

## Case Study 1: Estimating Click Probabilities

# Tackling an Unknown Number of Features with Sketching

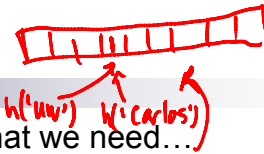
Machine Learning/Statistics for Big Data  
CSE599C1/STAT592, University of Washington

Emily Fox ~~Carlos Guestrin~~  
January 22<sup>nd</sup>, 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  $\mathbf{x}^{(t)}$ :

if mistake +  $\forall x_i^{(t)} \neq 0$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + y^{(t)} x_i^{(t)}$$

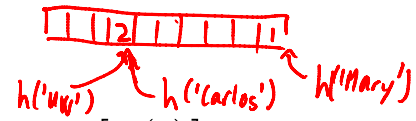
- Can we generalize the Bloom Filter?

©Carlos Guestrin 2013

2

# Count-Min Sketch: single vector

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



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

See 'Carlos' and 'uw'  $\Rightarrow \text{Count}[4] = 2$      $Q('Carlos') = \text{Count}[4] = 2 > 1$

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

$$\text{Count}[j] = \sum_{i: h(i)=j} a_i$$

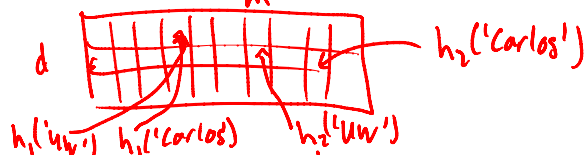
*over-est.*

*true counts*

$$Q(i) \rightarrow \text{return } \hat{a}_i = \text{Count}[h(i)] \geq a_i$$

# Count-Min Sketch: general case

- Keep  $d$  by  $m$  Count matrix



- $d$  hash functions:  $h_1, \dots, h_d$ 
  - Just like in Bloom Filter, decrease errors with multiple hashes
  - Every time see string  $i$ :

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

## Querying the Count-Min Sketch

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

- Query  $Q(i)$ ?

- What is in  $\text{Count}[j, k]$ ?

$$\text{Count}[j, k] = \sum_{i: h_j(i)=k} a_i$$

- Thus:  $Q(i)$

$$\text{each } \text{Count}[j, h_j(i)] \geq a_i$$

- Return:

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

← tightest upper bound

©Carlos Guestrin 2013

5

## Analysis of Count-Min Sketch

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

- Set:

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

- Then, after seeing  $n$  elements:

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

← not bigger by more than  $\epsilon n$

- 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  $j$ :

$$(i \neq k) \wedge (h_j(i) = h_j(k))$$

- Bounding expected value:

$$E[I_{i,j,k}] = P(h_j(i) = h_j(k)) = \frac{1}{m} \leq \frac{\epsilon}{e}$$

- $X_{i,j}$  = total colliding mass on estimate of count of  $i$  in hash  $j$ :

$$X_{i,j} = \sum_{k \neq i} I_{i,j,k} a_k$$

$$\text{Count}[j, h_j(i)] = a_i + X_{i,j}$$

- Bounding colliding mass:

$$E[X_{i,j}] = \sum_{k \neq i} a_k E[I_{i,j,k}] \leq \frac{\epsilon}{e} \sum_{k \neq i} a_k \leq \frac{\epsilon}{e} n$$

- 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:  $\text{Count}[j, h_j(i)] = a_i + X_{i,j}$      $E[X_{i,j}] \leq \frac{\epsilon}{e} n$

- Markov inequality: For  $z_1, \dots, z_k$  positive iid random variables

$$P(\forall z_i : z_i > \alpha E[z_i]) < \alpha^{-k} \quad P(z_i > a) < \frac{E[z_i]}{a}$$

$\underbrace{\hspace{10em}}_{\text{ind. of } z_i} \quad \underbrace{\hspace{10em}}_{a = \alpha E[z_i]}$

- Applying to the Count-Min sketch:

$$\begin{aligned} P(\hat{a}_i > a_i + \epsilon n) &= P(\forall j, \text{count}[j, h_j(i)] > a_i + \epsilon n) \\ &= P(\forall j, a_i + X_{i,j} > a_i + \epsilon n) \\ &= P(\forall j, X_{i,j} > \epsilon n) < e^{-d} \leq \delta \end{aligned}$$

$d \text{ hash fens} \quad \uparrow \quad \text{"d"} \quad \uparrow \quad \text{Small prob.}$

©Carlos Guestrin 2013

8

## But Our Updates may be positive or Negative

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

- Count-Min sketch for positive & negative case

- $a_i$  no longer necessarily positive

- Update the same: Observe change  $\Delta_i$  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?

$$\hat{a}_i = \text{median}_j \text{Count}[j, h_j(i)]$$

$$a_i = \text{Count}[1, h_1(i)] + \text{Count}[2, h_2(i)] + \text{Count}[3, h_3(i)]$$

- Bound:  $|\hat{a}_i - a_i| \leq \underline{3\epsilon \|\mathbf{a}\|_1}$

- With probability at least  $1 - \delta^{1/4}$ , where  $\|\mathbf{a}\| = \sum_i |a_i|$

pos. or neg.  
non-int.

©Carlos Guestrin 2013

9

## Finally, Sketching for Perceptron

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

- Never need to know size of vocabulary!  $\star$

- Make a mistake, update Count-Min matrix:

$$\forall i \quad x_i^{(t)} \neq 0$$

$$\forall j \quad \text{Count}[j, h_j(i)] += y^{(t)} x_i^{(t)}$$

- Making a prediction:

$$\text{sign}(\mathbf{w}^{(t)} \cdot \mathbf{x}) = \text{sign} \left( \sum_i w_i^{(t)} x_i \right) \approx w_i^{(t)}$$

$$\sum_i w_i^{(t)} x_i \approx \sum_{i \neq 0} x_i \text{Median count}[j, h_j(i)]$$

- 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 22<sup>nd</sup>, 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 **X**
  - **Output:** Set of k similar articles



X



©Emily Fox 2013

14

# Document Representation

- Bag of words model



ignore order  
of the words

$$x = \begin{bmatrix} w_{c_1} \\ w_{c_2} \\ \vdots \\ w_{c_d} \end{bmatrix}$$

word count

$|V| = \text{size of vocab.}$   
 $= d$

©Emily Fox 2013

15

# 1-Nearest Neighbor

- Articles  $X = \{x^1, \dots, x^N\}$ ,  $x^i \in \mathbb{R}^d$
- Query:  $x$
- 1-NN
  - Goal: find article in  $X$  "closest" to  $x$   
★ need distance metric★
  - Formulation:  $d(u, v)$

$$x^{NN} = \arg \min_{x^i \in X} d(x^i, x)$$

©Emily Fox 2013

16



# 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 x
  - Formulation:

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

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 2013

17

# Distance Metrics – Euclidean

$$d(u, v) = \sqrt{\sum_{i=1}^d (u_i - v_i)^2} = \|u - v\|_2$$

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

Equivalently,

*weight for dim i*

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

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

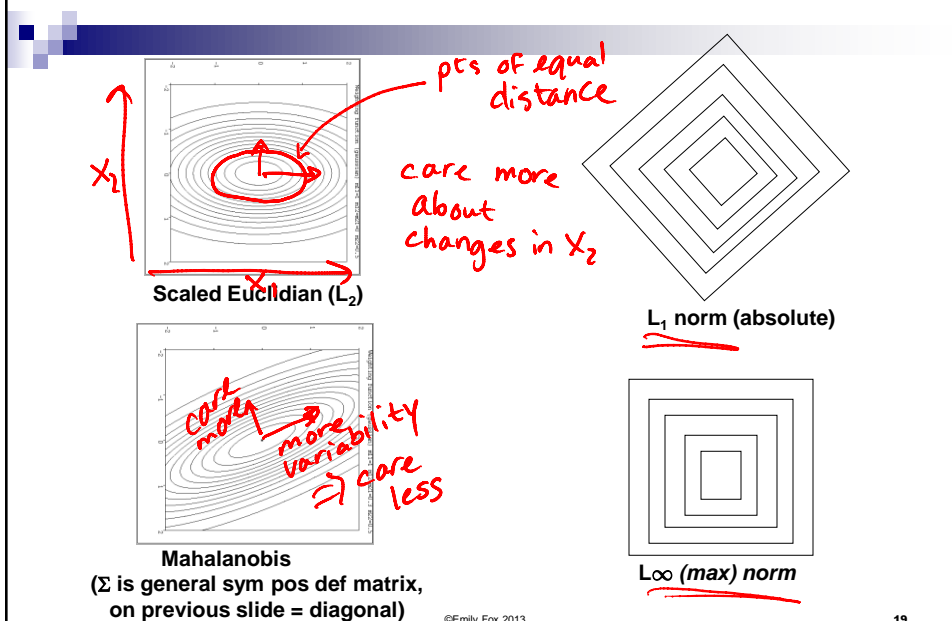
Other Metrics...

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

©Emily Fox 2013

18

## Notable Distance Metrics (and their level sets)



## Euclidean Distance + Document Retrieval

- Recall distance metric

$$d(u, v) = \sqrt{\sum_{i=1}^d (u_i - v_i)^2} = \|u - v\|_2$$

- What if each document were  $\alpha$  times longer?

- Scale word count vectors

$$u \leftarrow \alpha u$$

$$v \leftarrow \alpha v$$

- What happens to measure of similarity?

$$\|\alpha u - \alpha v\|_2 = \alpha \|u - v\|_2 > \|u - v\|_2$$

$\alpha > 1$

now less similar

- Good to normalize vectors

$$\|u\|_2 = \|v\|_2 = 1$$

# Issues with Document Representation

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

"The koala hugs the tree..."

"Tree huggers..."

rare words swamped by common words

common words like "tree" and "hug" are dominating (and "the")

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

# TF-IDF

- Term frequency:

$$tf(t, d) = \# \text{ of occur of } t \text{ in } d \triangleq f(t, d)$$

term doc

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

- Inverse document frequency:

$$idf(t, X) = \log \frac{|X|}{1 + |\{d \in X : t \in d\}|} \rightarrow 0 \text{ if many docs } d$$

$\frac{f(t, d)}{\max\{f(t, d); \text{word}\}}$  ← prevent bias towards long docs

- tf-idf:

$$tfidf(t, d, X) = tf(t, d) \times idf(t, X) > 0 \text{ on}$$

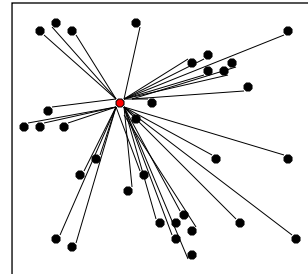
- High for document  $d$  with high frequency of term  $t$  (high "term frequency") and few documents containing term  $t$  in the corpus (high "inverse doc frequency")

# 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!
- $O(M \log k)$  per k-NN query!



33 Distance Computations

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

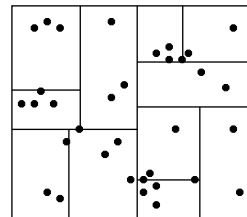
©Emily Fox 2013

23

# KD-Trees

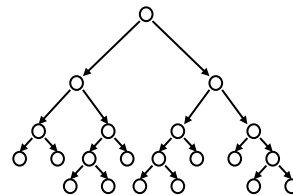
## Smarter approach: *kd-trees*

- 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...

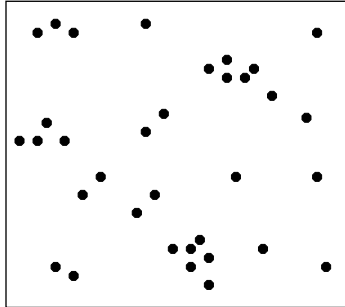


©Emily Fox 2013

24

# KD-Tree Construction

*2d example*

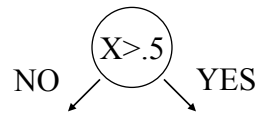
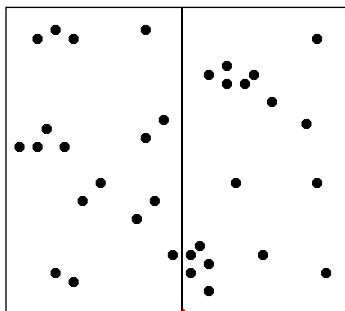


Pt	X	Y
1	0.00	0.00
2	1.00	4.31
3	0.13	2.85
...	...	...

*x<sup>1</sup>  
x<sup>2</sup>  
x<sup>3</sup>  
⋮*

- Start with a list of  $d$ -dimensional points.

# KD-Tree Construction



Pt	X	Y
1	<u>0.00</u>	0.00
3	<u>0.13</u>	2.85
...	...	...

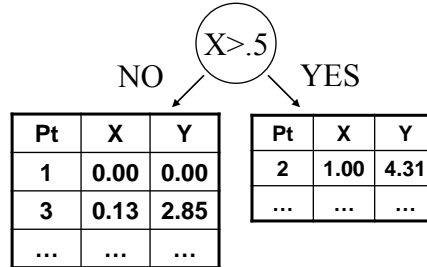
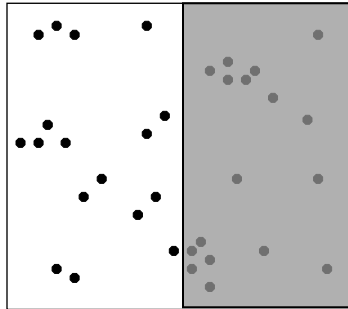
Pt	X	Y
2	<u>1.00</u>	4.31
...	...	...

*$x \leq V$     $x > V$*

*X in this case  
in this case  $V=0.5$*

- 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  $x_{d_j}^i \leq V$ .

# KD-Tree Construction

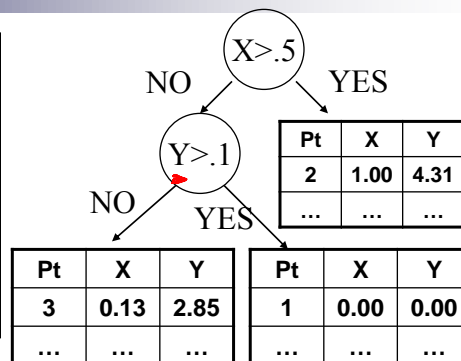
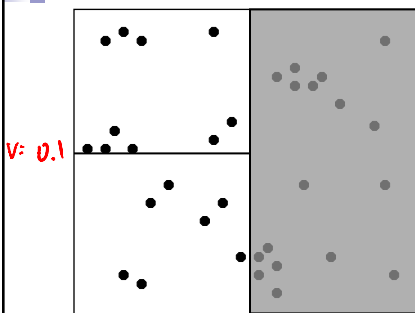


- Consider each group separately and possibly split again (along same/different dimension).

□ Stopping criterion to be discussed...

*How to choose dimension...*

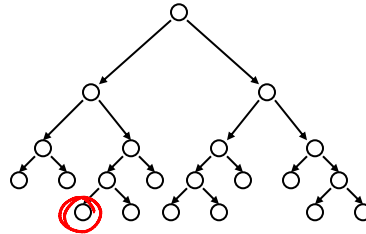
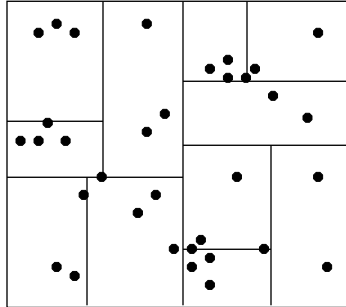
# KD-Tree Construction



- Consider each group separately and possibly split again (along same/different dimension).

□ Stopping criterion to be discussed...

# KD-Tree Construction



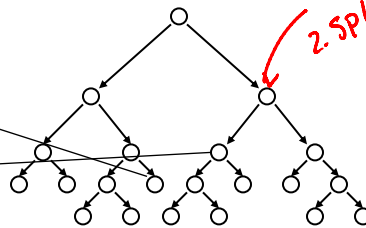
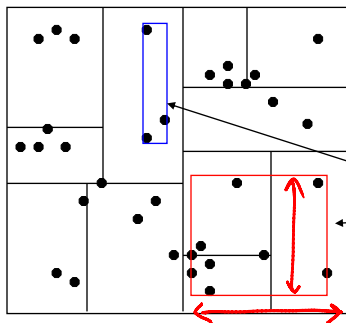
- Continue splitting points in each set
  - creates a binary tree structure
- Each leaf node contains a list of points

*satisfying all conditions down the tree to that point*

©Emily Fox 2013

29

# KD-Tree Construction



*Store:  
1. which dim?  
2. split value*

- 3. ■ 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

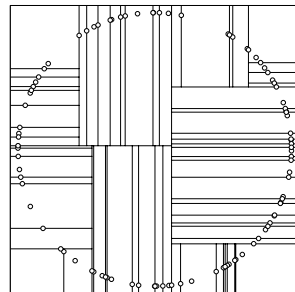
Use heuristics to make splitting decisions:

- Which dimension do we split along?  
*widest (or alternate)*
- Which value do we split at?  
*median of chosen split dim (or center)*
- When do we stop?  
*fewer than m pts left*  
*OR*  
*box hits minimum width*

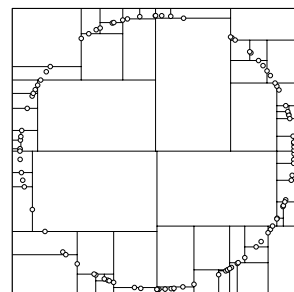
©Emily Fox 2013

31

# Many heuristics...



median heuristic



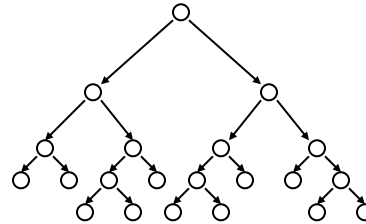
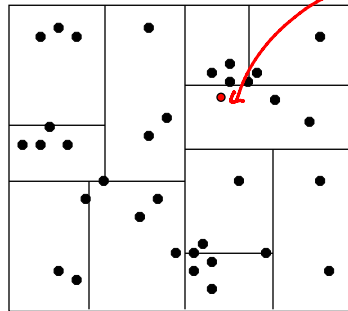
center-of-range heuristic

©Emily Fox 2013

32



# Nearest Neighbor with KD Trees

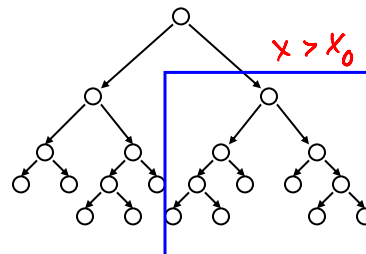
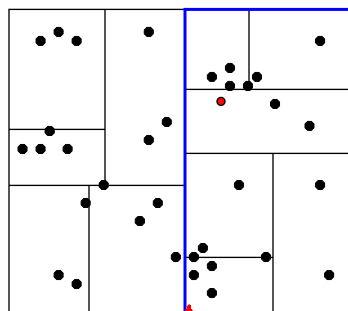


- Traverse the tree looking for the nearest neighbor of the query point.

©Emily Fox 2013

33

# Nearest Neighbor with KD Trees

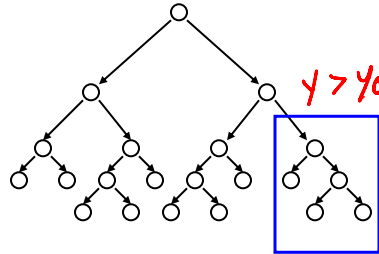
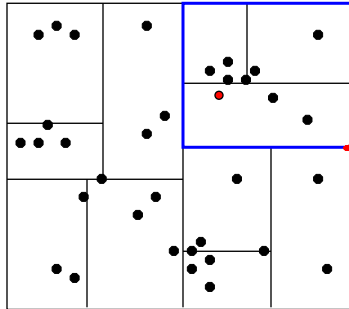


- Examine nearby points first:
  - Explore branch of tree closest to the query point first.

©Emily Fox 2013

34

# Nearest Neighbor with KD Trees

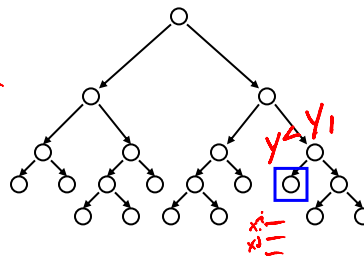
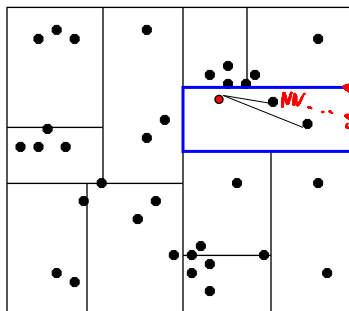


- Examine nearby points first:
  - Explore branch of tree closest to the query point first.

©Emily Fox 2013

35

# Nearest Neighbor with KD Trees

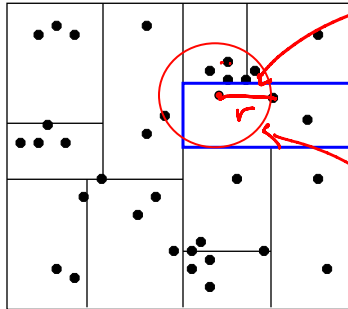


- When we reach a leaf node:
  - Compute the distance to each point in the node.

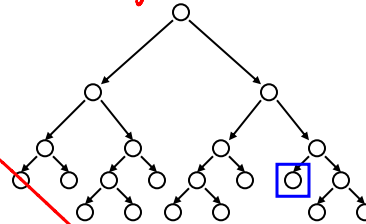
©Emily Fox 2013

36

# Nearest Neighbor with KD Trees



distance to closest neighbor found so far

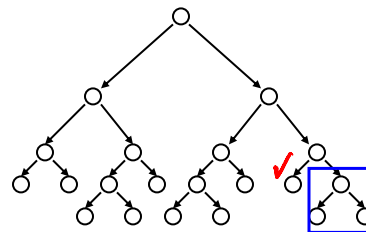
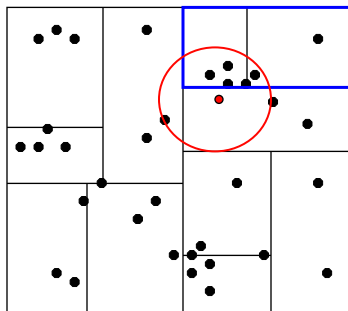


does NN have to be in this box?

NO

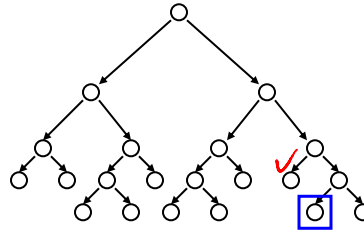
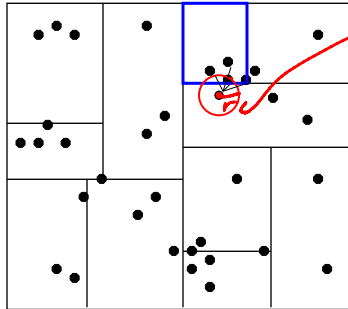
- When we reach a leaf node:
  - Compute the distance to each point in the node.

# Nearest Neighbor with KD Trees



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

# Nearest Neighbor with KD Trees

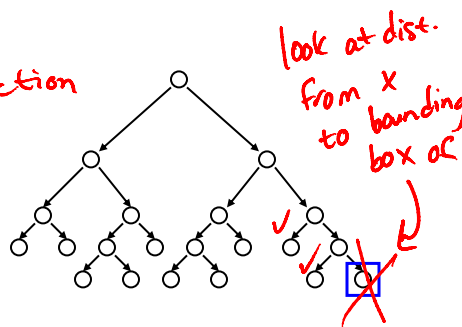
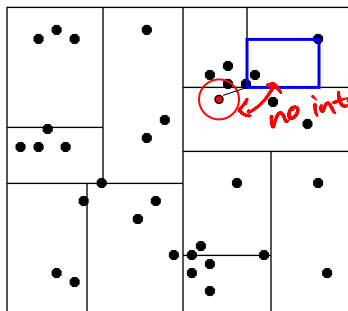


- Each time a new closest node is found, update the distance bound

©Emily Fox 2013

39

# Nearest Neighbor with KD Trees

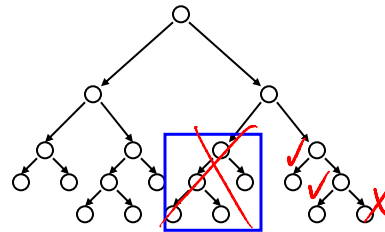
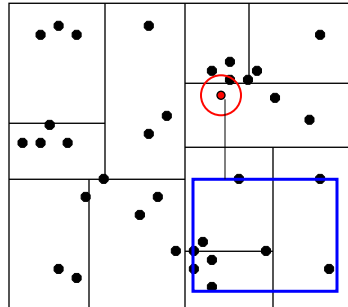


- Using the distance bound and bounding box of each node:
  - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2013

40

# Nearest Neighbor with KD Trees

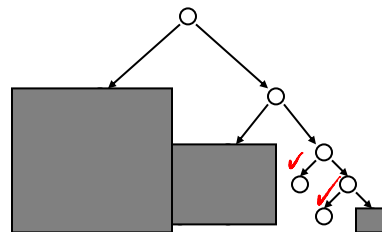
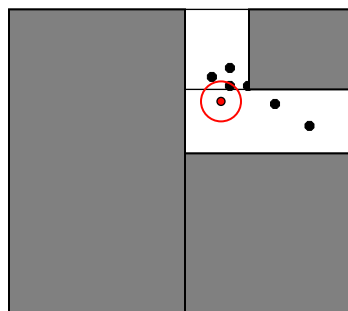


- Using the distance bound and bounding box of each node:
  - Prune parts of the tree that could NOT include the nearest neighbor

©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

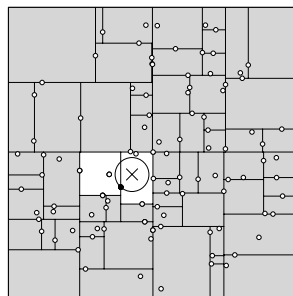
# Complexity

- For (nearly) balanced, binary trees...
- Construction
  - Size:  $2N-1 \rightarrow O(N)$
  - Depth:  $O(\log N)$
  - Median + send points left right:  $O(N)$  at every tree level
  - Construction time:  $O(N \log N)$  (smart)
- 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)

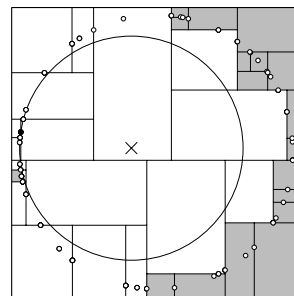
©Emily Fox 2013

43

# Complexity



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



pruned only a few  
 $O(N)$

©Emily Fox 2013

44

## Complexity for N Queries

- Ask for nearest neighbor to each document

*N queries*

- Brute force 1-NN:  $O(N^2)$

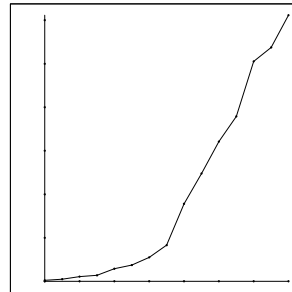
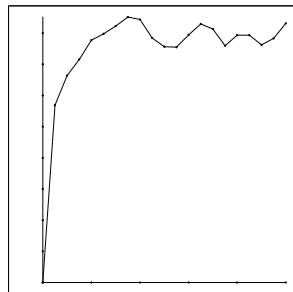
- kd-trees:  $O(N \log N) \rightarrow O(N^2)$

*↑  
potentially  
large savings!*

©Emily Fox 2013

45

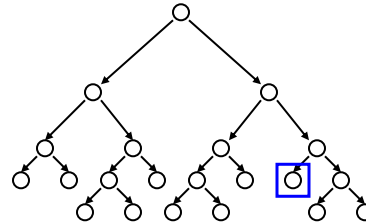
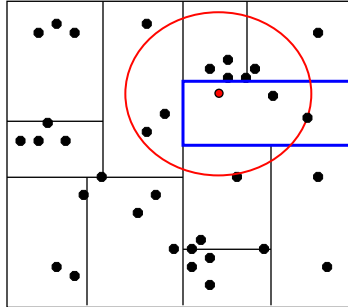
## Inspections vs. $N$ and $d$



©Emily Fox 2013

46

# K-NN with KD Trees

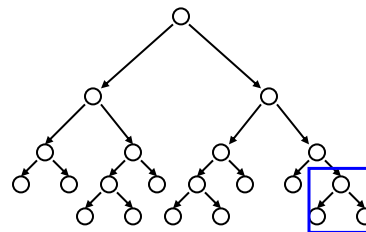
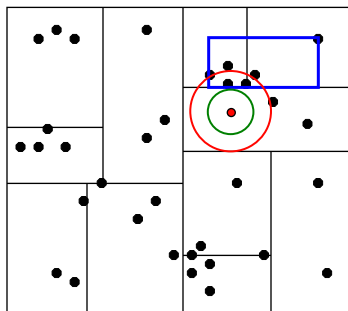


- Exactly the same algorithm, but maintain distance as distance to furthest of current  $k$  nearest neighbors
- Complexity is:

©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  $r/\alpha$ .
- 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](http://grist.caltech.edu/sc4devo/.../files/sc4devo_scalable_datamining.ppt)