



Universidade Federal de Santa Maria – Centro de Tecnologia

Peso: 5

Disciplina: Mineração de Dados (ELC1098)

Professor: Dr. Joaquim Assunção

Atenção: Leia atentamente as questões e escreva de forma clara e legível.

Nome: Lana Romato

Curso: Ciência da Computação Data: 06/09/18

- 1) Mineração de dados está inserida no contexto de um processo maior, KDD, cujas etapas são prévias e posteriores a mineração em si. Embora no dia-a-dia sejam tratados como sinônimos, o processo de KDD envolve tarefas como seleção e pré-processamento de dados, que ainda não são, tecnicamente, considerados como mineração. Descreva o porquê da existência dessas etapas prévias. Quais são os problemas e como essas etapas são úteis para a mineração em si.

(3)

A etapa de Seleção serve para selecionar os dados de um banco de dados, como o MySQL por exemplo. A etapa de Pré-processamento serve para transformar dados brutos em dados mais simples de se processar futuramente. Nessa última etapa, existe a fase de limpeza, onde os dados faltantes não completados para que exista uma conformidade entre eles, mas muitas vezes pode trazer dúvida quanto ao método a ser utilizado: média, valor mínimo, máximo, mediana dos dados.

- 2) Quais são as três grandes categorias de dados (quanto a estruturação)? Destas, qual é a mais próxima de um formato comumente minerado e por quê? (2)

Estruturados - formato mais tabulado. Semi-estruturado - Não estruturado - fotos e vídeos. XML e HTML. O formato mais próximo é o estruturado, pois tabular facilita a distribuição e mineração dos dados, mas os não estruturados também podem ser minerados.

- 3) Considerando conjuntos de dados estruturados "tidy". Por que não é comum armazenar dados nesse formato e porquê pode ser útil transformar os dados para esse formato? (2)

Não é comum pois a tabela fica muito grande se tiver muitos dados, mas é útil porque não apresentados com mais clareza e na hora de juntar tabelas com os mesmos tipos de dados, eles não ficam separados ou confundidos em linhas/columnas diferentes.



4) Analise o seguinte conjunto e responda:

4.1) Caso nosso projeto de mineração tenha como objetivo encontrar correlação entre tamanho de tela, valores e lucro. Quais informações podemos descartar? (1,5)

A coluna Produto e o resto da descrição que não fala sobre a tela 0.5

4.2) Considerando esse objetivo (4.1), devo alterar (adicionar ou editar) algo nesse conjunto? Por que? (1,5)

Deixar a coluna ID, retirar a coluna Produto, retirar as informações que não falam sobre tela da coluna Descrição (se não for possível, não deixar), juntar as colunas Valor compra e Valor venda substituindo-as e deixar a coluna de Unidades vendidas. Assim a tabela ficaria muito mais clara e objetiva, deixando apenas as informações necessárias 0.3

ID	Produto	Descricao	Valor compra	Valor venda	Unidades vendidas
1	Monitor LG 29p	Entradas 1,2,3. 29p. Full HD IPS 29UM68	250	290	29
...	...	...	...	...	...
999	Monitor Samsung 30p	Entr. 1,2,3,4. 30p. 4K	500	600	45

(Descrição)

ID	Tela	Det. Venda	Unidades Vendidas
----	------	------------	-------------------