



Exact Inference as Optimization

Assume we have a factored distribution

$$P(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(U_i)$$

In exact inference, we find a set of calibrated beliefs that represent $P(x)$

That is, we find beliefs that match the distribution represented by given set of initial potentials

We can view exact inference as searching over the set of distributions \mathcal{Q} that are representable by the cluster tree to find a distribution Q^* that matches P



Searching for a calibrated distribution that is as close as possible to P

the KL-divergence (relative-divergence)

$$D(P \parallel Q) = \int_x P(x) \log \frac{P(x)}{Q(x)} dx$$

Search for a distribution Q that minimizes $D(Q \parallel P)$

Suppose we are given a clique tree structure \mathcal{T} for P
(i.e. \mathcal{T} satisfies the running intersection property and family preserving property)
and a set of beliefs

$$\mathcal{Q} = \{ \beta_i \mid i \in V \} \cup \{ M_{ij} \mid (i,j) \in E_{\mathcal{T}} \}$$

the set of beliefs in \mathcal{T} defines a distribution Q

$$Q(x) = \frac{\prod_{i \in V} \beta_i}{\prod_{(i,j) \in E_{\mathcal{T}}} M_{ij}}$$

if \mathcal{Q} is a set of calibrated beliefs for \mathcal{T} , then
$$\begin{cases} \beta_i(c_i) = Q(c_i) \\ M_{ij}(S_{ij}) = Q(S_{ij}) \end{cases}$$

$$\text{find } \mathcal{Q} = \{ \beta_i \mid i \in V \} \cup \{ M_{ij} \mid (i,j) \in E_{\mathcal{T}} \}$$

$$\text{minimize } D(Q \parallel P)$$

$$\text{s.t. } \begin{cases} M_{ij}(S_{ij}) = \sum_{c_i \in S_{ij}} \beta_i(c_i) & \forall (i,j) \in E_{\mathcal{T}} \quad \forall S_{ij} \in \mathcal{V}(S_{ij}) \\ \sum_{c_i} \beta_i(c_i) = 1 & \forall i \in V \end{cases}$$

The Energy Functional

The objective $D(\alpha \| P_0)$ is unwieldy for direct optimization, but we can rewrite the relative entropy in a simpler form

$$D(\alpha \| P_0) = \int \alpha(x) \cdot \ln \frac{\alpha(x)}{P_0(x)} dx = \ln Z - F[P_0, \alpha]$$

where $F[P_0, \alpha]$ is the energy functional

$$F[P_0, \alpha] = E_{\alpha}[\ln \psi(x)] + H_{\alpha}(x) = \sum_{\mu} E_{\alpha}[\ln \phi_{\mu}] + H_{\alpha}(x)$$

and $H_{\alpha}(x) = -E_{\alpha}[\ln \alpha(x)] = -\int \alpha(x) \cdot \ln \alpha(x) dx$ is the entropy of α

$$\begin{aligned} D(\alpha \| P_0) &= \int \alpha(x) \cdot \ln \frac{\alpha(x)}{P_0(x)} dx \\ &= \int \alpha(x) \cdot \ln \alpha(x) dx - \int \alpha(x) \cdot \ln P_0(x) dx \\ &= E_{\alpha}[\ln \alpha(x)] - E_{\alpha}[\ln P_0(x)] \end{aligned}$$

$$\begin{aligned} P_0(x) &= \frac{1}{Z} \tilde{P}_0(x) = \frac{1}{Z} \int \phi_{\mu}(U_{\mu}) \quad U_{\mu} \text{ is the projection of } \mathcal{U} \text{ on } \text{State}[\mu] \\ &= E_{\alpha}[\ln \alpha(x)] - E_{\alpha}[\sum_{\mu} \ln \phi_{\mu}(U_{\mu}) - \ln Z] \\ &= -H_{\alpha}(x) - \sum_{\mu} E_{\alpha}[\ln \phi_{\mu}(U_{\mu})] + \ln Z \\ &= -F[P_0, \alpha] + \ln Z \end{aligned}$$

Since Z is the normalizing constant, doesn't depend on α .

\therefore minimize KL-divergence \rightarrow maximize $F[P_0, \alpha]$

$$F[P_0, \alpha] = \underbrace{-E_{\alpha}[\ln \alpha(x)]}_{\text{entropy term}} + \underbrace{\sum_{\mu} E_{\alpha}[\ln \phi_{\mu}(U_{\mu})]}_{\text{energy term}}$$

Optimizing the Energy Functional

$$\begin{aligned} D(\alpha \| P_0) &= \int \alpha(x) \cdot \ln \frac{\alpha(x)}{P_0(x)} dx = -H_{\alpha}(x) - \sum_{\mu} E_{\alpha}[\ln \phi_{\mu}(U_{\mu})] + \ln Z \\ &= -F[P_0, \alpha] + \ln Z \end{aligned}$$

$$D(\alpha \| P_0) \geq 0 \quad \ln Z \geq F[P_0, \alpha]$$

the energy functional is a lower bound of the log of the partition function

Variational method: want to solve a problem by introducing new variational parameters that increase the degrees of freedom over which we optimize. Each choice of those parameters gives an approximate answer.

Exact Inference as Optimization

Given a cluster tree τ with a set of beliefs \mathcal{Q} , define the factored energy functional

$$\tilde{F}[\mathbb{P}_\theta, \mathcal{Q}] = \sum_{i \in V_\tau} E_{\text{comp}}[\ln \psi_i] + \sum_{i \in V_\tau} H_{\beta_i}(C_i) - \sum_{(i,j) \in E_\tau} H_{\alpha_j}(S_j)$$

where ψ_i is the initial potential assigned to C_i : $\psi_i = \prod_{A \in \mathcal{A}_i} \phi$

$$\tilde{F}[\mathbb{P}_\theta, \mathcal{Q}] = \sum_{i \in V_\tau} E_{\text{comp}}[\ln \prod_{A \in \mathcal{A}_i} \phi] - \sum_{i \in V_\tau} \int_{C_i} \beta_i(C_i) \ln \beta_i(C_i) dC_i + \sum_{(i,j) \in E_\tau} \int_{S_j} \alpha_j(S_j) \ln \alpha_j(S_j) dS_j$$

$$\begin{aligned} \text{I}^\circ \sum_i E_{\text{comp}}[\ln \psi_i] &= \sum_i \int_{C_i} \beta_i(C_i) \ln \prod_{A \in \mathcal{A}_i} \phi(u_A) dC_i \\ &= \sum_i \sum_{A \in \mathcal{A}_i} \int_{C_i} \alpha_i(C_i) \ln \phi(u_A) dC_i \\ &= \sum_{A \in \mathcal{A}} E_{\text{comp}}[\ln \phi] \end{aligned}$$

$$\begin{aligned} \text{I}^\circ H_{\alpha_j}(S_j) &= - \int_{S_j} \alpha_j(S_j) \ln \alpha_j(S_j) dS_j = - \int_{S_j} \alpha_j(S_j) \cdot \ln \frac{\prod_{i \in V_\tau} \beta_i(C_i)}{\prod_{(i,j) \in E_\tau} \alpha_j(S_j)} dS_j \\ &= - \int_{S_j} \alpha_j(S_j) \cdot \sum_{i \in V_\tau} \ln \beta_i(C_i) dS_j + \int_{S_j} \alpha_j(S_j) \sum_{(i,j) \in E_\tau} \ln \alpha_j(S_j) dS_j \\ &= - \sum_{i \in V_\tau} \int_{S_j} \alpha_j(S_j) \ln \beta_i(C_i) dS_j + \sum_{(i,j) \in E_\tau} \int_{S_j} \alpha_j(S_j) \ln \alpha_j(S_j) dS_j \\ &= - \sum_{i \in V_\tau} \int_{C_i} \alpha_j(S_j) \ln \beta_i(C_i) dC_i + \sum_{(i,j) \in E_\tau} \int_{S_j} \alpha_j(S_j) \ln \alpha_j(S_j) dS_j \\ &= \sum_{i \in V_\tau} H_{\beta_i}(C_i) - \sum_{(i,j) \in E_\tau} H_{\alpha_j}(S_j) \end{aligned}$$

$$F[\mathbb{P}_\theta, \mathcal{Q}] = \sum_{A \in \mathcal{A}} E_{\text{comp}}[\ln \phi] + H_{\alpha_j}(S_j) = \sum_i E_{\text{comp}}[\ln \psi_i] + \sum_{i \in V_\tau} H_{\beta_i}(C_i) - \sum_{(i,j) \in E_\tau} H_{\alpha_j}(S_j) = \tilde{F}[\mathbb{P}_\theta, \mathcal{Q}]$$

$$\text{Find } \mathcal{Q} = \{ \beta_i : i \in V_\tau \} \cup \{ \alpha_j : (i,j) \in E_\tau \}$$

$$\text{max. } F[\mathbb{P}_\theta, \mathcal{Q}] = \tilde{F}[\mathbb{P}_\theta, \mathcal{Q}]$$

$$\text{s.t. } \alpha_j(S_j) = \prod_{i \in V_\tau} \beta_i(C_i)$$

$$\sum_{i \in V_\tau} \beta_i(C_i) = 1$$

$$\beta_i(C_i) \geq 0$$

$$\forall (i,j) \in E_\tau, \forall S_j \in \text{Val}(S_j)$$

$$\forall i \in V_\tau$$

$$\forall i \in V_\tau, \forall C_i \in \text{Val}(C_i)$$

} marginal consistency

} joint probability

Fixed-Point Characterization

If a clique tree τ is an I-map of P_0 , then there is a unique solution to

$$\begin{aligned} \text{Find } Q &= \{\beta_i : i \in U\} \cup \{\mu_{ij} : (ij) \in E_\tau\} \\ \text{min. } D(Q \| P_0) &= F[P_0, Q] + \ln(z) = -\tilde{F}[P_0, Q] \\ \text{s.t. } \mu_{ij}(S_{ij}) &= \sum_{C \in S_{ij}} \beta_i(C) \quad \forall (ij) \in E_\tau \quad \forall S_{ij} \in \text{Val}(S_{ij}) \end{aligned}$$

$$D(Q \| P_0) = \int_{\mathcal{X}} Q(x) \ln \frac{Q(x)}{P_0(x)} dx$$

$$= \int_{\mathcal{X}} Q(x) \ln Q(x) dx - \int_{\mathcal{X}} Q(x) \ln P_0(x) dx$$

$$= -H(Q) - E_{Q, \tau}[\ln P_0(x)]$$

$$P_0(z) = (\prod_i \phi_i)(x) \cdot \frac{1}{z} = \frac{1}{z} \tilde{P}_0(z)$$

$$= -H(Q) - E_{Q, \tau}[\ln(\prod_i \phi_i)(x)] + \ln z$$

$$= -H(Q) - \sum_{i \in U} E_{Q, \tau}[\ln \phi_i(U)] + \ln z$$

$$= -F[P_0, Q] + \ln z$$

Energy functional: $F[P_0, Q] = \sum_i E_{Q, \tau}[\ln \phi_i] + H(Q)$

$$= -\sum_{i \in U} E_{Q, \tau}[\ln \phi_i] - \sum_{i \in U} H(\beta_i) + \sum_{(ij) \in E_\tau} H(\mu_{ij}) + \ln z$$

$$= -\tilde{F}[P_0, Q] + \ln z$$

Factorial energy functional $\tilde{F}[P_0, Q] = \sum_{i \in U} E_{Q, \tau}[\ln \phi_i] + \sum_{i \in U} H(\beta_i) - \sum_{(ij) \in E_\tau} H(\mu_{ij})$

$$\text{Find } Q = \{\beta_i : i \in U\} \cup \{\mu_{ij} : (ij) \in E_\tau\}$$

$$\text{min. } -\sum_{i \in U} \int_{\mathcal{C}_i} \beta_i(C_i) \ln \mu_{ij}(C_i) dC_i + \sum_{i \in U} \int_{\mathcal{C}_i} \beta_i(C_i) \ln \beta_i(C_i) dC_i - \sum_{(ij) \in E_\tau} \int_{S_{ij}} \mu_{ij}(S_{ij}) \ln \mu_{ij}(S_{ij}) dS_{ij}$$

$$\text{s.t. } \mu_{ij}(S_{ij}) = \left(\prod_{C_i \in S_{ij}} \beta_i(C_i) \right) \mu_{ij}(S_{ij}) \quad \forall (ij) \in E_\tau \quad \forall S_{ij} \in \text{Val}(S_{ij}) \quad \text{dual variable } \lambda_{j \rightarrow i}[S_{ij}]$$

$$\sum_{C_i} \beta_i(C_i) = 1 \quad \forall i \in U \quad \text{dual variable } \lambda_i$$

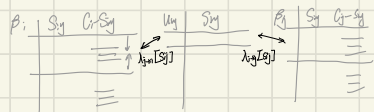
$$\beta_i(C_i) > 0 \quad \forall i \in U, C_i \in \text{Val}(C_i) \quad // \text{implicitly satisfied}$$



$$\mathcal{L} = -\tilde{F}[P_0, Q] + \sum_{i \in U} \lambda_i \left(\sum_{C_i} \beta_i(C_i) - 1 \right) + \sum_{(ij) \in E_\tau} \sum_{S_{ij}} \lambda_{j \rightarrow i}[S_{ij}] \left(\mu_{ij}(S_{ij}) - \left(\prod_{C_i \in S_{ij}} \beta_i(C_i) \right) \mu_{ij}(S_{ij}) \right)$$

$$\frac{\partial \mathcal{L}}{\partial \beta_i(C_i)} = -\ln \mu_{ij}(C_i) + \ln \beta_i(C_i) + 1 + \lambda_i - \sum_{(ij) \in E_\tau} \lambda_{j \rightarrow i}[S_{ij}] := 0$$

$$\beta_i(C_i) = \exp(+\lambda_i) \cdot \psi_i(C_i) \cdot \prod_{(ij) \in E_\tau} \exp(\lambda_{j \rightarrow i}[S_{ij}])$$



$$\frac{\partial \mathcal{L}}{\partial \mu_{ij}(S_{ij})} = -\ln \mu_{ij}(S_{ij}) - 1 + \lambda_{j \rightarrow i}[S_{ij}] + \lambda_{i \rightarrow j}[S_{ij}] := 0$$

$$\mu_{ij}(S_{ij}) = \exp(-1) \cdot \exp(\lambda_{j \rightarrow i}[S_{ij}]) \cdot \exp(\lambda_{i \rightarrow j}[S_{ij}])$$

$$\text{let } \hat{\mu}_{ij}[S_{ij}] = \exp(\lambda_{j \rightarrow i}[S_{ij}] - \frac{1}{2})$$

$$\mu_{ij}(S_{ij}) = \hat{\mu}_{i \rightarrow j}(S_{ij}) \cdot \hat{\mu}_{j \rightarrow i}(S_{ij})$$

$$\beta_i(C_i) = \exp(-\lambda_i) \cdot \psi_i(C_i) \cdot \prod_{(ij) \in E_\tau} \exp(\lambda_{j \rightarrow i}[S_{ij}] - \frac{1}{2}) \cdot \exp(\frac{\lambda_i}{2})$$

$$\beta(C_i) = \exp(+\lambda_i + \frac{\|N_i\|}{2}) \cdot \psi_i(C_i) \cdot \prod_{j \in N_i} \beta_j(S_j)$$

$$\begin{aligned} \delta_{i,j}(S_{ij}) &= \frac{U_{i,j}(S_{ij})}{\beta_{i,j}(S_{ij})} = \frac{(\sum_{C_{i,j}} \beta_i)(S_{ij})}{\beta_{i,j}(S_{ij})} \\ &= \exp(+\lambda_i + \frac{\|N_i\|}{2}) \cdot (\sum_{C_{i,j}} \psi_i)(S_{ij}) \prod_{k \in N_i, k \neq j} \delta_{k,i}(S_{kj}) \end{aligned}$$

$$\delta_{i,j}(S_{ij}) = C_i \cdot (\sum_{C_{i,j}} \psi_i) \prod_{k \in N_i, k \neq j} \delta_{k,i}(S_{kj}) \quad \text{belief propagation update rule}$$

A set of beliefs \mathcal{Q} is a stationary point of the optimization problem iff

$$\exists \{ \delta_{i,j}(S_{ij}) \mid (i,j) \in E \}$$

st.

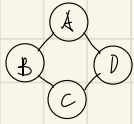
$$\delta_{i,j} \propto \sum_{C_{i,j}} \psi_i \left(\prod_{k \in N_i, k \neq j} \delta_{k,i} \right)$$

Moreover, we have

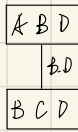
$$\beta_i \propto \psi_i \left(\prod_{j \in N_i} \delta_{j,i} \right)$$

$$U_{ij} = \delta_{i,j} \cdot \delta_{j,i}$$

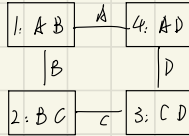
Propagation-Based Approximation: Example



Markov network



A clique tree



A cluster graph

the cluster graph contains loops (loop), we can still apply belief-update propagation (nothing in the algorithm relies on the fact that it's a tree)

$$1. \delta_{1,2}(B) = \sum_C \psi_1(A, B, C)$$

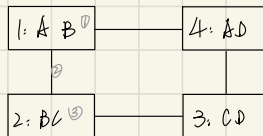
$$\beta_2(B, C) = \beta_2(B, C) \cdot \delta_{1,2}(B) / U_{2,1}(B) = \psi_2(B, C) \cdot \sum_A \psi_1(A, B)$$

$$U_{2,1}(B) = \delta_{2,1}(B) = \sum_C \psi_2(B, C)$$

$$2. \delta_{2,3}(C) = \sum_B \beta_2 = \sum_B \psi_2(B, C) \cdot \sum_A \psi_1(A, B)$$

$$\beta_3(C, D) = \beta_3(C, D) \cdot \delta_{2,3}(C) / U_{3,2}(C) = \psi_3(C, D) \cdot \sum_B \psi_2(B, C) \cdot \sum_A \psi_1(A, B)$$

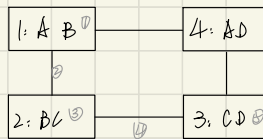
$$U_{3,2}(C) = \delta_{3,2}(C) = \sum_B \psi_3(B, C) \cdot \sum_A \psi_1(A, B)$$



$$1. \beta_1(A, B) = \psi_1(A, B)$$

$$2. U_{1,2}(B) = \delta_{1,2}(B) = \sum_C \psi_1(A, B, C) = \sum_A \psi_1(A, B)$$

$$3. \beta_2(B, C) = \psi_2(B, C) \cdot \sum_A \psi_1(A, B)$$



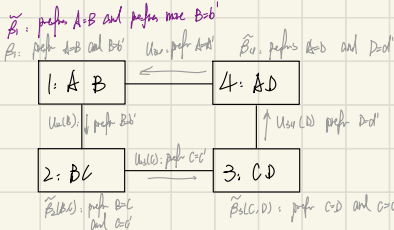
$$1. \beta_1(A, B) = \psi_1(A, B)$$

$$2. U_{1,2}(B) = \delta_{1,2}(B) = \sum_C \psi_1(A, B, C) = \sum_A \psi_1(A, B)$$

$$3. \beta_2(B, C) = \psi_2(B, C) \cdot \sum_A \psi_1(A, B)$$

$$4. U_{2,3}(C) = \delta_{2,3}(C) = \sum_B \psi_2(B, C) \cdot \sum_A \psi_1(A, B)$$

if all clusters favor consensus joint assignment (i.e. $\beta_i(a^*b^*)$ and $\beta_i(a^*b)$) $\Rightarrow \beta_i(a^*,b^*)$ and $\beta_i(a^*,b)$



the process may not converge
for cluster

Cluster-Graph Belief Propagation

A (loop) cluster graph \mathcal{U} satisfies the **running intersection property** if

whenever there is a variable such that $X \in C_i$ and $X \in C_j$,

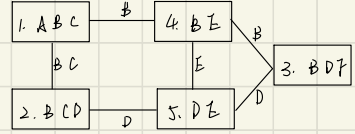
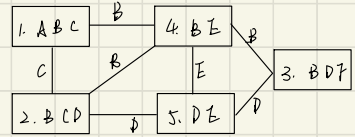
then there is a single path between C_i and C_j for which $X \in S_e \forall e$ in the path

running intersection property:

$\mathcal{P} \ni$ a path: information about x flow between all clusters that contain it

so that in a calibrated cluster graph, all clusters must agree about the marginal of x

\times only one path: prevents information about x from cycling endlessly in a loop



A cluster graph is **calibrated** if

$$\forall (i,j) \in \mathcal{E}_{\mathcal{U}}, \sum_{C \in \mathcal{S}_{ij}} \beta_i = \sum_{C \in \mathcal{S}_{ij}} \beta_j$$

(two clusters agrees on marginal of variables in \mathcal{S}_{ij})

not necessarily all vars \geq cliques have in common

def initialize_cluster_graph(\mathcal{U}):

for each cluster C_i :

$$\beta_i = \prod_{C \in \mathcal{C}_i} \phi$$

for each $(i,j) \in \mathcal{E}_{\mathcal{U}}$

$$\mathcal{S}_{ij} \neq \emptyset \quad \mathcal{S}_{j,i} \neq \emptyset$$

def sum_product_message (i : sending clique, j : receiving clique):

$$\psi_i(C_i) = \psi_i \cdot \prod_{k \in \mathcal{K}(C_i) \setminus \{j\}} \mathcal{S}_{i,k}$$

$$\tau_i(\mathcal{S}_{ij}) = \sum_{C \in \mathcal{S}_{ij}} \psi_i(C)$$

return $\tau_i(\mathcal{S}_{ij})$

def cluster_graph_sum_product_calibrate (\mathcal{F} : set of factors, \mathcal{U} : cluster graph):

initialize_cluster_graph(\mathcal{U})

while graph is not calibrated

select $(i,j) \in \mathcal{E}_{\mathcal{U}}$

$$\mathcal{S}_{ij}(\mathcal{S}_{ij}) = \text{sum_product_message}(i,j)$$

for each clique i

$$\beta_i \leftarrow \psi_i \cdot \prod_{k \in \mathcal{K}(C_i) \setminus \{j\}} \mathcal{S}_{i,k}$$

return $\{\beta_i\}$

other than the fact that the algorithm is applied to graphs rather than trees,
the algorithm is identical to sum-product calibration of clique trees (initialize all messages to 1)

can use belief-update messages to define belief-update calibration for cluster graphs

def initialize_cluster_graph():

for each cluster C_i :

$$\beta_i = \prod_{\phi \in C_i} \phi$$

for each edge $(i, j) \in E_u$

$$u_{ij} = 1$$

def belief_update_message(i: sending clique, j: receiving clique):

$$o_{i \rightarrow j} = \prod_{C_k \in \mathcal{C}_i} \beta_k$$

$$\beta_j = \beta_j \cdot \frac{o_{i \rightarrow j}}{u_{ij}}$$

$$u_{ij} = o_{i \rightarrow j}$$

def cluster_graph_belief_update_calibration(\mathcal{E} : set of factors, U : cluster graph over \mathcal{E}):

initialize_cluster_graph()

while graph is not calibrated:

select $(i, j) \in E_u$

belief_update_message(i, j)

return $\{\beta_i\}$, $\{u_{ij}\}$

Properties of cluster-graph belief propagation: Re-parameterization

Let U be a generalized cluster graph over a set of factors \mathcal{E} .

Consider the set of beliefs $\{\beta_i\}$ and separators $\{u_{ij}\}$ as any iteration of CB-CP-calibration

$$\tilde{P}_{\mathcal{E}}(Z) = \frac{\prod_{i \in \mathcal{C}_u} \beta_i(C_i)}{\prod_{(i, j) \in E_u} u_{ij}(S_{ij})}$$

where $\tilde{P}_{\mathcal{E}}(Z) = \prod_{\phi \in \mathcal{E}} \phi$ is the unnormalized distribution defined by \mathcal{E}

At each iteration: $\beta_i = \psi_i \prod_{j \in \mathcal{C}_i} \beta_j$

$u_{ij}(S_{ij}) = \delta_{ij} \cdot \beta_j$

$$\begin{aligned} \frac{\prod_{i \in \mathcal{C}_u} \beta_i(C_i)}{\prod_{(i, j) \in E_u} u_{ij}(S_{ij})} &= \frac{\prod_{i \in \mathcal{C}_u} \psi_i(C_i) \prod_{j \in \mathcal{C}_i} \beta_j(S_{ij})}{\prod_{(i, j) \in E_u} \delta_{ij}(S_{ij}) \beta_j(S_{ij})} \\ &= \prod_{i \in \mathcal{C}_u} \psi_i(C_i) \\ &= \tilde{P}_{\mathcal{E}}(Z) \end{aligned}$$

cluster-graph belief propagation preserves all the information about the original distribution
it does not "distill" the original factors by performing propagation only loops.

This process kind of tries to represent the original factors in a more useful form

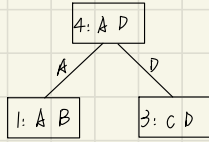
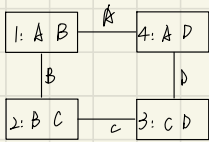
Properties of Cluster-Graph Belief Propagation: Tree Consistency

In a calibrated cluster tree, the belief over a cluster is the marginal of the distribution.

To characterize the beliefs we get by calibrating a cluster graph.

We can use the cluster tree invariance property applied to subtrees of a cluster graph.

A subtree T of U is a subset of clusters and edges that together form a tree that satisfies the running intersection property.

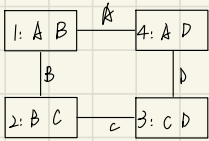


Removing edges from a cluster graph may result in violating the running intersection property.

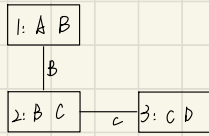
Once we select a tree T , we can think of it as defining a distribution $P_T(X) = \frac{\prod_{C \in T} \beta_C(C_i)}{\prod_{i \in \text{nodes}} \psi_i(S_i)}$ if the cluster graph is calibrated. So is the sub cluster tree, which satisfies the running intersection property.

$$\therefore \beta_C(C_i) = P_T(C_i)$$

Tree consistency: The beliefs over C_i in the tree are the marginals of P_T .



cluster graph U



cluster tree T (also a (calibrated) cluster graph)

$$P_U(A, B, C, D) = \frac{\beta_1(A, B) \cdot \beta_2(B, C) \cdot \beta_3(C, D) \cdot \beta_4(A, D)}{\psi_1(B) \cdot \psi_2(C) \cdot \psi_3(D) \cdot \psi_4(A)}$$

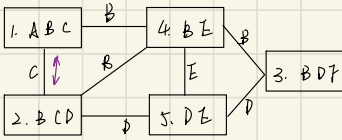
$$P_T(A, B, C, D) = P_U(A, B, C, D) \cdot \frac{\psi_1(B) \cdot \psi_2(C)}{\beta_1(A, B)}$$

Constructing Cluster Graphs

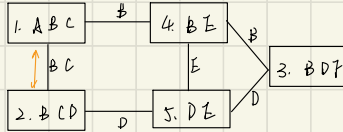
In the case of clique trees, different clique trees of a graphical model give the same answer

In the case of cluster graphs, different cluster graphs can lead to different answers

We have to consider the accuracy-cost trade-offs



message above C



message above BC

Assume $G(A, B, C)$ strongly prefer $B=C$

U_1 : the correlation is already conveyed from G_1 to G_2

U_2 : marginal on C is conveyed on (1,2), marginal on B is conveyed on (1-4-3)

the spurious dependency is lost

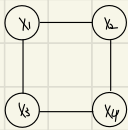
if marginal of C in $G(A, B, C)$ is uniform, then $S_{\text{var}}(C)$ is uniform, in U_1 , B and C seems to be two independent uniform to G_2

eg. pairwise Markov Network

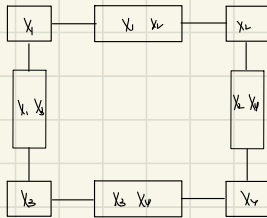
Univariate potential $\phi_i(x_i)$ for each variable x_i

pairwise potential $\phi_{ij}(x_i, x_j)$ over some pairs of variables

(any distribution can be reformulated as a pairwise Markov network with transformation of variables)



pairwise Markov network



cluster graph

Variational Analysis

for some cluster graph U and a distribution \mathcal{P} , define

$$\mathcal{Q}_{\mathcal{P}} = \{P(c_i) \mid i \in U\} \cup \{P(s_{ij}) \mid ij \in E_U\}$$

define the marginal polytope

$$\text{Marg}[U] = \{Q_{\mathcal{P}} : \mathcal{P} \text{ is a distribution over } \mathcal{Z}\}$$

The marginal polytope is the set of all cluster and separator beliefs that can be obtained by marginalizing an actual \mathcal{P}
 (there are calibrated cluster beliefs that do not represent the marginals of any single coherent distribution over \mathcal{Z})

The marginal polytope has exponentially many faces, cannot optimize over $\text{Marg}[U]$
 instead, we optimize the

$$\text{Local}[U] = \left\{ \begin{array}{l} \{\beta_i \mid i \in U\} \cup \\ \{u_{ij} \mid ij \in E_U\} \end{array} \right\} \left| \begin{array}{l} u_{ij}(s_{ij}) = \sum_{c_i, c_j} \beta_i(c_i) \beta_j(c_j) \quad \forall ij \in E_U, s_{ij} \in \text{val}(s_{ij}) \\ \sum_i \beta_i(c_i) = 1 \quad \forall i \in U \\ \beta_i(c_i) \geq 0 \quad \forall i \in U, c_i \in \text{val}(c_i) \end{array} \right\}$$

$$\begin{aligned} D(Q \| P_0) &= \int_{\mathcal{X}} Q(x) \cdot \ln \frac{Q(x)}{P_0(x)} dx \\ &= \int_{\mathcal{X}} Q(x) \cdot \ln Q(x) dx - \int_{\mathcal{X}} Q(x) \cdot \ln P_0(x) dx \\ &= -H(Q) - \mathbb{E}_{\text{true}}[\ln \prod_i \phi_i] + \ln Z \\ &= -H(Q) - \sum_i \mathbb{E}_{\text{true}}[\ln \phi_i] + \ln Z \\ &= -\tilde{F}[\tilde{P}_0, Q] + \ln Z \end{aligned} \quad \left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} \begin{array}{l} -\tilde{F}[\tilde{P}_0, Q] + \ln Z = D(Q \| P_0) \geq 0 \\ \tilde{F}[\tilde{P}_0, Q] \leq \ln Z \end{array}$$

Unlike for clique tree, $\tilde{F}[\tilde{P}_0, Q]$ is no longer a reformulation of the energy function, but an approximation of the energy function

Find Q

max. $\tilde{F}[\tilde{P}_0, Q]$

s.t. $Q \in \text{Local}[U]$

same optimization problem as for clique tree calibration,

but with \geq approximations (relaxations)

1° use the factored energy as an approximation of the true energy function

2° optimize over space of pseudo-marginals instead of the space of all coherent distributions

A set of beliefs Q is a stationary point of the optimization problem iff

$$\forall (ij) \in E_U, \text{ there are auxiliary factors } \delta_{i \rightarrow j}(s_{ij}) \text{ and } \delta_{j \rightarrow i}(s_{ij}) \text{ st. } \delta_{i \rightarrow j} \propto \sum_{c_j} \psi_i \cdot \prod_{k \in \text{neigh}(j)} \delta_{k \rightarrow j}$$

and we have

$$\beta_i \propto \psi_i \cdot \prod_{j \in \text{neigh}(i)} \delta_{j \rightarrow i}$$

$$u_{ij} = \delta_{i \rightarrow j} \cdot \delta_{j \rightarrow i}$$

Structured Variational Approximations

the structured variational approach aims to optimize the energy functional over a family \mathcal{Q} of reference distributions Q . This family is chosen to be computationally tractable, hence it's generally not sufficiently expressive to capture all of the information in P_0

$$\begin{aligned} \text{Find } Q \in \mathcal{Q} \\ \text{max. } F[P_0, Q] \end{aligned}$$

\mathcal{Q} is a given family of distributions.

$F[P_0, Q]$ is the exact energy functional, thus maximizing $F[Q]$ \leftrightarrow minimizing $D(Q||P_0)$

choose simple \mathcal{Q} :

- i: \mathcal{Q} can be described by a BN or MN with small tree-width, more efficient inference
- ii: poor approximation of P_0

Structured Variational Approximation: Mean Field Approximation

The mean field approximation finds the distribution Q that's closest to P_0 in terms of $D(Q||P_0)$ within the class of distributions representable as a product of independent marginals

$$Q(x) = \prod_i Q(x_i)$$

the approximation of P_0 as a fully factorized distribution is likely to lose a lot of information in P_0

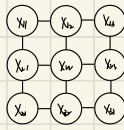
the energy functional $F[P_0, Q] = H_0(x) + \sum_{i \sim j} E_{x_{i,j}}[\ln \phi_{ij}]$

$$\begin{aligned} E_{x_{i,j}}[\ln \phi_{ij}] &= E_{\prod_{k \in \{i,j\}} Q(x_k)}[\ln \phi_{ij}] = \int_{\prod_{k \in \{i,j\}} Q(x_k)} Q(x_i) \cdot \ln \phi_{ij}(x_i, x_j) d x_i d x_j \\ &= \int_{\prod_{k \in \{i,j\}} Q(x_k)} \prod_{k \in \{i,j\}} Q(x_k) \ln \phi_{ij}(x_i, x_j) d x_i d x_j \end{aligned}$$

$$\begin{aligned} H_0(x) &= - \int_{\prod_{i \in \mathcal{V}} Q(x_i)} \ln Q(x) d x = - \int_{\prod_{i \in \mathcal{V}} Q(x_i)} \sum_i Q(x_i) \cdot \ln \prod_i Q(x_i) d x \\ &= - \sum_i \int_{\prod_{i \in \mathcal{V}} Q(x_i)} \prod_i Q(x_i) \ln Q(x_i) d x \\ &= - \sum_i \int_{\prod_{i \in \mathcal{V}} Q(x_i)} \prod_{i \neq i} Q(x_i) \cdot \int_{x_i} Q(x_i) \ln Q(x_i) d x_i \dots d x_i \quad // Q(x_i) \text{ is a marginal and integrates to 1} \\ &= - \sum_i \int_{x_i} Q(x_i) \ln Q(x_i) d x_i \\ &= \sum_i H_0(x_i) \end{aligned}$$

$$\text{if } Q(x) = \prod_i Q(x_i), \text{ then } H_0(x) = \sum_i H_0(x_i)$$

eg. 4x4 grid Markov network



$$\begin{aligned} \mathbb{E}_{\mathbb{P}_\theta}[\mathcal{Q}] &= H_\theta(x_i) + \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{I}} E_{\mathcal{Q}}[\ln \phi(y_{ij})] \\ &= \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{I}} E_{\mathcal{Q}}[\ln \phi(x_{ij}, x_{i,j})] \quad // \text{vertical potentials} \\ &\quad + \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{I}} E_{\mathcal{Q}}[\ln \phi(x_{ij}, x_{i,j+1})] \quad // \text{horizontal potentials} \\ &\quad + \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{I}} H_\theta(x_{ij}) \quad // \text{node potentials} \end{aligned}$$

Find $\mathcal{Q}(x_i)$
 min. $-\mathbb{E}_{\mathbb{P}_\theta}[\mathcal{Q}] = \sum_i \int_{\mathcal{X}_i} \mathcal{Q}(x_i) \cdot \ln \mathcal{Q}(x_i) dx_i - \sum_{\text{edges}(ij)} \int_{\mathcal{X}_{ij}} \prod_{i \in \text{edge}(ij)} \mathcal{Q}(x_i) \ln \phi(y_{ij}) dx_{ij}$
 s.t. $\sum_i \mathcal{Q}(x_i) = 1 \quad \forall i$

$$D(\mathcal{Q} \| \mathbb{P}_\theta) = \int_{\mathcal{X}} \mathcal{Q}(x) \cdot \ln \frac{\mathcal{Q}(x)}{\mathbb{P}_\theta(x)} dx \quad \text{is convex in } \mathcal{Q}(x) \text{ if } x \in \text{val}(x)$$

$$\mathcal{Q}(x) = \prod_{i \in \text{edge}(x)} \mathcal{Q}(x_i) \quad \text{is jointly convex in } \mathcal{Q}(x_i) \text{ and increasing if } \mathcal{Q}_i$$

$$D(\mathcal{Q} \| \mathbb{P}_\theta) = -\mathbb{E}_{\mathbb{P}_\theta}[\mathcal{Q}] + \ln Z \quad \therefore -\mathbb{E}_{\mathbb{P}_\theta}[\mathcal{Q}] \text{ is convex}$$

$$\mathcal{L} = \sum_i \int_{\mathcal{X}_i} \mathcal{Q}(x_i) \cdot \ln \mathcal{Q}(x_i) dx_i - \sum_{\text{edges}(ij)} \int_{\mathcal{X}_{ij}} \prod_{i \in \text{edge}(ij)} \mathcal{Q}(x_i) \ln \phi(y_{ij}) dx_{ij} + \sum_i \lambda_i (\sum_i \mathcal{Q}(x_i) - 1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathcal{Q}(x_i)} = 1 + \ln \mathcal{Q}(x_i) - \sum_{\substack{\text{edge}(ij) \\ x_i \in \text{edge}(ij)}} \int_{\mathcal{X}_{ij}} \prod_{j \in \text{edge}(ij)} \mathcal{Q}(x_j) \ln \phi(y_{ij}) dx_{ij} + \lambda$$

$$= 1 + \ln \mathcal{Q}(x_i) - \sum_{\mathbb{P}_\theta} E_{\mathcal{Q}}[\ln \phi | x_i = x_i] + \lambda \quad \therefore$$

$$\ln \mathcal{Q}(x_i) = -\lambda - 1 + \sum_{\mathbb{P}_\theta} E_{\mathcal{Q}}[\ln \phi | x_i = x_i]$$

The distribution \mathcal{Q} is a stationary point of mean-field optimization iff

$$\mathcal{Q}(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\mathbb{P}_\theta} E_{\mathcal{Q}}[\ln \phi | x_i = x_i] \right\}$$

$$\begin{aligned} \sum_{\mathbb{P}_\theta} E_{\mathcal{Q}}[\ln \phi | x_i = x_i] &= E_{\mathcal{Q}} \left[\sum_{\mathbb{P}_\theta} \ln \phi | x_i = x_i \right] = E_{\mathcal{Q}}[\ln \tilde{\mathbb{P}}_\theta(z) | x_i = x_i] \\ &= E_{x_{i-1} \sim \mathcal{Q}}[\ln \tilde{\mathbb{P}}_\theta(x_i, x_i)] \quad // \mathcal{Q} \text{ is product of marginals, } x_i = \sum_{\mathcal{I}} x_i \\ &= E_{x_{i-1} \sim \mathcal{Q}}[\ln \mathbb{P}_\theta(x_i | x_{i-1})] + E_{x_{i+1} \sim \mathcal{Q}}[\ln \mathbb{P}_\theta(x_i | x_{i+1})] \end{aligned}$$

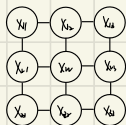
$$\mathcal{Q}(x_i) = \frac{1}{Z_i} \exp \left\{ E_{x_{i-1} \sim \mathcal{Q}}[\ln \mathbb{P}_\theta(x_i | x_{i-1})] \right\} \underbrace{\exp \left\{ E_{x_{i+1} \sim \mathcal{Q}}[\ln \mathbb{P}_\theta(x_i | x_{i+1})] \right\}}_{\text{does not depend on } x_i}$$

In the mean field approximation, $Q(x_i)$ is locally optimal only if

$$Q(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi: X_i \in \text{scope}[\phi]} E_{\text{rest } X_i \sim Q} [\ln \phi(U_\phi, x_i)] \right\}$$

$Q(x_i)$ has to be consistent with the expectation of the potential in which it appears

$$Q(x_{ij}) = \frac{1}{Z_{ij}} \exp \left\{ \begin{aligned} & \sum_{x_{ij}} Q(x_{ij}) \cdot \ln \phi(x_{ij}, x_{ij}) && \text{①} \\ & + \sum_{x_{ij+1}} Q(x_{ij+1}) \cdot \ln \phi(x_{ij}, x_{ij+1}) && \text{②} \\ & + \sum_{x_{i-1j}} Q(x_{i-1j}) \cdot \ln \phi(x_{i-1j}, x_{ij}) && \text{③} \\ & + \sum_{x_{ij+1}} Q(x_{ij+1}) \cdot \ln \phi(x_{ij}, x_{ij+1}) && \text{④} \end{aligned} \right.$$



each term is a geometric mean of one potential involving x_{ij}

def mean field approximation (\tilde{Q}, Q_0) :

$$Q = Q_0$$

$$\text{Unprocessed} = \mathcal{S}$$

while $\text{Unprocessed} \neq \emptyset$:

choose x_i from Unprocessed :

$$Q_{\text{old}}(x_i) = Q(x_i)$$

for $x_i \in \text{Val}(x_i)$

$$Q(x_i) = \exp \left\{ \sum_{\phi: X_i \in \text{scope}[\phi]} E_{\text{rest } X_i \sim Q} [\ln \phi(U_\phi, x_i)] \right\}$$

Coordinate ascent

Normalize $Q(x)$

if $Q_{\text{old}}(x_i) \neq Q(x_i)$

$$\text{Unprocessed} = \text{Unprocessed} \cup \left(\bigcup_{\phi: X_i \in \text{scope}[\phi]} \text{scope}[\phi] \right)$$

$$\text{Unprocessed} = \text{Unprocessed} - \{x_i\}$$

return Q