

SCC0173 – Mineração de Dados Biológicos

Aula 1 – Conceituação

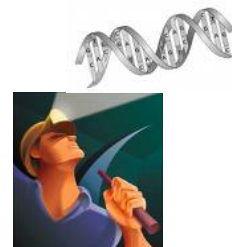
Prof. Ricardo J. G. B. Campello


Estagiário PAE: Pablo A. Jaskowiak

SCC / ICMC / USP

Créditos

- O material a seguir consiste de adaptações e extensões dos originais gentilmente cedidos:
 - pelo Prof. André C. P. L. F. de Carvalho
 - e seu estagiário PAE Ricardo Cerri
 - pelo Prof. João Luis Garcia Rosa






Resumo

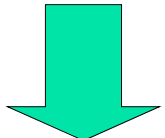
- Introdução
- Descoberta de Conhecimento em Bases de Dados (KDD)
- Etapas de KDD
- Mineração de Dados
- Aplicações

André Ponce de Leon F de Carvalho 3



Introdução

- Avanços recentes nas tecnologias de aquisição, transmissão e armazenamento de dados



Bases de dados cada vez maiores

André Ponce de Leon F de Carvalho 4

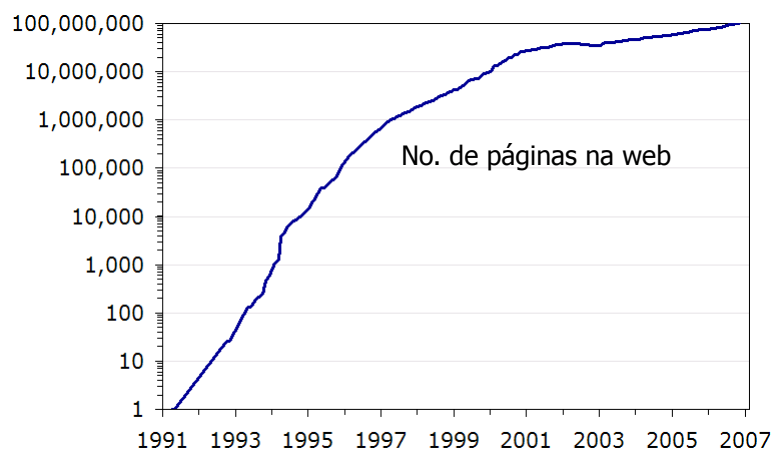
Introdução

- Estima-se que a quantidade de dados em Bases de Dados mundiais dobra a cada 20 meses
 - Transações bancárias
 - Utilização de cartões de crédito
 - Dados governamentais
 - Medições ambientais
 - Dados clínicos
 - Informações disponíveis na web
 - Projetos genoma

André Ponce de Leon F de Carvalho

5

Introdução

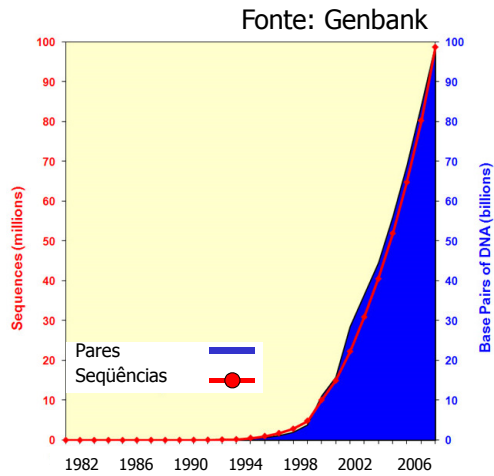


André Ponce de Leon F de Carvalho

6

Introdução

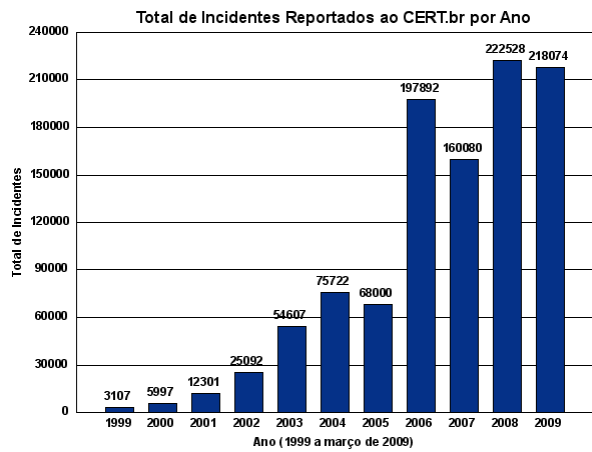
Crescimento do GenBank
1982-2009



André Ponce de Leon F de Carvalho

Introdução

Incidentes de
Segurança na
Internet Brasileira



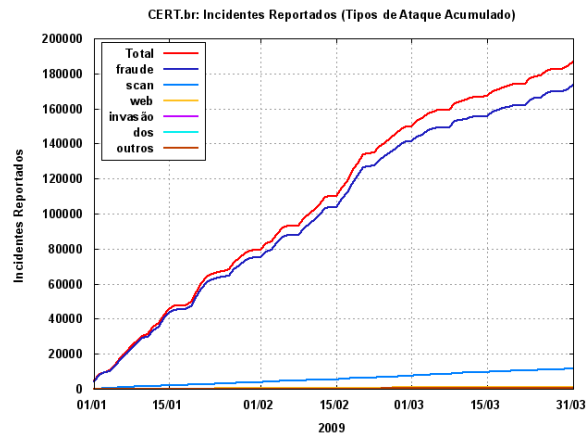
André Ponce de Leon F de Carvalho

8

Introdução

Incidentes de Segurança na Internet Brasileira

Janeiro a Março 2009



André Ponce de Leon F de Carvalho

9

Introdução

- Grandes bases de dados não são incomuns:
 - Informação de pontos de venda, registros governamentais, registros médicos e dados de cartões de crédito, etc.
 - Instrumentos científicos podem produzir terabytes (10^{12} ou 2^{40}) e petabytes (10^{15} ou 2^{50}) a taxas de gigabytes (10^9 ou 2^{30}) por hora
 - Capacidades de armazenamento aumentam
 - Bases de dados mais baratas e maiores, crescendo tanto em tamanho de campo e registro como em número de registros
- Limites humanos!

Slide adaptado do original do Prof. João L. G. Rosa

10



Introdução

- Dois exemplos:
 - Transações eletrônicas
 - BD Wal-Mart: 7.2 bilhões de diferentes transações de compra em 2008
 - Controle e monitoramento
 - BD NASA: recebe de satélites ~ 50 GB / hora

André Ponce de Leon F de Carvalho

11



Introdução

- Bases de Dados muito grandes podem conter (esconder) dados preciosos
- Existe um interesse crescente em explorar esses dados armazenados
 - Descobrir conhecimento novo e útil
 - Suporte a decisão

André Ponce de Leon F de Carvalho

12



Introdução

- Técnicas tradicionais de análise de dados permitem apenas consultas simples
 - P. ex., “Quantos itens de um produto em particular foram vendidos em um dado dia?”
 - Não conseguem responder consultas do tipo:
 - Que tecidos podem estar com tumor?
 - Qual a estrutura terciária de uma nova proteína
 - Técnicas mais sofisticadas, capazes de extrair **conhecimento** de grandes BD, são necessárias

André Ponce de Leon F de Carvalho

13

Conhecimento

Compreensão ou Modelo de um determinado Domínio ou Área

Dados – Informação – Conhecimento

- **Dados:**
 - quantificações de um determinado evento
- **Informação:**
 - dados contextualizados
- **Conhecimento:**
 - (modelo das) relações existentes entre as informações contidas nos dados e do seu significado

14

KDD

- **Descoberta de Conhecimento em Bases de Dados (KDD)**
 - Extração de **conhecimento** a partir de:
 - Registros de compras em grandes supermercados
 - Registros de empréstimos financeiros
 - Registros de transações financeiras
 - Registros médicos
 - Sequências de expressão-gênica
 - ...
 - Engloba a **Mineração de Dados** (Data Mining – DM)

André Ponce de Leon F de Carvalho 15

KDD e DM

- DM se refere especificamente às ferramentas (descritivas ou preditivas) de processamento de dados para extração de conhecimento:

```

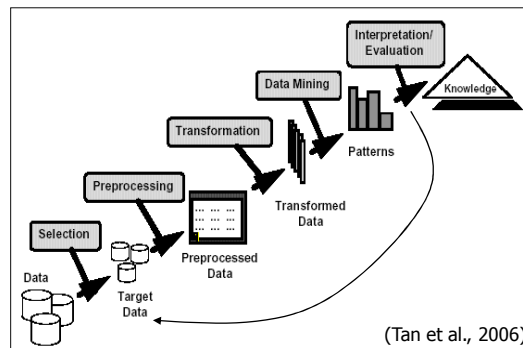
graph TD
    A[Aprendizado de Máquina] --> B[Data Mining e KDD]
    C[Estatística] --> B
    D[Banco de Dados Data Warehouse Information Retrieval] --> B
    E[Visualização de Dados] --> B
    B <--> F[Computação Paralela e Distribuída]
  
```

16

KDD e DM

- KDD é usualmente definido como um **ciclo** que envolve DM:
 - Formalmente, é o processo *não trivial* de identificar padrões **válidos, novos, potencialmente úteis e compreensíveis** em dados

Fayyad, U. M., Piatsky-Shapiro, G. and Smyth, P., "From Data Mining to Knowledge Discovery: An Overview" In **Advances In Knowledge Discovery and Data Mining**. AAAI/MIT Press, p. 1-34, 1996.



KDD e DM

- DM e KDD podem ser vistos como uma evolução dos processos tradicionais de análise de BDs:

Passo evolutivo	Tecnologias existentes	Provedores de produtos	Características
<i>Coleção de Dados</i> (1960s)	Computadores, fitas, discos.	IBM, CDC.	Retrospectivo, fornecimento estático de dados.
<i>Acesso aos Dados</i> (1980s)	RDBMS, SQL, ODBC.	Oracle, IBM, Microsoft.	Retrospectivo, fornecimento dinâmico de dados no nível do registro.
<i>Data Warehousing</i> (1990s)	Relational Data Warehouse, OLAP, MDDB.	NCR, Business Objects, COGNOS, Hyperion.	Retrospectivo, fornecimento dinâmico de dados em vários níveis.
<i>Data Mining</i> (2000s)	Algoritmos avançados, Bases de Dados muito grandes.	SAS, SPSS, IBM, Oracle, NCR.	Prospectivo, informação proativa e fornecimento de conhecimento.

Slide adaptado do original do Prof. João L. G. Rosa

18




KDD e DM



- Podem tratar com terabytes, ou mais...
- Hoje, petabytes não são incomuns...
- Quão grande é um petabyte?
 - 250 bilhões de páginas de texto,
 - 20 milhões de arquivos de aço de 4 gavetas,
 - 1 torre de mais de 3.000 km de altura de 1 bilhão de disquetes:
 - O Pico da Neblina, a mais alta montanha do Brasil, tem quase 3 km
 - foto: <http://www.ecotour.nl/Reizen/Expeditie%20Pico%20da%20Neblina/neb4.jpg>.

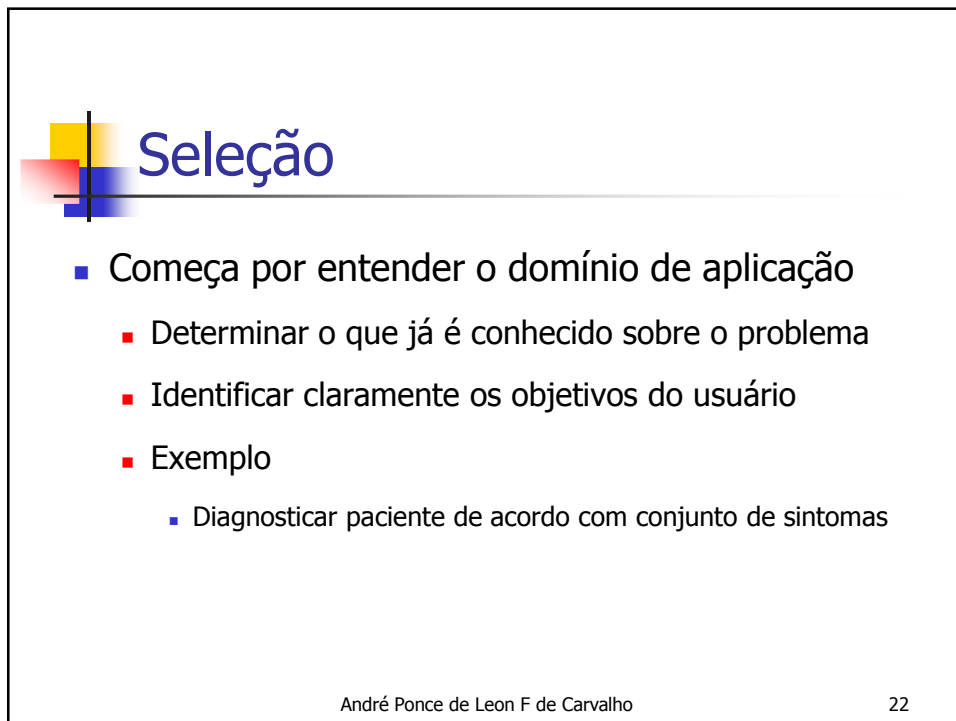
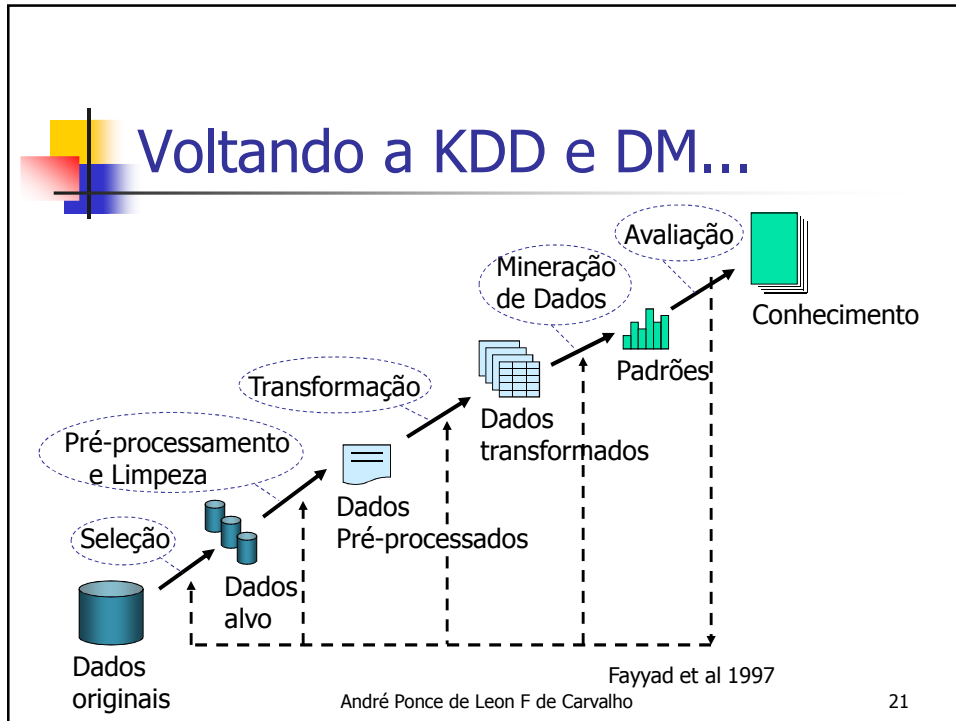
Slide adaptado do original do Prof. João L. G. Rosa 19



Tamanho dos Dados

- Tamanhos de conjuntos de dados
 - Pequeno
 - Conjunto de dados pode ser gerenciado pela ferramenta de KDD sozinha, geralmente em um arquivo e computador único
 - Médio
 - Necessária a integração do ambiente de KDD com Sistemas Gerenciadores de BDs (SGBDs), que gerenciam os dados
 - Grande
 - Quando o volume de dados é grande demais para ser gerenciado pelas ferramentas de um SGBD
 - Necessário sistemas mais sofisticados

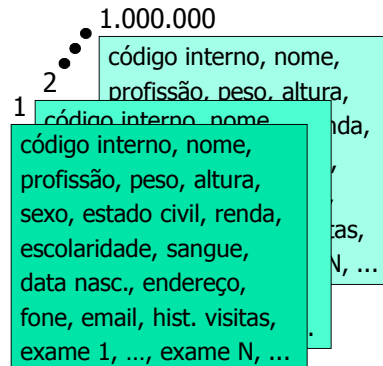
André Ponce de Leon F de Carvalho 20



Exemplo

- BD de um hospital
 - Composto por conjunto de registros de pacientes
 - Cada registro é composto de atributos
 - Informações pessoais
 - Sintomas

BD com registros de pacientes



André Ponce de Leon F de Carvalho

23

Seleção

- Toma-se então um subconjunto de interesse dentre os dados disponíveis
 - Subconjunto de registros
 - exemplos, instâncias ou objetos
 - Subconjunto de atributos
 - campos, variáveis ou características
 - considerados relevantes para o problema
 - os demais são claramente irrelevantes e descartados "manualmente"

André Ponce de Leon F de Carvalho

24

Exemplo

1.000.000

1 2

código interno, nome, profissão, peso, altura, ...

Conjunto com dados clínicos dos pacientes

986

1 17

código interno, nome, profissão, peso, altura, sexo, estado civil, renda, escolaridade, sangue, data nasc., endereço, fone, email, hist. visitas, exame 1, ..., exame N, ...

André Ponce de Leon F de Carvalho 25

Relembrando...

```

    graph TD
      A[Dados originais] -- Seleção --> B[Dados Pré-processados]
      B -- "Pré-processamento e Limpeza" --> B
      B -- Transformação --> C[Dados transformados]
      C -- "Mineração de Dados" --> D[Padrões]
      D -- Avaliação --> E[Conhecimento]
      A -.-> C
      B -.-> D
      C -.-> E
  
```

Fayyad et al 1997
André Ponce de Leon F de Carvalho 26



Pré-Processamento e Limpeza

- Melhorar a qualidade dos dados e facilitar sua posterior utilização
- As principais operações são
 - Eliminar dados duplicados
 - Lidar com valores ausentes
 - Lidar com ruído
 - p. ex. em imagens de fenômenos biológicos

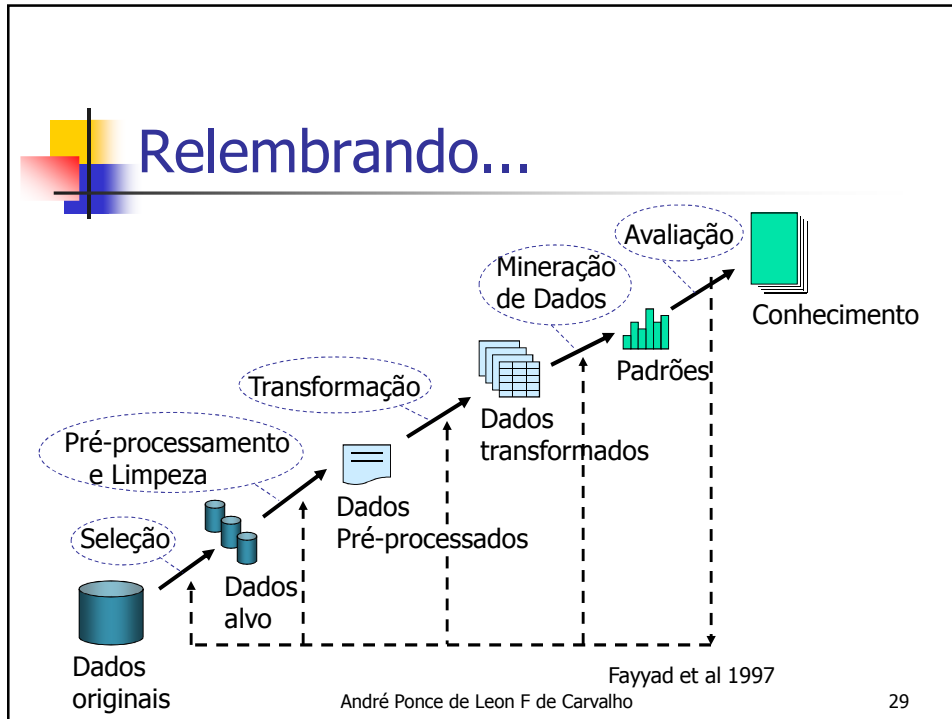
27



Transformação

- Inclui diversas operações, tais como:
 - Seleção "automática" de atributos
 - para eliminar atributos irrelevantes e/ou redundantes
 - Extração de características
 - para obter um menor no. de atributos mais relevantes a partir dos dados brutos
 - Discretizações e conversões
 - entre atributos de diferentes naturezas
 - Normalizações
 - para que os atributos ou mesmo os objetos apresentem determinadas propriedades de interesse
 - tipicamente propriedades estatísticas

28



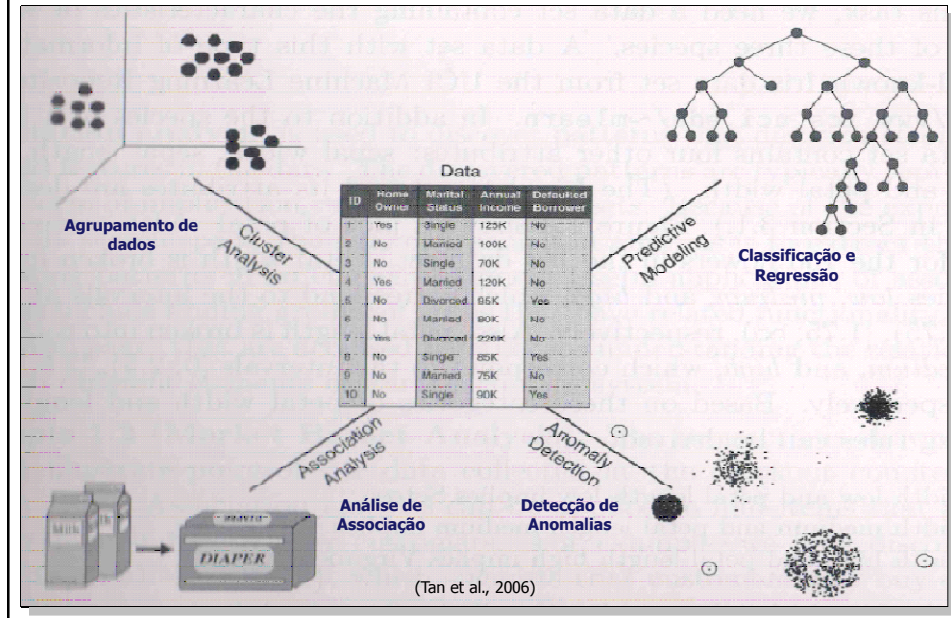
Mineração de Dados (DM)

- Principal passo no processo de KDD
 - DM e KDD são frequentemente utilizados como sinônimos
- Fronteiras da etapa de DM no processo de KDD são de difícil identificação
 - Pré-processamento e transformação de dados são frequentemente vistos como uma parte de DM

André Ponce de Leon F de Carvalho

30

Tarefas Fundamentais de DM



Interpretação / Avaliação

- Interpretação dos padrões obtidos na etapa de DM
 - Verifica se os padrões encontrados são:
 - válidos, novos, úteis e compreensíveis
 - Uma das etapas mais complexas
 - Em geral inclui análise estatística
 - Mas pode envolver elevado grau de subjetividade em algumas aplicações
 - Normalmente requer muito conhecimento de domínio
 - Possível retorno a qualquer uma das etapas anteriores
- Ferramentas de visualização de dados podem ter um papel de suporte importante

<http://www.crisp-dm.org>



CRISP-DM

- Projeto CRISP-DM
 - *Cross-Industry Standard Process for Data Mining*
 - Concebido em 1996 por 3 empresas:
 - Daimler-Chrysler
 - Aplicava DM em suas operações de negócios
 - SPSS
 - Prestava serviço de DM desde 1990
 - Desenvolveu primeira ferramenta comercial de DM (*Clementine*)
 - NCR
 - Tinha o propósito de adicionar valor aos seus produtos e negócios em Data Warehousing

André Ponce de Leon F de Carvalho

33



CRISP-DM

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

André Ponce de Leon F de Carvalho

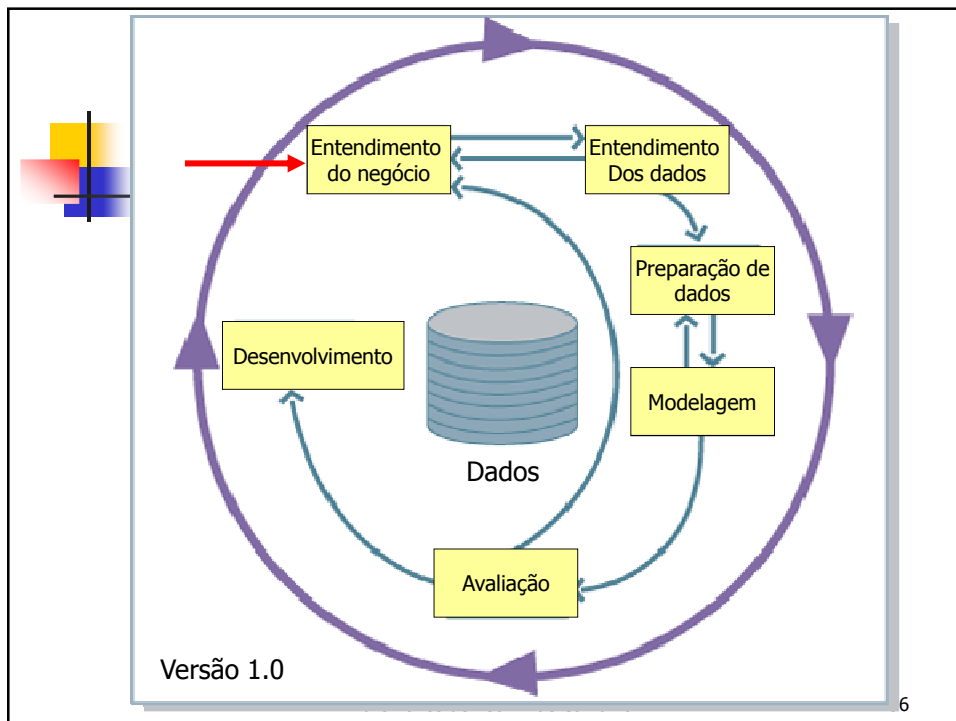
34

CRISP-DM

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

André Ponce de Leon F de Carvalho

35





Entendimento do Negócio

- Entender os objetivos do projeto
 - e requisitos do negócio (problema)
 - definir critérios para medir sucesso
- Converter o entendimento em definição de um problema de DM
- Traçar um planejamento preliminar para atingir objetivos

André Ponce de Leon F de Carvalho

37



Entendimento dos Dados

- Tem início com uma coleta dos dados
- Segue com atividades para:
 - Familiarizar-se com os dados
 - Calcular estatísticas básicas
 - Explorar os dados
 - Investigar atributos
 - Identificar problemas de qualidade nos dados
 - Realizar descobertas iniciais sobre os dados
 - Detectar subconjuntos interessantes

André Ponce de Leon F de Carvalho

38



Preparação dos Dados

- Cobre todas as atividades necessárias para construir o conjunto de dados final
 - Provavelmente executadas várias vezes
 - Sem uma ordem pré-definida
 - Inclui:
 - Seleção de tabelas, instâncias e atributos
 - Limpeza de dados
 - Transformação de dados

André Ponce de Leon F de Carvalho

39



Modelagem

- Escolha e aplicação de diferentes técnicas
 - bem como o ajuste de seus parâmetros
- Nota:
 - Existem várias técnicas para o mesmo tipo de problema de DM...
 - Algumas têm necessidades específicas para o formato dos dados
 - Frequentemente é necessário voltar à fase de preparação de dados

André Ponce de Leon F de Carvalho

40



Avaliação

- Protótipo do modelo estará construído
- Antes do seu desenvolvimento final:
 - É importante avaliá-lo e revisar os passos executados para sua construção
 - Para ter certeza de que atende adequadamente aos objetivos do problema
 - Verificar se algum aspecto importante não foi suficientemente considerado
- Decide se os resultados da DM serão usados

André Ponce de Leon F de Carvalho

41



Desenvolvimento

- Modelo em geral não atende às necessidades dos usuários finais, por diferentes razões, p. ex.:
 - Conhecimento precisa ser organizado e apresentado em uma forma que o usuário possa utilizá-lo
 - interfaces, ferramentas de geração de relatórios, etc.
 - usuário muitas vezes não é analista de DM !
 - Conhecimento em geral precisa ser explorado sob diferentes perspectivas por diferentes usuários e/ou em diferentes instantes de tempo
 - ferramentas de análise exploratória
 - ...

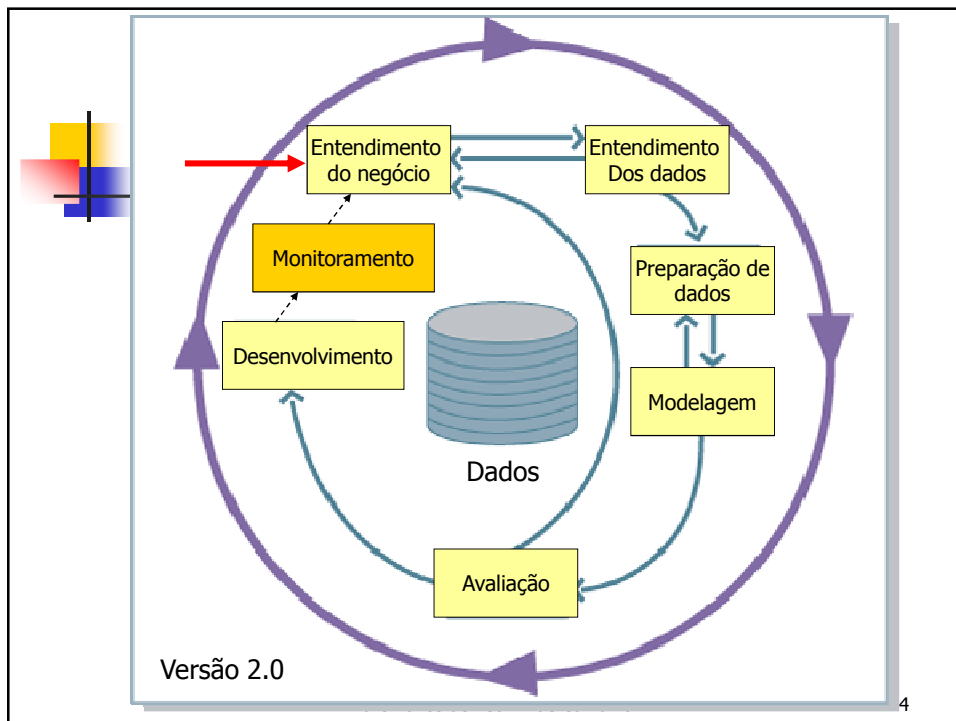
42

CRISP-DM

- Versão 2.0 em elaboração
 - *Special Interest Group* está sendo consultado
 - Possíveis mudanças
 - Divisão da fase de preparação de dados
 - Métodos de avaliação preliminares dentro da fase de modelagem
 - Inclusão de fase de monitoramento
 - ...

André Ponce de Leon F de Carvalho

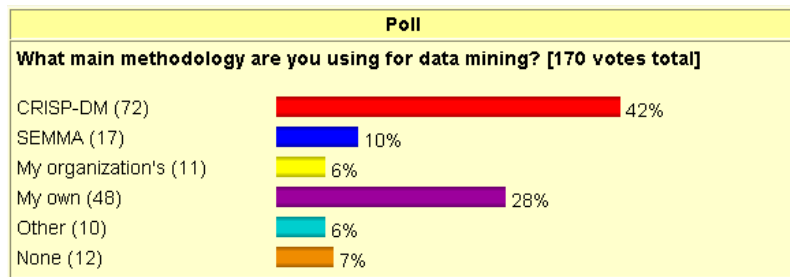
43



<http://www.kdnuggets.com>

Pesquisas KDnuggets

■ Pesquisa realizada em 2004

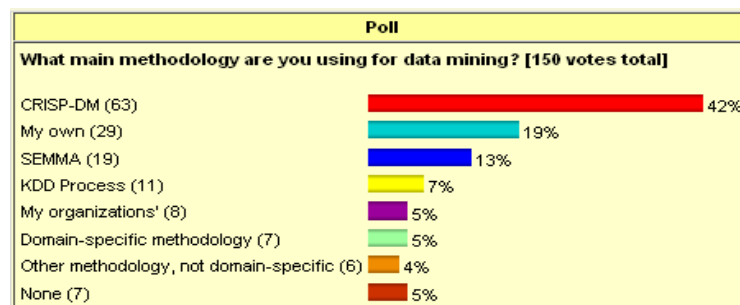


André Ponce de Leon F de Carvalho

45

Pesquisas KDnuggets

■ Pesquisa realizada em 2007



André Ponce de Leon F de Carvalho

46



Aplicações

- Número crescente de aplicações
 - Ciência e Medicina: descoberta de padrões de risco, diagnóstico de pacientes, análise de genoma, ...
 - Finanças: análise de risco (p. ex. de crédito), detecção de fraudes, gerenciamento de carteiras, ...
 - Internet: algoritmos de busca, marketing na web, ...
 - Indústrias: previsão de falhas, diagnóstico de qualidade, ...
 - Marketing: segmentação de mercado, ...

André Ponce de Leon F de Carvalho

47



Exemplos de Aplicações

- Siemens Medical
 - Ferramenta de DM para o Tratamento de Ataques cardíacos
 - Combina informações médicas de diversas fontes
 - Busca automática em registros combinados de 6 milhões de pacientes

André Ponce de Leon F de Carvalho

48



Exemplos de Aplicações

- Siemens Medical
 - Descobriu centenas de casos onde os melhores procedimentos médicos não haviam sido seguidos
 - Mas ainda havia tempo para intervir
 - Identificou pacientes elegíveis para estudos médicos
 - Ganhou o 2005 ICDM Data Mining Practice Prize

André Ponce de Leon F de Carvalho

49



Exemplos de Aplicações

- The Mitre Corporation
 - Ferramenta de DM para detecção de fraudes no IR
 - Sistema de DM usa Aprendizado de Máquina e Análise Estatística para descobrir sonegações
 - Resultados
 - Encontrou casos não descobertos por auditores
 - Reduziu tempo de análise
 - 2 semanas para poucas horas (dados de 2001)
 - 2º lugar no 2005 ICDM Data Mining Practice Prize

André Ponce de Leon F de Carvalho

50

Aplicações em Dados Biológicos

- Análise de expressão gênica
- Reconhecimento de genes
- Previsão de estruturas de proteínas
- Alinhamento de sequências
- Reconstrução de árvores filogenéticas
- ...

51

Análise de Expressão Gênica

- Medição do nível de expressão de grande quantidade de genes sob diversas condições
- Aplicações
 - Agrupar genes com expressões similares
 - Atribuir funções a genes não anotados
 - Identificar novas subclasses de doenças
 - Classificar tecidos
- Algoritmos de
 - Agrupamento
 - Classificação
 - Seleção de atributos

52

Reconhecimento de Genes

- Encontrar genes em sequências de DNA
- Análise em laboratório
 - Difícil e cara
- Análise computacional
 - Duas abordagens
 - Busca por sinal
 - Busca por conteúdo
 - Ambas utilizam algoritmos de Classificação

55

Previsão de Estruturas de Proteínas

- Predição da estrutura espacial de proteínas
 - Importante uma vez que sua estrutura determina ou provê indícios de sua função
- Métodos tradicionais
 - Cristalografia
- Métodos computacionais
 - Se baseiam em estruturas já conhecidas
 - Aplicação de métodos de classificação

56

Alinhamento de Sequências

- Busca por sequências similares
 - Possivelmente possuem ancestral comum
 - História evolutiva dos organismos
 - Diferenças entre sequências
 - Deleções
 - Inserções
 - Substituições
- Algoritmos
 - Programação dinâmica (BLAST)
 - Algoritmos genéticos

57

Reconstrução de Árvores Filogenéticas

- Árvores filogenéticas
 - Ancestralidade, relacionamento entre espécies ou grupos de espécies
 - Construídas com base em similaridades, por exemplo, entre sequências de bases
- Métodos computacionais
 - Método Neighbor-Joining, baseado no conceito de vizinho mais próximo.
 - Algoritmos de agrupamento hierárquicos

58

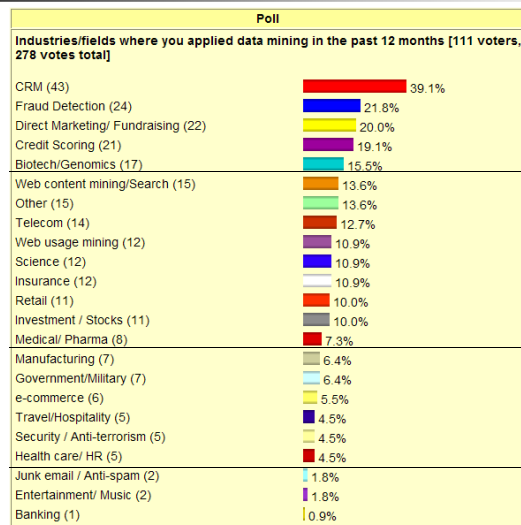
Pesquisas KDnuggets

- Aplicações de DM
- Em que indústrias / áreas você está atualmente aplicando DM
 - Fonte:
 - http://www.kdnuggets.com/polls/2006/data_mining_applications_industries.htm
 - Data: junho de 2006
 - 278 votos de 111 votantes

André Ponce de Leon F de Carvalho

59

Pesquisas KDnuggets



60

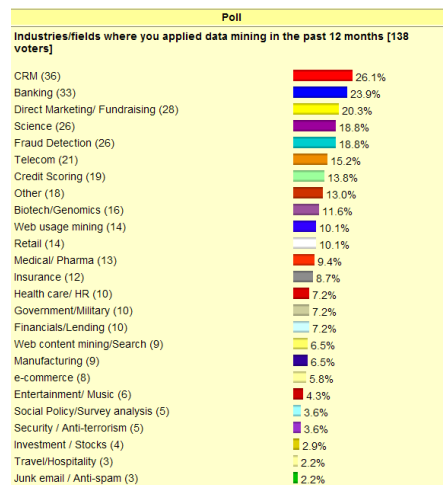
Pesquisas KDnuggets

- Aplicações de DM
- Em que indústrias / áreas você está atualmente aplicando DM
 - Fonte:
 - http://www.kdnuggets.com/polls/2007/data_mining_applications.htm
 - Data: junho de 2007
 - 138 votantes

André Ponce de Leon F de Carvalho

61

Pesquisas KDnuggets



André Ponce de Leon F de Carvalho

62



Pesquisas KDnuggets

- 2006 x 2007

Industry/Area	2006 Count	2007 Count	Change Y6 to Y7
Banking	1	33	3200%
Entertainment/ Music	2	6	200%
Science	12	26	117%
Health care/ HR	5	10	100%
Medical/ Pharma	8	13	63%
Junk email / Anti-spam	2	3	50%
Telecom	14	21	50%
Government/Military	7	10	43%
e-commerce	6	8	33%
Manufacturing	7	9	29%
Direct Marketing/ Fundraising	22	28	27%
Retail	11	14	27%
Other	15	18	20%
Web usage mining	12	14	17%
Fraud Detection	24	26	8%
Insurance	12	12	0%
Security / Anti-terrorism	5	5	0%
Biotech/Genomics	17	16	-6%
Credit Scoring	21	19	-10%
CRM	43	36	-16%
Travel/Hospitality	5	3	-40%
Web content mining/Search	15	9	-40%
Investment / Stocks	11	4	-64%

André Ponce de Leon F de Carvalho

63

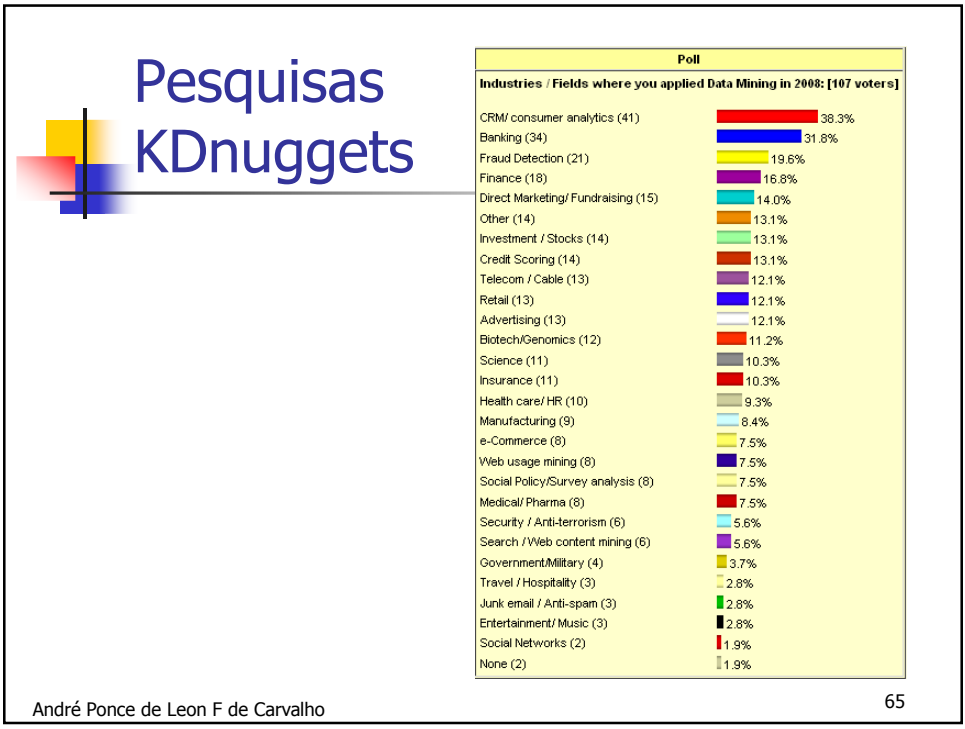


Pesquisas KDnuggets

- Aplicações de DM
- Em que indústrias / áreas você está atualmente aplicando DM
 - Fonte:
 - <http://www.kdnuggets.com/polls/2008/data-mining-applications.htm>
 - Data: dezembro de 2008
 - 107 votantes

André Ponce de Leon F de Carvalho

64



Produtos de DM























































André Ponce de Leon F de Carvalho 66

Mais Produtos


DataEngine DataMIND MarketMiner Inc. pcOLPARS
 DARWIN MINESET Urban Science GainSMARTS
 PRW ModelQuest Enterprise Affinium The Data Mining Suite
 Intelligent Miner for data comotion TextSense DBMiner Insight
 NETEZZA Partek Pro 5.0

André Ponce de Leon F de Carvalho 67

Alguns Mitos (Padhraic Smith)

- “Análise de dados pode ser completamente automatizada”
 - Julgamento humano é crítico na maioria das aplicações
 - De qualquer modo, semi-automatização é muito útil
- “Com uma quantidade massiva de dados, não é necessário estatística”
 - Grande volume leva a heterogeneidade
 - Precisa ainda mais de estatística


André Ponce de Leon F de Carvalho 68



Pacotes e Conjuntos de Dados

- WEKA
 - http://directory.google.com/Top/Computers/Artificial_Intelligence/Machine_Learning
- R
 - <http://lancet.mit.edu/ga>
- Machine Learning Data Repository UC Irvine
 - <http://www.ics.uci.edu/~mlearn/ML/Repository.html>


05/08/2010 André de Carvalho - ICMC/USP 69




Considerações Finais

- Expansão do volume de **dados**
- Necessidade de extrair **conhecimento**
- KDD é cada vez mais usado
- Cuidado com promessas exageradas
 - Sistemas Especialistas...

André Ponce de Leon F de Carvalho 70



Perguntas



André Ponce de Leon F de Carvalho

71