

Algorithms For Data Science: Gibbs Sampling

Héctor Corrada Bravo

University of Maryland, College Park, USA CMSC 644: 2020-05-04



Documents as *mixtures* of topics (Hoffman 1999)



46

We have a set of documents D

Each document modeled as a bag-of-words (bow) over dictionary W.

 $x_{w,d}$: the number of times word $w \in W$ appears in document $d \in D$.

Let's start with a simple model based on the frequency of word occurrences.

Each document is modeled as n_d draws from a *Multinomial* distribution with parameters $\theta_d = \{\theta_{1,d}, \dots, \theta_{W,d}\}$

Note $heta_{w,d} \geq 0$ and $\sum_w heta_{w,d} = 1$.

Probability of observed corpus D

$$Pr(D|\{ heta_d\}) \propto \prod_{d=1}^D \prod_{w=1}^W heta_{w,d}^{x_{w,d}}$$

Problem 1:

Prove MLE
$$\hat{ heta}_{w,d} = rac{x_{w,d}}{n_d}$$

We have a problem of type

$$egin{array}{lll} \min_x & f_0(x) \ ext{s.t.} & f_i(x) \leq 0 \; i=1,\ldots,m \ & h_i(x)=0 \; i=1,\ldots,p \end{array}$$

Note: This discussion follows Boyd and Vandenberghe, *Convex Optimization*

To solve these type of problems we will look at the *Lagrangian* function:

$$L(x,\lambda,
u)=f_0(x)+\sum_{i=1}^m\lambda_if_i(x)+\sum_{i=1}^p
u_ig_i(x)$$

We'll see these in more detail later, but there is a beautiful result giving *optimality conditions* based on the Lagrangian:

Suppose \tilde{x} , $\tilde{\lambda}$ and $\tilde{\nu}$ are *optimal*, then

$$egin{aligned} f_i(ilde{x}) &\leq 0 \ h_i(ilde{x}) &= 0 \ & ilde{\lambda}_i &\geq 0 \ & ilde{\lambda}_i f_i(ilde{x}) &= 0 \ & ilde{
aligned} &
aligned
aligned$$

We can use the gradient and feasibility conditions to prove the MLE result.

Let's change our document model to introduce topics.

The key idea is that the probability of observing a *word* in a *document* is given by two pieces:

- The probability of observing a *topic* in a document, and
- The probability of observing a *word* given a *topic*

$$Pr(w,d) = \sum_{t=1}^T Pr(w|t) Pr(t|d)$$

So, we rewrite corpus probability as

$$Pr(D|\{p_d\}\{ heta_t\}) \propto \prod_{d=1}^D \prod_{w=1}^W \left(\sum_{t=1}^T p_{t,d} heta_{w,t}
ight)^{x_{w,d}}$$

So, we rewrite corpus probability as

$$Pr(D|\{p_d\}\{ heta_t\}) \propto \prod_{d=1}^D \prod_{w=1}^W \left(\sum_{t=1}^T p_{t,d} heta_{w,t}
ight)^{x_{w,d}}$$

Mixture of topics!!

A fully observed model

Assume you know the *latent* number of occurences of word w in document d generated from topic t:

$$\Delta_{w,d,t}$$
 , such that $\sum_t \Delta_{w,d,t} = x_{w,d}$.

In that case we can rewrite corpus probability:

$$Pr(D|\{p_d\},\{ heta_t\}) \propto \prod_{d=1}^D \prod_{w=1}^W \prod_{t=1}^T (p_{t,d} heta_{w,t})^{\Delta_{w,d,t}}$$

Problem 2 Show MLEs given by

$$\hat{p}_{t,d} = \frac{\sum_{w=1}^{W} \Delta_{w,d,t}}{\sum_{t=1}^{T} \sum_{w=1}^{W} \Delta_{w,d,t}}$$
$$\hat{\theta}_{t,d} = \frac{\sum_{d=1}^{D} \Delta_{w,d,t}}{\sum_{w=1}^{W} \sum_{d=1}^{D} \Delta_{w,d,t}}$$

Since we don't observe $\Delta_{w,d,t}$ we use the EM algorithm

At each iteration (given current parameters $\{p_d\}$ and $\{\theta_d\}$ find *responsibility*

$$\gamma_{w,d,t} = E[\Delta_{w,d,t}|\{p_d\},\{ heta_t\}]$$

and maximize fully observed likelihood plugging in $\gamma_{w,d,t}$ for $\Delta_{w,d,t}$

Problem 4: Show

$$\gamma_{w,d,t} = x_{w,d} imes rac{p_{t,d} heta_{w,t}}{\sum_{t'=1}^T p_{t',d} heta_{w,t'}}$$

Ultimately, what we are interested in is learning topics

Perhaps instead of finding parameters θ that maximize likelihood

Sample from a distribution $Pr(\theta|D)$ that gives us topic estimates

Ultimately, what we are interested in is learning topics

Perhaps instead of finding parameters θ that maximize likelihood

Sample from a distribution $Pr(\theta|D)$ that gives us topic estimates

But, we only have talked about $Pr(D|\theta)$ how can we sample parameters?

Like EM, the trick here is to expand model with *latent* data Z^m

And sample from distribution $Pr(\theta, Z^m | Z)$

Like EM, the trick here is to expand model with *latent* data Z^m

And sample from distribution $Pr(heta, Z^m | Z)$

This is challenging, but sampling from $Pr(\theta|Z^m,Z)$ and $Pr(Z^m|\theta,Z)$ is easier

The Gibbs Sampler does exactly that

Property: After some rounds, samples from the conditional distributions $Pr(\theta|Z^m,Z)$

Correspond to samples from marginal $Pr(heta|Z) = \sum_{Z^m} Pr(heta, Z^m|Z)$

Quick aside, how to simulate data for pLSA?

- Generate parameters $\{p_d\}$ and $\{\theta_t\}$
- Generate $\Delta_{w,d,t}$

Let's go backwards, let's deal with $\Delta_{w,d,t}$

Let's go backwards, let's deal with $\Delta_{w,d,t}$

$$\Delta_{w,d,t} \sim \mathrm{Mult}_{\mathrm{x}_{\mathrm{w},\mathrm{d}}}(\gamma_{w,d,1},\ldots,\gamma_{w,d,T})$$

Where $\gamma_{w,d,t}$ was as given by E-step

Let's go backwards, let's deal with $\Delta_{w,d,t}$

$$\Delta_{w,d,t} \sim \operatorname{Mult}_{\mathrm{x}_{\mathrm{w},\mathrm{d}}}(\gamma_{w,d,1},\ldots,\gamma_{w,d,T})$$

Where $\gamma_{w,d,t}$ was as given by E-step

```
for d in range(num_docs):
    delta[d,w,:] = np.random.multinomial(doc_mat[d,w],
        gamma[d,w,:])
```

Hmm, that's a problem since we need $x_{w,d}$...

But, we know $Pr(w,d) = \sum_t p_{t,d} heta_{w,t}$ so, let's use that to generate each $x_{w,d}$ as

$$x_{w,d} \sim \operatorname{Mult}_{n_d}(Pr(1,d),\ldots,Pr(W,d))$$

Hmm, that's a problem since we need $x_{w,d}$...

But, we know $Pr(w,d) = \sum_t p_{t,d} \theta_{w,t}$ so, let's use that to generate each $x_{w,d}$ as

$$x_{w,d} \sim \operatorname{Mult}_{n_d}(Pr(1,d),\ldots,Pr(W,d))$$

for d in range(num_docs):

doc_mat[d,:] = np.random.multinomial(nw[d], np.sum(p[:,d] * theta), axis=0)

Now, how about p_d ? How do we generate the parameters of a Multinomial distribution?

Now, how about p_d ? How do we generate the parameters of a Multinomial distribution?

This is where the Dirichlet distribution comes in...

If $p_d \sim \operatorname{Dir}(lpha)$, then

$$Pr(p_d) \propto \prod_{t=1}^T p_{t,d}^{lpha_t-1}$$

Some interesting properties:

$$E[p_{t,d}] = rac{lpha_t}{\sum_{t'} lpha_{t'}}$$

So, if we set all $\alpha_t = 1$ we will tend to have uniform probability over topics (1/t each on average)

If we increase $\alpha_t = 100$ it will also have uniform probability but will have very little variance (it will almost always be 1/t)

Approximate Inference by Sampling So, we can say $p_d \sim {
m Dir}(lpha)$ and $heta_t \sim {
m Dir}(eta)$

So, we can say $p_d \sim \mathrm{Dir}(lpha)$ and $heta_t \sim \mathrm{Dir}(eta)$

And generate data as (with $\alpha_t = 1$)

```
for d in range(num_docs):
```

p[:,d] = np.random.dirichlet(1. * np.ones(num_topics))

So what we have is a *prior* over parameters $\{p_d\}$ and $\{\theta_t\}$: $Pr(p_d|\alpha)$ and $Pr(\theta_t|\beta)$

And we can formulate a distribution for missing data $\Delta_{w,d,t}$:

$$Pr(\Delta_{w,d,t}|p_d, heta_t,lpha,eta) =
onumber \ Pr(\Delta_{w,d,t}|p_d, heta_t)Pr(p_d|lpha)Pr(heta_t|eta)$$

However, what we care about is the *posterior* distribution $Pr(p_d|\Delta_{w,d,t}, \theta_t, \alpha, \beta)$

What do we do???

Another neat property of the Dirichlet distribution is that it is *conjugate* to the Multinomial

If $heta|lpha\sim {
m Dir}(lpha)$ and $X| heta\sim {
m Multinomial}(heta)$, then

```
	heta|X, lpha \sim \mathrm{Dir}(X+lpha)
```

That means we can sample p_d from

$$p_{t,d} \sim \mathrm{Dir}(\sum_w \Delta_{w,d,t} + lpha)$$

and

$$heta_{w,t} \sim \mathrm{Dir}(\sum_d \Delta_{w,d,t} + eta)$$

Coincidentally, we have just specified the **Latent Dirichlet Allocation** method for topic modeling.

This is the most commonly used method for topic modeling



Blei, Ng, Jordan (2003), JMLR 37 / 46

We can now specify a full Gibbs Sampler for an LDA mixture model.

Given:

- Word-document counts $x_{w,d}$
- Number of topics K
- Prior parameters α and β

Do: Learn parameters $\{p_d\}$ and $\{\theta_t\}$ for K topics

Step 0: Initialize parameters $\{p_d\}$ and $\{\theta_t\}$

 $p_d \sim \mathrm{Dir}(lpha)$

and

 $heta_t \sim {
m Dir}(eta)$

Step 1:

Sample $\Delta_{w,d,t}$ based on current parameters $\{p_d\}$ and $\{\theta_t\}$

$$\Delta_{w,d,.} \sim \operatorname{Mult}_{x_{w,d}}(\gamma_{w,d,1},\ldots,\gamma_{w,d,T})$$

Step 2:

Sample parameters from

$$p_{t,d} \sim \mathrm{Dir}(\sum_w \Delta_{w,d,t} + lpha)$$

and

$$heta_{w,t} \sim \mathrm{Dir}(\sum_d \Delta_{w,d,t} + eta)$$

Step 3:

Get samples for a few iterations (e.g., 200), we want to reach a stationary distribution...

Step 4:

Estimate $\hat{\Delta}_{w,d,t}$ as the average of the estimates from the last m iterations (e.g., m=500)

Step 5:

Estimate parameters p_d and $heta_t$ based on estimated $\hat{\Delta}_{w,d,t}$

$$\hat{p}_{t,d} = \frac{\sum_{w} \hat{\Delta}_{w,d,t} + \alpha}{\sum_{t} \sum_{w} \hat{\Delta}_{w,d,t} + \alpha}$$
$$\hat{\theta}_{w,t} = \frac{\sum_{d} \hat{\Delta}_{w,d,t} + \beta}{\sum_{w} \sum_{d} \hat{\Delta}_{w,d,t} + \beta}$$

44 / 46

Mixture models

We have now seen two different mixture models: soft k-means and topic models

Mixture models

We have now seen two different mixture models: soft k-means and topic models

Two inference procedures:

- Exact Inference with Maximum Likelihood using the EM algorithm
- Approximate Inference using Gibbs Sampling