

Taller Minería de datos aplicados a la educación



1ª parte

Introducción a la minería de datos

27 de junio de 2011

Mercedes Torrado

*Departamento Métodos de Investigación y
Diagnóstico en Educación (MIDE)*

Este trabajo cuenta con licencia de Creative Commons:

Minería de datos aplicados a la educación está sujeta a una licencia de Reconocimiento-
NoComercial-SinObraDerivada 3.0 (CC BY-NC-ND 3.0)

Para citar la obra:

Torrado, M. (2011) *Minería de datos aplicados a la educación*. Barcelona: Universidad de
Barcelona. Deposito Digital <http://hdl.handle.net/2445/19862>



Introducción a la Minería de datos

- *INTRODUCCIÓN*
- *CONCEPTO*
- *PROCESO DE MINERÍA DE DATOS*
- *EJEMPLOS*
- *PROGRAMAS INFORMÁTICOS*
- *OTRAS APLICACIONES*



Introducción a la Minería de datos

Los avances tecnológicos en las últimas décadas nos han facilitado enormemente el acceso a grandes volúmenes de datos.

La cantidad de información que manejamos hoy en día nos obliga a abordar el estudio de los datos/información desde una **perspectiva global y no parcelada**

La preocupación por disponer de información suficiente para la toma de decisiones.



¿Cómo podemos analizar tal cantidad de información e identificar aquella que nos permita tomar decisiones y mejorar?



Introducción a la Minería de datos

En los años 90 apareció el concepto **DATA MINING**.

Esta técnica se vinculó estrechamente con la dirección de empresas y en concreto al marketing.

La minería de datos o Data Mining puede definirse como una extracción de información desconocida no trivial y potencialmente útil de una gran cantidad de información



Bajo este término se engloban un conjunto de técnicas de análisis cuyo objetivo es extraer conocimiento implícito de la base de datos.



Introducción a la Minería de datos

DEFINICIÓN

El **Data mining** también es considerado como una tecnología emergente que parte, por un lado de las técnicas estadísticas y por otro de las técnicas de inteligencia artificial Aluja, 2001 ⁽¹⁾

Estadística se ha preocupado más por la posible generalización de los resultados

Inteligencia artificial – ofrece soluciones algorítmicas a los datos

La **Minería de datos** comprende un conjunto de técnicas para la descripción y predicción a partir de grandes masas de datos

(Viera et al., 2009: 12)

[1] Tomàs Aluja en su artículo *La minería de datos, entre la estadística y la inteligencia artificial* publicado en el 2001 en la revista **QÜESTIIO** (vol 25, 3, p 479-498) hace todo un repaso de los orígenes de la Minería de datos en cuanto a los elementos correspondientes de la Estadística y de la inteligencia artificial



Equivalencias de nomenclaturas entre la Estadística y la Inteligencia Artificial

(Aluja, 2001: 482)

Inteligencia artificial	Estadística
Red (network)	Modelo
Ejemplos (patterns)	Observaciones, individuos
Inputs, outputs, features	Variables
Inputs	Variables explicativas
Outputs, targets	Variables de respuesta
Errores	Residuos
Training, learning	Estimación
Función de error, coste	Criterio de ajuste
Pesos, coef. sinápticos	Parámetros
Aprendizaje supervisado	Regresión, discriminación
Aprendizaje no supervisado	Clasificación

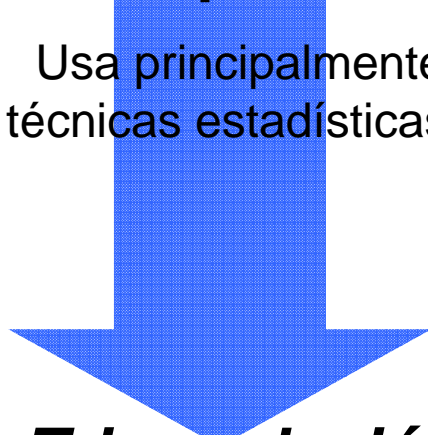


Introducción a la Minería de datos

Data mining se asocia a:

Minería predictiva

Usa principalmente técnicas estadísticas



Triangulación estadística

(Torrado, 2007)

Minería de datos para el descubrimiento del conocimiento

Usa principalmente técnicas de inteligencia artificial



Knowledge Discovery in Databases

(KDD)

Distintos tipos de conocimiento



SQL: Structured Query Language

OLAP: Online Analytical Processing

KDD: Knowledge Discovery on Databases



Introducción a la Minería de datos

Proceso de Descubrimiento de Conocimiento de Bases de Datos(KDD)

Las siglas **KDD** fue creada en 1995 para designar el conjunto de procesos, técnicas que propician el contexto en el cual la minería de datos tendrá lugar”

Una posible definición:

(Viera et al., 2009)

La integración de un conjunto de áreas que tienen como propósito la identificación de un conocimiento obtenido a partir de las bases de datos que aporten información hacia la toma de decisiones

(Riquelme et al., 2006)



Introducción a la Minería de datos

La finalidad del KDD es:

- Procesar automáticamente grandes cantidades de datos “brutos” **DATO**
- Identificar los patrones más significativos y relevantes **INFORMACIÓN**
- Presentarlos como conocimiento apropiado para satisfacer las metas del usuario **CONOCIMIENTO**

Fuente de Datos

Pre-procesamiento

Exploración y
transformación

Reconocimiento
de patrones

Evaluación e
interpretación

Fuente de
datos
Bases de
datos

Muestreo
Selección

Limpieza de
datos

- Datos que no existen
- Datos no clasificados
- Identificación de extremos

Transformación
de datos

- Reducir variables / dimensionalidad
- Creación de características
- Generación de variables nuevas

Modelado

- Clasificación
- Regresión
- Agrupamiento
- Asociación

Informes

PREPARACIÓN DE LOS DATOS

EXPLOTACIÓN

EVALUACIÓN

PROCESO KDD y de MINERÍA DE DATOS



Fuente de Datos

Pre-procesamiento

Exploración y
transformación

Reconocimiento
de patrones

Evaluación e
interpretación

CLASIFICACIÓN

La finalidad es crear un modelo para poder predecir la pertenencia a un determinado grupo

EJEMPLO:

Diagnosticar alumnos de éxito académico

PRUEBAS:

Árboles de decisiones, análisis discriminantes, etc..

Modelado

Informes

•Clasificación /
asociación

•Regresión

•Agrupamiento

EXPLORACIÓN

PROCESO KDD y de MINERÍA DE DATOS

Fuente de Datos

Pre-procesamiento

Exploración y
transformación

Reconocimiento
de patrones

Evaluación e
interpretación

REGRESIÓN

La finalidad es crear un modelo para poder predecir el valor de una variable dependiente a partir de otras independientes

EJEMPLO:

Estimar el rendimiento académico del primer año de carrera

PRUEBAS:

Regresión lineal, redes neuronales, regresión logística, etc...

Modelado

Informes

- Clasificación / asociación

- **Regresión**

- Agrupamiento

EXPLOTACIÓN

PROCESO KDD y de MINERÍA DE DATOS



Fuente de Datos

Pre-procesamiento

Exploración y
transformación

Reconocimiento
de patrones

Evaluación e
interpretación

AGRUPAMIENTO / SEGMENTACIÓN

La finalidad es crear un modelo para poder agrupar con características similares

EJEMPLO:

Identificar perfiles de alumnos

PRUEBAS:

K-medias, Bietápico, etc..

Modelado

Informes

•Clasificación /
asociación

•Regresión

•Agrupamiento

EXPLOTACIÓN

PROCESO KDD y de MINERÍA DE DATOS



Fuente de Datos

Pre-procesamiento

Exploración y
transformación

Reconocimiento
de patrones

Evaluación e
interpretación

LOS DATOS DEBEN SER (Viera et al., 2009)

Precisión – sin errores de medición

Consistencia – datos coherentes

Completos – sin falta de atributos

Relevancia – Concernientes al problema

No redundancia – Sin duplicar la misma información



PROCESO KDD y de MINERÍA DE DATOS



Introducción a la Minería de datos

Algunas aplicaciones (Riquelme, 2006)

Comercio y banca

Segmentación de clientes, previsión de ventas, análisis de riesgos

Medicina y farmacia

Diagnóstico de enfermedades y la efectividad de los tratamientos

Seguridad y detección de fraude

Reconocimiento facial, acceso a redes no permitidas,...

Astronomía

Identificación de nuevas estrellas y galaxias

Geología, minería, agricultura y pesca

Identificación de áreas de uso para distintos cultivos o pesca, explotación minera en base de datos de imágenes de satélites

Ciencias ambientales

Identificación de modelos de funcionamiento de ecosistemas naturales o artificiales

Ciencias sociales

Estudio de los flujos de opinión, identificar barrios con conflicto en función de valores socio-demográficos



Introducción a la Minería de datos

Algunas aplicaciones en Educación MDE

En el ámbito educativo la aplicación de la minería de datos como técnica de análisis **se ubica en el entorno del sistema educativo y en concreto en Educación superior**. Las base de datos que se utilizan en los sistemas educativos permiten disponer de una gran cantidad de información, tanto de los estudiantes, trabajadores, departamentos, universidades, etc...., por ejemplo la base de datos UNEIX

- La MDE tiene como objetivo obtener una mejor comprensión del proceso de aprendizaje de los estudiantes y de su participación global en el proceso, orientado a la mejora de la calidad y rentabilidad del sistema educativo

(Winters, T, 2006)



Introducción a la Minería de datos

Algunas aplicaciones en Educación MDE

- R. Alcover, J. Benlloch, P. Blesa, M. A. Calduch1, M. Celma, C. Ferri, J. Hernández-Orallo, L. Iniesta, J. Más, M. J. Ramírez-Quintana, A. Robles, J. M. Valiente, M. J. Vicent, L. R. Zúnica. (2007) **Análisis del rendimiento académico en los estudios de informática de la Universidad Politécnica de Valencia aplicando técnicas de minería de datos.** XIII Jornadas de enseñanza universitaria de la informática. Teruel. España (disponible internet)
- Quiroga, E. (2008) **Minería de datos en educación superior aplicada a un modelo de alerta académica.** Chile
- Valero, S. (2009) **Aplicación de la minería de datos para predecir la deserción.** Universidad tecnológica de Izúcar de Matamoros

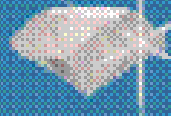
Algunos ejemplos visuales



EJEMPLOS

Cerón, M.A. y Gómez, H. (2010) Minería de datos

(<http://www.slideshare.net/04071977/mineria-de-datos>)



Análisis de Canasta

Ejemplo

Un ejemplo tradicional de minería de datos es el relacionado con una búsqueda en una bodega de datos, de un negocio de cadena, de hechos comunes y relevantes: Luego del proceso se dio como resultado la siguiente:

Si edad < 35;
y sexo = masculino;
y día = jueves
entonces compras incluyen
pañales;
y cerveza

Esto sirvió para que empresa tomara medidas relacionada con la ubicación de ciertos productos en sitios comunes.



EJEMPLOS

Cerón, M.A. y Gómez, H. (2010) Minería de datos

(<http://www.slideshare.net/04071977/mineria-de-datos>)



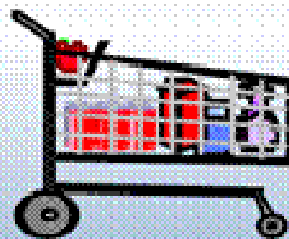
Análisis de Canasta (Market Basket Analysis)

Los hábitos de compra de los clientes pueden ser representados a través de asociaciones o correlaciones entre los diferentes productos que compran en sus "canastas".



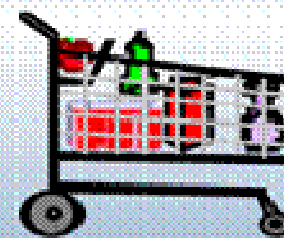
Cliente 1:

Arroz, puré, bebida



Cliente 2:

Arroz, helado,
pan



Cliente 1:

Arroz, bebida,
cerveza



EJEMPLOS

Cerón, M.A. y Gómez, H. (2010) Minería de datos

(<http://www.slideshare.net/04071977/mineria-de-datos>)

Ejemplo

Gestión de personal de una empresa: ¿Qué clases de empleados hay contratados?

Datos:

Id	Salary	Married	Car	Children	Rent/Owner	Union	Off sick/year	Work years	Gender
1	10000	yes	no	0	Rent	no	7	15	M
2	20000	no	yes	1	Rent	yes	3	3	F
3	15000	yes	yes	2	Owner	yes	5	10	M
4	30000	yes	yes	1	Rent	no	15	7	F
5	10000	yes	yes	0	Owner	yes	1	6	M
6	40000	no	yes	0	Rent	yes	3	16	F
7	25000	no	no	0	Rent	yes	0	8	M
8	20000	no	yes	0	Owner	yes	2	6	F
15	8000	no	yes	0	Rent	no	3	2	M
...



Modelo generado:

Minería de datos

Grupo 1: Sin niños y en una casa alquilada. Bajo número de uniones. Muchos días enfermos

Grupo 2: Sin niños y con coche. Alto número de uniones. Pocos días enfermos. Más mujeres y en una casa alquilada

Grupo 3: Con niños, casados y con coche. Más hombres y normalmente propietarios de casa. Bajo número de uniones





EJEMPLOS

Cerón, M.A. y Gómez, H. (2010) Minería de datos

(<http://www.slideshare.net/04071977/mineria-de-datos>)



Ejemplo

Tienda de TV: ¿Cuántas televisiones planas se venderán el próximo mes?

Datos:

PRODUCT	Month-12	...	Month-4	Month-3	Month-2	Month-1	Month
Flat TV 30"	20	...	52	14	139	74	?
Video-dvd-recorder	11	...	43	32	26	59	?
Discman	50	...	61	14	5	28	?
Five star fridge	3	...	21	27	1	49	?
Three star fridge	14	...	27	2	25	12	?
...

Modelo generado:

Minería de datos

Modelo lineal: número de televisiones para el próximo mes

$$V(\text{month})_{flatTV} = 0.62 V(\text{Month-1})_{flat-TV} + 0.33 V(\text{Month-2})_{flat-TV} + 0.12 V(\text{Month-1})_{DVD-Recorder} - 0.05$$



Introducción a la Minería de datos

PROGRAMAS INFORMÁTICOS

En esta última década han aparecido una serie de programas informáticos que nos han permitido analizar un gran volumen de datos

Sus diferencias radican en la presentación e implementación. Pasan por las mismas etapas: colecta de datos, depuración y análisis

SAS System

The screenshot displays the SAS Enterprise Miner interface. The main window title is "SAS - [SAS Enterprise Miner - mheyten [eCRM Fraud Demo]]". The menu bar includes File, Edit, View, Options, Actions, Window, and Help. The Explorer pane on the left shows the "Contents of 'SAS Environment'" with categories like Libraries, File Shortcuts, Favorite Folders, and My Computer. The central pane shows a hierarchical tree structure under "Sample", including "Input Data Source", "Sampling", "Data Partition", "Explore", "Modify", and "Model". The right pane displays a workflow diagram with the following steps: "FRDDemo.CASESET2" (Input Data Source) flows to "Data Partition", which then branches to "Tree" and "Regression". "Tree" flows to "Assessment". The bottom status bar shows "Item(s) deleted from diagram." and the file path "C:\Program Files\SAS\SAS System\9".

```
graph LR; A[FRDDemo.CASESET2] --> B[Data Partition]; B --> C[Tree]; B --> D[Regression]; C --> E[Assessment]
```

Diagram illustrating a workflow in SAS Enterprise Miner:

- Input Data Source: FRDDemo.CASESET2
- Process: Data Partition
- Process: Tree
- Process: Regression
- Process: Assessment

Workflow flow: FRDDemo.CASESET2 → Data Partition → Tree → Assessment; Data Partition → Regression.

SQL Server Data Mining

The screenshot displays the Microsoft Visual Studio interface for an Analysis Services project. The main window shows a cluster diagram with seven clusters (Cluster 1 through Cluster 7) connected by lines. The clusters are arranged in a roughly circular pattern, with Cluster 6 at the top, Cluster 4 on the left, Cluster 5 at the bottom left, Cluster 2 at the bottom right, Cluster 3 in the middle right, Cluster 7 in the middle, and Cluster 1 on the right. The diagram is titled "Sales Customer Territory Cus" and is viewed using the "Microsoft Cluster Viewer". The "Shading Variable" is set to "Population" and the "Density" is "None" with a "23%" indicator. The "State" is set to "None". The interface includes a menu bar (File, Edit, View, Project, Build, Debug, Database, Mining Model, Tools, Window, Community, Help), a toolbar, and a sidebar with "Toolbox", "Server Explorer", "Properties", "Deployment Progress", and "Solution Explorer". The status bar at the bottom shows "Deploy succeeded" and the system tray with the date "Wednesday" and time "3:39 PM".

ORACLE DATA MINING

Oracle Data Miner - Table: CD_BUYERS

Result Viewer: "DM4JSCD_BUYER19890_TM"

Result Viewer: CD_BUYERS20881_DT

Histogram for selected attribute

The screenshot displays the Oracle Data Miner interface with four main windows:

- Table: CD_BUYERS:** Shows a data table with columns: CUST_ID, CD_BUYER, AGE, MARITA, ANNUAL_INCOME. The table contains 20 rows of data.
- ROC Curve:** A Receiver Operating Characteristic (ROC) curve plot showing True Positive Rate vs. False Positive Rate. The Area Under the Curve (AUC) is 0.974251.
- Decision Tree:** A tree structure showing nodes and their associated rules. The root node is Node 0 (True). The tree splits based on various attributes like RELATIONSHIP, PAYROLL_DEDUCTION, AVE_CHECKING_BALANCE, CAPITAL_GAIN, and OCCUPATION.
- Histogram for AGE:** A histogram showing the distribution of the AGE attribute. The x-axis is Bin Count (0-700) and the y-axis is Bin Height. The distribution is roughly bell-shaped, centered around 38-40.

Confusion Matrix (Top Right):

	Others	1
Others	816	87
1	107	186

Confusion Matrix (Bottom Right):

	Others	1
Others	0	1
1	112500	0

Statistics (Bottom Right):

- Sample count: 3000
- Minimum value: 17
- Maximum value: 90
- Average value: 38.5
- Variance: 166.88
- Skewness: 0.61
- Kurtosis: -0.04

Decision Tree Rules (Middle):

- Node 0: True
- Node 1: RELATIONSHIP is in { Husband }
- Node 2: PAYROLL_DEDUCTION <= 97.5
- Node 3: PAYROLL_DEDUCTION > 97.5
- Node 4: AVE_CHECKING_BALANCE <= 0
- Node 5: AVE_CHECKING_BALANCE > 0
- Node 6: CAPITAL_GAIN <= 5715.5
- Node 7: CAPITAL_GAIN > 5715.5
- Node 8: OCCUPATION is in { Cleric, Cr. }
- Node 9: CAPITAL_GAIN <= 5483.0
- Node 10: OCCUPATION is in { Armed-F. En. }

Clementine





Introducción a la Minería de datos

Text Mining

Si bien es cierto que existe una gran cantidad de información almacenada en bases de datos, **la existencia de un gran volumen de documentos hace necesario aplicar algún tipo de sistema de análisis.**

El análisis presenta un mayor nivel de complejidad y de dimensiones en cuanto a la categorización de texto y procesamiento de lenguaje natural. La minería de texto o text mining permite la extracción y recuperación de la información

PASW Text Mining (2010)

Internet Web Mining

Otra de las aplicaciones de la minería de datos consiste en aplicar sus técnicas a documentos y servicios Web (minería de Web) en concreto, **el análisis de datos por Internet y on line.**

¿cuáles son las páginas web más visitadas?



Bibliografía

- Aluja, T (2001) *La minería de datos, entre la estadística y la inteligencia artificial*. **QÜESTIÓ**, vol 25,3, p 479-498
- Han, J. y Kamber, M. (2006) **Data mining, concepts and techniques**. USA
- Hernandez Orallo J.L. (2004) **Introducción a la minería de datos**. New York: Pearson Prentice Hall
- Pérez, C. Santín, D. (2007) **Minería de datos: técnicas y herramientas**. Madrid: Paraninfo
- Riquelme, J.C.; Ruiz, R y Gilbert, K. (2006) *Minería de datos: conceptos y tendencias*. **Revista Iberoamericana de Inteligencia artificial**, 29, pp 11-18
- Vieira Braga, L.P.; Ortiz Valencia, L.I.; Ramírez Carvajal, S.S. (2009) **Introducción a La Minería de Datos**. Rio de Janeiro: E-papers servicios editoriales
- Winters, T (2006) Educational Data Mining: Collection and Analysis of Score Matrices for Outcomes- Based Assessment .USA, University of California: Riverside



¿PASAMOS A LA SEGUNDA PARTE ? O ¿HACEMOS DESCANSO?