

Extracção de Conhecimento em Bases de Dados (ECBD ou KDD)

• O <u>ECBD</u> ou <u>KDD</u> é muitas vezes denominado de apenas Data Mining, ainda que, este seja, mais propriamente, uma das fases do processo (KDD conference, 1995, Montreal).

<u>Relação</u> do <u>ECBD</u> com <u>outras ferramentas</u> de exploração de informação:

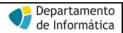
- Com ferramentas até agora descritas (capítulo anterior), poderse-á responder a questões como:
 - as vendas do produto X cresceram em Novembro?
 - as vendas do produto X diminuem quando há uma promoção do produto Y?
- Com ferramentas no domínio do ECBD/Data Mining, poderemos colocar a questão:
 - · Quais são os factores que determinam as vendas do produto X?



Análise Inteligente de Dados

2





DO/Reporting e OLAP x Data Mining

Relembrando o que já atrás foi focado (capítulo 1):

- Com as <u>ferramentas tradicionais</u>, o analista <u>coloca uma questão</u>, ou suposição ou talvez só uma inclinação e explora os dados. <u>Cria um modelo</u>, passo-apasso, trabalhando para <u>provar ou negar uma teoria</u>.
- É da <u>responsabilidade do analista propor cada hipótese</u>, <u>testá-la</u>, propor uma hipótese substituta ou adicional, testá-la <u>e assim sucessivamente</u>, e desta forma interactiva, criar o modelo.

Esta responsabilidade não desaparece inteiramente com data mining, mas,

- <u>muito do trabalho</u>, encontar o modelo apropriado, <u>é deslocado do analista para o computador</u>.
- O sistema toma a iniciativa da análise de dados, não o utilizador.

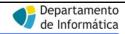
Benefícios:

- gerar o modelo requer menor esforço manual (mais eficiente);
- podem avaliar-se muito mais modelos, aumentando assim a possibilidade de encontar melhor modelo;
- o analista necessita de muito menor habilidade, dado que muitos dos procedimentos passo-a-passo são automáticos.



Análise Inteligente de Dados





Definição de Data Mining (1)

- Já no 1º capítulo, estabelecemos algumas diferenças entre Data Mining e outras ferramentas utilizadas no domínio da extracção de informação de uma base de dados (data query, reporting e OLAP).
- O Data Mining, ou mais genericamente o ECBD, pode ser visto segundo diversas perspectivas:
- 1. Numa perspectiva de negócio, será:
 - O processo de identificação de padrões e relacionamentos escondidos numa base de dados

Data Mining: Building Competitive Advantage

• Extracção de informação de negócio útil a partir de grandes bases de dados.

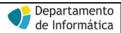
Data Warehousing, Data Mining and OLAP



Análise Inteligente de Dados

5





Definição de Data Mining (2)

- 2. Numa perspectiva funcional:
- É a procura de informação valiosa em grandes volumes de dados, resultado da cooperação de esforços humanos e de computadores. Os humanos desenham as bases de dados, descrevem problemas e estabelecem objectivos. Os computadores peneiram os dados, procurando padrões que correspondam aos objectivos.

Predictive Data Mining: a practical guide, Weiss S.M, and Indurkhya N.

- 3. Numa perspectiva mais académica:
- A <u>extracção implícita</u>, <u>não trivial</u> de <u>conhecimentos úteis</u>, <u>previamente desconhecidos</u>, dos <u>dados</u>.

Data Mining, Pieter Adrians, Dolf Zantige

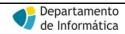
• O <u>processo</u> <u>não trivial</u> de identificação de <u>padrões válidos</u>, <u>novos, potencialmente úteis</u> e <u>compreensíveis</u> nos <u>dados</u>.

Frawley, Piatetsky e Matheus, 1991



Análise Inteligente de Dados





Definição de Data Mining (3)

Analisemos esta última definição:

Padrão - descrição mais simples do que a enumeração de todos os factos.

Processo - O processo de ECBD compreende, em geral, várias fases, envolvendo: (1) Definição do problema, (2) preparação dos dados, (3) procura de padrões, (4) avaliação dos resultados e (5) refinamento iterativo dos resultados.

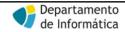
Não trivial - O processo deve envolver um certo grau de procura de padrões úteis. (Ex. calcular uma remuneração média dos clientes de uma base de dados sobre empréstimos, embora possa ser útil, não poderá ser entendido como extracção).



Análise Inteligente de Dados

7





Definição de Data Mining (4)

Validade - Os padrões extraídos devem, com um determinado grau de certeza, ser válidos para novos dados.

Novidade - A novidade pode ser medida com **referência aos dados** (comparação dos valores correntes com valores prévios ou esperados) ou **ao conhecimento** (comparação de uma nova descoberta com as anteriores).

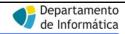
Utilidade potencial - Os padrões de detectados **devem conduzir potencialmente a acções úteis**.

Ex. num exemplo de empréstimos bancários, seria uma medida do aumento de lucros esperados para o banco em resultado da aplicação da regra de decisão decorrente do padrão obtido.



Análise Inteligente de Dados





Definição de Data Mining (5)

Compreensibilidade / Sensibilidade - Um dos objectivos da extracção de conhecimento é tornar os padrões gerados compreensíveis com vista a possibilitar uma melhor compreensão dos dados.

Como veremos, há técnicas de DM que são inerentemente mais potentes quanto a esta característica (ex. árvores de decisão - transparentes) do que outras (ex. redes neuronais - opacas).

Medida de Interesse - medida do valor de um padrão, combinando:

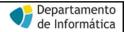
- validade
- novidade
- utilidade
- simplicidade



Análise Inteligente de Dados

9





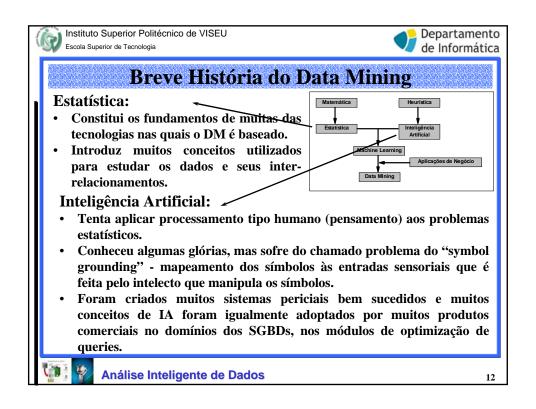
Poder do Data Mining

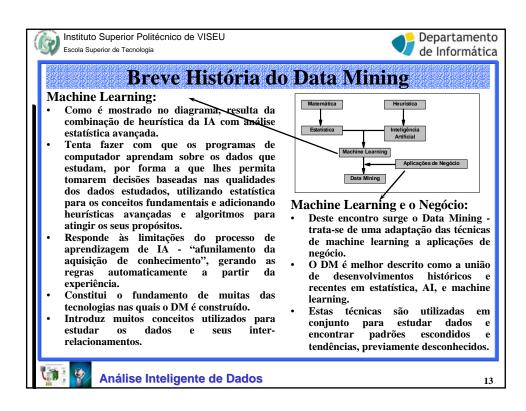
- O poder do data mining é devido ao facto de ele não depender das vistas humanas estreitas, para produzir os seus resultados, mas, em seu lugar, procura e identifica relacionamentos de que os humanos nunca teriam percepção.
- Uma boa forma de identificar esta realidade é avaliar o modo como um mestre de xadrez distingue um opositor humano de um cibernético.
 - Um computador faz muitas vezes jogadas que um humano nunca executaria, pois que este último "não olhou bem".
 - O que se passa é que a capacidade humana para explorar um grande número de movimentações, num tempo exíguo, é limitada. Tem assim que minimizar a "árvore de pesquisa", limitando o número de caminhos possíveis, baseados na pré-concepção do que entendemos como "estar ou não estar certo".



Análise Inteligente de Dados

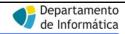












Distinção entre IA e Data Mining

Pode parecer que o Data Mining seja uma parte da IA, mas:

- Os sistemas IA lidam com a codificação do pensamento humano num programa de computador tentando simular a inteligência.
- Os sistema IA são conduzidos pelo conhecimento humano.

Data Mining, Estatística e Machine Learning, são sistemas:

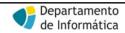
- Conduzidos por dados aprendem com exemplos da vida real e não com ideias pré-processadas. Utiliza informação histórica (experiência) para aprender.
- Desta forma:
 - estes sistemas são criados de forma automática e fácil,
 - podem ser actualizados rapidamente,
 - · são muito menos onerosos na sua construção,
 - não obrigam à existência de um perito intimamente ligado à criação do sistema.



Análise Inteligente de Dados

14





Distinção entre DM e Estatística

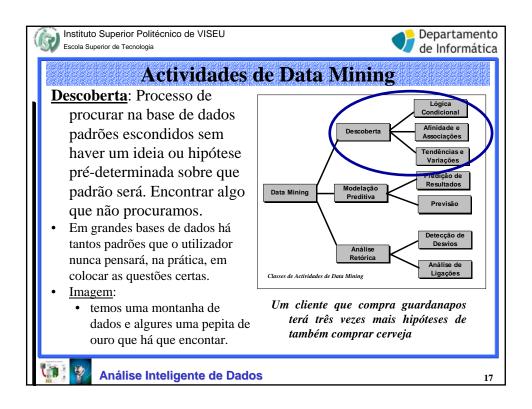
Ex. A regressão é utilizada para criar modelos capazes de predizer o comportamento de clientes, baseados em grandes volumes de dados.

Mas:

- o DM é capaz de ser utilizado pelo utilizador final
- já a estatística, terá de o ser por um perito na matéria
- "se a maioria da estatística se traduz num processo de estabelecer uma hipótese e e depois verificá-la, porque não deixar o computador fazer essas tentativas e testá-las automaticamente?"



Análise Inteligente de Dados

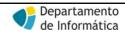












Fases do Processo de ECBD

Em princípio, o processo de ECBD consiste em seis estágios:

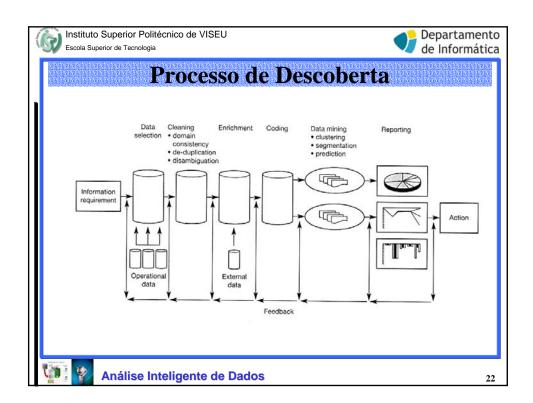
- Selecção dos Dados
- Depuração
- Enriquecimento
- Codificação
- Extracção de Conhecimento (data mining)
- Relato

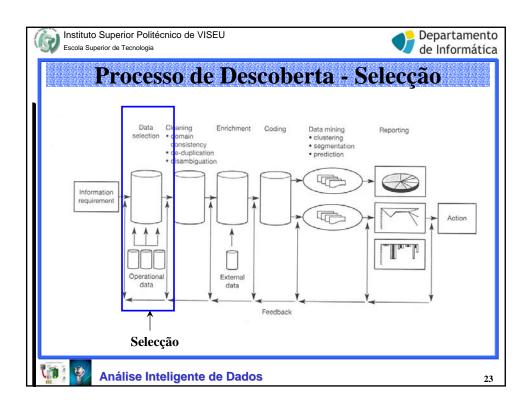
Observações:

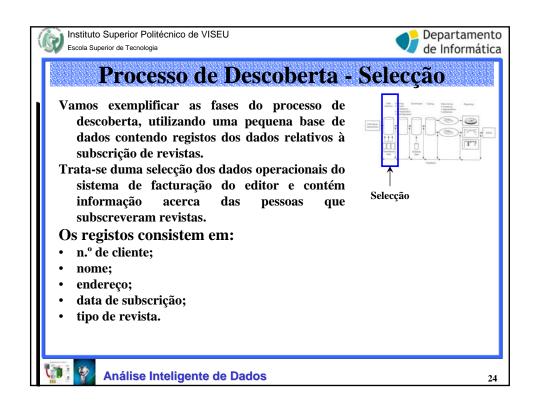
- Embora pareça que há uma trajectória linear, não é o caso. Em qualquer fase, pode ser necessário recuar uma ou mais fases; por exemplo, na fase da codificação ou de extracção de conhecimento, pode suceder apercebermo-nos que a fase de purificação está incompleta, ou descobrir novos dados e utilizá-los para novo enriquecimento.
- O processo é contínuo, devendo as organizações trabalhar continuamente os seus dados, identificando constantemente nova necessidade de informação, melhorando os dados para melhor atingir os objectivos.

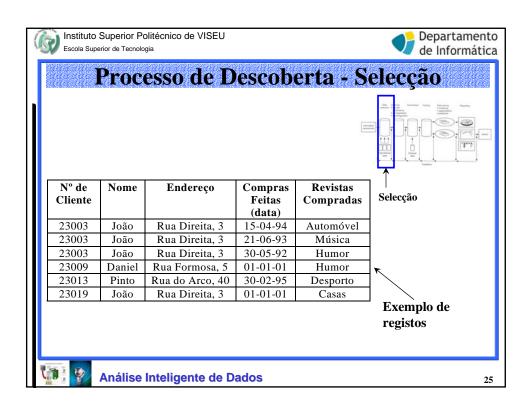


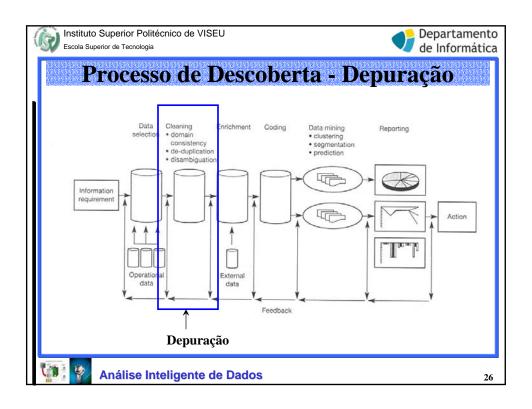
Análise Inteligente de Dados

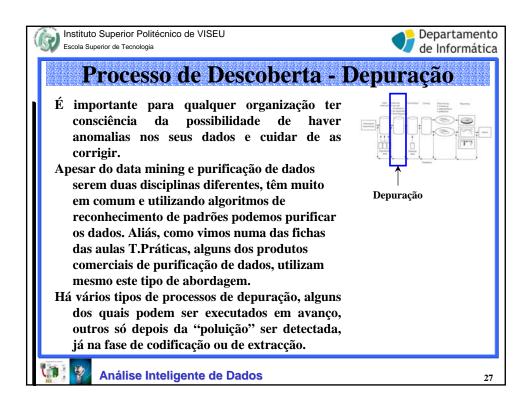


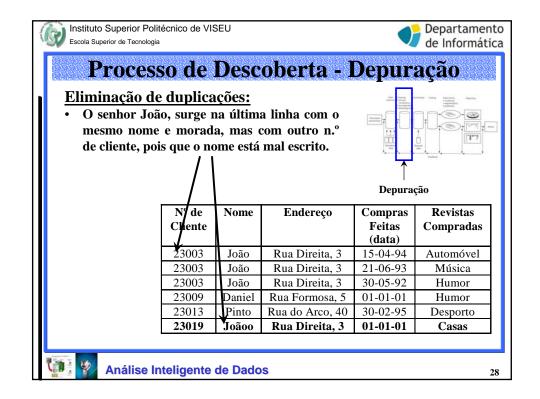


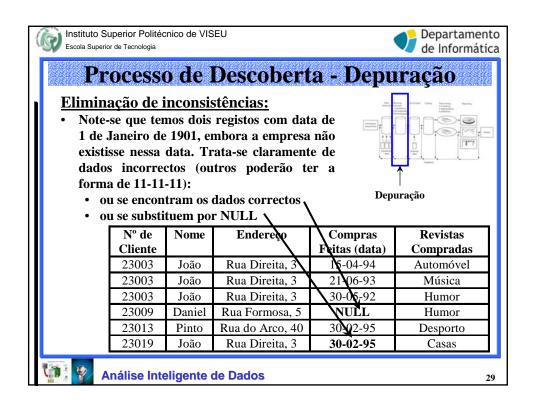


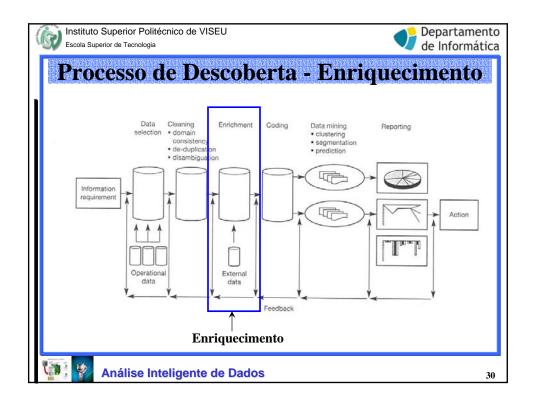


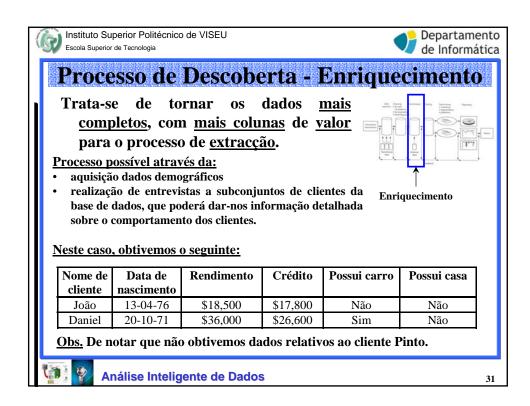


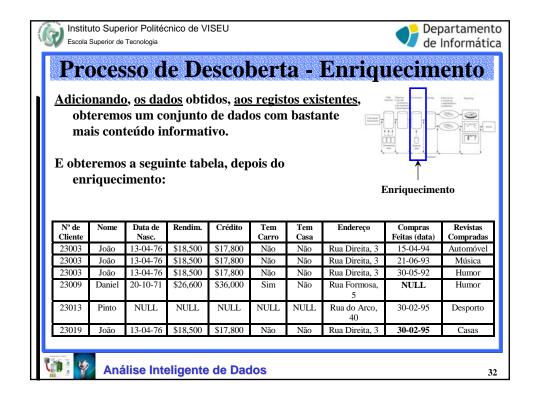


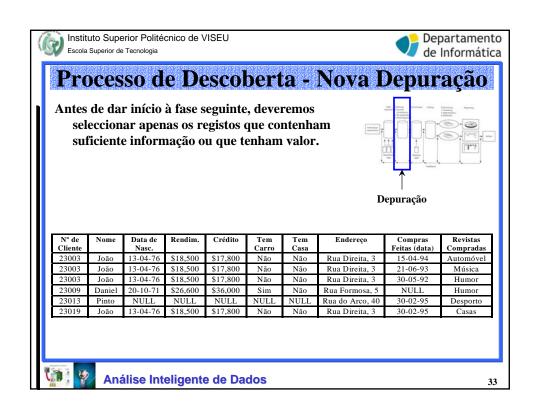


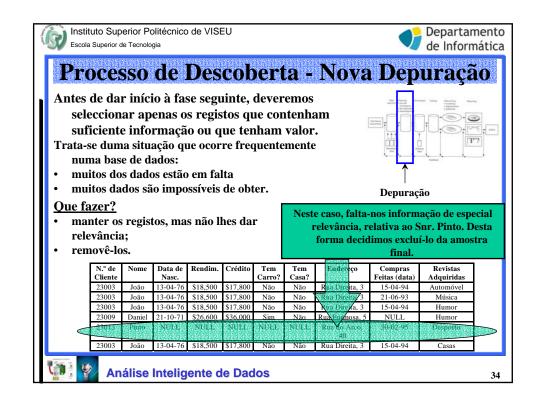


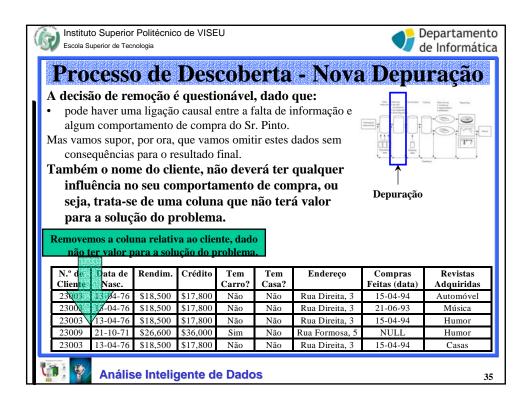


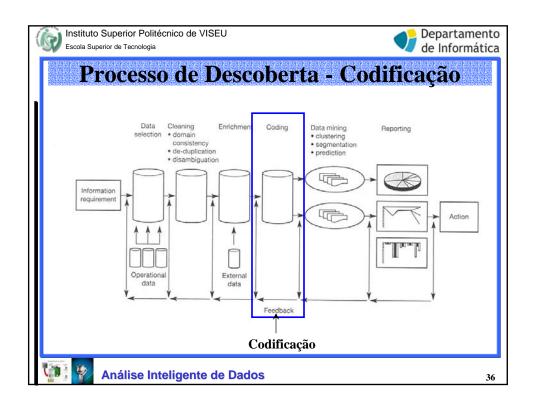


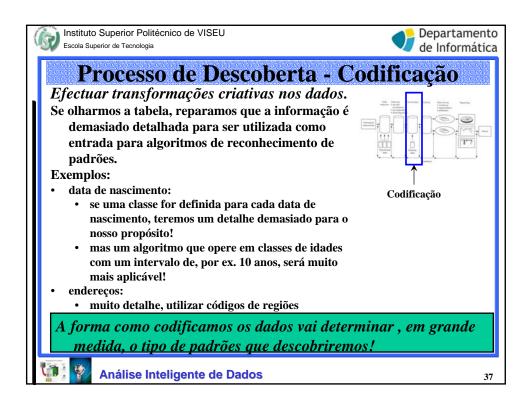


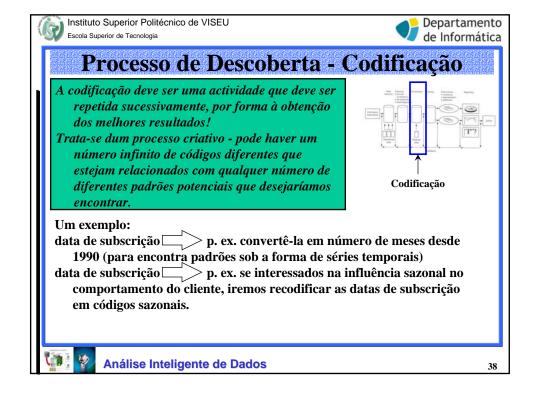


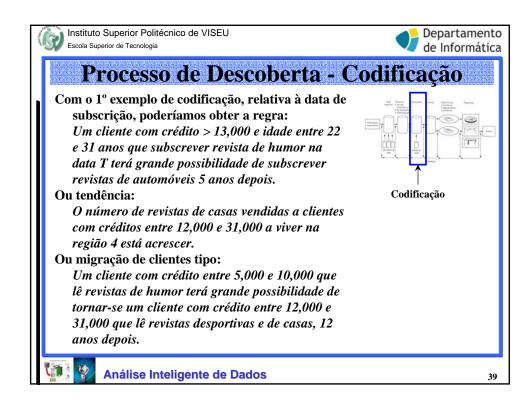


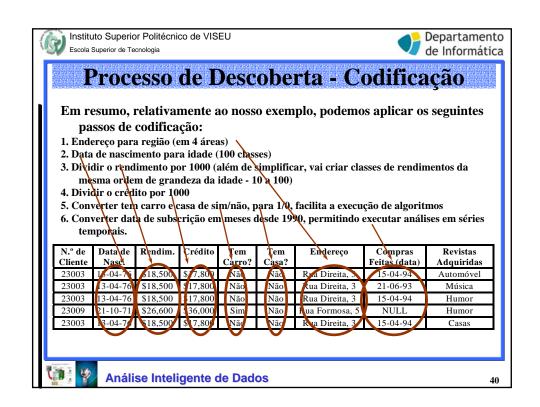




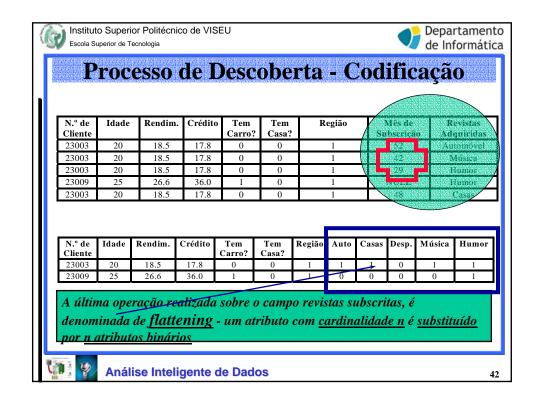


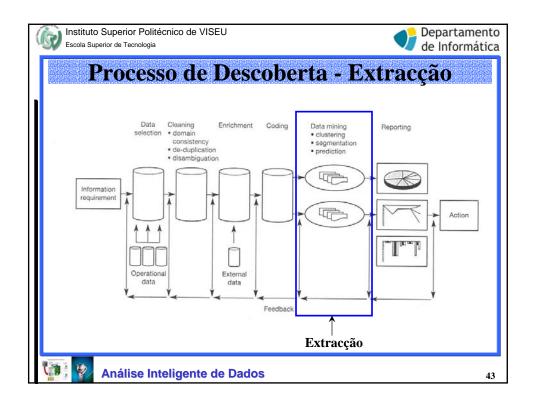


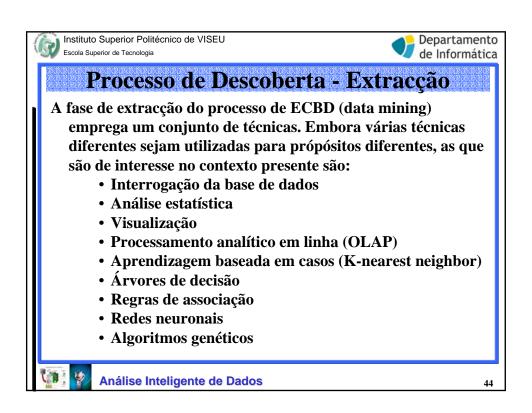




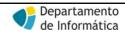












Extracção - Análise preliminar dos dados

Recurso a uma linguagem de interrogação (query tool)

- Trata-se de efectuar uma análise grosseira do conjunto de dados
- Através da utilização de SQL poderemos obter informação valiosa relativamente ao data set
 - já foi dito atrás que cerca de 80% da informação valiosa pode retirar-se através de query
 - só os restantes 20% obrigam à utilização de técnicas mais avançadas.

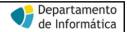
Para os exemplos seguintes, vamos utilizar uma base de dados relativa à subscrição de revistas (1000 clientes), de onde foi retirado o set até aqui utilizado.



Análise Inteligente de Dados

45





Extracção - Análise preliminar dos dados

Os valores médios, apresentam-se na tabela abaixo:

	Média
Idade	46.9
Rendimento	20.8
Crédito	34.9
Carro próprio	0.59
Casa própria	0.59
revista de carros	0.329
revista de casas	0.702
revistas de desporto	0.447
revista de música	0.146
revista de humor	0.081

Os números relativos à média são muito importantes, pois que nos dão uma norma que permite ajuizar do desempenho dos algoritmos de reconhecimento de padrões.

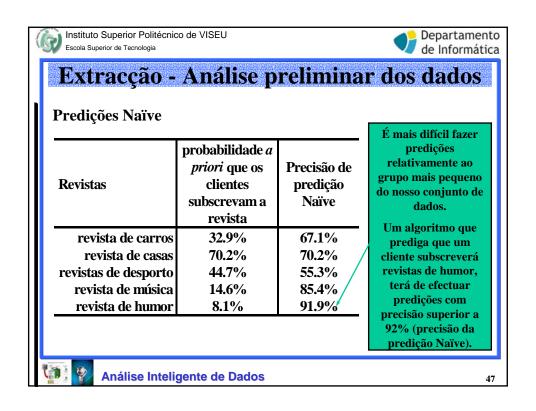
Um algoritmo que prediga sempre a <u>não</u> <u>subscrição</u> de revistas de automóveis, estará correcto em 671 por 1000 casos, cerca de 70%.

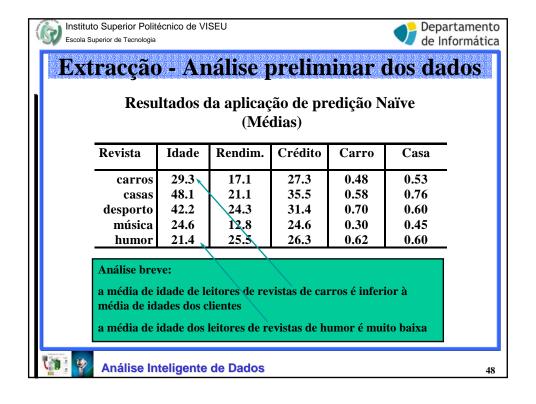
Qualquer algoritmo que reclame obter alguma visão sobre estes dados, permitindo alguma predição real, deverá melhorar este valor.

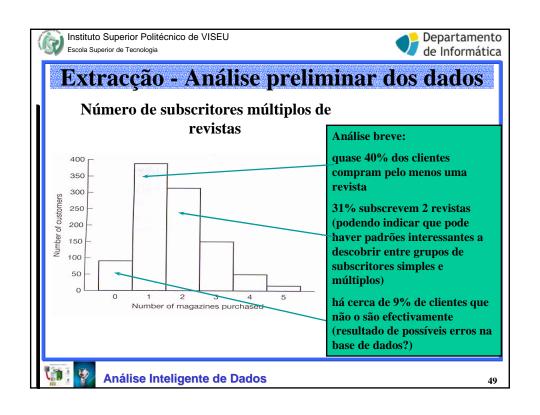
Predição Naïve - Resultado trivial que é obtido através dum método extremamente simples. O algoritmo de aprendizagem deve fazer melhor do que isto.

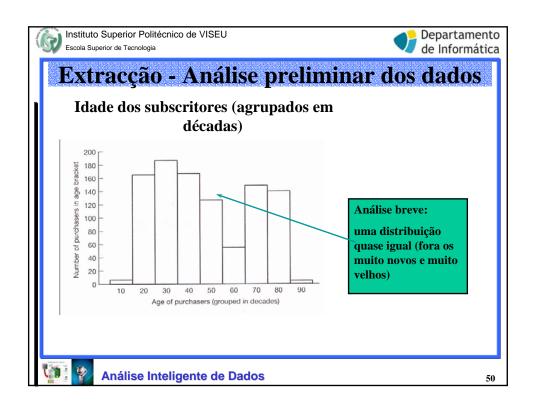


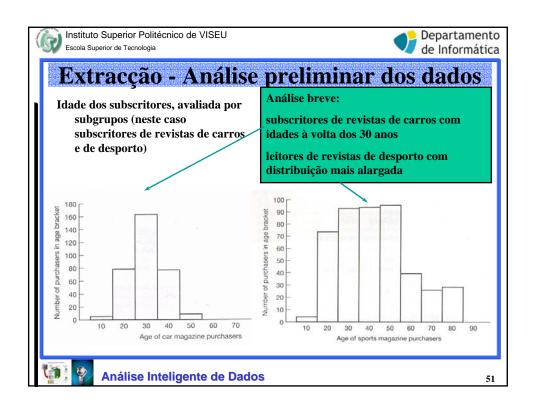
Análise Inteligente de Dados

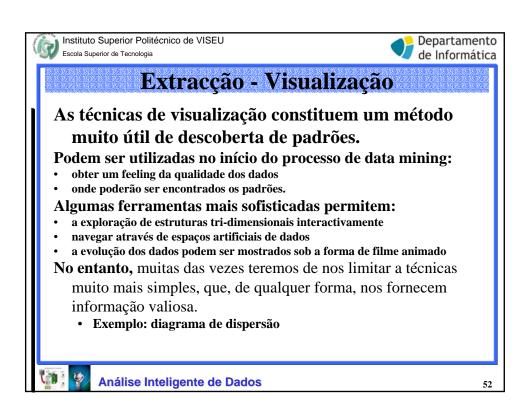


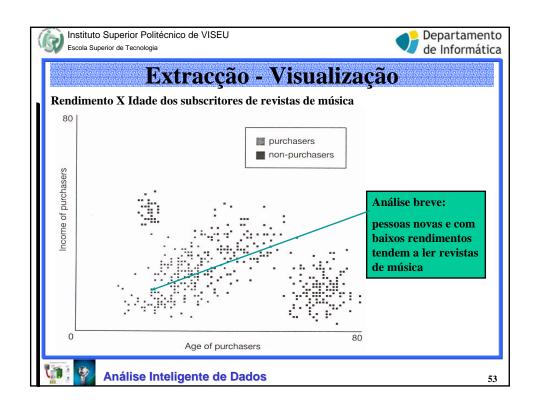


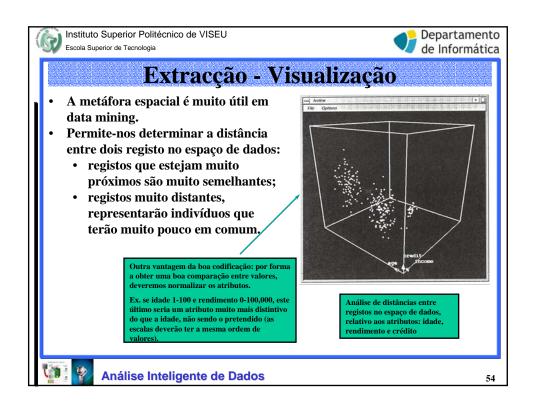


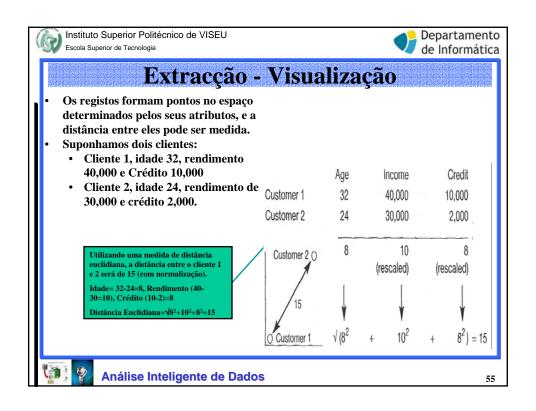


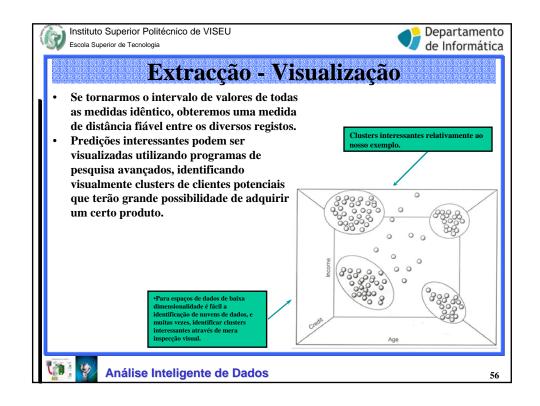




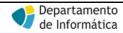












Extracção - Ferramentas OLAP

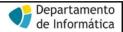
- Como já vimos no capítulo 2, estas ferramentas permitem o acesso a diversas formas de análise e visualização multidimensional e interactiva da informação.
- Aproximam a análise à forma de ver o negócio um determinado valor é visto segundo um conjunto de perspectivas lhe que fornecem caracterização.
- Funcionalidades dum sistema OLAP:
 - cálculo e modelação multidimensional
 - análise de tendências em sequências temporais
 - · análise de agrupamento de dados
 - movimentação e comparações ao longo das dimensões em consideração
- Observação Importante: as ferramentas OLAP não aprendem, não criam qualquer novo conhecimento, não pesquisam novas soluções. Obrigam, em regra, a motores multidimensionais intermédios ou a novas formas de armazenamento dos dados.



Análise Inteligente de Dados

57





Aplicar Data Mining (1)

Classificação e Regressão

- Representam a maioria dos problemas a que são aplicadas as técnicas de data mining, criando modelos capazes de predizer o valor ou classe de uma variável dependente, utilizando técnicas de indução supervisionada.
- São criados modelos :
 - de classificação se capazes de predizer a classe a que um determinado membro pertence;
 - de regressão se capazes de predizer um valor.
- Em classificação, a resposta é simplesmente verdade ou falso; já em regressão, a resposta é um número, como, por exemplo, os lucros ou perdas relativos a um empréstimo.

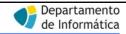
"O objectivo desta operação é utilizar o conteúdo da base de dados (dados sobre o passado) para gerar automaticamente um modelo que possa predizer o comportamento futuro."

"IBM's Data Mining Technology"



Análise Inteligente de Dados





Aplicar Data Mining (2)

Classificação e Regressão (continuação)

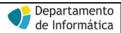
- Exemplos
 - predizer se um determinado empréstimo será ou não um bom risco de crédito;
 - predizer a rentabilidade de um cliente
 - predizer a probabilidade de que um determinado paciente tenha uma dada doença
- As séries temporais são apenas um tipo especial de problema relativo a regressão ou classificação, onde as medidas são obtidas a intervalos de tempo, como será o caso dos pagamentos relativos a um empréstimo.



Análise Inteligente de Dados

50





Aplicar Data Mining (3)

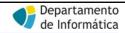
Associação e Sequência (também chamados de análise de cesto de compras)

- O <u>objectivo</u> é o <u>estabelecimento de relações</u> entre os <u>registos de uma base de dados</u>, não a criação de um modelo que caracterize o conteúdo de uma base de dados, como no processo anterior.
- Geram modelos descritivos que descobrem regras como:
 - os clientes que compram espaguete têm três vezes mais possibilidade de adquirirem queijo do que aqueles que não compram.
 - Exemplo, um gestor de vendas de um supermercado está muito interessado em conhecer que produtos se vendem conjuntamente, por forma a adquiri-los para a loja e a dispô-los fisicamente perto, implicando porventura um reordenamento das respectivas secções, mas decerto visando um incremento das vendas, pois que o cliente será automaticamente lembrado da conveniência da compra adicional.
- Trata-se duma operação que é suportada por técnicas de descoberta de associações e sequências.
 - A resposta deste tipo de operação aos problemas que lhe são colocados tem um formato lógico como "Um empréstimo é pago com 90% de confiança, quando o proprietário tem um história pessoal de pagamento atempado dos débitos efectuados com cartão de crédito".



Análise Inteligente de Dados





Aplicar Data Mining (4)

Clustering (segmentação da base de dados)

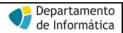
- Técnica descritiva que agrupa entidades similares em conjunto, colocando entidades dissimilares em grupos diferentes.
- Pode ser utilizada em marketing para encontrar grupos de clientes com afinidades e, em cuidados de saúde, para encontar pacientes com perfis semelhantes.
- Esta operação surge como resultado de necessidade de obter-se um resumo de cada base de dados ou antes de dar início a uma das operações de data mining.
- O Clustering é muito subjectivo:
 - Dado que se emprega uma medida de distância, como a técnica de vizinho mais próximo (nearest neighbor), os clusters estão completamente dependentes da medida da distância que é utilizada.
 - Normalmente é de interesse o envolvimento de um perito no domínio de análise pretendido, para propor a medida de distância apropriada.



Análise Inteligente de Dados

61





Aplicar Data Mining (5)

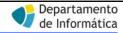
Clustering (segmentação da base de dados)

Um exemplo de utilização será o de uma cadeia de armazéns que mantém o registo das compras efectuadas pelos seus clientes, conhecendo as compras efectuadas, em qualquer visita de um seu cliente a cada uma das lojas. Neste caso, será muito útil segmentar a base de dados, baseando a divisão em períodos significativos de análise, como: período de "regresso às aulas", "antes do Natal", etc. Sobre cada um destes segmentos, poderá depois ser aplicada análise de ligações, para identificar que produtos são vendidos em conjunto. Para executar esta operação, são utilizadas técnicas de agrupamento "clustering".



Análise Inteligente de Dados





Aplicar Data Mining (6)

Detecção de Desvios

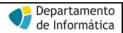
- Trata-se de uma operação no quadrante oposto da anterior, mas relacionada com ela.
 - Atrás, era importante identificar registos relacionados e estabelecer, com isso, grupos e as correspondentes divisões;
 - Agora, normalmente depois da segmentação efectuada, há que "identificar pontos que caem fora de um conjunto de dados particulares, explicando se serão devidos a ruído ou outras impurezas, ou por razões casuais";
 - É especialmente devido a estas últimas que esta operação é importante: em muitos casos, a identificação e explicação de um desvio é a fonte de descobertas puras, pois que expressam desvios em relação a expectativas e normas conhecidas previamente, podendo indicar o surgir de novas tendências ou oportunidades. Esta operação é suportada por técnicas estatísticas, tais como o teste de significância, onde sumarizações de estatísticas (média e desvio standard) são utilizadas para medir as diferenças.



Análise Inteligente de Dados

63





Aplicar Data Mining (7)

- Text Mining.
 - Também conhecido por Text Data Mining ou Descoberta de Conhecimento em Bases de Dados Textuais
 - Refere-se ao processo de extracção de padrões de interesse e não triviais ou conhecimento a partir de documentos de texto não estruturados.
 - Pode ser visto como uma extensão ao Data Mining ou KDD.
 - Dada a abundância de documentos (é a mais normal forma de arquivar informação – um estudo recente indica que cerca de 80% da informação das empresas reside sob a forma de documentos) possui um enorme potencial (maior mesmo do que o DMining).



Análise Inteligente de Dados

