**EM** $P(\mathbf{d},\mathbf{h}\,|\,\lambda)=\dfrac{e^{\lambda\cdot\varphi(\mathbf{d},\mathbf{h})}}{\displaystyle\sum_{\mathbf{d},\mathbf{h}}e^{\lambda\cdot\varphi(\mathbf{d},\mathbf{h})}}$

Example: Graph with no closed loops

$\varphi(\mathbf{d},\mathbf{h})=\begin{cases}\varphi(d_i,h_i):i=1,\ldots,N,\\\psi(h_i,h_{i+1}):i=1,\ldots,N-1\end{cases}$ $\lambda=\begin{cases}\lambda_i:i=1,\ldots,N,\\\mu_i:i=1,\ldots,N-1\end{cases}$

$$\lambda\cdot\varphi(\mathbf{d},\mathbf{h})=\sum_{i=1}^{N}\lambda_i\cdot\varphi(d_i,h_i)+\sum_{i=1}^{N-1}\mu_i\cdot\psi(h_i,h_{i+1})$$

$$Z[\lambda]=\sum_{\mathbf{d},\mathbf{h}}e^{\lambda\cdot\varphi(\mathbf{d},\mathbf{h})}$$

Note In this special case: $Z[\lambda]$ can be computed by Dynamic Programming because there are no closed loops.

Problem 1. Compute $P(\mathbf{d}\mid\lambda) = \sum_{\mathbf{h}} P(\mathbf{d},\mathbf{h}\mid\lambda)$

This can be computed by Dynamic Programming (sum rule)

If closed loops, approximated by Belief Propagation (sum-product)

Problem 2. Compute $\hat{\mathbf{h}} = \arg\max P(\mathbf{d},\mathbf{h}\mid\lambda)$

Again, this can be computed by DP (max)

If closed loops, approximated by BP (max-product)

Problem 3. Learning $\lambda$

Data $D = \left\{\mathbf{d}^m : m = 1,\ldots,M\right\}$ $\implies$ $\hat{\lambda} = \arg\max_{\lambda} \prod_{m=1}^{M} P(\mathbf{d}^m \mid \lambda)$

$$= \arg\max_{\lambda} \prod_{m=1}^{M} \sum_{\mathbf{h}^m} P(\mathbf{d}^m, \mathbf{h}^m \mid \lambda)$$

## EM with Free Energy

Introduce distribution $Q_m(\mathbf{h}^m)$

$$F\left[\boldsymbol{\lambda}:\{Q_m(\mathbf{h}^m)\}\right] = \sum_{m=1}^{M}\left\{-\log P(\mathbf{d}^m \mid \boldsymbol{\lambda}) + \sum_{\mathbf{h}^m} Q_m(\mathbf{h}^m)\log\frac{Q_m(\mathbf{h}^m)}{P(\mathbf{h}^m \mid \mathbf{d}^m, \boldsymbol{\lambda})}\right\}$$

Re-express this as:

$$F\left[\boldsymbol{\lambda}:\{Q_m(\mathbf{h}^m)\}\right] = \sum_{m=1}^{M}\left\{\sum_{\mathbf{h}^m} Q_m(\mathbf{h}^m)\log Q_m(\mathbf{h}^m) - \sum_{\mathbf{h}^m} Q_m(\mathbf{h}^m)\log P(\mathbf{h}^m, \mathbf{d}^m \mid \boldsymbol{\lambda})\right\}$$

The EM algorithm minimizes $F[\boldsymbol{\lambda}:\{Q_m(\cdot)\}]$ w.r.t. $\boldsymbol{\lambda}$ and the $\{Q_m(\cdot)\}$ alternatively

$$\boldsymbol{\lambda}^{t+1} = \arg\min_{\boldsymbol{\lambda}} F\left[\boldsymbol{\lambda}:\{Q_m^t(\cdot)\}\right]$$

$$Q_m^{t+1}(\cdot) = \arg\min_{Q_m} F\left[\boldsymbol{\lambda}^{t+1}:\{Q_m(\cdot)\}\right]$$

(B) $\lambda^{t+1} = \arg\min_{\lambda}\left\{-\sum_{\mathbf{h}^m} Q_m^t(\mathbf{h}^m)\log P(\mathbf{h}^m,\mathbf{d}^m\mid\lambda)\right\}$

(A) $Q_m^{t+1}(\mathbf{h}^m) = P(\mathbf{h}^m\mid\mathbf{d}^m,\lambda^t)$

(A) How to compute these update rules if $P(\mathbf{h}^m,\mathbf{d}^m\mid\lambda^t) = \dfrac{e^{\lambda\cdot\varphi(\mathbf{d},\mathbf{h})}}{Z[\lambda]}$ ?

$$P(\mathbf{h}^m\mid\mathbf{d}^m,\lambda^t) = \frac{P(\mathbf{h}^m,\mathbf{d}^m\mid\lambda^t)}{P(\mathbf{d}^m\mid\lambda^t)}$$

where $P(\mathbf{d}^m\mid\lambda^t) = \dfrac{1}{Z[\lambda]}\sum_{\mathbf{h}^m} e^{\lambda\cdot\varphi(\mathbf{d}^m,\mathbf{h}^m)}$

Hence, $P(\mathbf{h}^m\mid\mathbf{d}^m,\lambda^t) = \dfrac{e^{\lambda\cdot\varphi(\mathbf{d}^m,\mathbf{h}^m)}}{\boxed{\sum_{\mathbf{h}^m} e^{\lambda\cdot\varphi(\mathbf{d}^m,\mathbf{h}^m)}}}$ ⟵ This term can be directly computed by DP(sum) if the graph has no closed loop

Hence, $Q_m^{t+1}(\mathbf{h}^m) = P(\mathbf{h}^m\mid\mathbf{d}^m,\lambda^t)$ can be computed (no closed loops)

(B) How to compute $\lambda^{t+1} = \arg\min_{\lambda}\left\{-\sum_{\mathbf{h}^m} Q_m^t(\mathbf{h}^m)\log P(\mathbf{h}^m,\mathbf{d}^m\,|\,\lambda)\right\}$

Substitute $P(\mathbf{h},\mathbf{d}\,|\,\lambda) = \dfrac{e^{\lambda\cdot\varphi(\mathbf{h},\mathbf{d})}}{Z[\lambda]}$

Want to minimize: $G(\lambda) = -\displaystyle\sum_{m=1}^{M} Q_m^t(\mathbf{h}^m)\cdot\lambda\cdot\varphi(\mathbf{h}^m,\mathbf{d}^m) - \sum_{m=1}^{M}\log Z[\lambda]$

It can be shown that $G(\lambda)$ is a convex function of $\lambda$ (because $\log Z[\lambda]$ is convex)

The global minimum $\hat{\lambda}$ occurs where

$$\frac{\partial}{\partial\lambda} G(\hat{\lambda}) = 0$$

when $\dfrac{1}{M}\displaystyle\sum_{m=1}^{M} Q_m^t(\mathbf{h}^m)\varphi(\mathbf{h}^m,\mathbf{d}^m) = \sum_{\mathbf{h},\mathbf{d}} \varphi(\mathbf{h},\mathbf{d})P(\mathbf{h},\mathbf{d}\,|\,\lambda)$

i.e. when the expected statistics w.r.t. data $\mathbf{d}^m$ and $Q_m(\cdot)$
= the expected statistics of the model

**Note** Deriving the last equation as follows:

$$\frac{\partial}{\partial \boldsymbol{\lambda}} \log Z[\boldsymbol{\lambda}] = \frac{\partial}{\partial \boldsymbol{\lambda}} \log \sum_{\mathbf{h,d}} e^{\boldsymbol{\lambda} \cdot \varphi(\mathbf{h,d})} = \frac{\sum_{\mathbf{h,d}} \varphi(\mathbf{h,d}) \cdot e^{\boldsymbol{\lambda} \cdot \varphi(\mathbf{h,d})}}{\sum_{\mathbf{h,d}} e^{\boldsymbol{\lambda} \cdot \varphi(\mathbf{h,d})}} = \sum_{\mathbf{h,d}} P(\mathbf{h,d} \mid \boldsymbol{\lambda}) \cdot \boldsymbol{\varphi}(\mathbf{h,d})$$

(recall learning notes for
learning exponential distributions)

Hence, the update rule for $\boldsymbol{\lambda}^{t+1}$ requires finding the value of $\boldsymbol{\lambda}$
so that the expected statistics (w.r.t. the data $\mathbf{d}^m$ & $Q_m(\cdot)$) are
equal to the statistics of the model

Note This is a generalization of the result for Hidden Markov Models.