

Primeira versão do projeto da disciplina

Comparação entre os
algoritmos de ordenação
elementar

Wilton Lima Silva Saturnino, João Victor Barbosa

1 Introdução

Este trabalho foi desenvolvido para aplicação dos conceitos teóricos de programação estudados na disciplina Estrutura de Dados em paralelo com Laboratório de Estrutura de Dados. Espera-se com os resultados aplicar os algoritmos de ordenação e comparalos conforme os diferentes casos de complexidade estudados na disciplina.

Com base no que foi dito, os dados devem ser baixados na plataforma indicada pelo professor e a partir dos filtros disponibilizados deve-se baixarr os arquivos de dados para o sistema de compartilhamento de bicicletas na cidade de Los Angeles, California e região metropolitana da empresa Los Angeles Bike Share.

O tratamento dos dados está compreendiada entre 2016 e 2021 podem ser realizados a partir da aplicação dos algoritmos para ordenar os dados desenvolvidos em linfuagem JAVA e aplicando sobre arquivos CSV e para melhor visualização os resultados de tempo e desempenho podem ser plotados para verificar qual o método aplicado é mais eficiente para os diferentes casos de complexidade (melhor caso, pior caso, caso médio).

As tranformações sugeridas são realizadas a partir da identificação do id da estação substituindo pelo nome da estação contida no arquivo CSV, filtrar apenas as viagens que estão nas estações de maneira que seja considerado como prioridade as viagens que possuem duração maior que a média geral.

As ordenação são feias com base no resultado obtido pelas transformações e tratados em três casos. Ordenar o arquivo completo pelo nomes das estações em ordem alfabética, ordenar pelo campo de duração da viagem (campo duration) do menor para o maior, ordenar pela data de início da viagem mais recente para o mais antigo.

Para cada caso, melhor, médio e pior foram aplicados os seguintes algoritmos:

Insertion Sort; Selection Sort; Merge Sort; Quick Sort; Quick Sort com Mediana de 3; Counting Sort; e Heap Sort.

2 Descrição geral sobre o método utilizado

Neste trabalho foi necessário a utilização de um computador pessoal Windows e a teoria ministrada na disciplina a partir de reuniões em aulas presenciais. Para tratar os dados nos foi realizado testes de complexidade, os dados foram modificados sendo que no melhor caso modifica-se a ordem de forma crescente, no pior caso os dados são ajustados de forma invertida, e no caso médio eles são ajustados de forma randômica.

Executando o programa os dados contidos na pasta de dados csv são utilizados com upload no programa a partir do caminho aonde se encontra a pasta no computador e em seguida esses dados são filtrados pelo programa aplicando as ordenações e realizado os testes de complexidade para o melhor, médio e pior caso além da verificação do tempo esforço da máquina nos diferentes algoritmos de ordenação. Os dados gerados como resultados são separados e disponibilizado em uma pasta de dados filtrados em csv ao qual pode-se utilizar para entender o comparativo de tempo de trabalho dos diferentes métodos sendo possível utilizar softwares como o visualvm para visualizá-los.

A depender da CPU, Memória RAM, tamanho de arquivo e tempo de execução, pode-se se fazer diferentes interpretações sobre o custo benefício da utilização dos métodos e transformações, por isso foi feita uma análise que se refere aos recursos de máquina disponível para o processamento de cada algoritmo em cada caso de complexidade.

2.1 Descrição geral do ambiente de testes

(processamento, memória, sistema operacional)

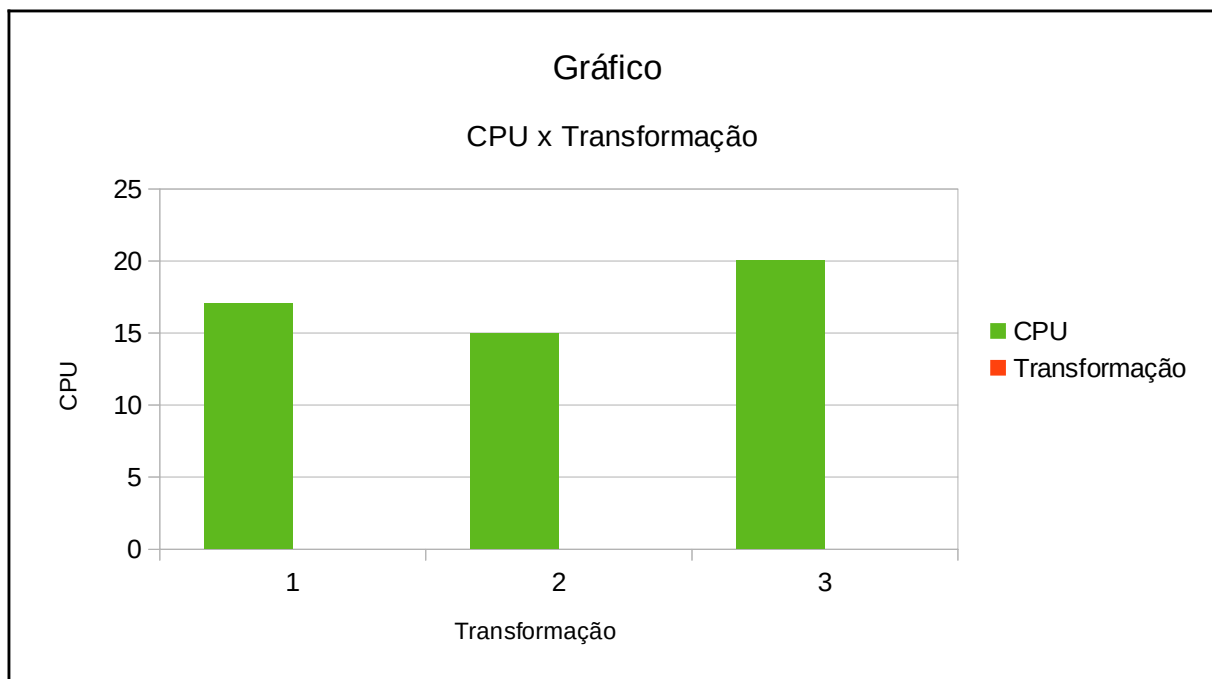
- Computador da empresa Dell;

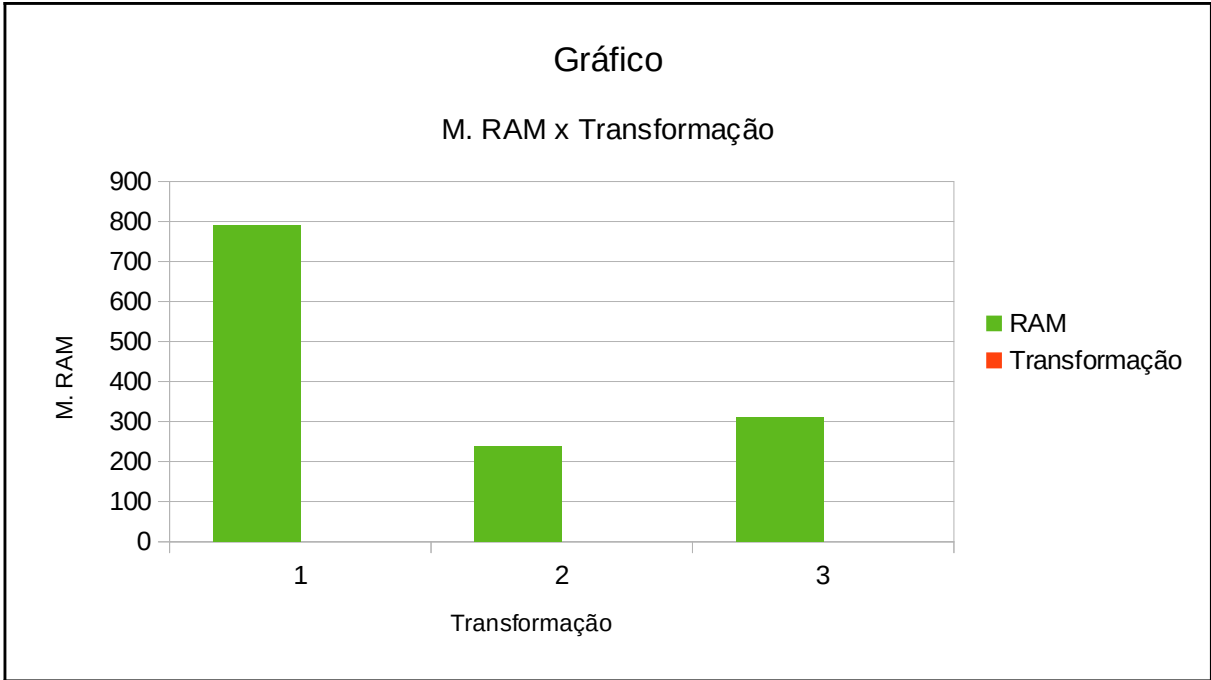
-
- Memória 16gb Ddr4 P/ Notebook Dell Inspiron 14 5000 ;
 - Processador Intel Core i5 10210U;
 - Eclipse software;
 - Sistema operacional Microsoft Windows;
 - Microsoft Excel;

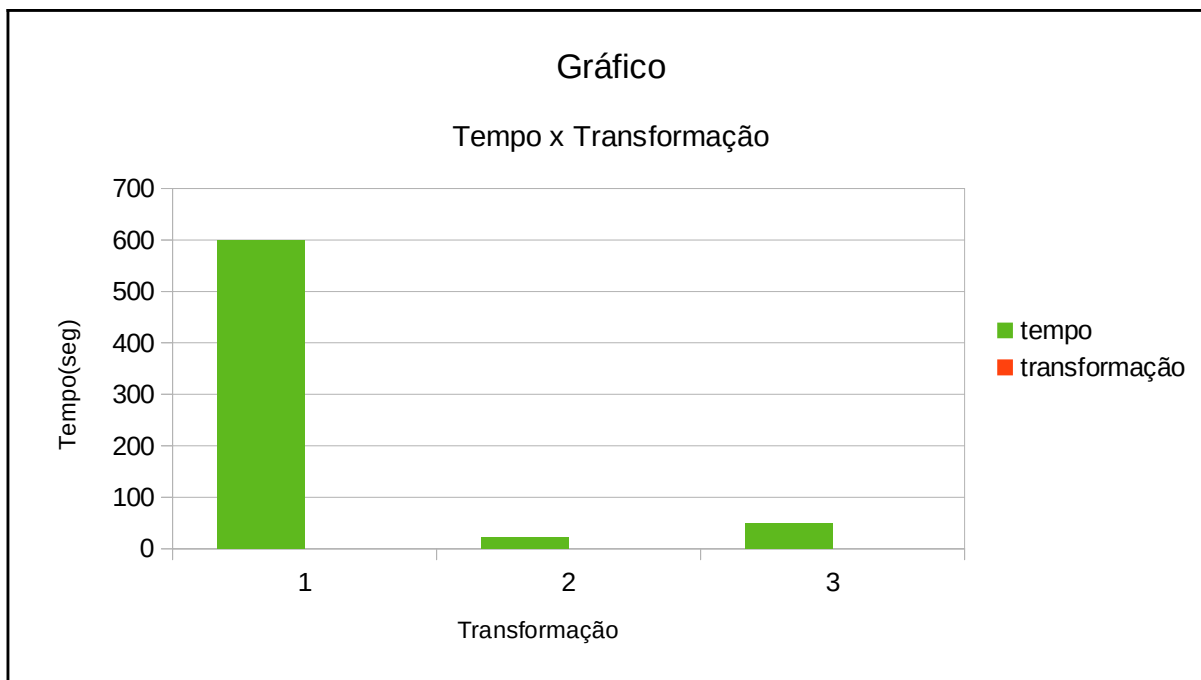
3 Resultados e Análise

Resultados dos testes usando tabelas e gráficos

Os resultados obtidos podem ser observados por meio da visualização do comportamento da máquina e desempenho como no gráfico abaixo:







No gráfico 01 é apresentado a primeira transformação substituindo os dados dos arquivos `id_star_station` e `id_end_station` como solicitado pelos nomes gerando o arquivo `station_name` do “`station.csv`”. Com esse procedimento o arquivo dá um salto de tamanho original de 173.035 kb para **235,567KB**, um acréscimo de aproximadamente 62,532KB.

O próximo passo será utilizar um arquivo chamado “`station_pasadena`” filtrar apenas as viagens que estão nas estações de Pasadena. Da mesma forma o arquivo é processado de forma mais ágil do processo por se tratar de um filtro de forma que a terceira transformação se porta de forma intermediária no tempo pela necessidade de percorrer todo o arquivo para a realização da filtragem.

Partindo para os métodos de ordenações utilizando o algoritmo Selection Sort, é perceptível sua inviabilidade de maneira que pelo modelo de comparação utilizado baseado em se passar sempre o menor valor do vetor para a primeira posição, ajustando o menor valor para a segunda posição e assim sucessivamente, o processo se mostrou trabalhoso e pela quantidade de dados a serem processados seriam necessários processar **XXX (milhões)** vezes o código gastando para isso **xx seg**, o que gastariam **xx horas** para serem ordenados seus casos de complexidade.

Comparado ao Selection Sort o insertion sort mostrou estrategicamente mais inteligente a ser utilizado neste projeto porque é um dos algoritmos mais eficientes em modelo de ordem quadrática, porém ele tem um comportamento melhor perante pequenas entradas, o que acabou sendo inviável nesse problema, necessitando de mais de 700 ordenações por min, extrapolando 40 hr para resolução de todas as ordenações.

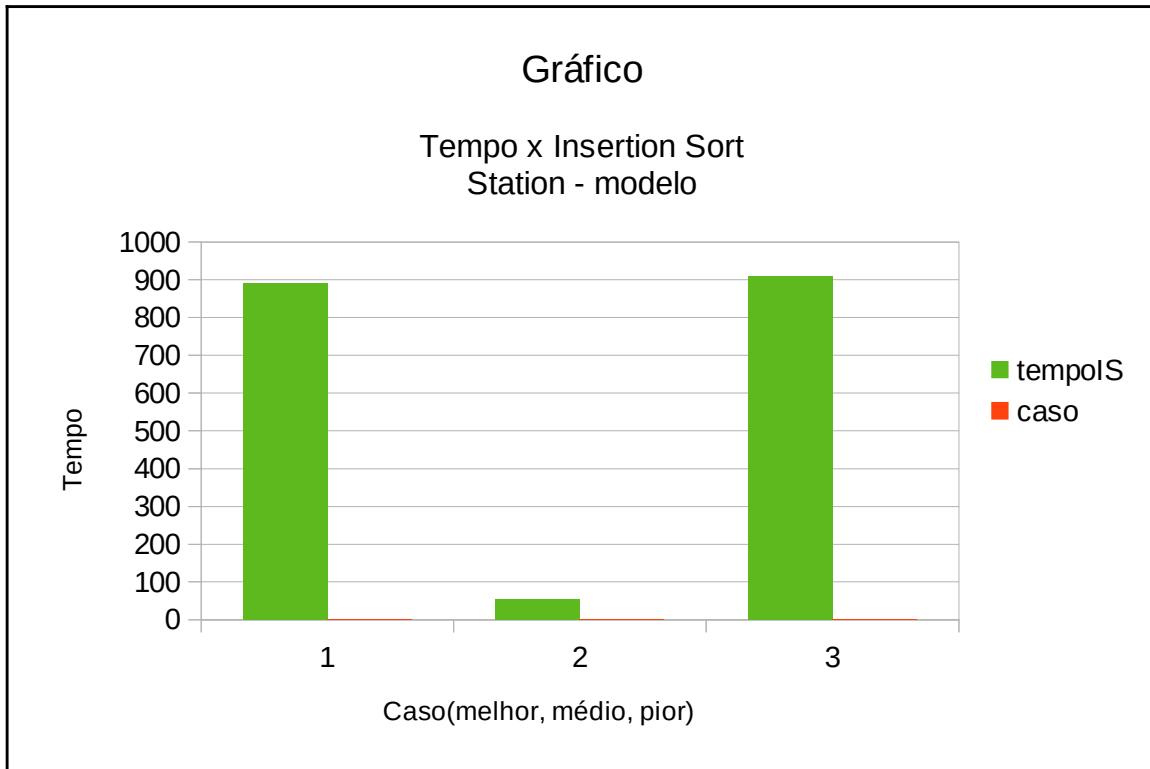
Sendo assim seria necessário entender o tradeoff porque apesar de necessitar de menos código para seu funcionamento é necessário percorrer cada linha e realizando a ordenação por inserção começando pelo índice um e cada nova posição deve ser ordenada procurando a posição correta no subarray ordenado à esquerda daquela posição.

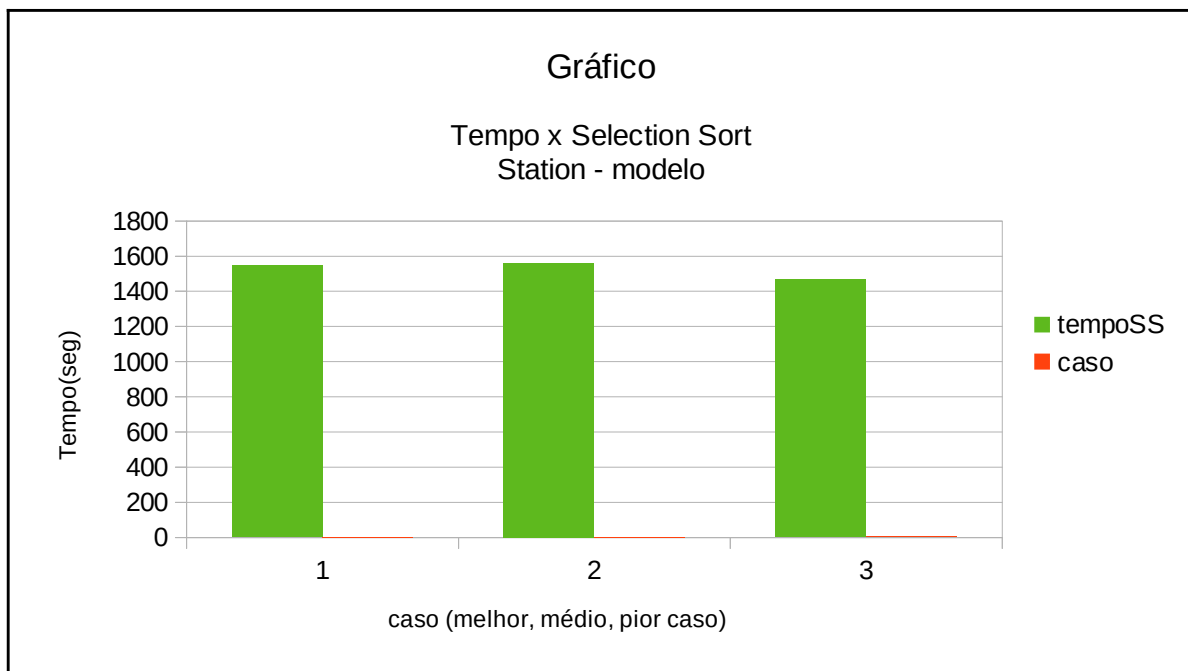
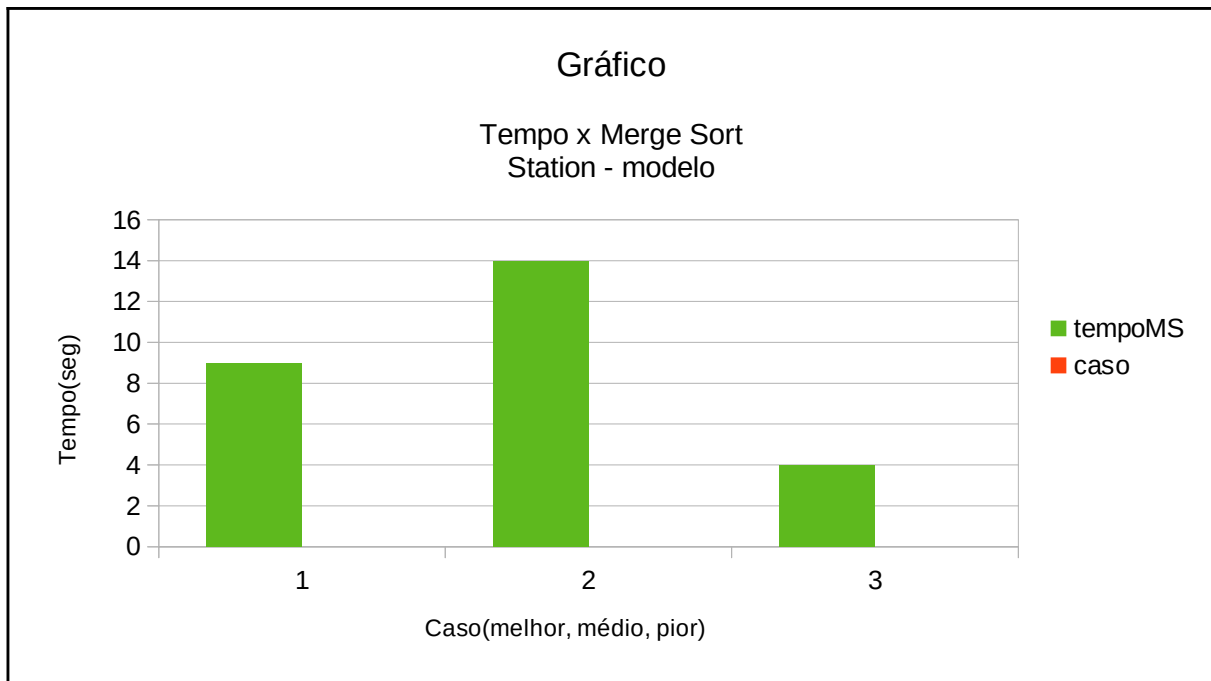
Alguns algoritmos testados não foram efetivos ou acessíveis a esse problema de maneira que foi inviável adaptar a ordenação ao tipo string. O Quick Sort, Count Sort e Mediana de 3

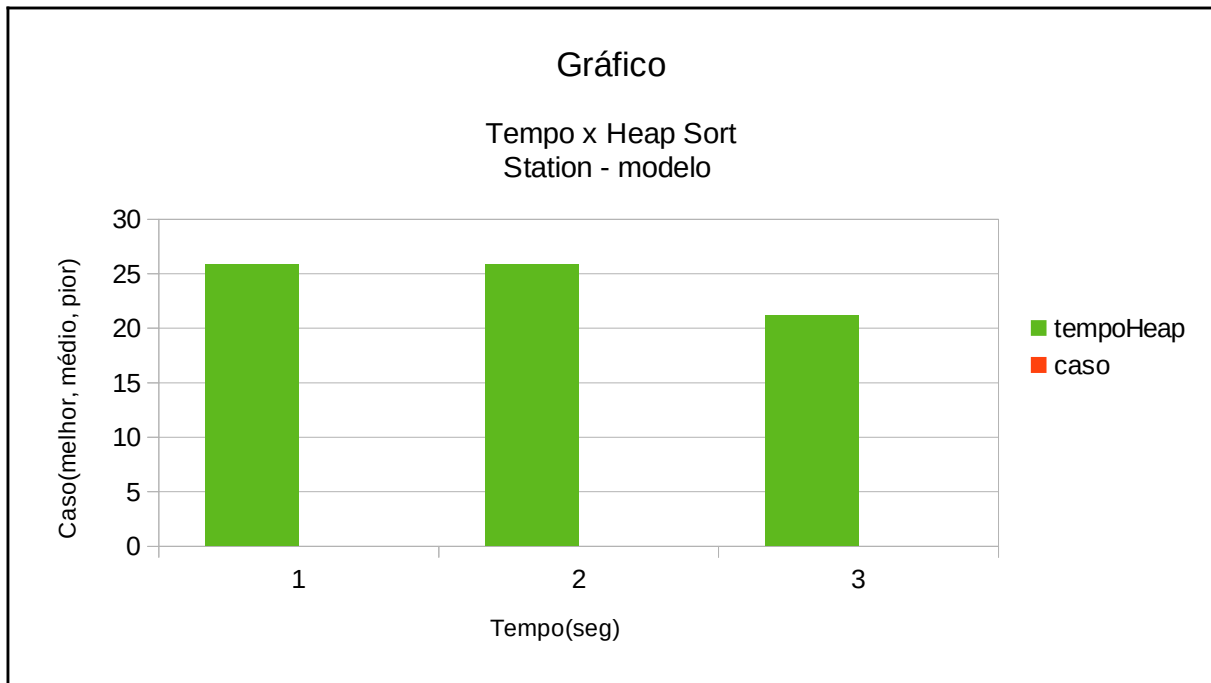
Na figura abaixo o gráfico das ordenações abordadas anteriormente:

Como resultado visivelmente o algoritmo selection sort se mostrou mais inviável pois leva um tempo elevado em todos os tipos de complexidade de

ordenação, por ser composto de dois laços um interno e outro externo, a checagem pra esse tipo de problema exige comparação a cada iteração um elemento com todos os outros exigindo mais esforço da máquina.

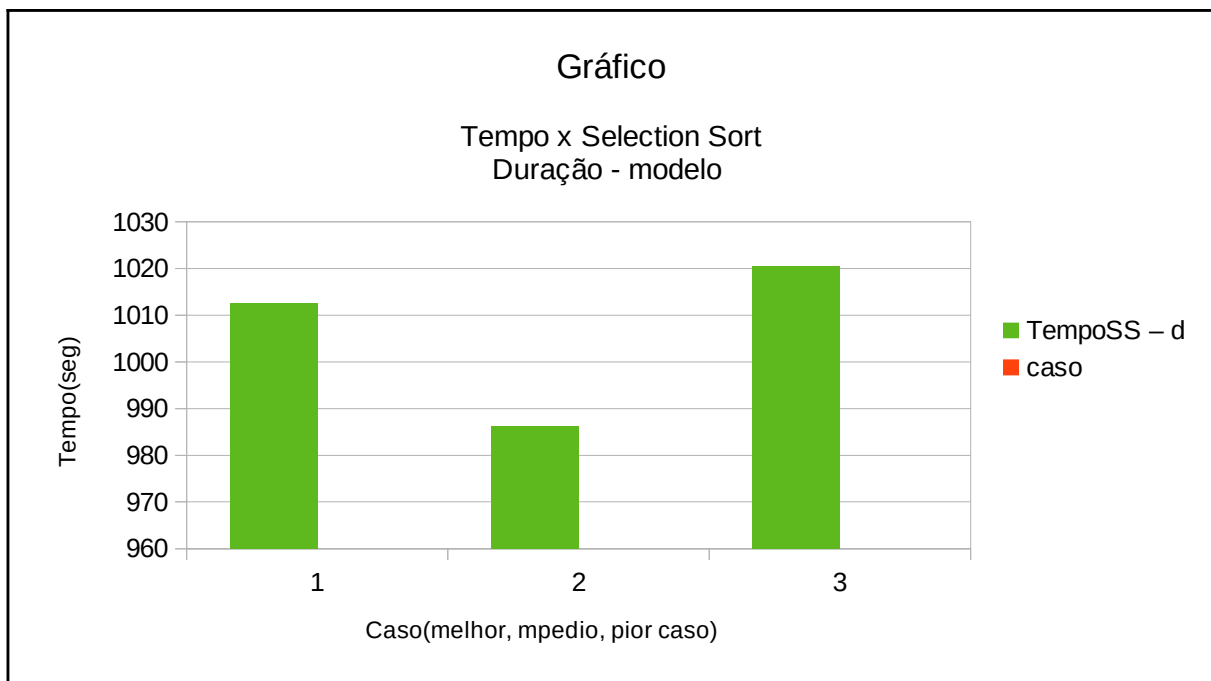
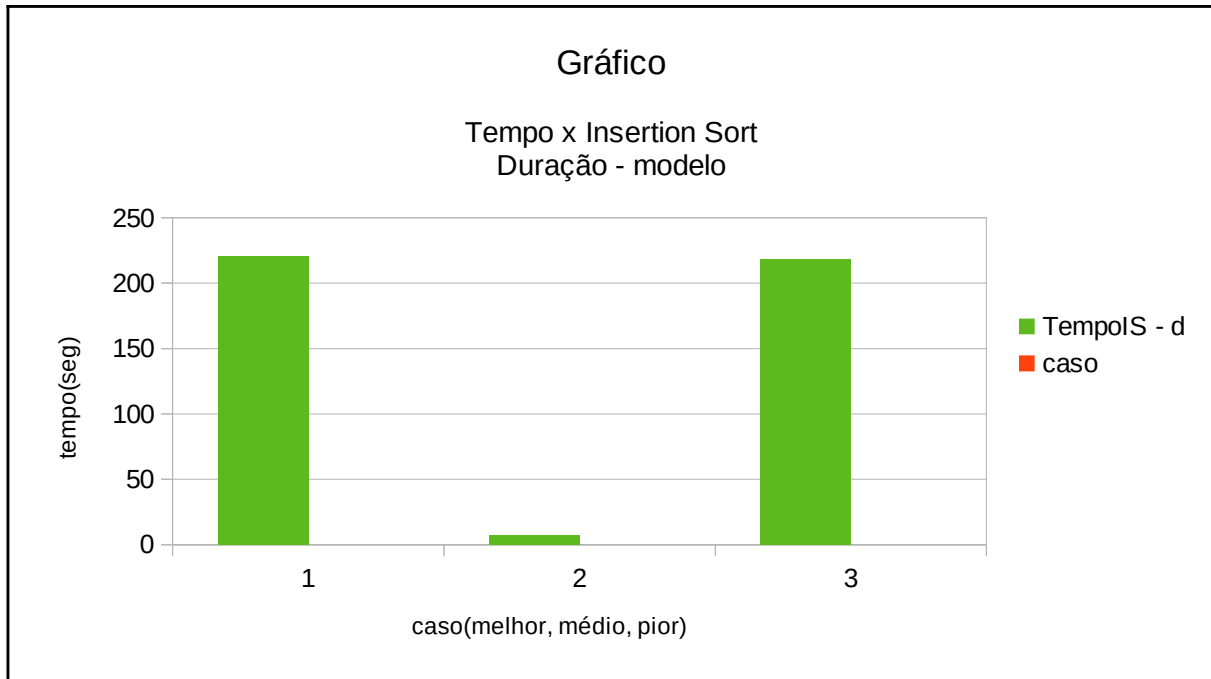


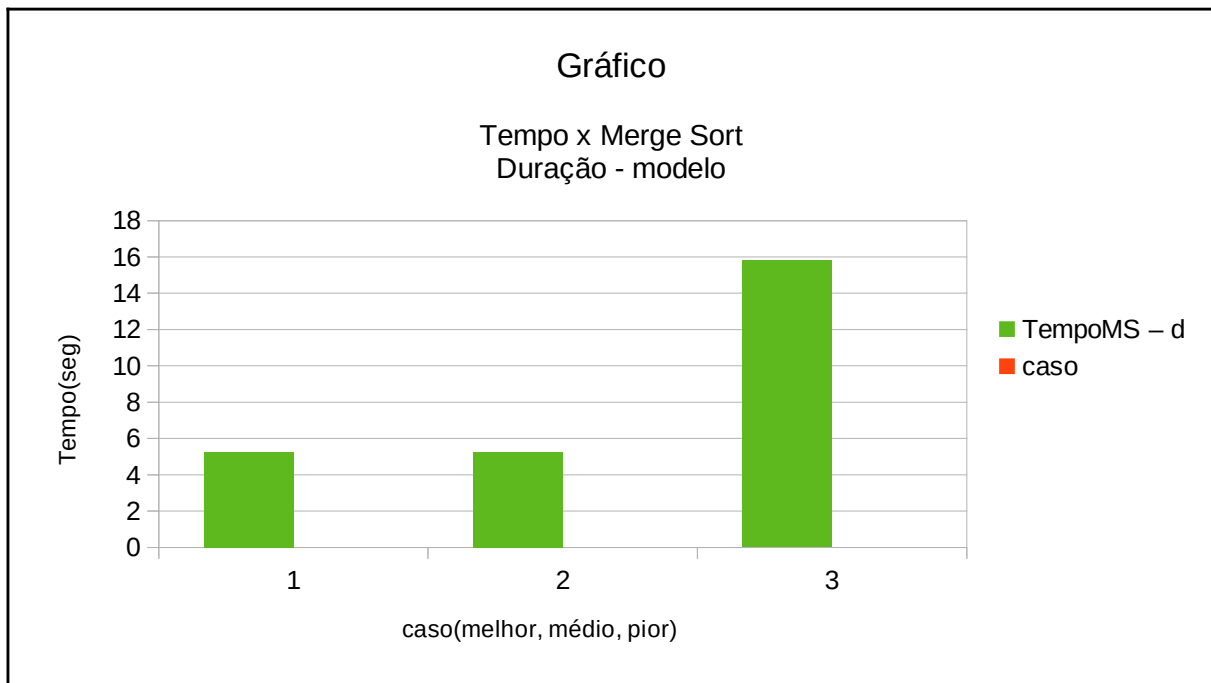
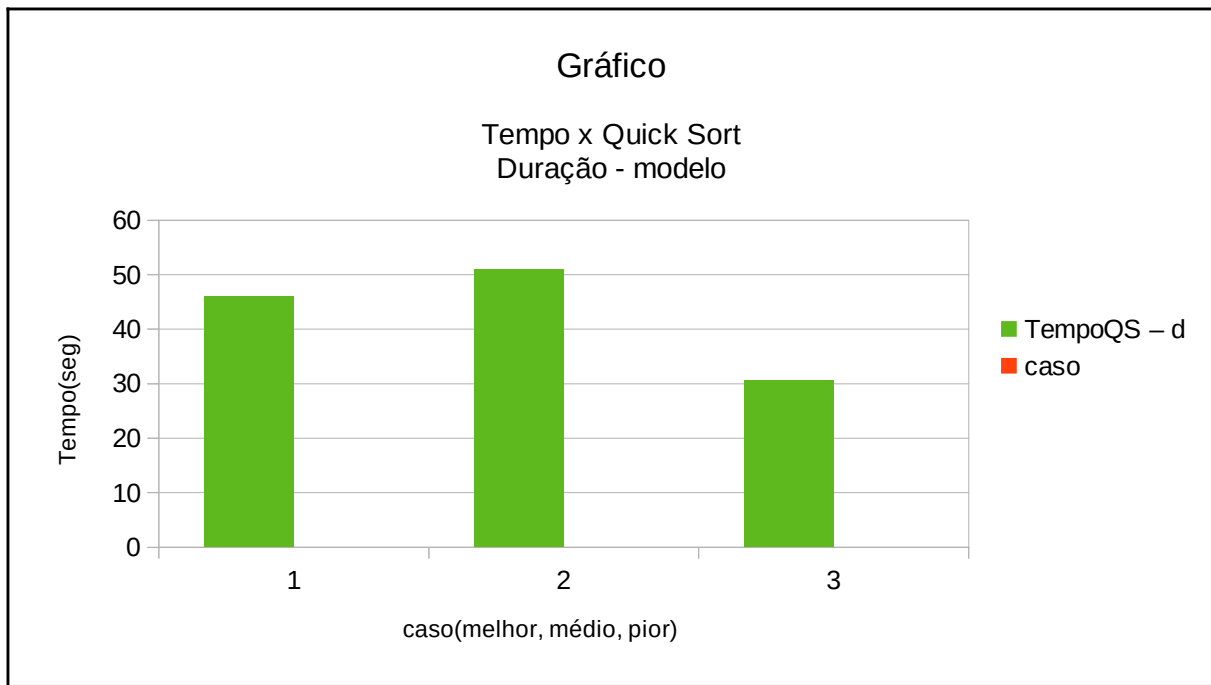


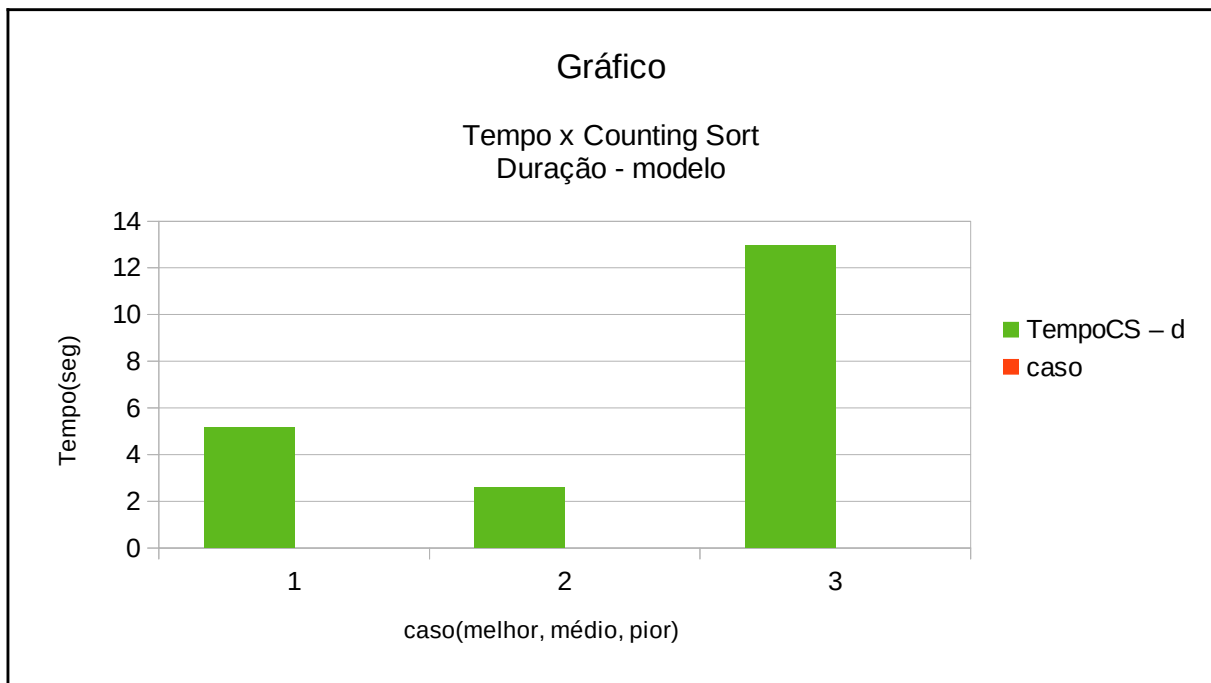
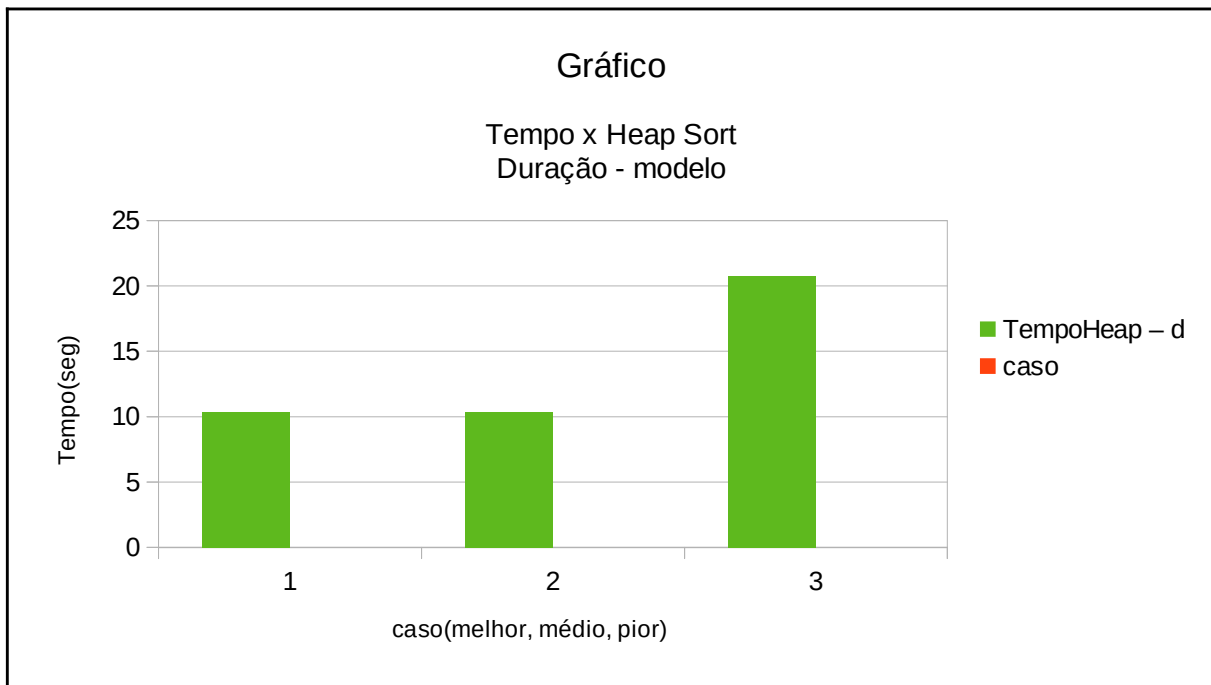


A Ordenação por Station (Nome), o algoritmo que executou mais rapidamente foi o Merge Sort, executando no menor tempo de ordenação em aproximadamente 3 seg.

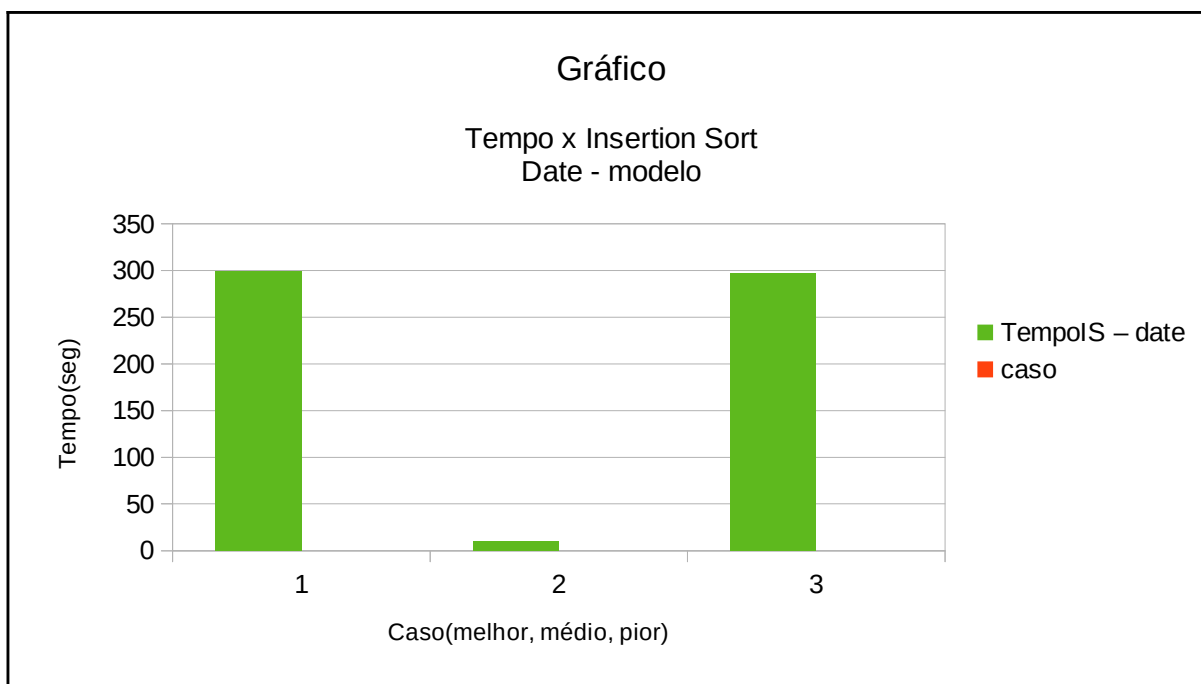
Ordenações por Duração:

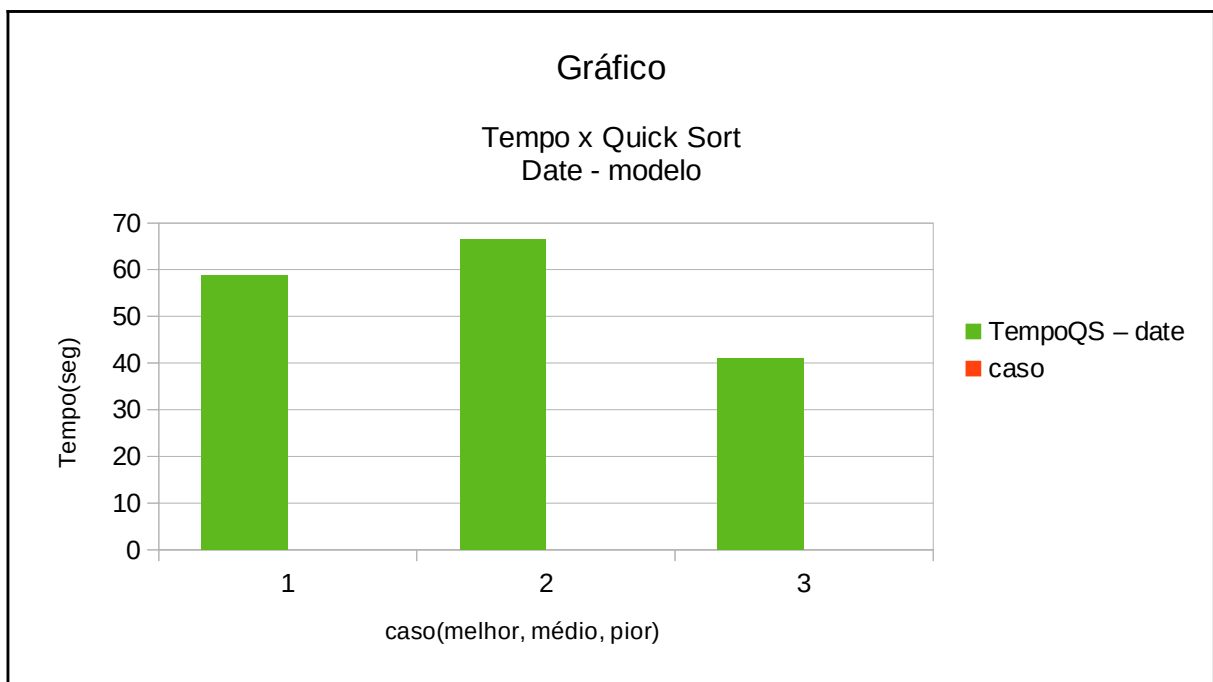
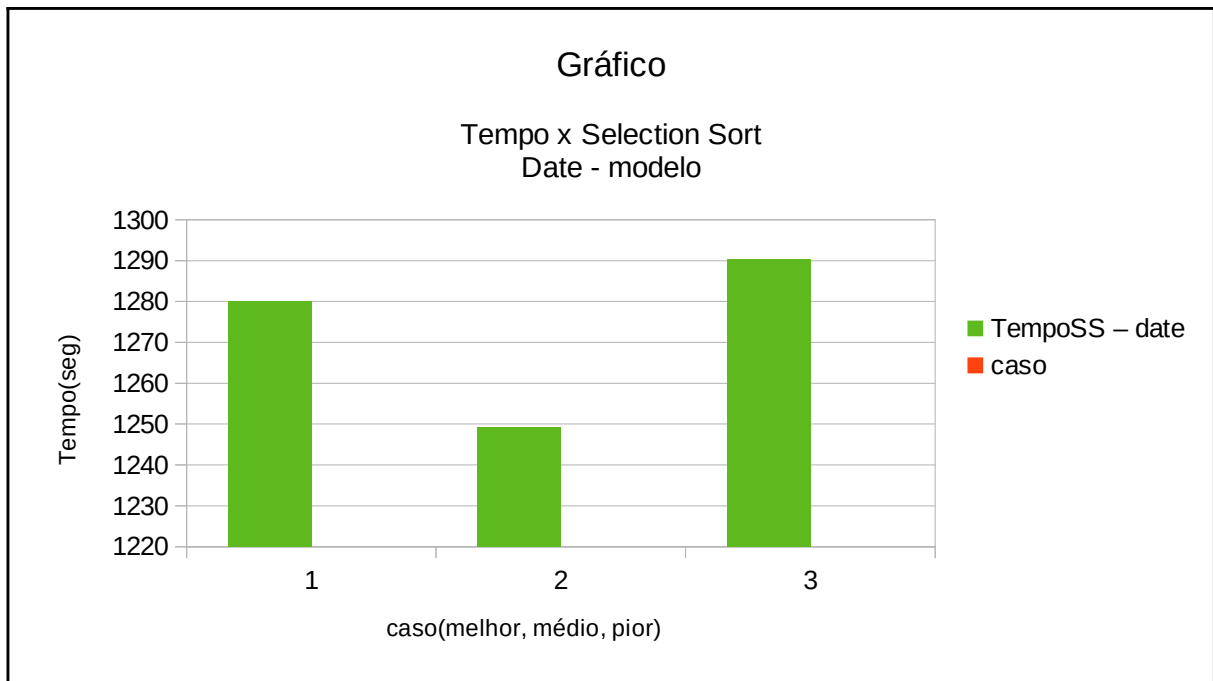


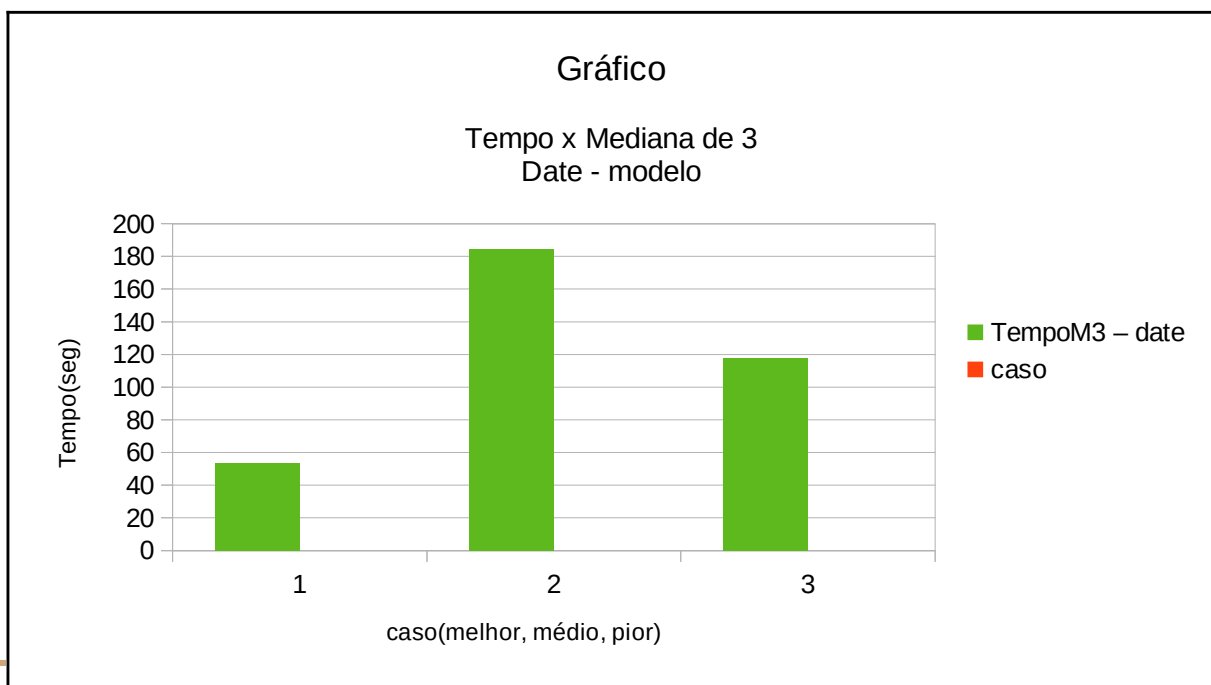
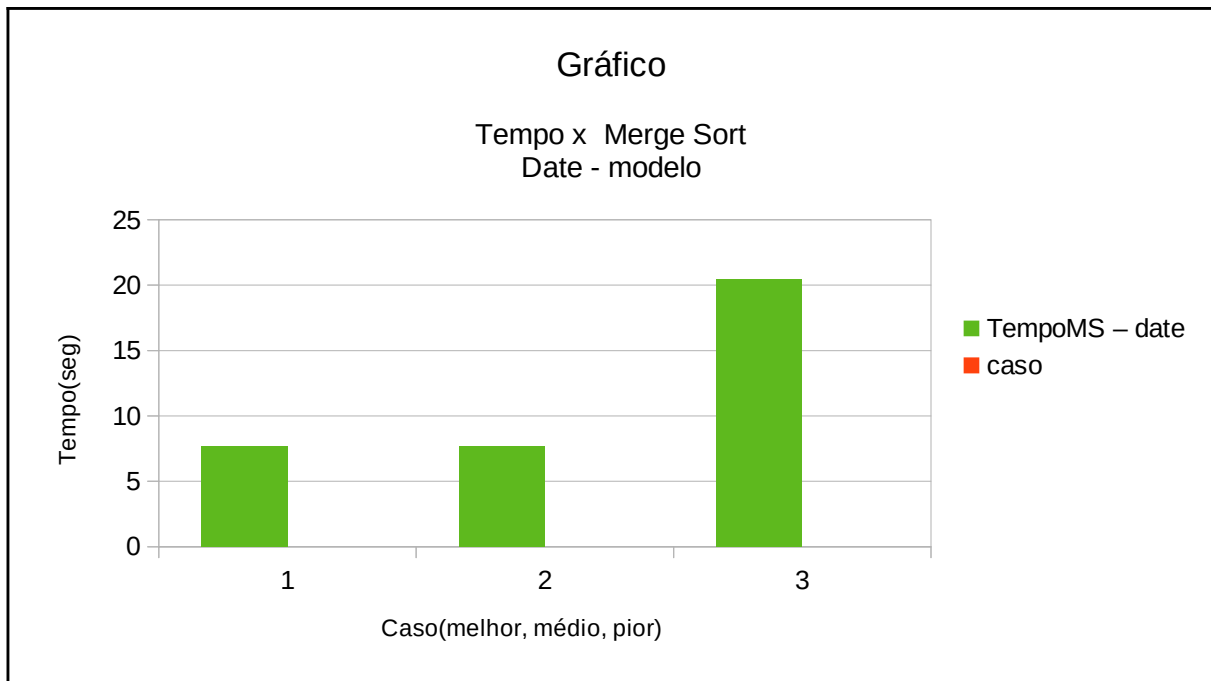




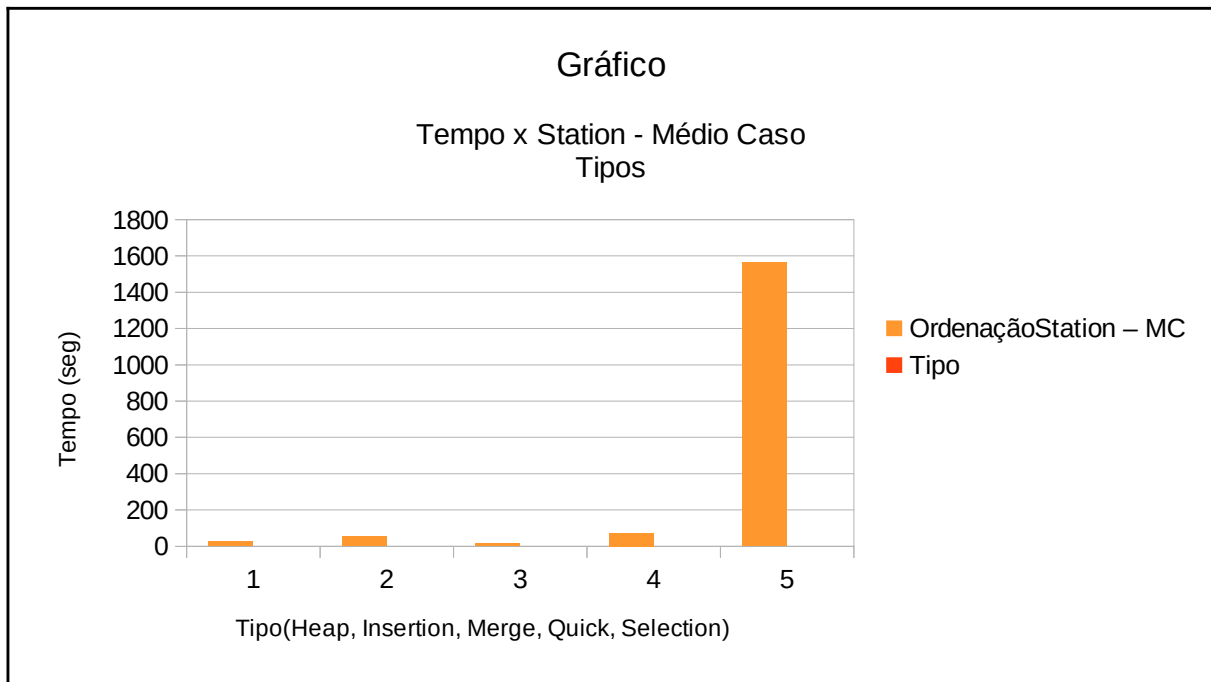
Ordenação Start_Date por algoritmos:

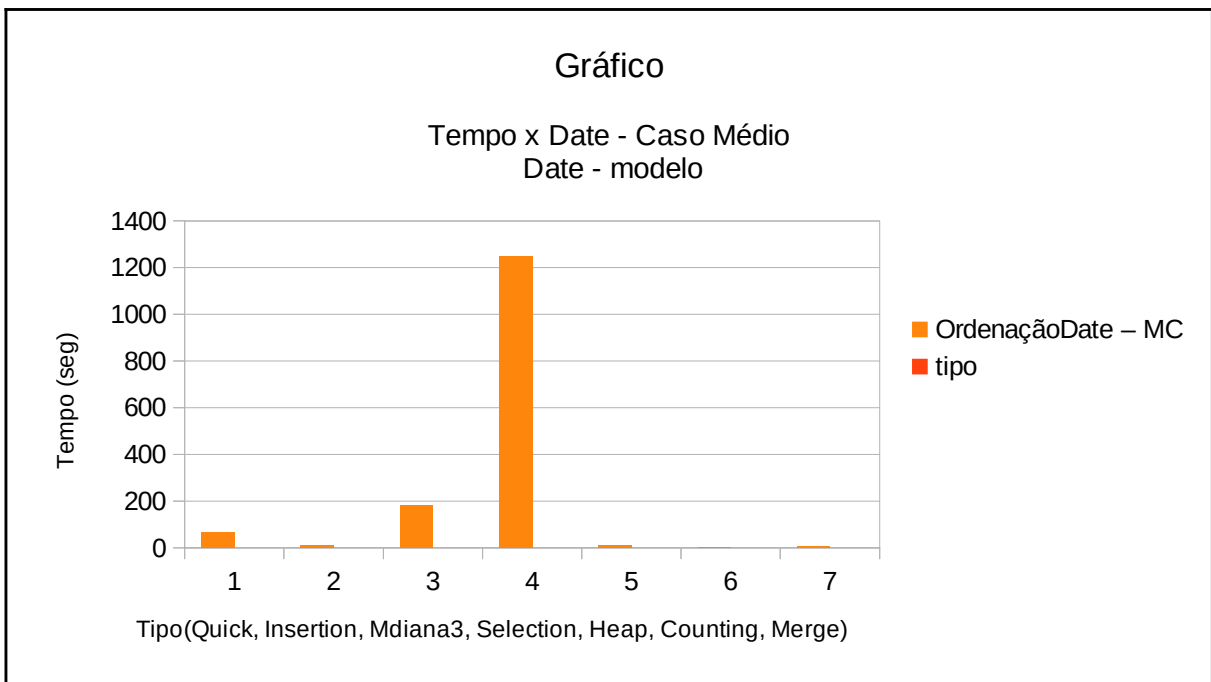
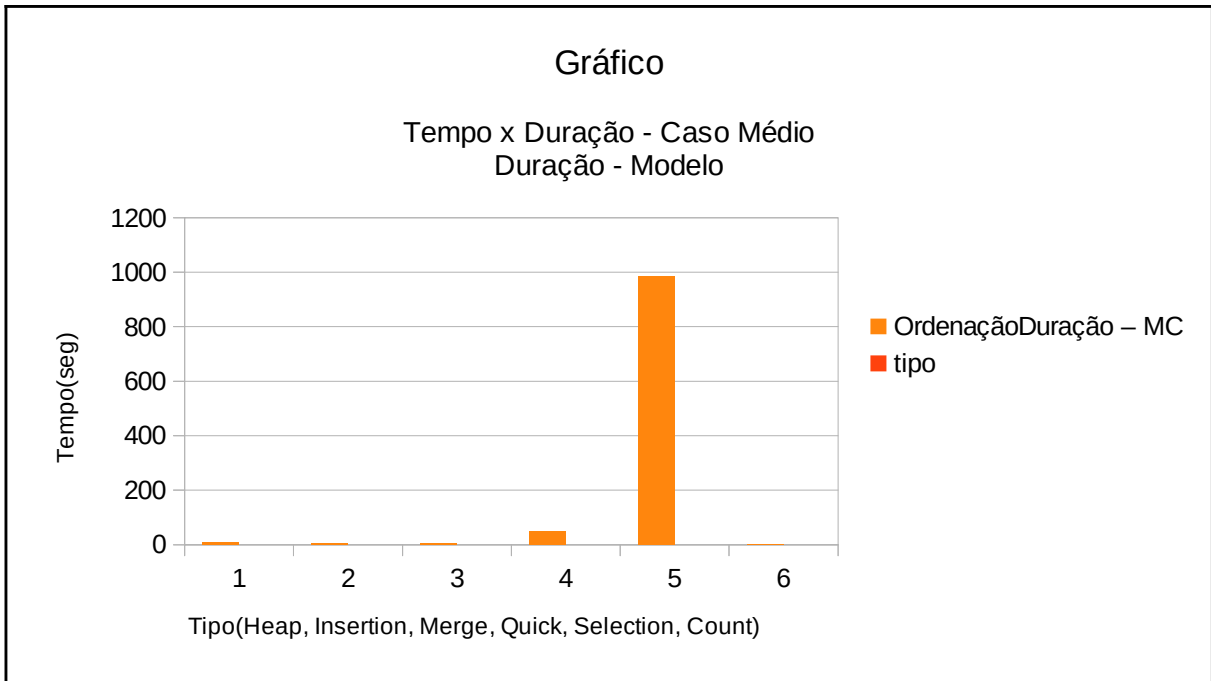






Comparar ordenação x caso médio:





O algoritmo mais eficiente foi o Merge Sort porque foi o que apresentou melhor performance em cada caso (melhor, médio e pior caso) e que conseguiu se ajustar adequadamente aos diferentes modelos de ordenação sugeridos ou seja independente da disposição dos dados a ordenação é eficiente de maneira que ele se comporta como $n \cdot \log n$ nos 3 casos.