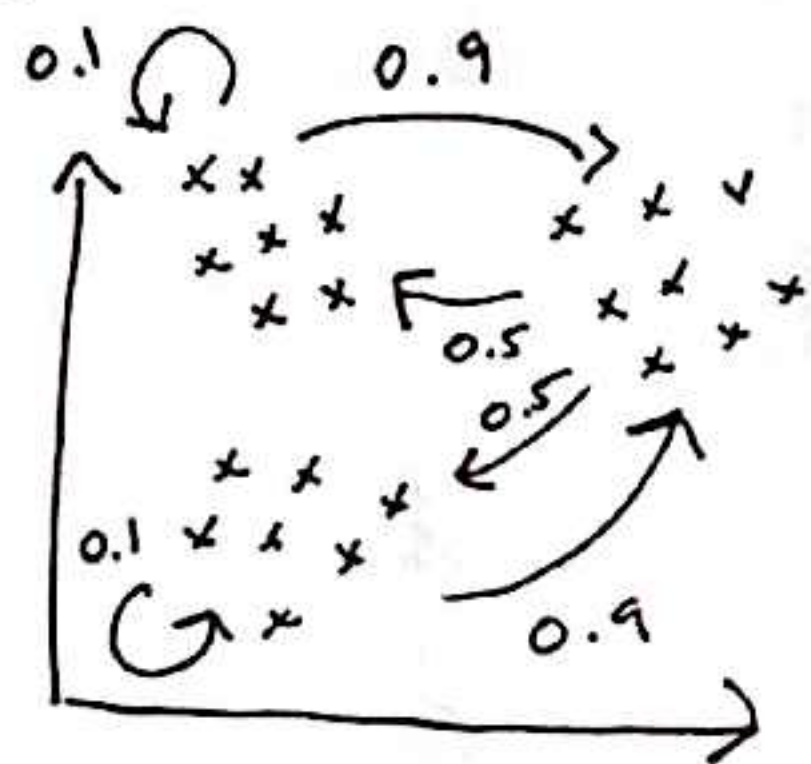


Sequential models

- in many settings data are sequential
 - language data
 - time indexed data
 - DNA
- we will study 2 sequential models - HMM & Kalman filter (further make EM concrete...)
- the HMM is a sequential generalization of the mixture model
- idea - the chosen mixture component ^{at time t} depends on the component at time t-1.

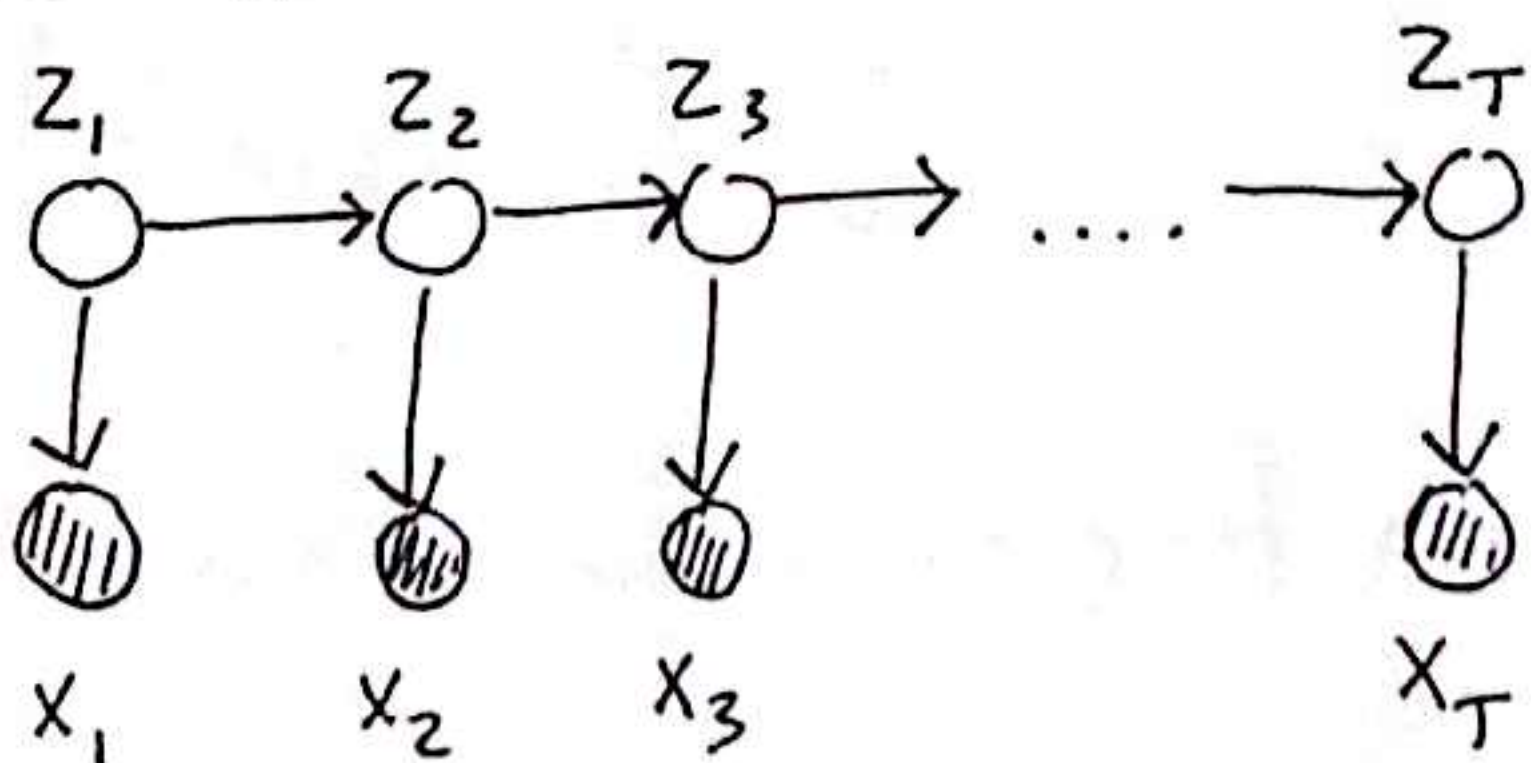
(contrast to a mixture, where the components are assumed \perp .)



The HMM models transitions b/w the components.

[show data generated from this picture.]

HMM modeling assumptions



z_t : discrete (K values)
 x_t : anything (e.g., \mathbb{R}^2)

$$p(z_{1:T}, x_{1:T}) = p(z_1) p(x_1 | z_1) \prod_{t=2}^T p(z_t | z_{t-1}) p(x_t | z_t)$$

(or group the x's all together...)

Parameters are —

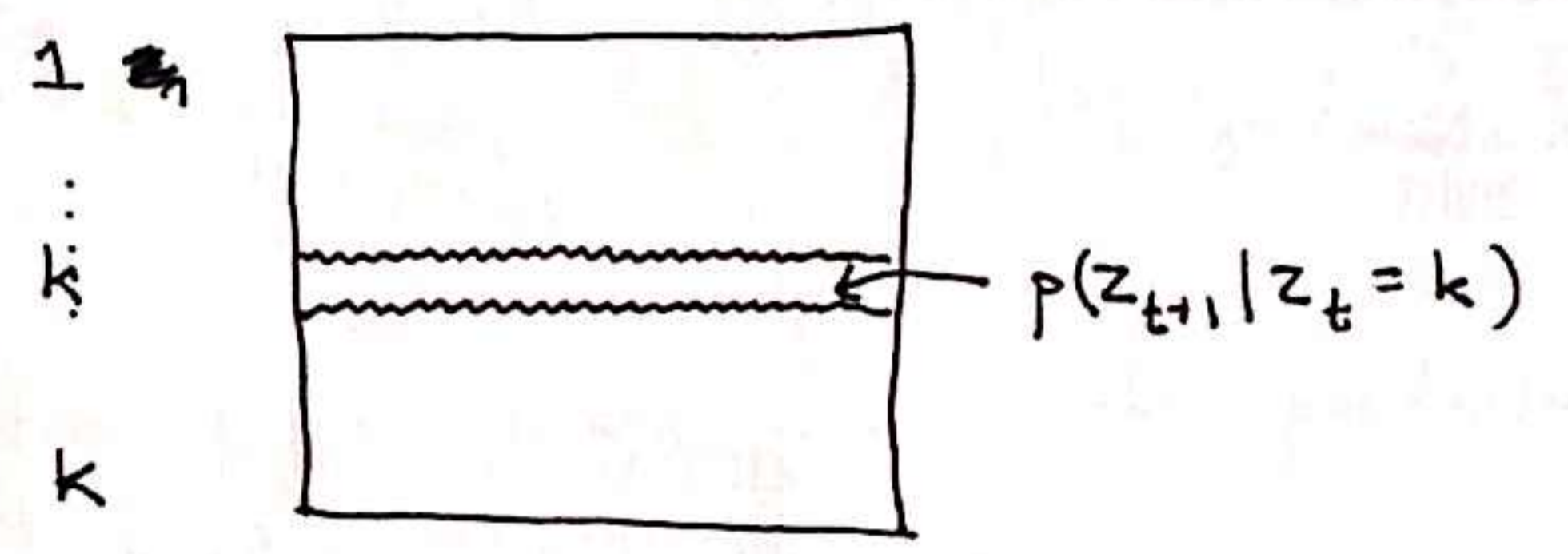
- Data generating distributions θ_k (also called emission probs)

$p(x_t | z_t) = \prod_{k=1}^K z_t^k p(x_t | \theta_k)$ ←

- As for mixtures, this can be anything.
- To be concrete, make θ_k parameters to 2D Gaussians (note - MN in genetics)

- Transition matrix A

- Probability of the next state given the current state.



- Recall that z_t is an indicator vector

- Thus,

$$p(z_t | z_{t-1}) = \prod_{k=1}^K \prod_{j=1}^K a_{jk}^{z_{t-1}^j z_t^k}$$

- Probability of the first state (to start things off) π

$$p(z_1) = \prod_{k=1}^K \pi_k^{z_1^k}$$

- Show a state transition diagram and emphasize it is not a G.M.

- Applications

- speech recognition
- handwriting recognition
- DNA analysis
- others...
- part of speech tagging

Fitting an HMM

EM algorithm.

① Write down the expected complete log likelihood

$$\begin{aligned} \mathbb{E}[\log p(z_{1:T}, x_{1:T})] &= \mathbb{E}\left[\log p(z_1) \prod_{t=2}^T p(z_t | z_{t-1}) \prod_{t=1}^T p(x_t | z_t)\right] \\ &= \mathbb{E}\left[\log \prod_{k=1}^K \pi_k^{z_1^k} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K a_{jk}^{z_{t-1}^j z_t^k} \prod_{t=1}^T \prod_{k=1}^K z_t^k p(x_t | \theta_k)\right] \\ &= \sum_{k=1}^K \mathbb{E}[z_1^k] \log \pi_k + \sum_{t=2}^T \sum_{j=1}^K \sum_{k=1}^K \mathbb{E}[z_{t-1}^j z_t^k] \log a_{jk} + \sum_{t=1}^T \mathbb{E}[z_t^k] \log p(x_t | \theta_k) \end{aligned}$$

All expectations are taken w/r/t the posterior. $p(z_{1:T} | x_{1:T})$

And, since z_t is an indicator vector —

$$\mathbb{E}_x[z_1^k] = p(z_1 = k | x_{1:T}), \quad \mathbb{E}_x[z_t^k] = p(z_t = k | x_{1:T})$$

$$\mathbb{E}_x[z_{t-1}^j z_t^k] = p(z_{t-1} = j, z_t = k | x_{1:T}) \triangleq \xi(z_{t-1}^j, z_t^k)$$

$$\Delta = \gamma(z_t^k)$$

② M-step Note as usual that the log complete likelihood decomposes.

$$\pi_k = \mathbb{E}[z_1^k] / \sum_j \mathbb{E}[z_1^j] \quad \text{normalized expected counts}$$

$$a_{jk} = \sum_{t=2}^T \mathbb{E}[z_{t-1}^j z_t^k] / \sum_{t=2}^T \sum_{l=1}^K \mathbb{E}[z_{t-1}^j z_t^l]$$

note: this stays at j.

θ_k : weighted maximum likelihood, as for mixture models.

$$\text{Gaussian} - \mu_k = \frac{\sum_{t=2}^T \mathbb{E}[z_t^k] x_t}{\sum_{t=2}^T \mathbb{E}[z_t^k]}$$

$$\Sigma_k = \frac{\sum_{t=2}^T \mathbb{E}[z_t^k] (x_t - \mu_k)(x_t - \mu_k)^T}{\sum_{t=2}^T \mathbb{E}[z_t^k]}$$

Discrete (e.g., biology) -

$$p(x_t | \theta_k) = \prod_{i=1}^V \theta_{ki}^{x_t^i}$$

note: this is how to do a mixture of \mathbb{R}^d Gaussians

$$\theta_{ki} = \frac{\sum_{t=2}^T \mathbb{E}[z_t^k] x_t^i}{\sum_{t=2}^T \mathbb{E}[z_t^k]}$$

Others... Poisson, etc.

E-step - also used for prediction.

This is the interesting part

$$\mathbb{E}[z_t | x_{1:T}] = p(z_t | x_{1:T})$$

$$= p(z_t, x_{1:T}) / p(x_{1:T})$$

$$= \cancel{p(x_{1:t}, z_t)} p(x_{1:t}, z_t) p(x_{t+1:T} | z_t) / p(x_{1:T})$$

We will compute these terms recursively -

$$\alpha(z_t) \triangleq p(x_{1:t}, z_t) \quad [K\text{-vector}]$$

$$\beta(z_t) \triangleq p(x_{t+1:T} | z_t) \quad [K\text{-vector}]$$

$$= \alpha(z_t) \beta(z_t) / p(x_{1:T})$$

- Also lets us compute:

$$\begin{aligned}
 \mathbb{E}[z_{t-1} z_t | x_{1:T}] &= p(z_{t-1}, z_t | x_{1:T}) \\
 &= p(x_{1:T}, z_{t-1}, z_t) / p(x_{1:T}) \\
 &= p(x_{1:t-1}, z_{t-1}) p(z_t | z_{t-1}) p(x_t | z_t) p(x_{t+1:T} | z_t) / p(x_{1:T}) \\
 &= \alpha(z_{t-1}) p(z_t | z_{t-1}) p(x_t | z_t) \beta(z_t) / p(x_{1:T})
 \end{aligned}$$

(Note—reursively applying this gives us any subsequence. But, we only need pairs for EM.)

- The likelihood is easy from these quantities, from \star

$$\sum_{z_t} \alpha(z_t) \beta(z_t) / p(x_{1:T}) = 1 \rightarrow p(x_{1:T}) = \sum_{z_t} \alpha(z_t) \beta(z_t)$$

- We are left to compute α, β : Recursive algorithm.

Sometimes called forward-backward, α - β , is an instance of propagation.

$$\begin{aligned}
 \alpha(z_1) &= p(z_1, x_1) \\
 &= \pi_{z_1} p(x_1 | z_1) \quad \text{for slight abuse of notation.}
 \end{aligned}$$

$$\begin{aligned}
 \alpha(z_{t+1}) &= p(x_{1:t+1}, z_{t+1}) \\
 &= p(x_{1:t+1} | z_{t+1}) p(z_{t+1}) \\
 &= p(x_{1:t} | z_{t+1}) p(x_{t+1} | z_{t+1}) p(z_{t+1}) \quad [x_{t+1} \perp\!\!\!\perp x_{1:t} | z_{t+1}] \\
 &= \cancel{\prod_{t+1}} p(x_{1:t}, z_{t+1}) p(x_{t+1} | z_{t+1}) \\
 &= \sum_{z_t} p(x_{1:t+1}, z_t, z_{t+1}) p(x_{t+1} | z_{t+1}) \\
 &= \sum_{z_t} p(x_{1:t}, z_{t+1} | z_t) p(z_{t+1}) p(x_{t+1} | z_{t+1})
 \end{aligned}$$

$$= \sum_{z_t} p(x_{1:t} | z_t) p(z_{t+1} | z_t) p(z_t) p(x_{t+1} | z_{t+1}) \quad [z_{t+1} \perp\!\!\!\perp x_{1:t} | z_t] \quad (6)$$

$$= \sum_{z_t} \underbrace{p(x_{1:t} | z_t)}_{\alpha(z_t)} \underbrace{p(z_{t+1} | z_t)}_{a_{z_t, z_{t+1}}} \underbrace{p(x_{t+1} | z_{t+1})}_{\text{emission prob.}}$$

Complexity for each step is $O(k^2)$:

for each value of z_{t+1} , we sum over k elements.

Complexity for the recursion is $O(Tk^2)$

Turn to the β recursion, i.e., the backward step.

$$\beta(z_t) = p(x_{t+1:T} | z_t)$$

$$= \sum_{z_{t+1}} p(x_{t+1:T} | z_{t+1}, z_t)$$

$$= \sum_{z_{t+1}} p(x_{t+1:T} | z_{t+1}, z_t) p(z_{t+1} | z_t) \quad [\text{chain rule}]$$

$$= \sum_{z_{t+1}} \underbrace{p(x_{t+2:T} | z_{t+1})}_{\beta(z_{t+1})} \underbrace{p(x_{t+1} | z_{t+1})}_{\text{emission}} \underbrace{p(z_{t+1} | z_t)}_{a_{z_t, z_{t+1}}} \quad \left[\begin{array}{l} x_{t+1} \perp\!\!\!\perp x_{t+2:T} | z_{t+1} \\ x_{t+1:T} \perp\!\!\!\perp z_t | z_{t+1} \end{array} \right]$$

Setting $\beta(z_T) \stackrel{\Delta}{=} 1$ makes the recursion work.

(note $\beta(z_T)$ is meaningless because x_{T+1} does not exist.)

$$\text{or } \beta(z_{T-1}) = p(x_T | z_{T-1})$$

$$= \sum_{z_T} p(x_T | z_T) p(z_T | z_{T-1})$$

Can also compute the predictive dist -

$$p(x_{T+1} | x_{1:T}) = \sum_{z_{T+1}} p(x_{T+1} | z_{T+1}) p(z_{T+1} | x_{1:T})$$

(7)

Many variants of HMMs -

~~Hidden Markov Models~~

Hierarchical HMMs - state transitions at diff. levels

Generalized HMMs - random # of emissions per time point

Shadower - HMM respecting a phylogenetic tree

Factorial HMMs - $\binom{K}{k}$ Factorial states.

"Infinite HMMs" - K is unbounded