

Pedagogical goal: more practice w/ developing generative stories

Outline:

- ① review mixture of multinomials ~~notation~~ (model, notation)
- ② new setting: "longitudinal" data [analogy w/ documents], where same person can have different types?
  - ②a choice of multinomials: from fixed set  $\rightarrow$  distribution over multinomials: the Dirichlet
  - the LDA story

LDA  
(Latent Dirichlet Allocation) -  
An important generative story

(For handout:)

Last time: data  $X$  has  $n$  rows like  $x_t = \begin{matrix} \text{pepsi} & \text{coke} & \text{sprite} \\ 3 & 0 & 6 \end{matrix} \dots$   
 $x_t = \begin{matrix} x_t[1] & x_t[2] & x_t[3] \end{matrix}$

Generative story: mixture of multinomials

For a given customer,  $x_t$  (corresponding to the  $t$ th row) =  $(x_t[1], x_t[2], \dots, x_t[d])^T$   
(transpose of the)

~~they~~ picks among given customer types w/ prob  $\pi[1], \dots, \pi[k]$   
mother nature  $\rightarrow$  # of customer types

A customer type  $j$  is represented by  $d$  values of  $\varphi_j$  ~~the parameters of multinomial~~, a multinomial.

- $\varphi_j[1]$  (prob of picking a pepsi)
- $\varphi_j[2]$  (prob of picking a coke)
- $\vdots$
- $\varphi_j[d]$

ex: an anything-but-coke type might have  $\varphi[1]=.6, \varphi[2]=.05, \varphi[3]=.35$ .

$$\sum_{l=1}^d \varphi_j[l] = 1, \text{ all } \varphi_j[l] \geq 0.$$

We let  $c_t$  = the type that customer  $t$  picked.

$x_t$  then picks their  $m$  purchases according to the  $\varphi_j$  they picked.

The prob of  $x_t$  ~~according to the  $\varphi_j$  they picked~~  $\left( \frac{\pi[c_t]}{x_t[1]! x_t[2]! \dots x_t[d]!} \prod_{l=1}^d \varphi_{c_t}[l]^{x_t[l]} \right)$

make this a clickerg:  
 swap  $\varphi$  and  $\pi$   
 swap  $\varphi$  and  $x$   
 swap  $\varphi$  and  $\pi$   
 don't swap anything

Task: given  $X$ , recover the best:  $\pi$  ( $k$  unknowns,  $k-1$  degrees of freedom)

$\varphi_j$  ( $1 \leq j \leq k$ ) (each has  $d$  unknowns,  $d-1$  d.o.f.)

ex: "20% of my population are sprite fans, ~~of~~ type (0.05, 0.05, .9)"

~~that's not a~~

our mixture-of-multinomials story is a pretty good one. <sup>① seems reasonable</sup>

<sup>②</sup> we can use EM to ~~recover parameters~~

choose "good" values of our unknown parameters.

(so, not that computationally

"complicated")

hidden vars:  $\{T, \theta_j\}$

and the  $c_j$ 's:

which type customer  $t$  picked ~~hidden var~~  $\theta_j, c_j$  customer  $t$ 's purchases over an extended time

new scenario:  $x_t$  represents ~~purchases over~~ period ... (so as to ~~get~~ <sup>cor tastes</sup> independence)

- ~~that's~~ <sup>car</sup> needs can change.

$\Rightarrow$  one  $\varphi_j$  isn't a good fit <sub>just</sub> <sup>convenient</sup>

{ one shopping trip, they need soft drinks  
but on a different trip, they want cereal,  
and on a third trip, they're choosing fruit.

we probably don't want to have one multinomial for:

people who hate-coke AND love Froot Loops AND love either bananas or oranges.

- instead of one  $\varphi_j$  for = coke-haters who love Froot Loops and either bananas or oranges!

have a  $\varphi$  for: when picking soft drinks, anything but coke <sub>(but occasionally)</sub>

<sup>plus</sup> ~~plus~~ a diff  $\varphi$  for: when picking cereal, favorite is Froot Loops, etc.

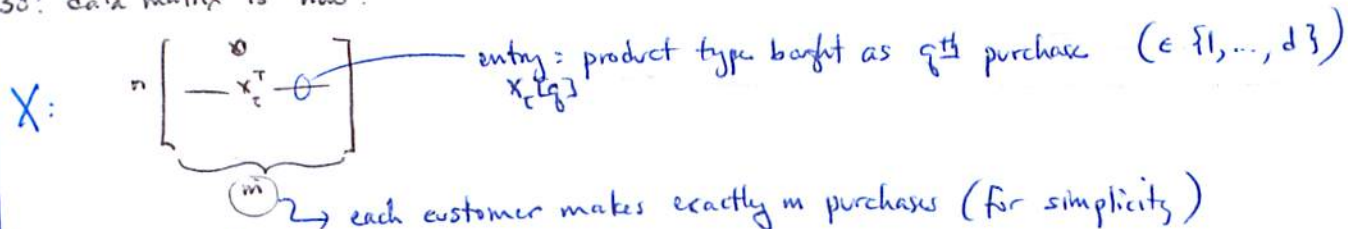
$\Rightarrow$  re-select a (possibly different)  $\varphi$  for each purchase customer  $t$  makes.

• have an index on sequence in which purchases were made

multinomials idea?   
 w/ 1. + for, same # of ~~params~~

in the likelihood   
 ~~what the  $\varphi_j$  is, is hidden variable, but it's ~~not~~ <sup>latent</sup>~~

so: data matrix is now:



task: recover values for:

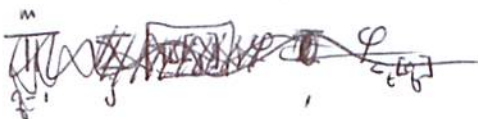
- hidden variables: the  $K$  different preference types  $\varphi_j$  ( $d$ -dimensional)
- (-almost the same story) the prob  $\pi$  of picking preference types ( $K$ -dimensional)
- the choices  $c_c [q]$  that customer  $c$  was using to pick their  $q^{\text{th}}$  purchase. ( $m \rightarrow c_c$  is  $m$ -dimensional) entries run from 1 to  $K$ .

definitely want to learn: what kind of preferences are there, and how likely is each of these preferences.

might want to know what 'frame of mind' (need) the user was in when making those purchases.

clicker of this:

the likelihood for vector  $x_c$  now looks something like this:  $\rightarrow$  given the hidden variables



$$\prod_{q=1}^m \boxed{\pi_j [c_c [q]]} \varphi_j [x_c [q]]$$

<ignoring normalization>

altho' if we wanted to marginalize over all the possible  $c_c [q]$  (assignments), we'd get:

$$\prod_{q=1}^m \sum_j \boxed{\pi_j [c_c [q]]} \varphi_j [x_c [q]]$$

- hidden structure makes likelihood simpler.

• we have <sup>many</sup> hidden variables  $c_t, \varphi_j$  = what is st. the  $g^{\text{th}}$  purchase by customer  $t$  was made according to ~~customer~~ preference type  $\varphi_j$ ?

• how should we model this choice?

W'd like to say, this person is mostly a pepsi-buyer, but sometimes buys bananas or oranges.

W'd like to say, this person is mostly a pepsi-buyer, Prof Snider but sometimes is a bananas-or-oranges buyer.

That is, their prob of ~~picking the pepsi~~ being a pepsi-buyer is higher than their prob of being a bananas-or-oranges buyer.

But this <sup>other</sup> person, say, prof Lee, never buys soft drinks @ all, and is only a bananas-or-oranges buyer.

Further

Another design choice: each customer has their own probabilities of choosing among the  $\varphi_j$  for each purchase.

ex: One ~~person~~ <sup>shops for soft drinks</sup> much more often than ~~another~~ <sup>shops for</sup> ~~fruit~~ <sup>soft drinks</sup> @ each ~~purchase~~ <sup>purchase</sup>.  
 another ~~picks~~ <sup>picks</sup> ~~bananas~~ <sup>bananas</sup> and ~~oranges~~ <sup>oranges</sup> with prob 50%, ~~cereal~~ <sup>cereal</sup> w/ prob 50% ~~at each purchase~~ <sup>at each purchase</sup>.

⇒ instead of a fixed  $\pi_j$ , we have a  $\pi_t$  for each customer  $t$ 's probability,

$\pi_t = [\pi_{t,j}]$ ,  $j \in \{1, \dots, k\}$  ; prob of picking preference-type  $\varphi_j$ .  
 ↑  
 customer  $t$

Think of  $\pi_t$  as user's profile over pref types  $\varphi_j$ . But we've ~~increased~~ <sup>now</sup> ~~mean~~ <sup>mean</sup> our # of hidden variables significantly.

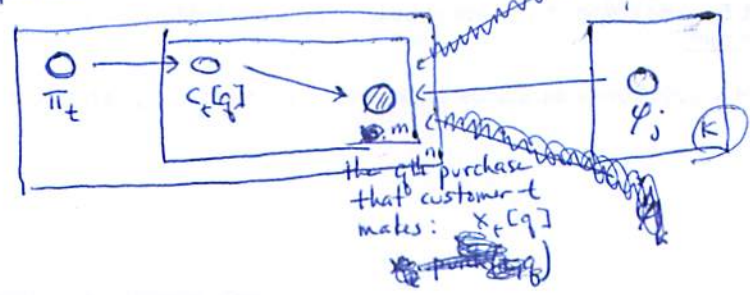
task so far: (just say, don't need on the board).

Given  $X$ , recover: the various preference types  $\varphi_j$  (d-dimensional vectors)  
 the various user profiles  $\pi_t$  (k-dimensional n k-dim vectors)  
 the preferences responsible for each purchase  $c_t, \varphi_j$  (n m-dim vectors)  
 each in  $\{1, \dots, k\}$  (we've fixed the # of purchases each user makes.)

skip

~~This is a lot of unknowns.~~

A (graphical-model-like) summary:



shading = observed.

# in corner = # of versions

Should we skip this? It kind of interrupts the flow

SKIP

So, if we've significantly expanded the search space of hidden variables.  
(for instance, we now have roughly  $n \times k$   $\pi$  values to fill in)  
we might expect that even as powerful a tool as EM might get in trouble  
(maybe too many stationary points or local maxima),

⊙ that is: greater model expressivity, @ cost of making search for hidden values harder.

- this kind of tradeoff is an important factor in developing your own generative stories.

What have we done in other such situations, when things seem complex or even impossible?

- in EM, we guessed what the right distribution would be
- when clustering proved impossible, we introduced constraints.

→ assume some prior info is given

- so, if we could a priori restrict the space of  $\pi_t$ 's to search over, or know that certain  $\pi_t$ 's are more likely, this would make things easier for us.

what prior can we put over multinomials  $\pi_t$ ?

<at least one student flipped the handout over to look >  
The probabilists & statisticians know a conjugate prior  
<you don't have to know what that is> for multinomials!

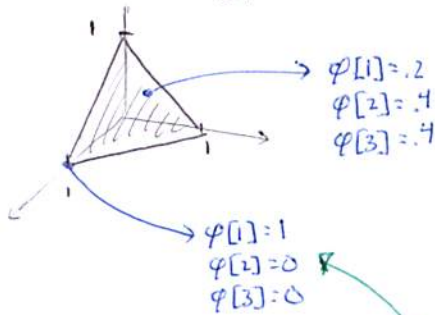
~~the~~

the Dirichlet prior - we're given you some ~~intuition~~  
geometric intuitions on the handout.

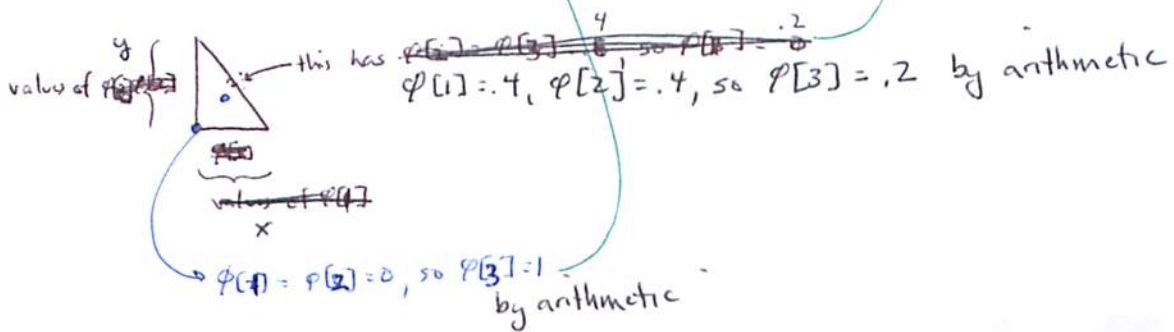
## About the Dirichlet distribution

Each  $\varphi$  can be thought of as a point on the probability simplex in  $\mathbb{R}^d$ . (i.e., vectors  $\varphi$  s.t.  $\sum_{k=1}^d \varphi[k] = 1, \varphi[k] \geq 0$ .)

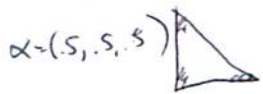
Ex:  $d=3$



Given that  $\varphi[k] = 1 - (\varphi[l] + \varphi[m])$ , we can "project" this simplex into 2-d by looking at it from the  $z$ -axis: (put your eye right there), like you're sighting down a pool cue!



<show 2nd slide 1st>:



heat map. bottom-right corner = multinomials that put most of their probability mass on the first product.  
The "redness" though = this Dirichlet prior puts high prob on those kinds of multinomials.

- the Dirichlet  $(1, 1, 1)$  is a uniform prior on multinomials.

- the Dirichlet  $(.5, .5, .5)$  isn't "symmetric" like the other two, and you see this in the fact that "the red splotch" isn't "symmetric".

(q: why are  $(.5, .5, .5)$  and  $(1, 1, 1)$  different?)

<then do 1st slide; go over it>