

Pesquisa e Ordenação de Dados

Unidade 1:

Contextualização

Pesquisa e Ordenação de Dados

- Ementa (PPC):
 - Métodos de ordenação. Busca linear e binária. Organização de arquivos. Persistência de dados. Ordenação externa. Índices (árvores B+ e hashing). Compactação de dados. Implementações com linguagem imperativa estruturada.
- Objetivo geral (PPC):
 - Utilizar estruturas de dados avançadas para ordenação e pesquisa de informações. Construir algoritmos para persistir dados e tratar dados persistidos. Analisar comparativamente os métodos de ordenação quadráticos e os métodos MergeSort, HeapSort e QuickSort.

Por que estudar pesquisa?



search

- Uma das funcionalidades mais úteis de um computador é sua capacidade de armazenar e recuperar dados
- Dados podem ser armazenados:
 - Na memória principal
 - Na memória secundária
 - Arquivos de texto, arquivos binários ou por intermédio de um SGBD
- Uma **pesquisa** ou (busca) visa **recuperar** um ou mais itens de um arquivo ou conjunto de dados
- É interessante que a busca não precise percorrer exaustivamente todo o conjunto de dados
 - Para isto, diversos algoritmos e estruturas de dados foram desenvolvidos, conforme veremos adiante.

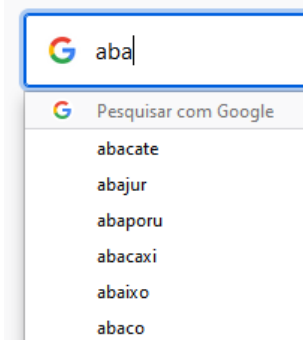
Por que estudar pesquisa?

- Descobrir se um item existe em um conjunto de dados



*existe o login
asebben123 no
cadastro de
usuários?*

- Listar valores que correspondem a um determinado critério de busca



*quais termos de
busca começam
com a string aba?*

- Recuperar dados associados a uma chave de busca única



*qual a descrição e o
preço do produto com
código de barras
123456789?*

- Listar todos os itens de um conjunto



SERVIÇO PÚBLICO FEDERAL
UNIVERSIDADE FEDERAL DA PÁRAMO, PA
PROFESSORIA DE EDUCAÇÃO
Rua 40, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000



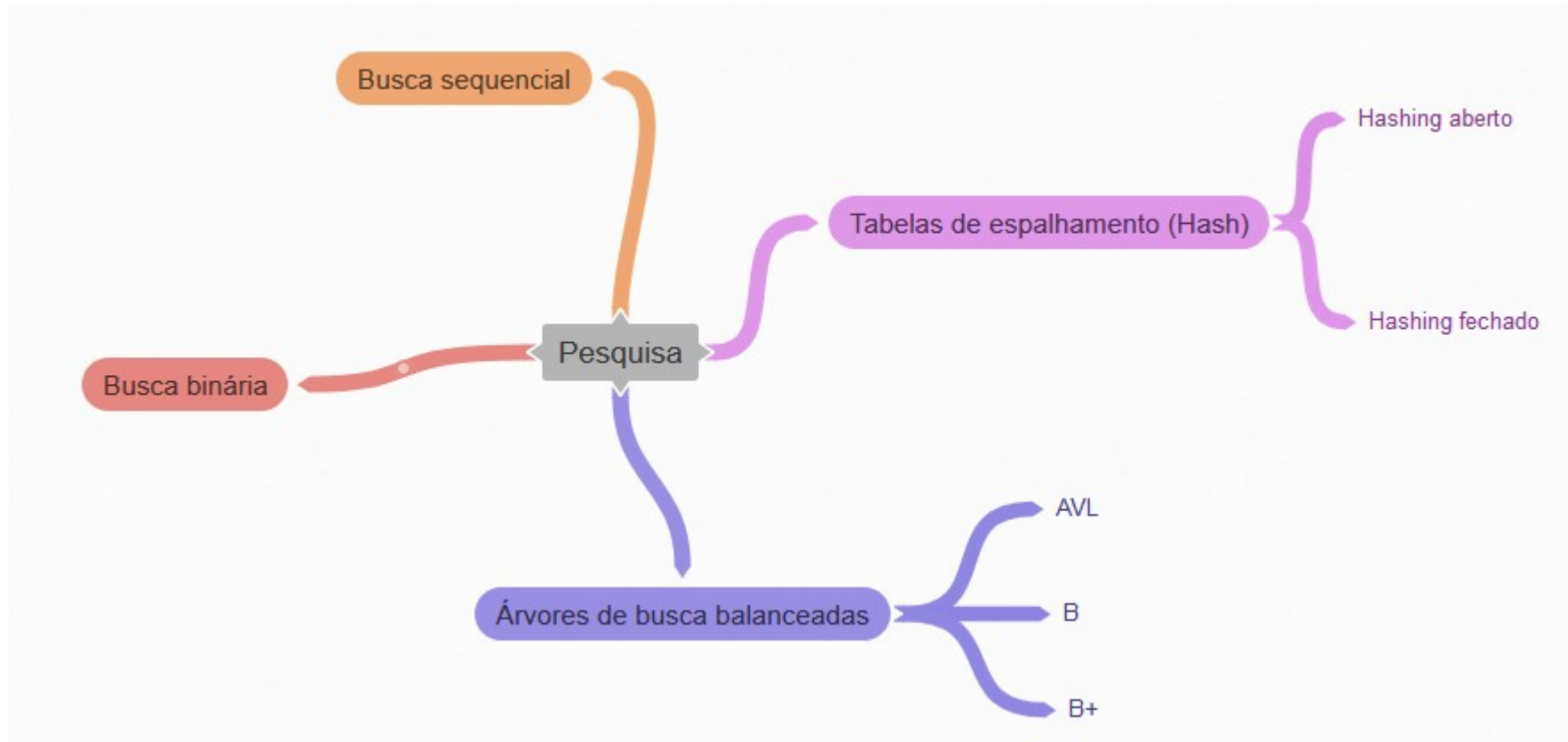
DIÁRIO DE CLASSE

Professor(a): ANTONIO DA SILVA
Componente Curricular: GEMAS - Jogos e brincadeiras de dados
Tema: 1 - Jogos de Comparação/Tema 1 - 7 (Jogos - Jogos de Comparação) - Jogos de Comparação
Tema 1 - Jogos de Comparação/Tema 1 - 7 (Jogos - Jogos de Comparação) - Jogos de Comparação
Tema 1 - Jogos de Comparação/Tema 1 - 7 (Jogos - Jogos de Comparação) - Jogos de Comparação

Nº	Matrícula	Nome	Total de Notas	Frequência	Nota Final	Situação
1	10000001	ALICE VANDER LINS DE LIMA PEREIRA	0	100,00%	0,00	CR
2	10000002	ALDENALDO DE LIMA PEREIRA	0	100,00%	0,00	CR
3	10000003	ANDRÉ LUCAS NUNES OLIVEIRA	0	100,00%	0,00	CR
4	10000004	BARBARA DE LIMA PEREIRA	0	100,00%	0,00	CR
5	10000005	BRAUN CARRELLA PEREIRA	0	100,00%	0,00	CR
6	10000006	BREUNA GABRIELA PEREIRA	0	100,00%	0,00	CR
7	10000007	CECÍLIA DE LIMA PEREIRA	0	100,00%	0,00	CR
8	10000008	DANIELA DE LIMA PEREIRA	0	100,00%	0,00	CR
9	10000009	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
10	10000010	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
11	10000011	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
12	10000012	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
13	10000013	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
14	10000014	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
15	10000015	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
16	10000016	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
17	10000017	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
18	10000018	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
19	10000019	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR
20	10000020	EDUARDO DE LIMA PEREIRA	0	100,00%	0,00	CR

*quem são os alunos
matriculados na
turma 27352?*




Pesquisa



Por que estudar ordenação?



- **Rearranjar um conjunto** de itens em ordem (ascendente ou descendente) para agilizar sua recuperação posterior é uma tarefa bastante comum em nosso cotidiano
- Esta atividade de colocar as coisas em ordem está presente na maioria das aplicações nas quais objetos armazenados precisam ser pesquisados e recuperados, tais como dicionários, índices de livros, tabelas e arquivos (ZIVIANI, 2004).

Selecionar	Foto do usuário	Nome / Sobrenome	Endereço de email	Cidade/Município	País	Último acesso ao curso
<input type="checkbox"/>		ANDRESSA SEBBEN	asebben@uffs.edu.br	CHAPECÓ/SC	Brasil	4 segundos
<input type="checkbox"/>		[REDACTED]	[REDACTED]	CHAPECÓ/SC	Brasil	1 dia 7 horas
<input type="checkbox"/>		[REDACTED]	[REDACTED]	CHAPECÓ/SC	Brasil	4 dias 6 horas

9345 Produtos | ■ ■ ■ ■ ■ ■ ■ ■

-33%

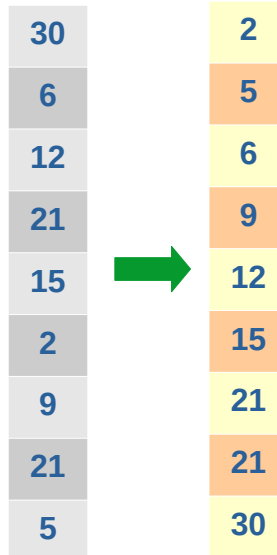


- Mais Populares ▾
- Mais Populares
 - Novidades
 - Menor Preço
 - Maior Preço
 - Maior desconto

- A ordenação é um problema largamente estudado na Ciência da Computação.
 - Veremos diversos algoritmos, com diferentes características.

Por que estudar ordenação?

- Quando uma lista de valores está ordenada, diversos problemas tornam-se fáceis - ou, pelo menos, mais fáceis (VISUALGO.NET):
 - Buscar por um valor específico na lista (busca binária)
 - Encontrar o valor mínimo, o valor máximo ou o k-ésimo menor/maior valor da lista
 - Verificar a unicidade de um valor e deletar duplicatas
 - Contar quantas vezes um valor específico aparece na lista
 - Calcular a mediana, bem como identificar os valores discrepantes ou “*outliers*”
 - Obter a intersecção/união entre a lista ordenada A e outra lista ordenada B
 - Encontrar valores x e y pertencentes à lista tal que $x+y$ resulta em um valor alvo z
- Aplicações menos óbvias incluem: compressão de dados, computação gráfica, biologia computacional, balanceamento de carga em sistemas paralelos, etc (SEDGEWICK, 2016).

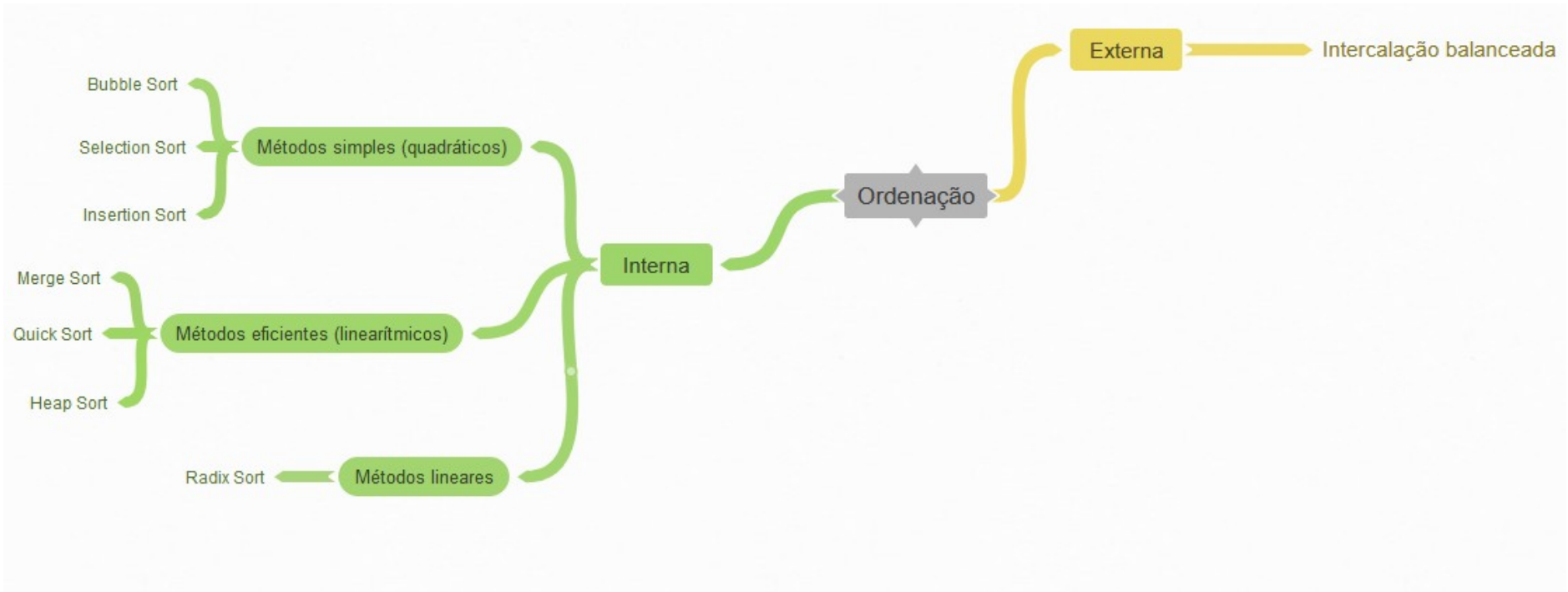


30	2
6	5
12	6
21	9
15	12
2	15
9	21
21	21
5	30

Ordenação (*Sorting*)

- Ordenação interna:
 - quando o conjunto de dados a ser ordenado encontra-se integralmente na memória principal;
 - qualquer registro pode ser imediatamente acessado;
- Ordenação externa:
 - quando o conjunto a ser ordenado não cabe na memória principal e, por isso, precisa ser armazenado em um arquivo em memória secundária (fita ou disco);
 - os registros são acessados sequencialmente ou em grandes blocos.

Ordenação – O que veremos



Análise de Algoritmos

- Dada a relevância da pesquisa e da ordenação, é altamente desejável que os métodos utilizados sejam **eficientes**
 - Não basta resolver o problema, é necessário ter um desempenho aceitável na prática!
 - Um algoritmo simples e rápido para um conjunto pequeno de dados pode se tornar inviável para um conjunto com milhões de registros
- Analisar a **complexidade** de um algoritmo = estimar os recursos necessários para que a tarefa seja completada
 - **Complexidade de tempo** → medida do tempo de execução
 - **Complexidade de espaço** → medida da quantidade de memória

Medida mais comum

Análise de Algoritmos

- Formas de computar:
 - **Análise empírica** (*benchmark*)
 - Implementar o algoritmo e realizar **experimentos** com diversos conjuntos de dados, de diferentes tamanhos e características, **medindo** o tempo de execução/consumo de memória;
 - Os resultados são dependentes da linguagem de programação utilizada, do compilador, da capacidade do hardware, da quantidade de processos sendo executados, etc;
 - Computa custos não aparentes (como alocação de memória);
 - Permite comparar computadores, linguagens, etc;
 - Depende da habilidade do programador;
 - Cenários de teste limitados.

Análise de Algoritmos

- Formas de computar:
 - **Análise matemática** (modelo teórico)
 - Estudo **formal** das propriedades do algoritmo;
 - Considera um computador idealizado onde cada operação executa em tempo constante e de forma sequencial;
 - Operações simples (comparações, atribuições, cálculos, incrementos, acesso a um elemento em um array) possuem mesmo custo
 - Realiza simplificações buscando considerar apenas o custo dominante;
 - Objetivo: **estimar como o algoritmo se comporta em função do tamanho do conjunto de dados** (entrada);
 - Independe de aspectos específicos de hardware, linguagem, compilador e ambiente de execução;

Complexidade

- Ao olhar para complexidade dos algoritmos, podemos:
 - prever seu desempenho;
 - comparar diferentes algoritmos que resolvem o mesmo problema;
 - avaliar sua viabilidade, provendo algumas garantias de que sua execução será completada no tempo esperado.
- Ao invés de olhar para o tempo exato de execução (o qual depende de fatores relacionados ao ambiente de execução), a noção de complexidade possibilita analisar e entender o **comportamento** do algoritmo, isto é,
 - **como o tempo de execução cresce à medida que o tamanho da entrada cresce.**

Análise de Complexidade

Exemplo

```
int menor(int *A, int n) {  
    int i, menor;  
    menor = A[0];  
    for(i = 1; i < n; i++) {  
        if(A[i] < menor)  
            menor = A[i];  
    }  
    return menor;  
}
```

1

2 da inicialização do laço
n-1 comparações
n-1 incrementos

$2+(n-1)+(n-1)$

n-1

n-1

no pior caso

1

$$\begin{aligned} &= 1 + 2 + n - 1 + n - 1 + n - 1 + n - 1 + 1 \\ &= 4n \end{aligned}$$

$$\cancel{4n} = n$$

Comportamento assintótico

60	50	40	30	20	10
----	----	----	----	----	----

0	1	2	3	4	5
---	---	---	---	---	---

Notação Assintótica

- Quando trabalhamos com N objetos, algumas operações tomarão tempo proporcional a N
 - Para valores pequenos de N , até mesmo um algoritmo ruim pode ter desempenho aceitável
 - Nos interessa saber a ordem de crescimento do algoritmo para **valores grandes** de N
 - À medida que N cresce, as constantes aditivas e multiplicativas têm cada vez menos impacto, podendo até mesmo se tornar irrelevantes

- Por exemplo, para valores de N suficientemente grandes, as funções

$$n^2$$

$$(3/2)n^2$$

$$9999n^2$$

$$n^2/1000$$

$$n^2+100n$$

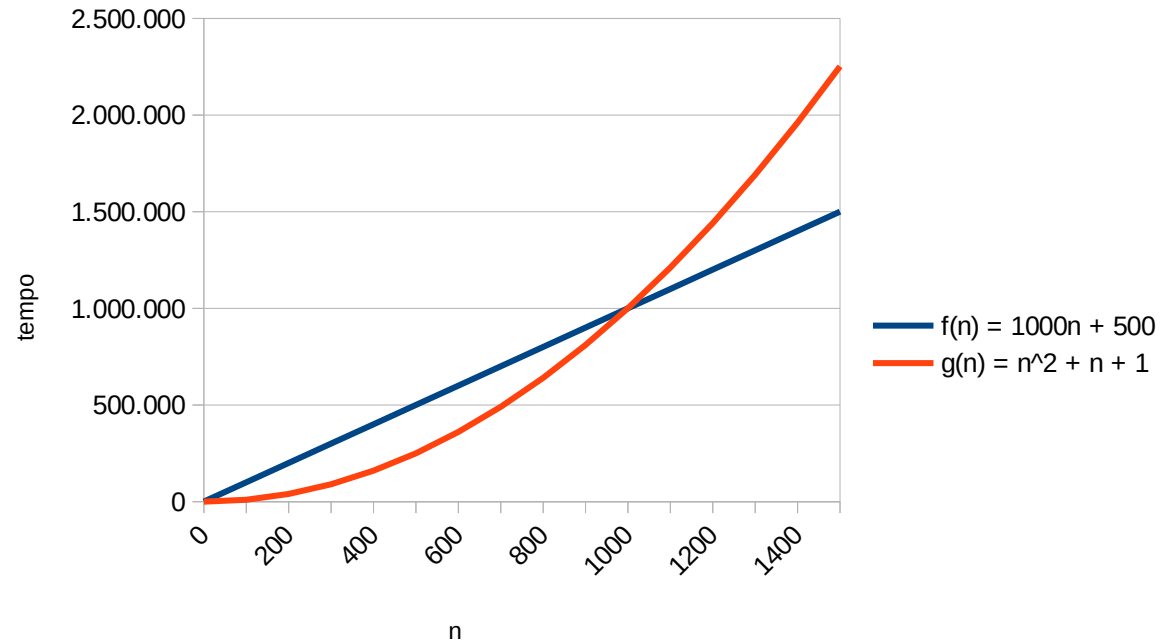
têm todas a mesma taxa de crescimento e, portanto, são todas "equivalentes".

Notação Assintótica

- Ex: considere as funções:

- $f(n) = 1000n + 500$

- $g(n) = n^2 + n + 1$



- Existe um valor de n a partir do qual $g(n)$ é sempre maior do que $f(n)$, tornando os demais termos e constantes pouco significativos.

Notação Assintótica

- Na notação assintótica, considera-se apenas o **termo dominante** da equação (ou seja, o termo de maior grau). Os termos de grau menor e as constantes aditivas e multiplicativas são desconsiderados.

Ex:

$$f(n) = 75$$

=

$$f(n) = 1$$

$$f(n) = 2n + 1$$

=

$$f(n) = n$$

$$f(n) = n^2 + n$$

=

$$f(n) = n^2$$

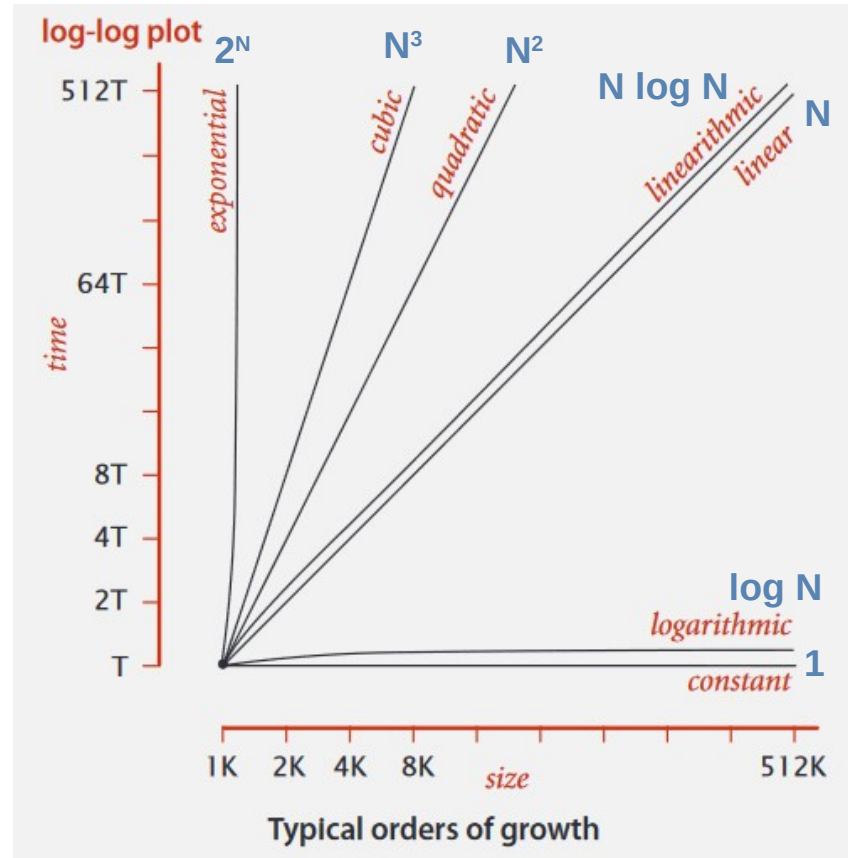
$$f(n) = 10n^3 + 4n - 1$$

=

$$f(n) = n^3$$

quando função não possui n , então o comportamento assintótico é 1 (constante)

Classes de Complexidade (Order of growth)



order of growth	name	typical code framework	description	example	$T(2N) / T(N)$
1	constant	<code>a = b + c;</code>	statement	add two numbers	1
$\log N$	logarithmic	<code>while (N > 1) { N = N / 2; ... }</code>	divide in half	binary search	~ 1
N	linear	<code>for (int i = 0; i < N; i++) { ... }</code>	loop	find the maximum	2
$N \log N$	linearithmic	[see mergesort lecture]	divide and conquer	mergesort	~ 2
N^2	quadratic	<code>for (int i = 0; i < N; i++) for (int j = 0; j < N; j++) { ... }</code>	double loop	check all pairs	4
N^3	cubic	<code>for (int i = 0; i < N; i++) for (int j = 0; j < N; j++) for (int k = 0; k < N; k++) { ... }</code>	triple loop	check all triples	8
2^N	exponential	[see combinatorial search lecture]	exhaustive search	check all subsets	$T(N)$

Complexidade

- **Melhor caso** (*best case*)
 - Representado pela letra grega Ω (ômega)
 - Consiste na entrada mais “fácil” (e a que deve preferencialmente ser buscada)
 - Constitui o limite inferior de custo (não pode ser mais rápido)
- **Pior caso** (*worst case*)
 - Representado pela letra O (ômicron), também chamada de “**Big O**”
 - Representa a entrada mais difícil (a que faz o algoritmo executar o maior número de operações)
 - Constitui o limite superior de custo, servindo como uma garantia para as demais entradas (não pode ser mais lento)
- **Caso médio** (*average case*)
 - Representado pela letra grega θ (theta)
 - Custo esperado para uma entrada aleatória
 - Difícil de determinar na maioria dos casos

Notação “Big O”:
a mais utilizada
para expressar a
complexidade do
algoritmo

Exemplo

Busca sequencial em um vetor de tamanho n

- Melhor caso: $\Omega(1)$
 - O valor procurado é o primeiro do vetor
- Pior caso: $O(n)$
 - O valor procurado é o último ou não faz parte do vetor
 - Será necessário visitar todos os n elementos do vetor até encontrar o valor procurado
- Caso médio: $\theta(n)$
 - Será necessário visitar na média $n/2$ elementos do vetor até encontrar o valor procurado

Dizemos que o algoritmo é $O(n)$, ou seja, linear