

last time, we started exploring the idea of clustering, a form of dim. reduction

(write down our formal defn, b/c might be hard to read on handout)

- looked @ functions for describing how good a clustering is, and explored relationships between them.
- but we didn't talk about how you get a good clustering (i.e., didn't describe any algorithms).

~~There's a good reason for this~~

Note that we're talking about a huge search space

<see "Stirling numbers of the second kind":

bounded below by $\frac{1}{2}(k^2+k+2)k^{n-k-1}$, according to wikipedia

~~flatgrass~~ Hatjivassiloglou; McKeown '93: $n=21, k=9, 1.23 \times 10^4$

... altho' that ~~doesn't~~ fact by itself doesn't imply that clustering is hard.

<sorting also has a huge search space>

nb: Srebro et al, "when is clustering hard?"

Today, two algs.

<follow handout 2nd page>

(1) seems to be $O(n^2)$ just to compute the criterion

- note that (2) requires being able to compute centroids - if your data isn't given as points in \mathbb{R}^d , might not be able to do this. Example: only given point-to-point distances.

visualization in D3; www.naftaliharris.com/blog/visualizing-k-means-clustering

should say about the alg: r_j^i is not necessarily r_j for C_j !! Maybe \hat{r}_j^i would have been better notation.

→ start w/ packed circles

explain colors = closest, \odot = centroid, colors = closest pts whose closest centroid is that color's.

(needed someone to 'put finger' on $\odot r_j^{i-1}$) reassign

(in class, accidentally picked an r_j^i that became "empty". oops!)

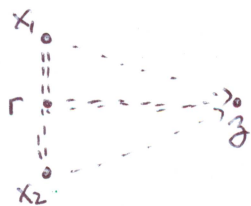
for smiling ~~total~~??

* specifically: can say k-means opt fn inspired by k, but it is actually:

min $\sum_j \sum_{x \in C_j} \|x - r_j\|_2^2$ so \hat{r}_j may not always be centroid;

- now: pt of (8). rewrite: $\left(\sum_{x \in C_j} \|x - r_j\|_2^2 \right) + n_j \|r_j - z\|_2^2 = \sum_{x \in C_j} \|x - z\|_2^2$ for arbitrary z .

last time: note that it seems like this is a violation of $\Delta \neq$: <click results>



that circle around the exponent was meaningful!

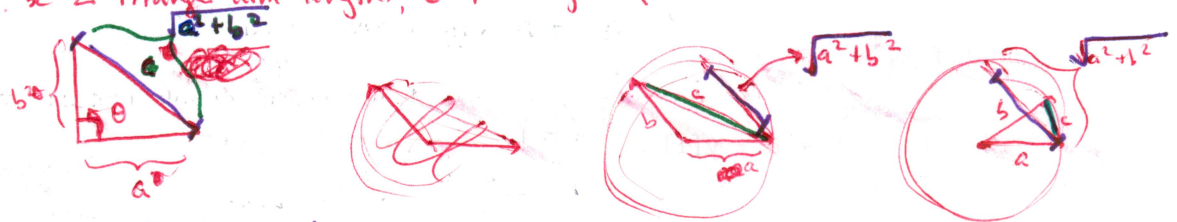
(9/2/21)

we resolve the paradox by noting that this is the Euclidean distance squared

$\| \cdot \|_2^2 \rightarrow$ square! not 1. $\Delta \nabla$ doesn't apply: $\langle \cdot \cdot \rangle$

again, sometimes the little things matter.

let a, b be 2 triangle arm-lengths, c the length of "connector".



here $a^2 + b^2 \neq c^2 = \|\cdot\|^2$

for right angle, $c = \sqrt{a^2 + b^2}$

$a^2 + b^2 < c^2$

$a^2 + b^2 > c^2$

so now that we @ least don't think ~~the lemma~~ that (8) is false a priori, why don't we go ahead & prove it?

$\sum_{x \in C_j} \|x - z\|_2^2$ need to introduce r_j . let's brute-force it.

$= \sum_{x \in C_j} \|x - r_j + r_j - z\|_2^2$

$\|a + b\|_2^2 = \sum_i a_i^2 + 2 \sum_i a_i b_i + \sum_i b_i^2$
 two vectors
 $= \|a\|_2^2 + 2 \sum_i a_i b_i + \|b\|_2^2$
 $= \|a\|_2^2 - 2a \cdot b + \|b\|_2^2$

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

$- 2 \sum_{x \in C_j} (x - r_j) \cdot (r_j - z) + \sum_{x \in C_j} \|r_j - z\|_2^2$

\downarrow
 $- 2(r_j - z) \cdot \sum_{x \in C_j} (x - r_j) + n_j \|r_j - z\|_2^2$
 second term in lemma

$- 2(r_j - z) \cdot \underbrace{\sum_{x \in C_j} (x - r_j)}_{=0}$

$\left(\sum_{x \in C_j} x \right) - n_j r_j$
 $= \frac{1}{n_j} \sum_{x \in C_j} x' - \frac{1}{n_j} \sum_{x \in C_j} x'$

making it a good choice for \hat{r}_j

- corollary: the z that minimizes ~~the LHS~~ either LHS/RHS is: the centroid

and, again, we assert w/out proof that this lemma yields that (1) \equiv (2). See Hopcroft; Kannan (link on webpage) for pf.
 So: considering RHS:

each $\|x - \tilde{r}_j\|$ is being reduced, where \tilde{r}_j 's are fixed.

<K-means> step 1 reduces distance to the "representative":

step 2 ~~minimizes w/in cluster~~ minimizes per cluster terms.

b/c we're choosing the best \tilde{r}_j 's for the ~~given~~ fixed C_j 's.

convergence given by ~~of possible~~ improvements in LHS having to be in discrete increments (if we assume no cycling btwn clusterings).

(this alternating minimization)

does it optimize? no (NP-hard problem).

partitioning problems are "pathy much always" NP-complete.

so that was K-means.

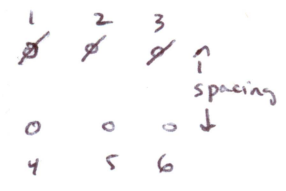
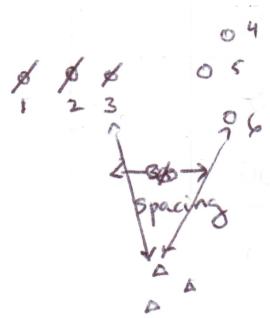
except when they aren't: lead-in to single-link clustering: an optimal alg. (for a certain clustering criteria)

note: "best friend" criteria ~~turns out to not~~ is not the single-link criterion.

~~probable~~ - probable hw exercise! Interesting to understand these subtleties.

instead: take ~~inter~~ inter-cluster spacing (front of handcat).

max the minimum of $d(C_j, C_{j'}) = \min_{x \in C_j, x' \in C_{j'}} \|x - x'\|_2^2$



note that ~~this~~ not all pairwise intra-cluster distances are small.

- idea of single-linkage: "knock out" the minimum ^{2-cluster} spacing
 \Rightarrow the max increases

but what about optimality?