

Algorithmique et structures de données : projet

Hamza BOUIHI, Mamadou THIONGANE

7 mai 2023

Table des matières

1	Algorithme naïf	2
1.1	Principe de l'algorithme	2
1.2	Implémentation	2
1.3	Complexité	3
2	Algorithme efficace : approche unidimensionnelle	3
2.1	Principe de l'algorithme	3
2.2	Implémentation	3
2.3	Complexité et performances	4
3	Algorithme efficace : approche bidimensionnelle	6
3.1	Algorithme utilisant un dictionnaire	6
3.1.1	Structure de dictionnaire	6
3.1.2	Principe de l'algorithme	6
3.1.3	Implémentation	7
3.1.4	Complexité et performances	8
3.2	Algorithme utilisant un arbre binaire de recherche	9
3.2.1	Structure d'arbre bidimensionnel	9
3.2.2	Principe de l'algorithme	10
3.2.3	Implémentation	10
3.2.4	Complexité et performances	11
4	Conclusion	12

On s'intéresse dans ce projet à l'identification de structures dans un nuage de points. On considère en entrée un ensemble de points distincts du plan ainsi qu'une distance seuil δ . Deux points p_i et p_j sont proches si $d(p_i, p_j) \leq \delta$.

À partir de cette information, il est possible d'extraire un graphe un graphe $G = (V, E)$ tel que

- tout point p_i correspond à un sommet $v_i \in V$ du graphe
- tout couple (p_i, p_j) de points proches correspond à une arête $(v_i, v_j) \in E$ du graphe

Dans ce projet, on cherche à calculer la distribution des tailles des composantes connexes de G .

1 Algorithme naïf

1.1 Principe de l'algorithme

L'algorithme se découpe de la manière suivante :

- la fonction `explore_connexe_comp` qui renvoie la taille d'une composante connexe. Pour cela, on crée une variable `pile` qui contient la liste des points à traiter dans la composante étudiée. Tant qu'elle est non vide, on itère sur TOUS les points en vérifiant que le point n'est pas dans la pile ou déjà traité avant de vérifier s'il est suffisamment proche pour le rajouter à la composante connexe.
- la fonction `print_components_sizes` qui se contente d'itérer sur la liste et d'appliquer `explore_connexe_comp` lorsqu'un point n'a pas été affecté à une composante connexe.

1.2 Implémentation

On peut réaliser l'algorithme suivant.

```

1  def explore_connexe_comp(dist, points, visited, first):
2      size, pile = 0, [first]
3      # Tant que la pile des voisins du sommet traité est non vide faire :
4      while pile:
5          indice_pts = pile.pop()
6          # Rend le point courant traité
7          visited[indice_pts] = True
8          size += 1
9          # Recherche de voisins
10         for i, point in enumerate(points):
11             if not visited[i] and i not in pile:
12                 if points[indice_pts].distance_to(points[i]) <= dist:
13                     pile.append(i)
14     return size
15
16 def print_components_sizes(dist, points):
17     n_pts, components_sizes, grid = len(points), [], {}
18     visited = [False] * n_pts
19     # Tant que les sommets n'ont pas tous été traités, appeler la fonction
20     for k in range(n_pts):
21         if not visited[k]:
22             size = explore_connexe_comp(dist, points, visited, k)
23             components_sizes += [size]
24     components_sizes.sort(reverse=True)
25     print(components_sizes)

```

1.3 Complexité

Dans toute la suite, on admettra que la complexité moyenne de la méthode `sort` est en $O(n \log(n))$.

On note $C(n)$ la complexité de la fonction `print_components_sizes` et $Q_i(n)$ la complexité de la fonction `explore_connexe_comp`.

$$C(n) = O\left(\sum_{i=1}^n Q_i(n)\right)$$

En notant k_i la taille de la i -ème composante connexe (éventuellement nulle),

$$C(n) = O\left(n \sum_{i=1}^n k_i\right) = O(n^2)$$

$$\text{car } \sum_{i=1}^n k_i = n.$$

2 Algorithme efficace : approche unidimensionnelle

2.1 Principe de l'algorithme

Théorème 1. Soit $p_i = (x_i, y_i)$ et $p_j = (x_j, y_j)$ deux points. Si $|x_j - x_i| > \delta$, alors $(p_i, p_j) \notin E$.

Démonstration. Si $(p_i, p_j) \in E$, alors $|x_j - x_i| \leq \sqrt{(x_j - x_i)^2 + (y_j - y_i)^2} = d(p_i, p_j) \leq \delta$, d'où le résultat par contraposition. \square

On peut alors trier la liste de points dans l'ordre lexicographique afin de réduire le nombre de calculs. En effet, pour chercher les points p_j proches d'un point p_i donné, il suffit de parcourir la liste à partir du premier point p non traité jusqu'à ce que la condition $x_j - x_i \leq \delta$ ne soit plus vérifiée car les points suivants ne peuvent pas appartenir à la composante connexe de p_i d'après le théorème 1. Il faut également réaliser un parcours décroissant entre p et p_i pour vérifier l'existence de points non traités proches de p mais pas de p_i .

2.2 Implémentation

On obtient l'algorithme suivant.

```
1 def print_components_sizes(distance, points):
2     n_points = len(points)
3     points.sort()
4     components_sizes = []
5
6     def component_size(i, j):
7         if j_min < j < n_points:
8             if not points[j].processed:
9                 point, other = points[i], points[j]
10                delta = abs(other.coordinates[0] - point.coordinates[0])
11                if delta <= distance:
12                    if point.distance_to(other) <= distance:
13                        points[j].processed = True
14                        components_sizes[-1] += 1
15                        component_size(j, j+1)
```

```

16         component_size(j, j-1)
17     if i < j:
18         component_size(i, j+1)
19     else:
20         component_size(i, j-1)
21 else:
22     if i < j:
23         component_size(i, j+1)
24     else:
25         component_size(i, j-1)
26
27 for i in range(n_points):
28     if not points[i].processed:
29         points[i].processed = True
30         components_sizes.append(1)
31         j_min = i
32         component_size(i, i+1)
33
34 components_sizes.sort(reverse=True)
35 print(components_sizes)

```

2.3 Complexité et performances

Dans le meilleur cas, le graphe n'admet qu'une composante connexe et la distribution de points est telle que chaque point n'admet qu'un unique point proche dans sa zone de recherche, auquel cas la complexité de la fonction `connectes_list` est en $O(n \log(n))$ (en raison de la complexité de la méthode `sort`).

Dans le pire cas, la distance seuil est telle que la zone de recherche regroupe un trop important nombre de points sans qu'ils soient suffisamment proches et la complexité est alors en $O(n^2)$.

Pour mesurer les performances de nos algorithmes, on génère aléatoirement une série d'instances en faisant varier le nombre de points ainsi que la distance seuil puis on exécute nos algorithmes avec les instances générées.

Pour la fonction `connectes_list`, on obtient les performances suivantes en figures 1 et 2.

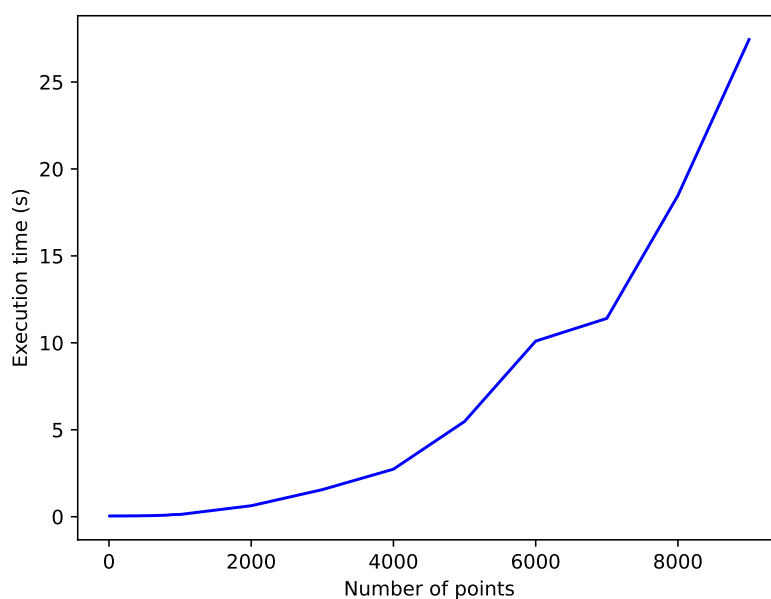


FIGURE 1 – Courbe de performances temporelles de `connectes_list` en fonction du nombre de points (distance = 0.05)

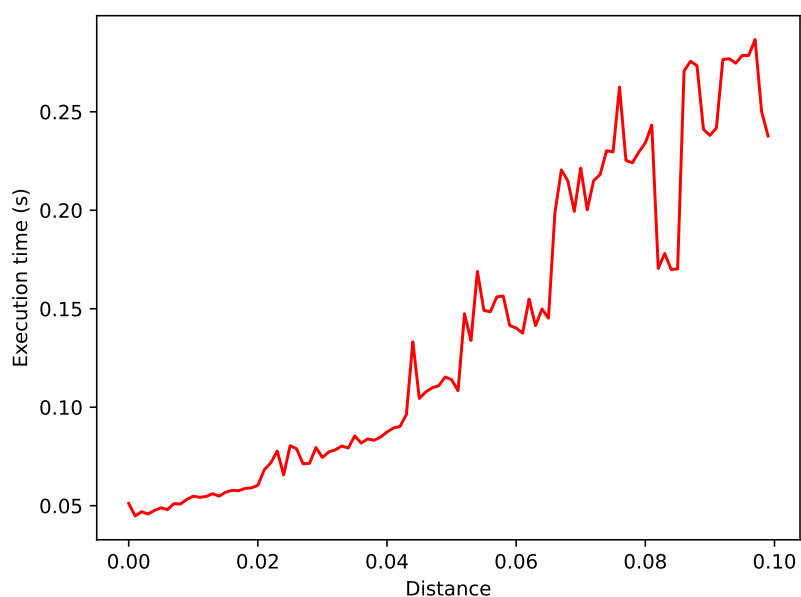


FIGURE 2 – Courbe de performances temporelles de `connectes_list` en fonction de la distance (pour 1000 points)

Sur la figure 2, on observe que le temps d'exécution augmente considérablement avec la distance. Ce résultat était prévisible : en effet, plus la distance est petite, plus la zone de recherche est réduite, ce qui diminue le nombre de calcul.

3 Algorithme efficace : approche bidimensionnelle

3.1 Algorithme utilisant un dictionnaire

3.1.1 Structure de dictionnaire

Un dictionnaire est une collection d'objets non-ordonnée. Chaque élément qu'il possède se compose d'une paire clé-valeur.

Comme les listes, les dictionnaires sont des objets muables et dynamiques. Néanmoins, il possède une propriété particulièrement intéressante en plus : deux éléments différents peuvent avoir la même clé (contrairement aux listes, si on considère que l'indice est une clé).

Cette propriété s'avère être très utile lorsque l'on est amené à regrouper des éléments qui ont une propriété commune : la clé.

3.1.2 Principe de l'algorithme

L'idée principale de l'algorithme est simple : réduire la liste des points qui pourraient appartenir à une composante connexe.

Pour cela, on quadrille le plan afin que le test d'appartenance à la composante connexe porte uniquement sur les carrés voisins du point que l'on est en train de traiter (cf. figure 3).

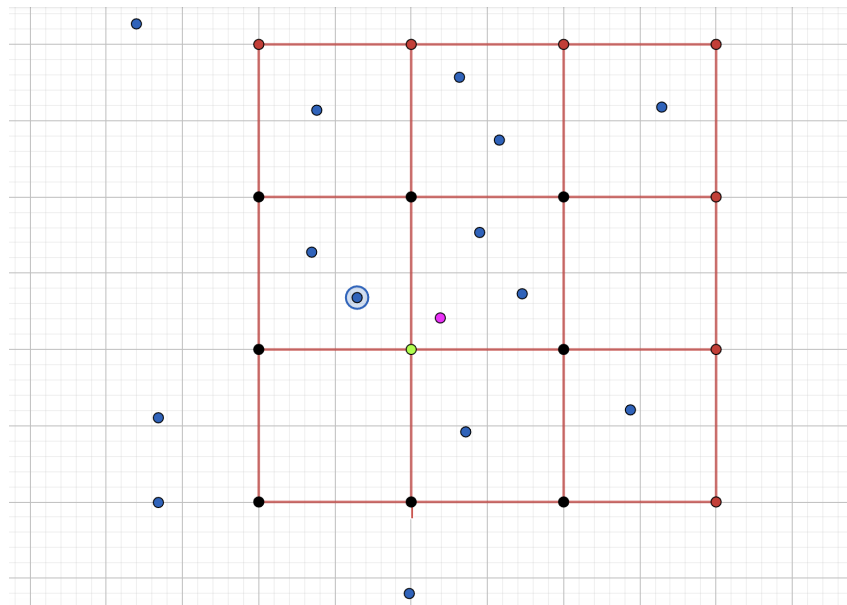


FIGURE 3 – Illustration de l'algorithme

La structure implémentée est un dictionnaire, nommé `grid`, et est implémentée de la manière suivante :

- **clé** : approximation à l'entier près des points de sorte à créer des paquets de points (correspond au point en vert fluo sur la figure 3)
- **valeur** : dictionnaire nommé `square` ayant pour clés les indices des points appartenant au carré et pour valeur les points qui correspondent aux indices (points bleus dans les carrés)

Ainsi, pour énumérer tous les points qui se situent à une distance inférieure au **seuil** du point courant que l'on traite (le violet sur la figure 3), il suffit de tester avec ceux qui appartiennent au même carré que le point courant et aux carrés voisins de celui-ci. Ceci est possible grâce à la variable `square_neighbor`.

3.1.3 Implémentation

On peut implémenter l'algorithme suivant.

```
1 def explore_connexe_comp(dist, points, visited, grid, first, square_neighbor):
2     size, pile = 0, set([first])
3
4     # Tant que la pile des voisins du sommet traité est non vide faire:
5     while pile:
6         indice_pts = pile.pop()
7         current_pts = points[indice_pts]
8         approx_pts = int(current_pts.coordinates[0]/dist),
9         int(current_pts.coordinates[1]/dist)
10
11         # Rend le point courant traité
12         visited[indice_pts] = True
13         grid[approx_pts].pop(indice_pts)
14         size += 1
15
16         # Recherche optimisée de possibles voisins
17         for coord in square_neighbor:
18             key = (approx_pts[0] + coord[0], approx_pts[1] + coord[1])
19             for index_possible_neighbor in grid.get(key, {}):
20                 if not visited[index_possible_neighbor]:
21                     if points[indice_pts].distance_to(
22                         grid[key][index_possible_neighbor]
23                     ) <= dist :
24                         pile.add(index_possible_neighbor)
25
26     return size
27
28 def print_components_sizes(dist, points):
29     n_pts, components_sizes, grid = len(points), [], {}
30     visited = [False] * n_pts
31     square_neighbor = [(-1, -1), (0, -1), (1, -1), (-1, 0),
32                        (0, 0), (1, 0), (-1, 1), (0, 1), (1, 1)]
33
34     # On ordonne les points par paquets qui appartiennent
35     # au même carré de côté dist grâce à l'approximation
36     # à l'entier le plus proche + structure de dictionnaire
37     for k in range(n_pts):
38         pts_k = points[k]
39         approx_pts = int(pts_k.coordinates[0]/dist),
40         int(pts_k.coordinates[1]/dist)
41         square = grid.get(approx_pts, {})
42         square[k] = pts_k
43         grid[approx_pts] = square
44
45     # Tant que les sommets n'ont pas tous été traités, appeler la fonction
46     for k in range(n_pts):
47         if not visited[k]:
48             size = explore_connexe_comp(
49                 dist, points, visited, grid, k, square_neighbor
50             )
51             components_sizes.append(size)
52
53     components_sizes.sort(reverse=True)
54     print(components_sizes)
```

3.1.4 Complexité et performances

Pour la fonction `connectes_dict`, on obtient les performances suivantes en figures 4 et 5.

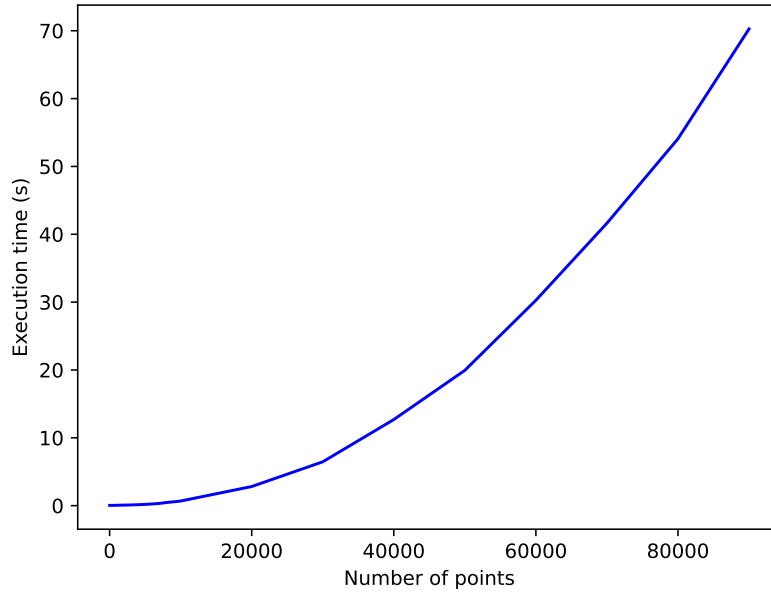


FIGURE 4 – Courbe de performances temporelles de `connectes_dict` en fonction du nombre de points (distance = 0.05)

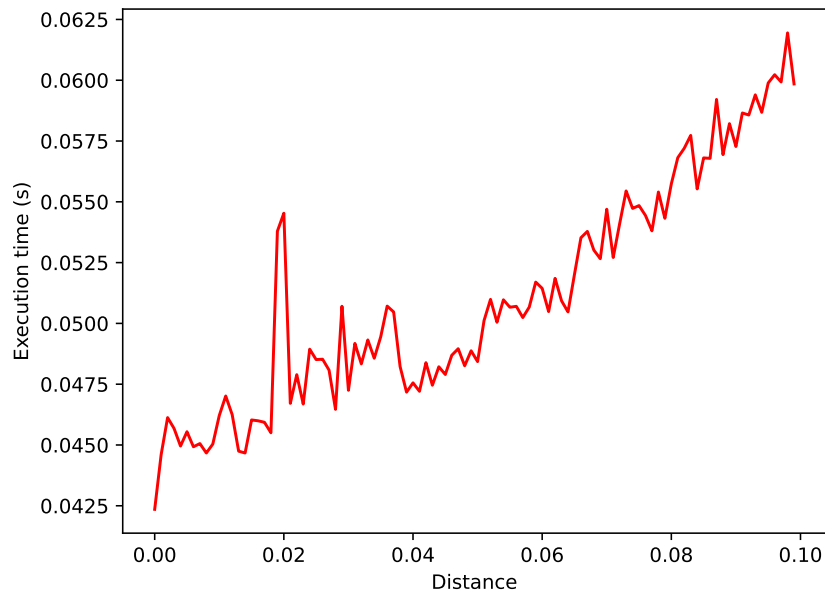


FIGURE 5 – Courbe de performances temporelles de `connectes_dict` en fonction de la distance (pour 10000 points)

3.2 Algorithme utilisant un arbre binaire de recherche

3.2.1 Structure d'arbre bidimensionnel

Un arbre bidimensionnel est une structure de données de partition du plan permettant de stocker des points.

Un arbre bidimensionnel est un arbre binaire dont chaque nœud contient un point. Si $p_i = (x_i, y_i)$ est un point stocké dans un nœud à profondeur paire (resp. impaire), alors, pour tout point $p_j = (x_j, y_j)$:

- si p_j appartient au sous-arbre gauche du nœud, $x_j < x_i$ (resp. $y_j < y_i$)
- si p_j appartient au sous-arbre droit du nœud, $x_j > x_i$ (resp. $y_j > y_i$)

On implémente une classe `Tree` correspondant à cette structure de données.

```
1 class Node:
2     def __init__(self, point, axis, left=None, right=None):
3         self.point = point
4         self.axis = axis
5         self.left = left
6         self.right = right
7
8
9 class Tree:
10     def __init__(self, points):
11         self.root = self.construct(points)
12
13     def construct(self, points, depth=0):
14         if not points:
15             return None
16
17         axis = depth % 2
18         median = mediane(points, axis)
19         m = len(points) // 2
20
21         return Node(
22             median,
23             axis,
24             self.construct(points[:m], depth+1),
25             self.construct(points[m+1:], depth+1)
26         )
```

Cette classe utilise une fonction `mediane` qui reprend le principe du tri rapide afin d'obtenir une complexité minimale en $O(n)$.

```
1 def partition(points, axis, g, d):
2     pivot = points[g]
3     m = g
4     for i in range(g+1, d):
5         if points[i].coordinates[axis] < pivot.coordinates[axis]:
6             m += 1
7             if i > m:
8                 points[i], points[m] = points[m], points[i]
9     if m > g:
10         points[m], points[g] = points[g], points[m]
11     return m
12
13
```

```

14 def mediane(points, axis):
15     k = len(points) // 2
16
17     def aux(g, d):
18         m = partition(points, axis, g, d)
19         if m == k:
20             return points[m]
21         elif m < k:
22             return aux(m+1, d)
23         else:
24             return aux(g, m)
25
26     return aux(0, len(points))

```

La construction de l'arbre est alors un algorithme "diviser pour régner" dont la complexité vérifie $C(n) = 2C(\frac{n}{2}) + O(n)$ soit $C(n) = O(n \log(n))$, d'après le master theorem.

3.2.2 Principe de l'algorithme

Pour chaque point p_i de la liste de points, on réalise un parcours en profondeur depuis la racine de l'arbre en utilisant les propriétés de l'arbre bidimensionnel pour réduire le nombre de calculs :

- on considère le point p_j stocké à la racine et on calcule la différence $d_x = x_j - x_i$
- si $|d_x| \leq \delta$, alors on vérifie que p_j n'a pas été traité puis on vérifie si $d(p_i, p_j) \leq \delta$. Dans ce cas, on marque p_j comme étant traité puis on parcourt de nouveau l'arbre depuis la racine avec p_j . On continue ensuite la recherche en profondeur avec p_i dans les sous-arbres enracinés en p_j selon la coordonnée y
- sinon, on a $d_x > \delta$ (resp. $d_x < -\delta$) et il suffit de continuer le parcours en profondeur uniquement dans le sous-arbre enraciné gauche (resp. droit) en p_j

3.2.3 Implémentation

L'algorithme correspondant est le suivant.

```

1 def print_components_sizes(distance, points):
2     tree = Tree(points)
3     components_sizes = []
4
5     def search(node, point):
6         if node is not None:
7             other = node.point
8             axis = node.axis
9             delta = other.coordinates[axis] - point.coordinates[axis]
10            if abs(delta) <= distance:
11                if not other.processed and point.distance_to(other) <= distance:
12                    other.processed = True
13                    components_sizes[-1] += 1
14                    search(tree.root, other)
15                    search(node.left, point)
16                    search(node.right, point)
17            elif delta > 0:
18                search(node.left, point)
19            else:
20                search(node.right, point)
21
22     for point in points:
23         if not point.processed:

```

```

24         point.processed = True
25         components_sizes.append(1)
26         search(tree.root, point)
27
28     components_sizes.sort(reverse=True)
29     print(components_sizes)

```

3.2.4 Complexité et performances

La fonction de construction de l'arbre bidimensionnel est de complexité en $O(n \log(n))$.

Dans le meilleur cas, le parcours en profondeur dans l'arbre ne se réalise que dans un sous-arbre à chaque appel et la complexité est en $O(\log(n))$. Dans le pire cas, le parcours en profondeur a une complexité en $O(n)$.

Ainsi, la fonction `connectes_tree` a une complexité quasi-linéaire.

Pour la fonction `connectes_tree`, on obtient les performances suivantes en figures 6 et 7.

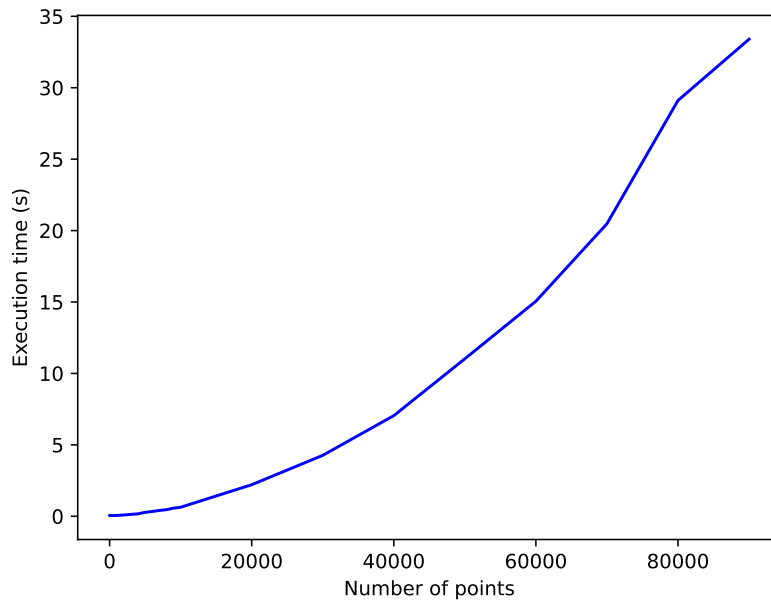


FIGURE 6 – Courbe de performances temporelles de `connectes_tree` en fonction du nombre de points (distance = 0.05)

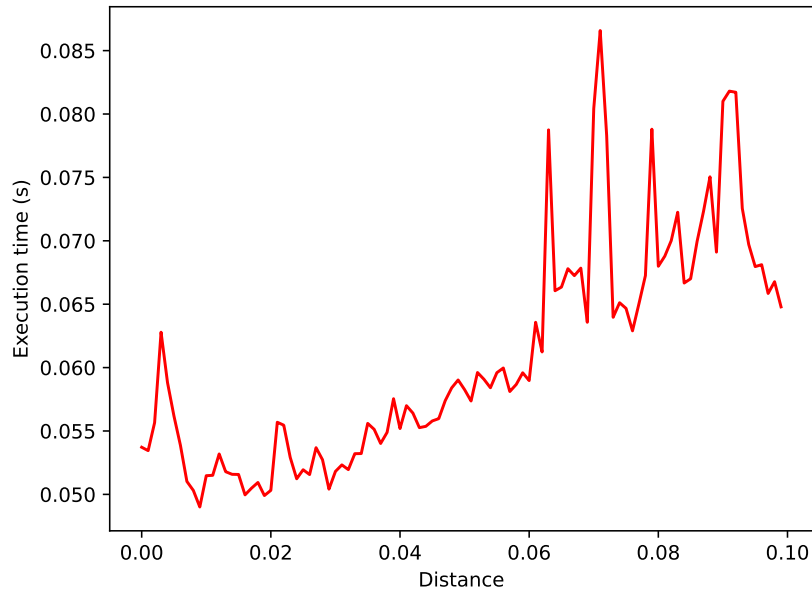


FIGURE 7 – Courbe de performances temporelles de `connectes_tree` en fonction de la distance (pour 10000 points)

Sur la figure 7, on remarque que `connectes_tree` est plus performant lorsque le nombre de composantes connexes est important.

4 Conclusion

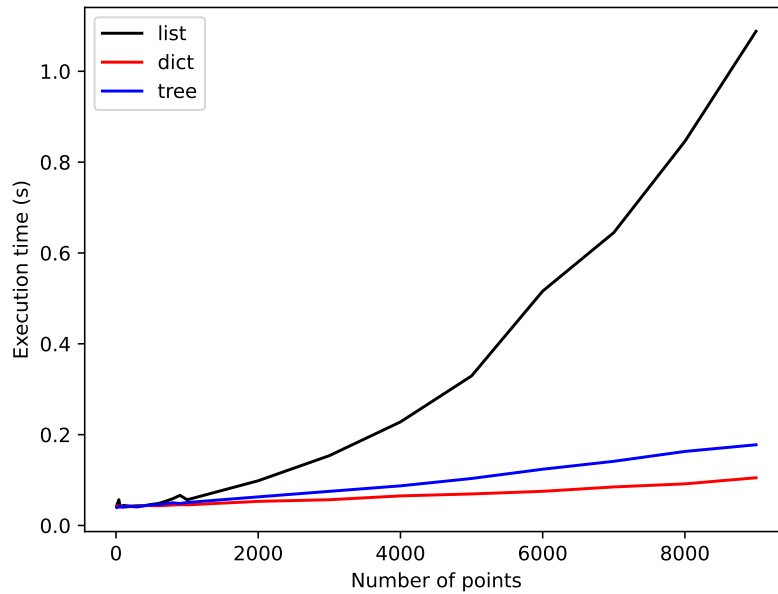


FIGURE 8 – Courbe de performances temporelles des différents algorithmes en fonction du nombre de points (distance = 0.01)

D'après la figure 8, la fonction `connectes_list` à un coût plus élevé que les deux autres fonctions, notamment lorsque que le nombre de points devient important. On peut conclure qu'une approche bidimensionnelle est préférable à une approche unidimensionnelle.

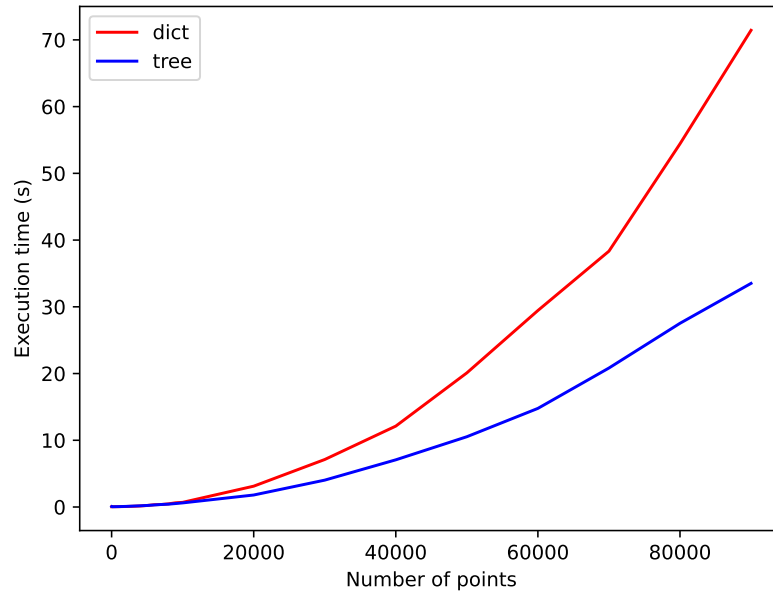


FIGURE 9 – Courbe de performances temporelles des différents algorithmes en fonction du nombre de points (distance = 0.05)

Malgré son efficacité pour un très grand nombre de points par rapport à `connect_dict`, un inconvénient de `connect_tree` est le nombre important d'appels récursifs qu'il nécessite, ce qui peut empêcher l'exécution du programme si ce nombre dépasse la limite autorisée. Si nous avions eu plus de temps, nous aurions tenté réaliser une version itérative de ce programme.