

Case Study 1: Estimating Click Probabilities

Tackling an Unknown Number of Features with Sketching

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

Emily Fox

January 16th, 2014

©Emily Fox 2014

1

Problem 1: Complexity of Update Rules for Logistic Regression

- Logistic regression update:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y = 1 | \mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

stoch. grad. asc. ↓

- Complexity of updates:

- Constant in number of data points ✓
- In number of features? *O(d)*
 - Problem both in terms of computational complexity and sample complexity

What if we have 1B features??

- What can we with very high dimensional feature spaces?
 - Kernels not always appropriate, or scalable ← *"kernel trick"*
 - What else?

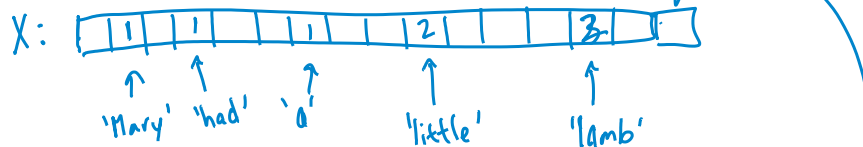
©Emily Fox 2014

2

Problem 2: Unknown Number of Features

- For example, bag-of-words features for text data:

- "Mary had a little lamb, little lamb..."



- What's the dimensionality of x ? *size of vocabulary ... millions*
- What if we see new word that was not in our vocabulary?
 - Obamacare
 - Theoretically, just keep going in your learning, and initialize $w_{\text{Obamacare}} = 0$
 - In practice, need to re-allocate memory, fix indices, ... A big problem for Big Data

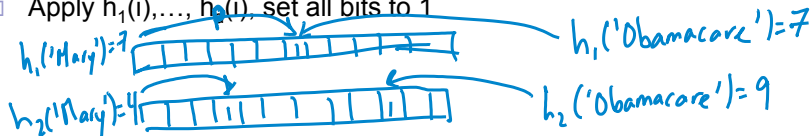
©Emily Fox 2014

3

Bloom Filter: Multiple Hash Tables

- Single hash table -> Many false positives
- Multiple hash tables with independent hash functions

- Apply $h_1(i), \dots, h_n(i)$, set all bits to 1



- Query $Q(i)$?

if $\forall j \quad h_j(i) = 1$
 $Q(i) = \text{very probably yes}$
 else $Q(i) = \text{no}$

'Mary' + 'Obamacare'
 collide in h_1 but not h_2

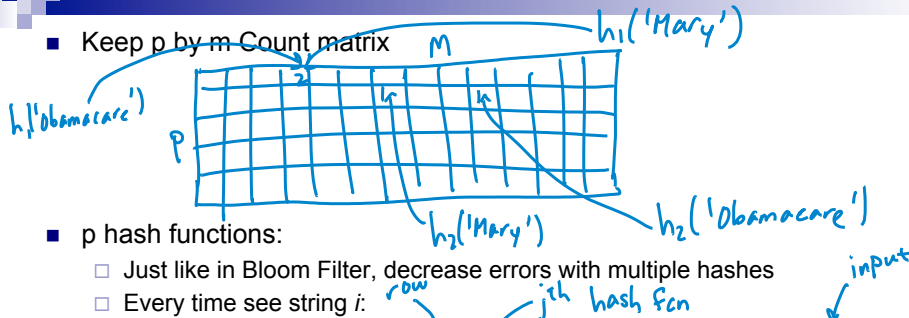
- Significantly decrease probability of false positives

©Emily Fox 2014

4

Count-Min Sketch: general case

- Keep p by m Count matrix



- p hash functions:

- Just like in Bloom Filter, decrease errors with multiple hashes
- Every time see string i :

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

©Emily Fox 2014

5

Finally, Sketching for LR

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y = 1 | \mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

- Never need to know size of vocabulary!
- At every iteration, update Count-Min matrix:

$$\forall j, k \quad \text{Count}[j, k] = (1 - \eta_t \lambda) \text{count}[j, k]$$

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

$$\forall j \quad \text{Count}[j, h_j(i)] += x_i^{(t)} \cdot \text{const} \quad \leftarrow \eta_t (y^{(t)} - P(Y=1|\dots))$$

- Making a prediction:

$$\text{Remember our est. of } w_i^{(t)} : \text{median}_j \text{Count}[j, h_j(i)]$$

Make pred:

$$-\log \text{odds} = w_0^{(t)} + \sum_{i: x_i \neq 0} \text{median}_j \text{Count}[j, h_j(i)] x_i^{(t)}$$

- Scales to huge problems, great practical implications...

©Emily Fox 2014

6

Hash Kernels

- Count-Min sketch not designed for negative updates
- Biased estimates of dot products
- Hash Kernels:** Very simple, but powerful idea to remove bias
- Pick 2 hash functions:
 - h : Just like in Count-Min hashing $h: X \rightarrow \{1, \dots, m\}$
 - ξ : Sign hash function $\xi: X \rightarrow \{+1, -1\}$
 - Removes the bias found in Count-Min hashing (see homework)

can think of as random projection of X

- Define a "kernel", a projection ϕ for \mathbf{x} :



$$\phi_i(\mathbf{x}) = \sum_{j:h(j)=i} f(j)x_j$$

IF $x_j = 7$
 $h(j) = 4$
 $f(j) = -1$
 add \downarrow -7
 to bin 4

©Emily Fox 2014

7

Hash Kernels Preserve Dot Products

$$\phi_i(\mathbf{x}) = \sum_{j:h(j)=i} \xi(j)x_j$$

- Hash kernels provide unbiased estimate of dot-products!

$$E_{h, \xi} [\phi(\mathbf{x}) \cdot \phi(\mathbf{y})] = \mathbf{x} \cdot \mathbf{y} \quad \text{pf: by homework}$$

- Variance decreases as $O(1/m)$ ← gets better w/ more dims

- Choosing m ? For $\epsilon > 0$, if

$$m = O\left(\frac{\log \frac{N}{\delta}}{\epsilon^2}\right) \quad \text{log in data size}$$

- Under certain conditions...
- Then, with probability at least $1 - \delta$:

$$(1 - \epsilon) \|\mathbf{x} - \mathbf{x}'\|_2^2 \leq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2^2 \leq (1 + \epsilon) \|\mathbf{x} - \mathbf{x}'\|_2^2$$

©Emily Fox 2014

8

Learning With Hash Kernels

- Given hash kernel of dimension m , specified by h and ξ
 - Learn m dimensional weight vector
- Observe data point \mathbf{x}
 - Dimension does not need to be specified a priori!
- Compute $\phi(\mathbf{x})$:
 - Initialize $\phi(\mathbf{x})$
 - For non-zero entries j of \mathbf{x}_j :
- Use normal update as if observation were $\phi(\mathbf{x})$, e.g., for LR using SGD:
$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + \phi_i(\mathbf{x}^{(t)}) [y^{(t)} - P(Y = 1 | \phi(\mathbf{x}^{(t)}), \mathbf{w}^{(t)})] \right\}$$

©Emily Fox 2014

9

Interesting Application of Hash Kernels: Multi-Task Learning

- Personalized click estimation for many users:
 - One global click prediction vector \mathbf{w} :
 - But...
 - A click prediction vector \mathbf{w}_u per user u :
 - But...
- Multi-task learning: Simultaneously solve multiple learning related problems:
 - Use information from one learning problem to inform the others
- In our simple example, learn both a global \mathbf{w} and one \mathbf{w}_u per user:
 - Prediction for user u :
 - If we know little about user u :
 - After a lot of data from user u :

©Emily Fox 2014

10

Problems with Simple Multi-Task Learning

- Dealing with new user is annoying, just like dealing with new words in vocabulary
- Dimensionality of joint parameter space is HUGE, e.g. personalized email spam classification from Weinberger et al.:
 - 3.2M emails
 - 40M unique tokens in vocabulary
 - 430K users
 - 16T parameters needed for personalized classification!

©Emily Fox 2014

11

Hash Kernels for Multi-Task Learning

- Simple, pretty solution with hash kernels:
 - Very multi-task learning as (sparse) learning problem with (huge) joint data point \mathbf{z} for point \mathbf{x} and user u :
- Estimating click probability as desired:
- Address huge dimensionality, new words, and new users using hash kernels:
 - Desired effect achieved if j includes both
 - just word (for global \mathbf{w})
 - word,user (for personalized \mathbf{w}_u)

©Emily Fox 2014

12

Simple Trick for Forming Projection $\phi(\mathbf{x}, u)$

- Observe data point \mathbf{x} for user u
 - Dimension does not need to be specified a priori and user can be unknown!
- Compute $\phi(\mathbf{x}, u)$:
 - Initialize $\phi(\mathbf{x}, u)$
 - For non-zero entries j of \mathbf{x}_j :
 - E.g., $j = \text{'Obamacare'}$
 - Need two contributions to ϕ :
 - Global contribution
 - Personalized Contribution
 - Simply:
- Learn as usual using $\phi(\mathbf{x}, u)$ instead of $\phi(\mathbf{x})$ in update function

©Emily Fox 2014

13

Results from Weinberger et al. on Spam Classification: Effect of m

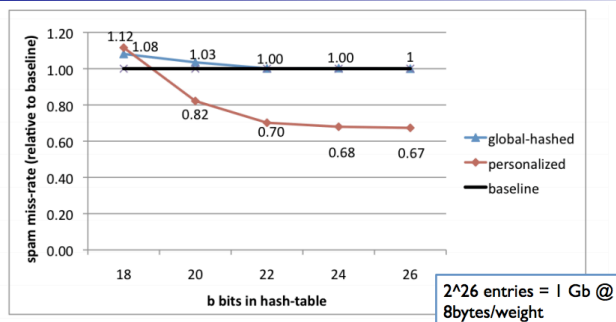


Figure 2. The decrease of uncaught spam over the baseline classifier averaged over all users. The classification threshold was chosen to keep the not-spam misclassification fixed at 1%. The hashed global classifier (*global-hashed*) converges relatively soon, showing that the distortion error ϵ_d vanishes. The personalized classifier results in an average improvement of up to 30%.

©Emily Fox 2014

14

Results from Weinberger et al. on Spam Classification: Illustrating Multi-Task Effect

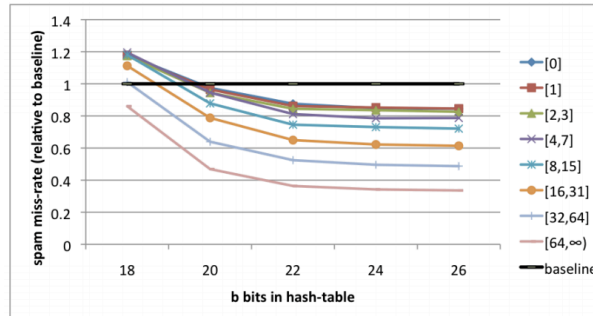


Figure 3. Results for users clustered by training emails. For example, the bucket [8, 15] consists of all users with eight to fifteen training emails. Although users in buckets with large amounts of training data do benefit more from the personalized classifier (up to 65% reduction in spam), even users that did not contribute to the training corpus at all obtain almost 20% spam-reduction.

©Emily Fox 2014

15

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 logistic regression
- Hash kernels:
 - Sparse representation for feature vectors
 - Very simple, use two hash function (Can use one hash function...take least significant bit to define ξ)
 - Quickly generate projection $\phi(\mathbf{x})$
 - Learn in projected space
- Multi-task learning:
 - Solve many related learning problems simultaneously
 - Very easy to implement with hash kernels
 - Significantly improve accuracy in some problems (if there is enough data from individual users)

©Emily Fox 2014

16

Case Study 2: Document Retrieval

Task Description: Finding Similar Documents

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

Emily Fox

January 16th, 2014

©Emily Fox 2014

17

Document Retrieval

- **Goal:** Retrieve documents of interest

- **Challenges:**

- Tons of articles out there
- How should we measure similarity?



©Emily Fox 2014

18

Task 1: Find Similar Documents

- To begin...

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



©Emily Fox 2014

19

Document Representation

- Bag of words model



©Emily Fox 2014

20

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:

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

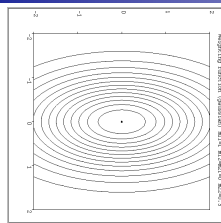
Other Metrics...

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

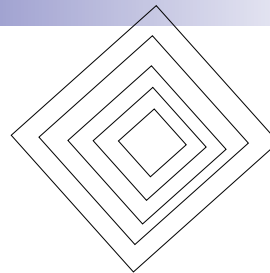
©Emily Fox 2014

23

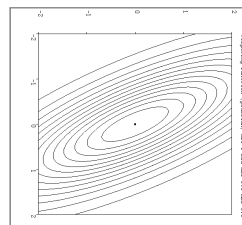
Notable Distance Metrics (and their level sets)



Scaled Euclidean (L_2)

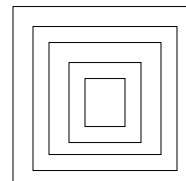


L_1 norm (absolute)



Mahalanobis

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



L_1 (max) norm

©Emily Fox 2014

24

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 2014

25

Issues with Document Representation

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



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

©Emily Fox 2014

26

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”)

©Emily Fox 2014

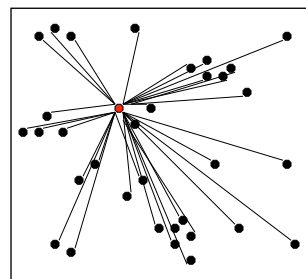
27

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)

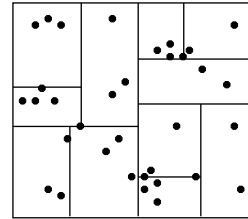
©Emily Fox 2014

28

KD-Trees

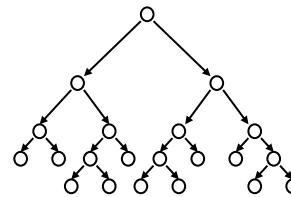
- Smarter approach: **kd-trees**

- 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

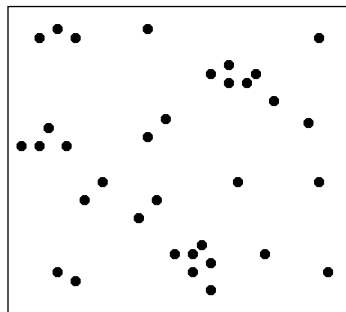
- We’ll get back to this...



©Emily Fox 2014

29

KD-Tree Construction



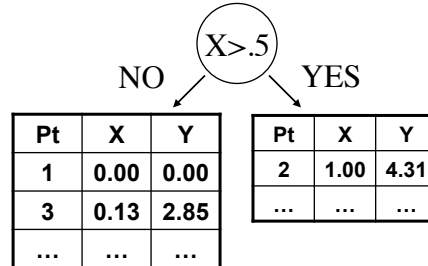
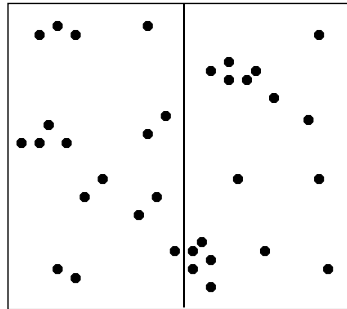
Pt	X	Y
1	0.00	0.00
2	1.00	4.31
3	0.13	2.85
...

- Start with a list of d -dimensional points.

©Emily Fox 2014

30

KD-Tree Construction

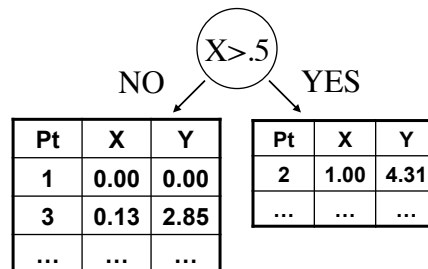
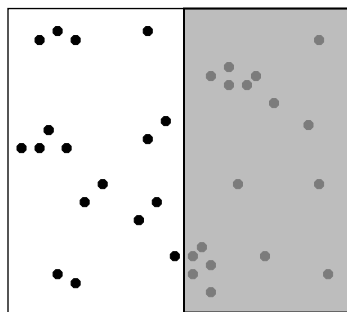


- Split the points into 2 groups by:
 - Choosing dimension d_j and value V (methods to be discussed...)
 - Separating the points into $x_{d_j}^i > V$ and $x_{d_j}^i \leq V$.

©Emily Fox 2014

31

KD-Tree Construction

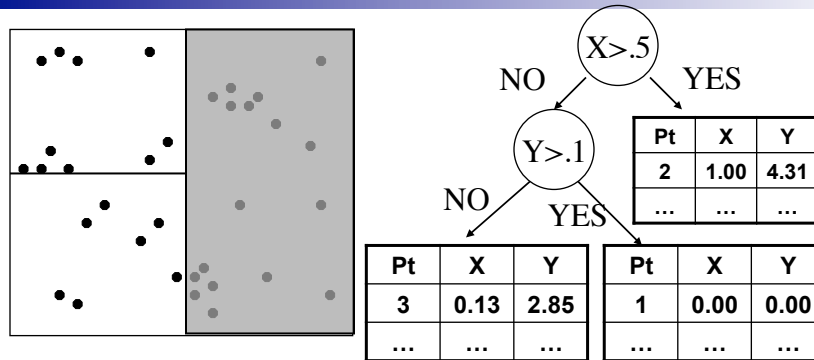


- Consider each group separately and possibly split again (along same/different dimension).
 - Stopping criterion to be discussed...

©Emily Fox 2014

32

KD-Tree Construction

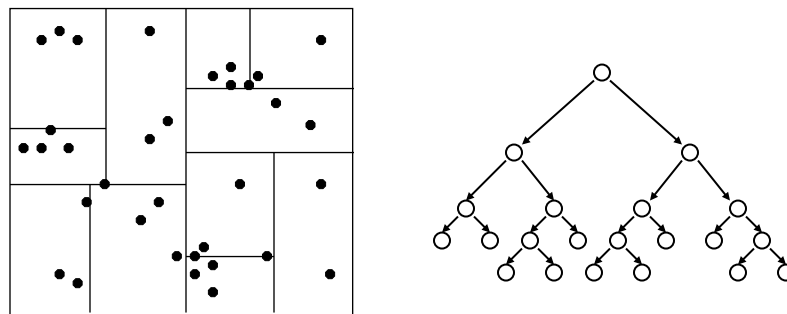


- Consider each group separately and possibly split again (along same/different dimension).
 - Stopping criterion to be discussed...

©Emily Fox 2014

33

KD-Tree Construction

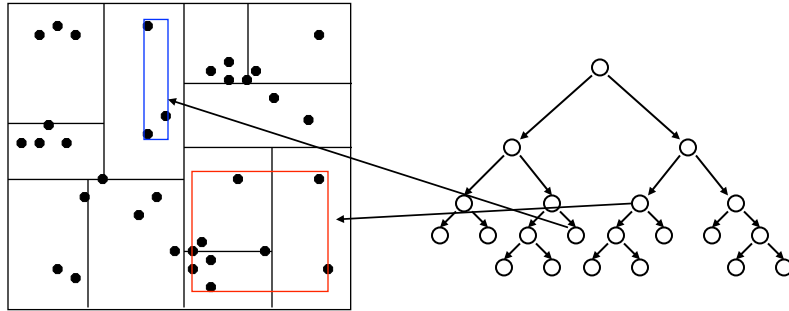


- Continue splitting points in each set
 - creates a binary tree structure
- Each leaf node contains a list of points

©Emily Fox 2014

34

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 2014

35

KD-Tree Construction

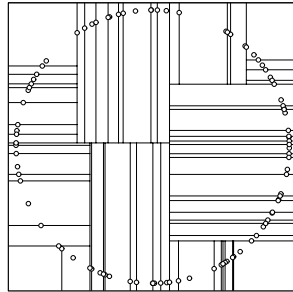
Use heuristics to make splitting decisions:

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

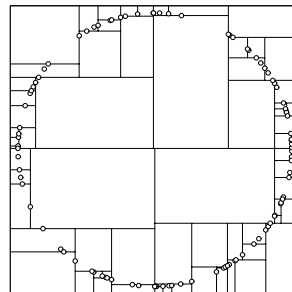
©Emily Fox 2014

36

Many heuristics...



median heuristic

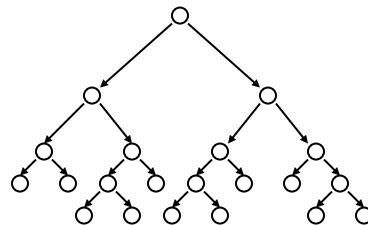
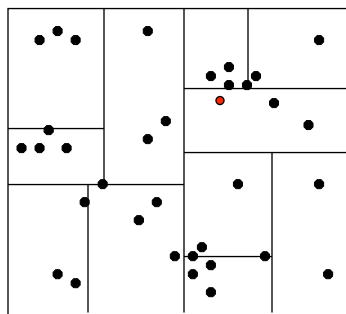


center-of-range heuristic

©Emily Fox 2014

37

Nearest Neighbor with KD Trees

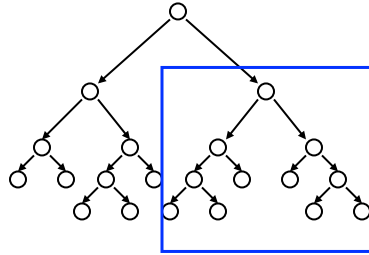
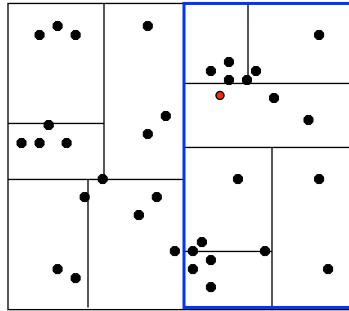


- Traverse the tree looking for the nearest neighbor of the query point.

©Emily Fox 2014

38

Nearest Neighbor with KD Trees

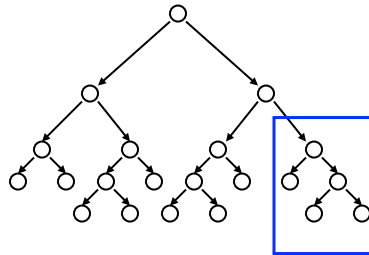
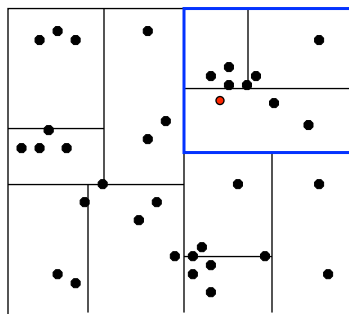


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2014

39

Nearest Neighbor with KD Trees

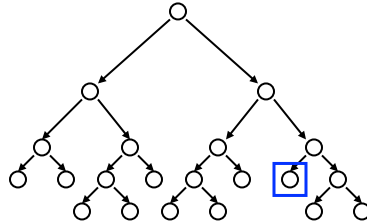
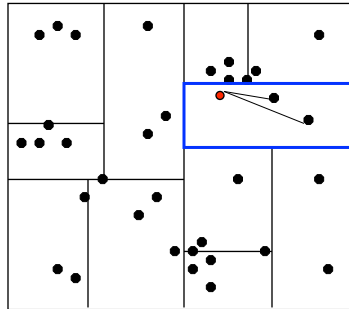


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2014

40

Nearest Neighbor with KD Trees

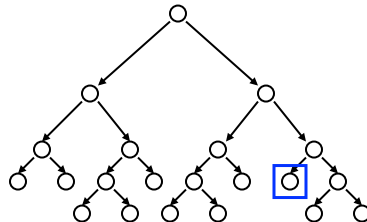
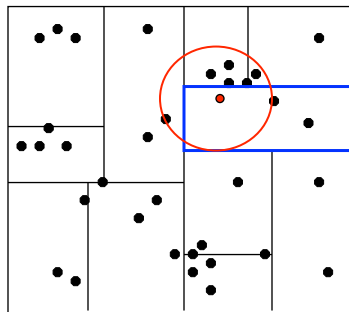


- When we reach a leaf node:
 - Compute the distance to each point in the node.

©Emily Fox 2014

41

Nearest Neighbor with KD Trees

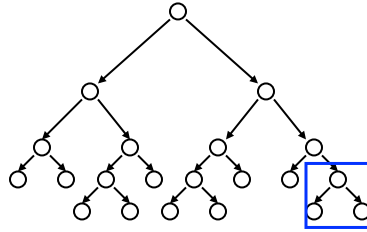
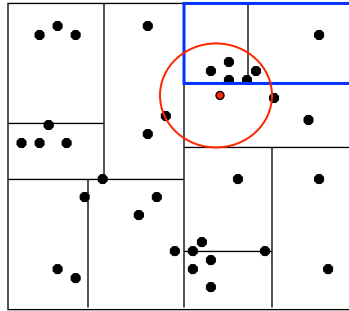


- When we reach a leaf node:
 - Compute the distance to each point in the node.

©Emily Fox 2014

42

Nearest Neighbor with KD Trees

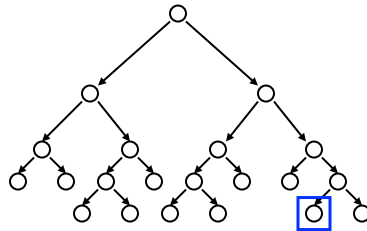
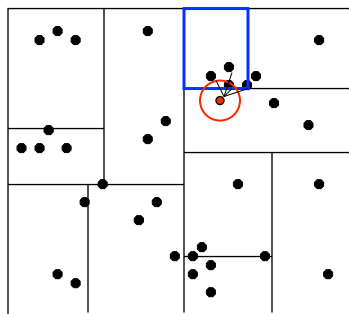


- Then backtrack and try the other branch at each node visited

©Emily Fox 2014

43

Nearest Neighbor with KD Trees

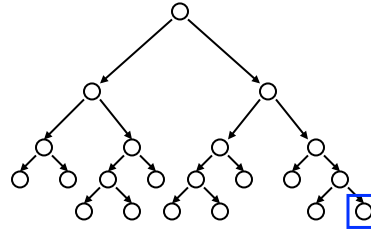
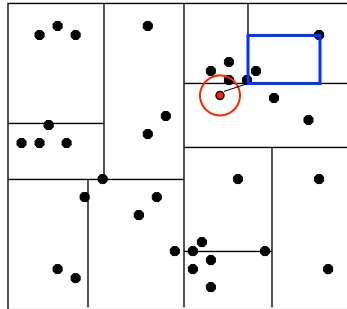


- Each time a new closest node is found, update the distance bound

©Emily Fox 2014

44

Nearest Neighbor with KD Trees

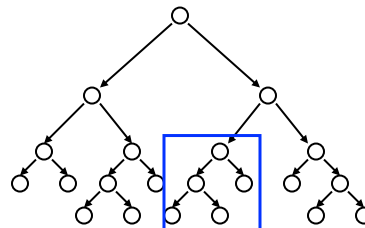
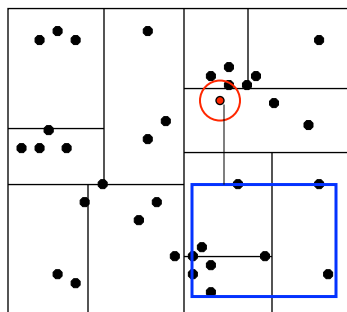


- Using the distance bound and bounding box of each node:
 - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2014

45

Nearest Neighbor with KD Trees

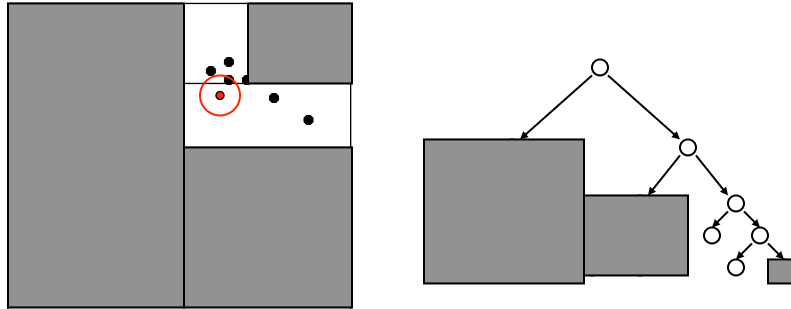


- Using the distance bound and bounding box of each node:
 - Prune parts of the tree that could NOT include the nearest neighbor

©Emily Fox 2014

46

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 2014

47

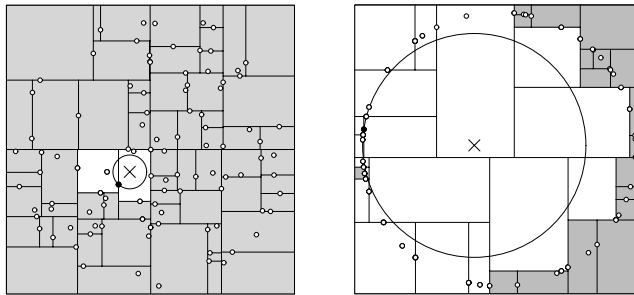
Complexity

- 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 2014

48

Complexity



©Emily Fox 2014

49

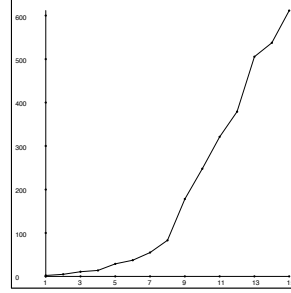
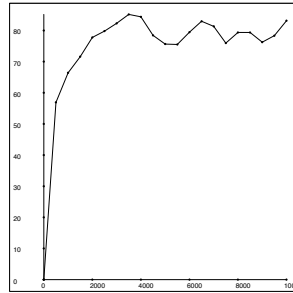
Complexity for N Queries

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

©Emily Fox 2014

50

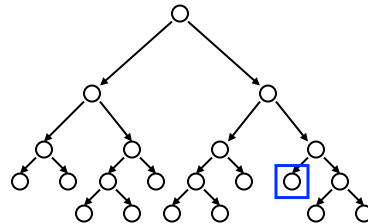
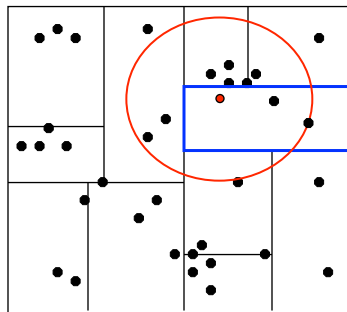
Inspections vs. N and d



©Emily Fox 2014

51

K-NN with KD Trees

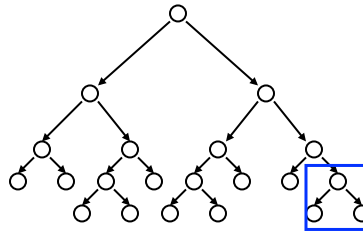
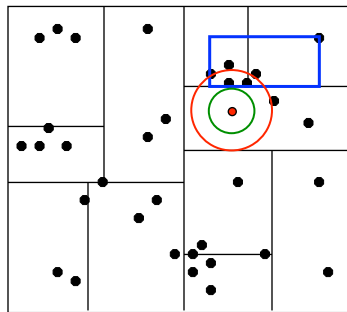


- Exactly the same algorithm, but maintain distance as distance to furthest of current k nearest neighbors
- Complexity is:

©Emily Fox 2014

52

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 2014

53

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 \rightarrow Everything equidistant (i.e., everything is far away)
 - Need technique to learn what features are important for your task

©Emily Fox 2014

54

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 2014

55

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 2014

56