

# Introducción a la Minería de Datos

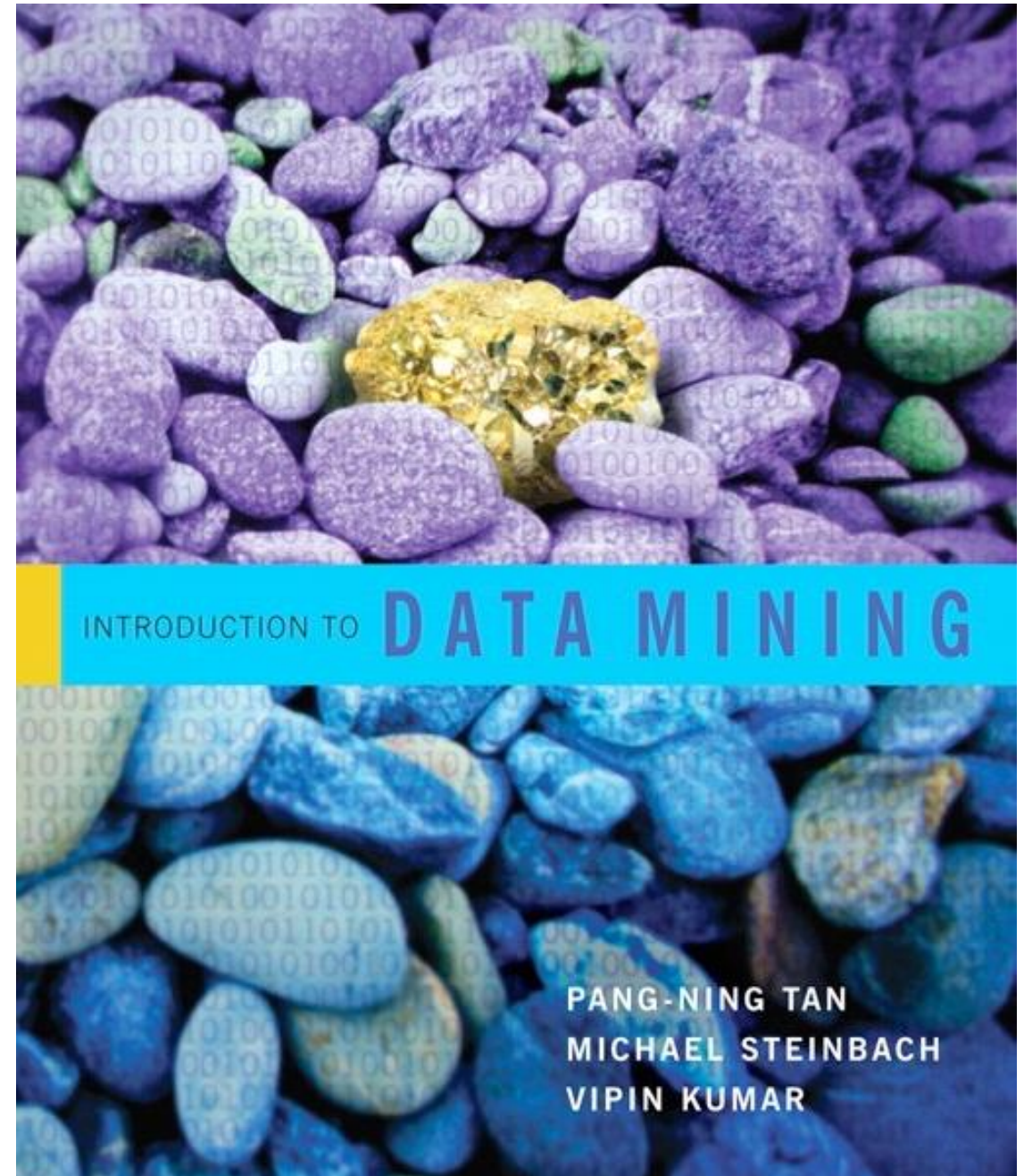
Basado en las slides de Bárbara Poblete

# Objetivos del curso

- **Curso introductorio**
- Aprender a aplicar el proceso de DM a datos reales
- Conocer, seleccionar y utilizar las técnicas básicas de DM
- Aprender a interpretar los resultados de estos procesos
- Proveer la base para adquirir conocimiento más avanzado

# Libro del Curso

- Introduction to Data Mining
- Autores: Pang-Ning Tan, Michael Steinbach, Vipin Kumar

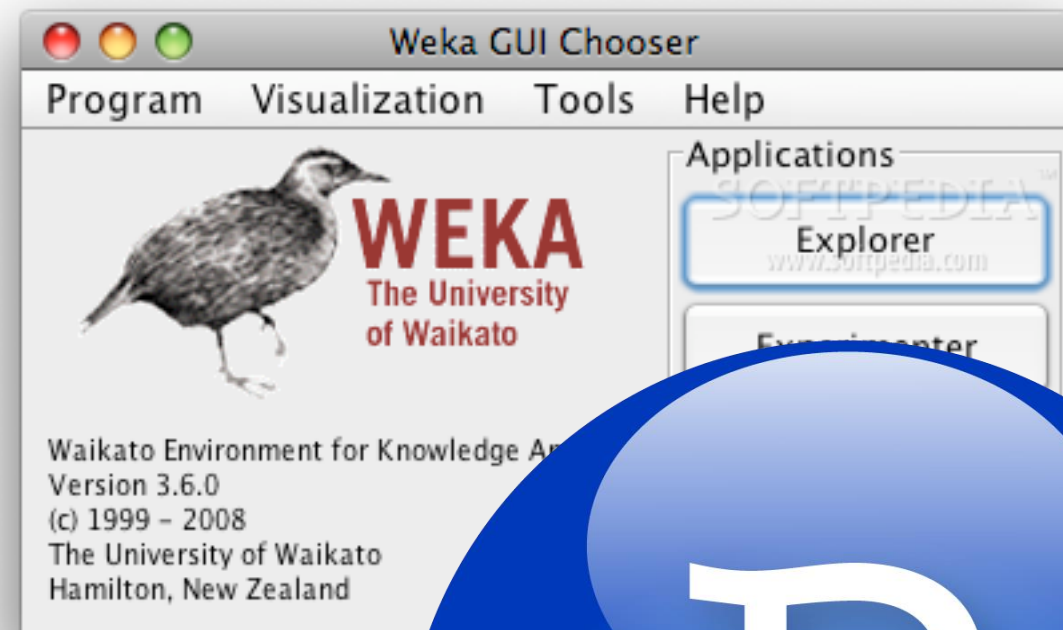


# Herramientas del curso

• ~~WEKA~~

• R (R Studio)

• Python



python™



¿Qué significa **Minar**?







# ¿Qué significa **Minar**?

Según la RAE:

*“Hacer grandes diligencias para conseguir algo”*

# ¿Qué es la Minería de Datos?

- Descubrir automáticamente información útil en grandes repositorios de datos



# ¿Qué es la Minería de Datos?

- Descubrir **automáticamente** información útil en grandes repositorios de datos

# ¿Qué es la Minería de Datos?

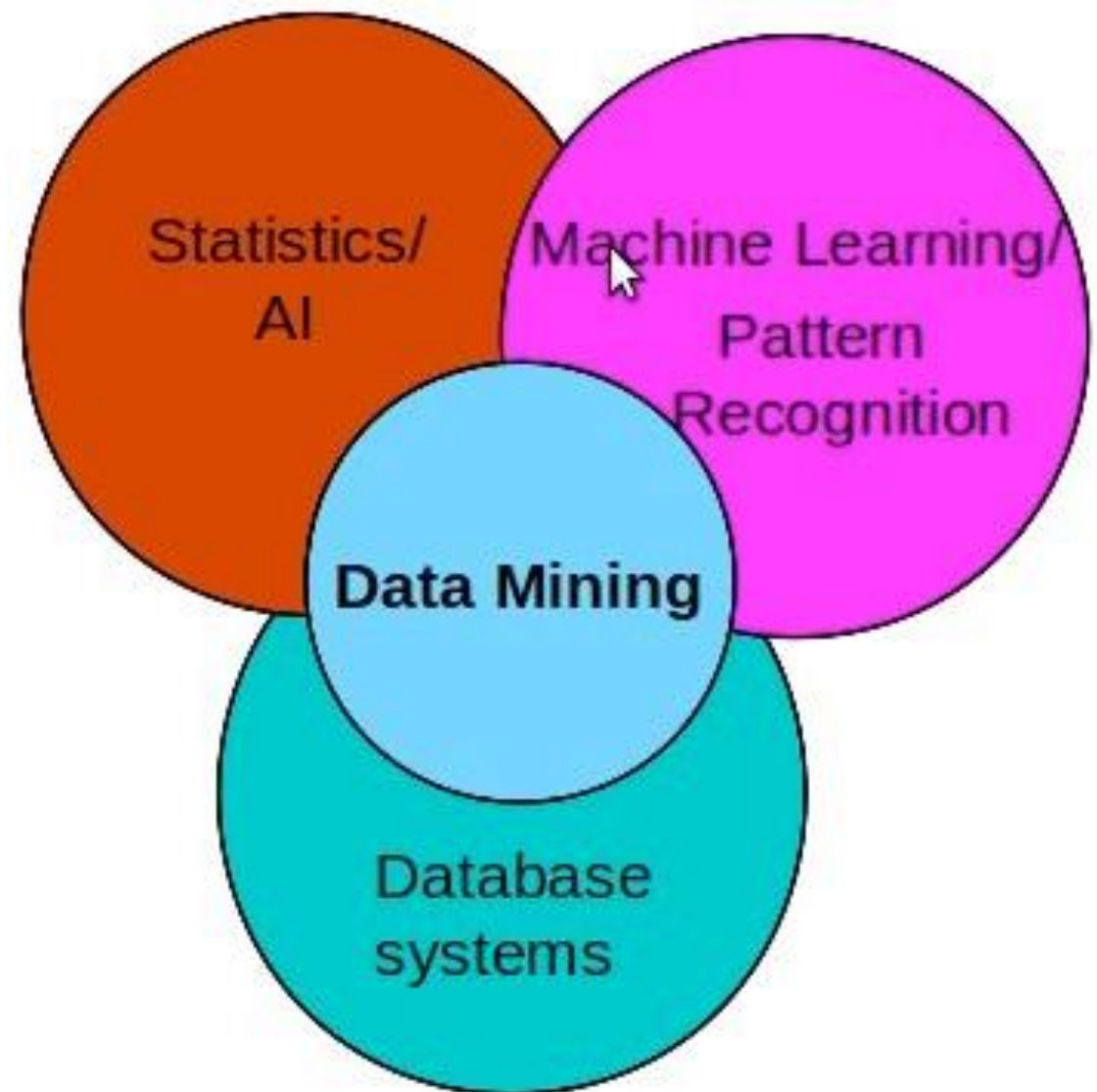
- Descubrir automáticamente información **útil** en grandes repositorios de datos

# ¿Qué es la Minería de Datos?

- Descubrir automáticamente información útil en **grandes repositorios** de datos

# Orígenes de la MD

- Une ideas de ML/AI, reconocimiento de patrones, estadística y BD
- Enfoques tradicionales fallan con datos masivos (alta dim., datos heterogéneos y distribuidos)





# ¿Cuál es la diferencia entre **Data Science**, **Machine Learning** e **Inteligencia Artificial**?

- Están de moda, pero no son lo mismo, ni son intercambiables
- **Data Science** es el nombre reciente para algo mucho más antiguo: **Data Mining (90's)**
- Definición (sobre) simplista:
  - **Data mining** genera **entendimiento**.
  - **Machine learning** genera **predicciones**.
  - **Artificial intelligence** genera **acciones**.

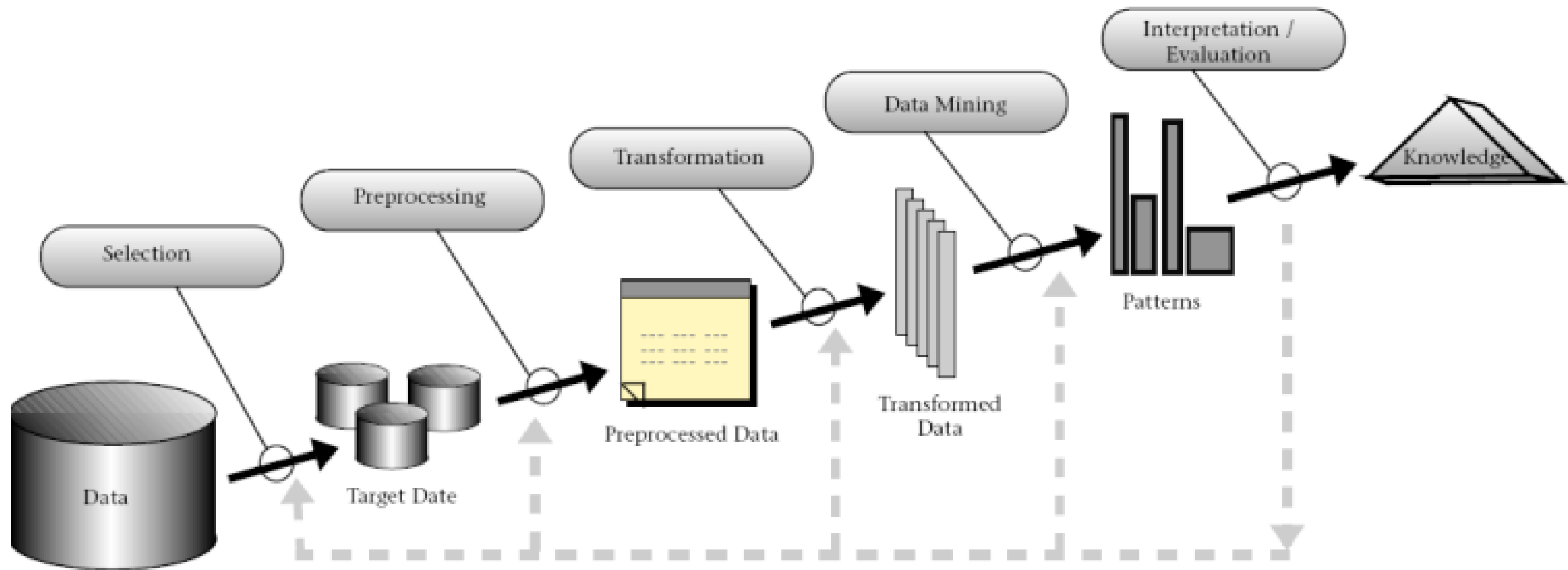
# ¿Cuál es la diferencia entre Data Science, Machine Learning e Inteligencia Artificial?

- **Artificial Intelligence:** auto reconoce una señal de STOP y toma la **acción** de frenar.
- **Machine Learning:** auto reconoce señales de STOP usando cámaras y **predice** en base a un entrenamiento cuando debe parar.
- **Data Mining:** auto transita por las calles y nos damos cuenta que su rendimiento no es el esperado. Luego, **entendemos** que esto se debe a varios factores externos.



# ¿Por qué es importante entender estas diferencias?

- Porque este **no es un curso de Machine Learning**, es un **curso de Minería de Datos**.
- **ML: Estudio, diseño y desarrollo de algoritmos** que permiten a los computadores aprender sin ser explícitamente programados (Arthur Samuel). Técnicas genéricas, aplicables a varios dominios.
- **Minería de Datos:** El enfoque está en **extraer conocimiento**, o patrones previamente desconocidos, a partir de (grandes) volúmenes de datos (en su mayoría no estructurados). Para esto se pueden utilizar técnicas de ML, entre otras. Requiere conocimiento de los datos mismos y su dominio.



## Knowledge Discovery in Databases (KDD)



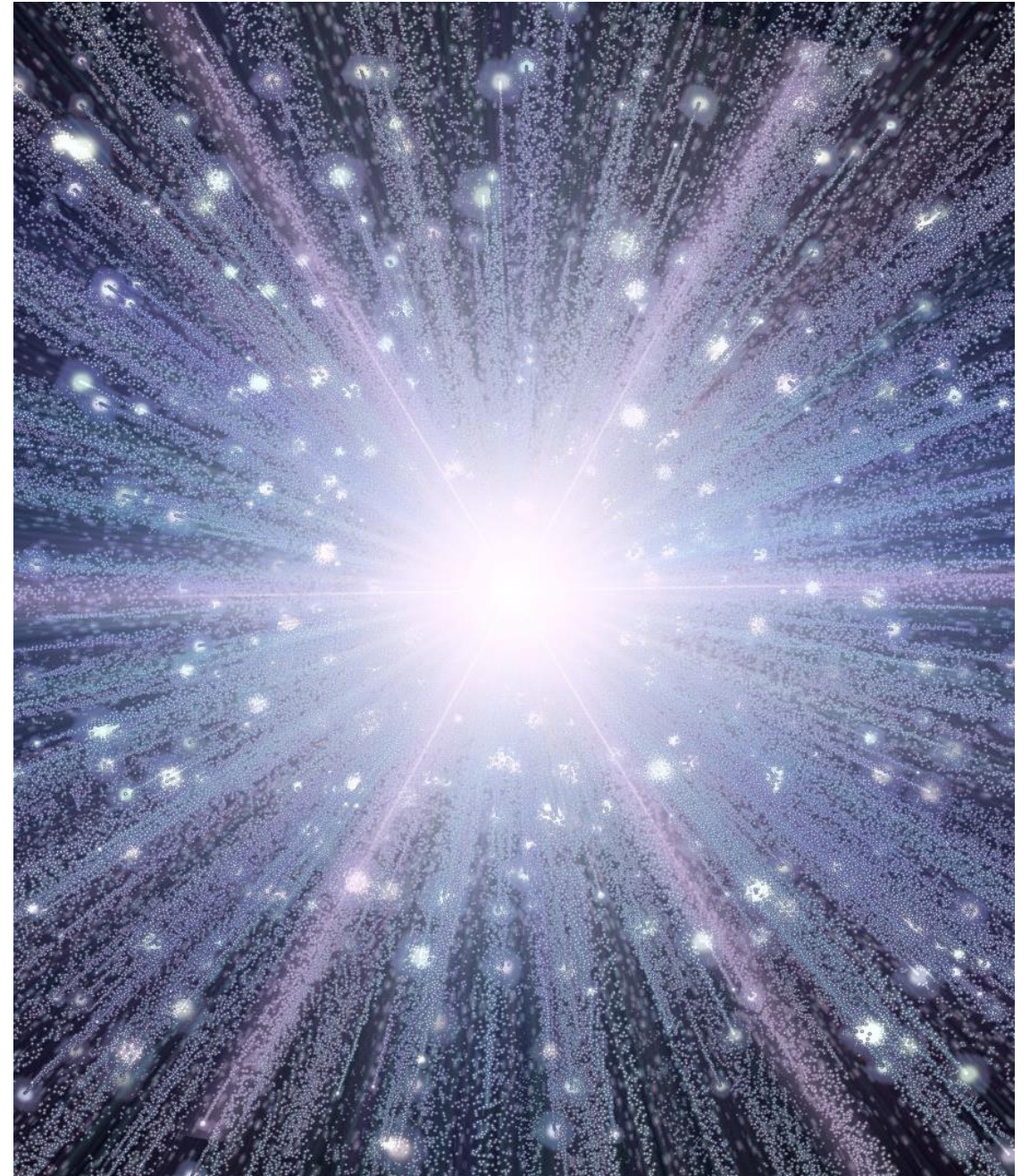
# Introducción a la Minería de Datos

## Parte 2 - Aplicaciones

Basado en las slides de Bárbara Poblete

# BIG BANG

- 2006 Hadoop
- Análisis de datos masivos al alcance de todos (cientos de start-ups)





# ¿Por qué hacer minería de datos?

- Aspecto comercial
- Aspecto científico



# ¿Por qué hacer minería de datos?

# ¿Motivación Comercial?

- Recolección de MUCHOS datos comerciales:
- Datos Web, e-commerce
- Compras en tiendas
- Transacciones en Bancos/ Tarjetas de Crédito







**Barbara's Amazon**

ON ORDER

0 items

AMAZON PRIME

Try Prime

View benefits

AUDIBLE AUDIOBOOKS

Try Audible

Get 2 free audiobooks

CUSTOMER SINCE

2008

## Recommended for you, Barbara



Literature & Fiction  
100 ITEMS



Science Fiction & Fantasy Books  
41 ITEMS



Prime Video – Unlimited Streaming for Prime Members  
27 ITEMS



Mystery, Thriller & Suspense Books  
55 ITEMS



Personal Care Products  
81 ITEMS



Recommended Based On *Sketching User Experiences: Getting the Design...*  
16 ITEMS



Office & School Supplies  
20 ITEMS



Cell Phones & Accessories  
10 ITEMS



Roll over image to zoom in

## 1 X Disney Frozen Pencil Case

by Innovative Designs, LLC

★★★★★ 4 customer reviews

Price: **\$5.30**

**In Stock.**

This item ships to **Santiago, Chile**. Want it **Friday, March 11**? Order within **9 hrs 49 mins** and choose **Amazon Global Priority Shipping** at checkout. [Learn more](#)

Sold by **JACOB'S** and **Fulfilled by Amazon**. Gift-wrap available.

Package Quantity: **1**

Style Name: **Purple**

- 1 Disney Frozen Pencil Case

15 new from **\$1.50**



### Frequently Bought Together



Total price: **\$22.55**

Add both to Cart

Add both to List

✓ **This item:** 1 X Disney Frozen Pencil Case **\$5.30**

✓ Thermos 12 Ounce Funtainer Bottle, Frozen Purple **\$17.25**

### Customers Who Bought This Item Also Bought

Page



Disney Frozen Light Blue Stationery Set Pack with Case (13 Pcs)

★★★★★ 39

\$7.40 ✓Prime



Disney Frozen Rolling 16" Backpack and Lunch Bag Lunchbox 2pc

★★★★★ 12

\$49.95 ✓Prime



Disney Frozen 1 Subject Wide Ruled Notebook - (Colors/Graphics Vary)

★★★★★ 14

\$4.67



Disney Frozen Elsa and Anna Kids Stationery Set (17 Pcs)

★★★★★ 19

\$8.95 ✓Prime



American Greetings Frozen Party Accessories, Pencils, 12 Count

★★★★★ 89

\$5.26 ✓Prime



Disney Frozen Hot Pink Elsa Anna and Olaf Stationery Set Pack with Case (13 Pcs)

★★★★★ 34



Thermos 12 Ounce Funtainer Bottle, Frozen Purple

★★★★★ 4

\$17.25 ✓Prime



## TV Thrillers &amp; Mysteries



## Romantic Movies



## Continue Watching for Barbara



## Watch It Again



## Top Picks for Barbara





Top Picks for Barbara



House, M.D.



★★★★★ 2014 TV-14 9 Seasons  
An awkward forensic anthropologist. An arrogant FBI agent. Together, they find justice in the dead.



★★★★★ 2010 TV-14 3 Seasons  
His deception detection is second to none. But his social skills? Well, they could use a little work.



★★★★★ 2015 18 11 Seasons  
Neither their patients' problems nor their own relationships are black-and-white. It's all shades of grey.



★★★★★ 2013 14 1 Season  
Elite FBI profilers play minds games to catch serial killers. Getting into murderers' heads can also get into yours.



★★★★★ 2015 TV-14 3 Seasons  
The legendary detective needs a doctor to keep him clean -- and maybe help round up a few murderers.

OVERVIEW

EPISODES

MORE LIKE THIS

DETAILS

Because you watched Cooked

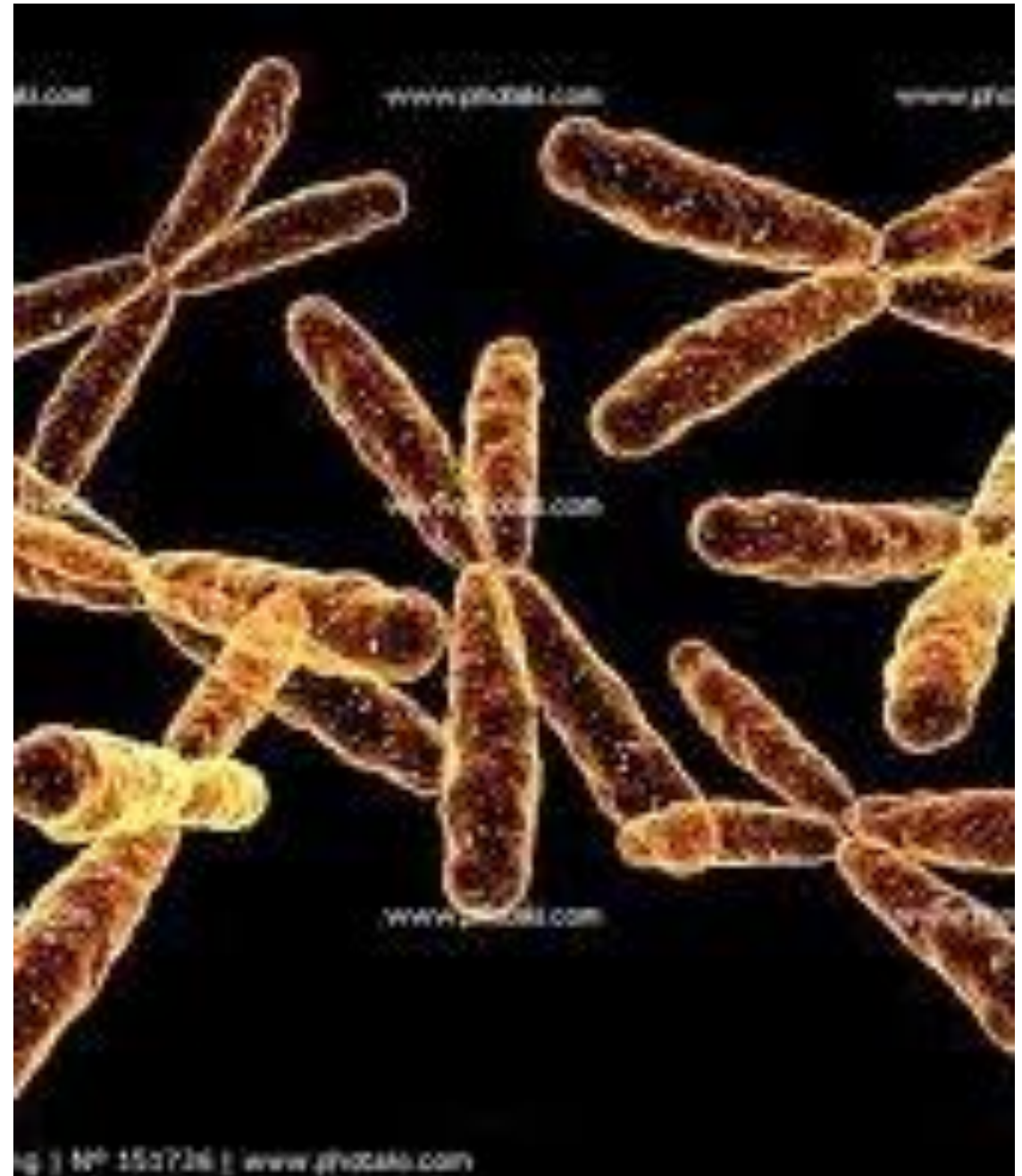




# ¿Por qué hacer minería de datos?

## ¿Motivación Científica?

- Datos (observaciones) recolectadas a gran velocidad (GB/hr, Tb/día)
- Telescopios, Satélites, Requerimientos Web, ADN, etc ([Google Flu Trends](#))



[Google.org home](#)

## Flu Trends

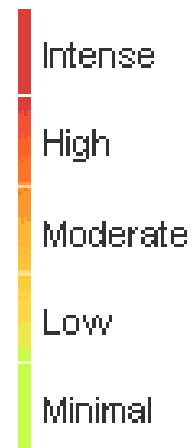
Select country 

[Home](#)

[How does this work?](#)

[FAQ](#)

### Flu activity



## Explore flu trends around the world

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. [Learn more »](#)



[Download world flu activity data](#)

# Introducción a la Minería de Datos

Parte 3 – Métodos y Técnicas

Basado en las slides de Bárbara Poblete

# Métodos utilizados en DM

- **Métodos predictivos:** Usar variables para predecir variables desconocidas o valores futuros de otras variables
- **Métodos descriptivos:** Encontrar patrones interpretables por humanos que permitan describir los datos



# Métodos utilizados en DM

- **Clasificación (Predictivo)**
- **Clustering (Descriptivo)**
- **Descubrimiento de Reglas de Asociación (Descriptivo)**
- **Descubrimiento de Patrones Secuenciales (Descriptivo)**
- **Regresión (Predictivo)**
- **Detección de Desviación (Predictivo)**



# Clasificación

- Set de Entrenamiento (atributos incluyendo clase)
- Busca modelar en atributo clase
- Objetivo: asignar la clase más correcta a records nuevos
- Set de Evaluación

*Categorica*

*Categorica*

*Categorica*

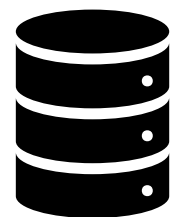
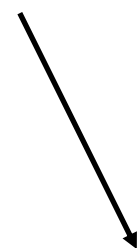
*Discreta*

*Categorica*

*Clase*

Nombre	Tipo sangre	Puede volar	Patas	Vive en el agua	Especie
Humano	Caliente	No	2	No	Mamífero
Rana	Fría	No	4	A veces	Anfibio
Paloma	Caliente	Si	2	No	Ave
Delfín	Caliente	No	0	Si	Mamífero

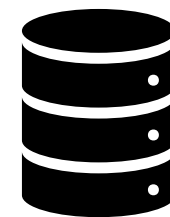
Nombre	Tipo sangre	Puede volar	Patas	Vive en el agua	Especie
Tortuga	Fría	No	4	A veces	?
Búho	Caliente	Si	2	No	?



**Conjunto de  
entrenamiento**



**Entrenar  
Clasificador**



**Conjunto  
de prueba**



**Modelo**

# Clasificación:

## Aplicación 1

- Marketing directo
- Meta: Reducir costos de publicidad apuntando directamente a potenciales compradores.
- ¿Cómo?

# CLUSTERING

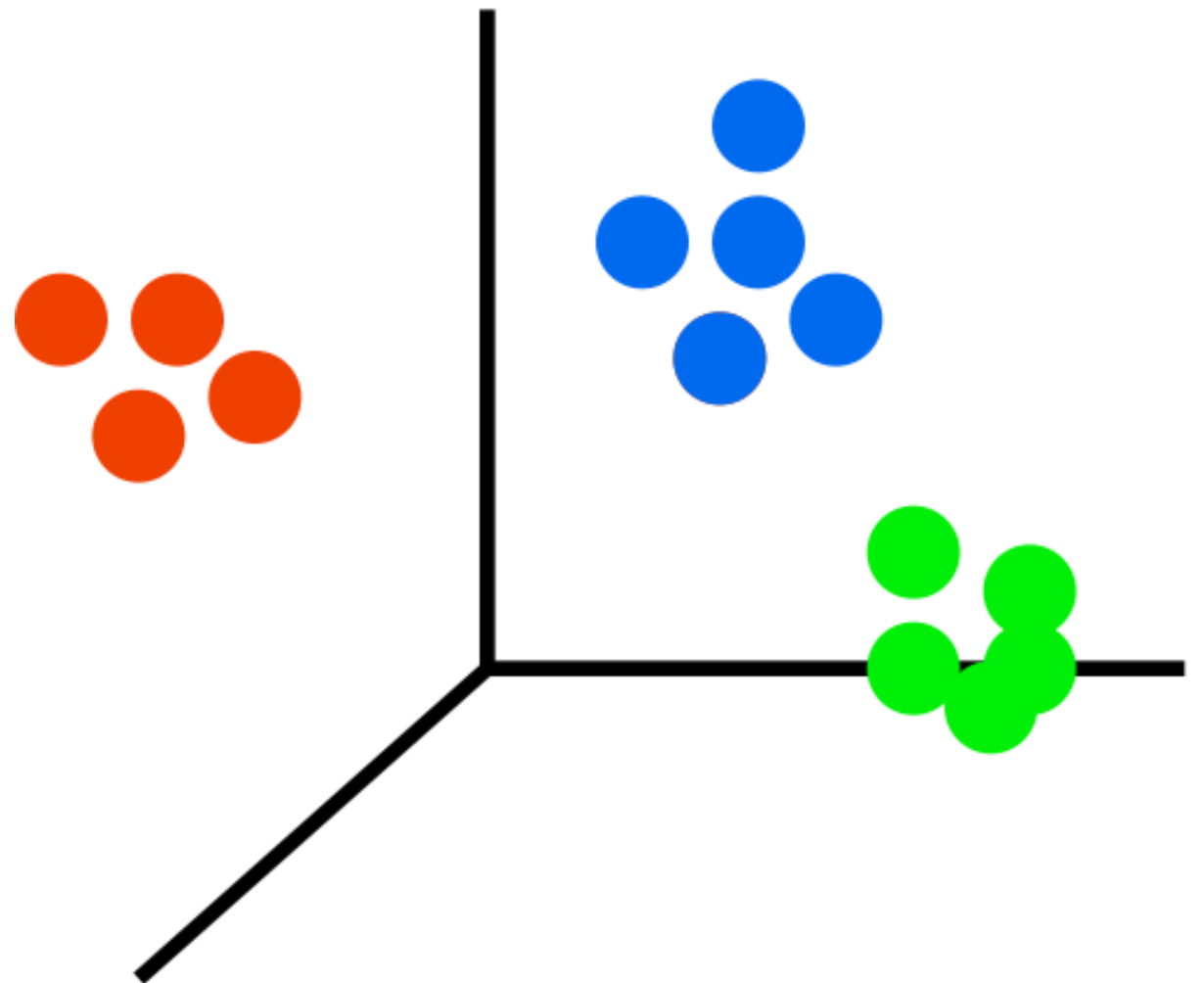
- Conjunto de puntos (datos), cada uno con un set de atributos y una medida de similitud
- Encontrar conjuntos tales que:
  - Puntos en un *cluster* sean más similares entre sí
  - Puntos en conjuntos diferentes sean menos similares entre sí

<i>Categorica</i>	<i>Categorica</i>	<i>Categorica</i>	<i>Discreta</i>	<i>Categorica</i>	<i>Clase</i>
Nombre	Tipo sangre	Puede volar	Patatas	Vive en el agua	Especie
Humano	Caliente	No	2	No	Mamífero
Rana	Fría	No	4	A veces	Anfibio
Paloma	Caliente	Si	2	No	Ave
Delfín	Caliente	No	0	Si	Mamífero
Tortuga	Fría	No	4	A veces	Reptil
Búho	Caliente	Si	2	No	Ave



# Visualización de clustering

- Clustering 3D basado basado en distancia Euclidiana
- Distancia intra-cluster es minimizada
- Distancia inter-cluster es maximizada



# Clustering

## Aplicación 1

- Segmentación de mercado
  - Meta: Subdividir un mercado en subconjuntos de clientes en donde cualquier conjunto es un potencial objetivo de marketing (ej: Netflix, Amazon)
  - ¿Cómo?

# Clustering

## Aplicación 2

- Clustering de documentos
  - Meta: Encontrar grupos de documentos que son similares entre sí, basándose en las palabras más importantes que contienen. (Directorios, Wikipedia)
  - ¿Cómo?

# Ejemplo

- Clustering de puntos: 3204 artículos del L.A. Times
- Medida de similitud: cuántas palabras tienen en común estos documentos (después de filtrar algunas palabras).

<i><b>Category</b></i>	<i><b>Total Articles</b></i>	<i><b>Correctly Placed</b></i>
<i><b>Financial</b></i>	555	364
<i><b>Foreign</b></i>	341	260
<i><b>National</b></i>	273	36
<i><b>Metro</b></i>	943	746
<i><b>Sports</b></i>	738	573
<i><b>Entertainment</b></i>	354	278

# Reglas de Asociación

- Dado un conjunto de records, cada uno contiene un número de elementos de una colección determinada
- Objetivo: Producir reglas de dependencia que predecirán la ocurrencia de un elemento (ítem) basándose en ocurrencias de otros ítems.



# Reglas de Asociación

TID

Items

1

Pan, Coca-cola, Pañales, Leche

2

Cerveza, Pan

3

Cerveza, Coca-cola, Pañales, Leche

4

Cerveza, Pan, Pañales, Leche

5

Coca-cola, Pañales, Leche

# Reglas de Asociación

## Aplicación 1

- Promoción de Marketing y Ventas
  - Sea la regla encontrada del tipo  
 $\{\text{Queso, ...}\} \longrightarrow \{\text{PapasFritas}\}$

# Patrones secuenciales

- Dado un set de objetos asociados a una línea de tiempo de eventos, encontrar los elementos que tengan fuertes dependencias secuenciales entre ellos
- Se forman reglas descubriendo patrones y luego se aplican restricciones de tiempo

# Regresión

- Predecir el valor de una variable continua, en base a valores de otras variables, asumiendo modelo de dependencia lineal o no-lineal.
- Estadística y redes neuronales



# Detección de desviación/anomalía

- Detectar desviaciones significativas de los valores normales

# Próxima Clase

- Leer reglas del curso: ver que no haya problemas con los requisitos de asistencia, se entiende que Ud. puede cumplirlos si sigue en el curso.
- Bonus track ver el video de [Hans Rosling](#).