

clustering possibilities  
: impossibilities.  
Intro to Gaussian  
Mixture Models

Outline:

big picture

implications of Kleinberg's impossibility theorem  
▶ what's "wrong" with ~~partitions~~ allowing partitions that refine each other?

▶ how should we think about the clustering problem, given the impossibility result?

ended up skipping

sidebar: proving a subpart of the impossibility theorem.  
reason: ~~the~~ kind of reasoning that can help in A2:

a souped-up version of

creating a dataset that forces a clustering algorithm to output a desired partition

see previous lecture's notes for proof. (or the actual paper)

• ((Gaussian) mixture) models: ~~can information abo~~

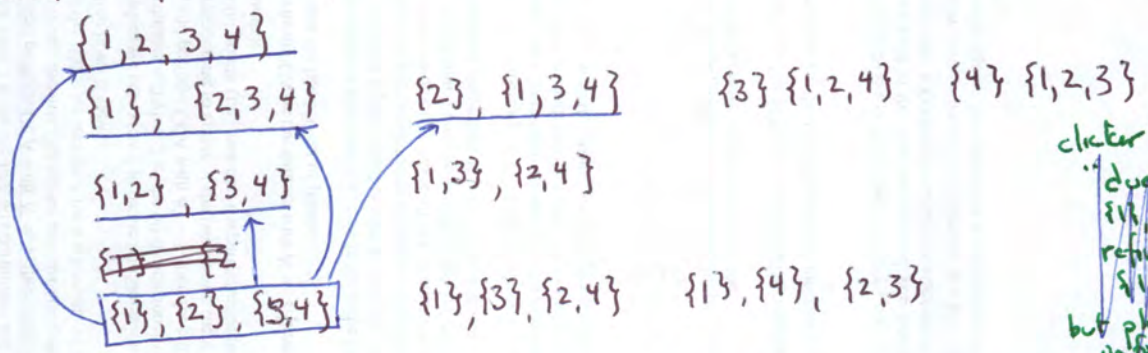
how does having prior information or assumptions change how we think about clustering?

pg 2 of handout: three properties of clustering functions we discussed last time.

not stated explicitly either last lecture or this one which was a mistake...  
placing here as a refresher/reference, altho ~~not~~ shown in lecture should be presented when covering the theorem's.

Ex:  $X = \{1, 2, 3, 4\}$

~~Set of all~~ <sup>some</sup> possible partitions of  $X$ :  
"clustering" = "partition"  
"cluster" = "a bin in a partition"



cluster of rc:  
"does {1}, {2}, {3}, {4} refine {1, 2, 3, 4}";  
but phrased informally, was instructive: 75/25 split.

Partition  $P_2$  refines partition  $P_1$  if ~~the~~ sets  $s_2 \in P_2$ ,  
 $\exists$  set  $s_1 \in P_1$ , s.t.  $s_1 \subseteq s_2$ .

In blue, we've ~~shown~~ let  $P_2 = \{1, 2, 3, 4\}$  and shown all the  $P_1$ 's that  $P_2$  refines.

"Refines" = ~~splits up some~~ "achieved by breaking up some clusters".

Recall: clustering functions  $f$  take as input a ~~data~~ distance function  $d$  and output a partition of  $X$ .

Review the three properties on the handout:

Richness: ~~there is no division of  $X$~~   
"there's no clustering/partitioning of  $X$  that can't be achieved by  $f$  by getting the right  $d$ ."

or,  
"for any partition you might want a priori (w/out looking at the data) there's some distance function  $d$  that causes  $f$  to output that target."

• so, an "expressiveness" property

Scale-invariance:

~~"the distance func"~~  
"output shouldn't depend on the units" (miles vs. mm) that the distances represent"

or,

~~"you do"~~  
"the output doesn't change if you scale all the distances by the same amount"

• so, a "stability" property

Consistency:

~~"if you change the  $d_i$~~

~~"suppose you have a  $d$  that gives you~~

"suppose  $f$  produces a partition  $P$  when given  $d$ .

If we alter  $d$  so that ~~the~~ <sup>each</sup> clusters of  $P$  gets even more "packed together"

and the ~~the~~ clusters of  $P$  grow even farther apart from each other, then  $f$  doesn't change its output"

• so, a "stability" property.



thm 3.1 + thm 3.2 = (scale-invariance + consistency)  $\leq$  "pseudo-richness"  
 $\neq$  richness.

pseudo-rich: possible output sets are the antichains for  $X$ .

sets of partitions ~~where~~ where no partition refines another.

implications:

clicker q: can we have a fn  $f$  whose range is an anti-chain  $A$  where  $|A| > 1$  and  $X \in \text{range}(f)$ ?

↑  
all points in one big cluster

class went ~~70/20~~ 80/20 on this.

No. Every other ~~part~~ partition of  $X$  refines  $X$ .

$\{1,2\}, \{3\}, \{4,5\}$  refines  $\{1,2,3,4,5\}$ .

[skip, but for those reading these notes:

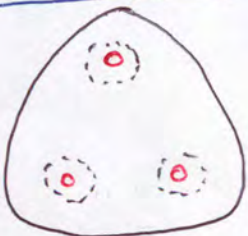
You also couldn't have  $|A| > 1$  and the partition  $P$  where every point is in its own cluster,

because  $P_{\text{singletons}}$  refines any other partition.

You can have  $A \subseteq$  set of all ~~single~~  $k$ -partitionings  
 (2 partitions of  $X$  that both consist of  $k$  clusters  
 can't possibly refine each other).  
 ↑  
different

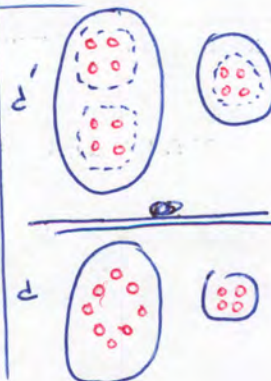
... and hey, that ~~helps~~ coincides with the fixed- $k$ -means; fixed- $k$  single-link algs that are so famous!

\* Why is refinement such a problem for scale-invariance and consistency?  
 Some examples for intuition.



$\square$  = a partition  
 $\square$  (with dashed border) = another partition, refining  $\square$ .

whether all-in-one or each-in-own seems to be a judgment about whether the distances are "big" or "little" in an absolute sense.



given  $\square$  = a partition  
 $\square$  (with dashed border) = another partition, refining  $\square$ ,  
 whether to ~~pick~~ pick  $\square$  or  $\square$  seems to depend on whether the <sup>top</sup> left-hand ~~subcluster~~ distances cause "subclusters" to arise.

So, should we just give up on clustering?  
Don't think so, it's just too useful a task - it tells you about the structure of your data.

It's just: you either have to:

- give up on some of the properties
- formulate your own properties
- break the formulation

~~xxxxxx~~ ... ~~including~~, ~~also~~  
... [perhaps what the theorem says we need is more information than just a distance fn.]  
(an example: ~~what k~~ predetermine k.)  
↑ important  
abbreviation of "function"

After all, when we are clustering, is our goal really:  
"just put points that are close together in the same group" ?

or is it  
xx determine the underlying structure of the data" xx  
↑ seems like a more fulfilling thing to do.

perhaps it has become time to consider a big question:

Where does ~~my~~ the data come from?





~~lecture~~

When does our data come from? (the birds & bees talk)

What kind of generative story or account can we make up that's reasonable for how clusters underly the data?

ex generative story: ~~percentage~~ <sup>parentage</sup> ("the apple doesn't fall far from the tree")

- pick ~~Mother Nature selects~~

We have tree species: oaks, maples, and apple trees. ~~Model of population~~  
⊗ Explain spatial location of the trees.

Mother Nature picks a species

- picks a <sup>pre-existing</sup> tree of that species to be a parent ←

- plants a seed around the location of that parent,   
 this intermediate variable is what distinguishes this from GMMs.

~~prob~~  
prob of ~~location~~ seed's location falls off w/ distance from parent.

ex: generative story: Gaussian sources < Gaussian mixture model >

~~Assume two types~~

We have species, oaks, maples, apple trees.

~~Have features regarding~~

Explain distribution of height, width, leaf lengths, leaf widths.   
 ↪ 4 #s

M.N. picks a species; grows a tree of that species, chooses the height, width values as from ~~approx~~ multivariate Gaussian for that species - i.e., each species has a template (with variance)

think as both mixture models:  
clustering { given observed data, w/out identifying labels (so, you don't get to see the actual trees, just the lat-longs or the height/weight data) (don't see the species!!)  
"recover the underlying parameters of the "source" model," go  
- in this case, to figure out which trees belong to the same species.  
?

But there's a problem with thinking of our quest as to recover "the" source model, and that is: ~~that~~ it's an impossible problem.

Ex: top ~~figure~~ figure of pg 4:

assume we're dealing with one-dimensional data.

black dots = data, all on one line.

and we tell you that the data was generated by two Gaussians, so one cluster should be the points generated by ~~the leftmost Gaussian~~ one of the Gaussians, and the other cluster should be the points generated by the other Gaussian.

So, one possibility is that we have ~~one Gaussian~~

{ one Gaussian centered under the left bump of the green curve,  
one Gaussian centered " " right " " " " }

green = resulting combined density fn.

this would be a plausible way the black dots got generated; it does look like there's a rightmost clump of points and a leftmost clump.

But, it really could be that the true source had one Gaussian a mile to the left, and one Gaussian a mile to the right, <so, both totally off the screen>.

We cannot say that that's not the source model.

it could be that that was the source model, and we just got a really unlikely draw of the data.

So, we ~~can~~ can't say that we can select the true source.

However, if we have to pick some choice of how to decide what model to settle on in order to create output,

then one idea is to use the ~~maximum-~~ maximum-likelihood principle (is on handout)

Assume fixed <sup>generating</sup> model class, where elements are indexed by parameter setting  $\theta$ . (a vector)  
Each possible source  $\uparrow p_\theta$

Ex: Gaussian mixture model w/ 2 Gaussians w/ fixed variances, one chosen w/ fixed prob  $2/3$ , one w/ fixed prob  $1/3$   $\uparrow$  1-d

$\Rightarrow$  possible  $\theta$ 's are the: the means of the two Gaussians.

ex: green ~~curve~~ curve:  $\theta = (-2, 2)$   $\leftarrow$  right mean  
 $\leftarrow$  left mean

resulting densities shown on handout



another possible choice would be our "miles-away" alternative

$$\theta = (-10,000, 10,000)$$

→ mean of 1<sup>st</sup> Gaussian → mean of second Gaussian.

~~ML principle~~  
So for given dataset  
pick

do write the English down first.

for given data  $X_{\text{given}}$ , pick model ~~with the  $\theta$  that maximizes  $P(X_{\text{given}} | \theta)$~~

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} P(X_{\text{given}} | \theta)$$

function of  $\theta$ .

thus, since the green-curve model thinks  $X_{\text{given}}$  (on handout, as black dots) is pretty likely, whereas the  $X_{\text{given}}$  is very unlikely according to the "miles-away" model, MLP says ~~to~~ prefer the green-curve model over the miles-away model.

How to find the  $\hat{\theta}_{\text{ML}}$ ? "Derivatize": set ~~to 0~~ param.

let  ~~$\theta = (a, b, c, \dots)$  for notational sake~~  
let:  $\theta = (\theta[1], \theta[2], \dots, \theta[l], \dots)$

"Derivatize":  $\frac{\partial}{\partial \theta[l]} P(X_{\text{given}} | \theta)$  and set to 0, solve for  $\theta[l]$ .

but, a trick:  
- almost always a good idea to use the log of the likelihood instead (and since log is monotone increasing in its argument, makes no diff. in finding the argmax).

Handout: the log-likelihood function for our two-Gaussian model,  $\frac{2}{3}$  vs.  $\frac{1}{3}$  prob:

$\theta = (\mu_1, \mu_2)$  →  $\mu_1$  is shown as x axis,  $\mu_2$  shown as y axis.

$$p(x | \theta) = p(x | \mu_1, \mu_2) = \frac{1}{3} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_1)^2} + \frac{2}{3} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_2)^2}$$

prob of choosing 1<sup>st</sup> Gaussian

prob of  $x$  if generated by 1<sup>st</sup> Gaussian.

$\log P(x | \mu_1, \mu_2)$  is the surface plot.

in orange: the application of an iterative method for trying to do gradient ascent to find a local max of the log-likelihood.

→ then are two trajectories, corresponding to two diff. starting points.

$\mu_a$  ;  $\mu_b$  are the two places this search converged.

↳ on the contour plot below the surface. (in an notation,  $\mu_a$  would be written  $\theta_a$ , and  $\mu_b$  as  $\theta_b$ )

on the figure above, the ~~distributions for~~ densities for  $\theta_a$  ;  $\theta_b$  are shown:

~~$\mu_a$  has~~

$\theta_a$  has  ~~$\mu_1 \approx 2$ ,  $\mu_2 \approx 2$~~   $\mu_1 \approx -2$ ,  $\mu_2 \approx 2$

$\theta_b$  has  $\mu_1 \approx -1$ ,  $\mu_2 \approx 2$ .

Note: the iterative method mentioned above and whose trajectories are depicted in the handout does not make use of latent variables. It instead uses the hideous equation whose existence is mentioned in the next lecture.