

Constrained Relative Entropy Minimization with Applications to Multitask Learning

Dissertation Defense

Oluwasanmi Koyejo

PhD Candidate
Department of Electrical and Computer Engineering
The University of Texas at Austin

April 29, 2013

Probabilistic Inference

Quantify uncertainty about x as $P(x)$.

Principle of maximum entropy

Choose $P(x)$ that captures known information and is as difficult as possible to discriminate from uniform^a (Jaynes, 1957).

^aUniform distribution represents complete ignorance.

Principle of minimum discrimination information

Given $P(x)$, choose $Q(x)$ that captures new information and is as difficult to discriminate from $P(x)$ as possible (Kullback, 1959).

Comparison with Bayesian Inference

Why do we need (yet) another inference rule?

Bayesian Inference:

- Solid conceptual framework derived from axioms of probability.
- Provides rules for updating $P(x)$ given data samples \mathcal{D} .
- Can impose support constraints by defining appropriate $P(x)$.
- Support constraints are preserved by $P(x) \mapsto P(x|\mathcal{D})$.

But:

- Other kinds of structure (e.g. moment constraints) generally do not propagate.
- Bayesian inference does not provide a formal approach for updating P with new *constraint* information.

Prior Work

Relative entropy inference

- MAXENT for natural language processing (Berger et al., 1996).
- l_1 and l_2 norm ball constraints (Dudík et al., 2007).
- Banach norm constraints (Altun and Smola, 2006).
- Margin constraints (“probabilistic support vector machines”):
 - Maximum entropy discrimination (Jaakkola et al., 1999).
 - MED-LDA (Zhu et al., 2009).
 - Link prediction (Zhu, 2012).

Multitask Learning

An overview

Given multiple (related) learning tasks:

$$f_1 : x \mapsto y_1, \quad f_2 : x \mapsto y_2, \quad f_3 : x \mapsto y_3, \quad \dots$$

Single task learning:

- Learn f_1, f_2, f_3, \dots separately.

Multiple task learning:

- Learn $F : x \mapsto Y$ jointly.
- Exploits inter-task relationships to improve performance
- When F is a bilinear function, MTL can be posed as a *matrix* estimation task.

Table of Contents

- 1 Introduction
 - Main Contributions
 - Background
- 2 Constrained Relative Entropy Minimization
 - Constrained Bayesian Inference
- 3 Applications
 - Multitask Bipartite Ranking
- 4 Conclusion

Main Contributions

Constrained Inference

A framework for probabilistic inference subject to expectation constraints:

- Constraints may be any combination of domain knowledge and data.
- Constraint sets may be nonconvex.
- Generalization of Bayesian inference that can incorporate additional structural constraints.

Main Contributions

Applications

- Constrained multitask learning and parameter estimation:
 - Applied to recommender systems and fMRI.
 - Rank constraints on the postdata mean matrix.
 - Kronecker constraints on the postdata covariance matrix.
 - Sparsity constraints on prior precision matrix.
- Nonparametric multitask learning:
 - Applied to predictive modeling of transposable data.
 - Nuclear norm constraints on the mean function.
- Multitask bipartite ranking:
 - Applied to disease-gene prioritization.
 - Ordering constraints on ranking variables.
 - Nuclear norm constraints on mean function of regression variables.

Table of Contents

- 1 Introduction
 - Main Contributions
 - Background
- 2 Constrained Relative Entropy Minimization
 - Constrained Bayesian Inference
- 3 Applications
 - Multitask Bipartite Ranking
- 4 Conclusion

Gaussian Process (GP)

- Collection of random variables $\{X(i), i \in \mathbb{I}\}$ completely described by:

$$\begin{aligned}\mu(i) &= \mathbb{E} [X(i)], \\ \mathcal{C}(i, i') &= \mathbb{E} [(X(i) - \mu(i))(X(i') - \mu(i'))].\end{aligned}$$

- The collection $\mathbf{x} = \{x(i), \forall i \in \mathbb{J}\}$ is jointly Gaussian for any finite index set $\mathbb{J} \subset \mathbb{I}$.
- Matrix-variate Gaussian process (MV-GP) $\mathcal{MG}\mathcal{P}(\phi, \mathcal{C}_{\mathbf{N}}, \mathcal{C}_{\mathbf{M}})$ is a special case:
 - Double index $(m, n) \in \mathbb{M} \times \mathbb{N} = \mathbb{J}$.

$$\mathcal{C}((m, n), (m', n')) = \mathcal{C}_{\mathbf{M}}(m, m')\mathcal{C}_{\mathbf{N}}(n, n').$$

Relative Entropy (Kullback-Leibler divergence)

For probability measures P, Q

Let P be absolutely continuous with respect to Q , there exist probability density functions p, q so that:

$$\text{KL}(p\|q) = \int_{\mathbb{X}} p(x) \log \frac{p(x)}{q(x)} dx$$

Some useful optimization-theoretic properties:

- Strictly convex wrt p .
- $\text{KL}(p\|q) \geq 0$.
- $\text{KL}(p\|q) = 0$ iff. $p = q$.

Exponential families

$$p_{\boldsymbol{\theta}}(x) = h(x)e^{\langle \boldsymbol{\eta}(\boldsymbol{\theta}), \mathbf{t}(x) \rangle - A(\boldsymbol{\theta})}$$

- $\boldsymbol{\eta}(\boldsymbol{\theta})$ natural parameter
- $\mathbf{t}(x)$ sufficient statistics
- $h(x)$ base measure
- $A(\boldsymbol{\theta})$ log-partition function

The domain of the parameter $\boldsymbol{\theta}$ is a convex set defined as:

$$\Theta = \left\{ \boldsymbol{\theta} \left| \int_{\mathbb{X}} h(x)e^{\langle \boldsymbol{\eta}(\boldsymbol{\theta}), \mathbf{t}(x) \rangle} dx < \infty \right. \right\}.$$

Constraints

We study constraints of the form $E [\beta(x)] \in C$.

- β is called the *feature function*.
- C is a constraint set.

Examples:

- Empirical average (Jaynes, 1957): $E [\beta(x)] = \frac{1}{N} \sum_n \beta(\hat{x}_n)$.
- Structural constraints (this dissertation):

$$E [\beta(x)] \geq 0, \quad \|E [\beta(x)]\|_0 \leq \epsilon.$$

Constrained Relative Entropy Minimization

Problem Statement

$$\inf_{q \in \mathcal{P}} \text{KL}(q||p) \text{ s.t. } E_q [\beta] \in \mathcal{C}$$

Assumptions:

- $\mathcal{C} \subset \mathcal{B}$ is closed.
- if q_* is a solution, then $E_{q_*} [\beta] = \mathbf{a}_*$ is bounded.

Representation result

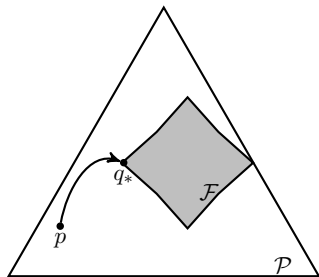
Every q_* can be represented as:

$$q_*(x) = p(x)e^{\langle \lambda_*, \beta(x) \rangle - G(\lambda_*)}$$

where λ_* is a function of \mathbf{a}_* .

Main Idea

Find solution using exponential family representation



$$\min_{\mathbf{c} \in A} \left[\min_{q \in \mathcal{P}} \text{KL}(q \| p) \text{ s.t. } E_q[\boldsymbol{\beta}] = \mathbf{c} \right]$$

- Inner problem is “easy” to solve.
- Find compact $A \subset C$ so that $\mathbf{a}_* \in A$.
- Define $\mathcal{F} = \{q_{\mathbf{c}} \mid \mathbf{c} \in A\} \subset \mathcal{P}$.
- Show that $q_* \in \mathcal{F}$.

Generalize to any $\mathcal{Q} \subset \mathcal{P}$ iff. $q_* \in \mathcal{Q}$.

Conjugate Priors for Relative Entropy

Motivation

- Exponential family representation requires normalization constant $G(\boldsymbol{\lambda})$.
- In many cases, computing $G(\boldsymbol{\lambda})$ is a challenging high dimensional problem.

Definition

The distribution $p \in \mathcal{Q}$ is a relative entropy conjugate prior distribution if any solution $q_* \in \mathcal{Q}$.

Conjugate Priors for Relative Entropy

Some properties

No need to compute $G(\lambda)$ explicitly if $p \in \mathcal{Q}$ is conjugate:

$$f_* = \arg \min_{f \in \mathcal{Q}} \text{KL}(f \| p) \text{ s.t. } E_f[\beta] \in \mathcal{C}$$

and $f_* \equiv q_*$.

Proposition

Let

$$\mathcal{Q} = \left\{ p \left| p_{\eta, \nu}(x) = h(x, \nu) e^{\langle \eta, \beta(x) \rangle - D(\eta, \nu)} \right. \right\},$$

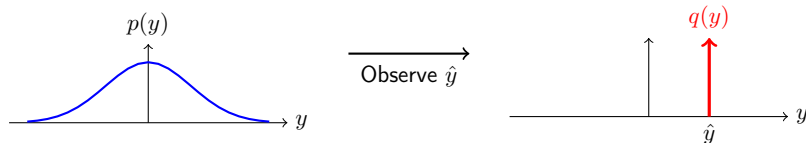
\mathcal{Q} are a family of relative entropy conjugate priors i.e. any $q_* \in \mathcal{Q}$.

Constrained Bayesian Inference

Bayesian Inference \subset Constrained Relative Entropy

Sample space: \mathbb{X} latent variables and \mathbb{Y} observations.

- Once you observe $y = \hat{y}$, there is no longer any uncertainty!
Williams (1980)



- Hence $q(y) = \delta_{\hat{y}}$
- Equivalent constraint set $E_q [\delta_{\hat{y}}] = 1$

Bayesian Inference \subset Constrained Relative Entropy

Solve:

$$\min_{q \in \mathcal{P}} \left[\text{KL}(q(x, y) \| p(x, y)) \text{ s.t. } E_q [\delta_{\hat{y}}] = 1 \right]$$

Solution is the Bayesian posterior!

$$q_*(x, y) = q_*(x|y = \hat{y})\delta_{\hat{y}} = p(x|y = \hat{y})\delta_{\hat{y}}$$

In other words:

$$q(x|y = \hat{y}) \equiv p(x|y = \hat{y})$$

Bayesian Inference with Additional Constraints

Our approach - incorporate further constraints on $q(x)$.

Constrained Bayesian inference

Solve:

$$\min_{q(y), q(x|y) \in \mathcal{P}} \left[\text{KL}(q(x|y)q(y) \| p(x, y)) \quad \text{s.t.} \quad \begin{array}{l} \mathbb{E}_q[\delta_{\hat{y}}] = 1 \\ \mathbb{E}_q[\beta(x)] \in \mathcal{C} \end{array} \right]$$

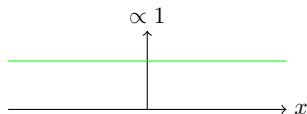
or the equivalent problem:

$$\min_{q \in \mathcal{P}} \left[\text{KL}(q(x) \| p(x|y)) \quad \text{s.t.} \quad \mathbb{E}_q[\beta(x)] \in \mathcal{C} \right]$$

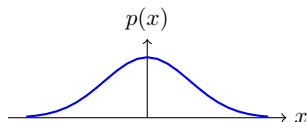
Table of Contents

- 1 Introduction
 - Main Contributions
 - Background
- 2 Constrained Relative Entropy Minimization
 - Constrained Bayesian Inference
- 3 Applications
 - Multitask Bipartite Ranking
- 4 Conclusion

Overview of the Approach

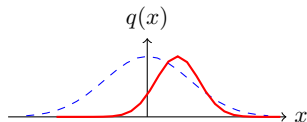


Domain \mathbb{X}
Complete ignorance



Maximize entropy given constraints

$$\min_p \left[H(p) \text{ s.t. } E_q [\gamma(x)] \in D \right]$$



Minimize relative entropy
given data and constraints

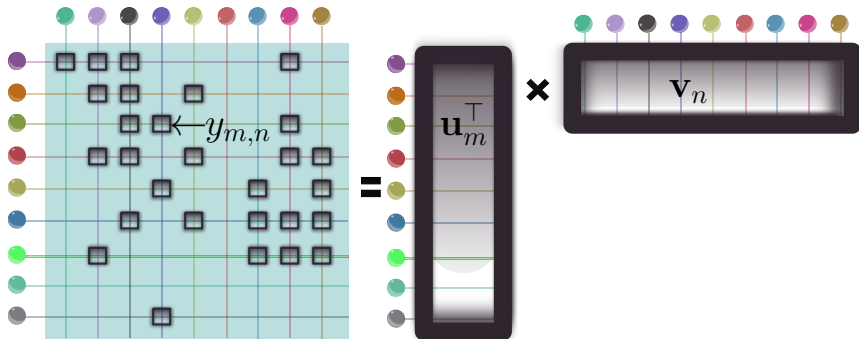
$$\min_q \left[\text{KL}(q||p) \text{ s.t. } E_q [\beta(x)] \in C \right]$$

Illustrative Example

Low Rank Expectation Constraints

Matrix Factorization

- Estimate rank R factors $\mathbf{u}_m \in \mathbb{R}^R$ and $\mathbf{v}_n \in \mathbb{R}^R$ so that $y_{m,n} \approx \mathbf{u}_m^\top \mathbf{v}_n$.



Initial constraints and Prior Distribution

Domain:

- $\mathbf{y} \in \mathbb{R}^K$ with $(m, n) \in K$
- $\mathbf{W} \in \mathbb{R}^{M \times N}$, $\mathbf{w} = \text{vec}(\mathbf{W})$, let $\mathbf{P} : M \times N \mapsto K$.

Initial information:

- $E[\mathbf{y}|\mathbf{w}] = \mathbf{P}\mathbf{w}$, $E[\mathbf{y}\mathbf{y}^\top] = \sigma^2\mathbf{I}$.
- $E[\mathbf{w}] = \mathbf{0}$, $E[\mathbf{w}\mathbf{w}^\top] = \mathbf{C} \otimes \mathbf{R}$.

Prior distribution

- Solve for maximum entropy distribution given constraints.
- Unique solution:

$$p(\mathbf{y}|\mathbf{w}) = \mathcal{N}(\mathbf{P}\mathbf{w}, \sigma^2), \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{C} \otimes \mathbf{R}).$$

Constrained Relative Entropy Approach

- Observe $\mathbf{y} = \hat{\mathbf{y}}$
- Additional constraints $\text{rank}(\mathbf{E}[\mathbf{W}]) \leq R$.

Solution

- Let $\mathcal{C} = \{\mathbf{B} \mid \text{rank}(\mathbf{B}) \leq R\}$
- $p(\mathbf{W}|\mathcal{D}) = \mathcal{N}(\mathbf{A}, \mathbf{\Sigma})$ is the Bayesian posterior.
- p is a relative entropy conjugate prior for $\beta(\mathbf{W}) = \mathbf{W}$.
- Hence q is Gaussian $q(\mathbf{W}) = \mathcal{N}(\mathbf{M}_*, \mathbf{S}_*)$ and:

$$\mathbf{M}_*, \mathbf{S}_* = \arg \min_{\mathbf{M}, \mathbf{S}} \left[\text{KL}(\mathcal{N}(\mathbf{M}, \mathbf{S}) \parallel \mathcal{N}(\mathbf{A}, \mathbf{\Sigma})) \text{ s.t. } \mathbf{M} \in \mathcal{C} \right]$$

Bayesian Approach

Factor Model

Prior Distribution:

$$\mathbf{y}|\mathbf{w} \sim \mathcal{N}(\mathbf{P}\mathbf{w}, \sigma^2),$$
$$\mathbf{W} = \mathbf{U}\mathbf{V}^\top, \quad \mathbf{u}_m = \mathcal{N}(\mathbf{0}, \mathbf{R}), \quad \mathbf{v}_n = \mathcal{N}(\mathbf{0}, \mathbf{C}).$$

Inference:

- Easy to compute some marginals e.g. $p(\mathbf{y}|\mathbf{U})$ (Lawrence and Urtasun, 2009),
- but $p(\mathbf{U}, \mathbf{V}|\mathbf{y})$ is analytically intractable!

In practice:

- Compute MAP estimate (equivalent to matrix factorization)
- Variational approximation.
- Sampling: $\mathbb{E}_{p(\mathbf{U}, \mathbf{V}|\mathbf{y})}[\mathbf{W}]$ is unlikely to be low rank.

Application to Multitask Learning

- Constrained multitask learning and parameter estimation:
 - Applied to recommender systems and fMRI.
 - Rank constraints on the postdata mean matrix.
 - Kronecker constraints on the postdata covariance matrix.
 - Sparsity constraints on prior precision matrix.
- Nonparametric multitask learning:
 - Applied to predictive modeling of transposable data.
 - Nuclear norm constraints on the mean function.
- Multitask bipartite ranking:
 - Applied to disease-gene prioritization.
 - Ordering constraints on ranking variables.
 - Nuclear norm constraints on mean function of regression variables.

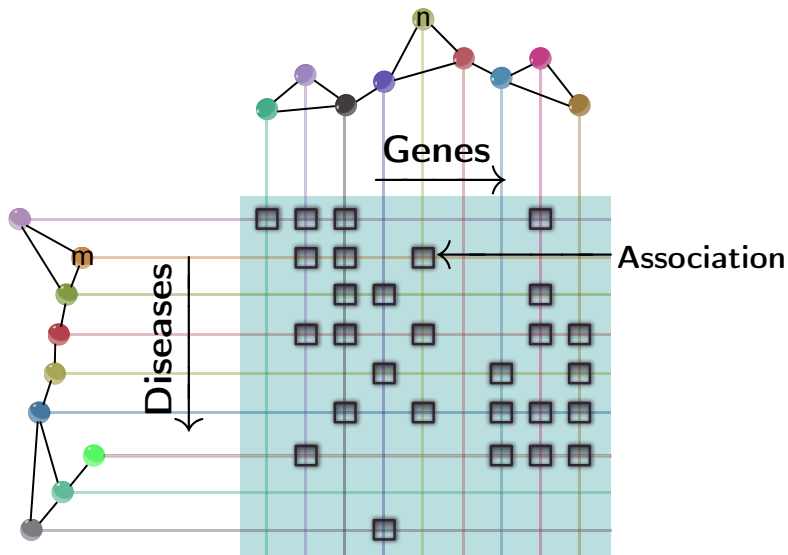
Multitask Bipartite Ranking

for Disease-Gene Prioritization

Genetic Diseases

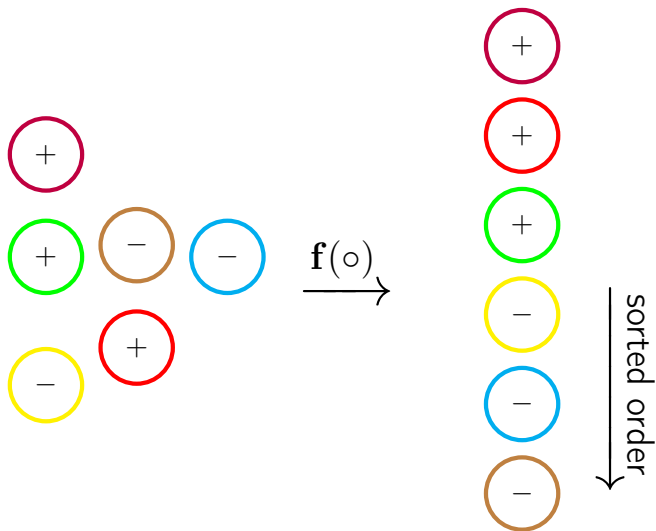
- Asthma, Cancer, Alzheimer's and many other diseases are associated with genetic mutations.
- Medical interest in unknown disease gene associations for targeted treatment, gene research, etc.
- 20,000 – 25,000 human genes and thousands of diseases so enumerative search is computationally intractable.
- Computational models can be used to suggest most likely genetic associations.

Gene Disease Data



Bipartite Ranking

Inferring order from binary examples



Problem Setup

Sample Space:

- Associate $x_{m,n}$ to each observation $y_{m,n}$ (*ranking variable*).
- Associate $z_{m,n}$ to each $x_{m,n}$ (*regression variable*).

Prior Distribution

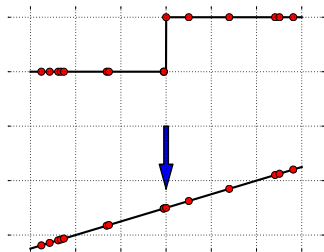
$$x_{m,n} \sim \mathcal{N}(z_{m,n}, \sigma^2), \quad Z \sim \mathcal{MG}\mathcal{P}(0, \mathcal{C}_M, \mathcal{C}_N).$$

Regression variable constraint

Nuclear norm constraints for parsimony and performance:

$$\mathbb{E}[Z] \in \mathcal{T}_\epsilon = \{\Phi \mid \|\Phi\|_{1-\mathcal{H}_C} \leq \epsilon\}$$

Ranking Variable Constraint



- Adapted from monotone retargeting (Acharyya et al., 2012)
- Key idea: Match ordering in \mathbf{x}_m with ordering in \mathbf{y}_m
- Encoded as $\mathbb{E}[\mathbf{x}] \in \mathcal{R}_y$

The resulting problem is given by:

$$q(\mathbf{x}, Z) = \arg \min_{q \in \mathcal{P}} \left[\begin{array}{l} \text{KL}(q(\mathbf{x}, Z) \| p(\mathbf{x}, Z)) \\ \text{s.t. } E_q[\mathbf{x}] \in \mathbf{R}_y, E_q[Z] \in \mathbf{T}_\epsilon \end{array} \right]$$

Solution

- $q(\mathbf{x}, Z)$ is a Gaussian process.
- Solution is unique!
- The marginal distribution $q(Z)$ is a matrix-variate GP with constrained (usually low rank) mean function.

Key Steps

- $q(\mathbf{x}, Z)$ is Gaussian for any finite index set via conjugacy.
- $\text{KL}(q(\mathbf{x}, Z) \| p(\mathbf{x}, Z))$ reduces to $\text{KL}(q(\mathbf{x}, \mathbf{Z}) \| p(\mathbf{x}, \mathbf{Z}))$
 - Proof uses representation theorem for spectral regularization (Abernethy et al., 2009).
 - \mathbf{Z} is the finite “hidden” matrix corresponding to observations.
- Solve for $q(\mathbf{x}, \mathbf{Z})$.
- “Expand” finite $q(\mathbf{x}, \mathbf{Z})$ to random process $q(\mathbf{x}, Z)$
 - Using Kolmogorov’s extension theorem (Kolmogorov, 1933).

Experiments

Datasets:

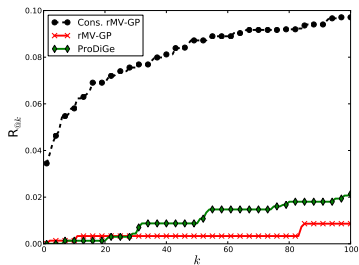
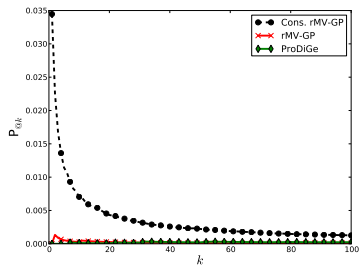
- **OMIM**: Set of manually-generated review articles linking inheritable diseases to one or a few genes ($M = 3,210$, $N = 13,614$, $J = 3,636$, $d = 0.0083\%$).
- **Medline**: Set of disease gene associations generated by processing annotations in the NIH PubMed/Medline literature ($M = 4,496$, $N = 21,243$, $J = 250,190$, $d = 0.36\%$).

Experimental Setup:

- “Known” diseases: randomly hidden associations.
- “New” diseases: randomly hidden diseases (*cold start*).

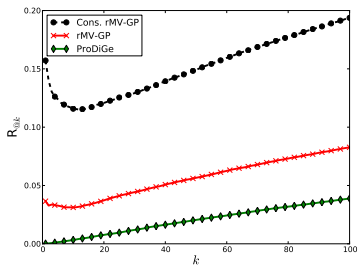
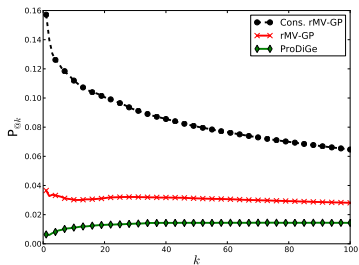
Results

OMIM



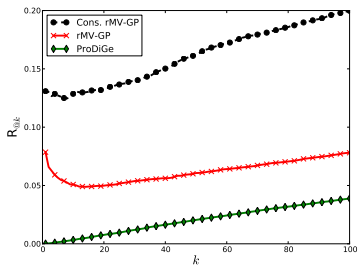
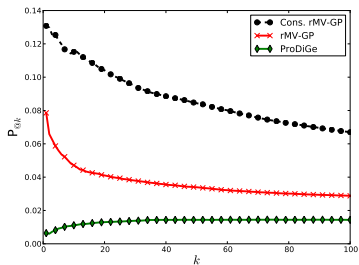
Results

Medline "Known" Diseases



Results

Medline "New" Diseases



Predicted Candidate Genes

Alzheimer's Disease

| Rank | Gene | Description |
|------|-------|---|
| 3 | IGF1 | Insulin-signaling protein linked to aging |
| 9 | MTHFR | Enzyme linked to cognitive impairments |
| 15 | IL1RN | Inflammatory marker linked to AD |
| 16 | APOE | Lipoprotein component linked to AD |

Prostate Cancer

| Rank | Gene | Description |
|------|-------|--|
| 2 | CRP | Marker linked to onset age and survival |
| 7 | AGT | Downstream products linked to metastasis |
| 9 | VEGFA | Promotes blood vessel growth in tumors |
| 13 | MMP9 | Protein linked to metastatic process |

Table of Contents

- 1 Introduction
 - Main Contributions
 - Background
- 2 Constrained Relative Entropy Minimization
 - Constrained Bayesian Inference
- 3 Applications
 - Multitask Bipartite Ranking
- 4 Conclusion

Conclusion

Developed a framework for probabilistic inference subject to expectation constraints:

- constraints may be any combination of domain knowledge and data.
- constraint set may be nonconvex.

Applied the results to multitask learning:

- Constrained multitask learning and parameter estimation.
- Nonparametric multitask learning for transposable data.
- Multitask bipartite ranking for disease gene prioritization.

Future Work

- Constrained relative entropy inference is difficult to optimize using sampling techniques. These may be necessary for complicated problems.
- The presented inference approach can be extended to other divergence metrics such as Csiszár divergences and Bregman divergences.
- Hyperparameter optimization is the most challenging aspect of the current approach. More research is required to improve available methods for large scale problems.
- Further investigation is required to fully understand the biological implications of the multitask bipartite ranking model.

Thank You

Related Publications

- Oluwasanmi Koyejo and Joydeep Ghosh. Constrained Bayesian inference for low rank multitask learning. In *Proceedings of the 29th conference on Uncertainty in artificial intelligence (UAI)*, 2013b. (Oral)
- Oluwasanmi Koyejo and Joydeep Ghosh. A representation approach for relative entropy minimization with expectation constraints. In *Workshop on Divergences and Divergence Learning (WDDL)*, 2013a
- Cheng Lee, Oluwasanmi Koyejo, and Joydeep Ghosh. Identifying candidate disease genes using a trace norm constrained bipartite raking model. In *Engineering in Medicine and Biology Society (EMBC)*, 2013
- Sreangsu Acharyya*, Oluwasanmi Koyejo*, and Joydeep Ghosh. Learning to rank with Bregman divergences and monotone retargeting. In *Proceedings of the 28th conference on Uncertainty in artificial intelligence (UAI)*, 2012
- Oluwasanmi Koyejo and Joydeep Ghosh. A kernel-based approach to exploiting interaction-networks in heterogeneous information sources for improved recommender systems. In *HetRec '11*, pages 9–16, 2011

Related Publications

Submitted Publications

- Oluwasanmi Koyejo, Cheng Lee, and Joydeep Ghosh. Constrained relative entropy minimization with application to disease gene prioritization. 2013d. Under review
- Oluwasanmi Koyejo, Cheng Lee, and Joydeep Ghosh. A constrained matrix-variate Gaussian process for transposable data. 2013c. Under review
- Sreangsu Acharyya, Oluwasanmi Koyejo, and Joydeep Ghosh. Learning to rank with Bregman divergences and monotone retargeting. 2013. Under review

Other Publications

- Mijung Park*, Oluwasanmi Koyejo*, Joydeep Ghosh, Russell R. Poldrack, and Jonathan W. Pillow. Bayesian structure learning for functional neuroimaging. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2013
- O. Koyejo, P. Patel, J. Ghosh, and R. A. Poldrack. Learning predictive cognitive structure from fMRI using supervised topic models. In *International Workshop on Pattern Recognition in Neuroimaging (PRNI)*, 2013a
- Oluwasanmi Koyejo, Cheng Lee, and Joydeep Ghosh. The trace norm constrained matrix-variate Gaussian process for multitask bipartite ranking. *CoRR*, abs/1302.2576, 2013e
- Oluwasanmi Koyejo and Joydeep Ghosh. MiPPS; a generative model for multi-manifold clustering. In *AAAI Fall Symposium on Manifold Learning and Its Applications*. AAAI Press, 2009
- Oluwasanmi Koyejo. Manifold learning and its applications: Reports of the AAAI 2010 fall symposia. *AI Magazine*, 32(1):93–100, 2011
- Richard Souvenir and Oluwasanmi Koyejo. Manifold learning and its applications: Reports of the AAAI 2009 fall symposia. *AI Magazine*, 31(1):88–94, 2010

Other Publications

Submitted Publications

- R. A. Poldrack, D. M. Barch, J. P. Mitchell, T. D. Wager, A. D. Wagner, J. T. Devlin, C. Cumba, O. Koyejo, and M. P. Milham. Towards open sharing of task-based fMRI data: The OpenfMRI project. 2013. Under review
- Oluwasanmi Koyejo, Sreangsu Acharyya, and Joydeep Ghosh. Retargeted matrix factorization. 2013b. Under review

References I

- Jacob Abernethy, Francis Bach, Theodoros Evgeniou, and Jean-Philippe Vert. A new approach to collaborative filtering: Operator estimation with spectral regularization. *JMLR*, 10:803–826, 2009. ISSN 1532-4435.
- Sreangsu Acharyya, Oluwasanmi Koyejo, and Joydeep Ghosh. Learning to rank with bregman divergences and monotone retargeting. In *Proceedings of the 28th conference on Uncertainty in artificial intelligence*, UAI '12, 2012. ISBN 0-9749039-0-6.
- Sreangsu Acharyya*, Oluwasanmi Koyejo*, and Joydeep Ghosh. Learning to rank with Bregman divergences and monotone retargeting. In *Proceedings of the 28th conference on Uncertainty in artificial intelligence (UAI)*, 2012.
- Sreangsu Acharyya, Oluwasanmi Koyejo, and Joydeep Ghosh. Learning to rank with Bregman divergences and monotone retargeting. 2013. Under review.
- Yasemin Altun and Alexander J. Smola. Unifying divergence minimization and statistical inference via convex duality. In *COLT*, 2006.
- Adam L. Berger, Vincent J. Della Pietra, and Stephen A. Della Pietra. A maximum entropy approach to natural language processing. *Comput. Linguist.*, 22(1):39–71, March 1996. ISSN 0891-2017.
- Miroslav Dudík, Steven J. Phillips, and Robert E. Schapire. Maximum entropy density estimation with generalized regularization and an application to species distribution modeling. *J. Mach. Learn. Res.*, 8:1217–1260, December 2007. ISSN 1532-4435.
- Tommi Jaakkola, Marina Meila, and Tony Jebara. Maximum entropy discrimination. In *NIPS*. MIT Press, 1999.
- E. T. Jaynes. Information Theory and Statistical Mechanics. *Physical Review Online Archive (Prola)*, 106(4):620–630, 1957.
- A. N. Kolmogorov. *Foundations of the theory of probability*. Chelsea, New York, 1933. Trans. N. Morrison (1956).
- O. Koyejo, P. Patel, J. Ghosh, and R. A. Poldrack. Learning predictive cognitive structure from fMRI using supervised topic models. In *International Workshop on Pattern Recognition in NeuroImaging (PRNI)*, 2013a.

References II

- Oluwasanmi Koyejo. Manifold learning and its applications: Reports of the AAAI 2010 fall symposia. *AI Magazine*, 32(1):93–100, 2011.
- Oluwasanmi Koyejo and Joydeep Ghosh. MiPPS; a generative model for multi-manifold clustering. In *AAAI Fall Symposium on Manifold Learning and Its Applications*. AAAI Press, 2009.
- Oluwasanmi Koyejo and Joydeep Ghosh. A kernel-based approach to exploiting interaction-networks in heterogeneous information sources for improved recommender systems. In *HetRec '11*, pages 9–16, 2011.
- Oluwasanmi Koyejo and Joydeep Ghosh. A representation approach for relative entropy minimization with expectation constraints. In *Workshop on Divergences and Divergence Learning (WDDL)*, 2013a.
- Oluwasanmi Koyejo and Joydeep Ghosh. Constrained Bayesian inference for low rank multitask learning. In *Proceedings of the 29th conference on Uncertainty in artificial intelligence (UAI)*, 2013b. (Oral).
- Oluwasanmi Koyejo, Sreangsu Acharyya, and Joydeep Ghosh. Retargeted matrix factorization. 2013b. Under review.
- Oluwasanmi Koyejo, Cheng Lee, and Joydeep Ghosh. A constrained matrix-variate Gaussian process for transposable data. 2013c. Under review.
- Oluwasanmi Koyejo, Cheng Lee, and Joydeep Ghosh. Constrained relative entropy minimization with application to disease gene prioritization. 2013d. Under review.
- Oluwasanmi Koyejo, Cheng Lee, and Joydeep Ghosh. The trace norm constrained matrix-variate Gaussian process for multitask bipartite ranking. *CoRR*, abs/1302.2576, 2013e.
- Solomon Kullback. *Information Theory and Statistics*. Dover, 1959.
- Neil D. Lawrence and Raquel Urtasun. Non-linear matrix factorization with gaussian processes. In *ICML*, pages 601–608. ACM, 2009.
- Cheng Lee, Oluwasanmi Koyejo, and Joydeep Ghosh. Identifying candidate disease genes using a trace norm constrained bipartite raking model. In *Engineering in Medicine and Biology Society (EMBC)*, 2013.

References III

- Mijung Park*, Oluwasanmi Koyejo*, Joydeep Ghosh, Russell R. Poldrack, and Jonathan W. Pillow. Bayesian structure learning for functional neuroimaging. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2013.
- R. A. Poldrack, D. M. Barch, J. P. Mitchell, T. D. Wager, A. D. Wagner, J. T. Devlin, C. Cumba, O. Koyejo, and M. P. Milham. Towards open sharing of task-based fMRI data: The OpenfMRI project. 2013. Under review.
- Richard Souvenir and Oluwasanmi Koyejo. Manifold learning and its applications: Reports of the AAAI 2009 fall symposia. *AI Magazine*, 31(1):88–94, 2010.
- Peter M Williams. Bayesian conditionalisation and the principle of minimum information. *The British Journal for the Philosophy of Science*, 31(2):131–144, 1980.
- Jun Zhu. Max-margin nonparametric latent feature models for link prediction. In *ICML*, 2012.
- Jun Zhu, Amr Ahmed, and Eric P. Xing. Medlda: maximum margin supervised topic models for regression and classification. In *ICML*, 2009.