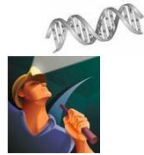


# Mineração de Dados em Biologia Molecular

KDD



Docente: André C. P. L. F. de Carvalho  
PAE: Victor Hugo Barella

## Tópicos do Módulo

- Introdução
- Descoberta de Conhecimento em Bases de Dados
- Etapas de KDD
- Mineração de Dados
- Aplicações

© André de Carvalho - ICMC/USP

2

## Introdução

- Bases de Dados podem conter (esconder) dados preciosos
- Existe um interesse crescente em explorar esses dados armazenados
  - Descobrir conhecimento novo
  - Apoio à tomada de decisão

© André de Carvalho - ICMC/USP

3

## O que é um dado?

- Dado
  - Exemplo, objeto, registro
  - Estrutura fundamental sobre a qual um sistema de informação é construído
  - Não tem um significado associado a ele
    - Ex.: 345,43

© André de Carvalho - ICMC/USP

4

## Exemplo

muito alto,muito alto,2,2,pequeno,baixo,ruim  
 muito alto,alto,3,5 ou mais,grande,baixo,ruim  
 muito alto,baixo,3,4,grande,baixo,ruim  
 médio,baixo,4,2,pequeno,alto,ruim  
 médio,baixo,3,4,pequeno,médio,médio  
 alto,alto,2,4,grande,médio,médio  
 baixo,baixo,5 ou mais,4,pequeno,médio,médio  
 baixo,médio,4,4,pequeno,médio,médio  
 baixo,médio,4,4,grande,médio,bom  
 baixo,baixo,4,5 ou mais,grande,médio,bom  
 médio,baixo,2,4,pequeno,alto,bom  
 baixo,médio,4,4,grande,alto,muito bom  
 médio,médio,2,4,grande,alto,muito bom  
 baixo,baixo,5 ou mais,5 ou mais,grande,alto,muito bom

© André de Carvalho - ICMC/USP

5

## Informação

- Dados com um significado associado
  - Tornam os dados úteis em uma tomada de decisão
  - Transformação de dados em informação
    - Geralmente pela apresentação dos dados em uma forma compreensível para o usuário
  - Informação é criada quando é associado um significado a um conjunto de dados

© André de Carvalho - ICMC/USP

6

## Exemplo - Carros

- Preço
  - Compra: muito-alto, alto, médio, baixo
  - Manutenção: muito-alto, alto, médio, baixo
- Características técnicas
  - Conforto
    - # portas: 2, 3, 4, 5-5 ou mais
    - # pessoas: 2, 4, 5 ou mais
    - Espaço do porta malas: pequeno, médio, grande
  - Segurança: baixo, médio, alto
- Aval. do carro: ruim, médio, bom, muito bom

© André de Carvalho - ICMC/USP 7

## Exemplo - Carros

- Preço
  - Compra: muito-alto, alto, médio, baixo
  - Manutenção: muito-alto, alto, médio, baixo
- Características técnicas
  - Conforto
    - # portas: 2, 3, 4, 5-5 ou mais
    - # pessoas: 2, 4, 5 ou mais
    - Espaço porta malas: pequeno, médio, grande
  - Segurança: baixo, médio, alto
- Aval. do carro: ruim, médio, bom, muito bom

© André de Carvalho - ICMC/USP 8

## Exemplo

```
tactagcaatacgttgcggtcggttaagtgtataatgacggtgctgtcgt
fgctatcctgacagttgcaagctgattgggtgcgttaacatcaacgcatcgcaa
gtactagagaactagtgacattatctttttttatcatgtagcggcg
aattgtgatgtatcgaagtggtgtagagtagatgtagaatacaaaactc
Tcgataaatactattgacgaaagctgaagactagaatgcctccggtgtag
aggggcaaggaggtggaagaggtgtagcgtataaagaactagagctcgttagt
agggatgaaactctctgtagcactagctcttctactgtgagtagcggag
cggatgagcttagagagcatgtagcctgacaactgcataaagtcttctg
cgctaggacttctgtgatttccatgctggttttgcgcaatgtaacgcttt
tatgggaacgagtcactcagggtctgactctgttactgtgaacattatt
agaggggtgactccaagaaggaagatgaggctagacgtctctgcatggatgga
gagagcatgtagcctcgacaactgcataaagtcttctgtagacgtgcctcag
```

© André de Carvalho - ICMC/USP 9

## Exemplo - Promotores

```
+ ,S10,tactagcaatacgttgcggtcggttaagtgtataatgacggtgctgtcgt
+ ,AMPc,tgctatcctgacagttgcaagctgattgggtgcgttaacatcaacgcatcgcaa
+ ,AROH,gtactagagaactagtgacattatctttttttatcatgtagcggcg
+ ,DEOP2,aattgtgatgtatcgaagtggtgtagagtagatgtagaatacaaaactc
+ ,LEU1_TRNA,tcgataaatactattgacgaaagctgaagactagaatgcctccggtgtag
+ ,MALEFG,aggggcaaggaggtggaagaggtgtagcgtataaagaactagagctcgttagt
-, 296,agggatgaaactctctgtagcactagctcttctactgtgagtagcggag
-, 648,ccgagtaggcttagagagcatgtagcctgacaactgcataaagtcttctg
-, 230,cgctaggacttctgtgatttccatgctggttttgcgcaatgtaacgcttt
-, 1163,tatgggaacgagtcactcagggtctgactctgttactgtgaacattatt
-, 1321,agaggggtgactccaagaaggaagatgaggctagacgtctctgcatggatgga
-, 663,gagagcatgtagcctcgacaactgcataaagtcttctgtagacgtgcctcag
```

© André de Carvalho - ICMC/USP 10

## KDD

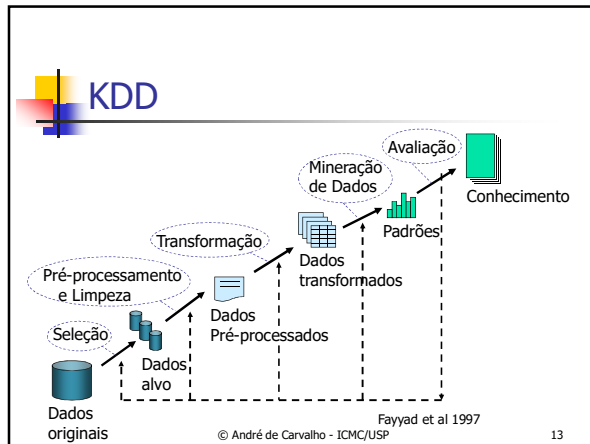
- Descoberta de conhecimento em Bancos de dados (BDs)
  - Knowledge Discovery in Databases
- Área de pesquisa em expansão
- Teorias e ferramentas computacionais para extrair informação útil de BDs
  - Informação útil = conhecimento

© André de Carvalho - ICMC/USP 11

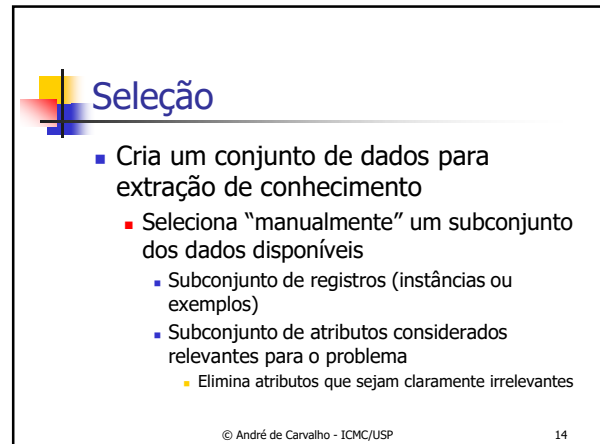
## KDD

- Processo de encontrar em dados padrões que são
  - Úteis
  - Novos
  - Válidos
  - Potencialmente compreensíveis
- Processo iterativo e iterativo
  - Várias etapas

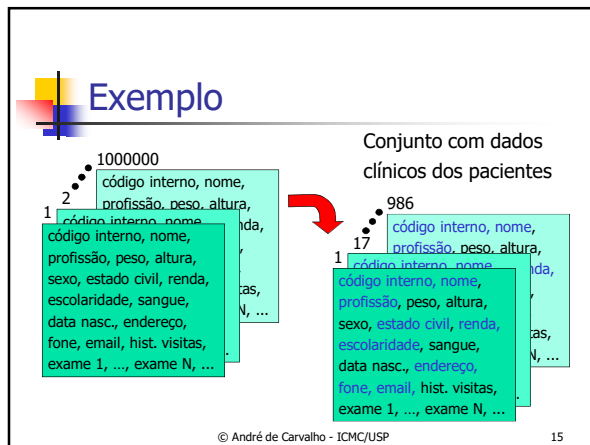
© André de Carvalho - ICMC/USP 12



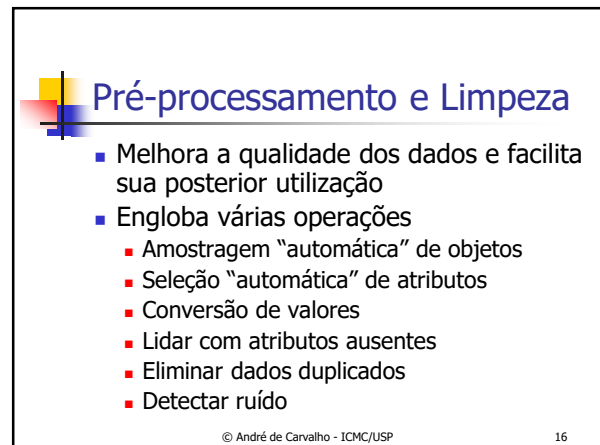
13



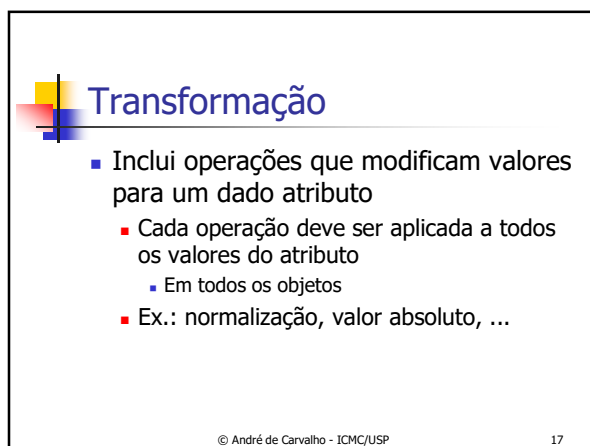
14



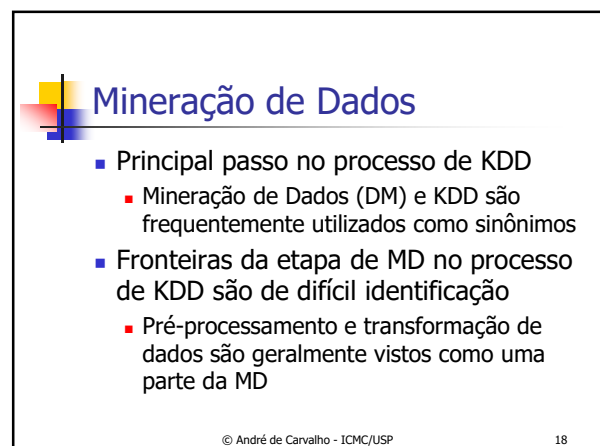
15



16



17



18

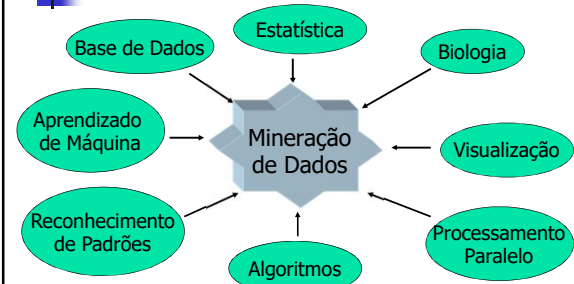
## Mineração de Dados

- MD X KDD
  - MD: ferramentas básicas utilizadas para extrair padrões de dados
  - KDD: processo que engloba o uso dessas ferramentas, além de:
    - Pré-processamento, seleção e transformação dos dados
    - Interpretação dos padrões
      - Geração de conhecimento
      - Suporte à tomada de decisão

© André de Carvalho - ICMC/USP

19

## Mineração de Dados



© André de Carvalho - ICMC/USP

20

## Mineração de Dados

- Outros termos utilizados para MD e KDD
  - Extração de conhecimento
  - Descoberta de informação
  - Extração de padrões
  - Análise exploratória de dados
  - Analítica (Data analytics ou analytics)

© André de Carvalho - ICMC/USP

21

## Analítica

- Ciência que analisa dados crus para extrair padrões desses dados
  - Pode englobar coleta e organização dos dados
- Analítica preditiva (predictive analytics)
  - Extraí modelos a partir de dados para fazer previsões futuras
- Analítica descritiva (descriptive analytics)
  - Sumariza ou condensa dados para extrair informações

© André de Carvalho - ICMC/USP

22

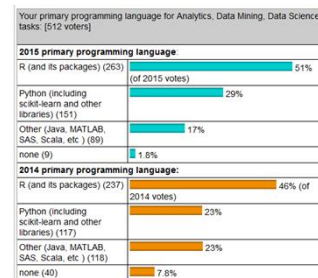
## Interpretação / Avaliação

- Interpretação dos padrões encontrados na etapa de MD
  - Possível retorno a qualquer uma das etapas anteriores para iteração adicional
- Valida padrões encontrados
  - Importante consulta a um especialista
- Inclui análise estatística
- Ferramentas de visualização fornecem um suporte importante

© André de Carvalho - ICMC/USP

23

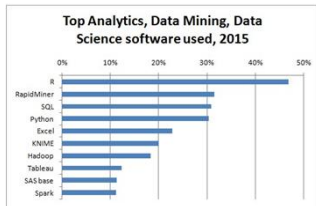
## Ferramentas



© André de Carvalho - ICMC/USP

24

## Ferramentas



© André de Carvalho - ICMC/USP

25

## Custos em MD Preditivo

- 15% - coleta de dados
- 60% - limpeza e pré-processamento de dados
- 15% - construção e análise de modelos
- 5% - aplicação
- 5% - melhorias contínuas

© André de Carvalho - ICMC/USP

26

## CRISP-DM

- Projeto CRISP-DM
  - *C*Ross-*I*ndustry *S*tandard *P*rocess for *D*ata *M*ining
  - Concebido em 1996 por:
    - Daimler-Chrysler
      - Aplicava MD em suas operações de negócios
    - SPSS
      - Prestava serviço de MD desde 1990
      - Desenvolveu primeira ferramenta comercial de MD (*Clemetine*)
    - NDR
      - Tinha o propósito de adicionar valor a sua enorme BD

© André de Carvalho - ICMC/USP

27

## CRISP-DM

- Projeto CRISP-DM
  - Desenvolveu um novo fluxo de processo para descoberta de conhecimento
    - A partir do processo anterior
      - Fayyad, Piatetsky-Shapiro and Smyth
  - Em resposta a requisitos de usuários
    - Definiu e validou processo de MD utilizado em vários setores industriais

© André de Carvalho - ICMC/USP

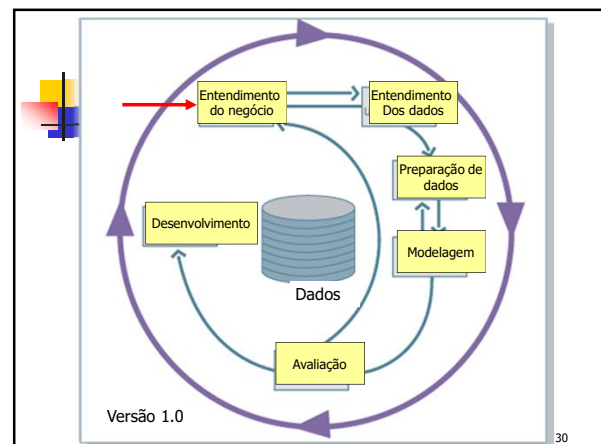
28

## CRISP-DM

- Metodologia procura tornar os projetos
  - Mais rápidos
  - Mais baratos
  - Mais confiáveis
  - Mais facilmente gerenciáveis
- Pode ser aplicada a pequenos projetos
- Metodologia padrão da indústria

© André de Carvalho - ICMC/USP

29

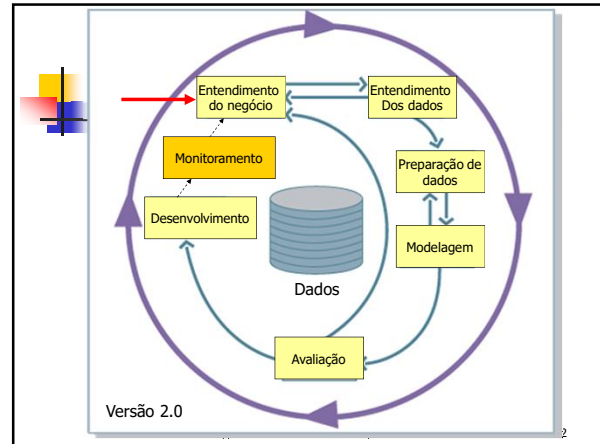


30

## CRISP-DM 2.0

- SIG (*special interest group*) formado entre 2006 e 2008
- Mudanças estudadas
  - Divisão da fase de preparação de dados
  - Métodos de avaliação dentro da fase de modelagem
  - Fase de avaliação associada a avaliação na empresa
    - Centro de pesquisa, laboratório, hospital,...
  - Inclusão de fase de monitoramento

© André de Carvalho - ICMC/USP 31



## CRISP-DM 2.0

- Não foi lançado até 2005
- SIG foi desfeito
  - Website CRISP-DM.org não esta mais ativo
- IBM criou nova metodologia
  - Refina e estende CRISP-DM
  - Analytics Solutions Unified Method for Data Mining/Predictive Analytics (ASUM-DM)

© André de Carvalho - ICMC/USP 33

## Produtos para MD

© André de Carvalho - ICMC/USP 34


## Mais Produtos

© André de Carvalho - ICMC/USP 35


## Considerações Finais

- Expansão do volume de dados armazenados
  - Necessidade de extrair conhecimento dos dados
- KDD é cada vez mais usado por órgãos privados e públicos
- Cuidado com promessas exageradas
- Leitura
  - Knowledge Discovery and Data Mining: Towards a Unifying Framework, U. Fayyad, P. Smyth, and G. Piatetsky-Shapiro, .2nd International Conference on Knowledge Discovery and Data Mining, 1996

© André de Carvalho - ICMC/USP 36

 Perguntas

---



© André de Carvalho - ICMC/USP 37