



# Lecture 11: Kernel Methods; Locality Sensitive Hashing

🕒 Created @July 29, 2024 5:03 PM

## Locality Sensitive Hashing (LSH)

### How to choose lines that divide bins?

- Randomly
  - choose a slope between 0 and 90 degrees because 2 objects that have a high cosine similarity will have a high chance of being grouped together
  - You can look at adjacent bins to check if something did get micategorized

### Bin indexing

- 0 means that the bin is above that particular line
- 1 means that the bin is below that particular line
- the line is based on the 1-indexed index value
- algorithm
  - draw  $h$  lines randomly
  - find the bin index of all points
  - figure out which point you want to find a match for and find its bin index
  - do exact nearest neighbor search only in that bin for a similarity match

- look at nearby bins if you have time

## Optimization/fine-tuning

- We can look at more bins to overcome the tradeoff that comes with better efficiency

Pick a random hyperplane that separates the points

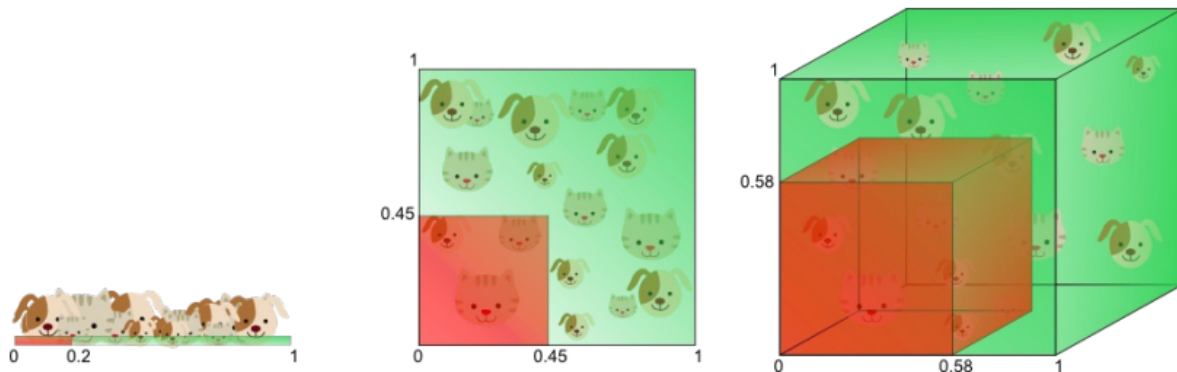
$h$  is the number of lines

$h = \ln(\# \text{ of dimensions})$

## Curse of Dimensionality

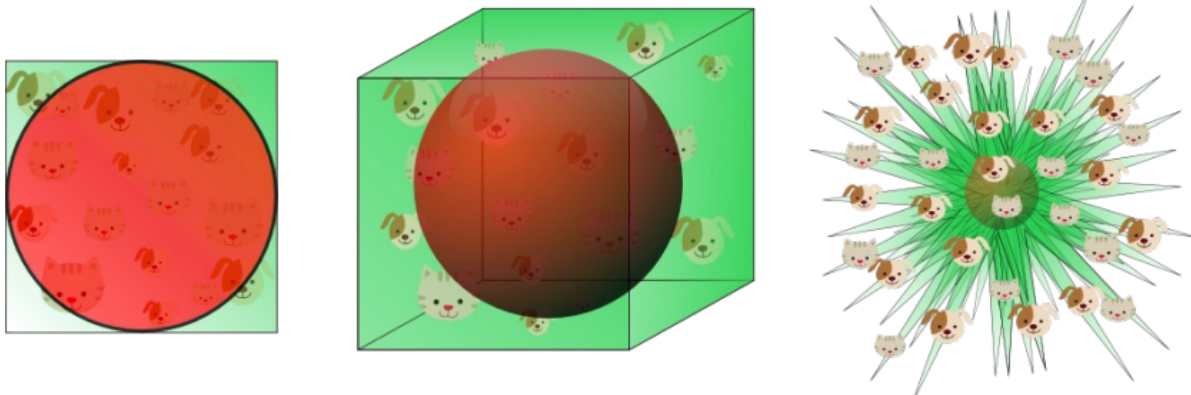
### Higher Dimensions

The sparsity of the data increases as the number of dimensions increases



So you need more data to reduce the sparsity (more space, so you need more data points to fill that space)

Most of the points become clustered at the corners of the shape as the number of dimensions increases



so your data becomes more sparse

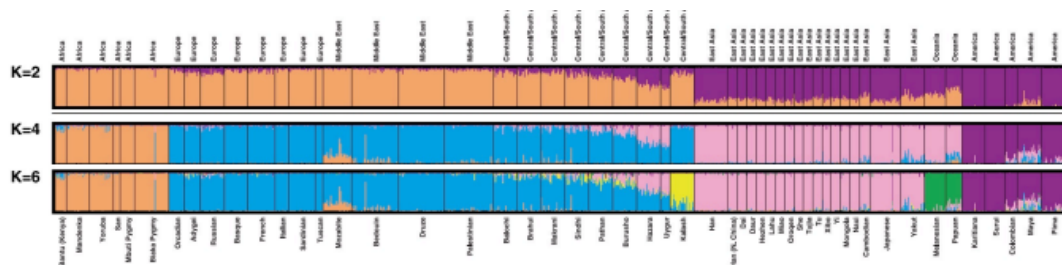
and also your nearest neighbors start becoming your farthest neighbors

Nearest neighbors start becoming farther and farther because they are on different dimensions

Lasso Regularization can be used to reduce the number of dimensions

## Clustering Overview

- unsupervised learning
  - CV, overfitting, bias-variance tradeoff, accuracy, error, etc don't apply because you don't have any labeled inputs or outputs
  - clustering gene sequences



- detects patterns in data that is not labeled
- quality metric in the ML pipeline is harder to evaluate for unsupervised learning because the outputs aren't labeled/set in stone
- Document retrieval case study

- given that someone read a sports article, what similar articles would you recommend
- if the data is labeled, you could use multi class classification
- but usually there isn't a defined boundary/difference between different categories or they might not be labeled
- clustering finds groups in a dataset
- a cluster is defined by
  - centroid: location of the center
  - spread: shape and size of the data
- 2 parts
  - find the clusters
  - assign each example to a cluster based on how close it is
  - closer distance means stronger similarity

## K-Means Clustering Algorithm

- algorithm to cluster different points together
- algorithm
  - first randomly initialize  $h$  points
  - map each point to the closest point
  - place the point in the average of the points in that cluster
  - repeat until convergence
- results heavily depend on the initialization
- k-means++ is more smart and is more likely to reach a local optima