

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? *old*
 - Problem both in terms of computational complexity and sample complexity

What if we have 1B features??

d very large

- What can we with very high dimensional feature spaces?

- Kernels not always appropriate, or scalable ← *"kernel trick"*
- What else?

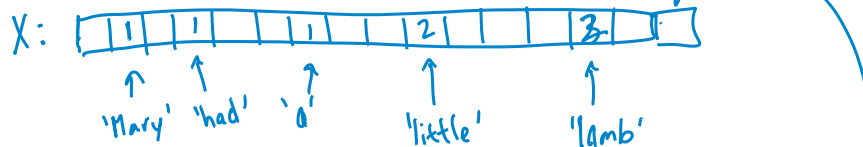
©Emily Fox 2014

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 *d can change*
 - 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

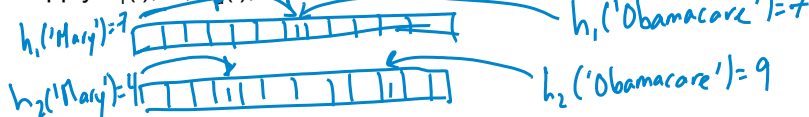
Bloom Filter: Multiple Hash Tables

- Single hash table -> Many false positives

just keep track of binary bit vector of length m ($m \ll d$)

- Multiple hash tables with independent hash functions

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



- Query $Q(i)$?

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

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

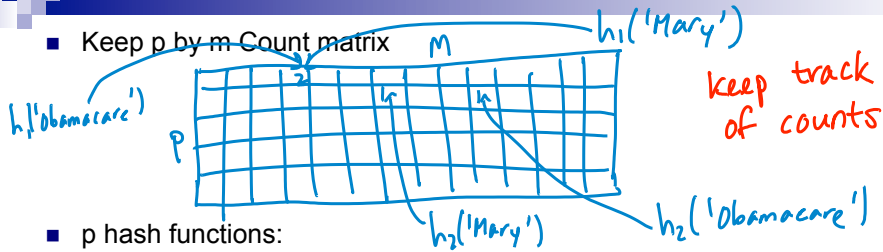
- Significantly decrease probability of false positives

©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:

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

$$\forall x_i^{(t)} \neq 0 \quad \forall j \quad \text{Count}[j, h_j(i)] += x_i^{(t)} \cdot \text{const}$$

"soft counts" instead

- Making a prediction:

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 $\rightarrow h: X \rightarrow \{1, \dots, m\}$
- ξ : Sign hash function $\rightarrow \xi: X \rightarrow \{+1, -1\}$
 - 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

Hash Kernels Preserve Dot Products

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

- 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 $O(1/m)$ \leftarrow 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 $1 - \delta$:

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

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 $\phi(\mathbf{x}) = \mathbf{0}$

- For non-zero entries j of \mathbf{x} :

$$\phi_{h(j)} += \xi(j) X_j$$

e.g. $j = \text{'Mary'}$ $h(\text{'Mary'}) = 7$
 $\xi(\text{'Mary'}) = -1$
 $\phi_7 += -X_{\text{'Mary'}}$

- Use normal update as if observation were $\phi(\mathbf{x})$, e.g., for LR using SGD:

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

w length m
 $i = 0, 1, \dots, m-1$

$$P(Y=1 | \phi(\mathbf{x}^{(t)}), \mathbf{w}^{(t)}) = \frac{\exp(\phi(\mathbf{x}^{(t)}) \cdot \mathbf{w}^{(t)})}{1 + \exp(\phi(\mathbf{x}^{(t)}) \cdot \mathbf{w}^{(t)})}$$

©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} :

predict using $\mathbf{w} \cdot \mathbf{x} \leftarrow \frac{\exp(\mathbf{w} \cdot \mathbf{x})}{1 + \exp(\mathbf{w} \cdot \mathbf{x})}$

- But... people are unique

- A click prediction vector \mathbf{w}_u per user u :

predict with $\mathbf{w}_u \cdot \mathbf{x}$

- But... people don't each provide much data (label)

- 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 :

$$(\mathbf{w} + \mathbf{w}_u) \cdot \mathbf{x} = \mathbf{w} \cdot \mathbf{x} + \mathbf{w}_u \cdot \mathbf{x}$$

- If we know little about user u :

basically $\mathbf{w} \cdot \mathbf{x}$

- After a lot of data from user u :

using $\mathbf{w} + \mathbf{w}_u$ as your vector

©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 :

$$\mathbf{z}_{(x,u)} = (\underbrace{x_1, \dots, x_d}_{\text{global}}, \underbrace{0, \dots, 0}_{\text{length } d}, \underbrace{x_1, \dots, x_d}_{\text{user } u}, \underbrace{0, \dots, 0}_{\text{\# of users}})$$

- Estimating click probability as desired:

$$\mathbf{w} = (\underbrace{w}_{\text{global}}, w_1, \dots, w_u, \dots, w_{\# \text{ users}}) \rightarrow \mathbf{z}_{(x,u)} \cdot \mathbf{w} = w \cdot \mathbf{x} + w_u \cdot \mathbf{x} = (w + w_u) \cdot \mathbf{x}$$

- Address huge dimensionality, new words, and new users using hash kernels:

$$\phi(\mathbf{z}_{(x,u)}) \quad \text{just like } w / \text{ hash kernels}$$

$$\phi_i = \sum_{j: h(j)=i} f(j) z_{(x,u),j}$$

- Desired effect achieved if j includes both
 - just word (for global w)
 - word,user (for personalized 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~~ ^{new}
- Compute $\phi(\mathbf{x}, u)$:
 - Initialize $\phi(\mathbf{x}, u) = \mathbf{0}$
 - For non-zero entries j of \mathbf{x} :
 - E.g., $j = \text{'Obamacare'}$
 - Need two contributions to ϕ :
 - Global contribution $\rightarrow \phi_h(\text{'Obamacare'}) + = f(\text{'Obamacare'}) \cdot \underline{X_j}$
 - Personalized Contribution
 - Simply: $u=17$ $\rightarrow \phi_h(\text{'Obamacare-user17'}) + = f(\text{'Obamacare-user17'}) \cdot \underline{X_j}$
- Learn as usual using $\phi(\mathbf{x}, u)$ instead of $\phi(\mathbf{x})$ in update function

new
 \uparrow dim $\phi(\mathbf{x}, u)$
 can change!

©Emily Fox 2014

13

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

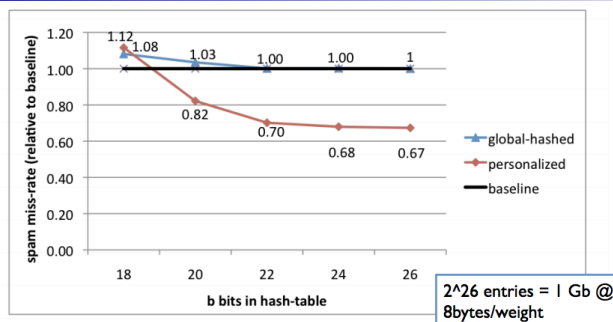


Figure 2. The decrease of uncaught spam over the baseline classifier averaged over all users. The classification threshold was chosen to keep the not-spam misclassification fixed at 1%. The hashed global classifier (*global-hashed*) converges relatively soon, showing that the distortion error ϵ_d vanishes.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

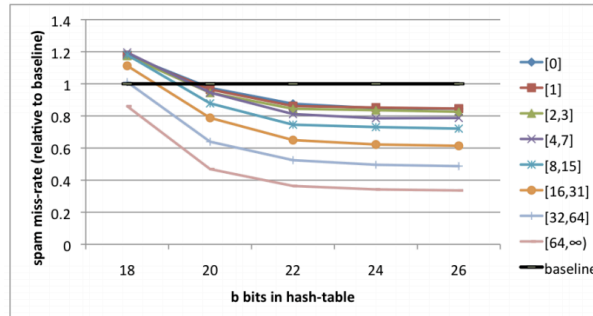


Figure 3. Results for users clustered by training emails. For example, the bucket [8, 15] consists of all users with eight to fifteen training emails. Although users in buckets with large amounts of training data do benefit more from the personalized classifier (up to 65% reduction in spam), even users 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 X
- Output: Set of k similar articles



X



©Emily Fox 2014

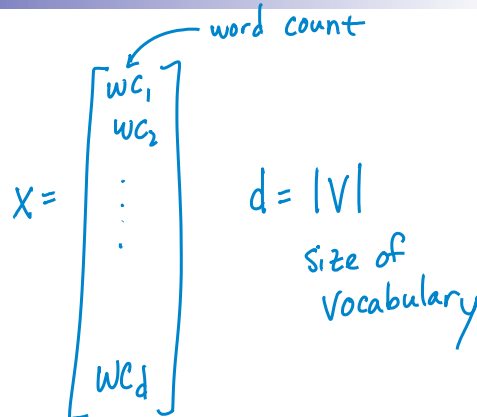
19

Document Representation

■ Bag of words model



ignore order of the words



©Emily Fox 2014

20

1-Nearest Neighbor

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

\swarrow query article

©Emily Fox 2014

21

k-Nearest Neighbor

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

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

$$\text{s.t. } \forall x^i \in \mathcal{X} \setminus \mathcal{X}^{NN} \leftarrow \text{"not in nearest neighbor set"}$$

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

©Emily Fox 2014

22

Distance Metrics – Euclidean

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

Or, more generally, $d(u, v) = \sqrt{\sum_{i=1}^d \sigma_i^2 (u_i - v_i)^2}$ *"scaled L2"*
 Equivalently, \leftarrow *weight for dim i*

$$d(u, v) = \sqrt{(u - v)' \Sigma (u - v)} \quad \text{where} \quad \Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & \sigma_d^2 \end{bmatrix}$$

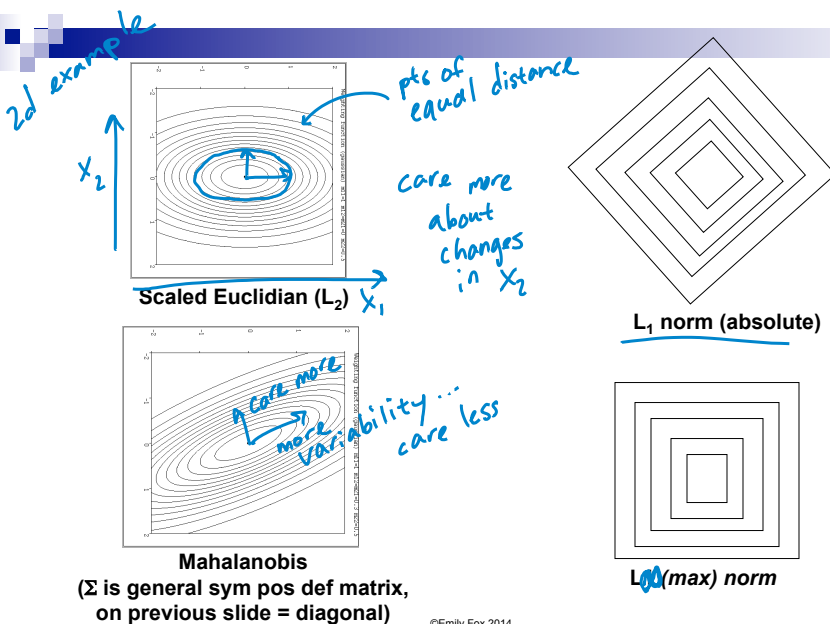
Other Metrics...

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

©Emily Fox 2014

23

Notable Distance Metrics (and their level sets)



©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

$$u \leftarrow \alpha u$$

$$v \leftarrow \alpha v$$

- 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

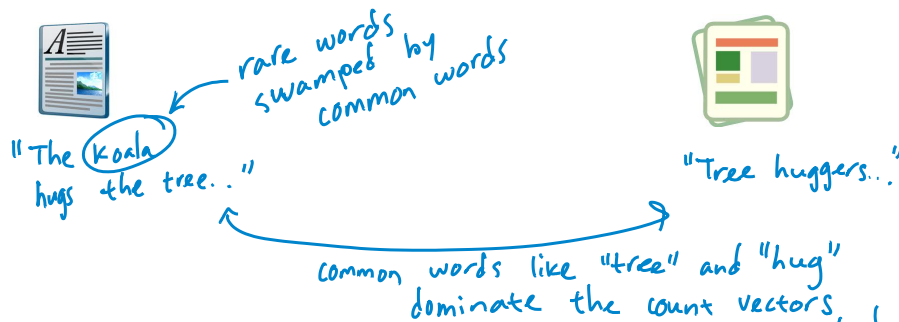
$$\|u\|_2 = \|v\|_2 = 1$$

©Emily Fox 2014

25

Issues with Document Representation

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



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

©Emily Fox 2014

26

TF-IDF

- Term frequency:

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

term ↑ *doc* ↑

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

$$\frac{f(t, d)}{\max\{f(w, d); w \in \mathcal{X}\}}$$

← prevent bias towards long doc

- Inverse document frequency:

$$idf(t, \mathcal{X}) = \log \frac{|\mathcal{X}|}{1 + |\{d \in \mathcal{X} : t \in d\}|}$$

→ 0 *t ∈ many docs d*
 > 0 *ow*

- tf-idf:

$$tfidf(t, d, \mathcal{X}) = tf(t, d) \times idf(t, \mathcal{X})$$

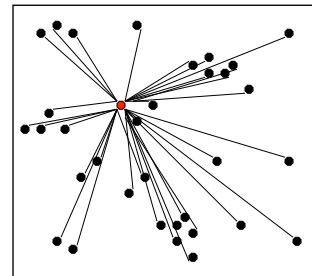
- High for document d with high frequency of term t (high "term frequency") and few documents containing term t in the corpus (high "inverse doc frequency")

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

↑ keep priority queue of top k + inserting into queue logk

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

KD-Trees

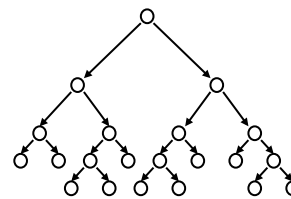
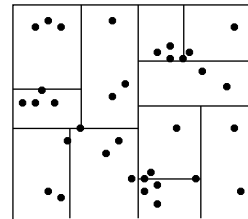
Smarter approach: *kd-trees*

- 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...

2d examples

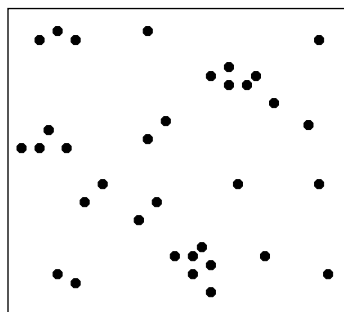


©Emily Fox 2014

29

KD-Tree Construction

dim 1 + dim 2



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

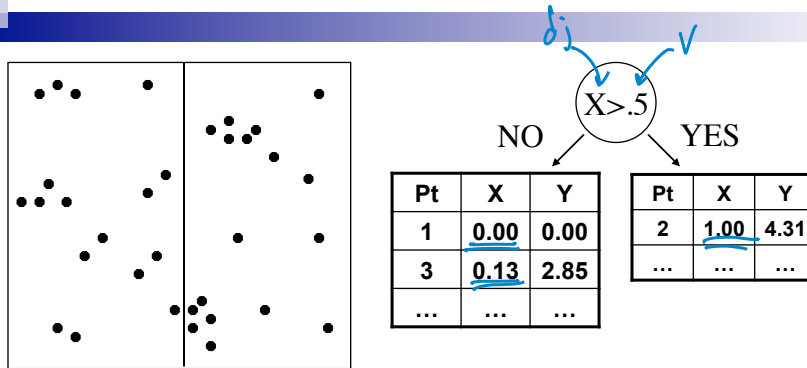
*x¹
x²
x³
⋮*

- Start with a list of d -dimensional points.

©Emily Fox 2014

30

KD-Tree Construction

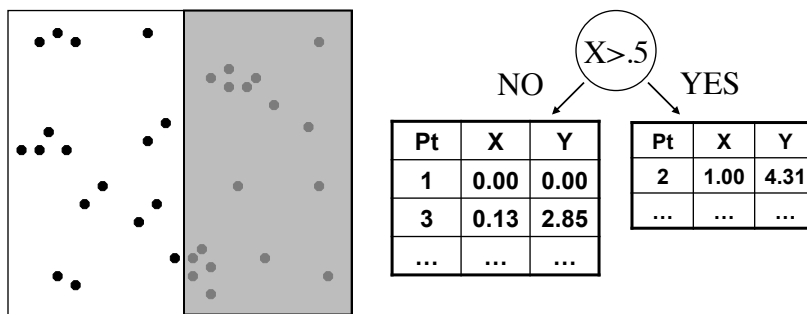


- Split the points into 2 groups by:
 - Choosing dimension d_j and value V (methods to be discussed...)
 - Separating the points into $x_{d_j}^i > V$ and $x_{d_j}^i \leq V$.
- Handwritten notes: $x \leq V$, $x > V$, X in this case, in this case, $V=0.5$*

©Emily Fox 2014

31

KD-Tree Construction

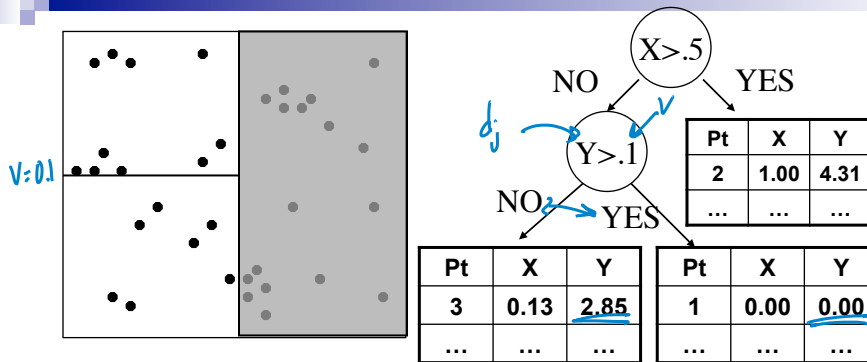


- Consider each group separately and possibly split again (along same/different dimension).
 - Stopping criterion to be discussed...
 - Handwritten note: how to choose dimension*

©Emily Fox 2014

32

KD-Tree Construction

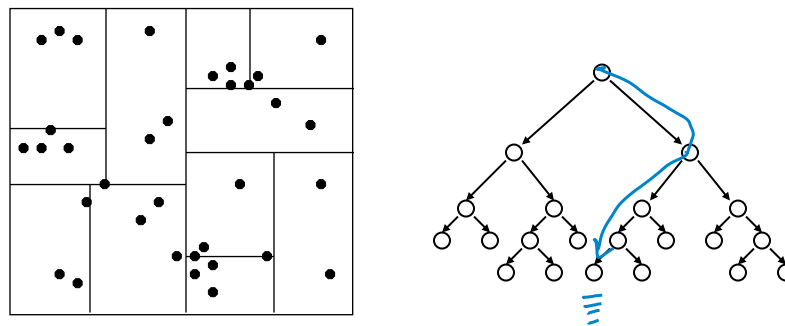


- Consider each group separately and possibly split again (along same/different dimension).
 - Stopping criterion to be discussed...

©Emily Fox 2014

33

KD-Tree Construction

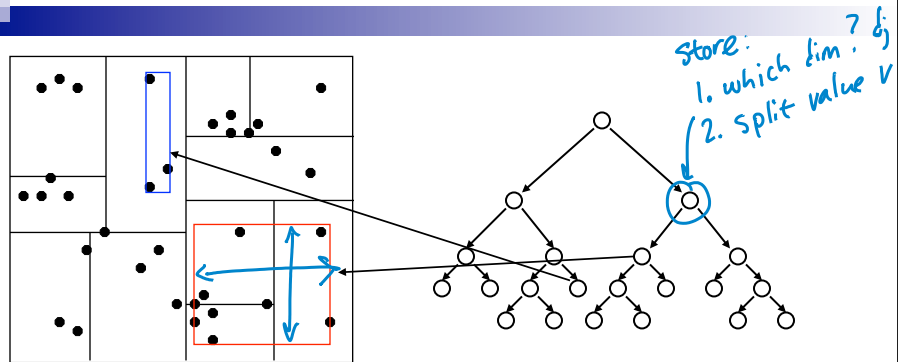


- Continue splitting points in each set
 - creates a binary tree structure
- Each leaf node contains a list of points

©Emily Fox 2014

34

KD-Tree Construction



- 3.
- 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

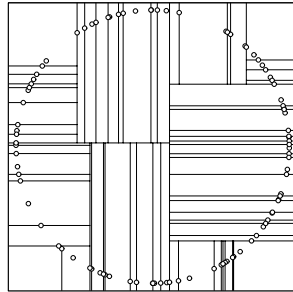
Use heuristics to make splitting decisions:

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

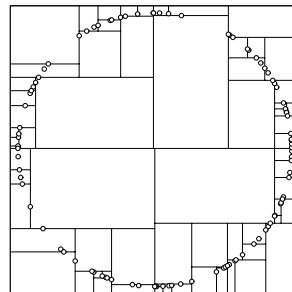
©Emily Fox 2014

36

Many heuristics...



median heuristic

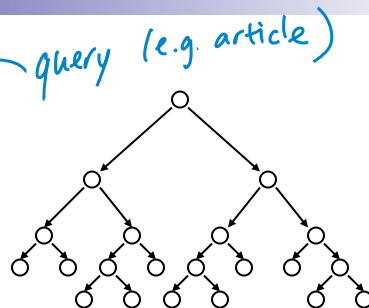
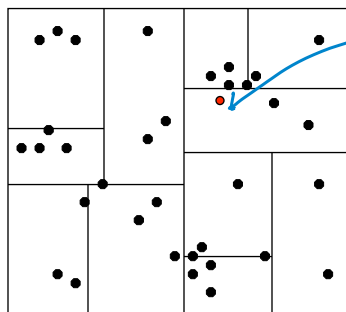


center-of-range heuristic

©Emily Fox 2014

37

Nearest Neighbor with KD Trees

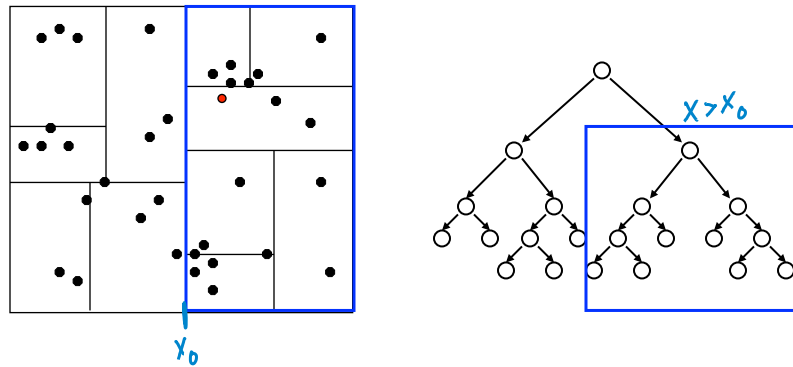


- Traverse the tree looking for the nearest neighbor of the query point.

©Emily Fox 2014

38

Nearest Neighbor with KD Trees

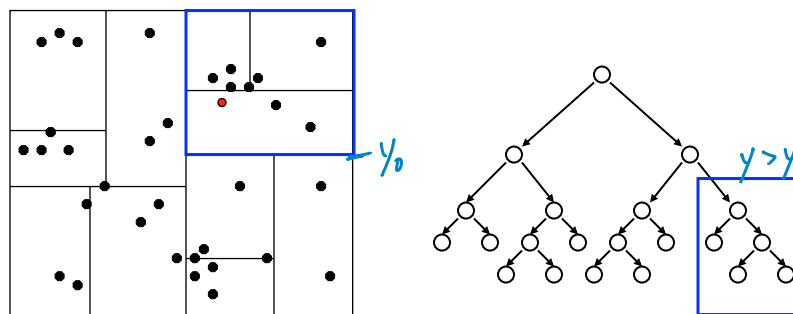


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2014

39

Nearest Neighbor with KD Trees

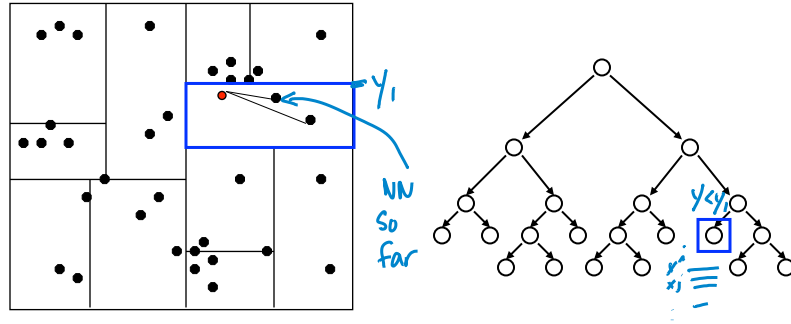


- Examine nearby points first:
 - Explore branch of tree closest to the query point first.

©Emily Fox 2014

40

Nearest Neighbor with KD Trees

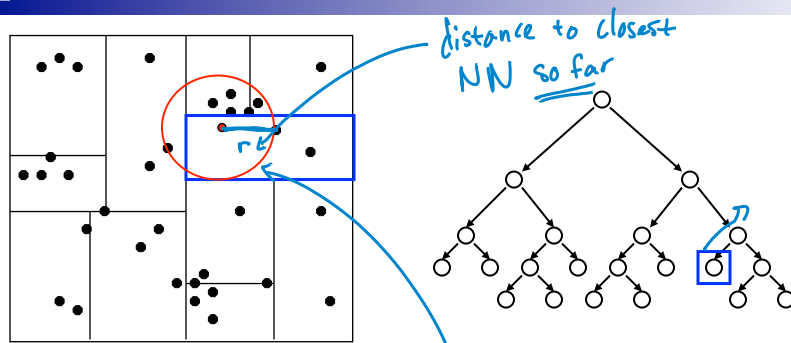


- When we reach a leaf node:
 - Compute the distance to each point in the node.

©Emily Fox 2014

41

Nearest Neighbor with KD Trees

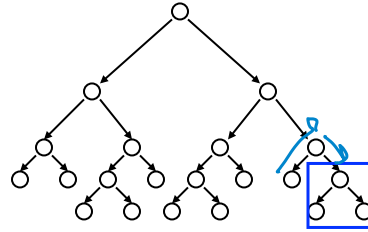
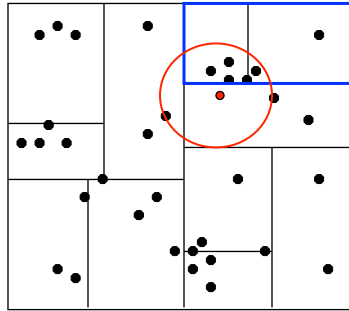


- When we reach a leaf node:
 - Compute the distance to each point in the node.
- does NN have to be in this box (blue)? NO

©Emily Fox 2014

42

Nearest Neighbor with KD Trees

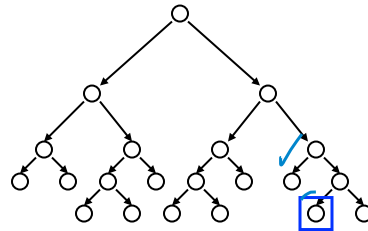
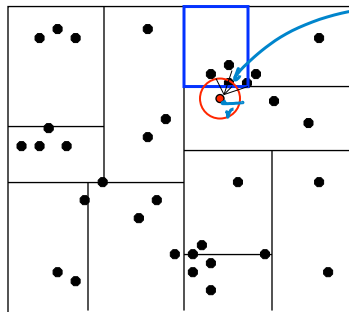


- Then backtrack and try the other branch at each node visited

©Emily Fox 2014

43

Nearest Neighbor with KD Trees

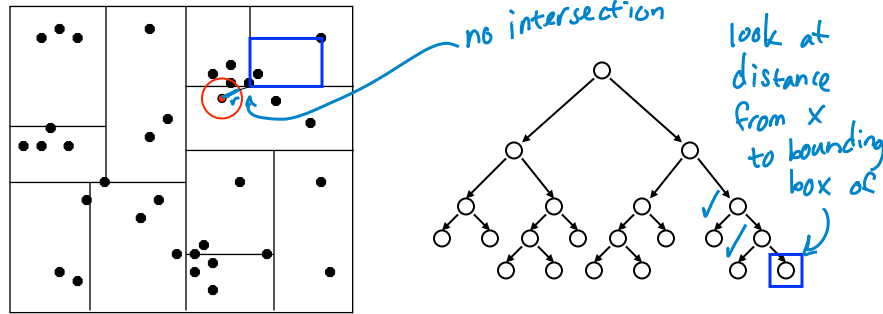


- Each time a new closest node is found, update the distance bound

©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

No article in this box could be the NN.