- Recap:
  - HMMs with known params $\odot = \{\pi_0, \wedge, \underline{\Phi}\}$ are trees.
  - (Draw model from last time)
  - Belief propagation therefore allows us to compute all singleton marginals $P(z_t \mid \bar{x}_{1:T})$ AKA "beliefs".
  - BP in HMMs $\equiv$ the "forwards-backwards" algo.
  - To learn parameters, we need inference.
  - EM: iterates between updating params, doing inference.

# Expectation - Maximization (EM) (1977)

- complete likelihood $P_\theta(x, z)$
- marginal likelihood $P_\theta(x)$ aka "evidence"
- Goal: Type-II MLE

$$\hat{\theta} = \underset{\theta}{\text{argmax}} \; P_\theta(x)$$

- Setting up the evidence lower bound (ELBO)

$$\log P_\theta(x) = \log \underset{z \sim Q(z)}{\mathbb{E}} \left[ \frac{P(x,z)}{Q(z)} \right]$$

trivially true for any density $Q(z)$ s.t.

$$Q(z) > 0 \quad \text{if} \quad P(z|x) > 0.$$

$$\geq \underset{Q}{\mathbb{E}} \left[ \log \frac{P_\theta(x,z)}{Q(z)} \right] \overset{\Delta}{=} B(Q, \theta)$$
$$(\text{ELBO})$$

$$= \underset{Q}{\mathbb{E}} \left[ \log P(x,z) \right] - \underbrace{\underset{Q}{\mathbb{E}} \left[ Q(z) \right]}_{= H(Q(z))}$$

$$= \underbrace{\underset{Q}{\mathbb{E}} \left[ \log P(x) \right]}_{= P(x)} + \underbrace{\underset{Q}{\mathbb{E}} \left[ \log \frac{P(z|x)}{Q(z)} \right]}_{= -kL\left( Q(z) \| P(z|x) \right)}$$

Picture: $\log P_\theta(x)$ $\quad\begin{array}{c} kL(Q, P_\theta) \\[20pt] B(\theta, Q) \end{array}$

The EM algorithm

- initialize $\theta_0$

for $m = 1, 2, \ldots$ until convergene in $\theta_m$

- $Q_m = \underset{Q}{\text{argmax}} \, B(\theta_{m-1}, Q)$    "E-step"

- $\theta_m = \underset{\theta}{\text{argmax}} \, B(\theta, Q_m)$    "M-step"

This is a <u>minorize-maximize</u> algorithm



$\theta_{m-1}$  $\theta_m$

- converges to <u>local</u> mode of $P(x; \theta)$
- random restarts for $\theta_0$ good idea
- ELBO is tight after E-step and always increases

## E-step

$$Q_m = \underset{Q}{\text{argmax}} \; B(\theta_{m-1}, Q)$$

$$= \underset{Q}{\text{argmax}} \; -kl(Q \| P_{\theta_{m-1}}(z|x))$$

$$= P_{\theta_{m-1}}(z|x)$$

## M-step

$$\theta_m = \underset{\theta}{\text{argmax}} \; B(\theta, Q_m)$$

$$= \underset{\theta}{\text{argmax}} \; \underset{Q_m}{\mathbb{E}}\left[\log P_\theta(x,z)\right]$$

$$= \underset{\theta}{\text{argmax}} \; \sum_z P_{\theta_{m-1}}(z|x) \log P_\theta(x,z)$$

Notice that the E-step seems redundant...

The entire algorithm could be stated as:

$$\boxed{\begin{array}{l} \text{until convergence} \\ \theta^{new} = \underset{\theta}{\text{argmax}} \; \underset{z \sim P(z|x; \theta^{old})}{\mathbb{E}}\left[\log P(x,z; \theta)\right] \end{array}}$$

... assuming that $\underset{P(z|x)}{\mathbb{E}}\left[\log P(x,z)\right]$ is tractable.

Is it? Often. But why? $\sum_z P(z|x; \theta^{old}) \log P(x,z; \theta)$

seems just as bad as $\sum_z P(x,z; \theta) = P(x; \theta) \cdots$

Recall last time:

$$P_\theta(x, z) = P_\theta(\bar{x}_{1:T}, z_{1:T}) = \prod_t P(z_t | z_{t-1}) \, P(\bar{x}_t | z_t)$$

$$= \prod_t \Lambda(z_{t-1}, z_t) \prod_k \text{Pois}(\bar{x}_t ; M_k)^{1(z_t = k)}$$

Think about the M-step for just $M_k$...

$$\text{argmax}_{M_k} \; \mathbb{E}\left[ \log P(\bar{x}_{1:T}, z_{1:T} ; \Lambda, M) \right]$$

$$= \text{argmax}_{M_k} \; \mathbb{E}\left[ \log \prod_t \text{Pois}(\bar{x}_t ; M_k)^{1(z_t = k)} \right]$$

$$= \text{argmax}_{M_k} \; \sum_{z_1} \cdots \sum_{z_T} P(z_1 \cdots z_T | \bar{x}_{1:T} ; M^{old}, \Lambda^{old})$$
$$\log \prod_t \text{Pois}(\cdots)^{1(\cdots)}$$

$$= \text{argmax}_{M_k} \; \sum_{z_1} \cdots \sum_{z_T} P(z_1 \cdots z_T | \bar{x}_{1:T}, M^{old}, \Lambda^{old})$$
$$\sum_t 1(z_t = k) \log \text{Pois}(\bar{x}_t ; M_k)$$

Note that:

$$\sum_{z_1} \cdots \sum_{z_T} P(z_1 \cdots z_T | \cdots) \, 1(z_t = k)$$

$$= \sum_{z_t} P(z_t | \cdots) \, 1(z_t = k)$$

$$= P(z_t = k | \cdots) \quad \text{"belief"}$$

$$= \underset{\mu_k}{\text{argmax}} \quad \sum_t \underbrace{P(z_t = k \mid \bar{x}_{1:T}, \mu_k^{old})}_{\overset{\circ}{=} q_{tk}} \log \text{Pois}(\bar{x}_t; \mu_k)$$

So in this case we only need the beliefs $q_{tk}$.

$$= \underset{\mu_k}{\text{argmax}} \quad \sum_t q_{tk} \left[ \bar{x}_t \log \mu_k - \mu_k \right]$$

$$= \underset{\mu_k}{\text{argmax}} \quad \left( \sum_t \bar{x}_t q_{tk} \right) \log \mu_k - \left( \sum_t q_{tk} \right) \mu_k$$

$$\hookrightarrow \quad \mu_k = \frac{\sum_t \bar{x}_t q_{tk}}{\sum_t q_{tk}}$$

This is why it is called the "Expectation" - Step.

① compute all <u>expectations</u> required by M-step:

$$\text{e.g.,} \quad q_{tk} = P(z_t = k \mid \bar{x}_{1:T}, \mu^{old}, \Lambda^{old}) \quad \forall \, t, k$$

$$\equiv \underset{z_{1:T} \sim P(z_1, \dots z_T \mid \bar{x}_{1:T}, \mu^{old}, \Lambda^{old})}{\mathbb{E} \left[ 1(z_t = k) \right]}$$

② <u>Maximize</u> $\mu, \Lambda$ ...

Notice also that the M-step was "nice".
why? Exponential family conditionals.

$$\underset{\mu_k}{\text{argmax}} \quad \sum_t \underbrace{P(z_t = k \mid \bar{x}_{1:T}, \mu_k^{old})}_{\overset{\Delta}{=} q_{tk}} \log \text{Pois}(\bar{x}_t; \mu_k)$$

Rewrite in terms of natural parameter $\eta_k = \log \mu_k$.

$$\underset{\eta_k}{\text{argmax}} \quad \sum_t q_{tk} \left( \eta_k^T t(\bar{x}_t) - a(\eta_k) \right)$$

$$\underset{\eta_k}{\text{argmax}} \quad \eta_k^T \left( \sum_t t(\bar{x}_t) q_{tk} \right) - \left( \sum_t q_{tk} \right) a(\eta_k)$$

$$\nabla_{\eta_k} " \qquad " = \sum_t q_{tk} t(\bar{x}_t) - \mu_k \left( \sum_t t(\bar{x}_t) \right)$$

$$\left( \text{Note for all expfam } \nabla_\eta a(\eta) = \mu \right)$$

$$= 0 \implies \eta_k = \dots \quad \text{s.t.} \quad \mu_k = \frac{\sum_t q_{tk} t(\bar{x}_t)}{\sum_t t(\bar{x}_t)} .$$

Returning to the question of why

$$\sum_z P(z|x; \theta^{old}) \log P(x, z; \theta)$$

is often a "nicer" objective than

$$\sum_q P(x, z; \theta)$$

At a very high level, for the following reason:

$$\mathbb{E} \log \prod \exp(\Sigma \cdots) = \Sigma \Sigma \mathbb{E}[\cdots]$$

we will see this motif again...

e.g. ...

$$\mathbb{E}\left[\log \prod_t \prod_k \exp\left(\eta_k^T t(x) - a(\eta_k)\right)^{f_{tk}} \right] \qquad \xleftarrow{\ \triangleq 1(z_t = k)}$$

$$= \mathbb{E}\left[\log \exp\left(\sum_k \eta_k^T \left(\sum_t f_{tk} t(x_t)\right) - \sum_k \left(\sum_t f_{tk}\right) a(\eta_k)\right)\right]$$

$$\cdots \qquad q_{tk} \stackrel{\triangle}{=} \mathbb{E}[f_{tk}]$$

$$= \sum_k \eta_k^T \left(\sum_t q_{tk} t(x_t)\right) - \sum_k \left(\sum_t q_{tk}\right) a(\eta_k)$$

Models with expfam conditionals and lots of conditional independence tend to lead to "nice" EM.

Do we only ever need singleton beliefs? No...

$$B(\theta, Q) \propto_\Lambda \mathbb{E}_Q[\log P(x, z; \vartheta)]$$

$$\propto_\Lambda \mathbb{E}_Q \left[ \log \prod_{t=1}^{T} \prod_{k=1}^{k} \prod_{j=1}^{k} P(z_t = k \mid z_{t-1} = j)^{\mathbb{1}(z_t = k)\mathbb{1}(z_{t-1} = j)} \right]$$

$$= \sum_t \sum_k \sum_j \mathbb{E}_Q \left[ \mathbb{1}(z_t = k, z_{t-1} = j) \right] \log \Lambda(j, k)$$

$$= \sum_k \sum_j \left[ \sum_t Q(z_t = k, z_{t-1} = j) \right] \log \Lambda(j, k)$$

$$\frac{\partial}{\partial \Lambda(j, k)} \left[ B(\theta, Q) + \eta_0 \left( \sum_k \Lambda(j, k) - 1 \right) \right] \quad \longleftarrow \text{ lagrange multiplier to enforce that } \Lambda_j \text{ sums to 1}$$

$$= \frac{\sum_t Q(z_t = k, z_{t-1} = j)}{\Lambda(j, k)} - \eta_0$$

$$0 = \cdots$$

$$\Lambda(j, k) = \frac{\sum_t Q(z_t = k, z_{t-1} = j)}{\eta_0}$$

Setting $\eta_0 = \sum_t \sum_{k'} Q(z_t = k', z_{t-1} = j) = \sum_t Q(z_{t-1} = j)$

satisfies the constraint that $\eta_0$ enforces.

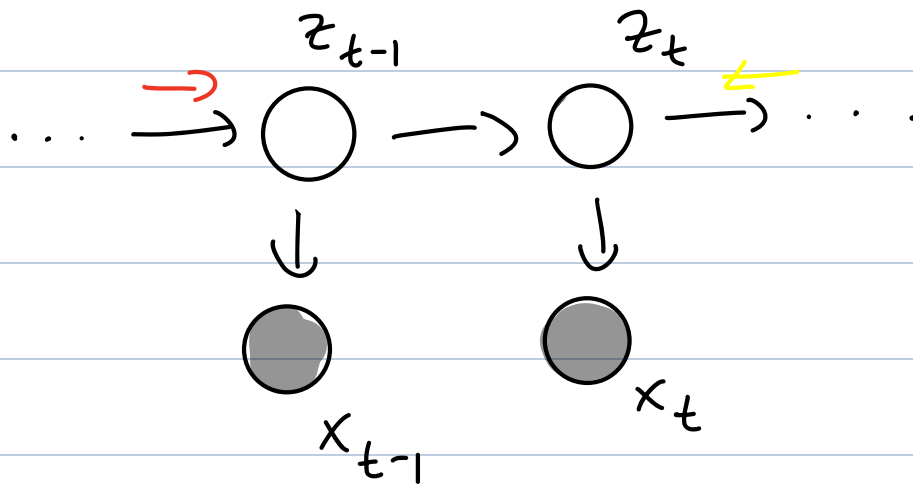$$\Lambda(j, k) = \frac{\sum_t Q(z_t = k, z_{t-1} = j)}{\sum_t Q(z_{t-1} = j)}$$

Summary: To maximize the ELBO with respect to $\Lambda$, we need the pairwise marginals $Q(z_t, z_{t-1})$.

# Pairwise marginals in HMMs

The optimal $Q^*(z_t, z_{t-1}) = P(z_t, z_{t-1} \mid \bar{x}_{1:T})$

$\rightarrow$ How do we compute the pairwise marginals?

Think of $(z_{t-1}, z_t)$ as a "super node"...



Then $p(z_{t-1}, z_t \mid \bar{x}_{1:T}) \propto p(z_{t-1}, z_t, \bar{x}_{1:T})$

$$\propto \alpha_{t-1}(z_{t-1}) \beta_t(z_t) \ell_{t-1}(z_{t-1}) \ell(z_t) \Lambda(z_{t-1}, z_t)$$

$$= p(z_{t-1}, z_t, \bar{x}_{1:T})$$

Note that this is still $O(k^2)$.

More generally for $N$-way marginal $O(k^N)$.

# What is or is not tractable in HMMs

- Evaluating the joint at given $z_1 \dots z_T$:

$$P(z_1 = z_1 \dots z_T = z_t \mid \bar{x}_{1:T}, \ominus) \qquad \mathcal{O}(k^2 T) \text{ to run BP.} \checkmark$$

- Evaluating the evidence

$$P(\bar{x}_{1:T} ; \ominus) \qquad \mathcal{O}(k^2 T) \checkmark$$

- Evaluating a gradient with backprop

$$\nabla_{\ominus} P(\bar{x}_{1:T} ; \ominus) \qquad \mathcal{O}(k^2 T) \checkmark$$

- Storing the joint:

$$P(z_1 \dots z_T \mid \bar{x}_{1:T}, \ominus) \quad \forall z_1 \dots z_T \qquad \mathcal{O}(k^T) \text{ values } \ddot{\frown}$$

- Storing the $N$-marginals:

$$\mathcal{O}(k^N) \text{ values} \qquad \ddot{/} \text{ depends on } N$$

---

Reasoning about the joint $P(z_1 \dots z_T \mid \dots)$ is often important.

e.g. speech recognition. Tractable options:

- MAP: $\underset{z_{1:T}}{\text{argmax}} \; P(z_1 \dots z_T \mid \bar{x}_{1:T} ; \theta)$
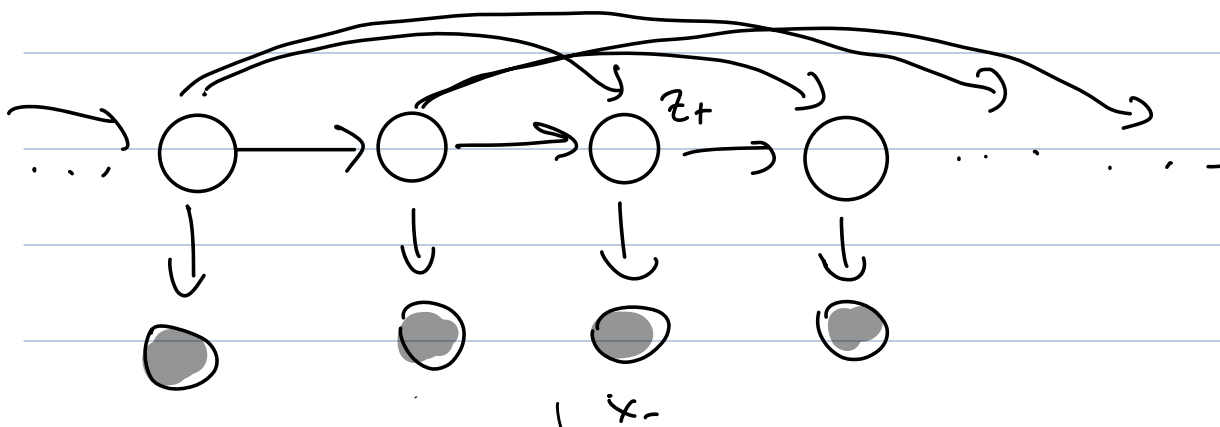
  "Viterbi algorithm"

- Sample: $z_1 \dots z_T \sim P(z_1 \dots z_T \mid \bar{x}_{1:T}, \ominus)$

  "forward filtering backward sampling" (FFBS)

Return to last lecture to cover:

Viterbi, FFBS, extensions of HMMs

# Monte Carlo (MC) - EM

Say we have an __m-order__ HMM



- $z_{t+1} \sim p(z_{t+1} \mid z_t \ldots z_{t-m})$

- $\Lambda$ is a $(KM \times K)$ matrix
  ( or a $\underbrace{K \times \cdots \times K}_{m+1}$ tensor)

- Exact EM would require all
  $(M+1)$-marginals $p(z_{t+1} \ldots z_{t-m} \mid \bar{x}_{1:T})$

- Large $m \to$ intractable

- Instead, replace the E-Step with:

  - $z_1^s \ldots z_T^s \sim p(z_1 \ldots z_T \mid \bar{x}_{1:T}, \Theta^{old})$    $s = 1 \ldots S$

  - $Q(z_{1:T}) = \frac{1}{S} \sum_{s=1}^{S} \mathbb{1}(z_{1:T} = z_{1:T}^s)$
    "atomic measure", each unique $z_{1:T}^s$ is an "atom".

- $\Theta^{new} = \operatorname*{argmax}_{\Theta} \frac{1}{S} \sum_{s} \log p(\bar{x}_{1:T}, z_{1:T}^s ; \Theta)$

- $S = 1$ is called "Stochastic EM".

# Variational EM

- Constrain $Q(z_{1:T})$ to be from a _tractable_ family $\mathcal{F}$

- e.g. $Q(z_{1:T}) = \prod_t Q(z_t)$ "factorized"

- ## Variational E-step:

$$Q^{new} = \operatorname*{argmin}_{Q \in \mathcal{F}} KL\left(Q(z) \| P(z|x; \Theta)\right)$$