

Probabilistic Graphical Models

Lecture 15 – Inference as Optimization

CS/CNS/EE 155

Andreas Krause

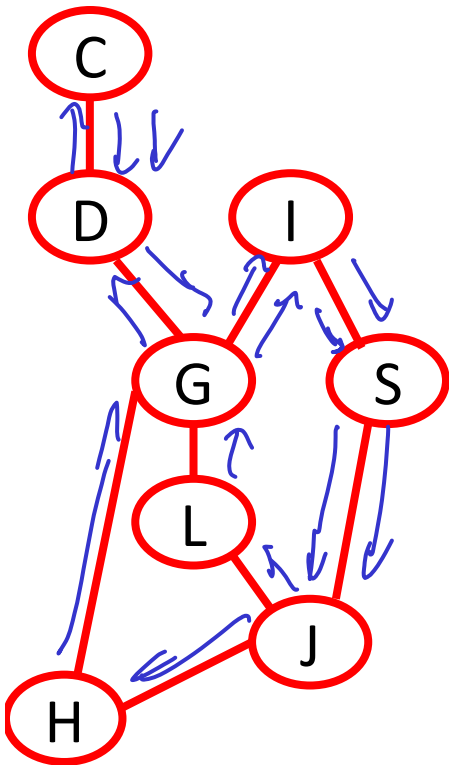
Announcements

- Homework 3 due next Monday (Nov 23)
- Project poster session on Friday December 4 (tentative)
- Final writeup (8 pages NIPS format) due Dec 9

Approximate inference

- Three major classes of general-purpose approaches
- **Message passing**
 - E.g.: Loopy Belief Propagation
- **Inference as optimization**
 - Approximate posterior distribution by simple distribution
 - Mean field / structured mean field
- **Sampling based inference**
 - Importance sampling, particle filtering
 - Gibbs sampling, MCMC
- Many other alternatives (often for special cases)

Loopy BP on arbitrary pairwise MNs



- What if we apply BP to a graph with loops?
 - Apply BP and hope for the best..

$$\delta_{i \rightarrow j}(X_j) = \sum_{x_i} \pi_i(x_i) \pi_{i,j}(x_i, X_j) \prod_{s \in N(i) \setminus \{j\}} \delta_{s \rightarrow i}(x_i)$$

- Will not generally converge.. ☹
- If it converges, will not necessarily get correct marginals ☹
- However, in practice, answers often still useful!

Approximate inference

- Three major classes of general-purpose approaches
- **Message passing**
 - E.g.: Loopy Belief Propagation (today!)
- **Inference as optimization**
 - Approximate posterior distribution by simple distribution
 - Mean field / structured mean field
 - Assumed density filtering / expectation propagation
- **Sampling based inference**
 - Importance sampling, particle filtering
 - Gibbs sampling, MCMC
- Many other alternatives (often for special cases)

Variational approximation

- Graphical model with intractable (high-treewidth) joint distribution $P(X_1, \dots, X_n)$

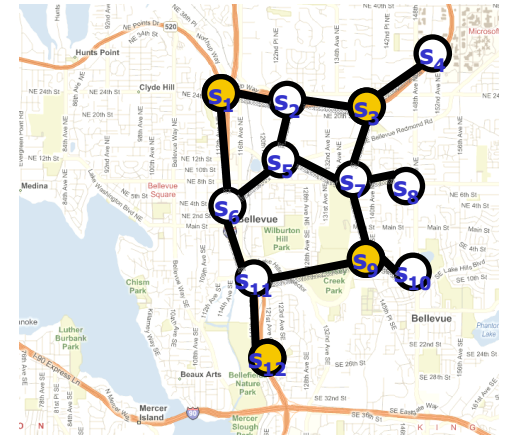
- Want to compute posterior distributions

$$\begin{aligned} & P(X_3, X_4 \mid X_1 = x_1, X_7 = x_7) \\ & \propto P(X_3, X_4, X_1 = x_1, X_7 = x_7) = \sum_{x_1, x_2, x_5, \dots} P(x_1, x_2, \dots, x_n) \end{aligned}$$

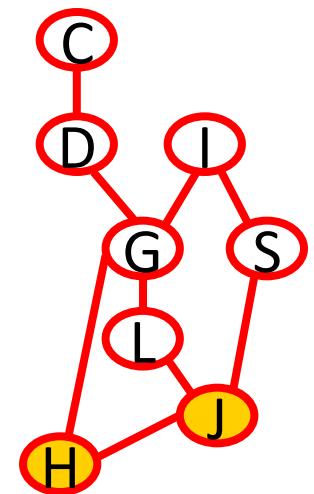
- Computing posterior exactly is intractable
- **Key idea:** Approximate posterior with *simpler* distribution that's *as close* to P as possible

Why should we hope that we can find a simple approximation?

- Prior distribution is complicated
 - Need to describe all possible states of the world (and relationships between variables)
- Posterior distribution is often simple:
 - Have made many observations → less uncertainty
 - Variables can become “almost independent”
- For now: Represent posterior as undirected model (and instantiate observations)



$$P(X_1, \dots, X_n \mid obs) = \frac{1}{Z} \prod_j \Psi_j(\mathbf{C}_j)$$



Variational approximation

- **Key idea:** Approximate posterior with simpler distribution that's as close as possible to P
 - What is a “simple” distribution?
 - What does “as close as possible” mean?
- **Simple** = efficient inference
 - Typically: factorized (fully independent, chain, tree, ...)
 - Gaussian approximation
- **As close as possible** = KL divergence (typically)
 - Other distance measures can be used too, but more challenging to compute

Kullback-Leibler (KL) divergence

- Distance between distributions

$$D(P||Q) = \int P(\mathbf{x}) \log \frac{P(\mathbf{x})}{Q(\mathbf{x})} d\mathbf{x}$$

- Properties:

- $D(P || Q) \geq 0$

- $P(x)=Q(x)$ almost everywhere $\Leftrightarrow D(P || Q) = 0$

- In general, $D(P || Q) \neq D(Q || P)$

- P determines when difference is important

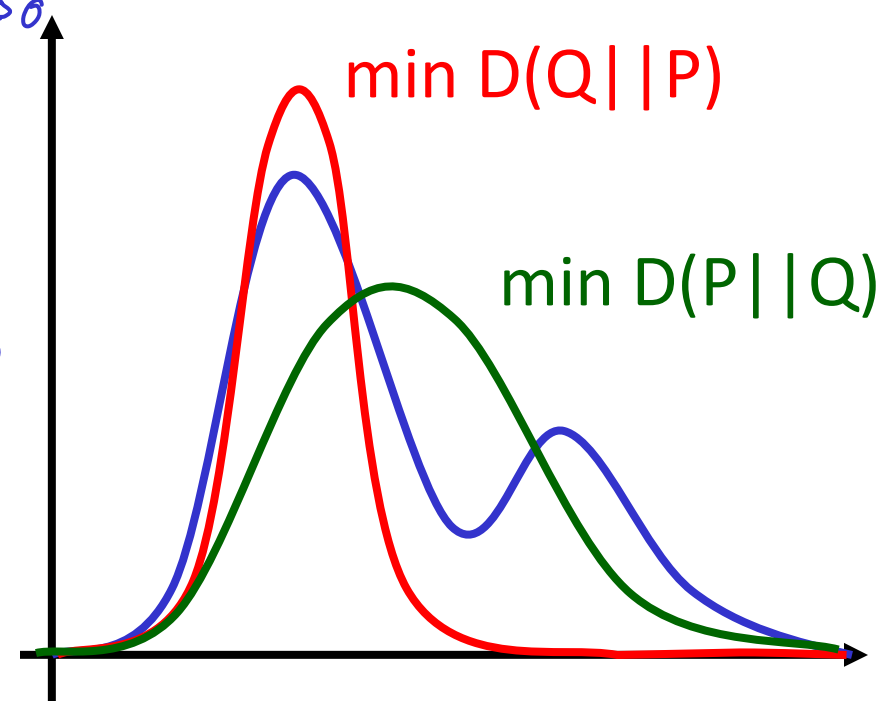
$$P(x) = 0, Q(x) \neq 0 \Rightarrow 0 \cdot \log \frac{0}{\varepsilon} = 0$$

$$P(x) = \varepsilon, Q(x) = 0 \Rightarrow \varepsilon \cdot \log \frac{\varepsilon}{0} = \infty$$

Finding simple approximate distributions

- KL divergence not symmetric; need to choose directions
- P : true distribution; Q : our approximation
- $D(P || Q)$
 - The “right” way
 - Q chosen to “support” P
 - Often intractable to compute
- $D(Q || P)$
 - The “reverse” way
 - Underestimates support (overconfident)
 - Will be tractable to compute
- Both special cases of α -divergence

$$P(x) > 0 \Rightarrow Q(x) > 0$$



“Simple” distributions

- Simplest distribution: Q fully factorized

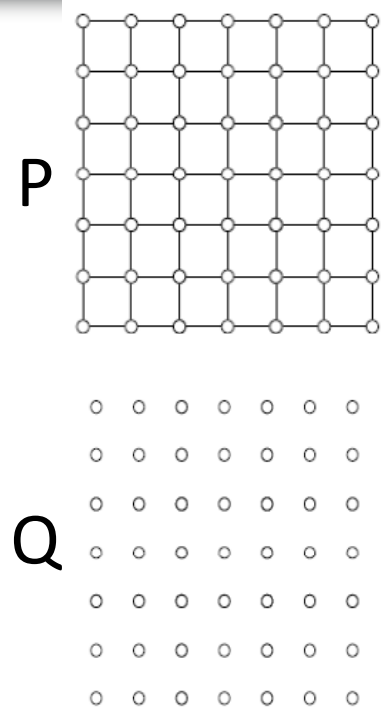
- $Q(X_1, \dots, X_n) = \prod_i Q_i(X_i)$

- \mathcal{M} = { Q : Q fully factorized}
= { Q : $Q(X) = \prod_i Q_i(X_i)$ }

$$Q^* = \operatorname{argmin}_{Q \in \mathcal{M}} D(Q || P)$$

$D(P || Q)$

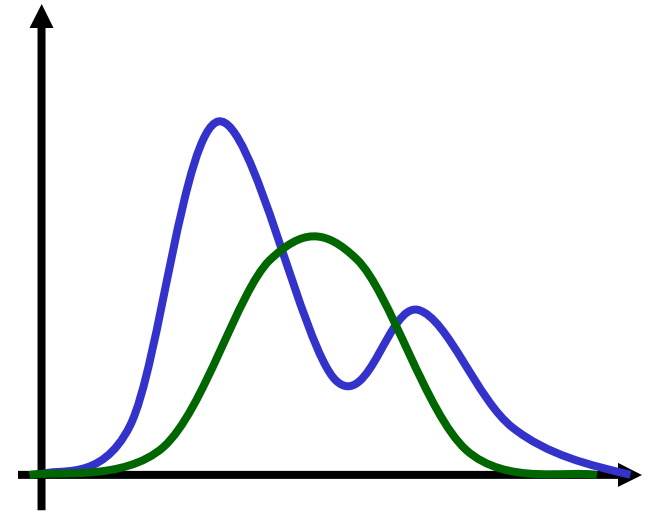
- Can also find more structured approximations
 - Chains: $Q(X_1, \dots, X_n) = \prod_i Q_i(X_i | X_{i=1})$
 - Trees
 - Any distributions one can do efficient inference on



Mean field approximation the “right way”

$$\begin{aligned} D(P \parallel Q) &= \sum_x P(x) \log \frac{P(x)}{Q(x)} \\ &= \underbrace{\sum_x P(x) \log P(x)}_{\text{constant}} - \underbrace{\sum_x P(x) \log Q(x)}_{(*)} \end{aligned}$$

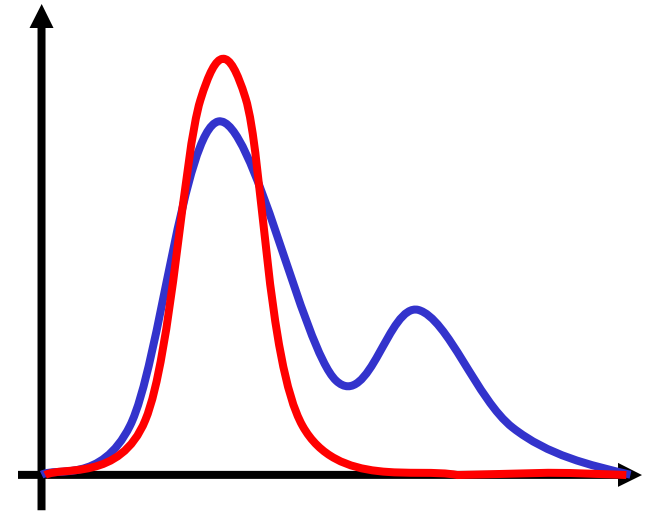
$$\begin{aligned} (*) &= \sum_x P(x) \log \prod_i Q_i(x_i) \\ &= \sum_x P(x) \sum_i \log Q_i(x_i) \\ &= \sum_i \sum_x P(x) \log Q_i(x_i) \\ &= \sum_i \left(\underbrace{\sum_{x_i} P(x_i) \log Q_i(x_i)}_{\text{Need } P(x_i) \text{ intractable!}} \right) \underbrace{\left(\sum_{x_{1:i-1}, i+1:m} P(x) \right)}_1 \end{aligned}$$



Mean field approximation the reverse way

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

$$= \underbrace{\sum_x Q(x) \log Q(x)}_{(1)} - \underbrace{\sum_x Q(x) \log P(x)}_{(2)}$$



$$(1) = -H(Q) = \sum_x Q(x) \log \prod_i Q_i(x_i)$$

$$= \sum_i \sum_x Q(x) \log Q_i(x_i)$$

$$= \sum_i \sum_{x_i} Q_i(x_i) \log Q_i(x_i) \left(\sum_{x_{-i}} Q(x_{-i}) \right) = - \sum_i H(Q_i)$$

$$(2) = \sum_x Q(x) \log \frac{1}{Z} \prod_i \psi_i(C_i) = \underbrace{- \sum_x Q(x) \log Z}_{-\log Z} + \underbrace{\sum_x Q(x) \sum_i \log \psi_i(C_i)}_{(*)}$$

$$(*) = \sum_i \sum_x Q(x) \log \underbrace{\psi_i(C_i)}_{\text{dep on small \# of vars}} = \sum_i \sum_{C_i} \underbrace{Q(C_i)}_{\text{Need marginal for } C_i} \log \psi_i(C_i)$$

Reverse KL for fully factorized case

$$D(Q||P) = - \sum_i \underbrace{\sum_x Q(x) \log \Psi_i(x)}_{\mathbb{E}_Q[\log \Psi_i]} - \sum_i H(Q_i) + \underbrace{\ln Z}_{\text{constant}}$$

$$\underbrace{\ln Z}_{\text{constant}} = \underbrace{D(Q||P)}_{\rightarrow \min} + \underbrace{\sum_i H(Q_i) + \sum_i \mathbb{E}_Q[\log \Psi_i]}_{\rightarrow \max}$$

$F(Q; P)$

KL and the partition function

Suppose $P(X_1, \dots, X_n) = Z^{-1} \prod_i \Psi_i(C_i)$ is Markov Network

Theorem: *For any Q (not necessarily fully factorized)*

$$\ln Z = \uparrow \underline{F[P; Q]} + \downarrow D(Q || P)$$

Hereby, $F[P; Q]$ is the following energy functional

$$F[P; Q] = \sum_i \mathbb{E}_Q[\ln \Psi_i] + H(Q)$$

Reverse KL vs. log-partition function

$$\ln Z = \underbrace{F[P; Q]}_{\geq 0} + \underbrace{D(Q||P)}_{\geq 0} \quad F[P; Q] = \sum_i \mathbb{E}[\ln \Psi_i] + H(Q)$$

Maximizing energy functional \Leftrightarrow Minimizing reverse KL

Corollary:

Energy function is lower bound on log partition function

$$P(x) = \frac{1}{Z} \prod_i \psi_i(x_i)$$

$$P(x) \leq \frac{1}{F(P; Q)} \prod_i \psi_i(x_i)$$

Implies upper bound on ~~event~~ probabilities!

Optimizing for mean field approximation

- Want to solve $\max_Q F[P; Q] = \max_Q \sum_j \mathbb{E}_Q[\ln \Psi_j] + \sum_i H(Q_i)$

$$\text{s.t. } \sum_{x_i} Q_i(x_i) = 1$$

- Solved via Lagrange multipliers: *there exist $\lambda_1 \dots \lambda_n$ s.t. optimization of (*) is equivalent to*

$$\max_Q \sum_j \mathbb{E}_Q[\ln \Psi_j] + \sum_j H(Q_j) + \sum_j \lambda_j \left[\sum_{x_j} Q_j(x_j) - 1 \right]$$

Differentiate and set to 0!

Minimum, Maximum or saddle point

Theorem: Q stationary point iff for each i and x_i :

$$Q_i(x_i) = \frac{1}{Z_i} \exp\left(\sum_j \mathbb{E}[\ln \Psi_j \mid x_i]\right)$$

Fixed point iteration for MF

- Initialize factors $Q^{(0)}_i$ arbitrarily; $t=0$
- Until converged, do
 - $t \leftarrow t+1$
 - For each variable i and each assignment x_i do

$$Q_i(x_i)^{(t+1)} = \frac{1}{Z_i} \exp\left(\sum_j \underbrace{\mathbb{E}_{Q^{(t)}} [\ln \Psi_j \mid x_i]}\right)$$

$$Z_i = \sum_{x_i} Q_i(x_i)$$

- Guaranteed to converge! 😊
- Gives both approx. distribution Q and lower bound on $\ln Z$
- Can get stuck in local optimum ☹️

Computing updates

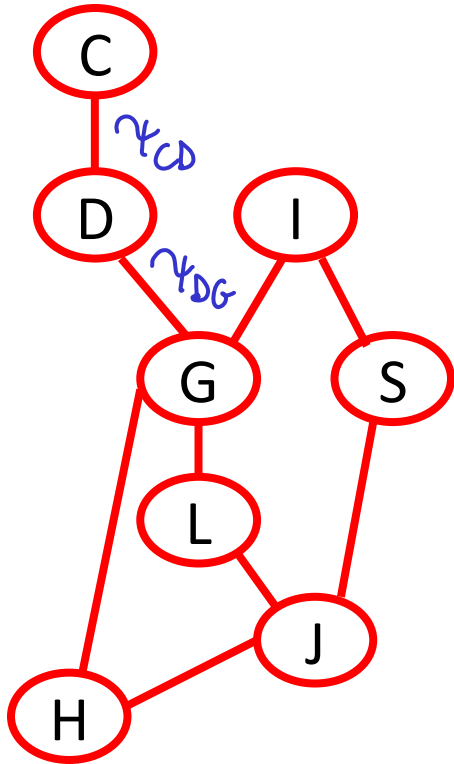
Need to compute

$$Q_i(x_i)^{(t+1)} = \frac{1}{Z_i} \exp\left(\sum_j \underbrace{\mathbb{E}_{Q^{(t)}}[\ln \Psi_j \mid x_i]}_{\text{Must compute expected log potentials: } (*)}$$

Must compute expected log potentials: $(*) = \mathbb{E}_Q[\ln \Psi_j \mid x_i]$

$$(*) = \sum_x Q^{(t)}(x \mid x_i) \ln \underbrace{\Psi_j(x)}_{=\Psi_j(c_j)} = \sum_{c_j \sim x_i} \underbrace{Q^{(t)}(c_j \mid x_i)}_{=\prod_{k \in C_j \setminus \{i\}} Q_k(x_k)} \ln \Psi_j(c_j)$$

Example iteration



$$\mathbb{E}_Q[\ln \psi_{DG} | X_D=1]$$

$$= \sum_{X_G} \underbrace{Q(X_G | X_D=1)}_{= Q_G(X_G)} \ln \psi_{DG}(X_D=1, X_G)$$

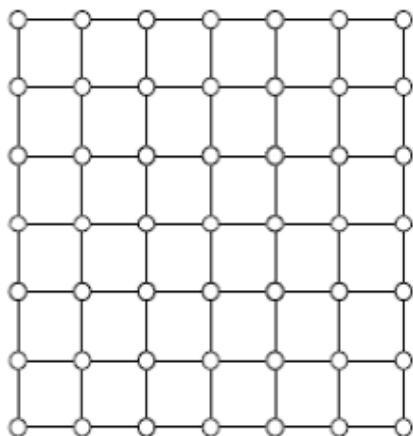
$$Q_D(X_D=1) = \frac{1}{Z_D} \exp(\mathbb{E}_Q[\ln \psi_{DG} | X_D=1] + \mathbb{E}_Q[\ln \psi_{CD} | X_D=1])$$

Structured mean field

Goal of variational inference:

Approximate complex distribution by simple distribution

True dist.



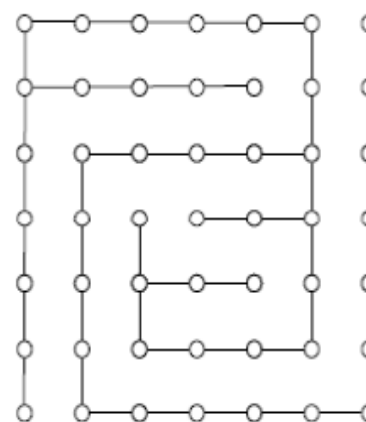
$$p(x) \propto \prod_c \phi_c(x_c)$$

Fully-factorized mean field



$$q(x) \propto \prod_i q_i(x_i)$$

Structured mean field



$$q(x) \propto q_A(x_A) q_B(x_B)$$

Structured mean-field approximations

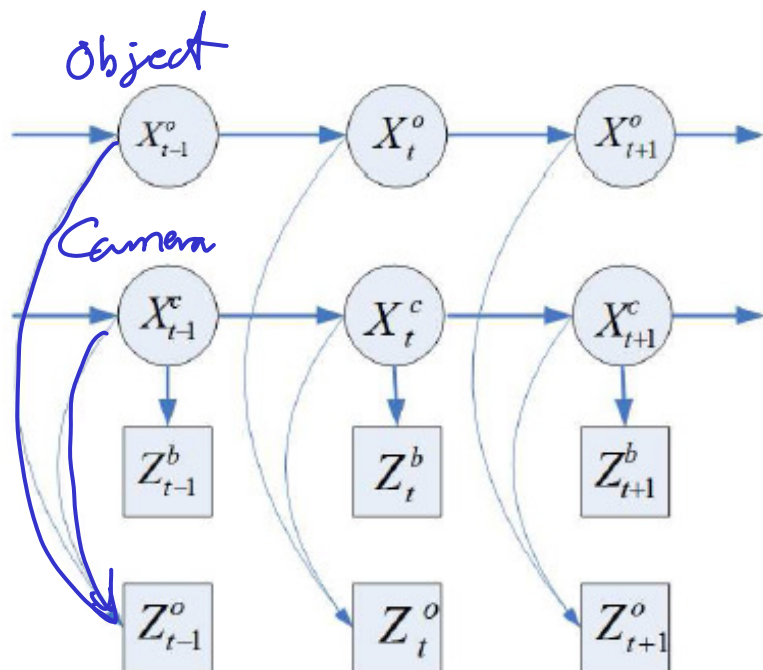
- Can get better approximations using **structured approximations**:

$$\max_{Q \in \mathcal{M}} F[P; Q] = \max_{Q \in \mathcal{M}} \sum_j \mathbb{E}_Q[\ln \Psi_j] + H(Q)$$

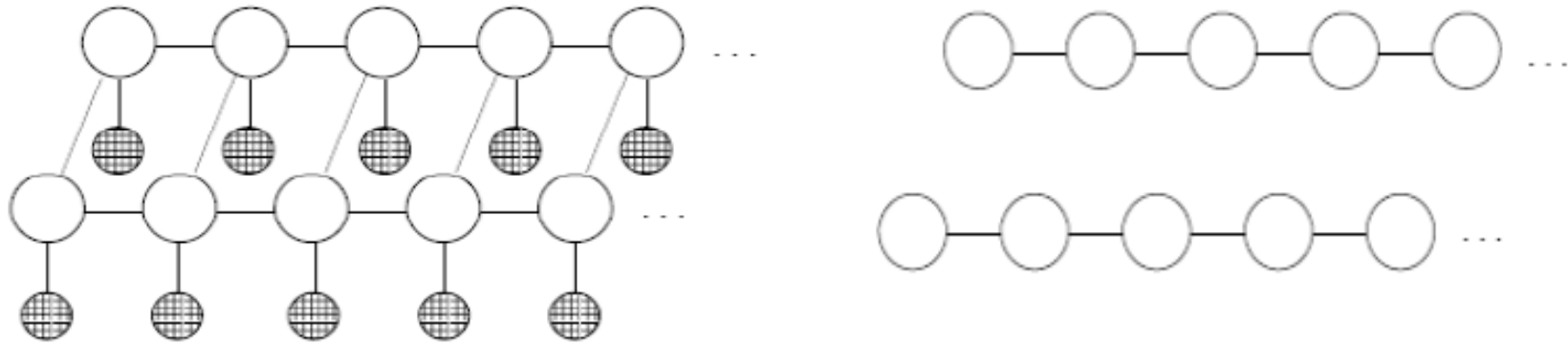
- Only need to be able to compute energy functional
- Can do whenever we can perform efficient inference in Q (e.g., chains, trees, low-treewidth models)
 - Update equations look similar as for fully-factorized case (see reading)

Example: Factorial HMM

- Simultaneous tracking and camera registration
- State space decomposed into object location and camera parameters



Variational approximations for FHMMs



$$\max_{Q \in \mathcal{M}} F[P; Q] = \max_{Q \in \mathcal{M}} \sum_j \mathbb{E}_Q[\ln \Psi_j] + H(Q)$$

- Approximate posterior by independent chains

$$\mathcal{M} = \left\{ Q : Q(\mathbf{X}) = \prod_c \prod_t Q_{c,t}(X_{c,t} | X_{c,t-1}) \right\}$$

Summary: Variational inference

- Approximate complex (intractable) distribution by simpler distribution that is “as close as possible”
- **Simple** = tractable (efficient inference)
- **Closeness** = Reverse KL (efficient to compute)
- Interpretation: Optimize **lower bound** on the log-partition function
 - Implies upper bound on event probabilities
- Efficient algorithm that's guaranteed to converge (in contrast to Loopy BP..), but possibly to local optimum

Approximate inference

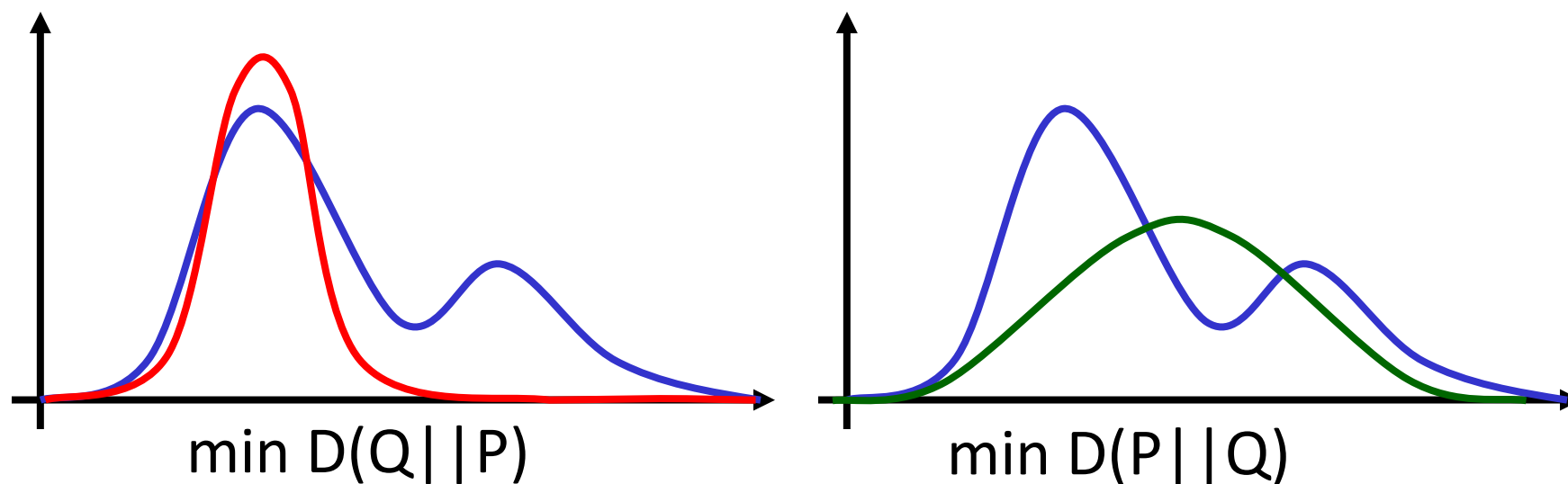
- Three major classes of general-purpose approaches
- **Message passing**
 - E.g.: Loopy Belief Propagation (today!)
- **Inference as optimization**
 - Approximate posterior distribution by simple distribution
 - Mean field / structured mean field
 - Assumed density filtering / expectation propagation
- **Sampling based inference**
 - Importance sampling, particle filtering
 - Gibbs sampling, MCMC
- Many other alternatives (often for special cases)

KL-divergence the “right” way:

- Find distribution $Q^* \in \mathcal{M}$:

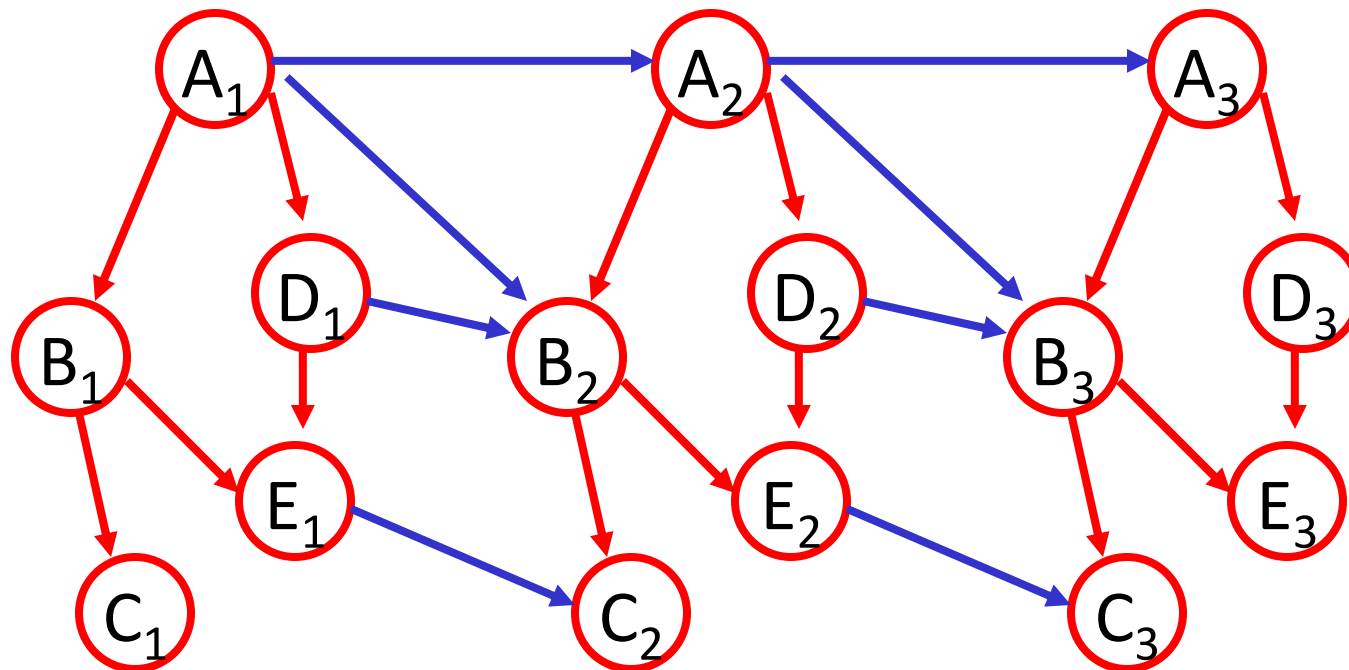
$$Q^* = \operatorname{argmin}_{Q \in \mathcal{M}} D(P || Q)$$

- In some applications, can compute $D(P || Q)$
 - Important example: Assumed density filtering in DBNs



Recall: Dynamic Bayesian Networks

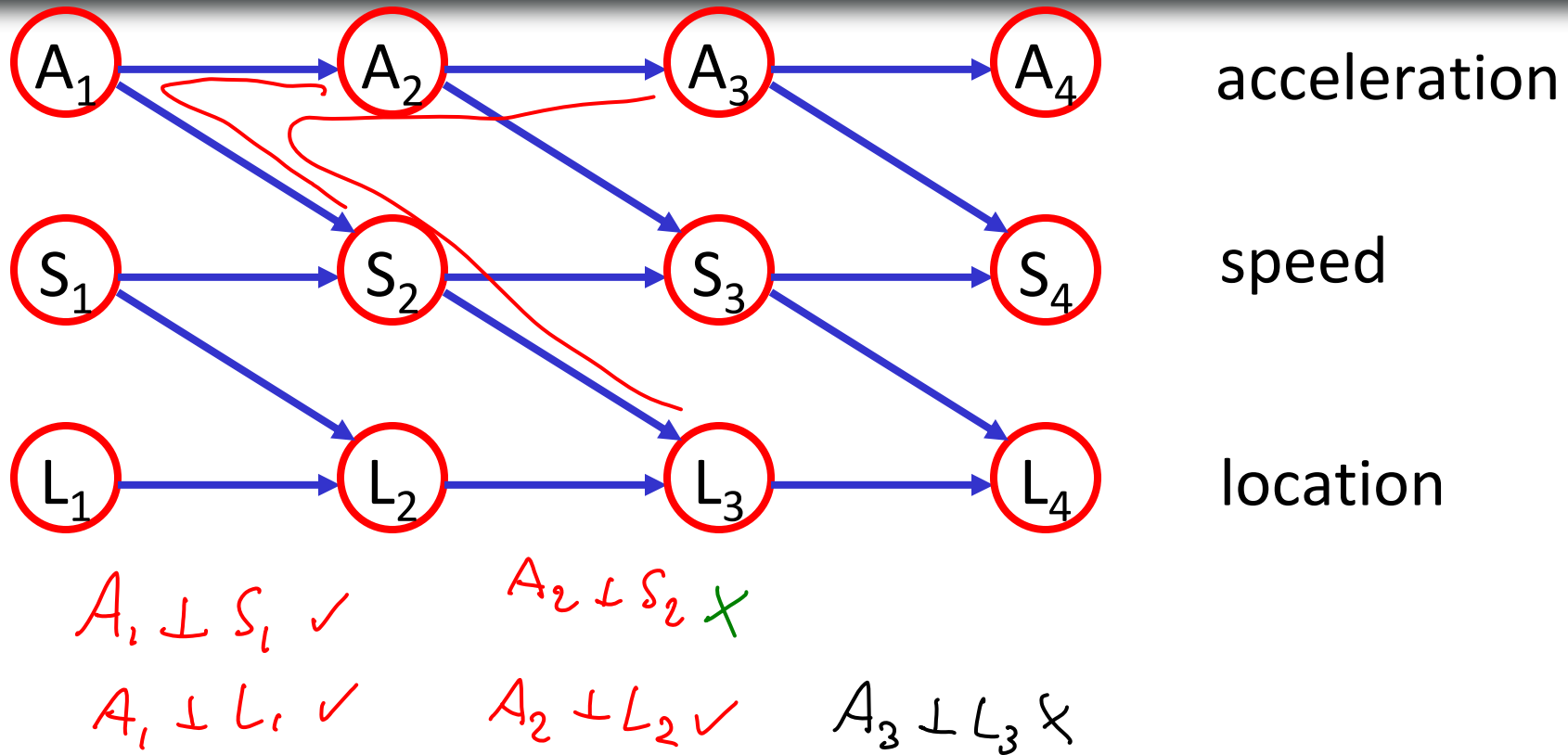
- At every timestep have a Bayesian Network



$$S_t = \{A_t, B_t, \dots, E_t\}$$

- Variables at each time step t called a “slice” S_t
- “Temporal” edges connecting S_{t+1} with S_t

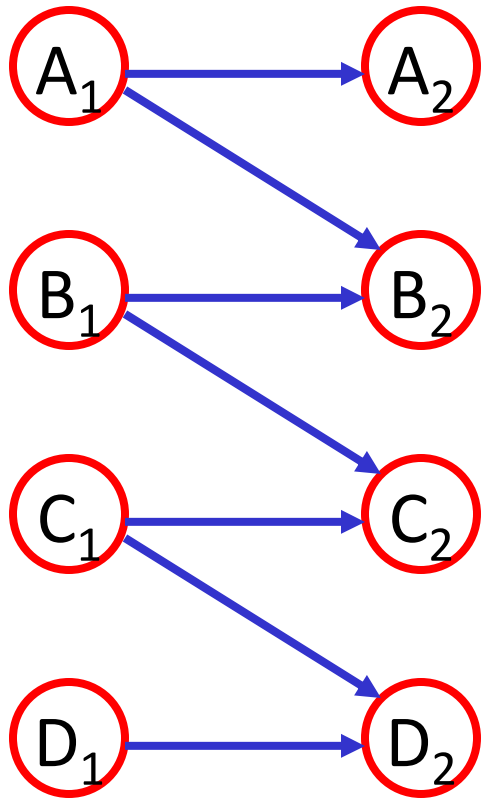
Flow of influence in DBNs



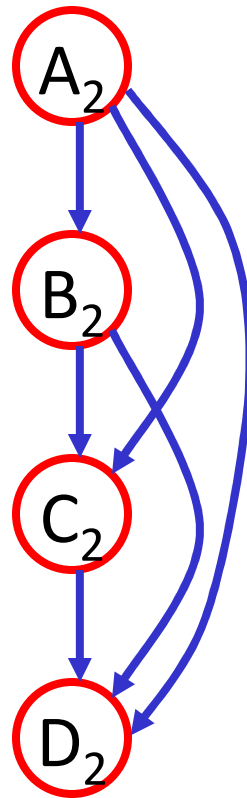
- Can we do efficient filtering in BNs?

Approximate inference in DBNs?

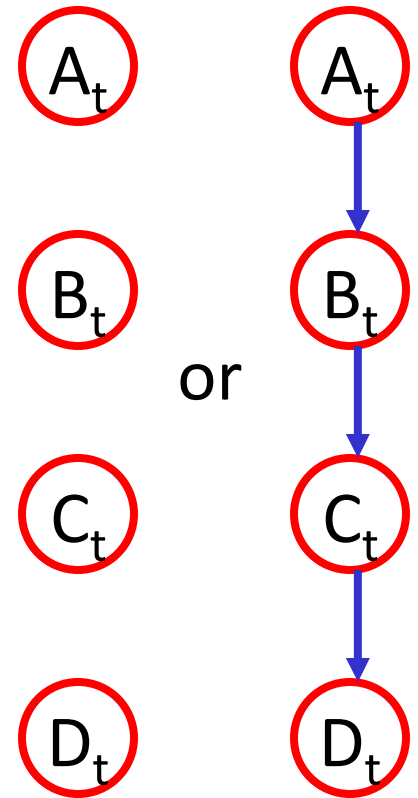
DBN



Marginals

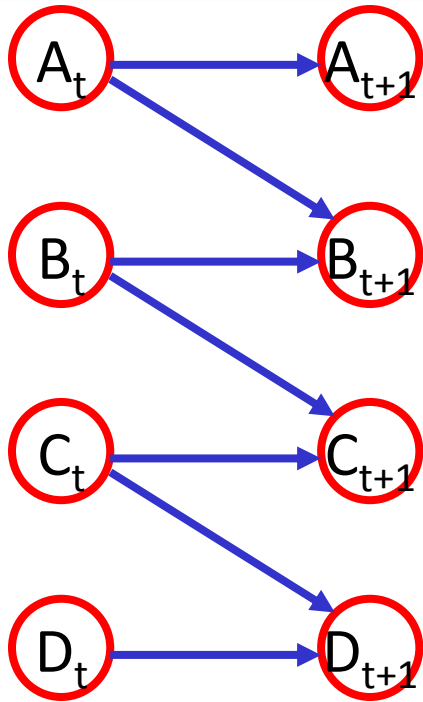


Approx. marginals



Want to find **tractable** approximation to marginals
that's **as close** to true marginals as possible

Assumed Density Filtering

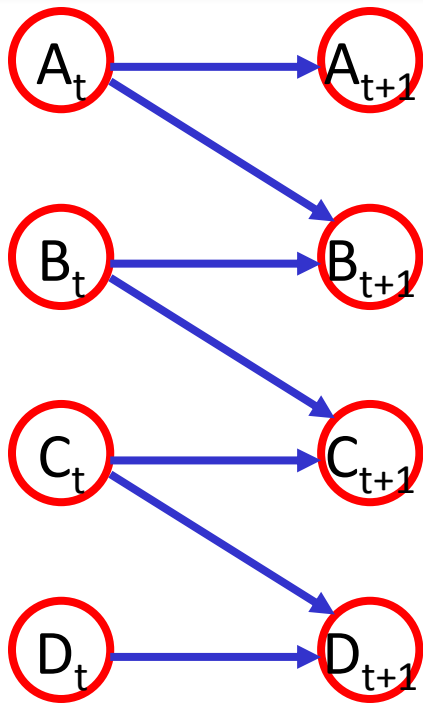


$$\begin{aligned}
 D(P(S_{t+1}) \parallel Q(S_{t+1})) &= \sum_{S_{t+1}} P(S_{t+1}) \log \frac{P(S_{t+1})}{Q(S_{t+1})} \\
 &= \text{const} - \underbrace{\sum_{S_{t+1}} P(S_{t+1}) \log Q(S_{t+1})}_{(*)}
 \end{aligned}$$

- Assume distribution $P(\mathbf{S}_t)$ for slice t factorizes
- $P(\mathbf{S}_{t+1})$ is fully connected ☹️
- Want to compute best-approximation Q^* for $P(\mathbf{S}_{t+1})$

$$Q^* = \operatorname{argmin} D(P \parallel Q)$$

Assumed Density Filtering



$$\sum_{s_{t+1}} P(s_{t+1}) \log \underbrace{Q(s_{t+1})}_{= \prod_i Q_i(s_{i,t+1})}$$

$$= \sum_i \sum_{s_{t+1}} P(s_{t+1}) \log Q_i(s_{i,t+1})$$

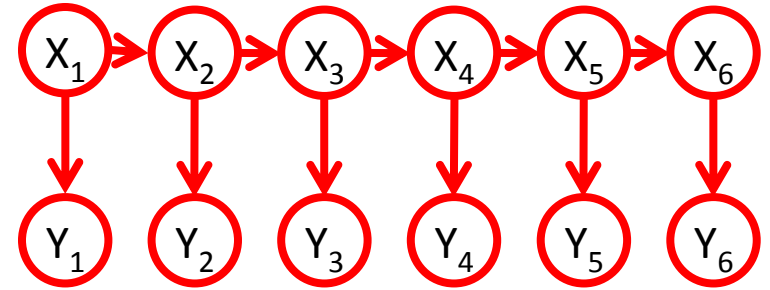
$$= \sum_i \sum_{s_{i,t+1}} \underbrace{P(s_{i,t+1})}_{\text{E.g., } P(A_{t+1})} \log Q_i(s_{i,t+1})$$

Can compute expectations efficiently

Get optimal Q^* by setting $Q_i^*(s_{i,t+1}) = P(s_{i,t+1})$

Recall: Bayesian filtering

- Start with $P(X_1)$
- At time t
 - Assume we have $P(X_t | y_{1..t-1})$
 - Condition: $P(X_t | y_{1..t})$



$$P(X_t | y_{1..t}) \propto P(X_t | y_{1..t-1}) \underbrace{P(y_t | X_t, y_{1..t-1})}_{\text{cond. ind. } P(y_t | X_t)}$$

- Prediction: $P(X_{t+1}, X_t | y_{1..t})$

$$P(X_{t+1}, X_t | y_{1..t}) = P(X_t | y_{1..t}) \cdot \underbrace{P(X_{t+1} | X_t, y_{1..t})}_{= P(X_{t+1} | X_t)}$$

- Marginalization: $P(X_{t+1} | y_{1..t})$

$$P(X_{t+1} | y_{1..t}) = \sum_{X_t} P(X_{t+1}, X_t | y_{1..t})$$

Assumed Density Filtering

- Start with $P(\mathbf{S}_1)$
- At every time step t : tractable approximation Q_t
 $Q_t(\mathbf{S}_t) \approx P(\mathbf{S}_t \mid \mathbf{O}_{1:t-1})$
- **Condition** on observation $\mathbf{O}_t \subseteq \mathbf{S}_t$: $Q_t(\mathbf{S}_t \mid \mathbf{O}_t)$
- **Predict**: multiply transition model to get $Q_t(\mathbf{S}_{t+1}, \mathbf{S}_t \mid \mathbf{O}_t)$
$$Q_t(\mathbf{S}_{t+1}, \mathbf{S}_t \mid \mathbf{O}_t) = Q_t(\mathbf{S}_t \mid \mathbf{O}_t) P(\mathbf{S}_{t+1} \mid \mathbf{S}_t)$$
- **Marginalize \mathbf{S}_t**
 - This is intractable (connects all variables in \mathbf{S}_{t+1})
 - Approximate $Q_t(\mathbf{S}_{t+1} \mid \mathbf{O}_t)$ by Q^* s.t.
 $Q^* = \operatorname{argmin}_Q D(Q_t(\mathbf{S}_{t+1}) \parallel Q(\mathbf{S}_{t+1}))$
 - This is done by matching moments:
for discrete models, ensure that $Q_{t+1}(\mathbf{s}_{t+1}) = Q_t(\mathbf{s}_{t+1} \mid \mathbf{o}_t)$

Summary of Assumed Density Filtering

- Variational inference technique for dynamical Bayesian Networks
- Find tractable approximation for each time slice that minimizes KL divergence (in the “right” way)
- Can show that errors don’t add up too much
- Examples:
 - Tractable inference in DBNs
 - Unscented Kalman Filter

Summary: Inference as optimization

- Approximate intractable distribution by a **tractable** one
- **Optimize parameters** of the distribution to make approximation as tight as possible
- Common distance measure: KL-divergence (both ways)
 - Special case of α -divergence
- Can get upper bounds on event probabilities, etc.