

Introdução

Mineração de Dados

Universidade Federal do ABC

Introdução

Mineração de Dados

DADOS, DADOS E MAIS DADOS

“Estamos nos afogando em informação, mas famintos por conhecimento” - John Naisbitt

- ▶ Progresso na coleta e armazenamento de dados tornaram comuns bases de dados enormes
 - ▶ Desde o início dos anos 90
 - ▶ Não há indícios de que isso irá parar tão cedo
 - ▶ Como tirar o melhor proveito possível dos dados

DADOS, DADOS E MAIS DADOS



Figura 1: Comércio Eletrônico

- ▶ O que você procura
- ▶ O que você compra
- ▶ O que você critica

DADOS, DADOS E MAIS DADOS



Figura 2: Twitter

- ▶ O que você escreve
- ▶ De onde você escreve
- ▶ Com quem você conversa
- ▶ Quem você lê

DADOS, DADOS E MAIS DADOS



Figura 3: Facebook

- ▶ Quem é você
- ▶ O que você faz
- ▶ Com quem e quando você faz
- ▶ O que você acha sobre o que os outros fazem

DADOS, DADOS E MAIS DADOS



Figura 4: *Smartphones*

- ▶ Com quem você fala
- ▶ Por onde você anda
- ▶ Como você se locomove
- ▶ Sem entrar em detalhes sobre *apps*

DADOS, DADOS E MAIS DADOS



Figura 5: Internet das coisas

- ▶ Acompanhamento 24/7
- ▶ Diversas questões preocupam nessa área
 - ▶ Não iremos abordar esses assuntos *nesta* disciplina

DADOS, DADOS E MAIS DADOS

- ▶ Esses são casos extremos, mas a situação é comum:
 - ▶ Concessionárias
 - ▶ Hospitais
 - ▶ Escolas
 - ▶ Bancos
 - ▶ ...
- ▶ Empresas ainda tomam decisões importantes considerando apenas em intuição
 - ▶ “HiPPO” *the highest-paid person’s opinion*
- ▶ Pessoas se apoiam muito em experiência e intuição mas pouco em dados

Introdução

Mineração de Dados

MINERAÇÃO DE DADOS

- ▶ Amadureceu conforme as bases de dados cresceram em tamanho e complexidade.
- ▶ Disciplina interdisciplinar relacionada a:
 - ▶ Estatística
 - ▶ Bancos de dados
 - ▶ Aprendizado de Máquina
 - ▶ Reconhecimento de Padrões
- ▶ Fronteiras entre as áreas **não** são rígidas
- ▶ Cervejas e fraldas (*fun fact*):
<https://www.kdnuggets.com/news/2000/n14/8i.html>

MINERAÇÃO DE DADOS

- ▶ Inicialmente era conhecida como a etapa de extração de padrões dentro do Processo de Descoberta de Conhecimento em Bases de Dados (*Knowledge Discovery in Databases*)
 - ▶ Com o tempo alguns pesquisadores começaram a utilizar como sinônimos

MINERAÇÃO DE DADOS

Definição (Hand, Manilla & Smith) Mineração de dados é a análise de bases de dados observacionais (frequentemente grandes) para encontrar relações desconhecidas e para sumarizar os dados em formas que sejam compreensíveis e úteis para o dono dos dados.

MINERAÇÃO DE DADOS

- ▶ Alguns termos na definição chamam atenção:
 - ▶ dados observacionais
 - ▶ relações desconhecidas
 - ▶ formas compreensivas
 - ▶ formas úteis

ETAPAS

- ▶ Identificação e formalização do problema
 - ▶ possivelmente a mais difícil
- ▶ Pré-processamento
 - ▶ Costuma ser a que leva mais tempo
- ▶ Extração de Padrões
- ▶ Pós-processamento

CRISP-DM

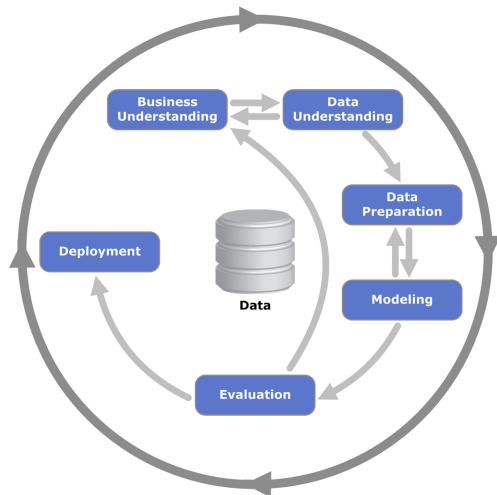


Figura 6: Cross-industry standard process for data mining

IDENTIFICAÇÃO E FORMALIZAÇÃO DO PROBLEMA

- ▶ Quais são as principais metas do processo?
- ▶ Como o desempenho será avaliado?
- ▶ A compreensibilidade do modelo deve ter peso maior do que seu desempenho?
- ▶ Qual o custo de um erro?

PRÉ-PROCESSAMENTO

- ▶ Integração de dados: múltiplas fontes e formatos
- ▶ Transformação: normalização/adequação
- ▶ Limpeza: dados inválidos são frequentes
 - ▶ Por exemplo, pressão sanguínea = 0 [*Pima Indians Diabetes*]
- ▶ Seleção e Redução: formas de “simplificar” os dados, discutiremos em outra aula

EXTRAÇÃO DE PADRÕES

- ▶ Onde a mágica acontece
- ▶ Aplicação de algoritmos de AM, estatística etc.
 - ▶ Tarefas preditivas
 - ▶ Tarefas descritivas
- ▶ Não existe algoritmo que seja melhor para todos os problemas
 - ▶ *No free-lunch theorem* (Wolpert)

PÓS-PROCESSAMENTO

- ▶ O que o modelo achou que é interessante?
- ▶ Apresentação do modelo de forma compreensível

ÁREAS RELACIONADAS

Business Intelligence

Data Science

Big Data

Predictive Analytics

...

O QUE NÃO IREMOS COBRIR

- ▶ Extração de dados
 - ▶ *Crawlers*
- ▶ Armazenamento e organização de dados
 - ▶ *Data Warehouse*
 - ▶ *Online Analytic Processing* (OLAP)
- ▶ Pós-processamento

REFERÊNCIAS

- ▶ D. Hand, H. Manilla e P. Smith. Principles of Data Mining. **Capítulo 1**
- ▶ S. Rezende. Sistemas Inteligentes: Fundamentos e Aplicações. **Capítulo 12**
- ▶ P. Tan, M. Steinbach e V. Kumar, Introduction to Data Mining. **Capítulo 1**
- ▶ Breiman, Leo. Statistical Modeling: The Two Cultures (with comments and a rejoinder by the author). Statist. Sci. 16 (2001), no. 3, 199–231. doi:10.1214/ss/1009213726. <https://projecteuclid.org/euclid.ss/1009213726>