



National University of Engineering (UNI)

School of Computer Science
Syllabus 2024-II

1. COURSE

CS370. Big Data (Mandatory)

2. GENERAL INFORMATION

- 2.1 Course : CS370. Big Data
- 2.2 Semester : 9th Semester.
- 2.3 Credits : 3
- 2.4 Horas : 1 HT; 4 HP;
- 2.5 Duration of the period : 16 weeks
- 2.6 Type of course : Mandatory
- 2.7 Learning modality : Face to face
 - CS272. Databases II. (5th Sem)
- 2.8 Prerequisites :
 - CS3P1. Parallel and Distributed Computing . (8th Sem)

3. PROFESSORS

Meetings after coordination with the professor

4. INTRODUCTION TO THE COURSE

Nowadays, knowing scalable approaches to processing and storing large volumes of information (terabytes, petabytes and even exabytes) is fundamental in computer science courses. Every day, every hour, every minute generates a large amount of information which needs to be processed, stored, analyzed.

5. GOALS

- That the student is able to create parallel applications to process large volumes of information
- That the student is able to compare the alternatives for the processing of big data
- That the student is able to propose architectures for a scalable application

6. COMPETENCES

- 1) Analyze a complex computing problem and apply principles of computing and other relevant disciplines to identify solutions. (Usage)
- 6) Apply computer science theory and software development fundamentals to produce computing-based solutions. (Usage)

7. TOPICS

Unit 1: Introducción a Big Data (15 hours)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none">• Overview on Cloud Computing• Distributed File System Overview• Overview of the MapReduce programming model	<ul style="list-style-type: none">• Explain the concept of Cloud Computing from the point of view of Big Data[Familiarizarse]• Explain the concept of Distributed File System [Familiarizarse]• Explain the concept of the MapReduce programming model[Familiarizarse]
Readings : [Cou+11]	

Unit 2: Hadoop (15 hours)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none"> • Hadoop overview. • History. • Hadoop Structure. • HDFS, Hadoop Distributed File System. • Programming Model MapReduce 	<ul style="list-style-type: none"> • Understand and explain the Hadoop suite [Familiarizarse] • Implement solutions using the MapReduce programming model. [Usar] • Understand how data is saved in the HDFS. [Familiarizarse]
Readings : [HDF11], [BVS13]	

Unit 3: Procesamiento de Grafos en larga escala (10 hours)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none"> • Pregel: A System for Large-scale Graph Processing. • Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud. • Apache Giraph is an iterative graph processing system built for high scalability. 	<ul style="list-style-type: none"> • Understand and explain the architecture of the Pregel project. [Familiarizarse] • Understand the GraphLab project architecture. [Familiarizarse] • Understand the architecture of the Giraph project. [Familiarizarse] • Implement solutions using Pregel, GraphLab or Giraph. [Usar]
Readings : [Low+12], [Mal+10], [Bal+08]	

8. WORKPLAN

8.1 Methodology

Individual and team participation is encouraged to present their ideas, motivating them with additional points in the different stages of the course evaluation.

8.2 Theory Sessions

The theory sessions are held in master classes with activities including active learning and roleplay to allow students to internalize the concepts.

8.3 Practical Sessions

The practical sessions are held in class where a series of exercises and/or practical concepts are developed through problem solving, problem solving, specific exercises and/or in application contexts.

9. EVALUATION SYSTEM

***** EVALUATION MISSING *****

10. BASIC BIBLIOGRAPHY

- [Bal+08] Shumeet Baluja et al. "Video Suggestion and Discovery for Youtube: Taking Random Walks Through the View Graph". In: *Proceedings of the 17th International Conference on World Wide Web*. WWW '08. Beijing, China: ACM, 2008, pp. 895–904. DOI: 10.1145/1367497.1367618. URL: <http://doi.acm.org/10.1145/1367497.1367618>.
- [Mal+10] Grzegorz Malewicz et al. "Pregel: A System for Large-scale Graph Processing". In: SIGMOD '10 (2010), pp. 135–146. DOI: 10.1145/1807167.1807184. URL: <http://doi.acm.org/10.1145/1807167.1807184>.
- [Cou+11] George Coulouris et al. *Distributed Systems: Concepts and Design*. 5th. USA: Addison-Wesley Publishing Company, 2011.
- [HDF11] Kai Hwang, Jack Dongarra, and Geoffrey C. Fox. *Distributed and Cloud Computing: From Parallel Processing to the Internet of Things*. 1st. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2011.

- [Low+12] Yucheng Low et al. “Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud”. In: *Proc. VLDB Endow.* 5.8 (Apr. 2012), pp. 716–727. DOI: 10 . 14778 / 2212351 . 2212354. URL: <http://dx.doi.org/10.14778/2212351.2212354>.
- [BVS13] Rajkumar Buyya, Christian Vecchiola, and S. Thamarai Selvi. *Mastering Cloud Computing: Foundations and Applications Programming*. 1st. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2013.