

Ciencia de Datos Educativos

Sebastián Ventura Soto

Grupo de Investigación *Knowledge Discovery
and Intelligent Systems*

Universidad de Córdoba

Contenidos

- **Introducción**
- **Ciencia de Datos Educativos**
 - **Perspectiva Histórica**
 - **Disciplinas asociadas**
 - **Ciclo de vida en CDE**
- **Casos de Éxito**
 - **Predicción del rendimiento académico**
 - **Personalización de la enseñanza**
 - **Formación de grupos y/o equipos de trabajo**
 - **Evaluación por pares y autoevaluación**
 - **Recomendación de itinerarios**
- **Nuevos retos**
 - **Problemas abiertos**
- **Conclusiones**

Introducción



Datos, datos... y más datos

- Vivimos en la era de los datos
- Los datos han cobrado una enorme importancia.
- La explotación de estos datos es aun más importante. Entre otras cosas permite:
 - Descubrir nuevos patrones
 - Comprobar/validar teorías científicas
 - Construir modelos para predecir el comportamiento de un sistema
 - Categorizar objetos
 - ...



Información generada en los Sistemas Educativos

- Instituciones educativas
 - Información de los estudiantes
 - Registros académicos
- Educación presencial
 - Múltiples datos generados por el alumno y su interacción con el profesor
 - Tamaño pequeño
- Sistemas tutores inteligentes
 - Acciones realizadas por los estudiantes
 - Resultados de la evaluación
 - ...
- Otros sistemas on-line (incluidos MOOCs)
 - Interacción con las plataformas de educación:
 - Páginas visitadas
 - Actividades realizadas
 - Comunicación con otros estudiantes
 - Resultados de la evaluación
- Simuladores y/o sistemas basados en gamificación
 - Acciones realizadas en el sistema
 - Progreso del aprendiz
 - Caminos recorridos
 - ...

Extracción de conocimiento en los Sistemas Educativos

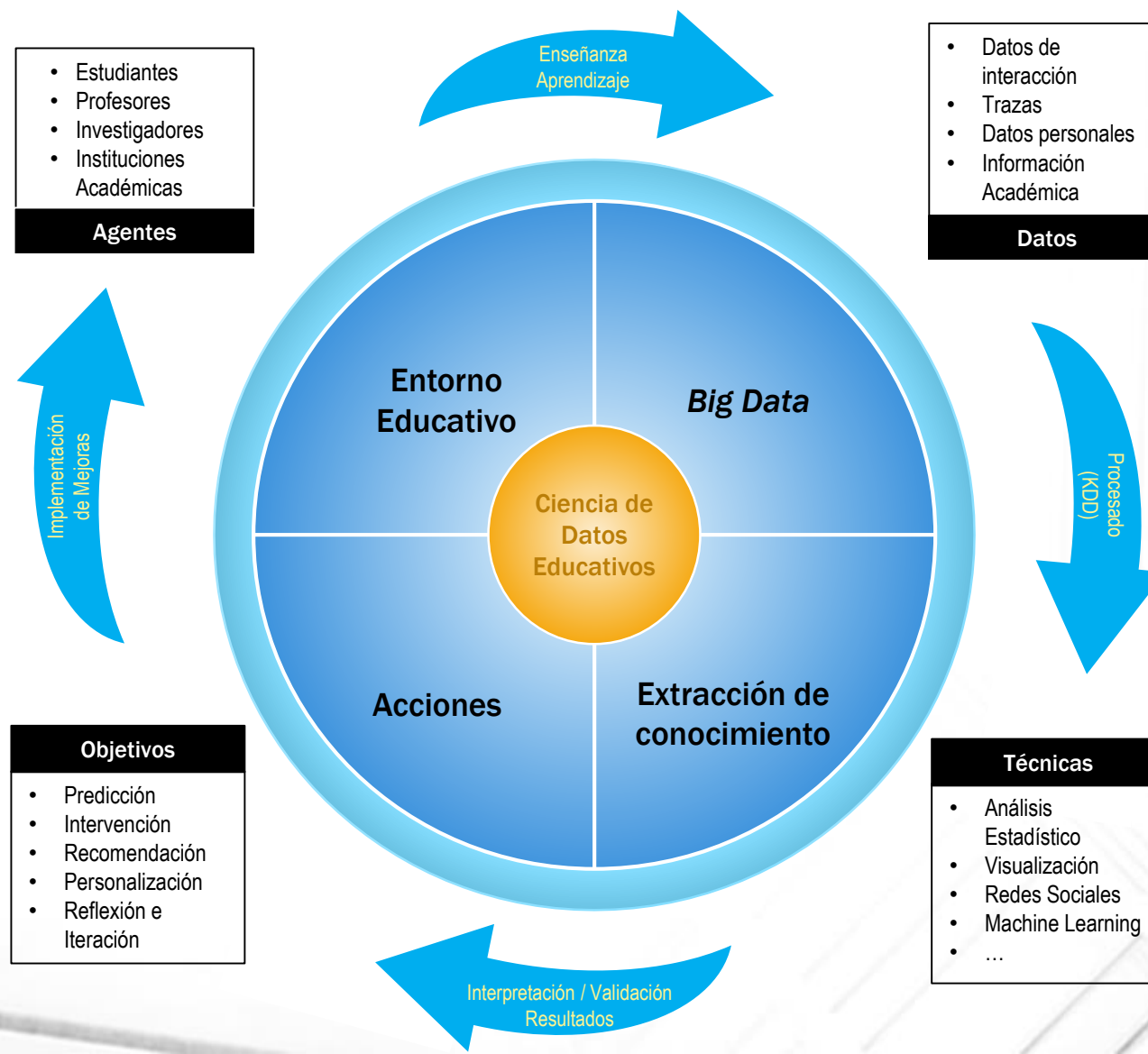
Toda la información generada por los sistemas educativos puede utilizarse en la extracción de nuevo conocimiento útil

El término **Ciencia de Datos Educativos** (*Educational Data Science*) abarca, de forma genérica, a multitud de disciplinas cuyo objetivo es aplicar las distintas técnicas de la Ciencia de Datos sobre la información generada en los sistemas educativos.

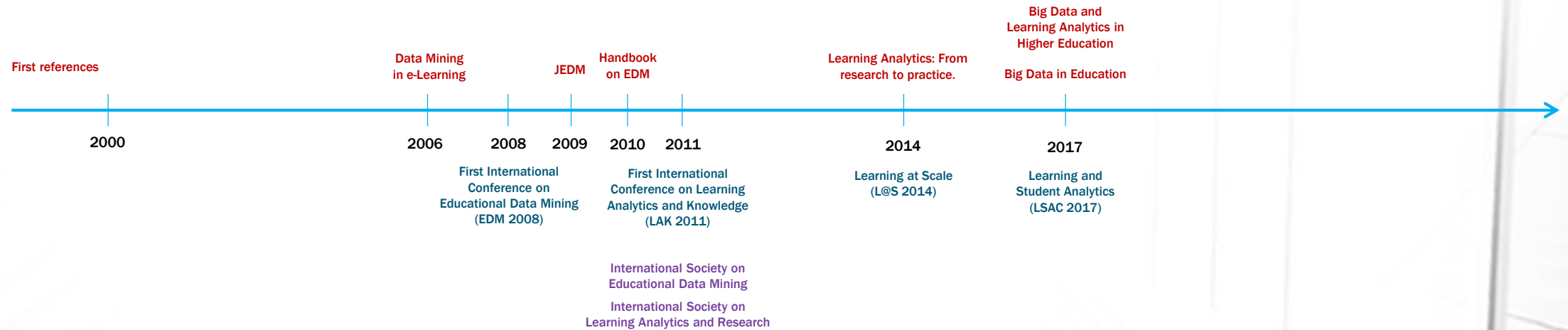
Ciencia de Datos Educativos



Ciclo de Vida en EDS



Ciencia de Datos Educativos. Una perspectiva histórica



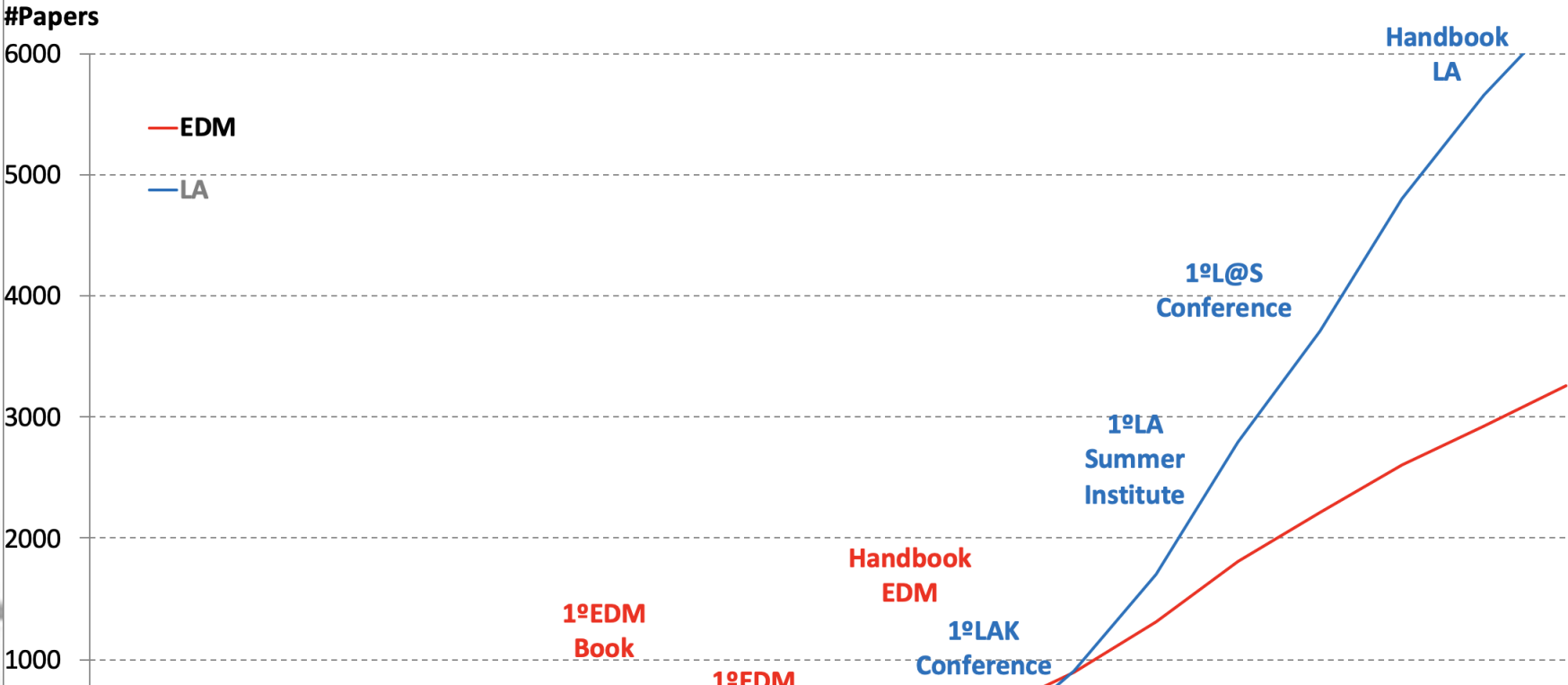
Ciencia de Datos Educativos: Disciplinas... ¿Incluidas?

El término Ciencia de Datos Educativos no está demasiado extendido. De hecho, existen términos mucho más populares que se refieren a la misma o a disciplinas parecidas

- **Minería de Datos Educativos** (*Educational Data Mining, EDM*). Aplicación de las técnicas de extracción de conocimiento a la mejora de los sistemas educativos.
- **Analítica de Aprendizaje** (*Learning Analytics, LA*). Aplicación de técnicas de analítica de datos para la mejora del proceso de enseñanza-aprendizaje.
- **Analítica Académica** (*Academic Analytics, AA*). Aplicación de las técnicas de analítica de datos e inteligencia de negocio a datos de las instituciones académicas.
- **Analítica de la Enseñanza** (*Teaching analytics, TA*). Analisis de datos de enseñanza y rendimiento académico.
- **Big Data en Educación** (*Big Data in Education, BDE*). Aplicación de técnicas de extracción de conocimiento en bases de datos educativas masivas, como las que se generan en los MOOCs



Ciencia de Datos Educativos. Evolución del Número de Publicaciones



Casos de éxito



Predicción del rendimiento académico

Concepto

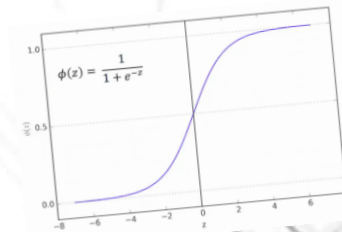
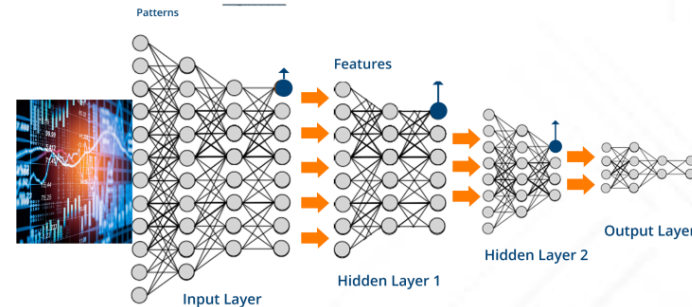
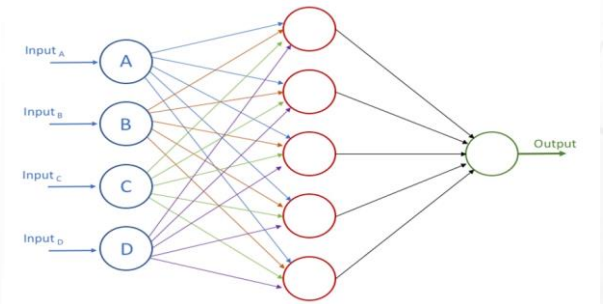
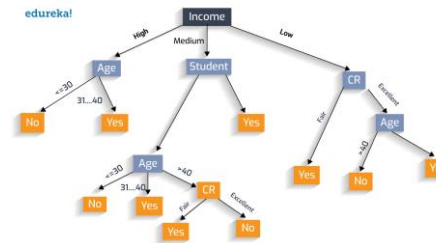
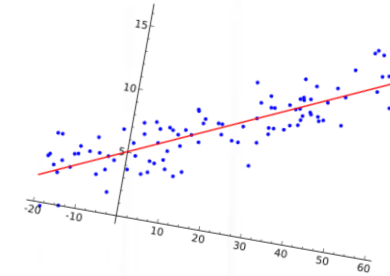
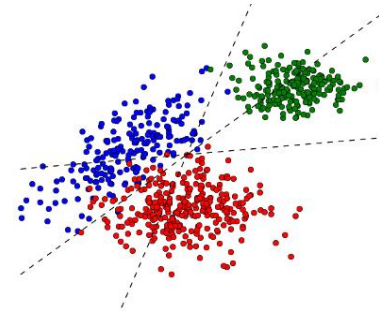
- Es una de las aplicaciones más estudiadas
- Se trata de obtener modelos para predecir el rendimiento del estudiante a partir de datos históricos
- Objetivos:
 - Detectar fracaso y/o abandono escolar
 - Mejor comprensión de ambos fenómenos.
 - Desarrollar medidas correctoras basadas en el conocimiento extraído de los modelos



Predicción del rendimiento académico

Metodologías empleadas

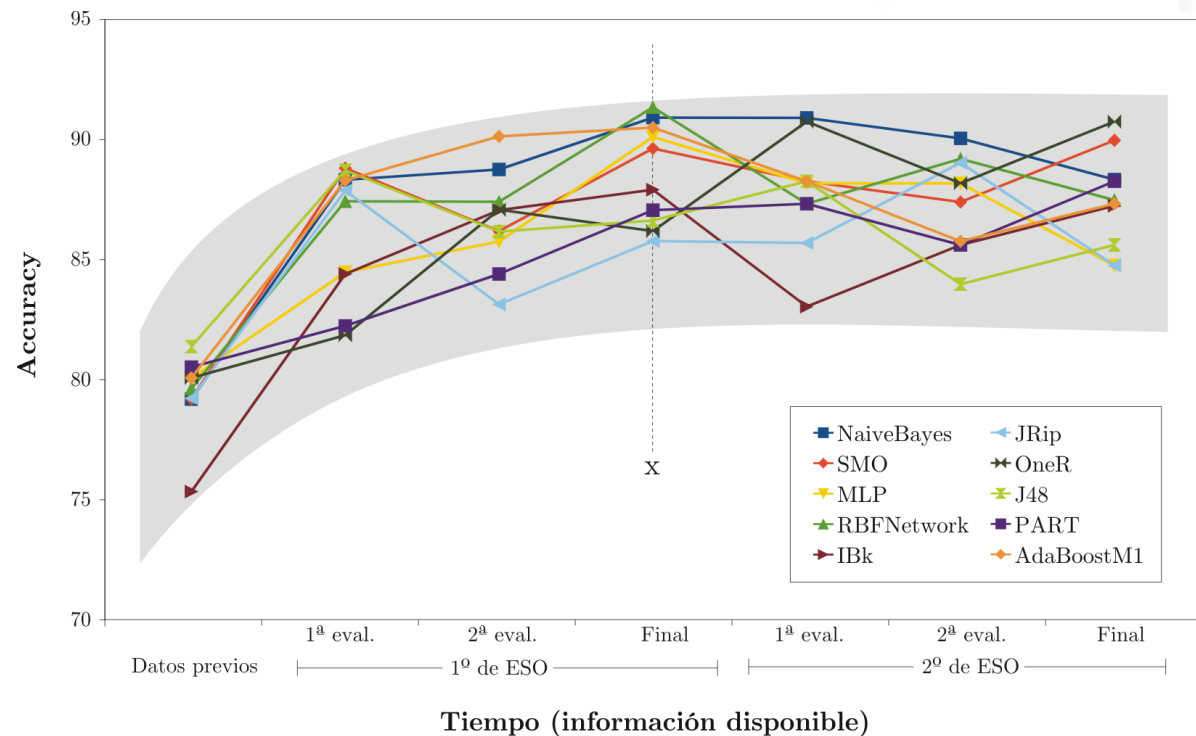
- Clasificación (cualitativa) vs regresión (numéricas).
- Modelos caja blanca (interpretables) vs caja negra (precisos)
- Modelos clásicos (estadística tradicional, *machine learning*) vs actuales (*deep learning*)



Predicción del rendimiento académico

Retos abiertos

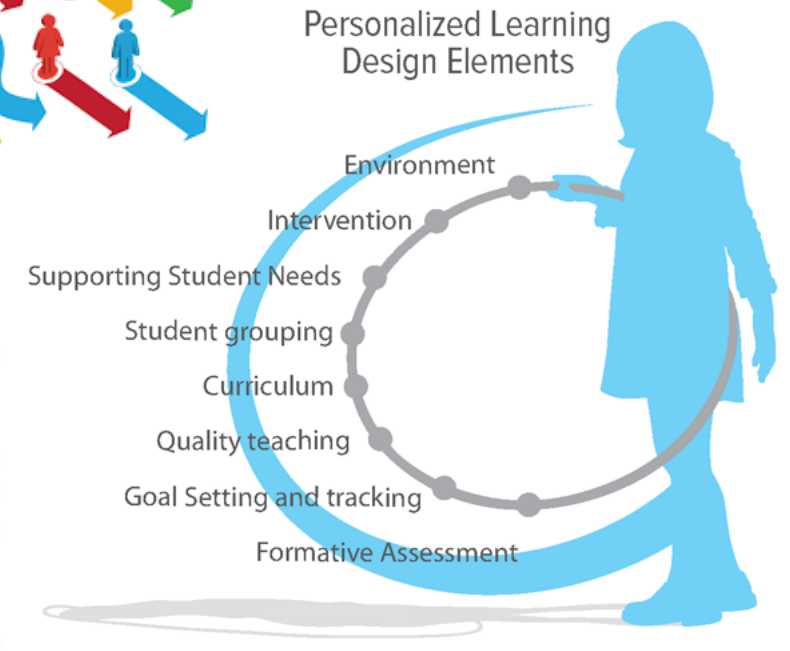
- Ventana temporal
- Datos procedentes de distintas fuentes
 - Datos tabulares vs. relacionales
 - Representaciones de datos flexibles
- Validez de los modelos
 - Caducidad de los modelos?
 - Big Data vs. Small Data



Personalización de la enseñanza

Concepto

- Término muy amplio:
 - Personalización de actividades
 - Personalización de contenidos
 - Personalización de caminos de aprendizaje
- Objetivos:
 - Mejorar la experiencia de los estudiantes en el sistema
 - Mejorar la eficiencia del proceso de enseñanza-aprendizaje:
 - Mejora del rendimiento
 - Metas perseguidas



Personalización de la enseñanza

Metodologías empleadas

Soluciones basadas en los principios de los Sistemas de Recomendación:

- Filtrado colaborativo:
 - Agrupamiento (*clustering*) basado en características
 - Análisis de redes sociales
- Modelos basados en contenido:
 - Modelos de clasificación convencionales
 - Análisis de sentimientos
- Sistemas híbridos



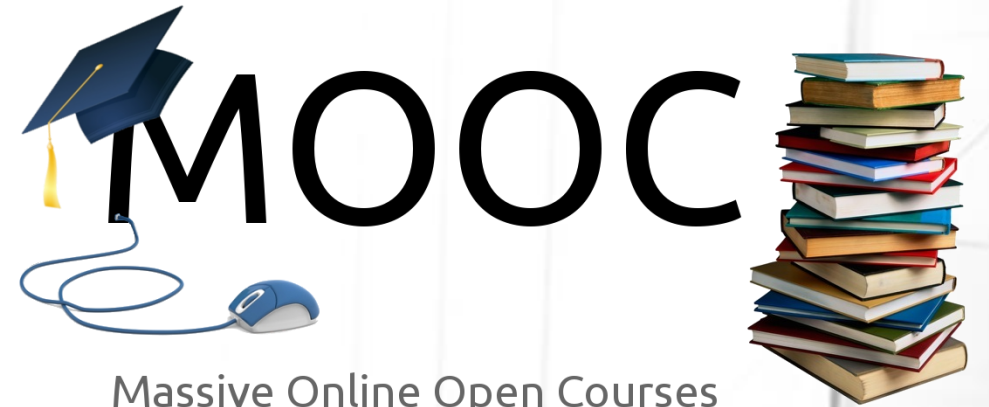
Personalización de la enseñanza

Recomendación de itinerarios / asignaturas en alumnos universitarios

- Problema reciente
- Doble objetivo:
 - Obtener un modelo que permita recomendar a los estudiantes un **itinerario académico** (qué asignaturas optativas elegir y en qué orden elegir las) en función de determinados objetivos: calificación media, especialización, esfuerzo realizado...
 - Establecer **índices de dificultad de asignaturas** basados en la experiencia previa. Estos índices pueden ser utilizados en la construcción del modelo indicado anteriormente
- Información disponible:
 - Registros académicos de una Universidad durante los **10 últimos años**. El número de registros es muy grande (más de **1000000** en este período de tiempo).

Modelado de la evaluación por pares y de la autoevaluación

- La aparición de los MOOCs (*Massive Open Online Courses*) plantea un gran número de retos asociados al número de estudiantes.
- La evaluación automática es un gran problema en algunas materias:
 - Ensayos,
 - Traducciones
 - Programas de ordenador
 - ...
- Nuevos modelos de evaluación:
 - Evaluación por pares
 - Autoevaluación



Massive Online Open Courses
Cursos online masivos y abiertos

Modelado de la evaluación por pares y de la autoevaluación

- Tanto en los sistemas de evaluación por pares como en los de autoevaluación es frecuente el uso de rúbricas (plantillas de evaluación).
- En ambos sistemas hay desviaciones entre las evaluaciones de los compañeros o del propio alumno y las evaluaciones del profesor.
- Se establecen modelos matemáticos que modelan dichas desviaciones:
 - Definición de una puntuación que mida la “reputación” de cada evaluador.

ORAL PRESENTATIONS RUBRIC

Criteria	Levels of Performance			
	1	2	3	4
Content	Energy assessment not clear; information included that does not support energy assessment in any way	There is a great deal of information that is not clearly connected to the energy assessment	Sufficient information that relates to energy assessment; many points made but uneven and little	An abundance of material clearly related with energy assessment; points are clearly made and all evidence supports energy assessment; varied use of materials
Coherence and Organization	Very Good	Good	Developing	
Research and collecting information 6	I collected <i>lots</i> of information from various places, such as books, the internet etc. 6 - 5	I collected <i>some</i> information from a few places. 4 - 3	I only collected a <i>little</i> information from few places. 2 - 1	
Sharing 8	I <i>always</i> shared my information or ideas with <i>all</i> my team members. 8 - 7	I <i>sometimes</i> shared information or ideas with my team members. 6 - 5 - 4	I shared <i>little</i> information or ideas with my team members. 3 - 2 - 1	
Completing tasks 8	I met <i>all</i> deadlines and I was not late for meetings or to complete work. 8 - 7	I met <i>most</i> deadlines and was only late for <i>some</i> meetings and to complete work. 6 - 5 - 4	I missed <i>many</i> deadlines and was <i>often</i> late for meetings or to complete work. 3 - 2 - 1	
Contribution 8	I <i>always</i> helped <i>every</i> team member with all tasks, such as gathering information, editing work. 8 - 7	I helped <i>some</i> of my team members, but not all to gather information and edit work. 6 - 5 - 4	I <i>didn't</i> help my team mates to gather information, edit work etc. 3 - 2 - 1	
Listening to other group members 5	I <i>always</i> listened to the ideas and suggestions from my team. 5 - 4	I <i>sometimes</i> listened to ideas and suggestions from my team. 3 - 2	I <i>didn't</i> listen to my other team members. I often did it my own way. 1	
Co-operating with my team 5	I <i>never</i> argued with my team members. I <i>always</i> talked about ideas and got everyone's opinion. 5 - 4	I <i>sometimes</i> argued with my team. I <i>sometimes</i> talked about ideas and thought about some opinions. 3 - 2	I <i>often</i> argued with my team mates. I <i>never</i> listened to their ideas and didn't think about their opinions. 1	

Feedback Rubric

	1
sentences are and well-d.	Many sentence fragments or run-on sentences OR paragraphing needs lots of work.
g needs	Writer makes more than 4 errors in grammar and/or spelling.
on and more errors in capitalization and punctuation.	Salutation and/or closing are missing.
using have 3 or more errors in capitalization and punctuation.	
The letter contains 3-4 accurate facts about the topic.	The letter contains 1-2 accurate facts about the topic.
The letter contains 3-4 accurate facts about the topic.	The letter contains no accurate facts about the topic.
Ideas were expressed in a clear and organized fashion. It was easy to figure out what the letter was about.	Ideas were expressed in a pretty clear manner, but could have been better.
Ideas were expressed in a clear and organized fashion. It was easy to figure out what the letter was about.	Ideas were somewhat organized, but were not very clear. It took more than one reading to figure out what the letter was about.
Ideas were expressed in a clear and organized fashion. It was easy to figure out what the letter was about.	The letter seemed to be a collection of unrelated sentences. It was very difficult to figure out what the letter was about.

Williams, S (2013) Project Rubrics Self-Assessment (V2)

Nuevos Retos



Problemas aún activos en EDS (EDM/LA)

- Desarrollo de herramientas específicas para entornos educativos y profesionales de la educación
- Datos dinámicos: minería de procesos y flujos de datos
- Análisis de video y sensores
- Interés del estudiante: emoción, afecto y elección
- Integración de las técnicas de ciencia de datos con la teoría educativa
- Mejorando el soporte al profesorado: efectividad de los materiales, metodología empleada, ...
- Análisis de datos no estructurados: documentos, foros, chats, sesiones de tutoría...
- Análisis de las redes sociales: Twitter, Facebook,...

Nuevos desafíos

Ryan Baker propone en un trabajo reciente, los “Premios Baker de Analítica de Aprendizaje”. Se trata de buscar soluciones a los siguientes problemas:

- Evaluación de los modelos desarrollados.
- Desarrollo de modelos interpretables.
- Aplicabilidad de los modelos
- Transferibilidad de los modelos
- Generalidad de los modelos

Conclusiones



Conclusiones

- La Ciencia de Datos puede utilizarse con éxito a la mejora de los sistemas educativos
- Se trata de una disciplina ya consolidada, pero que progresa muy deprisa, planteando problemas cada vez más complejos, de mayor interés
- La aparición de datos masivos (*big data*) en los sistemas educativos, y la incorporación de tecnologías emergentes (realidad virtual y aumentada, internet de las cosas, gemelos digitales, etc.) están dando una importancia creciente a la formación permanente.
- Las tecnologías de enseñanza online, soportadas por la inteligencia artificial y la computación cognitiva deben jugar un papel clave dentro de este nuevo escenario.

Ciencia de Datos en Educación

Sebastián Ventura Soto

Grupo de Investigación *Knowledge Discovery and Intelligent Systems*

Universidad de Córdoba