



- 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 L19 (12/3): on junction trees
- L9 (10/27): inference on junction trees, semirings,
- L10 (11/3): conditioning, hardness, LBP

Prof. Jeff Bilmes

- 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):
- Final Presentations: (12/10):

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

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F3/39 (pg.3/39)





Review

Möbius Inversion Lemma and Inclusion-Exclusion

- For any $A \subseteq V$, define two functions $\Omega: 2^V \to \mathbb{R}$ and $\Upsilon: 2^V \to \mathbb{R}$.
- Then the above inclusion-exclusion principle is one instance of the more general Möbius Inversion lemma, namely that each of the below two equations implies the other.

$$\forall A \subseteq V : \Upsilon(A) = \sum_{B:B \subseteq A} \Omega(B)$$
(16.13)

$$\forall A \subseteq V : \Omega(A) = \sum_{B:B \subseteq A} (-1)^{|A \setminus B|} \Upsilon(B)$$
 (16.14)

- Möbius Inversion lemma is also used to prove the Hammersley-Clifford theorem (that factorization and Markov property definitions of families are identical for positive distributions).
- We use it here to come up with alternative expressions for the entropy and for the marginal polytope.

Möbius Inversion Lemma for posets

- Let \mathcal{P} be a partially ordered set with binary relation \preceq .
- A zeta function of a poset is a mapping $\zeta : \mathcal{P} \times \mathcal{P} \to \mathbb{R}$ defined by

$$\zeta(g,h) = \begin{cases} 1 & \text{ if } g \preceq h, \\ 0 & \text{ otherwise.} \end{cases}$$
(16.23)

- The Möbius function $\omega : \mathcal{P} \times \mathcal{P} \to \mathbb{R}$ is the multiplicative inverse of this function. It is defined recursively:
- $\omega(g,g) = 1$ for all $g \in \mathcal{P}$
- $\omega(g,h) = 0$ for all $h : h \not\preceq g$.
- Given $\omega(g, f)$ defined for f such that $g \preceq f \prec h$, we define

$$\omega(g,h) = -\sum_{\{f \mid g \leq f \prec h\}} \omega(g,f)$$
(16.24)

• Then, ω and ζ are multiplicative inverses, in that

$$\sum_{f \in \mathcal{P}} \omega(g, f) \zeta(f, h) = \sum_{\{f | g \leq f \leq h\}} \omega(g, f) = \delta(g, h)$$
(16.25)

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

Review

F7/39 (pg.7/39)

General Möbius Inversion Lemma for Posets

Lemma 16.2.8 (General Möbius Inversion Lemma)

Given real valued functions Υ and Ω defined on poset \mathcal{P} , then $\Omega(h)$ may be expressed via $\Upsilon(\cdot)$ via

$$\Omega(h) = \sum_{g \leq h} \Upsilon(g) \quad \text{for all } h \in \mathcal{P}$$
 (16.23)

iff $\Upsilon(h)$ may be expressed via $\Omega(\cdot)$ via

$$\Upsilon(h) = \sum_{g \leq h} \Omega(g) \omega(g, h) \quad \text{for all } h \in \mathcal{P}$$
(16.24)

When $\mathcal{P} = 2^V$ for some set V (so this means that the poset consists of sets and all subsets of an underlying set V) this can be simplified, where \preceq becomes \subseteq ; and \succeq becomes \supseteq , like we saw above. (see Stanley, "Enumerative Combinatorics" for more info.)

Prof. Jeff Bilmes

Back to Kikuchi: Möbius and expressions of factorization

Suppose we are given marginals that factor w.r.t. a hypergraph G = (V, E), so we have μ = (μ_h, h ∈ E), then we can define new functions φ = (φ_h, h ∈ E) via Möbius inversion lemma as follows

$$\log \varphi_h(x_h) \triangleq \sum_{g \leq h} \omega(g, h) \log \mu_g(x_g)$$
(16.23)

• From Möbius inversion lemma, this then gives us a new way to write the log marginals, i.e., as

$$\log \mu_h(x_h) = \sum_{g \preceq h} \log \varphi_g(x_g) \tag{16.24}$$

• Key, when φ_h is defined as above, and G is a hypertree we have

$$p_{\mu}(x) = \prod_{h \in E} \varphi_h(x_h) \tag{16.25}$$

 \Rightarrow general way to factorize a distribution that factors w.r.t. a hypergraph.

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F9/39 (pg.9/39)

• Using Möbius, and Eqn. (16.23) we can write

$$I_h(\mu_h) = \sum_{x_h} \mu_h(x_h) \log \varphi_h(x_h) = \sum_{x_h} \mu_h(x_h) \left(\sum_{g \leq h} \omega(g, h) \log \mu_g(x_g) \right)$$

$$= \sum_{g \leq h} \omega(g, h) \left\{ \sum_{x_h} \mu_h(x_h) \log \mu_g(x_g) \right\}$$

$$= \sum_{f \leq h} \sum_{e \geq f} \omega(f, e) \left\{ \sum_{x_f} \mu_f(x_f) \log \mu_f(x_f) \right\} = -\sum_{f \leq h} c(f) H_f(\mu_f)$$
where we define overcounting numbers (~ shattering coefficient)

$$c(f) \triangleq \sum_{e \geq f} \omega(f, e)$$
(16.31)

$$H_{\text{hyper}}(\mu) = \sum_{h \in E} c(h) H_h(\mu_h)$$
(16.32)

Review

Logistics

Usable to get Kikuchi variational approximation

• Sum to one constraint:

$$\sum_{x_h} \tau_h(x_h) = 1 \tag{16.33}$$

• Local agreement via the hypergraph constraint. For any $g \preceq h$ must have marginalization condition

$$\sum_{x_{h\setminus g}} \tau_h(x_h) = \tau_g(x_g) \tag{16.34}$$

• Define new polyhedral constraint set $\mathbb{L}_t(G)$

 $\mathbb{L}_t(G) = \{ \tau \ge 0 | \text{ Equations (16.3) } \forall h, \text{ and (16.34) } \forall g \le h \text{ hold} \}$ (16.35)

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F11/39 (pg.11/39)

Review

Kikuchi variational approximation, entropy approx

• Generalized approximate (app) entropy for the hypergraph:

$$H_{\mathsf{app}} = \sum_{g \in E} c(g) H_g(\tau_g) \tag{16.33}$$

where H_g is hyperedge entropy and overcounting number defined by:

$$c(g) = \sum_{f \succeq g} \omega(g, f)$$
(16.34)

Variational Approach Amenable to Approximation Variational Approximations we cover

• Original variational representation of log partition function

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$
(16.1)

where dual takes form:

$$A^{*}(\mu) = \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^{\circ} \\ +\infty & \text{if } \mu \notin \overline{\mathcal{M}} \end{cases}$$
(16.2)

• 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 ______EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014 ______F13/39 (pg.13/39)

variational approximation, message passing on hypergraphs.

- Solution 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.

Kikuchi and Hypertree-based Methods

EP like variants

Kikuchi variational approximation

• This at last gets the Kikuchi variational approximation

$$A_{\mathsf{Kikuchi}}(\theta) = \max_{\tau \in \mathbb{L}_t(G)} \left\{ \langle \theta, \tau \rangle + H_{\mathsf{app}}(\tau) \right\}$$
(16.1)

- For a graph, this is exactly $A_{\text{Bethe}}(\theta)$.
- Also, if hypergraph is junction tree (r.i.p. holds, tree-local consistency implies global consistency), then also exact (although expensive, exponential in the tree-width to compute H_{app}).
- We can define message passing algorithms on the hypertree, and show that if it converges, it is a fixed point of the associated Lagrangian.









- The node sets that communicate with each other represented using hypergraph (hyperedges are the ndoe sets)
- Standard LBP algorithm is merely a special case of GBP
- Different choices of node sets/hyperedges and message passings give different GBP algorithms.
- This gives the user a gradual tradeoff between the most expensive, intractable, and accurate junction tree algorithm, and the least expensive but possibly quite inaccurate LBP algorithm.
- Allows a trade-off between complexity for accuracy!
- In many cases, convergence of GBP will be at fixed points of the Lagrangian for the generalized variational approximation

$$A_{\mathsf{Kikuchi}}(\theta) = \max_{\tau \in \mathbb{L}_t(G)} \left\{ \langle \theta, \tau \rangle + H_{\mathsf{app}}(\tau) \right\}$$
(16.2)

Prof. Jeff Bilmes

Kikuchi and Hypertree-based Method

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F17/39 (pg.17/39)

GBP examples: parent-to-child

In hypergraph Hasse-like diagram,

arrows point from parent (superset) to child (subset). Ex: on the right, set {1,2,4,5} is the parent of both {2,5} and {4,5}.



- For h ∈ E, let Par(h) be the set of parents. Also define descendants as D(h) = {g ∈ E|g ≺ h} and ancestors as A(h) = {g ∈ E|g ≻ h}.
- Also define $\mathcal{D}^+(h) = \mathcal{D}(h) \cup \{h\}$ and $\mathcal{A}^+(h) = \mathcal{A}(h) \cup \{h\}$
- If $f \succ g$ then x_f has more variables than x_g and one can perform a message of the form $M_{f \rightarrow g}(x_g) = \sum_{f \setminus g} \tau(x_f) = \sum_{f \setminus g} \tau(x_g, x_{f \setminus g})$





EP like variants

Conjugate Duality, Maximum Likelihood, Negative Entropy

Theorem 16.4.3 (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:

$$A^{*}(\mu) = \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^{\circ} \\ +\infty & \text{if } \mu \notin \overline{\mathcal{M}} \end{cases}$$
(16.3)

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

 $A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$ (16.4)

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

$$\mu = \int_{\mathsf{D}_X} \phi(x) p_{\theta}(x) \nu(dx) = \mathbb{E}_{\theta}[\phi(X)] = \nabla A(\theta)$$
 (16.5)

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

Kikuchi and Hypertree-based Metho

Expectation Propagation: basic idea

- Came from a method called "assumed density filtering" (ADF).
- Doing full inference involves exponential computation.
- We do a bit of inference, involving reasonable computation, and getting us a new distribution that is a bit more complex but not too much more complex.
- Before going further, we "project" this new distribution back down to a class of simple distributions.
- We then repeat the above step with a bit more of inference, different than what we did above.
- We keep repeating: do a bit of inference, and project, until all inference has been done.
- The difference between ADF and EP is that, with ADF at this stage we're done. With EP we can keep repeating the process of inference, projection.
- EP can be seen as a generalization of BP.
- Interestingly, EP is instance of our variational framework, Equation

F21/39 (pg.21/39)



$$\phi \triangleq (\phi_1, \phi_2, \dots, \phi_{d_T}) \tag{16.5}$$

• Intractable component

$$\Phi \triangleq (\Phi^1, \Phi^2, \dots, \Phi^{d_I}) \tag{16.6}$$

- ϕ_i are typically univariate, while Φ^i are typically multivariate (*b*-dimensional we'll assume), although this need not always be the case (but will be for our exposition).
- Consider exponential families associated with subcollection (ϕ, Φ) .



EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F23/39 (pg.23/39)









EP like variants

Associated Distributions: key points

The basic premises in the tractable-intractable partitioning between ϕ and Φ are:

- It is possible to compute marginals exactly in polynomial time for distributions of the base form (any member of the φ-exponential family).
- For each i = 1,..., d_I, exact polynomial-time computation is still possible for any Φⁱ-augmented form (any member of the (φ, Φⁱ)-exponential family).
- Intractable to perform exact computations with the full (ϕ, Φ) -exponential family.

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F27/39 (pg.27/39)

Kikuchi and Hypertree-based Methods EP like variants Example: Mixture models

- Let $X \in \mathbb{R}^m$ be Gaussian with distribution $N(0, \Sigma)$.
- Let $\varphi(y; \mu, \Lambda)$ be Gaussian with mean μ covariance Λ .
- Suppose y conditioned on x is a two-component Gaussian mixture model taking the form:

$$p(y|X = x) = (1 - \alpha)\varphi(y; 0, \sigma_0^2 I) + \alpha\varphi(y; x, \sigma_1^2 I)$$
 (16.13)

- Assume we have obtained n i.i.d. samples y^1, \ldots, y^n from mixture density, and goal is to produce posterior $p(x|y^1, \ldots, y^n)$, similar to Bayes-rule inverting a Naive-Bayes model.
- Using Bayes rule, we get mixture model with 2^n components!

$$p(x|y^{1},\ldots,y^{n}) \propto \exp\left(-\frac{1}{2}x^{\mathsf{T}}\Sigma^{-1}x\right) \prod_{i=1}^{n} p(y^{i}|X=x)$$
(16.14)
$$= \exp\left(-\frac{1}{2}x^{\mathsf{T}}\Sigma^{-1}x\right) \exp\left\{\sum_{i=1}^{n} \log p(y^{i}|X=x)\right\}$$
(16.15)

EP like variants

Example: Mixture models

- We equate $\exp\left(-\frac{1}{2}x^{\mathsf{T}}\Sigma^{-1}x\right)$ with $\exp(\langle\theta,\phi(x)\rangle)$, with $d_T = m$.
- Such a distribution is multivariate Gaussian, and getting marginals (say p(x_A) for A ⊆ [m]) from it is relatively "cheap" O(m³).
- $\exp\left\{\sum_{i=1}^{n}\log p(y^{i}|X=x)\right\}$ equates to $\prod_{i=1}^{d_{I}}\exp\left(\left\langle \tilde{\theta}^{i}, \Phi^{i}(x)\right\rangle\right)$, with b=1. These are the intractable factors.
- Base distribution $p(x; \theta, \vec{0}) \propto \exp\left(-\frac{1}{2}x^{\intercal}\Sigma^{-1}x\right)$ which is a Gaussian and easy as mentioned above.
- If we multiply in only one intractable term, complexity to produce marginal still not so bad (quite easy in fact).
- I.e., Φ^i -augmented distribution is proportional to

$$\exp\left(-\frac{1}{2}x^{\mathsf{T}}\Sigma^{-1}x\right)\left[(1-\alpha)\varphi(y^{i};0,\sigma_{0}^{2}I)+\alpha\varphi(y^{i};x,\sigma_{1}^{2}I)\right] \quad \textbf{(16.16)}$$

- Computing marginals is easy (mixture of only 2 components)
- If we multiply in all Φ^i , becomes intractable (2ⁿ potentially distinct components each of which requires marginalization).

```
Prof. Jeff Bilmes
```

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

29/39 (pg.29/39

EP like variants

Polytope and Base case

- We can partition the mean parameters $(\mu, \tilde{\mu}) \in \mathbb{R}^{d_T + d_I imes b}$
- Marginal polytope associated with these means

$$\mathcal{M}(\phi, \Phi) = \{(\mu, \tilde{\mu}) | (\mu, \tilde{\mu}) = \mathbb{E}_p[(\phi(X), \Phi(X))] \text{ for some } p\}$$
(16.17)

along with negative dual of cumulant, or entropy $H(\mu,\tilde{\mu})=-A^*(\mu,\tilde{\mu}).$

• We also have polytope associated with only base distribution

$$\mathcal{M}(\phi) = \left\{ \mu \in \mathbb{R}^{d_T} | \mu = \mathbb{E}_p(\phi(X)) \right\}$$
(16.18)

• Recall thm: any mean in the interior is realizable via an exponential family model, and associated entropy $H(\mu)$ is tractable.

Augmented Base case

• For each $i = 1 \dots d_I$ we have a Φ^i -augmented exp. model and polytope

$$\mathcal{M}(\phi, \Phi^i) = \left\{ (\mu, \tilde{\mu}^i) \in \mathbb{R}^{d_T + b} | (\mu, \tilde{\mu}^i) = \mathbb{E}_p[(\phi(X), \Phi^i(X))] \text{ for some } p \right\}$$
(16.19)

- Thus, any such mean parameters has instance for associated exponential family, and also $H(\mu, \tilde{\mu}^i)$ is easy to compute.
- Goal, variational approximation: Need outer bounds on $\mathcal{M}(\phi, \Phi)$ and expression for entropy (as is now normal).
- Turns out we can do this, and an iterative algorithm to find fixed points of associated Lagrangian, that correspond to EP.

Prof. Jeff Bilmes

Kikuchi and Hypertree-based Methods

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F31/39 (pg.31/39)

New EP-based outer bound

• For any mean parms $(\tau, \tilde{\tau})$ where $\tilde{\tau} = (\tilde{\tau}^1, \tilde{\tau}^2, \dots, \tilde{\tau}^{d_I})$, define coordinate "projection operation"

$$\Pi^{i}(\tau,\tilde{\tau}) \to (\tau,\tilde{\tau}^{i})$$
(16.20)

like variants

This operator simply removes all but $\tilde{\tau}^i$ from $\tilde{\tau}$.

• Define outer bound on true means $\mathcal{M}(\phi, \Phi)$ (which is still convex)

$$\mathcal{L}(\phi, \Phi) = \left\{ (\tau, \tilde{\tau}) | \tau \in \mathcal{M}(\phi), \Pi^{i}(\tau, \tilde{\tau}) \in \mathcal{M}(\phi, \Phi^{i}), \forall i \right\}$$
(16.21)

- Note, based on a set of projections onto $\mathcal{M}(\phi, \Phi^i)$.
- Outer bound, i.e., $\mathcal{M}(\phi, \Phi) \subseteq \mathcal{L}(\phi, \Phi)$, since:

$$\tau \in \mathcal{M}(\phi) \Leftrightarrow \exists p \text{ s.t. } \tau = E_p[\phi(X)]$$
 (16.22)

$$(\tau, \tilde{\tau}) \in \mathcal{L}(\phi, \Phi) \Leftrightarrow \tau \in \mathcal{M}(\phi) \& \exists p \text{ s.t. } (\tau, \tilde{\tau}^i) = E_p[\phi(X), \Phi^i(X)]$$
(16.23)

$$(\tau, \tilde{\tau}) \in \mathcal{M}(\phi, \Phi) \Leftrightarrow \exists p \text{ s.t. } (\tau, \tilde{\tau}) = E_p[\phi(X), \Phi(X)]$$
 (16.24)

• If Φ^i are edges of a graph (i.e. local consistency) then we get standard \mathbb{L} outer bound we saw before with Bethe approximation

EP outer bound entropy and opt

- For any mean parms (τ, τ̃) ∈ L(φ, Φ): A) There is a member of the φ-exponential family which mean parameters τ with entropy H(τ); B) Also, for i = 1...d_I, there is a member of the (φ, Φⁱ)-exponential family with mean parameters (τ, τ̃ⁱ) with entropy H(τ, τ̃ⁱ).
- Both entropy forms are easy to compute, and so is a new entropy approximation:

$$H(\tau, \tilde{\tau}) \approx H_{ep}(\tau, \tilde{\tau}) \triangleq H(\tau) + \sum_{\ell=1}^{d_I} \left[H(\tau, \tilde{\tau}^l) - H(\tau) \right]$$
(16.25)

• With outer bound and entropy expression, we get new variational form

$$\max_{(\tau,\tilde{\tau})\in\mathcal{L}(\phi,\Phi)}\left\{\left\langle \tau,\theta\right\rangle + \left\langle \tilde{\tau},\tilde{\theta}\right\rangle + H_{\mathsf{ep}}(\tau,\tilde{\tau})\right\}$$
(16.26)

- This characterizes the EP algorithms.
- Given graph G = (V, E) when we take ϕ to be unaries V and Φ to be edges E, we exactly recover Bethe approximation.

```
Prof. Jeff Bilr
```

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F33/39 (pg.33/39)

kuchi and Hypertree-based Methods

EP like variants

Lagrangian optimization setup

- Make d_I duplicates of vector $\tau \in \mathbb{R}^{d_T}$, call them $\eta^i \in \mathbb{R}^{d_T}$ for $i \in [d_T]$.
- This gives large set of pseudo-mean parameters

$$\left\{\tau, (\eta^i, \tilde{\tau}^i), i \in [d_I]\right\} \in \mathbb{R}^{d_T} \times (\mathbb{R}^{d_T} \times \mathbb{R}^b)^{d_I}$$
(16.27)

• We arrive at the optimization:

$$\max_{\left\{\tau, \left\{(\eta^{i}, \tilde{\tau}^{i})\right\}_{i}\right\}} \left\{ \langle \tau, \theta \rangle + \sum_{i=1}^{d_{I}} \left\langle \tilde{\tau}^{i}, \tilde{\theta}^{i} \right\rangle + H(\tau) + \sum_{i=1}^{d_{I}} \left[H(\eta^{i}, \tilde{\tau}^{i}) - H(\eta^{i}) \right] \right\}$$
(16.28)

subject to $\tau \in \mathcal{M}(\phi)$, and for all i that $\tau = \eta^i$ and that $(\eta^i, \tilde{\tau}^i) \in \mathcal{M}(\phi, \Phi^i)$.

• Use Lagrange multipliers to impose constraint $\eta^i = \tau$ for all i, and for the rest of the constraints too.

To Lagrangian optimization

• We get a Lagrangian version of the objective

$$L(\tau;\lambda) = \langle \tau, \theta \rangle + \sum_{i=1}^{d_I} \left\langle \tilde{\tau}^i, \tilde{\theta}^i \right\rangle + F(\tau; (\eta^i, \tilde{\tau}^i)) + \sum_{i=1}^{d_I} \left\langle \lambda^i, \tau - \eta^i \right\rangle + \dots$$
(16.29)

where

$$F(\tau; (\eta^{i}, \tilde{\tau}^{i})) = H(\tau) + \sum_{i=1}^{d_{I}} \left[H(\eta^{i}, \tilde{\tau}^{i}) - H(\eta^{i}) \right]$$
(16.30)

and where λ^i are the Lagrange multipliers assocaited with the constraint $\eta^i = \tau$ for all *i* (other multipliers not shown).

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 16 - Nov 24th, 2014

F35/39 (pg.35/39

1 Considering optimality conditions on what must hold for a solution $\{\tau, (\eta^i, \tilde{\tau}^i), i \in [d_I]\}$ to the above Lagrangian, must have properties: • τ belongs to relative interior, i.e., $\tau \in \mathcal{M}^{\circ}(\theta)$ of the base model. • $(\eta^i, \tilde{\tau}^i)$ belongs to relative interior of extended model, so $(\eta^i, \tilde{\tau}^i) \in \mathcal{M}^{\circ}(\phi, \Phi^i)$. • Means must agree, i.e., $\tau = \eta^i$ for all *i*. • First condition means we're a member of the ϕ -exponential family, and (it can be shown) has form: $q(x; \theta, \lambda) \propto \exp\left\{\left\langle \theta + \sum_{i=1}^d \lambda^i, \phi(x) \right\rangle\right\}$ (16.31) • Second condition means we're a member of the (ϕ, Φ^i) -exponential family, and (it can be shown) has form: $q^i(x, \theta, \tilde{\theta}^i, \lambda) \propto \exp\left(\left\langle \theta + \sum_{\ell \neq i} \lambda^\ell, \phi(x) \right\rangle + \left\langle \tilde{\theta}^i, \Phi^i(x) \right\rangle\right)$ (16.32)

≺ets I



