

# Predicción de abandonos en cursos online masivos y abiertos

**Manuel Alejandro Bacallado López**

Trabajo Final de Grado

II Congreso de Estudiantes de Informática de la Universidad de  
La Laguna

Escuela Superior de Ingeniería y Tecnología

30 de Noviembre de 2016

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

# Introducción

- El proyecto planteado está basado en un problema real:
  - Los cursos online masivos y abiertos se realizan por todo el mundo.
  - Muchos alumnos se inscribirán a estos cursos debido a su bajo coste y gran cantidad de materiales disponibles, pero no llegarán a terminarlos.
- Se analizarán los datos del problema para realizar una clasificación de los alumnos matriculados en los cursos, prestando atención a su abandono o finalización con éxito.
- La predicción se llevará a cabo mediante técnicas de minería de datos, ya que este proceso manual es lento y subjetivo.

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

## Antecedentes y estado actual del tema

- Los datos provienen del desafío KDD Cup 2015 (Predicting dropouts in MOOC)
- KDD Cup es una competición a nivel mundial de minería de datos y descubrimiento de conocimiento.
- Organizado por la ACM (Association of Computing Machinery).
- Vigencia desde 1997-Actualidad.
- Cada año se realiza el mismo proceso:
  - Habilitan los datasets con la información a tratar.
  - Instrucciones a seguir (Plazos de entrega, objetivos a cumplir).
  - Premio en \$ para los primeros puestos.

## Antecedentes y estado actual del tema

- Los datos se estructuran en las siguientes tablas:
  - date.csv
  - object.csv
  - enrollment\_train.csv - enrollment\_test.csv
  - log\_train.csv - log\_test.csv
  - truth\_train.csv

Introducción  
Antecedentes y estado actual del tema  
**Objetivos**  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos**
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía



## Objetivos

- Los objetivos de este proyecto se constituyen en:
  - Realizar una clasificación sobre los alumnos matriculados en los cursos online masivos y abiertos.
  - Utilizar herramientas de software libre para:
    - Crear una aplicación en Java que use los operadores internos de la aplicación RapidMiner Studio 7.0.
    - Almacenamiento de los datos iniciales en un base de datos y creación de tablas con conocimiento nuevo generado(Ingeniería de características).
    - Técnicas de ingeniería del software para un desarrollo profesional de la aplicación.

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
**Fundamentos de minería de datos**  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos**
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

# Fundamentos de minería de datos

- Se define como el proceso de extraer conocimiento útil y comprensible desde grandes cantidades de datos almacenados en distintos formatos.



# Fundamentos de minería de datos

- Tareas
  - Predictivas: Predecir uno o más valores para uno o más casos. Los casos van acompañados de una etiqueta(clase, categoría o valor numérico).
    - Clasificación: Cada registro en la base de datos pertenece a una clase, la cual se establece mediante el valor de un atributo denominado clase de la instancia.
  - Descriptivas: Los casos constituyen un conjunto sin etiquetas. El objetivo es describir los datos existentes.
- Técnicas
  - Árboles de decisión
  - Casos y vecindad(K-nn)
  - Técnicas probabilistas(Naive Bayes)

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
**Herramientas utilizadas**  
Aplicación software  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas**
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

## Herramientas utilizadas

- SGBD Maria DB
- Apache Maven
- Biblioteca de clases de RapidMiner Studio 7.0
- Java
- Eclipse
- Git
- Github
- Doxygen
- JUnit

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
**Aplicación software**  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software**
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

## Aplicación software

- La aplicación software ha sido realizada utilizando:
  - Java y Swing.
  - Los operadores internos de RapidMiner Studio 7.0.
  - Patrones de diseño:
    - Patrón Observador(Observer Pattern).
    - Patrón Estrategia(Strategy Pattern).
    - Patrón Modelo-Vista-Controlador(MVC Pattern).



Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
**Aplicación software**  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

## Aplicación software

- La aplicación software presenta la siguiente estructura:

The screenshot displays the 'Data Mining' application window. At the top, there are three buttons: 'Start Process', 'Reset Process', and 'Stop Process'. Below these, the interface is divided into three main sections for configuring operators:

- Training set operators:** Includes a 'Clear Training' button and a 'Training' button. The operator list contains: Select Attributes, Normalize, Numerical to Binominal, Set Role, and Naive Bayes.
- Test set operators:** Includes a 'Clear Test' button and a 'Test' button. The operator list contains: Read CSV, Select Attributes, Normalize, Numerical to Binominal, and Set Role.
- Apply set operators:** Includes a 'Clear Apply' button and an 'Apply' button. The operator list contains: Apply Model, Performance, Performance to Data, and Write CSV.

At the bottom of the window, there is a progress bar with six steps: Step 1: Data Access, Step 2: Blending, Step 3: Cleansing, Step 4: Modeling, Step 5: Scoring, and Step 6: Validation. Each step is represented by a small square icon.

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
**Caso de estudio**  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio**
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

## Caso de estudio

- Previamente a la generación de nuevas características se han realizado varias tareas:
  - Se han renombrado las tablas iniciales al castellano para una mejor comprensión.
  - Se eliminó la duplicidad de los registros en las tablas.
  - Se eliminaron atributos en las tablas que no aportaban información relevante.
  - Se analizó la tabla resultados\_train.csv(truth\_train.csv) para comprobar el número de abandonos sin realizar minería de datos, teniendo el valor 1 si abandona y valor 0 si continúa:
    - El número total de registros es de: 120.542.
    - El número de alumnos que abandonan es de: 95.581.
    - El número de alumnos que no abandonan es de: 24.691.

## Caso de estudio

- La generación de nuevo conocimiento se ha producido mediante una batería de cuestiones. A continuación se expondrán las más relevantes:
  - ¿Cuántos alumnos hay por curso?
  - ¿Cuántos cursos tiene un alumno?
  - ¿Cuántos alumnos de un curso aprobaron?
  - ¿Cuántos cursos aprobó un alumno?
  - ¿Número de días en los que un alumno accede a un curso?
  - ¿Número de cursos simultáneos que tiene un alumno en un determinado mes?
  - ¿Número de días entre el primero y último acceso por inscripción?

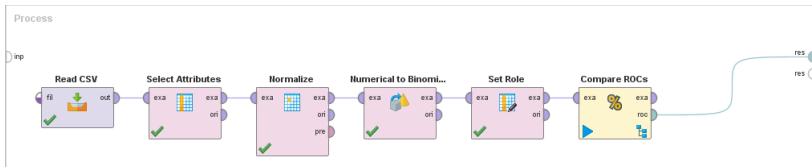
## Caso de estudio

- Los atributos seleccionados para la creación de la vista minable son:

Usuario	Curso_Id	NEDV	Rango	RangoCurso	NMV	NM	NTA	CS	Resultado
1qXC7Fjbp66GPQc6pHLfEu08WKOzxG4	7GRHBDsir1GKRZBTSMEnTYDr2JQm4xx	9	28	29	99	699	6	1	0
1ELMTIXpjncZU4WKxvrVri8AJR2gf	AXUJZGmZ0xaYSWazu8RQ1G5c76ECT1Kd	5	15	29	49	264	6	1	0
0K6JPqivQzicY4EV4nqMLCL3a08A97	DPnLzkJJqOOPRJfBxIHbQEERIYHu5la	11	25	29	77	398	7	1	0
088SASUPOVYUGhoEly8vlnkGvRBJoNwp	7GRHBDsir1GKRZBTSMEnTYDr2JQm4xx	18	29	29	101	699	7	2	0
0XSMd0GWzvML1r2AHrvvzWxbczqfFP	DPnLzkJJqOOPRJfBxIHbQEERIYHu5la	13	22	29	59	398	7	1	1
1h4cVFontLTW8vs6Jg6kDELHITwYJukub	TAYxoh39I2LZnfbpL0Lff2NcxzrCKplox	1	0	29	2	333	1	1	1
0U6Ls9kSIXfs9NGu5KLvjyV4KJOnn33	AXUJZGmZ0xaYSWazu8RQ1G5c76ECT1Kd	8	17	29	66	264	7	1	0
1T8ttUcZNn4nU0rZp09wSozsY7mAk229	DPnLzkJJqOOPRJfBxIHbQEERIYHu5la	8	24	29	39	398	7	1	0
0fEoLC14nqwTVtr21LrdveqXhccEIQPz	TAYxoh39I2LZnfbpL0Lff2NcxzrCKplox	10	19	29	45	333	7	2	1
0SIWaSfGyO3je3Wq0u3MOT4Rr4grcql3	DPnLzkJJqOOPRJfBxIHbQEERIYHu5la	10	16	29	55	398	7	1	0

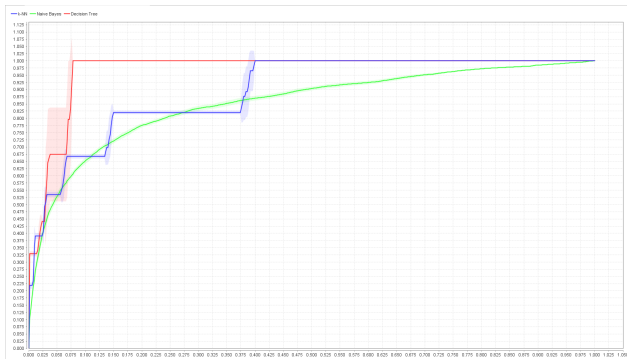
## Caso de estudio

- El proceso resultante en RapidMiner Studio 7.0 para resolver el caso de estudio es el siguiente:



## Caso de estudio

- El resultado utilizando RapidMiner Studio 7.0 es el siguiente:



## Caso de estudio

- Los resultados utilizando los algoritmos en la aplicación software son los siguientes:
  - Árbol de decisión:

Criterion	Value	Standard Deviation	Variance
accuracy	0.854968448789959		

- K-nn:

Criterion	Value	Standard Deviation	Variance
accuracy	0.8950488870397337		

- Naive Bayes:

Criterion	Value	Standard Deviation	Variance
accuracy	0.8474100270438943		



Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
**Conclusiones**  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones**
- 9 Líneas futuras
- 10 Bibliografía

## Conclusiones

- Con este trabajo se ha profundizado en el fantástico mundo de la minería de datos, realizando un caso práctico para los cursos online masivos y abiertos y desarrollando una aplicación en Java utilizando los operadores internos de RapidMiner Studio 7.0.

Ambos son temas de actualidad, los cursos online por la facilidad de inscripción y su bajo coste y la minería de datos por su continuo crecimiento y resolución de problemas enfocados a cualquier ámbito.

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
Conclusiones  
**Líneas futuras**  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras**
- 10 Bibliografía

## Líneas futuras

- La inclusión en la aplicación de más operadores de RapidMiner Studio 7.0 o versiones más actualizadas.
- La visualización de resultados gráficamente, ya sea utilizando las facilidades proporcionadas por las librerías de clases de RapidMiner Studio 7.0 o creándola desde cero o con ayuda de otras herramientas de visualización.
- Plantear la creación de una aplicación web que utilice los operadores de RapidMiner Studio 7.0 o versiones más actualizadas.
- Modificar la interfaz gráfica de usuario para que sea del estilo "Drag and Drop".

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Índice

- 1 Introducción
- 2 Antecedentes y estado actual del tema
- 3 Objetivos
- 4 Fundamentos de minería de datos
- 5 Herramientas utilizadas
- 6 Aplicación software
- 7 Caso de estudio
- 8 Conclusiones
- 9 Líneas futuras
- 10 Bibliografía

## Bibliografía

- SGBD Maria DB. <https://mariadb.org/>
- Apache Maven. <https://maven.apache.org/>
- Biblioteca de clases de RapidMiner Studio 7.0.  
<https://github.com/rapidminer/rapidminer-studio>
- Java. <https://www.java.com/es/>
- Eclipse. <https://eclipse.org/>
- Git. <https://git-scm.com/>
- Github. <https://github.com/>

## Bibliografía

- Javadoc.  
<http://www.oracle.com/technetwork/articles/java/index-jsp-135444.html>
- Doxygen. <http://www.stack.nl/~dimitri/doxygen/>
- JUnit. <http://junit.org/junit4/>
- Observer Pattern. E. Freeman, E. Freeman. Head First Design Pattern, O'Reilly, 2004.
- Strategy Pattern. E. Freeman, E. Freeman. Head First Design Pattern, O'Reilly, 2004.
- MVC Pattern. E. Freeman, E. Freeman. Head First Design Pattern, O'Reilly, 2004.

Introducción  
Antecedentes y estado actual del tema  
Objetivos  
Fundamentos de minería de datos  
Herramientas utilizadas  
Aplicación software  
Caso de estudio  
Conclusiones  
Líneas futuras  
Bibliografía

# Fin

Gracias por su atención.

E-mail: [manuelbacallado89@gmail.com](mailto:manuelbacallado89@gmail.com)