

Università degli Studi di Udine

Scienze Informatiche

Laboratorio di Algoritmi e Strutture Dati

Anno: 2019/2020

Riccardo Ferrarese

mat. 140491

ferrarese.riccardo@spes.uniud.it

Contents

1	Introduzione	2
2	Quick Select	2
	Algoritmo	2
	Pseudocodice	2
	Complessità e analisi dei tempi	3
	Grafici	4
3	Heap Select	7
	Algoritmo	7
	Pseudocodice	7
	Complessità e analisi dei tempi	7
	Grafici	8
4	Median of Median Select	12
	Algoritmo	12
	Pseudocodice	12
	Complessità e analisi dei tempi	14
	Grafici	15

1 Introduzione

Il seguente elaborato propone tre diverse soluzioni per il problema del calcolo del k -esimo elemento più piccolo in un vettore non ordinato di interi.

Le implementazioni sono state realizzate in linguaggio *C++*. Grafici, varianza e tempi medi totali sono stati calcolati con l'ausilio di codice scritto in *R*.

Nelle seguenti sezioni verranno presentate ed analizzate le soluzioni.

L'analisi dei tempi cerca di ridurre al minimo le variazioni nei tempi calcolati facendo eseguire più volte l'algoritmo su un input fissato per raggiungere una determinata precisione nella misurazione.

2 Quick Select

Algoritmo

Si tratta di una variante dell'algoritmo di ordinamento *quickSort*, in cui ogni chiamata ricorsiva su un intervallo $[i,j]$ del vettore fornito in input termina in tempo costante ogniqualvolta il parametro k non sia contenuto nell'intervallo $[i,j]$.

Dato in input un vettore di numeri *numbers* che va da $[1..dim]$ e la sua dimensione, la procedura **quick_select** utilizza **partition** (con pivot casuale) per partizionare il vettore attorno al pivot.

Così facendo se il pivot si trova ad un indice più grande di k , si continua la ricerca nella partizione di destra.

Viceversa se k è più grande dell'indice del pivot si continua la ricerca nella partizione di sinistra.

Pseudocodice

Algorithm 1 Quick Select - first solution

```
1: function QUICK_SELECT(numbers, start, end, k)
2:   if  $k \leq start$  and  $k \geq end$  then
3:      $pivot\_index \leftarrow partition(numbers, dim)$ 
4:     if  $pivot\_index == k$  then
5:       return numbers[k]
6:     else if  $pivot\_index > k$  then
7:       return quick_select(numbers, start,  $pivot\_index-1$ , k)
8:     else return quick_select(numbers, start,  $pivot\_index-1$ , k)
9:   end if
10: end if
11: return null
12: end function
```

Di seguito lo pseudocodice della procedura **Partition**, il quale preso un

vettore di numeri interi e un pivot randomico (in questo caso l'ultimo elemento del vettore) partiziona gli elementi più piccoli del pivot a sinistra e quelli più grandi a destra.

Come dimostrato a lezione *Partition* ha complessità $\mathcal{O}(n)$.

Algorithm 2 Partition_random

```

1: function PARTITION(numbers, start, end)
2:   pivot  $\leftarrow$  numbers[end]
3:   for j  $\leftarrow$  start, j  $\leq$  end, j ++ do
4:     if numbers[j]  $\leq$  pivot then
5:       i  $\leftarrow$  i + 1
6:       swap(numbers[i], numbers[j])
7:     end if
8:   end for
9:   swap(numbers[i + 1], numbers[end])
10:  return i + 1                                      $\triangleright$  posizione del pivot
11: end function

```

Complessità e analisi dei tempi

Definiamo l'equazione di complessità dell'algoritmo **quick_select**.

Sia $m \in [1..n-1]$ il numero di numeri su cui svolgo la ricorsione, sia $[i, j]$ la porzione di array in cui mi sto muovendo:

$$T(n) = \begin{cases} \mathcal{O}(1) & \text{se } k \notin [i, j], \\ T(m) + \mathcal{O}(n) & \text{con } m = j - i \text{ e } k \in [i, j] \end{cases}$$

La complessità dipende dalla dimensione di m , che nel caso peggiore corrisponde a $n-1$, cioè quando il pivot cade come primo o ultimo elemento ed effettuo la ricorsione sul resto dell'array.

Se ad ogni chiamata ricorsiva ci ritroviamo nel caso descritto precedentemente, al massimo verranno effettuate n chiamate su array di dimensione $n-1$.

Quindi l'algoritmo nel caso peggiore risulta avere complessità pari a:

$$\begin{aligned}
T(n) &= T(n-1) + \mathcal{O}(n) \\
&= \sum_{i=1}^n \mathcal{O}(i) \\
&= \mathcal{O}(n^2)
\end{aligned} \tag{1}$$

Nel caso medio, il pivot ha valore tale da dividere quasi equamente l'array in parti uguali. L'equazione di complessità diventa quindi:

$$\begin{aligned}
T(n) &= T\left(\frac{n}{2}\right) + \theta(n) && \text{sia } n = 2^k \\
T(2^k) &= T(2^{k-1}) + \theta(2^k) \\
G(k) &= G(k-1) + \theta(2^k) \\
\sum_{i=0}^{\log(n)} \mathcal{O}(2^i) &= \mathcal{O}\left(\sum_{i=0}^{\log(n)} 2^i\right) = \theta(2^i)
\end{aligned} \tag{2}$$

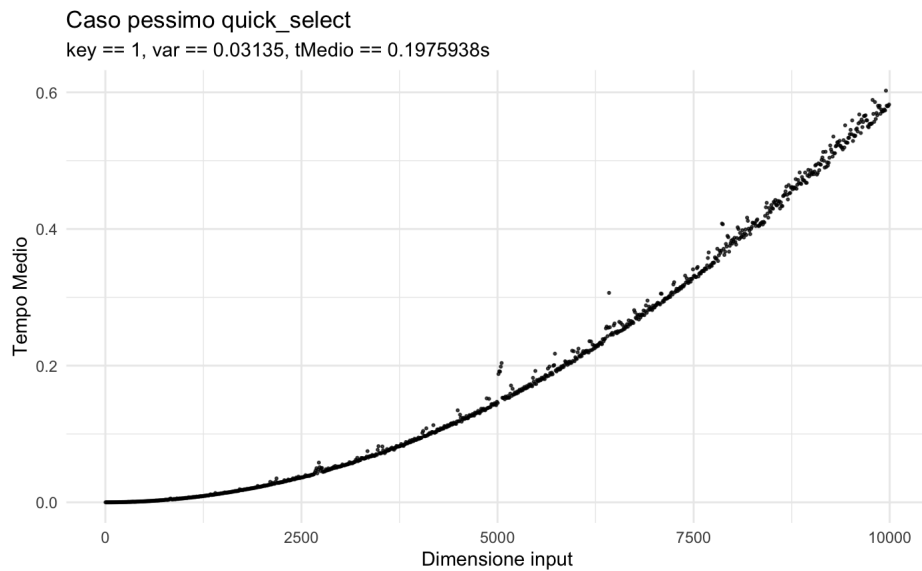
$$T(2^k) = \theta(2^i) \iff T(n) = \theta(n) \tag{3}$$

Grafici

Si mostra l'andamento dei tempi al variare dell'input, su dimensione crescente, e al variare di k .

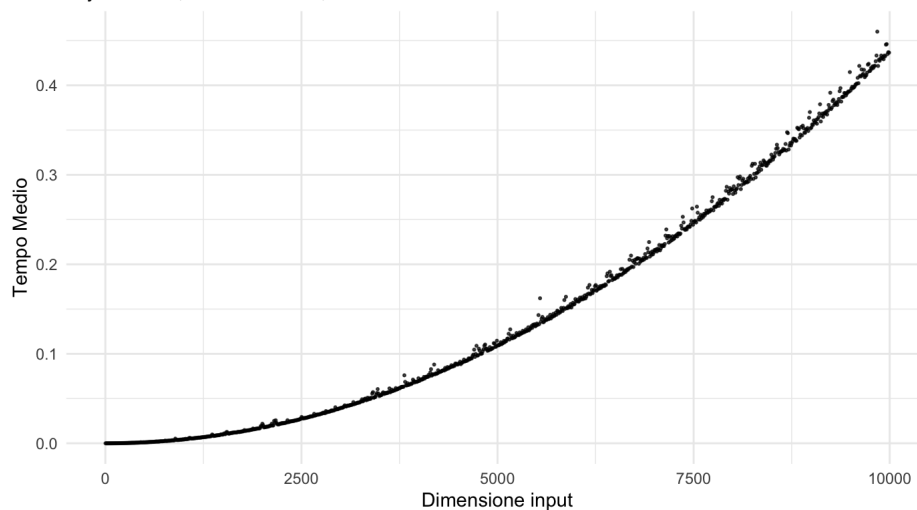
Per il caso pessimo si è scelto di generare delle stringhe di input in cui tutti i valori sono uguali, in modo da partizionare l'array sempre in maniera sbilanciata ed effettuare la ricorsione su $n - 1$ elementi.

Si danno due esempi per k diverso.



Caso pessimo quick_select

key == dim/2, var == 0.01729, tMedio == 0.1470228s

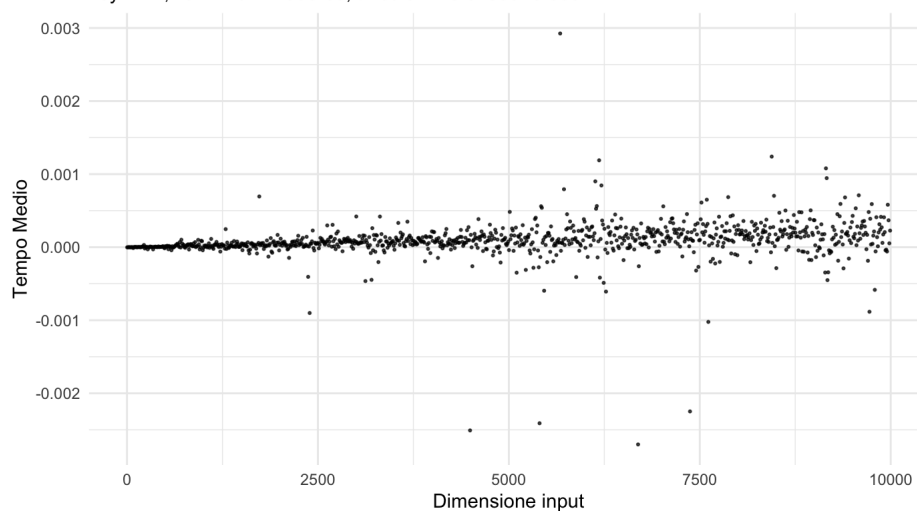


Come da ipotesi il tempo medio ha crescita quadratica rispetto alla dimensione dell'input.

Di seguito si mostrano dei casi medi, in cui i valori dell'input sono scelti randomicamente. Nel misurare i tempi si riscontrano dei valori negativi, si ipotizza a causa della scorporazione dal tempo di esecuzione della preparazione dell'input, che potrebbe impiegare un tempo maggiore rispetto al tempo totale nel caso in cui il valore da trovare viene trovato subito.

Caso medio quick_select

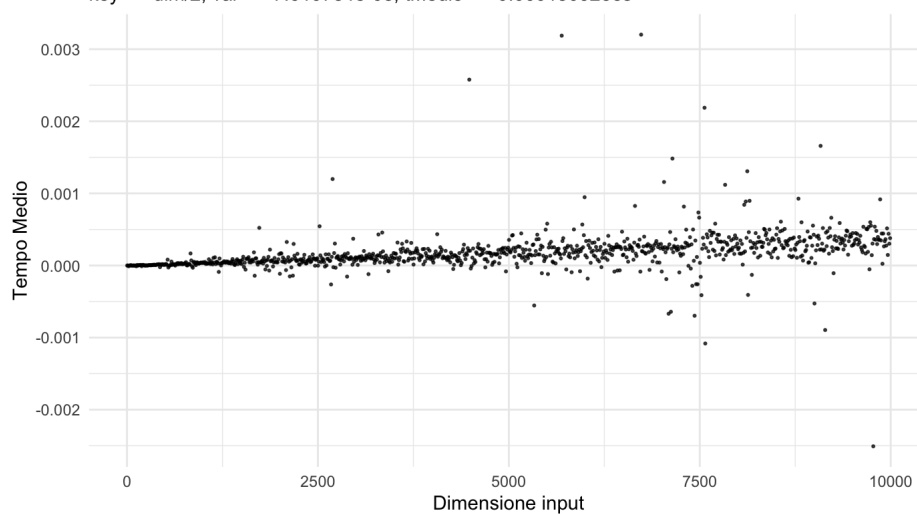
key == 1, var == 6.77295e-08, tMedio == 9.329392e-05s



Già da questo grafico, oltre a vedere che il tempo medio ha crescita più o meno lineare rispetto la dimensione dell'input, si può notare come i casi non siano uniformi, questo a causa della scelta randomica del pivot per la partizione e della distribuzione dei valori nel vettore. Si danno altri esempi al variare di k .

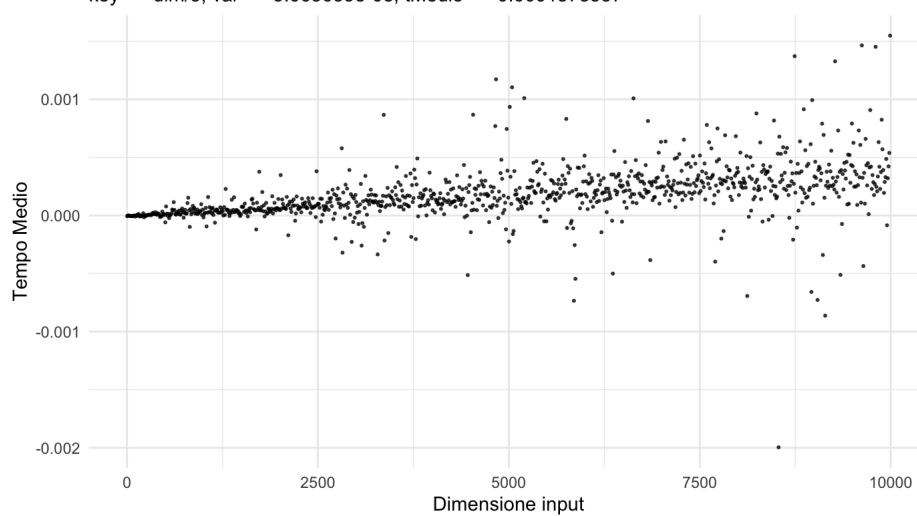
Caso medio quick_select

key == dim/2, var == 7.610784e-08, tMedio == 0.0001809288s



Caso medio quick_select

key == dim/3, var == 5.663655e-08, tMedio == 0.0001873337



3 Heap Select

Algoritmo

Questo algoritmo di selezione risolve il problema del k -esimo elemento più piccolo di una lista non ordinata di numeri utilizzando la struttura dati Heap. Se il valore di k è minore della metà, scegliamo di utilizzare la min-heap, viceversa, se è più grande della metà, possiamo sfruttare la max-heap per riportarci al problema duale, cioè cercando il $(dim - k + 1)$ -esimo elemento più grande.

La prima heap H1 é costruita a partire dal vettore fornito in input in tempo lineare e non viene modificata. La seconda heap H2 contiene solamente indici di valori corrispondenti alla prima heap. Inizialmente H2 contiene un solo nodo, corrispondente alla radice di H1. All' i -esima iterazione, per i che va da 1 a $k-1$ (con k opportuno se scelgo di utilizzare la max-heap) l'algoritmo estrae la radice di H2, corrispondente a un nodo x_i in H1 e reinserisce in H2 gli indici del figlio destro e sinistro di x_i (se esistono) nella heap H1. Dopo $k-1$ iterazioni, la radice di H2 corrisponderà all'indice del k -esimo elemento più piccolo del vettore fornito in input.

Pseudocodice

Algorithm 3 Heap Select - second solution

```
1: function HEAP_SELECT( $numbers, start, end, k$ )
2:   if  $k \geq \frac{dim}{2}$  then
3:     return  $max\_heap\_select(numbers, start, end, k)$ 
4:   else
5:     return  $min\_heap\_select(numbers, start, end, k)$ 
6:   end if
7:   return 1
8: end function
```

Di seguito lo pseudocodice della procedura **max_heap_select**, analoga sarà la procedura **min_heap_select** con valori di iterazione $[1..k-1]$ fatta con il k preso in input.

Per implementare la precedente procedura si è scelto di utilizzare una struttura per la **minHeap** e una per la **maxHeap**, inoltre vengono implementate due strutture di heap ausiliare (max e min) in cui i nodi hanno come valore gli indici che corrispondono ai valori dei nodi nella heap iniziale. Il metodo $successoriDi(x)$ inserisce nella seconda heap gli indici del figlio sinistro e destro del nodo avente valore x .

Algorithm 4 Heap Select - second solution

```
1: function MAX_HEAP_SELECT(numbers, dim, k)
2:   MaxHeap_1  $\leftarrow$  numbers
3:   MaxHeap_2  $\leftarrow$  MaxHeap_1.getMax()
4:   i  $\leftarrow$  0
5:   new_k = HeapSize(MaxHeap_1) - k + 1
6:   while i  $\leq$  new_k - 1 do
7:     x  $\leftarrow$  MaxHeap_2.extractMax()
8:     MaxHeap_2  $\leftarrow$  successoriDi(x)     $\triangleright$  Aggiungo a 2 i successori di x
9:     i ++
10:  end while
11:  return MaxHeap_1[MaxHeap_2.getMax()]
12: end function
```

Complessità e analisi dei tempi

Definiamo l'equazione di complessità dell'algoritmo **heap_select_min**, analogamente si può procedere per la soluzione con le *maxHeap*.

Sappiamo che la complessità di *Build_Heap* è pari a $\mathcal{O}(n)$. Il ciclo while viene eseguito k volte e il corpo ha complessità pari alla complessità di *extractMin* sulla heap ausiliaria. Quest'ultima conterrà al massimo k elementi essendo che ad ogni iterazione viene rimosso un elemento e successivamente ne vengono inseriti due.

Facciamo un esempio per $k = 4$, con $k - 2$ iterazioni del corpo del while:

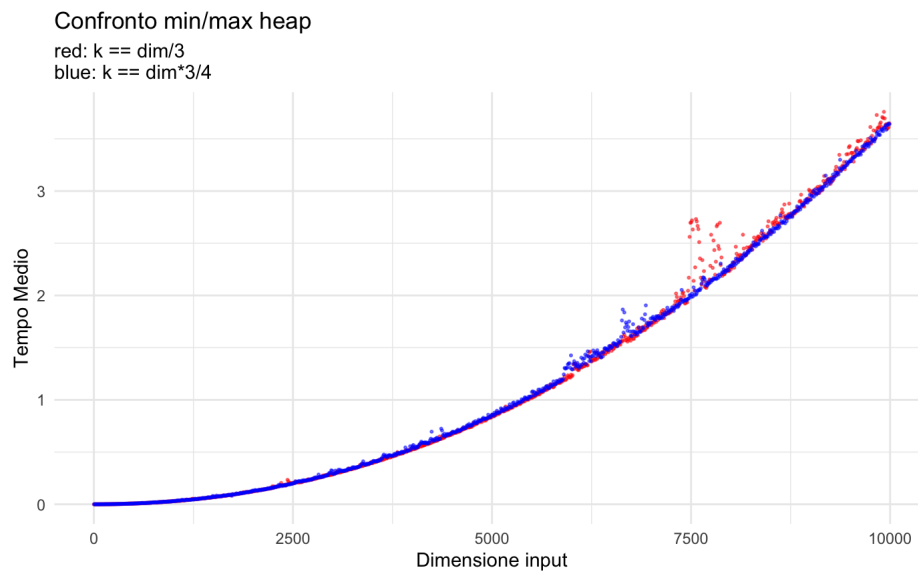
- $i = 0 \rightarrow$ un elemento iniziale, due finali.
- $i = 1 \rightarrow$ due elemento iniziale, tre finali.
- $i = 2 \rightarrow$ tre elemento iniziale, quattro finali.

L'equazione di complessità risulta quindi:

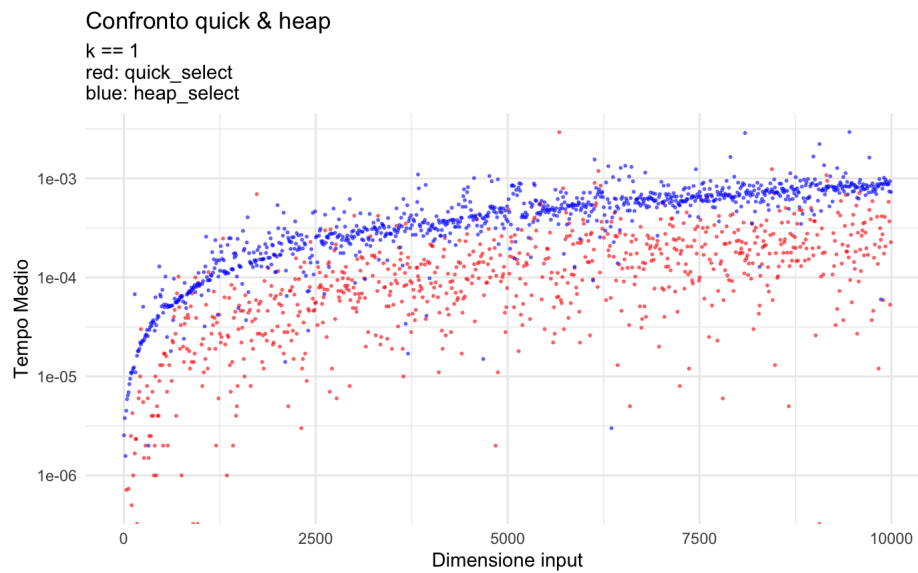
$$\begin{aligned} T(n) &= \mathcal{O}(n) + \mathcal{O}(k \log(k)) \\ &= \mathcal{O}(n + k \log(k)) \end{aligned}$$

Grafici

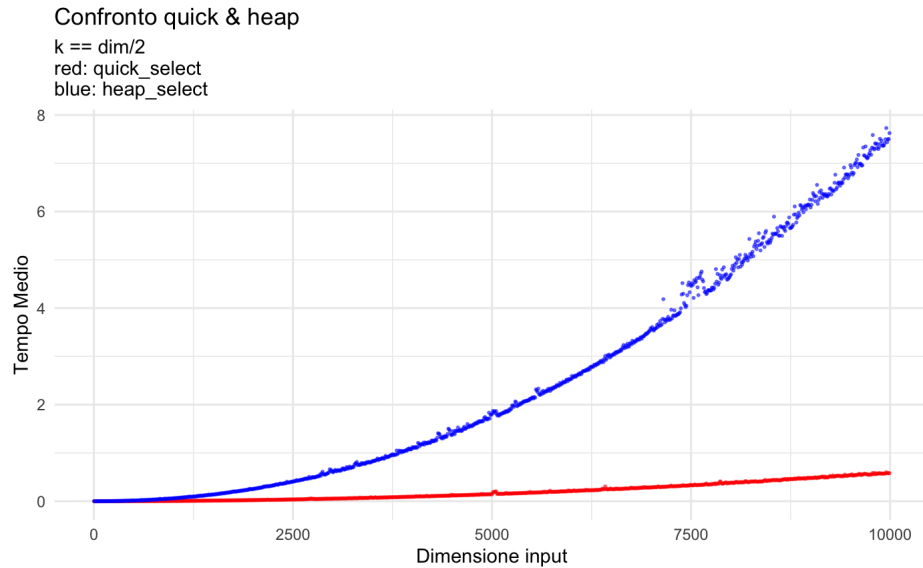
Di seguito viene proposto un grafico comparativo per k rispettivamente nella prima e seconda metà del vettore. Si mostra che gli andamenti dei tempi di esecuzione corrispondono, avendo utilizzato nel primo caso una *min-heap* e nel secondo una *max-heap*.



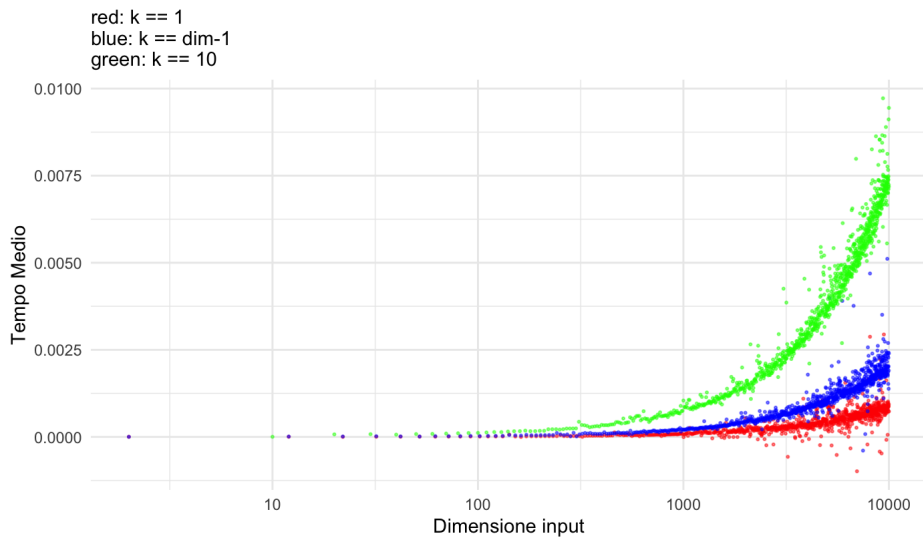
Il seguente grafico mostra come in caso di k molto piccolo (o rispettivamente molto grande) è preferibile l'algoritmo *heap_select* rispetto a *quick_select* avendo tempo medio più uniforme, non dipendendo dalla disposizione degli elementi nel vettore. L'asse y è in scala logaritmica per visualizzare meglio la dispersione dei dati rispetto il tempo.

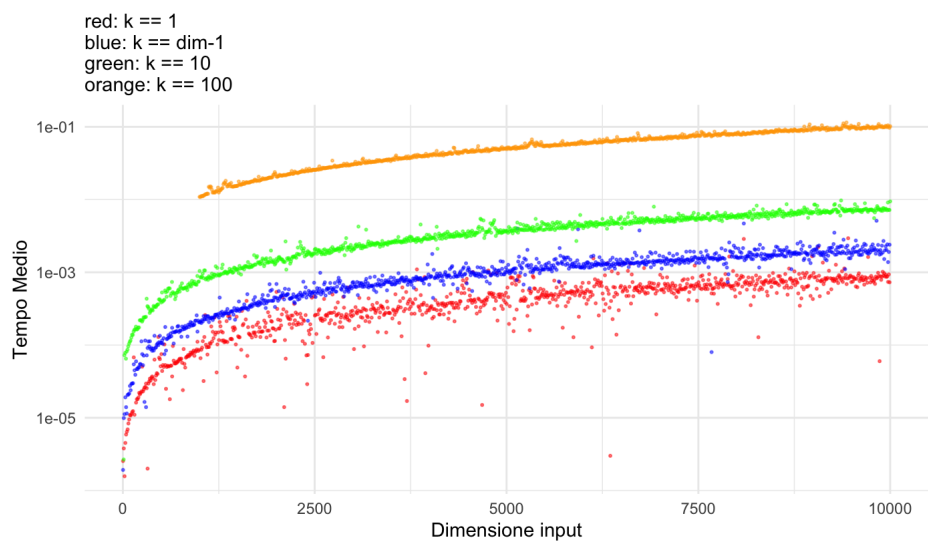
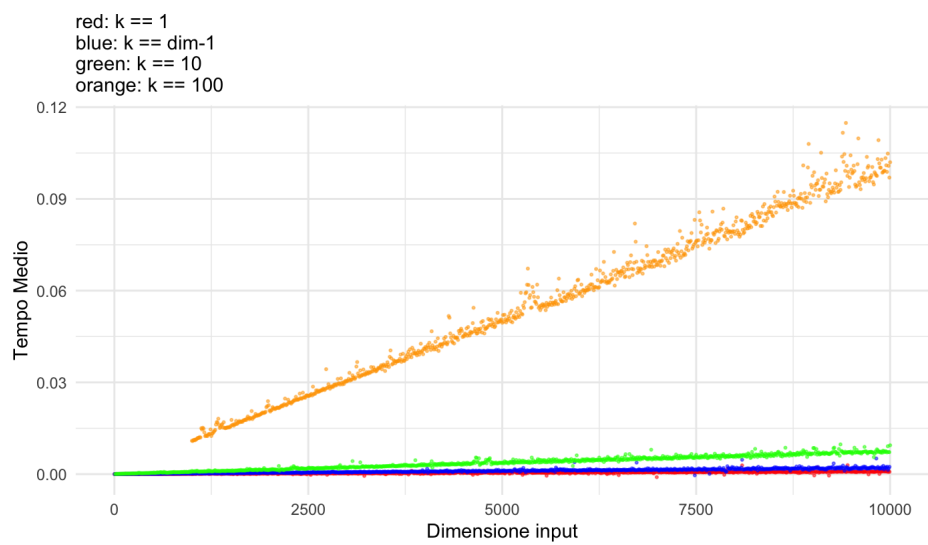


Considerando invece il caso di un k centrale, l'algoritmo *heap_select* risulta pessimo rispetto il worst case di *quick_select*.



Il seguente grafico, con asse delle x in scala logaritmica, mostra come per vettori di dimensione relativamente ridotta non si ha grande differenza di tempo rispetto alla scelta di k . Per vettori di dimensione abbastanza grande la differenza tra $k = 10$ e $k = 1$ risulta importante per il tempo di esecuzione.





4 Median of Median Select

Algoritmo

Questo algoritmo ci permette di trovare la mediana dei mediani ed utilizzarla come pivot nella partizione degli elementi, in modo da avere il pivot compreso tra $\frac{1}{4}$ e $\frac{3}{4}$ dell'array. Più precisamente, l'algoritmo esegue le seguenti operazioni:

- divisione dell'array in blocchi di 5 elementi, escluso eventualmente l'ultimo blocco che potrà contenere meno di 5 elementi;
- ordinamento e calcolo della mediana di ciascun blocco;
- calcolo della mediana M delle mediane dei blocchi, attraverso chiamata ricorsiva allo stesso algoritmo;
- partizionamento dell'intero array attorno alla mediana M, attraverso una variante della procedura *Partition* dell'algoritmo *quick select*;
- chiamata ricorsiva nella parte di array che sta a sinistra o a destra della mediana M, in funzione del valore k fornito in input.

L'algoritmo è in place muovendosi nell'array attraverso gli indici, *InsertionSort* copia la porzione da ordinare per non modificare l'array su cui successivamente fare partition. Supponendo n numero di elementi nel vettore, il numero di gruppetti che si verranno a formare è uguale a $\frac{n+4}{5}$, aggiungiamo 4 per far sì che il numero minimo di gruppetti sia 1 e non 0.

Pseudocodice

La funzione ha come argomenti l'array di numeri interi, il valore k , $start$ e end che sono gli indici di inizio e fine della porzione di array su cui si esegue la chiamata ricorsiva.

La funzione *find_median* copia la porzione di array su cui si esegue la chiamata di funzione (per non modificare l'array iniziale) e la ordina per trovare l'elemento mediano.

In corrispondenza della linea 11 dell'algoritmo 5, se la dimensione dell'array dei mediani è minore o uguale a 5, si proceda direttamente a trovare l'elemento mediano, altrimenti si faccia una chiamata ricorsiva alla procedura *Median of Median Select* con argomenti l'array dei mediani e k con $dim/2$.

Algorithm 5 Median of Median Select - third solution

```
1: function MEDIAN_OF_MEDIAN_SELECT(numbers, start, end, k)
2:   dim = end - start + 1
3:   if k > 0 and k ≤ dim then
4:     mediani $[\frac{dim+4}{5}]$ 
5:     for i = 0 to  $\frac{dim}{5}$  do
6:       mediani[i] = findMedian(numbers, start + i * 5, left + i * 5 + 4)
7:     end for
8:     if (i * 5 < dim) then ▷ Per gli elementi fuori dai gruppetti
9:       mediani[i] = findMedian(numbers, start + i * 5, dim%5)
10:    end if
11:    if (i ≤ 5) then
12:      mediano = findMedian(numbers, 0, i - 1)
13:    else
14:      mediano = median_of_median_select(median, 0, i - 1, i/2)
15:    end if
16:    pos = partition(numbers, start, end, medianOfmedian)
17:    if (pos - left == k - 1) then
18:      return numbers[pos]
19:    else if (pos - left > k - 1) then
20:      return median_of_median_select(numbers, start, pos - 1, k)
21:    else
22:      new_k ← k - pos + left - 1
23:      return median_of_median_select(numbers, pos + 1, end, new_k)
24:    end if
25:  end if
26: end function
```

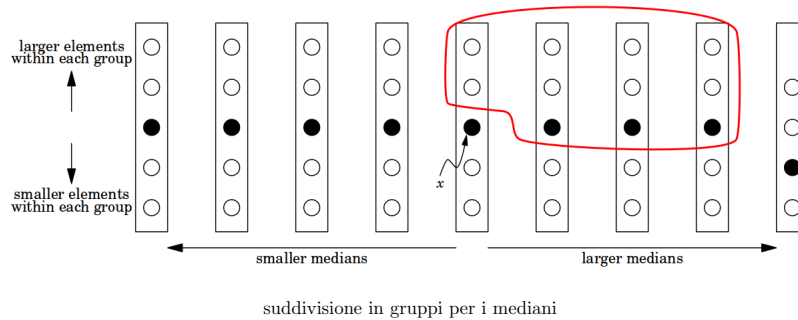
Algorithm 6 Partition

```
1: function PARTITION(numbers, start, end, medianOfmedian)
2:   i ← start
3:   while numbers[i] != medianOfmedian do ▷ cerco indice del pivot
4:     i ++
5:   end while
6:   swap(numbers[i], numbers[end]) ▷ spostato il pivot alla fine
7:   for j ← start to end do
8:     if numbers[j] ≤ medianOfmedian then
9:       i ← i + 1
10:      swap(numbers[i], numbers[j])
11:    end if
12:  end for
13:  swap(numbers[i + 1], numbers[end])
14:  return i + 1 ▷ posizione del pivot
15: end function
```

Complessità e analisi dei tempi

Definiamo l'equazione di complessità dell'algoritmo **Median of Median Select**. Di seguito si elencano le complessità dei passi svolti durante l'esecuzione:

- $\mathcal{O}(n)$ divisione degli elementi in gruppi di 5.
- $\mathcal{O}(n)$ ordinamento e selezione dell'elemento mediano di ogni gruppo.
- $T(\frac{n}{5})$ esecuzione della chiamata ricorsiva sulla lista degli elementi medi ottenuti al passo precedente.
- $\mathcal{O}(n)$ costo di partition su tutto l'array in input utilizzando come pivot l'elemento ottenuto nel passo precedente.
- $T(\frac{3}{4}n)$ ricorsione su uno dei sottoarray in base al valore di k .



Per analizzare il tempo di esecuzione dell'algoritmo dobbiamo prima determinare quanti elementi sono più grandi del pivot. Chiamiamo p il pivot trovato. Almeno la metà delle mediane trovate sono maggiori della mediana delle mediane e quindi almeno metà degli $\lfloor \frac{n}{5} \rfloor$ gruppi contribuiscono con 3 elementi che sono più grandi di x , tranne il gruppo di x stesso che ha meno elementi e quindi si hanno almeno:

$$3(\lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil - 2) \geq \frac{3}{10}n - 6$$

elementi più grandi di x .

Un ragionamento speculare può esser fatto per il numero di elementi minori di x che quindi sono $\geq \frac{3}{10}n - 6$.

Nel caso peggiore quindi la chiamata verrà effettuata su di un array di $\frac{7}{10}n + 6$ elementi. L'equazione di complessità risulta quindi:

$$T(n) = \begin{cases} \theta(1) & \text{con } n \leq c, \\ T(\lceil \frac{n}{5} \rceil) + T(\frac{7}{10}n + 6) + \theta(n) & \end{cases}$$

Dimostriamo che per qualche costante $cT(n) \leq cn$

$$\begin{aligned}
T(n) &\leq c \left\lceil \frac{n}{5} \right\rceil + c \left\lceil \frac{7}{10} + 6 \right\rceil + \mathcal{O}(n) \\
&\leq c \frac{n}{5} + c + \frac{7}{10}cn + 6c + \mathcal{O}(n) \\
&\leq \frac{9}{10} + 7c + \mathcal{O}(n) \\
&\leq cn
\end{aligned}$$

Quindi abbiamo:

$$c \frac{9}{10}n + 7c + kn \leq cn \iff kn \leq c(\frac{1}{10}n - 7) \iff k \leq c(\frac{1}{10} - \frac{7}{n}) \iff c \geq \frac{k}{(\frac{1}{10} - \frac{7}{n})}$$

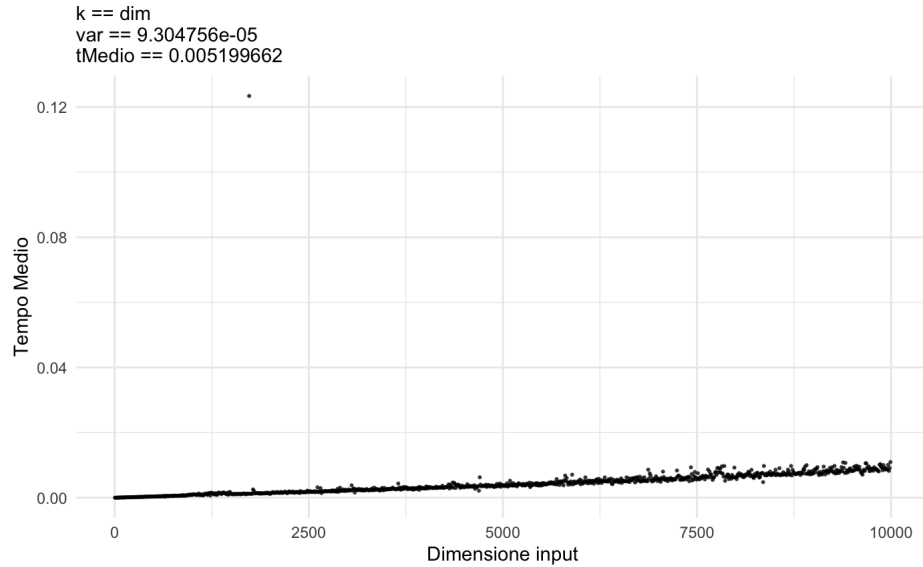
Se prendiamo $n \geq 80$ $c = \frac{k}{\frac{1}{10} - \frac{7}{80}}$ allora c è sicuramente più grande.

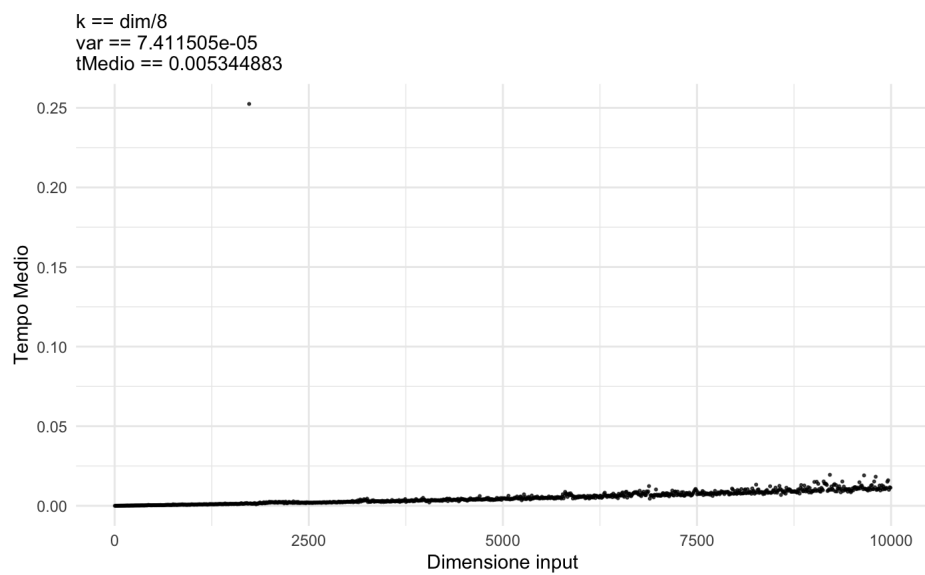
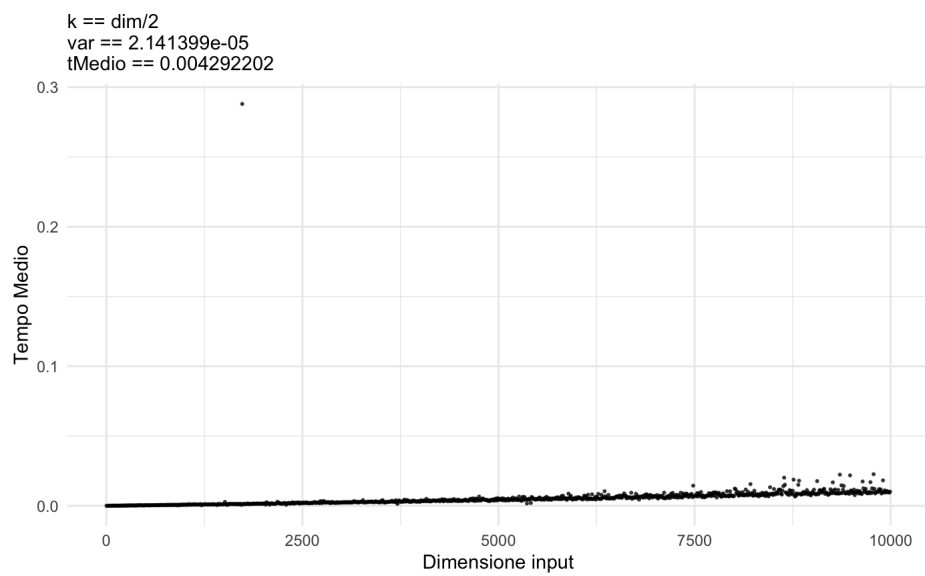
Concludiamo che la procedura **Median of median Select** sia nel caso pessimo che in quello medio ha complessità $\mathcal{O}(n)$.

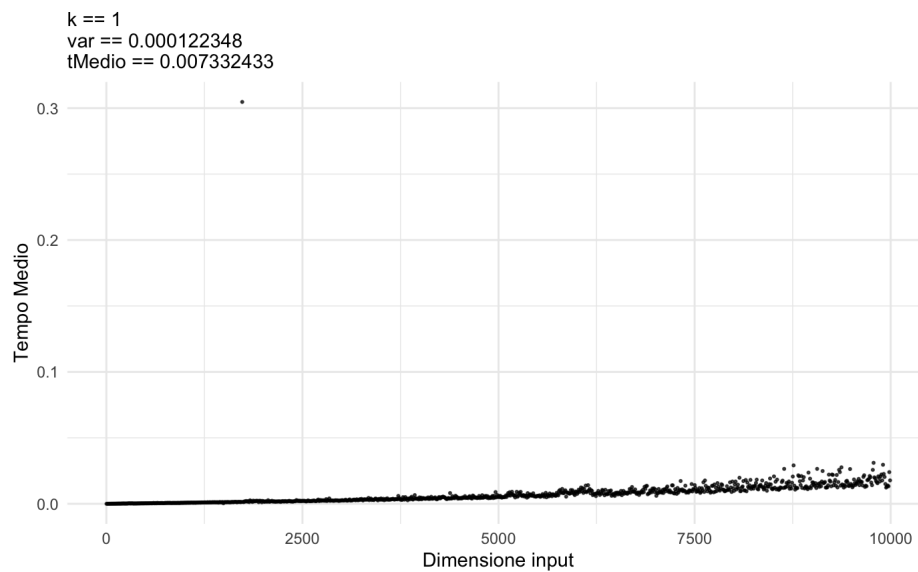
Il grafico dei tempi risultante conferma l'analisi teorica fatta, risultando la soluzione più efficiente per il problema proposto.

Grafici

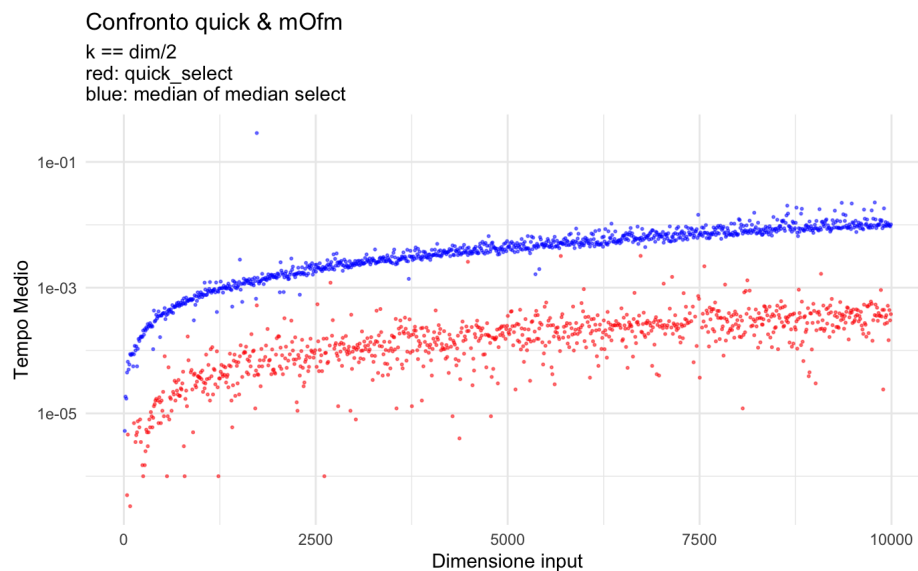
Si mostra al variare di k i grafici ottenuti dall'analisi del tempo di esecuzione. Si può notare come per qualsiasi k scelto, i tempi al variare della dimensione crescono in maniera lineare e uniforme rispetto a input con valori diversi, rispecchiando l'analisi teorica con di complessità $\mathcal{O}(n)$ in tutti i casi.



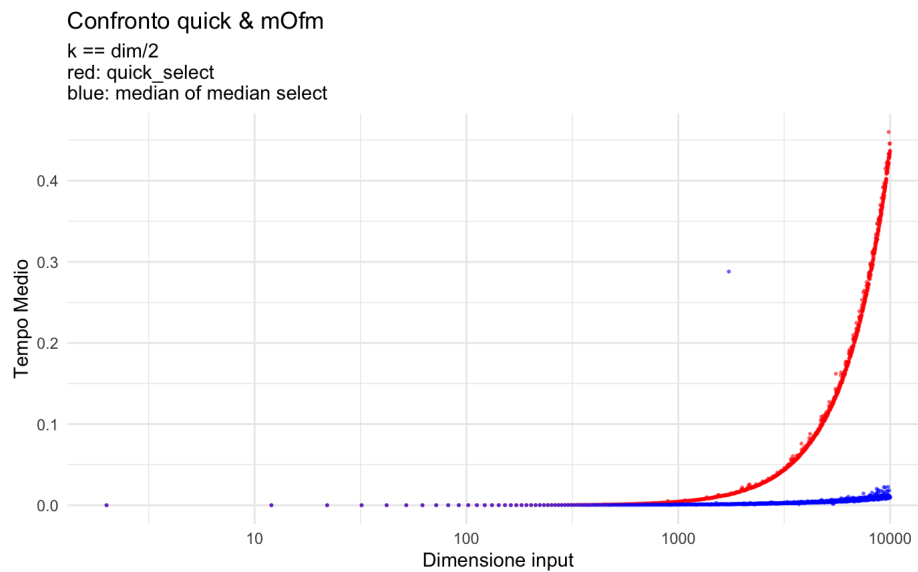




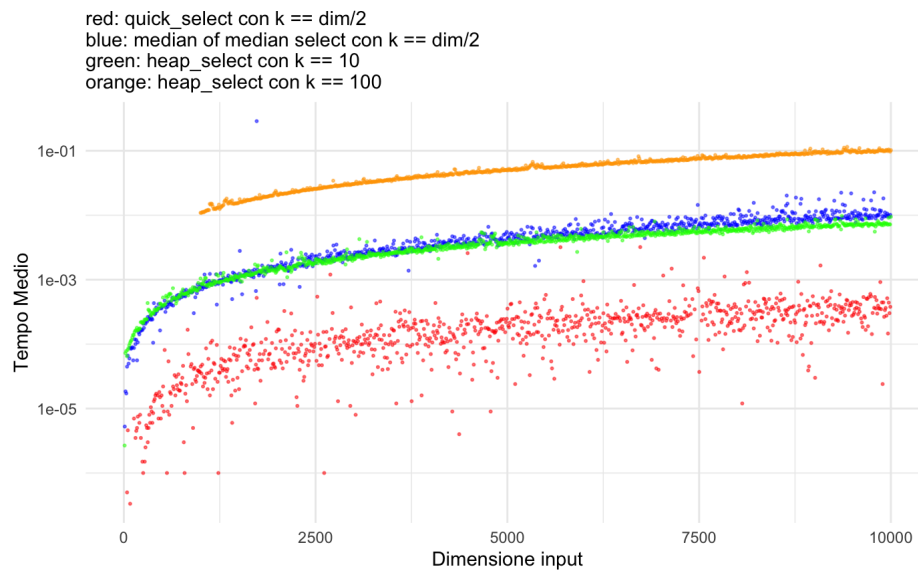
Solamente guardano il valore del tempo medio totale, ma a confermare l'ipotesi anche il grafico sottostante, sembrerebbe che l'algoritmo *quick_select* sia migliore. Ciò non è del tutto vero, essendo che è molto più facile ricadere nel caso pessimo rispetto a quello medio scegliendo un pivot casuale, a differenza di *mOfm_select* che mantiene il tempo medio più o meno costante qualsiasi sia la distribuzione dell'input e la k scelta.



Conviene comunque l'utilizzo di *mOfm_select* essendo che, nel caso in cui si ricada nel caso pessimo di *quick_select*, su input di dimensioni rilevanti, viene risparmiato un tempo di esecuzione importante.



Ugualmente in confronto a *heap_select*, avendo quest'ultimo relativamente pochi best case su input di dimensioni rilevanti, rispetto alla costanza di *mOfm_select*. Di seguito un grafico comparativo fra tutti gli algoritmi presentati.



Si può notare come l'uso di *heap_select* dipenda molto dalla k scelta, l'uso di *quick_select* dalla distribuzione dell'input. In casi con k che varia molto, o senza conoscere la distribuzione dell'input, è preferibile l'utilizzo dell'algoritmo *median of median select*.