

Inference as Optimization

Sargur Srihari

srihari@cedar.buffalo.edu

Topics in Inference as Optimization

- Overview
- Exact Inference revisited
- The Energy Functional
- Optimizing the Energy Functional

Exact and Approximate Inference

- PGMs represent probability distributions $P_{\Phi}(\chi)$
 - Where χ is a set of variables and Φ is a set of factors
- Inference is the task of answering queries
 - e.g., compute conditional probability $P_{\Phi}(Y|E=e)$, $Y, E \subseteq \chi$
- Problem of inference in PGMs is \mathcal{NP} -hard
 - Worst case is exponential
- *Exact Inference* is often efficient using
 - Variable Elimination or Clique tree Algorithms
 - But complexity is exponential in tree width of network
 - In such cases exact algorithms become infeasible
- This motivates *approximate inference*

Approximate Target Distribution

- We consider approximate inference methods where the approximation arises from constructing an approximation to target distribution P_{Φ}
- This approximation takes a simpler form that allows inference
- Simpler approximating form exploits factorization structure of PGM

Principles of Approximate Algorithms

- Approximate inference methods share common conceptual principles:
 1. Find target class \mathcal{Q} of “easy” distributions Q and
 2. Then Search for an instance within that class that best approximates P_{Φ}
 3. Answer queries using inference on Q instead of P_{Φ}
 4. All methods *optimize* the same target function for measuring similarity between Q and P_{Φ}
- This reformulates inference problem as:
 - Optimizing an objective function over class \mathcal{Q}

Reformulated Inference Problem

- This problem is one of constrained optimization
 - i.e., find distribution Q that *minimizes* $D(Q|| P_{\Phi})$
 - Such problems can be solved by variety of different optimization techniques
- Technique most often used for PGMs is based on Lagrange multipliers
- Constrained optimization and Lagrange solution is discussed next

What is constrained optimization?

- Ex: find the maximum entropy distribution over X with $Val(X) = \{x^1, \dots, x^K\}$ where entropy is

$$H(X) = -\sum_{k=1}^K p(x^k) \log P(x^k)$$

– Unconstrained Optimization

- Use gradient method treating each $P(x^k)$ as a parameter θ_k
 - Compute gradient of $H_P(X)$ wrt parameters: $\frac{\partial}{\partial \theta_k} H(X) = -\log(\theta_k) - 1$
 - Setting partial derivative to 0 we get $\log(\theta_k) = 1$, or $\theta_k = 1/2$
 - But nos do not add up to 1, and hence not a distribution
 - Flaw in analysis: we want constraints that $\sum_k \theta_k = 1$, and $\theta_k \geq 0$

– Constrained Optimization

- Maximizing a function f under equality constraints
- Find θ
- Maximizing $f(\theta)$
- Subject to $c_1(\theta) = 0$
.....
 $c_m(\theta) = 0$

Method of Lagrange multipliers allows us to solve constrained optimization problems using tools for unconstrained optimization. Lagrangian is

$$J(\theta, \lambda) = f(\theta) - \sum_{j=1}^m \lambda_j c_j(\theta)$$

Lagrange leads to Message Passing

- Method of Lagrange multipliers produces a set of equations that characterize the optima of the objective
- It produces a set of fixed-point equations that define each variable in terms of others
- Fixed point equations derived from constrained energy optimization can be viewed as passing messages over a graph object

Categories of methods in this class

1. Message passing on Clique Tree

- Loopy belief propagation
 - Optimize approximate versions of the energy functional

2. Message passing on Clique Trees with approximate messages

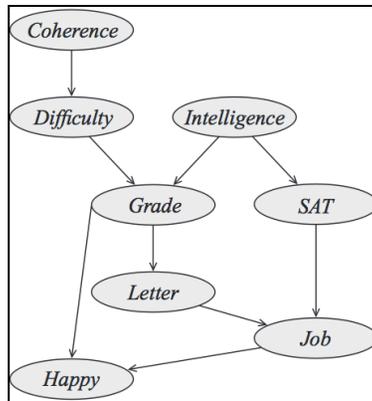
- Called expectation propagation
 - Maximize exact energy functional but with relaxed constraints on Q

3. Mean-field method

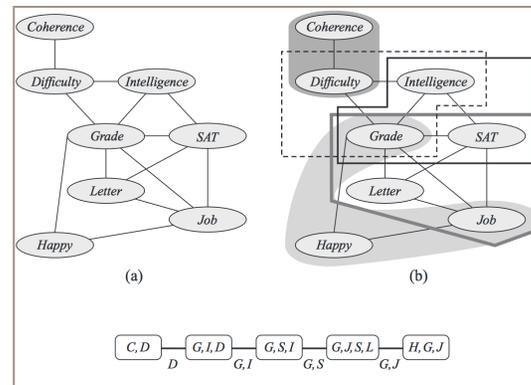
- Originates in statistical physics
 - Focus on Q that has simple factorization

Examples of Clique Tree

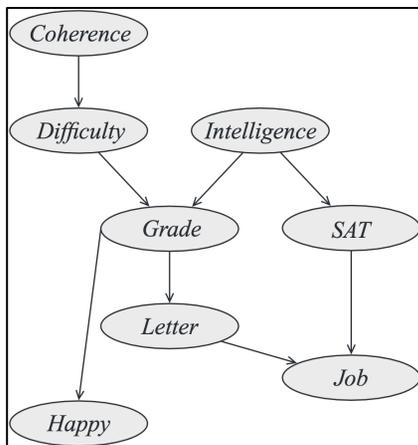
Bayesian Network 1



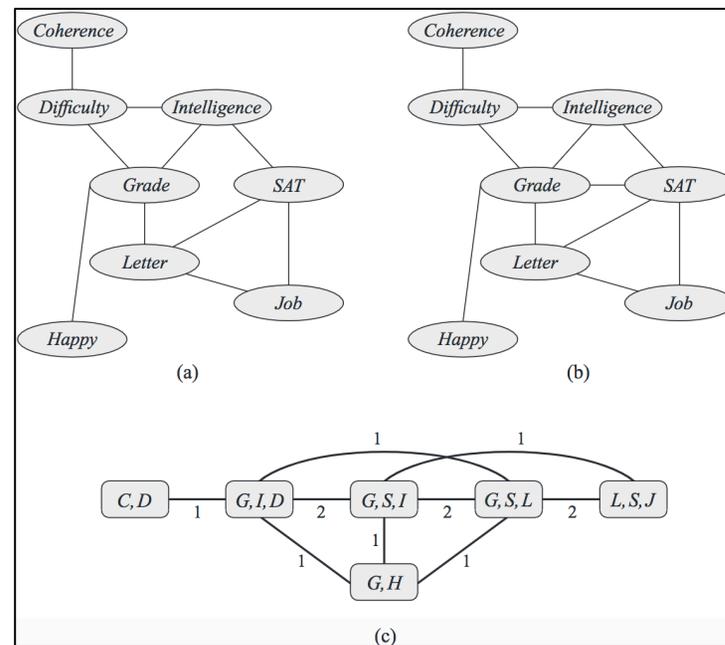
Moralized Graph and Clique tree



Bayesian Network 2



Moralized Graph



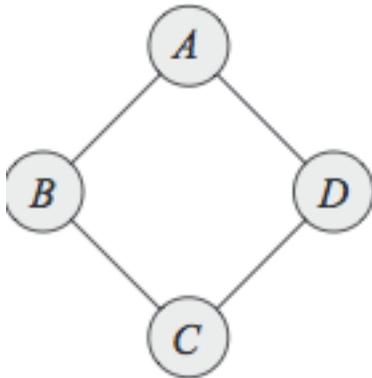
Triangu-
lation

Cluster
Graph

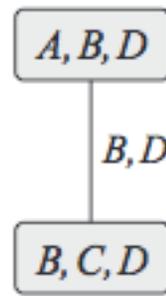
Clique tree

Belief Propagation

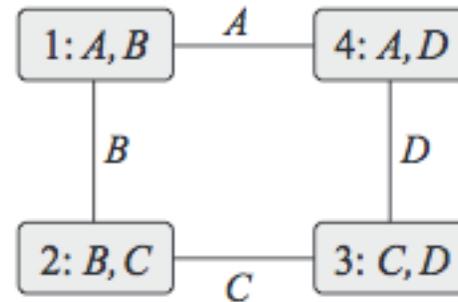
A simple network



A Clique Tree



A Cluster Graph



- Clique tree and cluster graph are alternative ways of doing inference
- Cluster graph may contain loops
- Inference is Called *Loopy Belief Propagation*
- Clusters are smaller than in Clique Tree

Exact Inference Revisited

- We have a factorized distribution of the form

$$P_{\Phi}(X) = \frac{1}{Z} \prod_{\phi \in \Phi} \phi(U_{\phi})$$

- where $U_{\phi} = \text{Scope}(\phi)$
- Factors are:
 - CPDs in a BN or
 - potentials in a MN
- We are interested in answering queries:
 - about marginal probabilities of variables and
 - about the partition function

Cluster Graph Representation

- End-product of Belief Propagation is a calibrated cluster tree
 - A calibrated set of beliefs represents a distribution
- We view exact inference as searching over the set of distributions Q that are representable by the cluster tree to find a distribution Q^* that matches P_{Φ}

Cluster graph U for factors Φ over χ is an undirected graph

Each of whose nodes i is associated with a subset

Each edge between pair of clusters C_i and C_j is associated with a sepset $S_{i,j} \subseteq C_i \cap C_j$

A tree T is a clique tree for graph \mathcal{H} if $C_i \subseteq \chi$

Each node in T corresponds to a clique in \mathcal{H} and each maximal clique in \mathcal{H} is a node in T

Each sepset $S_{i,j}$ separates $W_{\langle(j,j)} and $W_{\langle(j,i)}$ in $\mathcal{H}$$

Distance between Q and P_{Φ}

- We need to optimize distance between Q and P_{Φ} without answering hard queries about P_{Φ}
- *Relative entropy* (or K-L divergence) allows us to exploit the structure of P_{Φ} without performing reasoning with it
 - *Relative entropy* of P_1 and P_2 defined as

$$D(P_1 \parallel P_2) = E_{P_1} \left[\frac{\ln P_1(\chi)}{\ln P_2(\chi)} \right]$$

- It is always non-negative
- Equal to 0 if and only if $P_1 = P_2$
- We search for distribution Q that *minimizes* $D(Q \parallel P_{\Phi})$

Specifying the set Q

- We need to specify objects to optimize over

- Suppose we are given:

- a clique tree structure T for P_{Φ} ,
- a set of beliefs $Q = \{\beta_i: i \in V_T\} \cup \{\mu_{i,j}: (i,j) \in E_T\}$

where C_i are clusters in T ,

β_i denote beliefs over C_i and

$\mu_{i,j}$ denotes beliefs $S_{i,j}$ of edges in T

- Set of beliefs in T defines a distribution Q by

$$Q(x) = \frac{\prod_{i \in V_T} \beta_i}{\prod_{(i,j) \in E_T} \mu_{i,j}}$$

The beliefs correspond to marginals of Q

$$\beta_i[c_i] = Q(c_i)$$

$$\mu_{ij}[s_{ij}] = Q(s_{ij})$$

- We are now searching over a set of distributions Q
 - that are representable by a set of beliefs Q over the cliques and sepsets in a particular clique tree structure T

Statement of Inference as Optimization

- Exact inference is one of maximizing $-D(Q \parallel P_\Phi)$ over the space of calibrated sets Q

Ctree-Optimize-KL

- Find** $Q = \{\beta_i : i \in V_T\} \cup \{\mu_{i,j} : (i,j) \in E_T\}$

- Maximizing** $-D(Q \parallel P_\Phi)$

- Subject to**

$$\begin{aligned} \mu_{i,j}[s_{i,j}] &= \sum_{c_i \sim s_{i,j}} \beta_i(c_i) \quad \forall (i,j) \in E_T, \forall s_{i,j} \in \text{Val}(S_{i,j}) \\ \sum_{c_i} \beta_i(c_i) &= 1 \quad \forall i \in V_T \end{aligned}$$

- Theorem:* If T is an I-map of P_Φ then there is a unique solution to Ctree-Optimize-KL

Possible approach

- Examine different configurations of beliefs that satisfy marginal consistency constraints
 - Select the configuration that maximizes the objective
 - Such as exhaustive examination is impossible to perform
- Instead of searching over a space of all calibrated trees we can search over a space of simpler distributions
 - We will not find a distribution equivalent to P_{Φ} but one that is reasonably close

The Energy Functional

- Directly evaluating $D(Q \parallel P_{\Phi})$ is unwieldy

$$D(P_1 \parallel P_2) = E_{P_1} \left[\frac{\ln P_1(\chi)}{\ln P_2(\chi)} \right] = \sum_{\chi} P_1[\chi] \left[\frac{\ln P_1(\chi)}{\ln P_2(\chi)} \right]$$

- Because summation over all χ is infeasible in practice

- Instead use equivalent form $D(Q \parallel P_{\Phi}) = \ln Z - F(\tilde{P}_{\Phi}, Q)$

- Where F is the energy functional

- *Theorem:* $F[\tilde{P}_{\Phi}, Q] = E_Q[\ln \tilde{P}(\chi)] + H_Q(\chi) = \sum_{\phi \in \Phi} E_Q[\ln \phi] + H_Q(\chi)$

- Since the term $\ln Z$ does not depend on Q ,

- minimizing relative entropy $D(Q \parallel P_{\Phi})$ is equivalent to maximizing the energy functional $F(\tilde{P}_{\Phi}, Q)$

- Energy functional $F[\tilde{P}_{\Phi}, Q] = \sum_{\phi \in \Phi} E_Q[\ln \phi] + H_Q(\chi)$ has two terms:

- energy term (expectation of logs of factors in Φ) and entropy term

Optimizing the Energy Functional

- From here onward we pose the problem of finding a good Q as one of maximizing the energy functional
 - Equivalently minimizing the relative entropy
 - Importantly energy functional involves expectations in Q
 - By choosing Q that allow efficient inference we can evaluate/optimize the energy functional
- Moreover, energy Functional is a lower bound on partition function
 - Since $D(Q||P_{\Phi}) \geq 0$ we have $\ln Z \geq F[\tilde{P}_{\Phi}, Q]$
 - Useful since partition function is usually the hardest part of inference
 - Plays important role in learning

Strategies for optimizing energy functional

- Methods are referred to as *Variational Methods*
- Refers to a strategy in which we introduce new parameters that increase the degrees of freedom
- Each choice of these parameters gives a different approximation
- We attempt to optimize the variational parameters to get the best approximation
- *Variational calculus*: finding optima of a functional
 - E.g., distribution that maximizes entropy

Further Topics in Variational Methods

- Exact Inference
- Propagation-Based Approximations
- Propagation with Approximate Messages
- Structured Variational Approximations