

\* Equivalence of inference tasks.

• Suppose we have a black-box tool that computes the partition function  $Z_G$  of any graphical model  $G$ .

then we can compute Marginal  $P(X_i=0) \propto Z_{G(X_i=0)}$

$\cong$  is the graphical model conditioned on  $x_i$  being 0.

\* We will focus on the task of computing  $Z$  given a graphical model.

for

$$P(x) = \frac{1}{Z} \cdot \underbrace{\prod_{(i,j) \in E} f_{ij}(x_i, x_j)}_{f^{\text{total}}(x)}$$

• We are given  $f^{\text{total}}(x)$  but not  $Z$ .

• We focus on computing (approximately) the **log partition function**

$$\Phi \cong \log Z = \log \left\{ \sum_{x \in \mathcal{X}^n} f^{\text{total}}(x) \right\}$$

• Def. Variational characterization of log partition function.

is to represent as a solution of an optimization problem.

$$\Phi = \sup Q(b)$$

$b$ : distribution over  $\mathcal{X}^n$   
PMP

• Def. **Gibbs free energy**  $Q_{f^{\text{total}}}(b)$ .

$$Q_{f^{\text{total}}}(b) = \sum_{x \in \mathcal{X}^n} b(x) \cdot \log f^{\text{total}}(x) - \sum_{x \in \mathcal{X}^n} b(x) \log b(x)$$

$$= - \underbrace{\mathbb{E}_b \left[ \underbrace{-\log f^{\text{total}}(x)}_{H(x) \text{ energy of state } x} \right]}_{\text{expected energy}} + \underbrace{\mathbb{E}_b \left[ \underbrace{-\log b(x)}_{\text{energy of } b} \right]}_{\text{energy of } b}$$

\*  $P(x) = \frac{1}{Z} \exp(-\underbrace{H(x)}_{\text{energy}})$ : low energy state is more likely.

claim:  $\left\{ \begin{array}{l} \mathcal{G}_{\text{free}}(b) \text{ is strictly concave} \\ \sup_b \mathcal{G}_{\text{free}}(b) = \Phi \\ P(x) = \text{arg max}_b \mathcal{G}_{\text{free}}(b). \end{array} \right.$

Interpretation:  $b^*(x) = P(x)$  minimizes average energy while maximizing entropy.  
 (= more likely for each state)  
 (= more # of states)

Proof: 
$$\begin{aligned} \mathcal{G}_{\text{free}}(b) &= \sum_x b(x) \log f_{\text{total}}(x) - \sum_x b(x) \log b(x) \\ &= \sum_x b(x) \left[ \log Z + \log \underbrace{\frac{1}{Z} f_{\text{total}}(x)}_{P(x)} \right] - \sum_x b(x) \log b(x) \\ &= \log Z + \sum_x (b(x) \log P(x) - b(x) \log b(x)) \\ &= \Phi - \underbrace{D_{\text{KL}}(b \parallel P)}_{\text{Kullback-Leibler Divergence}} \leq \Phi \text{ \& equal if } b=P. \end{aligned}$$

From information theory, we know

$$D_{\text{KL}}(b \parallel P) \geq 0$$

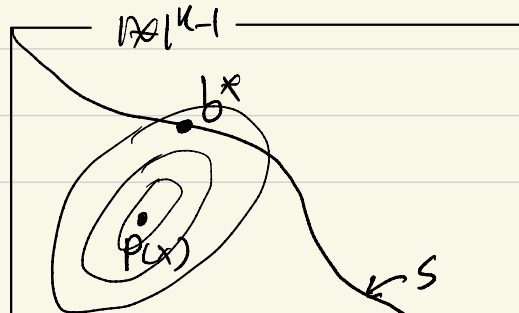
$D_{\text{KL}}(b \parallel P) = 0$  iff  $b=P$ .

$D_{\text{KL}}(b \parallel P)$  is convex in  $b$ .

then, 
$$\Phi = \sup_{b: \Delta_{|\mathcal{X}|^n-1}} \mathcal{G}_{\text{free}}(b)$$

is a concave maximization, but in  $|\mathcal{X}|^n - 1$  dimensions.

Strategy



\* we instead search over a small dimensional space, and find approximate solution.

$$\Phi \geq \sup_{b \in S} \mathcal{G}_{\text{free}}(b)$$

this provides a lower bound

\* Def. Naive Mean Field approach.

Consider Naive Mean Field factorization

$$S_{MF} = \{ b \in \Delta_{|\mathcal{X}| \times |\mathcal{Y}|} : b(x) = b_1(x_1) \times b_2(x_2) \times \dots \times b_n(x_n) \}$$

with a slight abuse of notation, let  $b = b_1 \times b_2 \times \dots \times b_n$ ,  
and plug in to Gibbs free energy to get

$$F_{MF}(b) = G_{free}(b_1 \times b_2 \times \dots \times b_n)$$

$$= \sum_{(i,j) \in E} \sum_{x_i, x_j} b_i(x_i) b_j(x_j) \log J_{ij}(x_i, x_j) - \sum_{i \in V, x_i} b_i(x_i) \log b_i(x_i)$$

Mean Field Variational Inference problem.

$$\max_{b \in S_{MF}} F_{MF}(b)$$

$$\text{s.t. } b_i(x_i) \geq 0 \quad \forall i, x_i$$

$$\sum_{x_i} b_i(x_i) = 1 \quad \forall i$$

\*  $b_i(\cdot)$ 's play the role of approximate marginal distribution

$$b_i(x_i) \cong P(x_i)$$

\* Dimension  $b \in \mathcal{R}^{n(|\mathcal{X}|-1)}$ ,  $b_i \in \Delta_{|\mathcal{X}|-1}$

\* but now the objective is no longer concave:  $b_i(x_i) \cdot b_j(x_j)$  is bilinear.

\* we look for a local maxima

Stationary Point of Naive Meanfield is characterized by Lagrangian

$$\mathcal{L}(b, \lambda) = F_{MF}(b) - \sum_{i \in V} \lambda_i \left\{ \sum_{x_i \in \mathcal{X}} b_i(x_i) - 1 \right\}$$

(we don't include posteriors as they appear in  $\log(\cdot)$  in the objective)

$$= \frac{1}{2} \sum_{i \in V, x_i \in \mathcal{X}} b_i(x_i) \left\{ \sum_{j \in \partial i, x_j \in \mathcal{X}} b_j(x_j) \cdot \log J_{ij}(x_i, x_j) \right\}$$

$$- \sum_{i \in V, x_i \in \mathcal{X}} b_i(x_i) \log b_i(x_i) - \sum_{i \in V} \lambda_i \sum_{x_i \in \mathcal{X}} (b_i(x_i) - 1)$$

taking the derivative, w.r.t  $b_i(x_i)$

$$\frac{\partial \mathcal{L}(b, \lambda)}{\partial b_i(x_i)} = \sum_{j \in \mathcal{N}_i} \sum_{x_j \in \mathcal{X}} b_j(x_j) \log f_{ij}(x_i, x_j) - 1 - \log b_i(x_i) - \lambda_i = 0$$

Def. Naive Mean field equation

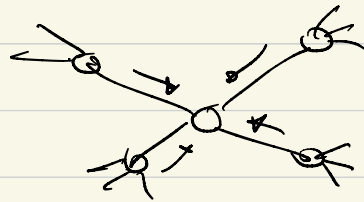
$$b_i(x_i) \propto \exp \left\{ \sum_{j \in \mathcal{N}_i} \sum_{x_j \in \mathcal{X}} \log f_{ij}(x_i, x_j) b_j(x_j) \right\}$$

is a fixed point of:  $b = F_{MF}(b)$

which can be <sup>solved</sup> by

$$b^{(t+1)} = F_{MF}(b^{(t)})$$

One can think of this as having a belief at each node (as opposed to edges) and update by aggregating its neighbors.



similar to gossip algorithms

this gives a very poor approximation.

ex)  $x_1, x_2 \in \{0, 1\}$ .

$$P(x) = \frac{1}{2} \mathbb{I}(x_1 = x_2) \quad \text{vs.} \quad P(x) = \frac{1}{2} \mathbb{I}(x_1 \neq x_2)$$

$$x_1 \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{matrix} \vdots \frac{1}{2} \\ \vdots \frac{1}{2} \end{matrix}$$

$$x_1 \begin{bmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{bmatrix} \begin{matrix} \vdots \frac{1}{2} \\ \vdots \frac{1}{2} \end{matrix}$$

[ account for marginal  $P(x_i)$   
 [ exact only if  $P(x) = P(x_1)P(x_2) \dots P(x_n)$ .



# \* Bethe free energy

[ account for pairwise correlations induced by edges/factors  
 ] exact on distribution  $P$  on trees

- Parameters:  $b_i(x_i)$ : approximates  $P(x_i)$   
 $b_{ij}(x_i, x_j)$ : approximates  $P(x_i, x_j)$

use  $b = \{ b_i, b_{ij} \}$  notation

w.r.t  $G$ .

Def. Ideally we want to search over **Globally Consistent Marginals**  
 $MARG(G) \triangleq \{ b = \{ \{ b_i \}_{i \in V}, \{ b_{ij} \}_{i,j \in E} \} \text{ s.t. } \exists P(x) \text{ with } \left. \begin{array}{l} b_i(x_i) = \sum_{x_{-i}} P(x), \forall i \in V \\ b_{ij}(x_i, x_j) = \sum_{x_{-i,j}} P(x), \forall i,j \in E \end{array} \right\}$

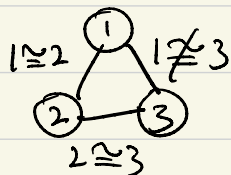
but checking if  $b \in MARG$  is NP-hard.

Def. Instead we propose searching over **Locally Consistent Marginals**  
 $LOC(G) \triangleq \{ b \text{ s.t. } \left. \begin{array}{l} \sum_{x_i} b_i(x_i) = 1, \forall i \in V \\ \sum_{x_j} b_{ij}(x_i, x_j) = b_i(x_i), \forall i,j \in E \end{array} \right\}$

Local consistency does not imply Global consistency

Counter example  $\mathcal{X} = \{0, 1\}$ ,  $b_1 = b_2 = b_3 = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$

$$b_{12} = b_{23} = \begin{bmatrix} 0.49 & 0.01 \\ 0.01 & 0.49 \end{bmatrix}, \quad b_{13} = \begin{bmatrix} 0.01 & 0.49 \\ 0.49 & 0.01 \end{bmatrix}$$



we can check that  $b_i(x_i) = \sum_{x_j} b_{ij}(x_i, x_j)$

• for general graph  $G$ .



• for  $G$  that is a tree.

$\forall$  locally consistent  $b = \{b_{ij}, b_{ji}\}$ ,  $\exists P(x)$  that is globally consistent for  $b$ .

$$\tilde{P}(x) = \prod_{i \in V} b_i(x_i) \cdot \prod_{(i,j) \in E} \frac{b_{ij}(x_i, x_j)}{b_i(x_i) b_j(x_j)}$$

claim:  $\tilde{P}(x)$  is globally consistent if  $G = \text{tree}$ .

$$\text{that is } \begin{cases} b_i(x_i) = \tilde{P}(x_i), \forall i \in V \\ b_{ij}(x_i, x_j) = \tilde{P}(x_i, x_j), \forall (i,j) \in E \end{cases}$$

proof: by induction

$$n=1 \quad \tilde{P}(x_1) = b(x_1) \quad \text{is globally consistent.}$$

suppose it is true for a tree with  $n$  nodes, then we consider a tree with 1 more node that is connected to node  $n$ .



$$\tilde{P}(x_1^n, x_{n+1}) = \tilde{P}(x_1^n) \cdot \frac{b_{n,n+1}(x_n, x_{n+1})}{b_n(x_n) b_{n+1}(x_{n+1})} b_{n+1}(x_{n+1})$$

$$\tilde{P}(x_n, x_{n+1}) = \sum_{x_1, \dots, x_n} \tilde{P}(x_1^n, x_{n+1}) = \tilde{P}(x_n) \cdot \frac{b_{n,n+1}}{b_n \cdot b_{n+1}} \cdot b_{n+1}$$

$$\text{induction } \rightarrow = b(x_n) \cdot \frac{b_{n,n+1}}{b_n \cdot b_{n+1}} \cdot b_{n+1}$$

$$= b_{n,n+1}(x_n, x_{n+1}).$$

Def. **Bethe free energy** on a tree.

Recall, Gibbs free energy for general  $b(x)$  is defined as

$$\mathcal{G}_{\text{total}}(b) = \underbrace{-\mathbb{E}_b[-\log f_{\text{total}}(X)]}_{\text{Energy}} + \underbrace{\mathbb{E}_b[-\log b(x)]}_{\text{Entropy}}$$

We evaluate it on  $b(x) = \prod_{i \in V} b_i(\alpha_i) \prod_{(i,j) \in E} \frac{b_{ij}(X_i, X_j)}{b_i(\alpha_i) b_j(\alpha_j)}$

$$\text{Energy} = - \sum_{(i,j) \in E} \sum_{X_i, X_j} b_{ij}(X_i, X_j) \cdot \log f_{ij}(X_i, X_j)$$

$$\begin{aligned} \text{Entropy} &= \mathbb{E}_b \left[ -\log \left( \prod_i b_i(\alpha_i) \prod_{(i,j) \in E} \frac{b_{ij}(X_i, X_j)}{b_i(\alpha_i) b_j(\alpha_j)} \right) \right] \\ &= \sum_i \sum_{X_i} \left( -b_i(\alpha_i) \log b_i(\alpha_i) \right) - \underbrace{\sum_{(i,j) \in E} \sum_{X_i, X_j} \left( \log b_{ij} - \log b_i - \log b_j \right) b_{ij}}_{\sum_{(i,j) \in E} \sum_{X_i, X_j} b_{ij} \log b_{ij} - \sum_{i \in V} \sum_{X_i} \text{deg}(i) b_i \log b_i} \end{aligned}$$

**Bethe free energy**

$$\mathbb{F}(b) \stackrel{\text{def}}{=} \mathcal{G}_{\text{total}}(b) = \underbrace{-\text{Energy}}_{\text{Tree}} + \text{Entropy}$$

$$= \sum_{(i,j) \in E} \sum_{X_i, X_j} b_{ij}(X_i, X_j) \left( \log f_{ij}(X_i, X_j) - \log b_{ij}(X_i, X_j) \right) + \sum_{i \in V} (\text{deg}(i)-1) \sum_{X_i} b_i(\alpha_i) \log b_i(\alpha_i)$$

claim. if  $G$  is a tree then,

$$\sup_{b \in \text{LOC}(G)} \mathbb{F}(b) = \sup_{b \in \text{MARG}(G)} \mathcal{G}(b) = \Phi \quad : \log \text{ partition function.}$$

this is called **Bethe variational problem**.

if you apply this formula  $\mathbb{F}(b)$  and optimize for a non-tree  $G$ , then it is called **Bethe approximation**.

\*From Bethe free energy to belief propagation.

claim. fixed points of BP are one-to-one correspondence with stationary points of Bethe variational problem.

Further, BP messages  $\{m_{i \rightarrow j}(x_i)\}$  are simple exponentials of Lagrangian variables  $\{\lambda_{i,j}^*(x_i)\}$ .

proof.

Define Lagrangian multipliers  $\lambda_i$  for  $\sum_{x_i} b(x_i) = 1$

$\lambda_{i \rightarrow j}(x_i)$  for  $\sum_{x_j} b_{ij}(x_i, x_j) = b(x_i)$

$$\mathcal{L}(b, \lambda) = \mathbb{F}(b) - \sum_i \lambda_i \left\{ \sum_{x_i} b(x_i) - 1 \right\} - \sum_{(i,j)} \sum_{x_i} \lambda_{i \rightarrow j}(x_i) \left\{ \sum_{x_j} b_{ij}(x_i, x_j) - b_i(x_i) \right\}$$

taking the derivative,

$$\nabla_{b_{ij}(x_i, x_j)} \mathcal{L}(b, \lambda) = -1 - \log b_{ij}(x_i, x_j) + \log f_{ij}(x_i, x_j) - \lambda_{i \rightarrow j}(x_i) - \lambda_{j \rightarrow i}(x_j)$$

$$\nabla_{b_i(x_i)} \mathcal{L}(b, \lambda) = -(1 - \deg(i)) \log [b_i(x_i) \cdot e] - \lambda_i + \sum_{j \in \partial i} \lambda_{i \rightarrow j}(x_i)$$

Setting them to zero,

$$b_{ij}^*(x_i, x_j) = f_{ij}(x_i, x_j) \exp \left\{ -1 - \lambda_{i \rightarrow j}(x_i) - \lambda_{j \rightarrow i}(x_j) \right\}$$

$$b_i^*(x_i) \propto \exp \left\{ -\frac{1}{\deg(i)} \sum_{j \in \partial i} \lambda_{i \rightarrow j}(x_i) \right\}$$

$$\sum_{x_j} b_{ij}^*(x_i, x_j) = b_i^*(x_i)$$

We change variables as:  $m_{i \rightarrow j}(x_i) \propto e^{-\lambda_{i \rightarrow j}(x_i)}$

$$b_{ij}^*(x_i, x_j) \propto m_{i \rightarrow j}(x_i) f_{ij}(x_i, x_j) m_{j \rightarrow i}(x_j)$$

$$b_i^*(x_i) \propto \prod_{j \in \partial i} \left\{ (m_{i \rightarrow j}(x_i))^{\frac{1}{\deg(i)-1}} \right\}$$

$$\sum_{x_j} b_{ij}^*(x_i, x_j) = b_i^*(x_i)$$

To show equivalence b/w BP vs. Bethe free  
we start with.

$$\prod_{k \in \partial(i)} \left\{ \sum_{X_k} b_{ik}^*(X_i, X_k) \right\} \stackrel{\substack{\uparrow \\ \text{Local} \\ \text{causality}}}{=} \prod_{k \in \partial(i)} b_i^*(X_i) = b_i^*(X_i)^{\deg(i)-1}$$

$\downarrow$  Stationarity  $\rightarrow S$ 
 $\leftarrow$  stationary  $F(b)$

$$\prod_{k \in \partial(i)} \left\{ m_{i \rightarrow k}(X_i) \cdot \sum_{X_k} m_{k \rightarrow i}(X_k) f_{ik}(X_i, X_k) \right\} \propto \prod_{k \in \partial(i)} m_{i \rightarrow k}(X_i)$$

$\downarrow$ 
 $\downarrow$  cancelling

$$\prod_{k \in \partial(i)} \sum_{X_k} m_{k \rightarrow i}(X_k) f_{ik}(X_i, X_k) \propto$$

$\downarrow$   
 we arrive at BP update