

CS181: Bayes Nets - Part III: Hidden Markov Models (Prob)

- Unsupervised Learning: goal model $p(x)$
- Specifically: structure in $p(x)$: x has multiple dims that have interesting conditional independence relationships



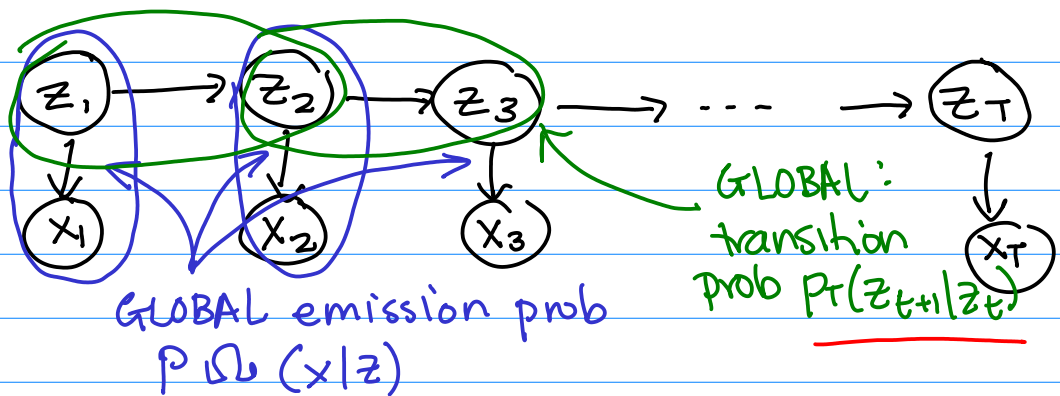
- Noted: relationships can be used to perform inference more efficiently: $p(\text{query}_{\text{var}} | \text{evidence})$
e.g. $p(A|C)$ or $p(B)$

TODAY: Specific type of BN: Hidden Markov Model
(Timeseries model)



(LOCAL)
hidden/
unobserved
structure

observed
vector \vec{x}_n



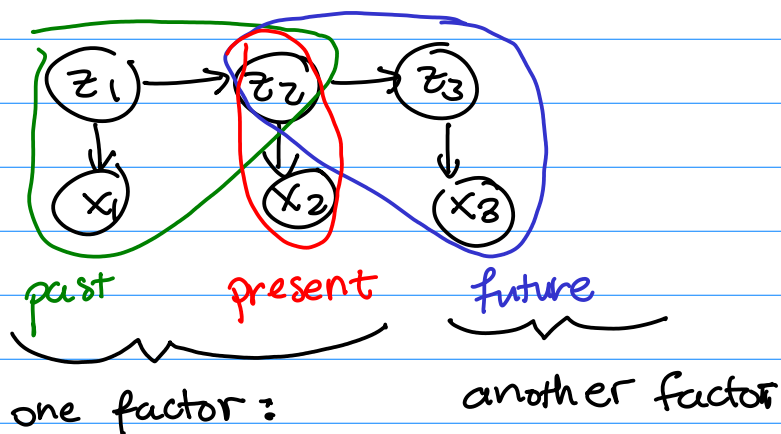
$z =$ hallway	:	$x =$ router1	router2	router3
		.7	.2	.1
rec room	:	.1	.8	.1
dining room	:	.1	.1	.8

Why do we want this structure? What can we do with it?

1. Filtering: where am I now? $p(z_t | x_1 \dots x_t)$
2. Smoothing: where was I at time t ? $p(z_t | x_1 \dots x_T)$
3. Prediction: what's my next meas? $p(x_{t+1} | x_1 \dots x_t)$
4. prob(seq): is this a good model? $p(x_1 \dots x_T)$
5. Best path: what path was most likely taken? $\arg\max_{z_1 \dots z_T} p(z_1 \dots z_T | x_1 \dots x_T)$

→ we've seen this type of problem! all have form $p(-|-)$ that requires various conditionals & marginals of the joint: $p(x_1 \dots x_t, z_1 \dots z_T)$

intuition:



Formalizing:

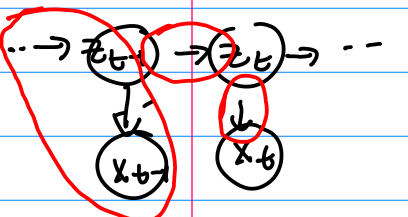
Lets' consider $p(z_t, x_1 \dots x_t)$

(aside: $p(\bar{z}_t^k | x_1 \dots x_t) \propto p(\bar{z}_t^k, x_1 \dots x_t) \quad p(x_1 \dots x_t) =$

$$p(z_t = k | x_1 \dots x_t) = \frac{p(z_t = k, x_1 \dots x_t)}{\sum_j p(z_t = j, x_1 \dots x_t)} \quad \checkmark$$

$$\alpha_t: [K] \rightarrow \alpha_t(z_t = k)$$

$$\alpha_t(z_t) = p(z_t, x_1 \dots x_t) = \sum_{z_{t-1}=1}^K p(z_{t-1}, z_t, x_1 \dots x_t)$$



$$\begin{aligned}
 &= \sum_{z_{t-1}=1}^K p(x_t | z_t) p(z_t | z_{t-1}) p(z_{t-1}, x_1 \dots x_{t-1}) \\
 &= \underbrace{p(x_t | z_t)}_{\substack{K: \\ x_t \text{ obs,} \\ p(x_t | z_t = k)}} \sum_{z_{t-1}=1}^K \underbrace{p(z_t | z_{t-1})}_{\substack{K \times K \\ \text{transition} \\ \text{matrix}}} \underbrace{\alpha_{t-1}(z_{t-1})}_K
 \end{aligned}$$

$$\alpha_t(z_t) = \sum_{z_1} \sum_{z_2} \dots \sum_{z_{t-1}} p(z_1 \dots z_{t-1}, z_t, x_1 \dots x_t)$$

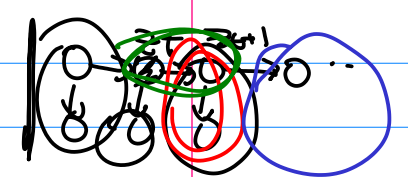
exponentially many combinations

FORWARD PASS
(considers the past)

but given structure, efficient dynamic program

Next: BACKWARD PASS (future)

for smoothing problem: $p(z_t, x_1 \dots x_T)$



$$\begin{aligned}
 &= p(x_1 \dots x_t, z_t) p(x_{t+1} \dots x_T | z_t, x_1 \dots x_t) \\
 &= \underbrace{p(x_1 \dots x_t, z_t)}_{\substack{\text{FORWARD PASS} \\ \alpha}} \underbrace{p(x_{t+1} \dots x_T | z_t)}_{\substack{\text{BACKWARD PASS} \\ \beta}}
 \end{aligned}$$

true in general using cond. ind.

$$\beta_t(z_t) = p(x_{t+1} \dots x_T | z_t)$$

$$\underbrace{\text{vector size } K}_{\substack{\text{vec size} \\ K}} = \sum_{z_{t+1}=1}^K p(x_{t+1} \dots x_T, z_{t+1} | z_t)$$

generally allowed, explicitly marginalizing

$$= \sum \underbrace{p(x_{t+1} | z_{t+1})}_{\text{red}} \underbrace{p(x_{t+2} \dots x_T | z_{t+1})}_{\text{blue}} \underbrace{p(z_{t+1} | z_t)}_{\text{green}} \text{ using str.}$$

$$= \sum_{z_{t+1}=1}^K p(x_{t+1} | z_{t+1}) p(z_{t+1} | z_t) \beta_{t+1}(z_{t+1})$$

convention: $\beta_T(z_T) = \begin{bmatrix} 1 \\ \vdots \end{bmatrix}$

Return to our tasks:

1. filtering: $p(z_t | x_1 \dots x_t) \propto \alpha_t(z_t)$

2. smoothing: $p(z_t | x_1 \dots x_T) \propto \alpha_t(z_t) \beta_t(z_t)$

3. prob seq: $p(x_1 \dots x_T) = \sum_{z_T=1}^K \alpha_T(z_T)$

$p(z_1 \dots z_T, x_1 \dots x_T)$

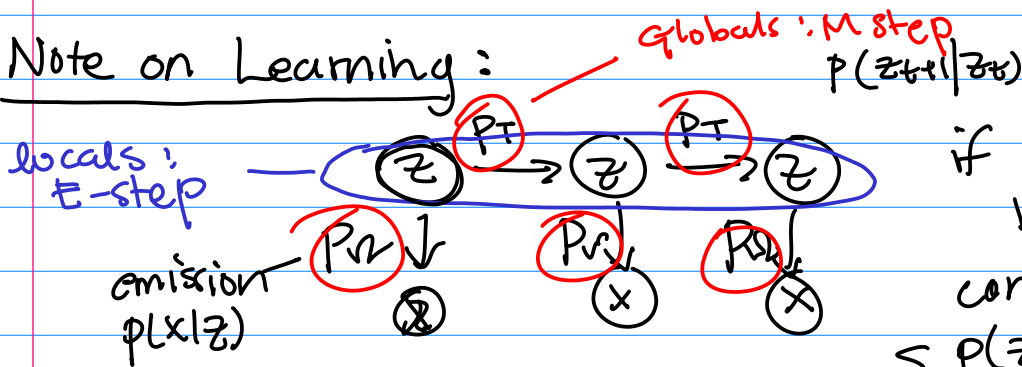
$K^T \dots \sum_{z_1} z_1 z_2 \dots \sum_{z_T} p(z_1 \dots z_T, x_1 \dots x_T)$

(recall: $\alpha_T(z_T) = \frac{p(z_T, x_1 \dots x_T)}{p(x_1 \dots x_T)}$)

4. prediction: $p(x_{t+1} | x_1 \dots x_t)$

$$= \frac{p(x_1 \dots x_{t+1})}{p(x_1 \dots x_t)} = \frac{\sum_{z_{t+1}} \alpha_{t+1}(z_{t+1})}{\sum_{z_t} \alpha_t(z_t)}$$

Note on Learning:



if p_T, p_R are known, we

can compute

$$p(z_t | x_1 \dots x_T)$$

$$p(z_t, z_{t+1} | x_1 \dots x_T)$$

if we have these, we can update p_T, p_R

(recall: EM opt. lower bound on $p(x)$)

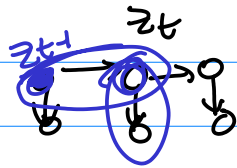
Note on finding the most likely path:

$$\delta_t(z_t) = \max_{z_1 \dots z_{t-1}} p(x_1 \dots x_t, \underbrace{z_1 \dots z_{t-1}}_{\text{vec of length } k}, \underbrace{z_t}_{\text{length } k})$$

vec of
length k

→ intuitively: ~~what~~ if we end up at $z_t = k$, what was the most likely path there?

$$\delta_t(z_t) = \max_{z_{t-1}} \delta_{t-1}(z_{t-1}) p(z_t | z_{t-1}) p(x_t | z_t)$$



could you
get to $z_{t-1} = k$
in a likely
way

could
you
transition
to z_t

would
you have
seen
 x_t ?

