

Case Study 2: Document Retrieval

Task Description: Finding Similar Documents

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

April 14th, 2015

©Emily Fox 2015

1

Task 1: Find Similar Documents

■ To begin...

- **Input:** Query article ✗
- **Output:** Set of k similar articles



✗



©Emily Fox 2015

2

k-Nearest Neighbor

- Articles $X = \{x^1, \dots, x^N\}$, $x^i \in \mathbb{R}^d$
- Query: $x \in \mathbb{R}^d$
- k-NN
 - Goal: find k articles in X closest to x

□ Formulation:

$$X^{NN} = \{x^{NN_1}, \dots, x^{NN_k}\} \subseteq X$$

$$\text{s.t. } \forall x^i \in X \setminus X^{NN}$$

$$d(x^i, x) \geq \max_{x^{NN_i} \in X^{NN}} d(x^{NN_i}, x)$$

©Emily Fox 2015

3

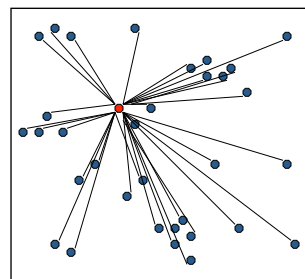
Issues with Search Techniques

- Naïve approach:

Brute force search

- Given a query point x
- Scan through each point x^i
- $O(N)$ distance computations per 1-NN query!
- $O(N \log k)$ per k -NN query!

↑ keep priority queue
of top k
+ inserting into
queue is $\log k$



33 Distance Computations

- What if N is huge???
(and many queries)

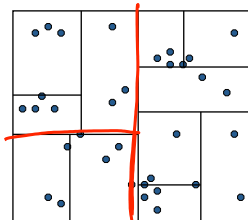
©Emily Fox 2015

4

KD-Trees

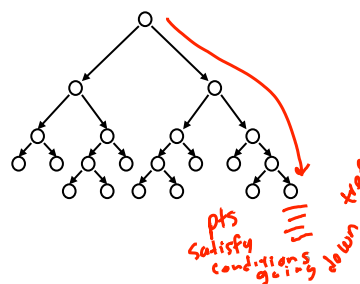
Smarter approach: **kd-trees**

- Structured organization of documents
 - Recursively partitions points into axis aligned boxes.
- Enables more efficient pruning of search space
 - Examine nearby points first.
 - Ignore any points that are further than the nearest point found so far.



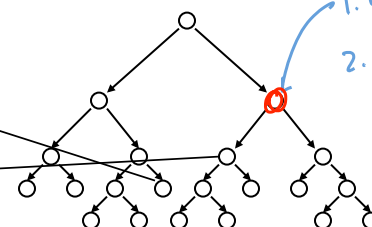
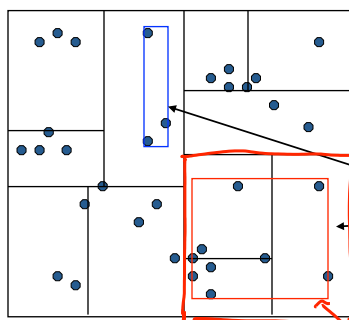
■ **kd-trees** work “well” in “low-medium” dimensions d

- We'll get back to this...



©Emily Fox 2015

KD-Tree Construction



3. Keep one additional piece of information at each node:
- ★ □ The (tight) bounds of the points at or below this node.

©Emily Fox 2015

6

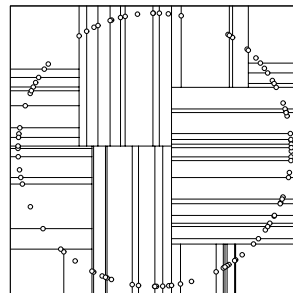
KD-Tree Construction

- Use heuristics to make splitting decisions:
- Which dimension do we split along?
widest (or alternate)
- Which value do we split at?
median of chosen split dim (or center)
- When do we stop?
fewer than m pt left
or
box hits minimum width

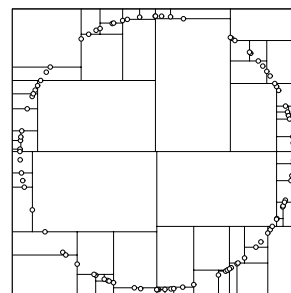
©Emily Fox 2015

7

Many heuristics...



median heuristic

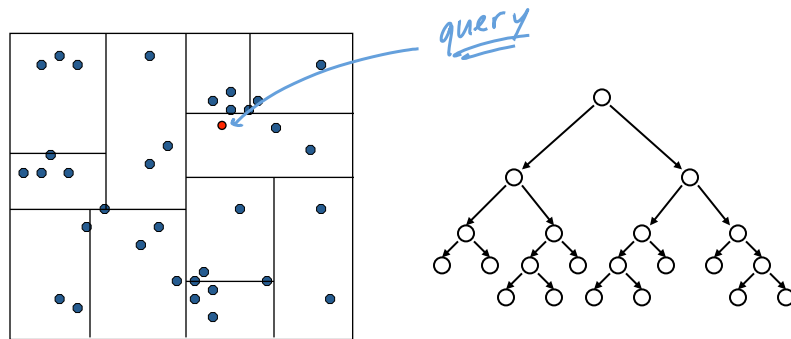


center-of-range heuristic

©Emily Fox 2015

8

Nearest Neighbor with KD Trees

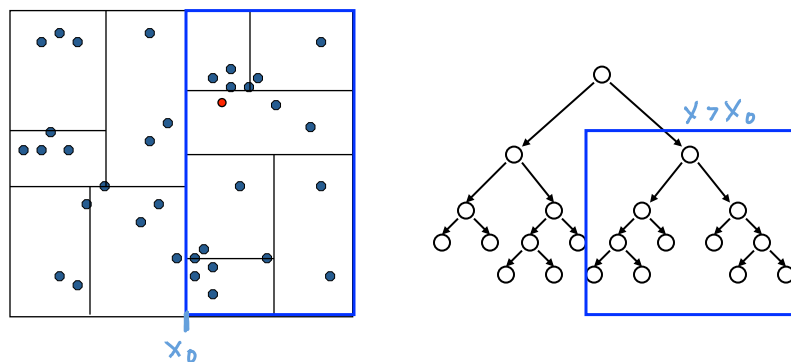


- Traverse the tree looking for the nearest neighbor of the query point.

©Emily Fox 2015

9

Nearest Neighbor with KD Trees

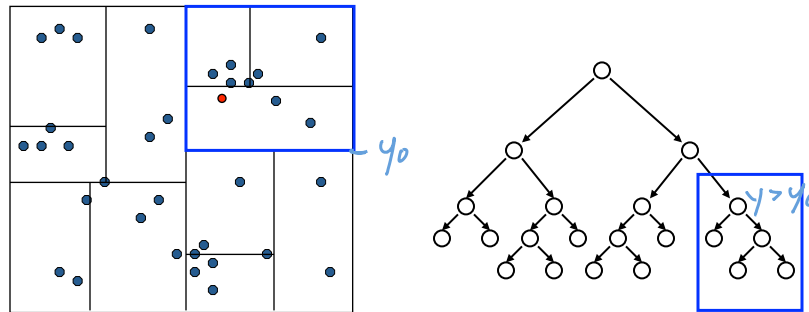


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2015

10

Nearest Neighbor with KD Trees

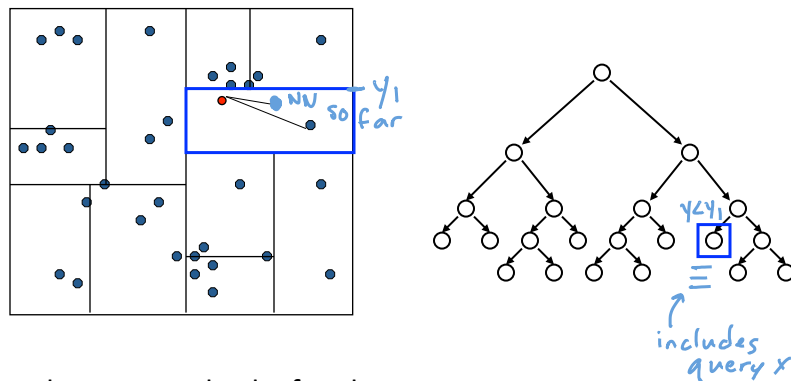


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2015

11

Nearest Neighbor with KD Trees

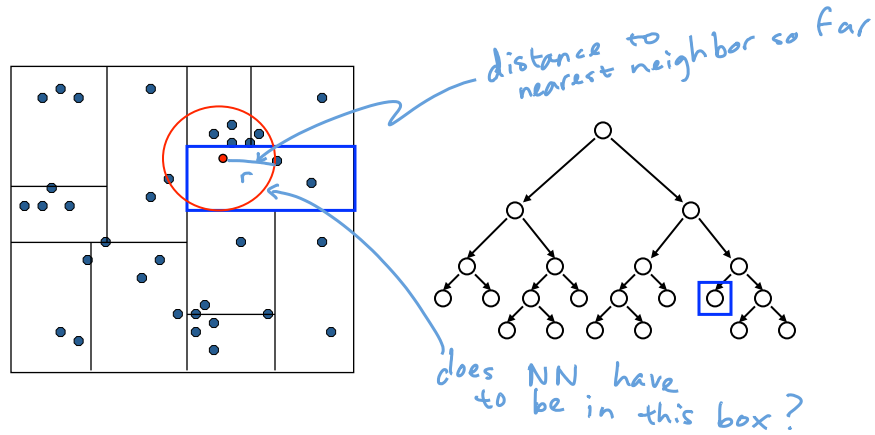


- When we reach a leaf node:
 - Compute the distance to each point in the node.

©Emily Fox 2015

12

Nearest Neighbor with KD Trees



- When we reach a leaf node:

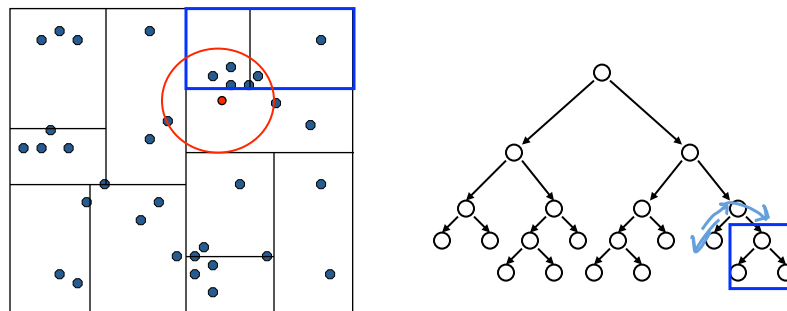
□ Compute the distance to each point in the node.

NO

©Emily Fox 2015

13

Nearest Neighbor with KD Trees

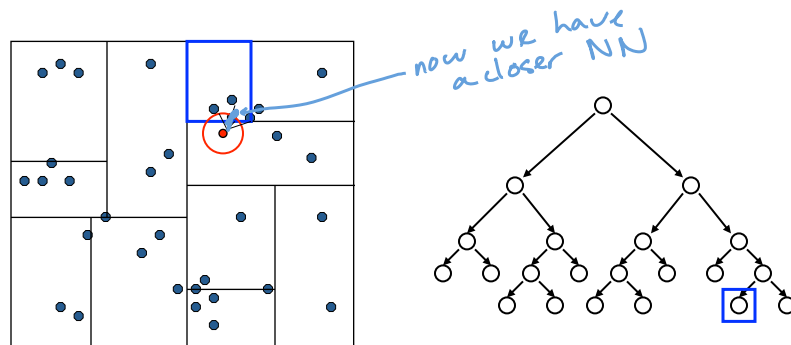


- Then backtrack and try the other branch at each node visited

©Emily Fox 2015

14

Nearest Neighbor with KD Trees

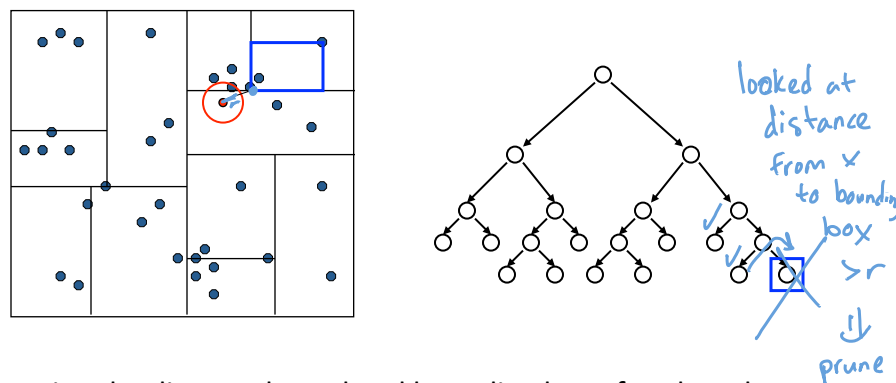


- Each time a new closest node is found, update the distance bound

©Emily Fox 2015

15

Nearest Neighbor with KD Trees

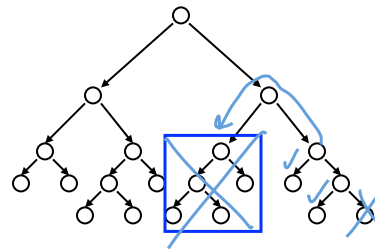
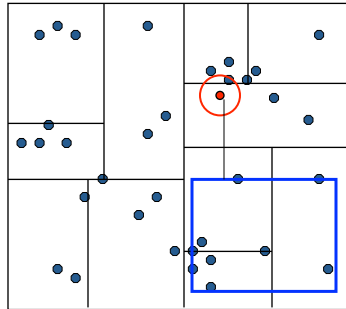


- Using the distance bound and bounding box of each node:
 - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2015

16

Nearest Neighbor with KD Trees

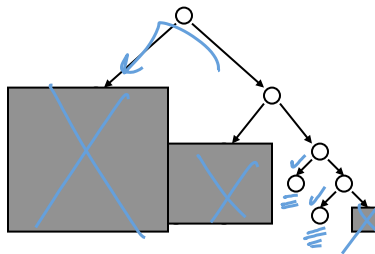
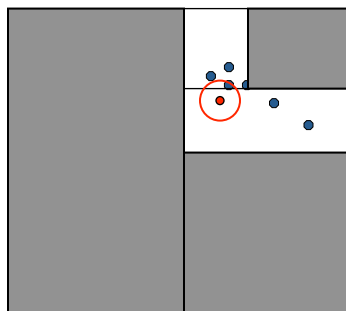


- Using the distance bound and bounding box of each node:
 - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2015

17

Nearest Neighbor with KD Trees



- Using the distance bound and bounding box of each node:
 - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2015

18

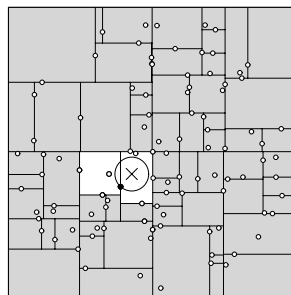
Complexity

- For (nearly) balanced, binary trees...
- Construction
 - Size: $2N-1 \rightarrow O(N)$
 - Depth: $O(\log N)$
 - Median + send points left right: $O(N)$ at every tree level (smart)
 - Construction time: $O(N \log N)$
- 1-NN query
 - Traverse down tree to starting point: $O(\log N)$
 - Maximum backtrack and traverse: $O(N)$ worst case
 - Complexity range: $O(\log N) \rightarrow O(N)$
- Under some assumptions on distribution of points, we get $O(\log N)$ but exponential in d (see citations in reading)

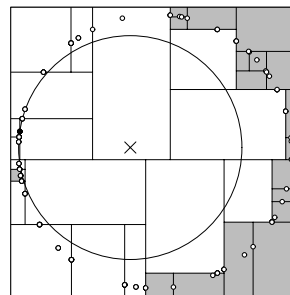
©Emily Fox 2015

19

Complexity



pruned many
(closer to $O(\log N)$)



pruned few
(closer to $O(N)$)

©Emily Fox 2015

20

Complexity for N Queries

- Ask for nearest neighbor to each document

N queries

- Brute force 1-NN:

$O(N^2)$

- kd-trees:

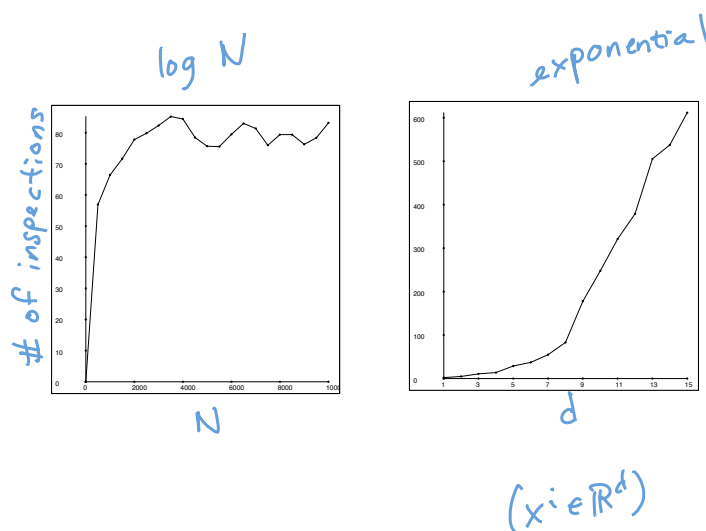
$O(N \log N) \rightarrow O(N^2)$

↑
potentially large
savings

©Emily Fox 2015

21

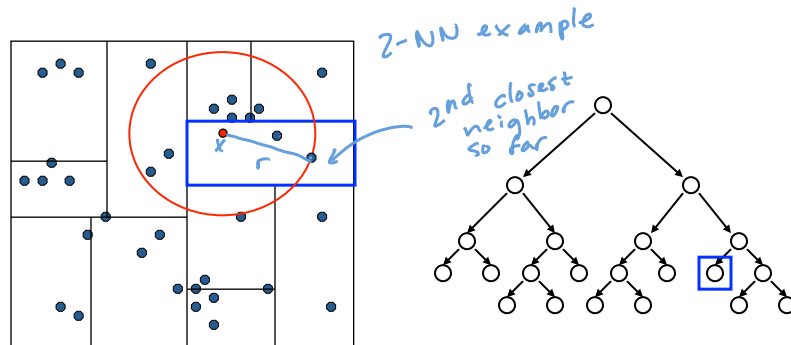
Inspections vs. N and d



©Emily Fox 2015

22

K-NN with KD Trees

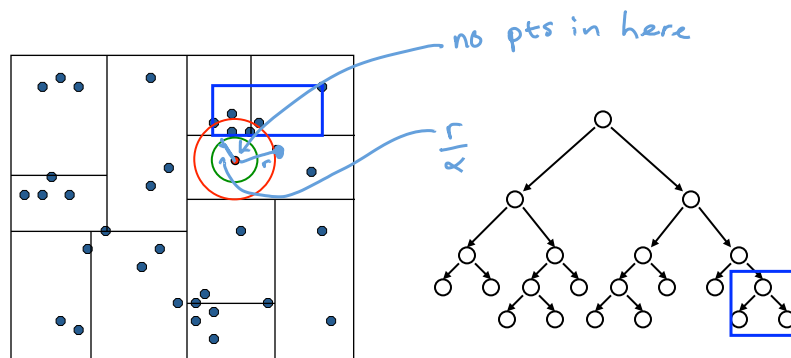


- Exactly the same algorithm, but maintain distance as distance to furthest of current k nearest neighbors
- Complexity is: $O(k \log N)$

©Emily Fox 2015

23

Approximate K-NN with KD Trees



- Before:** Prune when distance to bounding box $> r$
- Now:** Prune when distance to bounding box $> \frac{r}{\alpha}$ $\alpha > 1$
- Will prune more than allowed, but can guarantee that if we return a neighbor at distance r , then there is no neighbor closer than r/α .
- In practice this bound is loose...Can be closer to optimal.
- Saves lots of search time at little cost in quality of nearest neighbor.

©Emily Fox 2015

24

Wrapping Up – Important Points

kd-trees

- Tons of variants
 - On construction of trees (heuristics for splitting, stopping, representing branches...)
 - Other representational data structures for fast NN search (e.g., ball trees,...)

Nearest Neighbor Search

- Distance metric and data representation are crucial to answer returned

★ For both...

- High dimensional spaces are hard! *large d*
 - Number of kd-tree searches can be exponential in dimension
 - Rule of thumb... $N \gg 2^d$... Typically useless.
 - Distances are sensitive to irrelevant features
 - Most dimensions are just noise → Everything equidistant (i.e., everything is far away)
 - Need technique to learn what features are important for your task

©Emily Fox 2015

25

What you need to know

- Document retrieval task
 - Document representation (bag of words)
 - tf-idf
- Nearest neighbor search
 - Formulation
 - Different distance metrics and sensitivity to choice
 - Challenges with large N
- kd-trees for nearest neighbor search
 - Construction of tree
 - NN search algorithm using tree
 - Complexity of construction and query
 - Challenges with large d

©Emily Fox 2015

26

Acknowledgment

- This lecture contains some material from Andrew Moore's excellent collection of ML tutorials:
 - <http://www.cs.cmu.edu/~awm/tutorials>
- In particular, see:
 - http://grist.caltech.edu/sc4devo/.../files/sc4devo_scalable_datamining.ppt

©Emily Fox 2015

27

Case Study 2: Document Retrieval

Locality-Sensitive Hashing Random Projections for NN Search

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox
April 14th, 2015

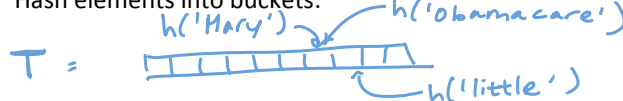
©Emily Fox 2015

28

Using Hashing to Find Neighbors

- KD-trees are cool, but...
 - Non-trivial to implement efficiently
 - Problems with high-dimensional data
- Approximate neighbor finding...
 - Don't find exact neighbor, but that's OK for many apps, especially with Big Data

- What if we could use hash functions: $h: \mathcal{X} \rightarrow \{1, \dots, m\}$
 - Hash elements into buckets:



typical
version
where we
keep list
of words
in every bin

- Look for neighbors that fall in same bucket as x :

$h(x)=i$, for all $y \in T[h(x)=i]$ look for neighbors there

- But, by design...

$$P(h(x)=h(x')) = \frac{1}{m} \quad \forall x'$$

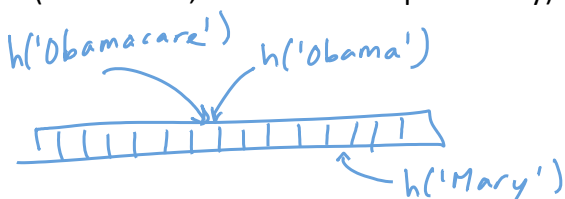
even if $d(x, x')$ is low $\nRightarrow h(x) \neq h(x')$

©Emily Fox 2015

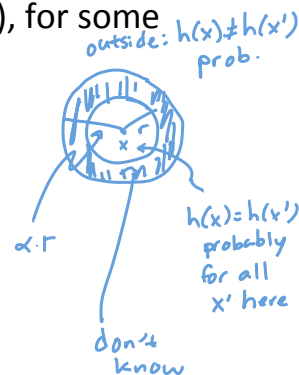
29

Locality Sensitive Hashing (LSH)

- A LSH function h satisfies (for example), for some similarity function d , for $r > 0$, $\alpha > 1$:
 - $d(x, x') \leq r$, then $P(h(x)=h(x'))$ is high
 - $d(x, x') > \alpha \cdot r$, then $P(h(x)=h(x'))$ is low
 - (in between, not sure about probability)



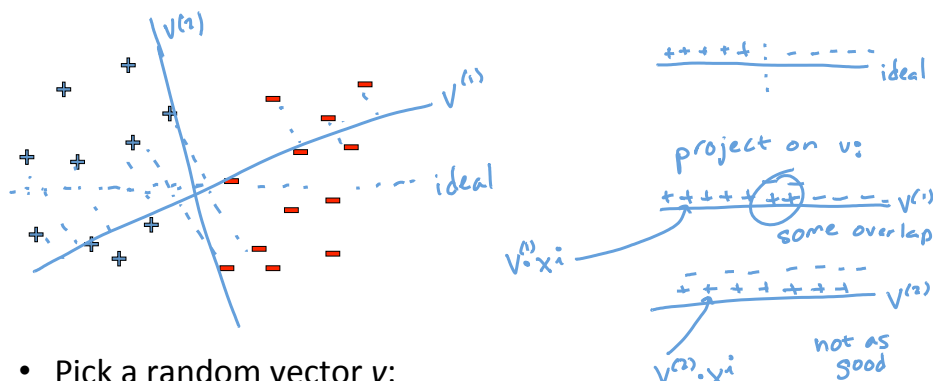
now look for pts hashing to same bin



©Emily Fox 2015

30

Random Projection Illustration



- Pick a random vector v :
 - Independent Gaussian coordinates

$$v_i \stackrel{iid}{\sim} N(0,1)$$

define d -dim vector $[v_1, \dots, v_d]$

- Preserves separability for most vectors
 - Gets better with more random vectors

$$V^{(1)} = [v_1^{(1)}, \dots, v_d^{(1)}]$$

©Emily Fox 2015

31

Multiple Random Projections: Approximating Dot Products

- Pick m random vectors $v^{(i)}$:
 - Independent Gaussian coordinates
- Approximate dot products:
 - Cheaper, e.g., learn in smaller m dimensional space
- Only need logarithmic number of dimensions!
 - N data points, approximate dot-product within $\epsilon > 0$:

$$x \cdot y \approx \frac{1}{m} \phi(x) \cdot \phi(y) = \frac{1}{m} \sum_{i=1}^m (v^{(i)} \cdot x) \cdot (v^{(i)} \cdot y)$$

$$m = O\left(\frac{\log N}{\epsilon^2}\right)$$

- But all sparsity is lost

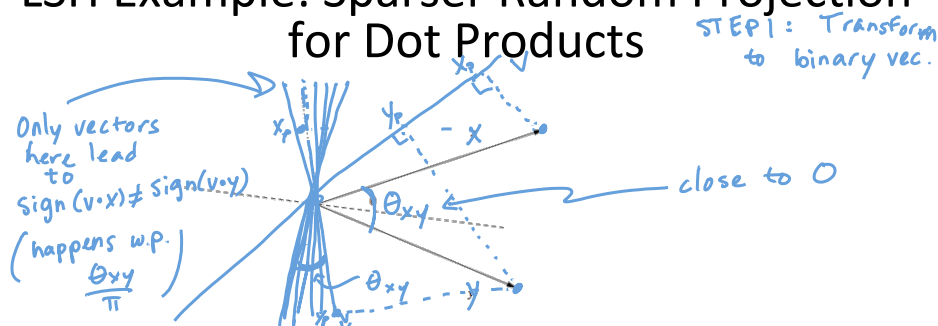
$$v^{(i)} \text{ are dense wpl} \Rightarrow v^{(i)} \cdot x \neq 0 \text{ wpl}$$

If I have big data
 $N \rightarrow$ very large,
 but only need $\log N$
 random vecs.
 even if x is sparse,
 $\phi(x)$ is not sparse

©Emily Fox 2015

32

LSH Example: Sparser Random Projection for Dot Products



- Pick random vectors $v^{(i)} \sim N(0, I)$ ← m in total
- Simple 0/1 projection: $\phi_i(x) = \begin{cases} 1 & \text{if } \text{sign}(v^{(i)} \cdot x) \geq 0 \\ 0 & \text{if } \text{sign}(v^{(i)} \cdot x) < 0 \end{cases}$
- Now, each vector is approximated by a bit-vector
 $\phi(x) = (0, 0, 1, 0, 1, 1, 1, 0)$
 m -dim
- Dot-product approximation:
 $\frac{x \cdot y}{\|x\| \|y\|} = \cos \theta_{xy} \approx \cos \left(\frac{\pi}{m} \text{HamDist}(\phi(x), \phi(y)) \right)$

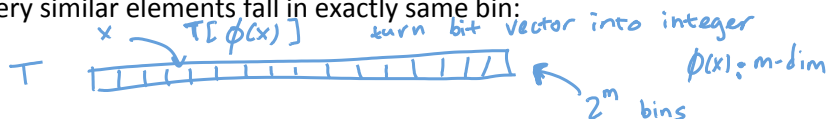
©Emily Fox 2015

33

LSH for Approximate Neighbor Finding

STEP 2: Use $\phi(x)$ as index in hash table

- Very similar elements fall in exactly same bin:



- And, nearby bins are also nearby:

in terms of Hamming dist.

- Simple neighbor finding with LSH:
 - For bins b of increasing hamming distance to $\phi(x)$:
 - Look for neighbors of x in bin b

$\forall y$ in $T[b]$, compute $d(x, y)$
 keep closest

- Stop when run out of time

- Pick m such that $N/2^m$ is "smallish" (in practice)

©Emily Fox 2015

similar $x \rightarrow \phi(x) = (0, 1, 1, 0)$
 $y \rightarrow \phi(y) = (0, 1, 1, 0)$ or $(1, 1, 1, 0)$
 but not $(1, 0, 0, 1)$

similarity, not dist.
 $S(x, y) = \cos \theta_{xy}$
 "cosine similarity"
 $= \cos \left(\frac{\pi}{m} \text{HamDist}(\phi(x), \phi(y)) \right)$
 $S(x, y)$ high
 $\Rightarrow \cos \theta_{xy} \approx 1$
 $\Rightarrow \theta_{xy} \approx 0$
 $\Rightarrow \text{HamDist} \approx 0$
 $\Rightarrow x, y$ same bin

What you need to know

- **Locality-Sensitive Hashing (LSH):** nearby points hash to the same or nearby bins
- LSH uses **random projections**
 - Only $O(\log N/\epsilon^2)$ vectors needed
 - But vectors and results are **not sparse**
- **Use LSH for nearest neighbors by mapping elements into bins**
 - Bin index is defined by bit vector from LSH
 - Find nearest neighbors by going through bins
- **Hash kernels:**
 - Sparse representation for feature vectors
 - Very simple, use two hash functions
 - Can even use one hash function, and take least significant bit to define ξ
 - Quickly generate projection $\phi(x)$
 - **Learn in projected space**