

Relatório 3 de Laboratório de Introdução à Ciências da Computação 2

Vítor Amorim Fróis

28 de dezembro de 2021

Resumo

Durante o último módulo da matéria de Laboratório de ICC2, foram estudados métodos de ordenação que não utilizam comparação: Counting, Bucket e Radix, obtendo assim resultados melhores em diversas situações. Apesar disso, existem limitações para esses algoritmos, que nem sempre funcionam tão bem.

Esse relatório busca analisar isso de modo a detalhar quando cada método pode ou não ser útil, além de comparar o que foi estudado ao longo do semestre.

1 Introdução

Ao longo do último semestre, foram estudados diversos métodos de ordenação, alguns com complexidade melhor ou pior de acordo com as regras da notação assintótica.

É fato que dependendo da situação, cada algoritmo pode performar melhor ou pior. Até agora, entretanto, nenhum dos métodos estudados ultrapassa a barreira de $\Theta(n) = n \log n$. Entretanto, ao olhar o problema por outra forma é possível melhorar os métodos. Esse relatório possui como objetivo analisar tais métodos, trazendo suas análises assintóticas, os melhores e piores casos, além de gráficos comparativos com o que foi estudado nos últimos relatórios.

2 Metodologia e desenvolvimento

2.1 A barreira logarítmica

Ao longo dos últimos módulos foram estudados métodos que se utilizavam de comparações para ordenar. Tais algoritmos possuem um limite inferior que pode ser demonstrado a partir de uma árvore de decisões binária. [2]

Existem $n!$ resultados possíveis para um algoritmo de ordenação, isto é, a permutação de todos os elementos do vetor. Uma árvore de decisão do modelo que usa somente comparações deve possuir $n!$ folhas, e portanto uma altura de no mínimo $\Omega(\log(n!))$. É possível aplicar a aproximação de Stirling para encontrar que $n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$ ou considerar que $n! > (n/2)^{n/2}$. Assim:

$$\log(n!) = \log(n/2)^{n/2}$$

$$\log(n!) = \frac{n}{2} \log(n/2)$$

$$\log(n!) = n \log(n)$$

Isso nos leva a uma complexidade mínima de $\Theta(n) = n \log(n)$ [1]. Entretanto, ao quebrar a barreira das comparações, é possível melhorar os algoritmos.

2.2 Counting Sort

Como o próprio nome acusa, o Counting Sort é um algoritmo que se vale da contagem a fim de ordenar sequências de números. Para isso, é necessário um vetor auxiliar de tamanho $k = \max - \min$, que incrementa o valor de das chaves de acordo com os índices do original. Logo, é necessário $O(n)$ para inserir no vetor auxiliar e $O(k)$ ordenar no final.

Complexidade de tempo: $O(n + k)$

Complexidade de memória: $O(k)$

Espera-se que o Counting Sort performe bem com números inteiros e chaves em pequenos intervalos, dado que para um k muito grande, a complexidade do algoritmo cresce muito. Diferente dos métodos estudados nos módulos 1 e 2, a taxa de desordem do vetor não importa.

Para o código a seguir e dos próximos métodos, considere que a representa operações aritméticas, b atribuições, e c comparações.

```
1 void countingsort(int* v, int n) {
2     // 0 - encontrar chaves min e max
3     int final[n];
4     int max, min, i;
5     max = min = v[0];
6     for (int i=1; i < n; i++) {
7         if (v[i] > max) max = v[i]; // c+b
8         if (v[i] < min) min = v[i]; // c+b
9     }
10    //N(2b+2c)
11
12    // 1 - cria vetor de contagem
13    int k = max-min+1;
14    int *contagem = (int*) calloc(k, sizeof(int));
15    // Kb
16
17    // 2 - conta a frequencia de cada chave
18    for (i = 0; i < n; i++)
19        contagem[v[i]-min]++; // a + a
20    // 2aN
21
22    contagem[0] = 0;
23    for(i = 1; i < k; i++)
24        contagem[i] += contagem[i-1]; //a + b
25    //K(a+b)
26
27    // 3 - recria o vetor usando as frequencias para cada chave no vetor de contagem
28    for(i = 0; i < n; i++)
29        final[contagem[v[i]-min]++] = v[i]; //a+a+b
30    //(2a+b)
31
32    for(i = 0; i < n; i++)
33        v[i] = final[i]; //b
34    //Nb
35
36    // total = 2cN+2bN+Kb+2aN+aK+bK+2aN+bN+bN
37    //         = 4Na + 4Nb + 2Nc + Ka + 2Kb
38    //         = 10N + 3K
39    free(contagem);
40 }
```

2.3 Bucket Sort

Ainda que a implementação do Counting Sort simples seja instável e seja possível resolver isso através do uso de registros, outro método existe para solucionar tal problema. Através do uso de "baldes", o Bucket Sort realiza ordenação por contagem abaixo do limite anteriormente conhecido.

De forma similar ao Counting, é necessário alocar um vetor de listas ligadas de tamanho $k = \max - \min$, que representarão os baldes. Cada balde pode armazenar um valor ou intervalo.¹ Com complexidade $O(n)$ os elementos são inseridos no vetor. Então, basta esvaziar os baldes, ordenados, a partir da cabeça, inserindo no vetor original $(k + n)$.

$$\text{Totalizando: } n + k + n + k \rightarrow 2(n + k)$$

$$O(n + k) \text{ em tempo}$$

Em termos de memória gasta, porém, devido ao uso de ponteiros nas listas ligadas, existe uma constante ϵ [3], tal que a complexidade de espaço é

$$O(n + k + \epsilon)$$

Isso significa que o algoritmo possui complexidade muito semelhante ao Counting Sort, porém usando mais memória auxiliar com o objetivo de manter a estabilidade.

```
1 void bucketsort(int *v, int n) {
2     // 1- encontra min e max
3     int max, min, i;
4     max = min = v[0];
5     for (int i=1; i < n; i++) {
6         if (v[i] > max) max = v[i]; // c+b
7         if (v[i] < min) min = v[i]; // c+b
8     }
9     // (2c+2b)n
10
11    // 2- criar um vetor auxiliar contendo listas (buckets)
12    // cada bucket possui um ponteiro para o inicio e outro
13    // para o fim da lista
14    int k = max-min+1;
15    Bucket *B = (Bucket *) calloc(k, sizeof(Bucket));
16    for(int i = 0; i < k; i++)
17        B[i].begin = NULL;
18    // Kb
19
20    // 3 - preenche os buckets com as chaves do vetor
21    // de entrada
22    for (i = 0; i < n; i++) { //(c+a)
23        int key = v[i]; //b
24        Node *newnode = malloc(sizeof(Node));
25        newnode->elem = v[i]; //b
26        newnode->next = NULL; //b
27        //3a
28
29        if (B[key].begin == NULL) B[key].begin = newnode; //b
30        else (B[key].end->next = newnode; //b
31
32        B[key].end = newnode; //b
33        //3b+c+2b+a+c
34    }
35    //n(a+5b+2c)
36
37    // 4 - percorre cada bucket, removendo os elementos do inicio da fila e inserindo na posicao
38    // correta
39    int j; // percorre buckets
40    i = 0; // percorre vetor de entrada //b
41    for (j = 0; j < k; j++) { //a+c
42        Node *pos;
```

¹Ao receber um intervalo, é possível, de acordo com a aplicação, aplicar outro método de ordenação que funcione bem com pequenos vetores ou intervalos para ordenar cada Bucket.

```

42 pos = B[j].begin; //b
43 while (pos != NULL) { //c
44     v[i] = pos->elem; //b
45     i++; //a+b
46     Node *del = pos; //b
47     pos = pos->next; //b
48     B[j].begin = pos; //b
49     free(del);
50 }
51 }
52 // (n+K)(a+5b+2c)
53
54 // total = 2Nc + 2Nb + Kb + Na + 5Nb + 2Nc + Na + 5Nb + 2Nc + Ka + 5Kb + 2Kc
55 //         = 6Nc + 12Nb + 6Kb + 2Na + Ka + 2Kc
56 //         = n(2a+12b+6c) + K(a+6b+2c)
57 //         = n + K
58 free(B);

```

2.4 Radix Sort

Essencialmente, os algoritmos vistos até agora são muito bons, mas ainda é possível fazer melhor. Para um k muito grande, Counting e Bucket não são suficientes. Um método de ordenação mais enxuto é o Radix Sort. A ideia do algoritmo é dividir cada número em potências da base utilizada e ordenar do menos significativo para o mais significativo com ajuda de algum dos dois métodos anteriores para alcançar complexidade de tempo linear.

Utilizando Counting Sort: Para um número qualquer na base b , o número de dígitos máximo d pode ser representado por $\log_b c + 1$ e a complexidade será $d \times \text{Counting Sort}$. Como todos os algarismos estarão entre 0 e b , a complexidade do counting será $O(n + b)$.

$$O((n + b) \times d) \text{ ou}$$

$$O((n + b) \times (\log_b c + 1))$$

Aqui, é importante encontrar uma boa base para diminuir a complexidade de tempo do algoritmo. Ao resolver, encontra-se $b = n$.

$$O((n + n) \times (\log_n c + 1)) \rightarrow O(n \log_n c)$$

Se $c \leq n^p$, o algoritmo se torna $O(np)$. [1] Para o Radix funcionar, é necessário fazer algumas modificações no algoritmo de subrotina.

```

1 void counting(int* original, int n, int exp, int base) {
2     // Nao e preciso calcular min e max, pois min sera 0 e max o exp
3     // 1 - cria vetor de contagem
4     int contagem[base];
5     for(int i = 0; i < base; i++) contagem[i] = 0; //
6     int final[n];
7     // Kb
8
9     // 2 - conta a frequencia de cada chave
10    int i;
11    for (i = 0; i < n; i++)
12        contagem[(original[i]/exp)%base]++; // 4a
13    // 4aN
14    print(contagem, base);
15
16    for(i = 1; i < base; i++)
17        contagem[i] += contagem[i-1]; //K(a+b)
18
19    // 3 - recria o vetor usando as frequencias para cada chave no vetor de contagem
20    for (i = n-1; i >= 0; i--)
21        final[--contagem[(original[i]/exp)%base]] = original[i]; //b + a+a+a
22    //N(3a+b)
23
24    for(i = 0; i < n; i++)
25        original[i] = final[i]; //Na

```

```

26
27 //Total = Kb + 4Na + Ka + Kb + N3a + Nb
28 //      = N(7a+b) + K(a+2b)
29 //      = N + K
30 }

```

Então, simples chamadas do método ordenam um vetor.

```

1 void radixsort(int v[], int n){
2     int max = v[0];
3     for (int i=1; i < n; i++)
4         if (v[i] > max) max = v[i]; // c+b
5         //n(c+b)
6     int base = n;
7
8     for (int exp = 1; max / exp > 0; exp *= base) //log_n(c)
9         counting(v, n, exp, base); //counting
10    //log_n(c) * (N+K)
11 }

```

Veja que a complexidade ainda é logarítmica. Espera-se que com a escolha correta de base, o algoritmo se torne linear.

3 Resultados

Para a seção de resultados, é importante repetir as mesmas experiências que foram realizadas com outros métodos afim de fazer boas comparações. Portanto, houveram testes de ordenação com vetores de 5000 a 100000 elementos. Todos os vetores utilizados foram aleatórios, pois a falta de comparação nos algoritmos faz com que ordenação prévia não seja um fator importante na análise.

Entretanto, como destacado na parte de Metodologia, a complexidade de tempo depende do tamanho máximo do vetor k . Assim, para comparação desses 3 métodos, haverão 3 gráficos inicialmente, gerados a partir de vetores uniforme, denso e esparso.

3.1 Vetor Uniforme

No vetor uniforme, os elementos são distribuídos de maneira uniforme pelo vetor. Ou seja, $k = n$, e consequentemente, $O(n + k) \rightarrow O(2n)$. É esperado que todos vetores sejam lineares, incluindo o Radix Sort que conta com a escolha adequada de n , descrita anteriormente, para alcançar complexidade $O(n)$.

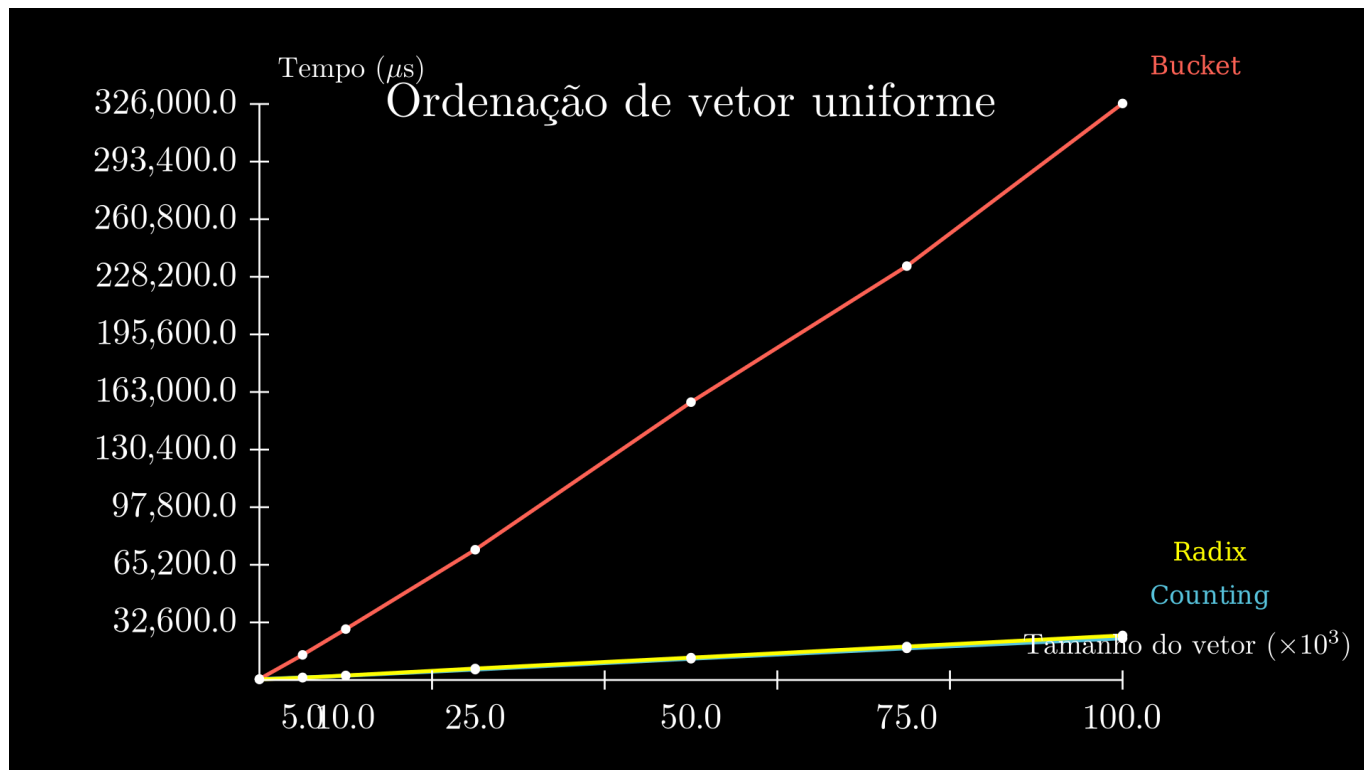


Figura 1: Ordenação de vetores uniformes com até 100000 elementos

De fato, os vetores cumprem a linearidade, contudo, o Bucket Sort performa muito pior que os outros. Isso acontece devido as constantes excluídas na análise assintótica, fator que interferiu em análises passadas também. Além disso, é possível perceber que Counting e Radix ficam bem próximos. E isso se deve ao fato de que para vetores uniformes, o Radix é apenas um Counting Sort.

3.2 Vetor Denso

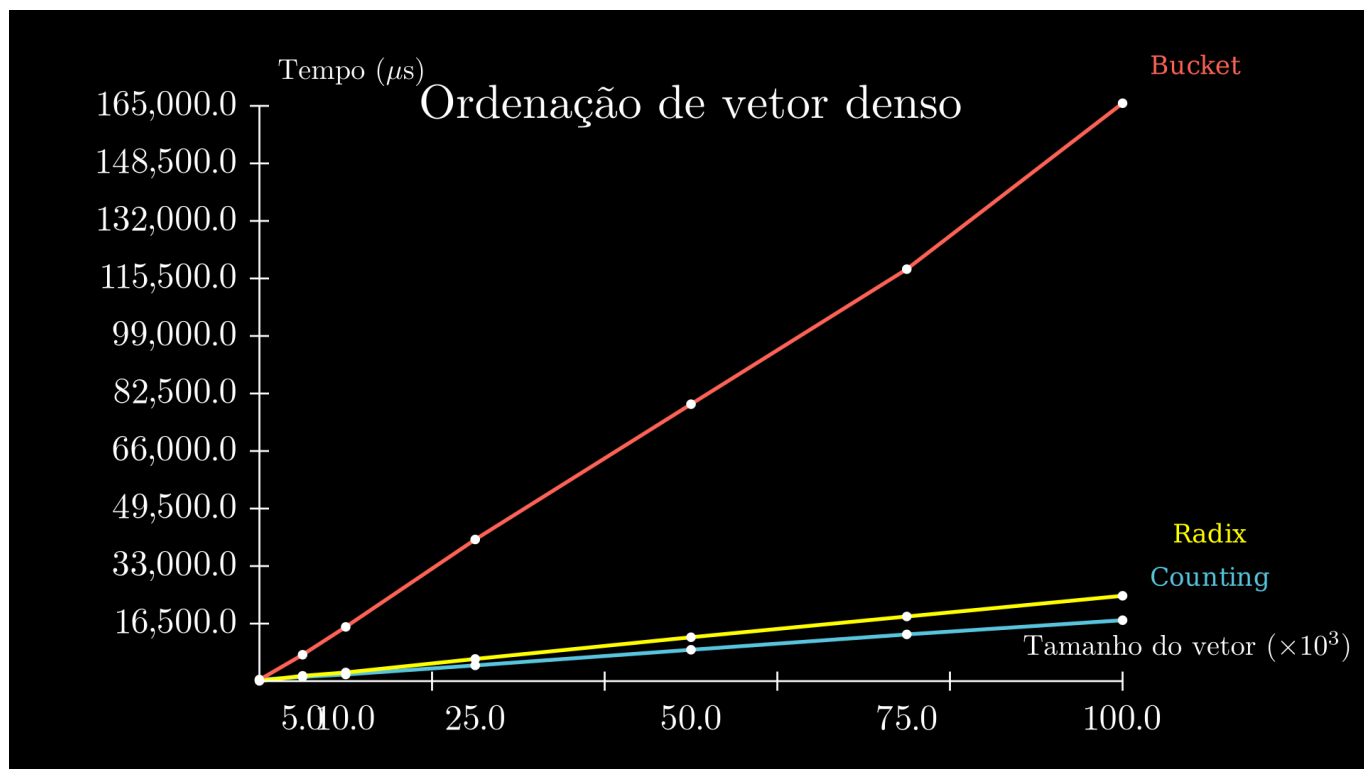
Já no vetor denso, os elementos possuem tamanho máximo de $\frac{1}{100}$ do número de elementos. Isso significa que há diversas repetições, cenário em que os algoritmos analisados melhor devem atuar, pois considerando que a complexidade média é $O(n + k)$ e quando $k = \frac{n}{100}$,

$$O(n + \frac{n}{100}) < O(n + n). \text{ Veja 3.1}$$

A fração é desprezível: $O(n) < O(2n)$.

Assim, os métodos devem apresentar desempenho 2 vezes mais rápido que no vetor uniforme.

Figura 2: Ordenação de vetores densos com até 100000 elementos

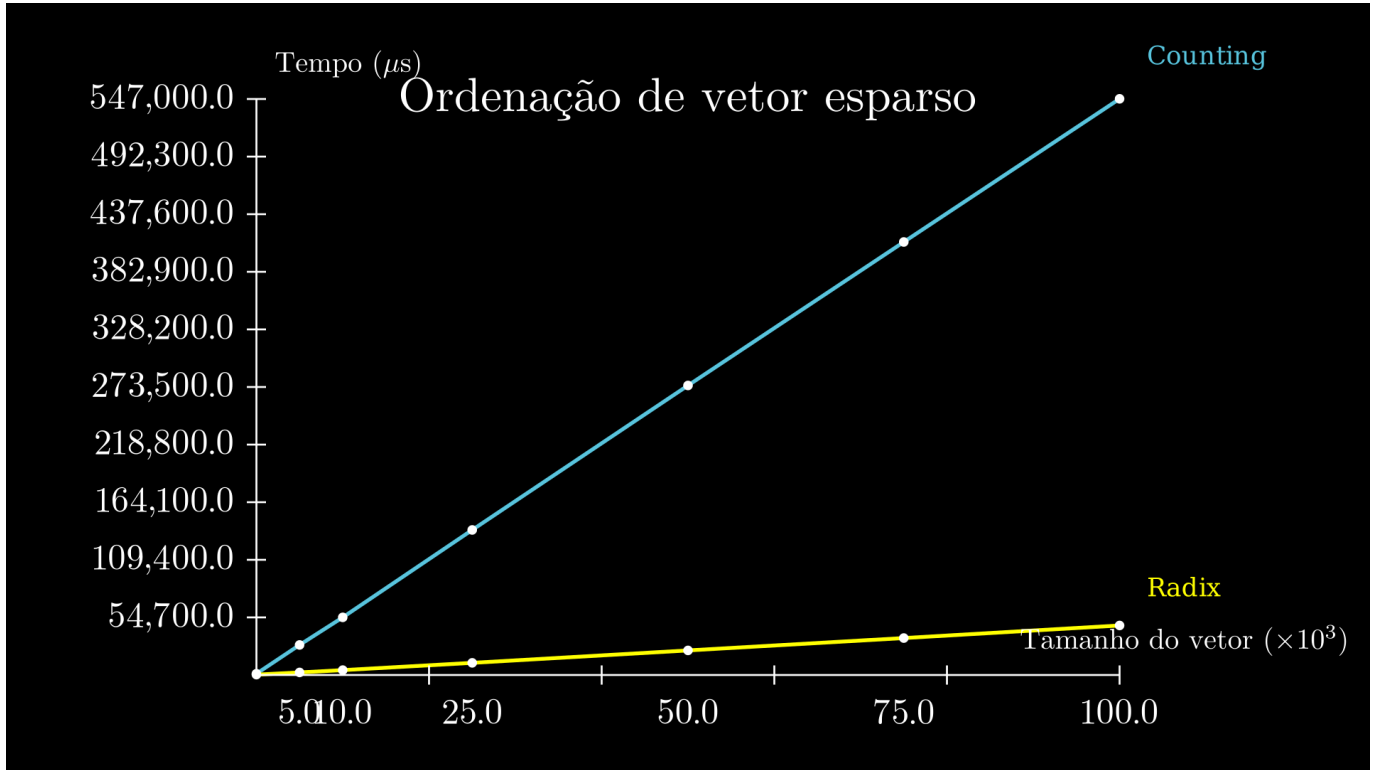


Ao comparar com o gráfico anterior (1), é possível observar que isso realmente acontece. O tempo médio passa de aproximadamente 32000 μs para 16000 μs .

3.3 Vetor Esperso

Ao contrário do denso, o vetor esperso possui os elementos espalhados e pouca repetição. O valor máximo é até $100\times$ o tamanho do vetor. Logo, a complexidade é $O(n + k) \rightarrow O(n + 100n)$, isto é, $55.5\times$ pior que no vetor uniforme para Counting e Bucket 1. Entretanto, com a escolha correta de base para o Radix, é possível manter baixo tempo de execução, como mostrado em 3.

Figura 3: Ordenação de vetores esparsos com até 100000 elementos



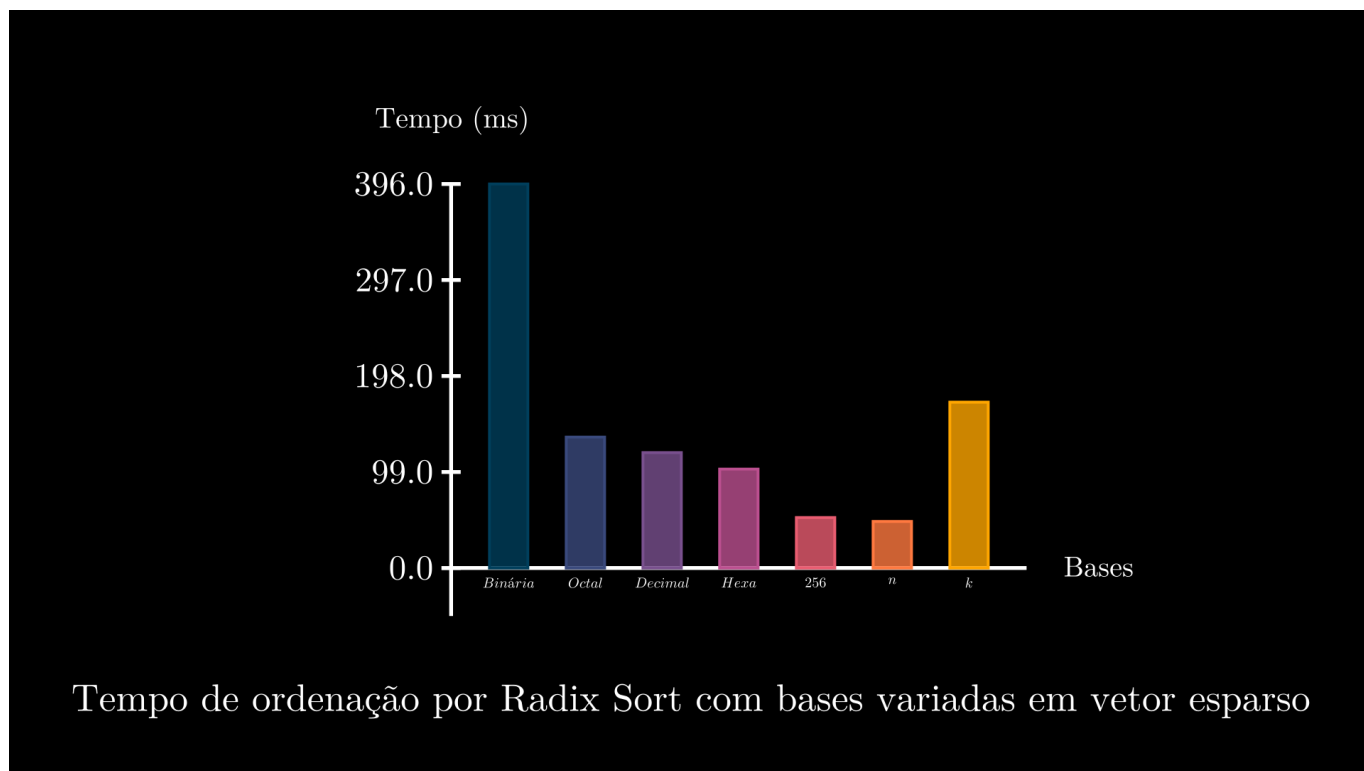
O computador não conseguiu suportar o alto custo de espaço do Bucket Sort, com milhares de listas encadeadas, excedendo a memória reservada para o funcionamento do algoritmo e assim sendo retirado do gráfico.

Nos algoritmos restantes, entretanto, Counting funciona muito mal, devido ao alto valor de k , enquanto com um bom b , o Radix continua com pequena constante.

3.4 Base do Radix

O diferencial dos dois outros métodos para o Radix Sort é a escolha de base, que permite ao algoritmo quebrar um número em partes para ordenar com tempo linear e baixo tempo de execução em qualquer situação. O gráfico 4 mostra porque a melhor escolha de base é $b = n$.

Figura 4: Eficiência de algumas bases comuns no Radix Sort ordenando vetores esparsos



Nesse gráfico foi utilizado um vetor de tamanho 100000 e máximo de 10000000. De todas as bases, binária, decimal, octal, k , a que melhor performa é n . Assim, confirmam-se os cálculos feitos na parte de metodologia sobre o Radix Sort. 2.4 apesar de ainda linear, é bem pior que o Radix.

3.5 Comparação com outros métodos

Também é preciso comparar os métodos estudados nesse relatório com os passados. Ainda que os lineares tendam a ser melhores no geral, cada um tem suas nuances e especializações, de forma que é importante analisar a situação antes de escolher o método com o qual se deseja trabalhar. Para a construção do gráfico 5 foram utilizados vetores de tamanho 10000 a 50000, randomicamente gerados.

Figura 5: Os métodos de ordenação estudados ao longo do semestre

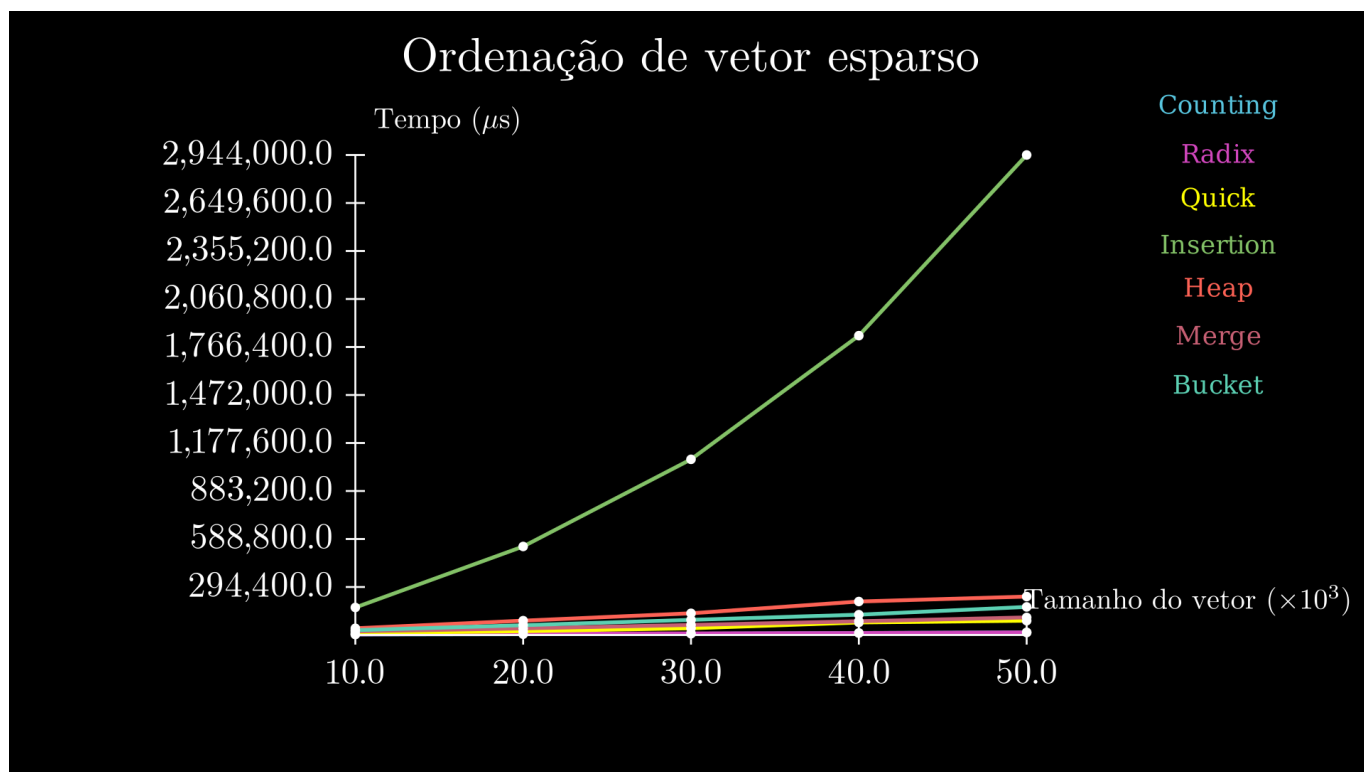
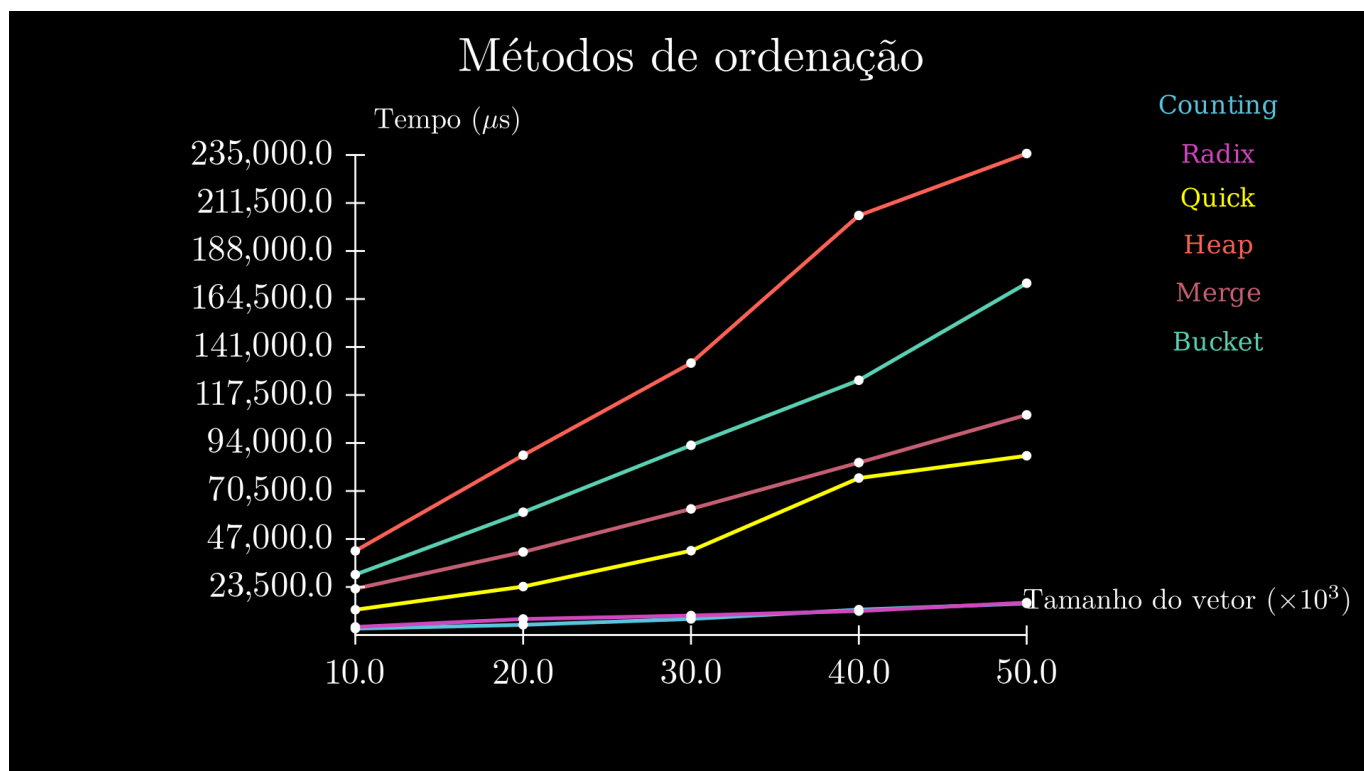


Figura 6: Exclui-se o método com eficiência quadrática para permitir melhor visualização



O Insertion, de complexidade $O(n^2)$, atrapalha a análise dos outros métodos. Ao retirá-lo, obtemos o gráfico 6, no qual é possível observar os métodos de contagem se sobressaindo, enquanto aqueles que utilizam comparação se mostram mais demorados.

4 Conclusão

Durante a elaboração desse relatório, além de aprender mais sobre a ferramenta LaTeX, foi possível constatar a eficiência dos métodos estudados na teoria, bem como comparar com o que havia sido estudado ao longo do ano. Com o Bucket Sort sofrendo por *Stack Overflow*, foi possível entender quão grave é a super utilização de ponteiros. Também pude constatar que a utilização da base $n = k$ é a mais eficiente dentre todas, no método mais eficiente. Por fim, da comparação com outros métodos que estudamos ao longo do ano, foi possível comparar a linearidade e linear-logaritimicidade entre eles.

Referências

- [1] Erik Demaine. *7. Counting Sort, Radix Sort, Lower Bounds for Sorting*. MIT OpenCourseWare, 2013. URL: <https://www.youtube.com/watch?v=Nz1KZXbghj8>.
- [2] Jason Ku Erik Demaine e Justin Solomon. *Introduction to Algorithms. Recitation 5*. URL: https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-006-introduction-to-algorithms-spring-2020/lecture-notes/MIT6_006S20_r05.pdf.
- [3] Moacir Antonelli Ponti. *ICC2 (3) 3 - Bucketsort: algoritmo e análise*. Universidade de São Paulo, 2021. URL: <https://www.youtube.com/watch?v=8ZyFKHJJJIA>.