

Dual Decomposition for Marginal Inference

Justin Domke

Rochester Institute of Technology

AAAI 2011

Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

Experiments

Conclusions

Conclusions

Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

Experiments

Conclusions

Conclusions

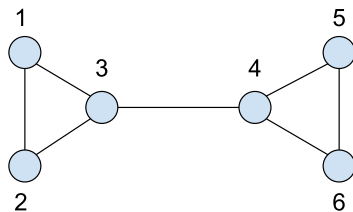
Graphical Models

- Markov Random Field / Factor Graph:

$$p(\mathbf{x}) \propto \prod_c \psi(\mathbf{x}_c)$$

Graphical Models

$$c_1 = \{1, 2, 3\}, \quad c_2 = \{3, 4\}, \quad c_3 = \{4, 5, 6\}$$



$$\begin{aligned} p(\mathbf{x}) &\propto \prod_c \psi(\mathbf{x}_c) \\ &= \psi(x_1, x_2, x_3) \psi(x_3, x_4) \psi(x_4, x_5, x_6) \end{aligned}$$

Marginal Inference

- Want to recover $p(X_i = x_i)$.
- Brute-force sum: Define $\hat{p}(\mathbf{x}) = \prod_c \psi(x_c)$

$$P(X_i = x_i) = \frac{1}{Z} \sum_{x_1} \dots \sum_{x_{i-1}} \sum_{x_{i+1}} \dots \sum_{x_M} \hat{p}(\mathbf{x})$$

$$Z = \sum_{x_1} \dots \sum_{x_M} \hat{p}(\mathbf{x})$$

- On trees, can do sums quickly by dynamic programming.
 - Sum-product algorithm / belief propagation
- #P-hard
 - Approximate: Tree-reweighted belief propagation (TRW)
 - This paper: Same approximation as TRW, different algorithm.

Marginal Inference

- Want to recover $p(X_i = x_i)$.
- Brute-force sum: Define $\hat{p}(\mathbf{x}) = \prod_c \psi(\mathbf{x}_c)$

$$P(X_i = x_i) = \frac{1}{Z} \sum_{\mathbf{x}_1} \dots \sum_{\mathbf{x}_{i-1}} \sum_{\mathbf{x}_{i+1}} \dots \sum_{\mathbf{x}_M} \hat{p}(\mathbf{x})$$

$$Z = \sum_{\mathbf{x}_1} \dots \sum_{\mathbf{x}_M} \hat{p}(\mathbf{x})$$

- On trees, can do sums quickly by dynamic programming.
 - Sum-product algorithm / belief propagation
- #P-hard
 - Approximate: Tree-reweighted belief propagation (TRW)
 - This paper: Same approximation as TRW, different algorithm.

Marginal Inference

- Want to recover $p(X_i = x_i)$.
- Brute-force sum: Define $\hat{p}(\mathbf{x}) = \prod_c \psi(\mathbf{x}_c)$

$$P(X_i = x_i) = \frac{1}{Z} \sum_{\mathbf{x}_1} \dots \sum_{\mathbf{x}_{i-1}} \sum_{\mathbf{x}_{i+1}} \dots \sum_{\mathbf{x}_M} \hat{p}(\mathbf{x})$$

$$Z = \sum_{\mathbf{x}_1} \dots \sum_{\mathbf{x}_M} \hat{p}(\mathbf{x})$$

- On trees, can do sums quickly by dynamic programming.
 - Sum-product algorithm / belief propagation
- #P-hard
 - Approximate: Tree-reweighted belief propagation (TRW)
 - This paper: Same approximation as TRW, different algorithm.

Marginal Inference

- Want to recover $p(X_i = x_i)$.
- Brute-force sum: Define $\hat{p}(\mathbf{x}) = \prod_c \psi(\mathbf{x}_c)$

$$P(X_i = x_i) = \frac{1}{Z} \sum_{\mathbf{x}_1} \dots \sum_{\mathbf{x}_{i-1}} \sum_{\mathbf{x}_{i+1}} \dots \sum_{\mathbf{x}_M} \hat{p}(\mathbf{x})$$

$$Z = \sum_{\mathbf{x}_1} \dots \sum_{\mathbf{x}_M} \hat{p}(\mathbf{x})$$

- On trees, can do sums quickly by dynamic programming.
 - Sum-product algorithm / belief propagation
- #P-hard
 - Approximate: Tree-reweighted belief propagation (TRW)
 - This paper: Same approximation as TRW, different algorithm.

Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

Experiments

Conclusions

Conclusions

Motivation

- TRW Convergence rates can be very slow.
 - If lucky, TRW = block coordinate ascent on dual.
- TRW may fail to converge.
 - Damping converges in practice, slower.
 - Recent alternatives guarantee convergence.
[Hazan & Shashua 2009, Globerson & Jaakkola 2007b]
 - Not claimed faster than TRW. TRW-S [Meltzer et al. 2009] is an exception.
- This paper: use a quasi-newton method on dual.
 - Line searches guarantee convergence.
 - Hopefully, faster convergence.

Motivation

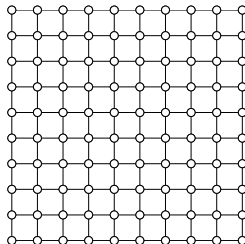
- TRW Convergence rates can be very slow.
 - If lucky, TRW = block coordinate ascent on dual.
- TRW may fail to converge.
 - Damping converges in practice, slower.
 - Recent alternatives guarantee convergence.
[Hazan & Shashua 2009, Globerson & Jaakkola 2007b]
 - Not claimed faster than TRW. TRW-S [Meltzer et al. 2009] is an exception.
- This paper: use a quasi-newton method on dual.
 - Line searches guarantee convergence.
 - Hopefully, faster convergence.

Motivation

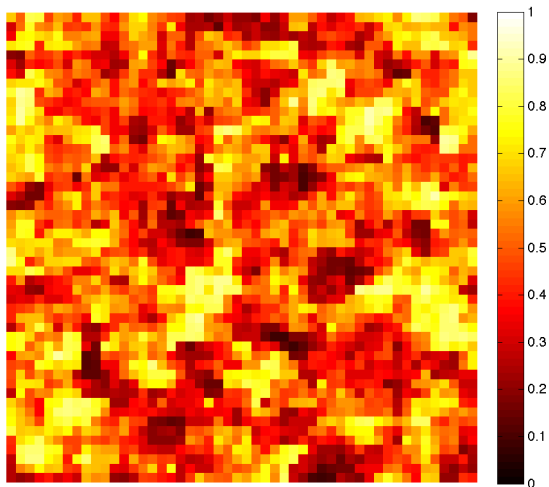
- TRW Convergence rates can be very slow.
 - If lucky, TRW = block coordinate ascent on dual.
- TRW may fail to converge.
 - Damping converges in practice, slower.
 - Recent alternatives guarantee convergence.
[Hazan & Shashua 2009, Globerson & Jaakkola 2007b]
 - Not claimed faster than TRW. TRW-S [Meltzer et al. 2009] is an exception.
- This paper: use a quasi-newton method on dual.
 - Line searches guarantee convergence.
 - Hopefully, faster convergence.

Ising Model

- $x_i \in \{-1, +1\}$
- $p(\mathbf{x}) \propto \prod_{ij} \exp(\theta(x_i, x_j)) \prod_i \exp(\theta(x_i))$
- $\theta(x_i) = \alpha_F x_i, \quad \alpha_F \in [-1, +1]$
- $\theta(x_i, x_j) = \alpha_I x_i x_j, \quad \alpha_I \in [0, T]$ for various T

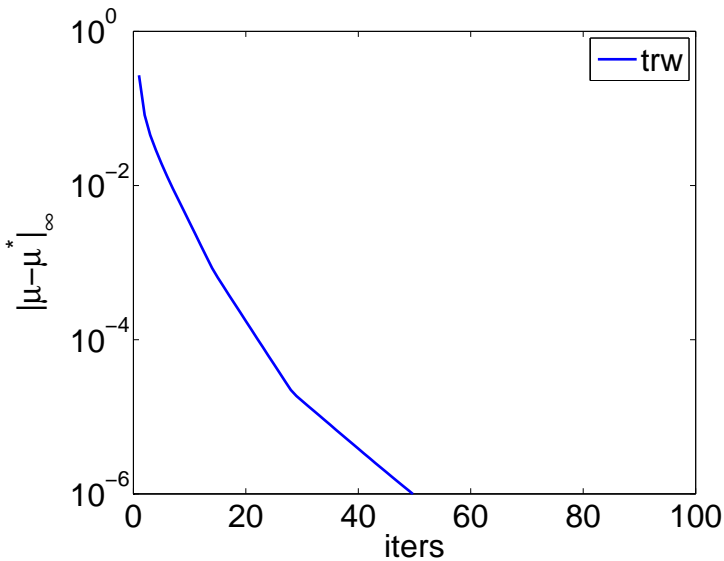


$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 1]$$



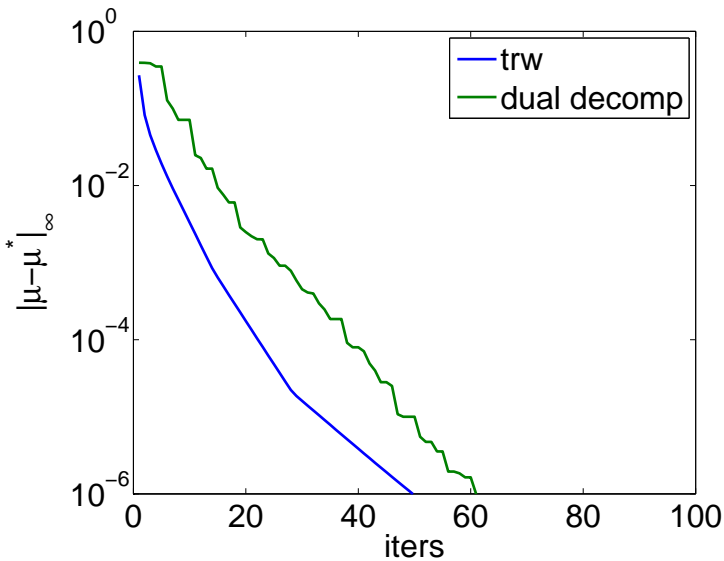


$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 1]$$



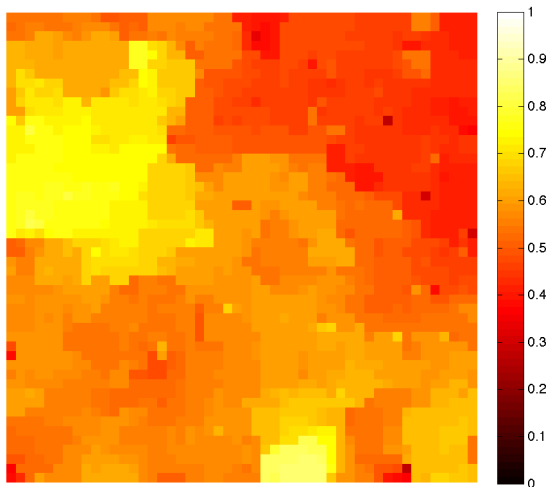


$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 1]$$

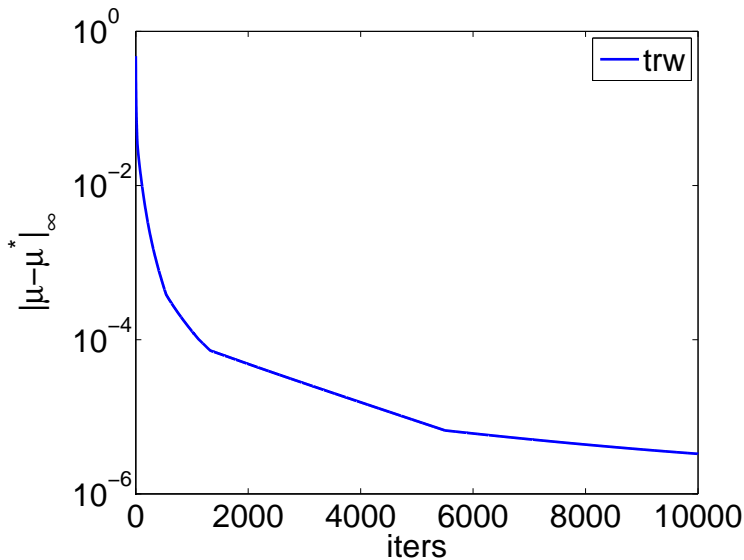


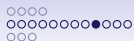


$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 3]$$

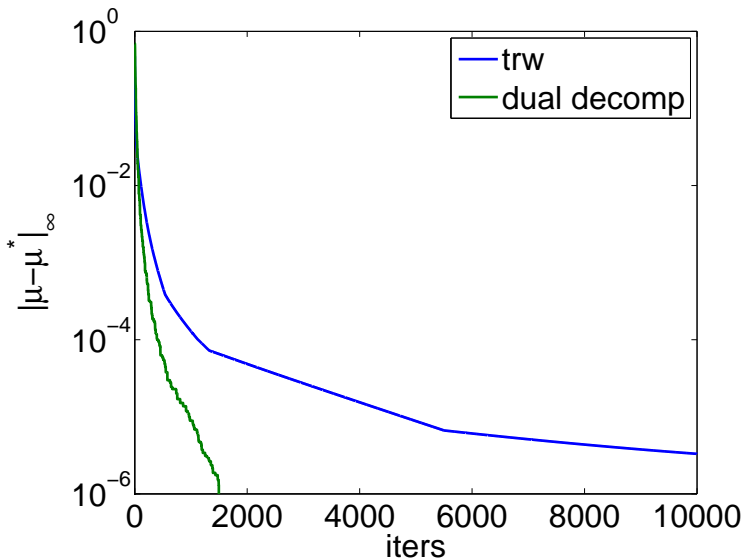


$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 3]$$

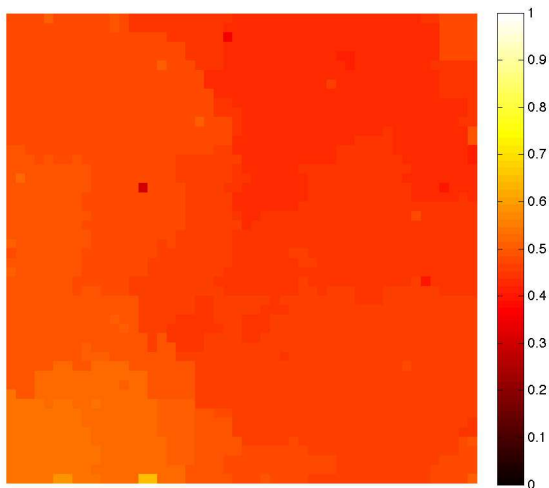




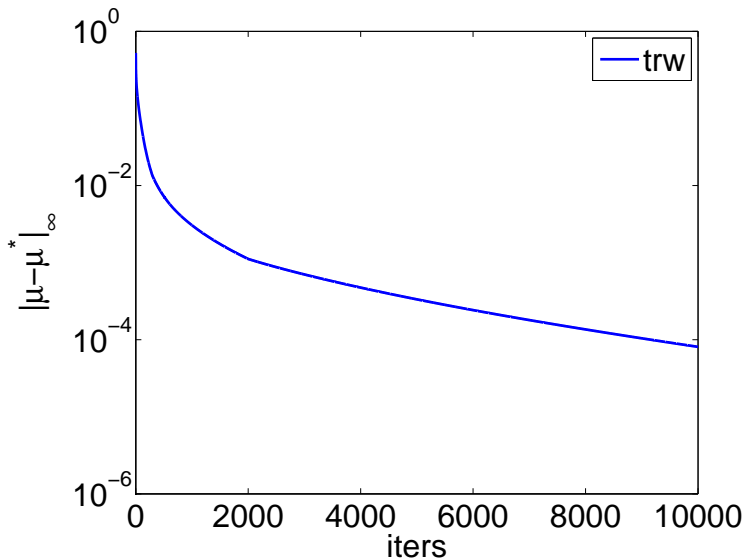
$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 3]$$

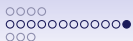


$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 5]$$

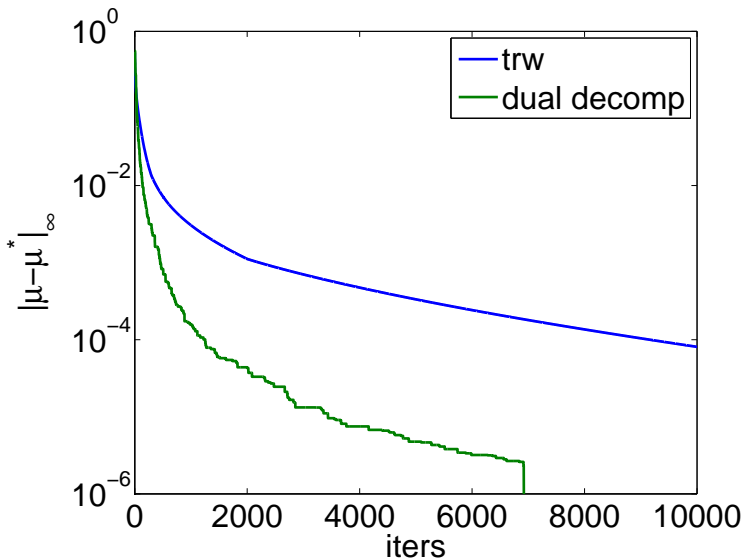


$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 5]$$





$$\theta(x_i, x_j) = \alpha_l x_i x_j, \quad \alpha_l \in [0, 5]$$



Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

Experiments

Conclusions

Conclusions

Wait a Second

Question: Why should I care about very accurately computing approximate marginals!?

Answer: You might not.

One reason to care:

- Number of iterations TRW needs for reasonable results is not easy to predict.

Wait a Second

Question: Why should I care about very accurately computing approximate marginals!?

Answer: You might not.

One reason to care:

- Number of iterations TRW needs for reasonable results is not easy to predict.

Why I Care

Want to fit a CRF with some loss $L(\theta) = M(\mu(\theta))$.

Algorithm (Domke, 2010):

1. Get μ by running TRW with parameters θ .
2. Compute $\frac{dM(\mu)}{d\mu}$
3. Get μ^+ by running TRW with parameters $\theta + r \frac{dM}{d\mu}$
4. $\frac{dL}{d\theta} \approx \frac{1}{r} (\mu^+ - \mu)$

Strong convergence needed for difference $\mu^+ - \mu$ to be meaningful.

Why I Care

Want to fit a CRF with some loss $L(\theta) = M(\mu(\theta))$.

Algorithm (Domke, 2010):

1. Get μ by running TRW with parameters θ .
2. Compute $\frac{dM(\mu)}{d\mu}$
3. Get μ^+ by running TRW with parameters $\theta + r \frac{dM}{d\mu}$
4. $\frac{dL}{d\theta} \approx \frac{1}{r} (\mu^+ - \mu)$

Strong convergence needed for difference $\mu^+ - \mu$ to be meaningful.

Why I Care

Want to fit a CRF with some loss $L(\theta) = M(\mu(\theta))$.

Algorithm (Domke, 2010):

1. Get μ by running TRW with parameters θ .
2. Compute $\frac{dM(\mu)}{d\mu}$
3. Get μ^+ by running TRW with parameters $\theta + r \frac{dM}{d\mu}$
4. $\frac{dL}{d\theta} \approx \frac{1}{r} (\mu^+ - \mu)$

Strong convergence needed for difference $\mu^+ - \mu$ to be meaningful.

Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

Experiments

Conclusions

Conclusions

Dual Decomposition with Two subproblems

$$\max_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x})$$

- Can quickly and exactly maximize $f(\mathbf{x}) + \mathbf{a} \cdot \mathbf{x}$.
- Can quickly and exactly maximize $g(\mathbf{x}) + \mathbf{b} \cdot \mathbf{x}$.

Dual Decomposition with Two subproblems

- Transform $\max_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x})$ to a constrained problem:

$$\begin{aligned} \max_{\mathbf{x}, \mathbf{y}} \quad & f(\mathbf{x}) + g(\mathbf{y}) \\ \text{s.t.} \quad & \mathbf{x} = \mathbf{y} \end{aligned}$$

- Leads to dual problem:

$$\begin{aligned} \min_{\lambda} h(\lambda), \quad h(\lambda) = & \max_{\mathbf{x}} f(\mathbf{x}) + \lambda \cdot \mathbf{x} \\ & + \max_{\mathbf{y}} g(\mathbf{y}) - \lambda \cdot \mathbf{y} \end{aligned}$$

Dual Decomposition with Two subproblems

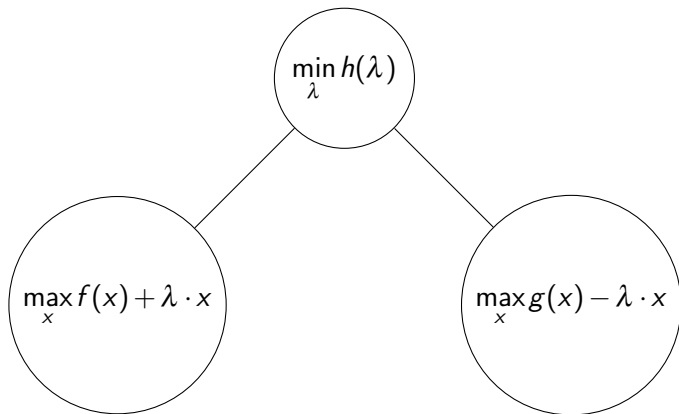
- Transform $\max_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x})$ to a constrained problem:

$$\begin{array}{ll} \max_{\mathbf{x}, \mathbf{y}} & f(\mathbf{x}) + g(\mathbf{y}) \\ \text{s.t.} & \mathbf{x} = \mathbf{y} \end{array}$$

- Leads to dual problem:

$$\begin{aligned} \min_{\lambda} h(\lambda), \quad h(\lambda) &= \max_{\mathbf{x}} f(\mathbf{x}) + \lambda \cdot \mathbf{x} \\ &+ \max_{\mathbf{y}} g(\mathbf{y}) - \lambda \cdot \mathbf{y} \end{aligned}$$

Dual Decomposition with Two subproblems



Dual Decomposition with N subproblems

$$\max_{\mathbf{x}} \sum_{i=1}^N f_i(\mathbf{x})$$

- Can quickly and exactly maximize $f_i(\mathbf{x}) + \mathbf{a}_i \cdot \mathbf{x}$, for all i .

Dual Decomposition with N subproblems

- Transform $\max_{\mathbf{x}} \sum_i f_i(\mathbf{x}_i)$ to a constrained problem:

$$\begin{aligned} \max_{\{\mathbf{x}_i\}} \quad & \sum_i f_i(\mathbf{x}_i) \\ \text{s.t.} \quad & \mathbf{x}_i = \frac{1}{N} \sum_j \mathbf{x}_j \end{aligned}$$

- Leads to dual problem:

$$\min_{\lambda} h(\lambda), \quad h(\lambda) = \sum_i h_i(\lambda)$$

$$h_i(\lambda) = \max_{\mathbf{x}_i} f_i(\mathbf{x}_i) + (\lambda_i - \frac{1}{N} \sum_j \lambda_j) \cdot \mathbf{x}_i$$

Dual Decomposition with N subproblems

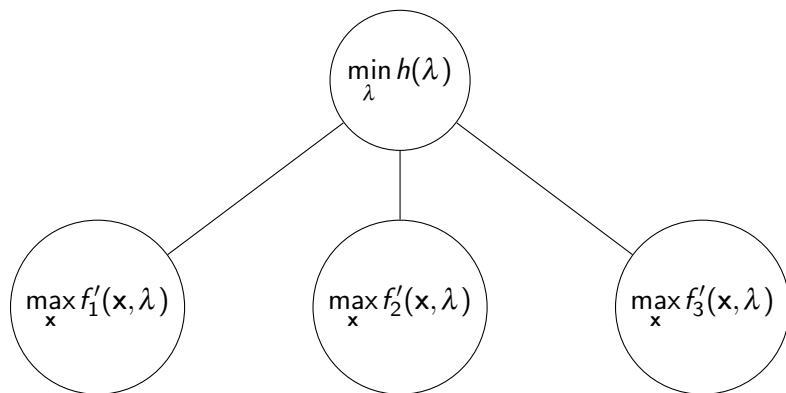
- Transform $\max_{\mathbf{x}} \sum_i f_i(\mathbf{x}_i)$ to a constrained problem:

$$\begin{aligned} \max_{\{\mathbf{x}_i\}} \quad & \sum_i f_i(\mathbf{x}_i) \\ \text{s.t.} \quad & \mathbf{x}_i = \frac{1}{N} \sum_j \mathbf{x}_j \end{aligned}$$

- Leads to dual problem:

$$\begin{aligned} \min_{\lambda} h(\lambda), \quad h(\lambda) &= \sum_i h_i(\lambda) \\ h_i(\lambda) &= \max_{\mathbf{x}_i} f_i(\mathbf{x}_i) + (\lambda_i - \frac{1}{N} \sum_j \lambda_j) \cdot \mathbf{x}_i \end{aligned}$$

Dual Decomposition with N subproblems



$$f'_i(\mathbf{x}, \lambda) = f_i(\mathbf{x}_i) + (\lambda_i - \frac{1}{N} \sum_j \lambda_j) \cdot \mathbf{x}_i$$

Dual Decomposition with N subproblems

- Has been used extensively for MAP inference.
 - $h(\lambda)$ is non-differentiable.
- For marginal inference, $h(\lambda)$ is differentiable, convex.

Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

Experiments

Conclusions

Conclusions

Variational Inference

Can represent a graphical model in exponential family:

$$p(\mathbf{x}; \theta) = \exp(\mathbf{f}(\mathbf{x}) \cdot \theta - A(\theta)), \quad A(\theta) = \log \sum_{\mathbf{x}} \exp(\mathbf{f}(\mathbf{x}) \cdot \theta)$$

Can compute A as [Wainwright and Jordan]

$$A(\theta) = \max_{\mu \in \mathcal{M}} \theta \cdot \mu + H(\mu)$$

- \mathcal{M} is marginal polytope (hard).
- H is entropy (hard).

Variational Inference

Exact inference: $A(\theta) = \max_{\mu \in \mathcal{M}} \theta \cdot \mu + H(\mu)$

TRW approximation: $B(\theta) = \max_{\mu \in \mathcal{L}} \theta \cdot \mu + \sum_T \rho_T H(\mu(T))$

- \mathcal{L} - is marginal polytope (easy)
- $H(\mu(T))$ - entropy of marginals projected onto tree T (easy)

Our problem: how to compute B ?

Variational Inference

Exact inference: $A(\theta) = \max_{\mu \in \mathcal{M}} \theta \cdot \mu + H(\mu)$

TRW approximation: $B(\theta) = \max_{\mu \in \mathcal{L}} \theta \cdot \mu + \sum_T \rho_T H(\mu(T))$

- \mathcal{L} - is marginal polytope (easy)
- $H(\mu(T))$ - entropy of marginals projected onto tree T (easy)

Our problem: how to compute B ?

Variational Inference

Exact inference: $A(\theta) = \max_{\mu \in \mathcal{M}} \theta \cdot \mu + H(\mu)$

TRW approximation: $B(\theta) = \max_{\mu \in \mathcal{L}} \theta \cdot \mu + \sum_T \rho_T H(\mu(T))$

- \mathcal{L} - is marginal polytope (easy)
- $H(\mu(T))$ - entropy of marginals projected onto tree T (easy)

Our problem: how to compute B ?

Variational Inference

Exact inference: $A(\theta) = \max_{\mu \in \mathcal{M}} \theta \cdot \mu + H(\mu)$

TRW approximation: $B(\theta) = \max_{\mu \in \mathcal{L}} \theta \cdot \mu + \sum_T \rho_T H(\mu(T))$

- \mathcal{L} - is marginal polytope (easy)
- $H(\mu(T))$ - entropy of marginals projected onto tree T (easy)

Our problem: how to compute B ?

Dual Decomposition for Marginal Inference

TRW approximation: $B(\theta) = \max_{\mu \in \mathcal{L}} \theta \cdot \mu + \sum_T \rho_T H(\mu(T))$

Theorem (main result):

$$B(\theta) = \min_{\{\theta^T\}} h(\{\theta^T\}) \quad \text{s.t.} \quad \sum_{T:a \in T} \theta_a^T = \theta_a$$

$$h(\{\theta^T\}) = \sum_T B_T(\theta^T)$$

$$B_T(\theta^T) = \max_{\mu^T \in \mathcal{M}_T} \theta^T \cdot \mu^T + \rho_T H_T(\mu^T)$$

$B_T(\theta^T)$ is computable by running regular sum-product algorithm.

Dual Decomposition for Marginal Inference

TRW approximation: $B(\theta) = \max_{\mu \in \mathcal{L}} \theta \cdot \mu + \sum_T \rho_T H(\mu(T))$

Theorem (main result):

$$B(\theta) = \min_{\{\theta^T\}} h(\{\theta^T\}) \quad \text{s.t.} \quad \sum_{T: a \in T} \theta_a^T = \theta_a$$

$$h(\{\theta^T\}) = \sum_T B_T(\theta^T)$$

$$B_T(\theta^T) = \max_{\mu^T \in \mathcal{M}_T} \theta^T \cdot \mu^T + \rho_T H_T(\mu^T)$$

$B_T(\theta^T)$ is computable by running regular sum-product algorithm.

Dual Decomposition for Marginal Inference

TRW approximation: $B(\theta) = \max_{\mu \in \mathcal{L}} \theta \cdot \mu + \sum_T \rho_T H(\mu(T))$

Theorem (main result):

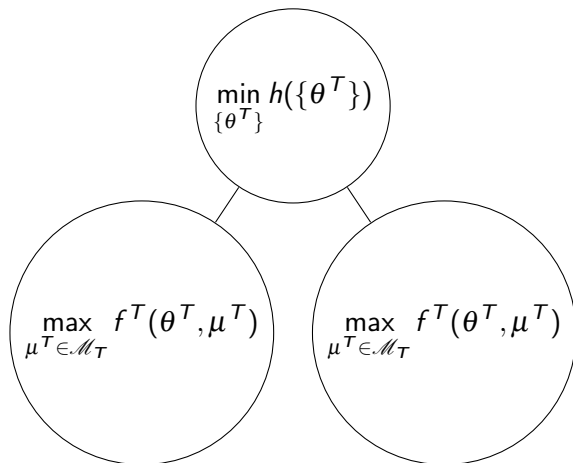
$$B(\theta) = \min_{\{\theta^T\}} h(\{\theta^T\}) \quad \text{s.t.} \quad \sum_{T: a \in T} \theta_a^T = \theta_a$$

$$h(\{\theta^T\}) = \sum_T B_T(\theta^T)$$

$$B_T(\theta^T) = \max_{\mu^T \in \mathcal{M}_T} \theta^T \cdot \mu^T + \rho_T H_T(\mu^T)$$

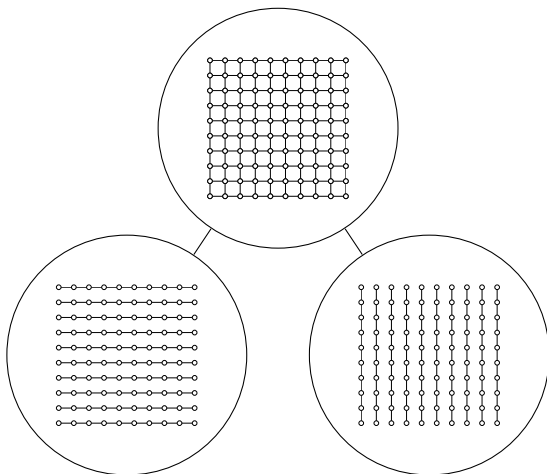
$B_T(\theta^T)$ is computable by running regular sum-product algorithm.

Dual Decomposition for Marginal Inference



$$f^T(\theta^T, \mu^T) = \theta^T \cdot \mu^T + \rho_T H_T(\mu^T)$$

Dual Decomposition for Marginal Inference



Dual Decomposition for Marginal Inference

Inference: Plug $\min_{\{\theta^T\}} \sum_T B_T(\theta^T)$ into L-BFGS.

- Guarantees convergence. (Line searches)
- Fast convergence rates.

Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

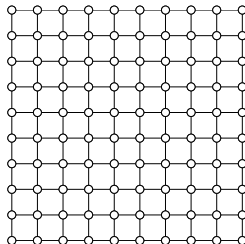
Experiments

Conclusions

Conclusions

Ising Model

- $x_i \in \{-1, +1\}$
- $p(\mathbf{x}) \propto \prod_{ij} \exp(\theta(x_i, x_j)) \prod_i \exp(\theta(x_i))$
- $\theta(x_i) = \alpha_F x_i$
- $\theta(x_i, x_j) = \alpha_I x_i x_j$



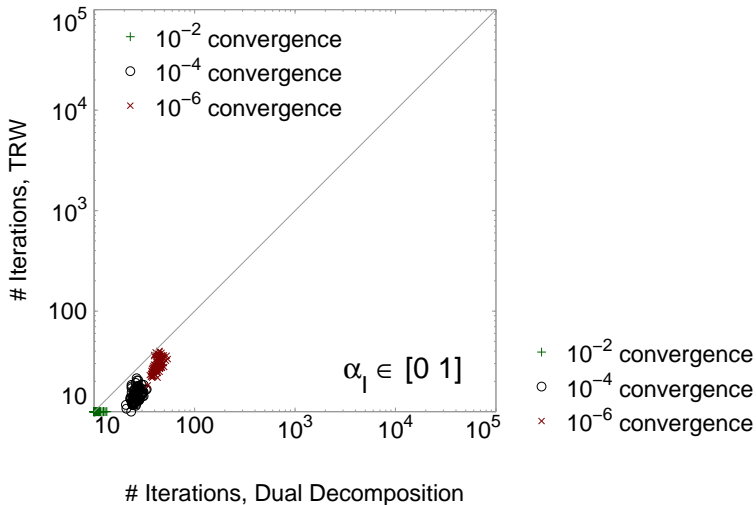
Algorithms

Algorithms Compared:

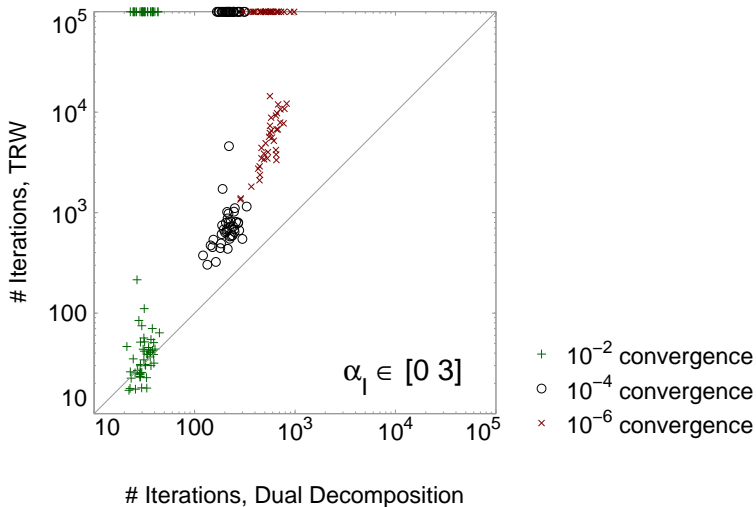
- Dual Decomposition + L-BFGS
- TRW
- TRW with damping of 1/2 in the log-domain.
- TRW-S [Meltzer et al. 2009]

Max of 10^5 iterations allowed.

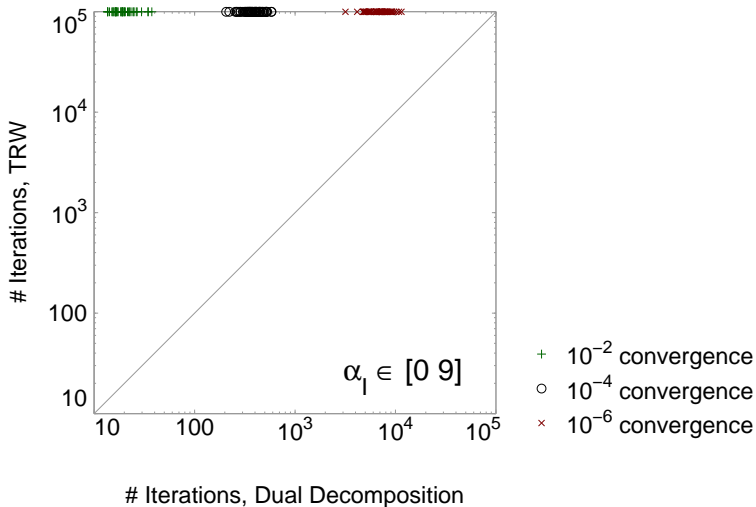
Dual Decomposition vs. TRW



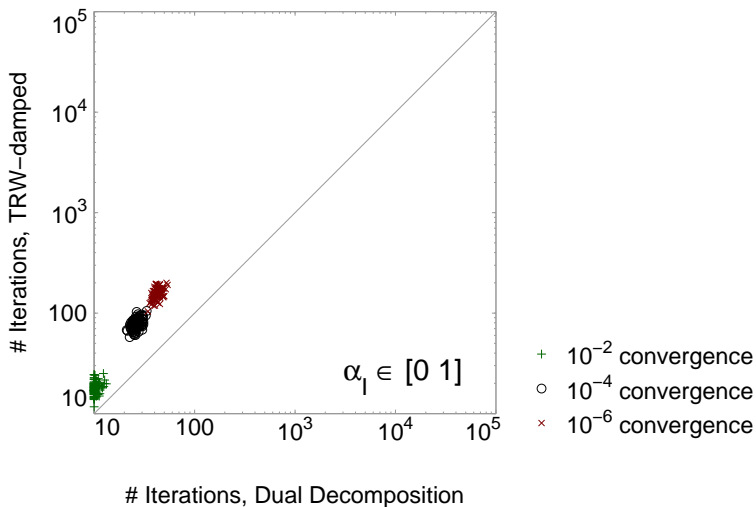
Dual Decomposition vs. TRW



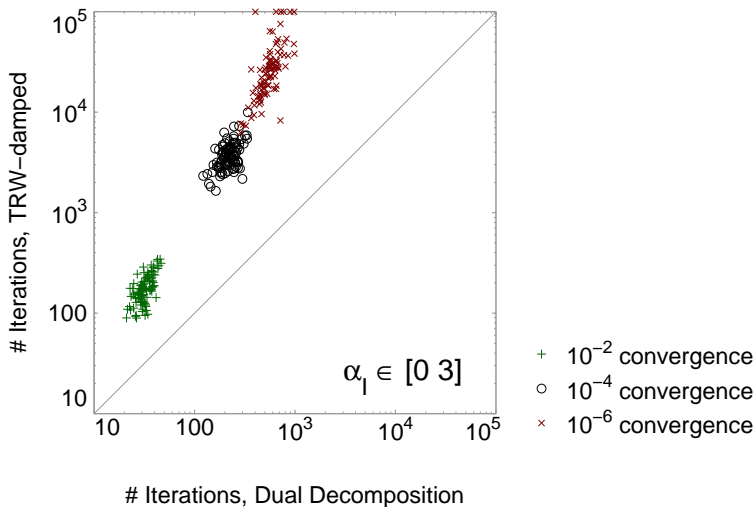
Dual Decomposition vs. TRW



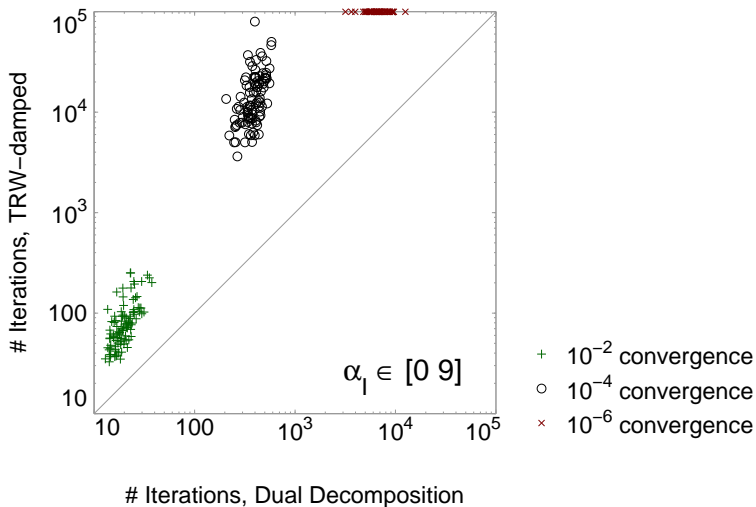
Dual Decomposition vs. TRW-damped



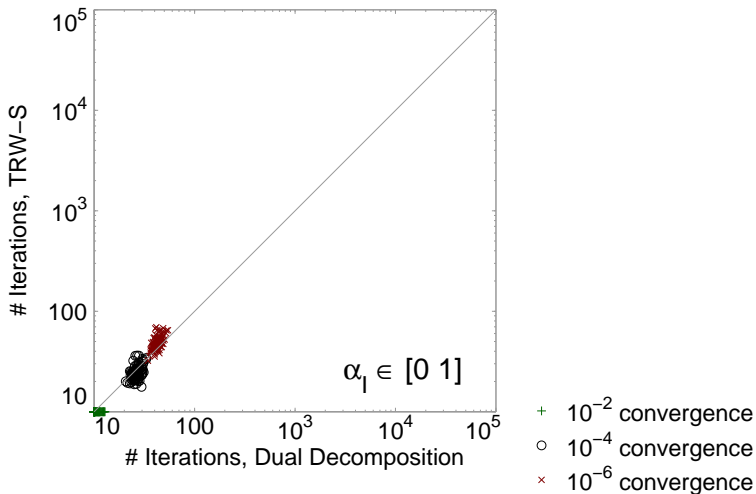
Dual Decomposition vs. TRW-damped



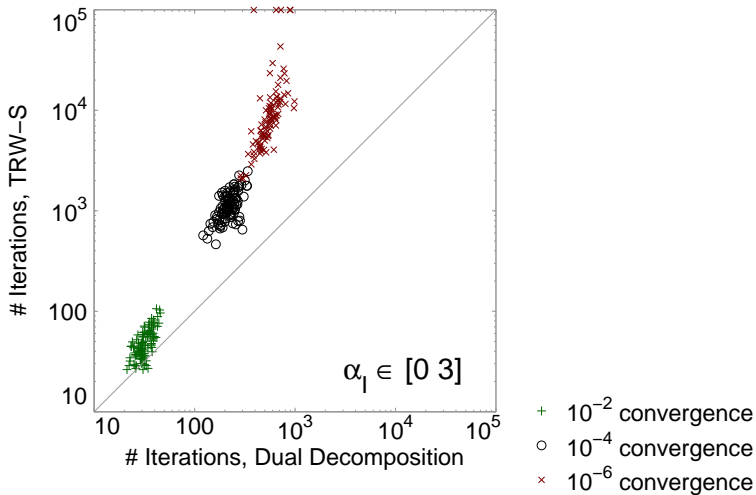
Dual Decomposition vs. TRW-damped



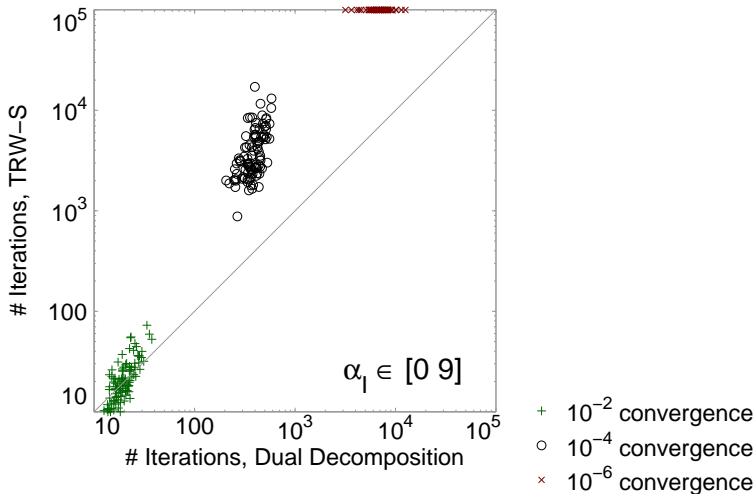
Dual Decomposition vs. TRW-S



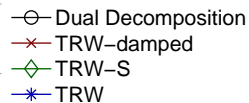
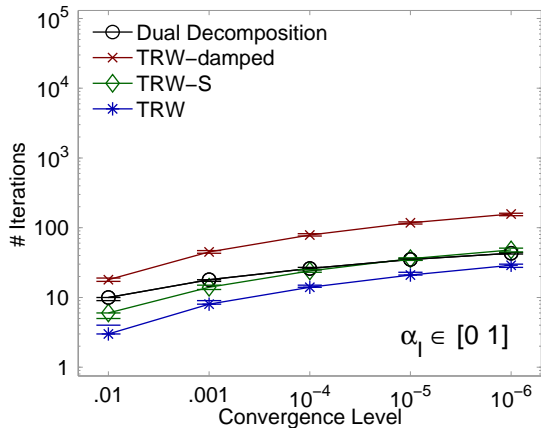
Dual Decomposition vs. TRW-S



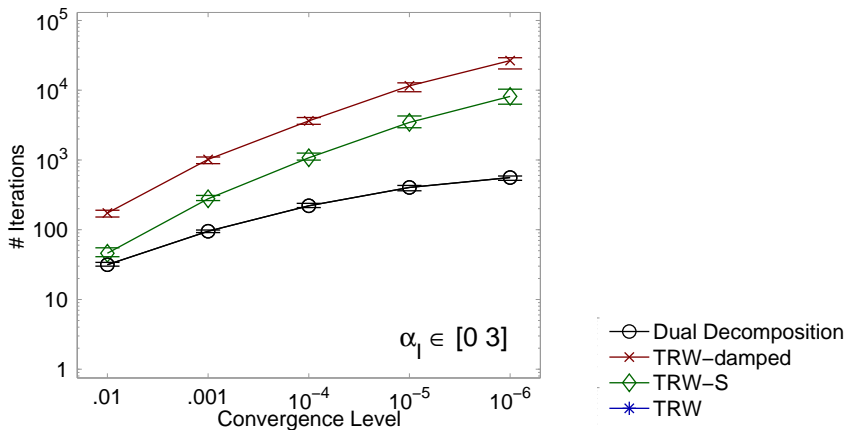
Dual Decomposition vs. TRW-S



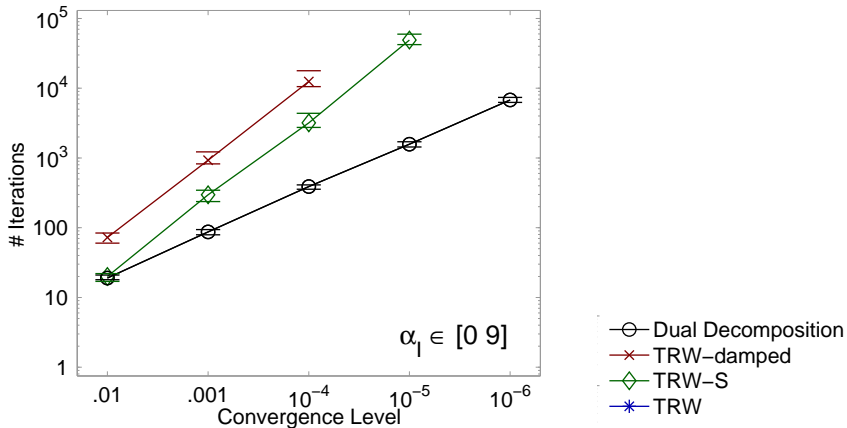
Convergence



Convergence



Convergence



Outline

Introduction

Graphical Models

Motivation

Wait a Second

Dual Decomposition

Dual Decomposition in General

Dual Decomposition for Marginal Inference

Experimental Results

Experiments

Conclusions

Conclusions

Conclusions

- Dual Decomposition
 - Faster on “hard” problems or if strong convergence needed.
- Caveats
 - Not really faster on “easy” problems.
 - Restriction on tree distribution $P(T)$.

Conclusions

- Dual Decomposition
 - Faster on “hard” problems or if strong convergence needed.
- Caveats
 - Not really faster on “easy” problems.
 - Restriction on tree distribution $P(T)$.

Conclusions

Thank you

Graphical models toolbox: people.rit.edu/jcdicsa/
(Dual decomposition coming soon.)