

Big Data Analytics

Session 4

- *Online Session* -

2019 - 2020

By **PhD. Samia Chehbi Gamoura**
Associate Professor - EM Strasbourg
France

Dear students welcome to my online MOOC platform. I am the avatar of your professor Samia Gamoura of this course, entitled "Big Data Analytics." This online session is also available on the website of your professor: www.samiagamoura.com. The last work provided by your professor was Exercise 3 about Machine Learning algorithms in Business cases. Your professor gave you enough time to work on the solution in the previous session. Some of you already sent me their complete exercises. Thanks for everyone.



Big Data Analytics → Advanced Analytics → Machine Learning

ML Algorithm	Type	Business Case	Area	With Big Data?	With Big Data, Applicable?
Association Rule	Unsupervised	Market basket Analysis	E-marketing	No	Yes
Naive Bayes Algorithm	Supervised	Housing Price Prediction	Real Estate	Yes	No
AdaBoost	Ensemble	Corporate Bankruptcy	Risk Management	Yes	Yes
K-Means Clustering	Unsupervised	Market Segmentation	Marketing Strategy	No	Yes
Logistic Regression	Supervised	Predicts Business Failure	Risk Management	Yes	Yes
Artificial Neural Network	Supervised	Tourism Demand Forecasts	Tourism and hospitality industry	Yes	Yes
Decision tree	Supervised	Financial crisis prediction	Financial industry and budgeting	Yes	Yes
Version space learning	Supervised	Recommender system	Marketing	Yes	No
Rough set rules	Unsupervised	Product classification	Production in Supply chain management	Yes	Yes
K-Means	Unsupervised	Stock market forecasting	Finance, investments	Yes	Yes



Please check this list of the selected correct answers in your solutions. Don't worry, you can upload this solution afterward. You can find it in PDF format in Moodle in the Exercises part, as you can see in the screenshot in the subsequent slide.

Session 4 - Online - Exercises and Case Studies

Session 4 - Online - Exercises and Case Studies

- Your solutions reviewed by the Professor
 - Group 1_Big Data Exercise Revised_By_Gamoura.docx
 - Group 2_Applicative Exercise 3_revised_by_Gamoura.xlsx
 - Group 3_Big data_revised_by_Gamoura.docx
 - Group 4_Case-study 02.03.2020_revised_by_Gamoura.docx
 - Course_EM054M8K_Big Data Analytics Content Exercise 3_Sol.pdf
 - Scoring_List_Students_Big Data Analytics_2020_Exercise_3.pdf



Only 68% of students provided their solutions. All of them have +0.5 point as an additional reward in the practical part of the exam. To check the list of names, please report to the PDF file of scoring uploaded in the Exercises part in Moodle. Your professor reviewed all your solutions and uploaded them in a separated directory as you can see in the screenshot here.



Big Data Analytics → Advanced Analytics → Machine Learning

Group 1_Big Data Exercise_ Revised_By_Gamoura.docx - Word

ML Algorithm	Type	Business Case	Area	With Big Data?	With Big Data, Applicable?
Naive Bayes Algorithm	Supervised	Housing Price Prediction	Real Estate	Yes	No
Multifactor Linear Regression Model	Supervised	Industrial Gas Business	Supply Chain Management	Yes	Yes
Ada Boost	Supervised Ensemble	Corporate Bankruptcy	Risk Management	Yes	Yes
PCA principal component analysis	Unsupervised	Reengineering Design of Purchasing System	Accounting	Yes	Yes
K-Means	Unsupervised	Market Segmentation	Marketing Strategy	No	No
KNN K-nearest neighbors	Supervised	Stock Price Prediction	Investments	No	No
CART Classification and Regression Trees	Supervised	Algorithm and Blood Donor Classification	Corporate Healthcare	Yes	Yes
Logistic Regression	Supervised	Predicts Business Failure	Risk Management	Yes	Yes

Comments and corrections by Samia Gamoura:

- OK
- Please specify. SCM is large field. Here what was the case? Prediction of what?
- OK
- It's not a machine learning technique but comes to support an ML
- Why not in this case?
- OK. Good as we need proximity here so no need for Big Data
- I don't understand!
- OK
- Good

Samia Gamoura
Mis en forme: Barré

When you download and open your file in your given solution, you find the comments and corrections of your professor. Please get the corrections and accept or reject the modification by using the track tool in your word file. For the Excel files, please open your comments. Now, let's switch to the course to continue with the Machine Learning part.



Analytics → Types → Advanced Analytics → Machine Learning

4 categories (meta – approaches) following the input Data:

- **Information-Based Machine Learning**
- **Erro-Based Machine Learning**
- **Similarity-Based Machine Learning**
- **Probability-Based Machine Learning**

There are four meta-categories if we base the categorization on the input data. You do not need to know all the technical aspects of these categories and algorithms as you do not have the technical background of algorithmic and computer science in this class. However, you need to know, at least, these approaches of Machine Learning are used in Business and management applications in enterprises.



■ Information-Based Machine Learning

- Shannon measurement theory presented in 1948 by Claude Shannon
- The extraction of a distinctive measurement from the Data Descriptive Features through the information contents.
- Need of use of information content measurement and metrics

There are four categories of Machine Learning meta algorithms if we base the categorization on the input Data. Information-based category of algorithms in Machine Learning are founded on Shannon measurement Theory of 1948. In this class of algorithms, the distinction is based on measurement from descriptive features in the information content. Keep in mind here that you need measurement metrics to perform in this category.



Information-Based Machine Learning

Approach	Types	Sub-Types
Decision Trees (DTR)	Supervised	Classification & Prediction
Conditional Random Fields (CRF)	Supervised	Classification & Prediction
Linear Regression (LRE)	Supervised	Classification & Prediction
Non Linear Regression (NLR)	Supervised	Classification & Prediction
Polynomial regression (PRE)	Supervised	Classification & Prediction
Maximum Entropy (MEN)	Supervised	Classification & Prediction
Frequent Pattern Growth (FRG)	Unsupervised	Association Rule
Apriori Algorithm (AAL)	Unsupervised	Association Rule
Eclat Algorithm (EAL)	Unsupervised	Association Rule
Random Decision Forest (RDF)	Ensemble	Bagging
AdaBoost Algorithm (ABA)	Ensemble	Bagging
LogitBoost Algorithm (LBA)	Ensemble	Bagging

In Business application of Machine Learning. These are some of the well-known ML algorithms in information-based learning like Decision Trees and Linear Regression that are supervised. Others like Eclat, and Apriori are unsupervised. We can find also some Ensemble meta algorithms like Ada Boost and Random Forest.



■ Error-Based Machine Learning

- Mathematical models that are founded on the idea of getting performance through error minimization
- Need error measure and the error surface (space)

The second category gathers the Error-Based Machine Learning approaches. They perform by computing the error as the measurement of the difference between the desired class and the obtained class. Backpropagation is known as the widely used method of error computing and correcting in the learning model. The algorithms of this category always need the mathematical measurement of the error.



■ Error-Based Machine Learning

Approach	Types	Sub-types
Artificial Neural Network (ANN)	Supervised	Classification & prediction
Fisher Linear Discriminant (FLD)	Supervised	Classification & prediction
Logistic Regression (LRE)	Supervised	Classification & prediction
Support Vector Machine (SVM)	Supervised	Classification & prediction



Algorithms in Error-Based models could be Support Vector Machines (SVM) and Artificial Neural Networks including a Deep Learning algorithms such as Conv Networks. Error-based algorithms are mainly supervised.



■ Similarity-Based Machine Learning

- The best way of predicting the future is to simply compare and find similarities of features from the past
- Need techniques of metrics in determining similarities degrees in the space of descriptive features

Similarity-based Machine Learning algorithms are in the third category. They are based on the hypothesis that the best way to predict the future is to simply compare and find similarities with the past. Similarities key features and comparison metrics are then needed for these approaches.



■ Similarity-Based Machine Learning

Approach	Types	Sub-Types
Case-Bases Reasoning (CBR)	Supervised	Classification & Prediction
Hidden Markov Models (HMM)	Supervised	Classification & Prediction
K-Means (KME)	Unsupervised	Clustering
Mixture Models (MMO)	Unsupervised	Clustering
Hierarchical Cluster Analysis (HCA)	Unsupervised	Clustering
Hebbian Learning (HLE)	Unsupervised	Unsupervised Artificial Neural Networks
Generative Adversarial Networks (GAN)	Unsupervised	Unsupervised Artificial Neural Networks
Sarsa (SAR)	Réinforcement	Model-free Reinforcement
Q-Learning (QLE)	Réinforcement	Model-free Reinforcement
Divergence de Kullback-Leibler (DKL)	Réinforcement	Model-based Reinforcement

In similarity-based approaches, we find the different modes of learning classes: supervised learning with hidden Markov models for example, unsupervised learning with clustering, and reinforcement learning with Q-learning and Sarsa algorithm.



■ Probability-Based Machine Learning

- Based on Bayes theorem
- The future is a random event based on relative frequencies with calculation of conditional probabilities based on the sequence of present and the past

The fourth category is probability-based learning technique. The key idea in these approaches is the use of Bayes theorem. In Bayes, the future is defined as a random event based on relative frequencies with calculation of conditional probabilities based on the sequence of present and the past.



■ Probability-Based Machine Learning

Approach	Types	Sub-types
Naive Bayes Classifier (NBC)	Supervised	Classification & prediction
Bayesian Networks (BNE)	Supervised	Classification & prediction

In business problems using probability-based approaches in Machine Learning, we find the majority of use in supervised algorithms such as Naive Bayes classifiers and Bayesian networks in classification.



Case Study 2

Case Study 2

Exigency Mandatory

Type Individual Work

Mode Homework

Estimated duration 120 minutes

Evaluation included in the practical part of the Exam

Repository **Session 4 - Case Study 2 to do (Mandatory - Part of the exam evaluation)**

Deadline **20th March**

This online is now finished. Now, you have homework to do. Now. You have this case study number 2 to do. This case is individual. It's a homework to perform by yourself. The estimated duration to do it is about two hours (120 minutes). This work will be evaluated by your professor afterward. The evaluation is included in the practical part of the exam. The deadline to provide this work is 20th March.



Case Study 2

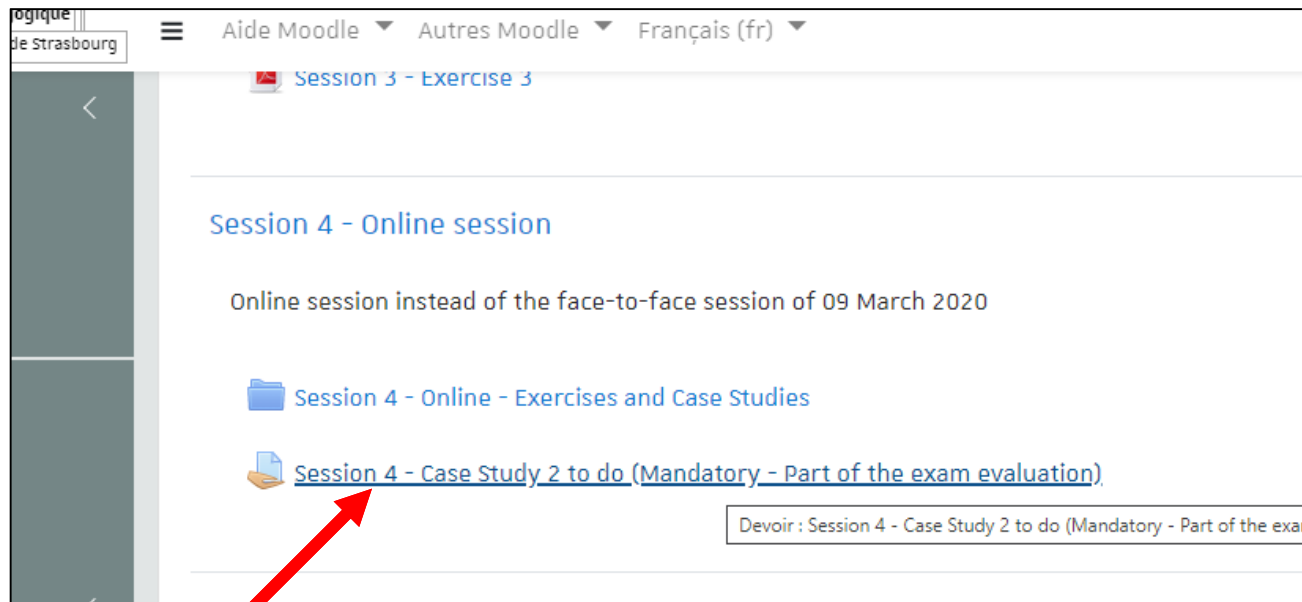
- Realize a research investigation in literature on one of the following topics:
(* please select only one topic in the list):
 1. Chatbot in Marketing
 2. Opinion mining in Customer Relationship Management
 3. Additive Manufacturing
 4. Recommender system in web Marketplace
 5. Social Computing in customer prospection
 6. Augmented product in Sales Management
 7. Business Process Automation
- by following the structure here below.
 1. Context and Definitions
 2. Examples and cases from literature
 3. Your criticism(s) and your opinion,
 4. Perspectives of use and conclusion

By using Internet documentation and academic references such as Google Scholar, try to realize a research investigation in literature. You have to select one topic to develop in a report (MS Word for example), by following the structure here below. Content and Definitions. Examples and cases from literature. Your criticism(s) and your opinion. Perspectives of use and conclusion. You have a list of topics proposed by your professor. Each student has to select one topic to develop. There is no issue if two students select the same topic. But this work still individual and each student should provide his own work.



Case Study 2

- Repository: Session 4 - Case Study 2 to do (Mandatory - Part of the exam evaluation)



As already mentioned, you have to upload your files in the repository named Session 4 – Case Study 2.



Case Study 2

- **Deadline: 20th March**

édagogique
sité de Strasbourg

Aide Moodle ▾ Autres Moodle ▾ Français (fr) ▾

Session 4 - Case Study 2 to do (Mandatory - Part of the exam evaluation)

Session 4 - Case Study 2 to do (Mandatory - Part of the exam evaluation)

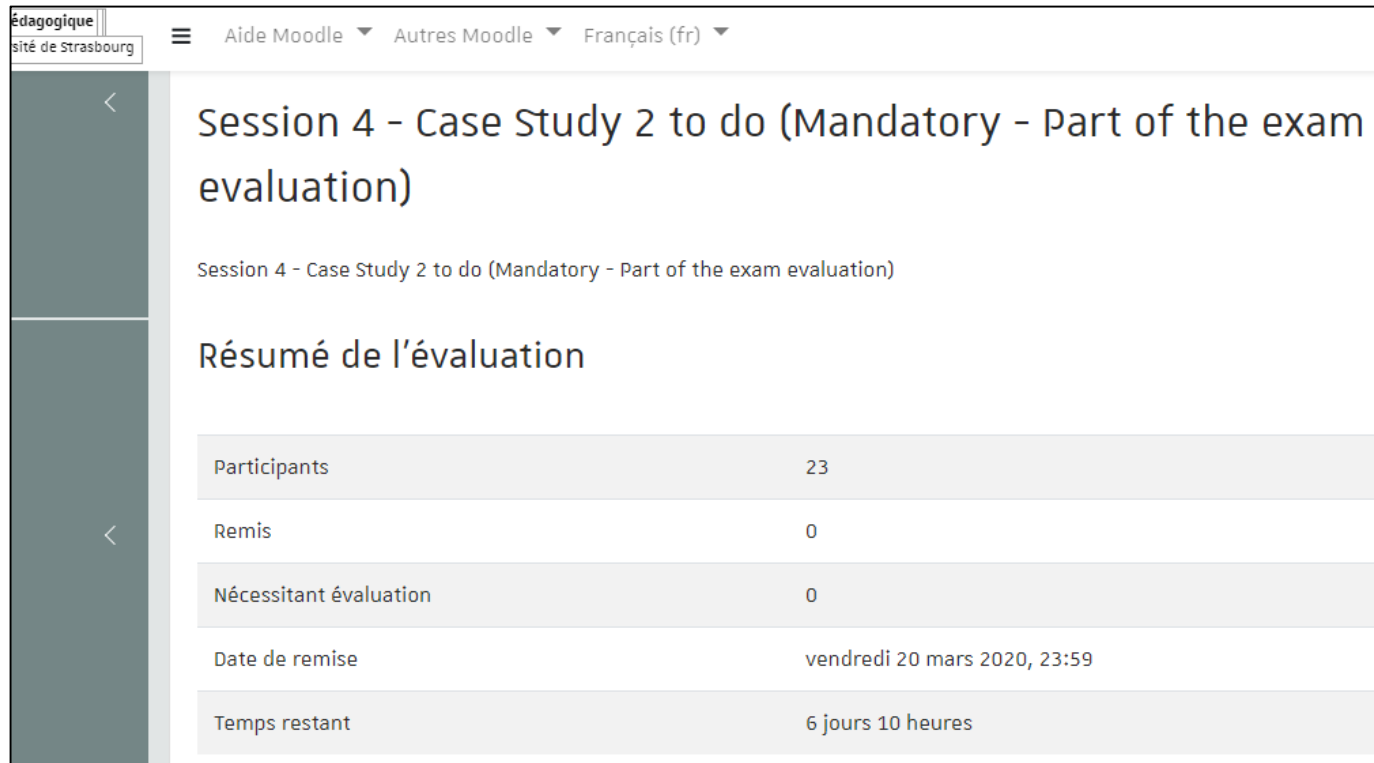
Résumé de l'évaluation

Participants	23
Remis	0
Nécessitant évaluation	0
Date de remise	vendredi 20 mars 2020, 23:59
Temps restant	6 jours 10 heures

The deadline is Friday 20th March. The system will close the repository automatically when deadline is over. No file will be accepted afterward. Files sent by emails are not accepted.



Case Study 2



Edagogique
Université de Strasbourg

Aide Moodle ▾ Autres Moodle ▾ Français (fr) ▾

Session 4 - Case Study 2 to do (Mandatory - Part of the exam evaluation)

Session 4 - Case Study 2 to do (Mandatory - Part of the exam evaluation)

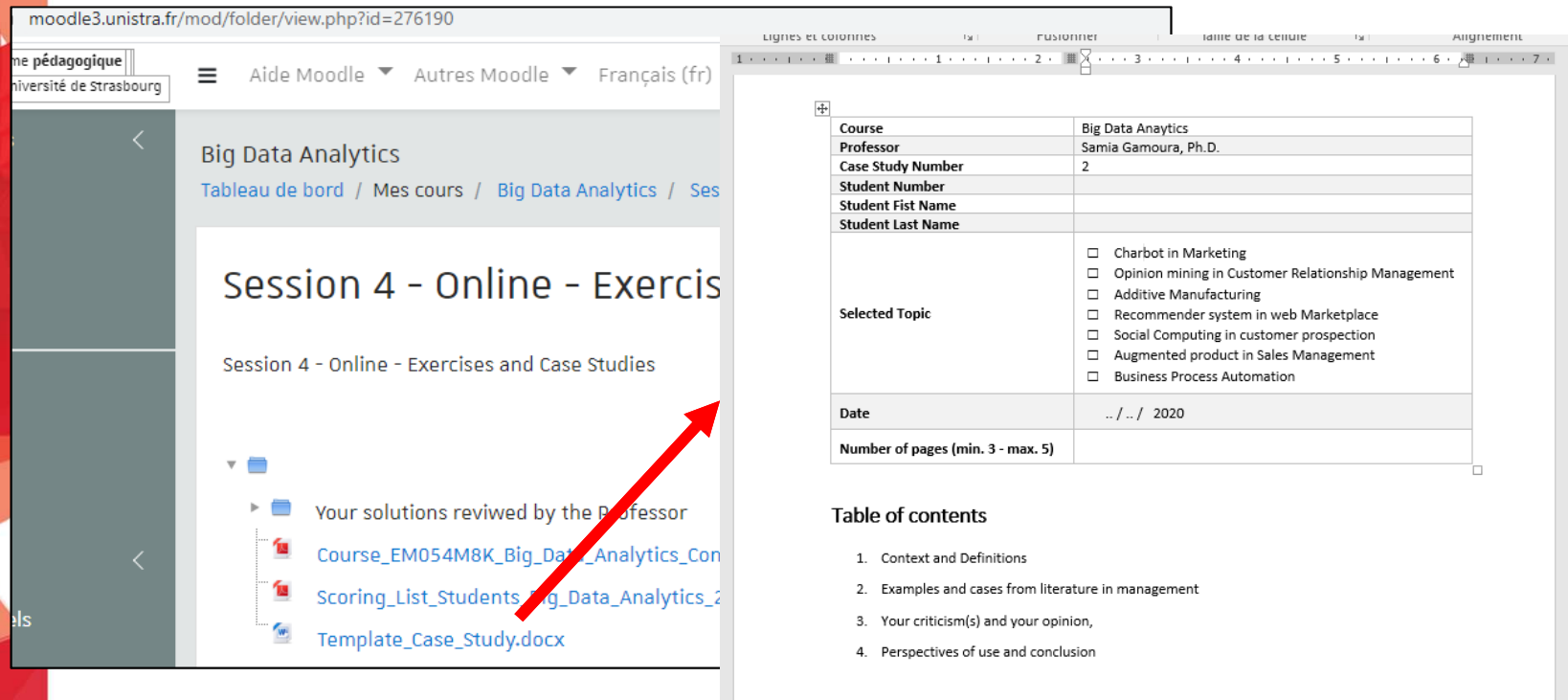
Résumé de l'évaluation

Participants	23
Remis	0
Nécessitant évaluation	0
Date de remise	vendredi 20 mars 2020, 23:59
Temps restant	6 jours 10 heures

The file you provide must be in Word format (not PDF) and must be named with your full name. For example, the file of the student **TOM BRYAN** must have the name: **TOM BRYAN _Case_Study_2.docx (or TOM BRYAN _Case_Study_2.doc)**



Case Study 2



The screenshot shows a Moodle course page for 'Big Data Analytics' at the University of Strasbourg. The page title is 'Session 4 - Online - Exercis'. Below the title, there is a sub-section 'Session 4 - Online - Exercises and Case Studies'. A red arrow points from a file named 'Template_Case_Study.docx' in the course content area to a form on the right side of the page. The form is titled 'Table of contents' and contains the following fields:

Course	Big Data Analytics
Professor	Samia Gamoura, Ph.D.
Case Study Number	2
Student Number	
Student First Name	
Student Last Name	
Selected Topic	<input type="checkbox"/> Charbot in Marketing <input type="checkbox"/> Opinion mining in Customer Relationship Management <input type="checkbox"/> Additive Manufacturing <input type="checkbox"/> Recommender system in web Marketplace <input type="checkbox"/> Social Computing in customer prospection <input type="checkbox"/> Augmented product in Sales Management <input type="checkbox"/> Business Process Automation
Date	.. / .. / 2020
Number of pages (min. 3 - max. 5)	

Below the form, there is a 'Table of contents' section with the following items:

1. Context and Definitions
2. Examples and cases from literature in management
3. Your criticism(s) and your opinion,
4. Perspectives of use and conclusion

You find also a template word file to use for your report. You have to fill the head and develop the content as required.



Case Study 2

moodle3.unistra.fr/course/view.php?id=9662#section-4

Université de Strasbourg

Aide Moodle Autres Moodle Français (fr)

Session 3 - Exercice 3

Ajouter une act

Session 4 - Online session

Online session instead of the face-to-face session of 09 March 2020

Session 4 - Online - Exercises and Case Studies

Session 4 - Case Study 2 to do (Mandatory - Part of the exam evaluation)

Session 4_Forum_OnlineSession

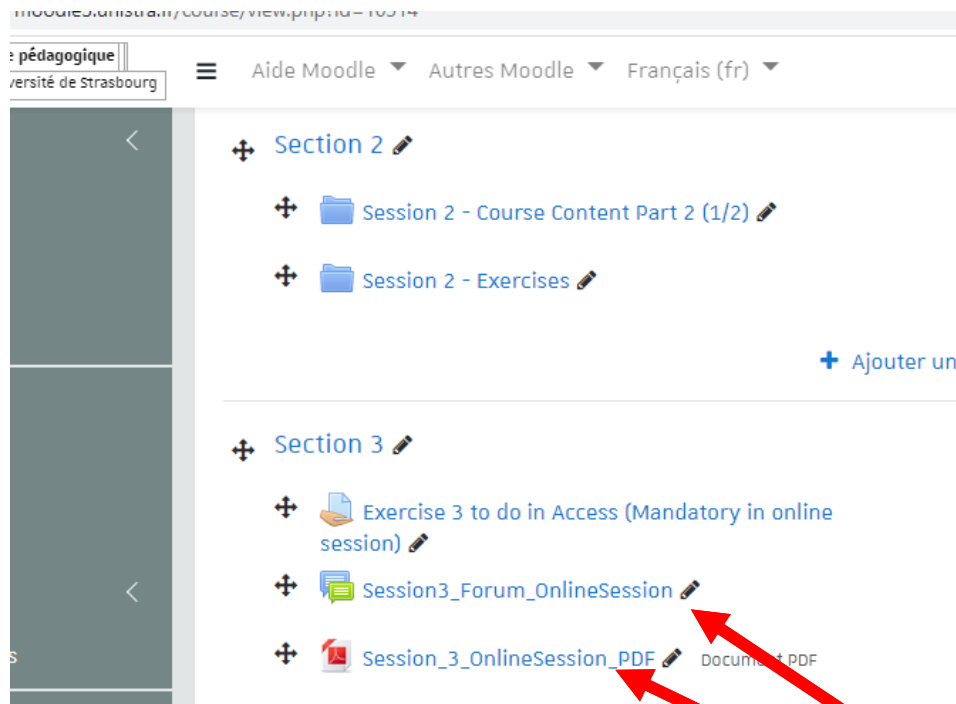
Ajouter une act

Section 5

Now you can download and review the PDF file of this online session in Moodle. If you have any questions, please use the Forum platform in Moodle to follow exchanges and ask your questions. Your professor will exchange with you through this Forum to provide you with answers about this session.



- Individual Work. Duration 60 minutes. Mandatory
Each student must create the Database Cinemas under MS Acces



Now you can download and review the PDF file of this online session in Moodle. If you have any questions, please use the Forum platform in Moodle to follow exchanges and ask your questions. Your professor will exchange with you through this Forum to provide you with answers about this session.



Bibliography

- Book: 'Data Analytics Made Accessible'. 2018. by Anil Maheshwari
- Book: 'Too Big to Ignore: The Business Case for Big Data'. by award-winning
- Book: 'Data Smart: Using Data Science to Transform Information into Insight', by J. W. Foreman'.
- Paper: 'Almeida, F. (2018). Big Data: Concept, Potentialities and Vulnerabilities'. Emerging Science Journal, 2(1).
- McAfee, A., Brynjolfsson, E., Davenport, T. H., Patil, D. J., & Barton, D. (2012). Big data: the management revolution. Harvard business review, 90(10), 60-68.
- Zikopoulos, P., & Eaton, C. (2011). Understanding big data: Analytics for enterprise class hadoop and streaming data. McGraw-Hill Osborne Media.
- Kwon, O., Lee, N., & Shin, B. (2014). Data quality management, data usage experience and acquisition intention of big data analytics. International Journal of Information Management, 34(3), 387-394.

10:00



End of this online session...

