

Supervised:  $Y \in \mathbb{R}$  or  $Y \in \mathcal{C}$  and  $X \in \mathbb{R}^p$

If  $(X, Y)$  are random then they have some joint dist.  $p(x, y)$ .

Supervised Learning is basically density estimation where we want to estimate  $p(y|x)$ .

prob. of seeing a  $y$  given a  $x$

often we just want

$$\hat{f} = \underset{f}{\operatorname{argmin}} E_{Y|X} L(f)$$

some mean we estimating

Bayes Theorem says that

$$p(x, y) = p(y|x) p(x)$$

supervised interest in this

not so much this

Unsupervised Learning

Just have  $X \in \mathbb{R}^p$   
(no  $Y$ )

We want to estimate  $p(x)$

→ summarize

## Quirks of Unsupervised

→ no direct measure of success

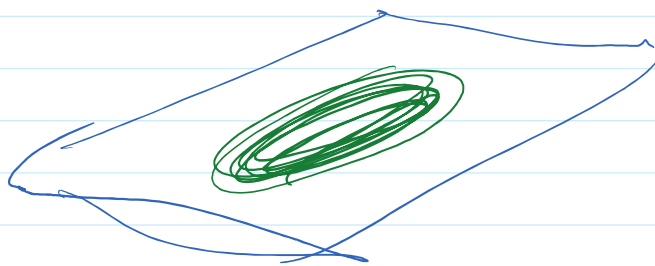
→ so many methods & opinions

## How to summarize $p(x)$

→ if  $P$  is low (low-dim'l problem)  
directly estimate  $p(x)$  / visualize

→ if  $P$  is high  $\Rightarrow$  ??? summary stats?

or PCA: find a low-dim'l subspace  
that contains high density parts of  $p(x)$



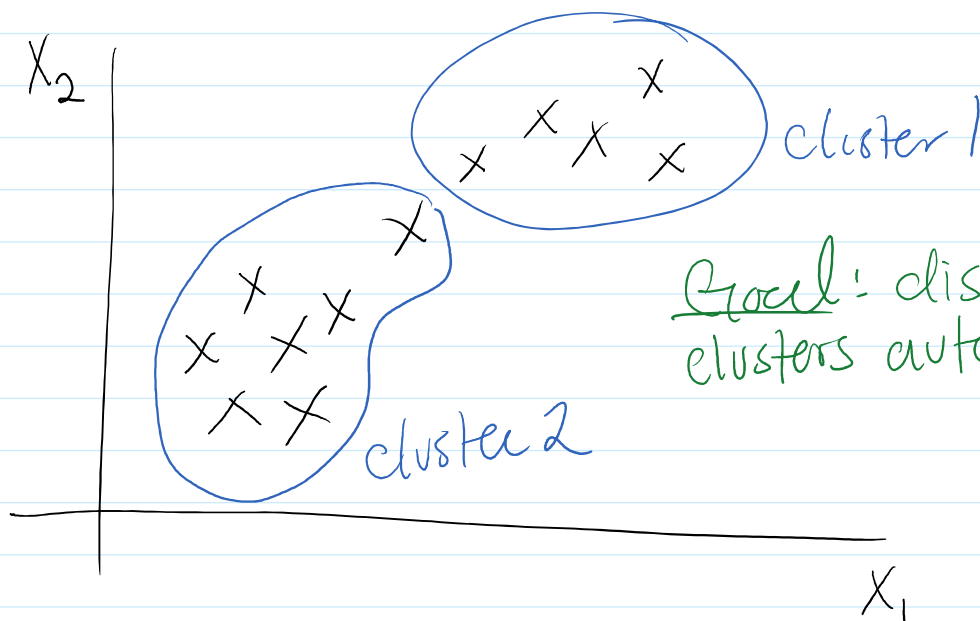
---

## Clustering:

Basically tries to find convex sets of

high density parts of  $p(x)$ .

Classic: want to create "clusters" w/in our data so that pts w/in a cluster are more similar to each other than points in different clusters.



Goal: discover clusters automatically

Sim/Dissim

To do clustering, we need some measure of either similarity or dissimilarity.

Ultimately (almost) all clustering algo just need a measure of sim/dissim b/w all pts.

If I have  $N$  objects then we can form a matrix

$D_{N \times N}$

1st 2nd ...  $i$ th ...  $j$ th ...  $N$

so that  $D_{ii} = \text{dissim btwn } i^{\text{th}} \text{ and } i^{\text{th}} \text{ object}$

[If I have a measure of sim. btwn objects,  
I can create  $D$  using a decreasing  
transformation.]

Many/most cluster algs only need  $D$ , don't  
even need some  $N \times P$  matrix  $X$ .

Properties req. for  $D$

①  $D = \begin{bmatrix} 0 & & \\ & \ddots & \\ & & 0 \end{bmatrix}$  ( $D_{ii} = 0$ )  
zeros down main diag

②  $D_{ii} \geq 0$  (non-neg.)

③  $D = D^T$  (symmetric)

---

If we have features  $X_{N \times P}$  we can form  
attribute-based dissims

i.e. use diff btwn feats to calc  $D$ .

$$D_{ii} = d(x_i, x_i) = \sum^P d_{ij}(x_i, x_i)$$

*(Red arrows point from the text "dissim for" to the  $d_{ij}$  term, and from "dissim meas for feat j" to the  $j$  index in the summation.)*

$$D_{ii'} = d(x_i, x_{i'}) = \sum_{j=1}^P d_j(x_i, x_{i'})$$

$\swarrow$   $\nearrow$   $\nearrow$   
 rows  $i$  and  $i'$   
 of  $X$

$\nwarrow$   $\nearrow$   $\nearrow$   
 for feat  $j$

Caveat: careful w/ scale

Ex. ① Numeric Feature  $j$  Euclidean (Squared) dist.

$$d_j(x_i, x_{i'}) = (x_{ij} - x_{i'j})^2$$

② Ordinal Feature  $j$

transform level  $i$  to  $\frac{i - 1/2}{M}$  for  $i = 1, \dots, M$

and treat as numeric (eg. use Euclidean)

③ Categorical Feature  $j$

$$d_j(x_i, x_{i'}) = \mathbb{I}(x_{ij} \neq x_{i'j})$$

$$\begin{cases} 1 & \text{if same} \\ 0 & \text{else} \end{cases}$$

---

Combinatorial Clustering Algos

Assume data comes from one of  $K$  clusters

$$G_1, \dots, G_K$$

want to do: assign each  $x_i$  to some  $G_k$

How: I do this assignment to minimize some measurement ("Loss") of not being clustered.

Classic measurement

$$W = \begin{matrix} \text{total} \\ \text{w/in cluster} \\ \text{dissim} \end{matrix} = \sum_{k=1}^K \sum_{i, i' \in G_k} D_{ii'}$$

= total dissim measure  
across elements in  
some cluster.

Should be small if clustered well  
// large if not clustered well

Note that if

$$T = \sum_{i, i'} D_{ii'} = \text{total dissim}$$

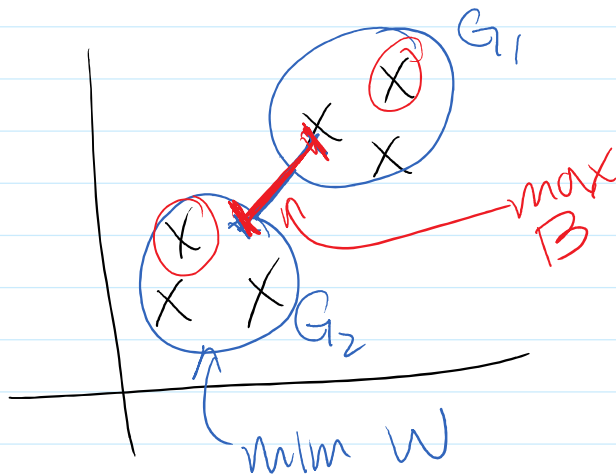
and  $\begin{matrix} \text{total} \\ \text{w/in cluster} \end{matrix} \sum_{k=1}^K T_k = T$

$$B = \begin{matrix} \text{total} \\ \text{btwn cluster} \\ \text{dissim} \end{matrix} = \sum_{k=1}^K \sum_{i \in G_k} \sum_{i' \notin G_k} D_{ii'} \\ = \text{total dissim btwn pairs} \\ \text{in diff clusters}$$

Note:  $T = W + B$

To find good  $G_s$  we can

- ① minimize  $W$
- or ② Since  $W = T - B$  we can maximize  $B$



How to find  $G_s$ ?

Ideally, consider all possible  $G_s$

Ideally, consider all possible  $\gamma$ s

Not computationally tractable

Ex,  $N=19$  and  $K=4$

then # of possible assignments is  $\sim 10^{10}$

Soln: Greedy.

K-Means:

Today: all features are numeric, and

$$D_{ii'} = \|x_i - x_{i'}\|^2$$

$\hookrightarrow$  squared euclidean dist.

In this case we can show that

$$W = \sum_{k=1}^K N_k \sum_{i \in G_k} \|x_i - \bar{x}_k\|^2$$

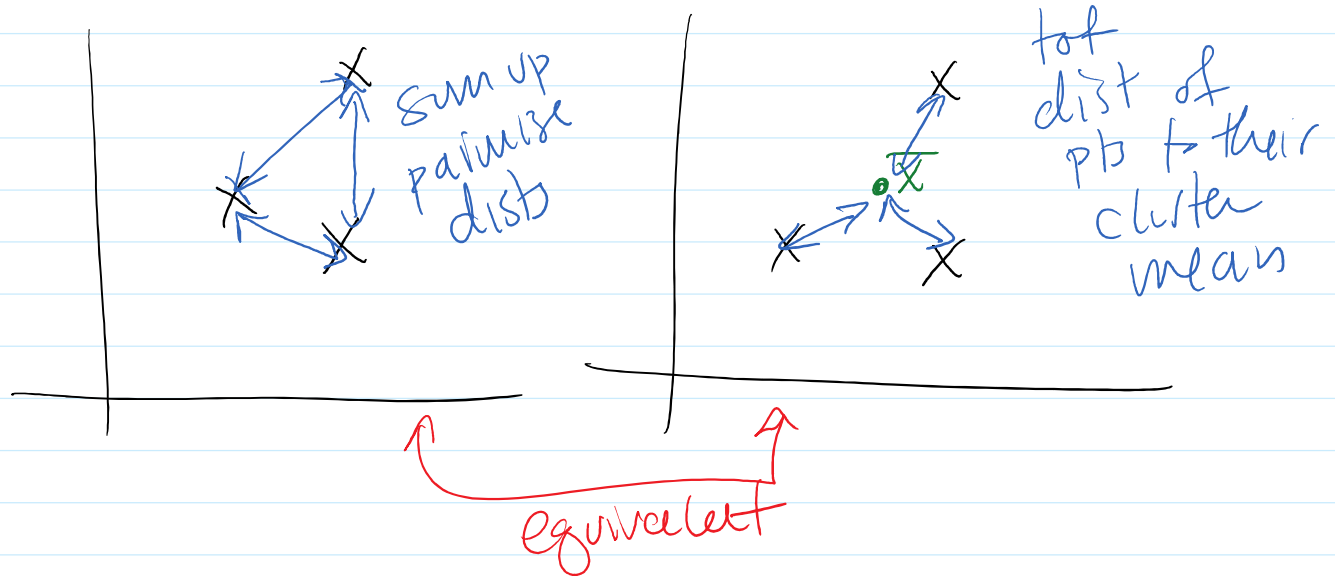
$N_k = \#$  in  $k^{\text{th}}$  cluster

$\bar{x}_k = \text{mean of } k^{\text{th}} \text{ cluster}$

$$\left[ \text{Recall! } W = \sum_{k=1}^K \sum_{i, i' \in G_k} D_{ii'} = \sum_{k=1}^K \sum_{i, i' \in G_k} \|x_i - x_{i'}\|^2 \right]$$

tot. pairwise dist. of pts in a cluster





Of course Goal: choose  $G_1, \dots, G_K$  so that we minimize  $W$ .

Classic Way of doing: Lloyd's Algorithm

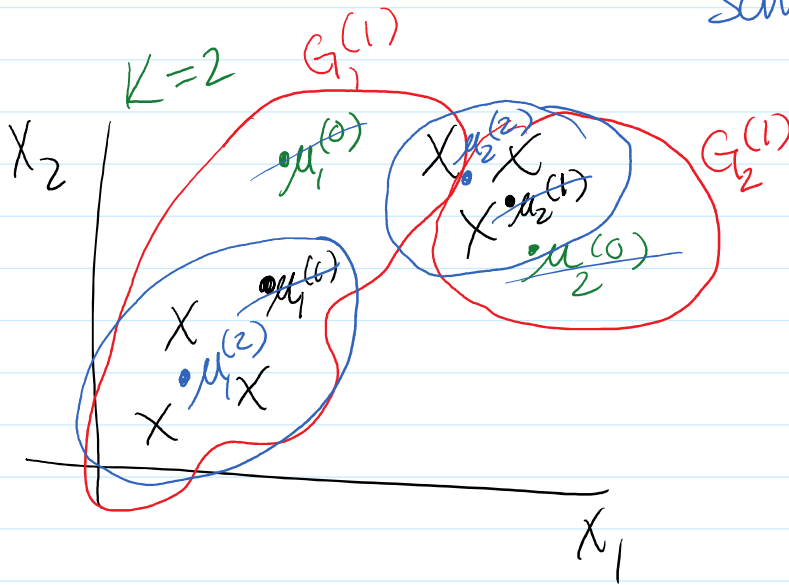
① Initialize Step:

Make initial (random) guesses of group means:  $\mu_1^{(0)}, \dots, \mu_K^{(0)}$

Loop over some iterations  $t=1, 2, 3, 4, \dots$   
where at the  $t^{\text{th}}$  iteration

① Assignment Step assign each  $x_i$  to the Group w/ the closest mean  $\mu_k^{(t-1)}$

② Update Step: re-compute  $\mu_s$  as sample means of new assignments



Do this until means stop moving.

---