# An Empirical Bayes Approach to Topic Modeling

#### Anirban Gangopadhyay<sup>1</sup>

<sup>1</sup>Department of Computer Science Columbia University

ICPR 2020

Anirban Gangopadhyay

An Empirical Bayes Approach to Topic Model

ICPR 2020 1 / 47

## Introduction

### 2 Hierarchical Bayesian Models

- Generative Process
- Inference
- Statistical Assumptions
- 3 Graph Clustering
  - Algorithm & Generative Process
  - Inference
  - Statistical Assumptions
- 4 Factor Analysis
  - Algorithmic and Generative Process
  - Statistical Assumptions
  - Inference
- 5 Empirical Bayes Approach
- 6 Applications & Results

- Consider the problem of modeling text corpora
  - Finding short descriptions of members of a collection (ie derive topics, communities, clusters, etc)
  - Scoring similarities across members (ie finding similar documents, topics)
- Many disparate models exist
- Derived from different mathematical disciplines that make different assumptions about data population

- We give an Empirical Bayes framework for choosing optimal algorithm
- Map each algorithm to a graphical model
- Choose optimal model given the observed variables in our data (based on the assumptions each model makes)

## Introduction

## 2 Hierarchical Bayesian Models

#### Generative Process

- Inference
- Statistical Assumptions
- 3 Graph Clustering
  - Algorithm & Generative Process
  - Inference
  - Statistical Assumptions
- Factor Analysis
  - Algorithmic and Generative Process
  - Statistical Assumptions
  - Inference
- Empirical Bayes Approach
- Applications & Results

We outline the generative process for each document  $\overrightarrow{w}$  in a corpus D:

- Choose  $N \sim Poisson(\xi)$
- 2 Choose  $\Theta \sim Dir(\alpha)$
- **③** For each of the *N* words  $w_n$ :
  - Choose a topic  $z_n \sim Mult(\Theta)$
  - **2** Choose a word  $w_n$  from  $p(w_n|z_n, \beta)$ , a multinomial probability conditioned on topic  $z_n$ .



・ロト ・日 ・ ・ ヨ ・ ・

- Joint distribution:  $p(\theta, z, \vec{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^{N} p(z_n | \theta) p(w_n | z_n, \beta).$
- To obtain the marginal distribution of a document, we integrate over  $\theta$  and sum over z:
- $p(\overrightarrow{w}|\alpha,\beta) = \int p(\theta|\alpha) \cdot (\prod_{n=1}^{N} \sum_{z_n} p(z_n|\theta) \cdot p(w_n|z_n,\beta)) d\theta.$
- Since we assume documents are i.i.d., we obtain the probability of the corpus by taking  $\prod_{d=1}^{M} p(\vec{w} | \alpha, \beta)$ .
- we wish to perform a maximum likelihood of
- $p(D|\alpha,\beta) = \prod_{i=1}^{M} \int p(\theta_d|\alpha) \cdot (\prod_{n=1}^{N_d} \sum_{z_{d_n}} p(z_{d_n}|\theta_d) \cdot p(w_{d_n}|z_{d_n},\beta)) d\theta_d$

## Introduction

### 2 Hierarchical Bayesian Models

Generative Process

#### Inference

• Statistical Assumptions

## 3 Graph Clustering

- Algorithm & Generative Process
- Inference
- Statistical Assumptions

## Factor Analysis

- Algorithmic and Generative Process
- Statistical Assumptions
- Inference
- Empirical Bayes Approach
- Applications & Results

• Posterior distribution given a document is intractable to compute:

• 
$$p(\theta, z | w, \alpha, \beta) = \frac{p(\theta, z, w | \alpha, \beta)}{p(w | \alpha, \beta)}.$$

- We consider variational inference as an approximate inference method
- Variational Inference Derive parameters of a tractable model that has minimal KL divergence to LDA model

• Using variational inference, we obtain a family of distributions on the latent variables  $\gamma,\phi_1,...,\phi_m$ 

• 
$$q(\theta, z|\gamma, \phi) = q(\theta|\gamma) \cdot \prod_{n=1}^{N} q(z_n|\phi_n)$$

- Dirichlet parameter  $\gamma$ , multinomial parameters  $\phi_1, ..., \phi_m$  are free variational parameters
- The family of distributions are a family of lower bounds derived via Jensen's inequality.



э **ICPR 2020** 12/47

-

• • • • • • • • • • • •

Continued

- We drop the edges between  $\theta, z, w$  and the w nodes in the above model
- To find the optimal  $\gamma, \phi$ , we minimize the KL divergence
- $(\gamma^*, \phi^*) = \operatorname{argmin}_{\gamma, \phi} D(q(\theta, z | \gamma, \phi) || p(\theta, z | w, \alpha, \beta))$
- $\gamma, \phi$  are conditioned on  $\overrightarrow{w}$  which implies that:
- $(\gamma^*, \phi^*) = \operatorname{argmin}_{\gamma, \phi} D(q(\theta, z | \overrightarrow{w}(\gamma), \overrightarrow{w}(\phi)) || p(\theta, z | \overrightarrow{w}, \alpha, \beta))$

Since  $\alpha,\beta$  are hidden, we can use Expectation-Maximization to derive  $\gamma^*,\phi^*.$ 

- (E-step) For each document, find the optimizing values of the variational parameters  $\{\gamma_d^*, \phi_d^* : d \in D\}$
- **(**M-step) Maximize the resulting lower bound on the log likelihood w.r.t.  $\alpha, \beta$  (corresponds to finding the MLE for each document under approximate posterior)

## Introduction

## 2 Hierarchical Bayesian Models

- Generative Process
- Inference

#### Statistical Assumptions

## 3 Graph Clustering

- Algorithm & Generative Process
- Inference
- Statistical Assumptions

#### Factor Analysis

- Algorithmic and Generative Process
- Statistical Assumptions
- Inference
- Empirical Bayes Approach
- Applications & Results

- **Graphical Model**: LDA allows each document to be generated from a distribution of topics whereas the two level Dirichlet-Multinomial clustering model would restrict a document to be associated with a single topic.
- **Exchangeability**: DeFinetti's theorem states that the joint distribution of an infinitely exchangeable sequence of random variables is as if a random parameter were drawn from some distribution and then the random variables in question were i.i.d. conditioned on the parameter.
- Allows us to assume a bag of words model for LDA (words are i.i.d. conditioned on the topic it is generated from)



Μ

・ロト ・ 日 ト ・ ヨ ト ・

▲ ■ ► ■ つへの ICPR 2020 17/4

## Introduction

#### 2 Hierarchical Bayesian Models

- Generative Process
- Inference
- Statistical Assumptions

## 3 Graph Clustering

## • Algorithm & Generative Process

- Inference
- Statistical Assumptions

## Factor Analysis

- Algorithmic and Generative Process
- Statistical Assumptions
- Inference
- Empirical Bayes Approach
- Applications & Results

We outline the algorithmic process for deriving topics (communities) given a corpus:

- Construct a network graph where each node represents a word w<sub>i</sub> in our overall dictionary (constructed from D) and the edge weight e<sub>w<sub>i</sub>,w<sub>j</sub></sub> between two words represents the number of documents w<sub>i</sub>, w<sub>j</sub> appear in together.
- We consider two metrics as a measure of goodness of community structure for the weighted, undirected graph constructed above: modularity Q and modularity density Q<sub>ds</sub>.
  - Modularity:  $Q(G) = \sum_{C_i \in C} \left[\frac{W_{c_i}^{in}}{W} \left(\frac{W_{c_i}}{2W}\right)^2\right]$  where W represents the sum of all edge weights in the entire graph,  $W_{c_i}^{in}$  is the sum of weights of edges within community  $C_i$ ,  $W_{C_i} = 2W_{C_i}^{in} + W_{C_i}^{out}$  where  $W_{C_i}^{out}$  is the sum of weights of edges with exactly one endpoint inside  $C_i$ .

- We wish to maximize Q (or Q<sub>ds</sub>) to obtain the optimal set of communities C = {C<sub>i</sub>}. These communities (set of words) may represent topics.
- A document may be broken down into a set of topics based on the span of words in  $\vec{w}$ .

We outline the generative process (and underlying statistical assumptions) for generating a network graph, its underlying communities and the edges (based on the assumptions made in the case where modularity is maximized).

- Choose the number of nodes N ~ Dist(α) where α is a hyperparameter and Dist(·) is a distribution TBD
- Por i = 1 to n:
  - Choose d<sub>i</sub> ~ Dist2(β) where d<sub>i</sub> represents the degree of node i, β is a hyperparameter, Dist2(·) is TBD.
  - **2** Choose the cluster  $z_n \sim Mult(\Theta)$
- Solution Let us denote A<sub>ij</sub> to be the edge weight between node i and node j. Then A<sub>ij</sub> ∼ Dist3(·) as a function of d<sub>i</sub>, d<sub>j</sub>, z<sub>i</sub>, z<sub>j</sub>

• • = • • = •

- We derive full joint probability, graphical model for the case where all priors are point parameters and only the edges are generative.
- **Goal**: choose  $\operatorname{argmax}_{\overrightarrow{z}} p(A|N, \overrightarrow{d}, \alpha, \beta; \overrightarrow{z})$
- Wish to choose  $p(A|N, \overrightarrow{d}, \alpha, \beta; \overrightarrow{z})$  such that modularity Q is maximized.
- Can be shown that choosing  $\arg\max_z L(\overrightarrow{z})$  is equivalent to choosing  $\arg\max_z Q(\overrightarrow{z})$
- Key insight: By choosing Dist(·), Dist2(·) to be point parameters, Dist3(·) = δ<sub>zi=zj</sub>N(μ<sub>1</sub>, σ) + δ<sub>zi≠zj</sub>N(μ<sub>2</sub>, σ), we have derived a one-layer generative process whose ML parameters will maximize modularity

• • = • • = •



(4) (5) (4) (5)

Image: A matrix

Derive joint probability, graphical model for the case where we presume a generative process for each derived parameter (only hyperparameters are point parameters).

- Choose  $Dist(\cdot)$  appropriately
- Dist2(·) is traditionally Poisson and is a fixed observed variable in the configuration model. We derive the Poisson assumption below based on reasonable base assumptions.
- Dist3(·) was assumed to be a GMM in the single layer model and Poisson(d<sub>i</sub>, d<sub>j</sub>, ω<sub>z<sub>i</sub>,z<sub>j</sub></sub>) in the stochastic block model (which we use for inference).

- Assume the probability of an edge between two nodes is represented by p.
- Furthermore the probability of an edge between two nodes is assumed to be independent of the existence of any other edges in the graph
- So Let z be the average number of edges a given node is connected to. Then  $p = \frac{z}{N-1}$  where N is the total number of nodes.
- Let us denote  $p_k$  to be the probability a given node has degree k. Then  $p_k$  can be derived as follows:  $p_k = {N \choose k} p^k \cdot (1-p)^{N-k} \approx \frac{z^k \cdot e^{-z}}{k!}$ . Note the equality becomes exact in the limit of large N.
- We see  $p_k$  follows a Poisson distribution.

# Graphical Model

Three Layer Model



## Introduction

### 2 Hierarchical Bayesian Models

- Generative Process
- Inference
- Statistical Assumptions

## 3 Graph Clustering

• Algorithm & Generative Process

#### Inference

Statistical Assumptions

#### Factor Analysis

- Algorithmic and Generative Process
- Statistical Assumptions
- Inference
- Empirical Bayes Approach
- Applications & Results

- Consider deriving communities that maximize modularity via a spectral algorithm
- Can show equivalence of this inference method to the single layer model for the case of two communities [Newman]
- **Key Insight:** Can directly compare the single layer model to other derived graphical models when choosing the optimal model for topic modeling.

- We use the degree corrected block model as our derived graphical model for approximate inference



★ Ξ >

- Construct  $L = D^{-1/2}AD^{-1/2}$  where A is the adjacency matrix, D is the diagonal matrix of vertex degrees.
- Derive the eigenvalues and eigenvectors λ<sub>1</sub>, λ<sub>2</sub>, ..., λ<sub>k</sub>, s<sub>1</sub>, ..., s<sub>k</sub> of L (using Lanczos algorithm)
- Consider s<sub>2</sub> and partition s<sub>2</sub> into k groups clustered by distance, for predetermined cluster size k
- $\overrightarrow{z}$  can be represented as the cluster assignment for node *i*

## Introduction

## 2 Hierarchical Bayesian Models

- Generative Process
- Inference
- Statistical Assumptions

## 3 Graph Clustering

- Algorithm & Generative Process
- Inference
- Statistical Assumptions

#### Factor Analysis

- Algorithmic and Generative Process
- Statistical Assumptions
- Inference
- Empirical Bayes Approach
- Applications & Results

- We explore the relevance of the bag of words assumption for our given graph based model
  - Our model makes the assumption that  $d_i$  is i.i.d. (for i = 1 to n). Furthermore  $A_{ij}$  is i.i.d.  $\forall i, j$  conditioned on  $\alpha, \beta, \vec{d}$ .
  - Oi Finetti's theorem can be applied with β being the hyperparameter, w<sub>i</sub>, w<sub>j</sub> co-occuring i.i.d. of other co-occurences
  - $\bigcirc$   $A_{ij}$  can be treated as an infinitely exchangeable sequence
- Hence, we see words don't co-occur in any specific order across our corpus of documents. Therefore we can assume our model follows the bag of words assumption.

## Introduction

## 2 Hierarchical Bayesian Models

- Generative Process
- Inference
- Statistical Assumptions
- 3 Graph Clustering
  - Algorithm & Generative Process
  - Inference
  - Statistical Assumptions
- 4 Factor Analysis
  - Algorithmic and Generative Process
  - Statistical Assumptions
  - Inference
  - Empirical Bayes Approach
  - Applications & Results

- Assume we are given a term document matrix  $D(n \times m)$  where n is the size of the vocabulary and m is the number of documents.
- We assume there exists a matrix D̂ that can be decomposed as such: D̂ = AW where A is a n × r matrix with columns representing the topics and w is the r × m matrix representing topic weights for the given m documents.
- We attempt to find A, W such that  $||D \hat{D}||_F < \epsilon$  for some predefined  $\epsilon$ . We assume A, W are non-negative matrices.
- Hence we recover our topic matrix A and given a new document  $\hat{d}_i$ , A we can recover the topic weights  $\hat{w}_i$  (since  $\hat{d}_i = A \cdot \hat{w}_i$ ).

Let A, n (size of the vocabulary), r (number of topics),  $\alpha$  (Dirichlet parameter) be fixed.

• for 
$$i = 1$$
 to  $m$ :

• 
$$w_i \sim \text{Dir}(\alpha)$$
  
•  $\epsilon_i \sim N(0, \sigma^2)$   
•  $d_i = A \cdot w_i + \epsilon_i$ 

The reason we add noise in our generative process is because the observed term document matrix D may not be factorizable.



Anirban Gangopadhyay

An Empirical Bayes Approach to Topic Model

· ▲ ≣ ト ■ - ∽ へ ICPR 2020 37 / 4

A D N A B N A B N A B N

## Introduction

## 2 Hierarchical Bayesian Models

- Generative Process
- Inference
- Statistical Assumptions
- 3 Graph Clustering
  - Algorithm & Generative Process
  - Inference
  - Statistical Assumptions

## 4 Factor Analysis

- Algorithmic and Generative Process
- Statistical Assumptions
- Inference
- Empirical Bayes Approach
- Applications & Results

- We assume our topic matrix A is p-seperable
- **Seperability** *r* rows out of *n* have only one non-zero value such that if permuted form the *r* × *r* identity matrix (as a submatrix within *A*)
- Discuss how a non-negative matrix A is invertible (relevant for Algorithm, Inference algorithm)
- Discuss how bag of words assumption holds here

## Introduction

## 2 Hierarchical Bayesian Models

- Generative Process
- Inference
- Statistical Assumptions
- 3 Graph Clustering
  - Algorithm & Generative Process
  - Inference
  - Statistical Assumptions

## 4 Factor Analysis

- Algorithmic and Generative Process
- Statistical Assumptions

## Inference

- Empirical Bayes Approach
- Applications & Results

- Recovering non-negative matrices A, W is called NMF and is NP-hard.
- A given observed document  $d_i$  is a noisy approximation to the true document  $A \cdot w_i$  which represents the weighted sum of topics.

## Inference Algorithm Given non-noisy D

- **()** We are given D, for which we assume there exists A, W s.t. D = AW
- 2 Let us consider  $Q = \frac{1}{m} D \cdot D^T \rightarrow Q = \frac{1}{m} AW \cdot W^T A \rightarrow A \cdot \frac{1}{m} WW^T A = AR(\tau) A$
- Note A is p-seperable and Q is a product of two non-negative matrices namely  $A, R(\tau)A^{T}$
- We can rewrite Q = AB where  $B = R(\tau)A^{T}$ .
- Let us assume we figure out the *r* anchor words of *A*. Then a simple transformation yields *Q̂* where we permute the the rows of *A* s.t. the top *r* rows form the identity matrix, we permute the columns of *R*(τ)*A*<sup>T</sup> to match.
- This allows us to derive  $\hat{Q}, R(\hat{\tau})A^T$  which are permuted versions of the original matrices.
- We can hence recover A from these matrices.

• • = • • = •



▲ ■ ト ■ 少 Q ペ
ICPR 2020 43 / 47

Image: A matrix

< ⊒ >

Here we outline how to choose the right algorithm (including the inference process) given the observed variables in our data.

- LDA & variational inference graphical model
- Modularity maximization & spectral clustering graphical model
- Plain LDA graphical model (for theoretical purposes)
- One layer Modularity maximization generative model (for theoretical purposes)
- NMF
- NMF & noise & seperability assumption

- Use model selection framework to evaluate which of the three models is most suitable given a corpus of data D.
- Model evaluation is done on labeled data using a random forest
- Data
  - Twitter
  - 20 newsgroup
  - NIPs

# Applications & Results



イロト イヨト イヨト イ

# Applications & Results



▲ 注 ▶ 注 少 Q C
ICPR 2020 47 / 47

イロト イヨト イヨト イヨ