

Lect 19: K-d tree for fast Nearest Neighbor Search

In our previous lecture, we introduced the K-nearest neighbor decision rule. It has a Bayesian interpretation and an asymptotic error rate no more than twice the Bayesian error, as the number of training examples goes to infinity. This sounds promising. The drawback are

- 1, it needs to memorize all n-samples (n could be very large)
- 2, it needs to search the nearest neighbors on-line, complexity is $O(n)$.

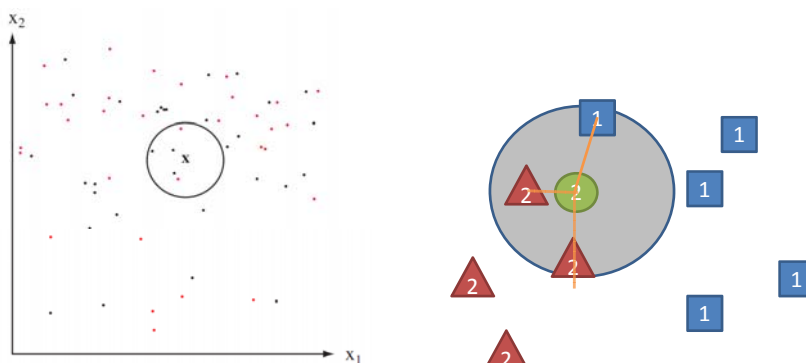
In this lecture, we introduce the Kd-tree technique for fast approximate nearest neighbor search:

Kd-tree: finding K-nearest neighbors in a d-dimension feature space.

This lecture note is based on a short tutorial (a chapter from a Ph.D. thesis) and a conference paper. Both are provided at the course webpage.

k-Nearest Neighbor Decision Rule

$$\omega(x) = \operatorname{argmax}_i \{k_1, k_2, \dots, k_c\}$$



Nearest Neighbor Search is Time-Consuming

Database:

$$X_1 = \{x_{11}, \dots, x_{1d}\}$$

$$X_2 = \{x_{21}, \dots, x_{2d}\}$$

.....

.....

.....

$$X_n = \{x_{n1}, \dots, x_{nd}\}$$

Query word:

$$X_q = \{x_{q1}, \dots, x_{qd}\}$$

Exactly nearest neighbor search is : $O(dn)$

If the dimension d is over 10 to 20, any methods for building search/index structures (including the Kd-trees) do no better than brute-force linear search. So people seek for approximate nearest neighbor. One may argue that approximate search is sufficient because the distance is often not precisely defined anyway.

Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

What is k-d tree

- k-d tree is a multidimensional **binary** search tree [Bentley '75]
- Structuring the codebook(exemplar-set) more intelligently
- Nearest Neighbor Search : **Given an input vector, find the r closest codeword in the codebook and output its index**

Input vector

$$X_q = \{x_{q1}, \dots, x_{qk}\}$$

Codebook

$$X_1 = \{x_{11}, \dots, x_{1k}\}$$

$$X_2 = \{x_{21}, \dots, x_{2k}\}$$

.....

.....

.....

$$X_n = \{x_{n1}, \dots, x_{nk}\}$$

r nearest neighbor

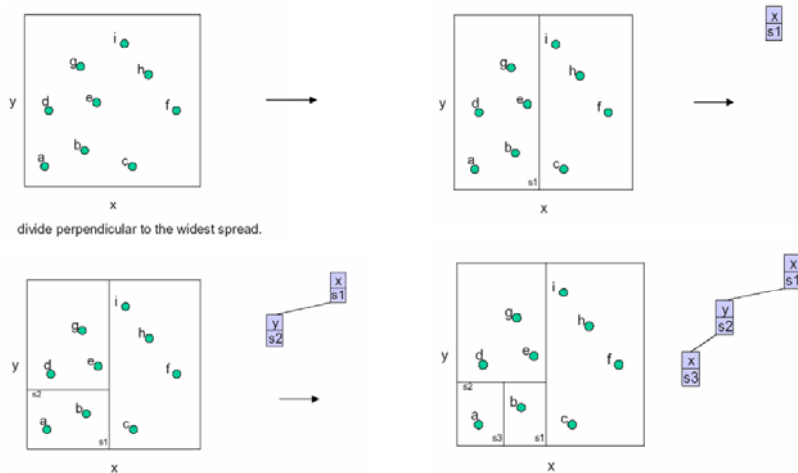
$$\begin{matrix} X'_1 \\ \vdots \\ X'_k \end{matrix}$$

Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

K-d tree construction

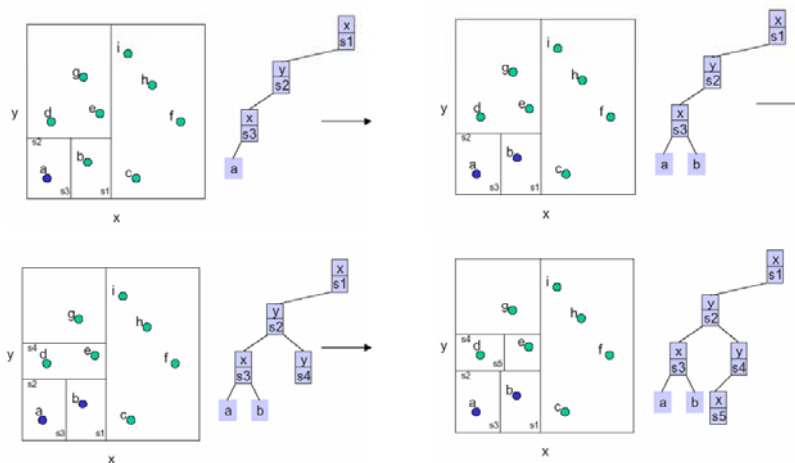
Each time, it chooses a dimension x that has the widest spread (variance), and divide the axis at the mean (or median) value.



Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

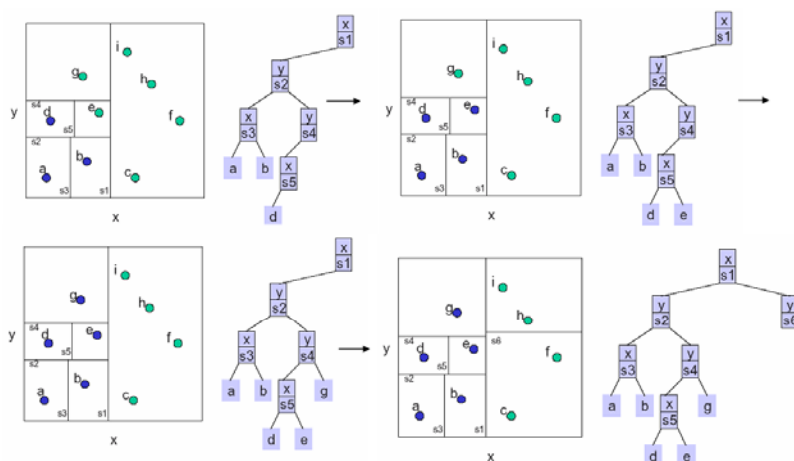
K-d tree construction



Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

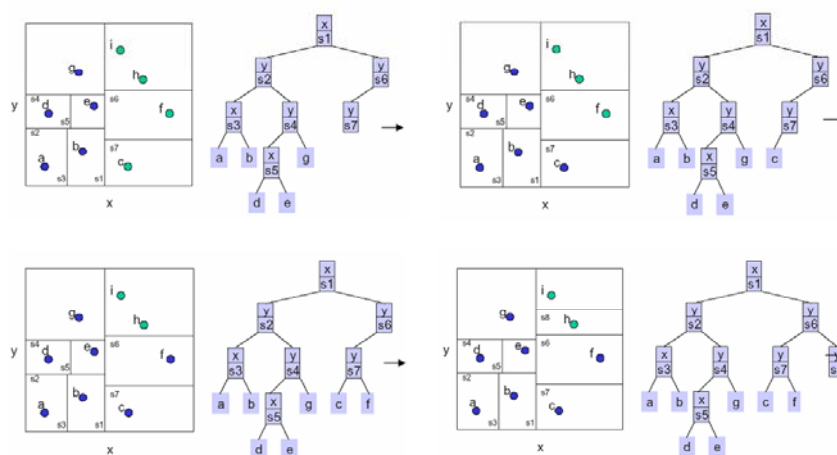
K-d tree construction



Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

K-d tree construction

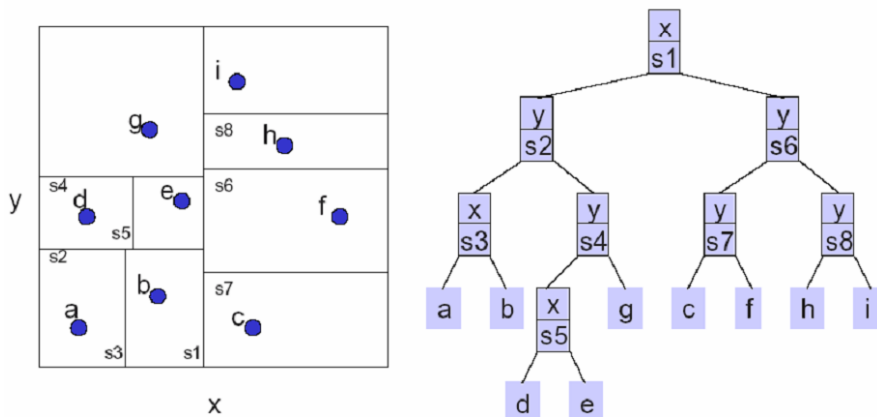


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Final K-d tree constructed ($d=2$)

Bifurcation based on mean values, the samples are on the leaf nodes. The search follows the tree.

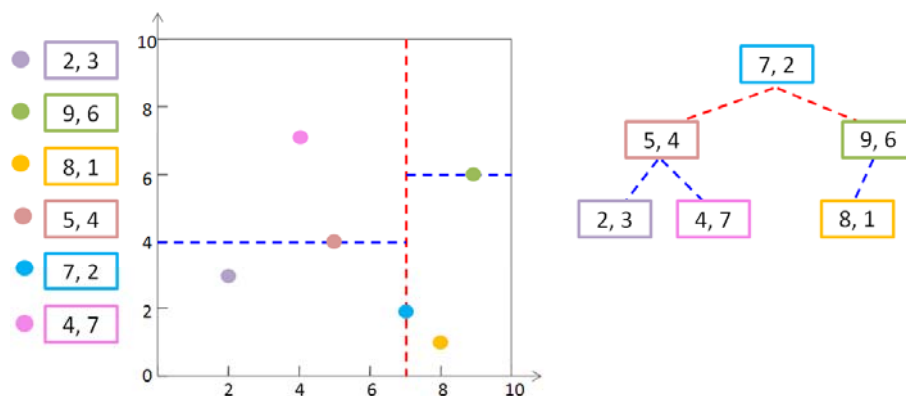


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

An example of Kd-tree based on median values

The inner nodes are also samples

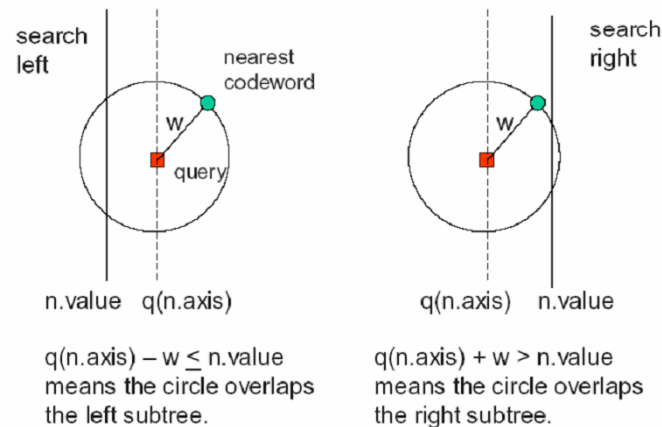


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Approximate Nearest Neighbor Search

Using a circle with radius w to tolerate errors. If the query point is within w -distance to the boundary, then it will search for both sides.

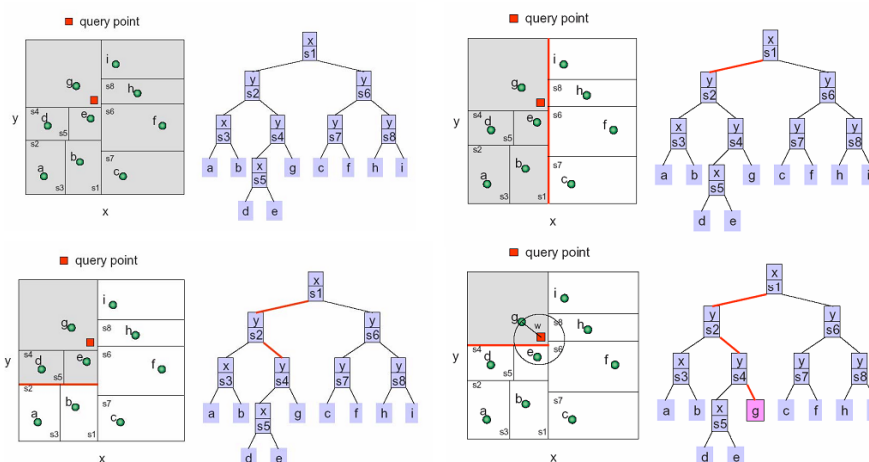


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Approximate Nearest Neighbor Search (example)

This is a depth-first-search: it goes all the way down the tree to get a neighbor (leaf-node g) that contains the query point, and then it back-traces to other leaf-nodes.

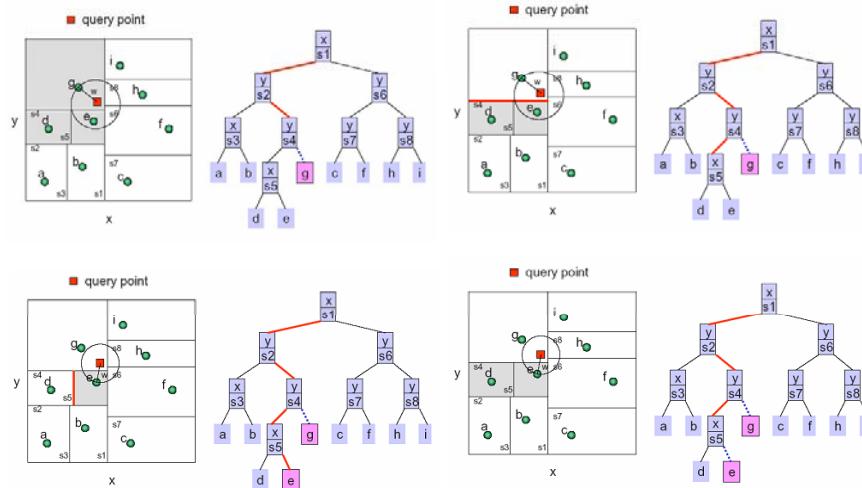


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Approximate Nearest Neighbor Search

But the real nearest point is **e** not **g**. It searches all branches that overlaps the w-circle.

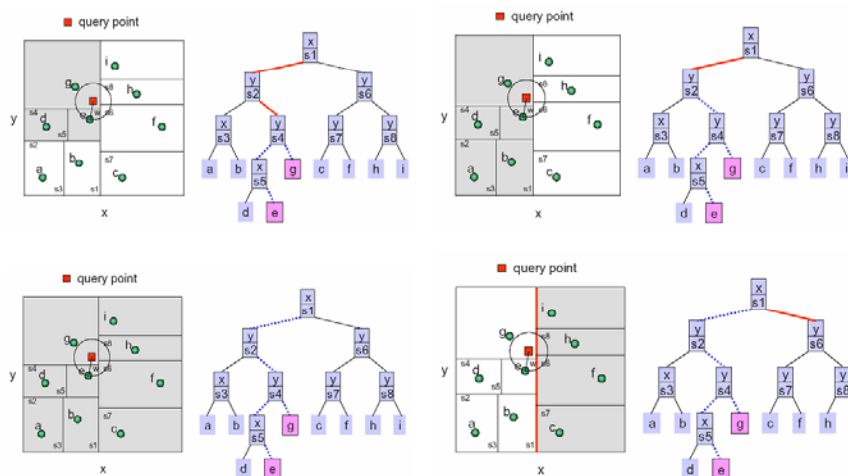


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Approximate Nearest Neighbor Search

Trace to the left lobe of the tree because the w-circle touch the boundary at the root node.



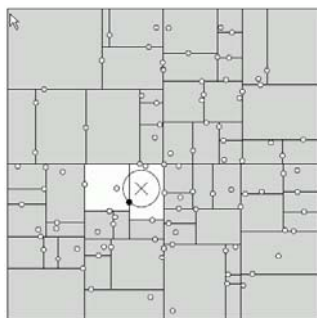
Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

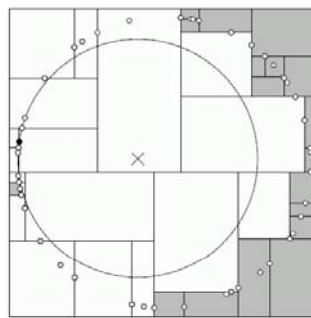
Complexity analysis

The number of searched nodes in the Kd-tree is between $O(\log n)$ and $O(n)$.

$$O(\log n) \leq c \leq O(n)$$



A nice case: 2 leaf nodes are searched.



A bad case: majority nodes are searched.
High-dimensional data are just like this.

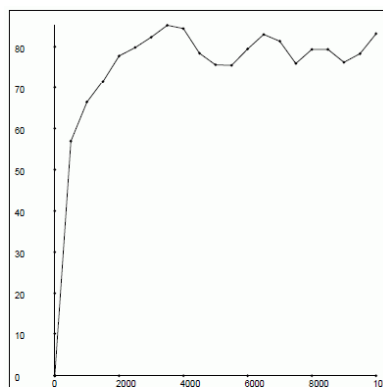
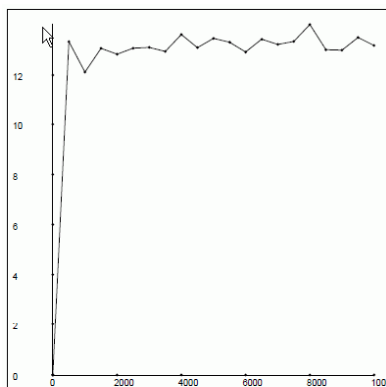
Examples illustrated by Andrew Moore, 1991 (see handout)

Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Complexity analysis: empirical results

Plot the number of leaf nodes inspected against the total leaf node number in the tree.



(Left) $K=4$ (vector dimension), but data distributed in a 3-dim space (intrinsic dimen).
(right) $K=8$ true dimension is 8.

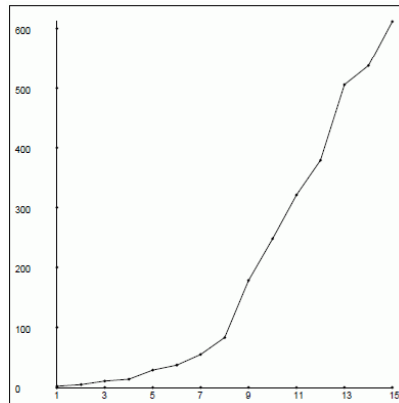
Each point is averaged over 500 searches in the tree.

Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Complexity analysis: empirical results

Plot the number of leaf nodes inspected against the dimension K .



When dimension is over 10-20, it will become more and more like the bad cases.

Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Randomized k-d trees

The Kd-tree search does return all the nearest neighbors within a tolerance radius w .

But such distance measure may not be robust. For example, if certain entry in the query vector has noise perturbation, then the search results can be largely affected.

To resolve this problem, people construct many Kd-trees by choosing different dimensions for bifurcations. During query, all these trees are searched.

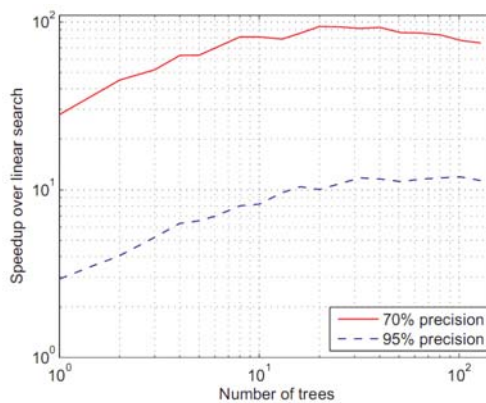
Other techniques for randomized search: Hashing.

Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

Multiple randomized k-d trees

Multiple trees are built, instead of only one in traditional k-d tree



Query point is q ,
True nearest neighbor is p^* ,
you only search it once and return p :

$$d(q, p) \leq (1 + \epsilon) d(q, p^*)$$

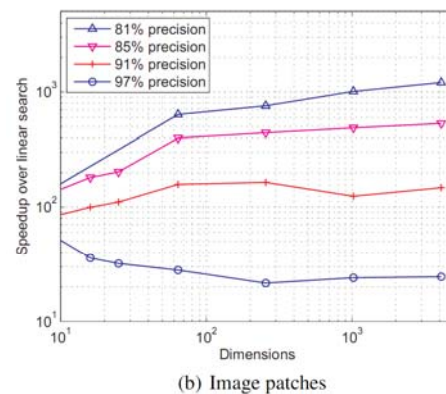
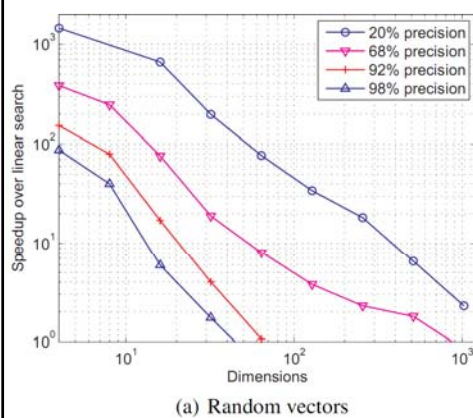
Marius Muja and David G. Lowe, "Fast Approximate Nearest Neighbors with Automatic Algorithm Configuration", in International Conference on Computer Vision Theory and Applications, 2009

Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

k-d tree

- Experiments: Data dimensionality

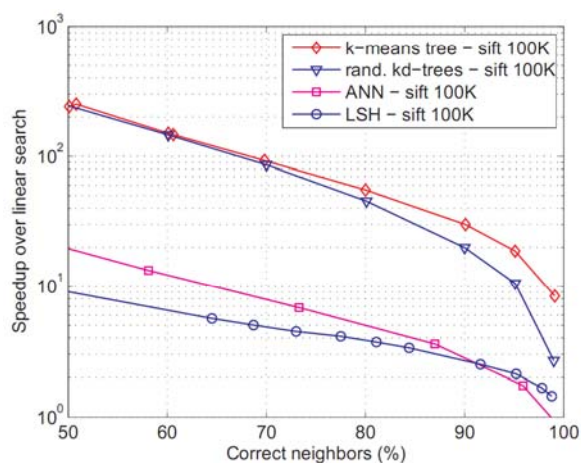


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

k-d tree

- Experiments: Search precision

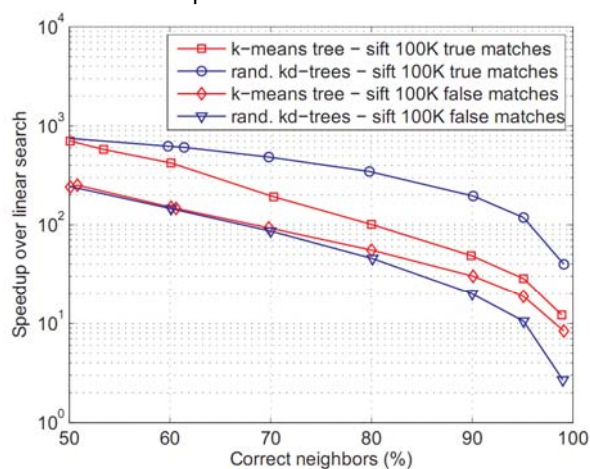


Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu

k-d tree

- Experiments: Search precision



Stat 231-CS276A: Pattern Recognition and Machine Learning

© S.C. Zhu