

## k-d-Trees

Verallgemeinerung von Quadtrees,  
bzw. Verallgemeinerung von binären Suchbäumen  
auf  $k$  Dimensionen.

(Letzteres war der Satz von Friedman, Bentley & Finkel  
1977; das "k" im Namen war ursprünglich für die  
Dim. gedacht, also 2-d-Trees; 3-d-Trees, etc.; inzwischen  
aber einfach nur k-d-Tree, oder k-d-Tree (ältere Schreibw.))

Konstruktions-algo:

Bsp.:  $P = \{p^1, \dots, p^n\} \in \mathbb{R}^d$

wähle Splitachse  $i$  (d.h.,  $i$ -te Koord.)

bestimme Median  $m$  von  $\{p_i^1, \dots, p_i^n\}$

bilde Teilmengen

$$P^- := \{p \in P \mid p_i \leq m\}$$

$$P^+ := \{p \in P \mid p_i > m\}$$

Rekursion mit  $P^-$ ,  $P^+$

Knoten  $v$  speichert

- Kinder zu  $P^-$  und  $P^+$
- Median  $m$  und Splitachse  $i$
- evtl. BBox ( $P$ )

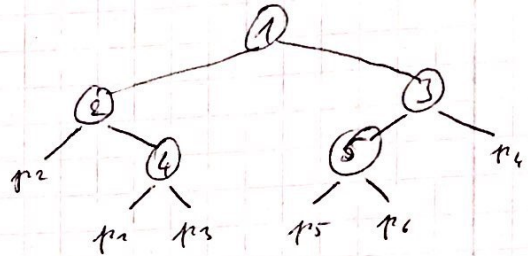
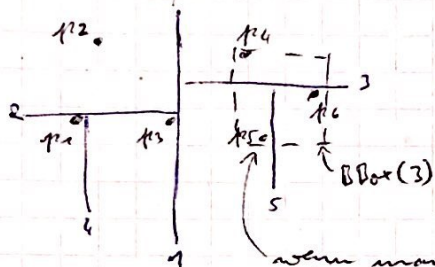
Abbruch der Rekursion bei  $|P| = 1$

Terminologie:

- $i$  heißt auch "Discriminator" (insbes. bei Datenbanken)
- $m$  und  $i$  definieren eine "Splitting-Ebene"

Genauso kann man kd-Trees über Objekten aufbauen, z.B. Records einer Datenbank die nicht unmittelbar Eltern im Raum sind; man benötigt nur auf jeder Komponente eine totale Ordnung.

Beispiel:



wenn man den Median nimmt, dann liegt natürlich der Median immer genau auf der Split-Ebene

Tiefe:  $O(\log n)$   
Bemerkung:

Ein kd-Tree induziert eine Partitionierung des Raumes.

Bezeichne mit  $R(v) \subseteq \mathbb{R}^d$  die Teilregion des Knotens  $v$ .  
 $R(\text{root}) = [-\infty, +\infty]^d$ .  
(nicht verwechseln mit  $BBox(P(v))$ !)

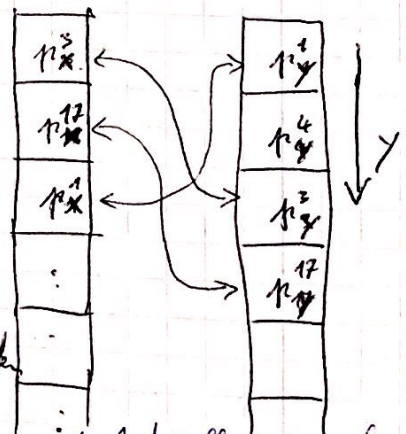
(In der Praxis startet man mit der  $BBox(P)$  an der Wurzel und aktualisiert  $R(v)$  bei der Traversierung.)

Implementierung (in 2D, for sake of simplicity):

sortiere  $P$  1x entlang  $x$  und 1x entlang  $y \Rightarrow$  "x-Liste"

versetze die Listen

Median finden ist jetzt  $O(1)$   
Speicherbedarf  $O(d \cdot n)$  während des Aufbaus  
Listen aufteilen ist  $O(d \cdot n)$   
(aufteilen braucht man doch gar nicht mehr!)



Laufzeit für Aufbau:

Reprocessing  $\rightarrow O(d \cdot n \log n)$  (splitter aller sortierter Listen)

Rekursion:  $T(n) = O(d \cdot n) + 2T(\frac{n}{2})$   
 $\Rightarrow$  Gesamtlaufzeit  $O(d \cdot n \log n)$   
(not für Listen selbst, sondern für die Teilung der Intervalle in größer)

Platz:  $O(n)$   
Der Term  $O(n)$  kommt von der Medianfindung: beim ersten Mal findet man den Median, z.B. in der x-Liste, in  $O(1)$ ; dann teilt man die Liste in 2 Hälften und macht die Rekursion: jetzt wird z.B. entlang  $y$  aufgeteilt; man mischt die ganze x-Liste neu, um herauszufinden, wie weit nach dabei sind; es reicht man bei jeder Rekursion einen Satz von  $d$  Indexintervallen nach unten.

Variationen:

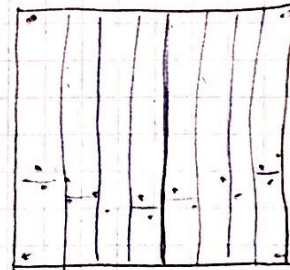
Freidase der Median-Plt im Knoten ( $P = P^- \cup P^+ \cup \{\bar{x}\}$ )

1. Split-Ebene:

a) "longest spread kd tree": wähle diejenige Plt  $i$ , entlang der die Plt die breiteste Verteilung haben, so also  $BBox(P(v))$  größte Ausdehnung hat.

Problem: falls Plt auf einer Mannigfaltigkeit (ungefähr) liegen, deren Dim.  $< d$  ist, führt diese Heuristik oft zu langen, schmalen Zellen.

Beispiel:



b) "longest side kd tree":

splitte entlang der Koord., wo

$R(v)$  am längsten ist.

(Ist also <sup>(fast)</sup> völlig unabhängig von  $P$ )

[Dickerson, Duncan, Goodrich: K-d trees are better when cut on the longest side; ESA 2000]

c) Zyklisch: erst  $x$ , dann  $y$ , dann  $z$ , wieder  $x$ , etc...

2. "Binning" ("bucketing"):

Abbruch bei  $|P| \leq b$ .

(Asymptotische Laufzeiten bleiben gleich.)

$b \approx 20 \dots 30$

[Greenspan, Gohin, Talbot: Acceleration of Binning Nearest Neighbor Methods]

Bemerkung:

Bei kd-Trees fällt die Anzahl Plt assoziiert mit  $R(v)$

auf dem Weg nach unten exponentiell; bei Quadrees

fällt die Größe der Knoten exponentiell (Größe = Seitenlänge).

hier passt "Stackless kd Tree Traversal" (Cropped kd-trees)

Das Nearest-Neighbor-Problem ("closest point problem"):

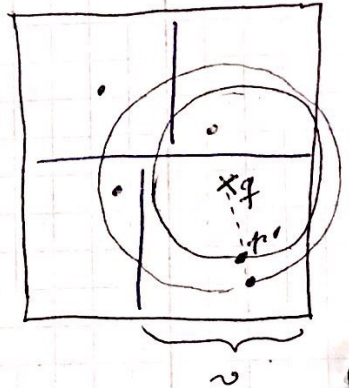
input: Gegeben Ptmenge  $P \subseteq \mathbb{R}^d$  und  $q \in \mathbb{R}^d$ .

output: Gesucht  $p^* \in P$ , so dass  $\forall p \in P: \|p^* - q\| \leq \|p - q\|$ .  
 $p^*$  heißt "NN".

Idee: bestimme das Blatt  $v$  im kd-Tree mit  $q \in R(v)$ ;

das liefert einen Kandidaten  $p'$ ;

man muss man nur noch Blätter besuchen, die die Kugel  $K(q, r)$  um  $q$  mit Radius  $\underbrace{\|q - p'\|}_r$  schneiden.



Algorithmus:

$nn(v, p, r)$ :  
input  $v$ , input  $p$ ,  $r$ :

input:  $p$  = aktueller Kandidat als NN

output:  $r = \|p - q\|$   
 new candidate  $p, p'$  // possibly same as inp.  
 precondition:  $K(q, r)$  überlappt  $R(v)$

[Friedman, Bentley, Finkel: An algo for finding best matches; '77]

//  $K$  = Kugel mit Radius  $r$  um Punkt  $q$

if  $v$  ist Blatt:

$p' :=$  nearest neighbor aus  $P(v)$  zu  $q$

$r' := \|q - p'\|$

if  $r' < r$  then:  $(p, r) := (p', r')$

else:

// Rekursion in das nähere Kind

if  $q_i \leq m_i$ :

$p, r := nn(v_l, p, r)$

else:

$p, r := nn(v_r, p, r)$

// hier kann  $K(q, r)$  schon kleiner geworden sein

// sei  $m, i$  die Split-Ebene in  $v$

// seien  $v_l, v_r$  = linkes / rechtes Kind von  $v$

// hier braucht man keinen bounds overlap test! (?)

// Rekursion in das entfernere Kind

if  $q_i \leq m$ :

if  $K(q, r)$  überlappt  $R(v_r)$ :

$$p, r := \text{min}(v_r, p, r)$$

else:

if  $K(q, r)$  überlappt  $R(v_l)$ :

$$p, r := \text{min}(v_l, p, r)$$

end if

if  $K(q, r) \subseteq R(v)$ :

$p^* := p$  ist die Lsg.; Rekursion abbrechen

return  $p, r$

(\*\*)

Diese Teile könnte man doch sicher schon in die vorigen if-Zweige einbauen

(\*) braucht

man in einer rekursiven Impl. nicht, da es, sobald dieser Fall einsetzt, sicher ist, dass wir nie wieder in den Knoten eines Vater-Knoten zurückkehren werden, wegen (\*\*).

\*) heißt bei Friedman et al. "ball within bounds test"; könnte man auch weglassen, der Algo geht dann halt weiter hoch, aber nie wieder runter in einen anderen bst

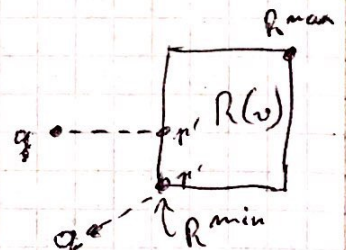
(\*\*) "bounds overlap ball test":

$K(q, r)$  überlappt  $R \Leftrightarrow d(q, R) < r \Leftrightarrow$

$\|q - p'\| < r$ , wobei  $p' \in R$  der nächste

Pkt zu  $q = (q_1, \dots, q_d)$  ist:

$$p'_j = \begin{cases} R_j^{\min} & ; q_j < R_j^{\min} \\ q_j & ; R_j^{\min} \leq q_j \leq R_j^{\max} \\ R_j^{\max} & ; R_j^{\max} < q_j \end{cases}$$



Aufruf:  $\text{min}(\text{root}, \text{NULL}, \infty)$

Analog macht man k-NN-Suche.

Und analog geht "farthest neighbor"-Suche.

Mit kd-Trees kann man sehr einfach auch sog.  
"orthogonal range queries" machen, also alle Pkte  
liefern, die in einer achsen-parallelen Box liegen  $\rightarrow$  Üb.aufgabe  
Ray-Tracing mit kd-Trees: siehe CG 2!

Laufzeit:

Klar ist  $T(n) \in \Omega(\log n)$  und  $T(n) \in O(n)$ ;

bessere Schranken für worst-case gibt's nicht (bis jetzt = 2000)  
(Erickson et al. zeigen in ihrem paper [Dickson, Duncan  
Sachrich]

eine expected Laufzeit von  $O(\log n)$ , aber nur für best Annahmen  
über die Verteilung der Pkte, und für "worst spread" kd-Trees;  
und komischerweise kommt kein  $d$  in der Laufzeit vor.)

# Exkurs: The Curse of Dimensionality

Vorab die  
4D-Videos!

(eine Erzählung ...)

Satz (o. Bew.):

1. Ein kd-Baum über einer Menge von  $n$  Pkten in  $\mathbb{R}^d$  erlaubt eine orthogonale Bereichsanfrage in Zeit  $O(n^{1-\frac{1}{d}} + k)$  und Platz  $O(n)$ .  
↑ # output

2. Jeder algo zur orthog. Range Query über irgend einer Datenstruktur, die nur  $O(n)$  viel Platz benötigt, hat mindestens Laufzeit  $\Omega(n^{1-\frac{1}{d}} + k)$ . [Rolf Klein's Book, S. 134]

(Wenn man also nur  $O(n)$  viel Platz spendiert, dann sind kd-Trees die opt. Datenstruktur für orthog. Range Queries.)

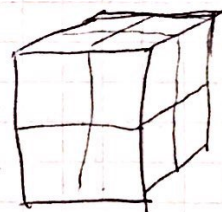
Ein einfaches Rechenbeispiel:

Bsp.  $N = 10^7$  uniform verteilte Pkte in einem Würfel in  $\mathbb{R}^d$ .

Unterteile den Würfel in  $c = 2^d$  Zellen (wie bei Octree)

Anzahl Zellen  $c = 2^d$

→ erwartete Anzahl Pkte pro Zelle  $\mu = \frac{N}{c}$



→ Mindest-Anteil leere Zellen  $e = \frac{c-N}{c}$  (konservative Annahme: ein Pkt pro Zelle)

$d$	$\mu$	$e$
10	$9,8 \cdot 10^5$	0%
30	0,009	99,1%
100	$8 \cdot 10^{-24}$	$\sim 100\%$

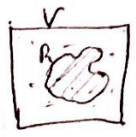
M.a.W.: je höher die Dimension, desto mehr Punkte können man in den Würfel mit fester Seitenlänge packen, ohne dass die Dichte steigt!

Zur Erinnerung:

Sei Volumen  $V \subseteq \mathbb{R}^d$  und Gebiet  $R \subseteq V$  gegeben.

Sei  $P \subseteq \mathbb{R}^d$  eine <sup>zufällig</sup> uniform in  $V$  verteilte Punktmenge.

Dann ist die erwartete Anzahl Pkte in  $R = \frac{\text{Vol}(R)}{\text{Vol}(V)} \cdot |P|$ .



Betrachte die Hyperkugel  $K_r \subseteq \mathbb{R}^d$ :

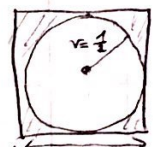
[Manis Kappen: The curse of dimensionality; Seite 1/2]

$$\text{Vol}(K_r) = r^d \cdot \frac{\pi^{d/2}}{(d/2)!}, \text{ falls } d \text{ gerade}$$

Setzt Einheitskugel in Einheitswürfel

$$\frac{\text{Vol}(K_1)}{1^d} = \frac{(\pi^{d/2})^{d/2}}{(d/2)!} = \frac{0,78^{d/2}}{(d/2)!} \rightarrow 0 \text{ für } d \rightarrow \infty$$

(gilt auch für beliebige  $r$ !)



(und das sogar sehr schnell, obwohl  $K$  alle Hyperecken berührt!)

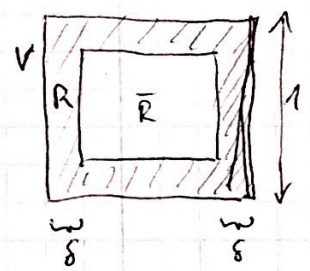
D.h., das meiste Volumen befindet sich in den Ecken!

Und die erwartete Anzahl Pkte in der Kugel  $\rightarrow 0$ !

(Wahrscheinl., da zufällig eines Pkt festleg. Pkt mit  $d$  steigt)

Betrachte den <sup>dünnen</sup> Rand des Hyperwürfels (Hypercube shell)

$$\text{Vol}(R) = 1 - \underbrace{(1 - 2\delta)^d}_{\text{Vol}(\bar{R})}$$

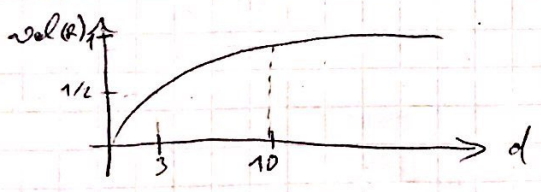


Wahrsch.keit, daß ein Pkt  $p \in P$

den Abstand  $\leq 0.1$  von der Oberfläche des  $1$ -Würfels hat (also in  $R$  liegt):

$d$	$P_r (= \text{Vol}(R) \cdot 100\%)$
3	49%
10	89%
30	99,8%

$\rightarrow$  Je höher die Dimension, desto mehr Punkte liegen an der Oberfläche; m.a.W.



$\rightarrow$  Darts im hoch-dimensionalen Raum ist noch viel schwieriger! ;)

Übungsaufgabe: Wie wächst die Anzahl der "approximativ Ns" mit  $d$ ?

Log.:  $V_\epsilon \sim (1 + \epsilon)^d - 1$





Kombinatorische Zahlen für einen Hyperwürfel:

# Ecken =  $2^d$  , # Kanten =  $d \cdot 2^{d-1}$

# Quadrate =  $\frac{d(d-1)}{2} 2^{d-2}$  (# k-dim Facetten =  $\binom{d}{k} 2^{d-k}$ )

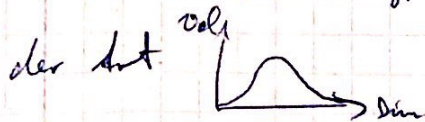
# Hyperecken =  $2d$  ← Hyperwürfel der Dim  $d-1$

Länge der Diagonale =  $\sqrt{d}$

(Man kann sich also einen Hyperwürfel ein wenig so vorstellen wie ein Gebirge mit einem kugelförmigen "Kern", dessen Volumen verschwindend klein ist/wird, und exponentiell vielen "Ecken" (wie ein Tegel), in denen fast das ganze Volumen enthalten ist.) [Körper 9.5]

← { Never compare  $\text{Vol}(K_d)$  with  $\text{Vol}(K_{d+1})$  !!  
(think "cm" versus "cm<sup>2</sup>" !)

Zum Volumen der Hyperkugel: man sieht da manchmal eine Kurve



Das ist aber Quatsch, denn es ist völlig sinnlos, den Inhalt in versch. Dimensionen zu vergleichen

z.B. Fläche in 2D mit Vol in 3D! → Wikipedia-Artikel "Hypersphäre"

Message:  
man darf nur  
 $\text{Vol}(Kug)$  oder  
 $\text{Vol}(Würfel)$   
vergleichen

Remember the Closest Pairs algorithm:

Phase 3 involved a "sparsity" argument.

Consider pts  $P \subseteq \mathbb{R}^d$  with min. dist.  $\delta$

Take any  $p \in P$ , place cube  $C$  centered at  $p$  with "radius"  $\delta$ .

Consider  $L \subseteq P$  inside  $C$ .

Question:  $|L| = ?$

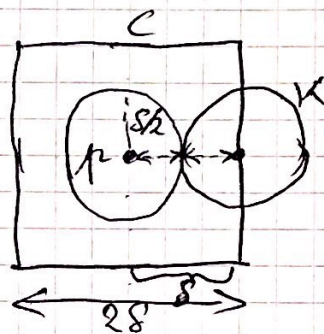
Consider balls  $K_d$  around each  $q \in L$  with radius  $\frac{\delta}{2}$ .  
No two balls can intersect!

Upper bound:

$$\begin{aligned} |L| &\leq \max \# \text{ balls in } C \\ &\approx \frac{\text{Vol}(C)}{\text{Vol}(K)} = \frac{(2\delta)^d}{\frac{\pi^{d/2}}{(d/2)!} \left(\frac{\delta}{2}\right)^d} \\ &= 4^d \cdot \frac{(d/2)!}{\pi^{d/2}} \end{aligned}$$

$$\approx 4^d \frac{(d/2)^{d/2}}{\pi^{d/2}} = 16^{d/2} \cdot \frac{(d/2)^{d/2}}{\pi^{d/2}} \approx (2.54d)^{d/2}$$

Stirling's approx.



Good news: in case  $d$  is fixed  $\rightarrow |L| \in O(1)$

Bad news:  $|L| \rightarrow \infty$  as  $d \rightarrow \infty$

## Approximate Nearest Neighbors

Wegen der Curse of Dimensionality und der unteren Schranke von  $\Omega(n^{1-1/d})$  für orthog. Range Queries glauben viele, daß es keinen Sinn macht, Algos für NN in hohen Dimen. zu suchen, die <sup>signifikant</sup> besser als  $O(n)$  sind. (und nur Platz  $O(n)$  brauchen).

Was kann man tun? Mehr Platz spendieren: etwas machen, nicht exakten NN suchen (reicht für die allermeisten Appl.)

Dickson, Duncan, Gao: kd-trees are better when cut on the longest side, ESA 2000.  
Duncan, Gao, Kobayashi: Balanced aspect ratio trees...  
Arya, Mount, Netanyahu, Silberman, Wu: An optimal algo for approx nearest neighbor searching - ; 1998, J.B. + 19

Def.:

Sei  $S \subseteq \mathbb{R}^d$  eine Menge von Pkten,  $q \in \mathbb{R}^d$  ein Query-Pkt,  $p^* \in S$  der NN, und  $\varepsilon > 0$ .

Dann heißt ein  $p \in S$  " $(1+\varepsilon)$ -<sup>approximate</sup> nearest neighbor"  $\Leftrightarrow$   
 $d(p, q) \leq (1+\varepsilon) \cdot d(p^*, q)$ .

Bezeichnungen:

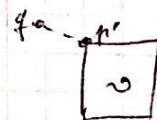
Als jetzt:  $v$  und  $R(v)$  bzw. "Knoten" und "Region des Knoten"

sind Synonyme.

Annahme: metric space, d.h.  $d()$  ist Metrik, z.B.  
 $d(p, q) := \|p - q\|$ ,  $d(v, q) := d(R(v), q)$

$p_v :=$  der in  $v$  liegende Pkt, falls  $v$  Blatt ist

(sonst hatten wir ja  $R(v) =$  Menge der in  $v$  liegenden Pkte)



Algo:

$Q$  = Priority-Queue mit Zeiger auf Knoten im  $k$ - $d$ -Tree, sortiert nach  $d(v, q) =$  Abstand zw.  $q$  und  $v$ ; (hier am tiefen kleinster Abstand vorne)

$p^0 :=$  Blatt unevtl. weit weg (= aktueller Kandidat des kNN)

$v :=$  Root des  $k$ - $d$ -Trees

while  $d(v, q) < \frac{1}{1+\epsilon} d(p^0, q)$ :

while  $v \neq$  Blatt:

seien  $v_1, v_2$  die Kinder von  $v$ , und  $d(v_1, q)$

füge  $v_2$  in  $Q$  ein

$\leq d(v_2, q)$

while  $v := v_1$

if  $d(p_v, q) < d(p^0, q)$ :

{  $v$  ist jetzt Blatt }

then

$p^0 := p_v$

$v \leftarrow$  extractMin( $Q$ )

$v :=$  minimales Element aus  $Q$ , lösche dieses aus  $Q$ ;

end while

return  $p^0$

(Das entscheidende ist, dass der Algo das  $p^*$  gar nicht kennen muss! Ist für  $\epsilon=0$  genau der ursprüngl. Algo, nur eben in iterativer Form)

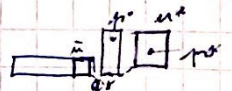
Zeit:  $O(l \cdot \log n)$ ,  $l = \#$  besuchter Blätter

Bew.:

1. Korrektheit: Sei  $u^*$  Blatt aus  $k$ - $d$ -Tree, das  $p^*$  enthält

a)  $u^*$  wurde besucht  $\Rightarrow$  Algo liefert  $p^* = p^0$

b)  $u^*$  wurde nicht besucht  $\Rightarrow p^0 \neq p^* \Rightarrow$



$d(p^*, q) \geq d(u^*, q) \geq d(\bar{u}, q) \geq \frac{1}{1+\epsilon} d(p^0, q)$

weil  $d()$  Metrik

nähere Knoten werden zuerst besucht

Abbruchzeit.

wobei  $\bar{u}$  der zuletzt besuchte Knoten sei. (Blatt oder innerer Knoten)

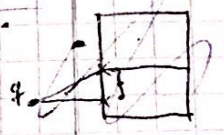
2. Zeit:

pro <sup>inner</sup> Schleifendurchlauf:  $k$  Knoten extrahieren,  $O(\log n)$  Knoten einfügen; <sup>inner loop</sup>  $O(1)$

Fibonacci-Heap verwenden  $\Rightarrow O(\log n)$  für Extrah.,  $O(1)$  amortisierte Kosten für Einfügen;  $d(v, q)$  braucht  $O(1) \Rightarrow$  Beh.

Bem.:

1. In der Praxis bleibt Heap relativ klein  $\rightarrow$  verwende normalen Heap.
2.  ~~$d(v, q)$  braucht nur  $O(1)$ , weil man Distanz inkrementell aus Dist. zum Vater berechnen kann.~~
3. Analog geht "(1- $\epsilon$ )-farthest neighbor".
4. gibt PD-Software "AVL Library", die <sup>gibt inzwischen bessere S-Felder</sup> verwendet aber eine wesentlich komplexere Datenstruktur, weil damals noch nicht bekannt war, daß auch die alten "longest-side b-d-trees" diese schöne worst-case-Laufzeit haben.
5. # besuchter Blätter  $\in O(\log^{d-1} n)$  [ESA 2000];



sieht man, indem man Schranke für Anzahl der Knoten des k-d-Tree findet, die einen Strahlus "durchstechen".  
 Genauer:  $O\left(\left(\frac{\log n}{\epsilon}\right)^{d-1}\right)$   
 gilt nur für "longest-side" kd-Tree!

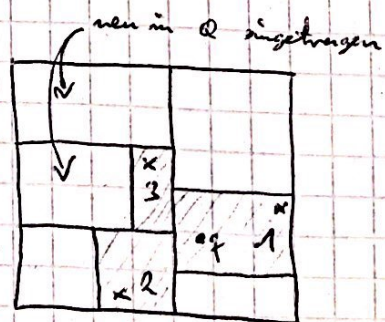


Satz:

Zu einer Menge  $S \subseteq \mathbb{R}^d$ ,  $|S|=n$ , und einem Query-Pkt  $q \in \mathbb{R}^d$  findet man den  $(1+\epsilon)$ -nearest neighbor in Zeit  $O\left(\frac{\log^d n}{\epsilon^{d-1}}\right)$ .

Beispiel zum Algorithmus:

k-d-Tree, soweit er in der Queue eingetragen ist/war:



in Q eingetragen  
 schon besuchtes Blatt, aus Q schon wieder entfernt.

## Beste ANN-Algos

Achtung: "beste" hängt von vielen Faktoren ab!

- Metrik oder nicht-metrischer Abstand
- Verteilung der Pkte (hohe Korrelation oder uniform)
- Dimension

Folgende zwei Algos funktionieren in praxi sehr gut für viele Daten

Alternatives Abbruchkriterium: besuche max  $L_{max}$  viele Blätter.

Alternatives Güte-Maß des ANN-Algo:

$$\text{precision} = \frac{\# \text{ exakter NN's}}{\# \text{ Queries}}$$

$$\text{error} = 1 - \text{precision} \quad (\text{eigtl. "error"})$$

Randomized kd-Tree (RKD):

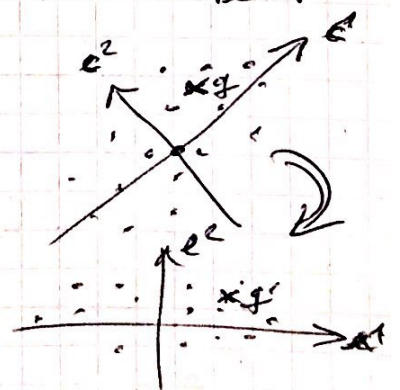
- Bestimme  $D$  Achsen (dimensionen) mit höchster Varianz der Pkte
- Wähle eine davon zufällig
- Splitte entlang Median dieser Achse

Experimente zeigen:  $D=5$  reicht

Muja & Lowe: Fast Approximate Nearest Neighbors with Automatic Algo Configuration; 2005.  
Silpa-Anan & Hartley: Optimized KD-trees for fast image descriptor matching; 2008.

PCA-RKD:

- Bestimme Principal Components der Pkte  $P$
- Transformiere Pkte  $\rightarrow P'$
- Baue RKD über  $P'$
- Transformiere Query-Pkt  $q \rightarrow q'$
- Weiter mit bisherigem ANN-Algo



RKD-Forest (auch mit PCA):

Bau mehrere (20-50) RKD-Trees über  $P$  (bzw  $P'$ )

ANN-Search mittels RKD-Forest:

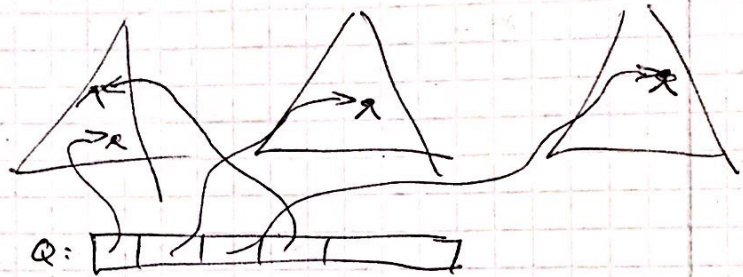
Maintain 1  $p$ -queue only (one for all)

Descend each RKD-tree down to leaf closest to  $q$   
(first iteration of outer while loop of ANN algo)

Choose closest  $p$  of all those leaves

Proceed with ANN-algo as before  
put others in  $p$ -queue

Queue contains pointers to nodes in different trees!



Warum funktioniert das besser?

Bsp.: 1 Mio Pkte  $\approx 2^{20} \Rightarrow$  Tiefe der kd-Trees = 20

Dimension = 100  $\Rightarrow$  ANN schaut 80 Einträge

von  $q$  gar nicht an!

$\Rightarrow q$  und  $p^o$  in Blatt sind nah zueinander  
in 20 Dim., aber in den 80 anderen  
nicht notwendigerweise

Hier hilft die Randomisierung

Andere Betrachtung:

- Beim klassischen ANN-Algo hängen die  
nacheinander (aus der  $P$ -Queue kommenden)  
besuchten Zellen stark voneinander ab

- Bei obigem Algo liegen die Zellen aus der  
 $P$ -Queue in verschiedenen, randomisierten  
kd-Trees  $\Rightarrow$  weniger Abhängigkeit dazwischen!

$k$ -Means-Tree :

1. Partitioniere  $P$  mittels  $k$ -means-Clustering in  $k$  Cluster
2. Erzeuge für jeden Cluster einen  $k$ -Means-Tree
3. Erzeuge Knoten mit  $k$  Kindern

ANN-Search mittels  $k$ -Means-Tree :

Transversiere Baum bis zum nächsten Blatt  
Schreibe nicht-besuchte Kinder in  $P$ -Aneane  
Sortiere  $P$ -Aneane bzgl.  $\text{dist}(q, \text{cluster-center})$   
Stoppe nach  $L_{\max}$  besuchten Blättern

Bem.:

- $k$ -means-Clustering benötigt auch NN  
(allerdings immer nur zu den  $k$  Mittelpunkten der Cluster)
- Randomisierung und Forest haben <sup>hier</sup> nichts gebracht (lt. Autoren)
- Aufbau ist rel. teuer  $\rightarrow$  führe nur wenige Iterationen in  $k$ -means-Clustering durch, nicht bis Konvergenz laufen lassen.  
ANN wird kaum beeinträchtigt.