

Recherche des plus proches voisins par des arbres kd

Motivation : Recherche des PPV \rightarrow Utile en apprentissage supervisé pour former des groupes de points.

Définitions : On appelle **point** un élément de \mathbb{N}^d , où $d \in \mathbb{N}^*$ est appelée la **dimension**. La **distance** entre deux points (x_1, \dots, x_d) et (y_1, \dots, y_d) est définie par : $d((x_1, \dots, x_d), (y_1, \dots, y_d)) = (x_1 - y_1)^2 + \dots + (x_d - y_d)^2$.

Problème : k -PPV

Entrée : $p_1, \dots, p_n \in \mathbb{N}^d$, $p \in \mathbb{N}^d$, $k \in \{1, \dots, n\}$

Sortie : k plus proches voisins de p parmi p_1, \dots, p_n .

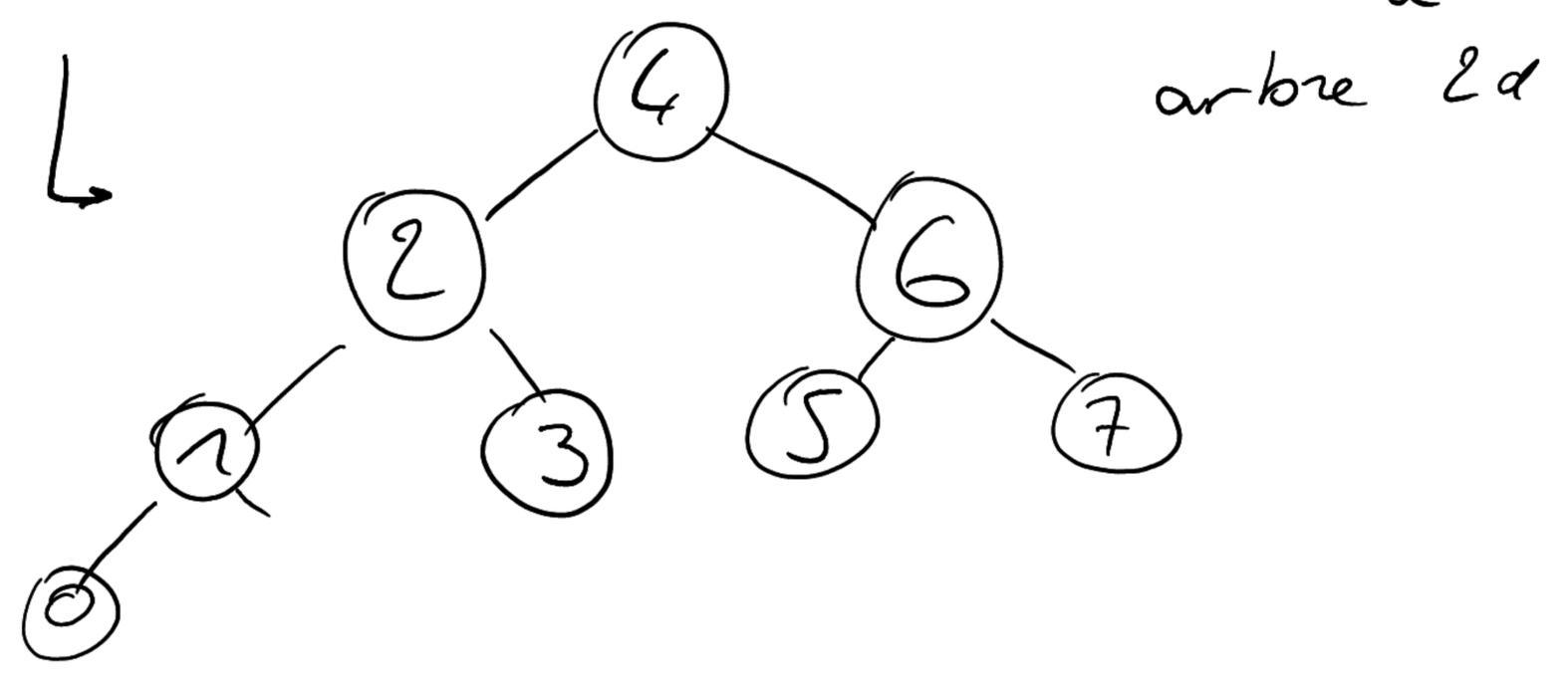
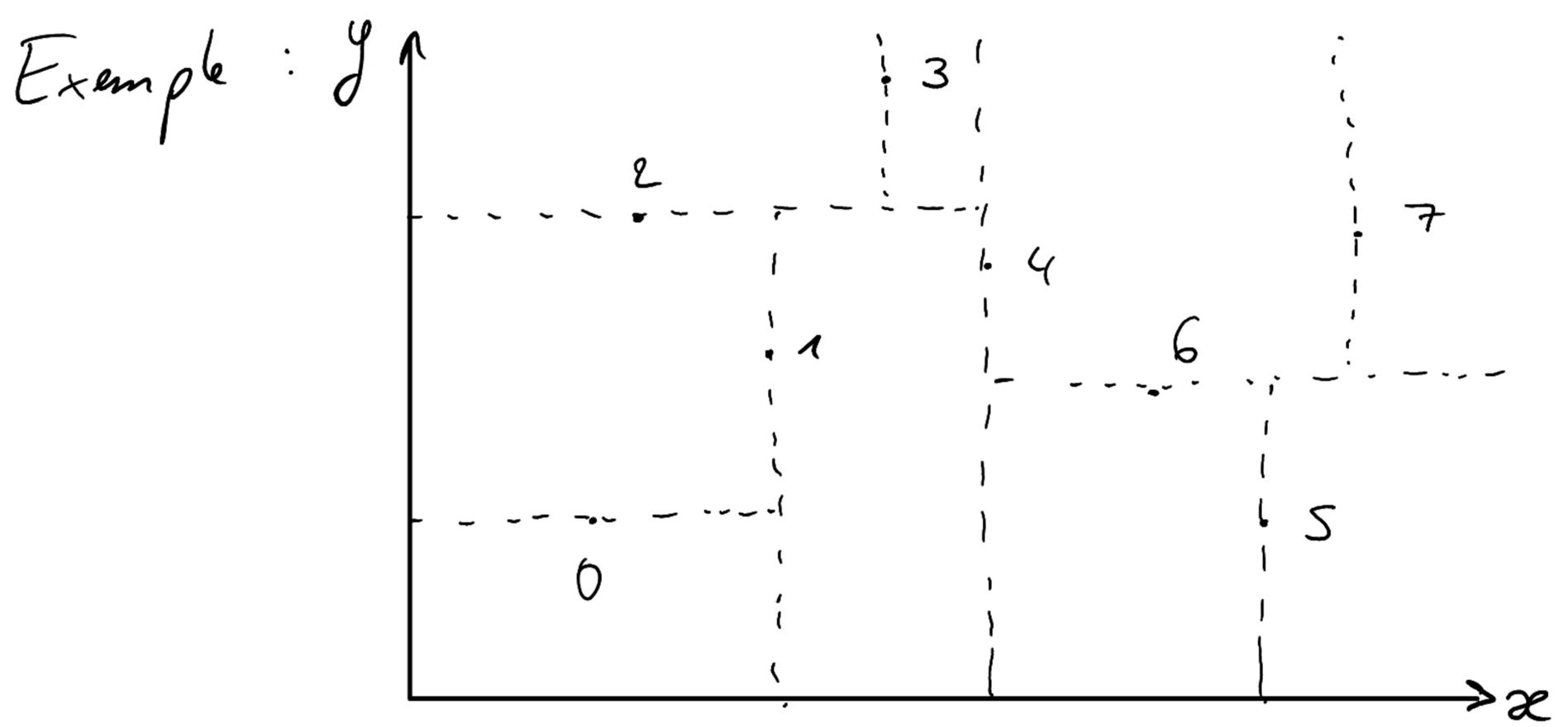
Solutions naïves :

- Tas min contenant les k plus proches voisins
 $\hookrightarrow O(n \log(k))$ en temps et $O(k)$ en mémoire

- Tas min contenant toutes les distances
 $\hookrightarrow O(n + n \log(n))$ en temps et $O(n)$ en mémoire
 \hookrightarrow peu efficace pour beaucoup de requêtes.

Idee : Utiliser une structure d'arbre binaire pour diviser l'espace de recherche

Définition: Un **arbre k -dimensionnel** est un arbre binaire de recherche équilibré où les éléments sont des éléments de \mathbb{N}^d comparés à chaque profondeur dans une dimension différente.



Fonction construction_arbre - rd(points)

| Retourner construction(0, points, 0, |points|)

Fonction construction(i, points, bas, haut)

Si: bas = haut alors

| Retourner Vide

Si non

mid := (bas + haut) / 2

idx_med := Indice du médian des points pour la i^e coordonnée

Echanger points[idx_med] et points[mid]

Placer dans points[bas:mid-1] les éléments \leq $points[idx_med]$ et dans points[mid+1:haut] les éléments \geq $points[idx_med]$.

Retourner Noeud (construction((i+1)/dim, points, bas, mid-1),
points[mid],
construction((i+1)/dim, points, mid+1, haut)

)

Privilégier la recherche pour le développement!
Pas le temps de faire les deux algorithmes.

Fonction $ppv(\kappa, \text{point}, \text{kd tree})$

voisins := creer_tas_max()

Fonction visiter(arbre, i)

Si arbre = Vide alors skip

Si non. arbre = Noeud(l, p, r) alors

t1, t2 := Si point0.i \leq p.i alors (l, r) sinon (r, l)

visiter(t1, (i+1) % dim)

Si |voisins| < κ ou

max(voisins).distance > $\overbrace{|p.i - \text{point0.i}|}^c$ alors

enfiler(voisins, {point = p; distance = $d(p, \text{point0})$ })

Si |voisins| > κ alors supprimer_max(voisins)

visiter(t2)

visiter(kdtree)

Retourner voisins

Complexité :

$O(n \log^2(n))$

pour la construction de l'arbre

$O(\kappa \log(n))$ (en moyenne) par appel à ppv (admis)

Exemple :

