



Data mining

Minería de datos

Autor: Alfredo Daza Vergaray

© Derechos de autor registrados:

Empresa Editora Macro EIRL

© Derechos de edición, arte gráfico y diagramación reservados:

Empresa Editora Macro EIRL

Coordinación de edición:

Magaly Ramon Quiroz

Diseño de portada:

Rudy Herrera Torres

Corrección de estilo:

Martín Vargas Canchanya

Diagramación:

Julissa Ventocilla Fernández

Edición a cargo de:

© Empresa Editora Macro EIRL

Av. Paseo de la República N.° 5613, Miraflores, Lima, Perú

☎ Teléfono: (511) 748 0560

✉ E-mail: proyectoeditorial@editorialmacro.com

🌐 Página web: www.editorialmacro.com

Primera edición: julio 2016

Tiraje: 1200 ejemplares

Impresión

Talleres gráficos de la Empresa Editora Macro EIRL

Jr. San Agustín N.° 612-624, Surquillo, Lima, Perú

ISBN N.° 978-612-304-417-6

Hecho el depósito legal en la Biblioteca Nacional del Perú N.° 2016-08276

Prohibida la reproducción parcial o total, por cualquier medio o método, de este libro sin previa autorización de la Empresa Editora Macro EIRL.

Índice

Introducción	11
CAPÍTULO 1: Conceptos básicos de minería de datos	13
1.1 Minería de datos	15
1.2 Procesos de minería de datos (KDD).....	16
1.3 Metodología CRISP	17
1.4 Modelo	20
1.5 Modelo híbrido	20
1.6 Predicción	21
1.7 Almacén de datos (<i>data warehouse</i>).....	21
Resumen	22
CAPÍTULO 2: Técnicas y aplicación de la minería de datos.....	25
2.1 Modelos de minería de datos	27
2.2 Métodos de minería de datos	27
2.2.1 Árboles de clasificación.....	28
2.2.2 Redes neuronales.....	33
2.3 Aplicación de la minería de datos.....	37
2.3.1 Minería de datos en la educación	39
Resumen	57
CAPÍTULO 3: Presentación general de SPSS Clementine.....	61
3.1 SPSS Clementine	63
3.1.1 Sector público	63
3.1.2 CRM.....	64
3.1.3 <i>Web mining</i>	64
3.1.4 Desarrollo de fármacos	65
Resumen	66
CAPÍTULO 4: Interfaz y categorías de SPSS Clementine	67
4.1 Elementos de la interfaz de SPSS Clementine	69
4.1.1 Clementine Stream Canvas	69
4.1.2 Nodos Palette	70
4.1.3 Clementine Managers	70
4.1.4 Clementine Projects	72

4.2 Categorías de SPSS Clementine	73
4.2.1 Categoría Source	74
4.2.2 Categoría Record Ops	75
4.2.3 Categoría Field Ops	76
4.2.4 Categoría Output	76
4.2.5 Categoría Graphs	77
4.2.6 Categoría Modeling.....	78
4.2.7 Categoría Export	79
Resumen	80
CAPÍTULO 5: Instalación de SPSS Clementine	81
5.1 Instalación del programa SPSS Clementine	83
5.1.1 Pasos para la instalación del programa SPSS Clementine	83
Resumen	92
CAPÍTULO 6: Aplicaciones con diferentes técnicas de minería de datos.....	93
6.1 Caso n.º 1: Predicción de juego de tenis (árboles de decisión).....	95
6.2 Caso n.º 2: Predicción de planta iris.....	109
6.3 Caso n.º 3: Predicción de fármacos.....	122
6.4 Caso n.º 4: Problemas de clúster (caso empleados Memolum Web).....	136
6.5 Caso n.º 5: Agrupamientos en relación a las ventas	141
6.6 Caso n.º 6: Datos erróneos y faltantes (caso empleados Memolum Web).....	147
6.7 Caso n.º 7: Obtener y transformar datos a través de ODBC (conexión de base de datos abierta)	166
6.8 Caso n.º 8: Catalog_forecast (series de tiempo).....	176
6.9 Caso n.º 9: <i>Computer hardware data set</i>	181
6.10 Caso n.º 10: Detección de fraude.....	187
6.11 Caso n.º 11: Validación del modelo Drug con datos nuevos	195
6.12 Caso n.º 12: Integración y partición de datos	200
6.13 Caso n.º 13: Columna vertebral (partición de datos).....	210
6.14 Caso n.º 14: Validación cruzada.....	220
6.15 Caso n.º 15: Trabajar con pocos registros	224
6.16 Caso n.º 16: Reglas de asociación y dependencia	233
6.17 Caso n.º 17: Regresión logística (telecomunicaciones <i>churn</i>)	243
6.18 Caso n.º 18: Predicción secuencial	254
6.19 Caso n.º 19: Exportación de modelos y resultados.....	261
6.20 Caso n.º 20: Series de tiempo (pronosticar).....	267
Resumen	279
BIBLIOGRAFÍA	282

1.1 Minería de datos

La minería de datos se ha definido de diferentes maneras. A continuación, se mencionarán algunos de estos conceptos para un mejor entendimiento, en especial, por aquellas personas que recién se estén iniciando en el maravilloso mundo de la extracción del conocimiento:

A. Primera definición

La minería de datos se define como aquel proceso que consiste en extraer conocimiento útil y comprensible, previamente desconocido, desde grandes cantidades de datos almacenados en distintos formatos. En ese sentido, la tarea fundamental de la minería de datos es encontrar modelos inteligibles a partir de los datos recogidos (Hernández *et al.*, 2004).

B. Segunda definición

Según Hernández *et al.* (2004), la minería de datos implica un proceso de descubrimiento de nuevas y significativas relaciones, patrones y tendencias que son observadas al examinar grandes cantidades de información.

C. Tercera definición

Según Han y Kamber (2006) la minería de datos se refiere a la extracción de conocimiento o «minería» de grandes cantidades de datos. Sin embargo, de acuerdo con su perspectiva el nombre sería, en realidad, inapropiado, ya que, por ejemplo, la actividad minera que extrae oro de las rocas o de la arena se conoce como «minería de oro» en vez de «minería de roca» o «de extracción de arena». Por ende, partiendo de una lógica similar, la minería de datos debería haber recibido el nombre más apropiado de «minería de datos del conocimiento», el cual, por desgracia, es un poco largo. Ahora bien, sucede que «minería» es una palabra que porta la idea de un proceso por el cual se extrae un pequeño conjunto de elementos (pepitas) poseedores de una cierta cantidad de materia prima (metales preciosos). Así, a pesar de ser un nombre poco apropiado, al vincular las ideas de datos y extracción, «minería de datos» se ha convertido en una opción más popular. Frente a esto, lo único que cabe advertir es que existen muchas otras expresiones similares a esta, las cuales, empero, tienen un diferente matiz de sentido, tales como la minería de datos de conocimiento, la extracción de conocimientos, análisis de datos, análisis de patrones, arqueología de datos y filtración de información.

D. Cuarta definición

Según González (2005), la minería de datos es el proceso por el cual se genera un modelo útil para la predicción. Dicho modelo se construye teniendo como fundamento los datos que se encuentran en una base de datos, a los cuales se le ha aplicado algún algoritmo justamente con el fin de plantear un modelo.

En conclusión, se podría decir que la minería de datos es un proceso que integra los datos de diferentes fuentes (SQL Server, Oracle, Excel, etc.) para, posteriormente, extraer un importante conocimiento, es decir, identificar información trascendente, valiosa y útil, a partir de lo cual las instituciones van a poder tomar alguna significativa decisión.

1.2 Procesos de minería de datos (KDD)

Las etapas para la realización de la minería de datos siempre son las mismas independientemente de la técnica específica a usar. El conjunto de las partes de este proceso, también conocido por sus siglas en inglés «KDD» (*knowledge discovery in databases* o «descubrimiento de conocimiento en bases de datos»), es descrito de la siguiente manera por Fayyad *et al.* (1996).

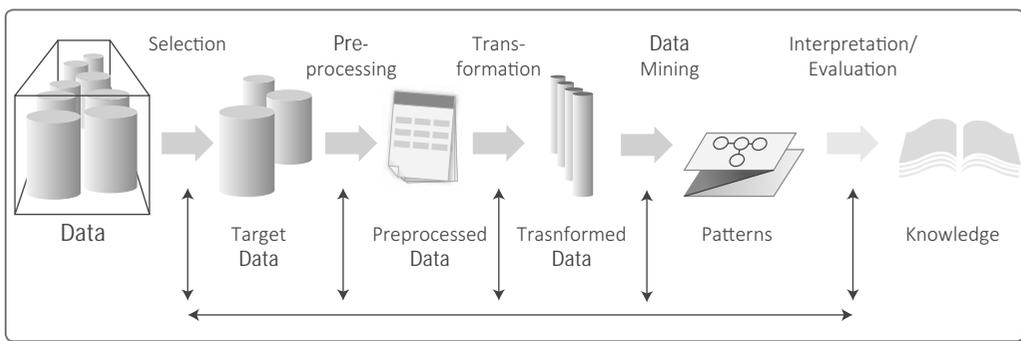


Figura 1.1 Descripción general de los pasos que constituyen el proceso KDD

Fuente: Fayyad *et al.* (1996).

El proceso de KDD es interactivo e iterativo (con muchas decisiones tomadas por el usuario) e implica numerosos pasos que se resumen así:

1. Aprendizaje del dominio de la aplicación: este paso incluye la adquisición del conocimiento previo relevante y el planteo de los objetivos de la aplicación.
2. Creación de un conjunto de datos de destino: por medio de esto se escoge el conjunto de datos o se elige el subconjunto de variables o muestras de datos en los cuales el descubrimiento se va a realizar.
3. Limpieza de datos y preprocesamiento: aquí se dan las operaciones básicas como la eliminación de ruido, la recogida de la información necesaria para modelar, la determinación de estrategias para el manejo de los campos de datos que faltan, la contabilidad de la información en tiempo y secuencia de los cambios conocidos, la decisión en torno al uso de DBMS (tales como tipos de datos y esquemas) y la asignación de valores desaparecidos y desconocidos.
4. Reducción de datos y de proyección: en este punto se lleva a cabo la búsqueda de las características útiles para representar los datos, dependiendo del objetivo y de la tarea. Para ello, se ejecuta la reducción de dimensionalidad o se aplican métodos de transformación que permitan acortar el número efectivo de las variables en estudio o encontrar las representaciones invariables para los datos.