.74

Add t

☐ I own it   ☐ Not interested   ☒ ☆☆☆☆☆ Rate this item
Recommended because you purchased **Data Science for Business** and more (Fix this)

---

2.    LEARNING FROM DATA

**Learning From Data**
by Yaser S. Abu-Mostafa (March 27, 2012)
Average Customer Review: ★★★★☆ ☑ (62)
Available from **these sellers**.

19 used & new from $28.00

See all bu

☐ I own it   ☐ Not interested   ☒ ☆☆☆☆☆ Rate this item
Recommended because you purchased **Applied Predictive Modeling** and more (Fix this)

---

3.    LOOK INSIDE!

An Introduction
to Statistical
Learning

**An Introduction to Statistical Learning: with Applications in R (Springer Texts in Statistics)**
by Gareth James (August 12, 2013)
Average Customer Review: ★★★★★ ☑ (32)
In Stock

**List Price:** $79.99
**Price: $61.62**
61 used & new from $57.18

Add t

☐ I own it   ☐ Not interested   ☒ ☆☆☆☆☆ Rate this item
Recommended because you purchased **Applied Predictive Modeling** and more (Fix this)

---

4.    Bayesian Data Analysis

**Bayesian Data Analysis, Third Edition (Chapman & Hall/CRC Texts in Statistical Science)**
by Andrew Gelman (November 1, 2013)

10

# Personalized discounts

# Plant based food

# Prediction of Crypto-Currency price

# Customer profiling



http://www.rdginsights.com.au/hero-products

# Customer Abandonment

# Employee Evaluation

# Text mining
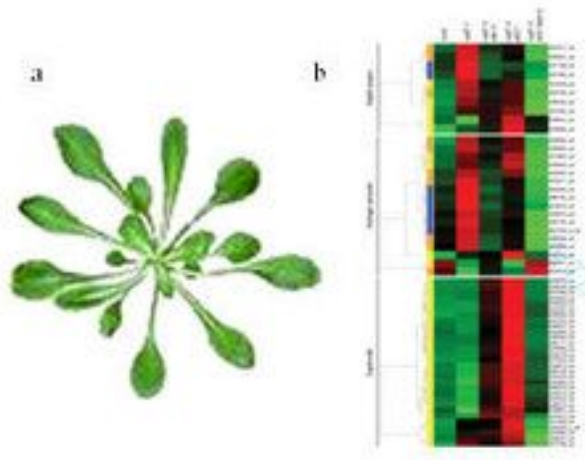
# Social Networks data mining

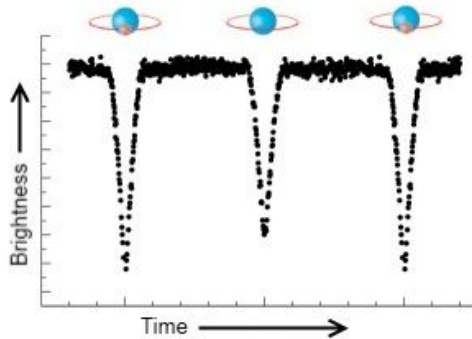# From offline to online world: Reporting and analysis

# Biotechnology

# Automatic classification of Variable stars
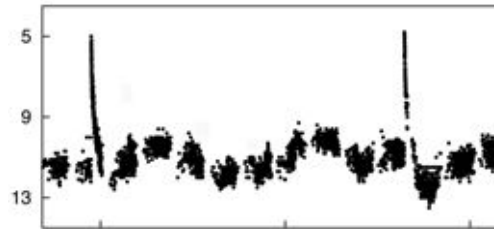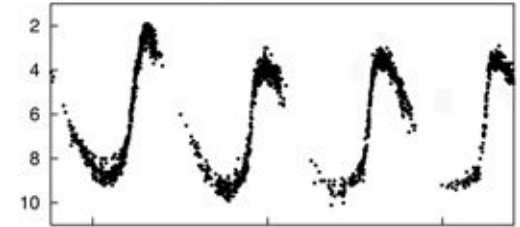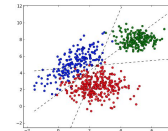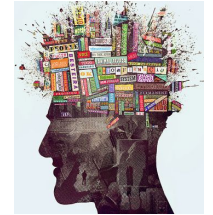
# Contenidos del Curso

**Knowledge**



Results

**4**. Evaluation and visualization

**3**. Data Mining
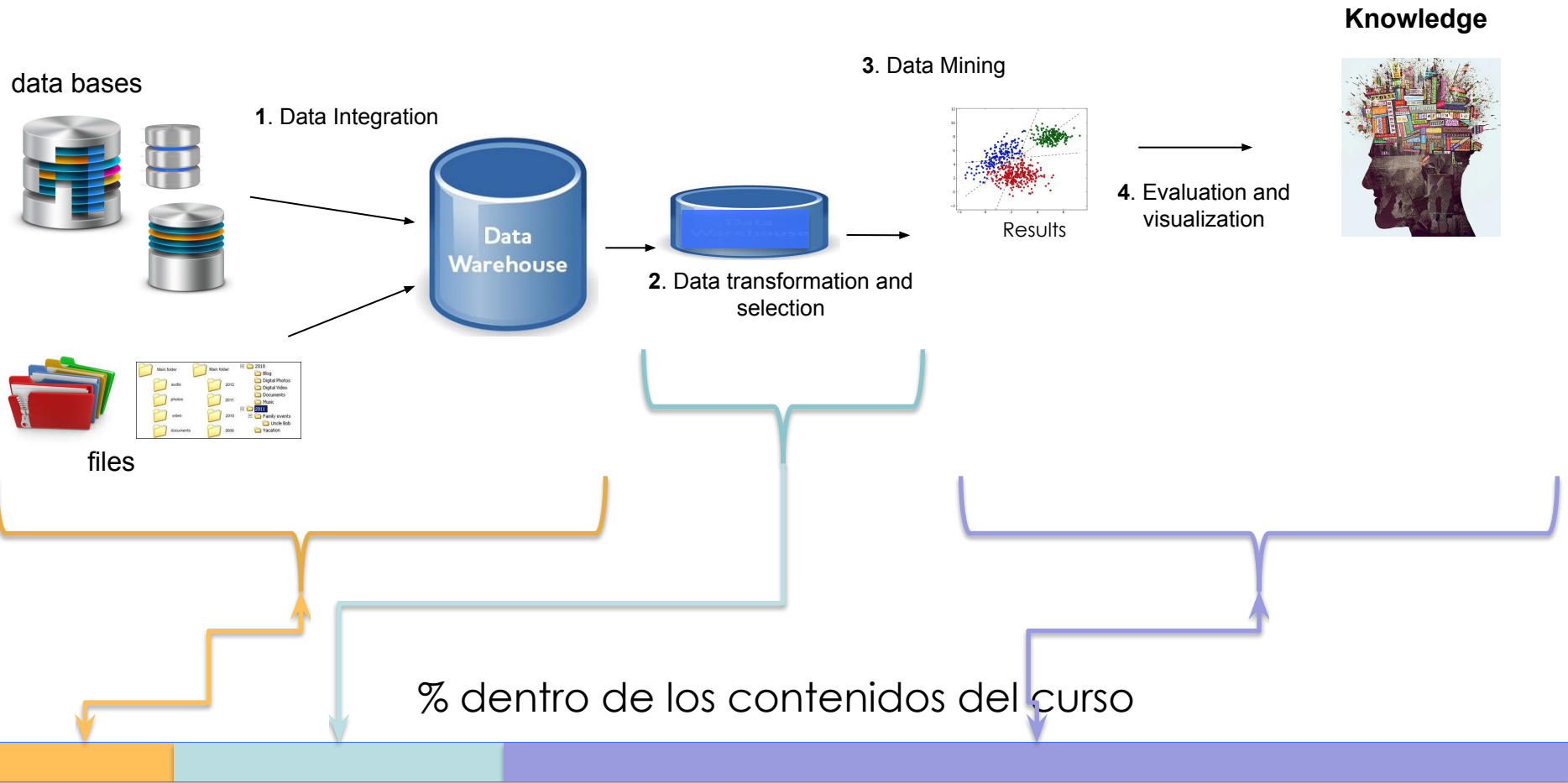
**2**. Data transformation and selection

Data Warehouse

**1**. Data Integration

data bases

files

Preprocessing

# Contenidos del Curso

data bases

**1**. Data Integration

**3**. Data Mining

**Knowledge**

Data Warehouse

**2**. Data transformation and selection

Results

**4**. Evaluation and visualization

files

% dentro de los contenidos del curso