
An Incremental Subgradient Algorithm for MAP Estimation in Graphical Models

Jeremy Jancsary **Gerald Matz**¹ **Harald Trost**²

`jeremy.jancsary@ofai.at`

December 8, 2010



¹Vienna University of Technology

²Medical University of Vienna

A Few Things to Take Away From This Talk

- Why you might be interested in our algorithm:
 - It is **efficient**, both computationally and memory-wise.
 - It finds **better solutions** than the methods we compared it to.
 - It **converges** to the global optimum of the first-order LP relaxation of the MAP problem (which is tight in some cases).
 - You get a **certificate** of optimality w.r.t. the discrete problem.
- How it all works:
 - Start out with tree-reweighted upper bound (Wainwright et al., 2005).
 - The upper bound is developed until it assumes a degenerate form involving a large number of easy problems.
 - The *tightest* bound can then be found very efficiently using incremental methods, solving one easy problem at a time.
 - Equivalent to maximizing the LP relaxation of the MAP problem.



Outline

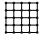
- 1 The Problem
- 2 Towards a Solution
- 3 Some Properties
- 4 Experimental Results
- 5 Conclusion



Outline

- 1 The Problem
- 2 Towards a Solution
- 3 Some Properties
- 4 Experimental Results
- 5 Conclusion

Maximum-a-Posteriori (MAP) Estimation

Consider an undirected graphical model (e.g. ) with vertex set \mathcal{V} and edge set \mathcal{E} defined over discrete random variables with pairwise interactions. Potential of a particular variable state $\mathbf{x} \in \mathcal{X}^n$:

$$P(\mathbf{x}; \boldsymbol{\theta}) = \sum_{s \in \mathcal{V}} \theta_s(x_s) + \sum_{(s,t) \in \mathcal{E}} \theta_{st}(x_s, x_t).$$

MAP Estimation (Discrete Problem)

(OP1)

$$\begin{aligned}\bar{P}(\boldsymbol{\theta}) &= \max_{\mathbf{x} \in \mathcal{X}^n} P(\mathbf{x}; \boldsymbol{\theta}), \\ \bar{\mathbf{x}} &= \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}^n} P(\mathbf{x}; \boldsymbol{\theta}).\end{aligned}$$

Computation of these quantities is **NP-hard** in general (notable exceptions: trees; binary variables + submodular energies).



Outline

- 1 The Problem
- 2 Towards a Solution**
- 3 Some Properties
- 4 Experimental Results
- 5 Conclusion

Relaxing the Discrete Problem

Tightest Tree-Reweighted Upper Bound

(OP2)

$$\min_{\{\theta^T\} \in \mathcal{C}(\theta)} \sum_T \rho^T \bar{P}(\theta^T),$$

$$\mathcal{C}(\theta) = \left\{ \{\theta^T\} \mid \sum_T \rho^T \theta^T = \theta \right\} \text{ and } \{\rho^T\} \in \left\{ \rho \geq \mathbf{0} \mid \sum_T \rho^T = 1 \right\}.$$

is connected through strong duality to

First-Order LP Relaxation

(OP3)

$$\max_{\mu \in \mathcal{L}} \sum_s \mu_s \cdot \theta_s + \sum_{(s,t)} \mu_{st} \cdot \theta_{st},$$

$$\mathcal{L} = \left\{ \mu \geq \mathbf{0} \mid \begin{array}{l} \sum_{x_s} \mu_s(x_s) = 1 \text{ for all } s \\ \sum_{x'_s} \mu_{st}(x'_s, x_t) = \mu_t(x_t), \sum_{x'_t} \mu_{st}(x_s, x'_t) = \mu_s(x_s) \end{array} \right\}.$$

Simplifying the Upper Bound (1)

- Curious fact #1 (Wainwright et al., 2005): Choice of ρ irrelevant as long as all edges are covered (otherwise, $\mathcal{C}(\theta)$ is empty).
Minimization of tree-reweighted upper bound \equiv maximization of LP relaxation, which does not depend on ρ .
- Can pick small set $S(\mathcal{T})$ of trees needed to cover all edges and set $\rho^T = \rho \stackrel{\text{def}}{=} 1/|S(\mathcal{T})|$ if $T \in S(\mathcal{T})$, and $\rho^T = 0$ otherwise.
- Exploit linearity of $P(\cdot)$, move ρ into params—viz. $\lambda^T = \rho\theta^T$:

“Dual Decomposition”-like Formulation (OP4)

$$\min_{\{\lambda^T\} \in \mathcal{S}(\theta)} \sum_{T \in \mathcal{S}(\mathcal{T})} \bar{P}(\lambda^T) \text{ with } \mathcal{S}(\theta) = \left\{ \{\lambda^T\} \mid \sum_{T \in \mathcal{S}(\mathcal{T})} \lambda^T = \theta \right\}.$$

Simplifying the Upper Bound (2)

- **Curious fact #2** (Kolmogorov, 2006): Trees in tree-reweighted upper bound need not be spanning (\rightarrow no impact on *tightest* bound).
- **Central idea #1**: Choose each tree T as single edge $E = (s, t)$.
- Determines almost all parameters:

$$\lambda_{st}^E \stackrel{!}{=} \begin{cases} \theta_{st} & \text{if } E = (s, t) \\ \mathbf{0} & \text{otherwise} \end{cases}, \quad \lambda_s^E \stackrel{!}{=} \mathbf{0} \text{ if } s \notin E.$$

- Remaining parameters of an edge $E = (s, t)$: $\lambda^E = \{\lambda_s^E, \lambda_t^E\}$.

Tightest Degenerate Upper Bound

(OP5)

$$\min_{\lambda \in \mathcal{Q}(\theta)} D(\lambda; \theta) \stackrel{\text{def}}{=} \sum_E \max_{(x_s, x_t)} \{\lambda_s^E(x_s) + \lambda_t^E(x_t) + \theta_{st}(x_s, x_t)\},$$

$$\mathcal{Q}(\theta) = \left\{ \{\lambda^E\} \mid \sum_{E:s \in E} \lambda_s^E = \theta_s \text{ for all } s \right\}.$$

Tightening the Upper Bound

Subgradient

Objective $D(\lambda; \theta)$ is non-differentiable, but $g \in \mathbb{R}^{2|\mathcal{X}^E|}$ is given by: $g_s^E(x_s) = [x_s = \bar{x}_s^E]$, $g_t^E(x_t) = [x_t = \bar{x}_t^E]$ for all $E = (s, t)$, x_s, x_t , where we use $(\bar{x}_s^E, \bar{x}_t^E)$ to refer to the **edge MAP state** (cf. OP5).

Projection

We need to solve $\operatorname{argmin}_{\lambda \in Q(\theta)} \|\lambda - \lambda'\|_2^2$. Solution obtained as:

$$\mathcal{P}_\theta(\lambda') = \left\{ \lambda_s^E(x_s) \leftarrow \lambda_s'^E(x_s) - \left(\sum_{E' \in \mathcal{E}_s} \lambda_s'^{E'}(x_s) - \theta_s(x_s) \right) / |\mathcal{E}_s| \right\},$$

which distributes amount of change uniformly over adjacent edges.

Central idea #2: Separable, non-diff. problem, cheap projection \rightarrow use incremental subgradient method (Nedić and Bertsekas, 2001).

The Algorithm

Input : Graph G , target parameters θ , initial feasible point λ

Output: Feasible primal solution \tilde{x} that is an approximation to \bar{x}

choose initial feasible primal solution \tilde{x} arbitrarily ;

repeat

pick next step size α and shuffle the set of edges \mathcal{E} ;

foreach $E = (s, t) \in \mathcal{E}$ **do**

find MAP state: $(\bar{x}_s^E, \bar{x}_t^E)$;

subtract subgradient: $\lambda_s^E(\bar{x}_s^E) \leftarrow \lambda_s^E(\bar{x}_s^E) - \alpha$, $\lambda_t^E(\bar{x}_t^E) \leftarrow \lambda_t^E(\bar{x}_t^E) - \alpha$;

foreach $E' \in \mathcal{E}_s$ **do** project: $\lambda_s^{E'}(\bar{x}_s^E) \leftarrow \lambda_s^{E'}(\bar{x}_s^E) + \alpha/|\mathcal{E}_s|$;

foreach $E' \in \mathcal{E}_t$ **do** project: $\lambda_t^{E'}(\bar{x}_t^E) \leftarrow \lambda_t^{E'}(\bar{x}_t^E) + \alpha/|\mathcal{E}_t|$;

foreach $s \in \mathcal{V}$ **do** build candidate \tilde{c} : $\tilde{c}_s \leftarrow$ at random from $\{\bar{x}_s^E \mid E \in \mathcal{E}_s\}$;

if $P(\tilde{c}; \theta) > P(\tilde{x}; \theta)$ **then**

accept best primal solution so far: $\tilde{x} \leftarrow \tilde{c}$;

if $D(\lambda; \theta) = P(\tilde{x}; \theta)$ **then**

optimal primal solution found: **return** \tilde{x} ;

until converged;

approximate primal solution found: **return** \tilde{x} ;



Outline

- 1 The Problem
- 2 Towards a Solution
- 3 Some Properties**
- 4 Experimental Results
- 5 Conclusion

Formal Guarantees

Proposition (Global Convergence)

For an appropriately chosen sequence of step sizes $\{\alpha^{(k)}\}$, convergence to the global optimum of (OP5) is guaranteed as $k \rightarrow \infty$.

The choice of $\{\alpha^{(k)}\}$ is discussed in detail in the paper.

Proposition (Optimality of Primal Solutions)

Assume that at an outer iteration, $P(\tilde{x}; \theta) = D(\lambda; \theta)$. It follows that \tilde{x} maximizes $P(\cdot)$ and λ minimizes $D(\cdot)$. This happens precisely if for each node, the edge MAP states agree on a common node MAP state.

We thus obtain a certificate of optimality for our primal solution.

Comparison to Related Approaches

Several methods have been devised that all aim at solving the first-order linear programming relaxation (OP3).

Method	Converg.	Global	Rate	Memory
INCMP (Our Method)	yes	yes	sublinear	$\mathcal{O}(\mathcal{X})$
DDSUB (Komodakis et al., 2007)	yes	yes	sublinear	$\mathcal{O}(\mathcal{X} ^2)$
TRWMP (Wainwright et al., 2005)	no	no	?	$\mathcal{O}(\mathcal{X})$
TRW-S (Kolmogorov, 2006)	yes	no	?	$\mathcal{O}(\mathcal{X})$
MPLP (Globerson et al., 2007)	yes	no	?	$\mathcal{O}(\mathcal{X})$
PROXLP (Ravikumar et al., 2010)	yes	yes	superlinear	$\mathcal{O}(\mathcal{X} ^2)$
DDPROX (Jojic et al., 2010)	yes	yes	linear	$\mathcal{O}(\mathcal{X} ^2)$

The convergence rates and working memory requirements listed above are **asymptotic** and do not necessarily reveal a lot about real-world performance (cost of an iteration is crucial).



Outline

- 1 The Problem
- 2 Towards a Solution
- 3 Some Properties
- 4 Experimental Results**
- 5 Conclusion

Experimental Setup

We compared three solvers (INCMP, DDSUB and TRWMP) on three different types of graphs, averaged over 20 runs.

GridsingUni: A 50×50 grid with binary variables ($\mathcal{X} = \{-1, +1\}$) and potentials given by $\theta_s(x_s) = \gamma x_s$ and $\theta_{st}(x_s, x_t) = \gamma x_s x_t$ with $\gamma \sim \mathcal{U}(-1, +1)$ drawn independently for each node and edge.

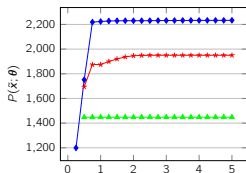
GridMultiGauss: A 20×20 grid with variables of arity $|\mathcal{X}| = 16$ and potentials chosen as $\theta_s(x_s) = 0$ and $\theta_{st}(x_s, x_t) \sim \mathcal{N}(0, 15)$ independently.

ComplsingUni: A complete graph of 50 binary variables with potentials chosen akin to GridsingUni.

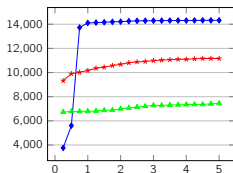
Results

- Measured the score $P(\tilde{x}; \theta)$ of the best primal solution \tilde{x} found so far, as a function of running time (seconds).
- For INCMP and DDSUB, constructed \tilde{x} randomly from the edge and tree MAP states, respectively (at each iteration).
- For TRWMP, used the maximizers of the node beliefs.

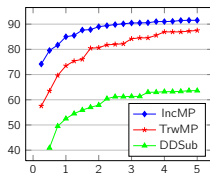
GridIsingUni 



GridMultiGauss 



ComplIsingUni 



Outline

- 1 The Problem
- 2 Towards a Solution
- 3 Some Properties
- 4 Experimental Results
- 5 Conclusion**

Future Work

- Step size α can be determined analytically so as not to increase the dual objective $P(\lambda, \theta)$.
 - Turns algorithm into a “dual descent method” (Bertsekas, 1999).
 - Open question: Can global convergence still be guaranteed?
 - Most likely, can get stuck in a “corner” (akin to MPLP).
- Use as computational core in a branch-and-bound scheme.
 - Low working memory requirements, ideal for parallelization.
 - Constraints like “ $X_s \stackrel{!}{=} x_s$ ” can easily be added, warm-starting should work rather well.
- Release source code as part of the **PhiWeave** package for approximate training of discriminative graphical models.

Some References



M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky.

MAP estimation via agreement on (hyper)trees: Message passing and linear-programming approaches.
IEEE Transactions on Information Theory, Vol. 51(11), 2005.



V. Kolmogorov.

Convergent tree-reweighted message passing for energy minimization.
Pattern Analysis and Machine Intelligence, Vol. 28(10), 2006.



A. Nedić and D. P. Bertsekas.

Incremental subgradient methods for nondifferentiable optimization.
SIAM Journal on Optimization, Vol. 12, 2001.



N. Komodakis, N. Paragios, and G. Tziritas.

MRF optimization via dual decomposition: Message-passing revisited.
IEEE 11th International Conference on Computer Vision, 2007.



A. Globerson and T. Jaakkola.

Fixing max-product: Convergent message passing algorithms for MAP-LP relaxations.
Advances in Neural Information Processing Systems, 2007.



P. Ravikumar, A. Argarwal, and M. J. Wainwright.

Message-passing for graph-structured linear programs: proximal methods and rounding schemes.
Journal of Machine Learning Research, Vol. 11, 2010.



V. Jojic, S. Gould, and D. Koller.

Accelerated dual decomposition for MAP inference.
27th International Conference on Machine Learning, 2010.

