

Announcements

- Wainwright and Jordan Graphical Models, Exponential Families, and Variational Inference http://www.nowpublishers.com/product. aspx?product=MAL&doi=2200000001
- Should have read chapters 1 through 5 in our book. Read chapter 7.
- Also read chapter 8 (integer programming, although we probably won't cover that chapter in class unfortunately).
- Also should have read "Divergence measures and message passing" by Thomas Minka, and "Structured Region Graphs: Morphing EP into GBP", by Welling, Minka, and Teh.
- Assignment due Wednesday (Dec 3rd) night, 11:45pm. Final project proposal final progress report (one page max).
- Update: For status update, final writeup, and talk, use notation as close as possible to that used in class!

- L1 (9/29): Introduction, Families, Semantics
- L2 (10/1): MRFs, elimination, Inference on Trees
- L3 (10/6): Tree inference, message passing, more general queries, non-tree)
- L4 (10/8): Non-trees, perfect elimination, triangulated graphs
- L5 (10/13): triangulated graphs, k-trees, the triangulation process/heuristics
- L6 (10/15): multiple queries, decomposable models, junction trees
- L7 (10/20): junction trees, begin intersection graphs
- L8 (10/22): intersection graphs, inference on junction trees
- L9 (10/27): inference on junction trees, semirings,
- L10 (11/3): conditioning, hardness, LBP

- L11 (11/5): LBP, exponential models,
- L12 (11/10): exponential models, mean params and polytopes,
- L13 (11/12): polytopes, tree outer bound, Bethe entropy approx.
- L14 (11/17): Bethe entropy approx, loop series correction
- L15 (11/19): Hypergraphs, posets, Mobius, Kikuchi
- L16 (11/24): Kikuchi, Expectation Propagation
- L17 (11/26): Expectation Propagation, Mean Field
- L18 (12/1): Structured mean field, Convex relaxations and upper bounds, tree reweighted case
- L19 (12/3):
- Final Presentations: (12/10):

Finals Week: Dec 8th-12th, 2014.

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014 Prof. Jeff Bilmes

F3/45 (pg.3/45)

Review Conjugate Duality, Maximum Likelihood, Negative Entropy

Theorem 18.2.3 (Relationship between A and A^*)

(a) For any $\mu \in \mathcal{M}^{\circ}$, $\theta(\mu)$ unique canonical parameter sat. matching condition, then conj. dual takes form:

$$A^{*}(\mu) = \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^{\circ} \\ +\infty & \text{if } \mu \notin \overline{\mathcal{M}} \end{cases}$$
(18.3)

(b) Partition function has variational representation (dual of dual)

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$
(18.4)

(c) For $\theta \in \Omega$, sup occurs at $\mu \in \mathcal{M}^{\circ}$ of moment matching conditions

$$\mu = \int_{\mathsf{D}_X} \phi(x) p_{\theta}(x) \nu(dx) = \mathbb{E}_{\theta}[\phi(X)] = \nabla A(\theta)$$
 (18.5)

Variational Approach Amenable to Approximation Variational Approximations we cover

• Original variational representation of log partition function

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$
(18.1)

where dual takes form:

$$A^{*}(\mu) = \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^{\circ} \\ +\infty & \text{if } \mu \notin \overline{\mathcal{M}} \end{cases}$$
(18.2)

• Given efficient expression for $A(\theta)$, we can compute marginals of interest.

 Above expression (dual of the dual) offers strategies to approximate or (upper or lower) bound A(θ). We either approximate M or -A*(μ) or (most likely) both.

Set $\mathcal{M} \leftarrow \mathbb{L}$ and $-A^*(\mu) \leftarrow H_{\mathsf{Bethe}}(\tau)$ to get Bethe variational approximation, LBP fixed point.

2 Set
$$\mathcal{M} \leftarrow \mathbb{L}_t(G)$$
 (hypergraph marginal polytope), $-A^*(\mu) \leftarrow H_{app}(\tau)$
Prof. Jeff Bilmes $\mathcal{L}_{g \in E}$ (g) Fig. (g) (via the Bill 2014 (pg.5/45)) (via the Bill 2014 (pg.5/45)

variational approximation, message passing on hypergraphs.

Service Partition τ into $(\tau, \tilde{\tau})$, and set $\mathcal{M} \leftarrow \mathcal{L}(\phi, \Phi)$ and set $-A^*(\mu) \leftarrow H_{ep}(\tau, \tilde{\tau})$ to get expectation propagation.

Logistics

Review

Variational Approach Amenable to Approximation Variational Approximations we cover

• Original variational representation of log partition function

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$
(18.1)

where dual takes form:

$$A^{*}(\mu) = \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^{\circ} \\ +\infty & \text{if } \mu \notin \overline{\mathcal{M}} \end{cases}$$
(18.2)

- Given efficient expression for $A(\theta)$, we can compute marginals of interest.
- Above expression (dual of the dual) offers strategies to approximate or (upper or lower) bound A(θ). We either approximate M or -A*(μ) or (most likely) both.
- Set $\mathcal{M} \leftarrow \mathbb{L}$ and $-A^*(\mu) \leftarrow H_{\mathsf{Bethe}}(\tau)$ to get Bethe variational approximation, LBP fixed point.
- Set $\mathcal{M} \leftarrow \mathbb{L}_t(G)$ (hypergraph marginal polytope), $-A^*(\mu) \leftarrow H_{app}(\tau)$ Prof. Jeff Bilmes Where Happ $-\frac{\text{EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014}{\sum g \in E \cup (9) \to g(\tau g)}$ (Via Wienerse) to get remember F6/45 (pg.6/45)
- variational approximation, message passing on hypergraphs.

EP as variational: Summary of key points

• Fixed points of EP exist assuming Lagrangian form has at least one optimum. No guarantees that EP will converge, but if it does it will be at a stationary point of the Lagrangian. • EP can be seen to be based on variational framework, using Bethe-like entropy and convex outer bound for the mean parameters. • When base distribution is unaries and Φ^i is the edges of a graph, we in fact get standard Bethe approximation, and standard sum-product LBP. • Moment matching of EP can be seen as striving for solution of associated Lagrangian. • Lost of flexibility here, depending on what the base distribution is (e.g., could be a k-tree, clusters, or many other structures as well). • Can also be done for Gaussian mixture and other distributions. • Many more details, variations, and possible roads to new research. See text and also see Tom Minka's papers. http://research.microsoft.com/en-us/um/people/minka/papers/ EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014 F7/45 (pg.7/45)

Logistics	Review
Mean Field	

- So far, we have been using an outer bound on \mathcal{M} .
- In mean-field methods, we use an "inner bound", a subset of \mathcal{M} constructed so as to make the optimization of $A(\theta)$ easier.
- Since subset, we get immediate bound on $A(\theta)$, all else (i.e., the entropy) being equal.
- Key: we based the inner bound on a "tractable family" like a 1-tree or even a 0-tree (all independent) so that the variational problem can be computed efficiently.
- Convexity of the optimization problem is often lost still, however, in the general case (due to the inner bound).
- Thus, in mean field, we will get a lower bound on $A(\theta)$ but not a convex procedure to find it (both good and bad news).

Tractable Families (for mean field approach)

- We have graph G = (V, E) which is intractable and we find a spanning subgraph (recall, spanning = all nodes, subgraph = subset of edges), i..e, F = (V, E_F) where E_F ⊆ E.
- Simplest example: $F = (V, \emptyset)$ all independence model.
- Tree example: F = (V, E_T) where edges E_T ⊂ E constitute a spanning tree.
- Exponential family, sufficient statistics $\phi = (\phi_{\alpha}, \alpha \in \mathcal{I})$ associated with this family $\mathcal{I}(F) \subseteq \mathcal{I}$. These are the statistics that need respect the Markov properties of only the subgraph F.
- Ω gets smaller too, canonical *F*-respecting parameters are of the form:

$$\mathbb{R}^{|\mathcal{I}|} \ni \Omega(F) \triangleq \{\theta \in \Omega | \theta_{\alpha} = 0 \ \forall \alpha \in \mathcal{I} \setminus \mathcal{I}(F)\} \subseteq \Omega.$$
(18.14)

Notice, all parameters associated with sufficient statistic not in $\mathcal{I}(F)$ are set to zero, those statistics are nonexistent in F.

• If parameter was not zero, model would not respect the familiy of F.

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014

Prof. Jeff Bilmes

Logistics

Inner bound Approximate Polytope

- Before, we had $\mathcal{M}(G; \phi) (= \mathcal{M}_G(G; \phi))$, all possible mean parameters associated with G and associated set of sufficient statistics ϕ .
- For a given subgraph *F*, we only consider those mean parameters possible under *F*-respecting models. I.e.,

$$\mathcal{M}_F(G;\phi) = \left\{ \mu \in \mathbb{R}^d | \mu = \mathbb{E}_{\theta}[\phi(x)] \text{ for some } \theta \in \Omega(F) \right\}$$
(18.18)

• Therefore, since $\theta \in \Omega(F) \subseteq \Omega$, we have that

$$\mathcal{M}_F^{\circ}(G;\phi) \subseteq \mathcal{M}^{\circ}(G;\phi) \tag{18.19}$$

and so $\mathcal{M}_F^\circ(G;\phi)$ is an $\mbox{ inner approximation of the set of realizable mean parameters.$

• Shorthand notation: $M_F^\circ(G) = M_F^\circ(G;\phi)$ and $M^\circ(G) = M^\circ(G;\phi)$

Review

F9/45 (pg.9/45)

Review

Tractable Dual

- Normally dual $A^*(\mu) = \sup_{\theta \in \Omega} (\langle \theta, \mu \rangle A(\theta))$ is intractable or unavailable, but key idea is that if $\mu \in \mathcal{M}_F(G)$ it will be possible to compute easily.
- Thus, goal of mean field (from variational approximation perspective) is to form $A_{MF}(\theta)$ where:

$$A(\theta) \ge \max_{\mu \in \mathcal{M}_F(G)} \left\{ \langle \mu, \theta \rangle - A_F^*(\mu) \right\} \triangleq A_{\mathsf{MF}}(\theta)$$
(18.23)

where $A_F^*(\mu)$ corresponds to dual function restricted to inner bound set $\mathcal{F}(G)$. I.e., when we expand $A_F^*(\mu)$, we can take advantage of the fact that μ is restricted in all cases, so $A_F^*(\mu)$ might be greatly simplified relative to $A^*(\mu)$.

• Note, for $\mu \in \mathcal{M}_F(G)$ and since $\mathcal{M}_F(G) \subseteq \mathcal{M}(G)$, $A_F^*(\mu)$ is not an approximation, rather it is just easy to compute.

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014

F11/45 (pg.11/45)

Mean field, KL-Divergence, Exponential Model Families

• Thus, solving the mean-field variational problem (see Eqn. (18.23)) of:

$$\max_{\mu \in \mathcal{M}_F(G)} \left\{ \langle \mu, \theta \rangle - A_F^*(\mu) \right\} = \max_{\mu \in \mathcal{M}_F(G)} \left\{ \langle \mu, \theta \rangle - A^*(\mu) \right\}$$
(18.34)

is identical to minimizing KL Divergence $D(\mu || \theta)$ subject to constraint $\mu \in \mathcal{M}_F(G)$.

• I.e., mean field can be seen as finding the best approximation, in terms of this particular KL-divergence, to p_{θ} , over a family of "nice" distributions $M_F(G)$.

Review

F13/45 (pg.13/45)

Logistics

Naïve Mean field for Ising Model: optimization

• We get variational lower bound problem

$$A(\theta) \ge \max_{(\mu_1,...,\mu_m)\in[0,1]^m} \left\{ \sum_{s\in V} \theta_s \mu_s + \sum_{(s,t)\in E} \theta_{st} \mu_s \mu_t + \sum_{s\in V} H_s(\mu_s) \right\}$$
(18.35)

- Have constrained form of edge mean parameters $\mu_{st} = \mu_s \mu_t$
- $(\mu_1, \ldots, \mu_m) \in [0, 1]^m$ is *m*-D hypercube.
- We have a non-convex problem, so while it is a bound, it might be hard to get as tight as possible.
- One way to optimize is to do coordinate ascent (given otherwise fixed vector, optimize one value at a time).
- If each coordinate optimization is optimal, we'll get a stationary point.

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014

• Fortunately, each coordinate optimization is concave!

```
6 Key idea: set of sufficient statistics that yield efficient inference need not be all independence. Could be a tree, or a chain, or a set of trees/chains.
6 "structured" in general means that it is not a monolithic single variable, but is a vector with some decomposability properties.
9 In Structured mean field, we exploit this and it again can be seen in our variational framework.
9 We first see a nice way that we can use fixed points of the mean field primal/dual equations to derive a general form of the mean field update.
```

Prof. Jeff Bilmes

eld Cnv

Tree Re-weighted (

Structured Mean Field

- Again, $\mathcal{I}(F)$ is set of suff. stats. corresponding to F, and we have corresponding mean vector $\mu(F) = (\mu_{\alpha}, \alpha \in \mathcal{I}(F))$.
- Define new quantity $\mathcal{M}(F)$, the set of realizable mean parameters associated with F, so that $\mu(F) \in \mathcal{M}(F)$. Thus, $\mathcal{M}(F) \subseteq \mathbb{R}^{|\mathcal{I}(F)|}$.
- Note also, $\mathcal{M}(F) \neq \mathcal{M}_F(G)$, their dimensions are entirely different.
- Key thing: in mean field, $\mu(F) \in \mathcal{M}(F)$ and there is no real need to mention the full $M_F(G)$. Also, the dual A_F^* depends on only $\mu(F)$ not μ (the other values are derivations from entries within $\mu(F)$.
- Other mean parameters μ_β for β ∈ I \ I(F) do play a role in the value of the mean field variational problem but their value is derivable from values μ(F), thus we can express the μ_β in functional form based on values μ(F).
- Thus, for each $\beta \in \mathcal{I} \setminus \mathcal{I}(F)$, we set $\mu_{\beta} = g_{\beta}(\mu(F))$ for function g_{β} .
- Ex: mean field Ising, edges $(s,t) \in E$, get $\mu_{st} = g_{st}(\mu(F)) = \mu_s \mu_t$.

```
Prof. Jeff Bilme
```

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014

.5/45 (pg.15/45)

Str. Mean Field Cnvx Relax/Up. Bounds Tree Re-weighted Case Refs Structured Mean Field Refs

• The mean field optimization problem becomes

$$\max_{\mu \in \mathcal{M}_{F}(G)} \left\{ \langle \mu, \theta \rangle - A_{F}^{*}(\mu) \right\}$$
(18.1)
$$= \max_{\mu(F) \in \mathcal{M}(F)} \left\{ \underbrace{\sum_{\alpha \in \mathcal{I}(F)} \theta_{\alpha} \mu_{\alpha} + \sum_{\alpha \in \mathcal{I}^{c}(F)} \theta_{\alpha} g_{\alpha}(\mu(F)) - A_{F}^{*}(\mu(F))}_{f(\mu(F))} \right\}$$
(18.2)

• With this, we can recover our sigmoid mean field coordinate update process by iterating fixed point equations of f, i.e., for $\beta \in \mathcal{I}(F)$,

$$\frac{\partial f}{\partial \mu_{\beta}}(\mu(F)) = \theta_{\beta} + \sum_{\alpha \in \mathcal{I}(G) \setminus \mathcal{I}(F)} \theta_{\alpha} \frac{\partial g_{\alpha}}{\partial \mu_{\beta}}(\mu(F)) - \frac{\partial A_{F}^{*}}{\partial \mu_{\beta}}(\mu(F))$$
(18.3)



- After each update of Eqn. (18.5), a mean parameter, say $\mu(F)_{\delta}$, that depends on any of the updated canonical parameter also needs to be updated before doing the next update.
- Since we're using a tractable sub-structure *F*, we can then update the out-of-date mean parameters using any exact inference algorithm (e.g., junction tree, possible since sub-structure is tractable), and then repeat Eqn. (18.5).

Mean Field Cnvx Relax/Up. Bounds

Structured Mean Field

- Alternatively, we can transform back to mean parameters right away using ∇A is the forward mapping, maping from canonical to mean.
- I.e., we can derive a mean field mean parameter to mean parameter update equation using A_F since $\nabla A_F(\gamma(F)) = \mu(F)$,
- We get update, for $\beta \in \mathcal{I}(F)$:

Prof leff Bilme

$$\mu_{\beta}(F) \leftarrow \frac{\partial A_F}{\partial \gamma_{\beta}} \left(\theta_{\beta} + \sum_{\alpha \in \mathcal{I}(G) \setminus \mathcal{I}(F)} \theta_{\alpha} \nabla g_{\alpha}(\mu(F)) \right)$$
(18.6)

- This generalizes our mean field coordinate ascent update from before, where in that case we would get $\frac{\partial A_F}{\partial \gamma_{\beta}}$ as being the sigmoid mapping.
- But here, we can use this for any tractable substructure (e.g., trees or chains or collections thereof).

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st. 2014



F19/45 (pg.19/45)



Convex Relaxations and Upper Bounds

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$
(18.7)

- Other than mean field (which gives lower bound on $A(\theta)$), none of the other approximation methods have been anything other than approximation methods.
- What about upper bounds?
- We would like both lower and upper bounds of A(θ) since that will allow us to produce upper and lower bounds of the probabilistic queries we wish to perform.
- If the upper and lower bounds between a given probably p is small, $p_L \leq p \leq p_U$, with $p_U - p_L \leq \epsilon$, we have guarantees, for a particular instance of a model.
- In this next chapter (Chap 7), we will "convexify" $H(\mu)$ and at the same time produce upper bounds.

```
Prof. Jeff Bilmes
```

 $\mathsf{EE512a}/\mathsf{Fall}$ 2014/Graphical Models - Lecture 18 - Dec 1st

3/45 (pg.23/45

Cnvx Relax/Up. Bounds Convex Relaxations and Upper Bounds - Relaxed Entropy • Recall sufficient stats $\phi = (\phi_{\alpha}, \alpha \in \mathcal{I})$ and canonical parameters $\theta = (\theta_{\alpha}, \alpha \in \mathcal{I}).$ • In general, inference (computing mean parameters) starting from canonical parameters is hard for a given G. • For a tractable subgraph F, it is not so hard, as we saw in the mean field case. Note in mean field case, we had one particular F. • Let \mathfrak{D} be a set of subfamilies that are tractable. • I.e., \mathfrak{D} might be all spanning trees of G, or some subset of spanning trees that we like. • As before, $\mathcal{I}(F) \subseteq \mathcal{I}$ are the subset of indices of the suff. stats. that abide by F, and $|\mathcal{I}(F)| = d(F) < d = |\mathcal{I}|$ suff. stats. • As before, $\mathcal{M}(F)$ is set of realizable mean parameters associated with F, and $\mu(F) \in \mathcal{M}(F)$. Thus, $\mathcal{M}(F) \subseteq \mathbb{R}^{|\mathcal{I}(F)|}$, and $\mathcal{M}(F) = \left\{ \mu \in \mathbb{R}^{|\mathcal{I}(F)|} | \exists p \text{ s.t. } \mu_{\alpha} = \mathbb{E}_p[\phi_{\alpha}(X)] \ \forall \alpha \in \mathcal{I}(F) \right\}$ (18.8) Note $\mathcal{M}_F(G) \neq \mathcal{M}(F)$. EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014 Prof leff Bilmes F24/45 (pg.24/45)



Str. Mean Field Cnvx Relax/Up. Bounds Tree Re-weighted Case Refs Convex Relaxations and Upper Bounds - Relaxed Entropy

Proposition 18.4.1 (Maximum Entropy Bounds)

Given any mean parameter $\mu \in \mathcal{M}$ and its projection $\mu(F)$ onto any subgraph F, we have the bound

$$A^*(\mu(F)) \le A^*(\mu)$$
 (18.9)

or alternatively stated, $H(\mu(F)) \ge H(\mu)$, entropy of projection is higher.

- Intuition: $H(\mu) = H(p_{\mu})$ is the entropy of the exponential family model with mean parameters μ .
- equivalently H(μ) = H(p_μ) is the entropy of the distribution that is the solution to the maximum entropy problem subject to the constraints that it has μ = E_{p_θ}[φ(X)].
- Fewer constraints when forming $\mu(F)$ (see Eqn. (18.8)), so entropy in corresponding maxent problem can only, if anything, get larger.
- Thus, $H(\mu(F)) \ge H(\mu)$.



Str. Mean Field Cnvx Relax/Up. Bounds Tree Re-weighted Case Refs Convex Relaxations and Upper Bounds - Relaxed Entropy

- Note that the upper bound is true for each $F \in \mathfrak{D}$, and thus would be true for mixtures of different $F \in \mathfrak{D}$.
- We can form a distribution over tractable structures, i.e., $\rho \in \mathbb{R}^{|\mathfrak{D}|}$, i.e., $\rho(F) \geq 0$ for $F \in \mathfrak{D}$ and $\sum_{F \in \mathfrak{D}} \rho(F) = 1$
- Convex combination over $F \in \mathfrak{D}$, gives more general upper bound

$$H(\mu) \le \mathbb{E}_{\rho}[H(\mu(F))] = \sum_{F \in \mathfrak{D}} \rho(F)H(\mu(F))$$
(18.13)

- This will be our convexified upper bound on entropy (lower bound on the dual).
- compared to mean field, we are not choosing only one structure, but many of them, and mixing them together in a certain way.
- This so far gives us an upper bound on $A(\theta)$, but we still need an outer bound. The combination will give us our uppper bound on $A(\theta)$.



Str. Mean Field	Cnvx Relax/Up. Bounds	Tree Re-weighted Case	Refs
Convex Upper	Bounds		

• Combining the upper bound on entropy, and the outer bound on \mathcal{M} , we get a new variational approximation to the cumulant function.

$$B_{\mathfrak{D}}(\theta;\rho) \stackrel{\Delta}{=} \sup_{\tau \in \mathcal{L}(G;\mathfrak{D})} \left\{ \langle \tau, \theta \rangle + \sum_{F \in \mathfrak{D}} \rho(F) H(\tau(F)) \right\}$$
(18.15)

- Objective is convex in θ since it is a max over a set of affine functions of θ (i.e., $g(\theta) = \max_{\tau} \langle \tau, \theta \rangle + c_{\tau}$)
- Evaluating the objective (optimization) is concave, so possible to get!
- Also, $\mathcal{L}(G;\mathfrak{D})$ is a convex outer bound on $\mathcal{M}(G)$
- Thus B_D(θ; ρ) is convex, has a global optimal solution, it approximates A(θ), and best of all is an upper bound, A(θ) ≤ B_D(θ; ρ)

Str. Mean Field

. 2. Relax/Up. Bounds

Tree Re-weighted (

Tree-reweighted sum-product and Bethe

- We can get convex upper bounds in the tree case, and a new style of sum-product algorithm.
- Consider MRF again

$$p_{\theta}(x) \propto \exp\left\{\sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t)\right\}$$
(18.16)

- Let \mathfrak{T} be a set of all spanning trees T of G, and let ρ be a distribution over them, $\sum_{T \in \mathfrak{T}} \rho(T) = 1$.
- Thus, we have $H(\mu) \leq \sum_{T \in \mathfrak{T}} \rho(T) H(\mu(T))$
- For any T, $H(\mu(T))$ has an easy form, i.e.,

$$H(\mu(T)) = \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E(T)} I_{st}(\mu_{st})$$
(18.17)

• We want to use this to see what happens when we take the expected value w.r.t. distribution ρ .

```
Prof. Jeff Bilmes
```

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014

F31/45 (pg.31/45)

 Str. Mean Field
 Cnvx Relax/Up. Bounds
 Tree Re-weighted Case
 Refs

 Tree-reweighted sum-product and Bethe

- Every tree is spanning, all tress have all nodes, so the probability, according to ρ of a given node is always 1. I.e., $\rho_s = 1, \forall s \in V$.
- Thus, in $\mathbb{E}_{\rho}[H(\mu(T))]$, we have a term of the form $\sum_{s \in V} H_s(\mu_s)$.
- For edges we need $\rho_{st} = \mathbb{E}_{\rho}[\mathbb{I}[(s,t) \in E(T)]]$, this indicates the probability of presence of an edge in the set \mathfrak{T} .
- The expression becomes

$$H(\mu) \le \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} \rho_{st} I_{st}(\mu_{st})$$
(18.18)

Note right hand sum is over all E (not just a given spanning tree) and terms are weighted by probability of the given edge ρ_{st} .

• ρ_{st} is edge appearance probability, $\rho = (\rho_{st}, (s, t) \in E)$ is spanning tree polytope.



Str. Mean Field	Cnvx Relax/Up. Bounds	Tree Re-weighted Case	Refs I
Tree-reweig	ghted sum-product a	nd Bethe	

- \bullet We also need outer bound on $\mathcal{M}.$
- For discrete case $\mathcal{M} = \mathbb{M}(G)$ is marginal polytope.
- M(T) is marginal polytope for tree, and for a tree is the same as L(T), the locally consistent pseudo-marginals (which recall are marginals for a tree).
- Thus, $\mu(T) \in \mathbb{M}(T)$ requires non-negativity, sum-to-one (at each node), and edge-to-node consistency (marginalization) on each edge. If G = T then we're done.
- For general G, If we ask for $\mu(T) \in \mathbb{M}(T)$ for all $T \in \mathfrak{T}$, this is identical to asking for local marginalization on every edge of G.
- Thus, in this case $\mathcal{L}(G; \mathfrak{I})$ is just the set of locally consistent pseudomarginals, and is the same as the outer bound we saw in the Bethe variational approximation $\mathbb{L}(G)$.
- In Bethe case, however, we did not have a bound on entropy, only an outer bound on the marginal polytope. Now, however, we also have a (convexification based) bound on entropy.



Cnvx Relax/Up. Bounds

Rets I

Tree-reweighted sum-product and Bethe

Theorem 18.5.1 (Tree-Reweighted Bethe and Sum-Product)

(a) For any choice of edge appearance vector $\rho = (\rho_{st}, (s, t) \in E)$ in the spanning tree polytope, the cumulant function $A(\theta)$ evaluated at θ is upper bounded by the solution of the tree reweighted Bethe variational problem (BVP):

$$B_{\mathfrak{T}}(\theta;\rho) = \max_{\tau \in \mathbb{L}(G)} \left\{ \langle \tau, \theta \rangle + \sum_{s \in V} H_s(\tau_s) - \sum_{(s,t) \in E} \rho_{st} I_{st}(\tau_{st}) \right\}$$

$$\geq A(\theta)$$
(18.20)

For any edge appearance vector such that $\rho_{st} > 0$ for all edges (s, t), this problem is strictly convex with a unique optimum.

Prof. Jeff Bilmes

. . .

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014

F35/45 (pg.35/45)

Str. Mean Field Cnvx Relax/Up. Bounds Tree Re-weighted Case Tree-reweighted sum-product and Bethe

Theorem 18.5.1 (Tree-Reweighted Bethe and Sum-Product)

(b) The tree-reweighted BVP can be solved using the tree-reweighted sum-product updates

$$M_{t \to s}(x_s) \leftarrow \kappa \sum_{x'_t \in \mathcal{X}_t} \varphi_{st}(x_s, x'_t) \frac{\prod_{v \in N(t) \setminus \{s\}} [M_{v \to t}(x'_t)]^{\rho_{vt}}}{[M_{s \to t}(x'_t)]^{(1-\rho_{ts})}} \quad (18.21)$$

where $\varphi_{st}(x_s, x'_t) = \exp\left(\frac{1}{\rho_{st}}\phi_{st}(x_s, x'_t) + \theta_t(x'_t)\right)$. The updates have a unique fixed point under assumptions given in (a).

Str. Mean Field

Cnvx Relax/Up. Bounds

Tree Re-weighted C

Tree-reweighted sum-product and Bethe

- Note that if $\rho_{st} \leftarrow 1$, for all $(s,t) \in E$, then we recover standard LBP and Bethe approximation.
- However, if $\rho_{st} = 1$ then edge (s,t) appears in all spanning trees. If this is indeed true for all spanning trees T, then G must be a tree, and we get back standard tree-based message passing we saw in lecture 2!!
- Thus, this is a true convex generalization, when $\rho_{st} < 1$ for many s, t.
- Note that ρ = (ρ_{st}, (s, t) ∈ E) must live in the "spanning tree polytope" ⊆ ℝ^E₊, i.e., a convex combination of vertices consisting of characteristic (indicator) functions of spanning trees (see example earlier). I.e., Let ℑ be the set of all spanning trees, and 1_T ∈ {0,1}^E be the characteristic vector of T ∈ ℑ. Then we must have that

$$\rho \in \operatorname{conv}(\{\mathbf{1}_T : T \in \mathfrak{T}\}) \tag{18.22}$$

where $conv(\cdot)$ is the convex hull of its argument.

Prof. Jeff Bilme

EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014

F37/45 (pg.37/45

 Str. Mean Field
 Cnvx Relax/Up. Bounds
 Tree Re-weighted Case
 Refs

 More on spanning tree polytope
 Image: Comparison of the state of

• Spanning tree polytope takes the form

$$\rho \in \operatorname{conv}(\{\mathbf{1}_T : T \in \mathfrak{T}\}) \tag{18.23}$$

where \mathfrak{T} is set of all spanning trees.

- Consider graphic matroid on G = (V, E) with rank function r(A) for any $A \subseteq E$.
- Then A is a spanning tree iff r(A) = |A| and |A| = m 1.
- Consider polytopes:

$$P_r = \left\{ x \in \mathbb{R}^E_+ : x(A) \le r(A), \forall A \subseteq E \right\}$$
(18.24)

$$B_r = P_r \cap \left\{ x \in \mathbb{R}^E_+ : x(E) = r(E) \right\}$$
(18.25)

- Then if T is a spanning tree, $\mathbf{1}_T \in B_r$, and $B_r = \operatorname{conv}(\{\mathbf{1}_T : T \in \mathfrak{T}\}).$
- Edmonds showed that a simple fast greedy procedure will maximize a linear function over this polytope, and this can be useful for finding good points in the spanning tree polytope.

Str. Mean Field	Cnvx Relax/Up. Bounds	Tree Re-weighted Cas	e Refs I
Tree-reweig	hted sum-product:	convex vs. upp	per bound
 In above ca upper boun 	se, we have both a conv d property.	exification of the c	umulant and an
 It should be can have co 	e pointed out that these nvex without upper bou	are not mutual req nd and vice verse.	uirements: one
Prof. Jeff Bilmes	EE512a/Fall 2014/Graphical Models	s - Lecture 18 - Dec 1st, 2014	F39/45 (pg.39/45)

Str. Mean Field Cnvx Relax/Up. Bounds Tree Re-weighted Case Refs Tree-reweighted sum-product fixed point Image: Compare the sum-product fixed point Refs

The fixed point we ultimately reach has following form:

$$\tau_s^*(x_s) = \kappa \exp\{\theta_s(x_s)\} \prod_{v \in N(s)} [M_{v \to s}^*(x_s)]^{\rho_{vs}}$$
(18.26)

$$\tau_{st}^{*}(x_{s}, x_{t}) = \kappa \varphi_{st}(x_{s}, x_{t}) \frac{\prod_{v \in N(s) \setminus t} [M_{vs}^{*}(x_{s})]^{\rho_{vs}} \prod_{v \in N(t) \setminus s} [M_{vt}^{*}(x_{t})]^{\rho_{vt}}}{[M_{ts}^{*}(x_{s})]^{(1-\rho_{st})} [M_{st}^{*}(x_{t})]^{(1-\rho_{ts})}}$$
(18.27)

with $\varphi_{st}(x_s, x_t) = \exp\left\{\frac{1}{\rho_{st}}\theta_{st}(x_s, x_t) + \theta_s(x_s) + \theta_t(x_t)\right\}$ where the * versions are the final (convergent) messages.

• In practice: damping of messages *M* appears in practice to help reach convergence, where each new message is a convex mixture of the previous version of itself and the new message according to the equations.

Str. Mean Field

 $\mathsf{Cnvx}\;\mathsf{Relax}/\mathsf{Up}.\;\mathsf{Bounds}$

hypertree-reweighted sum-product

- Why stop at trees, instead could use hypertrees and then deduce a hypertree version of the reweighted BP algorithm.
- Example in book considers k-trees, with tree width at most t. I.e. $\mathfrak{T}(t)$.
- Then we get the same form of bounds

$$H(\mu) \le E_{\rho}[H(\mu(T))] = \sum_{T \in \mathfrak{T}(t)} \rho(T)H(\mu(T))$$
 (18.28)

but here T is over all valid k-trees.

• This leads to a convexified Kikuchi variational problem

$$A(\theta) \le B_{\mathfrak{B}(t)}(\theta;\rho) = \max_{\tau \in \mathbb{L}(G)} \left\{ \langle \tau, \theta \rangle + \mathbb{E}_{\rho}[H(\tau(T))] \right\}$$
(18.29)

same form (but different than) before.

• Optimizing ρ over hypertree polytope is hard, unfortunately.

Prof. Jeff Bilmes EE512a/Fall 2014/Graphical Models - Lecture 18 - Dec 1st, 2014 F41/45 (pg.41/45)

Str. Mean Field	Cnvx Relax/Up. Bounds	Tree Re-weighted Case	Refs
Reweighted EP			

- Other variational variants have convexified version.
- Convexified forms of EP

$$H_{ep}(\tau, \tilde{\tau}; \rho) = H(\tau) + \sum_{\ell=1}^{d_I} \rho(\ell) [H(\tau, \tilde{\tau}^{\ell}) - H(\tau)]$$
(18.30)

where $\sum_{\ell} \rho(\ell) = 1$.

- In this case, reweighted entropy is concave!
- Lagrangian formulation leads to solutions that are a form of "reweighted" EP, ideas which also are sometimes called "power EP" (blending the above reweighted sum-product ideas and EP).



 Str. Mean Field
 Cnvx Relax/Up. Bounds
 Tree Re-weighted Case
 R

 Variational Approach Amenable to Approximation
 Variational Approximations we cover
 R

• Original variational representation of log partition function

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$
(18.1)

where dual takes form:

$$A^{*}(\mu) = \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^{\circ} \\ +\infty & \text{if } \mu \notin \overline{\mathcal{M}} \end{cases}$$
(18.2)

- Given efficient expression for $A(\theta)$, we can compute marginals of interest.
- Above expression (dual of the dual) offers strategies to approximate or (upper or lower) bound A(θ). We either approximate M or -A*(μ) or (most likely) both.
- Set $\mathcal{M} \leftarrow \mathbb{L}$ and $-A^*(\mu) \leftarrow H_{\mathsf{Bethe}}(\tau)$ to get Bethe variational approximation, LBP fixed point.
- Set $\mathcal{M} \leftarrow \mathbb{L}_t(G)$ (hypergraph marginal polytope), $-A^*(\mu) \leftarrow H_{app}(\tau)$ Prof. Jeff Bilmes Where $H_{app} = \frac{\mathsf{EE512a}/\mathsf{Fall} \, 2014/\mathsf{Graphical Models} - \mathsf{Lecture 18} - \mathsf{Dec 1st}, \, 2014}{\sum g \in E \cup \{9\}} + \frac{\mathsf{F44}/45 \, (\mathsf{pg.44}/45)}{2g \in E \cup \{9\}} + \frac{\mathsf{F44}/45 \, (\mathsf{pg.44}/45)}$
- variational approximation, message passing on hypergraphs.

