# Non-Parametric Bayesian Models

**Computational Linguistics** 

Alexander Koller

21 January 2020

with help from Christoph Teichmann and illustrations by Martín Villalba



# **Topic models**



learn: word probs. ← given: raw documents → learn: topic mixture for (abstract) *topics* in each document

# Examples



topic mixture for one article in *Science*   $\begin{array}{l} 15 \text{ words with highest } \varphi_{k,w} \\ \text{for each topic over whole corpus} \\ \text{(with made-up topic label)} \end{array}$ 

#### Examples

#### development of topics from Science over time (1880-2002)



#### Last time

Say you come across some people who have been stabbed or poisoned. You know that each of them was killed by a pirate or a ninja. You can tell how each person died, but not by whom they were killed.



## **Generative story**

• We assume deaths are generated as follows:

 $\begin{aligned} &(\theta_{pi}, \theta_{ni}) \sim \text{Dir}(\alpha, \alpha) \\ &(\phi_{st|pi}, \phi_{po|pi}), (\phi_{st|ni}, \phi_{po|ni}) \sim \text{Dir}(\beta, \beta) \\ &z_1, \dots, z_K \sim \text{Categorical}(\theta) \\ &w_i \sim \text{Categorical}(\phi_{zi}) \end{aligned}$ 

- That is:
  - $P(z_i = pi) = \theta_{pi}$ ,  $P(z_i = ni) = \theta_{ni}$
  - if  $z_i$  came out as "pi", then  $P(w_i = st) = \phi_{st|pi}$

I abbreviate  $\theta = (\theta_{pi}, \theta_{ni}), \phi_{pi} = (\phi_{st|pi}, \phi_{po|pi}), \phi_{ni} = (\phi_{st|ni}, \phi_{po|ni}).$  $\alpha, \beta$  are assumed given and are called *hyperparameters*.



# Supervised learning

If all killers are known,  $P(M \mid D)$  is easy to compute.



$$\begin{split} P(M) &= \operatorname{Dir}_{\alpha,\alpha}(\theta) \cdot \operatorname{Dir}_{\beta,\beta}(\phi_{\mathrm{pi}}) \cdot \operatorname{Dir}_{\beta,\beta}(\phi_{\mathrm{ni}}) \\ &\propto \theta_{\mathrm{pi}}^{\alpha-1} \cdot \theta_{\mathrm{ni}}^{\alpha-1} \cdot \phi_{\mathrm{st}|\mathrm{pi}}^{\beta-1} \cdot \phi_{\mathrm{po}|\mathrm{pi}}^{\beta-1} \cdot \phi_{\mathrm{st}|\mathrm{ni}}^{\beta-1} \cdot \phi_{\mathrm{po}|\mathrm{ni}}^{\beta-1} \\ P(D \mid M) &= P(z_1 = \mathrm{pi}, w_1 = \mathrm{st}, z_2 = \mathrm{ni}, w_2 = \mathrm{po}) \\ &= \theta_{\mathrm{pi}} \cdot \phi_{\mathrm{st}|\mathrm{pi}} \cdot \theta_{\mathrm{ni}} \cdot \phi_{\mathrm{po}|\mathrm{ni}} \\ P(M \mid D) \propto P(D \mid M) \cdot P(M) \\ &\propto \theta_{\mathrm{pi}}^{\alpha} \cdot \theta_{\mathrm{ni}}^{\alpha} \cdot \phi_{\mathrm{st}|\mathrm{pi}}^{\beta-1} \cdot \phi_{\mathrm{po}|\mathrm{pi}}^{\beta-1} \cdot \phi_{\mathrm{po}|\mathrm{ni}}^{\beta-1} \\ &\propto \operatorname{Dir}_{\alpha+1,\alpha+1}(\theta) \cdot \operatorname{Dir}_{\beta+1,\beta}(\phi_{\mathrm{pi}}) \cdot \operatorname{Dir}_{\beta,\beta+1}(\phi_{\mathrm{ni}}) \end{split}$$



# **Gibbs Sampling**

• Gibbs sampling is MCMC method for computing expected values under posterior distribution.



# Let's simplify

• To bring out today's point more clearly, consider a parade of observable pirates and ninjas (unlike last time, where they were latent):

$$\theta = (\theta_{pi}, \theta_{ni}) \sim Dir(\alpha_{pi}, \alpha_{ni})$$
  
 $z_1, ..., z_N \sim Categorical(\theta)$ 



• Posterior after observations  $D = z_1, ..., z_N$ :  $P(\theta \mid D) = \text{Dir}_{\alpha_{pi}+n_{pi},\alpha_{ni}+n_{ni}}(\theta_{pi}, \theta_{ni})$ 

#### **Predictive distributions**



## **Predictive distribution**

• We can determine the predictive distribution by marginalizing over the model:

$$P(z_{i} = pi \mid z_{1}, \dots, z_{i-1}) = \int P(z_{i} = pi \mid M, z_{1}, \dots, z_{i-1}) \cdot P(M \mid z_{1}, \dots, z_{i-1}) dM$$
$$= \int \theta_{pi} \cdot P(M \mid z_{1}, \dots, z_{i-1}) dM$$
$$= \dots$$
$$= \frac{\alpha_{pi} + n_{pi}^{-i}}{\alpha_{pi} + \alpha_{ni} + (i-1)}$$
# pirates in 1, ..., i-1

• Or equivalently, with  $\alpha = \alpha_{pi} + \alpha_{ni}$  and  $\pi_{pi} = \alpha_{pi} / \alpha$ :

$$P(z_i = pi \mid z_1, \dots, z_{i-1}) = \frac{n_{pi}^{-i}}{\alpha + (i-1)} + \frac{\alpha \cdot \pi_{pi}}{\alpha + (i-1)}$$

## **Teichmann Pocket Process**

- Illustrate this distribution as follows:
  - left pocket contains  $\alpha_{pi}$  pirates,  $\alpha_{ni}$  ninjas
  - right pocket contains α jokers
  - randomly draw a card from right pocket
  - if it is pirate or ninja, output that guy and put him *and a clone of him* in right pocket
  - if it is joker, randomly draw a guy X from left pocket;
    output X, put him back in left pocket, and put clone in right
- Officially called "Polya urn scheme", but I prefer Christoph Teichmann's pocket metaphor.

#### **Teichmann Pocket Process** $P(z_i = pi \mid z_1, ..., z_{i-1}) = \frac{n_{pi}^{-i}}{\alpha + (i-1)} + \frac{\alpha}{\alpha + (i-1)} \cdot \pi_{pi}$

![](_page_12_Picture_1.jpeg)

 $\mathbf{Z}_1$ 

 $\mathbf{Z}_2$ 

## Non-parametric models

- Key limitation of models so far: must specify number K of topics / of killer types.
- Will now generalize this to a class of Bayesian models that automatically pick as many killer types as needed to fit the data.
- Called *non-parametric models* because number of parameters not fixed in advance like in  $(\theta_1, ..., \theta_K)$ .

## Non-parametric models

- Idea: prob dist over infinite space of events.
  - assume some *base distribution* G over these events
  - add Polya-urn-style caching model on top of it
- Can simply adapt predictive distribution:

$$P(z_i = k \mid z_1, \dots, z_{i-1}) = \frac{n_k^{-i}}{\alpha + (i-1)} + \frac{\alpha}{\alpha + (i-1)} \cdot G(k)$$

• Earlier pirate-ninja distribution:  $G(pi) = \pi_{pi}, G(ni) = \pi_{ni} \text{ with } \pi_{pi} + \pi_{ni} = 1.$ 

#### **Teichmann Pocket Process** $P(z_i = k \mid z_1, \dots, z_{i-1}) = \frac{n_k^{-i}}{\alpha + (i-1)} + \frac{\alpha}{\alpha + (i-1)} \cdot G(k)$

![](_page_15_Picture_1.jpeg)

## **Chinese Restaurant Process**

- Alternative illustration (very popular in literature):
  - Chinese restaurant with infinite sequence of tables, each of which has infinite seating capacity.
  - With probability  $\frac{n_k^{-i}}{\alpha + (i-1)}$ , customer i chooses to sit at table k (which has  $n_k^{-i}$  other people sitting at it).
  - With probability  $\frac{\alpha}{\alpha + (i-1)}$  customer opens up new table, and label L for new table drawn at random from G(L).

![](_page_16_Figure_5.jpeg)

# Exchangeability

- The CRP is *exchangeable*: when computing
  P(z<sub>i</sub> | z<sup>-i</sup>), can pretend that z<sub>i</sub> is last event (and tables have all customers from z<sup>-i</sup> on them).
- Can therefore use predictive prob in two ways:
  - as predictive prob, to predict next unseen event
  - in Gibbs sampling, to resample  $z_i$  based on the others
- De Finetti's theorem: Exchangeable observations are independent given some latent variables.

 $\mathbf{Z}_2$ 

 $\mathbf{Z}_{\mathbf{N}}$ 

 for CRP, distribution over latent variables is *Dirichlet process*

## **Grammar induction**

- Tree substitution grammar (TSG) is a grammar formalism in which *elementary trees* are combined using the *substitution* operation.
- In the Penn Treebank, we can only observe the derived trees.
  - Unclear how they were constructed from elementary trees.
- How can we induce a (probabilistic) TSG grammar from the Penn Treebank trees?

## **PTSG: Example**

![](_page_19_Figure_1.jpeg)

#### **TSG induction**

![](_page_20_Figure_1.jpeg)

# Why not EM?

- Our default method for learning with latent variables L so far: Expectation Maximization (EM).
- EM tries to find maximum-likelihood estimate:

$$\max_{p} P(D \mid p)$$
$$= \max_{p} \sum_{L} P(L) \cdot P(D \mid L, p)$$

- This does not work for grammar induction: Max-likelihood estimate makes each derived tree a single elementary tree.
  - Need prior on L to avoid this: What does a "reasonable" grammar look like?

#### **Generative story**

![](_page_22_Figure_1.jpeg)

In full model (Cohn et al. 2010), there is a separate restaurant for each root symbol. Also uses extension of CRP called the *Pitman-Yor Process*.

## **Base distribution for e-trees**

• Prob dist over infinite set of e-trees; prob decays exponentially with size of e-tree (hence sum to 1).

![](_page_23_Figure_2.jpeg)

![](_page_23_Picture_3.jpeg)

• Also need to specify distribution P<sub>C</sub> over cf. rules.

## **Gibbs Sampling for TSGs**

Gibbs state: nodes marked as 1 (substitution site) or 0 (inside e-tree)

![](_page_24_Figure_2.jpeg)

## **Grammar induction**

- Using Gibbs sampler, we can sample from posterior, given PTB as observations.
- We want to learn how to parse new strings with PTSG. Can do this in various ways, e.g.:
  - estimate expected values of PTSG parameters  $\theta$  with Gibbs
  - include sentence to be parsed in Gibbs sampling and return most frequent tree (MPT, MPD in Cohn et al.)

#### Results

	$\leq 40$		all	
Parser	<b>F1</b>	EX	<b>F1</b>	EX
MLE PCFG	64.2	7.2	63.1	6.7
TSG PYP Viterbi	83.6	24.6	82.7	22.9
TSG PYP MPD	84.2	27.2	83.3	25.4
TSG PYP MPT	84.7	28.0	83.8	26.2
TSG PYP MER	85.4	27.2	84.7	25.8
DOP (Zuidema, 2007)			83.8	26.9
Berkeley parser (Petrov and Klein, 2007)	90.6		90.0	
Berkeley parser (restricted)	87.3	31.0	86.6	29.0
Reranking parser (Charniak and Johnson, 2005)	92.0		91.4	
Shindo et al., 2012 (single)	91.6		91.1	
Shindo et al., 2012 (multiple)	92.9		92.4	

## Conclusion

- Predictive probabilities:
  - integrate posterior distribution over models
  - yields intuitive stochastic processes (Polya)
- Extend to non-parametric models:
  - distributions over infinite domains with caching
  - Chinese Restaurant Process, Pitman-Yor
- Apply to grammar induction, e.g. for TSGs.