

Trabalho Prático 2: Biblioteca Digital de Arendelle

Victor Hugo Silva Moura

Universidade Federal de Minas Gerais (UFMG)

Belo Horizonte – MG – Brasil

victorhugosmoura@gmail.com

Abstract. *The following text describes the “Arendelle’s Digital Library” problem and its solution using different implementations of the QuickSort algorithm. For each implementation, a series of tests were conducted in order to analyze and compare their respective performances.*

Resumo. *O texto a seguir descreve o problema “Biblioteca Digital de Arendelle” e sua resolução usando diferentes implementações do algoritmo QuickSort. Para cada implementação é feita uma série de testes a fim de analisar e comparar suas respectivas performances.*

1. Introdução

O problema *Biblioteca Digital de Arendelle* consiste em ordenar o acervo de livros e pergaminhos que o reino de Arendelle possui em suas bibliotecas físicas. Por simplificação, cada item do acervo é considerado como um número inteiro e, deste modo, o real problema é ordenar um array de inteiros. Para realizar essa ordenação existem 7 algoritmos candidatos que devem ser comparados a fim de descobrir qual é a eficiência deles para a resolução desse problema. Os algoritmos candidatos são variações do algoritmo QuickSort, sendo eles:

1. **Quicksort clássico** – Seleção de pivô usando o elemento central.
2. **Quicksort mediana de três** – Seleção do pivô usando a “mediana de três” elementos, em que o pivô é escolhido usando a mediana entre a chave mais à esquerda, a chave mais à direita e a chave central (como no algoritmo clássico).
3. **Quicksort primeiro elemento** – Seleção do pivô como sendo o primeiro elemento do subconjunto.
4. **Quicksort inserção 1%** - O processo de partição é interrompido quando o subvetor tiver menos de $k = 1\%$ chaves. A partição então deve ser ordenada usando uma implementação especial do algoritmo de ordenação por inserção, preparada para ordenar um subvetor. Seleção de pivô usando a “mediana de três” elementos, descrita acima.
5. **Quicksort inserção 5%** - Mesmo que o anterior, com $k = 5\%$.
6. **Quicksort inserção 10%** - Mesmo que o anterior, com $k = 10\%$.
7. **Quicksort não recursivo** – Implementação que não usa recursividade. Utiliza pilha para simular as chamadas de função recursivas e identificar os intervalos a serem ordenados a cada momento. A seleção do pivô deve ser feita assim como no Quicksort clássico.

Para cada algoritmo são feitos testes com diferentes tipos de arrays de entrada (ordenado crescente, ordenado decrescente e aleatório) e diferentes tamanhos de array com intuito de obter aproximadamente o tempo de execução, número de comparações entre chaves e número de movimentações de registro médios (mediana no caso do

tempo de execução). Posteriormente deve ser feita uma análise dos três parâmetros citados anteriormente. A linguagem utilizada na solução desse problema foi a linguagem C.

2. Implementação

A implementação desse problema se deu em duas etapas distintas, sendo elas:

- Identificar os parâmetros passados para o programa (tipo de QuickSort, tamanho e organização dos arrays, impressão dos arrays) e organizar os dados seguindo esses parâmetros;
- Montar e executar o tipo certo de QuickSort de acordo com a entrada do programa.

Sendo assim, vamos dividir essa seção em duas partes para abordar cada etapa individualmente.

2.1. Identificação dos parâmetros

Os parâmetros para o programa são passados na forma de argumentos no seguinte formato:

./nomedoprograma <variação> <tipo> <tamanho> [-p]

onde *variação*, *tipo* e *tamanho* são o tipo de QuickSort a ser executado, o tipo de ordenação dos arrays e o tamanho do arrays, respectivamente. O parâmetro *-p*, quando presente, indica que é necessário imprimir os arrays originais utilizados na ordenação.

A linguagem C fornece uma forma simples de lidar com argumentos. Por meio da passagem de dois parâmetros para a função *main*, é possível saber a quantidade de argumentos e seus valores. Esses parâmetros são o *argc* e o *argv*, sendo o primeiro um valor inteiro que guarda a quantidade de argumentos e o segundo um array de strings contendo o valor dos argumentos. Dessa forma, a declaração do *main* fica no seguinte formato:

*int main(int argc, char *argv[]) { ... }*

Após o recebimento dos argumentos/parâmetros, a *variação* e o *tamanho* dos arrays são passados para uma função responsável por gerar os arrays de teste. Os arrays de ordem aleatória utilizam uma função que gera números aleatórios (*rand()* da biblioteca *stdlib.h*) enquanto que os arrays ordenados em ordem crescente e decrescente vão de 1 a *n* e de *n* a 1, respectivamente, sendo *n* o tamanho dos arrays.

Imediatamente após gerar os arrays, é conferido se o parâmetro *-p* estava presente na entrada. Se sim, é feita uma cópia dos arrays originais de forma a garantir que a ordem original de cada um deles não seja perdida. Por fim, é feita uma chamada para a função QuickSort indicando o tipo de ordenação que foi passado. Esta função, que será abordada na próxima seção, é responsável por identificar e aplicar o tipo certo de ordenação para os arrays.

2.2. QuickSort

O princípio básico do QuickSort consiste, em poucas palavras, em organizar os elementos de acordo com um pivô previamente escolhido. O método de partição, que é o método principal dessa ordenação, é dividido em três etapas simples:

- Escolher um elemento aleatório do vetor, que será o pivô;

- Colocar os elementos menores que o pivô à esquerda dele, no vetor;
- Colocar os elementos maiores que o pivô à direita dele.

Com isso temos duas partições: uma composta apenas de elementos menores que o pivô e outra composta apenas de elementos maiores que o pivô. A partir disso, é só aplicar o método novamente nas partições até que as subpartições geradas tenham apenas um elemento, o que significa que o vetor está ordenado.

Inicialmente cada variação de QuickSort foi tratada de forma individual, sendo criado um conjunto específico de funções para cada uma delas. Porém, após desenvolver e testar cada variação, o objetivo se tornou reduzir a quantidade de código repetido por meio da junção de algumas funções. Com esse processo, as variações de QuickSort foram divididas em três grupos principais: QuickSort Padrão, QuickSort com Inserção e QuickSort Não Recursivo. O primeiro grupo é composto pelas variações QuickSort Clássico, Mediana de Três e Primeiro Elemento, cuja variação se dá apenas pela escolha do pivô. O segundo grupo é composto por QuickSort Inserção 1%, 5% e 10%, cuja variação se dá apenas pela quantidade de elementos necessários para a chamada do método de ordenação por inserção. O último grupo é composto apenas pela variação QuickSort Não Recursivo, pois essa é a única variação que não utiliza recursão.

Para cada um dos dois primeiros conjuntos foi desenvolvida uma função Sort que coordena a escolha de pivô e as chamadas recursivas (além da chamada para o método de inserção, no caso do segundo grupo). O método de partição é único no código, ou seja, todas as variações de QuickSort utilizam ele, e a diferenciação entre as formas de escolha de pivô é indicada por um parâmetro passado para a função. Cada troca de elementos do vetor é considerada como 3 movimentações (cópia do primeiro elemento para um auxiliar, cópia do segundo para o primeiro, cópia do auxiliar para o segundo). Uma estrutura de dados Pilha é utilizada na função de QuickSort Não Recursivo para armazenar a ordem das chamadas para a função de partição, guardando os índices do array para as próximas chamadas.

As funções Sort dos grupos QuickSort Padrão e QuickSort com Inserção, e a função de QuickSort não recursivo são coordenadas por uma função geral chamada QuickSort (mencionada no final da seção anterior). Essa função, que recebe um array e sua variação, determina qual função de ordenação será chamada e quais parâmetros serão passados para a função. Por exemplo: para a variação QuickSort Clássico, o QuickSort chama a função Sort do grupo QuickSort Padrão passando como parâmetro de escolha de pivô o número 0, que indica que o pivô escolhido será o elemento central. Para um QuickSort Mediana de Três, a chamada seria a mesma, com exceção do parâmetro de escolha de pivô que agora seria o número 2, indicando uma escolha de pivô por mediana de três. Para o grupo de QuickSort com Inserção, é passado o número de elementos a partir do qual o método de inserção deve ser chamado.

As funções de inserção e partição possuem um contador de trocas de elementos do vetor e de comparações entre elementos do vetor que são retornadas a cada chamada. Isso possibilita a contagem desses parâmetros, que será necessária na saída do programa e na seção de Análise Experimental (seção 4). A saída do programa segue o seguinte formato:

<variação> <tipo> <tamanho> <n_comp> <n_mov> <tempo>

onde *variação*, *tipo* e *tamanho* são os parâmetros recebidos na entrada, *n_comp* e *n_mov* se referem ao número médio de comparações de chaves e de movimentações de registros efetuadas e *tempo* ao tempo mediano de execução, em microssegundos.

3. Instruções de Compilação e Execução

O compilador utilizado durante toda a implementação foi o GNU C Compiler (GCC) versão 7.3.0 e o sistema operacional foi um Linux Ubuntu 18.04 64-bit.

Para compilar o programa foi desenvolvido um arquivo makefile que é responsável por gerenciar os arquivos a serem compilados e a ordem de compilação. Para executar esse aqui, basta digitar o comando *make* no terminal.

A execução do programa (como mostrado na seção 2.1) é feita pelo comando:

./nomedoprograma <variação> <tipo> <tamanho> [-p]

Obs.: Para alguns casos de teste, o número de chamadas recursivas ultrapassa o limite da pilha de execução devido ao grande número de elementos somado ao pior caso do QuickSort. Para evitar isso, o comando *ulimit -s <tamanho em bytes>*, que aumenta o tamanho da pilha, deve ser usado antes da execução do programa.

4. Análise Experimental

As análises experimentais seguiram um modelo determinado pelas instruções do trabalho. O modelo consiste em fazer testes de arrays variando entre 50 mil e 500 mil elementos, em intervalos de 50 mil elementos, para cada variação de QuickSort e de ordenação do array. Sendo assim temos: *10 tamanhos x 7 quicksorts x 3 ordenações do array* = 210 testes. Porém, para cada teste foram usados 20 arrays de modo a obter a mediana do tempo de execução e, as médias de comparações e movimentações de elementos no array. Sendo assim, *4200* arrays foram durante os testes.

Especificações técnicas do computador de testes:

- **Processador:** Intel i7-3612QM CPU @ 2.10GHz
- **Memória RAM:** 4 GB
- **Sistema Operacional:** Linux Ubuntu (versão 18.04)

Os dados utilizados para análise e todos os gráficos estão na pasta *dados* enviada junto com esta documentação. Os arquivos dos gráficos estão nomeados de forma a facilitar a identificação de cada um deles.

4.1. Análise do Tempo

A subpasta *Time* contém todos os gráficos que comparam a performance (em tempo) dos métodos de ordenação para determinada ordenação inicial dos elementos e quantidade de elementos do array. Além dela, a subpasta *Full Plots* contém gráficos de comparação de tempo mais gerais entre os métodos. Entre eles, dois gráficos são bem interessantes de serem analisados. O primeiro (Figura 1) compara a mediana do tempo para arrays de 50 mil a 500 mil elementos ordenados previamente em ordem crescente. O gráfico à esquerda contém os dados em escala original, enquanto que o segundo retira uma das linhas, o que dá um zoom no gráfico e permite a visualização das outras linhas.

O que pode-se notar por meio desses gráficos é que o tempo de ordenação do QuickSort com pivô no primeiro elemento é consideravelmente maior do que o dos outros

métodos. Isso ocorre porque, como os dados já estão ordenados, na hora da partição o pivô sempre será o menor elemento do subvetor a ser ordenado, gerando assim a pior partição possível para os elementos. Assim, o tempo de ordenação por esse método entra no pior caso do QuickSort, que é $O(n^2)$. Já para os métodos que usam inserção, o que ocorre é o contrário. Como o array já está ordenado, a complexidade de ordenação entra no melhor caso do método de inserção, que é $O(n)$, o que explica eles serem mais rápidos em tempo execução para esses arrays. O gráfico *Time Plot – OrdD* (da pasta Full Plots) entra no mesmo caso, então a análise é semelhante com exceção de que o pivô, para o método QPE (QuickSort Primeiro Elemento), sempre será o maior elemento do subvetor a ser ordenado.

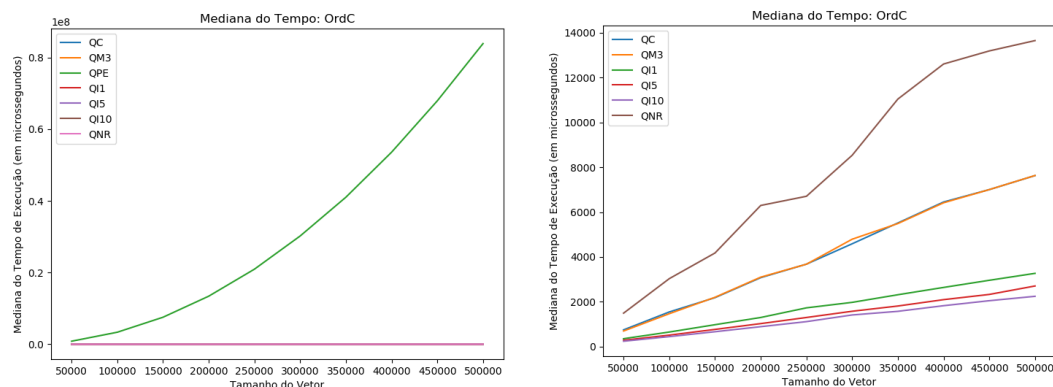


Figura 1. Mediana do Tempo para Vetores Previamente Ordenados

O segundo gráfico (Figura 2) compara a mediana do tempo para arrays de 50 mil a 500 mil elementos ordenados de forma aleatória. Nesse caso, os QuickSorts que utilizam inserção ao final da ordenação são os mais lentos. Isso ocorre porque, no pior caso e no caso médio, o algoritmo de inserção é $O(n^2)$, o que impacta bastante no tempo para vetores grandes, explicando o porquê o QuickSort com 10% inserção é mais lento que o de 5% e o de 1%.

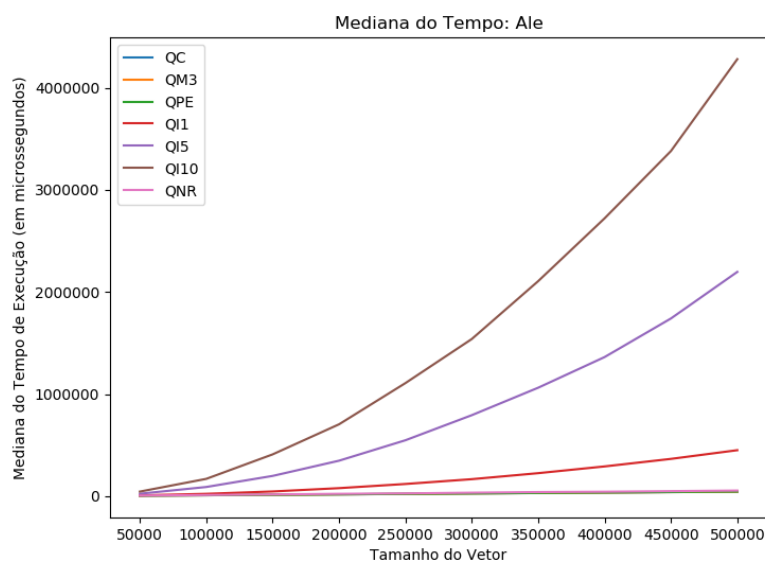


Figura 2. Mediana do Tempo para Vetores Aleatórios

4.2. Análise de Comparações e Trocas

As subpastas *Comparisons* e *Swap* contêm todos os gráficos que comparam o número médio de comparações e de trocas, respectivamente. Assim como no caso anterior, a subpasta *Full Plots* contém gráficos interessantes, que serão analisados.

Os gráficos de comparações seguem um padrão parecido com os gráficos de tempo. Para vetores previamente ordenados em ordem crescente (Figura 3), ou decrescente (Comp Plot – OrdD da pasta Full Plots), o QuickSort com pivô no primeiro elemento faz o maior número de comparações, pelo mesmo motivo citado na subseção anterior. As partições geradas são sempre as piores possíveis, levando ao pior caso no número de comparações também. O gráfico da esquerda novamente mostra os dados em escala real, enquanto que o gráfico da direita é um zoom que possibilita ver as outras curvas.

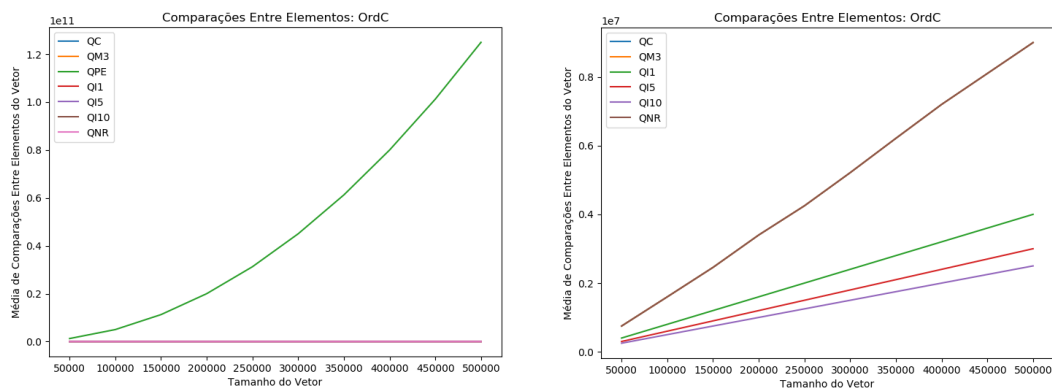


Figura 3. Média de Comparações para Vetores Previamente Ordenados

Para vetores aleatórios, os QuickSorts que utilizam inserção fazem o maior número de comparações (Figura 4), comportamento semelhante ao que foi visto na comparação de tempo, pois o pior caso do algoritmo de inserção também é $O(n^2)$ para comparações.

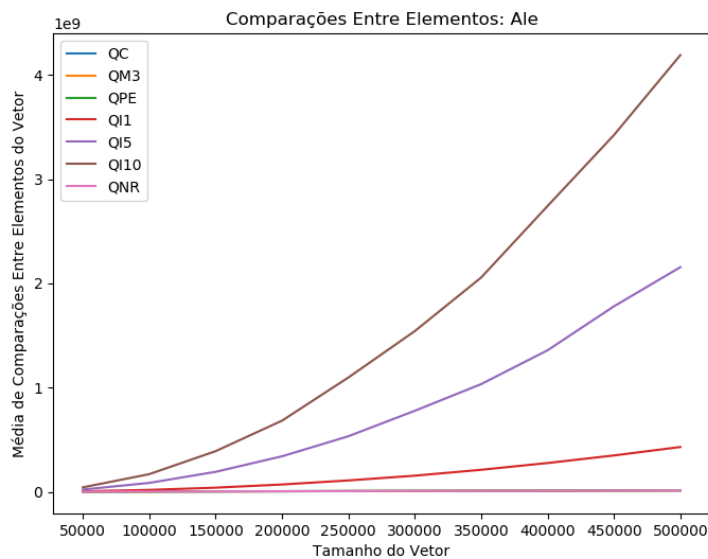


Figura 4. Média de Comparações para Vetores Aleatórios

Os gráficos de trocas tem um comportamento bem diferente. Para vetores com elementos em ordem aleatória (Figura 5), os QuickSorts com inserção ainda possuem uma performance pior que os outros métodos. Porém para vetores ordenados decrescentes (Figura 6), os QuickSorts Clássico, Mediana de Três e Não-Recursoivo fazem o maior número de trocas. Isso ocorre porque o pivô para esses três casos é o elemento central, com isso todos os elementos do array são trocados de posição pelo menos uma vez, devido à ordenação inicial.

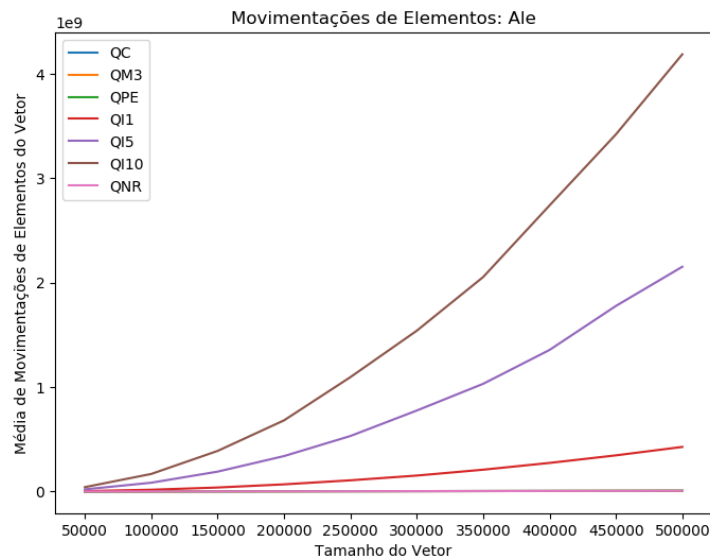


Figura 5. Movimentações de Elementos para Vetores Aleatórios

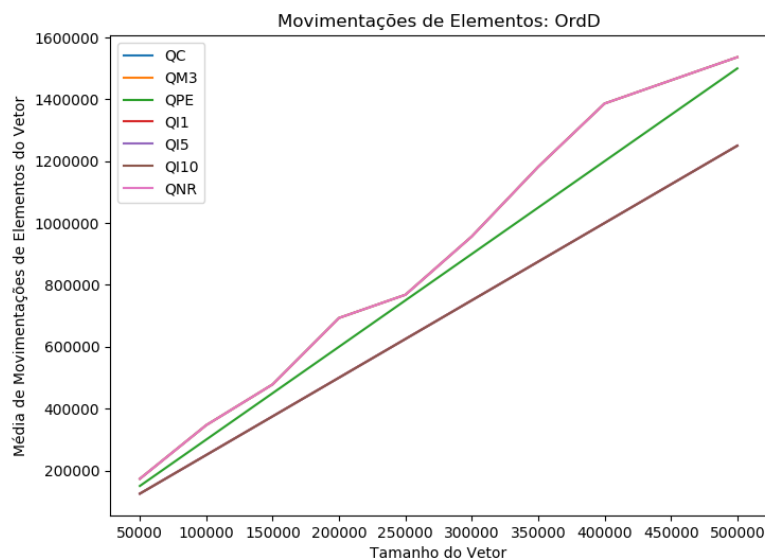


Figura 6. Movimentações de Elementos para Vetores Ordenados Decrescentes

Para vetores ordenados crescentes (Figura 7), o QuickSort com pivô no primeiro elemento continua a ser o destaque, desta vez possuindo o maior número de trocas. Isso decorre da implementação do QuickSort que, nesse caso, sempre troca o pivô de posição com ele mesmo. Sendo assim, temos um número de trocas linear igual ao $3(n-1)$ (com n sendo o número de elementos do vetor), porque são feitas $n-1$ chamadas do método de partição, onde cada uma troca o pivô com ele mesmo (3 movimentações).

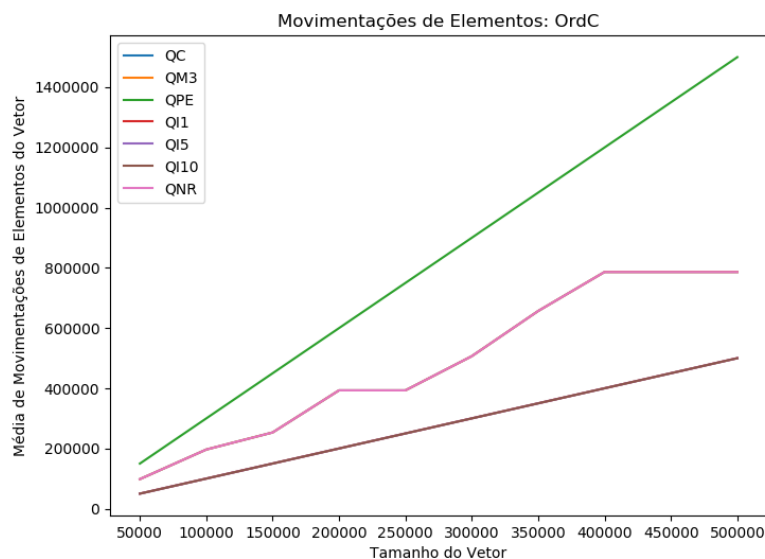


Figura 7. Movimentações de Elementos para Vetores Previamente Ordenados

5. Conclusão

Por meio dos experimentos, e da análise deles, foi possível observar a diferença de performance entre os diferentes tipos de implementação QuickSort e como cada uma se comporta em diferentes situações. Pode-se perceber que para certos tipos de ordenação dos dados, existem variações do QuickSort que obtém melhor desempenho, como por exemplo o QuickSort com Inserção para vetores previamente ordenados ou quase ordenados. Também existem casos em que o desempenho mostra qual variação do QuickSort não deve ser usada para certos tipos de dados, como é o caso do QuickSort com pivô no primeiro elemento para vetores ordenados crescentes ou decrescentes.

Como os dados foram obtidos através de experimentos, possíveis erros em algumas medidas são esperados. Mas, de maneira geral, o comportamento dos parâmetros medidos se mostrou condizente com a teoria.

A análise de algoritmos e a comparação de performance entre eles possibilita uma melhor escolha de qual algoritmo implementar para certo tipo de problema, o que pode ser fundamental em alguns casos.

6. Referências

Ziviani, N. *Projeto de Algoritmos: com implementações em PASCAL e C. 2 ed.* São Paulo: Thomson, 2004

Chaimowicz, L. e Prates, R. *Ordenação: Quicksort*, disponível na turma *ESTRUTURAS DE DADOS – METATURMA* do Moodle UFMG 2019/1. Acesso em maio/2019.

Ahuja, K. e Patade, A. *Command line arguments in C/C++*, <https://www.geeksforgeeks.org/command-line-arguments-in-c-cpp/>, acesso em maio/2019.