



**Soft Management of Internet and Learning** 



# Análisis de Sentimientos y otros retos del aprovechamiento inteligente de los datos masivos

José A. Olivas

UCM, mayo, 2018

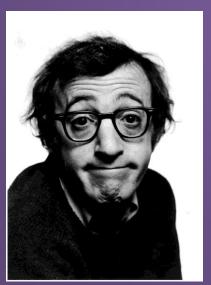




**Soft Management of Internet and Learning** 



# Desmontando a Harry, a Google y otros mitos de la era digital...



UCM, mayo, 2018





**Soft Management of Internet and Learning** 



# El aprendiz de Data Scientist...



UCM, mayo, 2018

# Comencemos por el final...

...los datos masivos

### ¿Qué son los datos?: DATOS / INFORMACIÓN /CONOCIMIENTO

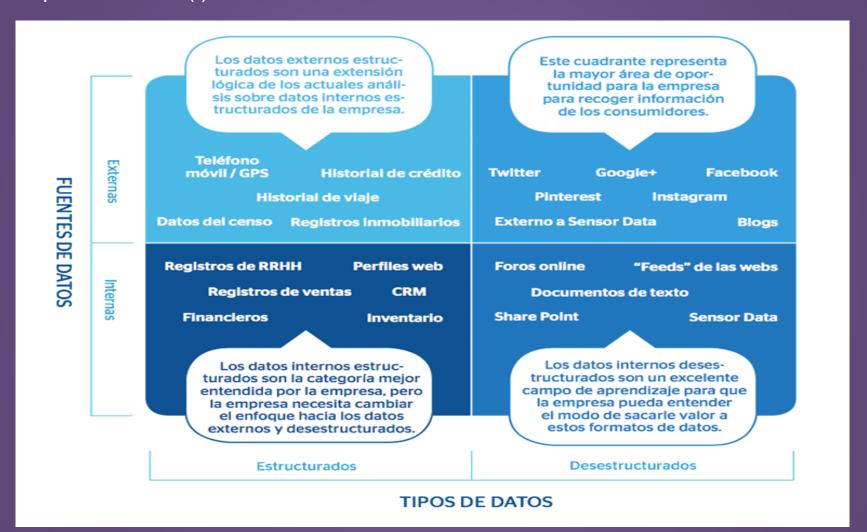








#### Tipos de datos (I)





# Tipos de datos (III)

Document Database	Graph Databases
Couchbase	Neo4j
■ MarkLogic mongoDB  Wide Column Stores	InfiniteGraph The Distributed Graph Database  Key-Value Databases
e redis	accumulo
amazon Dynamobb  in rick	HYPERTABLE™  Cassandra FIBASE  Amazon SimpleDB



¿ Dónde residen los datos ?

# What is a data lake?

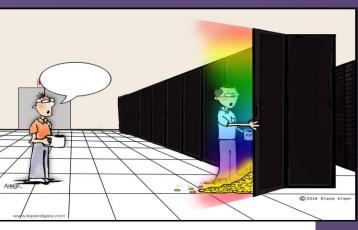
A repository for large quantities and varieties of data, both structured and unstructured.

Data generalists/ programmers can tap the stream data for real-time analytics.



The data lake accepts input from various sources and can preserve both the original data fidelity and the lineage of data

transformations Data



The lake can serve as a staging area for the data warehouse, the location of more carefully "treated" data for reporting and analysis in batch mode.



#### ¿ Cómo se consiguen ?



- Sistemas Transaccionales (operadores que recogen las peticiones a través de Call-Centers)
- Transacciones que se generan en las Webs (ficheros weblogs)
- Los sensores permiten capturar las magnitudes físicas o químicas y convertirlas en datos, por ejemplo temperatura, luz, distancia, aceleración, inclinación, desplazamiento, presión, fuerza, humedad, sonido, movimiento o el pH.
- IoT, Smart Cities...
- Redes sociales
- Etc, Etc...

#### ¿ Cómo se consiguen ?



- Sistemas Transaccionales (operadores que recogen las peticiones a través de Call-Centers)
- Transacciones que se generan en las Webs (ficheros weblogs)
- Los sensores permiten capturar las magnitudes físicas o químicas y convertirlas en datos, por ejemplo temperatura, luz, distancia, aceleración, inclinación, desplazamiento, presión, fuerza, humedad, sonido, movimiento o el pH.
- IoT, Smart Cities...
- Redes sociales
- Etc, Etc...



### El Business Intelligence (BI)



definición de business intelligence (BI)

La capacidad de transformar datos en información para ayudar a gestionar una empresa es el dominio de la inteligencia empresarial de negocios (BI), que consiste en los procesos, aplicaciones y prácticas que apoyen la toma de decisiones ejecutivas

#### **BI** operacional

- soporta funciones al nivel operacional
- capacidad en tiempo real o cerca de real-time
- comprende y cubre los procesos.



# Crítica...

### Demasiado restringido:

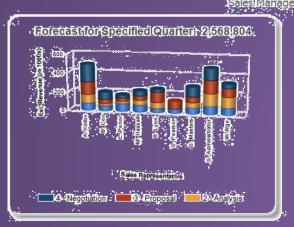
- …transformar datos en información…
- …apoyen la toma de decisiones…

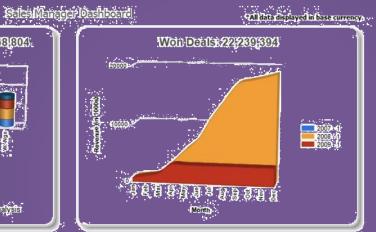
¡¡ Hay muchas otras cosas que se pueden hacer !!

Veamos las posibles salidas...



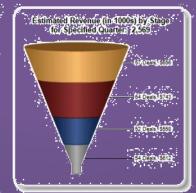
#### Outputs...





Top 10 Key Deals				
Account	Est. Close Date	Est. Revenue (in 1000s)	Recent Activity	
deWollie policies	<i>Piji</i> ws	\$74100	Q	
eyik eendi	E/[4]/20101	\$435570	Q.	
bikinine.	2//1/2010	\$434,30	øī.	
ilgii fioye	(8/25/2010)	\$431150	Ø.	
onsolidated essengen	5/j:(5/j:2010)	\$393,000	Ø:	
tycower@Ught	E[/E]/009	\$339.69	<b>(</b> )	
ontoso Partis	<b>E</b> [/2009	<b>\$338</b> (51	Ø	
akirkamilinay	4/20/2009	\$333512	Qi.	
ity Power (Allighi)	<i>5[19]</i> 2009	(55)263	Ø.	
adrona Solutions	13/9/x009	\$336,94	<b>6</b> 1	

Top 10 Sales Leaders in 2009			
Sales Representative	Actual Revenue (in 1000s)	Win Rate	
Anton Kiri by	<b>\$</b> 7,510,15821	4723	
System Administrator	<b>6</b> 3)640(8966	<b>493</b>	
eimon Pearson	<b>67/7</b> 6007613	4373	
Mark Hassall	<b>\$1</b> 1,867 <b>,</b> 92774	451%	
Brian Cox	<b>\$</b> 1129832422	461%	
William Ngo	<b>(</b> \$1\\$1515(\$228)	47%	
kokentiyon	<b>(51)</b> 454)(1014)	461%	
Lori Penor	<b>(\$</b> 1\\3372).6267	45%	
Steve Masters	(\$11,359,97723	443%	
aoini¢i∈n	\$799.0690	41125	







# Crítica...

#### De nuevo demasiado restringido:

- Esto es sólo visualización
- Conocimiento...
  - Sistemas de Ayuda a la Decisión (DSS).
  - Sistemas Recomendadores (Recommender Systems).
  - Análisis de series temporales (Predicción vs Pronóstico).
- Segmentación.

#### ii Patrones!!

- Las salidas condicionan todo el proceso.
- No se debe ir "a ciegas" hacia delante





# Big Data

- Aproximación ingenua y crítica.
- Definición abierta de Big Data.

"Big Data" es en el sector de tecnologías de la información y la comunicación una referencia a los sistemas que manipulan grandes conjuntos de datos. Las dificultades más habituales en estos casos se centran en la captura, el almacenamiento, búsqueda, compartición, análisis y visualización. (Wikipedia)

"Big data" es un término aplicado a <u>conjuntos</u> <u>de datos</u> que superan la capacidad del <u>software</u> <u>habitual</u> para ser <u>capturados</u>, <u>gestionados</u> y <u>procesados</u> en un <u>tiempo razonable</u>. Los tamaños del "big data" se hallan constantemente en <u>aumento</u>. (Wikipedia)

# Big Data

- Aproximación ingenua y crítica.
- Definición abierta de Big Data.

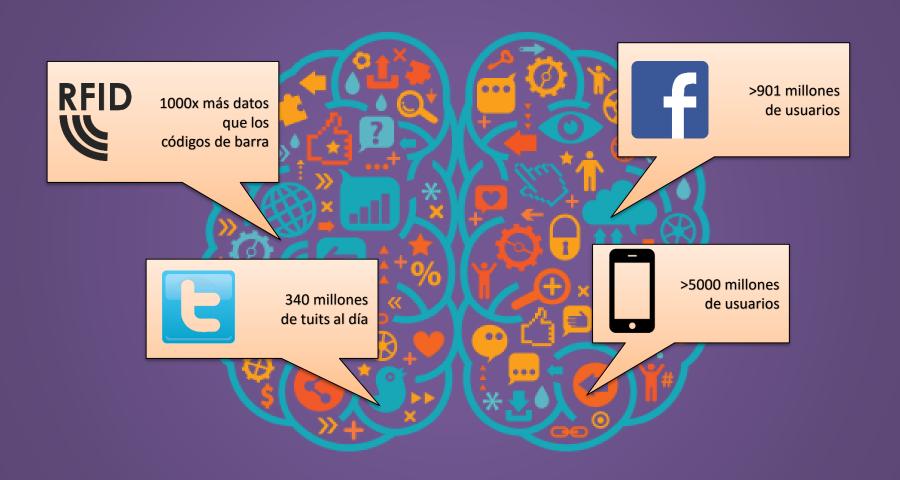
"Big Data" es en el sector de tecnologías de la información y la comunicación una referencia a los sistemas que manipulan grandes conjuntos de datos. Las dificultades más habituales en estos casos se centran en la captura, el almacenamiento, búsqueda, compartición, análisis y visualización. (Wikipedia)

"Big data" es un término aplicado a <u>conjuntos</u> <u>de datos</u> que superan la capacidad del <u>software habitual</u> para ser <u>capturados</u>, <u>gestionados</u> y <u>procesados</u> en un <u>tiempo razonable</u>. Los tamaños del "big data" se hallan constantemente en <u>aumento</u>.

(Wikipedia)

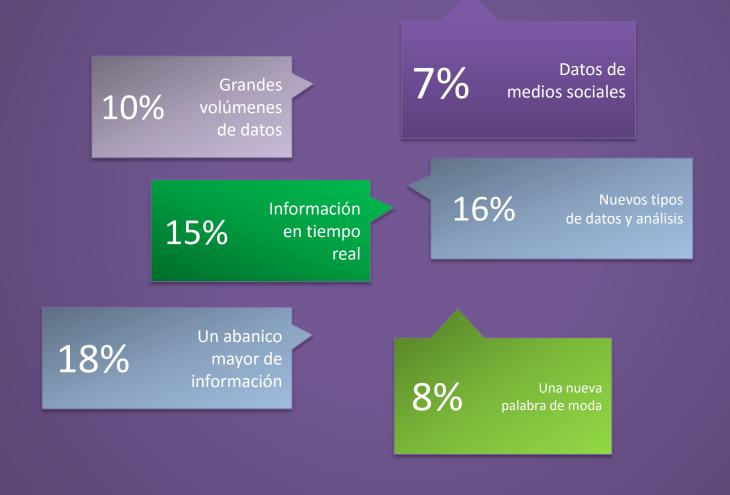
¡El nudo Gordiano de Google!



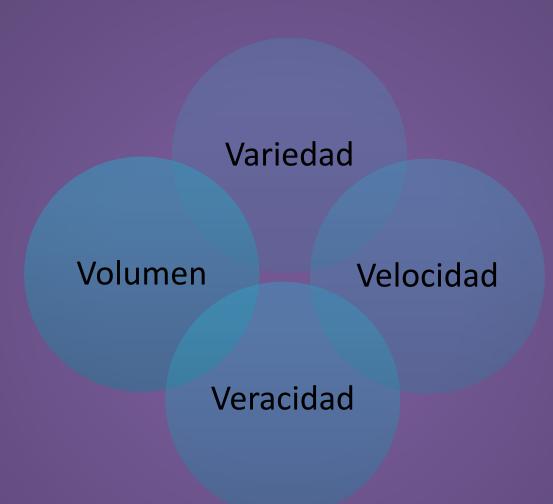




# ¿Cómo se percibe el Big Data?

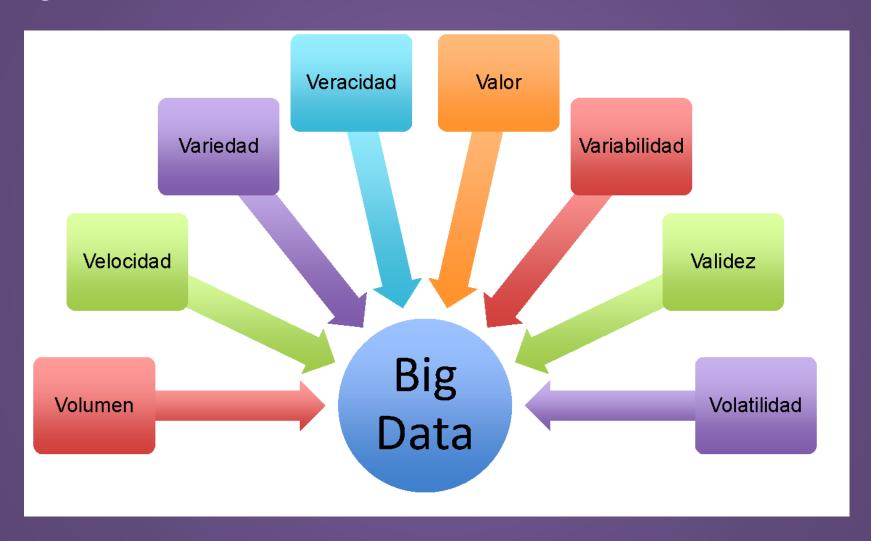








# ¿O las 8 V's?...



# Definición

- Datos...
- Información...
- ¿Conocimiento?
  - Abstracción-Patrones...
  - Dimensión humana...
  - La Web...
  - Google...

# ¡¡Explosión en la cantidad de datos!!

#### A380:

- Más de 1 billón de líneas de código.
- Cada motor genera 10 Tb cada media hora.
- Mas de 640 Tb de información por vuelo.
- Twitter genera más de 15 Tb de datos al día.
- Las principales bolsas generan más de 1 Tb al día.
- La capacidad de almacenamiento de ha doblado cada 3 años desde los 80s.

# ¡¡Explosión en la cantidad de datos!!

Historias Clínicas Electrónicas:

9.000.000.000 documentos sólo en España...

# ¡¡Problemas graves al gestionarlos!!

- A380 de Quantas (32-2009) ¡SATURACIÓN!
- A330 de Air France (447-2010) ¡INCONSISTENCIA!
- B777 Malayo (370-2014) ¡INCERTIDUMBRE!
- Twitter, ¡ANÁLISIS DE SENTIMIENTOS! PLN.
- No se usan las Historias Clínicas Electrónicas.

¡¡Explosión en la cantidad de datos!!

# ¿Habitualmente qué hacemos con todos estos datos?

¡¡Explosión en la cantidad de datos!!

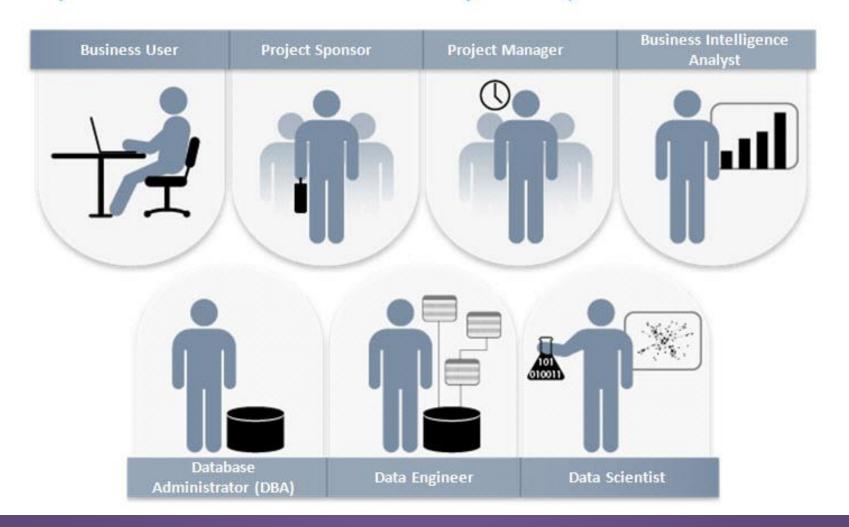
¿Habitualmente qué hacemos con todos estos datos?

ilGNORARLOS!

# ...el aprovechamiento inteligente

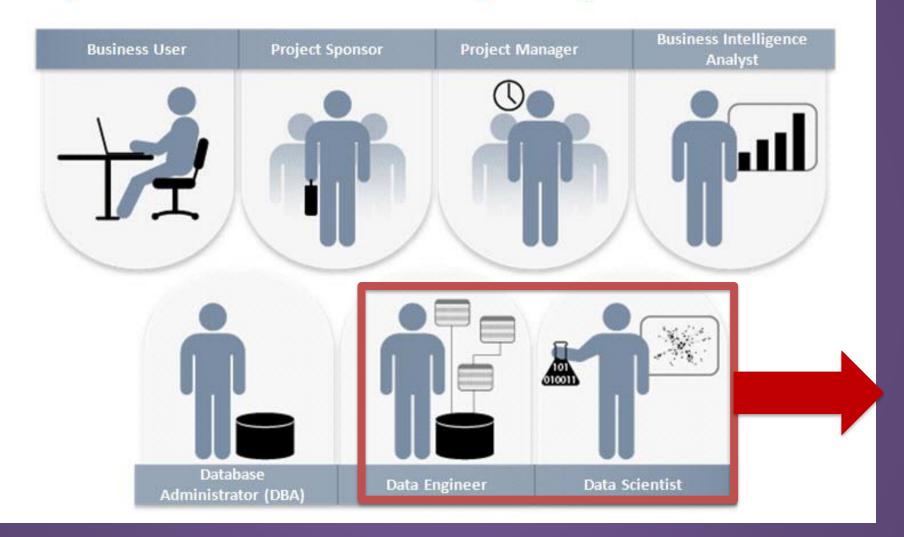


# Key Roles for a Successful Analytic Project

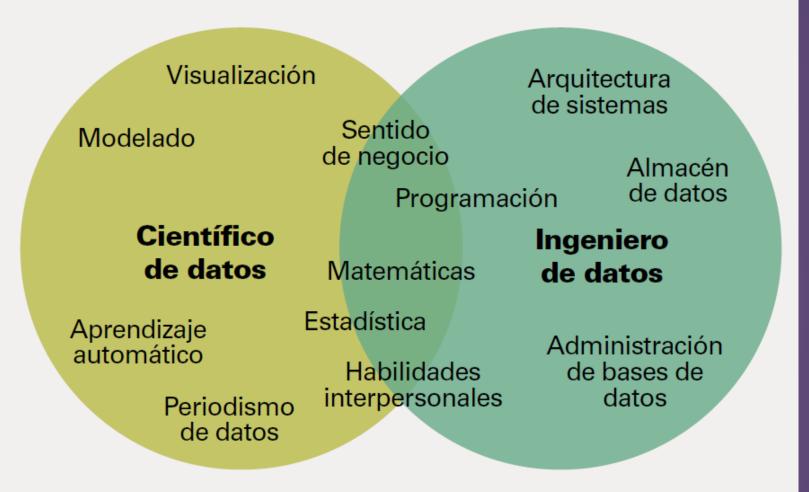




# Key Roles for a Successful Analytic Project





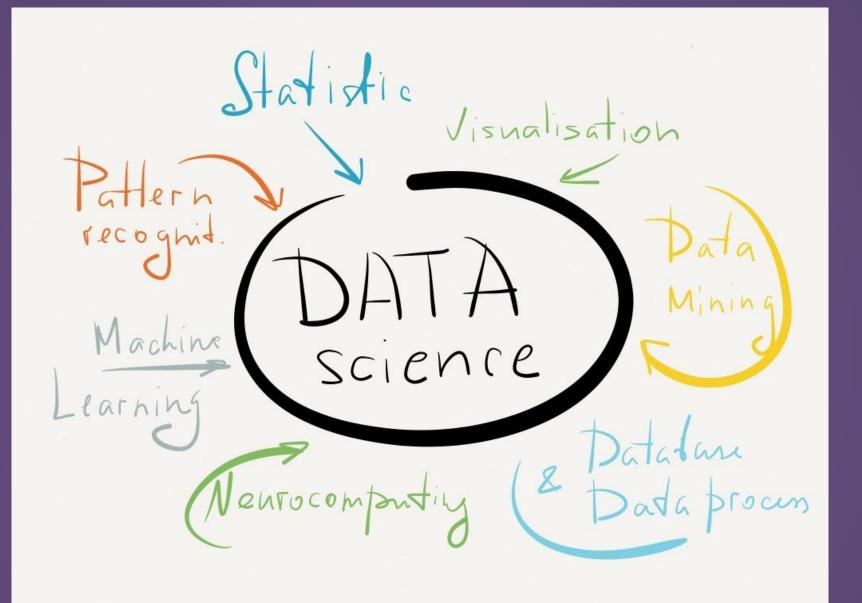


Fuente: Universitat Oberta de Catalunya. Máster en Business Intelligence y Big Data (2016)









# Métodos basados en la estadística

- Técnicas de Regresión y correlación
  - Lineal
  - Múltiple
  - Logística
  - CART (Classification And Regression Trees, Leo Breiman)
- Técnicas de extrapolación de funciones.
- Técnicas de aproximación y ajuste de funciones.
- Técnicas de agrupamiento basadas en medidas estadísticas (clustering).
- Etc.

# Métodos de Machine Learning (Aprendizaje automático)

### Aprendizaje por Analogía.

(Transformacional, derivacional, Razonamiento basado en Casos...)

#### Paradigma Inductivo.

Árboles de decisión, algoritmos de inducción pura...

### Paradigma Conexionista.

Redes Neuronales Artificiales...

#### Paradigma Evolutivo.

 Algoritmos Genéticos, otros métodos de optimización, colonias de insectos, descenso estocástico del gradiente...

### Modelos gráficos probabilistas.

Bayesianos, cadenas de Markov, Filtros de Kalman, redes de creencia,
 Máquinas de Soporte Vectorial (SVM), Metaheurísticas...

# Técnicas de Clustering: EJEMPLOS

- Clustering Jerárquico.
- Paradigma Conexionista.
  - Redes Neuronales Artificiales: SOM (Self Organized Maps,
     Mapas de Kohonen). Toolbox de Matlab SOM.
  - Etc.

### Técnicas de Clustering: EJEMPLOS

- Modelos estadísticos y probabilistas.
  - K-means, c-means,
  - K-nearest neighbours (KNN),
  - Mean shift (ventanas circulares con un centroide),
  - Dirichlet process (estocásticos basados en distribuciones de probabilidad). LDA (Latent Dirichlet Allocation),
  - Modelos Gaussianos,
  - Etc.

### Técnicas de Clustering: EJEMPLOS

- Extensiones basadas en Lógica Borrosa.
  - Fuzzy K-means,
  - Fuzzy c-means,
  - Isodata,
  - Etc.

### Técnicas de Clasificación: EJEMPLOS

- Paradigma Inductivo. Árboles de decisión:
  - ID3,
  - CART,
  - C4.5,
  - See5,
  - Random Forest (de moda en Big Data), Leo Breiman 2001,
  - Etc.

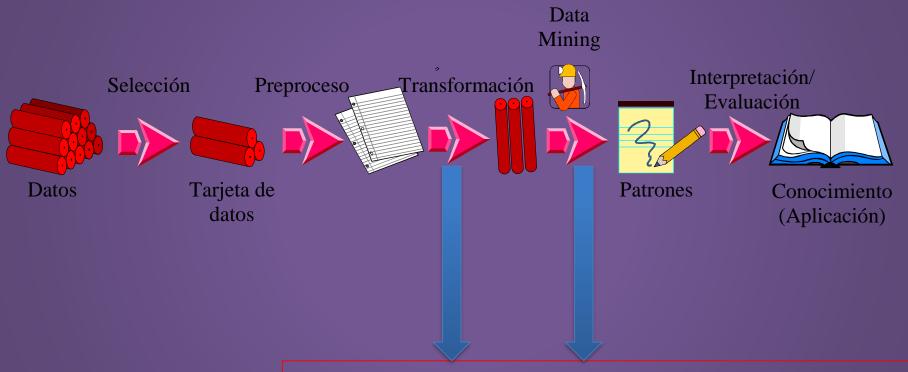
### Técnicas de Clasificación: EJEMPLOS

- Paradigma Conexionista. Redes Neuronales Artificiales:
  - Perceptrón Multicapa (con backpropagation),
  - Convolucionales,
  - Neocognitrones,
  - Redes de Hopfield,
  - Redes recurrentes,
  - Adaline,
  - Deep Learning (de moda en Big Data),
  - Etc.

### Técnicas de Clasificación: EJEMPLOS

- Modelos estadísticos y probabilistas.
  - Redes Bayesianas,
  - Naive-Bayes,
  - Máquinas de Soporte Vectorial (SVM),
  - Metaheurísticas,
  - Etc.

### La importancia del KDD



Consisten habitualmente en convertir un proceso de 'Clustering' en uno de 'clasificación'.

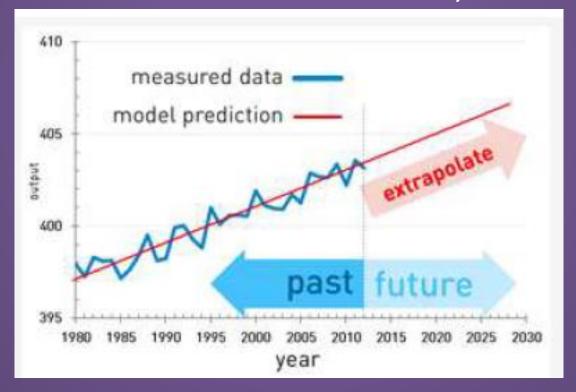
### Adecuación de los métodos a los problemas.

#### Análisis Predictivo

- Extrapolación de funciones (Tendencia para el futuro, pero no hay capacidad de pronóstico – hechos/cambios puntuales-).
- Correlaciones entre variables (Demasiado evidente, no suele funcionar de forma muy fina).
- Encontrar 'patrones' en los datos que puedan ser aplicados a situaciones futuras (KDD y Data Mining).
- Métodos de CLUSTERING Y CLASIFICACIÓN.

#### Análisis Predictivo

• Extrapolación de funciones (Por ejemplo Estimaciones o Líneas de Tendencia).



### Análisis Prescriptivo

- El análisis predictivo se centra en **un escenario** futuro.
- El prescriptivo se centra en múltiples alternativas.
- Por lo tanto, un modelo prescriptivo puede ser considerado como una combinación de modelos predictivos (uno por cada posible escenario), que se ejecutan en paralelo.
- El objetivo es encontrar la mejor opción posible: **OPTIMIZACIÓN**.

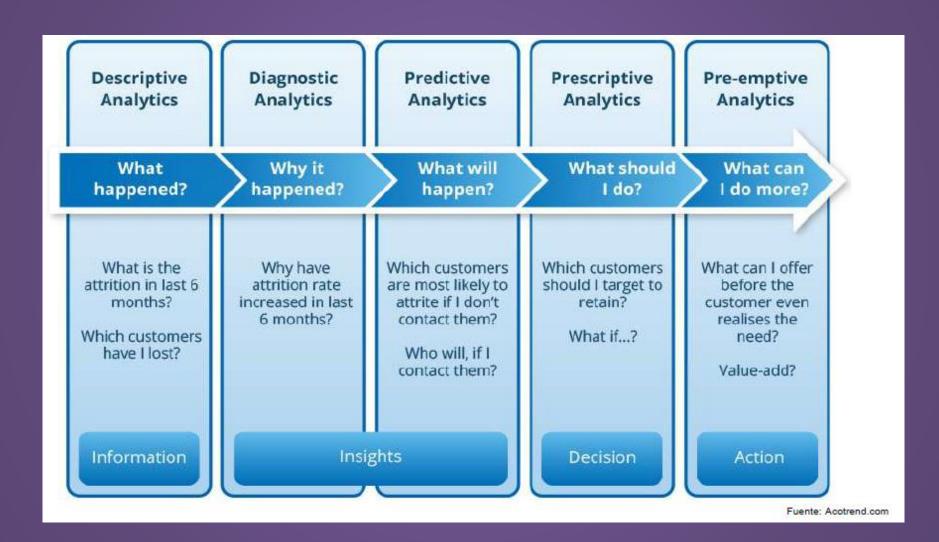
### Análisis Prescriptivo

### Técnicas:

- Técnicas de Investigación Operativa,
- Algoritmos Genéticos,
- Técnicas estocásticas,
- Metaheurísticas,
- Etc.

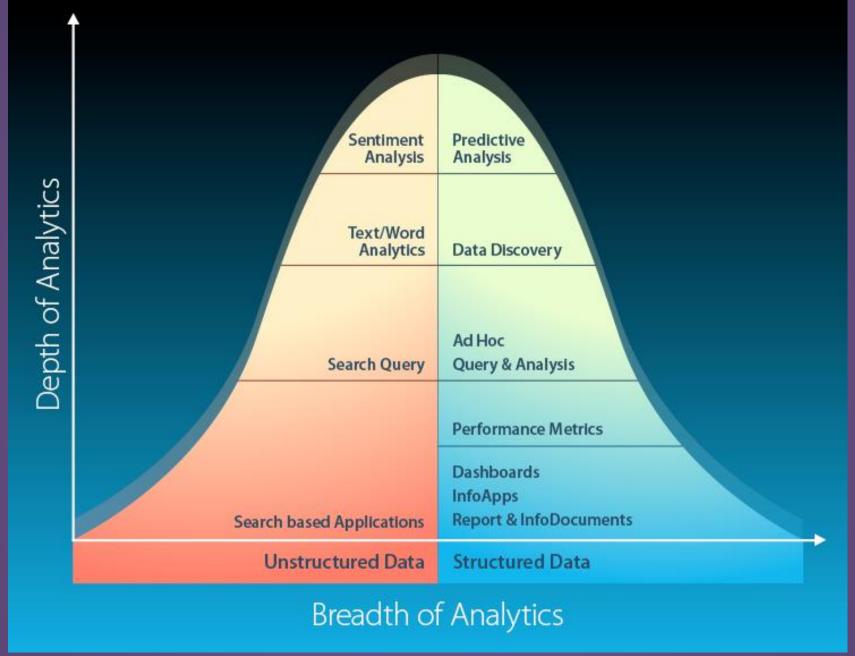


### Análisis Prescriptivo



### ...análisis de sentimientos y otros retos





## Acceso y la búsqueda de información digital.

- Los buscadores son eficientes, pero no eficaces.
- Ejemplos: sinonimia, veracidad (reputación), variedades diatópicas, operadores, tendencias...

"Búsqueda eficaz de información en la Web" (EDULP, 2011)

 El nuevo reto del análisis inteligente en Internet y las redes sociales:

"asíncrono" vs. "síncrono"

Reflexión/preparación vs. Inmediatez/visceralidad

• El nuevo reto del análisis inteligente en Internet y las redes sociales:

"asíncrono" vs. "síncrono"

Reflexión/preparación vs. Inmediatez/visceralidad

• ¡Dimensión Humana!

• El nuevo reto del análisis inteligente en Internet y las redes sociales:

"asíncrono" vs. "síncrono"

Reflexión/preparación vs. Inmediatez/visceralidad

¡Dimensión Humana! CONTEXTO

• El nuevo reto del análisis inteligente en Internet y las redes sociales:

"asíncrono" vs. "síncrono"

Reflexión/preparación vs. Inmediatez/visceralidad

¡Dimensión Humana! CONTEXTO

• El nuevo reto del análisis inteligente en Internet y las redes sociales:

"asíncrono" vs. "síncrono"

Reflexión/preparación vs. Inmediatez/visceralidad

- ¡Dimensión Humana! CONTEXTO
- El gran reto: PLN

• El nuevo reto del análisis inteligente en Internet y las redes sociales:

"asíncrono" vs. "síncrono"

Reflexión/preparación vs. Inmediatez/visceralidad

- ¡Dimensión Humana! CONTEXTO
- El gran reto: PLN
- Análisis de Sentimientos.

<sup>&</sup>quot;Sentiment analysis: A review and comparative analysis of web services" (Information Sciences, 2015)

• El nuevo reto del análisis inteligente en Internet y las redes sociales:

"asíncrono" vs. "síncrono"

Reflexión/preparación vs. Inmediatez/visceralidad

- ¡Dimensión Humana! CONTEXTO
- El gran reto: PLN
- Análisis de Sentimientos. ¿SIRI?

"Sentiment analysis: A review and comparative analysis of web services" (Information Sciences, 2015)





Lógica Borrosa Lotfi A. Zadeh 1921-2017







**Soft Management of Internet and Learning** 



# Análisis de Sentimientos y otros retos del aprovechamiento inteligente de los datos masivos

### iiMUCHAS GRACIAS!!

José A. Olivas

UCM, mayo, 2018

Joseangel.olivas@uclm.es



¡Hasta siempre Lotfi!