# Foundations of Bayesian NLP MSc Artificial Intelligence

Lecturer: Wilker Aziz Institute for Logic, Language, and Computation

2019

# The problem with $\ensuremath{\mathsf{MLE}}$

Motivating example from Liang and Klein (2007)

mixture of Gaussians trained via EM

Motivating example from Liang and Klein (2007)

mixture of Gaussians trained via EM



Example from Liang and Klein (2007): ACL tutorial on Structured Bayesian Nonparametric Models slides

Motivating example from Liang and Klein (2007)

mixture of Gaussians trained via EM



Example from Liang and Klein (2007): ACL tutorial on Structured Bayesian Nonparametric Models slides

Motivating example from Liang and Klein (2007)

mixture of Gaussians trained via EM



Example from Liang and Klein (2007): ACL tutorial on Structured Bayesian Nonparametric Models slides

Motivating example from Liang and Klein (2007)

mixture of Gaussians trained via EM



 as the capacity of the model increases (more clusters), training likelihood strictly improves

Example from Liang and Klein (2007): ACL tutorial on Structured Bayesian Nonparametric Models slides

Motivating example from Liang and Klein (2007)

mixture of Gaussians trained via EM



- as the capacity of the model increases (more clusters), training likelihood strictly improves
- but what happens with test likelihood?

Example from Liang and Klein (2007): ACL tutorial on Structured Bayesian Nonparametric Models slides

Motivating example from Liang and Klein (2007)

mixture of Gaussians trained via EM



- as the capacity of the model increases (more clusters), training likelihood strictly improves
- but what happens with test likelihood?

Example from Liang and Klein (2007): ACL tutorial on Structured Bayesian Nonparametric Models slides

# The problem with $\ensuremath{\mathsf{MLE}}$

That's why you were told to always do model selection

- on heldout set
- preferably via cross-validation

That's why you were told to always do model selection

- on heldout set
- preferably via cross-validation

Can you see limitations of this approach?

That's why you were told to always do model selection

- on heldout set
- preferably via cross-validation

#### Can you see limitations of this approach?

- availability of data
- representativeness of heldout set
- discrete optimisation: combinatorial search over models

# NLP1

#### Preliminaries

Bayesian modelling

Applications

#### • N observations $\mathbf{x} = \langle x_1, \dots, x_N \rangle$

- N observations  $\mathbf{x} = \langle x_1, \dots, x_N \rangle$
- *i*th observation  $x_i \in \{1, \ldots, K\}$

- N observations  $\mathbf{x} = \langle x_1, \dots, x_N \rangle$
- *i*th observation  $x_i \in \{1, \ldots, K\}$
- ▶ all but the *i*th observation  $\mathbf{x}_{-i}$

- N observations  $\mathbf{x} = \langle x_1, \dots, x_N \rangle$
- *i*th observation  $x_i \in \{1, \ldots, K\}$
- ▶ all but the *i*th observation  $\mathbf{x}_{-i}$
- $\blacktriangleright$  N cluster indicators
  - $\mathbf{z} = \langle z_1, \ldots, z_N \rangle$

- N observations  $\mathbf{x} = \langle x_1, \dots, x_N \rangle$
- *i*th observation  $x_i \in \{1, \ldots, K\}$
- ▶ all but the *i*th observation  $\mathbf{x}_{-i}$
- N cluster indicators
  - $\mathbf{z} = \langle z_1, \ldots, z_N \rangle$
- *i*th cluster indicator  $z_i \in \{1, \ldots, C\}$

- N observations  $\mathbf{x} = \langle x_1, \dots, x_N \rangle$
- *i*th observation  $x_i \in \{1, \ldots, K\}$
- ▶ all but the *i*th observation  $\mathbf{x}_{-i}$
- N cluster indicators  $\mathbf{z} = \langle z_1, \dots, z_N \rangle$
- *i*th cluster indicator  $z_i \in \{1, \ldots, C\}$
- ▶ all but the *i*th cluster assignment  $\mathbf{z}_{-i}$

- N observations  $\mathbf{x} = \langle x_1, \dots, x_N \rangle$
- *i*th observation  $x_i \in \{1, \ldots, K\}$
- ▶ all but the *i*th observation  $\mathbf{x}_{-i}$
- N cluster indicators  $\mathbf{z} = \langle z_1, \dots, z_N \rangle$
- *i*th cluster indicator  $z_i \in \{1, \ldots, C\}$
- ▶ all but the *i*th cluster assignment  $\mathbf{z}_{-i}$
- Parameter vector

 $\theta = \langle \theta_1, \dots, \theta_K \rangle$ 

- N observations  $\mathbf{x} = \langle x_1, \dots, x_N \rangle$
- *i*th observation  $x_i \in \{1, \ldots, K\}$
- ▶ all but the *i*th observation  $\mathbf{x}_{-i}$
- N cluster indicators  $\mathbf{z} = \langle z_1, \dots, z_N \rangle$
- *i*th cluster indicator  $z_i \in \{1, \ldots, C\}$
- $\blacktriangleright$  all but the *i*th cluster assignment  $\mathbf{z}_{-i}$
- Parameter vector

$$\theta = \langle \theta_1, \ldots, \theta_K \rangle$$

► Collection of parameter vectors  $\boldsymbol{\theta} = \langle \theta^{(1)}, \dots, \theta^{(C)} \rangle$ 



Let's assume x to be 1 of K, and z to be 1 of C

categorical likelihood



Let's assume x to be 1 of K, and z to be 1 of C

- categorical likelihood
- uniform prior over mixture components, i.e. mixing weights are fixed and uniform



Let's assume x to be 1 of K, and z to be 1 of C

- categorical likelihood
- uniform prior over mixture components, i.e. mixing weights are fixed and uniform

$$\bullet \ \theta^{(c)} \in \Delta_{K-1}$$



Let's assume x to be 1 of K, and z to be 1 of C

- categorical likelihood
- uniform prior over mixture components, i.e. mixing weights are fixed and uniform

$$\bullet \ \theta^{(c)} \in \Delta_{K-1}$$

For  $i = 1, \ldots, N$ 

 $Z_i \sim \mathcal{U}(C)$ 



Let's assume x to be 1 of K, and z to be 1 of C

- categorical likelihood
- uniform prior over mixture components, i.e. mixing weights are fixed and uniform

$$\bullet \ \theta^{(c)} \in \Delta_{K-1}$$

For  $i = 1, \ldots, N$ 

$$Z_i \sim \mathcal{U}(C)$$

$$X_i | \boldsymbol{\theta}, \mathbf{z}_{-i}, z_i = c \sim \operatorname{Cat}(\theta^{(c)})$$
(1)



For  $i = 1, \ldots, N$ 

Let's assume x to be 1 of K, and z to be 1 of C

- categorical likelihood
- uniform prior over mixture components, i.e. mixing weights are fixed and uniform

$$\bullet \ \theta^{(c)} \in \Delta_{K-1}$$

 $Z_i \sim \mathcal{U}(C)$   $X_i | \boldsymbol{\theta}, \mathbf{z}_{-i}, z_i = c \sim \operatorname{Cat}(\theta^{(c)})$ (1)

What is a sensible conditional distribution  $X|\theta^{(c)} \sim \operatorname{Cat}(\theta^{(c)})$ ?

c = 1 (the blue cluster), K = 4





$$c = 1$$
 (the blue cluster),  $K = 4$ 



$$c = 1$$
 (the blue cluster),  $K = 4$ 



c = 1 (the blue cluster), K = 4

Can you make any assumptions before observing data?

What does Bayes rule tell you?

$$\underbrace{P(h|d)}_{\text{posterior}} =$$

What does Bayes rule tell you?



What does Bayes rule tell you?



What does Bayes rule tell you?



the likelihood tells you how well a hypothesis h explains the observed data d;

What does Bayes rule tell you?



the likelihood tells you how well a hypothesis h explains the observed data d;

the prior tells you how much h conforms to expectations about what a good hypothesis looks like regardless of observed data;

What does Bayes rule tell you?



- the likelihood tells you how well a hypothesis h explains the observed data d;
- the prior tells you how much h conforms to expectations about what a good hypothesis looks like regardless of observed data;
- the evidence tells you how well your model *M* explains the data, i.e. *P(d)* is actually *P(d|M)*
#### Bayes rule

What does Bayes rule tell you?



- the likelihood tells you how well a hypothesis h explains the observed data d;
- the prior tells you how much h conforms to expectations about what a good hypothesis looks like regardless of observed data;
- ► the evidence tells you how well your model *M* explains the data, i.e. *P(d)* is actually *P(d|M)*
- the posterior updates our beliefs about hypotheses in light of observed data.

An optimisation problem based on the (log-)likelihood function

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)$$

An optimisation problem based on the (log-)likelihood function

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h) = \underset{h}{\operatorname{arg\,max}} \underbrace{\log P(d|h)}_{\mathcal{L}(h)}$$
(3)

An optimisation problem based on the (log-)likelihood function

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h) = \underset{h}{\operatorname{arg\,max}} \underbrace{\log P(d|h)}_{\mathcal{L}(h)}$$
(3)

all hypotheses are equally likely a priori;

An optimisation problem based on the (log-)likelihood function

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h) = \underset{h}{\operatorname{arg\,max}} \underbrace{\log P(d|h)}_{\mathcal{L}(h)}$$
(3)

- all hypotheses are equally likely a priori;
- can be approached by coordinate ascent methods;

An optimisation problem based on the (log-)likelihood function

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h) = \underset{h}{\operatorname{arg\,max}} \underbrace{\log P(d|h)}_{\mathcal{L}(h)}$$
(3)

- all hypotheses are equally likely a priori;
- can be approached by coordinate ascent methods;
- local optimality guarantees;

#### All the same a priori



#### Before data, MLE is equally happy with the hypotheses on the left

Maximum a posteriori

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)P(h)$$
$$= \underset{h}{\operatorname{arg\,max}} \log P(d|h) + \log P(h)$$

• perhaps fine if P(h) has a single narrow peak

(4)

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)P(h)$$
  
= 
$$\underset{h}{\operatorname{arg\,max}} \log P(d|h) + \log P(h)$$
 (4)

- perhaps fine if P(h) has a single narrow peak
- priors often indicate preference for a subset of hypotheses over another, multiple peaks make optimisation considerably harder

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)P(h)$$
  
= 
$$\underset{h}{\operatorname{arg\,max}} \log P(d|h) + \log P(h)$$
 (4)

- perhaps fine if P(h) has a single narrow peak
- priors often indicate preference for a subset of hypotheses over another, multiple peaks make optimisation considerably harder
- still a point estimate, teaches us very little about the overall model (set of assumptions)

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)P(h)$$
  
= 
$$\underset{h}{\operatorname{arg\,max}} \log P(d|h) + \log P(h)$$
 (4)

- perhaps fine if P(h) has a single narrow peak
- priors often indicate preference for a subset of hypotheses over another, multiple peaks make optimisation considerably harder
- still a point estimate, teaches us very little about the overall model (set of assumptions)
- "I read before that Bayesian priors are just like regularisers,

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)P(h)$$
  
= 
$$\underset{h}{\operatorname{arg\,max}} \log P(d|h) + \log P(h)$$
 (4)

- perhaps fine if P(h) has a single narrow peak
- priors often indicate preference for a subset of hypotheses over another, multiple peaks make optimisation considerably harder
- still a point estimate, teaches us very little about the overall model (set of assumptions)
- "I read before that Bayesian priors are just like regularisers, I even know that a Gaussian prior is just  $L_2$  regularisation"

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)P(h)$$
  
= 
$$\underset{h}{\operatorname{arg\,max}} \log P(d|h) + \log P(h)$$
 (4)

- perhaps fine if P(h) has a single narrow peak
- priors often indicate preference for a subset of hypotheses over another, multiple peaks make optimisation considerably harder
- still a point estimate, teaches us very little about the overall model (set of assumptions)
- "I read before that Bayesian priors are just like regularisers, I even know that a Gaussian prior is just  $L_2$  regularisation"
  - that only covers the specification of a prior

Maximum a posteriori

$$h^{\star} = \underset{h}{\operatorname{arg\,max}} P(d|h)P(h)$$
  
= 
$$\underset{h}{\operatorname{arg\,max}} \log P(d|h) + \log P(h)$$
 (4)

- perhaps fine if P(h) has a single narrow peak
- priors often indicate preference for a subset of hypotheses over another, multiple peaks make optimisation considerably harder
- still a point estimate, teaches us very little about the overall model (set of assumptions)

"I read before that Bayesian priors are just like regularisers, I even know that a Gaussian prior is just  $L_2$  regularisation"

- that only covers the specification of a prior
- Bayesian modelling does not end at prior specification you need the crucial part: posterior inference

NLP1

Preliminaries

#### Bayesian modelling Dirichlet-Multinomial model

Applications



In a Bayesian model, parameters are no different from data

they are random variables much like data



In a Bayesian model, parameters are no different from data

- they are random variables much like data
- only they are not observed



In a Bayesian model, parameters are no different from data

- they are random variables much like data
- only they are not observed

Bayesians do condition on deterministic quantities

 $\triangleright$   $\beta$  here are called *hyperparameters* 



In a Bayesian model, parameters are no different from data

- they are random variables much like data
- only they are not observed

Bayesians do condition on deterministic quantities

- $\blacktriangleright$   $\beta$  here are called *hyperparameters*
- but most Bayesians leave those fixed (no search!)



In a Bayesian model, parameters are no different from data

- they are random variables much like data
- only they are not observed

Bayesians do condition on deterministic quantities

- $\triangleright$   $\beta$  here are called *hyperparameters*
- but most Bayesians leave those fixed (no search!)

We will study an example that illustrates important concepts Dirichlet-Multinomial model

### Dirichlet distribution

A distribution over the open simplex of K-dimensional vectors we denote the simplex by

$$\Delta_{K-1} = \left\{ \theta \in \mathbb{R}_{>0}^{K} : \sum_{k=1}^{K} \theta_{k} = 1 \right\} \subseteq \mathbb{R}_{>0}^{K}$$
(5)







#### Count vector

For observations x, where  $x_i$  is 1 of K define  $n^{(\mathbf{x})}$  as the K-dimensional vector such that

$$n_k = \sum_{i=1}^{N} [x_i = k]$$
(6)

#### Count vector

For observations x, where  $x_i$  is 1 of K define  $n^{(x)}$  as the K-dimensional vector such that

$$n_k = \sum_{i=1}^N [x_i = k] \tag{6}$$

Example: for K = 3 and N = 6

$$\mathbf{x} = \langle x_1 = 2, x_2 = 3, x_3 = 1, x_4 = 2, x_5 = 2, x_6 = 3 \rangle$$
  
$$n^{(\mathbf{x})} =$$

#### Count vector

For observations x, where  $x_i$  is 1 of K define  $n^{(x)}$  as the K-dimensional vector such that

$$n_k = \sum_{i=1}^N [x_i = k] \tag{6}$$

Example: for K = 3 and N = 6

$$\mathbf{x} = \langle x_1 = 2, x_2 = 3, x_3 = 1, x_4 = 2, x_5 = 2, x_6 = 3 \rangle$$
$$n^{(\mathbf{x})} = \langle n_1 = 1, n_2 = 3, n_3 = 2 \rangle$$

#### Gamma function

A generalisation of the factorial function to  $\ensuremath{\mathbb{R}}$ 

$$\Gamma(z) = \int_0^\infty \epsilon^{z-1} \exp(-\epsilon) d\epsilon$$
(7)

Properties

• 
$$\Gamma(n) = (n-1)!$$
 for positive integer  $n$   
•  $\Gamma(z) = (z-1)\Gamma(z-1)$ 

# **Dirchlet-Multinomial**



#### Model

$$\begin{aligned} \theta &| \beta \sim \text{Dir}(\beta) \\ X_i &| \theta \sim \text{Cat}(\theta) \quad \text{for } i = 1, \dots, N \end{aligned}$$
 (8)

#### **Dirchlet-Multinomial**



Model

$$\begin{array}{l}
\theta|eta \sim \operatorname{Dir}(eta)\\
X_i|\theta \sim \operatorname{Cat}(\theta) \quad \text{for } i = 1, \dots, N
\end{array}$$
(8)

Joint distribution

$$P(\mathbf{x}, \theta|\beta) = P(\theta)P(\mathbf{x}|\theta)$$
  
= Dir(\theta|\beta) Mult(n<sup>(x)</sup>|\theta, N) (9)

# Multinomial likelihood

For  $\theta \in \Delta_{K-1}$ 

$$P(\mathbf{x}|\theta) = \text{Mult}(n^{(\mathbf{x})}|\theta, N)$$

# Multinomial likelihood

For  $\theta \in \Delta_{K-1}$ 

$$P(\mathbf{x}|\theta) = \text{Mult}(n^{(\mathbf{x})}|\theta, N)$$
$$= \frac{N!}{\prod_{k=1}^{K} n_k!} \prod_{k=1}^{K} \theta_k^{n_k}$$

#### Multinomial likelihood For $\theta \in \Delta_{K-1}$

$$P(\mathbf{x}|\theta) = \operatorname{Mult}(n^{(\mathbf{x})}|\theta, N)$$
$$= \frac{N!}{\prod_{k=1}^{K} n_k!} \prod_{k=1}^{K} \theta_k^{n_k}$$
$$= \frac{\Gamma(\sum_{k=1}^{K} n_k + 1)}{\prod_{k=1}^{K} \Gamma(n_k + 1)} \prod_{k=1}^{K} \theta_k^{n_k}$$

#### Multinomial likelihood For $\theta \in \Delta_{K-1}$

$$\in \Delta_{K-1}$$

$$P(\mathbf{x}|\theta) = \operatorname{Mult}(n^{(\mathbf{x})}|\theta, N)$$

$$= \frac{N!}{\prod_{k=1}^{K} n_k!} \prod_{k=1}^{K} \theta_k^{n_k}$$

$$= \frac{\Gamma(\sum_{k=1}^{K} n_k + 1)}{\prod_{k=1}^{K} \Gamma(n_k + 1)} \prod_{k=1}^{K} \theta_k^{n_k}$$

$$(10)$$

Example: for K = 3 and N = 6

$$\theta = \langle \theta_1 = 0.2, \theta_2 = 0.3, \theta_3 = 0.5 \rangle$$
  

$$\mathbf{x} = \langle x_1 = 2, x_2 = 3, x_3 = 1, x_4 = 2, x_5 = 2, x_6 = 3 \rangle$$
  

$$n^{(\mathbf{x})} = \langle n_1 = 1, n_2 = 3, n_3 = 2 \rangle$$

$$P(\mathbf{x}|\theta) = \frac{\Gamma(\ldots)}{\prod \cdots} \theta_1^1 \times \theta_2^3 \times \theta_3^2$$

Wilker Aziz

NLP1 2019

# Dirichlet prior

For  $\beta \in \mathbb{R}_{>0}^K$ 

$$\operatorname{Dir}(\boldsymbol{\theta}|\boldsymbol{\beta}) = \frac{\Gamma(\sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(\beta_k)} \prod_{k=1}^{K} \theta_k^{\beta_k - 1}$$

# Dirichlet prior

For  $\beta \in \mathbb{R}_{>0}^K$ 

$$\operatorname{Dir}(\theta|\beta) = \frac{\Gamma(\sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(\beta_k)} \prod_{k=1}^{K} \theta_k^{\beta_k - 1}$$

$$\propto \prod_{k=1}^{K} \theta_k^{\beta_k - 1}$$
(11)

## Dirichlet prior

For  $\beta \in \mathbb{R}_{>0}^K$ 

$$\operatorname{Dir}(\theta|\beta) = \frac{\Gamma(\sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(\beta_k)} \prod_{k=1}^{K} \theta_k^{\beta_k - 1}$$

$$\propto \prod_{k=1}^{K} \theta_k^{\beta_k - 1}$$
(11)

We call

$$\int_{\Delta_{K-1}} \prod_{k=1}^{K} \theta_k^{\beta_k - 1} = \frac{\prod_{k=1}^{K} \Gamma(\beta_k)}{\Gamma(\sum_{k=1}^{K} \beta_k)}$$

the Dirichlet normaliser

Posterior

#### $P(\boldsymbol{\theta}|\mathbf{x},\beta) \propto$

#### Posterior

 $P(\boldsymbol{\theta}|\mathbf{x},\boldsymbol{\beta}) \propto P(\mathbf{x}|\boldsymbol{\theta})P(\boldsymbol{\theta}|\boldsymbol{\beta})$








$$\begin{split} P(\boldsymbol{\theta}|\mathbf{x},\boldsymbol{\beta}) &\propto P(\mathbf{x}|\boldsymbol{\theta})P(\boldsymbol{\theta}|\boldsymbol{\beta}) \\ &\propto \underbrace{\frac{\Gamma(\sum_{k=1}^{K}n_{k}+1)}{\prod_{k=1}^{K}\Gamma(n_{k}+1)}\prod_{k=1}^{K}\theta_{k}^{n_{k}}}_{\mathrm{Mult}(n^{(\mathbf{x})}|\boldsymbol{\theta})} \times \underbrace{\frac{\Gamma(\sum_{k=1}^{K}\beta_{k})}{\prod_{k=1}^{K}\Gamma(\beta_{k})}\prod_{k=1}^{K}\theta_{k}^{\beta_{k}-1}}_{\mathrm{Dir}(\boldsymbol{\theta}|\boldsymbol{\beta})} \\ &\propto \prod_{k=1}^{K}\theta_{k}^{n_{k}} \times \prod_{k=1}^{K}\theta_{k}^{\beta_{k}-1} \\ &= \prod_{k=1}^{K}\theta_{k}^{n_{k}+\beta_{k}-1} \end{split}$$

$$P(\theta|\mathbf{x},\beta) \propto P(\mathbf{x}|\theta)P(\theta|\beta)$$

$$\propto \underbrace{\frac{\Gamma(\sum_{k=1}^{K} n_k + 1)}{\prod_{k=1}^{K} \Gamma(n_k + 1)} \prod_{k=1}^{K} \theta_k^{n_k}}_{\operatorname{Mult}(n^{(\mathbf{x})}|\theta)} \times \underbrace{\frac{\Gamma(\sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(\beta_k)} \prod_{k=1}^{K} \theta_k^{\beta_k - 1}}_{\operatorname{Dir}(\theta|\beta)}}_{\operatorname{Dir}(\theta|\beta)}$$

$$\propto \prod_{k=1}^{K} \theta_k^{n_k} \times \prod_{k=1}^{K} \theta_k^{\beta_k - 1}$$

$$= \prod_{k=1}^{K} \theta_k^{n_k + \beta_k - 1} \propto \operatorname{Dir}(\theta|n^{(\mathbf{x})} + \beta)$$

$$\begin{split} P(\boldsymbol{\theta}|\mathbf{x},\boldsymbol{\beta}) &\propto P(\mathbf{x}|\boldsymbol{\theta})P(\boldsymbol{\theta}|\boldsymbol{\beta}) \\ &\propto \underbrace{\frac{\Gamma(\sum_{k=1}^{K}n_{k}+1)}{\prod_{k=1}^{K}\Gamma(n_{k}+1)}\prod_{k=1}^{K}\theta_{k}^{n_{k}}}_{\mathrm{Mult}(n^{(\mathbf{x})}|\boldsymbol{\theta})} \times \underbrace{\frac{\Gamma(\sum_{k=1}^{K}\beta_{k})}{\prod_{k=1}^{K}\Gamma(\beta_{k})}\prod_{k=1}^{K}\theta_{k}^{\beta_{k}-1}}_{\mathrm{Dir}(\boldsymbol{\theta}|\boldsymbol{\beta})} \\ &\propto \prod_{k=1}^{K}\theta_{k}^{n_{k}} \times \prod_{k=1}^{K}\theta_{k}^{\beta_{k}-1} \\ &= \prod_{k=1}^{K}\theta_{k}^{n_{k}+\beta_{k}-1} \propto \mathrm{Dir}(\boldsymbol{\theta}|n^{(\mathbf{x})}+\boldsymbol{\beta}) \end{split}$$

Thus

$$P(\theta | \mathbf{x}, \beta) = \underbrace{\prod_{\substack{1 \text{ normaliser of Dir}(n^{(\mathbf{x})} + \beta)}} \prod_{k=1}^{K} \theta_k^{n_k + \beta_k - 1}$$
(12)

$$\begin{split} P(\boldsymbol{\theta}|\mathbf{x},\boldsymbol{\beta}) &\propto P(\mathbf{x}|\boldsymbol{\theta})P(\boldsymbol{\theta}|\boldsymbol{\beta}) \\ &\propto \underbrace{\frac{\Gamma(\sum_{k=1}^{K}n_{k}+1)}{\prod_{k=1}^{K}\Gamma(n_{k}+1)}\prod_{k=1}^{K}\theta_{k}^{n_{k}}}_{\mathrm{Mult}(n^{(\mathbf{x})}|\boldsymbol{\theta})} \times \underbrace{\frac{\Gamma(\sum_{k=1}^{K}\beta_{k})}{\prod_{k=1}^{K}\Gamma(\beta_{k})}\prod_{k=1}^{K}\theta_{k}^{\beta_{k}-1}}_{\mathrm{Dir}(\boldsymbol{\theta}|\boldsymbol{\beta})} \\ &\propto \prod_{k=1}^{K}\theta_{k}^{n_{k}} \times \prod_{k=1}^{K}\theta_{k}^{\beta_{k}-1} \\ &= \prod_{k=1}^{K}\theta_{k}^{n_{k}+\beta_{k}-1} \propto \mathrm{Dir}(\boldsymbol{\theta}|n^{(\mathbf{x})}+\boldsymbol{\beta}) \end{split}$$

Thus

$$P(\theta|\mathbf{x},\beta) = \underbrace{\frac{\Gamma(N+\sum_{k=1}^{K}\beta_k)}{\prod_{k=1}^{K}\Gamma(n_k+\beta_k)}}_{\frac{1}{\text{normaliser}} \text{ of } \text{Dir}(n^{(\mathbf{x})}+\beta)} \prod_{k=1}^{K} \theta_k^{n_k+\beta_k-1}$$
(12)

NLP1 2019

#### Posterior predictive distribution

Suppose a new data point  $x_{N+1} = j$  is available



$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} P(\theta, x_{N+1} | \mathbf{x}, \beta) d\theta$$

 $x_{N+1}$  is independent of  ${\bf x}$  given  $\theta$ 

#### Posterior predictive distribution

Suppose a new data point  $x_{N+1} = j$  is available



$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} P(\theta, x_{N+1} | \mathbf{x}, \beta) d\theta$$
$$= \int_{\Delta_{K-1}} \underbrace{P(x_{N+1} = j | \theta)}_{\text{likelihood}} \underbrace{P(\theta | \mathbf{x}, \beta)}_{\text{posterior}} d\theta$$

 $x_{N+1}$  is independent of  ${\bf x}$  given  $\theta$ 

Suppose a new data point  $x_{N+1} = j$  is available

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\text{likelihood}}}_{\text{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\text{posterior}}}_{\text{posterior}} d\theta$$

$$\begin{split} P(x_{N+1} = j | \mathbf{x}, \beta) &= \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}_{\text{likelihood}} \underbrace{P(\theta | \mathbf{x}, \beta)}_{\text{posterior}} d\theta}_{\text{posterior}} \end{split}$$
$$= \int_{\Delta_{K-1}} \theta_j \times d\theta \end{split}$$



$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\text{likelihood}}}_{\text{posterior}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\text{posterior}}}_{\text{posterior}} d\theta$$
$$= \int_{\Delta_{K-1}} \theta_j \times \underbrace{\frac{\Gamma(N + \sum_{k=1}^K \beta_k)}{\prod_{k=1}^K \Gamma(n_k + \beta_k)}}_{\text{constant wrt } \theta} \prod_{k=1}^K \theta_k^{n_k + \beta_k - 1} d\theta$$
$$= \underbrace{\frac{\Gamma(N + \sum_{k=1}^K \beta_k)}{\prod_{k=1}^K \Gamma(n_k + \beta_k)}}_{\text{constant wrt } \theta} \int_{\Delta_{K-1}} \theta_j \times d\theta$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\mathsf{likelihood}}}_{\mathsf{posterior}} \underbrace{P(\theta | \mathbf{x}, \beta)}_{\mathsf{posterior}} \mathrm{d}\theta$$
$$= \int_{\Delta_{K-1}} \theta_j \times \underbrace{\frac{\Gamma(N + \sum_{k=1}^K \beta_k)}{\prod_{k=1}^K \Gamma(n_k + \beta_k)}}_{\mathsf{constant wrt } \theta} \prod_{k=1}^K \theta_k^{n_k + \beta_k - 1} \mathrm{d}\theta$$
$$= \underbrace{\frac{\Gamma(N + \sum_{k=1}^K \beta_k)}{\prod_{k=1}^K \Gamma(n_k + \beta_k)}}_{\mathsf{constant wrt } \theta} \int_{\Delta_{K-1}} \theta_j \times \underbrace{\theta_j^{n_j + \beta_j - 1}}_{\prod_{k=1}^K \theta_k^{n_k + \beta_k - 1}} \operatorname{d}\theta$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\text{likelihood}}}_{\text{posterior}} \underbrace{P(\theta | \mathbf{x}, \beta)}_{\text{posterior}} d\theta$$
$$= \int_{\Delta_{K-1}} \theta_j \times \underbrace{\frac{\Gamma(N + \sum_{k=1}^K \beta_k)}{\prod_{k=1}^K \Gamma(n_k + \beta_k)}}_{\text{constant wrt } \theta} \prod_{k=1}^K \theta_k^{n_k + \beta_k - 1} d\theta$$
$$= \underbrace{\frac{\Gamma(N + \sum_{k=1}^K \beta_k)}{\prod_{k=1}^K \Gamma(n_k + \beta_k)}}_{\text{constant wrt } \theta} \int_{\Delta_{K-1}} \theta_j \times \underbrace{\theta_j^{n_j + \beta_j - 1}}_{\prod_{k=1}^K \theta_k^{n_k + \beta_k - 1}} d\theta$$
$$= \frac{\Gamma(N + \sum_{k=1}^K \beta_k)}{\prod_{k=1}^K \Gamma(n_k + \beta_k)} \int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} d\theta$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\mathsf{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\mathsf{posterior}}}_{\mathsf{posterior}} \mathrm{d}\theta$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} \mathrm{d}\theta}_{k \neq j}$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\text{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\text{posterior}}}_{\text{posterior}} d\theta$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{\substack{k \neq j}} \theta_k^{n_k + \beta_k - 1} d\theta}_{\text{Dir normaliser}}$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\mathsf{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\mathsf{posterior}}}_{\mathsf{posterior}} \mathrm{d}\theta$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} \mathrm{d}\theta}_{\mathsf{Dir normaliser}}$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\frac{1}{\sum_{k=1}^{K} \beta_k}}_{\mathsf{Dir normaliser}}$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\mathsf{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\mathsf{posterior}}}_{\mathsf{posterior}} \mathrm{d}\theta$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} \mathrm{d}\theta}_{\mathsf{Dir normaliser}}$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k + 1)}{\Gamma(N + \sum_{k=1}^{K} \beta_k + 1)}$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{P(x_{N+1} = j | \theta)}_{\text{likelihood}} \underbrace{P(\theta | \mathbf{x}, \beta)}_{\text{posterior}} d\theta$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} d\theta}_{\text{Dir normaliser}}$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \frac{\Gamma(n_j + \beta_j + 1) \prod_{k \neq j} \Gamma(n_k + \beta_k)}{\Gamma(N + \sum_{k=1}^{K} \beta_k + 1)}$$

$$\begin{split} P(x_{N+1} = j | \mathbf{x}, \beta) &= \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\mathsf{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\mathsf{posterior}}}_{\mathsf{posterior}} \mathrm{d}\theta \\ &= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} \mathrm{d}\theta}_{\mathsf{Dir normaliser}} \\ &= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \frac{\Gamma(n_j + \beta_j + 1) \prod_{k \neq j} \Gamma(n_k + \beta_k)}{\Gamma(N + \sum_{k=1}^{K} \beta_k + 1)} \\ &= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \frac{(n_j + \beta_j) \Gamma(n_j + \beta_j) \prod_{k \neq j} \Gamma(n_k + \beta_k)}{(N + \sum_{k=1}^{K} \beta_k) \Gamma(N + \sum_{k=1}^{K} \beta_k)} \end{split}$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\mathsf{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\mathsf{posterior}}}_{\mathsf{posterior}} \mathrm{d}\theta$$

$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} \mathrm{d}\theta}_{\mathsf{Dir normaliser}}$$

$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\frac{\Gamma(n_j + \beta_j + 1) \prod_{k \neq j} \Gamma(n_k + \beta_k)}{\Gamma(N + \sum_{k=1}^{K} \beta_k + 1)}}_{(N + \sum_{k=1}^{K} \beta_k + 1)}$$

$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\frac{(n_j + \beta_j)\Gamma(n_j + \beta_j) \prod_{k \neq j} \Gamma(n_k + \beta_k)}{(N + \sum_{k=1}^{K} \beta_k)\Gamma(N + \sum_{k=1}^{K} \beta_k)}}_{(N + \sum_{k=1}^{K} \beta_k)\Gamma(N + \sum_{k=1}^{K} \beta_k)}$$

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \int_{\Delta_{K-1}} \underbrace{\frac{P(x_{N+1} = j | \theta)}{\mathsf{likelihood}} \underbrace{\frac{P(\theta | \mathbf{x}, \beta)}{\mathsf{posterior}}}_{\mathsf{posterior}} \mathrm{d}\theta}$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\int_{\Delta_{K-1}} \theta_j^{n_j + \beta_j} \prod_{k \neq j} \theta_k^{n_k + \beta_k - 1} \mathrm{d}\theta}_{\mathsf{Dir normaliser}}$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\frac{\Gamma(n_j + \beta_j + 1) \prod_{k \neq j} \Gamma(n_k + \beta_k)}{\Gamma(N + \sum_{k=1}^{K} \beta_k + 1)}}_{(N + \sum_{k=1}^{K} \beta_k + 1)}$$
$$= \frac{\Gamma(N + \sum_{k=1}^{K} \beta_k)}{\prod_{k=1}^{K} \Gamma(n_k + \beta_k)} \underbrace{\frac{(n_j + \beta_j)\Gamma(n_j + \beta_j) \prod_{k \neq j} \Gamma(n_k + \beta_k)}{(N + \sum_{k=1}^{K} \beta_k)\Gamma(N + \sum_{k=1}^{K} \beta_k)}}_{(N + \sum_{k=1}^{K} \beta_k)\Gamma(N + \sum_{k=1}^{K} \beta_k)}$$
$$= \frac{n_j + \beta_j}{N + \sum_{k=1}^{K} \beta_k}$$

Wilker Aziz

NLP1 2019

Dirchlet-Multinomial (overview)



Joint distribution

$$P(\mathbf{x}, \theta | \beta) = P(\theta) P(\mathbf{x} | \theta)$$
  
= Dir(\theta|\beta) Mult(n<sup>(x)</sup>|\theta, N) (13)

Posterior

$$P(\theta|\mathbf{x},\beta) = \text{Dir}(\theta|n^{(\mathbf{x})} + \beta)$$
(14)

Predictive posterior

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \frac{n_j + \beta_j}{N + \sum_{k=1}^K \beta_k}$$
(15)

Random variables are called exchangeable under a model when all permutations of the set of outcomes have the same probability

Random variables are called exchangeable under a model when all permutations of the set of outcomes have the same probability

 in our Dirichlet-Multinomial model any re-ordering of the observations is equally likely to occur

Random variables are called exchangeable under a model when all permutations of the set of outcomes have the same probability

 in our Dirichlet-Multinomial model any re-ordering of the observations is equally likely to occur

Combine that fact with the predictive posterior result

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \frac{n_j + \beta_j}{N + \sum_{k=1}^K \beta_k}$$
(16)

Random variables are called exchangeable under a model when all permutations of the set of outcomes have the same probability

 in our Dirichlet-Multinomial model any re-ordering of the observations is equally likely to occur

Combine that fact with the predictive posterior result

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \frac{n_j + \beta_j}{N + \sum_{k=1}^K \beta_k}$$
(16)

and we can single out any observation, e.g.  $\mathbf{x}_i$ 

$$P(\mathbf{x}_i = j | \mathbf{x}_{-i}, \beta) = -----$$
(17)

Random variables are called exchangeable under a model when all permutations of the set of outcomes have the same probability

 in our Dirichlet-Multinomial model any re-ordering of the observations is equally likely to occur

Combine that fact with the predictive posterior result

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \frac{n_j + \beta_j}{N + \sum_{k=1}^K \beta_k}$$
(16)

and we can single out any observation, e.g.  $\mathbf{x}_i$ 

$$P(\mathbf{x}_i = j | \mathbf{x}_{-i}, \beta) = \frac{1}{N - 1 + \sum_{k=1}^{K} \beta_k}$$
(17)

Random variables are called exchangeable under a model when all permutations of the set of outcomes have the same probability

 in our Dirichlet-Multinomial model any re-ordering of the observations is equally likely to occur

Combine that fact with the predictive posterior result

$$P(x_{N+1} = j | \mathbf{x}, \beta) = \frac{n_j + \beta_j}{N + \sum_{k=1}^K \beta_k}$$
(16)

and we can single out any observation, e.g.  $\mathbf{x}_i$ 

$$P(\mathbf{x}_i = j | \mathbf{x}_{-i}, \beta) = \frac{n_j^{(\mathbf{x}_{-i})} + \beta_j}{N - 1 + \sum_{k=1}^K \beta_k}$$
(17)

Friends do not let friends optimise

no point estimates, we use all possible model parameters

- no point estimates, we use all possible model parameters
- this is called Bayesian inference, or simply, inference

- no point estimates, we use all possible model parameters
- this is called Bayesian inference, or simply, inference
- Bayesian models have memory: the posterior summarises what we learnt from data

- no point estimates, we use all possible model parameters
- this is called Bayesian inference, or simply, inference
- Bayesian models have memory: the posterior summarises what we learnt from data
- ► If we collect more data  $\mathbf{x}'$ , we can update the posterior,  $P(\theta|\mathbf{x}, \mathbf{x}', \beta) = \text{Dir}(\theta|n^{(\mathbf{x})} + n^{(\mathbf{x}')} + \beta)$
# Summary

Friends do not let friends optimise

- no point estimates, we use all possible model parameters
- this is called Bayesian inference, or simply, inference
- Bayesian models have memory: the posterior summarises what we learnt from data
- If we collect more data  $\mathbf{x}'$ , we can update the posterior,  $P(\theta|\mathbf{x}, \mathbf{x}', \beta) = \text{Dir}(\theta|n^{(\mathbf{x})} + n^{(\mathbf{x}')} + \beta)$
- MLE is memoryless: there is one fixed θ, no matter how much more data you see, θ will never change

# NLP1

Preliminaries

Bayesian modelling

Applications

Wilker Aziz NLP1 2019



Suppose the task is to label observations  $x_1, \ldots, x_N$  with cluster assignments  $z_1, \ldots, z_n$ .



Suppose the task is to label observations  $x_1, \ldots, x_N$  with cluster assignments  $z_1, \ldots, z_n$ .

We would need to explore our posterior  $p(\mathbf{z}|\mathbf{x}, \alpha, \beta)$ but this is a very complex object



Suppose the task is to label observations  $x_1, \ldots, x_N$  with cluster assignments  $z_1, \ldots, z_n$ .

We would need to explore our posterior  $p(\mathbf{z}|\mathbf{x}, \alpha, \beta)$ but this is a very complex object

- conditioning on x induces (undirected) dependencies amongst
   z and θ this is called *moralisation* in PGMs
- integrating θ out induces (undirected) dependencies amongst
   z and x this is called *variable elimination* in PGMs



Let's try to explore the posterior distribution one variable at a time, that is, let's try and characterise  $P(z_i, x_i | \mathbf{z}_{-i}, \mathbf{x}_{-i}, \alpha, \beta)$ 

This will lead to a class of approximate inference algorithms known as Markov Chain Monte Carlo (MCMC)

Wilker Aziz NLP1 2019



due to marginal independence of  $\theta_c$  (with respect to  $\theta_{c'\neq c}$ )



Define counts based on joint assignments to  $\mathbf{x}_{-i}, \mathbf{z}_{-i}$ 1/C  $n_{c,k} = \sum_{j \neq i} [z_j = c] [x_j = k]$   $n_c = \sum_{k=1}^K n_{c,k}$   $\mathbf{n}_c = [n_{c,1}, \dots, n_{c,K}]$ 







What does it mean to have uniform prior over components?

What does it mean to have uniform prior over components?

unlike it may seem, it does not mean to promote diversity! Let's see whether the posterior is *peaked* 

$$P(z|x) = \frac{\frac{1}{C} \times P(x|z)}{P(x)} \propto P(x|z)$$

What does it mean to have uniform prior over components?

unlike it may seem, it does not mean to promote diversity! Let's see whether the posterior is *peaked* 

$$P(z|x) = \frac{\frac{1}{C} \times P(x|z)}{P(x)} \propto P(x|z)$$

uniform prior leaves it up to the likelihood to control sparsity

What does it mean to have uniform prior over components?

unlike it may seem, it does not mean to promote diversity! Let's see whether the posterior is *peaked* 

$$P(z|x) = \frac{\frac{1}{C} \times P(x|z)}{P(x)} \propto P(x|z)$$

 ▶ uniform prior leaves it up to the likelihood to control sparsity
 ▶ luckily we are promoting sparse likelihoods X|z because θ<sup>(z)</sup> ~ Dir(β)

What does it mean to have uniform prior over components?

unlike it may seem, it does not mean to promote diversity! Let's see whether the posterior is *peaked* 

$$P(z|x) = \frac{\frac{1}{C} \times P(x|z)}{P(x)} \propto P(x|z)$$

uniform prior leaves it up to the likelihood to control sparsity

- ► luckily we are promoting sparse likelihoods X|z because θ<sup>(z)</sup> ~ Dir(β)
- ▶ but *P*(*z*) has nothing to do with it!

What does it mean to have uniform prior over components?

unlike it may seem, it does not mean to promote diversity! Let's see whether the posterior is *peaked* 

$$P(z|x) = \frac{\frac{1}{C} \times P(x|z)}{P(x)} \propto P(x|z)$$

uniform prior leaves it up to the likelihood to control sparsity

- ► luckily we are promoting sparse likelihoods X|z because θ<sup>(z)</sup> ~ Dir(β)
- but P(z) has nothing to do with it!

Is there really no preference we can express about P(z)?

What does it mean to have uniform prior over components?

unlike it may seem, it does not mean to promote diversity! Let's see whether the posterior is *peaked* 

$$P(z|x) = \frac{\frac{1}{C} \times P(x|z)}{P(x)} \propto P(x|z)$$

uniform prior leaves it up to the likelihood to control sparsity

- ► luckily we are promoting sparse likelihoods X|z because θ<sup>(z)</sup> ~ Dir(β)
- ▶ but *P*(*z*) has nothing to do with it!

Is there really no preference we can express about P(z)?

what about preferring to use fewer components?

### Sparse prior over mixing weights

Say we have 10 components, how do you want to use them?









$$P(x_i = k, z_i = c | \mathbf{x}_{-i}, \mathbf{z}_{-i}, \alpha, \beta)$$



$$P(x_{i} = k, z_{i} = c | \mathbf{x}_{-i}, \mathbf{z}_{-i}, \alpha, \beta)$$

$$= \int \int \underbrace{P(\phi | \mathbf{z}_{-i}, \mathbf{x}_{-i}, \alpha)}_{\text{Dir}(\phi | \alpha + \mathbf{t})} \underbrace{P(\theta_{c} | \mathbf{z}_{-i}, \mathbf{x}_{-i}, \beta)}_{\text{Dir}(\theta_{c} | \beta + \mathbf{n}_{c})} \underbrace{P(x_{i} = k, z_{i} = c | \theta_{c}, \phi)}_{\phi_{c} \times \theta_{c,k}} d\theta_{c} d\phi$$



$$\begin{split} P(x_i = k, z_i = c | \mathbf{x}_{-i}, \mathbf{z}_{-i}, \alpha, \beta) \\ &= \int \int \underbrace{P(\phi | \mathbf{z}_{-i}, \mathbf{x}_{-i}, \alpha)}_{\operatorname{Dir}(\phi | \alpha + \mathbf{t})} \underbrace{P(\theta_c | \mathbf{z}_{-i}, \mathbf{x}_{-i}, \beta)}_{\operatorname{Dir}(\theta_c | \beta + \mathbf{n}_c)} \underbrace{P(x_i = k, z_i = c | \theta_c, \phi)}_{\phi_c \times \theta_{c,k}} \mathrm{d}\theta_c \mathrm{d}\phi \\ &= \int_{\Delta_{C-1}} \phi_c \operatorname{Dir}(\phi | \alpha + \mathbf{t}) \mathrm{d}\phi \times \int_{\Delta_{K-1}} \theta_{c,k} \operatorname{Dir}(\theta_c | \beta + \mathbf{n}_c) \mathrm{d}\theta_c \end{split}$$

Wilker Aziz NLP1 2019



$$\begin{split} P(x_i = k, z_i = c | \mathbf{x}_{-i}, \mathbf{z}_{-i}, \alpha, \beta) \\ &= \int \int \underbrace{P(\phi | \mathbf{z}_{-i}, \mathbf{x}_{-i}, \alpha)}_{\operatorname{Dir}(\phi | \alpha + \mathbf{t})} \underbrace{P(\theta_c | \mathbf{z}_{-i}, \mathbf{x}_{-i}, \beta)}_{\operatorname{Dir}(\theta_c | \beta + \mathbf{n}_c)} \underbrace{P(x_i = k, z_i = c | \theta_c, \phi)}_{\phi_c \times \theta_{c,k}} \, \mathrm{d}\theta_c \mathrm{d}\phi \\ &= \int_{\Delta_{C-1}} \phi_c \operatorname{Dir}(\phi | \alpha + \mathbf{t}) \mathrm{d}\phi \times \int_{\Delta_{K-1}} \theta_{c,k} \operatorname{Dir}(\theta_c | \beta + \mathbf{n}_c) \mathrm{d}\theta_c \\ &= \frac{n_c + \alpha}{N - 1 + C\alpha} \times \frac{n_{c,k} + \beta}{n_c + K\beta} \end{split}$$





Define counts based on joint assignments to  $\mathbf{x}_{-i}, \mathbf{z}_{-i}$ 

$$m_{b,c} = \sum_{j \neq i} [z_{j-1} = b][z_j = c]$$
$$m_b = \sum_{c=1}^C m_{b,c}$$



 $P(z_i = c | \mathbf{x}, \mathbf{z}_{-i}, \alpha, \beta)$ 

Define counts based on joint assignments to  $\mathbf{x}_{-i}, \mathbf{z}_{-i}$ 

$$m_{b,c} = \sum_{j \neq i} [z_{j-1} = b][z_j = c]$$
$$m_b = \sum_{c=1}^C m_{b,c}$$



Define counts based on joint assignments to  $\mathbf{x}_{-i}, \mathbf{z}_{-i}$ 

$$m_{b,c} = \sum_{j \neq i} [z_{j-1} = b][z_j = c]$$

$$m_{b,c} = \sum_{j \neq i} [z_{j-1} = b][z_j = c]$$

$$m_b = \sum_{c=1}^C m_{b,c}$$

$$P(z_i = c | \mathbf{x}, \mathbf{z}_{-i}, \alpha, \beta) \text{ note that } \begin{cases} z_{i-1} = b \\ z_{i+1} = d \end{cases} \text{ is in } \mathbf{z}_{-i}$$



Wilker Aziz NLP1 2019



Wilker Aziz NLP1 2019

We draw from the posterior  $P(\mathbf{z}|\mathbf{x})$  via a Markov chain of random states  $Y_1, \ldots, Y_T$  where  $P(y_t|y_{< t}) = P(y_t|y_{t-1})$ 

We draw from the posterior  $P(\mathbf{z}|\mathbf{x})$  via a Markov chain of random states  $Y_1, \ldots, Y_T$  where  $P(y_t|y_{< t}) = P(y_t|y_{t-1})$ 

▶ the transition probability from y to y' is coded in a matrix  $\mathbf{P}$  $P_{ij}$  corresponds to P(Y = i, Y = j)

We draw from the posterior  $P(\mathbf{z}|\mathbf{x})$  via a Markov chain of random states  $Y_1, \ldots, Y_T$  where  $P(y_t|y_{< t}) = P(y_t|y_{t-1})$ 

- ► the transition probability from y to y' is coded in a matrix P P<sub>ij</sub> corresponds to P(Y = i, Y = j)
- under certain conditions the chain converges to a stationary distribution  $\pi$  such that  $\mathbf{P}\pi = \pi$

We draw from the posterior  $P(\mathbf{z}|\mathbf{x})$  via a Markov chain of random states  $Y_1, \ldots, Y_T$  where  $P(y_t|y_{< t}) = P(y_t|y_{t-1})$ 

- ► the transition probability from y to y' is coded in a matrix P P<sub>ij</sub> corresponds to P(Y = i, Y = j)
- under certain conditions the chain converges to a stationary distribution  $\pi$  such that  $\mathbf{P}\pi = \pi$
- possible states are assignments to the variables in the model
- ▶ the transition probability from y to y' is coded in a matrix P P<sub>ij</sub> corresponds to P(Y = i, Y = j)
- under certain conditions the chain converges to a stationary distribution  $\pi$  such that  $\mathbf{P}\pi = \pi$
- possible states are assignments to the variables in the model
- **b** by defining **P** properly we guarantee that  $\pi$  is the true posterior

- ▶ the transition probability from y to y' is coded in a matrix P P<sub>ij</sub> corresponds to P(Y = i, Y = j)
- under certain conditions the chain converges to a stationary distribution  $\pi$  such that  $\mathbf{P}\pi = \pi$
- possible states are assignments to the variables in the model
- by defining  ${f P}$  properly we guarantee that  $\pi$  is the true posterior
- $\blacktriangleright$  once the chain has converged each  $y_t$  will be a sample from the posterior

- the transition probability from y to y' is coded in a matrix P P<sub>ij</sub> corresponds to P(Y = i, Y = j)
- under certain conditions the chain converges to a stationary distribution  $\pi$  such that  $\mathbf{P}\pi = \pi$
- possible states are assignments to the variables in the model
- by defining  ${f P}$  properly we guarantee that  $\pi$  is the true posterior
- once the chain has converged each y<sub>t</sub> will be a sample from the posterior
- we can design P by decomposing it  $P_1 \cdots P_M$ where each component satisfies  $P_k(y, y')\pi(y) = P_k(y', y)\pi(y')$

- the transition probability from y to y' is coded in a matrix P P<sub>ij</sub> corresponds to P(Y = i, Y = j)
- under certain conditions the chain converges to a stationary distribution  $\pi$  such that  $\mathbf{P}\pi = \pi$
- possible states are assignments to the variables in the model
- by defining  ${f P}$  properly we guarantee that  $\pi$  is the true posterior
- once the chain has converged each y<sub>t</sub> will be a sample from the posterior
- we can design P by decomposing it  $P_1 \cdots P_M$ where each component satisfies  $P_k(y, y')\pi(y) = P_k(y', y)\pi(y')$
- applying each of P<sub>k</sub> in turn or choosing P<sub>k</sub> at random produces a P that satisfies the necessary conditions

We want to sample from  $P(\mathbf{z}|\mathbf{x})$  with a Markov chain a state  $y_t = \mathbf{z}^{(t)}$  is the *t*-th assignment to  $\mathbf{z}$ 

We want to sample from  $P(\mathbf{z}|\mathbf{x})$  with a Markov chain a state  $y_t = \mathbf{z}^{(t)}$  is the *t*-th assignment to  $\mathbf{z}$ 

To obtain a new state we

1. start a draft state  $\mathbf{z} = \mathbf{z}^{(t-1)}$ 

We want to sample from  $P(\mathbf{z}|\mathbf{x})$  with a Markov chain a state  $y_t = \mathbf{z}^{(t)}$  is the *t*-th assignment to  $\mathbf{z}$ 

To obtain a new state we

1. start a draft state  $\mathbf{z} = \mathbf{z}^{(t-1)}$ 

2. repeat for 
$$i = 1, \ldots, N$$

We want to sample from  $P(\mathbf{z}|\mathbf{x})$  with a Markov chain a state  $y_t = \mathbf{z}^{(t)}$  is the t-th assignment to  $\mathbf{z}$ 

To obtain a new state we

- 1. start a draft state  $\mathbf{z} = \mathbf{z}^{(t-1)}$
- 2. repeat for  $i = 1, \ldots, N$ 
  - ▶ resample Z<sub>i</sub> ~ P(z<sub>i</sub>|x<sub>-i</sub>, z<sub>-i</sub>) only variables in the Markov blanket of z<sub>i</sub> play a role that's why this is feasible

We want to sample from  $P(\mathbf{z}|\mathbf{x})$  with a Markov chain a state  $y_t = \mathbf{z}^{(t)}$  is the t-th assignment to  $\mathbf{z}$ 

To obtain a new state we

- 1. start a draft state  $\mathbf{z} = \mathbf{z}^{(t-1)}$
- 2. repeat for  $i = 1, \ldots, N$ 
  - ▶ resample Z<sub>i</sub> ~ P(z<sub>i</sub>|x<sub>-i</sub>, z<sub>-i</sub>) only variables in the Markov blanket of z<sub>i</sub> play a role that's why this is feasible
- 3. after complete pass over the data we have a new state  $\mathbf{z}^{(t)}$

We want to sample from  $P(\mathbf{z}|\mathbf{x})$  with a Markov chain a state  $y_t = \mathbf{z}^{(t)}$  is the t-th assignment to  $\mathbf{z}$ 

To obtain a new state we

- 1. start a draft state  $\mathbf{z} = \mathbf{z}^{(t-1)}$
- 2. repeat for  $i = 1, \ldots, N$ 
  - ▶ resample Z<sub>i</sub> ~ P(z<sub>i</sub>|x<sub>-i</sub>, z<sub>-i</sub>) only variables in the Markov blanket of z<sub>i</sub> play a role that's why this is feasible
- 3. after complete pass over the data we have a new state  $\mathbf{z}^{(t)}$

When we have collected a large number T of samples

we can summarise the distribution and/or make decisions

### More in NLP

- Topic modelling (Blei et al. 2003)
- Unsupervised POS tagging (Goldwater and Griffiths 2007)
- PCFGs (Johnson et al. 2007)
- ▶ IBM model 1 and 2 (Mermer and Saraclar 2011)
- ▶ Word alignments without NULL words (Schulz et al. 2016)
- HMM alignments (Schulz and Aziz 2016)

Friends don't let friends optimise

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification
- Bayesian modelling is about uncertainty quantification

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification
- Bayesian modelling is about uncertainty quantification
- Bayesians compare models (a set of assumptions) not point estimates

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification
- Bayesian modelling is about uncertainty quantification
- Bayesians compare models (a set of assumptions) not point estimates
- Comparing Bayesian models is easier

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification
- Bayesian modelling is about uncertainty quantification
- Bayesians compare models (a set of assumptions) not point estimates
- Comparing Bayesian models is easier
- Bayesian modelling requires some maths ;)

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification
- Bayesian modelling is about uncertainty quantification
- Bayesians compare models (a set of assumptions) not point estimates
- Comparing Bayesian models is easier
- Bayesian modelling requires some maths ;)
- Some families enjoy analytically available posteriors

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification
- Bayesian modelling is about uncertainty quantification
- Bayesians compare models (a set of assumptions) not point estimates
- Comparing Bayesian models is easier
- Bayesian modelling requires some maths ;)
- Some families enjoy analytically available posteriors
- Inference can be done by simulation (MCMC)

- Friends don't let friends optimise
- Bayesian modelling is not only about prior specification
- Bayesian modelling is about uncertainty quantification
- Bayesians compare models (a set of assumptions) not point estimates
- Comparing Bayesian models is easier
- Bayesian modelling requires some maths ;)
- Some families enjoy analytically available posteriors
- Inference can be done by simulation (MCMC)

# Beyond

For more on latent variable modelling, especially with structured data

- ▶ take NLP2
- though most of it will be *frequentist* (but for very good reasons!)

For more on Bayesian modelling, approximate inference, and probabilistic modelling with neural networks

- take DL4NLP
- though MCMC will not be the method of choice, instead we will look into variational inference
- and we will need to count on optimisation =0
- though with a nice twist ;)

#### References I

- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. J. Mach. Learn. Res., 3:993–1022, March 2003. ISSN 1532-4435. URL http://dl.acm.org/citation.cfm?id=944919.944937.
- Sharon Goldwater and Tom Griffiths. A fully bayesian approach to unsupervised part-of-speech tagging. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 744–751, Prague, Czech Republic, June 2007. Association for Computational Linguistics. URL

http://www.aclweb.org/anthology/P07-1094.

### References II

Mark Johnson, Thomas Griffiths, and Sharon Goldwater. Bayesian inference for PCFGs via Markov chain Monte Carlo. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pages 139–146, Rochester, New York, April 2007. Association for Computational Linguistics. URL http://www.aclweb.org/anthology/N/N07/N07-1018.

Coskun Mermer and Murat Saraclar. Bayesian word alignment for statistical machine translation. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 182–187, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL http://www.aclweb.org/anthology/P11-2032.

### References III

Philip Schulz and Wilker Aziz. Fast collocation-based bayesian hmm word alignment. In Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, pages 3146–3155, Osaka, Japan, December 2016. The COLING 2016 Organizing Committee. URL http://aclweb.org/anthology/C16-1296.

Philip Schulz, Wilker Aziz, and Khalil Sima'an. Word alignment without null words. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 169–174, Berlin, Germany, August 2016. Association for Computational Linguistics. URL http://anthology.aclweb.org/P16-2028.