**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24

1

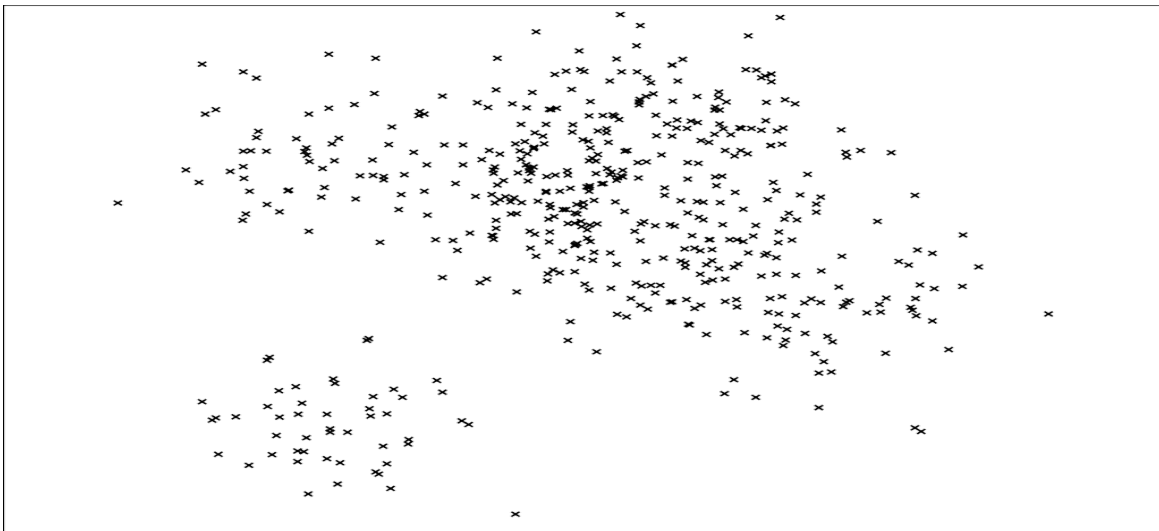# Approximate Nearest Neighbor Search
## Tree and Quantization based approaches



Source: Yusuke Matsui

# ANNOY – from O(nd) to O(log n)

Source: Erik Bernhardsson

## How do you build a tree from these vectors?

Source: Erik Bernhardsson

**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24

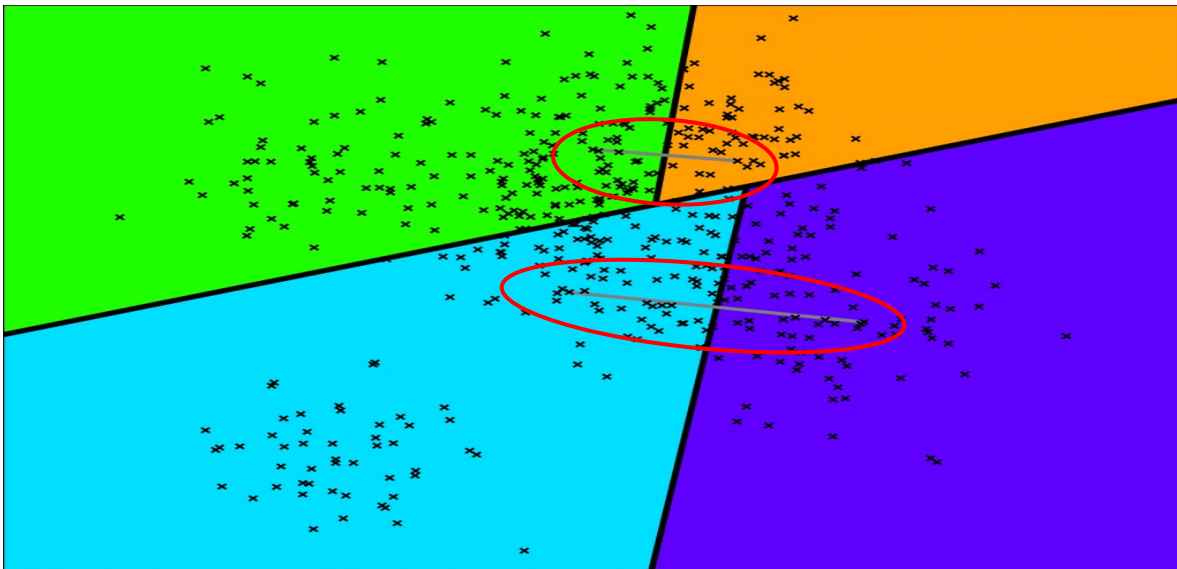## Pick two points randomly; split the feature space by the hyperplane equidistant from the two points.



Source: Erik Bernhardsson
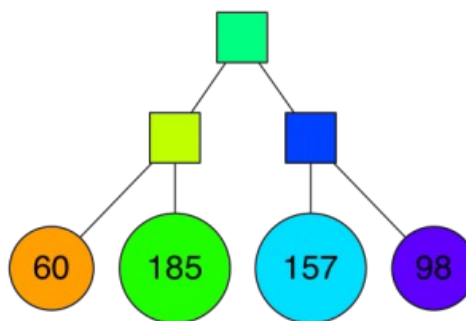
## Split each subspace recursively



Source: Erik Bernhardsson

**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24
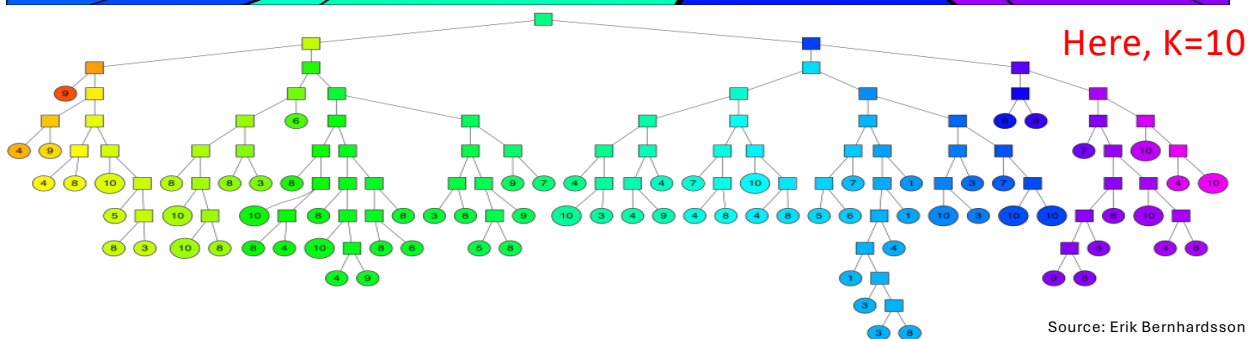
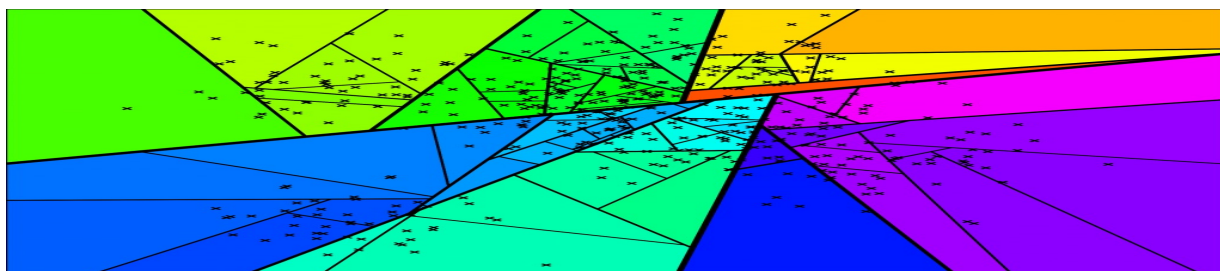# Split each subspace recursively



and the tree starts to evolve
(Intermediate node defines a hyperplane)



Source: Erik Bernhardsson

# Repeat until at most K items are left in each node



Here, K=10

Source: Erik Bernhardsson

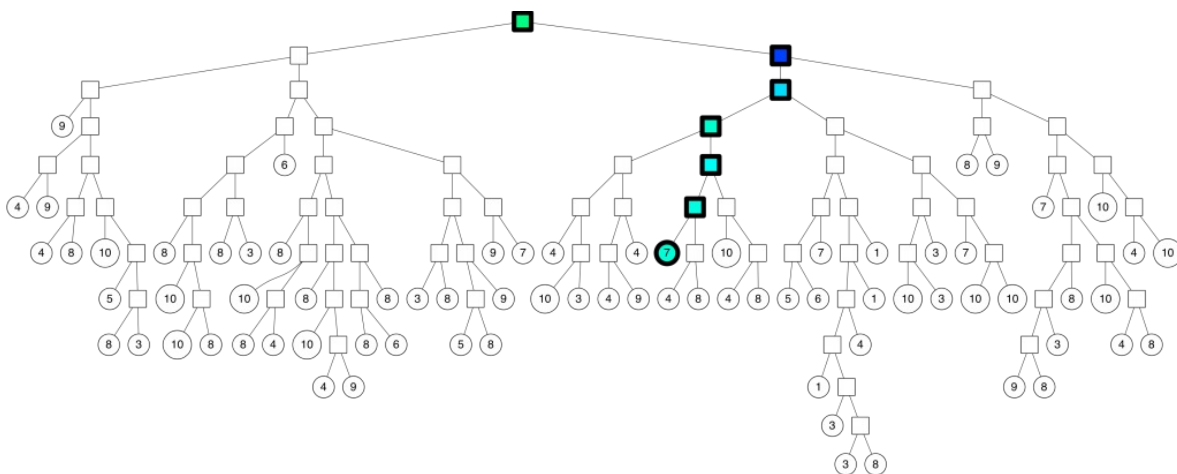**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24

# How do you find documents similar to a query?



Source: Erik Bernhardsson

# How do you search the tree?



Source: Erik Bernhardsson

**Strictly for class use. Do NOT share outside the class now or in future**
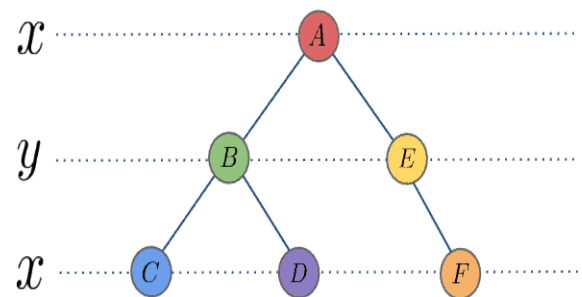
11/9/24

# A better solution: priority queue

Source: Erik Bernhardsson

# Even better solution: build a forest of trees

Source: Erik Bernhardsson

**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24

# Even better solution: build a forest of trees



Source: Erik Bernhardsson



## K-D Tree

Source: https://www.baeldung.com

# KD-Tree



### FLANN: Fast Library for Approximate Nearest Neighbors

Images are from [Muja and Lowe, TPAMI 2014]



Randomized KD Tree

k-means Tree

➤ Automatically select "Randomized KD Tree" or "k-means Tree"
  https://github.com/mariusmuja/flann

☺ Good code base. Implemented in OpenCV and PCL
☺ Very popular in the late 00's and early 10's
☹ Large memory consumption. The original data need to be stored
☹ Not actively maintained now

16

Source: Yusuke Matsui

**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24



Inverted File (IVF)

IVF : Intuition



1-NN: VORONOI DIAGRAM



Visualization of the Induced Decision Boundary

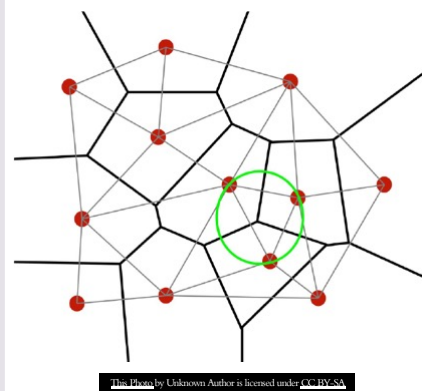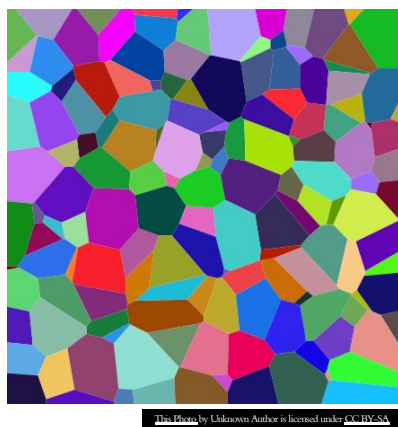**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24

# Inverted File Index (IVF): Centroids represent clusters



Source: Pinecone

# The inverted index is from centroids to the vectors in each cluster



Source: Pinecone

To search, find closest centroids and search in the corresponding clusters



Source: Pinecone

# Why IVF?

- Faster to build; the index size is much smaller
- However, search is slower than in HNSW
  $O(\log N) \; vs \; O(sqrt(N))$
- HNSW has a better recall as well
- A better approach: In IVF, search the nearest centroids using HNSW for better recall!

**Strictly for class use. Do NOT share outside the class now or in future**

11/9/24

# Product Quantization

## A Space and Computation saving technique

## Product Quantization

**How many vectors of floats are possible in a vector space?**

Infinite – each element of the vector can be any of the infinite floats

**How can we reduce the number of possible vectors (scope)?**

Approximate the floats by a representative finite range of integers

**In k-means, what is the representative vector in each cluster?**

Centroid – the "mean" for each of the k-clusters

**What if each of the element (dimension / feature) of the original vector is mapped to the number of the closest centroid in that dimension?**