

1 实验说明

1.1 实验目的

熟悉主要排序算法，对排序算法的时间复杂度及其影响因素有直观认识。

1.2 实验内容

- 随机生成一个包括 n 个整数的数组，元素取值范围是 $[1, 1000]$ ，利用插入排序，归并排序，快速排序，堆排序，基数排序，桶排序等算法对数组进行非降序排序，记录不同算法的运行时间。
 - 改变数组规模 $n = 5 \times 10^4, 1 \times 10^5, 2 \times 10^5, 3 \times 10^5, 5 \times 10^5$ ，记录不同规模下各个算法的排序时间。
 - 对固定规模 $n = 1 \times 10^5$ 的数组进行随机扰乱，对扰乱后的数组进行排序并记录各个算法的排序时间。本实验要求重复 5 次，观察输入数据分布和运行时间的关系。
 - 各个算法的时间复杂度理论分析，并与实验二进行对比。
- ### 1.3 检查代码正确性
- 产生一个长度为 20 随机数组，允许有重复元素，分别用各种算法进行排序并输出排序结果。
 - 产生一个长度为 1000 的随机数组，允许有重复元素，分别用各种算法进行排序并输出排序结果以及耗费时间，使用 `for` 循环比较其他算法的排序结果是否与插入排序的结果相同。
 - 固定随机数种子，方便复现实验结果。

2 排序算法及其时间复杂度

2.1 插入排序

2.1.1 伪代码

Algorithm 1: Insertion Sort
Input: An array A of n numbers. Output: The array A sorted in nondecreasing order.
<pre>1. for $i = 2$ to n: 2. $key \leftarrow A[i]$ 3. $j \leftarrow i$ 4. while $j > 0$ and $A[j-1] > key$: 5. $A[j] \leftarrow A[j-1]$ 6. $j \leftarrow j-1$ 7. $A[j] \leftarrow key$</pre>

2.1.2 时间复杂度

插入排序算法中，`for` 循环的迭代次数为 $n-1$ ，最坏情况下，`while` 内循环的迭代次数依次为 $n-1, n-2, \dots, 1$ ，因此插入排序的运行时间为

$$T(n) = \Theta(n) + \sum_{i=1}^{n-1} O(i) = O(n^2).$$

最好情况下，数组已有序，不执行交换，时间复杂度为 $O(n)$ 。平均情况下，比较次数可大致视为最坏情况的一半，时间复杂度仍为 $O(n^2)$ 。

2.2 归并排序

2.2.1 伪代码

Algorithm 2: Merge Sort
Input: An array A of n numbers. Output: The array A sorted in nondecreasing order, returned by <code>MERGE-SORT($A, 1, n$)</code> .
<pre>function MERGE-SORT(A, p, r): 1. if $p < r$: 2. $q \leftarrow \lfloor \frac{p+r}{2} \rfloor$ 3. MERGE-SORT(A, p, q) 4. MERGE-SORT($A, q+1, r$) 5. MERGE(A, p, q, r) function MERGE(A, p, q, r): 1. $n_1 \leftarrow q-p+1, n_2 \leftarrow r-q$ 2. Let $L[1..n_1+1]$ and $R[1..n_2+1]$ be new arrays. 3. for $i = 1$ to n_1: 4. $L[i] \leftarrow A[p+i-1]$ 5. for $j = 1$ to n_2: 6. $R[j] \leftarrow A[q+j]$ 7. $L[n_1+1] \leftarrow \infty, R[n_2+1] \leftarrow \infty$ 8. $i \leftarrow 1, j \leftarrow 1$ 9. for $k = p$ to r: 10. if $L[i] \leq R[j]$: 11. $A[k] \leftarrow L[i]$ 12. $i \leftarrow i+1$ 13. else: 14. $A[k] \leftarrow R[j]$ 15. $j \leftarrow j+1$</pre>

2.2.2 时间复杂度

MERGE 操作的时间复杂度为 $\Theta(n)$ ，故归并排序的运行时间 $T(n)$ 的递归式：

$$T(n) = 2T\left(\frac{n}{2}\right) + \Theta(n).$$

由主定理第二种情况可得， $T(n) = O(n \log n)$ 。

2.3 快速排序

2.3.1 伪代码

Algorithm 3: Quick Sort
Input: An array A of n numbers. Output: The array A sorted in nondecreasing order, returned by <code>QUICK-SORT($A, 1, n$)</code> .
<pre>function PARTITION(A, p, r): 1. while $p < r$: 2. $q \leftarrow \text{PARTITION}(\mathbf{A}, \mathbf{p}, \mathbf{r})$ 3. QUICK-SORT($A, p, q-1$) 4. $p \leftarrow q+1$ function PARTITION(A, p, r): 1. $x \leftarrow A[r]$ 2. $i \leftarrow p$ 3. for $j = p$ to $r-1$: 4. if $A[j] \leq x$: 5. SWAP($A[i], A[j]$) 6. $i \leftarrow i+1$ 7. SWAP($A[i], A[r]$) 8. return $i+1$</pre>

2.3.2 时间复杂度

PARTITION 操作的时间复杂度为 $\Theta(n)$ 。最坏情况下，划分产生的子数组规模为 $n-1$ 和 0，运行时间为

$$T(n) = T(n-1) + \Theta(n) = O(n^2).$$

最好情况下，划分产生的子数组规模平衡，运行时间为

$$T(n) = 2T\left(\frac{n}{2}\right) + \Theta(n) = O(n \log n).$$

当划分为常数比例时，运行时间仍为 $O(n \log n)$ 。

平均情况下，为方便分析，将数组 A 的各个元素重新命名为 a_1, a_2, \dots, a_n 其中 a_i 是第 i 小的元素。定义 $A_{ij} = \{a_i, a_{i+1}, \dots, a_j\}$ 。假定使用 RANDOMIZED-PARTITION 操作，主元的选取独立且随机。令随机变量 X_{ij} 指示 a_i 与 a_j 是否被比较， X 为总比较次数，则

$$\begin{aligned} \mathbb{E}[X] &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr\{a_i \text{ and } a_j \text{ are compared}\} \\ &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr\{a_i \text{ or } a_j \text{ is the first pivot chosen from } A_{ij}\} \\ &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} < n \sum_{k=1}^n \frac{2}{k} = O(n \log n). \end{aligned}$$

故平均情况的运行时间为 $O(n + n \log n) = O(n \log n)$ 。

2.4 堆排序

2.4.1 伪代码

Algorithm 4: Heap Sort
Input: An array A of n numbers. Output: The array A sorted in nondecreasing order.
<pre>function BUILD-MAX-HEAP(A, n): 1. $swap \leftarrow n$ to 2: 2. MAX-HEAPIFY($A, 1, i-1$) function BUILD-MAX-HEAP(A, n): 1. $h \leftarrow \lfloor n/2 \rfloor$ 2. for $i = h$ to 1: 3. MAX-HEAPIFY(A, i, n) function MAX-HEAPIFY(A, n, i): 1. largest $\leftarrow i$ 2. repeat: 3. $l \leftarrow 2i, r \leftarrow 2i+1$ 4. if $l \leq n$ and $A[l] > A[largest]$: 5. largest $\leftarrow l$ 6. if $r \leq n$ and $A[r] > A[largest]$: 7. largest $\leftarrow r$ 8. SWAP($A[i], A[largest]$) 9. until $i \neq largest$</pre>

2.4.2 时间复杂度

由于树高为 $O(\log n)$ ，MAX-HEAPIFY 的时间复杂度为 $O(\log n)$ ，故 BUILD-MAX-HEAP 的时间复杂度为 $O(n \log n)$ （事实上，更紧确的界为 $O(n)$ ，此处不赘述），进而 HEAP-SORT 的时间复杂度为 $O(n \log n)$ ， $\text{BUILD-MAX-HEAP} = O(n \log n)$ 。

最好情况下，建堆后数组已有序，MAX-HEAPIFY 花费 $\Theta(1)$ 时间，总时间复杂度为 $O(n)$ 。

2.5 基数排序

2.5.1 伪代码

Algorithm 5: Radix Sort
Input: An array A of n numbers, each of which has d digits and is in base k . Output: The array A sorted in nondecreasing order.
<pre>1. Let $C[0..k-1]$ be an array of size k, and $B[1..n]$ be an array of size n. 2. for $j = 1$ to d: 3. for $i = 0$ to $k-1$: 4. $C[i] \leftarrow 0$ 5. for $i = 1$ to n: 6. $C[\lfloor A[i]/k^{j-1} \rfloor \bmod k] \leftarrow C[\lfloor A[i]/k^{j-1} \rfloor \bmod k] + 1$ 7. for $i = 1$ to $k-1$: 8. $C[i] \leftarrow C[i] + C[i-1]$ 9. for $i = n$ to 1: 10. $B[C[\lfloor A[i]/k^{j-1} \rfloor \bmod k]] \leftarrow A[i]$ 11. $C[\lfloor A[i]/k^{j-1} \rfloor \bmod k] \leftarrow C[\lfloor A[i]/k^{j-1} \rfloor \bmod k] - 1$ 12. $A[1..n] \leftarrow B[1..n]$</pre>

2.5.2 时间复杂度

对每一位均使用 COUNTING-SORT 进行稳定排序，耗时 $\Theta(n+k)$ ，共进行 d 轮，故总时间复杂度为 $O(d(n+k))$ 。

2.6 桶排序

2.6.1 伪代码

Algorithm 6: Bucket Sort
Input: An array A of n numbers, each of which is in the interval $[0, 1]$. Output: The array A sorted in nondecreasing order.
<pre>1. Let $B[0..n-1]$ be an array of n linked lists. 2. for $i = 0$ to $n-1$: 3. $B[i] \leftarrow \emptyset$ 4. for $j = 1$ to n: 5. Insert $A[j]$ into $B[\lfloor nA[j] \rfloor]$. 6. for $i = 0$ to $n-1$: 7. Sort $B[i]$ using INSERTION-SORT. 8. $A[1..n] \leftarrow B[0] \oplus B[1] \oplus \dots \oplus B[n-1]$</pre>

2.6.2 时间复杂度

假设输入数据服从均匀分布，求期望时间复杂度。桶 $B[i]$ 中的元素位于区间 $[\frac{i}{n}, \frac{i+1}{n})$ 。记

$$X_{ij} = \mathbb{I}\left\{\frac{i}{n} \leq A_j < \frac{i+1}{n}\right\} \sim \text{Bernoulli}\left(\frac{1}{n}\right).$$

由于 $X_{i1}, X_{i2}, \dots, X_{in}$ 相互独立，故 $B[i]$ 的大小 n_i 满足

$$n_i = \sum_{j=1}^n X_{ij} \sim \text{Binomial}\left(n, \frac{1}{n}\right),$$

$$\mathbb{E}[n_i^2] = \mathbb{E}^2[n_i] + \text{Var}[n_i] = 2 - \frac{1}{n}.$$

因此，桶排序在平均情况下的运行时间为

$$\Theta(n) + nO\left(\mathbb{E}\left[n_i^2\right]\right) = O(n^2).$$

最坏情况下，所有元素都落在同一个桶中，此时桶排序退化为 $O(n^2)$ 的插入排序。

3 排序算法运行时间与数组规模的关系

随机生成一个包括 n 个整数的数组，元素取值范围是 $[1, 1000]$ 。改变数组规模 $n = 5 \times 10^4, 1 \times 10^5, 2 \times 10^5, 3 \times 10^5, 5 \times 10^5$ ，利用插入排序，归并排序，快速排序，堆排序，基数排序，桶排序等算法分别对数组进行非降序排序，记录不同规模下各个算法的排序时间，如下表所示。

Table 1. Time cost of different sorting algorithms with different array sizes. For each cell, the first number in parentheses is the mean value of 5 runs, and the second number is the standard deviation.

n	std :: sort	std :: qsort	INSERTION-SORT
50000	$(1.571 \pm 0.063) \times 10^{-3}$	$(2.460 \pm 0.085) \times 10^{-3}$	$(1.619 \pm 0.029) \times 10^0$
100000	$(3.129 \pm 0.039) \times 10^{-3}$	$(4.879 \pm 0.096) \times 10^{-3}$	$(6.364 \pm 0.018) \times 10^0$
150000	$(4.705 \pm 0.086) \times 10^{-3}$	$(7.179 \pm 0.072) \times 10^{-3}$	$(1.436 \pm 0.004) \times 10^1$
200000	$(6.164 \pm 0.138) \times 10^{-3}$	$(9.588 \pm 0.088) \times 10^{-3}$	$(2.563 \pm 0.004) \times 10^1$
250000	$(7.820 \pm 0.193) \times 10^{-3}$	$(1.198 \pm 0.029) \times 10^{-2}$	$(4.006 \pm 0.004) \times 10^1$
300000	$(9.043 \pm 0.310) \times 10^{-3}$	$(1.431 \pm 0.033) \times 10^{-2}$	$(5.771 \pm 0.011) \times 10^1$
350000	$(1.057 \pm 0.032) \times 10^{-2}$	$(1.658 \pm 0.009) \times 10^{-2}$	$(7.849 \pm 0.009) \times 10^1$
400000	$(1.201 \pm 0.016) \times 10^{-2}$	$(1.900 \pm 0.011) \times 10^{-2}$	$(1.026 \pm 0.001) \times 10^2$
450000	$(1.383 \pm 0.023) \times 10^{-2}$	$(2.131 \pm 0.025) \times 10^{-2}$	$(1.302 \pm 0.002) \times 10^2$
500000	$(1.519 \pm 0.015) \times 10^{-2}$	$(2.362 \pm 0.030) \times 10^{-2}$	$(1.610 \pm 0.003) \times 10^2$

n	HEAP-SORT	RADIX-SORT	BUCKET-SORT
50000	$(8.823 \pm 0.190) \times 10^{-3}$	$(5.210 \pm 0.147) \times 10^{-3}$	$(5.262 \pm 0.119) \times 10^{-3}$
100000	$(1.813 \pm 0.025) \times 10^{-2}$	$(1.365 \pm 0.008) \times 10^{-2}$	$(1.351 \pm 0.008) \times 10^{-2}$
150000	$(2.754 \pm 0.024) \times 10^{-2}$	$(2.590 \pm 0.057) \times 10^{-2}$	$(2.571 \pm 0.039) \times 10^{-2}$
200000	$(3.607 \pm 0.075) \times 10^{-2}$	$(4.044 \pm 0.085) \times 10^{-2}$	$(4.125 \pm 0.042) \times 10^{-2}$
250000	$(4.485 \pm 0.071) \times 10^{-2}$	$(5.938 \pm 0.035) \times 10^{-2}$	$(5.994 \pm 0.046) \times 10^{-2}$
300000	$(5.489 \pm 0.076) \times 10^{-2}$	$(8.154 \pm 0.214) \times 10^{-2}$	$(8.084 \pm 0.174) \times 10^{-2}$
350000	$(6.345 \pm 0.054) \times 10^{-2}$	$(1.057 \pm 0.013) \times 10^{-1}$	$(1.043 \pm 0.017) \times 10^{-1}$
400000	$(7.244 \pm 0.072) \times 10^{-2}$	$(1.342 \pm 0.004) \times 10^{-1}$	$(1.026 \pm 0.004) \times 10^{-1}$
450000	$(8.116 \pm 0.179) \times 10^{-2}$	$(1.655 \pm 0.031) \times 10^{-1}$	$(1.619 \pm 0.010) \times 10^{-1}$
500000	$(8.936 \pm 0.074) \times 10^{-2}$	$(1.977 \pm 0.010) \times 10^{-1}$	$(1.957 \pm 0.009) \times 10^{-1}$

在对数坐标系中，绘制不同算法数组规模与运行时间的折线图，如下图所示。

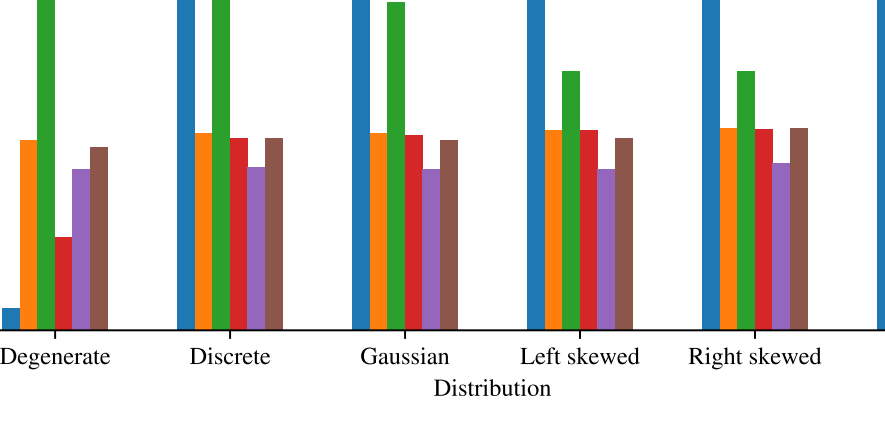


Figure 1. Time cost of different sorting algorithms with different array sizes.

- 由图 1 可以看出，平均情况下，各算法 $\log T(n)$ 随 $\log n$ 增长的速率：
- 根据前文理论分析，插入排序的期望运行时间为 $O(n^2)$ ，基数排序（ d 与 k 固定）及桶排序的期望运行时间为 $O(n)$ ，故图 1 中插入排序对应斜率为 2，基数排序及桶排序对应斜率为 1。
 - Gaussian：从均值 500，标准差为 5 的正态分布中随机抽取并四舍五入为整数，丢弃 $[1, 1000]$ 范围外的元素。
- 下面，针对每种排序算法，检验其时间复杂度。对于 INSERTION-SORT，以 n^2 为横坐标，绘制数组规模与运行时间的折线图，如下图所示。

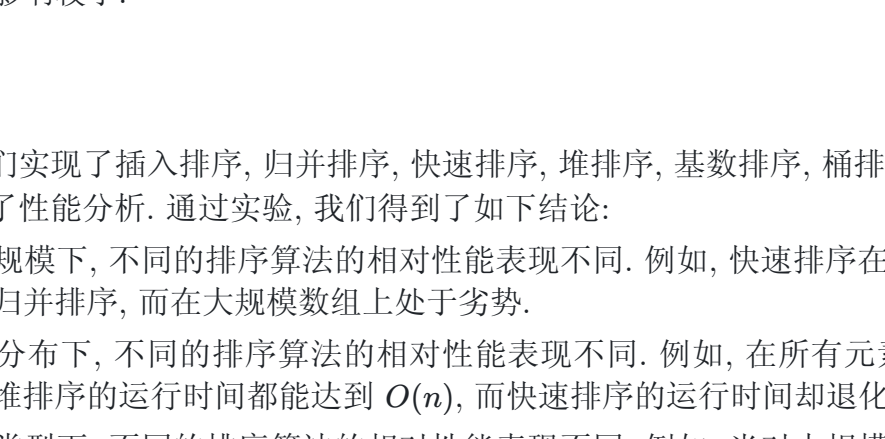


Figure 2. Time cost of INSERTION-SORT with different array sizes.

可见 $T(n)$ 与 n^2 成线性关系，验证了插入排序的平均时间复杂度为 $O(n^2)$ 。

对于 MERGE-SORT, QUICK-SORT 和 HEAP-SORT，以 $n \log n$ 为横坐标，绘制数组规模与运行时间的折线图，如下图所示。

Figure 3. Time cost of MERGE-SORT, QUICK-SORT and HEAP-SORT with different array sizes.

可见，当 n 充分大时，三种算法的 $T(n)$ 均与 $n \log n$ 大致成线性关系，验证了这三种排序算法的平均时间复杂度为 $O(n \log n)$ 。有如下结论：

- 当 n 较小时，快速排序的运行时间最短，优于归并排序和堆排序，可能有以下原因：
 - QUICK-SORT 存在尾递归，相较 MERGE-SORT 递归开销较小。
 - Gaussian：从均值 500，标准差为 5 的正态分布中随机抽取并四舍五入为整数，丢弃 $[1, 1000]$ 范围外的元素。
 - QUICK-SORT 总是对连续的子数组进行排序，Cache 命中率较高。
- 当 n 较大时，快速排序的运行时间最长，增长速度显著高于归并排序和堆排序。这说明，相较另外两种算法，QUICK-SORT 所需比较、交换等操作的次数更多，比值 $\frac{T(n)}{n \log n}$ 趋向于更大的常数。

对于 RADIX-SORT 和 BUCKET-SORT，以 n 为横坐标，绘制数组规模与运行时间的折线图，如下图所示。

Figure 4. Time cost of RADIX-SORT and BUCKET-SORT with different array sizes.

可见，平均情况下，两种算法的 $T(n)$ 均与 n 大致成线性关系，符合其时间复杂度。然而，由于实现方式的不同，两种算法的运行时间存在较大差异。RADIX-SORT 针对整数实现 $\Theta(d(n+k))$ 的排序，在所有算法中总是具有最短的运行时间。

4 排序算法运行时间与数据分布的关系

从不同的数据分布中抽样，生成一个包括 $n = 10^5$ 个整数的数组，元素取值范围是 $[1, 1000]$ 。针对每种数据分布，随机扰乱相应的数组并排序，统计六种排序算法的运行时间。现针对以下六种数据分布进行实验：

- Degenerate：退化分布，所有元素均为 500。
- Discrete：从集合 $\{100, 200, 300, 400, 500, 600, 700, 800, 900, 1000\}$ 中随机抽取。
- Gaussian：从均值为 500，标准差为 5 的正态分布中随机抽取并四舍五入为整数，丢弃 $[1, 1000]$ 范围外的元素。
- Left skewed：从均值为 50 的指数分布中随机抽取，向下取整后加 1，丢弃大于 1000 的元素。该分布具有左偏的特点。
- Right skewed：从均值为 50 的指数分布中随机抽取，向下取整后从 1000 中减去，丢弃小于 1 的元素。该分布具有右偏的特点。
- Uniform：随机抽取区间 $[1, 1000]$ 中的整数。

抽样所得数组元素的频数分布如下图所示。

Figure 5. Histogram of the data distributions.

实验得到每种排序算法的运行时间如下表所示。

Table 2. Time cost of different sorting algorithms with different data distributions.

Distribution	std :: sort	std :: qsort	INSERTION-SORT
Degenerate	$(6.694 \pm 0.031) \times 10^{-5}$	$(1.711 \pm 0.036) \times 10^{-4}$	$(2.345 \pm 0.043) \times 10^{-4}$
Discrete	$(1.027 \pm 0.142) \times 10^{-3}$	$(1.497 \pm 0.096) \times 10^{-3}$	$(5.819 \pm 0.026) \times 10^0$
Gaussian	$(1.347 \pm 0.088) \times 10^{-3}$	$(2.018 \pm 0.042) \times 10^{-3}$	$(6.269 \pm 0.104) \times 10^0$
Left skewed	$(2.212 \pm 0.057) \times 10^{-3}$	$(3.322 \pm 0.096) \times 10^{-3}$	$(6.399 \pm 0.013) \times 10^0$
Right skewed	$(2.364 \pm 0.387) \times 10^{-3}$	$(3.384 \pm 0.194) \times 10^{-3}$	$(6.486 \pm 0.200) \times 10^0$
Uniform	$(3.101 \pm 0.072) \times 10^{-3}$	$(4.841 \pm 0.139) \times 10^{-3}$	$(6.465 \pm 0.054) \times 10^0$

Distribution	MERGE-SORT	QUICK-SORT	RANDOMIZED-QUICK-SORT
Degenerate	$(1.304 \pm 0.031) \times 10^{-3}$	$(6.488 \pm 0.119) \times 10^0$	$(6.437 \pm 0.050) \times 10^0$
Discrete	$(1.542 \pm 0.036) \times 10^{-2}$	$(6.446 \pm 0.027) \times 10^{-1}$	$(6.426 \pm 0.026) \times 10^{-1}$
Gaussian	$(1.559 \pm 0.042) \times 10^{-2}$	$(3.721 \pm 0.060) \times 10^{-1}$	$(3.788 \pm 0.138) \times 10^{-1}$
Left skewed	$(1.700 \pm 0.030) \times 10^{-2}$	$(7.161 \pm 0.109) \times 10^{-2}$	$(7.115 \pm 0.092) \times 10^{-2}$
Right skewed	$(1.712 \pm 0.041) \times 10^{-2}$	$(7.551 \pm 0.882) \times 10^{-2}$	$(7.731 \pm 1.308) \times 10^{-2}$
Uniform	$(1.762 \pm 0.036) \times 10^{-2}$	$(1.398 \pm 0.021) \times 10^{-2}$	$(1.394 \pm 0.034) \times 10^{-2}$

绘制得不同算法，不同数据分布下的运行时间柱状图如下。其中，纵轴采用对数坐标。

Figure 6. Time cost of