

Efficient nearest neighbor search with KD trees and LSH

Andreas Knoblauch

HRI-EU Scientific Seminar, 21/03/2012

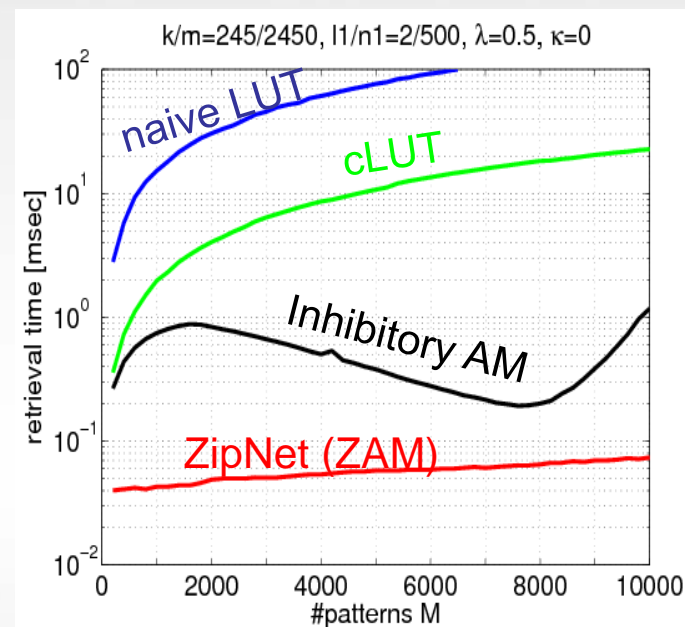
innovation through science

Applications:

- Similarity search in large databases
- prototype-based classification
- matching of local image features or image patches

Algorithms:

- look-up tables (naive)
- **kd-trees** (Bentley '90)
- multi-index hashing (Greene/Parnas/Yao '94)
- neural associative networks (Willshaw'69, Knoblauch '10)
- **locality-sensitive-hashing** (Andoni/Indyk '06)
- ...
- ideal: search time $O(1)$ independent of DB size N (as in brain)



Today: Look at **kd-trees** & LSH

↳ locality sensitive hashing
 ↳ ist alternative läuft über Hashing \Rightarrow sehr schnell

Given a database of key vectors $\{u^\mu\}$ and a query u^* , return the nearest neighbor (NN) key vector u^{μ^*} (e.g., w.r.t. Euclidean distance $d=d_2$)

Relevant problem parameters:

N : size of database (# stored key vectors) (Anzahl Datenvektoren)

D : dimensionality of key vectors

(K : # non-zero entries per key vector)

an wichtigsten!

Relevant complexity measures:

- t_q : query time (Anzahl Rechenschritte um NN zu finden) \rightarrow nearest neighbor zu u^* Datensatz finden \rightarrow ist u^3 wie an schnellsten nach?
- t_c : time to create index structure (wie lang geht es um KD-Tree aufbauen)
- t_u : update time (insert / delete keys) zu updaten und wieviel Speicher nötig?
- S : size of index structure

\rightarrow Datenbank linear durchsuchen.
suche den Datenvektor, der den kleinsten Abstand zu u^* hat

Naive algorithm (linear scan of database):

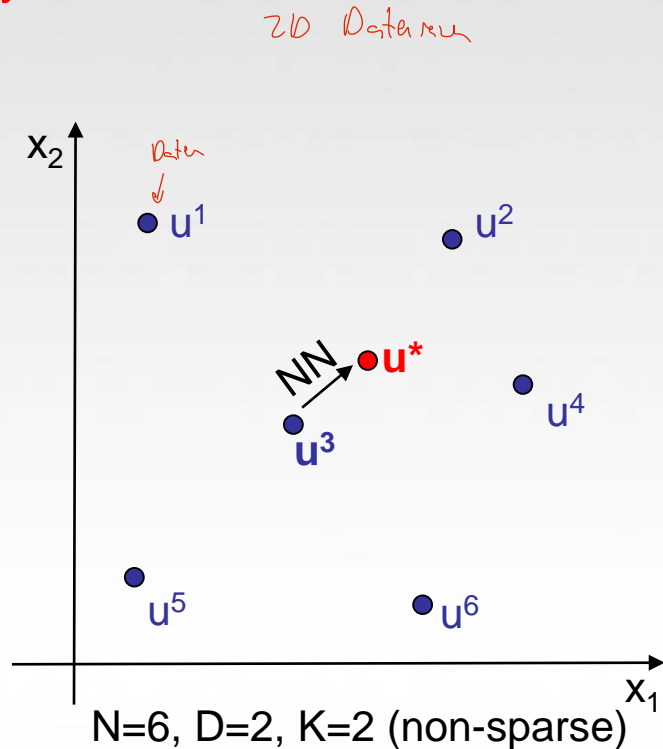
$\mu^* = 1$;
For $\mu=2$ TO N DO
 IF $d(u^*, u^\mu) < d(u^*, u^{\mu^*})$ $\mu^* = \mu$ \leftarrow linear scan DB durchsuchen und Abstand für jeden Datenvektor berechnen
RETURN μ^* \leftarrow neuer NN definieren der besser ist
 \leftarrow Index des finalen NN

$t_q = O(ND)$
 $t_c = O(ND)$
 $t_u = O(D)$
 $S = O(ND)$

(or $O(NK)$ for sparse data vectors)
(or $O(NK)$)
(or $O(K)$)
(or $O(NK)$)

\rightarrow prohibitive for large N, D

Laufzeitparameter

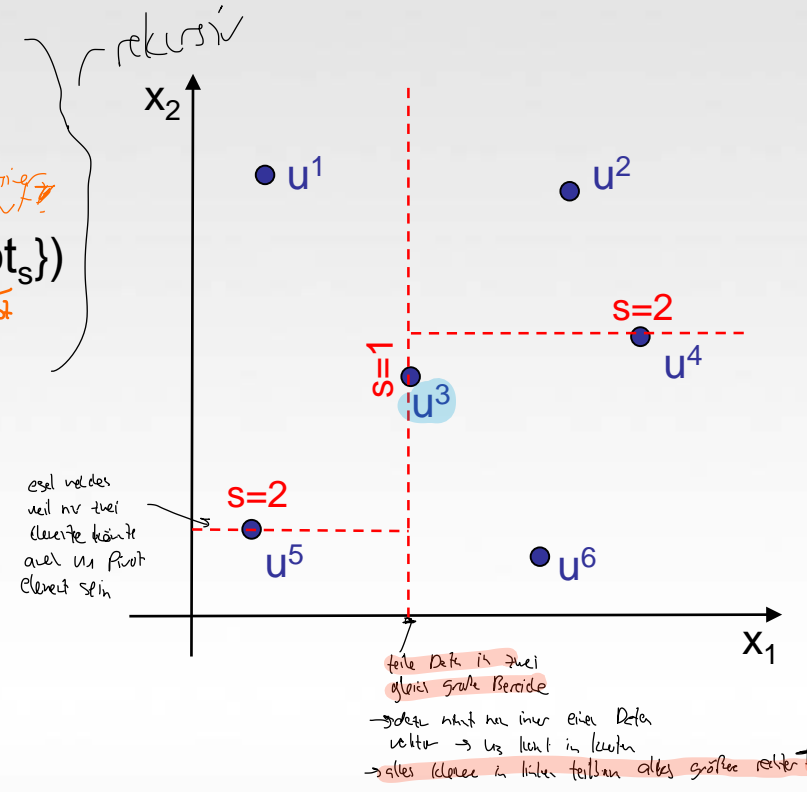
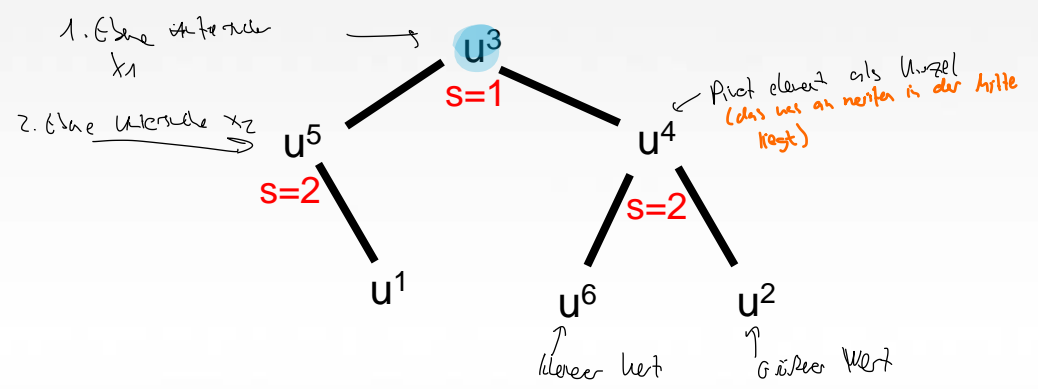


Binärer Baum: jeder Knoten hat max. 2 Zweige!
An jeder Stelle des Baums wird 1 D. ausgewählt (1. Schritt hier z.B. x_1)
→ dann sucht man schnellwert, sodass Daten gleichmäßig aufgeteilt werden
kann KD-Tree auf Menge von Daten (Matrix)
Master auf linken und rechten Baum
Split-Dimension
kürzester Knoten

KD-trees (1): Creation

function (root, left, right, s) = Create-KD-Tree(DB) // returns kd-tree for database DB:={ u^i }

- 1) IF isEmpty(DB) THEN RETURN 0 (Rekursions-Ende)
- 2) Select pivot vector root \in DB and split dimension $s \in \{1, \dots, D\}$
(siehe Split-Dimension und welche Split-Vektor)
(kürzer als root)
- 3) left := Create-KD-Tree($\{u \in DB : u_s < root_s\}$)
(kleiner als root) rekursiv auf!
- 4) right := Create-KD-Tree($\{u \in DB \setminus \{root\} : u_s \geq root_s\}$)
(größer als root)
- 5) RETURN (root, left, right, s)



- KD-tree defines **partitioning of key space into rectangular regions**
 - **various variants** to implement step 2 (choose root, s)
 - try to get **balanced trees** → short paths
Ziel: → dann belad sich in jedem Schritt überm Daten → benutzte ist dann $\log(x)$
 - try to get **square regions** (and **avoid elongated regions**)
→ maximize „cut-off opportunities“ (see retrieval...)
- ähnlich wie bei Dichte-schätzung? Raum wird in Rechtecke eingeteilt!

⇒ Ziel: Finde zu u^* den NN aus den KD-Tree

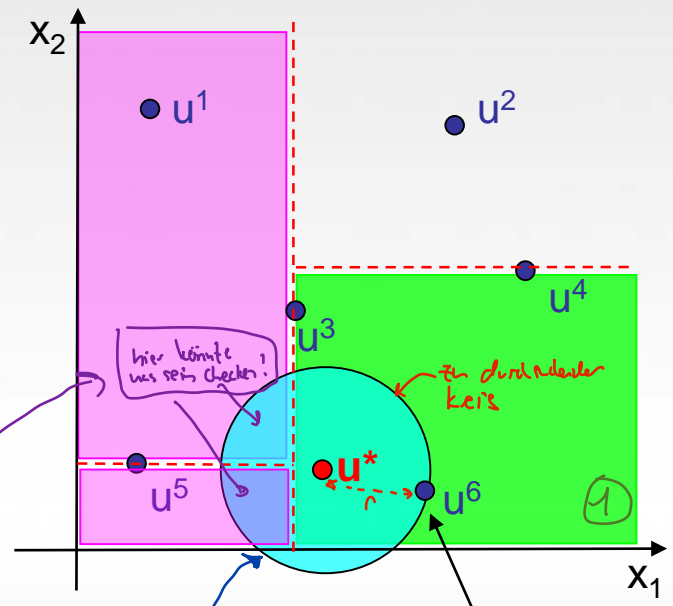
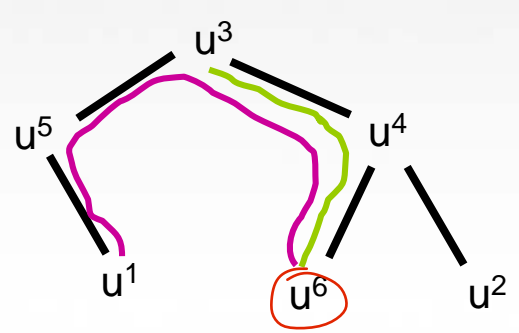
KD-trees (2): Retrieval

KD-Tree Aufspaltvektor

Function $u^{NN} = \text{Query-KD-Tree}(\text{KDT}, u^*)$ // Find NN of query u^* in KD-Tree

- 1) Find first approximation u^{NN1} in leaf region of u^*
- 2) Do backtracking to investigate keys in regions that overlap with sphere $\{x: d(u^*, x) \leq d(u^*, u^{NN1})\}$
- 3) Return true nearest neighbor u^{NN}

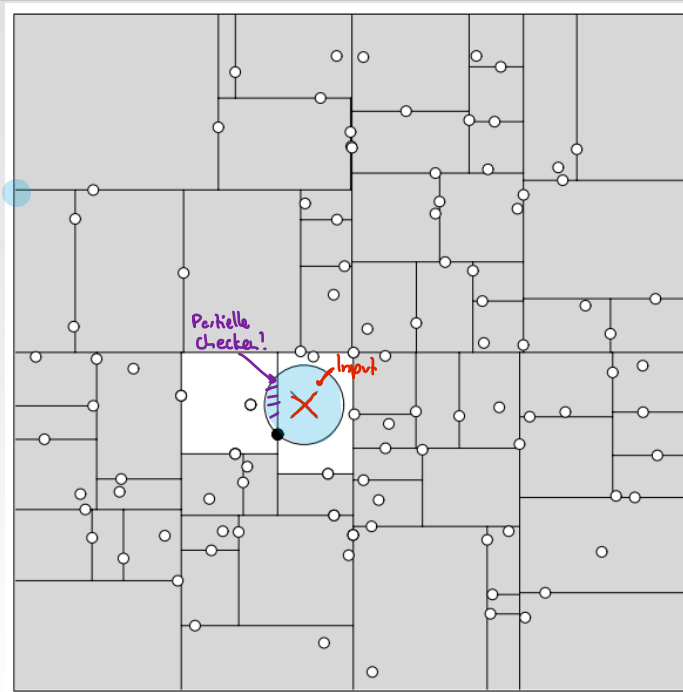
(siehe Approx. für NN)
↪ suche von oben Partitionen in der u^* liegt
↪ ist rechter im linken oder rechten Teilbaum?
 (kleiner oder größer als pivot?)



→ je mehr Nachbarn jeher desto schlechter!
→ Funktioniert falls 1. NN sehr nahe am inputvektor ist

1. Kandidat für NN → $= u^{NN1}$
Problem: $= u^{NN}$
→ geht noch ein NN in benachbarten Partitione
→ falls auch einer existiert muss der im Kreis sich befinden
→ prüfe diese Nachbarnregionen

KD-trees (3): Performance & Limitations



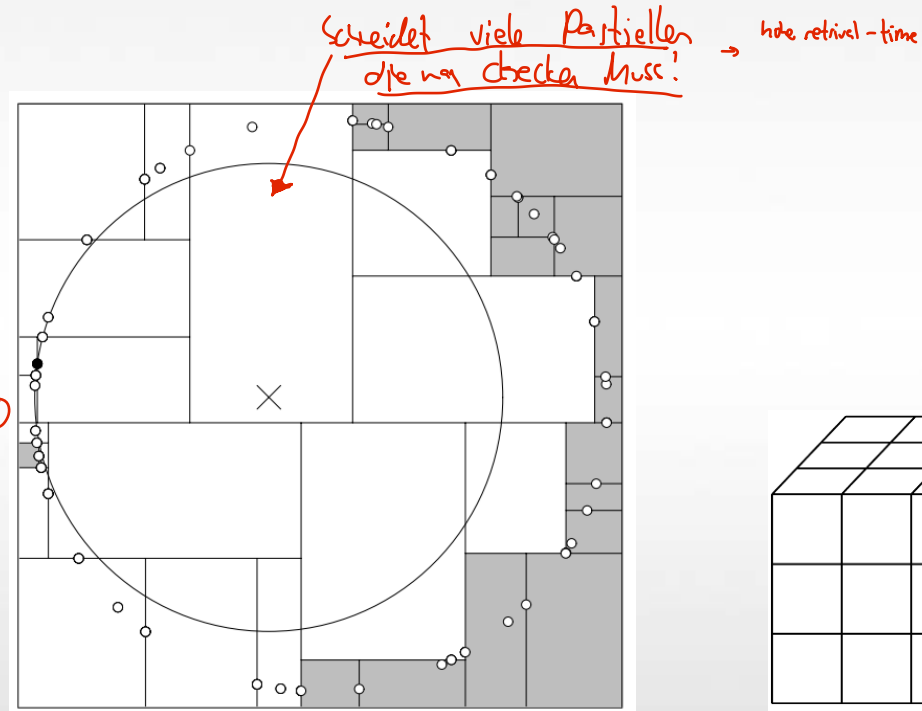
KD-tree algorithm works fine if the first approximative NN is close enough to the query to cut-off rest of the tree

⇒ je mehr Regionen der Kreis schneidet desto schlechter!

- path to leaf $\sim \log N$ steps (for balanced trees)
- constant steps for backtracking (for hypercubic regions)
- query time $t_q \sim \log N$ (for fixed D)

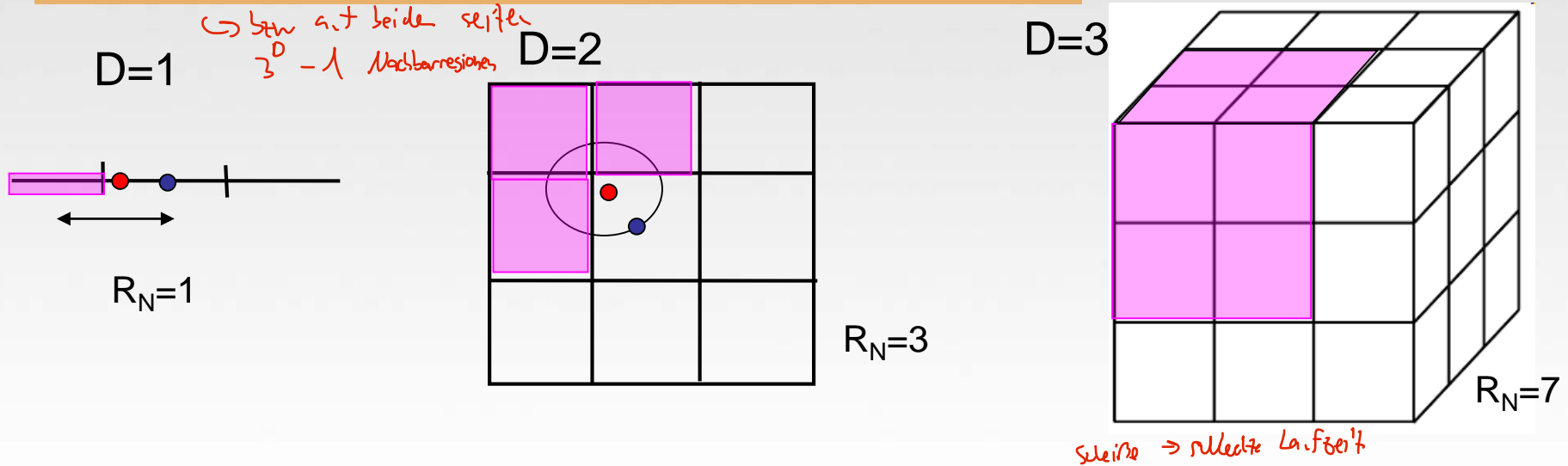
Potential problems:

- **unbalanced trees** *→ dann muss man alles durchlaufen*
- queries not drawn from data distribution *alle liegen außen, ich neh in der mitte... schlechter!*
- **non-hypercubic regions** *(Länglicher stellt quadrat → man schneidet oft!)*
- sparse data vectors
- curse of dimensionality *↳ wird bei viel D sehr kompliziert → testen viele nachbarregionen*



KD-trees (4): Curse of dimensionality

- # backtracking steps depends on $R_N := \#$ neighboring regions
- $R_N = 2^D - 1$ increases exponentially with dimension D



Curse of dimensionality: query time $t_q \sim \log N + 2^D$ exponential in D

- KD-trees are useful only for low-dimensional data, e.g., $1 < D < 20$
- other tree-based partitioning algorithms (KDB, R, SR, SS trees) suffer from similar problems
- there are no known efficient NN algorithms for large D
- next best solution: Approximate NN algorithms...

bzw. bis 50 noch okay

*(sind ähnliche Algorithmen...
liefern nur Näherung)*

Approximate Nearest Neighbor (ANN)

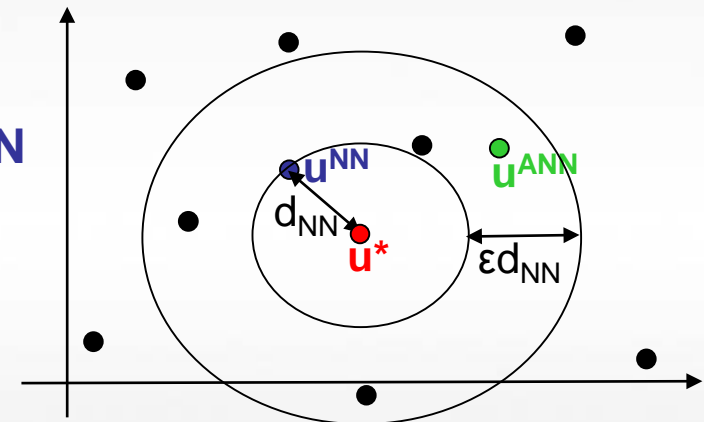
Reasoning: Choosing a particular metrics is a heuristics anyway
 → Therefore, for many applications, **solving ANN may be sufficient**

Many different methods for ANN problem:

- KD (or other) trees without backtracking
- K-means NN (for clustered data)
- Neural associative networks
- Locality sensitive hashing (LSH) (ähnlich efficient)

Related problem with bounded error: ϵ -NN

An algorithm solving ϵ -NN returns an **ANN** that is at most factor $1+\epsilon$ worse than **true NN**



What is actually hashing?

Given: large database, key k
 Goal : find content associated
 with key k

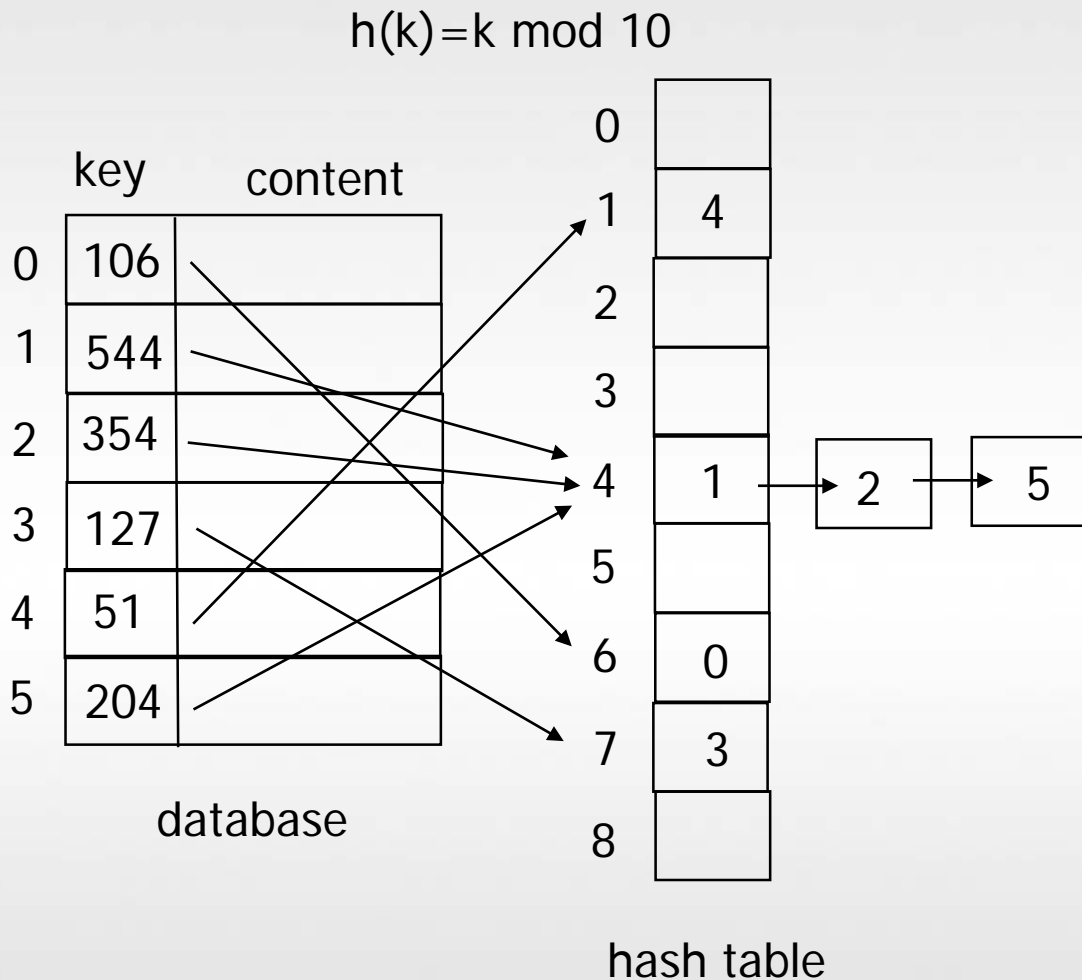
Linear search: $O(n)$

Hashing with $h(k)$: $O(1)$

Desired properties of $h(k)$:

- compute $h(k)$ in $O(1)$
- distribute keys uniformly

→ Not locality-sensitive!



(i) Build-up of index structure

Input: Key vectors $\{u_1, u_2, \dots, u_N\}$ and $l := \# \text{hash tables}$

Output: Hash tables T_1, \dots, T_l

Algorithm:

1) for each $i=1..l$

Generate random hash function $g_i: u_j \rightarrow \text{bucket index}$

2) for each $i=1..l$

for each $j=1..N$

store key vector u_j in bucket $g_i(u)$ of hash table T_i

For example: g_i are chosen randomly from set

$H = \{h_i: u \rightarrow u_i \mid i=1..D\}$ (projections onto i -th component)

(ii) Approximate Nearest-Neighbor Query

Input: Query vector u^* and $K := \#$ nearest neighbors

Output: K (or less) NN vectors

Algorithm:

1) $S = \{ \}$

2) for each $i=1..l$

$S = S + \{\text{vectors in bucket } g_i(u^*) \text{ of table } T_i\}$

3) Return K nearest neighbors of u^* found in set S
(e.g., by using linear search)

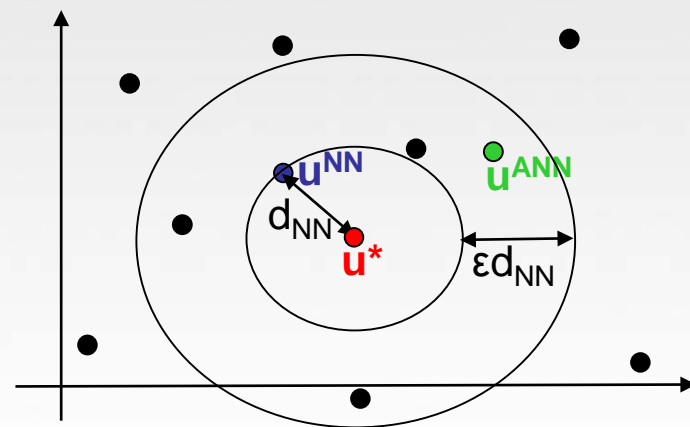
Note: S is typically much(!) smaller than DB size

One can show (Andoni/Indyk '06) for appropriate choices of H (family of hash functions) and l (# of hash tables):

- LSH solves ϵ -NN problem
- LSH is efficient:
 query time $t_q \sim DN^{1/(1+\epsilon)2}$
 space $\sim n^{1+(1+\epsilon)2}$

For example, for $\epsilon=1$:

$$t_q \sim DN^{0.25}$$



- Efficient solution for ANN problem in high-dim spaces
- recent variants
 - Coherency Sensitive Hashing (Korman/Avidan, 2011)
 - Complementary Hashing (Xu et al, 2011)
 - ...

Thank you for your attention!