

Announcements:

- Yavenu Paris' Friday 3/28 hours moved to ~~the~~ today 8-9pm, location T34.
- All spring break office hours are cancelled, but:
 - Mevlana Gemici will hold his Sunday hours if you request by email with 24 hours notice
 - Jack Hessel will hold his Monday hours if you request by email with 24 hours notice

Lec. 17, 3/27/1
• More generative stories
(and hence more EM) need for not?
• Easing into graphical models.

Coursework scheduling updates:

- Given time constraints, we are almost surely only going to assign one more assignment. As for the "competitions"...
- The first "competition" dataset is still being tested and documented. We hope to release the dataset before Spring Break, so students can take a look if they want...
 - ... but not release instructions, so that there's no pressure to do anything over Spring Break.
- The due date for the second competition is, as determined by the registrar, May 11th, 4:30 pm.

A2 pedagogical note (see also Prof Sridharan's email):

What we want you to learn from this assignment is:

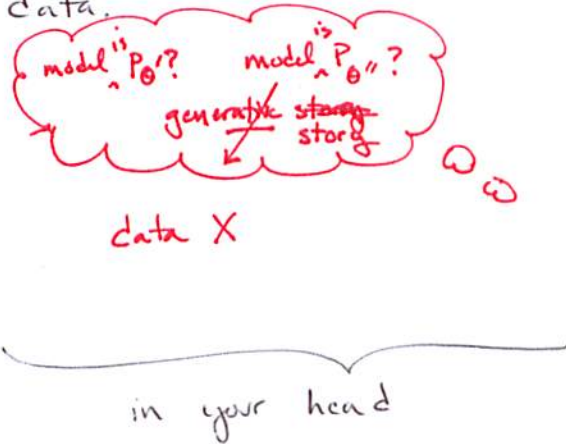
NOT that bizarre, degenerate, unrealistic datasets reveal instability for clustering algorithms,

BUT ~~that~~ ~~what~~ ~~kinds~~ ~~of~~ ~~dataset~~
skip

So far: we've moved from ~~just~~ considering our data as "just given" (projection, clustering)

to thinking (imagining) that the data is generated by some (unknown member of) probabilistic process a family of

and regarding our task as ~~picking~~ selecting a model based on the data.



The model P_θ* you pick is your "dimension-reduced" understanding of the data.

"These 10K datapoints are really just examples of 3 species of trees, where the oak trees look like are usually on the east coast, the orange trees tend to concentrate in Florida".

how to pick θ' vs. θ"?
maximum likelihood principle
maximum a posteriori principle

} either requires maxing
~~which could be hard~~
often hard to compute directly

EM often comes to the "rescue"

~~quotes b/c it often is~~
~~only going to get you~~
quotes b/c it may only get you to a stationary point.

~~* very general and powerful~~
* very general technique:

We're teaching it to you b/c it's important.

when your model's hidden structure
~~makes it easy~~
~~making~~ makes ~~taking~~ derivatives
solving "derivative=0" easy
if only it weren't hidden

~~Let's demonstrate that EM~~

let's further exemplify EM by applying it to another set of models.

Gaussian mixture models: apply to "numerical, spatially concentrated" data.

But there are lots of kinds of data that they don't apply to.

What about category-

Contrasting example: purchase data for market segmentation.
Mixture of multinomials

customer	# bags of chips	# boxes of	...
	product 1	product 2	...
customer:	3	0	6
	bottles of Coca Cola	bottles of Coke	bottles of Sprite

(don't normalize now - will compare w/ the ϕ_i 's)
(typically let's assume normalized)

this kind of 'count' data just doesn't look like Gaussian-generated.
You could go ahead & try to fit a GMM to it, but why not pick a more appropriate generative story?

(can't buy 6.2 bottles)

notation for multinomials: choice over D options:

ϕ_j : the set of D parameters for the j th multinomial.

$\phi_j[l]$: prob of picking the product l .
ex: ~~pepsi fans~~. anything-but-pepsi

ex for ~~pepsi fans~~: pepsi enemies:

$\phi_j[1] = 0.05$ \rightarrow desert island scenario
 $\phi_j[2] = 0.6$ \rightarrow equally likely to pick competing product.
 $\phi_j[3] = 0.35$

✓ generative story:

[there are K types of customers.
assume each customer makes exactly m purchases.]

probably not an essential assumption

pick the customer type j according to Π_j distribution over types.

(I don't want to deal w/ proper norm normalization or the whole multinomial coeff stuff.)

pick the m purchases by the distribution given by ϕ_j .

Task: given X , pick a Π_j ϕ_1, \dots, ϕ_K
 \rightarrow each is $\phi_i[l], \dots, \phi_i[D], m$

hope you might discover this like 'oh, 20% of my customer base is people who buy anything but pepsi'.

ML approach: find params maximizing ^(log) likelihood of data:

Recall: for a given φ_i and a given $x_t = x_t[1], x_t[2], \dots, x_t[d]$

$$P(x_t | \varphi_i) = \frac{m!}{x_t[1]! x_t[2]! \dots x_t[d]!} \cdot \varphi_i[1]^{x_t[1]} \cdot \varphi_i[2]^{x_t[2]} \dots \varphi_i[d]^{x_t[d]}$$

MLE ~~se~~ ~~PKXX~~

$$\text{loglik}(X) = \log \left(\prod_t P(x_t | \theta) \right)$$

↳ our π_i φ_i 's.

$$= \sum_t \log P(x_t | \theta)$$

= <generative story> use π box to not confuse w/ product.

$$\sum_t \log \left(\sum_j \pi_j \frac{m!}{x_t[1]! \dots x_t[d]!} \prod_{l=1}^d \varphi_j[l]^{x_t[l]} \right)$$

... ready

could be any customer type, so sum over them all.

... ready to take some derivatives & set to 0?

clicker question: how many derivatives ^{must we} derive ~~how~~ to compute?

(6) 29% ~~3%~~ (A) \bullet K (one for each φ_i) + K (for $\prod_j \pi_j$)

(10) 49% (B) \bullet $K \times d$ (one for each $\varphi_i[l]$) + K (one for each π_j)

(2) 10% ~~(C) almost A, but off concept~~
 (C) my answer is almost but not quite (A)

(3) 14% (D) my answer is almost but not quite (B).

↳ it's this one, but I only need $K \times (d-1) + (d-1)$

since $\sum_{l=1}^d \varphi_i[l] = 1$.

(I pick something w/ prob 1)

enforce w/ Lagrange multiplier λ_i or just directly normalize by sum (anti-Pepsi's chance of Spnk).

now, if we try to take the derivative

clicker q:

$$\sum_t \frac{\partial}{\partial \varphi_i[l]} \left(\sum_j \pi_j \frac{m!}{x_t[1]! \dots x_t[d]!} \prod_{l=1}^d \varphi_j[l]^{x_t[l]} \right)$$

OMG $\varphi_{i,l}[l']$ are everywhere, ~~omit~~

↳ all constant wrt i, l'

fracs: $\frac{d}{dy} \log(z y^n) = \frac{(n+1) z^{n-1} y^n}{z y^n} = \frac{n+1}{y}$

if $l \neq l'$ and $i \neq i'$, then $\varphi_{i,l}[l']$ looks like a constant to $\frac{\partial}{\partial \varphi_{i,l}[l]}$

clicker question: ~~do we need EM?~~ should we use EM, instead of directly solving for the max. likelihood soln?

(a) 43% (A) Yes, because we can ^{without} get the equation: take the derivatives but can't solve for $\varphi_j, [L']$ independently of the other params. (i.e., just in terms of the data)

(b) 14% (B) Yes, because we can't take derivatives if we need

(c) 29% (C) No, because we can take the derivatives and solve for $\varphi_j, [L']$

14% (D) No, because there's no hidden structure for EM to take advantage of.

so, what's the problem?

• sum over i in denom:
 (unknown)
 { keeping all the φ_i 's in the equation
 wouldn't be "necessary" if we knew that
 that x_t "actually" is of type "j"
 (no longer a mixture prob)

the c_t = which type x_t is.
 (e.g. "anti-pepsi").

EM idea: ~~initialize initial parameter~~

• we don't know the c_t 's,
 ... so just average (take expectation) over all their possible values
 • but if we ~~can estimate~~ know their distribution,
 we can take the expectation over possible c_t values:

⇒ instead of find $\frac{\partial \log L}{\partial \varphi_j, [L']}$ maximizing $\sum_t \sum_{k=1}^m p(c_t=k) \log P(x_t, c_t=k | \varphi_j)$

for $k=j$, $\log \frac{m!}{\dots!} \left(\prod_i \varphi_i^{x_t^{(j)}} \prod_j \varphi_j^{c_t^{(j)}} \right) \cdot \varphi_j^{x_t^{(j)}}$
 ↓
 constant wrt $\varphi_j, [L']$

this was our $Q_t(c_t)$ from last time lecture.

↳ Q = "like a P".

~~q~~ = $P(c_t=k | \text{some parameters}, x_t)$

where do we get $P(c_t=k | \text{old params}, x_t)$?
 Our previous guess of $\pi[k]!$

↑
 an initial guess!

EM: just keep imagining what would be useful, and guess when you run out of other possibilities!

$$\frac{\partial}{\partial \varphi_j[l']} \left[\log Q_t(c_t) + \log \left[\text{big constant} \cdot \varphi_j[l']^{x_t[l']} \right] \right]$$

$$= \frac{x_t[l']}{\varphi_j[l']}$$

set to 0

$$\Rightarrow \varphi_j[l'] \cdot \infty \cdot x_t[l'] \cdot Q_t(c_t)$$

how many times do we see
part being deriv.