Variational Methods for Discrete-Data Latent Gaussian Models

Mohammad Emtiyaz Khan

University of British Columbia Vancouver, Canada

March 6, 2012

The Big Picture

- Joint density models for data with mixed data types
- Bayesian models principled and robust approach
- Algorithms that are not only accurate and fast, but are also easy to tune, implement, and intuitive (speed-accuracy tradeoffs)

Variational Methods for Discrete-Data Latent Gaussian Models

Sources of Discrete Data

User rating data



Survey/voting data and blogs for sentiment analysis



Health data



tag correlation.



Consumer choice data



Sports/game data



Mohammad Emtiyaz Khan

Slide 3 of 46

Motivation: Recommendation system

Movie rating dataset - Missing values - Different types of data

	User1	User2	User3	User4	User5	User6	
Movie1	9	2	3		9		
Movie2	8			8		2	
Movie3		2		8			
Movie4	3	8	8			1	
Movie5	2		7		1		
Movie6		7		2		1	

Missing Ratings

From Wikipedia on Netflix-prize dataset

"The training set is such that the average user rated over 200 movies, and the average movie was rated by over 5000 users. But there is wide variance in the data—some movies in the training set have as few as 3 ratings, while one user rated over 17,000 movies."

Movielens Dataset



Missing ratings for movies

Sources of Discrete Data

User rating data



Survey/voting data and blogs for sentiment analysis



Health data



tag correlation.



Consumer choice data



Sports/game data



Mohammad Emtiyaz Khan

Slide 8 of 46

What we need!

For these datasets, we need a method of analysis which

- Handles missing values efficiently
- Makes efficient use of the data by weighting "reliable" data vectors more than the "unreliable" ones
- Makes efficient use of the data by "fusing" different types of data efficiently (binary, ordinal, categorical, count, text)

Factor Model



Bayesian Learning

$$\max_{\boldsymbol{\theta}} \sum_{n=1}^{N} \log \int p(\mathbf{y}_n | \mathbf{z}_n, \mathbf{W}) \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{z}_n$$





$$\geq \max_{\boldsymbol{\theta}} \sum_{n=1}^{N} \max_{\boldsymbol{\psi}_n} \underline{\mathcal{L}}_n(\boldsymbol{\theta}, \boldsymbol{\psi}_n)$$

This talk: Lower bound maximization

Variational Methods

$$\max_{\boldsymbol{\theta}} \sum_{n=1}^{N} \log \int p(\mathbf{y}_n | \mathbf{z}_n, \mathbf{W}) \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{z}_n \geq \max_{\boldsymbol{\theta}} \sum_{n=1}^{N} \max_{\boldsymbol{\psi}_n} \underline{\mathcal{L}}_n(\boldsymbol{\theta}, \boldsymbol{\psi}_n)$$

- Design tractable bounds to reduce approximation error
- Efficient optimization since lower bounds are concave : good convergence rates and easy convergence diagnostics
- Efficient expectation-maximization (EM) algorithms for parameter leaning
- Comparable performance to MCMC, but much faster
- Algorithms with a wide range of speed-accuracy trade-offs

Outline

- Latent Gaussian models
- Bounds for binary data
- Bounds for categorical data
- Results
- Future work and conclusions

Outline

- Latent Gaussian models
 - Definition and examples
 - Problem with parameter learning
- Bounds for binary data
- Bounds for categorical data
- Results
- Future work and conclusions

Latent Gaussian Model (LGM)

n=1:N

Likelihood Examples

Data type	Distribution	$p(y=1 \eta) = \frac{e^{\eta}}{1+e^{\eta}}$
Real	Gaussian	$1 + e^{\eta}$
Count	Poisson	0.9
Binary	Bernoulli-Logit	
Categorical	Multinomial-Logit	
Ordinal	Proportional-odds	η 2
		$\mathbf{v}(y = k \boldsymbol{\eta}) = \frac{e^{\eta_k}}{\sum_{j=1}^{K} e^{\eta_j}}$



Mohammad Emtiyaz Khan

Slide 16 of 46

Parameter Estimation

$$\max_{\boldsymbol{\theta}} \sum_{n=1}^{N} \log \int \prod_{d=1}^{D} p(y_{dn} | \mathbf{z}, \boldsymbol{\theta}) \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{z}$$





Jensen's Lower Bound

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}|\mathbf{y}) &= \log \int \prod_{d=1}^{D} p(y_d|\mathbf{z}, \boldsymbol{\theta}) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{z} \\ &= \log \int \frac{\prod_{d=1}^{D} p(y_d|\mathbf{z}, \boldsymbol{\theta}) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V})} \mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V}) d\mathbf{z} \end{aligned}$$

$$\mathcal{L}(\boldsymbol{\theta}|\mathbf{y}) \geq \max_{\mathbf{m},\mathbf{V}} \sum_{d=1}^{D} \int [\log p(y_d|\mathbf{z},\boldsymbol{\theta})] \mathcal{N}(\mathbf{z}|\mathbf{m},\mathbf{V}) d\mathbf{z} - KL [\mathcal{N}(\mathbf{m},\mathbf{V})||\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma})]$$

Variational Lower Bound

$$\max_{\boldsymbol{\theta}} \sum_{n=1}^{N} \max_{\mathbf{m}_{n}, \mathbf{V}_{n}} \sum_{d=1}^{D} \int \log p(y_{dn} | \boldsymbol{\eta}_{dn}) \mathcal{N}(\mathbf{z} | \mathbf{m}_{n}, \mathbf{V}_{n}) d\mathbf{z}$$
$$- KL[\mathcal{N}(\mathbf{m}_{n}, \mathbf{V}_{n}) | | \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})]$$

- Generalized EM algorithm
- E-step involves minimizing convex function
- Early stopping in E-step
- (Almost) no tuning parameters
- Easy convergence diagnostics

Outline

- Latent Gaussian models
- Bounds for binary data
 - Bernoulli-logistic likelihood
 - The Bohning bound (Khan, Marlin, Bouchard, Murphy, NIPS 2010)
 - Piecewise bounds (Marlin, Khan, Murphy, ICML 2011)
- Bounds for categorical data
- Results
- Future work and conclusions

Bernoulli-Logit Likelihood

$$p(y = 1|\eta) = \frac{e^{\eta}}{1 + e^{\eta}}$$

$$\log p(y = 1|\eta) = \eta - \log(1 + e^{\eta})$$

$$\mathcal{L}(\theta|\mathbf{y}) \geq \max_{\mathbf{m}, \mathbf{V}} \sum_{d=1}^{D} \int [\log p(y_d|\mathbf{z}, \theta)] \mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V}) d\mathbf{z}$$

$$-KL [\mathcal{N}(\mathbf{m}, \mathbf{V}) || \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})]$$

 $= \max_{\mathbf{m},\mathbf{V}} \sum_{d=1}^{D} \int \left[-\log(1+e^{\eta})\right] \mathcal{N}(\tilde{m}_{d}, \tilde{v}_{d}) d\eta + \text{tractable terms} \\ \inf_{\mathbf{m},\mathbf{N}} \inf_{\mathbf{M}} \operatorname{d} \mathbf{V} \right] \mathcal{N}(\tilde{m}_{d}, \tilde{v}_{d}) d\eta + \operatorname{tractable terms} \\ \operatorname{tractable terms} \\ \operatorname{in} \mathbf{m} \text{ and } \mathbf{V}$

Mohammad Emtiyaz Khan

Local Variational Bounds





- Bohning's bound (Khan, Marlin, Bouchard, Murphy 2010)
- Jaakola's bound (Jaakkola and Jordan1996)
- Piecewise quadratic bounds (Marlin, Khan, Murphy 2011)

Bohning Bound is Faster





Slide 25 of 46

Piecewise bounds are more accurate



Details of Piecewise bounds



- Find cut points and parameters of each piece by minimizing maximum error
- Linear pieces (Hsiung, Kim and Boyd, 2008)
- Quadratic Pieces (Nelder-Mead method)
- Fixed Piecewise Bounds!
- Increase accuracy by increasing the number of pieces

Outline

- Latent Gaussian models
- Bounds for binary data
- Bounds for categorical data
 - Multinomial-logistic likelihood and local variational bounds
 - Stick-breaking likelihood (Khan, Mohamed, Marlin, Murphy, AI-Stats 2012)
- Results
- Future work and conclusions

Multinomial-Logit Likelihood

$$p(y = k | \boldsymbol{\eta}) = \frac{e^{\eta_k}}{\sum_{j=1}^{K} e^{\eta_j}}$$

$$\log p(y = k | \boldsymbol{\eta}) = \eta_k - \log \sum_{j=1}^{K} e^{\eta_j}$$

$$\mathcal{L}(\boldsymbol{\theta} | \mathbf{y}) \geq \max_{\mathbf{m}, \mathbf{V}} \sum_{d=1}^{D} \int [\log p(y_d | \mathbf{z}, \boldsymbol{\theta})] \mathcal{N}(\mathbf{z} | \mathbf{m}, \mathbf{V}) d\mathbf{z}$$

$$- KL [\mathcal{N}(\mathbf{m}, \mathbf{V}) || \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})]$$

Local Variational bounds

- The Bohning bound
 - Fast and closed form updates
- The log bound (Blei and Lafferty 2006)
 - More accurate than the Bohning bound, but slower
- The product of sigmoid bound (Bouchard 2007)

Stick-Breaking Likelihood

$$p(y = 1|\eta) = \frac{e^{\eta_1}}{1 + e^{\eta_1}}$$

$$p(y = 2|\eta) = \left(1 - \frac{e^{\eta_1}}{1 + e^{\eta_1}}\right) \frac{e^{\eta_2}}{1 + e^{\eta_2}}$$

$$p(y = 3|\eta) = \left(1 - \frac{e^{\eta_1}}{1 + e^{\eta_1}}\right) \left(1 - \frac{e^{\eta_2}}{1 + e^{\eta_2}}\right) \frac{e^{\eta_3}}{1 + e^{\eta_3}}$$

$$\vdots$$

$$p(y = K|\eta) = \prod_{j=1}^{K-1} \left(1 - \frac{e^{\eta_j}}{1 + e^{\eta_j}}\right)$$

$$\log p(y = k | \boldsymbol{\eta}) = \eta_k - \sum_{j=1}^{K-1} \mathbf{I}(j \le k) \log(1 + e^{\eta_j})$$

Slide 32 of 46

Mohammad Emtiyaz Khan

Outline

- Latent Gaussian models
- Bounds for binary data
- Bounds for categorical data
- Results
- Future work and conclusions

Speed Accuracy Trade-offs

Binary FA : UCI voting dataset (D=15, N=435)



Comparison with EP

Binary Gaussian Process : Ionosphere dataset (D=200)

 $\Sigma_{ij} = \sigma \exp[-||x_i - x_j||^2/s]$



EP vs PW : Posterior Distribution



Slide 38 of 46

Mohammad Emtiyaz Khan

Comparison with EP

- Both methods give very similar results for GPs
- Our approach can be easily extended to factor models
- Variational EM objective function is well-defined and can be obtained by solving minimization of convex functions
- Numerically stable

MultiClass Gaussian Process







Categorical Factor Analysis

Glass dataset (D = 10, N = 958, sum of K = 29)

Outline

- Latent Gaussian models
- Bounds for binary data
- Bounds for categorical data
- Results
- Future work and conclusions

Future Work

- Large-scale collaborative filtering
- Use convexity to design approximate gradient methods
- Sparse Gaussian Posterior Distribution
- Tuning HMC using Bayesian optimization methods
- Latent Sparse-factor model
- Conditional models (e.g. to model for tag-image correlation)

Conclusions

- Variational methods show comparable performance with existing approaches
 - The main sources of errors is the bounding error
 - Design of piecewise bounds to control these errors
 - A good control over speed-accuracy trade-offs can be obtained
- Variational lower bounds can be optimized efficiently
 - Use of convex optimization methods to get fast convergence rates and easy convergence diagnostics
 - Design of efficient expectation-maximization (EM) algorithms

Collaborators

Kevin Murphy UBC

Guillaume Bouchard XRCE, France

Benjamin Marlin U. Mass-Amherst

Shakir Mohamed U. Cambridge, now at UBC

Variational Methods for Discrete-Data Latent Gaussian Models

Thank You