

## Résumé d'article - Cours d'Apprentissage et Fouille de données Alan Tucholka

### An introductory tutorial on kd-trees

Andrew W. Moore  
Carnegie Mellon University

Cet extrait de thèse est une introduction aux structures de type kd-tree, il présente également une implémentation explicite et détaillée d'un algorithme de recherche du plus proche voisin.

#### **Introduction aux kd-trees**

Un kd-tree est une structure de données en arbre binaire pour stocker un ensemble fini de points d'un espace de dimension  $k$ . Chaque nœud contient : (i) les coordonnées d'un point de l'espace (nommé *dom*) à  $k$  dimensions, (ii) une dimension de séparation (nommée *split*), (iii) une branche gauche (respectivement droite) pour les points à gauche (respectivement à droite) de l'hyperplan de séparation. Un point est la branche gauche du nœud si et seulement si sa valeur sur la dimension *split* est inférieure à la valeur de la même dimension du point *dom*.

#### **Recherche du plus proche voisin en utilisant les kd-trees**

Dans cette section, nous verrons une description de l'algorithme que propose l'auteur pour la recherche du plus proche voisin en utilisant un kd-tree.

Une première approximation consiste à trouver le point de la feuille du kd-tree qui contient la cible. S'il existe un point plus proche, il doit se trouver obligatoirement dans l'hypersphère centrée sur la cible et passant par ce point (cf. le cercle figure 1). Passons maintenant au nœud parent de la feuille. Il n'est pas nécessaire de calculer la distance entre ce nœud et tous les nœuds de la branche non explorée s'il n'y a pas d'intersection entre l'hypersphère et l'hyperrectangle non exploré (cf. rectangle gris figure 1). L'auteur présente dans l'article comment implémenter les opérations sur les hyperrectangles.

Regardons plus en détail l'algorithme présenté dans le tableau 1. L'algorithme est appelé récursivement avec 4 variables : le kd-tree, la cible, l'hyperrectangle contenant l'arbre, une distance maximum. L'appel initial prend en paramètre l'hyperrectangle englobant tous les points, et une distance maximum placée à l'infini.

- L'étape 1 gère le cas trivial ou l'arbre est vide.
- L'étape 4 coupe l'hyperrectangle actuel en deux suivant l'axe de coupe du nœud courant.
- Les étapes 5 à 7 déterminent quel hyperrectangle contient la cible.
- L'étape 8 appelle récursivement la fonction avec les paramètres correspondants à la branche de l'arbre contenant la cible.

A partir de maintenant, nous avons une approximation du point le plus proche. Tout point plus proche doit reposer dans l'hypersphère décrite dans le paragraphe précédent.

- L'étape 10 teste s'il y a une intersection entre l'hypersphère et l'hyperrectangle.

Si c'est le cas :

- L'étape 10.1 calcule la distance entre le nœud courant et la cible.
- L'étape 10.2 appelle récursivement la fonction sur la branche encore non explorée.

#### **Performances**

L'algorithme parcourt l'arbre en profondeur jusqu'à arriver à une feuille, il ne va inspecter un nœud qu'une seule fois. Le nombre de nœuds parcouru est donc compris entre  $\log(N)$  et  $N$ . Cependant,  $N$  est exponentiellement dépendant de  $k$ , le nombre de dimensions.

## ANNEXES

<b>Algorithm:</b>	Nearest Neighbour in a <i>kd</i> -tree
<b>Input:</b>	<b>kd</b> , of type <b>kdtree</b> <b>target</b> , of type <b>domain vector</b> <b>hr</b> , of type <b>hyperrectangle</b> <b>max-dist-sqd</b> , of type <b>float</b>
<b>Output:</b>	<b>nearest</b> , of type <b>exemplar</b> <b>dist-sqd</b> , of type <b>float</b>
<b>Pre:</b>	<i>Is-legal-kdtree(kd)</i>
<b>Post:</b>	Informally, the postcondition is that <b>nearest</b> is a nearest exemplar to <b>target</b> which also lies both within the hyperrectangle <b>hr</b> and within distance $\sqrt{\text{max-dist-sqd}}$ of <b>target</b> . $\sqrt{\text{dist-sqd}}$ is the distance of this nearest point. If there is no such point then <b>dist-sqd</b> contains infinity.
<b>Code:</b>	<pre> 1.   if <b>kd</b> is empty then set <b>dist-sqd</b> to infinity and exit. 2.   <b>s</b> := <b>split</b> field of <b>kd</b> 3.   <b>pivot</b> := <b>dom-elt</b> field of <b>kd</b> 4.   Cut <b>hr</b> into two sub-hyperrectangles <b>left-hr</b> and <b>right-hr</b>.       The cut plane is through <b>pivot</b> and perpendicular to the <b>s</b> dimension. 5.   <b>target-in-left</b> := <b>target</b><sub>s</sub> ≤ <b>pivot</b><sub>s</sub> 6.   if <b>target-in-left</b> then 6.1      <b>nearer-kd</b> := <b>left</b> field of <b>kd</b> and <b>nearer-hr</b> := <b>left-hr</b> 6.2      <b>further-kd</b> := <b>right</b> field of <b>kd</b> and <b>further-hr</b> := <b>right-hr</b> 7.   if <b>not target-in-left</b> then 7.1      <b>nearer-kd</b> := <b>right</b> field of <b>kd</b> and <b>nearer-hr</b> := <b>right-hr</b> 7.2      <b>further-kd</b> := <b>left</b> field of <b>kd</b> and <b>further-hr</b> := <b>left-hr</b> 8.   Recursively call <b>Nearest Neighbour</b> with parameters       (<b>nearer-kd</b>,<b>target</b>, <b>nearer-hr</b>,<b>max-dist-sqd</b>), storing the results       in <b>nearest</b> and <b>dist-sqd</b> 9.   <b>max-dist-sqd</b> := minimum of <b>max-dist-sqd</b> and <b>dist-sqd</b> 10.  A nearer point could only lie in <b>further-kd</b> if there were some       part of <b>further-hr</b> within distance <math>\sqrt{\text{max-dist-sqd}}</math> of <b>target</b>.       if this is the case then 10.1     if <math>(\text{pivot} - \text{target})^2 &lt; \text{dist-sqd}</math> then 10.1.1      <b>nearest</b> := (<b>pivot</b>, <b>range-elt</b> field of <b>kd</b>) 10.1.2      <b>dist-sqd</b> := <math>(\text{pivot} - \text{target})^2</math> 10.1.3      <b>max-dist-sqd</b> := <b>dist-sqd</b> 10.2     Recursively call <b>Nearest Neighbour</b> with parameters       (<b>further-kd</b>,<b>target</b>, <b>further-hr</b>,<b>max-dist-sqd</b>),       storing the results in <b>temp-nearest</b> and <b>temp-dist-sqd</b> 10.3     If <b>temp-dist-sqd</b> &lt; <b>dist-sqd</b> then 10.3.1      <b>nearest</b> := <b>temp-nearest</b> and <b>dist-sqd</b> := <b>temp-dist-sqd</b> </pre>
<b>Proof:</b>	Outlined in text

Tableau 1 - Algorithme de recherche du plus proche voisin.

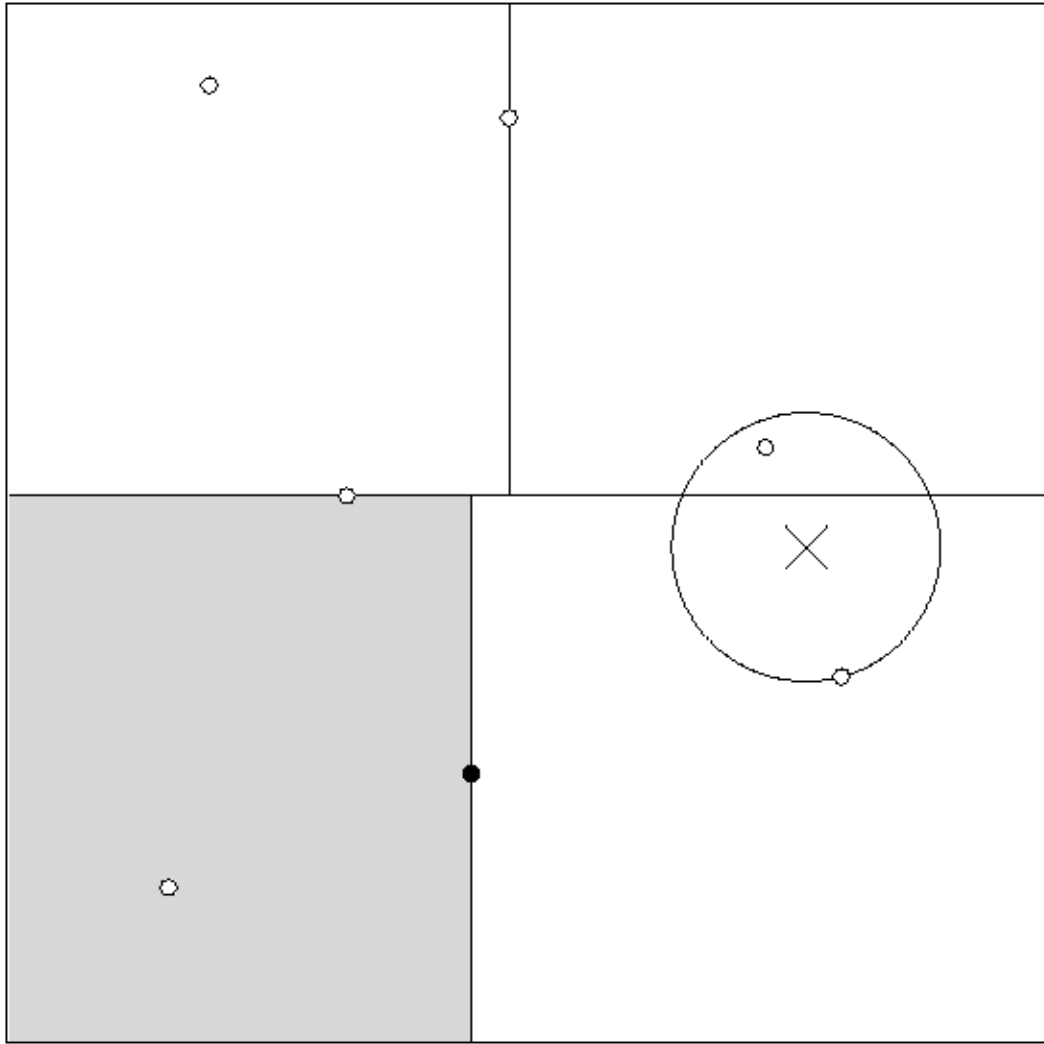


Figure 1 - Le point noir est le parent du premier point le plus proche.  
Il est inutile de parcourir l'hyperrectangle en gris.