

What if there is noise?

### Abstract Topic Model

① unknown topic matrix  $A$ , distribution  $M$  on simplex in  $\mathbb{R}^r$

② For each document  $i$

    a) Sample  $W_i$  from  $M$

    b) Generate  $L$  words i.i.d. from  $AW_i$

Can we recover  $A$ ? Notice that

$$\mathbb{E}[\tilde{M} | W] = M = AW$$

↑  
empirical term-by-document

but  $\tilde{M}$  and  $M$  are far apart. Why? sparsity

This is a rich model, containing  
supp on corners dirichlet lognormal hierarchical  
pure, LDA, CTM, Pachinko

as special cases

Thm [Anura, Ge, Moitra] there is a polynomial time algorithm for learning  $A$  when it is separable, under nondegeneracy conditions



Furthermore we can expand

$$P[w_1=j | w_2=j'] =$$

$$\sum_{i'} P[w_1=j | w_2=j', t_2=i'] P[t_2=i' | w_2=j']$$

By conditional independence

$$(*) = \sum_{i'} P[w_1=j | t_2=i'] P[t_2=i' | w_2=j']$$

Claim:  $P[w_1=j | t_2=i'] = P[w_1=j | w_2=\pi(i')]$

This implies

$$(*) = \sum_{i'} P[w_1=j | w_2=\pi(i')] \underbrace{P[t_2=i' | w_2=j']}_{\text{unknowns}}$$

(1)

This linear system has unique soln if  $R$  is full rank

Note: It uses all words, not just anchor words

Finally, by Bayes rule

$$(2) P[w_2=j' | t_2=i'] = \frac{P[t_2=i' | w_2=j'] P[w_2=j']}{P[t_2=i']}$$

and moreover

$$IP[t_2 = i'] = \sum_j IP[t_2 = i' | w_2 = j] IP[w_2 = j]$$

Thus our algorithm is

### Anchor-Bayes

Compute Gram matrix  $G$

Compute the anchor words (via separable NMF)

Solve (1) for  $IP[t_2 = i' | w_2 = j']$

Solve (2) to compute  $IP[w_2 = j' | t_2 = i']$

$$= A_{j' i'}$$

Proof of Claim let's expand

$$IP[w_1 = j | w_2 = \Pi(i')] =$$

$$\sum_{i''} IP[w_1 = j | w_2 = \Pi(i''), t_2 = i'']$$

$$IP[t_2 = i'' | w_2 = \Pi(i'')]$$

$$= \begin{cases} 1 & \text{if } i'' = 0 \\ 0 & \text{else} \end{cases}$$

does this use  
uniqueness of  
anchor words?  
no

$$= P[w_1 = j \mid w_2 = \pi(i'), t_2 = i']$$

by conditional independence

$$= P[w_1 = j \mid t_2 = i'] \quad \square$$

Are natural topic models separable?

UCI Dataset  $\rightsquigarrow$  MALLET  $\rightsquigarrow$   $\hat{A}$   
300k NYT articles

Findings: with  $r=200$ , about 0.9 fraction  
of topics had a near anchor word  
(posterior  $\geq 0.9$ )

How do algorithms based on separability  
perform?

$\hat{A}$   $\rightsquigarrow$  synthetic documents  $\rightsquigarrow$  MALLET  
 $\rightsquigarrow$  Anchor Bayes

Findings: More accurate, and a hundred  
times faster

what do we find on real data?

show figures