

Plano de Ensino

Inteligência na Web e Big Data

Fabricio Olivetti de França e Thiago Ferreira Covões
folivetti@ufabc.edu.br, thiago.covoes@ufabc.edu.br

Centro de Matemática, Computação e Cognição
Universidade Federal do ABC



Sobre os Docentes

Formado em Engenharia Elétrica com ênfase em Computação pela Universidade Católica de Santos.

Mestrado e Doutorado no mesmo curso na Universidade Estadual de Campinas.

- Co-agrupamento de Dados
- Regressão Simbólica
- Algoritmos Evolutivos
- Teoria das Categorias Aplicada
- Aprendizado de Máquina

Formado em Ciência da Computação pela Universidade Católica de Santos.

Mestrado e Doutorado no mesmo curso na Universidade de São Paulo - Campus São Carlos.

- Aprendizado de máquina
- Mineração de dados
- Algoritmos evolutivos

Sobre o Curso

O avanço da tecnologia nos permitiu obter dados em massa de diversas fontes:

- Transações Bancárias
- Dados de medidas de sensores
- Experimentos genéticos
- Uso de Redes Sociais
- Construção de conteúdo colaborativo

Grande quantidade de dados cria a necessidade de armazenamento.

- Maior HD tem cerca de 16TB por \$7.000,00
- Dados na ordem de PetaBytes!!

Grande quantidade de dados cria a necessidade de armazenamento.

- Maior HD tem cerca de 16TB por \$7.000,00
- Dados na ordem de PetaBytes!!

Solução: Sistema de Arquivos Distribuídos

Crescimento dos Dados



Figura 1: Fonte: datacenterknowledge.com

Desses dados surge a necessidade de extrair informações úteis:

- Possíveis dívidas bancárias
- Genes relacionados a doenças
- Previsão de desastres naturais

- Muitos algoritmos consolidados
- Foco em bases de dados pequenas e bem estruturadas
- Bases de dados pequenas em apenas uma máquina!

Temos:

- Bases de Dados grandes e distribuídas em diversas máquinas
- Dados brutos, não estruturados
- Sem a mínima ideia de como trabalhar com eles!

Objetivo do Curso

- Conhecimento básico sobre Ciência dos Dados
- Paradigma funcional / Pensamento Paralelo
- Enfoque em algoritmos escaláveis e distribuídos

<http://folivetti.github.io/teaching/2019-spring-teaching-2>

- Aulas expositivas
- Atividades de Programação
- Projetos - Seminários

No site.

Livros sobre Mineração de Dados:

Mining Massive Datasets

Introduction to Data Mining – P-N. Tan, M. Steinbach, V.Kumar
– Addison Wesley 2005.

Data mining : practical machine learning tools and techniques /
3. ed.

Livros sobre Apache Spark:

<https://www.gitbook.com/book/jaceklaskowski/mastering-apache-spark>

ODERSKY, Martin et al. Programming in Scala. Mountain View, USA: Artima Press, 2008. xxxix, 736. ISBN 9780981531601.

Livros sobre Python:

LUTZ, Mark; ASCHER, David; ALYH69. Aprendendo python.
2. ed. Porto Alegre, RS: Bookman, 2007. 566 p. ISBN
857780013x.

[https://spark.apache.org/docs/2.2.0/
rdd-programming-guide.html](https://spark.apache.org/docs/2.2.0/rdd-programming-guide.html)

- 02 provas cada uma
- 01 projeto em grupo

Cada avaliação valerá 10 ptos e a média final será calculada via média harmônica.

O projeto deve ser entregues via Github Classroom (link no site da disciplina)

Cada aluno deve escolher um tema de projeto relacionado com a disciplina.

O projeto deve ser entregue em forma de relatório científico com até 3 páginas e códigos.

Apresentação de 10 a 15 minutos.

Possíveis temas:

- Implementação distribuída de um algoritmo.
- Implementação de um algoritmo escalável.

Deve envolver uma quantidade suficiente de DADOS!

O conceito final será calculado da seguinte forma:

```
1 conceito :: Double -> Char
2 conceito nota
3   | nota >= 9 = 'A'
4   | nota >= 8 = 'B'
5   | nota >= 6 = 'C'
6   | nota >= 5 = 'D'
7   | otherwise = 'F'
```
