

Case Study 1: Estimating Click Probabilities

Tackling an Unknown Number of Features with Sketching

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Carlos Guestrin
January 22nd, 2013

©Carlos Guestrin 2013

1

Sketching Counts

- Bloom Filter is super cool, but not what we need...
 - We don't just care about whether a feature existed before, but to keep track of counts of occurrences of features!
- Recall Perceptron update:
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$
- Must keep track of counts of each feature (weighed by $y^{(t)}$):
 - E.g., with sparse data, for each non-zero dimension i in $\mathbf{x}^{(t)}$:
- Can we generalize the Bloom Filter?

©Carlos Guestrin 2013

2

Count-Min Sketch: single vector

- Simpler problem: Count how many times you see each string
- Single hash function:
 - Keep Count vector of length m
 - every time see string i :

$$Count[h(i)] \leftarrow Count[h(i)] + 1$$

- Again, collisions could be a problem:
 - a_i is the count of element i :

©Carlos Guestrin 2013

3

Count-Min Sketch: general case

- Keep d by m Count matrix
- d hash functions:
 - Just like in Bloom Filter, decrease errors with multiple hashes
 - Every time see string i :

$$\forall j \in \{1, \dots, d\} : Count[j, h(i)] \leftarrow Count[j, h(i)] + 1$$

©Carlos Guestrin 2013

4

Querying the Count-Min Sketch

$\forall j \in \{1, \dots, d\} : \text{Count}[j, h(i)] \leftarrow \text{Count}[j, h(i)] + 1$

- Query $Q(i)$?
 - What is in $\text{Count}[j, k]$?

 - Thus:

 - Return:

©Carlos Guestrin 2013

5

Analysis of Count-Min Sketch

$$\hat{a}_i = \min_j \text{Count}[j, h(i)] \geq a_i$$

- Set:
 $m = \left\lceil \frac{e}{\epsilon} \right\rceil \quad d = \left\lceil \ln \frac{1}{\delta} \right\rceil$

- Then, after seeing n elements:

$$\hat{a}_i \leq a_i + \epsilon n$$

- With probability at least $1 - \delta$

©Carlos Guestrin 2013

6

Proof of Count-Min for Point Query with Positive Counts: Part 1 – Expected Bound

- $I_{i,j,k}$ = indicator that i & k collide on hash j :
- Bounding expected value:
- $X_{i,j}$ = total colliding mass on estimate of count of i in hash j :
- Bounding colliding mass:
- Thus, estimate from each hash function is close in expectation

©Carlos Guestrin 2013

7

Proof of Count-Min for Point Query with Positive Counts: Part 2 – High Probability Bounds

- What we know: $Count[j, h(i)] = a_i + X_{i,j}$ $E[X_{i,j}] \leq \frac{\epsilon}{e}n$
- Markov inequality: For z_1, \dots, z_k positive iid random variables
$$P(\forall z_i : z_i > \alpha E[z_i]) < \alpha^{-k}$$
- Applying to the Count-Min sketch:

©Carlos Guestrin 2013

8

But Our Updates may be positive or Negative

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$

- Count-Min sketch for positive & negative case
 - a_i no longer necessarily positive
- Update the same: Observe change Δ_i to element i :
 $\forall j \in \{1, \dots, d\} : \text{Count}[j, h(i)] \leftarrow \text{Count}[j, h(i)] + \Delta_i$
 - Each $\text{Count}[j, h(i)]$ no longer an upper bound on a_i
- How do we make a prediction?

- Bound: $|\hat{a}_i - a_i| \leq 3\epsilon \|\mathbf{a}\|_1$
 - With probability at least $1 - \delta^{1/4}$, where $\|\mathbf{a}\| = \sum_i |a_i|$

©Carlos Guestrin 2013

9

Finally, Sketching for Perceptron

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$

- Never need to know size of vocabulary!
- Make a mistake, update Count-Min matrix:

- Making a prediction:

- Scales to huge problems, great practical implications... More next time

©Carlos Guestrin 2013

10

What you need to know

- Hash functions
- Bloom filter
 - Test membership with some false positives, but very small number of bits per element
- Count-Min sketch
 - Positive counts: upper bound with nice rates of convergence
 - General case
- Application to Perceptron Learning and Prediction

Case Study 2: Document Retrieval

Task Description: Finding Similar Documents

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Emily Fox

January 22nd, 2013

Document Retrieval

- **Goal:** Retrieve documents of interest
- **Challenges:**
 - Tons of articles out there
 - How should we measure similarity?



©Emily Fox 2013

13

Task 1: Find Similar Documents

- **To begin...**
 - **Input:** Query article
 - **Output:** Set of k similar articles



©Emily Fox 2013

14

Document Representation

- Bag of words model



1-Nearest Neighbor

- Articles
- Query:
- 1-NN
 - Goal:
 - Formulation:

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:
 - Formulation:

©Emily Fox 2013

17

Distance Metrics – Euclidean

$$d(u, v) = \sqrt{\sum_{i=1}^d (u_i - v_i)^2}$$

Or, more generally, $d(u, v) = \sqrt{\sum_{i=1}^d \sigma_i^2 (u_i - v_i)^2}$

Equivalently,

$$d(u, v) = \sqrt{(u - v)' \Sigma (u - v)}$$

where $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & \sigma_d^2 \end{bmatrix}$

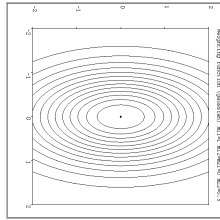
Other Metrics...

- Mahalanobis, Rank-based, Correlation-based, cosine similarity...

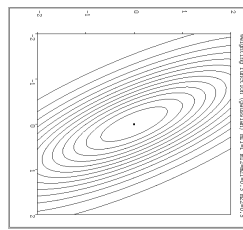
©Emily Fox 2013

18

Notable Distance Metrics (and their level sets)

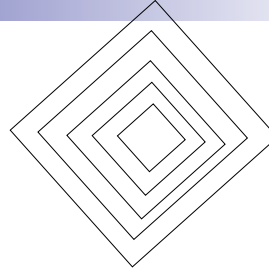


Scaled Euclidian (L_2)

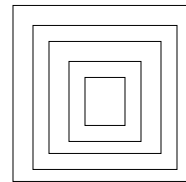


Mahalanobis

(Σ is general sym pos def matrix,
on previous slide = diagonal)



L_1 norm (absolute)



$L1$ (max) norm

©Emily Fox 2013

19

Euclidean Distance + Document Retrieval

- Recall distance metric

$$d(u, v) = \sqrt{\sum_{i=1}^d (u_i - v_i)^2}$$

- What if each document were α times longer?
 - Scale word count vectors
 - What happens to measure of similarity?
- Good to normalize vectors

©Emily Fox 2013

20

Issues with Document Representation

- Words counts are **bad** for standard similarity metrics



- Term Frequency – Inverse Document Frequency (tf-idf)
 - Increase importance of rare words

TF-IDF

- Term frequency:

$$tf(t, d) =$$

- Could also use $\{0, 1\}, 1 + \log f(t, d), \dots$

- Inverse document frequency:

$$idf(t, D) =$$

- tf-idf:

$$tfidf(t, d, D) =$$

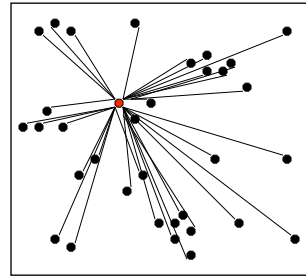
- High for document d with high frequency of term t (high “term frequency”) and few documents containing term t in the corpus (high “inverse doc frequency”)

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!



33 Distance Computations

■ What if N is huge???

(and many queries)

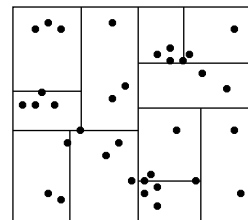
©Emily Fox 2013

23

KD-Trees

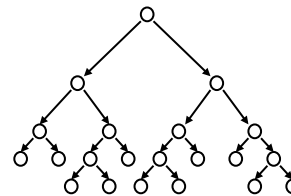
■ Smarter approach: *kd-trees*

- 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

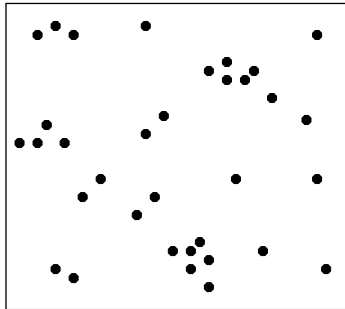
- We’ll get back to this...



©Emily Fox 2013

24

KD-Tree Construction



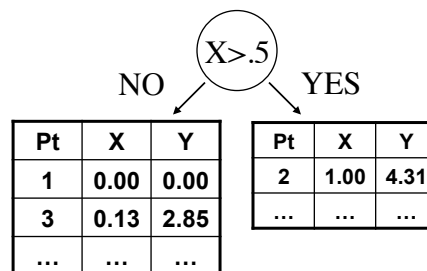
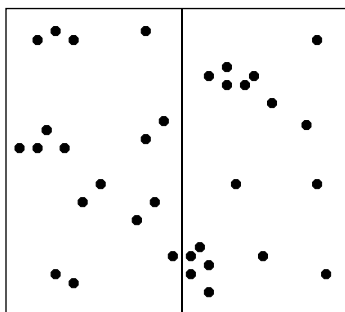
Pt	X	Y
1	0.00	0.00
2	1.00	4.31
3	0.13	2.85
...

- Start with a list of d -dimensional points.

©Emily Fox 2013

25

KD-Tree Construction

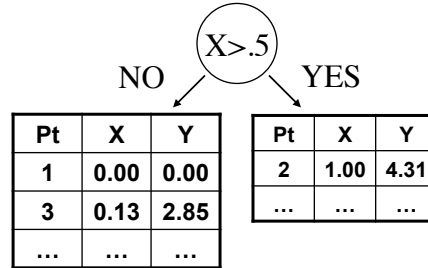
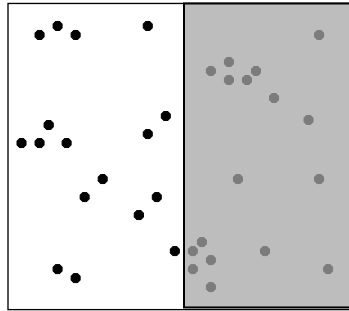


- Split the points into 2 groups by:
 - Choosing dimension d_j and value V (methods to be discussed...)
 - Separating the points into $x_{d_j}^i > V$ and $x_{d_j}^i \leq V$.

©Emily Fox 2013

26

KD-Tree Construction

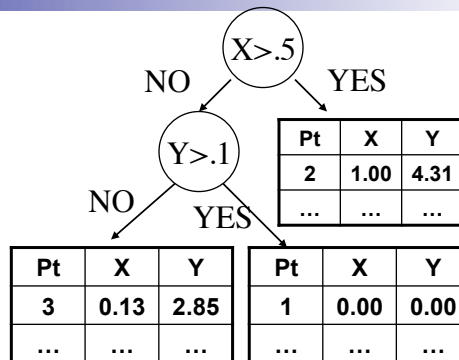
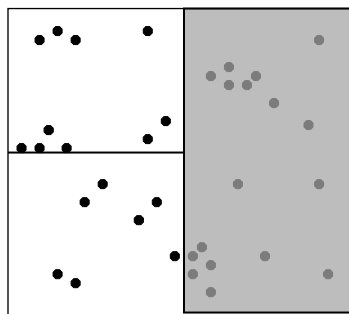


- Consider each group separately and possibly split again (along same/different dimension).
 - Stopping criterion to be discussed...

©Emily Fox 2013

27

KD-Tree Construction

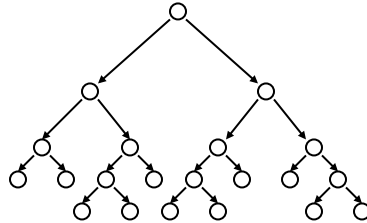
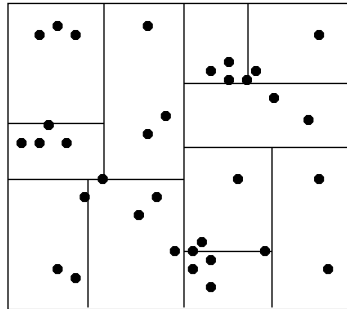


- Consider each group separately and possibly split again (along same/different dimension).
 - Stopping criterion to be discussed...

©Emily Fox 2013

28

KD-Tree Construction

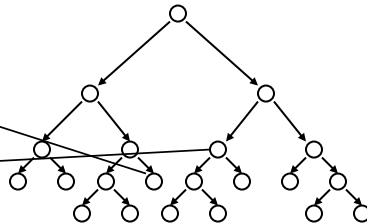
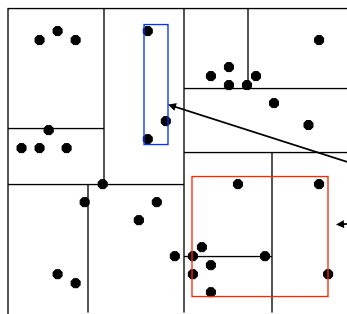


- Continue splitting points in each set
 - creates a binary tree structure
- Each leaf node contains a list of points

©Emily Fox 2013

29

KD-Tree Construction



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

©Emily Fox 2013

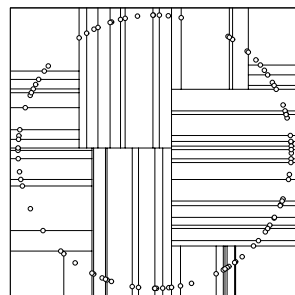
30

KD-Tree Construction

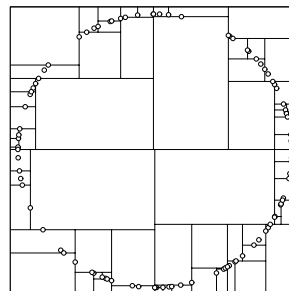
Use heuristics to make splitting decisions:

- Which dimension do we split along?
- Which value do we split at?
- When do we stop?

Many heuristics...

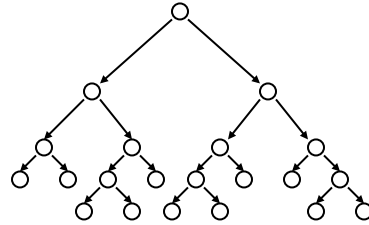
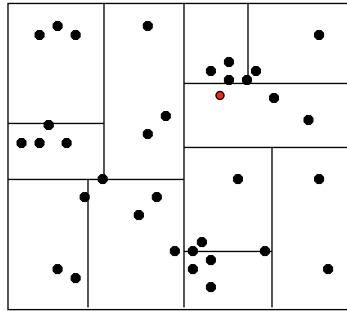


median heuristic



center-of-range heuristic

Nearest Neighbor with KD Trees

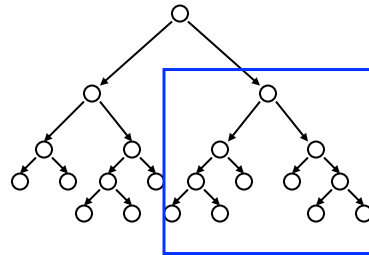
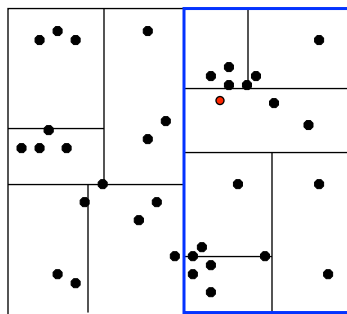


- Traverse the tree looking for the nearest neighbor of the query point.

©Emily Fox 2013

33

Nearest Neighbor with KD Trees

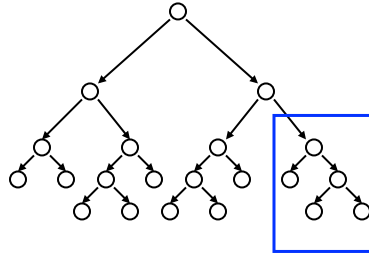
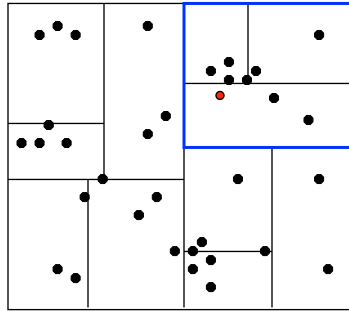


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2013

34

Nearest Neighbor with KD Trees

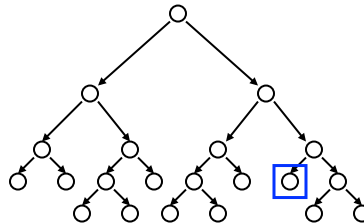
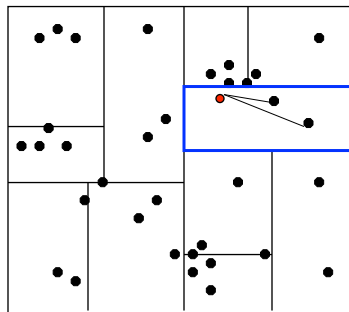


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2013

35

Nearest Neighbor with KD Trees

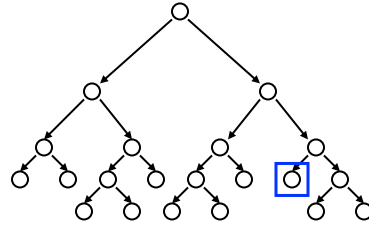
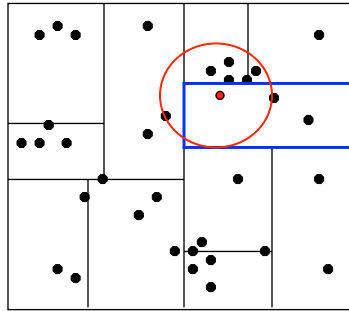


- When we reach a leaf node:
 - Compute the distance to each point in the node.

©Emily Fox 2013

36

Nearest Neighbor with KD Trees

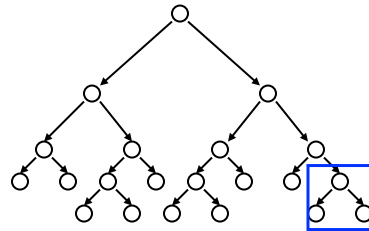
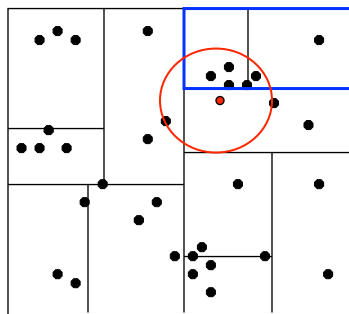


- When we reach a leaf node:
 - Compute the distance to each point in the node.

©Emily Fox 2013

37

Nearest Neighbor with KD Trees

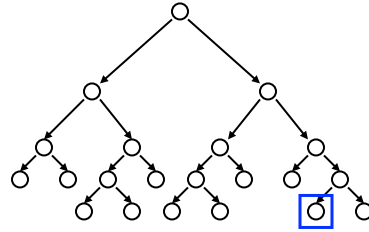
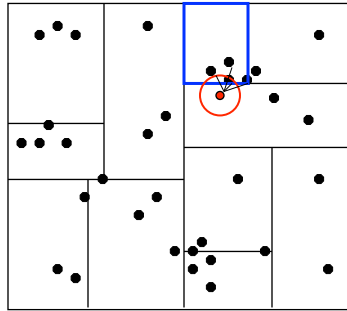


- Then backtrack and try the other branch at each node visited

©Emily Fox 2013

38

Nearest Neighbor with KD Trees

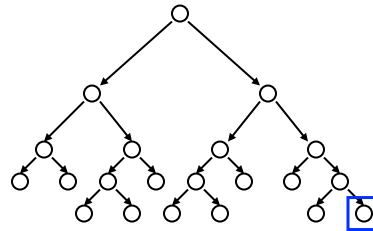
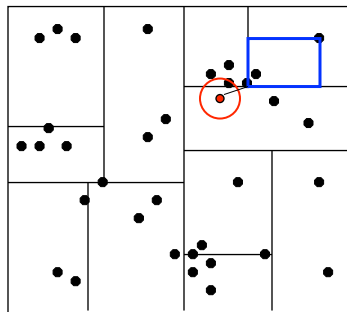


- Each time a new closest node is found, update the distance bound

©Emily Fox 2013

39

Nearest Neighbor with KD Trees

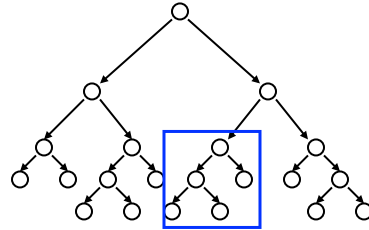
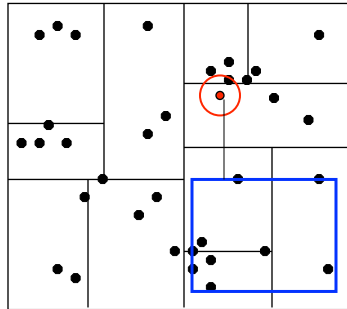


- Using the distance bound and bounding box of each node:
 - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2013

40

Nearest Neighbor with KD Trees

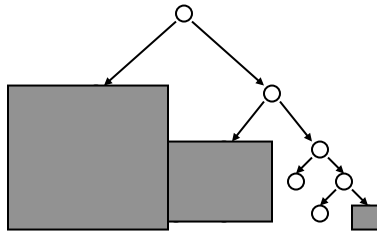
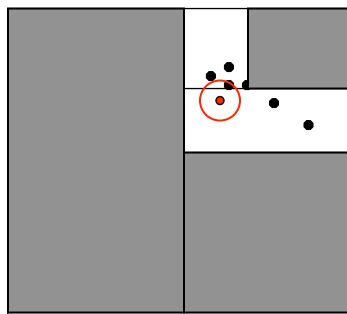


- Using the distance bound and bounding box of each node:
 - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2013

41

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 2013

42

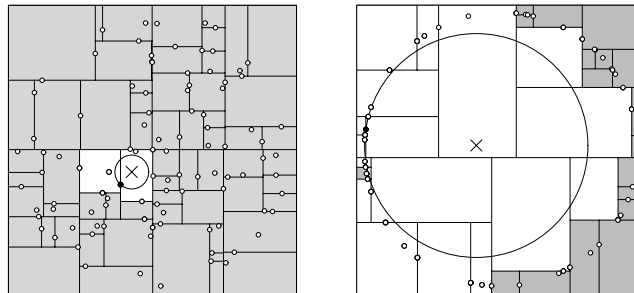
Complexity

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

©Emily Fox 2013

43

Complexity



©Emily Fox 2013

44

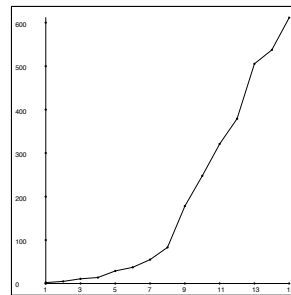
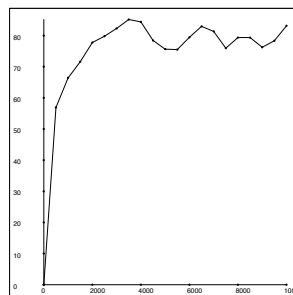
Complexity for N Queries

- Ask for nearest neighbor to each document
- Brute force 1-NN:
- kd-trees:

©Emily Fox 2013

45

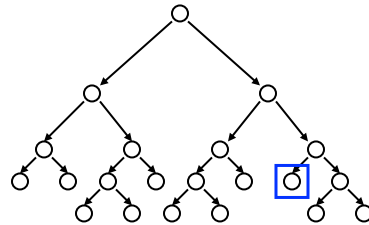
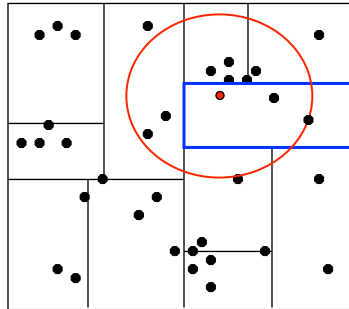
Inspections vs. N and d



©Emily Fox 2013

46

K-NN with KD Trees

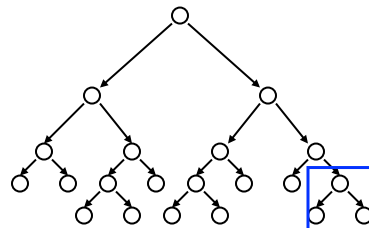
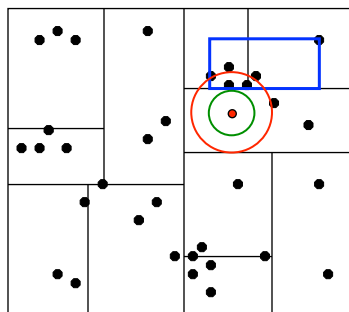


- Exactly the same algorithm, but maintain distance as distance to furthest of current k nearest neighbors
- Complexity is:

©Emily Fox 2013

47

Approximate K-NN with KD Trees



- **Before:** Prune when distance to bounding box $>$
- **Now:** Prune when distance to bounding box $>$
- 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 2013

48

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!
 - 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 2013

49

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 2013

50

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