

## Lecture November 5

### k-means Clustering

Data set  $\{x_0, x_1 \dots x_{n-1}\}$   
of  $n$ -observations of  
a  $D$ -dimensional random  
variable  $X$

The aim is to find a  
partition of the data  
into same number of  
clusters -  $k$  -

- Need to find assignments  
of points belonging to  
a given cluster -  $k$  -
- Define the cluster by  
its position  $\mu_k$  (center)

Optimize

$$C = \sum_{i=0}^{n-1} \sum_{k=0}^{K-1} a_{ik} \|x_i - \mu_k\|_2^2$$

$a_{ik} = \begin{cases} 1 & \text{if within} \\ 0 & \text{else,} \end{cases}$

$$a_{ik} = \begin{cases} 1 & \text{if } k = \arg \min \underbrace{\|x_i - \mu_k\|} \\ 0 & \text{else} \end{cases}$$

- First choose some initial value  $\mu_k$  ( $K$ )  
optimize  $a_{ik}$  while keeping  $\mu_k$  fixed.
- with  $a_{ik}$  fixed,  
optimize  $\mu_k$

Derivatives w.r.t  $\mu_k$

$$2 \sum_{i=0}^{n-1} a_{ik} (x_i - \mu_k) = 0$$

$$\mu_K = \frac{\sum_i n_{ik} x_i}{\sum_i n_{ik}}$$

$\nearrow$   
 # points in  
 a cluster.

The values of  $\mu_K$  are defined by the mean values defined by the data points in each cluster.

## Decision Trees

ensemble methods {
 

- Bagging (random forest)
- random forests (different trees)
- Voting
- Boosting, gradient boosting

Typical Regression case



