

# Documentação Técnica - Ordenação Externa com Merge Sort

## Visão Geral

Este documento apresenta a implementação de um algoritmo de **ordenação externa** baseado no **Merge Sort** para arquivos CSV grandes que excedem a capacidade da memória RAM disponível.

## Estruturas de Dados Utilizadas

### 1. Classe OrdenacaoExterna

- **Atributos principais:**
  - `tamanho_buffer`: Controla o número máximo de registros carregados na memória
  - `arquivos_temporarios`: Lista para rastrear arquivos temporários criados
  - `diretorio_temp`: Diretório temporário para armazenar runs intermediários
  - `coluna_chave_atual`: Referência à coluna utilizada para ordenação

### 2. Estruturas de Dados Auxiliares

- **Listas Python**: Utilizadas como buffer para armazenar registros em memória
- **CSV Readers/Writers**: Para manipulação eficiente de arquivos CSV
- **Arquivos temporários**: Runs ordenados armazenados em disco

## Divisão de Módulos

### Módulo Principal (`mergesort_externo.py`)

Contém a implementação completa do algoritmo dividida em métodos especializados:

#### Métodos Públicos:

- `ordenar_arquivo()`: Interface principal para ordenação externa
- `__init__()`: Inicialização da classe com configurações

## Métodos Privados de Divisão:

- `_dividir_em_runs()`: Divide o arquivo original em runs menores
- `_salvar_run()`: Ordena e salva um run individual
- `_merge_sort_interno()`: Implementa merge sort para dados em memória
- `_merge_interno()`: Mescla duas listas ordenadas em memória

## Métodos Privados de Merge Externo:

- `_merge_externo()`: Coordena o merge de todos os runs
- `_merge_dois_runs()`: Mescla dois arquivos de runs
- `_comparar_registros()`: Compara registros baseado na chave de ordenação

## Métodos Auxiliares:

- `_limpar_arquivos_temporarios()`: Remove arquivos temporários
- `_get_coluna_chave()`: Retorna a coluna chave atual

## Módulos de Teste e Demonstração

- `demo.py`: Demonstração simples com dados pequenos
- `exemplo_teste.py`: Testes abrangentes com dados maiores

# Descrição das Rotinas e Funções

## 1. Fase de Divisão

`_dividir_em_runs()`

**Propósito:** Divide o arquivo original em runs menores que cabem na memória.

**Algoritmo:**

1. Abre o arquivo CSV original
2. Lê o cabeçalho para identificar colunas
3. Carrega registros até atingir o limite do buffer
4. Quando o buffer está cheio, chama `_salvar_run()`
5. Repete até processar todo o arquivo

**Complexidade:**  $O(n)$  onde  $n$  é o número de registros

`_salvar_run()`

**Propósito:** Ordena um buffer de dados e salva como arquivo temporário.

**Algoritmo:**

1. Chama `_merge_sort_interno()` para ordenar o buffer
2. Cria arquivo temporário único
3. Escreve cabeçalho e registros ordenados
4. Adiciona arquivo à lista de temporários

**Complexidade:**  $O(k \log k)$  onde  $k$  é o tamanho do buffer

`_merge_sort_interno()`

**Propósito:** Implementa merge sort clássico para dados em memória.

**Algoritmo:**

1. Caso base: retorna se lista tem  $\leq 1$  elemento
2. Divide lista ao meio
3. Recursivamente ordena cada metade
4. Mescla as metades usando `_merge_interno()`

**Complexidade:**  $O(k \log k)$  onde  $k$  é o tamanho da lista

## 2. Fase de Merge Externo

`_merge_externo()`

**Propósito:** Coordena o merge de todos os runs até obter um arquivo final.

**Algoritmo:**

1. Enquanto houver mais de um run:
  - o Processa runs em pares
  - o Chama `_merge_dois_runs()` para cada par
  - o Substitui lista de runs pelos resultados mesclados
2. Retorna o único run restante

**Complexidade:**  $O(n \log r)$  onde  $r$  é o número de runs

`_merge_dois_runs()`

**Propósito:** Mescla dois arquivos de runs ordenados.

**Algoritmo:**

1. Abre os dois arquivos de runs e um arquivo de saída
2. Lê cabeçalhos e primeira linha de cada arquivo
3. Compara registros usando `_comparar_registros()`
4. Escreve o menor registro no arquivo de saída
5. Avança no arquivo correspondente
6. Repete até esgotar ambos os arquivos

**Complexidade:**  $O(m + n)$  onde  $m$  e  $n$  são os tamanhos dos runs

### 3. Funções Auxiliares

`_comparar_registros()`

**Propósito:** Compara dois registros baseado na chave de ordenação.

**Algoritmo:**

1. Extrai valores da coluna chave de ambos os registros
2. Tenta conversão numérica para comparação apropriada
3. Aplica lógica de ordenação (ascendente/descendente)
4. Retorna resultado da comparação

**Complexidade:**  $O(1)$

## Complexidades de Tempo e Espaço

### Complexidade de Tempo

Análise Geral:

- **Divisão em runs:**  $O(n)$  para leitura +  $O(n \log k)$  para ordenação interna
- **Merge externo:**  $O(n \log r)$  onde  $r$  é o número de runs
- **Total:**  $O(n \log n)$  - equivalente ao merge sort clássico

Detalhamento por Fase:

1. **Leitura e divisão:**  $O(n)$
2. **Ordenação interna dos runs:**  $O(n \log k)$  onde  $k = \text{tamanho\_buffer}$
3. **Merge externo:**  $O(n \log r)$  onde  $r = \lceil n/k \rceil$
4. **Resultado final:**  $O(n \log n)$

### Complexidade de Espaço

Memória RAM:

- **Espaço principal:**  $O(k)$  onde  $k = \text{tamanho\_buffer}$
- **Espaço adicional:**  $O(1)$  para variáveis auxiliares
- **Total em RAM:**  $O(k)$

Espaço em Disco:

- **Runs temporários:**  $O(n)$  - mesmo tamanho do arquivo original
- **Arquivo final:**  $O(n)$

- **Total em disco:**  $O(n)$  adicional durante o processo

# Problemas e Observações Encontrados

## 1. Problemas de Implementação

### Gerenciamento de Memória

- **Problema:** Controle preciso do uso de memória
- **Solução:** Implementação de buffer com tamanho configurável
- **Observação:** Tamanho do buffer afeta diretamente a performance

### Manipulação de Arquivos CSV

- **Problema:** Diferentes codificações e formatos
- **Solução:** Uso de encoding UTF-8 e tratamento de exceções
- **Observação:** Necessário validar formato dos dados de entrada

### Limpeza de Arquivos Temporários

- **Problema:** Possível acúmulo de arquivos temporários em caso de erro
- **Solução:** Uso de try/finally para garantir limpeza
- **Observação:** Implementado com diretório temporário do sistema

## 2. Desafios de Performance

### I/O Intensivo

- **Problema:** Muitas operações de leitura/escrita em disco
- **Impacto:** Performance dependente da velocidade do disco
- **Mitigação:** Uso de buffer adequado e minimização de operações

### Comparação de Tipos de Dados

- **Problema:** Dados podem ser numéricos ou texto
- **Solução:** Tentativa de conversão numérica com fallback para string
- **Observação:** Comparação consistente independente do tipo

## 3. Limitações Identificadas

### Escalabilidade

- **Limitação:** Espaço em disco deve ser pelo menos 2x o arquivo original
- **Impacto:** Não viável para sistemas com pouco espaço livre

### Tipos de Dados Suportados

- **Limitação:** Comparação básica entre números e strings
- **Extensão possível:** Suporte a datas, tipos customizados

## Estabilidade da Ordenação

- **Observação:** Algoritmo preserva ordem relativa de elementos iguais
- **Benefício:** Importante para ordenações secundárias

# Testes e Validação

## Casos de Teste Implementados

1. **Teste com dados pequenos:** Verificação de funcionamento básico
2. **Teste com buffer pequeno:** Simulação de limitação de memória
3. **Teste com diferentes colunas:** Validação de flexibilidade
4. **Teste com diferentes ordens:** Ascendente e descendente
5. **Teste com tipos mistos:** Numéricos e texto

## Resultados dos Testes

- **Funcionalidade:** Todos os testes passaram com sucesso
- **Performance:** Comportamento  $O(n \log n)$  confirmado
- **Robustez:** Tratamento adequado de casos extremos
- **Limpeza:** Arquivos temporários removidos corretamente

# Conclusão

## Resultados Obtidos

A implementação do algoritmo de ordenação externa baseado em Merge Sort atendeu com sucesso todos os requisitos especificados:

1. **Funcionalidade completa:** Ordena arquivos CSV grandes que não cabem na memória
2. **Flexibilidade:** Suporta ordenação por qualquer coluna e em ambas as direções
3. **Eficiência:** Mantém complexidade  $O(n \log n)$  do merge sort clássico
4. **Robustez:** Trata diferentes tipos de dados e cenários de erro
5. **Limpeza:** Gerencia adequadamente recursos temporários

## Vantagens da Solução

1. **Escalabilidade:** Processa arquivos maiores que a RAM disponível
2. **Estabilidade:** Preserva ordem relativa de elementos iguais
3. **Configurabilidade:** Tamanho do buffer ajustável conforme recursos

4. **Portabilidade:** Usa apenas bibliotecas padrão do Python

## Aplicações Práticas

Esta implementação é adequada para:

- Processamento de grandes datasets em análise de dados
- Ordenação de logs de sistema
- Preparação de dados para algoritmos que requerem entrada ordenada
- Sistemas com limitações de memória

## Trabalhos Futuros

Possíveis melhorias incluem:

- Suporte a múltiplas colunas de ordenação
- Otimização para SSDs (operações sequenciais vs. aleatórias)
- Interface gráfica para usuários não técnicos
- Suporte a diferentes formatos de arquivo (JSON, XML)
- Implementação paralela para sistemas multi-core

A solução desenvolvida demonstra com sucesso a aplicação prática dos conceitos de ordenação externa e algoritmos de divisão e conquista para resolver problemas reais de processamento de grandes volumes de dados.