

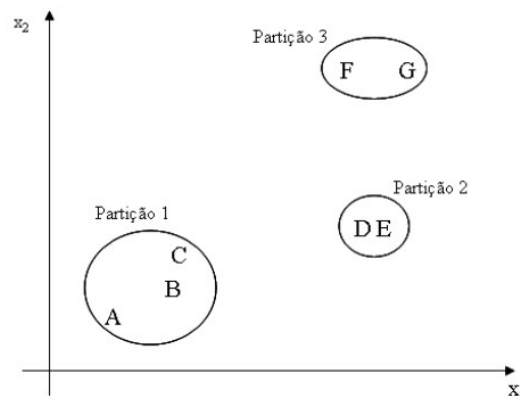
## Paradigmas de Programação Segundo Trabalho

Prof. Flávio Miguel Varejão

### I. Descrição do Problema

Agrupamento de dados multidimensionais é um dos problemas mais comuns na área de aprendizado de máquina. Esse problema consiste em dividir um conjunto de pontos em um espaço multidimensional em um determinado número pré-especificado de grupos de modo que os pontos pertencentes a um mesmo grupo estão mais relacionados entre si e menos relacionados em relação aos pontos associados aos outros grupos.

A figura abaixo ilustra um exemplo de agrupamento no qual os sete pontos {A, B, C, D, E, F, G} foram agrupados em três grupos, indicando que os padrões {A, B, C} são mais similares entre si do que em relação aos demais, assim como os padrões {D, E} e {F, G}.



Formalmente, dado um conjunto de dados  $X$  com  $N$  pontos  $\{x_1, \dots, x_N\}$ , sendo que cada ponto  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}]$  possui  $d$  coordenadas (dimensões), deseja-se encontrar  $K$  grupos  $\{C_1, \dots, C_K\}$  de tal forma que as seguintes condições sejam atendidas:

- $C_j \neq \emptyset, j = 1, \dots, K$
- $\bigcup_{j=1}^K C_j = X$
- $C_i \cap C_j = \emptyset, i \neq j, i, j = 1, \dots, K$

Uma forma de realizar agrupamento de dados multidimensionais envolve inicialmente a criação uma lista de ligações entre pares de pontos e em seguida cortar as ligações mais distantes de forma a criar os  $K$  grupos desejados. O corte da lista é sempre na maior ligação presente na lista.

O pseudo-código a seguir ilustra os passos para realização de agrupamento de dados usando uma álista de ligações mais próximas entre pares de pontos:

1. Escolher o primeiro ponto lido do arquivo de entrada como ponto corrente.
2. Escolher o ponto ainda não escolhido mais próximo ao ponto corrente para incluir a ligação do ponto corrente a esse ponto na lista.

3. Tornar o ponto recentemente escolhido como corrente.
4. Repetir os passos 2 e 3 até que todos os pontos tenham sido adicionados à lista de ligações.
5. Escolher a maior ligação da lista para dividi-la em duas listas.
6. Repetir o passo 5 para cada lista até que se tenham apenas K listas. A lista a ser dividida é sempre aquela que contém a maior ligação no estado corrente.
7. Ao final, os pontos pertencentes a cada lista compõe os grupos correspondentes.

Neste trabalho será usada a distância Euclidiana  $||x_i - x_j||$  como métrica de distância. Ela é calculada pela expressão:

$$||x_i - x_j|| = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{id} - x_{jd})^2}$$

Cada ponto do conjunto de dados a ser agrupado terá suas coordenadas (valor numérico em ponto flutuante) expressas em uma linha do arquivo csv de entrada. A linha em que os dados do ponto se encontra será o seu identificador único. Assim, o ponto na primeira linha será identificado pelo número 1, o ponto na segunda linha será identificado pelo número 2 e assim por diante.

Para uniformizar os resultados, o ponto inicial escolhido será o ponto 1, isto é, o primeiro ponto lido do arquivo.

Além disso, se houver igualdade na distância entre pontos, o critério de desempate levará em consideração o identificador único do ponto, isto é, o número de sua linha no arquivo. Por exemplo, se a distância entre o ponto 3 e o ponto 5 é igual a distância entre o ponto 3 e o ponto 17, será considerada como menor distância aquela entre 3 e 5.

$$(3, 5, 10) < (3, 17, 10)$$

De forma análoga, se a distância entre o ponto 7 e 11 é igual a distância entre o ponto 10 e 2, será considerada como menor distância aquela entre 7 e 11 porque 7 é menor do que 10.

$$(7, 11, 65) < (10, 2, 65)$$

## II. Especificação do Sistema

### Funcionalidades a serem implementadas:

1. Leitura do nome do arquivo de entrada, do nome do arquivo de saída e do número de grupos da entrada padrão.
2. Leitura da base de dados do arquivo csv de entrada.
3. Realização do agrupamento de dados.
4. Gravação dos identificadores dos pontos dos grupos no arquivo csv de saída. Cada linha do arquivo de saída corresponderá a um grupo.
5. Apresentação dos identificadores dos pontos dos grupos na saída padrão. Cada linha da saída corresponderá a um grupo.

Os exemplos seguintes são apenas ilustrativos dos formatos de entrada e saída e não existe correspondência entre os dados e o resultado. Em outras palavras, o arquivo de saída apresentado não contém os grupos corretos que deveriam ter sido gerados.

Exemplo de formato de arquivo de entrada:

```
7, 5.4, 6.32, 9
17, 32.3, 5, 9.99
33, 54, 5.6, 65.8
77.7, 33.4, 98, 7.56
8.9, 5.8, 6, 9
```

Exemplo de formato de arquivo de saída (com K = 2):

```
1, 3, 5
2, 4
```

Exemplo de formato de interação do programa com o usuário:

```
Forneça o nome do arquivo de entrada: base.csv
Forneça o nome do arquivo de saída: saida.csv
Forneça o número de grupos (K): 2
Agrupamentos:
1, 3, 5
2, 4
```

### **III. Condições de Entrega**

O trabalho deve ser feito individualmente e submetido pelo sistema da sala virtual até a data limite (30 de julho de 2025).

O trabalho deve ser submetido em um arquivo zip contendo todos os arquivos com código fonte em haskell. O arquivo zip deve possuir o nome Trab2\_Nome\_Sobrenome. Note que a data limite já leva em conta um dia adicional de tolerância para o caso de problemas de submissão via rede. Isso significa que o aluno deve submeter seu trabalho até no máximo um dia antes da data limite. Se o aluno resolver submeter o trabalho na data limite, estará fazendo isso assumindo o risco do trabalho ser cadastrado no sistema após o prazo. Em caso de recebimento do trabalho após a data limite, o trabalho não será avaliado e a nota será ZERO. Aluno que receber zero por este motivo e vier pedir para o professor considerar o trabalho não será nem respondido. Trabalhos em que se configure cópia receberão nota zero independente de quem fez ou quem copiou.

### **IV. Requisitos da implementação**

- a. Modularize seu código adequadamente.
- b. Crie códigos claros e organizados. Utilize um estilo de programação consistente, Comente seu código.
- c. Os arquivos do programa devem ser lidos e gerados na mesma pasta onde se encontram os arquivos fonte do seu programa.

## **V. Observação importante**

**Caso haja algum erro neste documento, serão publicadas novas versões e divulgadas erratas em sala de aula. É responsabilidade do aluno manter-se informado, freqüentando as aulas ou acompanhando as novidades na página da disciplina na sala virtual.**