



Class Road Map - EE512a

- 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

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

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

Prof. Jeff Bilmes

 L11 (11/5): LBP, exponential models, mean params and polytopes

- L13 (11/10):
- L14 (11/12):
- L15 (11/17):
- L16 (11/19):
- L17 (11/24):
- L18 (11/26):
- L19 (12/1):
- L20 (12/3):
- Final Presentations: (12/10):

F3/60 (pg.3/60)



Review

Belief Propagation: message definition

Generic message definition

$$\mu_{i \to j}(x_j) = \sum_{x_i} \psi_{i,j}(x_i, x_j) \prod_{k \in \delta(i) \setminus \{j\}} \mu_{k \to i}(x_i)$$
(11.5)

• If graph is a tree, and if we obey MPP order, then we will reach a point where we've got marginals. I.e.,

$$p(x_i) \propto \prod_{j \in \delta(i)} \mu_{j \to i}(x_i)$$
 (11.6)

and

$$p(x_i, x_j) \propto \psi_{i,j}(x_i, x_j) \prod_{k \in \delta(i) \setminus \{j\}} \mu_{k \to i}(x_i) \prod_{\ell \in \delta(j) \setminus \{i\}} \mu_{\ell \to j}(x_j) M \quad (11.7)$$

```
Prof. Jeff Bilmes
```

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

F5/60 (pg.5/60)

Review

Choices for dealing with higher order factors in MRFs

So, to deal with MRFs with higher order factors, we can:

- Itransform MRF to have only pairwise interactions, add more variables. we can keep using BP on MRF edges (as done above), makes the math a bit easier, does not change fundamental computational cost. Possible since for any given p, we know the interaction terms.
- Alternatively, we can define BP on factor graphs.
- Alternatively, could define BP directly on the maxcliques of the MRF (but maxcliques are not easy to get in a MRF when not triangulated).

For the remainder of this term, we'll assume we've done the pair-wise transformation (i.e., option 1 above).



Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

F7/60 (pg.7/60)

Review

Messages as matrix multiply

=

$$\mu_{i \to j}(x_j) \propto \sum_{x_i} \psi_{i,j}(x_i, x_j) \psi_i(x_i) \prod_{k \in \delta(i) \setminus \{j\}} \mu_{k \to i}(x_i)$$
(11.9)

$$= \sum_{x_i} \psi'_{i,j}(x_i, x_j) \mu_{\neg j \to i}(x_i)$$
(11.10)

$$= (\psi'_{i,j})^T \mu_{\neg j \to i}$$
 (11.11)

- Here, $\psi'_{i,j}$ is a matrix and $\mu_{\neg j \rightarrow i}$ is a column vector.
- Going from state μ^t to μ^{t+1} is like matrix-vector multiply group messages from μ^t together into one vector representing $\mu_{\neg j \rightarrow i}$ for each $(i, j) \in E$, do the matrix-vector update, and store result in new state vector μ^{t+1} .
- If G is tree, μ^t will converged after D steps.

What if graph has cycles?

- MPP causes deadlock since there is no way to start sending messages
- Like before, we can assume that messages have an initial state, e.g., $\mu_{i \to j}(x_j) = 1$ for all $(i, j) \in E(G)$ note this is bi-directional. This breaks deadlock.
- We can then start sending messages. Will we converge after D steps? What does D even mean here?
- No, in fact we could oscillate forever.

LBP Next phase of class exponential models μ Param./Marg. Polytope Refs Belief Propagation, Cycles, and Oscillation Refs Refs Refs Refs

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

- Consider odd length cycle (e.g., C_3 , C_5 , etc.), C_3 is sufficient i j k i
- Assume all messages start out at state $\mu_{i \to j} = [1, 0]^T$.
- Consider (pairwise) edge functions, for each i, j

$$\psi_{ij}(x_i, x_j) = \begin{bmatrix} 0 & 1\\ 1 & 0 \end{bmatrix}$$
(11.1)

• then we have

Prof. Jeff Bilmes

$$\mu_{j \to k}(x_k) = \sum_{x_j} \psi_{j,k}(x_j, x_k) \mu_{i \to j}(x_j)$$
(11.2)

• or in matrix form

$$\mu_{j \to k} = (\psi_{j,k})^T \mu_{i \to j} \tag{11.3}$$

F9/60 (pg.9/60)

Belief Propagation, Cycles, and Oscillation

I the second second

- Let $\mu_{i\to j}^t$ be the $t^{\rm th}$ formed message, with $\mu_{i\to j}^0$ being the starting state at $[1,0]^T$.
- Then $\mu_{i\to j}^1=[0,1]^T$, $\mu_{i\to j}^2=[1,0]^T$, $\mu_{i\to j}^3=[0,1]^T$, and so on, never converging. In fact,

$$\mu_{i \to j}^{t+1} = (\psi_{i,j})^T \mu_{k \to i}^t$$
(11.4)

$$= (\psi_{i,j})^T (\psi_{k,i})^T \mu_{j \to k}^t$$
(11.5)

$$= (\psi_{i,j})^T (\psi_{k,i})^T (\psi_{j,k})^T \mu_{i \to j}^t$$
(11.6)

$$= \begin{bmatrix} 0 & 1\\ 1 & 0 \end{bmatrix}^{3} \mu_{i \to j}^{t}$$
(11.7)

$$= \begin{bmatrix} 0 & 1\\ 1 & 0 \end{bmatrix} \mu_{i \to j}^t \tag{11.8}$$

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

F11/60 (pg.11/60

 LBP
 Next phase of class
 exponential models
 µ Param./Marg. Polytope
 Refs

 Belief Propagation, Cycles, and Oscillation

- Thus, each time we go around the loop in the cycle, the message configuration for each (i, j) will flip, thereby never converging.
- Damping the messages? I.e., Let $0 \le \gamma < 1$ and treat messages as

$$\mu_{i \to j}^t \leftarrow \gamma \mu_{i \to j}^t + (1 - \gamma) \mu_{i \to j}^{t-1}$$
(11.9)

- Empirical Folklore if we converge quickly without damping, the quality of the resulting marginals might be good. If we don't converge quickly, w/o damping, might indicate some problem.
- Ways out of this problem: Other message schedules, other forms of the interaction matrices, and other initializations.



LBP Next phase of class exponential models μ Param./Marg. Polytope Refs Belief Propagation, Single Cycle

- Consider a graph with a single cycle C_{ℓ} .
- It could be a cycle with trees hanging off of each node. We send messages from the leaves of those dangling trees to the cycle (root) nodes, leaving only a cycle remaining.
- Consider what happens to $\mu_{i \to j}^t$ as t increases. w.l.o.g. consider $\mu_{\ell \to 1}^t$
- Let the cycle be nodes $(1, 2, 3, \dots, \ell, 1)$

$$\mu_{\ell \to 1}^{t+1} = \left(\prod_{i=1}^{\ell-1} (\psi_{i,i+1})^T \right) \mu_{\ell \to 1}^t$$
(11.10)

$$= M \mu_{\ell \to 1}^t \tag{11.11}$$

• Will this converge to anything?



LBP Next phase of class exponential models μ Param./Marg. Polytope Refs Belief Propagation, Single Cycle Refs

From this, we the following theorem follows almost immediately.

Theorem 11.3.2

1. $\mu_{\ell \to 1}$ converges to the principle eigenvector of M.

2. $\mu_{2\rightarrow 1}$ converges to the principle eigenvector of M^T .

3. The convergence rate is determined by the ratio of the largest and second largest eigenvalue of M.

4. The diagonal elements of M correspond to correct marginal $p(x_1)$

5. The steady state "pseudo-marginal" $b(x_1)$ is related to the true marginal by $b(x_1) = \beta p(x_1) + (1 - \beta)q(x_1)$ where β is the ratio of the largest eigenvalue of M to the sum of all eigenvalues, and $q(x_1)$ depends on the eigenvectors of M.

Proof.

See Weiss2000.



- It still does not tell us that we end up with correct marginals, rather we get "pseudo-marginals", which are locally normalized, but might not be the correct marginals.
- Moreover, they might not be the correct marginals for any probability distribution.
- Also, we'd like a characterization of LBP's convergence (if it happens) for more general graphs, with an arbitrary number of loops.



LBP	Next phase of class	exponential models	μ Param./Marg. Polytope	Refs I
exponential fam	ily models			

- φ = (φ_α, α ∈ I) is a collection of functions known as potential functions, sufficient statistics, or features. I is an index set of size d = |I|.
- Each φ_α is a function of x, φ_α(x) but it usually does not use all of x (only a subset of elements). Notation φ_α(x_{C_α}) assumed implicitly understood, where C_α ⊆ V(G).
- θ is a vector of canonical parameters (same length, |*I*|). θ ∈ Ω ⊆ ℝ^d where d = |*I*|.
- We can define a family as

$$p_{\theta}(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$$
(11.12)

Note that we're using ϕ here in the exponent, before we were using it out of the exponent.

Note that φ(x) = (φ₁(x), φ₂(x),..., φ_{|I|}) where again each φ_i(x) might use only some of the elements in vector x. φ : D_X^m → ℝ^d.



LBP Next phase of class exponential models μ Param./Marg. Polytope Refs exponential family models and clique features

- Example: single scalar discrete random variable X ∈ {1, 2, ..., k} might have indicator feature for all possible values α_i(x) ≜ 1(x = i)
 — in this case |C_α| = 1 for all α ∈ I.
- Could even think of {C_α}_{α∈𝒯} as cliques of some graph, but not necessarily maxcliques.
- Likely not dealing with triangulated models. Could be based on cliques, or cliques and subsets of cliques (consider 4-cycle with edges and vertices).
- Key: $p \in \mathcal{F}(G, \mathcal{M}^{(f)})$ by Hammersley-Clifford theorem,
 - where G = (V, E) where V is the nodes corresponding to vector x,
 - and E is formed by using $\{C_{\alpha}\}_{\alpha \in \mathcal{I}}$ as an edge clique cover: \exists an $\alpha \in \mathcal{I}$ such that $u, v \in C_{\alpha}$ where $u, v \in V(G) \Leftrightarrow$ there is an edge $(u, v) \in E(G)$.



Prof. Jeff Bilmes

exponential family models

• Based on underlying set of parameters θ , we have family of models

$$p_{\theta}(x) = \frac{1}{Z(\theta)} \exp\left\{\sum_{\alpha \in \mathcal{I}} \theta_{\alpha} \phi_{\alpha}(x)\right\} = \exp(\langle \theta, \phi(x) \rangle - A(\theta)) \quad (11.13)$$

• To ensure normalized, we use log partition (cumulant) function

$$A(\theta) = \log \int_{\mathsf{D}_X} \exp\left(\langle \theta, \phi(x) \rangle\right) \nu(dx) \tag{11.14}$$

with $\theta\in\Omega\stackrel{\Delta}{=}\left\{\theta\in\mathbb{R}^{d}|A(\theta)<+\infty\right\}$

- $A(\theta)$ is convex function of θ , so Ω is convex.
- Exponential family for which Ω is open is called regular

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2

F25/60 (pg.25/60)

 LBP
 Next phase of class
 exponential models
 µ. Param./Marg. Polytope
 Refs

 exponential family models

 Refs

• Based on underlying set of parameters θ , we have family of models

$$p_{\theta}(x) = \frac{1}{Z(\theta)} \exp\left\{\sum_{\alpha \in \mathcal{I}} \theta_{\alpha} \phi_{\alpha}(x)\right\} = \exp(\langle \theta, \phi(x) \rangle - A(\theta)) \quad (11.15)$$

- family can arise for a number of reasons, e.g., distribution having maximum entropy but that satisfies certain (moment) constraints.
- Given data $\mathbf{D} = \{\bar{x}_E^{(i)}\}_{i=1}^M$, form the expected statistics (requirements) of a model, with $\bar{x}^{(i)} \sim p(x)$

$$\hat{\mu}_{\alpha} = \frac{1}{M} \sum_{i=1}^{M} \phi_{\alpha}(\bar{x}^{(i)})$$
(11.16)

Thus, $\lim_{M\to\infty}\hat{\mu}_{\alpha}=E_p[\phi_{\alpha}(X)]=\mu_{\alpha}$

Exponential family models

• Goal ("estimation", or "machine learning") is to find

$$p^* \in \operatorname*{argmax}_{p \in \mathcal{U}} H(p) \text{ s.t. } \mathbb{E}_p[\phi_\alpha(X)] = \hat{\mu}_\alpha \ \forall \alpha \in \mathcal{I}$$
 (11.17)

where $\forall \alpha \in \mathcal{I}$

$$\mathbb{E}_p[\phi_\alpha(X)] = \int_{\mathsf{D}_X} \phi_\alpha(x) p(x) \nu(dx)$$
(11.18)

- $\mathbb{E}_p[\phi_\alpha(X)]$ is mean value as measured by potential function, so above is a form of moment matching.
- Maximum entropy (MaxEnt) distribution is solved by taking distribution in form of Eq. 11.15, by finding θ that solves

$$E_{p_{\theta}}[\phi_{\alpha}(X)] = \hat{\mu}_{\alpha} \text{ for all } \alpha \in \mathcal{I}$$
(11.19)

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

BP Next phase of class exponential models μ Param./Marg. Polytope Refs Minimal Representation of Exponential Family

• Solution as form:

$$p_{\theta}(x) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$$
(11.20)

where
$$A(\theta) = \log \int_{\mathsf{D}_X} \exp\left(\langle \theta, \phi(x) \rangle\right) \nu(dx)$$
 (11.21)

Exercise: show that solution to Eqn (11.17) has this form.

- Minimal representation Does not exist a nonzero vector γ ∈ ℝ^d for which ⟨γ, φ(x)⟩ is constant ∀x (that are ν-measurable).
- I.e., guarantee that, for all $\gamma \in \mathbb{R}^D$, there exists $x_1 \neq x_2$, with $\nu(x_1), \nu(x_2) > 0$, such that $\langle \gamma, \phi(x_1) \rangle \neq \langle \gamma, \phi(x_2) \rangle$.
- essential idea: that for a set of sufficient stats *I*, there is not a lower-dimensional vector |*I*'| < |*I*| that is also sufficient (a min suf stat is a function of all other suf stats).
- We can't reduce the dimensionality d without changing the family.

F27/60 (pg.27/60

• We'll see later, this useful in understanding BP algorithm.

Prof. Jeff Bilmes EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014 F29/60 (pg.29/60)

 LBP
 Next phase of class
 exponential models
 μ Param./Marg. Polytope
 Refs

 Exponential family models

 </t

• Minimal representation of Bernoulli distribution is

$$p(x|\gamma) = \exp(\gamma x - A(\gamma)) \tag{11.28}$$

• overcomplete rep of Bernoulli dist.

$$p(x|\theta_0, \theta_1) = \exp(\langle \theta, \phi(x) \rangle)$$
(11.29)

$$= \exp(\theta_0(1-x) + \theta_1 x - A(\gamma))$$
 (11.30)

where $\theta = (\theta_0, \theta_1)$ and $\phi(x) = (1 - x, x)$.

- Is there a vector a s.t. $\langle a, \phi(x) \rangle = c$ for all x, ν -a.e.?
- If a = (1, 1) then $\langle a, \phi(x) \rangle = (1 x) + x = 1$
- This is overcomplete since there is a linear combination of feature functions that are constant.
- Since $\theta_0(1-x) + \theta_1 x = \theta_0 + x(\theta_1 \theta_0)$, any parameters of form $\theta_1 \theta_0 = \gamma$ gives same distribution.

LBP

Famous Example - Ising Model

• Famous example is the Ising model in statistical physics. We have a grid network with pairwise interactions, each variable is 0/1-valued binary, and parameters associated with pairs being both on. Model becomes

$$p_{\theta}(x) = \exp\left\{\sum_{v \in V} \theta_v x_v + \sum_{(s,t) \in E} \theta_{st} x_s x_t - A(\theta)\right\}, \quad (11.31)$$

with

$$A(\theta) = \log \sum_{x \in \{0,1\}^m} \exp\left\{\sum_{v \in V} \theta_v x_v + \sum_{(s,t) \in E} \theta_{st} x_s x_t - A(\theta)\right\}$$
(11.32)

• Note that this is in minimal form. Any change to parameters will result in different distribution

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

F31/60 (pg.31/60)

• Note, in this case \mathcal{I} is all singletons (unaries) and all pairs, so that $\{C_{\alpha}\}_{\alpha} = \{\{x_i\}_i, \{x_i x_j\}_{(i,j) \in E}\}.$ • We can easily generalize this via a set system. I.e., consider (V, \mathcal{V}) , where $\mathcal{V} = \{V_1, V_2, \dots, V_{|\mathcal{V}|}\}$ and where $\forall i, V_i \subseteq V.$ • We can form sufficient statistic set via $\{C_{\alpha}\}_{\alpha} = \{\{x_V\}_{V \in \mathcal{V}}\}.$ • Higher order factors/interaction functions/potential functions/sufficient statistics.

LBP Next phase of class exponential models μ Param./Marg. Polytope Refs
Multivariate Gaussian

• Usually, multivariate Gausisan is parameterized via mean and covariance matrix. For canonical exponential form, we use mean and correlation matrix. I.e.

$$p_{\theta}(x) = \exp\left\{ \langle \theta, x \rangle + \frac{1}{2} \left\langle \! \left\langle \Theta, x x^{\mathsf{T}} \right\rangle \! \right\rangle - A(\theta, \Theta) \right\}$$
(11.36)

- $\langle\!\langle \Theta, xx^{\intercal} \rangle\!\rangle = \sum_{ij} \Theta_{ij} x_i x_j$ is Frobenius norm.
- So sufficient statistics are $(x_i)_{i=1}^n$ and $(x_i x_j)_{i,j}$
- $\Theta_{s,t} = 0$ means identical to missing edge in corresponding graph (marginal independence).
- Any other constraints on Θ ? negative definite
- Mixtures of Gaussians can also be parameterized in exponential form (but note, key is that it is the joint distribution $p_{\theta_s}(y_s, x_s)$).

LBP	Next phase of class	exponential models	μ Param./Marg. Polytope	Refs I
Mean Paramete	ers, Convex	< Cores		

• Consider quantities μ_{α} associated with statistic ϕ_{α} defined as:

$$\mu_{\alpha} = \mathbb{E}_p[\phi_{\alpha}(X)] = \int \phi_{\alpha}(x)p(x)\nu(dx)$$
 (11.37)

- this defines a vector of "mean parameters" $(\mu_1, \mu_2, \ldots, \mu_d)$ with $d = |\mathcal{I}|$.
- Define all the possible such vectors

$$\mathcal{M}(\phi) = \mathcal{M} \stackrel{\Delta}{=} \left\{ \mu \in \mathbb{R}^d : \exists p \text{ s.t. } \mu_\alpha = \mathbb{E}_p[\phi_\alpha(X)], \ \forall \alpha \in \mathcal{I} \right\}$$
(11.38)

- We don't say p was necessarily exponential family
- \mathcal{M} is convex since expected value is a linear operator. So convex combinations of p and p' will lead to convex combinations of μ and μ'
- \mathcal{M} is like a "convex core" of all distributions expressed via ϕ .

Prof. Jeff Bilme

LBP

class exponential models

Mean Parameters and Gaussians

- Here, we have $\mathbb{E}[XX^{\intercal}] = C$ and $\mu = \mathbb{E}X$. Question is, how to define \mathcal{M} ?
- Given definition of C and μ , then $C \mu \mu^{\mathsf{T}}$ must be valid covariance matrix (since this is $\mathbb{E}[X \mathbb{E}X][X \mathbb{E}X]^{\mathsf{T}} = C \mu \mu^{\mathsf{T}}$).
- Thus, $C \mu \mu^{\mathsf{T}} \succeq 0$, thus p.s.d. matrix.
- On the other hand, if this is true, we can form a Gaussian using $C \mu \mu^{\mathsf{T}}$ as the covariance matrix.
- $\bullet\,$ Thus, for Gaussian MRFs, ${\cal M}$ has the form

$$\mathcal{M} = \left\{ (\mu, C) \in \mathbb{R}^m \times \mathcal{S}^m_+ | C - \mu \mu^{\mathsf{T}} \succeq 0 \right\}$$
(11.39)

where \mathcal{S}^m_+ is the set of symmetric positive semi-definite matrices.

```
Prof. Jeff Bilme
```

EE512a/Fall 2014/Graphical Models - Lecture 11

F37/60 (pg.37/60

"Illustration of the set \mathcal{M} for a scalar Gaussian: the model has two mean parameters $\mu = \mathbb{E}[X]$ and $\Sigma_{11} = \mathbb{E}[X^2]$, which must satisfy the quadratic contraint $\Sigma_{11} - \mu^2 \ge 0$. Notice that \mathcal{M} is convex, which is a general property."

Also, don't confuse the "mean parameters" with the means of a Gaussian. The typical means of Gaussians are means in this new sense, but those means are not all of the means. ©

Prof. Jeff Bilmes

<page-header><section-header><section-header><section-header><section-header><section-header><section-header>

 LBP
 Next phase of class
 exponential models
 μ Param./Marg. Polytope
 Refs

 Mean Parameters and Polytopes
 Refs
 Refs

• Polytopes can be represented as a set of linear inequalities, i.e., there is a $|J| \times d$ matrix A and |J|-element column vector b with

$$M = \left\{ \mu \in \mathbb{R}^d : A\mu \ge b \right\} = \left\{ \mu \in \mathbb{R}^d : \langle a_j, \mu \rangle \ge b_j, \forall j \in J \right\}$$
(11.42)

with A having rows a_j .

"Ising model with two variables $(X_1, X_2) \in \{0, 1\}^2$. Three mean parameters $\mu_1 = \mathbb{E}[X_1]$, $\mu_2 = \mathbb{E}[X_2]$, $\mu_{12} = \mathbb{E}[X_2X_2]$, must satisfy constraints $0 \le \mu_{12} \le \mu_i$ for i = 1, 2, and $1 + \mu_{12} - \mu_1 - \mu_2 \ge 0$. These constraints carve out a polytope with four facets, contained within the unit hypercube $[0, 1]^3$."

LBP Next phase of class exponential models μ Param./Marg. Polytope Refs Mean Parameters and Marginal Polytopes

• Mean parameters are now true (fully specified) marginals, i.e., $\mu_v(j) = p(x_v=j)$ and $\mu_{st}(j,k) = p(x_s=j,x_t=k)$ since

$$\mu_{v,j} = \mathbb{E}_p[\mathbf{1}(x_v = j)] = p(x_v = j)$$
(11.49)

$$\mu_{st,jk} = \mathbb{E}_p[\mathbf{1}(x_s = j, x_t = k)] = p(x_s = j, x_t = k)$$
(11.50)

- Such an \mathcal{M} is called the *marginal polytope*. Any μ must live in the polytope that corresponds to node and edge true marginals!!
- We can also associate such a polytope with a graph G, where we take only (s,t) ∈ E(G). Denote this as M(G).
- This polytope can help us to characterize when BP converges (there might be an outer bound of this polytope), or it might characterize the result of a mean-field approximation (an inner bound of this polytope) as we'll see.

Learning is the dual of Inference

• Ex: Estimate θ with $\hat{\theta}$ based on data $\mathbf{D} = \{\bar{x}_E^{(i)}\}_{i=1}^M$ of size M, likelihood function

$$\ell(\theta, \mathbf{D}) = \frac{1}{M} \sum_{i=1}^{M} \log p_{\theta}(\bar{x}^{(i)}) = \langle \theta, \hat{\mu} \rangle - A(\theta)$$
(11.51)

where empirical means given by

$$\hat{\mu} = \hat{\mathbb{E}}[\phi(X)] = \frac{1}{M} \sum_{i=1}^{M} \phi(\bar{x}^{(i)})$$
(11.52)

• By taking derivatives of the above, it is easy to see that solution is the point $\hat{\theta}$ such that (empirical matches expected means)

$$\mathbb{E}_{\hat{\theta}}[\phi(X)] = \hat{\mu} \tag{11.53}$$

this is the the backward mapping problem, going from μ to θ .

• This is identical to the maximum entropy problem.

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th, 2014

LBP Next phase of class exponential models μ Param./Marg. Polytope Refs Learning is the dual of Inference

- I.e., solution to the maximum likelihood problem is one that satisfies the moment constraints and has the exponential model form.
- The exponential model form arises when we find the maximum entropy distribution over distributions satisfying the moment constraints.
- If we do maximum entropy learning, where does the exp(·) function come from? From the entropy function. I.e., the exponential form is the distribution that has maximum entropy having those constraints.

Prof leff Bilmes

F47/60 (pg.47/60)

 LBP
 Next phase of class
 exponential models
 μ Param./Marg. Polytope
 Refs

 Log partition (or cumulant) function
 <t

$$A(\theta) = \log \int_{\mathsf{D}_X} \langle \theta, \phi(x) \rangle \,\nu(dx) \tag{11.54}$$

- If we know the log partition function, we know a lot for an exponential family model. In particular, we know
- $A(\theta)$ is convex in θ (strictly so if minimal representation).
- It yields cumulants of the random vector $\phi(X)$

$$\frac{\partial A}{\partial \theta_{\alpha}}(\theta) = \mathbb{E}_{\theta}[\phi_{\alpha}(X)] = \int \phi_{\alpha}(X)p_{\theta}(x)\nu(dx) = \mu_{\alpha}$$
(11.55)

in general, derivative of log part. function is expected value of featureAlso, we get

$$\frac{\partial^2 A}{\partial \theta_{\alpha_1} \partial \theta_{\alpha_2}}(\theta) = \mathbb{E}_{\theta}[\phi_{\alpha_1}(X)\phi_{\alpha_2}(X)] - \mathbb{E}_{\theta}[\phi_{\alpha_1}(X)]\mathbb{E}_{\theta}[\phi_{\alpha_2}(X)]$$
(11.56)

Proof given in book.

Prof. Jeff Bilme

- For minimal exponential family models, this mapping is one-to-one, that is there is a unique pairing between μ and θ .
- For non-minimal exponential families, more than one θ for a given μ (not surprising since multiple θ 's can yield the same distribution).
- For non-exponential families, other distributions can yield μ, but the exponential family one is the one that has maximum entropy. ex1: Gaussian, a distribution with maximum entropy amongst all other distributions with same mean and covariance. ex2: Consider the maximum entropy optimization problem, yields a distribution with exactly this property.
- Key point: all mean parameters are realizable by member of exp. family.

LBP Next phase of class exponential models μ Param./Marg. Polytope Refs			
Mappings - onto			
Moreover,			
Theorem 11.6.2			
In a minimal exponential family, the gradient map $ abla A$ is onto the interior			
of $\mathcal M$ (denoted $\mathcal M^\circ$). Consequently, for each $\mu\in\mathcal M^\circ$, there exists some			
$\theta = \theta(\mu) \in \Omega$ such that $\mathbb{E}_{\theta}[\phi(X)] = \mu$.			

- Example: consider, for example, a Gaussian.
- Any mean parameter (set of means $\mathbb{E}[X]$ and correlations $\mathbb{E}[XX^T]$) can be realized by a Gaussian having those same mean parameters (moments).
- The Gaussian won't nec. be the "true" distribution (in such case, the "true" distribution would not be an exponential family model with those moments).
- The theorem here is more general and applies for any set of sufficient statistics.

phase of class exponential models

Conjugate Duality

• Consider maximum likelihood problem for exp. family

$$\theta^* \in \operatorname*{argmax}_{\theta} \left(\langle \theta, \hat{\mu} \rangle - A(\theta) \right)$$
(11.57)

• Convex conjugate dual of $A(\theta)$ is defined as:

$$A^{*}(\mu) \stackrel{\Delta}{=} \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right)$$
(11.58)

- So dual is optimal value of the ML problem, when $\mu \in \mathcal{M}$
- Key: when $\mu \in \mathcal{M}$, dual is negative entropy of exp. model $p_{\theta(\mu)}$ where $\theta(\mu)$ is the unique set of canonical parameters satisfying this matching condition

$$\mathbb{E}_{\theta(\mu)}[\phi(X)] = \nabla A(\theta(\mu)) = \mu \tag{11.59}$$

• When $\mu \notin \mathcal{M}$, then $A^*(\mu) = +\infty$, optimization with dual need consider points only in \mathcal{M} .

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 11 - Nov 5th,

F55/60 (pg.55/60)

LBP Next phase of class exponential models μ Param./Marg. Polytope Refs Conjugate Duality

Theorem 11.6.3 (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:

$$A^{*}(\mu) = \sup_{\theta \in \Omega} \left(\langle \theta, \mu \rangle - A(\theta) \right) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^{\circ} \\ +\infty & \text{if } \mu \in \bar{\mathcal{M}} \end{cases}$$
(11.60)

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

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$
(11.61)

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

$$\mu = \int_{\mathsf{D}_X} \phi(x) p_\theta(x) \nu(dx) = \mathbb{E}_\theta[\phi(X)] = \nabla A(\theta)$$
(11.62)

• Wainwright and Jordan Graphical Models, Exponential Families, and Variational Inference http://www.nowpublishers.com/product. aspx?product=MAL&doi=2200000001