# Graphical models, message-passing algorithms, and variational methods: Part II

Martin Wainwright

Department of Statistics, and

Department of Electrical Engineering and Computer Science,

UC Berkeley, Berkeley, CA USA


*Email:* `wainwrig@{stat,eecs}.berkeley.edu`

For further information (tutorial slides, films of course lectures), see:

`www.eecs.berkeley.edu/~wainwrig/`

# Introduction

- graphical models are used and studied in various applied statistical and computational fields:

  - machine learning and artificial intelligence

  - computational biology

  - statistical signal/image processing

  - communication and information theory

  - statistical physics

  - .....

- based on correspondences between graph theory and probability theory

- important but difficult problems:

  - computing likelihoods, marginal distributions, modes

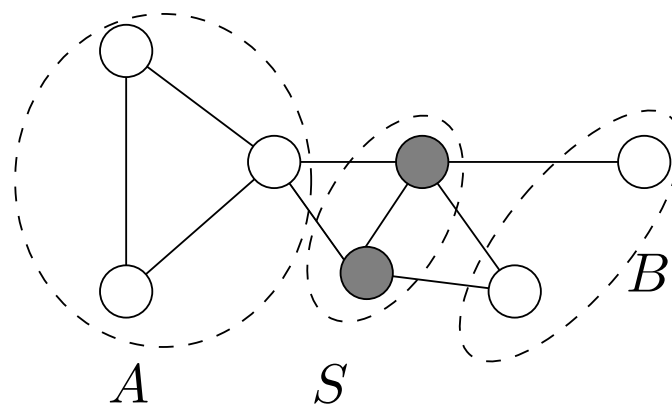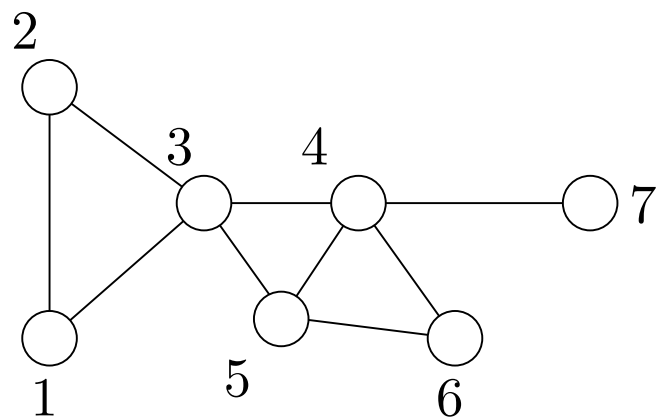  - estimating model parameters and structure from (noisy) data

# Outline

1. Recap

    (a) Background on graphical models

    (b) Some applications and challenging problems

    (c) Illustrations of some message-passing algorithms

2. Exponential families and variational methods

    (a) What is a variational method (and why should I care)?

    (b) Graphical models as exponential families

    (c) Variational representations from conjugate duality

3. Exact techniques as variational methods

    (a) Gaussian inference on arbitrary graphs

    (b) Belief-propagation/sum-product on trees (e.g., Kalman filter; $\alpha$-$\beta$ alg.)

    (c) Max-product on trees (e.g., Viterbi)

4. Approximate techniques as variational methods

    (a) Mean field and variants

    (b) Belief propagation and extensions on graphs with cycles

    (c) Convex methods and upper bounds

    (d) Tree-reweighted max-product and linear programs

# Undirected graphical models

Based on correspondences between graphs and random variables.

- given an undirected graph $G = (V, E)$, each node $s$ has an associated random variable $X_s$

- for each subset $A \subseteq V$, define $X_A := \{X_s, s \in A\}$.

g replacements



PSfrag replacements

Maximal cliques $(123), (345), (456), (47)$        Vertex cutset $S$

- a *clique* $C \subseteq V$ is a subset of vertices all joined by edges

- a *vertex cutset* is a subset $S \subset V$ whose removal breaks the graph into two or more pieces

# Factorization and Markov properties

The graph $G$ can be used to impose constraints on the random vector $X = X_V$ (or on the distribution $p$) in different ways.

**Markov property:** $X$ is *Markov w.r.t $G$* if $X_A$ and $X_B$ are conditionally indpt. given $X_S$ whenever $S$ separates $A$ and $B$.

**Factorization:** The distribution $p$ *factorizes according to $G$* if it can be expressed as a product over cliques:

$$p(\mathbf{x}) \quad = \quad \frac{1}{Z} \prod_{C \in \mathcal{C}} \underbrace{\exp\left\{\theta_C(x_C)\right\}}$$

compatibility function on clique $C$

**Theorem:** (Hammersley-Clifford) For strictly positive $p(\cdot)$, the Markov property and the Factorization property are equivalent.

# Challenging computational problems

Frequently, it is of interest to compute various quantities associated with an undirected graphical model:

(a) the log normalization constant $\log Z$

(b) local marginal distributions or other local statistics

(c) modes or most probable configurations

Relevant dimensions often grow rapidly in graph size $\implies$ major computational challenges.

**Example:** Consider a naive approach to computing the normalization constant for binary random variables:

$$Z = \sum_{\mathbf{x} \in \{0,1\}^n} \prod_{C \in \mathcal{C}} \exp\left\{\theta_C(x_C)\right\}$$

Complexity scales exponentially as $2^n$.

# Gibbs sampling in the Ising model

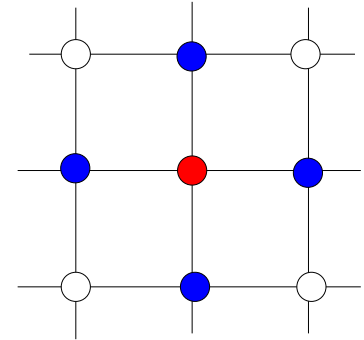- binary variables on a graph $G = (V, E)$ with pairwise interactions:

$$p(\mathbf{x}; \theta) \quad \propto \quad \exp\Big\{ \sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t \Big\}$$

- Update $x_s^{(m+1)}$ stochastically based on values $x_{\mathcal{N}(s)}^{(m)}$ at neighbors:

1. Choose $s \in V$ at random.

2. Sample $u \sim \mathcal{U}(0,1)$ and update

$$x_s^{(m+1)} = \begin{cases} 1 & \text{if } u \leq \{1 + \exp[-(\theta_s + \sum_{t \in \mathcal{N}(s)} \theta_{st} x_t^{(m)})]\}^{-1} \\ 0 & \text{otherwise} \end{cases},$$

- sequence $\{\mathbf{x}^{(m)}\}$ converges (in a stochastic sense) to a sample from $p(\mathbf{x}; \theta)$
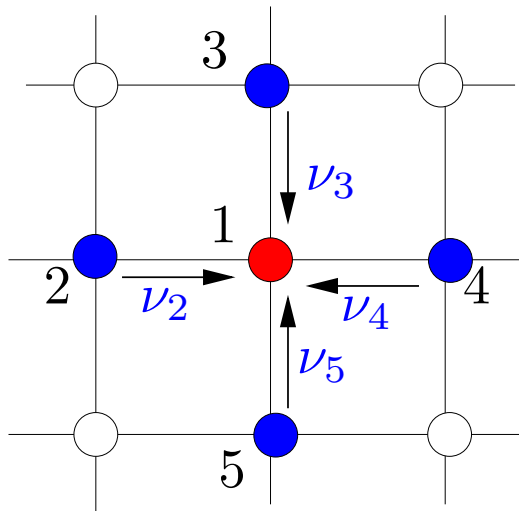
# Mean field updates in the Ising model

- binary variables on a graph $G = (V, E)$ with pairwise interactions:

$$p(\mathbf{x}; \theta) \quad \propto \quad \exp\left\{ \sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t \right\}$$

- simple (deterministic) message-passing algorithm involving
  *variational parameters* $\nu_s \in (0,1)$ at each node

lacements



1. Choose $s \in V$ at random.

2. Update $\nu_s$ based on neighbors
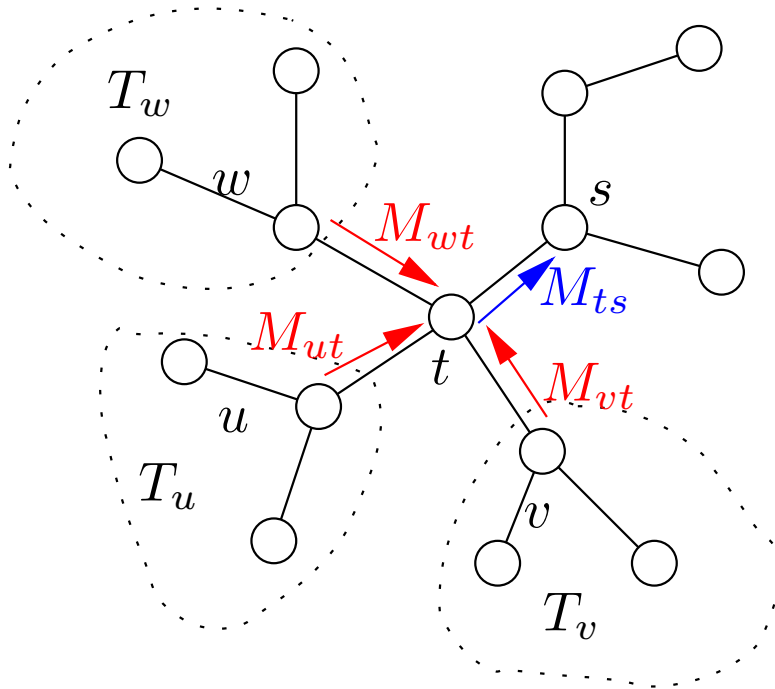   $\{\nu_t, \ t \in \mathcal{N}(s)\}$:

$$\nu_s \quad \longleftarrow \quad \left\{ 1 + \exp\left[ -(\theta_s + \sum_{t \in \mathcal{N}(s)} \theta_{st} \nu_t) \right] \right\}^{-1}$$

## Questions:

- principled derivation?       • convergence and accuracy?

8

# **Sum and max-product algorithms: On trees**

Exact for trees, but approximate for graphs with cycles.



$M_{ts} \quad \equiv \quad$ message from node $t$ to $s$

$\Gamma(t) \quad \equiv \quad$ neighbors of node $t$

Sum-product: for marginals

(generalizes $\alpha - \beta$ algorithm; Kalman filter)
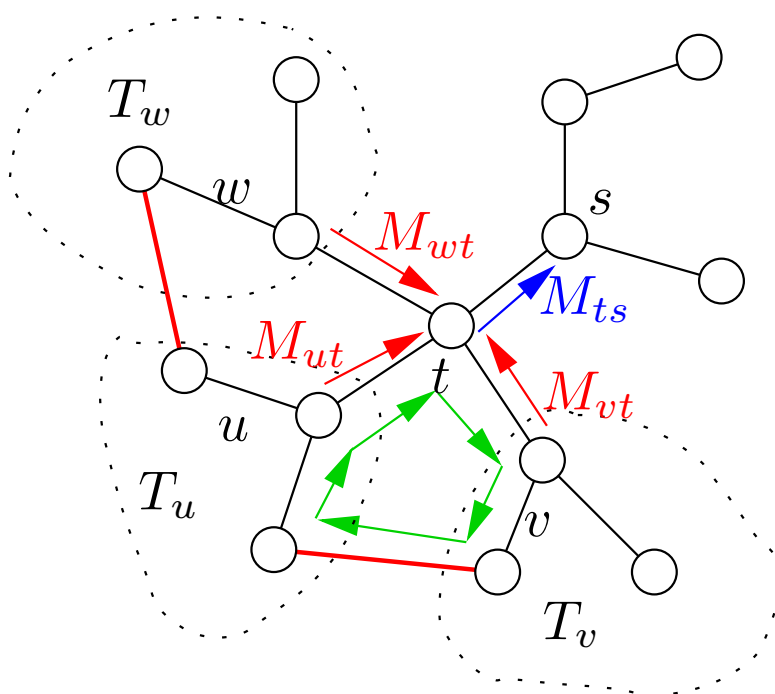
Max-product: for MAP configurations

(generalizes Viterbi algorithm)

Update: $\quad \mathbf{M_{ts}(x_s)} \leftarrow \sum_{x'_t \in \mathcal{X}_t} \left\{ \exp\left[ \theta_{st}(x_s, x'_t) + \theta_t(x'_t) \right] \prod_{v \in \Gamma(t) \setminus s} \mathbf{M_{vt}(x_t)} \right\}$

Marginals: $\quad p(x_s; \theta) \propto \exp\{\theta_t(x_t)\} \prod_{t \in \Gamma(s)} M_{ts}(x_s).$

# Sum and max-product: On graphs with cycles

- what about applying same updates on graph with cycles?

- updates need not converge (effect of cycles)

- seems naive, but remarkably successful in many applications



acements

$M_{ts} \equiv$ message from node $t$ to $s$

$\Gamma(t) \equiv$ neighbors of node $t$

Sum-product: for marginals

Max-product: for modes

**Questions:** • meaning of these updates for graphs with cycles?

• convergence? accuracy of resulting "marginals"?

# Outline

1. Recap

    (a) Background on graphical models

    (b) Some applications and challenging problems

    (c) Illustrations of some message-passing algorithms

2. Exponential families and variational methods

    (a) What is a variational method (and why should I care)?

    (b) Graphical models as exponential families

    (c) Variational representations from conjugate duality

3. Exact techniques as variational methods

    (a) Gaussian inference on arbitrary graphs

    (b) Belief-propagation/sum-product on trees (e.g., Kalman filter; $\alpha$-$\beta$ alg.)

    (c) Max-product on trees (e.g., Viterbi)

4. Approximate techniques as variational methods

    (a) Mean field and variants

    (b) Belief propagation and extensions

    (c) Convex methods and upper bounds

    (d) Tree-reweighted max-product and linear programs

# Variational methods

- *"variational"*: umbrella term for optimization-based formulation of problems, and methods for their solution

- historical roots in the calculus of variations

- modern variational methods encompass a wider class of methods (e.g., dynamic programming; finite-element methods)

**Variational principle:** Representation of a quantity of interest $\widehat{\mathbf{u}}$ as the solution of an optimization problem.

1. allows the quantity $\widehat{\mathbf{u}}$ to be studied through the lens of the optimization problem

2. approximations to $\widehat{\mathbf{u}}$ can be obtained by approximating or relaxing the variational principle

# Illustration: A simple variational principle

*Goal:* Given a vector $\mathbf{y} \in \mathbb{R}^n$ and a symmetric matrix $Q \succ 0$, solve the linear system $Q\mathbf{u} = \mathbf{y}$.

*Unique solution* $\widehat{\mathbf{u}}(\mathbf{y}) = Q^{-1}\mathbf{y}$ can be obtained by matrix inversion.

*Variational formulation:* Consider the function $J_{\mathbf{y}} : \mathbb{R}^n \to \mathbb{R}$ defined by

$$J_{\mathbf{y}}(\mathbf{u}) \quad := \quad \frac{1}{2}\mathbf{u}^T Q\mathbf{u} - \mathbf{y}^T\mathbf{u}.$$

It is strictly convex, and the minimum is uniquely attained:

$$\widehat{\mathbf{u}}(\mathbf{y}) \;=\; \arg\min_{\mathbf{u}\in\mathbb{R}^n} J_{\mathbf{y}}(\mathbf{u}) \;=\; Q^{-1}\mathbf{y}.$$

Various methods for solving linear systems (e.g., conjugate gradient) exploit this variational representation.

# Useful variational principles for graphical models?

Consider an undirected graphical model:

$$p(\mathbf{x}) \quad = \quad \frac{1}{Z} \prod_{C \in \mathbf{C}} \exp\{\theta_C(x_C)\}.$$

Core problems that arise in many applications:

(a) computing the log normalization constant $\log Z$

(b) computing local marginal distributions (e.g., $p(x_s) = \sum_{x_t, t \neq s} p(\mathbf{x})$)

(c) computing modes or most likely configurations $\widehat{\mathbf{x}} \in \arg\max_{\mathbf{x}} p(\mathbf{x})$

**Approach:** Develop variational representations of all of these problems by exploiting results from:

(a) exponential families                                                      (e.g.,Brown, 1986)

(b) convex duality                                                      (e.g., Rockafellar,1973)

# Maximum entropy formulation of graphical models

- suppose that we have measurements $\widehat{\mu}$ of the average values of some (local) functions $\phi_\alpha : \mathcal{X}^n \to \mathbb{R}$

- in general, will be many distributions $p$ that satisfy the measurement constraints $\mathbb{E}_p[\phi_\alpha(\mathbf{x})] = \widehat{\mu}$

- will consider finding the $p$ with maximum "uncertainty" subject to the observations, with uncertainty measured by entropy

$$H(p) = -\sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}).$$

**Constrained maximum entropy problem:** Find $\widehat{p}$ to solve

$$\max_{p \in \mathcal{P}} H(p) \qquad \text{such that} \qquad \mathbb{E}_p[\phi_\alpha(\mathbf{x})] = \widehat{\mu}$$

- elementary argument with Lagrange multipliers shows that solution takes the exponential form

$$\widehat{p}(\mathbf{x}; \theta) \propto \exp\Big\{ \sum_{\alpha \in \mathcal{I}} \theta_\alpha \phi_\alpha(\mathbf{x}) \Big\}.$$

# Exponential families

$$\phi_\alpha : \mathcal{X}^n \to \mathbb{R} \qquad \equiv \qquad \text{sufficient statistic}$$

$$\boldsymbol{\phi} = \{\phi_\alpha, \alpha \in \mathcal{I}\} \qquad \equiv \qquad \text{vector of sufficient statistics}$$

$$\theta = \{\theta_\alpha, \alpha \in \mathcal{I}\} \qquad \equiv \qquad \text{parameter vector}$$

$$\boldsymbol{\nu} \qquad \equiv \qquad \text{base measure (e.g., Lebesgue, counting)}$$

- parameterized family of densities (w.r.t. $\boldsymbol{\nu}$):

$$p(\mathbf{x}; \theta) \;=\; \exp\Big\{ \sum_\alpha \theta_\alpha \phi_\alpha(\mathbf{x}) \;-\; A(\theta)\Big\}$$

- cumulant generating function (log normalization constant):

$$A(\theta) = \log\Big(\int \exp\{\langle \theta,\, \boldsymbol{\phi}(\mathbf{x})\rangle\} \boldsymbol{\nu}(d\mathbf{x})\Big)$$

- set of valid parameters $\Theta := \{\theta \in \mathbb{R}^d \mid A(\theta) < +\infty\}$.

- will focus on *regular* families for which $\Theta$ is open.

# Examples: Scalar exponential families

| Family | $\mathcal{X}$ | $\nu$ | $\log p(\mathbf{x}; \theta)$ | $A(\theta)$ |
|--------|------|------|------------------|--------|
| Bernoulli | $\{0,1\}$ | Counting | $\theta x - A(\theta)$ | $\log[1 + \exp(\theta)]$ |
| Gaussian | $\mathbb{R}$ | Lebesgue | $\theta_1 x + \theta_2 x^2 - A(\theta)$ | $\frac{1}{2}[\theta_1 + \log \frac{2\pi e}{-\theta_2}]$ |
| Exponential | $(0, +\infty)$ | Lebesgue | $\theta(-x) - A(\theta)$ | $-\log \theta$ |
| Poisson | $\{0, 1, 2 \ldots\}$ | Counting $h(x) = 1/x!$ | $\theta x - A(\theta)$ | $\exp(\theta)$ |

# Graphical models as exponential families

- choose random variables $X_s$ at each vertex $s \in V$ from an arbitrary exponential family (e.g., Bernoulli, Gaussian, Dirichlet etc.)

- exponential family can be the same at each node (e.g., multivariate PSfrag replacements Gaussian), or different (e.g., mixture models).



**Key requirement:** The collection $\phi$ of sufficient statistics *must* respect the structure of $G$.

# Example: Discrete Markov random field

$$\theta_{st}(x_s, x_t)$$

$$\theta_t(x_t) \qquad \theta_s(x_s)$$

acements

Indicators:
$$\mathbb{I}_j(x_s) = \begin{cases} 1 & \text{if } x_s = j \\ 0 & \text{otherwise} \end{cases}$$

Parameters:
$$\theta_s = \{\theta_{s;j}, j \in \mathcal{X}_s\}$$
$$\theta_{st} = \{\theta_{st;jk}, (j,k) \in \mathcal{X}_s \times \mathcal{X}_t\}$$

Compact form:
$$\theta_s(x_s) := \sum_j \theta_{s;j} \mathbb{I}_j(x_s)$$
$$\theta_{st}(x_s, x_t) := \sum_{j,k} \theta_{st;jk} \mathbb{I}_j(x_s) \mathbb{I}_k(x_t)$$

Density (w.r.t. counting measure) of the form:

$$p(\mathbf{x}; \theta) \quad \propto \quad \exp\Big\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \Big\}$$

Cumulant generating function (log normalization constant):

$$A(\theta) \quad = \quad \log \sum_{\mathbf{x} \in \mathcal{X}^n} \exp\Big\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \Big\}$$

# Special case: Hidden Markov model

- Markov chain $\{X_1, X_2, \ldots\}$ evolving in time, with noisy observation $Y_t$ at each time $t$



- an HMM is a particular type of discrete MRF, representing the conditional $p(\mathbf{x} \mid \mathbf{y}; \theta)$

- exponential parameters have a concrete interpretation

$$
\begin{aligned}
\theta_{23}(x_2, x_3) &= \log p(x_3 \mid x_2) \\
\theta_5(x_5) &= \log p(y_5 \mid x_5)
\end{aligned}
$$

- the cumulant generating function $A(\theta)$ is equal to the log likelihood $\log p(\mathbf{y}; \theta)$

# Example: Multivariate Gaussian

$U(\theta)$: Matrix of natural parameters $\qquad$ $\phi(\mathbf{x})$: Matrix of sufficient statistics

$$
\begin{bmatrix}
0 & \theta_1 & \theta_2 & \ldots & \theta_n \\
\theta_1 & \theta_{11} & \theta_{12} & \ldots & \theta_{1n} \\
\theta_2 & \theta_{21} & \theta_{22} & \ldots & \theta_{2n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
\theta_n & \theta_{n1} & \theta_{n2} & \ldots & \theta_{nn}
\end{bmatrix}
\qquad
\begin{bmatrix}
1 & x_1 & x_2 & \ldots & x_n \\
x_1 & (x_1)^2 & x_1 x_2 & \ldots & x_1 x_n \\
x_2 & x_2 x_1 & (x_2)^2 & \ldots & x_2 x_n \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
x_n & x_n x_1 & x_n x_2 & \ldots & (x_n)^2
\end{bmatrix}
$$

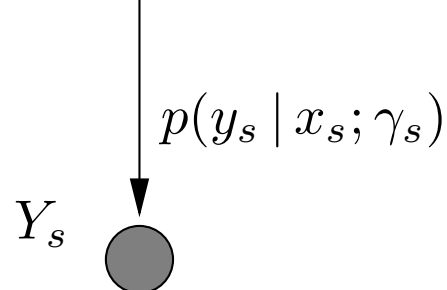Edgewise natural parameters $\theta_{st} = \theta_{ts}$ must respect graph structure:



$\quad$ (a) Graph structure $\quad$ (b) Structure of $[Z(\theta)]_{st} = \theta_{st}$.

# Example: Mixture of Gaussians

- can form *mixture models* by combining different types of random variables

- let $Y_s$ be conditionally Gaussian given the discrete variable $X_s$ with parameters $\gamma_{s;j} = (\mu_{s;j}, \sigma^2_{s;j})$:

$X_s \quad \bigcirc p(x_s; \theta_s)$

$p(y_s \,|\, x_s; \gamma_s)$

$Y_s \quad \bullet$

$X_s \quad \equiv \quad$ mixture indicator

$Y_s \quad \equiv \quad$ mixture of Gaussian

replacements

- couple the mixture indicators $\mathbf{X} = \{X_s, s \in V\}$ using a discrete MRF

- overall model has the exponential form

$$p(\mathbf{y}, \mathbf{x}; \theta, \gamma) \;\propto\; \prod_{s \in V} p(y_s \,|\, x_s; \gamma_s) \; \exp\Big\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t)] \Big\}.$$

# Conjugate dual functions

- conjugate duality is a fertile source of variational representations

- any function $f$ can be used to define another function $f^*$ as follows:

$$f^*(v) \quad := \quad \sup_{u \in \mathbb{R}^n} \big\{ \langle v, \, u \rangle - f(u) \big\}.$$

- easy to show that $f^*$ is always a convex function

- how about taking the "dual of the dual"? I.e., what is $(f^*)^*$?

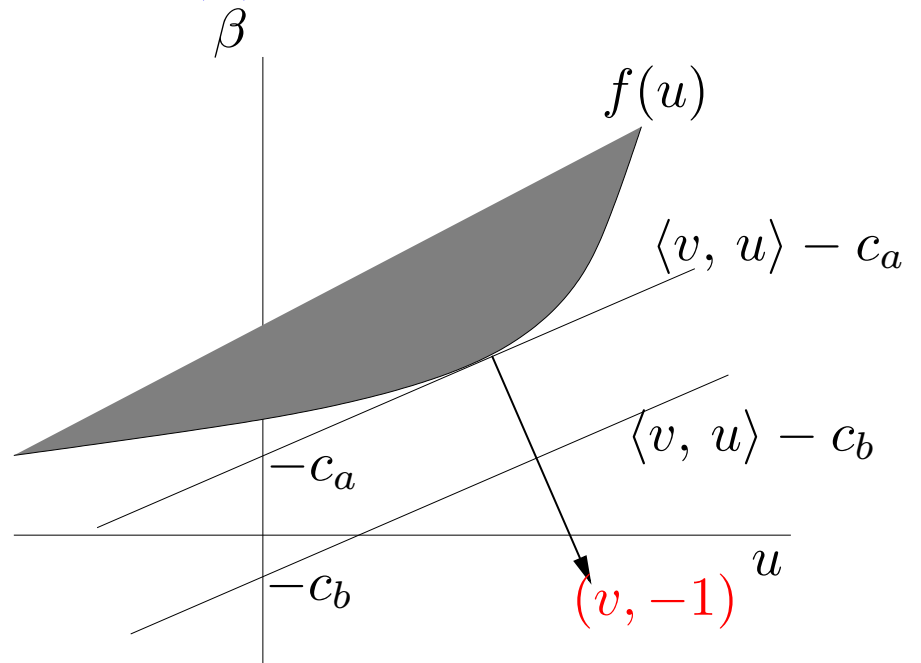- when $f$ is well-behaved (convex and lower semi-continuous), we have $(f^*)^* = f$, or alternatively stated:

$$f(u) \quad = \quad \sup_{v \in \mathbb{R}^n} \big\{ \langle u, \, v \rangle - f^*(v) \big\}$$

# Geometric view: Supporting hyperplanes

**Question:** Given all hyperplanes in $\mathbb{R}^n \times \mathbb{R}$ with normal $(v, -1)$, what is the intercept of the one that supports $\mathrm{epi}(f)$?



Epigraph of $f$: PSfrag replacements

$$\mathrm{epi}(f) := \{(u, \beta) \in \mathbb{R}^{n+1} \mid f(u) \leq \beta\}.$$

Analytically, we require the smallest $c \in \mathbb{R}$ such that:

$$\langle v, u \rangle - c \;\; \leq \;\; f(u) \quad \text{for} \;\; \text{all} \;\; u \in \mathbb{R}^n$$

By re-arranging, we find that this optimal $c^*$ is the dual value:

$$c^* \;\; = \;\; \sup_{u \in \mathbb{R}^n} \left\{ \langle v, u \rangle - f(u) \right\}.$$

# Example: Single Bernoulli

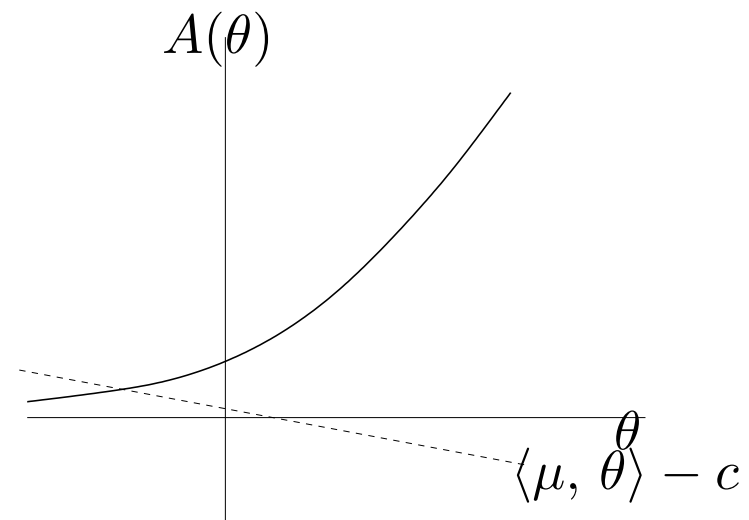Random variable $X \in \{0, 1\}$ yields exponential family of the form:

$$p(x; \theta) \propto \exp\{\theta x\} \quad \text{with} \quad A(\theta) = \log[1 + \exp(\theta)].$$

Let's compute the dual $A^*(\mu) := \sup_{\theta \in \mathbb{R}} \{\mu\theta - \log[1 + \exp(\theta)]\}$.

(Possible) stationary point: $\mu = \exp(\theta)/[1 + \exp(\theta)]$.



(a) Epigraph supported      (b) Epigraph *cannot* be supported

We find that:
$$A^*(\mu) = \begin{cases} \mu \log \mu + (1 - \mu) \log(1 - \mu) & \text{if } \mu \in [0, 1] \\ +\infty & \text{otherwise.} \end{cases}$$

Leads to the variational representation: $A(\theta) = \max_{\mu \in [0,1]} \{\mu \cdot \theta - A^*(\mu)\}$.

# More general computation of the dual $A^*$

- consider the definition of the dual function:

$$A^*(\mu) \quad = \quad \sup_{\theta \in \mathbb{R}^d} \{ \langle \mu, \, \theta \rangle - A(\theta) \}.$$

- taking derivatives w.r.t $\theta$ to find a stationary point yields:

$$\mu - \nabla A(\theta) \quad = \quad 0.$$

- <u>Useful fact:</u> Derivatives of $A$ yield *mean parameters:*

$$\frac{\partial A}{\partial \theta_\alpha}(\theta) \quad = \quad \mathbb{E}_\theta[\phi_\alpha(\mathbf{x})] \; := \; \int \phi_\alpha(\mathbf{x}) p(\mathbf{x}; \theta) \boldsymbol{\nu}(\mathbf{x}).$$

---

Thus, stationary points satisfy the equation:

$$\mu \quad = \quad \mathbb{E}_\theta[\boldsymbol{\phi}(\mathbf{x})] \tag{1}$$

# Computation of dual (continued)

- assume solution $\theta(\mu)$ to equation (1) exists

- strict concavity of objective guarantees that $\theta(\mu)$ attains global maximum with value

$$
\begin{aligned}
A^*(\mu) &= \langle \mu, \, \theta(\mu) \rangle - A(\theta(\mu)) \\
&= \mathbb{E}_{\theta(\mu)}\Big[ \langle \theta(\mu), \, \phi(\mathbf{x}) \rangle - A(\theta(\mu)) \Big] \\
&= \mathbb{E}_{\theta(\mu)}[\log p(\mathbf{x}; \theta(\mu))]
\end{aligned}
$$

- recall the definition of *entropy*:

$$
H(p(\mathbf{x})) \;:=\; -\int \big[\log p(\mathbf{x})\big] p(\mathbf{x}) \boldsymbol{\nu}(d\mathbf{x})
$$

- thus, we recognize that $A^*(\mu) = -H(p(\mathbf{x}; \theta(\mu)))$ when equation (1) has a solution

**Question:** For which $\mu \in \mathbb{R}^d$ does equation (1) have a solution $\theta(\mu)$?

# Sets of realizable mean parameters

- for any distribution $p(\cdot)$, define a vector $\mu \in \mathbb{R}^d$ of *mean parameters*:

$$\mu_\alpha \quad := \quad \int \phi_\alpha(\mathbf{x}) p(\mathbf{x}) \boldsymbol{\nu}(d\mathbf{x})$$

- now consider the set $\mathcal{M}(G; \boldsymbol{\phi})$ of all realizable mean parameters:

$$\mathcal{M}(G; \boldsymbol{\phi}) = \left\{ \mu \in \mathbb{R}^d \mid \mu_\alpha = \int \phi_\alpha(\mathbf{x}) p(\mathbf{x}) \boldsymbol{\nu}(d\mathbf{x}) \quad \text{for } some \ p(\cdot) \right\}$$

- for discrete families, we refer to this set as a *marginal polytope*, denoted by $\mathrm{MARG}(G; \boldsymbol{\phi})$

# Examples of $\mathcal{M}$: Gaussian MRF

$\phi(\mathbf{x})$ Matrix of sufficient statistics $\qquad$ $U(\mu)$ Matrix of mean parameters

$$
\begin{bmatrix}
1 & x_1 & x_2 & \ldots & x_n \\
x_1 & (x_1)^2 & x_1 x_2 & \ldots & x_1 x_n \\
x_2 & x_2 x_1 & (x_2)^2 & \ldots & x_2 x_n \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
x_n & x_n x_1 & x_n x_2 & \ldots & (x_n)^2
\end{bmatrix}
\qquad
\begin{bmatrix}
1 & \mu_1 & \mu_2 & \ldots & \mu_n \\
\mu_1 & \mu_{11} & \mu_{12} & \ldots & \mu_{1n} \\
\mu_2 & \mu_{21} & \mu_{22} & \ldots & \mu_{2n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
\mu_n & \mu_{n1} & \mu_{n2} & \ldots & \mu_{nn}
\end{bmatrix}
$$

- Gaussian mean parameters are specified by a single semidefinite constraint as $\mathcal{M}_{Gauss} = \{\mu \in \mathbb{R}^{n+\binom{n}{2}} \mid U(\mu) \succeq 0\}$.

**Scalar case:**

$$
U(\mu) = \begin{bmatrix} 1 & \mu_1 \\ \mu_1 & \mu_{11} \end{bmatrix}
$$

PSfrag replacements



$\mu_{11}$

$\mathcal{M}_{gauss}$

$\mu_1$

# Examples of $\mathcal{M}$: Discrete MRF

- sufficient statistics:
$$\mathbb{I}_j(x_s) \qquad \text{for } s = 1, \ldots n, \quad j \in \mathcal{X}_s$$
$$\mathbb{I}_{jk}(x_s, x_t) \quad \text{for}(s,t) \in E, \quad (j,k) \in \mathcal{X}_s \times \mathcal{X}_t$$
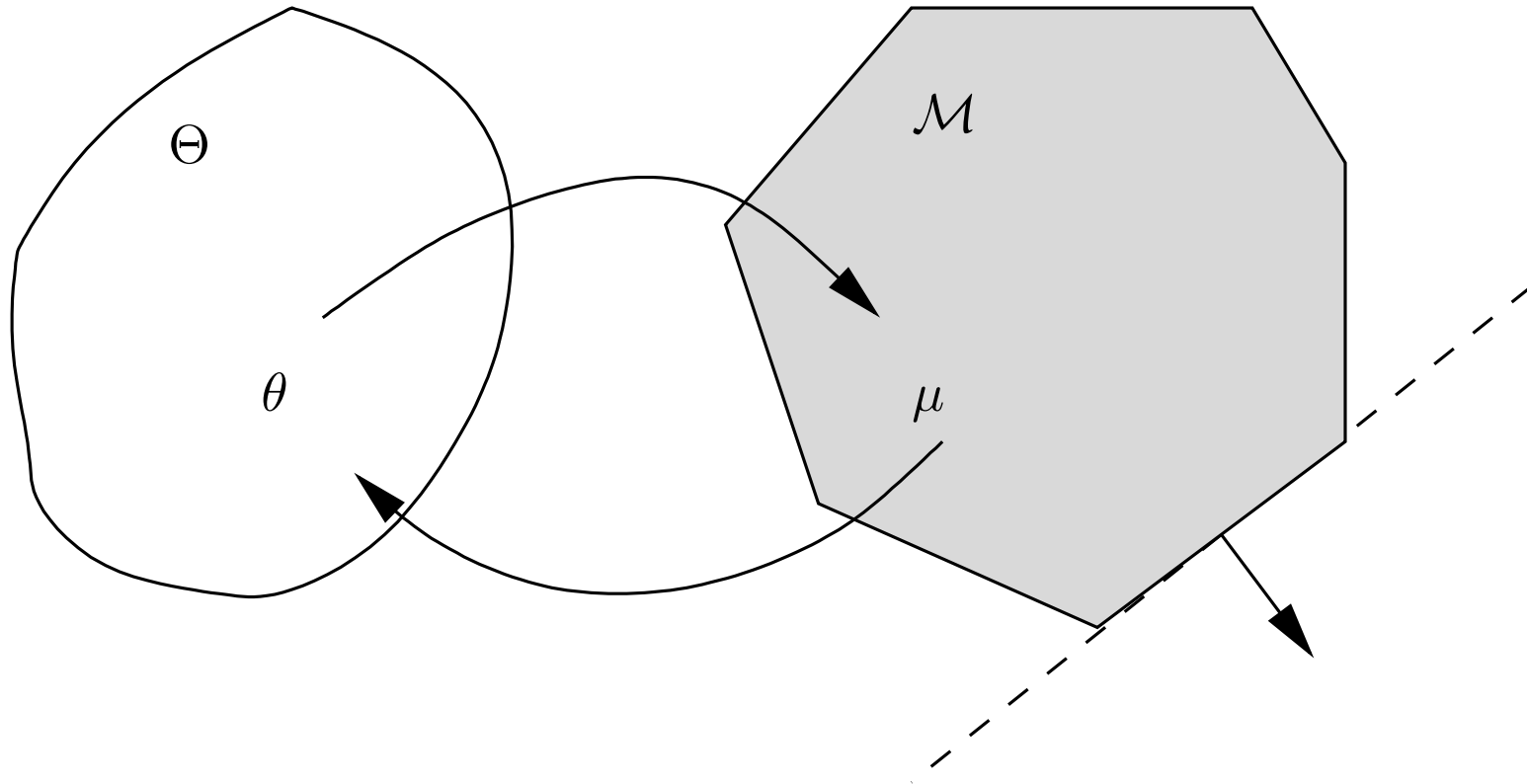
- mean parameters are simply marginal probabilities, represented as:

$$\mu_s(x_s) := \sum_{j \in \mathcal{X}_s} \mu_{s;j} \mathbb{I}_j(x_s), \qquad \mu_{st}(x_s, x_t) := \sum_{(j,k) \in \mathcal{X}_s \times \mathcal{X}_t} \mu_{st;jk} \mathbb{I}_{jk}(x_s, x_t)$$

$\mu_\mathbf{e}$

MARG$(G)$

acements

$a_j$

$\langle a_j, \mu \rangle = b_j$

- denote the set of realizable $\mu_s$ and $\mu_{st}$ by MARG$(G)$

- refer to it as the *marginal polytope*

- extremely difficult to characterize for general graphs

30

# Geometry and moment mapping



Θ

$\mathcal{M}$

θ

μ

eplacements

For suitable classes of graphical models in exponential form, the gradient map $\nabla A$ is a bijection between $\Theta$ and the interior of $\mathcal{M}$.

(e.g., Brown, 1986; Efron, 1978)

# Variational principle in terms of mean parameters

- The conjugate dual of $A$ takes the form:

$$A^*(\mu) \;=\; \begin{cases} -H(p(\mathbf{x}; \theta(\mu))) & \text{if } \mu \in \text{int}\, \mathcal{M}(G; \boldsymbol{\phi}) \\ +\infty & \text{if } \mu \notin \text{cl}\, \mathcal{M}(G; \boldsymbol{\phi}). \end{cases}$$

*Interpretation:*

- $A^*(\mu)$ is finite (and equal to a certain negative entropy) for any $\mu$ that is globally realizable

- if $\mu \notin \text{cl}\, \mathcal{M}(G; \boldsymbol{\phi})$, then the max. entropy problem is *infeasible*

- The cumulant generating function $A$ has the representation:

$$\underbrace{A(\theta)}_{\text{cumulant generating func.}} \;=\; \underbrace{\sup_{\mu \in \mathcal{M}(G;\phi)} \{ \langle \theta,\, \mu \rangle \;-\; A^*(\mu) \}}_{\text{max. ent. problem over } \mathcal{M}},$$

- in contrast to the "free energy" approach, solving this problem provides both the value $A(\theta)$ and the exact mean parameters $\widehat{\mu}_\alpha = \mathbb{E}_\theta[\phi_\alpha(\mathbf{x})]$

# Alternative view: Kullback-Leibler divergence

- Kullback-Leibler divergence defines "distance" between probability distributions:

$$D(p \,\|\, q) \quad := \quad \int \left[ \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \right] p(\mathbf{x}) \boldsymbol{\nu}(d\mathbf{x})$$

- for two exponential family members $p(\mathbf{x}; \theta^1)$ and $p(\mathbf{x}; \theta^2)$, we have

$$D(p(\mathbf{x}; \theta^1) \,\|\, p(\mathbf{x}; \theta^2)) \quad = \quad A(\theta^2) - A(\theta^1) - \langle \mu^1, \, \theta^2 - \theta^1 \rangle$$

- substituting $A(\theta^1) = \langle \theta^1, \, \mu^1 \rangle - A^*(\mu^1)$ yields a *mixed form*:

$$D(p(\mathbf{x}; \theta^1) \,\|\, p(\mathbf{x}; \theta^2)) \quad \equiv \quad D(\mu^1 \,\|\, \theta^2) \; = \; A(\theta^2) + A^*(\mu^1) - \langle \mu^1, \, \theta^2 \rangle$$

Hence, the following two assertions are equivalent:

$$A(\theta^2) \quad = \quad \sup_{\mu^1 \in \mathcal{M}(G; \boldsymbol{\phi})} \quad \{ \langle \theta^2, \, \mu^1 \rangle \; - \; A^*(\mu^1) \}$$

$$0 \quad = \quad \inf_{\mu^1 \in \mathcal{M}(G; \boldsymbol{\phi})} D(\mu^1 \,\|\, \theta^2)$$

# Challenges

1. In general, mean parameter spaces $\mathcal{M}$ can be very difficult to characterize (e.g., multidimensional moment problems).

2. Entropy $A^*(\mu)$ as a function of *only* the mean parameters $\mu$ typically lacks an explicit form.

**Remarks:**

1. Variational representation clarifies why certain models are tractable.

2. For intractable cases, one strategy is to solve an approximate form of the optimization problem.

# Outline

1. Recap

    (a) Background on graphical models

    (b) Some applications and challenging problems

    (c) Illustrations of some message-passing algorithms

2. Exponential families and variational methods

    (a) What is a variational method (and why should I care)?

    (b) Graphical models as exponential families

    (c) Variational representations from conjugate duality

3. Exact techniques as variational methods

    (a) Gaussian inference on arbitrary graphs

    (b) Belief-propagation/sum-product on trees (e.g., Kalman filter; $\alpha$-$\beta$ alg.)

    (c) Max-product on trees (e.g., Viterbi)

4. Approximate techniques as variational methods

    (a) Mean field and variants

    (b) Belief propagation and extensions

    (c) Convex methods and upper bounds

    (d) Tree-reweighted max-product and linear programs

# A(i): Multivariate Gaussian (fixed covariance)

Consider the set of all Gaussians with fixed *inverse* covariance $Q \succ 0$.

- potentials $\phi(\mathbf{x}) = \{x_1, \ldots, x_n\}$ and natural parameter $\theta \in \Theta = \mathbb{R}^n$.

- cumulant generating function:

$$
A(\theta) \quad = \quad \log \int_{\mathbb{R}^n} \overbrace{\exp\Big\{\sum_{s=1}^{n} \theta_s x_s\Big\}}^{\text{density}} \underbrace{\exp\Big\{-\frac{1}{2}\mathbf{x}^T Q \mathbf{x}\Big\} d\mathbf{x}}_{\text{base measure}}
$$

- completing the square yields $A(\theta) = \frac{1}{2}\theta^T Q^{-1}\theta + \text{constant}$

- straightforward computation leads to the dual
$$
A^*(\mu) = \tfrac{1}{2}\mu^T Q \mu - \text{constant}
$$

- putting the pieces back together yields the variational principle

$$
A(\theta) \quad = \quad \sup_{\mu \in \mathbb{R}^n} \Big\{\theta^T \mu - \frac{1}{2}\mu^T Q \mu\Big\} + \text{constant}
$$

- optimum is uniquely obtained at the familiar Gaussian mean $\widehat{\mu} = Q^{-1}\theta$.

# A(ii): Multivariate Gaussian (arbitrary covariance)

- matrices of sufficient statistics, natural parameters, and mean parameters:

$$\phi(\mathbf{x}) = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \begin{bmatrix} 1 & \mathbf{x} \end{bmatrix}, \quad U(\theta) := \begin{bmatrix} 0 & [\theta_s] \\ [\theta_s] & [\theta_{st}] \end{bmatrix} \quad U(\mu) := \mathbb{E}\left\{ \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \begin{bmatrix} 1 & \mathbf{x} \end{bmatrix} \right\}$$

- cumulant generating function:

$$A(\theta) \quad = \quad \log \int \exp\left\{ \langle\!\langle U(\theta), \ \phi(\mathbf{x}) \rangle\!\rangle \right\} d\mathbf{x}$$

- computing the dual function:

$$A^*(\mu) \quad = \quad -\frac{1}{2} \log \det U(\mu) - \frac{n}{2} \log 2\pi e,$$

- exact variational principle is a *log-determinant problem*:

$$A(\theta) \quad = \quad \sup_{U(\mu) \succ 0, \ [U(\mu)]_{11}=1} \left\{ \langle\!\langle U(\theta), \ U(\mu) \rangle\!\rangle + \frac{1}{2} \log \det U(\mu) \right\} + \frac{n}{2} \log 2\pi e$$

- solution yields the *normal equations* for Gaussian mean and covariance.

# B: Belief propagation/sum-product on trees

- discrete variables $X_s \in \{0, 1, \ldots, m_s - 1\}$ on a *tree* $T = (V, E)$

- sufficient statistics: indicator functions for each node and edge

$$\mathbb{I}_j(x_s) \quad \text{for} \quad s = 1, \ldots n, \quad j \in \mathcal{X}_s$$

$$\mathbb{I}_{jk}(x_s, x_t) \quad \text{for} \quad (s, t) \in E, \quad (j, k) \in \mathcal{X}_s \times \mathcal{X}_t.$$

- exponential representation of distribution:

$$p(\mathbf{x}; \theta) \quad \propto \quad \exp\Big\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \Big\}$$

where $\theta_s(x_s) := \sum_{j \in \mathcal{X}_s} \theta_{s;j} \mathbb{I}_j(x_s)$   (and similarly for $\theta_{st}(x_s, x_t)$)

- mean parameters are simply marginal probabilities, represented as:

$$\mu_s(x_s) := \sum_{j \in \mathcal{X}_s} \mu_{s;j} \mathbb{I}_j(x_s), \qquad \mu_{st}(x_s, x_t) := \sum_{(j,k) \in \mathcal{X}_s \times \mathcal{X}_t} \mu_{st;jk} \mathbb{I}_{jk}(x_s, x_t)$$

- the marginals must belong to the following *marginal polytope*:

$$\text{MARG}(T) \quad := \quad \{ \mu \geq 0 \mid \sum_{x_s} \mu_s(x_s) = 1, \ \sum_{x_t} \mu_{st}(x_s, x_t) = \mu_s(x_s) \},$$

# Decomposition of entropy for trees

- by the junction tree theorem, any tree can be factorized in terms of its marginals $\mu \equiv \mu(\theta)$ as follows:

$$p(\mathbf{x}; \theta) = \prod_{s \in V} \mu_s(x_s) \prod_{(s,t) \in E} \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s)\mu_t(x_t)}$$

- taking logs and expectations leads to an entropy decomposition

$$H(p(\mathbf{x}; \theta)) = -A^*(\mu(\theta)) = \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} I_{st}(\mu_{st})$$

where

Single node entropy: $\quad H_s(\mu_s) := -\sum_{x_s} \mu_s(x_s) \log \mu_s(x_s)$

Mutual information: $\quad I_{st}(\mu_{st}) := \sum_{x_s, x_t} \mu_{st}(x_s, x_t) \log \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s)\mu_t(x_t)}$.

- thus, the dual function $A^*(\mu)$ has an explicit and easy form

# Exact variational principle on trees

- putting the pieces back together yields:

$$A(\theta) \;=\; \max_{\mu \in \mathrm{MARG}(T)} \Big\{ \langle \theta, \, \mu \rangle + \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E(T)} I_{st}(\mu_{st}) \Big\}.$$

- let's try to solve this problem by a (partial) Lagrangian formulation

- assign a Lagrange multiplier $\lambda_{ts}(x_s)$ for each constraint
  $C_{ts}(x_s) := \mu_s(x_s) - \sum_{x_t} \mu_{st}(x_s, x_t) = 0$

- will enforce the normalization $(\sum_{x_s} \mu_s(x_s) = 1)$ and non-negativity constraints explicitly

- the Lagrangian takes the form:

$$\mathcal{L}(\mu; \lambda) = \langle \theta, \, \mu \rangle + \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E(T)} I_{st}(\mu_{st})$$

$$+ \sum_{(s,t) \in E} \Big[ \sum_{x_t} \lambda_{st}(x_t) C_{st}(x_t) + \sum_{x_s} \lambda_{ts}(x_s) C_{ts}(x_s) \Big]$$

# Lagrangian derivation (continued)

- taking derivatives of the Lagrangian w.r.t $\mu_s$ and $\mu_{st}$ yields

$$\frac{\partial \mathcal{L}}{\partial \mu_s(x_s)} = \theta_s(x_s) - \log \mu_s(x_s) + \sum_{t \in \mathcal{N}(s)} \lambda_{ts}(x_s) + C$$

$$\frac{\partial \mathcal{L}}{\partial \mu_{st}(x_s, x_t)} = \theta_{st}(x_s, x_t) - \log \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s)\mu_t(x_t)} - \lambda_{ts}(x_s) - \lambda_{st}(x_t) + C'$$

- setting these partial derivatives to zero and simplifying:

$$\mu_s(x_s) \propto \exp\{\theta_s(x_s)\} \prod_{t \in \mathcal{N}(s)} \exp\{\lambda_{ts}(x_s)\}$$

$$\mu_s(x_s, x_t) \propto \exp\{\theta_s(x_s) + \theta_t(x_t) + \theta_{st}(x_s, x_t)\} \times$$

$$\prod_{u \in \mathcal{N}(s) \setminus t} \exp\{\lambda_{us}(x_s)\} \prod_{v \in \mathcal{N}(t) \setminus s} \exp\{\lambda_{vt}(x_t)\}$$

- enforcing the constraint $C_{ts}(x_s) = 0$ on these representations yields the familiar update rule for the *messages* $M_{ts}(x_s) = \exp(\lambda_{ts}(x_s))$:

$$M_{ts}(x_s) \leftarrow \sum_{x_t} \exp\{\theta_t(x_t) + \theta_{st}(x_s, x_t)\} \prod_{u \in \mathcal{N}(t) \setminus s} M_{ut}(x_t)$$

# C: Max-product algorithm on trees

**Question:** What should be the form of a variational principle for computing modes?

**Intuition:** Consider behavior of the family $\{p(\mathbf{x}; \beta\theta) \mid \beta > 0\}$.



(a) Low $\beta$          (b) High $\beta$

**Conclusion:** Problem of computing modes should be related to limiting form $(\beta \rightarrow +\infty)$ of computing marginals.

# Limiting form of the variational principle

- consider the variational principle for a discrete MRF of the form $p(\mathbf{x}; \beta\theta)$:

$$\frac{1}{\beta} A(\beta\theta) \quad = \quad \frac{1}{\beta} \max_{\mu \in \text{MARG}} \left\{ \langle \beta\theta, \, \mu \rangle - A^*(\mu) \right\}.$$

- taking limits as $\beta \to +\infty$ yields:

$$\underbrace{\max_{\mathbf{x} \in \mathcal{X}^N} \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}}_{\text{computation of modes}} \quad = \quad \underbrace{\max_{\mu \in \text{MARG}(G)} \left\{ \langle \theta, \, \mu \rangle \right\}}_{\text{linear program}}.$$

- thus, computing the mode in a discrete MRF is equivalent to a *linear program over the marginal polytope*

# Max-product on tree-structured MRFs

- recall the max-product (belief revision) updates:

$$M_{ts}(x_s) \quad \leftarrow \quad \max_{x_t} \exp\left\{\theta_t(x_t) + \theta_{st}(x_s, x_t)\right\} \prod_{u \in \mathcal{N}(t) \setminus s} M_{ut}(x_t)$$

- for trees, the variational principle (linear program) takes the especially simple form

$$\max_{\mu \in \text{MARG}(T)} \left\{ \sum_{s \in V} \theta_s(x_s)\mu_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t)\mu_{st}(x_s, x_t) \right\}$$

- constraint set is the marginal polytope for trees

$$\text{MARG}(T) \quad := \quad \left\{ \mu \geq 0 \;\Big|\; \sum_{x_s} \mu_s(x_s) = 1, \;\; \sum_{x_t} \mu_{st}(x_s, x_t) = \mu_s(x_s) \right\},$$

- a similar Lagrangian formulation shows that max-product is an iterative method for solving this linear program (details in Wainwright & Jordan, 2003)

# Outline

1. Recap

    (a) Background on graphical models

    (b) Some applications and challenging problems

    (c) Illustrations of some message-passing algorithms

2. Exponential families and variational methods

    (a) What is a variational method (and why should I care)?

    (b) Graphical models as exponential families

    (c) Variational representations from conjugate duality

3. Exact techniques as variational methods

    (a) Gaussian inference on arbitrary graphs

    (b) Belief-propagation/sum-product on trees (e.g., Kalman filter; $\alpha$-$\beta$ alg.)

    (c) Max-product on trees (e.g., Viterbi)

4. Approximate techniques as variational methods

    (a) Mean field and variants

    (b) Belief propagation and extensions

    (c) Convex relaxations and upper bounds

    (d) Tree-reweighted max-product and linear programs

# A: Mean field theory

**Difficulty:** (typically) no explicit form for $-A^*(\mu)$ (i.e., entropy as a function of mean parameters) $\implies$ exact variational principle is intractable.

**Idea:** Restrict $\mu$ to a *subset* of distributions for which $-A^*(\mu)$ has a tractable form.

**Examples:**

(a) For product distributions $p(\mathbf{x}) = \prod_{s \in V} \mu_s(x_s)$, entropy decomposes as $-A^*(\mu) = \sum_{s \in V} H_s(x_s)$.

(b) Similarly, for trees (more generally, decomposable graphs), the junction tree theorem yields an explicit form for $-A^*(\mu)$.

**Definition:** A subgraph $H$ of $G$ is *tractable* if the entropy has an explicit form for any distribution that respects $H$.

# Geometry of mean field



- let $H$ represent a *tractable subgraph* (i.e., for which $A^*$ has explicit form)

- let $\mathcal{M}_{tr}(G; H)$ represent tractable mean parameters:

$$\mathcal{M}_{tr}(G; H) := \{\mu |\ \mu = \mathbb{E}_\theta[\phi(\mathbf{x})]\ \ \text{s.t.}\ \theta\ \text{respects}\ H\}.$$



$\mu_{\mathbf{e}}$

$\mathcal{M}_{tr}$

$\mathcal{M}$

- under mild conditions, $\mathcal{M}_{tr}$ is a non-convex *inner approximation* to $\mathcal{M}$

- optimizing over $\mathcal{M}_{tr}$ (as opposed to $\mathcal{M}$) yields *lower bound*:

$$A(\theta)\ \geq\ \sup_{\widetilde{\mu} \in \mathcal{M}_{tr}}\ \{\langle \theta,\ \widetilde{\mu} \rangle - A^*(\widetilde{\mu})\}.$$

47

# Alternative view: Minimizing KL divergence

- recall the *mixed form* of the KL divergence between $p(\mathbf{x}; \theta)$ and $p(\mathbf{x}; \widetilde{\theta})$:

$$D(\widetilde{\mu} \,\|\, \theta) \quad = \quad A(\theta) + A^*(\widetilde{\mu}) - \langle \widetilde{\mu},\, \theta \rangle$$

- try to find the "best" approximation to $p(\mathbf{x}; \theta)$ in the sense of KL divergence

- in analytical terms, the problem of interest is

$$\inf_{\widetilde{\mu} \in \mathcal{M}_{tr}} D(\widetilde{\mu} \,\|\, \theta) \quad = \quad A(\theta) + \inf_{\widetilde{\mu} \in \mathcal{M}_{tr}} \left\{ A^*(\widetilde{\mu}) - \langle \widetilde{\mu},\, \theta \rangle \right\}$$

- hence, finding the tightest lower bound on $A(\theta)$ is equivalent to finding the best approximation to $p(\mathbf{x}; \theta)$ from distributions with $\widetilde{\mu} \in \mathcal{M}_{tr}$

# Example: Naive mean field algorithm for Ising model

- consider completely disconnected subgraph $H = (V, \emptyset)$

- permissible exponential parameters belong to subspace

$$\mathcal{E}(H) = \{\theta \in \mathbb{R}^d \mid \theta_{st} = 0 \;\; \forall \;\; (s,t) \in E\}$$

- allowed distributions take product form $p(\mathbf{x}; \theta) = \prod_{s \in V} p(x_s; \theta_s)$, and generate

$$\mathcal{M}_{tr}(G; H) = \{\mu \mid \mu_{st} = \mu_s \mu_t, \;\; \mu_s \in [0, 1]\}.$$

- approximate variational principle:

$$\max_{\mu_s \in [0,1]} \left\{ \sum_{s \in V} \theta_s \mu_s + \sum_{(s,t) \in E} \theta_{st} \mu_s \mu_t - \Big[ \sum_{s \in V} \mu_s \log \mu_s + (1 - \mu_s) \log(1 - \mu_s) \Big] \right\}.$$

- **Co-ordinate ascent:** with all $\{\mu_t, t \neq s\}$ fixed, problem is strictly concave in $\mu_s$ and optimum is attained at

$$\mu_s \;\; \longleftarrow \;\; \left\{ 1 + \exp\Big[ -\big(\theta_s + \sum_{t \in \mathcal{N}(s)} \theta_{st} \mu_t\big)\Big] \right\}^{-1}$$

49

# Example: Structured mean field for coupled HMM



(a)                                                    (b)

- entropy of distribution that respects $H$ decouples into sum: one term for each chain.

- *structured mean field updates* are an iterative method for finding the tightest approximation (either in terms of KL or lower bound)

# B: Belief propagation on arbitrary graphs

**Two main ingredients:**

1. Exact entropy $-A^*(\mu)$ is intractable, so let's approximate it.

   The *Bethe approximation* $A^*_{Bethe}(\mu) \approx A^*(\mu)$ is based on the exact expression for trees:

   $$-A^*_{Bethe}(\mu) \quad = \quad \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} I_{st}(\mu_{st}).$$

2. The *marginal polytope* $\mathrm{MARG}(G)$ is also difficult to characterize, so let's use the following (tree-based) outer bound:

   $$\mathrm{LOCAL}(G) \quad := \quad \{\, \tau \geq 0 \mid \sum_{x_s} \tau_s(x_s) = 1, \; \sum_{x_t} \tau_{st}(x_s, x_t) = \tau_s(x_s) \,\},$$

**Note:** Use $\tau$ to distinguish these locally consistent *pseudomarginals* from globally consistent marginals.

# Geometry of belief propagation

- combining these ingredients leads to the *Bethe variational principle:*

$$\max_{\tau \in \mathrm{LOCAL}(G)} \Big\{ \langle \theta, \, \tau \rangle + \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} I_{st}(\tau_{st}) \Big\}$$

- belief propagation can be derived as an iterative method for solving a Lagrangian formulation of the BVP          (Yedidia et al., 2002)



- belief propagation uses a *polyhedral outer approximation* to $\mathcal{M}$

- for any graph, $\mathrm{LOCAL}(G) \supseteq \mathrm{MARG}(G)$.

- equality holds $\iff$ $G$ is a tree.

52

# Illustration: Globally inconsistent BP fixed points

Consider the following assignment of pseudomarginals $\tau_s, \tau_{st}$:



Locally consistent (pseudo)marginals

- can verify that $\tau \in \text{LOCAL}(G)$, and that $\tau$ is a fixed point of belief propagation (with all constant messages)

- however, $\tau$ is globally inconsistent

**Note:** More generally: for any $\tau$ in the interior of $\text{LOCAL}(G)$, can construct a distribution with $\tau$ as a BP fixed point.

# High-level perspective

- message-passing algorithms (e.g., mean field, belief propagation) are solving approximate versions of exact variational principle in exponential families

- there are two *distinct* components to approximations:

  (a) can use either inner or outer bounds to $\mathcal{M}$

  (b) various approximations to entropy function $-A^*(\mu)$

- <u>mean field:</u> non-convex inner bound and exact form of entropy

- <u>BP:</u> polyhedral outer bound and non-convex Bethe approximation

- <u>Kikuchi and variants:</u> tighter polyhedral outer bounds and better entropy approximations

  (e.g.,Yedidia et al., 2002)

# Generalized belief propagation on hypergraphs

- a *hypergraph* is a natural generalization of a graph

- it consists of a set of vertices $V$ and a set $E$ of hyperedges, where each *hyperedge* is a subset of $V$

- convenient graphical representation in terms of *poset diagrams*



(a) Ordinary graph     (b) Hypertree (width 2)     (c) Hypergraph

- *descendants* and *ancestors* of a hyperedge $h$:

$$\mathcal{D}^+(h) := \{\, g \in E \mid g \subseteq h \,\}, \qquad \mathcal{A}^+(h) := \{\, g \in E \mid g \supseteq h \,\}.$$

# Hypertree factorization and entropy

- hypertrees are an alternative way to describe junction trees

- associated with any poset is a Möbius function $\omega : E \times E \to \mathbb{Z}$

$$\omega(g, g) = 1, \quad \omega(g, h) = - \sum_{\{f \mid g \subseteq f \subset h\}} \omega(g, f)$$

  Example: For Boolean poset, $\omega(g, h) = (-1)^{|h| \setminus |g|}$.

- use the Möbius function to define a correspondence between the collection of marginals $\mu := \{\mu_h\}$ and new set of functions $\varphi := \{\varphi_h\}$:

$$\log \varphi_h(x_h) = \sum_{g \in \mathcal{D}^+(h)} \omega(g, h) \log \mu_g(x_g), \qquad \log \mu_h(x_h) = \sum_{g \in \mathcal{D}^+(h)} \log \varphi_g(x_g).$$

- any hypertree-structured distribution is guaranteed to factor as:

$$p(\mathbf{x}) = \prod_{h \in E} \varphi_h(x_h).$$

# Examples: Hypertree factorization

1. **Ordinary tree:**

$$\varphi_s(x_s) = \mu_s(x_s) \qquad \text{for any vertex } s$$

$$\varphi_{st}(x_s, x_t) = \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s)\,\mu_t(x_t)} \qquad \text{for any edge } (s,t)$$

2. **Hypertree:**

$$\varphi_{1245} = \frac{\mu_{1245}}{\frac{\mu_{25}}{\mu_5}\frac{\mu_{45}}{\mu_5}\mu_5}$$

$$\varphi_{45} = \frac{\mu_{45}}{\mu_5}$$

$$\varphi_5 = \mu_5$$



Combining the pieces:

$$p = \frac{\mu_{1245}}{\frac{\mu_{25}}{\mu_5}\frac{\mu_{45}}{\mu_5}\mu_5}\ \frac{\mu_{2356}}{\frac{\mu_{25}}{\mu_5}\frac{\mu_{56}}{\mu_5}\mu_5}\ \frac{\mu_{4578}}{\frac{\mu_{45}}{\mu_5}\frac{\mu_{58}}{\mu_5}\mu_5}\ \frac{\mu_{25}}{\mu_5}\ \frac{\mu_{45}}{\mu_5}\ \frac{\mu_{56}}{\mu_5}\ \frac{\mu_{58}}{\mu_5}\ \mu_5 = \frac{\mu_{1245}\,\mu_{2356}\,\mu_{4578}}{\mu_{25}\,\mu_{45}}$$

# Building augmented hypergraphs

Better entropy approximations via augmented hypergraphs.



(a) Original

(b) Clustering

(c) Full covering



(d) Kikuchi

(e) Fails single counting

# C. Convex relaxations and upper bounds

Possible concerns with the Bethe/Kikuchi problems and variations?

(a) lack of convexity $\Rightarrow$ multiple local optima, and substantial algorithmic complications

(b) failure to bound the log partition function

**Goal:** Techniques for approximate computation of marginals and parameter estimation based on:

(a) convex variational problems $\Rightarrow$ unique global optimum

(b) relaxations of exact problem $\Rightarrow$ upper bounds on $A(\theta)$

**Usefulness of bounds:**

(a) interval estimates for marginals

(b) approximate parameter estimation

(c) large deviations (prob. of rare events)

# Bounds from "convexified" Bethe/Kikuchi problems

**Idea:** Upper bound $-A^*(\mu)$ by convex combination of tree-structured entropies.



PSfrag replacementsPSfrag replacements    PSfrag replacements    PSfrag replacements

$$-A^*(\mu) \quad \leq \quad -\rho(T^1)A^*(\mu(T^1)) \quad - \quad \rho(T^2)A^*(\mu(T^2)) \quad - \quad \rho(T^3)A^*(\mu(T^3))$$

- given any spanning tree $T$, define the moment-matched tree distribution:

$$p(\mathbf{x}; \mu(T)) \quad := \quad \prod_{s \in V} \mu_s(x_s) \prod_{(s,t) \in E} \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s)\,\mu_t(x_t)}$$

- use $-A^*(\mu(T))$ to denote the associated tree entropy

- let $\boldsymbol{\rho} = \{\rho(T)\}$ be a probability distribution over spanning trees
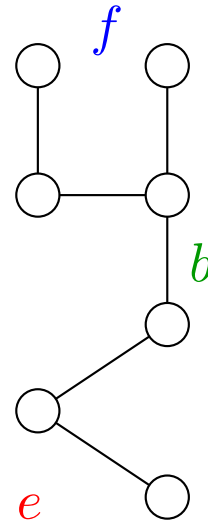
60

# Edge appearance probabilities

**Experiment:** What is the probability $\rho_e$ that a given edge $e \in E$ belongs to a tree $T$ drawn randomly under $\boldsymbol{\rho}$?
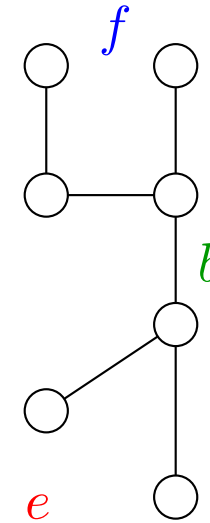


(a) Original     (b) $\rho(T^1) = \frac{1}{3}$     (c) $\rho(T^2) = \frac{1}{3}$     (d) $\rho(T^3) = \frac{1}{3}$

In this example:     $\rho_b = 1$;     $\rho_e = \frac{2}{3}$;     $\rho_f = \frac{1}{3}$.

The vector $\boldsymbol{\rho_e} = \{\, \rho_e \mid e \in E \,\}$ must belong to the *spanning tree polytope*, denoted $\mathbb{T}(G)$.        (Edmonds, 1971)

# Optimal bounds by tree-reweighted message-passing

Recall the constraint set of locally consistent marginal distributions:

$$\text{LOCAL}(G) \;=\; \{\, \tau \geq 0 \mid \underbrace{\sum_{x_s} \tau_s(x_s) = 1}_{\text{normalization}}, \;\; \underbrace{\sum_{x_s} \tau_{st}(x_s, x_t) = \tau_t(x_t)}_{\text{marginalization}} \,\}.$$

**Theorem:**     (Wainwright, Jaakkola, & Willsky, 2002; To appear in IEEE-IT)

(a) For any given edge weights $\boldsymbol{\rho_e} = \{\rho_e\}$ in the spanning tree polytope, the optimal upper bound over *all* tree parameters is given by:

$$A(\theta) \;\leq\; \max_{\tau \in \text{LOCAL}(G)} \{\, \langle \theta, \tau \rangle + \sum_{s \in V} H_s(\tau_s) - \sum_{(s,t) \in E} \rho_{st} I_{st}(\tau_{st}) \,\}.$$

(b) This optimization problem is strictly convex, and its unique optimum is specified by the fixed point of $\boldsymbol{\rho_e}$-reweighted message passing:

$$M_{ts}^*(x_s) = \kappa \sum_{x_t' \in \mathcal{X}_t} \left\{ \exp\left[ \frac{\theta_{st}(x_s, x_t')}{\rho_{st}} + \theta_t(x_t') \right] \frac{\prod\limits_{v \in \Gamma(t) \setminus s} \left[ M_{vt}^*(x_t) \right]^{\rho_{vt}}}{\left[ M_{st}^*(x_t) \right]^{(1-\rho_{ts})}} \right\}.$$

# Semidefinite constraints in convex relaxations

**Fact:** Belief propagation and its hypergraph-based generalizations all involve polyhedral (i.e., *linear*) outer bounds on the marginal polytope.

**Idea:** Use *semidefinite* constraints to generate more global outer bounds.

**Example:** For the Ising model, relevant mean parameters are $\mu_s = p(X_s = 1)$ and $\mu_{st} = p(X_s = 1, X_t = 1)$.
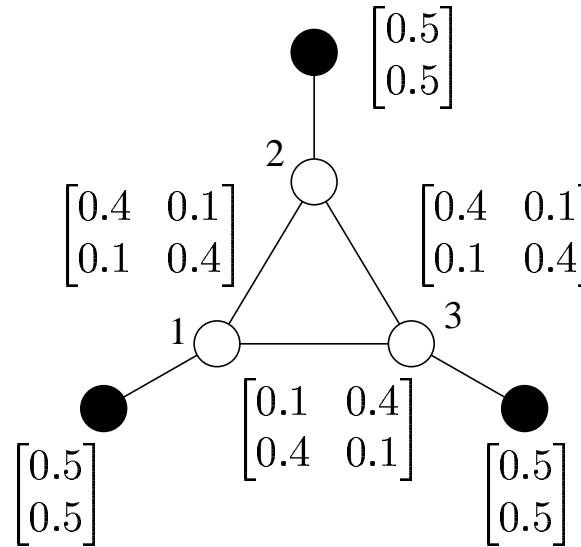
Define $\mathbf{y} = [1 \ \mathbf{x}]^T$, and consider the second-order moment matrix:

$$
\mathbb{E}[\mathbf{y}\mathbf{y}^T] = \begin{bmatrix}
1 & \mu_1 & \mu_2 & \cdots & \mu_n \\
\mu_1 & \mu_1 & \mu_{12} & \cdots & \mu_{1n} \\
\mu_2 & \mu_{12} & \mu_2 & \cdots & \mu_{2n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
\mu_n & \mu_{n1} & \mu_{n2} & \cdots & \mu_n
\end{bmatrix}
$$

It must be positive semidefinite, which imposes (an infinite number of) linear constraints on $\mu_s, \mu_{st}$.

# Illustrative example

$$\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

2

$$\begin{bmatrix} 0.4 & 0.1 \\ 0.1 & 0.4 \end{bmatrix} \qquad \begin{bmatrix} 0.4 & 0.1 \\ 0.1 & 0.4 \end{bmatrix}$$

3

Locally consistent
(pseudo)marginals

1

$$\begin{bmatrix} 0.1 & 0.4 \\ 0.4 & 0.1 \end{bmatrix}$$

$$\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \qquad \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

Second-order
moment matrix

$$\begin{bmatrix} \mu_1 & \mu_{12} & \mu_{13} \\ \mu_{21} & \mu_2 & \mu_{23} \\ \mu_{31} & \mu_{32} & \mu_3 \end{bmatrix} = \begin{bmatrix} 0.5 & 0.4 & 0.1 \\ 0.4 & 0.5 & 0.4 \\ 0.1 & 0.4 & 0.5 \end{bmatrix}$$

Not positive-semidefinite!

64

# Log-determinant relaxation

- based on optimizing over covariance matrices $M_1(\mu) \in \mathrm{SDEF}_1(K_n)$

**Theorem:** Consider an outer bound $\mathrm{OUT}(K_n)$ that satisfies:

$$\mathrm{MARG}(K_n) \quad \subseteq \quad \mathrm{OUT}(K_n) \subseteq \mathrm{SDEF}_1(K_n)$$

For any such outer bound, $A(\theta)$ is upper bounded by:

$$\max_{\mu \in \mathrm{OUT}(K_n)} \left\{ \langle \theta, \mu \rangle + \frac{1}{2} \log \det \left[ M_1(\mu) + \frac{1}{3} \mathrm{blkdiag}[0, I_n] \right] \right\} + \frac{n}{2} \log(\frac{\pi e}{2})$$

**Remarks:**

1. Log-det. problem can be solved efficiently by interior point methods.

2. Relevance for applications:

   (a) Upper bound on $A(\theta)$.

   (b) Method for computing approximate marginals.

<div align="right">(Wainwright & Jordan, 2003)</div>

# Results for approximating marginals



(a) Nearest-neighbor grid        (b) Fully connected

- average $\ell_1$ error in approximate marginals over 100 trials

- coupling types: repulsive $(-)$, mixed $(+/-)$, attractive $(+)$
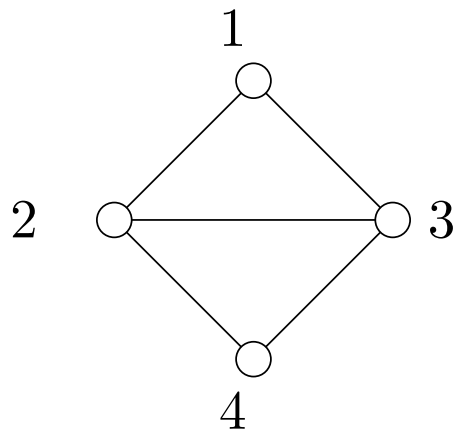
# D. Tree-reweighted max-product algorithm

- integer program (IP): optimizing cost function over a discrete space
  (e.g., $\{0,1\}^n$)

- IPs computationally intractable in general, but graphical structure
  can be exploited

- tree-structured IPs solvable in linear time with dynamic
  programming (max-product message-passing)

- standard max-product algorithm applied to graphs with cycles, but
  no guarantees in general

- a novel perspective:
  - "tree-reweighted" max-product for arbitrary graphs
  - optimality guarantees for reweighted message-passing
  - forges connection between max-product message passing and linear
    programming (LP) relaxations

# Some previous theory on ordinary max-product

- analysis of single-cycle graphs   (Aji & McEliece, 1998; Horn, 1999; Weiss, 1998)

- guarantees for attenuated max-product updates (Frey & Koetter, 2001)

- locally optimal on "tree-plus-loop" neighborhoods  (Weiss & Freeman, 2001)

- strengthened optimality results and computable error bounds on the gap                                                                 (Wainwright et al., 2003)

- exactness after finite # iterations for max. weight bipartite matching                                                                 (Bayati, Shah & Sharma, 2005)

# Standard analysis via computation tree

- standard tool: computation tree of message-passing updates

  (Gallager, 1963; Weiss, 2001; Richardson & Urbanke, 2001)



(a) Original graph        (b) Computation tree (4 iterations)

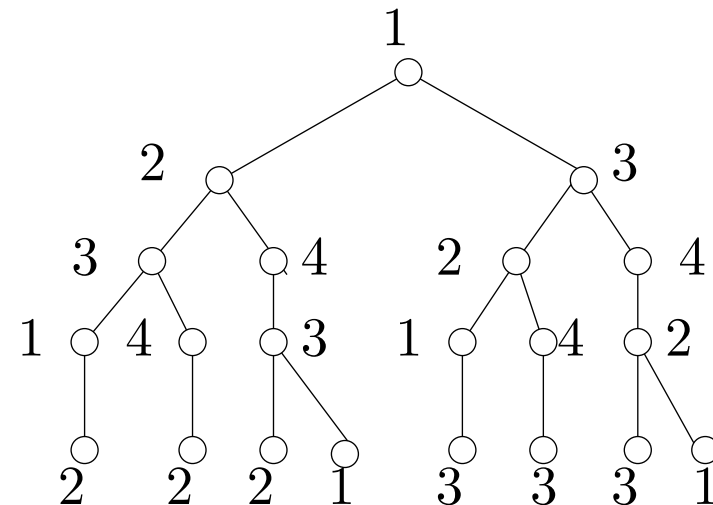- level $t$ of tree: all nodes whose messages reach the root (node 1) after $t$ iterations of message-passing

# Illustration: Non-exactness of standard max-product

**Intuition:**

- max-product is solving (exactly) a modified problem on the computation tree

- nodes *not equally weighted* in computation tree $\Rightarrow$ max-product can output an incorrect configuration
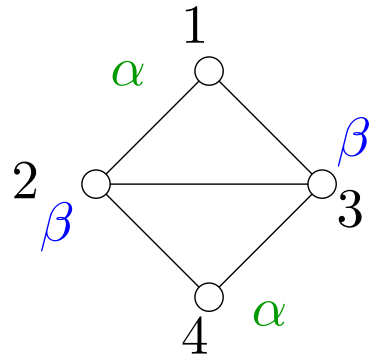


(a) Diamond graph $G_{\mathrm{dia}}$      (b) Computation tree (4 iterations)

- for example: asymptotic node fractions in this computation tree:

$$\begin{bmatrix} f(1) & f(2) & f(3) & f(4) \end{bmatrix} = \begin{bmatrix} 0.2393 & 0.2607 & 0.2607 & 0.2393 \end{bmatrix}$$

# A whole family of non-exact examples

- consider the following integer program on $G_{\text{dia}}$:

replacements



$$\theta_s(x_s) \quad \begin{cases} \alpha x_s & \text{if } s = 1 \text{ or } s = 4 \\ \beta x_s & \text{if } s = 2 \text{ or } s = 3 \end{cases}$$
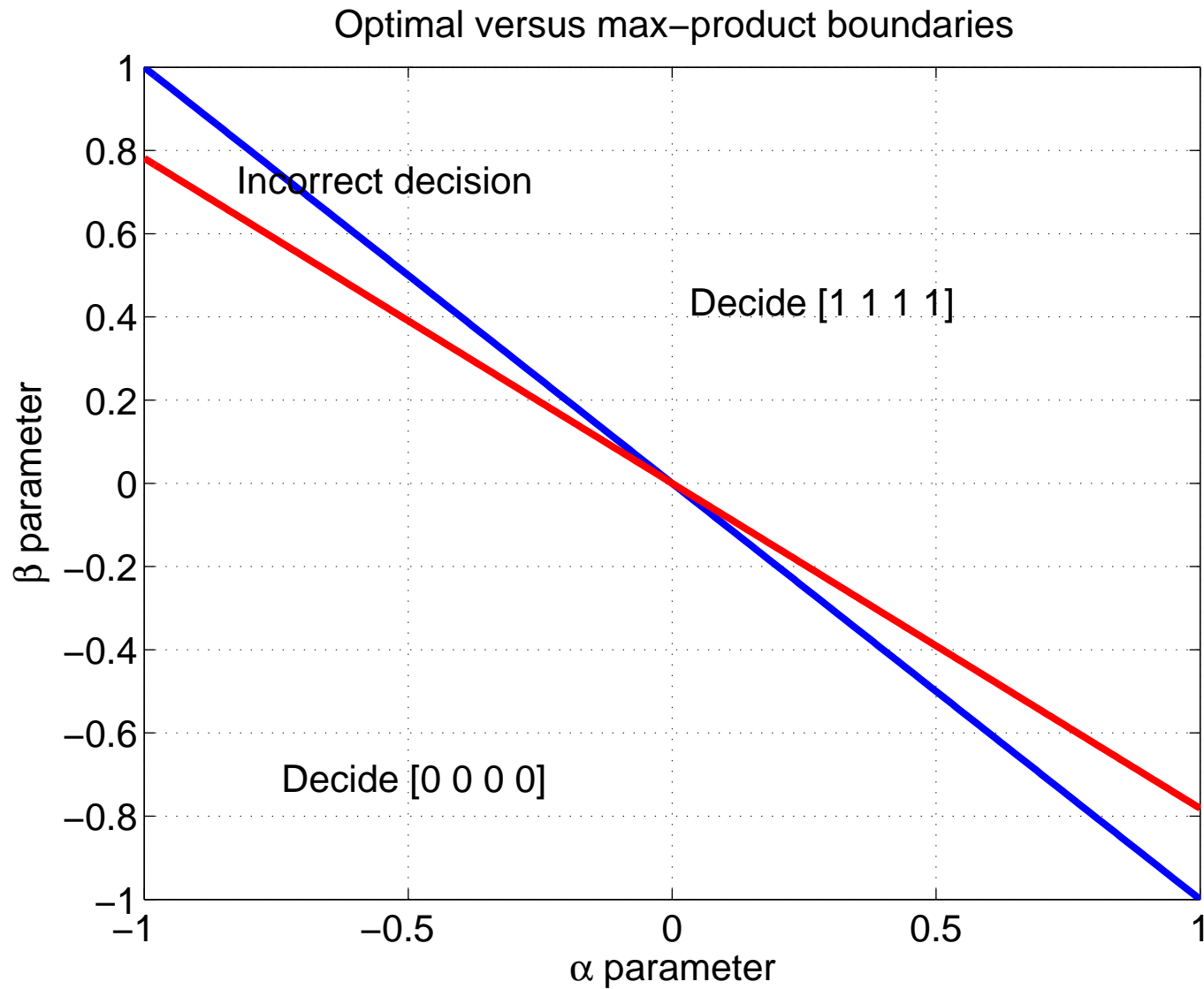
$$\theta_{st}(x_s, x_t) \quad = \quad \begin{cases} -\gamma & \text{if } x_s \neq x_t \\ 0 & \text{otherwise} \end{cases}$$

- for $\gamma$ sufficiently large, optimal solution is always either $0^4 = [0\,0\,0\,0]$ or $1^4 = [1\,1\,1\,1]$.

- max-product and optimum give *different* decision boundaries:

Optimum boundary: $\quad \widehat{\mathbf{x}} = \begin{cases} 1^4 & \text{if } 0.25\alpha + 0.25\beta \geq 0 \\ 0^4 & \text{otherwise} \end{cases}$

Max-product boundary: $\quad \widehat{\mathbf{x}} = \begin{cases} 1^4 & \text{if } 0.2393\alpha + 0.2607\beta \geq 0 \\ 0^4 & \text{otherwise} \end{cases}$

# Incorrect max-product decision boundary

# Tree-reweighted max-product algorithm

(Wainwright, Jaakkola & Willsky, 2002)

Message update from node $t$ to node $s$:

$$M_{ts}(x_s) \quad \leftarrow \quad \kappa \max_{x'_t \in \mathcal{X}_t} \left\{ \exp\left[ \underbrace{\frac{\theta_{st}(x_s, x'_t)}{\rho_{st}}}_{\text{reweighted edge}} + \theta_t(x'_t) \right] \frac{\displaystyle \prod_{v \in \Gamma(t) \backslash s} \overbrace{\left[M_{vt}(x_t)\right]^{\rho_{vt}}}^{\text{reweighted messages}}}{\underbrace{\left[M_{st}(x_t)\right]^{(1-\rho_{ts})}}_{\text{opposite message}}} \right\}.$$

**Properties:**

1. Modified updates remain *distributed* and *purely local* over the graph.

   - Messages are reweighted with $\rho_{st} \in [0,1]$.
2. Key differences:   - Potential on edge $(s,t)$ is rescaled by $\rho_{st} \in [0,1]$.
   - Update involves the reverse direction edge.

3. The choice $\rho_{st} = 1$ for all edges $(s,t)$ recovers standard update.

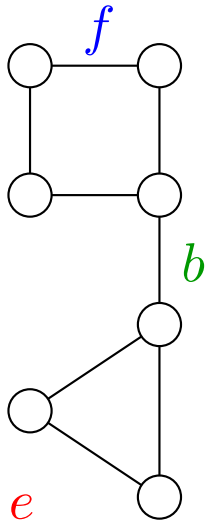# TRW max-product computes pseudo-max-marginals

- observe that the mode-finding problem can be solved via the exact *max-marginals*:

$$\mu_s(x_s) := \max_{\{\mathbf{x}' \,|\, x'_s = x_s\}} p(\mathbf{x}'), \qquad \mu_{st}(x_s, x_t) := \max_{\{\mathbf{x}' \,|\, x'_s = x_s, x'_t = x_t\}} p(\mathbf{x}')$$
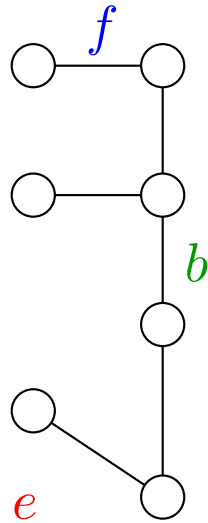
- any message fixed point $M^*$ defines a set of *pseudo-max-marginals* $\{\nu_s^*, \nu_{st}^*\}$ at nodes and edges

- will show these TRW message-passing never "lies", and establish optimality guarantees for suitable fixed points

- guarantees require the following conditions:

  (a) **Weight choice:** the edge weights $\rho = \{\rho_{st} \mid (s,t) \in E\}$ are "suitably" chosen (in the spanning tree polytope)

  (b) **Strong tree agreement:** There exists an $\mathbf{x}^*$ such that:
  
  - Nodewise optimal: $x_s^*$ belongs to $\arg\max_{x_s} \nu_s^*(x_s)$.
  - Edgewise optimal: $(x_s^*, x_t^*)$ belongs to $\arg\max_{x_s, x_t} \nu_{st}^*(x_s, x_t)$.
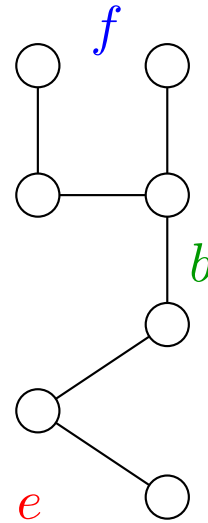
# Edge appearance probabilities

**Experiment:** What is the probability $\rho_e$ that a given edge $e \in E$ belongs to a tree $T$ drawn randomly under $\boldsymbol{\rho}$?
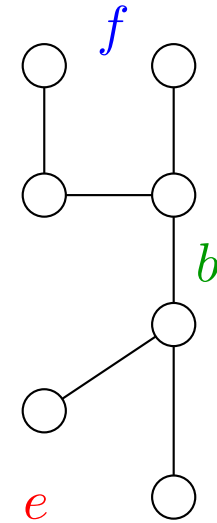


(a) Original    (b) $\rho(T^1) = \frac{1}{3}$    (c) $\rho(T^2) = \frac{1}{3}$    (d) $\rho(T^3) = \frac{1}{3}$

In this example:    $\rho_b = 1$;    $\rho_e = \frac{2}{3}$;    $\rho_f = \frac{1}{3}$.

The vector $\boldsymbol{\rho_e} = \{ \rho_e \mid e \in E \}$ must belong to the *spanning tree polytope*, denoted $\mathbb{T}(G)$. (Edmonds, 1971)

# TRW max-product never "lies"

**Set-up:** Any fixed point $\nu^*$ that satisfies the strong tree agreement (STA) condition defines a configuration $\mathbf{x}^* = (x_1^*, \ldots, x_n^*)$ such that

$$\underbrace{x_s^* \in \arg\max_{x_s} \nu_s^*(x_s),}_{\text{Node optimality}} \qquad \underbrace{(x_s^*, x_t^*) \in \arg\max_{x_s, x_t} \nu_{st}^*(x_s, x_t)}_{\text{Edge-wise optimality}}$$

**Theorem 1** (Arbitrary problems)               (Wainwright et al., 2003):

(a) Any STA configuration $\mathbf{x}^*$ is provably MAP-optimal for the graph with cycles.

(b) Any STA fixed point is a dual-optimal solution to a certain "tree-based" linear programming relaxation.

Hence, TRW max-product acknowledges failure by *lack of strong tree agreement*.

# Performance of tree-reweighted max-product

**Key question:** When can strong tree agreement be obtained?

- performance guarantees for particular problem classes:

  (a) guarantees for submodular and related binary problems
      (Kolmogorov & Wainwright, 2005)

  (b) LP decoding of error-control codes
      (Feldman, Wainwright & Karger, 2005; Feldman et al., 2004)

- empirical work on TRW max-product and tree LP relaxation:

  (a) LP decoding of error-control codes
      (Feldman, Wainwright & Karger, 2002, 2003, 2005; Koetter & Vontobel,
      2003, 2005)

  (b) data association problem in sensor networks
      (Chen et al., SPIE 2003)

  (c) solving stereo-problems using MRF-based formulation
      (Kolmogorov, 2005; Weiss, Meltzer & Chanover, 2005)

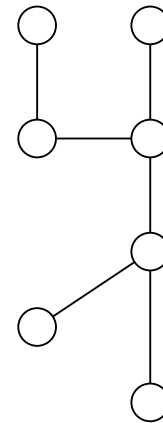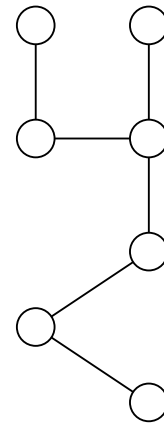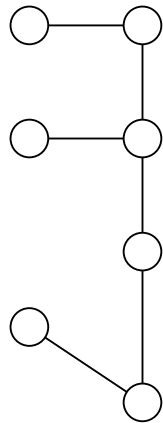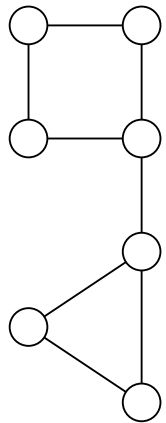# Basic idea: convex combinations of trees

**Observation:** Easy to find its MAP-optimal configurations on trees:

$$\mathrm{OPT}(\theta(T)) \quad := \quad \big\{ \mathbf{x} \in \mathcal{X}^n \mid \mathbf{x} \text{ is MAP-optimal for } p(\mathbf{x}; \theta(T)) \big\}.$$

**Idea:** Approximate original problem by a convex combination of trees.

$$\boldsymbol{\rho} = \{\rho(T)\} \quad \equiv \quad \text{probability distribution over spanning trees}$$

$$\theta(T) \quad \equiv \quad \text{tree-structured parameter vector}$$

$$* \qquad \theta^* \quad = \quad \rho(T^1)\theta(T^1) \quad + \quad \rho(T^2)\theta(T^2) \quad + \quad \rho(T^3)\theta(T^3)$$

$$\dagger \quad \mathrm{OPT}(\theta^*) \quad \supseteq \quad \mathrm{OPT}(\theta(T^1)) \quad \cap \quad \mathrm{OPT}(\theta(T^2)) \quad \cap \quad \mathrm{OPT}(\theta(T^3)).$$

# Tree fidelity and agreement

**Goal:** Want a set $\{\theta(T)\}$ of tree-structured parameters and probability distribution $\{\rho(T)\}$ such that:

*Fidelity to original:*
$$\overbrace{\theta^* = \sum_T \rho(T)\theta(T)}^{\text{Combining trees yields original problem.}}$$

*Tree agreement:* The set $\underbrace{\bigcap_T \mathrm{OPT}(\theta(T))}_{\text{Configurations on which all trees agree.}}$ is non-empty.

**Lemma:** Under the fidelity condition, the set $\bigcap_T \mathrm{OPT}(\theta(T))$ of configurations on which trees agree is contained within the optimal set $\mathrm{OPT}(\theta^*)$ for the original problem.

**Consequence:** If we can find a set of tree parameters that satisfy both conditions, then any configuration $\mathbf{x}^* \in \bigcap_T \mathrm{OPT}(\theta(T))$ is MAP-optimal for the *original problem* on the MRF with cycles.

# Message-passing as negotiation among trees

**Intuition:**

- reweighted message-passing $\Rightarrow$ negotiation among the trees

- ultimate goal $\Rightarrow$ adjust messages to obtain *tree agreement.*

- when tree agreement is obtained, the configuration $\mathbf{x}^*$ is MAP-optimal for the graph with cycles

- hence solving a sequence of tree problems (*very easy to do!*) suffices to find an optimum on the MRF (hard in general)
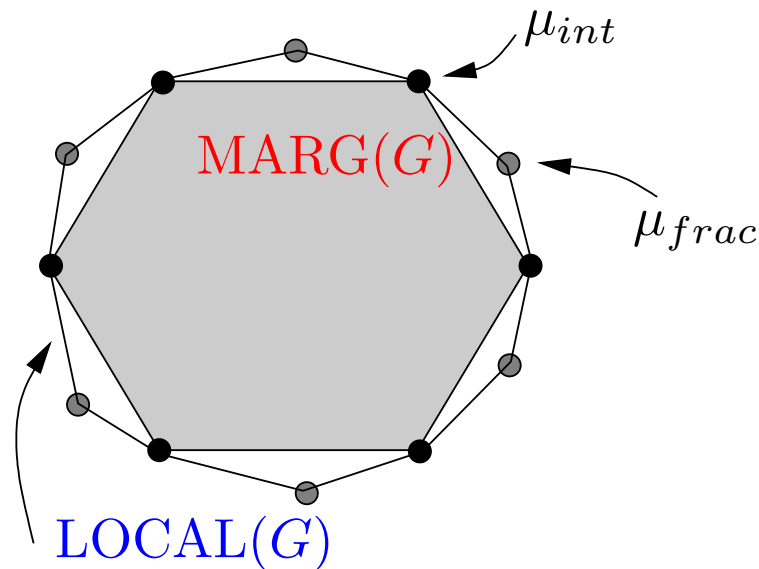
# Dual perspective: linear programming relaxation

- Reweighted message-passing is attempting to find a *tight upper bound* on the convex function $A_\infty(\theta^*) := \max_{\mathbf{x} \in \mathcal{X}^n} \langle \theta^*, \, \boldsymbol{\phi}(\mathbf{x}) \rangle$.

- The dual of this problem is a *linear programming relaxation* of the integer program.

- any unconstrained integer program (IP) can be re-written as a LP over the *convex hull of all possible solutions*      (Bertsimas & Tsitsiklis, 1997)

- the IP $\max_{\mathbf{x} \in \{0,1\}^n} \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}$ is exactly the same as the linear program (LP)

$$\max_{\mu_s, \mu_{st} \in \mathrm{MARG}(G)} \left\{ \sum_{s \in V} \sum_{x_s} \mu_s(x_s)\theta_s(x_s) + \sum_{(s,t) \in E} \sum_{x_s, x_t} \theta_{st}(x_s, x_t)\mu_{st}(x_s, x_t) \right\}$$

- here $\mathrm{MARG}(G)$ is the *marginal polytope* of all realizable marginal distributions $\mu_s$ and $\mu_{st}$

# Geometric perspective of LP relaxation



PSfrag replacements

- relaxation is based on replacing the exact marginal polytope (very hard to characterize!) with the tree relaxation

$$\text{LOCAL}(G) \quad := \quad \{\, \tau \geq 0 \mid \sum_{x_s} \tau_s(x_s) = 1, \ \sum_{x_t} \tau_{st}(x_s, x_t) = \tau_s(x_s) \,\},$$

- leads to a (poly-time solvable) linear programming relaxation

# Guarantees for submodular and binary problems

- stronger optimality assertions can be made for particular classes of problems

- important problem class: binary quadratic programs (i.e., modes of pairwise graphical model with binary variables)

- subclass of binary quadratic programs: *supermodular* interactions:

$$\underbrace{\theta_{st}(0,0) + \theta_{st}(1,1)}_{\text{agreement}} \quad \geq \quad \underbrace{\theta_{st}(1,0) + \theta_{st}(0,1).}_{\text{disagreement}}$$

- supermodular maximization can be performed in polynomial-time (max-flow)

**Theorem 2:** (Binary supermodular)      (Kolmogorov & Wainwright, 2005)
The TRW max-product algorithm always succeeds for binary supermodular problems.

# Partial information

**Question:** Can TRW message-passing still yield useful (partial) information when strong tree agreement does not hold?
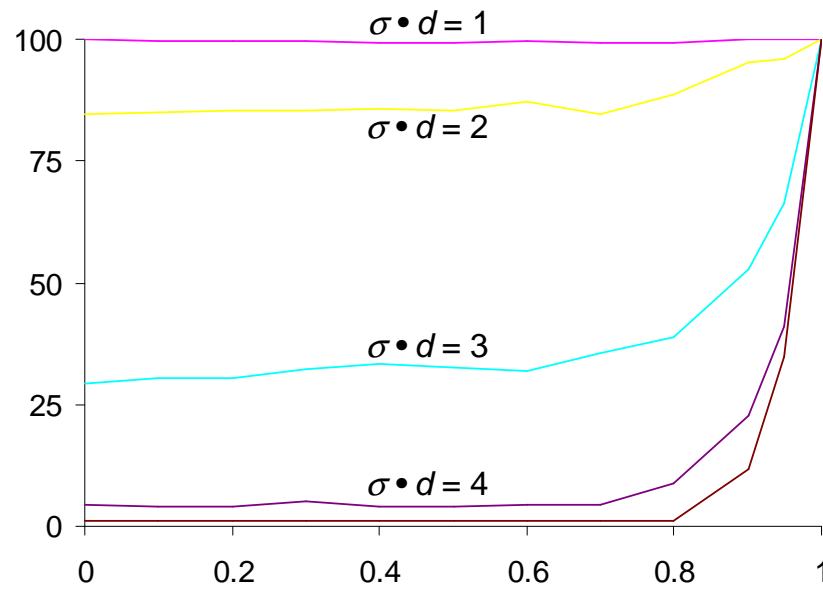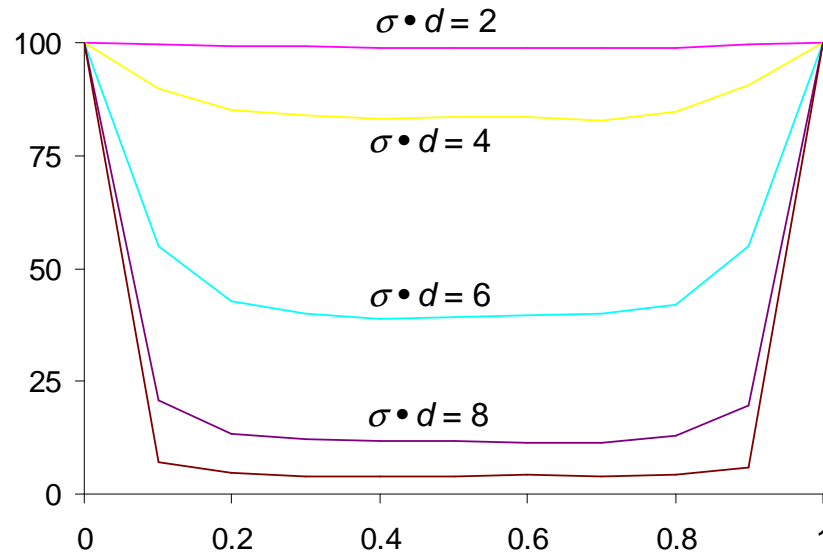
**Theorem 2:** (Binary persistence) (Kolmogorov & Wainwright, 2005; Hammer et al., 1984)

Let $S \subseteq V$ be the subset of vertices for which there exists a single point $x_s^* \in \arg\max_{x_s} \nu_s^*(x_s)$. Then for *any optimal solution* $\mathbf{y} \in \mathrm{OPT}(\theta^*)$, it holds that $y_s = x_s^*$.
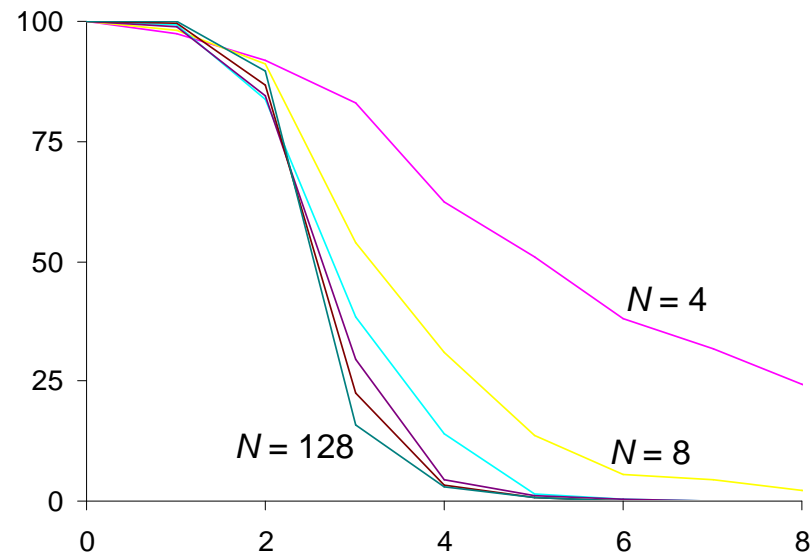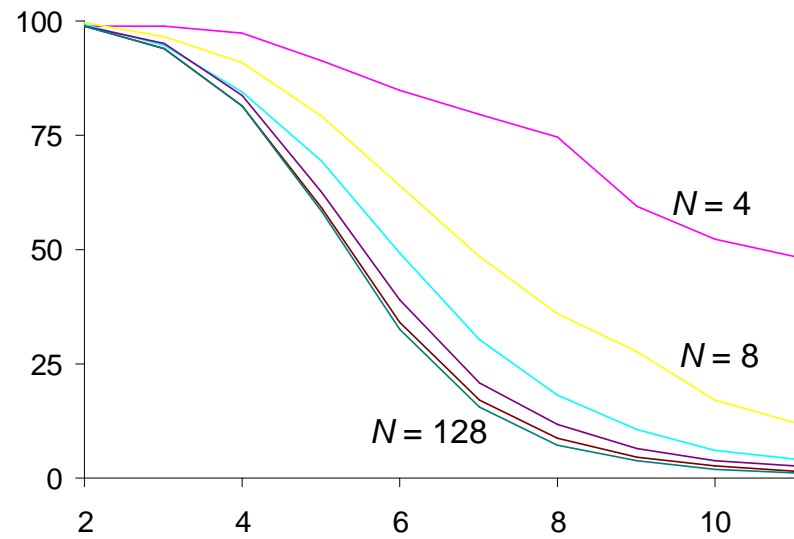
Some outstanding questions:

- what fraction of variables are correctly determined in "typical" problems?

- do similar partial validity results hold for more general (non-binary) problems?

# Some experimental results: amount of frustration

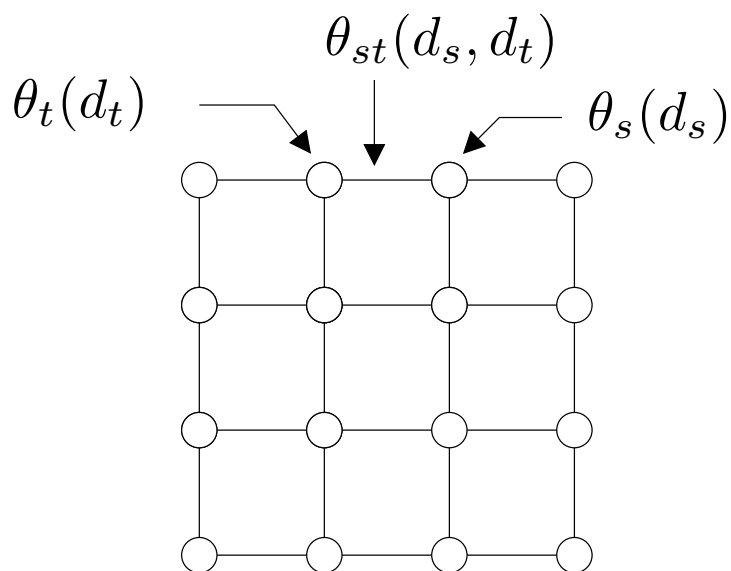# Some experimental results: coupling strength

# Disparity computation in stereo vision

- estimate depth in scenes based on two (or more) images taken from different positions

- one powerful source of stereo information:
  - biological vision: disparity from offset eyes
  - computer vision: disparity from offset cameras

- challenging (computational) problem: estimate the disparity at each point of an image based on a stereo pair

- broad range of research in both visual neuroscience and computer vision

# Global approaches to disparity computation

- wide range of approaches to disparity in computer vision (see, e.g., Scharstein & Szeliski, 2002)

- *global approaches*: disparity map based on optimization in an MRF

$$\theta_{st}(d_s, d_t)$$

$$\theta_t(d_t) \quad\quad\quad \theta_s(d_s)$$

acements

- grid-structured graph $G = (V, E)$
- $d_s \equiv$ disparity at grid position $s$
- $\theta_s(d_s) \equiv$ image data fidelity term
- $\theta_{st}(d_s, d_t) \equiv$ disparity coupling

- optimal disparity map $\widehat{\mathbf{d}}$ found by solving MAP estimation problem for this Markov random field

- computationally intractable (NP-hard) in general, but TRW max-product and other message-passing algorithms can be applied

# Middlebury stereo benchmark set

- standard set of benchmarked examples for stereo algorithms
  (Scharstein & Szeliski, 2002)

- Tsukuba data set: Image sizes $384 \times 288 \times 16$ ($W \times H \times D$)



(a) Original image      (b) Ground truth disparity

# Comparison of different methods



(a) Scanline dynamic programming

(b) Graph cuts

(c) Ordinary belief propagation

(d) Tree-reweighted max-product

(a), (b): Scharstein & Szeliski, 2002; (c): Sun et al., 2002 (d): Weiss, et al., 2005;

# Summary and future directions

- variational methods are based on converting statistical and computational problems to optimization:

  (a) complementary to sampling-based methods (e.g., MCMC)

  (b) a variety of new "relaxations" remain to be explored

- many open questions:

  (a) prior error bounds available only in special cases

  (b) extension to non-parametric settings?

  (c) hybrid techniques (variational and MCMC)

  (d) variational methods in parameter estimation

  (e) fast techniques for solving large-scale relaxations (e.g., SDPs, other convex programs)