

# Aula 13 – Mineração de Dados

## Regras de Associação - Parte2

Profa. Elaine Faria

UFU

# Material

- Este material foi construído por meio de traduções dos slides do do prof. Tan, disponíveis em:
  - <https://www-users.cse.umn.edu/~kumar001/dmbook/index.php>
  - Todas as figuras usadas foram retiradas dos slides do prof. Tan
- O material é baseado no livro
  - Tan P., SteinBack M. e Kumar V. Introduction to Data Mining, Pearson, 2006.
  - Faceli, K., Lorena, A. C., Gama, J., Carvalho, A. C. P. L. F., Inteligência Artificial: Uma abordagem de Aprendizado de Máquina, LTC, 2011.

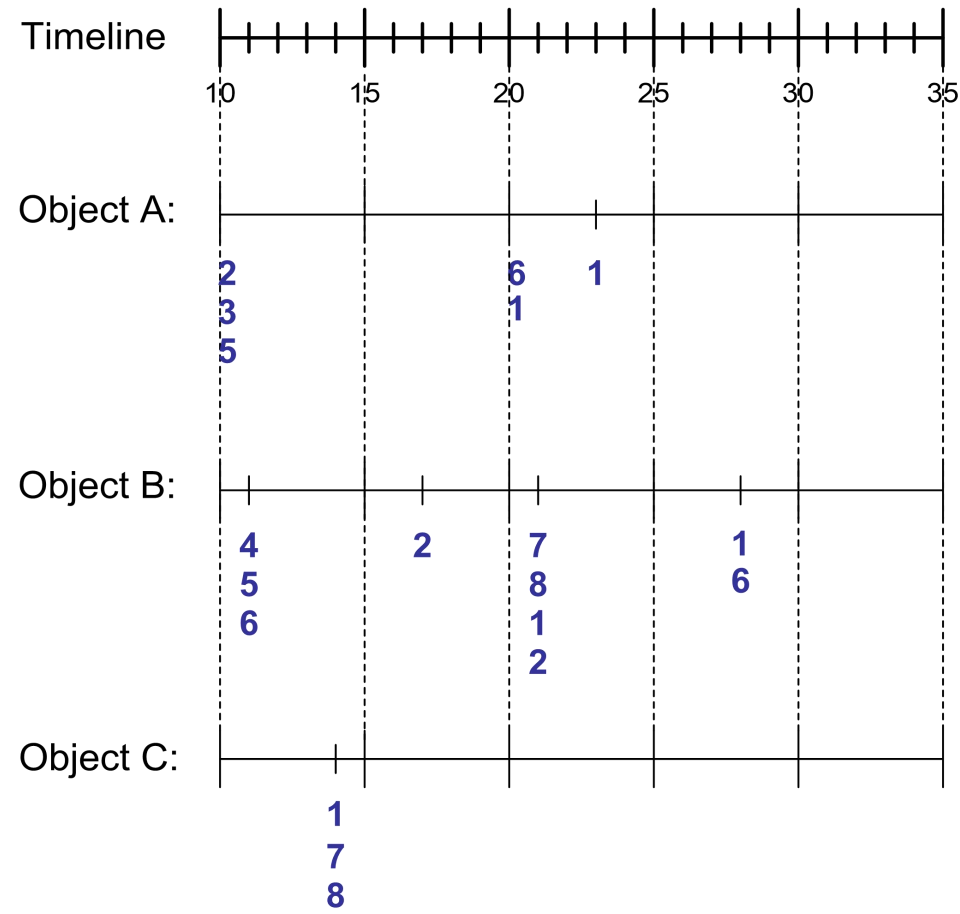
# Padrões Sequenciais

- As cesta de compras do supermercado contém uma informação temporal sobre quando um item foi comprado pelo cliente
  - Pode ser usado para colocar junto uma sequência de transações feitas por um cliente em um determinado período de tempo
- Conhecer a evolução das compras de seus clientes
  - EX: “que sequência de produtos são comprados por um mesmo cliente em momentos consecutivos”?
  - Se você pudesse descobrir que uma sequência de produtos  $\langle p1, p2, p3 \rangle$  é frequentemente comprada nesta ordem pelos clientes, você poderia enviar folhetos promocionais envolvendo os produtos  $p2$  ou  $p3$  para aqueles clientes que compraram o produto  $p1$ .

# Exemplo de uma base de dados sequências

Base de Dados de Sequência

Sequence ID	Timestamp	Events
A	10	2, 3, 5
A	20	6, 1
A	23	1
B	11	4, 5, 6
B	17	2
B	21	7, 8, 1, 2
B	28	1, 6
C	14	1, 8, 7



# Exemplo de sequências

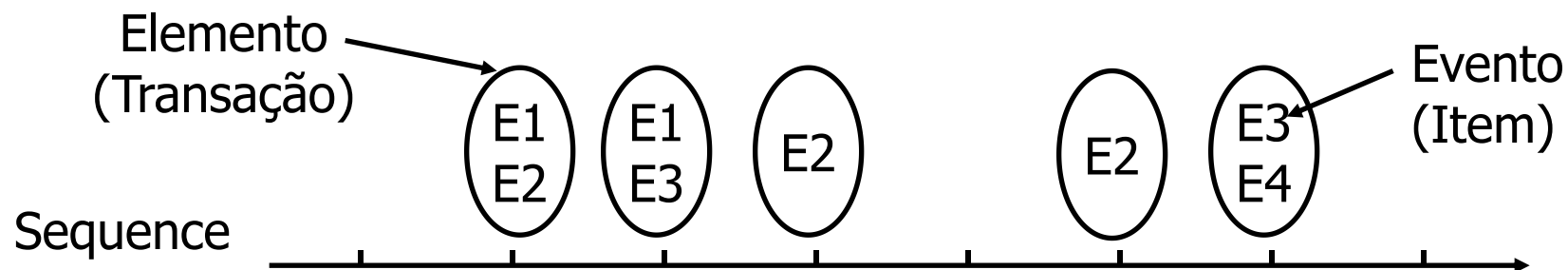
- Sequência de diferentes transações de um cliente em uma loja online:  
    < {Camera Digital,iPad} {cartão de memória} {headphone,capa de iPad} >
- Sequência de livros retirados em uma biblioteca:  
    <{Sociedade do Anel} {As Duas Torres} {Retorno do Rei}>

# Descoberta de padrões sequenciais: Exemplo

- Sequências de transações de ponto de vendas
  - Livraria de informática:  
(Intro\_To\_Visual\_C) (C++\_Primer) -->  
(Perl\_for\_dummies,Tcl\_Tk)
  - Loja de roupas esportivas:  
(Shoes) (Racket, Racketball) --> (Sports\_Jacket)

# Exemplo de elementos e eventos em uma base de dados sequencial

Bases de Dados	Sequência	Elemento (Transação)	Evento (Item)
Cliente	Histórico de compras de um determinado cliente	Um conjunto de itens comprados por um cliente no momento t	Livros, produtos diários, CDs, etc.
Dados da Web	Atividade de navegação de um determinado visitante da web	Uma coleção de arquivos visualizados por um visitante da web após um único clique do mouse	Página inicial, página de índice, informações de contato,
Dados de Evento	Histórico de eventos gerados por um determinado sensor	Eventos disparados por um sensor no tempo t	Tipos de alarmes gerados por sensores



# Dados de sequência vs. Dados de cesta de compras

Base de Dados de Sequência

Customer	Date	Items bought
A	10	2, 3, 5
A	20	1,6
A	23	1
B	11	4, 5, 6
B	17	2
B	21	1,2,7,8
B	28	1, 6
C	14	1,7,8

Dados de cesta de compras

Events
2, 3, 5
1,6
1
4,5,6
2
1,2,7,8
1,6
1,7,8



# Definição formal de uma sequência

- Uma sequência (ou padrão sequencial) é uma lista ordenada de elementos (itemsets)

$$s = \langle e_1, e_2, e_3 \dots \rangle$$

- Cada elemento contém uma coleção de eventos (items)

$$e_i = \{i_1, i_2, \dots, i_k\}$$

- Comprimento da sequência,  $|s|$ , é dada pelo nro de elementos na sequência
- Uma k-sequência é uma sequência que contém k eventos (items)
  - $\langle \{a, b\} \{a\} \rangle$  tem o comprimento 2 e é uma 3-sequência
    - Essa definição é usada pelo algoritmo GSP, mas é ligeiramente diferente no algoritmo AprioriAll

# Definição formal de uma sub-sequência

- Uma sequência  $t: \langle a_1 a_2 \dots a_n \rangle$  **está contida** em outra sequência  $s: \langle b_1 b_2 \dots b_m \rangle$  ( $m \geq n$ ) se existe inteiros :  $i_1 < i_2 < \dots < i_n$  tal que  $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$

- Exemplo ilustrativo:

$s:$ 
 $\quad b_1 \quad b_2 \quad b_3 \quad b_4 \quad b_5$   
 $t:$ 
 $\quad \quad a_1 \quad a_2 \quad \quad a_3$

$t$  é uma **subsequência** de  $s$  se  $a_1 \subseteq b_2, a_2 \subseteq b_3, a_3 \subseteq b_5$ .

Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{8\} \rangle$	$\langle \{2\} \{8\} \rangle$	Yes
$\langle \{1,2\} \{3,4\} \rangle$	$\langle \{1\} \{2\} \rangle$	No
$\langle \{2,4\} \{2,4\} \{2,5\} \rangle$	$\langle \{2\} \{4\} \rangle$	Yes
$\langle \{2,4\} \{2,5\} \{4,5\} \rangle$	$\langle \{2\} \{4\} \{5\} \rangle$	No
$\langle \{2,4\} \{2,5\} \{4,5\} \rangle$	$\langle \{2\} \{5\} \{5\} \rangle$	Yes
$\langle \{2,4\} \{2,5\} \{4,5\} \rangle$	$\langle \{2, 4, 5\} \rangle$	No

# Mineração de Padrão Sequencial: Definição

- **Dados:**

- uma base de sequências
- um limite mínimo de suporte especificado pelo usuário, *minsup*

- **Tarefa:**

- Encontrar todas as subsequências com suporte  $\geq minsup$

- O suporte de uma subsequência  $w$  é definido como a fração de todas as sequências de dados que contém  $w$
- Um *padrão sequencial* é uma subsequência frequente (i.e., uma subsequência cujo suporte é  $\geq minsup$ )

# Mineração de Padrão Sequencial: Exemplo

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

Examplos de Subsequências Frequentes:

< {1,2} >	s=60%
< {2,3} >	s=60%
< {2,4}>	s=80%
< {3} {5}>	s=80%
< {1} {2} >	s=80%
< {2} {2} >	s=60%
< {1} {2,3} >	s=60%
< {2} {2,3} >	s=60%
< {1,2} {2,3} >	s=60%

# Base de dados sequencial vs. Cesta de compra

Base de Dados de Sequências

Customer	Date	Items bought
A	10	2, 3, 5
A	20	1,6
A	23	1
B	11	4, 5, 6
B	17	2
B	21	1,2,7,8
B	28	1, 6
C	14	1,7,8

$\{2\} \rightarrow \{1\}$

$$\text{conf}(\{2\} \rightarrow \{1\}) = \frac{\sigma(\{2\} \{1\})}{\sigma(\{2\})}$$

Base de Dados de cesta de compras

Eventos
2, 3, 5
1,6
1
4,5,6
2
1,2,7,8
1,6
1,7,8

$(1,8) \rightarrow (7)$

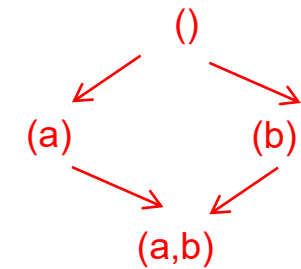
$$\text{conf}(1,8) \rightarrow (7)) = \frac{\sigma(1,7,8)}{\sigma(\{1,8\})}$$

# Extraindo padrões sequenciais

- Dados n eventos:  $i_1, i_2, i_3, \dots, i_n$
- 1-subsequência candidata  
 $\langle \{i_1\} \rangle, \langle \{i_2\} \rangle, \langle \{i_3\} \rangle, \dots, \langle \{i_n\} \rangle$
- 2-subsequências candidatas:  
 $\langle \{i_1, i_2\} \rangle, \langle \{i_1, i_3\} \rangle, \dots,$   
 $\langle \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_2\} \rangle, \dots, \langle \{i_n\} \{i_n\} \rangle$
- 3-subsequências candidatas:  
 $\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, \dots,$   
 $\langle \{i_1, i_2\} \{i_1\} \rangle, \langle \{i_1, i_2\} \{i_2\} \rangle, \dots,$   
 $\langle \{i_1\} \{i_1, i_2\} \rangle, \langle \{i_1\} \{i_1, i_3\} \rangle, \dots,$   
 $\langle \{i_1\} \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_1\} \{i_2\} \rangle, \dots$

# Extraindo padrões sequenciais: Exemplo simples

- Dados 2 eventos: a, b
- 1-subsequência candidata  
 $\langle \{a\} \rangle, \langle \{b\} \rangle$
- 2-subsequências candidatas  
 $\langle \{a\} \{a\} \rangle, \langle \{a\} \{b\} \rangle, \langle \{b\} \{a\} \rangle, \langle \{b\} \{b\} \rangle, \langle \{a, b\} \rangle.$
- 3-subsequências candidatas  
 $\langle \{a\} \{a\} \{a\} \rangle, \langle \{a\} \{a\} \{b\} \rangle, \langle \{a\} \{b\} \{a\} \rangle, \langle \{a\} \{b\} \{b\} \rangle,$   
 $\langle \{b\} \{b\} \{b\} \rangle, \langle \{b\} \{b\} \{a\} \rangle, \langle \{b\} \{a\} \{b\} \rangle, \langle \{b\} \{a\} \{a\} \rangle$   
 $\langle \{a, b\} \{a\} \rangle, \langle \{a, b\} \{b\} \rangle, \langle \{a\} \{a, b\} \rangle, \langle \{b\} \{a, b\} \rangle$



Item-set patterns