

Multidimensional Binary Search Trees

Martin Rau

December 14, 2021

Abstract

This entry provides a formalization of multidimensional binary trees, also known as k -d trees. It includes a balanced build algorithm as well as the nearest neighbor algorithm and the range search algorithm. It is based on the papers "Multidimensional binary search trees used for associative searching" [1] and "An Algorithm for Finding Best Matches in Logarithmic Expected Time" [2].

Contents

1	Definition of the k-d Tree	2
1.1	Definition of the k -d Tree Invariant and Related Functions . .	2
1.2	Lemmas adapted from <i>HOL-Library.Tree</i> to k -d Tree	4
1.3	Lemmas adapted from <i>HOL-Library.Tree-Real</i> to k -d Tree .	5
2	Building a balanced k-d Tree from a List of Points	6
2.1	Auxiliary Lemmas	7
2.2	Widest Spread Axis	7
2.3	Fast Axis Median	8
2.4	Building the Tree	9
2.5	Main Theorems	10
3	Range Searching	11
3.1	Rectangle Definition	11
3.2	Search Function	11
3.3	Auxiliary Lemmas	11
3.4	Main Theorem	12
4	Nearest Neighbor Search on the k-d Tree	12
4.1	Auxiliary Lemmas about <i>sorted-wrt</i>	12
4.2	Neighbors Sorted wrt. Distance	13
4.3	The Recursive Nearest Neighbor Algorithm	13
4.4	Auxiliary Lemmas	14
4.5	The Main Theorems	14
4.6	Nearest Neighbors Definition and Theorems	15

1 Definition of the k -d Tree

```
theory KD-Tree
imports
  Complex-Main
  HOL-Analysis.Finite-Cartesian-Product
  HOL-Analysis.Topology-Euclidean-Space
begin
```

A k -d tree is a space-partitioning data structure for organizing points in a k -dimensional space. In principle the k -d tree is a binary tree. The leafs hold the k -dimensional points and the nodes contain left and right subtrees as well as a discriminator v at a particular axis k . Every node divides the space into two parts by splitting along a hyperplane. Consider a node n with associated discriminator v at axis k . All points in the left subtree must have a value at axis k that is less than or equal to v and all points in the right subtree must have a value at axis k that is greater than v .

Deviations from the papers:

The chosen tree representation is taken from [2] with one minor adjustment. Originally the leafs hold buckets of points of size b . This representation fixes the bucket size to $b = 1$, a single point per Leaf. This is only a minor adjustment since the paper proves that $b = 1$ is the optimal bucket size for minimizing the running time of the nearest neighbor algorithm [2], only simplifies building the optimized k -d trees [2] and has little influence on the search algorithm [1].

```
type-synonym  $'k$  point = (real,  $'k$ ) vec
```

```
lemma dist-point-def:
  fixes  $p_0 :: ('k::finite)$  point
  shows  $dist\ p_0\ p_1 = sqrt\ (\sum\ k \in UNIV.\ (p_0\$k - p_1\$k)^2)$ 
  <proof>
```

```
datatype  $'k$  kdt =
  Leaf  $'k$  point
| Node  $'k$  real  $'k$  kdt  $'k$  kdt
```

1.1 Definition of the k -d Tree Invariant and Related Functions

```
fun set-kdt ::  $'k$  kdt  $\Rightarrow$  ( $'k$  point) set where
  set-kdt (Leaf  $p$ ) = {  $p$  }
| set-kdt (Node - -  $l$   $r$ ) = set-kdt  $l \cup$  set-kdt  $r$ 
```

```
definition spread :: ( $'k::finite$ )  $\Rightarrow$   $'k$  point set  $\Rightarrow$  real where
  spread  $k$   $P = (if\ P = \{\}$  then 0 else let  $V = (\lambda p.\ p\$k)$  '  $P$  in  $Max\ V - Min\ V)$ 
```

```
definition widest-spread-axis :: ( $'k::finite$ )  $\Rightarrow$   $'k$  set  $\Rightarrow$   $'k$  point set  $\Rightarrow$  bool where
```

widest-spread-axis $k K ps \longleftrightarrow (\forall k' \in K. \text{spread } k' ps \leq \text{spread } k ps)$

fun *invar* :: ('k::finite) *kdt* \Rightarrow *bool* **where**
invar (*Leaf* p) \longleftrightarrow *True*
| *invar* (*Node* $k v l r$) \longleftrightarrow $(\forall p \in \text{set-kdt } l. p\$k \leq v) \wedge (\forall p \in \text{set-kdt } r. v < p\$k)$
 \wedge
widest-spread-axis $k UNIV (\text{set-kdt } l \cup \text{set-kdt } r) \wedge \text{invar } l \wedge \text{invar } r$

fun *size-kdt* :: 'k *kdt* \Rightarrow *nat* **where**
size-kdt (*Leaf* $-$) = 1
| *size-kdt* (*Node* $- - l r$) = *size-kdt* l + *size-kdt* r

fun *height* :: 'k *kdt* \Rightarrow *nat* **where**
height (*Leaf* $-$) = 0
| *height* (*Node* $- - l r$) = *max* (*height* l) (*height* r) + 1

fun *min-height* :: 'k *kdt* \Rightarrow *nat* **where**
min-height (*Leaf* $-$) = 0
| *min-height* (*Node* $- - l r$) = *min* (*min-height* l) (*min-height* r) + 1

definition *balanced* :: 'k *kdt* \Rightarrow *bool* **where**
balanced *kdt* \longleftrightarrow *height* *kdt* - *min-height* *kdt* \leq 1

fun *complete* :: 'k *kdt* \Rightarrow *bool* **where**
complete (*Leaf* $-$) = *True*
| *complete* (*Node* $- - l r$) \longleftrightarrow *complete* $l \wedge \text{complete } r \wedge \text{height } l = \text{height } r$

lemma *invar-l*:
invar (*Node* $k v l r$) \Longrightarrow *invar* l
 $\langle \text{proof} \rangle$

lemma *invar-r*:
invar (*Node* $k v l r$) \Longrightarrow *invar* r
 $\langle \text{proof} \rangle$

lemma *invar-l-le-k*:
invar (*Node* $k v l r$) $\Longrightarrow \forall p \in \text{set-kdt } l. p\$k \leq v$
 $\langle \text{proof} \rangle$

lemma *invar-r-ge-k*:
invar (*Node* $k v l r$) $\Longrightarrow \forall p \in \text{set-kdt } r. v < p\k
 $\langle \text{proof} \rangle$

lemma *invar-set*:
set-kdt (*Node* $k v l r$) = *set-kdt* $l \cup \text{set-kdt } r$
 $\langle \text{proof} \rangle$

1.2 Lemmas adapted from *HOL–Library.Tree* to *k-d Tree*

lemma *size-ge0*[*simp*]:

$$0 < \text{size-kdt } kdt \\ \langle \text{proof} \rangle$$

lemma *eq-size-1*[*simp*]:

$$\text{size-kdt } kdt = 1 \longleftrightarrow (\exists p. \text{kdt} = \text{Leaf } p) \\ \langle \text{proof} \rangle$$

lemma *eq-1-size*[*simp*]:

$$1 = \text{size-kdt } kdt \longleftrightarrow (\exists p. \text{kdt} = \text{Leaf } p) \\ \langle \text{proof} \rangle$$

lemma *neq-Leaf-iff*:

$$(\nexists p. \text{kdt} = \text{Leaf } p) = (\exists k \ v \ l \ r. \text{kdt} = \text{Node } k \ v \ l \ r) \\ \langle \text{proof} \rangle$$

lemma *eq-height-0*[*simp*]:

$$\text{height } kdt = 0 \longleftrightarrow (\exists p. \text{kdt} = \text{Leaf } p) \\ \langle \text{proof} \rangle$$

lemma *eq-0-height*[*simp*]:

$$0 = \text{height } kdt \longleftrightarrow (\exists p. \text{kdt} = \text{Leaf } p) \\ \langle \text{proof} \rangle$$

lemma *eq-min-height-0*[*simp*]:

$$\text{min-height } kdt = 0 \longleftrightarrow (\exists p. \text{kdt} = \text{Leaf } p) \\ \langle \text{proof} \rangle$$

lemma *eq-0-min-height*[*simp*]:

$$0 = \text{min-height } kdt \longleftrightarrow (\exists p. \text{kdt} = \text{Leaf } p) \\ \langle \text{proof} \rangle$$

lemma *size-height*:

$$\text{size-kdt } kdt \leq 2 \wedge \text{height } kdt \\ \langle \text{proof} \rangle$$

lemma *min-height-le-height*:

$$\text{min-height } kdt \leq \text{height } kdt \\ \langle \text{proof} \rangle$$

lemma *min-height-size*:

$$2 \wedge \text{min-height } kdt \leq \text{size-kdt } kdt \\ \langle \text{proof} \rangle$$

lemma *complete-iff-height*:

$$\text{complete } kdt \longleftrightarrow (\text{min-height } kdt = \text{height } kdt) \\ \langle \text{proof} \rangle$$

lemma *size-if-complete*:

$complete\ kdt \implies size\text{-}kdt\ kdt = 2^{\wedge} height\ kdt$

<proof>

lemma *complete-if-size-height*:

$size\text{-}kdt\ kdt = 2^{\wedge} height\ kdt \implies complete\ kdt$

<proof>

lemma *complete-if-size-min-height*:

$size\text{-}kdt\ kdt = 2^{\wedge} min\text{-}height\ kdt \implies complete\ kdt$

<proof>

lemma *complete-iff-size*:

$complete\ kdt \longleftrightarrow size\text{-}kdt\ kdt = 2^{\wedge} height\ kdt$

<proof>

lemma *size-height-if-incomplete*:

$\neg complete\ kdt \implies size\text{-}kdt\ kdt < 2^{\wedge} height\ kdt$

<proof>

lemma *min-height-size-if-incomplete*:

$\neg complete\ kdt \implies 2^{\wedge} min\text{-}height\ kdt < size\text{-}kdt\ kdt$

<proof>

lemma *balanced-subtreeL*:

$balanced\ (Node\ k\ v\ l\ r) \implies balanced\ l$

<proof>

lemma *balanced-subtreeR*:

$balanced\ (Node\ k\ v\ l\ r) \implies balanced\ r$

<proof>

lemma *balanced-optimal*:

assumes $balanced\ kdt\ size\text{-}kdt\ kdt \leq size\text{-}kdt\ kdt'$

shows $height\ kdt \leq height\ kdt'$

<proof>

1.3 Lemmas adapted from *HOL–Library.Tree-Real* to *k-d Tree*

lemma *size-height-log*:

$\log\ 2\ (size\text{-}kdt\ kdt) \leq height\ kdt$

<proof>

lemma *min-height-size-log*:

$min\text{-}height\ kdt \leq \log\ 2\ (size\text{-}kdt\ kdt)$

<proof>

lemma *size-log-if-complete*:

$complete\ kdt \implies height\ kdt = \log\ 2\ (size\text{-}kdt\ kdt)$

<proof>

lemma *min-height-size-log-if-incomplete:*
 $\neg \text{complete kdt} \implies \text{min-height kdt} < \log 2 (\text{size-kdt kdt})$
<proof>

lemma *min-height-balanced:*
assumes *balanced kdt*
shows $\text{min-height kdt} = \text{nat}(\text{floor}(\log 2 (\text{size-kdt kdt})))$
<proof>

lemma *height-balanced:*
assumes *balanced kdt*
shows $\text{height kdt} = \text{nat}(\text{ceiling}(\log 2 (\text{size-kdt kdt})))$
<proof>

lemma *balanced-Node-if-wbal1:*
assumes *balanced l balanced r size-kdt l = size-kdt r + 1*
shows *balanced (Node k v l r)*
<proof>

lemma *balanced-sym:*
 $\text{balanced (Node k v l r)} \implies \text{balanced (Node k' v' r l)}$
<proof>

lemma *balanced-Node-if-wbal2:*
assumes *balanced l balanced r abs(int(size-kdt l) - int(size-kdt r)) ≤ 1*
shows *balanced (Node k v l r)*
<proof>

end

2 Building a balanced k -d Tree from a List of Points

theory *Build*
imports
 KD-Tree
 Median-Of-Medians-Selection.Median-Of-Medians-Selection
begin

Build a balanced k -d Tree by recursively partition the points into two lists. The partitioning criteria will be the median at a particular axis k . The left list will contain all points p with $p \ \$ \ k \leq \text{median}$. The right list will contain all points with median at axis $\text{median} < p \ \$ \ k$. The left and right list differ in length by one or none. The axis k will the widest spread axis.

2.1 Auxiliary Lemmas

lemma *length-filter-mset-sorted-nth*:

assumes *distinct xs n < length xs sorted xs*

shows $\{\# x \in \# \text{mset } xs. x \leq xs ! n \#\} = \text{mset } (\text{take } (n + 1) \text{ } xs)$

<proof>

lemma *length-filter-sort-nth*:

assumes *distinct xs n < length xs*

shows $\text{length } (\text{filter } (\lambda x. x \leq \text{sort } xs ! n) \text{ } xs) = n + 1$

<proof>

2.2 Widest Spread Axis

definition *calc-spread* :: $('k::\text{finite}) \Rightarrow 'k \text{ point list} \Rightarrow \text{real}$ **where**

calc-spread k ps = (case ps of [] \Rightarrow 0 | ps \Rightarrow

let ks = map ($\lambda p. p \$ k$) (tl ps) in

fold max ks ((hd ps) \$ k) - fold min ks ((hd ps) \$ k)

)

fun *widest-spread* :: $('k::\text{finite}) \text{ list} \Rightarrow 'k \text{ point list} \Rightarrow 'k \times \text{real}$ **where**

widest-spread [] - = undefined

| *widest-spread [k] ps = (k, calc-spread k ps)*

| *widest-spread (k # ks) ps = (*

let (k', s') = widest-spread ks ps in

let s = calc-spread k ps in

if s \leq s' then (k', s') else (k, s)

)

lemma *calc-spread-spec*:

calc-spread k ps = spread k (set ps)

<proof>

lemma *widest-spread-calc-spread*:

ks \neq [] \implies (k, s) = widest-spread ks ps \implies s = calc-spread k ps

<proof>

lemma *widest-spread-axis-Un*:

shows *widest-spread-axis k K P \implies spread k' P \leq spread k P \implies widest-spread-axis k (K \cup { k' }) P*

and *widest-spread-axis k K P \implies spread k P \leq spread k' P \implies widest-spread-axis k' (K \cup { k' }) P*

<proof>

lemma *widest-spread-spec*:

(k, s) = widest-spread ks ps \implies widest-spread-axis k (set ks) (set ps)

<proof>

2.3 Fast Axis Median

definition *axis-median* :: ('k::finite) ⇒ 'k point list ⇒ real **where**

axis-median k ps = (let n = (length ps - 1) div 2 in fast-select n (map (λp. p\$k) ps))

lemma *length-filter-le-axis-median*:

assumes 0 < length ps ∀ k. distinct (map (λp. p\$k) ps)

shows length (filter (λp. p\$k ≤ *axis-median* k ps) ps) = (length ps - 1) div 2 + 1

⟨*proof*⟩

definition *partition-by-median* :: ('k::finite) ⇒ 'k point list ⇒ 'k point list × real × 'k point list **where**

partition-by-median k ps = (
 let m = *axis-median* k ps in
 let (l, r) = partition (λp. p\$k ≤ m) ps in
 (l, m, r)
)

lemma *set-partition-by-median*:

(l, m, r) = *partition-by-median* k ps ⇒ set ps = set l ∪ set r

⟨*proof*⟩

lemma *filter-partition-by-median*:

assumes (l, m, r) = *partition-by-median* k ps

shows ∀ p ∈ set l. p\$k ≤ m

and ∀ p ∈ set r. ¬p\$k ≤ m

⟨*proof*⟩

lemma *sum-length-partition-by-median*:

assumes (l, m, r) = *partition-by-median* k ps

shows length ps = length l + length r

⟨*proof*⟩

lemma *length-l-partition-by-median*:

assumes 0 < length ps ∀ k. distinct (map (λp. p\$k) ps) (l, m, r) = *partition-by-median* k ps

shows length l = (length ps - 1) div 2 + 1

⟨*proof*⟩

corollary *lengths-partition-by-median-1*:

assumes 0 < length ps ∀ k. distinct (map (λp. p\$k) ps) (l, m, r) = *partition-by-median* k ps

shows length l - length r ≤ 1

and length r ≤ length l

and 0 < length l

and length r < length ps

⟨*proof*⟩

corollary *lengths-partition-by-median-2*:
assumes $1 < \text{length } ps \ \forall k. \text{ distinct } (\text{map } (\lambda p. p\$k) \ ps)$ $(l, m, r) = \text{partition-by-median } k \ ps$
shows $0 < \text{length } r$
and $\text{length } l < \text{length } ps$
<proof>

lemmas *length-partition-by-median =*
sum-length-partition-by-median length-l-partition-by-median
lengths-partition-by-median-1 lengths-partition-by-median-2

2.4 Building the Tree

function (*domintros, sequential*) *build* :: ('k::finite) list \Rightarrow 'k point list \Rightarrow 'k kdt
where
build - [] = *undefined*
| *build* - [p] = *Leaf* p
| *build* ks ps = (
 let (k, -) = *widest-spread* ks ps *in*
 let (l, m, r) = *partition-by-median* k ps *in*
 Node k m (*build* ks l) (*build* ks r)
)
<proof>

lemma *build-dominros3*:
assumes $(k, s) = \text{widest-spread } ks \ (x \# y \# zs)$ $(l, m, r) = \text{partition-by-median } k \ (x \# y \# zs)$
assumes *build-dom* (ks, l) *build-dom* (ks, r)
shows *build-dom* (ks, x # y # zs)
<proof>

lemma *build-termination*:
assumes $\forall k. \text{ distinct } (\text{map } (\lambda p. p\$k) \ ps)$
shows *build-dom* (ks, ps)
<proof>

lemma *build-psimp-1*:
 $ps = [p] \Longrightarrow \text{build } k \ ps = \text{Leaf } p$
<proof>

lemma *build-psimp-2*:
assumes $(k, s) = \text{widest-spread } ks \ (x \# y \# zs)$ $(l, m, r) = \text{partition-by-median } k \ (x \# y \# zs)$
assumes *build-dom* (ks, l) *build-dom* (ks, r)
shows *build* ks (x # y # zs) = *Node* k m (*build* ks l) (*build* ks r)
<proof>

lemma *length-xs-gt-1*:
 $1 < \text{length } xs \Longrightarrow \exists x \ y \ ys. \ xs = x \# y \# ys$

<proof>

lemma *build-psimp-3*:

assumes $1 < \text{length } ps$ $(k, s) = \text{widest-spread } ks$ ps $(l, m, r) = \text{partition-by-median } k$ ps

assumes $\text{build-dom } (ks, l)$ $\text{build-dom } (ks, r)$

shows $\text{build } ks$ $ps = \text{Node } k$ m $(\text{build } ks$ $l)$ $(\text{build } ks$ $r)$

<proof>

lemmas $\text{build-psimps}[simp] = \text{build-psimp-1 } \text{build-psimp-3}$

2.5 Main Theorems

theorem *set-build*:

$0 < \text{length } ps \implies \forall k. \text{distinct } (\text{map } (\lambda p. p\$k) ps) \implies \text{set } ps = \text{set-kdt } (\text{build } ks$

$ps)$

<proof>

theorem *invar-build*:

$0 < \text{length } ps \implies \forall k. \text{distinct } (\text{map } (\lambda p. p\$k) ps) \implies \text{set } ks = \text{UNIV} \implies \text{invar}$

$(\text{build } ks$ $ps)$

<proof>

theorem *size-build*:

$0 < \text{length } ps \implies \forall k. \text{distinct } (\text{map } (\lambda p. p\$k) ps) \implies \text{size-kdt } (\text{build } ks$ $ps) =$

$\text{length } ps$

<proof>

theorem *balanced-build*:

$0 < \text{length } ps \implies \forall k. \text{distinct } (\text{map } (\lambda p. p\$k) ps) \implies \text{balanced } (\text{build } ks$ $ps)$

<proof>

lemma *complete-if-balanced-size-2powh*:

assumes $\text{balanced } kdt$ $\text{size-kdt } kdt = 2 \wedge h$

shows $\text{complete } kdt$

<proof>

theorem *complete-build*:

$\text{length } ps = 2 \wedge h \implies \forall k. \text{distinct } (\text{map } (\lambda p. p\$k) ps) \implies \text{complete } (\text{build } k$ $ps)$

<proof>

corollary *height-build*:

assumes $\text{length } ps = 2 \wedge h$ $\forall k. \text{distinct } (\text{map } (\lambda p. p\$k) ps)$

shows $h = \text{height } (\text{build } k$ $ps)$

<proof>

end

3 Range Searching

```
theory Range-Search
imports
  KD-Tree
begin
```

Given two k -dimensional points p_0 and p_1 which bound the search space, the search should return only the points which satisfy the following criteria:

For every point p in the resulting set:

For every axis k :

$$p_0 \ \$ \ k \ \leq \ p \ \$ \ k \ \wedge \ p \ \$ \ k \ \leq \ p_1 \ \$ \ k$$

For a 2-d tree a query corresponds to selecting all the points in the rectangle that has p_0 and p_1 as its defining edges.

3.1 Rectangle Definition

```
lemma cbox-point-def:
  fixes  $p_0 :: ('k::finite) \ point$ 
  shows  $cbox \ p_0 \ p_1 = \{ p. \ \forall k. \ p_0 \ \$ \ k \ \leq \ p \ \$ \ k \ \wedge \ p \ \$ \ k \ \leq \ p_1 \ \$ \ k \}$ 
  <proof>
```

3.2 Search Function

```
fun search :: ('k::finite) point  $\Rightarrow$  'k point  $\Rightarrow$  'k kdt  $\Rightarrow$  'k point set where
  search  $p_0 \ p_1$  (Leaf  $p$ ) = (if  $p \in cbox \ p_0 \ p_1$  then {  $p$  } else {})
| search  $p_0 \ p_1$  (Node  $k \ v \ l \ r$ ) = (
  if  $v < p_0 \ \$ \ k$  then
    search  $p_0 \ p_1 \ r$ 
  else if  $p_1 \ \$ \ k < v$  then
    search  $p_0 \ p_1 \ l$ 
  else
    search  $p_0 \ p_1 \ l \ \cup \ search \ p_0 \ p_1 \ r$ 
)
```

3.3 Auxiliary Lemmas

```
lemma l-empty:
  assumes  $invar \ (Node \ k \ v \ l \ r) \ v < p_0 \ \$ \ k$ 
  shows  $set-kdt \ l \ \cap \ cbox \ p_0 \ p_1 = \{\}$ 
  <proof>
```

```
lemma r-empty:
  assumes  $invar \ (Node \ k \ v \ l \ r) \ p_1 \ \$ \ k < v$ 
  shows  $set-kdt \ r \ \cap \ cbox \ p_0 \ p_1 = \{\}$ 
  <proof>
```

3.4 Main Theorem

theorem *search-cbox*:
 assumes *invar kdt*
 shows $\text{search } p_0 \ p_1 \ kdt = \text{set-kdt } kdt \cap \text{cbox } p_0 \ p_1$
 $\langle \text{proof} \rangle$
end

4 Nearest Neighbor Search on the k -d Tree

theory *Nearest-Neighbors*
imports
 KD-Tree
begin

Verifying nearest neighbor search on the k -d tree. Given a k -d tree and a point p , which might not be in the tree, find the points ps that are closest to p using the Euclidean metric.

4.1 Auxiliary Lemmas about *sorted-wrt*

lemma
 assumes *sorted-wrt f xs*
 shows *sorted-wrt-take: sorted-wrt f (take n xs)*
 and *sorted-wrt-drop: sorted-wrt f (drop n xs)*
 $\langle \text{proof} \rangle$

definition *sorted-wrt-dist* :: $(k::\text{finite}) \ \text{point} \Rightarrow k \ \text{point list} \Rightarrow \text{bool}$ **where**
 $\text{sorted-wrt-dist } p \equiv \text{sorted-wrt } (\lambda p_0 \ p_1. \ \text{dist } p_0 \ p \leq \text{dist } p_1 \ p)$

lemma *sorted-wrt-dist-insort-key*:
 $\text{sorted-wrt-dist } p \ ps \Longrightarrow \text{sorted-wrt-dist } p \ (\text{insort-key } (\lambda q. \ \text{dist } q \ p) \ q \ ps)$
 $\langle \text{proof} \rangle$

lemma *sorted-wrt-dist-take-drop*:
 assumes *sorted-wrt-dist p ps*
 shows $\forall p_0 \in \text{set } (\text{take } n \ ps). \ \forall p_1 \in \text{set } (\text{drop } n \ ps). \ \text{dist } p_0 \ p \leq \text{dist } p_1 \ p$
 $\langle \text{proof} \rangle$

lemma *sorted-wrt-dist-last-take-mono*:
 assumes *sorted-wrt-dist p ps* $n \leq \text{length } ps$ $0 < n$
 shows $\text{dist } (\text{last } (\text{take } n \ ps)) \ p \leq \text{dist } (\text{last } ps) \ p$
 $\langle \text{proof} \rangle$

lemma *sorted-wrt-dist-last-insort-key-eq*:
 assumes *sorted-wrt-dist p ps* *insort-key* $(\lambda q. \ \text{dist } q \ p) \ q \ ps \neq ps \ @ \ [q]$
 shows $\text{last } (\text{insort-key } (\lambda q. \ \text{dist } q \ p) \ q \ ps) = \text{last } ps$
 $\langle \text{proof} \rangle$

lemma *sorted-wrt-dist-last*:
assumes *sorted-wrt-dist p ps*
shows $\forall q \in \text{set } ps. \text{dist } q \ p \leq \text{dist } (\text{last } ps) \ p$
<proof>

4.2 Neighbors Sorted wrt. Distance

definition *upd-nbors* :: $\text{nat} \Rightarrow ('k::\text{finite}) \text{point} \Rightarrow 'k \text{point} \Rightarrow 'k \text{point list} \Rightarrow 'k \text{point list}$ **where**
upd-nbors $n \ p \ q \ ps = \text{take } n \ (\text{insort-key } (\lambda q. \text{dist } q \ p) \ q \ ps)$

lemma *sorted-wrt-dist-nbors*:
assumes *sorted-wrt-dist p ps*
shows *sorted-wrt-dist p (upd-nbors n p q ps)*
<proof>

lemma *sorted-wrt-dist-nbors-diff*:
assumes *sorted-wrt-dist p ps*
shows $\forall r \in \text{set } ps \cup \{q\} - \text{set } (\text{upd-nbors } n \ p \ q \ ps). \forall s \in \text{set } (\text{upd-nbors } n \ p \ q \ ps). \text{dist } s \ p \leq \text{dist } r \ p$
<proof>

lemma *sorted-wrt-dist-last-upd-nbors-mono*:
assumes *sorted-wrt-dist p ps n ≤ length ps 0 < n*
shows *dist (last (upd-nbors n p q ps)) p ≤ dist (last ps) p*
<proof>

4.3 The Recursive Nearest Neighbor Algorithm

fun *nearest-nbors* :: $\text{nat} \Rightarrow ('k::\text{finite}) \text{point list} \Rightarrow 'k \text{point} \Rightarrow 'k \text{kdt} \Rightarrow 'k \text{point list}$ **where**
nearest-nbors $n \ ps \ p \ (\text{Leaf } q) = \text{upd-nbors } n \ p \ q \ ps$
| *nearest-nbors* $n \ ps \ p \ (\text{Node } k \ v \ l \ r) = (
 \text{if } p\$k \leq v \ \text{then}$
 $\text{let } \text{candidates} = \text{nearest-nbors } n \ ps \ p \ l \ \text{in}$
 $\text{if } \text{length } \text{candidates} = n \wedge \text{dist } p \ (\text{last } \text{candidates}) \leq \text{dist } v \ (p\$k) \ \text{then}$
 candidates
 else
 $\text{nearest-nbors } n \ \text{candidates } p \ r$
 else
 $\text{let } \text{candidates} = \text{nearest-nbors } n \ ps \ p \ r \ \text{in}$
 $\text{if } \text{length } \text{candidates} = n \wedge \text{dist } p \ (\text{last } \text{candidates}) \leq \text{dist } v \ (p\$k) \ \text{then}$
 candidates
 else
 $\text{nearest-nbors } n \ \text{candidates } p \ l$
 $)$

4.4 Auxiliary Lemmas

lemma *cutoff-r*:

assumes *invar* (*Node k v l r*)

assumes $p\$k \leq v \text{ dist } p \ c \leq \text{dist } (p\$k) \ v$

shows $\forall q \in \text{set-kdt } r. \text{ dist } p \ c \leq \text{dist } p \ q$

<proof>

lemma *cutoff-l*:

assumes *invar* (*Node k v l r*)

assumes $v \leq p\$k \text{ dist } p \ c \leq \text{dist } v \ (p\$k)$

shows $\forall q \in \text{set-kdt } l. \text{ dist } p \ c \leq \text{dist } p \ q$

<proof>

4.5 The Main Theorems

lemma *set-nns*:

$\text{set } (\text{nearest-nbors } n \ ps \ p \ \text{kdt}) \subseteq \text{set-kdt } \text{kdt} \cup \text{set } ps$

<proof>

lemma *length-nns*:

$\text{length } (\text{nearest-nbors } n \ ps \ p \ \text{kdt}) = \min n \ (\text{size-kdt } \text{kdt} + \text{length } ps)$

<proof>

lemma *length-nns-gt-0*:

$0 < n \implies 0 < \text{length } (\text{nearest-nbors } n \ ps \ p \ \text{kdt})$

<proof>

lemma *length-nns-n*:

assumes $(\text{set-kdt } \text{kdt} \cup \text{set } ps) - \text{set } (\text{nearest-nbors } n \ ps \ p \ \text{kdt}) \neq \{\}$

shows $\text{length } (\text{nearest-nbors } n \ ps \ p \ \text{kdt}) = n$

<proof>

lemma *sorted-nns*:

$\text{sorted-wrt-dist } p \ ps \implies \text{sorted-wrt-dist } p \ (\text{nearest-nbors } n \ ps \ p \ \text{kdt})$

<proof>

lemma *distinct-nns*:

assumes *invar kdt distinct ps set ps* $\cap \text{set-kdt } \text{kdt} = \{\}$

shows *distinct* (*nearest-nbors n ps p kdt*)

<proof>

lemma *last-nns-mono*:

assumes *invar kdt sorted-wrt-dist p ps* $n \leq \text{length } ps \ 0 < n$

shows $\text{dist } (\text{last } (\text{nearest-nbors } n \ ps \ p \ \text{kdt})) \ p \leq \text{dist } (\text{last } ps) \ p$

<proof>

theorem *dist-nns*:

assumes *invar kdt sorted-wrt-dist p ps set ps* $\cap \text{set-kdt } \text{kdt} = \{\}$ *distinct ps* $0 < n$

shows $\forall q \in \text{set-kdt } kdt \cup \text{set } ps - \text{set } (\text{nearest-nbors } n \text{ } ps \text{ } p \text{ } kdt)$. $\text{dist } (\text{last } (\text{nearest-nbors } n \text{ } ps \text{ } p \text{ } kdt)) \text{ } p \leq \text{dist } q \text{ } p$
 ⟨proof⟩

4.6 Nearest Neighbors Definition and Theorems

definition $\text{nearest-neighbors} :: \text{nat} \Rightarrow ('k::\text{finite}) \text{ point} \Rightarrow 'k \text{ kdt} \Rightarrow 'k \text{ point list}$
where

$\text{nearest-neighbors } n \text{ } p \text{ } kdt = \text{nearest-nbors } n \text{ } [] \text{ } p \text{ } kdt$

theorem $\text{length-nearest-neighbors}$:

$\text{length } (\text{nearest-neighbors } n \text{ } p \text{ } kdt) = \text{min } n \text{ } (\text{size-kdt } kdt)$
 ⟨proof⟩

theorem $\text{sorted-wrt-dist-nearest-neighbors}$:

$\text{sorted-wrt-dist } p \text{ } (\text{nearest-neighbors } n \text{ } p \text{ } kdt)$
 ⟨proof⟩

theorem $\text{set-nearest-neighbors}$:

$\text{set } (\text{nearest-neighbors } n \text{ } p \text{ } kdt) \subseteq \text{set-kdt } kdt$
 ⟨proof⟩

theorem $\text{distinct-nearest-neighbors}$:

assumes $\text{invar } kdt$
shows $\text{distinct } (\text{nearest-neighbors } n \text{ } p \text{ } kdt)$
 ⟨proof⟩

theorem $\text{dist-nearest-neighbors}$:

assumes $\text{invar } kdt \text{ } nns = \text{nearest-neighbors } n \text{ } p \text{ } kdt$
shows $\forall q \in (\text{set-kdt } kdt - \text{set } nns)$. $\forall r \in \text{set } nns$. $\text{dist } r \text{ } p \leq \text{dist } q \text{ } p$
 ⟨proof⟩

end

References

- [1] J. L. Bentley. Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18(9):509–517, 1975.
- [2] J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Softw.*, 3(3):209–226, 1977.