Midterm: Friday Mar 4 1:10–3pm
UC273
2-sided sheet & calculator

## Hierarchical clustering

Start with $n$ clusters ($n$ observations)
— sucessing group together observations/cluster that are close.

have 2 clusters $U, V$
$d(U,V) = ?$

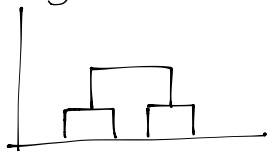Single linkage
$d(U,V) = \min(d(x_i, x_j) : x_i \in U, x_j \in V)$

Complete linkage
$d(U,V) = \max(d(x_i, x_j) : x_i \in U, x_j \in V)$

Average linkage
$d(U,V) = \text{ave}(d(x_i, x_j) : x_i \in U, x_j \in V)$

Result   clustering tree (dendrogram)



Blackboard
Clustering of
atheletes
record data

— dendrograms for 3 methods
— KOR clearly "different"

Other note:
— often do hierarchical clustering with "similarity" matrices
— high similarity = low distance
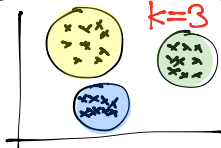e.g. single linkage define similarity b/w clusters = max pairwise similarity

## Towards model-based clustering

k-means clustering
$\overbrace{\phantom{xxxxxx}}^{\text{mixture model}}$

Informal model: $X$ has density $f(x) = \theta_1 f_1(x) + \cdots + \theta_k f_k(x)$
where $\theta_1, \ldots, \theta_k > 0$ with $\theta_1 + \cdots + \theta_k = 1$
— Also assume that $f_i(x) f_j(x) \doteq 0$ for $i \neq j$



k=3

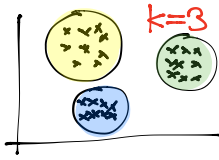— $f_1(x), \ldots, f_k(x)$ represent density of sub-populations
— $\theta_1, \ldots, \theta_k$ represent sub-pop'n proportions.
— in practice, $f_1(x), \ldots, f_k(x)$ and $\theta_1, \ldots, \theta_k$ are unknown
— If we know $f_1(x), \ldots f_k(x)$ then given an observation $x^*$, we
can predict very well which sub-pop'n it belongs to.

Problem: Given data $x_1, \ldots, x_n$, how to determine $k$ clusters (corresponding to $k$ sub-
pop'ns) of observations.
— Start by assuming $k$ is known.

$k=3$

— need to determine
- center points (centroids) of the $k$ clusters
- determine shape of clusters
- which observations belong to which clusters

For arbitrary centroids $\mu_1, \ldots, \mu_k$, find groups of observations $\underbrace{G_1, G_2, \cdots G_k}_{\text{disjoint}}$ to minimize $\underbrace{\sum_{j=1}^{k} \sum_{i \in G_j} d(x_i, \mu_j)}_{\text{minimize}} \Rightarrow$ fixed

$\Rightarrow G_1^*, G_2^*, \ldots, G_k^*$ depend on $\mu_1, \ldots, \mu_k$

Now minimize (w.r.t. $\mu_1, \ldots, \mu_k$),
$$\sum_{j=1}^{k} \sum_{i \in G_j^*} d(x_i, \mu_j) = g(\mu_1, \ldots, \mu_k)$$

In general, computationally very difficult!

$\underline{\text{Special case:}}$ $d(x_i, \mu_j) = \|x_i - \mu_j\|^2 = $ squared Euclidean norm

$\Rightarrow$ $k$-means

Re-express optimization problem:

$$\min_{\mu_1, \ldots, \mu_k} \min_{G_1, \ldots, G_k} \sum_{j=1}^{k} \sum_{i \in G_j} \|x_i - \mu_j\|^2$$

$$= \min_{G_1, \ldots, G_k} \min_{\mu_1, \ldots, \mu_k} \sum \sum \underbrace{\| \cdots \|^2}_{\text{same as above}} \longrightarrow \text{minimized at mean of points in the cluster } j$$

$$= \min_{G_1, \ldots, G_k} \sum_{j=1}^{k} \sum_{i \in G_j} \underbrace{\|x_i - \overline{x}_{G_j}\|^2}_{\text{within group (cluster) sum of squares}}$$

where $\overline{x}_{G_j} = \dfrac{1}{\text{number of pts in } G_j} \sum_{i \in G_j} x_i$

$\underline{\text{Comments}}$

① Still computationally very hard but good algorithms exist.

② In $k$-means clustering, shapes of clusters are all spheres (spheroids)

— can often transform variables (e.g. look at PC scores) s.t. shape assumption is not too severe.