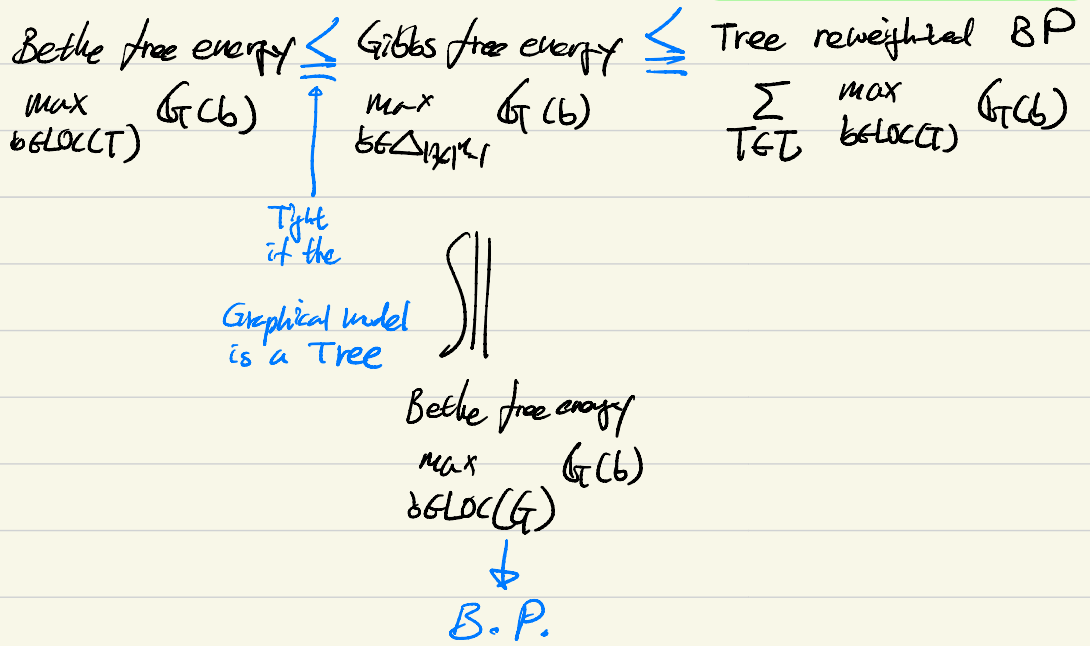
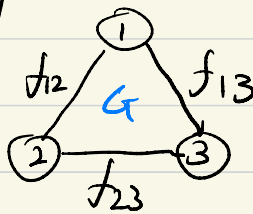


*Recap

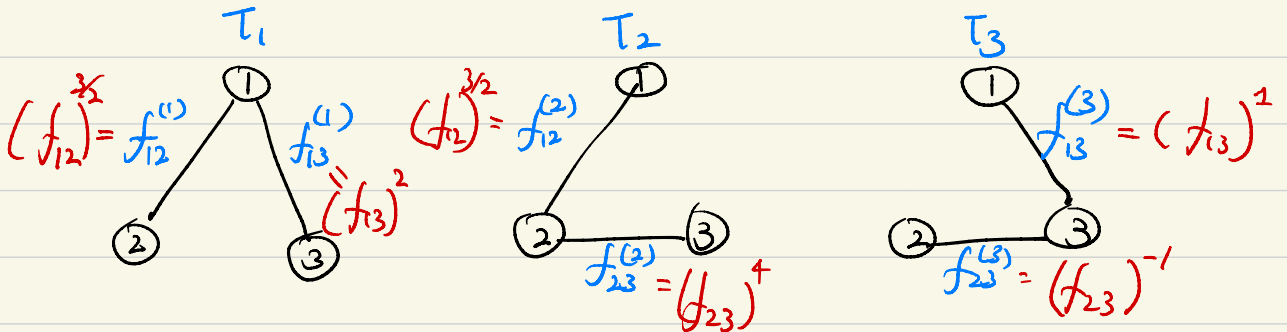


Tree reweighted belief propagation [more details on slides
 "A new class of upper bounds on the log partition function." 2005
 Wainwright, Jaakkola, Willsky

example \rightarrow given Graphical Model



Consider all spanning trees on G , and assign weights on these trees.



Rule #1: $\sum_k C_k = 1$

$C_1 = \frac{1}{3}$

$C_2 = \frac{1}{3}$

$C_3 = \frac{1}{3}$

Rule #2: for each $(i,j) \in E$ $\prod_k \left(f_{ij}^{(k)} \right)^{C_k} = f_{ij}$

Def. Tree reweighted belief propagation.

given $\mathcal{T} = \{T_k\}_{k=1}^K$ set of spanning trees, and corresponding weights $\{c_k\}$ s.t.

$$\sum c_k = 1$$

$$\log f_{ij}(x_i, x_j) = \sum_k c_k \log f_{ij}^{(k)}(x_i, x_j)$$

then the energy term in Gibbs free energy decomposes as

$$\begin{aligned} \mathbb{E}_b \left[-\sum_{i,j \in E} \log f_{ij}(x_i, x_j) \right] &= \mathbb{E}_b \left[-\sum_{(i,j) \in E} \sum_k c_k \log f_{ij}^{(k)}(x_i, x_j) \right] \\ &= \sum_k c_k \underbrace{\mathbb{E}_b \left[-\sum_{(i,j) \in E} \log f_{ij}^{(k)}(x_i, x_j) \right]}_{\text{energy of a one tree model.}} \end{aligned}$$

* Claim: $\log Z \leq \sum_k c_k \log Z_k$
 easy to compute with B.P. as it is a tree.

proof:

$$\log Z = \max_{b \in \Delta_{|\mathcal{X}|^2-1}} \mathbb{E}_b \left[\sum_{(i,j) \in E} \log f_{ij}(x_i, x_j) \right] + \text{Entropy}(b)$$

$$\sum c_k = 1 \rightarrow = \max_{b \in \Delta_{|\mathcal{X}|^2-1}} \sum_k c_k \left\{ \mathbb{E}_b \left[\sum_{(i,j) \in E} \log f_{ij}^{(k)}(x_i, x_j) \right] + \text{Entropy}(b) \right\}$$

$$\text{exchange sum \& max} \rightarrow \leq \sum_k c_k \max_{b \in \Delta_{|\mathcal{X}|^2-1}} \left\{ \mathbb{E}_b \left[\sum_{(i,j) \in E} \log f_{ij}^{(k)}(x_i, x_j) \right] + \text{Entropy}(b) \right\}$$

$$\text{graphical model is on tree} \rightarrow = \sum_k c_k \max_{b \in \text{LOC}(T_k)} \left\{ \mathbb{E}_b \left[\sum_{(i,j) \in E} \log f_{ij}^{(k)}(x_i, x_j) \right] + \text{Entropy}(b) \right\}$$

this can be solved exactly with B.P.

• Caveat $\left[\right.$ we want to select $c_k, f_{ij}^{(k)}$ that minimize $\sum c_k \log Z_k$
 $\left. \right]$ # of spanning trees can explode

* To solve inference problems $P(x)$

Variational Methods

↓
Belief Propagation

- Deterministic
- fast
- approximation

sampling

↓
Gibbs Sampling

- Randomized
- slower

• exact in the limit $N \rightarrow \infty$, but difficult to decide when to stop.

$$\hat{P}(X_i=a) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(X_i=a)$$

Def: Markov Chain Monte Carlo methods.

- construct a Markov chain $X^{(t)} \in \mathcal{X}^n \rightarrow X^{(t+1)} \in \mathcal{X}^n$ with transition matrix Q whose stationary distribution $P(x)$
- start with an arbitrary realization $X^{(0)}$ and run the Markov Chain until it converges close to its stationary distribution
- Repeat.

Q1: How do we construct Q ? → Metropolis-Hastings algorithm → Gibbs sampling

Q2: How long does it take for the Markov chain to converge?

- spectral analysis.
- path coupling.

Strategy

given a graphical model

$$G \Rightarrow P(x) = \frac{1}{Z} \prod_j \psi_j(x_i, x_j)$$

state $X^{(t)} \in \mathcal{X}^n$: realization

initial state $X^{(0)}$

construct transition matrix $Q \in \mathbb{R}^{|\mathcal{X}|^n \times |\mathcal{X}|^n}$

we will not write it in a memory, but conceptually define it sparse

Repeat: sample: $X^{(t+1)} \sim (X^{(t)})^T \cdot Q$



eventually $X^{(t)} \sim \pi$: stationary distribution of Q .
"P(x)"

Def. time-homogeneous finite-state Markov Chains.

Markov chain $\left\{ \begin{array}{l} \text{State space } \mathcal{X}^n \\ \text{Transition matrix } Q \in \mathbb{R}^{|\mathcal{X}^n| \times |\mathcal{X}^n|} \end{array} \right.$

$$Q_{xy} = \mathbb{P}(X_{t+1}=y \mid X_t=x)$$

Def. stationary distribution π of M.C. $\pi = \left\{ \right\} |\mathcal{X}^n|$

$$\pi^T \cdot Q = \pi^T$$

- might not be unique
 - might not exist
-) we do not put into detailed conditions on Q that ensures existence & uniqueness, but construct Q s.t. it is U&E.

Def. a Markov chain is Reversible if $\exists \pi$ s.t. detailed balance equation is satisfied

$$\pi_x \cdot Q_{xy} = \pi_y \cdot Q_{yx} \quad \text{for all } x, y. \quad (*)$$

Claim. π satisfying (*) is a stationary distribution of Q .

proof >

$$(\pi^T Q)_y = \sum_x \pi_x Q_{xy} \stackrel{(*)}{=} \sum_x \pi_y Q_{yx} \stackrel{Q \text{ is stochastic, i.e. row-sum to one}}{=} \pi_y.$$

$$\Rightarrow \pi^T Q = \pi^T$$

we will construct a Markov Chain Q that is reversible,

$\left\{ \begin{array}{l} \text{we can use } (*) \text{ to ensure } \pi = P(x) \\ \text{spectral analysis can be applied.} \end{array} \right.$

Def. Metropolis-Hastings Algorithm

- start with a candidate transition matrix K
- To ensure unique stationary distribution, it is sufficient to have
 - $K_{xx} > 0$, $\forall x \in \mathcal{X}^n$ [aperiodic]
 - for any $x, y \in \mathcal{X}^n$, \exists a path $(x_1=x, x_2, \dots, x_m=y)$ of positive probability transitions [irreducibility]

$$K_{x_i x_{i+1}} > 0, \forall i \in \{1, \dots, m-1\}$$

what we want

$$(*) \quad Q_{xy} \cdot P(x) = Q_{yx} \cdot P(y)$$

what we have

$$K_{xy} P(x) \underset{\substack{\uparrow \\ \text{w.l.o.g.}}}{>} K_{yx} P(y)$$

the main trick is to remove some "probability mass" from the larger one.

$$\text{Define: } R_{xy} \triangleq \min \left\{ 1, \frac{P(y) K_{yx}}{P(x) K_{xy}} \right\} \quad \text{for each } x, y$$

$$Q_{xy} = K_{xy} \cdot R_{xy}$$

$$Q_{xx} \underset{\substack{\uparrow \\ \text{s.t.} \\ \text{probability sums to one}}}{=} 1 - \sum_{y \neq x} Q_{xy}$$

with this rejection sampling ($\cdot R_{xy}$)

claim: $(Q, P(x))$ satisfy $(*)$.

$$\begin{aligned} \text{proof: } P(x) \cdot Q_{xy} &= P(x) \cdot R_{xy} \cdot K_{xy} \\ &= P(y) K_{yx} && \text{(if } R_{xy} \leq 1) \\ &= P(y) R_{yx} \cdot K_{yx} && \text{if } R_{yx} = 1 \end{aligned}$$

* Remark: as $P(x)$ is only used in $R_{xy} = \min\left\{1, \frac{P(y)K_{yx}}{P(x)K_{xy}}\right\}$

we only need $\frac{P(y)}{P(x)} = \prod_{(i,j) \in E} \frac{f_{ij}(y_i, y_j)}{f_{ij}(x_i, x_j)}$ } takes $O(|E|)$ computation
does not require Σ .

* Remark: Do we need to store K and $Q \in \mathbb{R}^{|\mathcal{X}|^n \times |\mathcal{X}|^n}$

- We can choose K to be simple, such as $K = \frac{1}{|\mathcal{X}|^n} \mathbb{1}\mathbb{1}^T$

Step 1 - At time t first generate candidate $x^{t+1/2}$ from $K(x^t, x^{t+1/2})$

Step 2 - Accept $x^{t+1/2}$ with probability $R_{x^t, x^{t+1/2}}$: $x^{t+1} = x^{t+1/2}$

Step 3 - Otherwise reject and keep current state : $x^{t+1} = x^t$

* Theorem: Metropolis-Hastings Algorithm finds L_1 -projection of K onto the space of reversible Markov chains with stationary distribution $P(x)$.

$$Q_{MH} = \arg \min_{\tilde{Q} \text{ s.t. } P^T \tilde{Q} = P^T} \sum_x \sum_{y \neq x} |P(x) \cdot K_{xy} - P(y) \tilde{Q}_{xy}|$$

* Are we done? - the art is in choosing K .

if the spread of K is too large, then acceptance rate is low

if the spread of K is too narrow, then mixing time can be large

example $\triangleright K = \frac{1}{|\mathcal{X}|^n} \mathbb{1}\mathbb{1}^T$, $R_{xy} = \min\left(1, \prod_{(i,j) \in E} \frac{f_{ij}(y_i, y_j)}{f_{ij}(x_i, x_j)}\right)$

all pairs first sampled with equal probability (as per K)

- but many candidates might be unlikely and be rejected.

* Def. Gibbs Sampling.

Q {
Step 1. sample $i \in \{1, \dots, n\}$ uniformly at random.
Step 2. set $y_{-i} = X_{-i}^{(t)}$
 $\stackrel{\text{[n] \setminus \{i\}}}{\text{[n] \setminus \{i\}}}$
Step 3. sample y_i from $P(y_i | X_{-i}^{(t)})$

* Remark. $P(y_i | X_{-i}^{(t)}) \propto \prod_{j \in \Theta_i} f_{ij}(y_i, X_j^{(t)})$ is efficient

* claim. $(Q, P(x))$ satisfy (*).

proof for y that differ at only i -th coordinate from x ,

$$P(x) \cdot Q_{xy} = P(x) \cdot \frac{1}{n} \cdot P(y_i | X_{-i})$$

$$\text{Bayes} \rightarrow = P(x_i | X_{-i}) P(x_{-i}) \cdot \frac{1}{n} P(y_i | X_{-i})$$

$$= \underbrace{P(x_i | X_{-i})}_{Q_{yx}} \cdot \frac{1}{n} \underbrace{P(x_{-i}) P(y_i | X_{-i})}_{P(y)}$$

otherwise $Q_{xy} = 0$ if x & y differ more than one coordinate.

* the resulting dynamics of the Markov chain is called **Glauber Dynamics**.