

Whats 'Natural' k to use for this data?
 Sometimes we don't know which k to use.

Purity

- intra-Class (intra-cluster) similarity is high
- inter-class similarity is low

in other words, data inside clusters should be consistent, while data among different clusters should differ.

Intra Cluster distance

Cluster i , distance $d(C)$

$$\text{avg. pairwise intra-cluster distance } G_i = \frac{\sum_{x, y \in C_i} d(x, y)}{(2|C_i|)}$$

do for all clusters in Cluster Set.

weighted intra-cluster diff for level j is

$$W_j = \sum_{i=1}^j |C_i| G_i$$

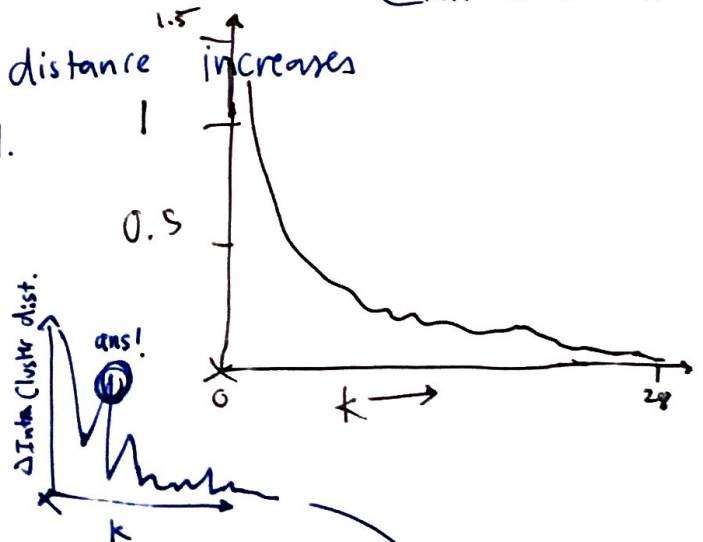
We want to minimize inter Cluster distance (& maximize intra Cluster similarity)

Graph-based Approach:

* when we merge, intra-cluster distance increases

- Choosing min. of this is not useful.
- well end up getting $k=N$,
- So, we look @ slope!

$$W_j' = \frac{(W_{j+1} - W_j)}{2}$$



1. Choose place where there's steepest absolute slope
 $k = \text{argmax}_j (|W_j'|)$

2. How about where is greatest Δ Slope? — 2nd derivative!
 (Knee of Graph)