

Outline: ~~partition into~~ (for now): partition into  $k$  non-empty disjoint sets.

- clustering: ~~partition into~~ (1)
- prove equivalence of within-cluster scatter & w/in cluster variation (2)

looks good, but computationally hard

<resolve  $\Delta$  inequality "paradox">

suggests an algorithm, but why is this optimizing a "good thing"?

Lecture 10  
3/3/15  
clustering algs:  
k-means,  
single-link

k-means

- give k-means alg, including examples.
- show k-means is improving its opt. criterion (no guarantee of optimality, tho')
- ~~show some successes & failures~~  $\leftarrow$  do we want an algorithm? or just point to alg's?

single-link

consider alternate criterion:  
- mention problem w/ (6) last time  
spacing: maximize between-cluster distance, defined as

$$d(C_i, C_j) = \min_{x \in C_i, x' \in C_j} \|x - x'\|_2^2$$

(actually, other functions ok here, too, only require  $f(x, x) = 0$ ,  $f(x, x) = f(x', x)$ ,  $f(x, x') > 0$  if  $x \neq x'$ .)

- single-link clustering finds an optimal solution
- behaviors

### Announcements:

- typo in AI Q2 re: error bars. OK if you used "typo" version - don't redo.
- updated lec 9 handout posted  $\leftarrow$  or just ~~be~~ print new copy?

From last time:

$n$  data points  $x_1, \dots, x_n$

$k$  clusters  $C_1, \dots, C_k$

for cluster  $C_j$ ,

$n_j = \#$  points in  $C_j$ . Req'd to be  $> 0$ .

$r_j = \frac{1}{n_j} \sum_{x \in C_j} x$ , the centroid.

think of "r" as standing for "representative".

(1) within-cluster scatter:

$$\sum_t \sum_{x_t \in C_j, t \leq s} \|x_t - x_s\|_2^2$$

(2) within-cluster variation

$$\sum_j n_j \sum_{x \in C_j} \|x - r_j\|_2^2$$

(8) lemma: for any point  $z$ ,

$$\sum_{x \in C_j} \|x - z\|_2^2 = \left( \sum_{x \in C_j} \|x - r_j\|_2^2 \right) + n_j \|r_j - z\|_2^2$$

argument follows Hopcroft k-means { Kannan 2014, section 8.2 }

Idea: (1) - general, intuitive, but computationally hard?

(ii) - if (8) were true, then ~~some algorithm~~ (1)  $\equiv$  (2)  
(we'll omit the pf, which is just algebra)

~~links easy~~  
fewer things to keep track of,  
@ any rate.

(iii) - ~~proof of (8)~~  
(2) suggests an algorithm

(iv) proof of (8) shows why the re-computation of centroids improves (2)

**k-means algorithm** Start with some initialization  $\hat{\mathbf{r}}_j^0$ ,  $j = 1, \dots, K$  (superscripts = iteration number, we start with  $i = 0$ ). Repeat until "convergence":

1. Assign each  $\mathbf{x}$  to its nearest *representative*  $\hat{\mathbf{r}}_j^i$ .
2. Set  $\hat{\mathbf{r}}_j^{i+1}$  to be the centroid of the  $\mathbf{x}$ s now assigned to it.
3. Increment  $i$

---

**single-link algorithm** We start with  $n$  clusters, each containing a single point. Here, think of clusters as connected components in an undirected graph. Repeat until there are only  $K$  clusters:

1. Add an edge between the two points  $x$  and  $x'$  in different clusters that have the minimum distance between them, thus merging their component clusters.