

Multi-Task Learning of Component Strengths in Non-Conjugate Admixture Models

Bryan Feeney

Read this presentation online at:

<http://tinyurl.com/bryansphd2014>

bryan.feeney@ucl.ac.uk

October , 2014

Research

Three Components

- ① Multi-Task Learning
- ② Admixture Modelling ("Topic Models")
- ③ Local Variational Bounds

Research

Three Components

- ① Multi-Task Learning
- ② Admixture Modelling ("Topic Models")
- ③ Local Variational Bounds

<http://tinyurl.com/bryansphd2014>

Research

Three Components

- ① **Multi-Task Learning**
- ② Admixture Modelling ("Topic Models")
- ③ Local Variational Bounds

Multi-Task Learning

Situation: Making many predictions from the same data

- 1 Predict exam scores in L subjects for several children[12]
- 2 Predict customers affinity to L observed aspects of a product[1]
- 3 Propose image captions by predicting $p(\text{word}|\text{image})$ from image features for L words (in this case $L > 10,000$)[2]

Multi-Task Learning

Situation: Making many predictions from the same data

- 1 Predict exam scores in L subjects for several children[12]
- 2 Predict customers affinity to L observed aspects of a product[1]
- 3 Propose image captions by predicting $p(\text{word}|\text{image})$ from image features for L words (in this case $L > 10,000$)[2]

Problem: Can we improve performance by transferring knowledge between tasks[15]

- 1 Learn correlations between tasks.
- 2 Learn a low-rank projection of the tasks themselves
- 3 Learn structure of features by how we use them - via regularization

Regularization

Learn L vectors w_l . How to *transfer* knowledge from inferring w_1, \dots, w_{l-1} to the task of inferring w_l

Regularization

Learn L vectors w_l . How to *transfer* knowledge from inferring w_1, \dots, w_{l-1} to the task of inferring w_l

$$y_{nl}|w_l \sim \mathcal{N}\left(w_l^\top x_n, \sigma^2 I\right) \quad w_l \sim \mathcal{N}\left(\mathbf{0}, \alpha^2 I\right)$$

Regularization

Learn L vectors w_l . How to *transfer* knowledge from inferring w_1, \dots, w_{l-1} to the task of inferring w_l

Bayesian approach - learn the prior[1]

$$y_{nl}|w_l \sim \mathcal{N}\left(w_l^\top x_n, \sigma^2 I\right)$$

$$m|\Sigma \sim \mathcal{N}\left(m_0, \frac{1}{\lambda}\Sigma\right)$$

$$w_l \sim \mathcal{N}(m, \Sigma)$$

$$\Sigma \sim \mathcal{W}^{-1}(\Sigma_0, \nu)$$

Regularization

Learn L vectors w_l . How to *transfer* knowledge from inferring w_1, \dots, w_{l-1} to the task of inferring w_l

Low-Rank Projections of the Feature Space

$$y_{nl}|w_l \sim \mathcal{N}\left(w_l^\top x_n, \sigma^2 I\right)$$

$$w_l|z_l \sim \mathcal{N}\left(Uz_l + m, \alpha^2 I\right)$$

$$z_l \sim \mathcal{N}(\mathbf{0}, I)$$

$$\Rightarrow w_l \sim \mathcal{N}\left(m, \alpha^2 I + UU^\top\right)$$

Regularization

Learn L vectors w_l . How to *transfer* knowledge from inferring w_1, \dots, w_{l-1} to the task of inferring w_l

Exploit sparsity using ARD priors

$$y_{nl}|w_l \sim \mathcal{N}\left(w_l^\top x_n, \sigma^2 I\right)$$

$$w_l \sim \mathcal{N}(\mathbf{0}, \text{diag}(\alpha))$$

$$\alpha_f \sim \mathcal{G}(a, b)$$

Regularization

Learn L vectors w_l . How to *transfer* knowledge from inferring w_1, \dots, w_{l-1} to the task of inferring w_l

Exploit sparsity using Clustered ARD priors[6]

$$y_{nl}|w_l \sim \mathcal{N}\left(w_l^\top x_n, \sigma^2 I\right)$$

$$z_l \sim \mathcal{M}(\theta, 1)$$

$$\theta \sim \mathcal{D}(\beta)$$

$$w_l|z_l \sim \mathcal{N}(\mathbf{0}, \text{diag}(\alpha_{z_l}))$$

$$\alpha_{fk} \sim \mathcal{G}(a, b)$$

Regularization

Learn L vectors w_I . How to *transfer* knowledge from inferring w_1, \dots, w_{I-1} to the task of inferring w_I

This is all just hierarchical Bayesian modelling

However analogous methods exist in error-optimisation approaches to machine learning which learn regularization functions instead of priors.

- Low-Rank projections in [3]
- Heterogeneous sparsity in [4]

Research

Three Components

- ① Multi-Task Learning
- ② **Admixture Modelling (“Topic Models”)**
- ③ Local Variational Bounds

Mixture Models versus Admixture Models

Classical Mixture Model of Text - One topic per document

$$\boldsymbol{\theta} \sim \mathcal{D}(\boldsymbol{\alpha}) \quad z_d \sim \mathcal{M}(\boldsymbol{\theta}, 1) \quad w_{dn} \sim \mathcal{M}(\boldsymbol{\phi}_{z_d}, 1)$$

Where each of the component vocabularies is drawn $\boldsymbol{\phi}_k \sim \mathcal{D}(\boldsymbol{\beta})$.

Mixture Models versus Admixture Models

Classical Mixture Model of Text - One topic per document

$$\theta \sim \mathcal{D}(\alpha) \quad z_d \sim \mathcal{M}(\theta, 1) \quad w_{dn} \sim \mathcal{M}(\phi_{z_d}, 1)$$

Where each of the component vocabularies is drawn $\phi_k \sim \mathcal{D}(\beta)$.

Mixture models struggle to generalise.

- Fewer clusters mean coarser estimates of cluster centroids
- But more clusters mean fewer datapoints per cluster, and thus sparser estimates of cluster centroids (for text), due to the assumption of one cluster per document

Mixture Models versus Admixture Models

The Secrets of Life and Death

by Rebecca Alexander

Members	Reviews	Popularity	Average rating	Mentions
46	10	251,994	☆☆☆ (3.2)	5

[Add to your library](#) [Add to wishlist](#)

Book Information

Member: Shuffy2

Title: The Secrets of Life and Death

Authors: Rebecca Alexander

Info: Broadway Books (2014), Paperback, 384 pages

Collections: Donated

Rating: ☆☆☆½ ☆ ☆

Tags: Early Reviewers, Fantasy

Work details

The Secrets of Life and Death by Rebecca Alexander

Members [all members](#)

Tags [numbers](#) [show all](#)

death demons Early Reviewers ebook English English
literature ER evil **fantasy** fiction historical fiction history
horror jännitys keskiaika medieval **occult** own police read read in 2014

Quick Links

- [Amazon.com \(direct\)](#)
- [Abebooks.com](#)
- [Amazon Kindle \(2 editions\)](#)
- [Audiible \(1 edition\)](#)
- [CD Audiobook \(0 editions\)](#)
- [Project Gutenberg \(0 editions\)](#)
- [Google Books — Loading...](#)
- [WorldCat](#)

Get this book

- [Local Book Search](#)
- [All sources](#)

Swap Ebooks Audio
— 2 pay 1 pay

Popular covers

(see all 4 covers)

Rating

Average: ☆☆☆ (3.2)

0.5

Mixture Models versus Admixture Models

The screenshot shows the Netflix interface for the TV show "The Thick of It". The browser address bar displays "www.netflix.com/WiMovie/70224481?trkid=13462100". The Netflix logo is at the top left, with "Browse" and "KIDS" links. A search bar and a user profile icon are at the top right. The main content area features a large video player with a "Resume" button. To the right of the player, the show's title "The Thick of It" is displayed, along with its release year (2005-2012), a rating icon, and the number of series (4). Below this, there are five stars and a "Not Interested" button. The text "Our best guess for Susannah: 5 stars" and "Average of 102,406 ratings: 4.4 stars" are shown. A description states: "This award-winning fictitious reality comedy cleverly pokes fun at the intricacies -- and ineptitude -- of the modern British government." At the bottom of the main area, there are buttons for "+ My List" and "Recommend to a friend". A progress bar indicates the user has watched 0m out of 29m, with a note "Series 3: Ep. 4 - You watched this title on 01/10/2014" and a "Report a problem" link. The "Member Reviews" section shows "33 member reviews" and a "Write a review" button. A review by "Susannah" is highlighted, showing a 5-star rating and the text: "The greatest political comedy of the 21st century. Witty, intelligent and prophetic with the most creative swearing ever committed to tape..". It also notes "293 out of 293 members found this review helpful". A sidebar on the right lists genres: "TV Programmes", "British TV Programmes", "TV Comedies", "British TV Comedies", and "Sitcoms". It also states "This programme is Witty, Quirky, Deadpan".

Watch The Thick of It Online

www.netflix.com/WiMovie/70224481?trkid=13462100

NETFLIX Browse KIDS Search

THE THICK OF IT Resume

The Thick of It
2005-2012 4 Series
★★★★★ Not Interested
Our best guess for Susannah: 5 stars
Average of 102,406 ratings: 4.4 stars
This award-winning fictitious reality comedy cleverly pokes fun at the intricacies -- and ineptitude -- of the modern British government.
+ My List Recommend to a friend

0m 29m
Series 3: Ep. 4 - You watched this title on 01/10/2014 Report a problem

Member Reviews 33 member reviews
Reviews voted most helpful
Write a review

★★★★★ The greatest political comedy of the 21st century. Witty, intelligent and prophetic with the most creative swearing ever committed to tape..
293 out of 293 members found this review helpful
You found this review
Helpful Not Helpful

Genres
TV Programmes
British TV Programmes
TV Comedies
British TV Comedies
Sitcoms
This programme is
Witty
Quirky
Deadpan

Mixture Models versus Admixture Models

The screenshot shows a web browser window displaying the Last.fm profile for the band The Strokes. The browser's address bar shows the URL www.last.fm/music/The+Strokes. The Last.fm interface includes a navigation bar with links for Music, Listen, Events, Charts, and Originals, along with Join and Login buttons. The main content area features a large profile picture of the band, a list of tags (Indie, rock, alternative, alternative rock), and a bio stating they are an American rock band from New York City, formed in 1998. To the right, there are statistics for 144,253,646 scrobbles and 2,992,315 listeners, a section for similar artists (Julian Casablancas and Albert Hammond, Jr.), and a list of tracks and albums. The bottom of the page shows a navigation bar with various icons for navigation and search.

The Strokes

Tracks Albums Pictures Videos Events Biography More...

144,253,646 scrobbles 2,992,315 listeners

Similar Artists

Julian Casablancas Albert Hammond, Jr.

Mixture Models versus Admixture Models

Classical Mixture Model of Text - One topic per document

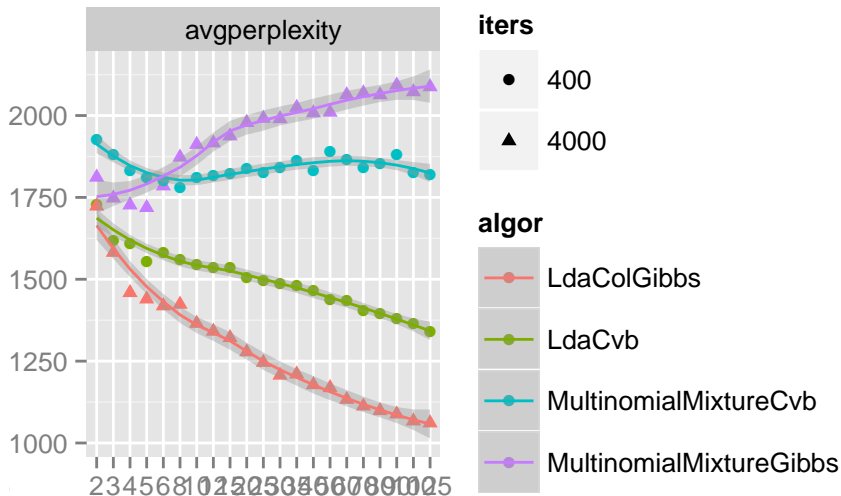
$$\theta \sim \mathcal{D}(\alpha) \quad z_d \sim \mathcal{M}(\theta, 1) \quad w_{dn} \sim \mathcal{M}(\phi_{z_d}, 1)$$

Where each of the component vocabularies is drawn $\phi_k \sim \mathcal{D}(\beta)$.

Admixture models assign a *mixture* of topics to each document, In the case of text the “topic-model” implementation[10] assigns a single topic to each word w_{dn} .

$$\theta_d \sim \mathcal{D}(\alpha) \quad z_{dn} \sim \mathcal{M}(\theta_d, 1) \quad w_{dn} \sim \mathcal{M}(\phi_{z_{dn}}, 1)$$

Mixture Models versus Admixture Models



Topic Models: Multinomial Admixtures

Five broad areas of research

- Richer Observation Models
 - Language Models[41][44][26]
 - Alternative distributions such as logistic-normal[9][17] or von Mises - Fisher[34]

Topic Models: Multinomial Admixtures

Five broad areas of research

- Richer Observation Models
- Alternatives to the Dirichlet prior
 - The “Logistic Normal” prior for the Correlated Topic Model[7]

$$\boldsymbol{\eta}_d \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\theta_{dk} = \sigma_k(\boldsymbol{\theta}_d) = \frac{\exp(\eta_{dk})}{\sum_j \exp(\eta_{dj})}$$

- Same quality of model fit with fewer topics

Topic Models: Multinomial Admixtures

Five broad areas of research

- Richer Observation Models
- Alternatives to the Dirichlet prior
- Bayesian Non-Parametrics to estimate Topic Counts
 - Hierarchical Dirichlet Processes[38]
 - Discrete Infinite Logistic Normal Distribution[32]

Topic Models: Multinomial Admixtures

Five broad areas of research

- Richer Observation Models
- Alternatives to the Dirichlet prior
- Bayesian Non-Parametrics to estimate Topic Counts
- Scalable Implementations:
 - Gibbs Sampling[33] and Collapsed Gibbs Sampling[19]
 - Variational[10] inference and Collapsed Variational[39][21] inference
 - MAP and other approximations[5]
 - Distributed "Big Data" approaches[37][31][16]
 - Optimised online approaches[22][23][29]

Topic Models: Multinomial Admixtures

Five broad areas of research

- Richer Observation Models
- Alternatives to the Dirichlet prior
- Bayesian Non-Parametrics to estimate Topic Counts
- Scalable Implementations:
- Use of Covariates x_d
 - Ad-hoc: Time[43], author[28], region[17]
 - “Downstream” Models[8][35][40]

$$p(w_d, x_d) = \int_{\theta_d} p(w_d | \theta_d) p(x_d | \theta_d) p(\theta_d) \quad (1)$$

- “Upstream” Models:[30]

$$p(w_d | x_d) = \int_{\theta_d} p(w_d | \theta_d) p(\theta_d | x_d) \quad (2)$$

Topic Models: Multinomial Admixtures

Five broad areas of research

- Richer Observation Models
- Alternatives to the Dirichlet prior
- Bayesian Non-Parametrics to estimate Topic Counts
- Scalable Implementations:
- Use of Covariates x_d

LDA also linked to multinomial PCA[14] and Non-Negative Matrix Factorization[18]

Research

Three Components

- ① Multi-Task Learning
- ② Admixture Modelling ("Topic Models")
- ③ **Local Variational Bounds**

Local Variational Bounds

The softmax transformation

$$\eta_d \sim \mathcal{N}(\mu, \Sigma) \qquad z_{dn} \sim \mathcal{M}(\sigma(\eta_d), 1) \qquad (1)$$

Problems

- Prior (a Gaussian) is not conjugate to the likelihood (a multinomial mixture)
- No analytic form for the posterior.

Approach is to approximate the likelihood using bounds

- Classic approach is a global bound over the entire likelihood (e.g. the Laplace approximation)
- A better fit, and more flexibility, can be obtained by using *local* bounds over the likelihood of each document.

Local Variational Bounds

- Global Bounds: Laplace, Delta Method[42]
- Quadratic Bounds: Bohning[11], Bouchard[13]
- Other more recent bounds
 - Piecewise Quadratic Bounds[27]
 - Tilted Bound[25]
 - Stick-breaking bound[24]

Research

Use-Cases: Microtexts

- Tweets: Predict text given user features (username, time tweet was posted). Notably predict absent words such as hashtags.
- Image captions: Predict caption given image features.

Incorporate multi-task learning into an upstream topic model

Partial Model

The Model

$$\eta_d \sim \mathcal{N}(A x_d, \Sigma)$$

$$z_{dn} \sim \mathcal{M}(\theta_d, 1)$$

$$\theta_d = \sigma(\eta_d)$$

$$w_{dn} \sim \mathcal{M}(\phi_{z_{dn}}, 1)$$

Partial Model

The Model

$$\eta_d \sim \mathcal{N}(A x_d, \Sigma)$$

$$\theta_d = \sigma(\eta_d)$$

$$z_{dn} \sim \mathcal{M}(\theta_d, 1)$$

$$w_{dn} \sim \mathcal{M}(\phi_{z_{dn}}, 1)$$

Prior over A?

Matrix-Variate Priors

Matrix-Variate Normal Distribution [20] :

$$A \sim \mathcal{N}(M, \Omega, \Sigma) \implies \text{vec}(A) \sim \mathcal{N}(\text{vec}(M), \Sigma \otimes \Omega)$$

Matrix-Variate Priors

Matrix-Variate Normal Distribution [20] :

$$A \sim \mathcal{N}(M, \Omega, \Sigma) \implies \text{vec}(A) \sim \mathcal{N}(\text{vec}(M), \Sigma \otimes \Omega)$$

$$\ln p(W) = -\frac{FL}{2} \ln 2\pi - \frac{F}{2} \ln |\Omega| - \frac{L}{2} \ln |\Sigma| - \frac{1}{2} \text{etr} \left(\Sigma^{-1} (W - M) \Omega^{-1} (W - M)^{\top} \right)$$

where $\text{etr}(X) = \exp(\text{tr}(X))$

Matrix-Variate Priors

Matrix-Variate Normal Distribution [20] :

$$A \sim \mathcal{N}(M, \Omega, \Sigma) \implies \text{vec}(A) \sim \mathcal{N}(\text{vec}(M), \Sigma \otimes \Omega)$$

$$V \sim \mathcal{N}(0, I, I)$$

$$A|V \sim \mathcal{N}(UZ, I, \Sigma)$$

$$A \sim \mathcal{N}(0, I + UU^T, \Sigma)$$

See also [2] for examples of matrix-variate normal priors.

Full Model

The Model

$$V \sim \mathcal{N}(0, I, I) \qquad A|V \sim \mathcal{N}(UV, I, \Sigma) \qquad (2)$$

$$\eta_d \sim \mathcal{N}(A x_d, \Sigma) \qquad \theta_d = \sigma(\eta_d) \qquad (3)$$

$$z_{dn} \sim \mathcal{M}(\theta_d, 1) \qquad w_{dn} \sim \mathcal{M}(\phi_{z_{dn}}, 1) \qquad (4)$$

Inference

- Bound the likelihood using quadratic bounds
- Use the variational bound with the mean-field approximation $q(V, A, \Theta, Z, \Phi) = q(V)q(A)q(\Theta)q(Z)q(\Phi)$
- Numerically eliminate $Z \in \mathbb{R}^{D \times N \times K}$ as it consumes huge amounts of computer memory.

Model

Two Datasets NIPS Papers from 1987 to 1999

- 682 Documents. Vocabulary of 12503 words
- Median document length of 1,532 words, total of 1,075,323 words across the corpus
- Features are: authors; citations; the year

Tweets from April to September 2013 (inclusive)

- 735,868 tweets from 572 users. Vocabulary of 82,698 words
- Median tweet length is 10 words, total of 7,272,228 word observations across the corpus
- Features are: authors; time at various granularities (hour, day, week, month)

Images are ongoing

Twitter Hashtags

- Our model predicts words according to $p(w_d|x_d)$
- So given a tweets features, we can generate words ourselves, instead of using the observed words
- We can generate hashtags
 - So given all tweets that do not have a given hashtag (like #eurozone)
 - In which tweet is it most likely to occur given the features

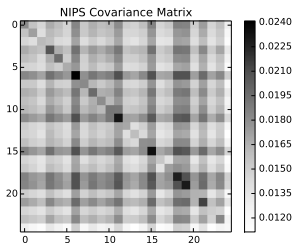
Twitter Hashtags

Hashtag	Tweet Words
#eurozone	#German finmin #Schaeuble argues #Karlsruhe may have no jurisdiction over #ECB measures.1st question by court's judges also focuses on this.
#usopen	Essa foi apenas a 2a vitoria de Gasquet nas oitavas de um Grand Slam. O frances alcanca sua primeira QFs em GS desde Wimbledon 2007 —
#f1	MT @f1paddockpass: ...a big shout out to the marshalls & volunteers here in Singapore. To all of you, heartfelt thanks @F1NightRace

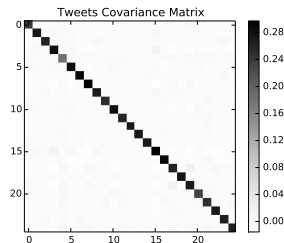
Twitter Hashtags

Hashtag	Tweet Words
#tcot	#Delaware Senate rejects bill to keep guns away from unstable people deemed danger to others http://hrld.us/110xnGt #guncontrol #NRA
#beer	No-Li Brewhouse on track for 150% growth, expanding annual capacity to 10,000 barrels ...
#obamacare	@michellemalkin #feded supporters simply believe you won't challenge them. #stopcommoncore @afpne

Covariances



(a) NIPS



(b) Tweets

Figure: Covariances over topics inferred using the Bohning implementation of the model for $K = 25$

What Next?

- Look at more complex matrix-variate priors, with low-rank approximations to both row and column covariances (the obvious approach doesn't work well)
- Large Scale inference (to 33 million tweets), potentially via SGD[23]
- Look at alternatives to the low-rank decomposition such as the use of matrix-variate Gaussian Scale models for sparsity[46]
- Look at richer downstream models for covariates such as collective-matrix[36] and tensor factorization[45].

Questions?

Reference II

- [4] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil.
Multi-Task feature learning.
In Advances in neural information processing systems,
volume 19, page 41, January 2005.
- [5] A Asuncion, M Welling, P Smyth, and YW Teh.
On smoothing and inference for topic models.
*In Proceedings of the Twenty-Fifth Conference on Uncertainty
in Artificial Intelligence*, pages 27–34. AUAI Press, 2012.
- [6] Bart Bakker and Tom Heskes.
Task Clustering and Gating for Bayesian Multitask Learning.
Journal of Machi, 4:83–99, 2003.

Reference III

- [7] D Blei and J Lafferty.
Correlated topic models.
Advances in neural information processing systems, 18:147,
2006.
- [8] David M Blei and Michael I Jordan.
Modeling annotated data.
*In Proceedings of the 26th annual international ACM SIGIR
conference on Research and development in informaion
retrieval.*, pages 127–134. ACM, 2003.
- [9] David M Blei and John D Lafferty.
Dynamic Topic Models.
*In Proceedings of the 23rd international conference on
Machine learning*, pages 113–120, 2006.

Reference V

- [13] Guillaume Bouchard.
Efficient Bounds for the Softmax Function and Applications
to Approximate Inference in Hybrid models.
In *NIPS*, pages 1–9, 2007.
- [14] Wray Buntine.
Variational Extensions to EM and Multinomial PCA.
Machine Learning: ECML 2002, pages 23–34, 2002.
- [15] Rich Caruana.
Multitask Learning.
Machine Learning, 28:41–75, 1997.
- [16] Jianfei Chen, Jun Zhu, Zi Wang, Xun Zheng, and Bo Zhang.
Scalable Inference for Logistic-Normal Topic Models.
In *NIPS*, pages 1–9, 2013.

Reference VIII

- [24] Mohamed Emtiyaz Khan, Shakir Mohamed, Benjamin M Marlin, and Kevin P Murphy.
A Stick-Breaking Likelihood for Categorical Data Analysis with Latent Gaussian Models.
In 15th International Conference on Artificial Intelligence and Statistics, volume XX, 2012.
- [25] David A Knowles and Thomas P Minka.
Non-conjugate Variational Message Passing for Multinomial and Binary Regression.
Advances in Neural Information Processing Systems, pages 1701–1709, 2011.

Reference XI

- [30] David Mimno and Andrew McCallum.
Topic Models Conditioned on Arbitrary Features with
Dirichlet-multinomial Regression.
*In Twenty-Fourth Conference on Uncertainty in Artificial
Intelligence*, June 2008.
- [31] David Newman and Max Welling.
Distributed algorithms for topic models.
The Journal of machine Learning research, 10:1801–1828,
2009.
- [32] John Paisley, Chong Wang, and David M. Blei.
The Discrete Infinite Logistic Normal Distribution.
Bayesian Analysis, 7(4):997–1034, December 2012.

Reference XII

- [33] J K Pritchard, M Stephens, and P Donnelly.
Inference of population structure using multilocus genotype data.
Genetics, 155(2):945–59, June 2000.
- [34] Joseph Reisinger, Austin Waters, and Raymond J Mooney.
Spherical Topic Models.
2010.
- [35] Konstantin Salomatin, Y Yang, and A Lad.
Multi-field correlated topic modeling.
In *Siam International Conference on Data Mining (SDM)*,
pages 628–637. SIAM, 2009.

Reference XIV

- [39] YW Teh, David Newman, and Max Welling.
A collapsed variational Bayesian inference algorithm for latent Dirichlet allocation.
Advances in Neural Information Processing Systems, 19(1353), 2007.
- [40] Seppo Virtanen, Yangqing Jia, Arto Klami, and Trevor Darrell.
Factorized multi-modal topic model.
In *Proceedings of the Twenty-Eighth Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI-12)*, pages 843–851. AUAI Press, 2012.

Reference XVI

- [44] Xuerui Wang, Andrew McCallum, and Xing Wei.
Topical n-grams: Phrase and topic discovery, with an application to information retrieval.
In Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on, pages 697–702. IEEE, 2007.
- [45] Liang Xiong, Xi Chen, Jeff Schneider, and Jaime G Carbonell.
Temporal Collaborative Filtering with Bayesian Probabilistic Tensor Factorization.
pages 211–222.
- [46] Ming Yang and Yingming Li.
Multi-Task Learning with Gaussian Matrix Generalized Inverse Gaussian Model.
In ICML, volume 1, pages 423–431, 2013.