

---

# Compressão de Arquivos Orientados a Coluna com PPM

Vinicius F. Garcia<sup>1</sup>, Sergio L. S. Mergen<sup>1</sup>

<sup>1</sup>Centro de Tecnologia  
Universidade Federal de Santa Maria, Brasil

Escola Regional de Banco de Dados, Londrina, Brasil

14 de abril de 2016



- 1 Introdução
  - Motivação
  - Investigação
- 2 Prediction by Partial Matching (PPM)
  - Metodologia do PPM
  - Análise do PPM
- 3 Avaliação
  - Algoritmos avaliados
  - Dados avaliados
  - Resultados
- 4 Conclusões
- 5 Trabalhos futuros



## Dados orientados a coluna

- Crescente popularização
- Abordagem NoSQL
- Consultas analíticas em poucas colunas

## Compressão de dados

- Utilização econômica de memória
- Melhor aproveitamento de servidores de BD
- Necessidade de algoritmos eficientes



# Dados orientados a coluna

---

<i>id</i>	<i>nome</i>	<i>sal.</i>	<i>inicio</i>
1	João Dias	2000	10-2014
2	Ana Galo	3500	05-2012
3	J. Andre	2800	12-2015

Sistema orientado a registros



$arq_1 \Rightarrow 1, 2, 3$

$arq_2 \Rightarrow$  João Dias, Ana Galo, J. Andre

$arq_3 \Rightarrow 2000, 3500, 2800$

$arq_4 \Rightarrow 10-2014, 05-2012, 12-2015$

Sistema orientado a coluna



## Objetivo

Determinar se algoritmos da família dos PPMs (Prediction by Partial Matching) são adequados e eficientes em relação aos demais testados para dados orientados a coluna.



## Informações gerais

- O método PPM codifica um símbolo de cada vez
- Utilização de um histórico chamado de contexto
- O primeiro símbolo após o contexto é o alvo da codificação
- A organização do contexto é tipicamente implementada com árvores
- A árvore gera um intervalo de probabilidade de aparição do símbolo

A C A B A \_ A \_ A B A \_ C  
Contexto



## Busca e codificação

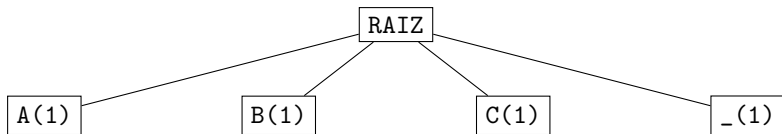
- Busca pelo símbolo alvo nível a nível que estejam disponíveis
- Níveis com símbolo alvo ausente resultam na codificação de um escape
- Entropia realizada por codificação aritmética

A C A B A \_ A \_ A B A \_ C



## Exemplificação

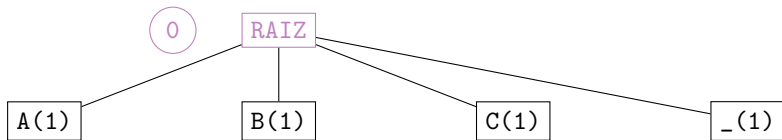
↓  
A C A B A \_ A \_ A B A \_ C





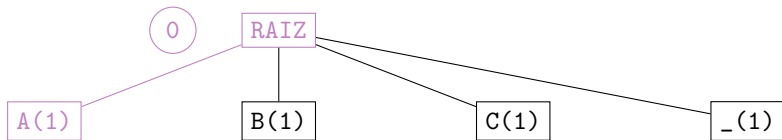
## Exemplificação

↓  
A C A B A \_ A \_ A B A \_ C



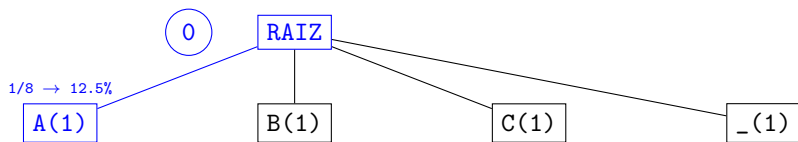
## Exemplificação

↓  
A C A B A \_ A \_ A B A \_ C



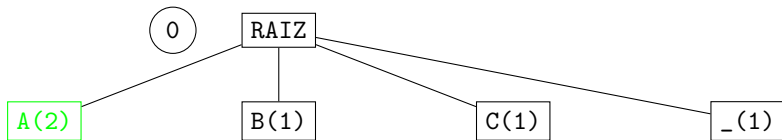
## Exemplificação

↓  
A C A B A \_ A \_ A B A \_ C

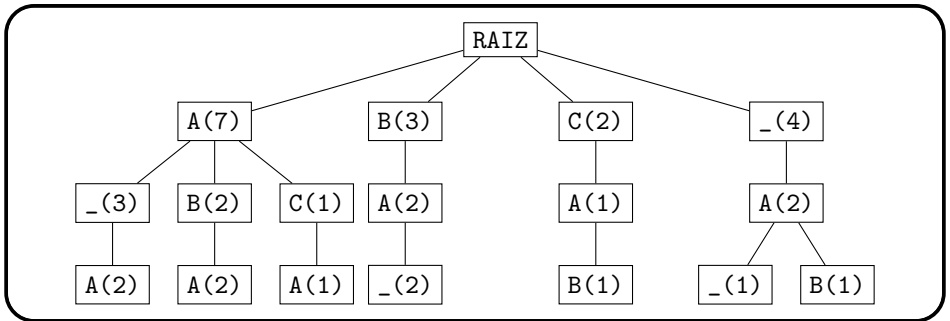


## Exemplificação

↓  
A C A B A \_ A \_ A B A \_ C



## Exemplificação

$$\underbrace{A}_{Contexto} \quad \downarrow \quad C$$


## Exemplificação

```

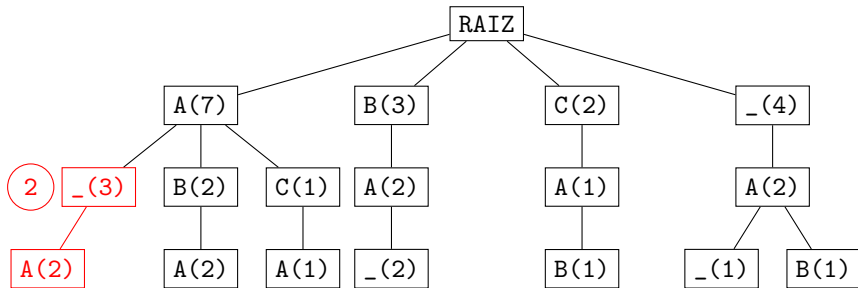
graph TD
    RAIZ[RAIZ] --> A7[A(7)]
    RAIZ --> B3[B(3)]
    RAIZ --> C2[C(2)]
    RAIZ --> U4[_ (4)]
    A7 --> 2((2))
    A7 --> U3[_ (3)]
    A7 --> B2[B(2)]
    2 --> A2_2[A(2)]
    U3 --> A2_U3[A(2)]
    B2 --> A2_B2[A(2)]
    B3 --> A2_B3[A(2)]
    A2_B3 --> U2[_ (2)]
    C2 --> A1[C(1)]
    A1 --> B1_1[B(1)]
    U4 --> A2_U4[A(2)]
    A2_U4 --> U1[_ (1)]
    A2_U4 --> B1_2[B(1)]
  
```

## Exemplificação

```

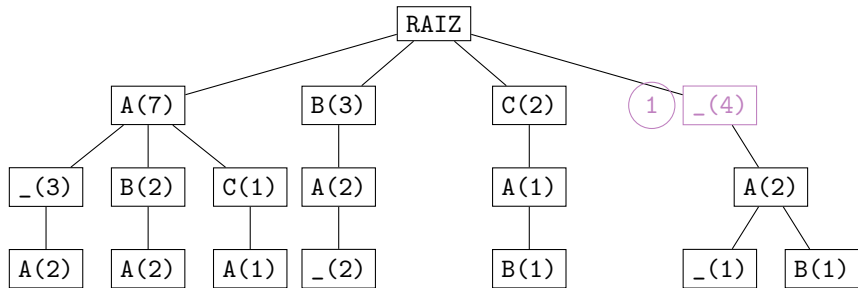
graph TD
    RAIZ[RAIZ] --> A7[A(7)]
    RAIZ --> B3[B(3)]
    RAIZ --> C2[C(2)]
    RAIZ --> U4[_ (4)]
    
    A7 --> 2((2))
    2 --> U3[_ (3)]
    U3 --> A2_1[A(2)]
    A7 --> B2[B(2)]
    B2 --> A2_2[A(2)]
    A7 --> C1[C(1)]
    C1 --> A1_1[A(1)]
    
    B3 --> A2_3[A(2)]
    A2_3 --> U2[_ (2)]
    
    C2 --> A1_2[A(1)]
    A1_2 --> B1_1[B(1)]
    
    U4 --> A2_4[A(2)]
    A2_4 --> U1[_ (1)]
    A2_4 --> B1_2[B(1)]
  
```

## Exemplificação

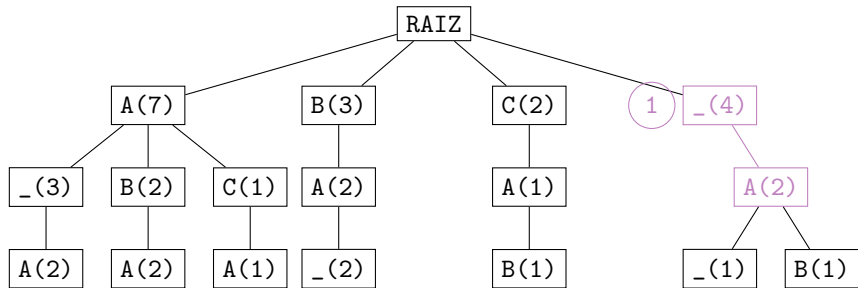
$$\underbrace{A \quad -}_{Contexto} \quad \downarrow \quad C$$




## Exemplificação

$$\underbrace{A \quad -}_{Contexto} \quad \downarrow \quad C$$


## Exemplificação

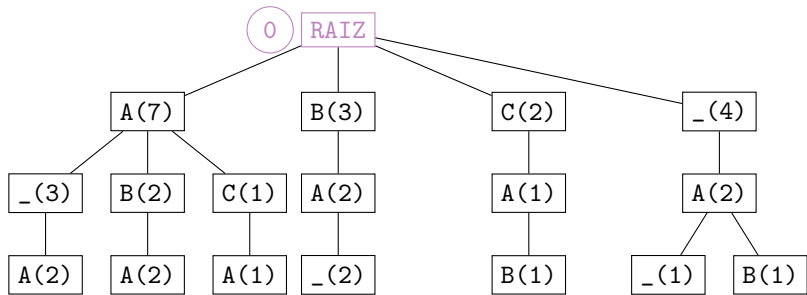
$$\underbrace{A \quad -}_{Contexto} \quad \downarrow \quad C$$


## Exemplificação

## Contexto



## Exemplificação

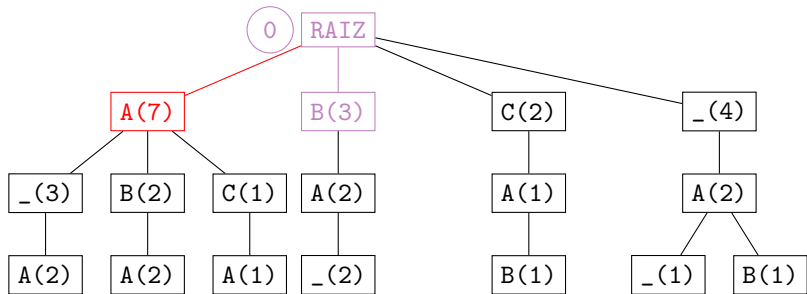
$$\underbrace{A \quad -}_{Contexto} \quad \downarrow \quad C$$


## Exemplificação

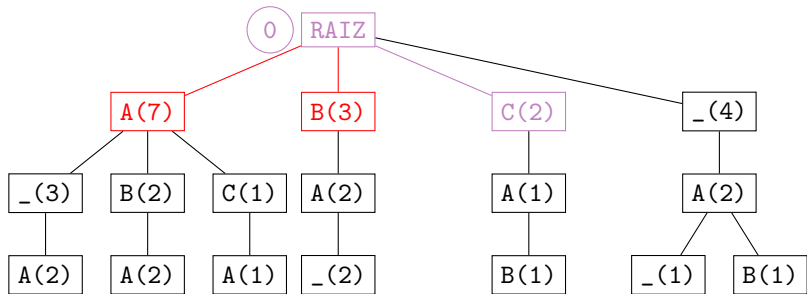
A C A B A \_ A \_ A B



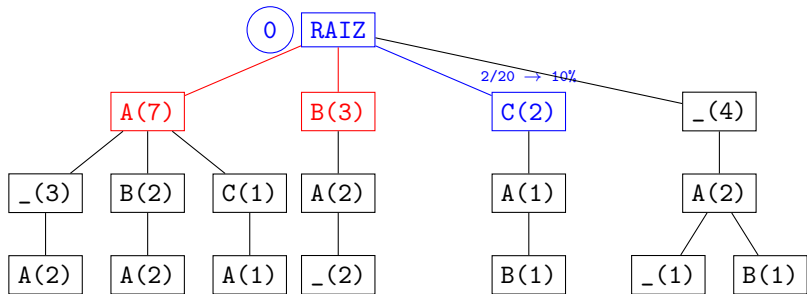
## Exemplificação

$$\underbrace{A \quad -}_{Contexto} \quad \downarrow \quad C$$


## Exemplificação

$$\underbrace{A \quad -}_{Contexto} \quad \downarrow \quad C$$


## Exemplificação

$$\underbrace{A \quad -}_{Contexto} \quad \downarrow \quad C$$




## Exemplificação

```

graph TD
    RAIZ[RAIZ] --> A7[A(7)]
    RAIZ --> B3[B(3)]
    RAIZ --> C3[C(3)]
    RAIZ --> M4[_(4)]
    A7 --> M3[_(3)]
    A7 --> B2[B(2)]
    A7 --> C1[C(1)]
    M3 --> A2_3[A(2)]
    M3 --> C1_3[C(1)]
    B2 --> A2_2[A(2)]
    C1 --> A1_1[A(1)]
    B3 --> A2_1[A(2)]
    A2_1 --> A1_2[A(1)]
    A2_1 --> M2[_(2)]
    C3 --> A1_3[A(1)]
    M4 --> A2_4[A(2)]
    M4 --> C1_4[C(1)]
    A2_4 --> M1[_(1)]
    A2_4 --> B1[B(1)]
  
```

## Taxa de compressão

- Forte dependência da existência de padrões
- Sequências que costumam aparecer juntas utilizarão poucos bits
- Situação comum em dados textuais

## Desempenho e uso de memória

- Quanto maior o número de padrões, mais rápida a compressão e descompressão
- Quanto maior o número de padrões, menor o custo em memória



# Algoritmos avaliados

---

<b>Algoritmo</b>	<b>Implementação</b>	<b>Programa</b>
BWT	C	BZip2
Cod. Aritimética	C	Próprio
LZ	C	GZip
PPM	C	Próprio



## TPC-H

- Benchmark para avaliação de performance em transações de BDs
- Banco de dados gerado de 1GB
- Adaptação do TPC-H para geração de dados orientados a coluna

## Calgary corpus

- Benchmark para avaliação de compressores de dados



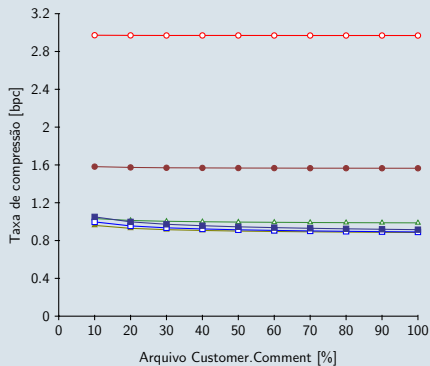
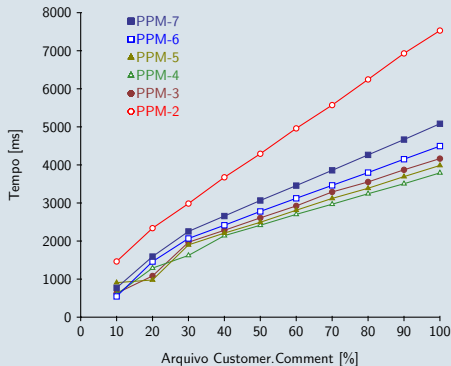
## Objetivo

Determinar o melhor tamanho de contexto para compressão de dados orientados a coluna.



# Determinação do tamanho de contextos

## Resultado



# Análise dos algoritmos sob dados orientados a coluna

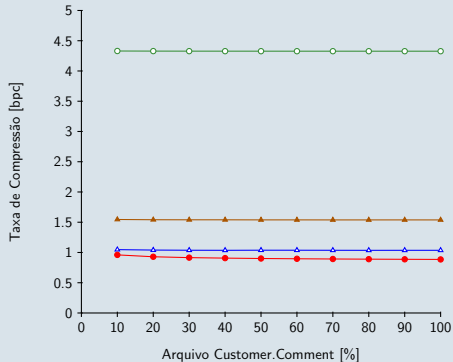
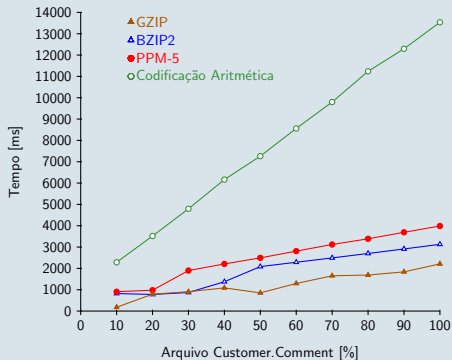
## Objetivo

Determinar o comportamento dos algoritmos propostos anteriormente com arquivos orientados a coluna.



# Análise dos algoritmos sob dados orientados a coluna

## Resultado





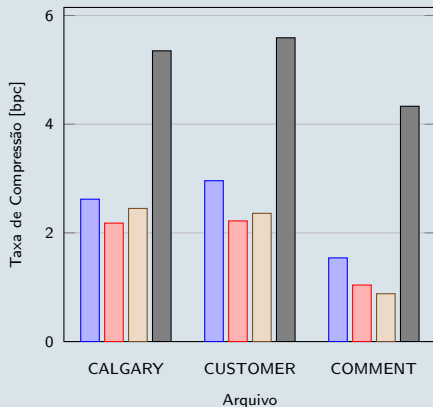
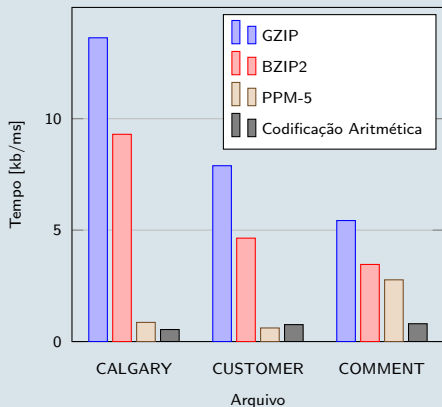
## Objetivo

Determinar o comportamento dos algoritmos propostos anteriormente com arquivos variados presentes no corpus Calgary e em dados orientados a registros.



# Análise dos algoritmos sob dados variados

## Resultado

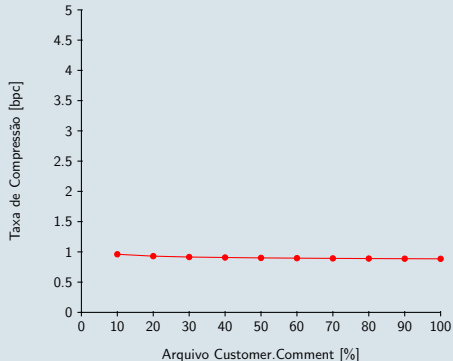
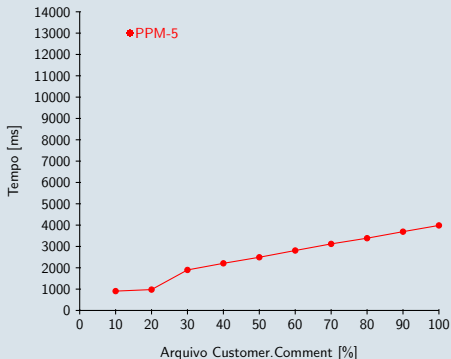


- Arquivos orientados a coluna são bem explorados por compressores baseados em padrões
- O método PPM mostrou um custo benefício aceitável na relação compressão/desempenho
- O método PPM mostra-se com aplicação viável sob dados orientados a coluna



## Atualização da árvore de contexto

- Analisar métricas de compressão (bpc) durante a execução do algoritmo
- Determinar pontos de não atualização da árvore
- Otimizar o tempo de execução do método PPM



# Compressão de Arquivos Orientados a Coluna com PPM

---

## Perguntas?

Universidade Federal de Santa Maria  
Santa Maria - RS

Vinícius Fülber Garcia  
vfulber@inf.ufsm.br

Sérgio L. S. Mergen  
mergen@inf.ufsm.br

