

Probabilistic & Unsupervised Learning

Convex Algorithms in Approximate Inference

Maneesh Sahani

maneesh@gatsby.ucl.ac.uk

Gatsby Computational Neuroscience Unit, and
MSc ML/CSML, Dept Computer Science
University College London

Term 1, Autumn 2014

Convexity and Approximate Inference

The theory of convex functions and convex spaces has long been central to optimisation. It has recently also found application in the theory of free energy and approximation:

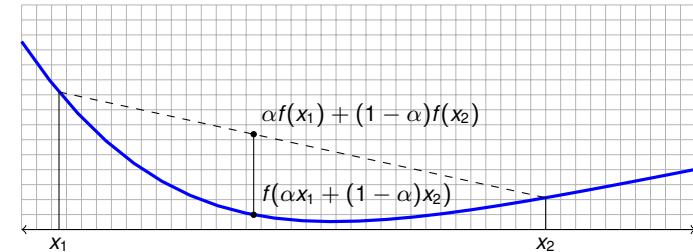
- ▶ Linear programming relaxation as an approximate method to find the MAP assignment in Markov random fields.
- ▶ Attractive Markov random fields: binary case exact and related to a maximum flow-minimum cut problem in graph theory (a linear program). Approximate otherwise.
- ▶ Unified view of approximate inference as optimization on the marginal polytope.
- ▶ Tree-structured convex upper bounds on the log partition function (convexified belief propagation).
- ▶ Learning graphical models using maximum margin principles and convex approximate inference.

Convexity

A convex function $f : X \rightarrow \mathbb{R}$ is one where

$$f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2)$$

for any $x_1, x_2 \in X$ and $0 \leq \alpha \leq 1$.



Convex functions have a global infimum (unless not bounded below) and there are efficient algorithms to find a minimum subject to convex constraints.

Examples: linear programs (LP), quadratic programs (QP), second-order cone programs (SOCP), semi-definite programs (SDP), geometric programs.

LP Relaxation for Markov Random Fields

Consider a discrete Markov random field (MRF) with pairwise interactions:

$$p(\mathbf{X}) = \frac{1}{Z} \prod_{(ij)} f_{ij}(X_i, X_j) \prod_i f_i(X_i) = \frac{1}{Z} \exp \left(\sum_{(ij)} E_{ij}(X_i, X_j) + \sum_i E_i(X_i) \right)$$

The problem is to find the most likely configuration \mathbf{X}^{MAP} :

$$\mathbf{X}^{\text{MAP}} = \underset{\mathbf{X}}{\operatorname{argmax}} \sum_{(ij)} E_{ij}(X_i, X_j) + \sum_i E_i(X_i)$$

Reformulate in terms of indicator variables:

$$b_i(x_i) = \delta(X_i = x_i)$$

$$b_{ij}(x_i, x_j) = \delta(X_i = x_i) \delta(X_j = x_j)$$

where $\delta(\cdot) = 1$ if argument is true, 0 otherwise. Each $b_i(x_i)$ is an indicator for whether variable X_i takes on value x_i . The indicator variables need to satisfy certain constraints:

$$b_i(x_i), b_{ij}(x_i, x_j) \in \{0, 1\} \quad \text{Indicator variables are binary variables.}$$

$$\sum_{x_i} b_i(x_i) = 1 \quad X_i \text{ takes on exactly one value.}$$

$$\sum_{x_j} b_{ij}(x_i, x_j) = b_i(x_i) \quad \text{Pairwise indicators are consistent with single-site indicators.}$$

LP Relaxation for Markov Random Fields

MAP assignment problem is equivalent to:

$$\operatorname{argmax}_{\{b_i, b_{ij}\}} \sum_{(ij)} \sum_{x_i, x_j} b_{ij}(x_i, x_j) E_{ij}(x_i, x_j) + \sum_i \sum_{x_i} b_i(x_i) E_i(x_i)$$

with constraints:

$$\forall i, j, x_i, x_j : \quad b_i(x_i), b_{ij}(x_i, x_j) \in \{0, 1\} \quad \sum_{x_i} b_i(x_i) = 1 \quad \sum_{x_j} b_{ij}(x_i, x_j) = b_i(x_i)$$

The linear programming relaxation for MRFs is:

$$\operatorname{argmax}_{\{b_i, b_{ij}\}} \sum_{(ij)} \sum_{x_i, x_j} b_{ij}(x_i, x_j) E_{ij}(x_i, x_j) + \sum_i \sum_{x_i} b_i(x_i) E_i(x_i)$$

with constraints:

$$\forall i, j, x_i, x_j : \quad b_i(x_i), b_{ij}(x_i, x_j) \in [0, 1] \quad \sum_{x_i} b_i(x_i) = 1 \quad \sum_{x_j} b_{ij}(x_i, x_j) = b_i(x_i)$$

Attractive Binary MRFs and Max Flow-Min Cut

Binary MRFs:

$$p(\mathbf{X}) = \frac{1}{Z} \exp \left(\sum_{(ij)} W_{ij} \delta(X_i = X_j) + \sum_i c_i X_i \right)$$

The binary MRF is attractive if $W_{ij} \geq 0$ for all i, j .

- ▶ Neighbouring variables 'prefer' to be in the same state.
- ▶ No loss of generality; any Boltzmann machines with positive interactions can be reparametrised to this form.
- ▶ Many practical MRFs are attractive, e.g. image segmentation, webpage classification.
- ▶ MAP \mathbf{X} can be found efficiently by converting problem into a maximum flow-minimum cut program.

LP Relaxation for Markov Random Fields

- ▶ The LP relaxation is a linear program which can be solved efficiently.
- ▶ If the solution is integral, i.e. each $b_i(x_i), b_{ij}(x_i, x_j) \in \{0, 1\}$, then the solution corresponds to **the MAP solution \mathbf{x}^{MAP}** .
- ▶ LP relaxation is a zero-temperature version of the Bethe free energy formulation of loopy BP, where the Bethe entropy term can be ignored.
- ▶ If the MRF is binary and attractive, then (a slightly different reformulation of LP relaxation) will **always give the MAP solution**.
- ▶ Next: we show how to find the MAP solution directly for binary attractive MRFs using network flow.

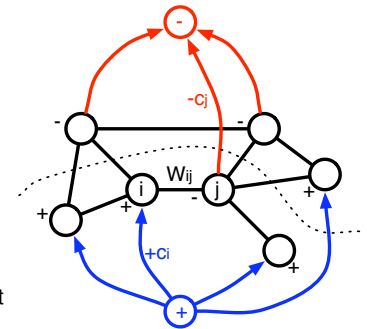
Attractive Binary MRFs and Max Flow-Min Cut

The MAP problem:

$$\operatorname{argmax}_{\mathbf{x}} \sum_{(ij)} W_{ij} \delta(x_i = x_j) + \sum_i c_i x_i$$

Construct a network as follows:

1. Edges (ij) are undirected with weight $\lambda_{ij} = W_{ij}$;
2. Add a **source** s and a **sink** t node;
3. $c_i > 0$: Connect the **source** node to variable i with weight $\lambda_{si} = c_i$;
4. $c_j < 0$: Connect variable j to the **sink** node with weight $\lambda_{jt} = -c_j$.



A **cut** is a partition of the nodes into S and T with $s \in S$ and $t \in T$. The weight of the cut is

$$\Lambda(S, T) = \sum_{i \in S, j \in T} \lambda_{ij}$$

The **minimum cut** problem is to find the cut with minimum weight.

Attractive Binary MRFs and Max Flow-Min Cut

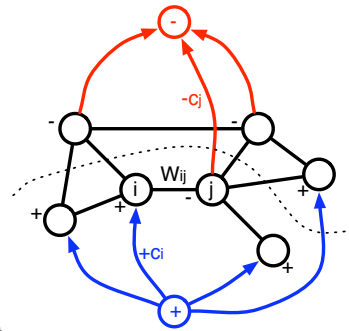
Identify an assignment $\mathbf{X} = \mathbf{x}$ with a cut:

$$S = \{s\} \cup \{i : x_i = 1\}$$

$$T = \{t\} \cup \{j : x_j = 0\}$$

The weight of the cut is:

$$\begin{aligned} \Lambda(S, T) &= \sum_{(ij)} W_{ij} \delta(x_i \neq x_j) \\ &\quad + \sum_i (1 - x_i) \max(0, c_i) \\ &\quad + \sum_j x_j \max(0, -c_j) \\ &= - \sum_{(ij)} W_{ij} \delta(x_i = x_j) - \sum_i x_i c_i + \text{constant} \end{aligned}$$



So finding the minimum cut corresponds to finding the MAP assignment.

How do we find the minimum cut? The minimum cut problem is dual to the **maximum flow problem**, i.e. find the maximum flow allowable from the source to the sink through the network. This can be solved extremely efficiently (see wikipedia entry).

The framework can be generalized to general attractive MRFs, but will not be exact anymore.

Exponential families: mean parameters and negative entropy

An exponential family distribution can also (almost always) be parameterised by the **means of the sufficient statistics**.

$$\mu(\theta) = \mathbb{E}_{\theta} [s(X)]$$

Consider the **negative entropy** of the distribution as a function of the mean parameter:

$$\Psi(\mu) = \mathbb{E}_{\theta} [\log p(X|\theta(\mu))] = \theta^T \mu - \Phi(\theta)$$

so

$$\theta^T \mu = \Phi(\theta) + \Psi(\mu)$$

The negative entropy is **dual** to the log-partition function. For example,

$$\begin{aligned} \frac{d}{d\mu} \Psi(\mu) &= \frac{\partial}{\partial \mu} (\theta^T \mu - \Phi(\theta)) + \frac{d\theta}{d\mu} \frac{\partial}{\partial \theta} (\theta^T \mu - \Phi(\theta)) \\ &= \theta + \frac{d\theta}{d\mu} (\mu - \mu) = \theta \end{aligned}$$

Exponential families: the log partition function

Consider an exponential family distribution with sufficient statistic $s(X)$ and natural parameter θ (and no base factor in X alone). We can write its probability or density function as

$$p(X|\theta) = \exp(\theta^T s(X) - \Phi(\theta))$$

where $\Phi(\theta)$ is the **log partition function**

$$\Phi(\theta) = \log \sum_x \exp(\theta^T s(x))$$

$\Phi(\theta)$ plays an important role in the theory of the exponential family. For example, it maps natural parameters to the moments of the sufficient statistics:

$$\frac{\partial}{\partial \theta} \Phi(\theta) = e^{-\Phi(\theta)} \sum_x s(x) e^{\theta^T s(x)} = \mathbb{E}_{\theta} [s(X)] = \mu(\theta) = \mu$$

$$\frac{\partial^2}{\partial \theta^2} \Phi(\theta) = e^{-\Phi(\theta)} \sum_x s(x)^2 e^{\theta^T s(x)} - e^{-2\Phi(\theta)} \left[\sum_x s(x) e^{\theta^T s(x)} \right]^2 = \mathbb{V}_{\theta} [s(X)]$$

The second derivative is thus positive semi-definite, and so $\Phi(\theta)$ is **convex in θ** .

Exponential families: duality

In fact, the log partition function and negative entropy are **conjugate dual** functions.

Consider the KL divergence between distributions with natural parameters θ and θ' :

$$\begin{aligned} \mathbf{KL}[\theta || \theta'] &= \mathbf{KL}[p(X|\theta) || p(X|\theta')] = \mathbb{E}_{\theta} [-\log p(X|\theta') + \log p(X|\theta)] \\ &= -\theta'^T \mu + \Phi(\theta') + \Psi(\mu) \geq 0 \\ \Rightarrow \Psi(\mu) &\geq \theta'^T \mu - \Phi(\theta') \end{aligned}$$

where μ are the mean parameters corresponding to θ .

Now, the minimum KL divergence of zero is reached iff $\theta = \theta'$, so

$$\Psi(\mu) = \sup_{\theta'} [\theta'^T \mu - \Phi(\theta')] \quad \text{and, if finite} \quad \theta(\mu) = \operatorname{argmax}_{\theta'} [\theta'^T \mu - \Phi(\theta')]$$

The left-hand equation is the definition of the conjugate dual of a convex function.

Continuous functions are reciprocally dual, so we also have:

$$\Phi(\theta) = \sup_{\mu'} [\theta^T \mu' - \Psi(\mu')] \quad \text{and, if finite} \quad \mu(\theta) = \operatorname{argmax}_{\mu'} [\theta^T \mu' - \Psi(\mu')]$$

Thus, duality gives us another relation between θ and μ .

Duality, inference and the free energy

Consider a joint exponential family distribution on observed \mathbf{x} and latent \mathbf{y} .

$$p(\mathbf{x}, \mathbf{y}) = \exp \left[\boldsymbol{\theta}^T s(\mathbf{x}, \mathbf{y}) - \Phi_{XY}(\boldsymbol{\theta}) \right]$$

The posterior on \mathbf{y} is also in the exponential family, with the **clamped** sufficient statistic $s_Y(\mathbf{y}; \mathbf{x}) = s_{XY}(\mathbf{x}^{\text{obs}}, \mathbf{y})$; the **same** (now possibly redundant) natural parameter $\boldsymbol{\theta}$; and partition function $\Phi_Y(\boldsymbol{\theta}) = \log \sum_{\mathbf{y}} \exp \boldsymbol{\theta}^T s_Y(\mathbf{y})$.

The likelihood is

$$\mathcal{L}(\boldsymbol{\theta}) = p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{\mathbf{y}} e^{\boldsymbol{\theta}^T s(\mathbf{x}, \mathbf{y}) - \Phi_{XY}(\boldsymbol{\theta})} = \sum_{\mathbf{y}} e^{\boldsymbol{\theta}^T s_Y(\mathbf{y}; \mathbf{x})} e^{-\Phi_{XY}(\boldsymbol{\theta})} = \exp[\Phi_Y(\boldsymbol{\theta}) - \Phi_{XY}(\boldsymbol{\theta})]$$

So we can write the log-likelihood as

$$\ell(\boldsymbol{\theta}) = \sup_{\boldsymbol{\mu}_Y} \left[\underbrace{\boldsymbol{\theta}^T \boldsymbol{\mu}_Y - \Phi_{XY}(\boldsymbol{\theta})}_{\langle \log p(\mathbf{x}, \mathbf{y}) \rangle_q} - \underbrace{\Psi(\boldsymbol{\mu}_Y)}_{\mathbf{H}[q]} \right] = \sup_{\boldsymbol{\mu}_Y} \mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\mu}_Y)$$

This is the familiar free energy with $q(\mathbf{y})$ represented by its mean parameters $\boldsymbol{\mu}_Y$!

Convexity and undirected trees

- We can parametrise a discrete pairwise MRF as follows:

$$\begin{aligned} p(\mathbf{X}) &= \frac{1}{Z} \prod_i f_i(X_i) \prod_{(ij)} f_{ij}(X_i, X_j) \\ &= \exp \left(\sum_i \sum_{x_i} \boldsymbol{\theta}_i(x_i) \delta(X_i = x_i) + \sum_{(ij)} \sum_{x_i, x_j} \boldsymbol{\theta}_{ij}(x_i, x_j) \delta(X_i = x_i) \delta(X_j = x_j) - \Phi(\boldsymbol{\theta}) \right) \end{aligned}$$

- So discrete MRFs are always exponential family, with natural and mean parameters:

$$\begin{aligned} \boldsymbol{\theta} &= [\boldsymbol{\theta}_i(x_i), \boldsymbol{\theta}_{ij}(x_i, x_j) \quad \forall i, j, x_i, x_j] \\ \boldsymbol{\mu} &= [p(X_i = x_i), p(X_i = x_i, X_j = x_j) \quad \forall i, j, x_i, x_j] \end{aligned}$$

In particular, the mean parameters are just the singleton and pairwise probability tables.

- If the MRF has tree structure \mathcal{T} , the negative entropy can be written in terms of the single-site entropies and mutual informations on edges:

$$\begin{aligned} \Psi(\boldsymbol{\mu}_{\mathcal{T}}) &= \mathbb{E}_{\boldsymbol{\theta}_{\mathcal{T}}} \left[\log \prod_i p(X_i) \prod_{(ij) \in \mathcal{T}} \frac{p(X_i, X_j)}{p(X_i)p(X_j)} \right] \\ &= - \sum_i H(X_i) + \sum_{(ij) \in \mathcal{T}} I(X_i, X_j) \end{aligned}$$

Inference with mean parameters

We have described inference in terms of the distribution q , approximating as needed. Can we describe it instead as an optimisation over $\boldsymbol{\mu}$ directly?

$$\boldsymbol{\mu}_Y^* = \underset{\boldsymbol{\mu}_Y}{\operatorname{argmax}} [\boldsymbol{\theta}^T \boldsymbol{\mu}_Y - \Psi(\boldsymbol{\mu}_Y)]$$

Concave maximisation(!), but two complications:

- The optimum must be found over **feasible** means. Interdependence of the sufficient statistics may prevent arbitrary sets of mean sufficient statistics being achieved
- Feasible means are convex combinations of all the single-configuration sufficient statistics.

$$\boldsymbol{\mu} = \sum_{\mathbf{x}} \nu(\mathbf{x}) s(\mathbf{x}) \quad \sum_{\mathbf{x}} \nu(\mathbf{x}) = 1$$

- Take a Boltzmann machine on two variables, x_1, x_2 .
- The sufficient stats are $s(\mathbf{x}) = [x_1, x_2, x_1 x_2]$.
- Clearly only the stats $S = \{[0, 0, 0], [0, 1, 0], [1, 0, 0], [1, 1, 1]\}$ are possible.
- Thus $\boldsymbol{\mu} \in \text{convex hull}(S)$.

- For a discrete distribution, this space of possible means is bounded by exponentially many hyperplanes connecting the discrete configuration stats: called the **marginal polytope**.
- Even when restricted to the marginal polytope, evaluating $\Psi(\boldsymbol{\mu})$ can be challenging.

The Bethe free energy again

We can see the Bethe free energy problem as a relaxation of the true free-energy optimisation:

$$\boldsymbol{\mu}_Y^* = \underset{\boldsymbol{\mu}_Y \in \mathcal{M}}{\operatorname{argmax}} [\boldsymbol{\theta}^T \boldsymbol{\mu}_Y - \Psi(\boldsymbol{\mu}_Y)]$$

where \mathcal{M} is the set of feasible means.

1. **Relax** $\mathcal{M} \rightarrow \mathcal{L}$, where \mathcal{L} is the set of **locally consistent** means (i.e. all nested means marginalise correctly).
2. **Approximate** $\Psi(\boldsymbol{\mu}_Y)$ by the tree-structured form

$$\Psi_{\text{Bethe}}(\boldsymbol{\mu}_Y) = - \sum_i H(X_i) + \sum_{(ij) \in \mathcal{T}} I(X_i, X_j)$$

\mathcal{L} is still a convex set (polytope for discrete problems). However Ψ_{Bethe} is not convex.

Convexifying BP

Consider instead an **upper bound** on $\Phi(\theta)$:

Imagine a set of spanning trees \mathcal{T} for the MRF, each with its own parameters $\theta_{\mathcal{T}}, \mu_{\mathcal{T}}$. By padding entries corresponding to off-tree edges with zero, we can assume that $\theta_{\mathcal{T}}$ has the same dimensionality as θ .

Suppose also that we have a distribution β over the spanning trees so that $\mathbb{E}_{\beta} [\theta_{\mathcal{T}}] = \theta$.

Then by the convexity of $\Phi(\theta)$,

$$\Phi(\theta) = \Phi(\mathbb{E}_{\beta} [\theta_{\mathcal{T}}]) \leq \mathbb{E}_{\beta} [\Phi(\theta_{\mathcal{T}})]$$

If we were to **tighten** the upper bound we might obtain a good approximation to Φ :

$$\Phi(\theta) \leq \inf_{\beta, \theta_{\mathcal{T}}: \mathbb{E}_{\beta} [\theta_{\mathcal{T}}] = \theta} \mathbb{E}_{\beta} [\Phi(\theta_{\mathcal{T}})]$$

Convex Upper Bounds on the Log Partition Function

$$\begin{aligned} \Phi(\theta) &\leq \sup_{\lambda} \inf_{\theta_{\mathcal{T}}} \mathbb{E}_{\beta} [\Phi(\theta_{\mathcal{T}})] - \lambda^{\top} (\mathbb{E}_{\beta} [\theta_{\mathcal{T}}] - \theta) \\ &= \sup_{\lambda} \lambda^{\top} \theta + \mathbb{E}_{\beta} \left[\inf_{\theta_{\mathcal{T}}} \Phi(\theta_{\mathcal{T}}) - \theta_{\mathcal{T}}^{\top} \Pi_{\mathcal{T}}(\lambda) \right] \\ &= \sup_{\lambda} \lambda^{\top} \theta + \mathbb{E}_{\beta} [-\Psi(\Pi_{\mathcal{T}}(\lambda))] \\ &= \sup_{\lambda} \lambda^{\top} \theta + \mathbb{E}_{\beta} \left[\sum_i H_{\lambda}(X_i) - \sum_{(ij) \in \mathcal{T}} l_{\lambda}(X_i, X_j) \right] \\ &= \sup_{\lambda} \lambda^{\top} \theta + \sum_i H_{\lambda}(X_i) - \sum_{(ij)} \beta_{ij} l_{\lambda}(X_i, X_j) \end{aligned}$$

This is a **convexified** Bethe free energy.

Convex Upper Bounds on the Log Partition Function

$$\Phi(\theta) \leq \inf_{\theta_{\mathcal{T}}: \mathbb{E}_{\beta} [\theta_{\mathcal{T}}] = \theta} \mathbb{E}_{\beta} [\Phi(\theta_{\mathcal{T}})]$$

Solve this constrained optimisation problem using Lagrange multipliers:

$$\mathcal{L} = \mathbb{E}_{\beta} [\Phi(\theta_{\mathcal{T}})] - \lambda^{\top} (\mathbb{E}_{\beta} [\theta_{\mathcal{T}}] - \theta)$$

Setting the derivatives wrt $\theta_{\mathcal{T}}$ to zero, we get:

$$\begin{aligned} \beta(\mathcal{T}) \lambda_{\mathcal{T}} - \beta(\mathcal{T}) \Pi_{\mathcal{T}}(\lambda) &= 0 \\ \lambda_{\mathcal{T}} &= \Pi_{\mathcal{T}}(\lambda) \end{aligned}$$

where $\Pi_{\mathcal{T}}(\lambda)$ are the Lagrange multipliers corresponding to vertices and edges on the tree \mathcal{T} .

Although there can be many $\theta_{\mathcal{T}}$ parameters, at optimum they are all constrained: their corresponding mean parameters are all consistent with each other and with λ .

References

- **Graphical Models, Exponential Families, and Variational Inference.** Wainwright and Jordan. **Foundations and Trends in Machine Learning**, 2008 1:1-305.
- Exact Maximum A Posteriori Estimation for Binary Images. Greig, Porteous and Seheult, Journal of the Royal Statistical Society B, 51(2):271-279, 1989.
- Fast Approximate Energy Minimization via Graph Cuts. Boykov, Veksler and Zabih, International Conference on Computer Vision 1999.
- MAP estimation via agreement on (hyper)trees: Message-passing and linear-programming approaches. Wainwright, Jaakkola and Willsky, IEEE Transactions on Information Theory, 2005, 51(11):3697-3717.
- Learning Associative Markov Networks. Taskar, Chatalbashev and Koller, International Conference on Machine Learning, 2004.
- A New Class of Upper Bounds on the Log Partition Function. Wainwright, Jaakkola and Willsky. IEEE Transactions on Information Theory, 2005, 51(7):2313-2335.
- MAP Estimation, Linear Programming and Belief Propagation with Convex Free Energies. Weiss, Yanover and Meltzer, Uncertainty in Artificial Intelligence, 2007.