



## Guia docent

# 270650 - DAKD - Anàlisi de Dades i Descobriment de Coneixement

Última modificació: 12/07/2021

**Unitat responsable:** Facultat d'Informàtica de Barcelona

**Unitat que imparteix:** 723 - CS - Departament de Ciències de la Computació.

**Titulació:** MÀSTER UNIVERSITARI EN INNOVACIÓ I RECERCA EN INFORMÀTICA (Pla 2012). (Assignatura optativa).

**Curs:** 2021

**Crèdits ECTS:** 6.0

**Idiomes:** Anglès

## PROFESSORAT

**Professorat responsable:** ALFREDO VELLIDO ALCACENA

**Altres:** Primer quadrimestre:

LUIS ANTONIO BELANCHE MUÑOZ - 10

ALFREDO VELLIDO ALCACENA - 10

## CAPACITATS PRÈVIES

Students are expected to have at least some basic background in the area of artificial intelligence and, more specifically, with the areas of Machine Learning and Computational Intelligence.

Some basic knowledge of probability theory and statistics would be beneficial.

Other than this, the course is open to students and researchers of all types of background.

## COMPETÈNCIES DE LA TITULACIÓ A LES QUALS CONTRIBUEIX L'ASSIGNATURA

### Específiques:

CEC1. Capacitat per aplicar el mètode científic en l'estudi i anàlisi de fenòmens i sistemes en qualsevol àmbit de la Informàtica, així com en la concepció, disseny i implantació de solucions informàtiques innovadores i originals.

CEC3. Capacitat per aplicar solucions innovadores i realitzar avanços en el coneixement que explotin els nous paradigmes de la Informàtica, particularment en entorns distribuïts.

### Genèriques:

CG3. Capacidad para el modelado matemático, cálculo y diseño experimental en centros tecnológicos y de ingeniería de empresa, particularmente en tareas de investigación e innovación en todos los ámbitos de la Informática.

### Transversals:

CTR4. ÚS SOLVENT DELS RECURSOS D'INFORMACIÓ: Gestionar l'adquisició, l'estructuració, l'anàlisi i la visualització de dades i d'informació de l'àmbit de l'enginyeria informàtica, i valorar de forma crítica els resultats d'aquesta gestió.

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

## METODOLOGIES DOCENTS

This course will build on different teaching methodology (TM) aspects, including:

TM1: Expositive seminars

TM2: Expositive-participative seminars

TM3: Orientation for individual assignments (essays)

TM4: Individual tutorization



## OBJECTIUS D'APRENENTATGE DE L'ASSIGNATURA

1. Presenting DM as a process that should involve a methodology id applied at its best.
2. Introducing the students to the new concept of DM for processes, called Process Mining.
3. Delving into some detail in one of the stages of DM: data exploration.
4. Dealing in detail with the problem of data visualization for exploration as a key issue in DM.
5. Introducing the students to the basics of probability theory as applied in Data Analysis and Knowledge Discovery (DAKD)
6. Introducing the students to the probabilistic variant of DAKD in the form of Statistical Machine Learning, both for supervised and unsupervised learning models.
7. Dealing in detail with different unsupervised models for data visualization, including case studies.
8. Approaching the multi-faceted concept of data mining (DM) from different perspectives.

## HORES TOTALS DE DEDICACIÓ DE L'ESTUDIANTAT

Tipus	Hores	Percentatge
Hores aprenentatge autònom	96,0	64.00
Hores grup gran	45,0	30.00
Hores activitats dirigides	9,0	6.00

**Dedicació total:** 150 h

## CONTINGUTS

### Introduction to the concept of data mining (DM).

#### Descripció:

DM is a multi-faceted concept that requires discussion and clarification. We will do this at the beginning of the course.

### DM as a methodology.

#### Descripció:

We argue that DM should not be focused on the concept of data analysis/modeling, but, instead, should be treated as a methodology with diverse inter-related stages.

### DM for processes: Process Mining.

#### Descripció:

A new development in DM methodologies is that which deals with one specifically suited for processes. It is called Process Mining and will be described and discussed in this course.

### Data exploration in DM.

#### Descripció:

One of the main stages of well-structures DM methodologies is Data exploration. It will be discussed as a preamble to data visualization.



### Data visualization for exploration.

#### Descripció:

One of the aspects of the problem of data exploration is data visualization. It has a research 'life' of its own as it involves not only computer-based mathematical models, but also natural perception and processing.

### Basics of probability theory in Data Analysis and Knowledge Discovery (DAKD)

#### Descripció:

For a long time in the last half-century, multivariate statistics and artificial intelligence (mostly in the field of machine learning) have developed in parallel without fully meeting. Statistical machine learning has bridged that field over the last two decades. We introduce it by first providing some basic principles of probability theory (Bayesian inference).

### Statistical Machine Learning for DAKD: supervised models.

#### Descripció:

Once the basics of Bayesian inference are set, we will delve into the field of Statistical Machine Learning for IDA, starting with supervised learning models, with an emphasis on feed-forward artificial neural networks.

### Statistical Machine Learning for DAKD: unsupervised models.

#### Descripció:

Once the basics of Bayesian inference and of Statistical Machine Learning for IDA in supervised models are set, we will continue with unsupervised models, focusing on self-organizing maps and related models.

### Unsupervised models for data visualization, with case studies.

#### Descripció:

In the final item of the contents of the course, we will bring statistical machine learning and data visualization together by discussing some probabilistic unsupervised learning models for data visualization, including some case studies as an example.



## ACTIVITATS

### Essay on DAKD for DM

**Descripció:**

Students will have to write a research essay on the topic of DAKD for DM, with different options:

1. State of the art on an specific DAKD-DM topic
2. Evaluation of an DAKD-DM software tool with original experiments
3. Pure research essay, with original experimental content

**Objectius específics:**

1, 2, 3, 4, 5, 6, 7, 8

**Competències relacionades:**

CG3. Capacidad para el modelado matemático, cálculo y diseño experimental en centros tecnológicos y de ingeniería de empresa, particularmente en tareas de investigación e innovación en todos los ámbitos de la Informática.

CEC1. Capacitat per aplicar el mètode científic en l'estudi i anàlisi de fenòmens i sistemes en qualsevol àmbit de la Informàtica, així com en la concepció, disseny i implantació de solucions informàtiques innovadores i originals.

CEC3. Capacitat per aplicar solucions innovadores i realitzar avanços en el coneixement que explotin els nous paradigmes de la Informàtica, particularment en entorns distribuïts.

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

CTR4. ÚS SOLVENT DELS RECURSOS D'INFORMACIÓ: Gestionar l'adquisició, l'estructuració, l'anàlisi i la visualització de dades i d'informació de l'àmbit de l'enginyeria informàtica, i valorar de forma crítica els resultats d'aquesta gestió.

**Dedicació:** 3h

Activitats dirigides: 3h

### Introduction to Data Mining and its Methodologies

**Descripció:**

Introduction to Data Mining as a general concept and to its methodologies for practical implementation

**Objectius específics:**

1

**Competències relacionades:**

CG3. Capacidad para el modelado matemático, cálculo y diseño experimental en centros tecnológicos y de ingeniería de empresa, particularmente en tareas de investigación e innovación en todos los ámbitos de la Informática.

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

CTR4. ÚS SOLVENT DELS RECURSOS D'INFORMACIÓ: Gestionar l'adquisició, l'estructuració, l'anàlisi i la visualització de dades i d'informació de l'àmbit de l'enginyeria informàtica, i valorar de forma crítica els resultats d'aquesta gestió.

**Dedicació:** 23h

Grup gran/Teoria: 9h

Activitats dirigides: 1h

Aprendentatge autònom: 13h



## Process Mining

### Descripció:

Introduction to the novel concept of Process Mining and its application within the DM framework.

### Objectius específics:

2

### Competències relacionades:

CG3. Capacidad para el modelado matemático, cálculo y diseño experimental en centros tecnológicos y de ingeniería de empresa, particularmente en tareas de investigación e innovación en todos los ámbitos de la Informática.

CEC3. Capacitat per aplicar solucions innovadores i realitzar avanços en el coneixement que explotin els nous paradigmes de la Informàtica, particularment en entorns distribuïts.

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

### Dedicació: 9h

Grup gran/Teoria: 3h

Activitats dirigides: 1h

Aprenentatge autònom: 5h

## Data Visualization

### Descripció:

As part of the DM stage of Data Exploration, we focus in the problem of Data Visualization.

### Objectius específics:

3, 4

### Competències relacionades:

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

CTR4. ÚS SOLVENT DELS RECURSOS D'INFORMACIÓ: Gestionar l'adquisició, l'estructuració, l'anàlisi i la visualització de dades i d'informació de l'àmbit de l'enginyeria informàtica, i valorar de forma crítica els resultats d'aquesta gestió.

### Dedicació: 16h

Grup gran/Teoria: 6h

Activitats dirigides: 1h

Aprenentatge autònom: 9h



## Basics of probability theory for intelligent data analysis

### Descripció:

Introduction to probability theory for intelligent data analysis, with a focus on Bayesian statistics

### Objectius específics:

5

### Competències relacionades:

CEC1. Capacitat per aplicar el mètode científic en l'estudi i anàlisi de fenòmens i sistemes en qualsevol àmbit de la Informàtica, així com en la concepció, disseny i implantació de solucions informàtiques innovadores i originals.

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

CTR4. ÚS SOLVENT DELS RECURSOS D'INFORMACIÓ: Gestionar l'adquisició, l'estructuració, l'anàlisi i la visualització de dades i d'informació de l'àmbit de l'enginyeria informàtica, i valorar de forma crítica els resultats d'aquesta gestió.

### Dedicació: 16h

Grup gran/Teoria: 6h

Activitats dirigides: 1h

Aprenentatge autònom: 9h

## Statistical Machine Learning methods

### Descripció:

The meeting of statistics and machine learning: Statistical Machine Learning methods, from the point of view of both supervised and unsupervised learning

### Objectius específics:

5, 6

### Competències relacionades:

CEC1. Capacitat per aplicar el mètode científic en l'estudi i anàlisi de fenòmens i sistemes en qualsevol àmbit de la Informàtica, així com en la concepció, disseny i implantació de solucions informàtiques innovadores i originals.

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

CTR4. ÚS SOLVENT DELS RECURSOS D'INFORMACIÓ: Gestionar l'adquisició, l'estructuració, l'anàlisi i la visualització de dades i d'informació de l'àmbit de l'enginyeria informàtica, i valorar de forma crítica els resultats d'aquesta gestió.

### Dedicació: 31h

Grup gran/Teoria: 12h

Activitats dirigides: 1h

Aprenentatge autònom: 18h



## SML in data visualization, with case studies

### Descripció:

We merge the topics of SML and data visualization, illustrating its use with some real case studies

### Objectius específics:

4, 7, 8

### Competències relacionades:

CG3. Capacidad para el modelado matemático, cálculo y diseño experimental en centros tecnológicos y de ingeniería de empresa, particularmente en tareas de investigación e innovación en todos los ámbitos de la Informática.

CTR6. RAONAMENT: Capacitat de raonament crític, lògic i matemàtic. Capacitat de resoldre problemes en la seva àrea d'estudi.

Capacitat d'abstracció: capacitat de crear i utilitzar models que reflecteixin situacions reals. Capacitat de dissenyar i realitzar experiments senzills, i analitzar-ne i interpretar-ne els resultats. Capacitat d'anàlisi, de síntesi i d'avaluació.

CTR4. ÚS SOLVENT DELS RECURSOS D'INFORMACIÓ: Gestionar l'adquisició, l'estructuració, l'anàlisi i la visualització de dades i d'informació de l'àmbit de l'enginyeria informàtica, i valorar de forma crítica els resultats d'aquesta gestió.

**Dedicació:** 25h

Grup gran/Teoria: 9h

Activitats dirigides: 1h

Aprendentatge autònom: 15h

## SISTEMA DE QUALIFICACIÓ

The course will be evaluated through a final essay that will take one of these three modalities:

1. State of the art on an specific IDA-DM topic
2. Evaluation of an IDA-DM software tool with original experiments
3. Pure research essay, with original experimental content

## BIBLIOGRAFIA

### Bàsica:

- MacKay, D.J.C. Information theory, inference, and learning algorithms. Cambridge University Press, 2003. ISBN 0521642981.
- Hand, D.; Mannila, H.; Smyth, P. Principles of data mining. MIT Press, 2001. ISBN 026208290X.
- Bishop, C.M. Pattern recognition and machine learning. New York: Springer, 2006. ISBN 0387310738.

### Complementària:

- Hand, D.J. Statistics: a very short introduction. Oxford University Press, 2008. ISBN 9780199233564.
- Spence, R. Information visualization: design for interaction. 2nd ed. Pearson/Prentice Hall, 2007. ISBN 9780132065504.
- Yau, N. Visualize this: the flowing data guide to design, visualization, and statistics. Wiley, 2011. ISBN 9780470944882.