

## K-MEANS CLUSTERING

Scheme to divide  $N$  data points into  $K$  clusters  
 $\{x_j\}$   $\{C_i\}$

$$r_{ji} = \begin{cases} 1, & \text{if } x_j \text{ is assigned to } C_i \\ 0, & \text{if } x_j \text{ is assigned to other clusters} \end{cases}$$

Distortion Measure:  $J = \sum_{j=1}^N \sum_{i=1}^K r_{ji} \|\vec{x}_j - \vec{\mu}_i\|^2$   
(Loss function)

Find  $r_{ji}$  &  $\vec{\mu}_i$  which minimize  $J$ .

### Iterative Process

step 0 - Choose  $\vec{\mu}_i$  randomly from given data points.

step 1 -  $r_{ji} = 1$  for  $(j, i)$  such that  $\vec{\mu}_i$  is the closest mean vector for given point  $\vec{x}_j$   
else  $r_{ji} = 0$

step 2 - For  $r_{ji}$  estimated in step 1, update  $\vec{\mu}_i$ :  
 $\nabla_{\vec{\mu}_i} J = 0 \Rightarrow 2 \sum_{j=1}^N r_{ji} (\vec{x}_j - \vec{\mu}_i) = 0$

$$\Rightarrow \vec{\mu}_i = \frac{\sum_j r_{ji} \vec{x}_j}{\sum_j r_{ji}}$$

Iterate till convergence.

(guaranteed)

could be local minima

Number of points in cluster  $C_i$

Convergence is slow & several variations available.

### Lossy Data Compression

For each of the  $N$  points, only store the identity of the cluster to which it belongs.  
& values of the  $K$  cluster centers,  $\vec{\mu}_i$   
uses significantly less data if  $K \ll N$ :  
also called vector quantization  
&  $\vec{\mu}_i$  code-book vectors

## Example of data compression

Each image has  $N$  pixels of  $\{R, G, B\}$  values stored with 8-bit precision.  
So each image needs  $24N$  bits

$K$ -classes needs  $\log_2 K$  bits

code-book vectors need  $24K$  bits

$\therefore$  Total no. of bits =  $24K + N \log_2 K$ .

$\therefore$  Compression Ratio =  $\frac{24K + N \log_2 K}{24N} \times 100\%$

$$= \left( \frac{K}{N} + \frac{1}{24} \log_2 K \right) \times 100\%$$

$$= 13.94\% \quad \text{for } K = 10 \text{ \& } N = 100 \times 100$$

### K-Means

Unsupervised

$K$  = no. of clusters into which the given data needs to be classified.

### KNN

Supervised

$K$  = no. of nearest neighbors to use for classification of a new data point

1

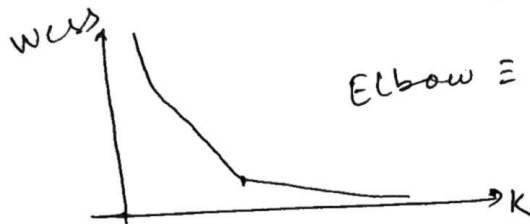
# K-Mean & K-NN

on K-Mean

- What will you do if number of clusters is not known?

WCSS [within cluster sum of squares]

$$WCSS(K) = \sum_{i=1}^K \sum_{j=1}^N (x_j - \mu_i)^2$$



Elbow  $\equiv$  Sudden change of behaviour

- ~~What~~ Can we do image clustering using this method?

- ~~Images~~ Similar images may have different orientation
- How to handle very large dataset?
  - ~~image set~~ - use minibatches.
- How to classify a new data point?

K-NN

~~when k is large what happens to bias of~~

small k - more noisy prediction  
large k - more stable prediction

- How to use K-NN for regression?

## K-Mean

- Unsupervised
- K - need to be determined by elbow method
- Does not work well with clusters of different sizes
- Used for clustering
- Eager learner

## K-NN

- Supervised
- $K \sim \sqrt{N}$
- Size difference is not a problem
- Used for classification or regression
- Lazy learner