

# Data Mining

Nom : **EUDES**

Prénom : **Richard**

Année : **M2**

Semestre : **10**

Nature : **CM**

Volume horaire : **24 H**

ECTS / Coef : **2**

Prérequis	Une connaissance de niveau M1 des statistiques et de l'optimisation (recherche opérationnelle). Une expérience antérieure du logiciel SAS est utile (la programmation SAS base est un plus).
Résumé	Ce cours couvre les compétences attendues d'un data miner / data scientist pour structurer un protocole d'analyse de données à l'aide de workflows (SAS Enterprise Miner) dans la création de modèles non supervisés (segmentations, analyses d'association et de séquence) et supervisés (arbre de décision, régression et modèles de réseaux neuronaux). Il résume l'ensemble des étapes statistiques nécessaires à la préparation de données et à la construction de modèles pertinents et interprétables. En outre, ce cours permet d'introduire des éléments de gestion de projets propres aux démarches analytiques modernes (CRISP, SEMMA).
Objectifs	<ul style="list-style-type: none"><li>- définir un projet et explorer les données graphiquement (dataviz)</li><li>- modifier les données pour de meilleurs résultats d'analyse (transformer les variables, les données manquantes,...)</li><li>- construire et comprendre des modèles prédictifs tels que des arbres de décision, des modèles de régression, des réseaux de neurones (et d'autres algorithmes comme SVM, PLS Regression, Knn, ..)</li><li>- comparer et expliquer des modèles complexes</li><li>- générer et utiliser le code de score (scoring).</li><li>- appliquer l'analyse d'association et de séquence aux données de transaction (techniques de clustering)</li><li>- utiliser d'autres outils de modélisation tels que l'induction de règles, le Gradient Boosting et les Support Vector Machines</li></ul> <p>Ce cours peut vous aider à vous préparer aux examens de certification suivants: Modélisation prédictive avec SAS Enterprise Miner.</p>
Bibliographie	<p>Beck, A. 1997. "Herb Edelstein discusses the usefulness of data mining." <i>DS Star</i>. Vol. 1, N0. 2. Available <a href="http://www.tgc.com/dsstar/">www.tgc.com/dsstar/</a>.</p> <p>Bishop, C. M. 1995. <i>Neural Networks for Pattern Recognition</i>. New York: Oxford University Press.</p> <p>Breiman, L. et al. 1984. <i>Classification and Regression Trees</i>. Belmont, CA: Wadsworth International Group.</p> <p>Hand, D. J. 1997. <i>Construction and Assessment of Classification Rules</i>. New York: John Wiley &amp; Sons, Inc.</p> <p>Hand, D. J. 2005. "What you get is what you want? – Some dangers of black box data mining." <i>M2005 Conference Proceedings</i>, Cary, NC: SAS Institute Inc.</p> <p>Hand, D. J. 2006. "Classifier technology and the illusion of progress." <i>Statistical Science</i> 21:1-14.</p> <p>Hand, D. J. and W. E. Henley. 1997. "Statistical classification methods in consumer credit scoring: a review." <i>Journal of the Royal Statistical Society A</i> 160:523-541.</p> <p>Hand, David, Heikki Mannila, and Padraic Smyth. 2001. <i>Principles of Data Mining</i>. Cambridge, Massachusetts: The MIT Press.</p> <p>Harrell, F. E. 2006. <i>Regression Modeling Strategies</i>. New York: Springer-Verlag New York, Inc.</p>

- Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2001. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York: Springer-Verlag New York, Inc.
- Hoaglin, D. C., F. Mosteller, and J. W. Tukey. 1983. *Understanding Robust and Exploratory Data Analysis*. New York: John Wiley & Sons, Inc.
- Kass, G. V. 1980. "An exploratory technique for investigating large quantities of categorical data." *Applied Statistics* 29:119-127.
- Mosteller, F. and J. W. Tukey. 1977. *Data Analysis and Regression*. Reading, MA: Addison-Wesley.
- Piatetsky-Shapiro, G. 1998. "What Wal-Mart might do with Barbie association rules." *Knowledge Discovery Nuggets*, 98:1. Available <http://www.kdnuggets.com/>.
- Ripley, B. D. 1996. *Pattern Recognition and Neural Networks*. New York: Cambridge University Press.
- Rubinkam, M. 2006. "Internet Merchants Fighting Costs of Credit Card Fraud," *AP Worldstream*. The Associated Press.
- Rud, Olivia Parr. 2001. *Data Mining Cookbook: Modeling Data, Risk, and Customer Relationship Management*. New York: John Wiley & Sons, Inc.
- Sarle, W. S. 1983. *Cubic Clustering Criterion*. SAS Technical Report A-108. Cary, NC: SAS Institute Inc.
- Sarle, W. S. 1994a. "Neural Networks and Statistical Models," *Proceedings of the Nineteenth Annual SAS® Users Group International Conference*. Cary: NC, SAS Institute Inc., 1538-1550.
- Sarle, W. S. 1994b. "Neural Network Implementation in SAS® Software," *Proceedings of the Nineteenth Annual SAS® Users Group International Conference*. Cary: NC, SAS Institute Inc., 1550-1573.
- Sarle, W. S. 1995. "Stopped Training and Other Remedies for Overfitting." *Proceedings of the 27th Symposium on the Interface*.
- SAS Institute Inc. 2002. *SAS® 9 Language: Reference, Volumes 1 and 2*. Cary, NC: SAS Institute Inc.
- SAS Institute Inc. 2002. *SAS® 9 Procedures Guide*. Cary, NC: SAS Institute Inc.
- SAS Institute Inc. 2002. *SAS/STAT® 9 User's Guide, Volumes 1, 2, and 3*. Cary, NC: SAS Institute Inc.
- Weiss, S. M. and C. A. Kulikowski. 1991. *Computer Systems That Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems*. San Mateo, CA: Morgan Kaufmann.

---

# PLAN

---

Le support de cours étant en anglais, il se structure de la manière suivante :

**Chapter 1 Introduction**

- 1.1 Introduction to data mining (history) and data science

**Chapter 2 Accessing and Assaying Prepared Data**

- 2.1 Introduction
- 2.2 Designing a Project, Library, and Diagram
- 2.3 Defining a Data Source
- 2.4 Exploring a Data Source

**Chapter 3 Introduction to Predictive Modeling: Decision Trees**

- 3.1 Introduction
- 3.2 Cultivating Decision Trees
- 3.3 Optimizing the Complexity of Decision Trees
- 3.4 Understanding Additional Diagnostic Tools
- 3.5 Autonomous Tree Growth Options

**Chapter 4 Introduction to Predictive Modeling: Regressions**

- 4.1 Introduction
- 4.2 Selecting Regression Inputs
- 4.3 Optimizing Regression Complexity
- 4.4 Interpreting Regression Models
- 4.5 Transforming Inputs
- 4.6 Categorical Inputs
- 4.7 Polynomial Regressions

**Chapter 5 Introduction to Predictive Modeling: Neural Networks and Other Modeling Tools**

- 5.1 Introduction
- 5.2 Input Selection
- 5.3 Stopped Training
- 5.4 Other Modeling Tools

**Chapter 6 Model Assessment**

- 6.1 Model Fit Statistics
- 6.2 Statistical Graphics
- 6.3 Adjusting for Separate Sampling
- 6.4 Profit Matrices

**Chapter 7 Model Implementation**

- 7.1 Introduction
- 7.2 Internally Scored Data Sets
- 7.3 Score Code Modules

**Chapter 8 Introduction to Pattern Discovery**

- 8.1 Introduction
- 8.2 Cluster Analysis
- 8.3 Market Basket Analysis

**Chapter 9 Special Topics**

- 9.1 Introduction
- 9.2 Ensemble Models
- 9.3 Variable Selection
- 9.4 Categorical Input Consolidation
- 9.5 Surrogate Models
- 9.6 SAS Rapid Predictive Modeler

**Appendix A Case Studies**

- A.1 Banking Segmentation Case Study
- A.2 Web Site Usage Associations Case Study
- A.3 Credit Risk Case Study
- A.4 Enrollment Management Case Study