

Introdução ao agrupamento de dados

Mineração de Dados

Ronaldo C. Prati

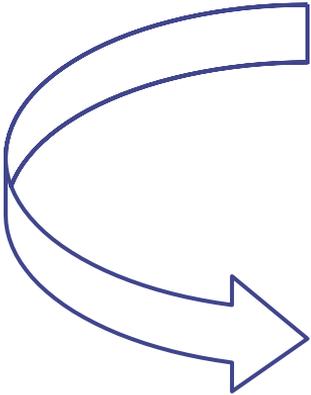
Aprendizado não supervisionado

- No aprendizado supervisionado, temos dados associados a uma variável de interesse
 - classificação: variável discreta
 - regressão: variável numérica
- No aprendizado não-supervisionado, não temos essa variável de interesse
- Objetivos é encontrar padrões de interesse nos dados
 - Análise de agrupamentos
 - Regras de associação

1. Motivação

Diversas ciências se baseiam na organização de objetos de acordo com suas similaridades;

- Biologia: Reino: *Animalia*
Ramo: *Chordata*
Classe: *Mammalia*
Ordem: *Primatas*
Família: *Hominidae*
Gênero: *Homo* (homem moderno e parentes)
Espécie: *Homo sapiens*



Humanos se interessam por *categorizações*:



stk325153rkn
www.fotosearch.com.br

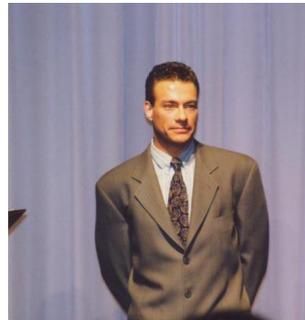
□ Música: erudita, popular, religiosa, etc..

□ Filmes:

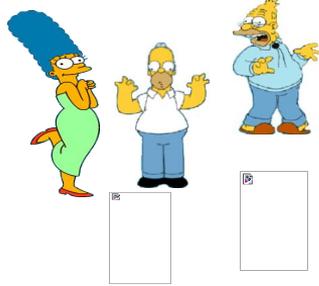
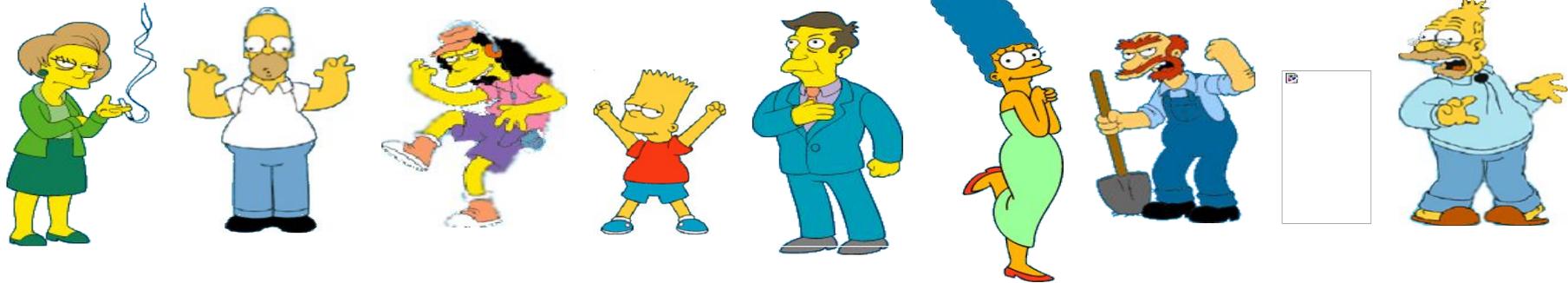
- Animação, Aventura, Comédia, Drama, Musical, etc...



- Agrupamento de atores:



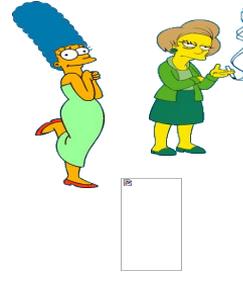
Como agrupar naturalmente os seguintes objetos?



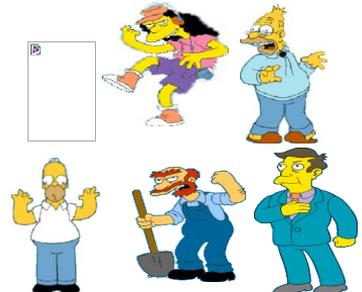
Família



Empregados



Mulheres



Homens

→ *Grupo* é um conceito subjetivo!

O que é um grupo ?

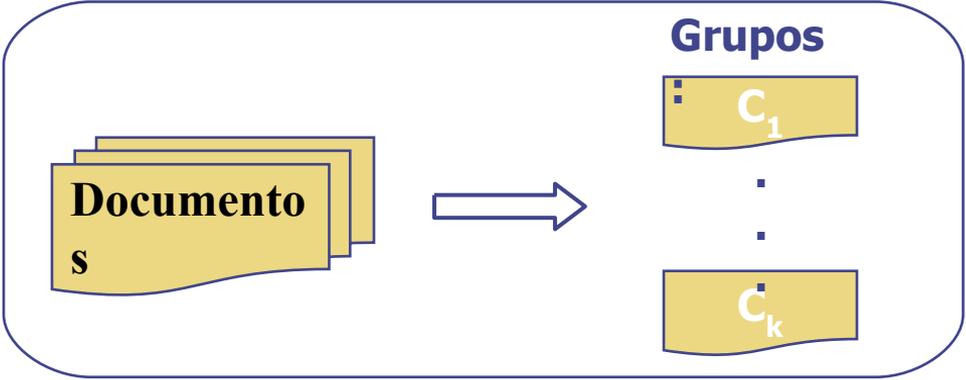
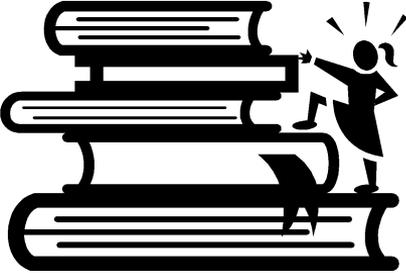
- Definições subjetivas:
 - “Semelhanças entre objetos”.
 - Quais atributos devemos considerar para computar similaridades?



- Numa “abordagem matemática”, critérios numéricos usualmente consideram:
 - Homogeneidade (coesão interna);
 - Heterogeneidade (separação entre grupos);

- Apesar das dificuldades apresentadas, a literatura sobre ADs é rica e bem estabelecida;
- Há medidas de dis(similaridade) bem estudadas e fundamentadas para diversos tipos de dados (e também para diversos domínios de aplicação):
 - Dados Numéricos;
 - Dados Categóricos / Nominais;
 - Dados Binários;
 - Etc.

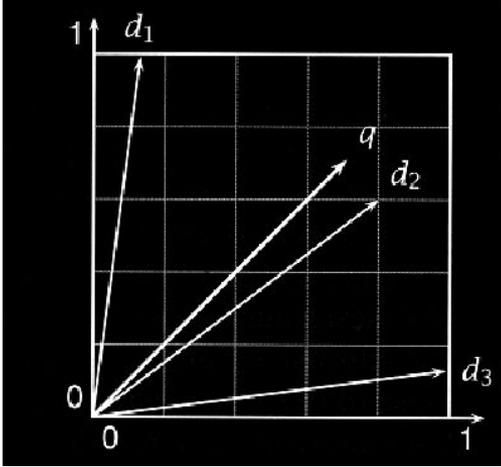
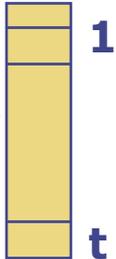
- Mineração de Textos:



Como?

Documento:
Bag-of-words

Vetor
de
Palavras



Aplicações?

Agrupamento de Dados (ADs) é uma técnica importante para *Análise Exploratória de Dados* :

- Engenharia;
- Biologia;
- Psicologia;
- Medicina;
- Administração (*Marketing* , Finanças,...);
- Ciência da Computação:
 - Bioinformática;
 - Dados coletados via sensores;
 - Componentes de sistemas inteligentes;
 - Componentes de algoritmos para aprendizado de máquina, ...

Data science/machine learning methods you used in the past 12 months: [732 voters]

(<https://www.kdnuggets.com/2017/12/top-data-science-machine-learning-methods.html>)

Regression 60%

Clustering 55%

Decision Trees/Rules 51%

K-NN 39%

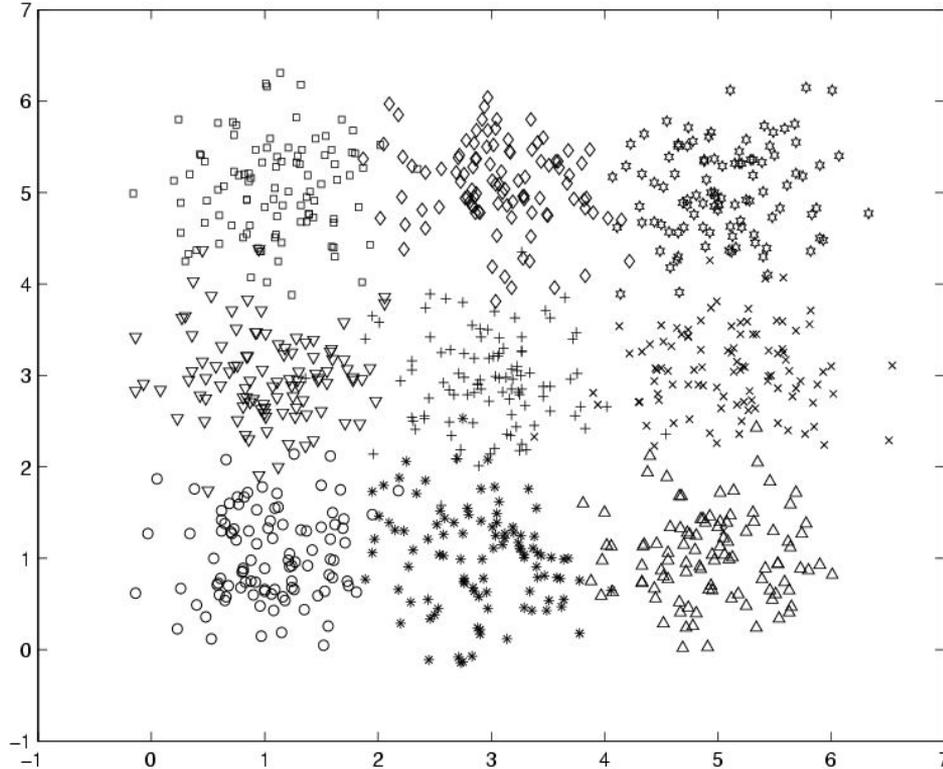
...

2. Conceitos Básicos

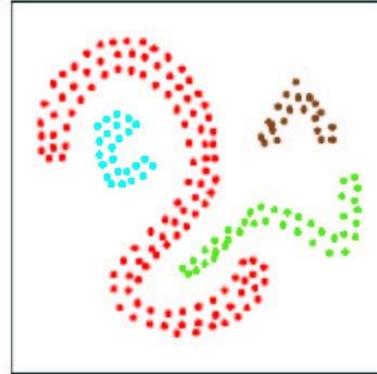
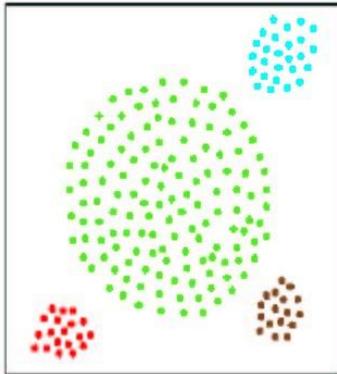
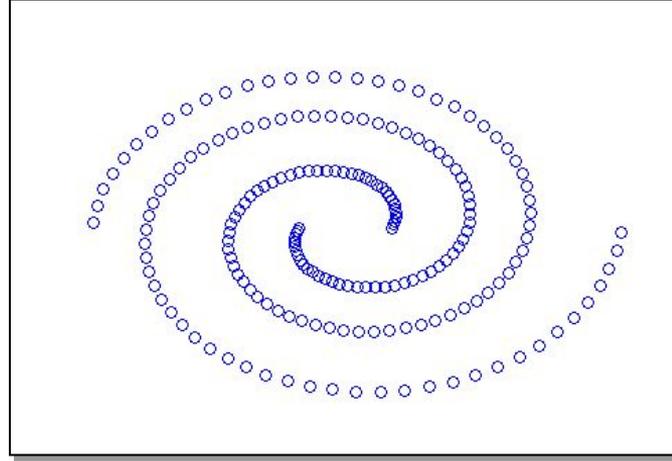
Algumas Definições (Everitt, 1974):

- *Um cluster (grupo) é um conjunto de entidades semelhantes, e entidades pertencentes a diferentes clusters não são semelhantes.*
 - *Um grupo é uma aglomeração de pontos no espaço tal que a distância entre quaisquer dois pontos no grupo é menor do que a distância entre qualquer ponto no grupo e qualquer ponto fora deste.*
 - *Grupos podem ser descritos como regiões conectadas de um espaço multidimensional contendo uma densidade de pontos relativamente alta, separada de outras tais regiões por uma região contendo uma densidade relativamente baixa de pontos.*
- **Humanos reconhecem clusters no plano quando os vêem, sem saber explicar exatamente porquê (Jain & Dubes, 1988),...**

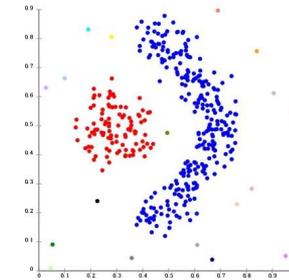
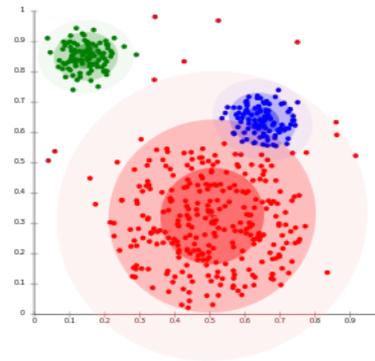
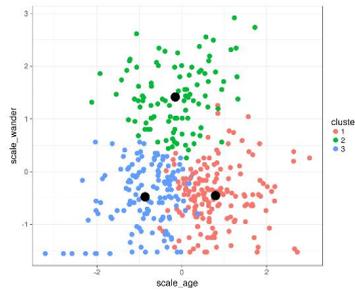
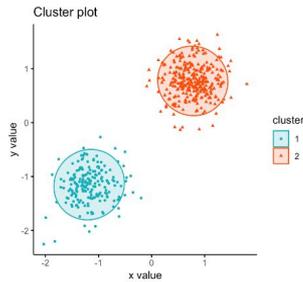
Quais são os grupos ?



Quais são os grupos ? ...



- Grupos podem ter diferentes tamanhos, formas e densidades
- Grupos podem formar uma hierarquia
- Grupos podem ter sobreposição ou serem disjuntos

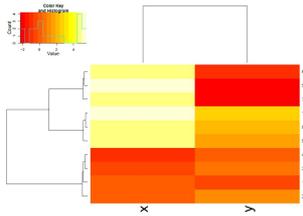


- Algoritmos para agrupamento de dados induzem *clusters*;
- Conceito semelhante ao de tendência (*bias*) indutiva estudado em aprendizado de máquina;
- Medidas de dis(similaridade), índices de validade relativos, parâmetros definidos pelo usuário, etc. (dependente do domínio/problema)
- Sob o ponto de vista de AM: *projetista define o que o computador pode aprender.*
- *Existem centenas de algoritmos...*

Tipos de agrupamento

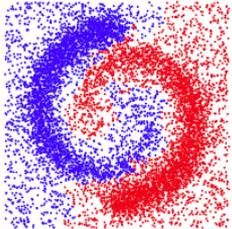
Baseado em ligações

e.g. Agrupamento hierárquico



Baseado em densidade

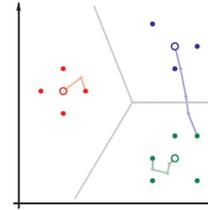
e.g. DBSCAN



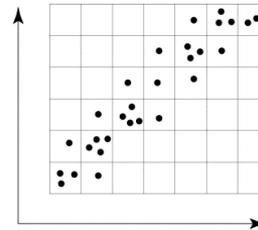
Agrupamento por

partição

e.g. k-Means



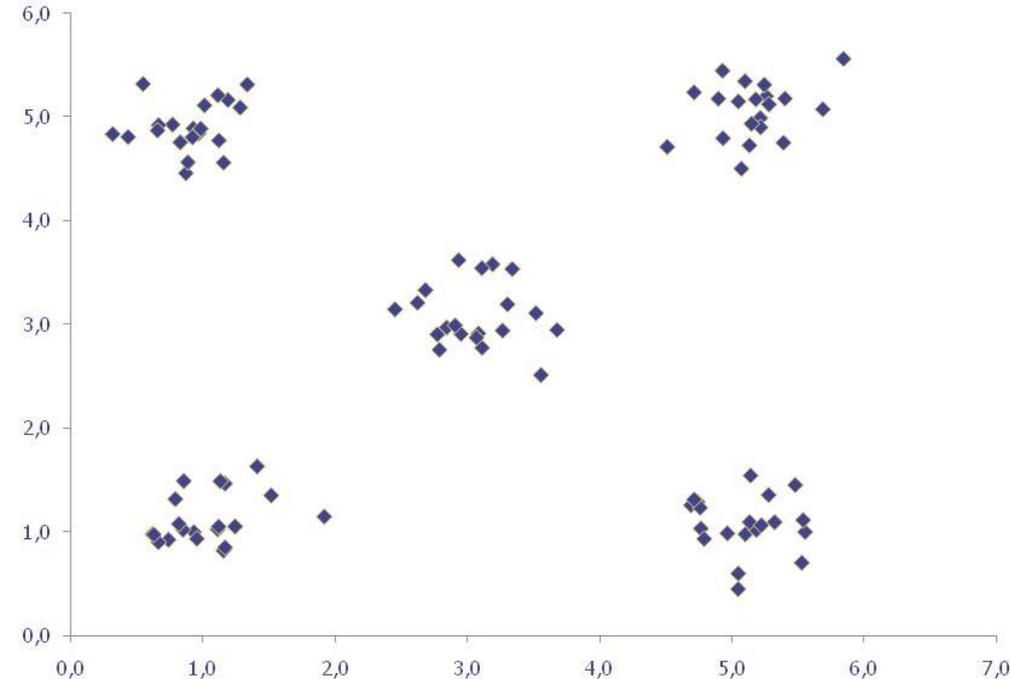
Baseado em Grid



Independendentemente do método o principal objetivo do agrupamento de dados é:

- Maximizar a homogeneidade interna ao *cluster* e a heterogeneidade entre diferentes *clusters*.
- Objetos que pertencem ao mesmo *cluster* devem ser mais semelhantes entre si do que em relação a objetos de outros clusters;
- Medidas de (Dis)similaridade possuem polarizações (biases):
 - vantagens/desvantagens dependentes do domínio: Distância Euclidiana, Correlação de Pearson, Coseno, etc.

Agrupamento X Classificação?

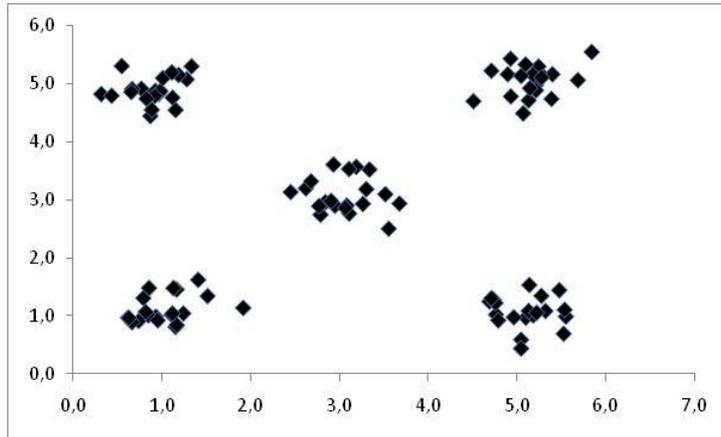
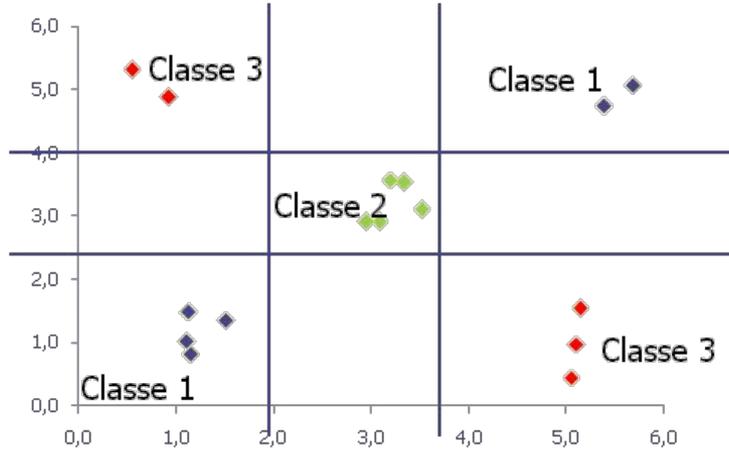


Agrupamento:
Indução de grupos
a partir da base de
dados...

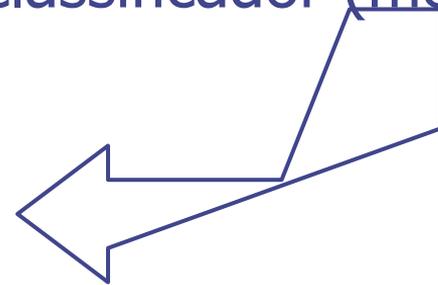


□ Grupos obtidos serão então cuidadosamente estudados.

Agrupamento X Classificação?



Base de treinamento com dados rotulados:
□ classificador (modelo)



Rotular dados de teste em função do modelo obtido.