

A.I. in health informatics

lecture 8 structured learning

kevin small &
byron wallace

today

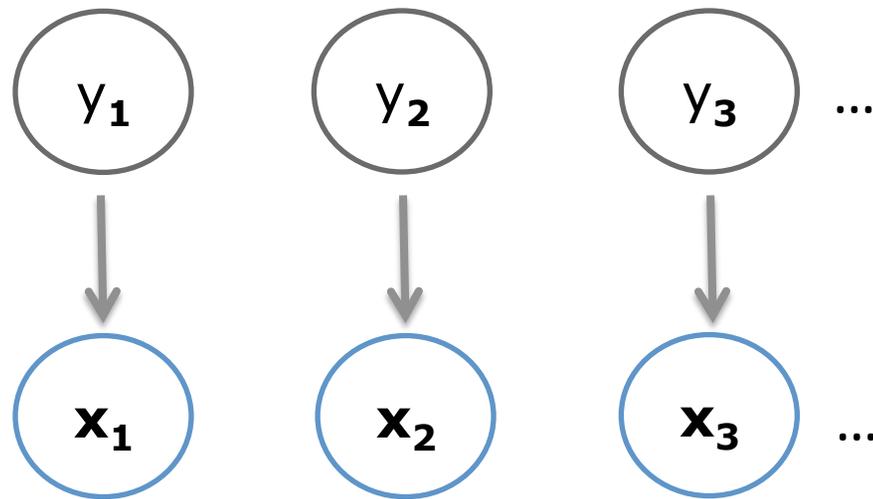
- models for structured learning: HMMs and CRFs
 - structured learning is particularly useful in biomedical applications: parsing (clinical) text; genetic data, etc.
 - we'll cover this in more detail; but need the basics first

unstructured learning

assumptions:

- we're given a set of i.i.d. instances $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ and their (univariate) labels $\{y_1, y_2, \dots, y_N\}$
- no order or *sequence* to the data

unstructured learning: graphically



structured learning

assumptions:

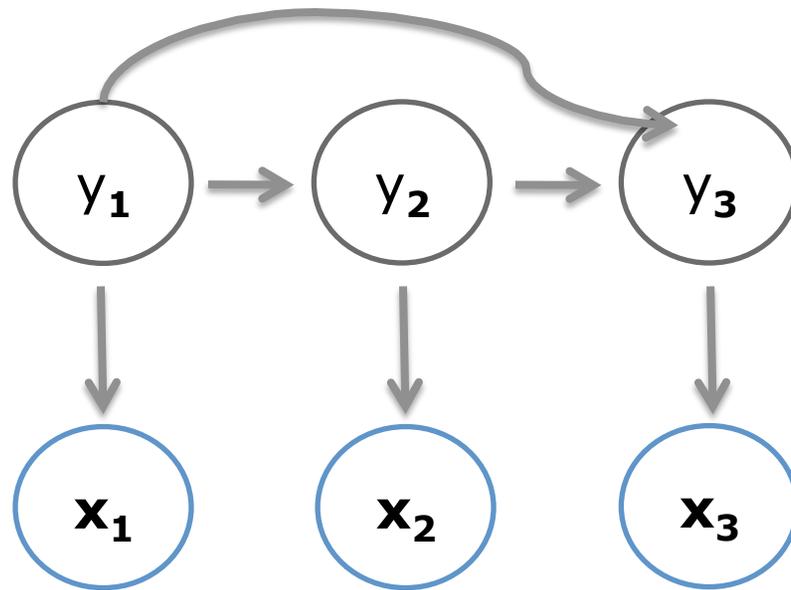
- there is some correlation between a label \mathbf{y}_i and the *preceding* labels, in other words,
- \mathbf{y}_i is **structured**, ie., y_{i+1} is affected by the previous labels $y_i, y_{i-1} \dots$

structured learning



- consider the task of part-of-speech (POS) tagging sentences
 - *nouns* tend to follow *verbs*

structured learning: graphically



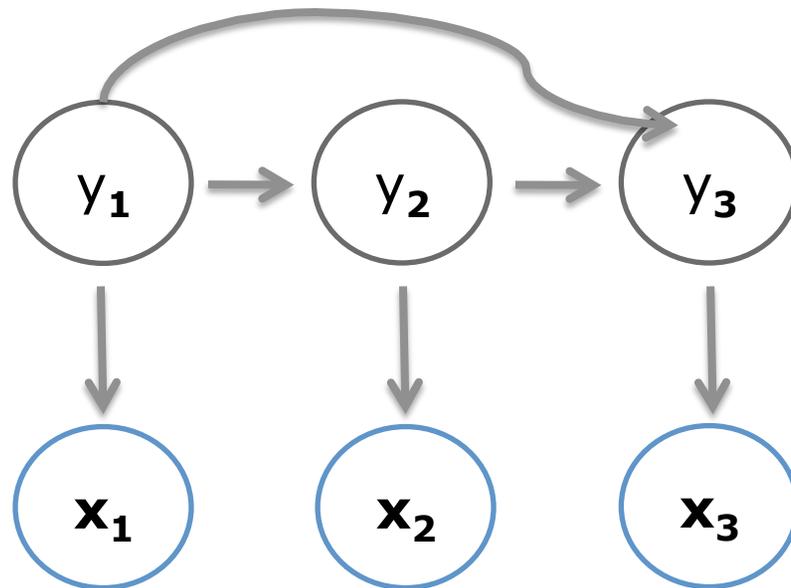
probability & structured learning

$$p(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$



intractable!

(first-order) HMM



* usually assume states (y_s) are *latent*; but not always (see weather example, upcoming)

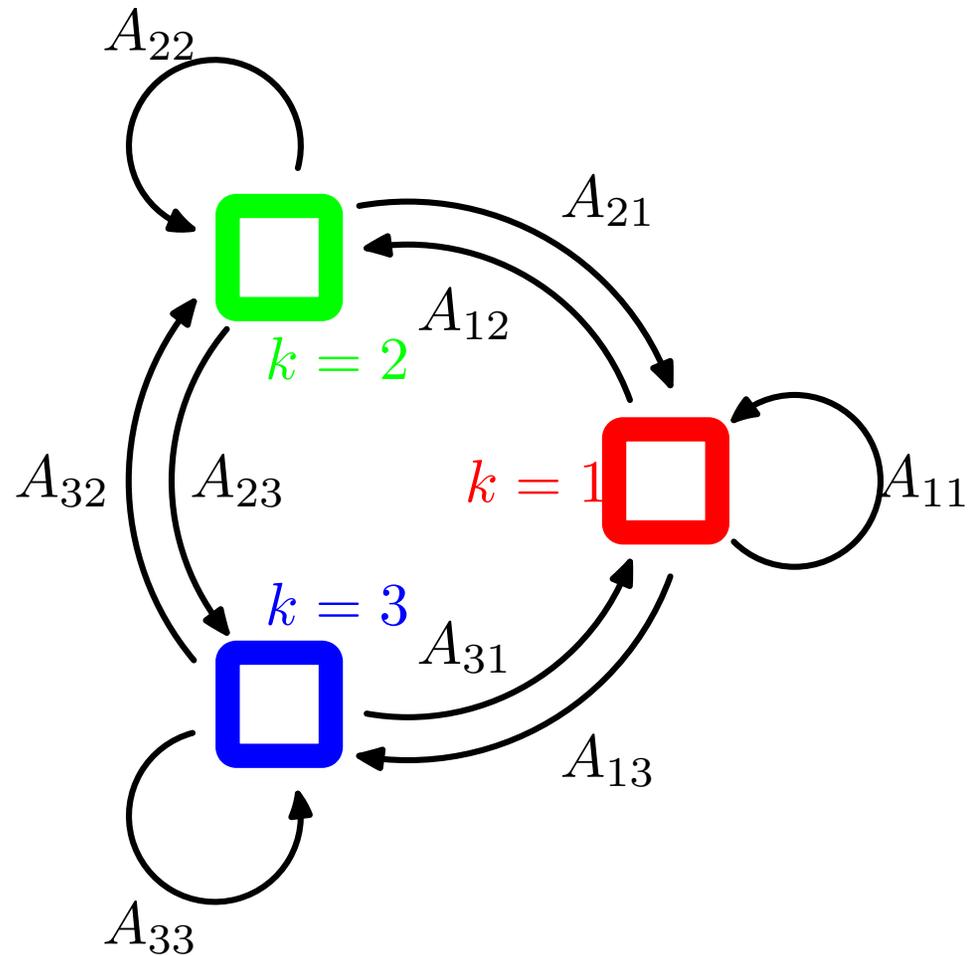
(first-order) MM

$$p(x_1, x_2, \dots, x_n) = p(x_1) \prod_{i=2}^n p(x_i | x_{i-1})$$

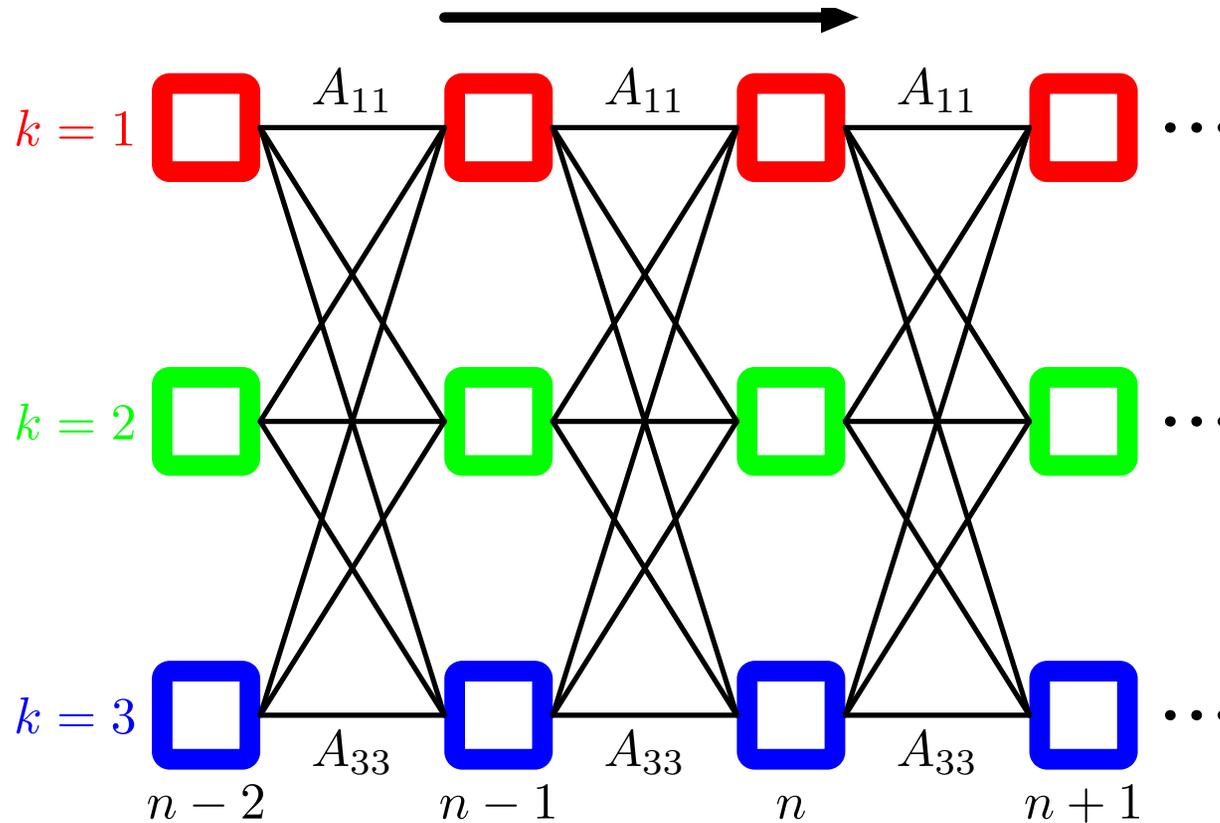
tractable!



markov model



markov model: unfolded



markov model

- \mathbf{a}_{ij} – probability of transitioning from state i to j
 - $\sum_j a_{ij} = 1$ (we have to go somewhere!)
 - $a_{ij} \geq 0$

let's talk about the weather

- the world has three states:

1 rainy, 2 cloudy, 3 sunny

- (in the weather case the markov model is not hidden)

- our transition matrix is:

$$A = \{a_{ij}\} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

note that this system
likes to stay where it is!

the weather

- probability that of observing the weather
{sunny, sunny, rainy}?

$$A = \{a_{ij}\} = \begin{matrix} & \begin{matrix} R & C & S \end{matrix} \\ \begin{matrix} - \\ \\ \\ \end{matrix} & \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix} \end{matrix}$$
$$= 1.0 * P(S|S)*P(R|S) = .8*.1$$

the weather

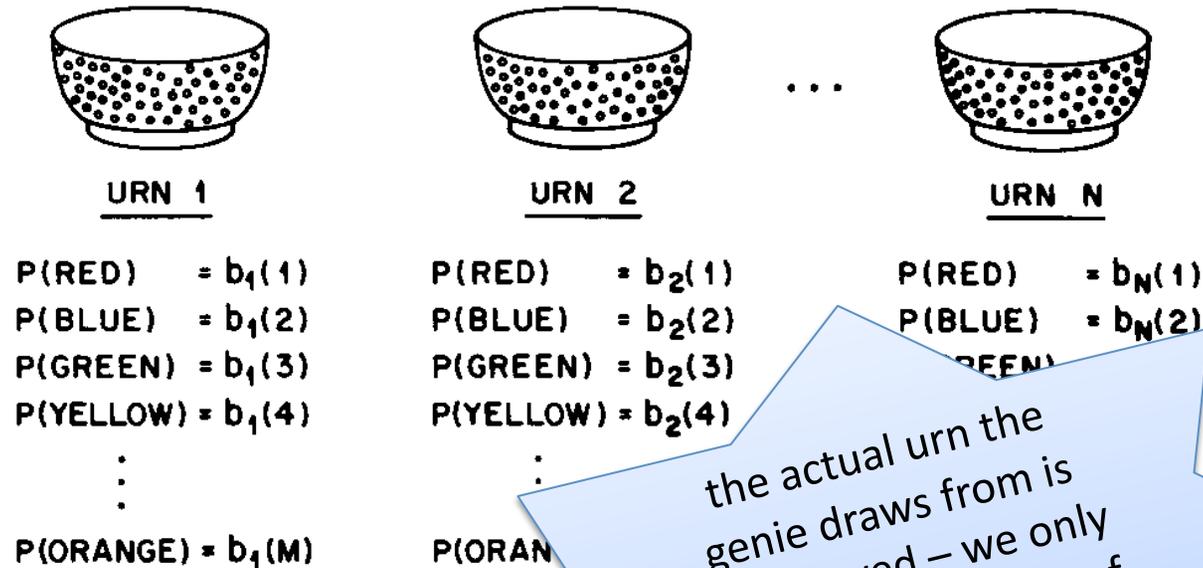
- it's sunny today. what's the probability that it remains so for k days?

$$A = \{a_{ij}\} = \begin{matrix} & \begin{matrix} R & C & S \end{matrix} \\ \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix} \end{matrix}$$
$$= (.8)^k$$

hidden markov models

- in weather example, we only care about *transitions*
 - the (weather) states were the observations
- often we need to model data in which an observation is generated conditional on some latent, underlying state
 - e.g., bias coin example; urn-genie example
- enter the HMM

the urn example



$O = \{\text{GREEN, GREEN, BLUE, RED}\}$

the actual urn the genie draws from is unobserved – we only know the sequence of draws!

a **genie (!!!)** is in the room choosing which urn to draw balls from – each urn contains different proportions of the various colored balls

hidden markov models: sufficient parameters

N states (latent), M symbols (observed): symbols are observed conditioned on the current state

A - transition probabilities (from urn to urn)

B - symbol emission probabilities (color proportions in each urn)

π - initial state distribution (initial urn likelihood)

$\lambda = (A, B, \pi)$ specifies our model

hidden markov models: the three problems

- 1** given a set of observations $\mathbf{o} = \{o_1, o_2 \dots o_T\}$ and a model λ , compute $P(\mathbf{o}|\lambda)$
- 2** given \mathbf{o} , λ , calculate most likely latent states (usually thought of as labels) $\mathbf{q} = \{q_1, q_2 \dots q_T\}$
- 3** given \mathbf{o} , calculate λ

hidden markov models: problem 1

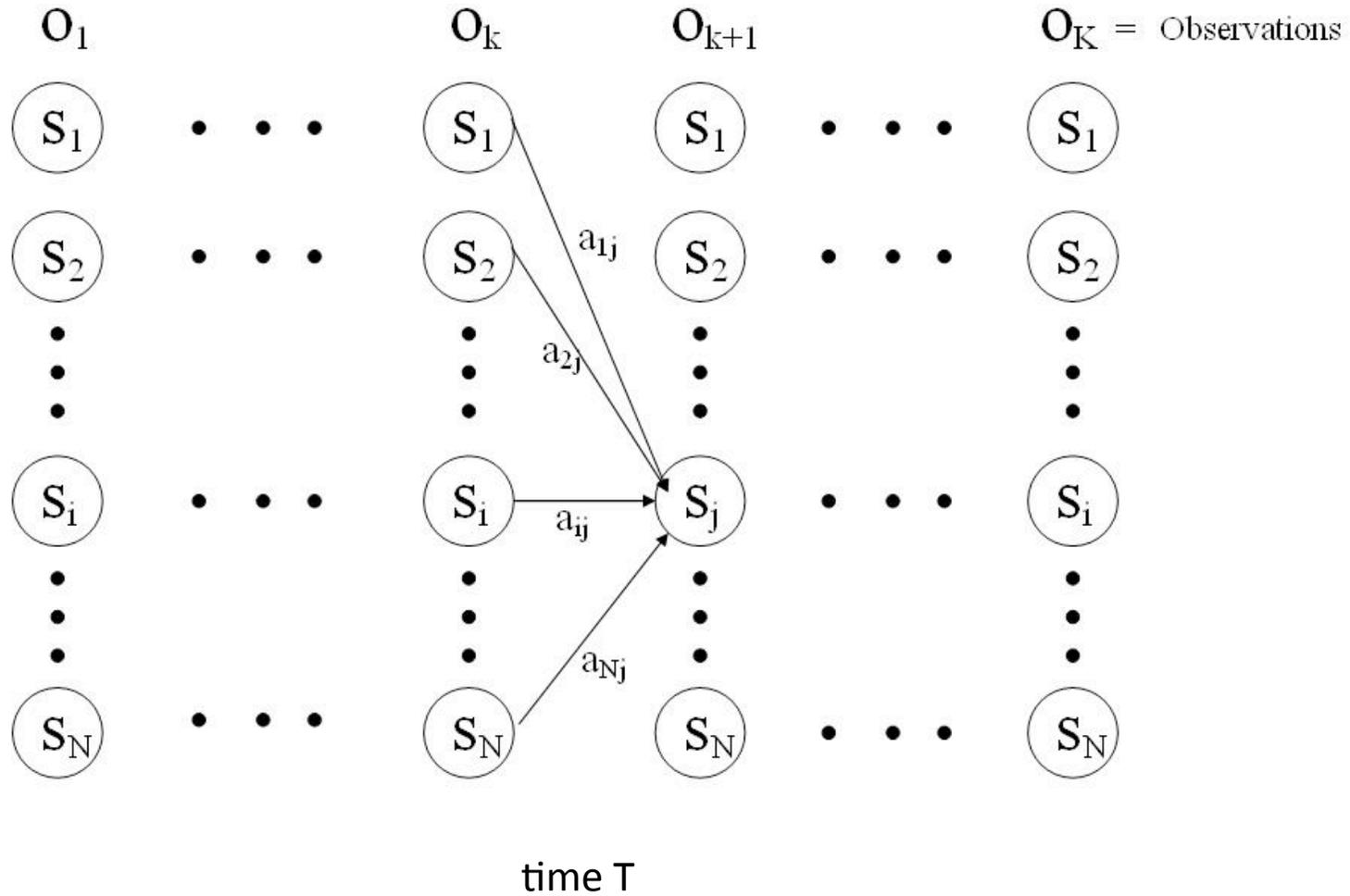
1 given a set of observations $\mathbf{o} = \{o_1, o_2, \dots, o_T\}$ and a model λ , compute $P(\mathbf{o} | \lambda)$

$$P(\mathbf{o} | \lambda) = \sum_{\text{all } \mathbf{q}} P(\mathbf{o} | \mathbf{q}, \lambda) P(\mathbf{q} | \lambda) =$$

... cool. so we're done?

not quite. this will require $O(N^T)$ calculations

dynamic programming to the rescue!



dynamic programming to the rescue!

1) Initialization:

$$\alpha_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N.$$

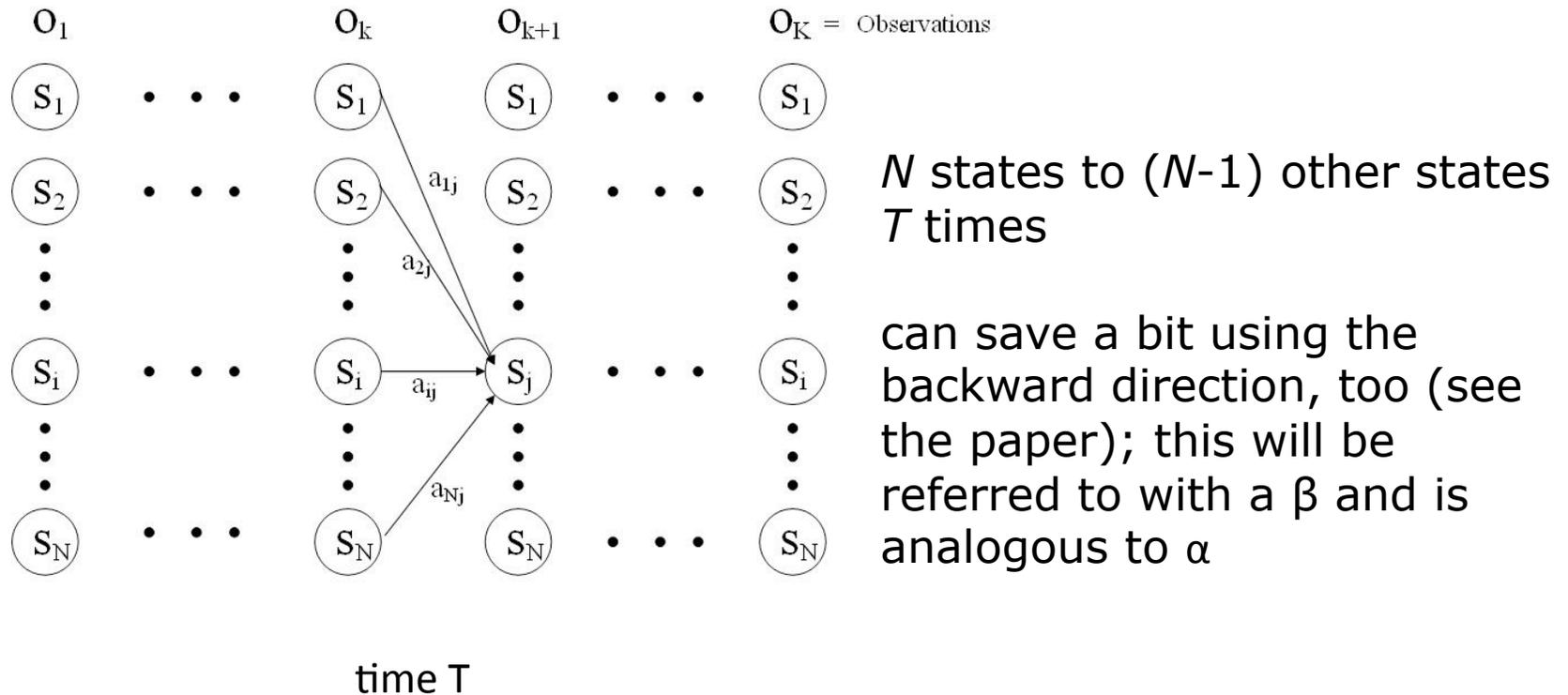
2) Induction:

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1}), \quad 1 \leq t \leq T - 1$$
$$1 \leq j \leq N.$$

3) Termination:

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i).$$

about that runtime



hidden markov models: the three problems

- 1** given a set of observations $\mathbf{o} = \{o_1, o_2 \dots o_T\}$ and a model λ , compute $P(\mathbf{o}|\lambda)$
- 2** given \mathbf{o} , λ , calculate most likely latent states (usually thought of as labels) $\mathbf{q} = \{q_1, q_2 \dots q_T\}$
- 3** given \mathbf{o} , calculate λ

hidden markov models: problem 2

- 2** given a set of observations $\mathbf{o} = \{o_1, o_2 \dots o_T\}$
and a model λ , find most likely latent states
(urns, say) $\mathbf{q}^* = \{q_1, q_2 \dots q_T\}$
- a unique solution need not exist! what are
we even optimizing here?
 - individual most likely states q_i ?
 - joint probability \mathbf{q} ?

hidden markov models: problem 2

problem optimizing for the most likely state at *any given time* can lead to *impossible* sequences under λ

so let's consider the joint instead

solving problem 2; the joint solution

want to solve for \mathbf{q}^* :

$$\mathbf{q}^* \leftarrow \underset{\mathbf{q}}{\operatorname{argmax}} P(\mathbf{q}, 0 \mid \lambda)$$

- obviously enumerating all possible \mathbf{q} sequences is infeasible
- dynamic programming to the rescue, again!

the viterbi algorithm

define:

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P[q_1 q_2 \dots q_t = i, O_1 O_2 \dots O_t | \lambda]$$

(most likely sequence up to time t-1). the inductive step

$$\delta_{t+1}(j) = [\max_i \delta_t(i) a_{ij}] \cdot b_j(O_{t+1})$$

i.e., from the best path so far, we find the most likely transition/emission probability



that's it! we just need to keep track of the states we visit!

hidden markov models: the three problems

- 1** given a set of observations $\mathbf{o} = \{o_1, o_2 \dots o_T\}$ and a model λ , compute $P(\mathbf{o}|\lambda)$
- 2** given \mathbf{o} , λ , calculate most likely latent states (usually thought of as labels) $\mathbf{q} = \{q_1, q_2 \dots q_T\}$
- 3** given \mathbf{o} , calculate λ

EM for λ

define

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O|\lambda)}$$

emit observation in state j

forward to state i

backward to state j

from i to j

(the probability of being in state i at time t and state j at time $t+1$)

EM for λ

probability of being in state i is:

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j).$$

(we have to transition *somewhere*)

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions from } S_i$$

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{expected number of transitions from } S_i \text{ to } S_j.$$

EM for λ

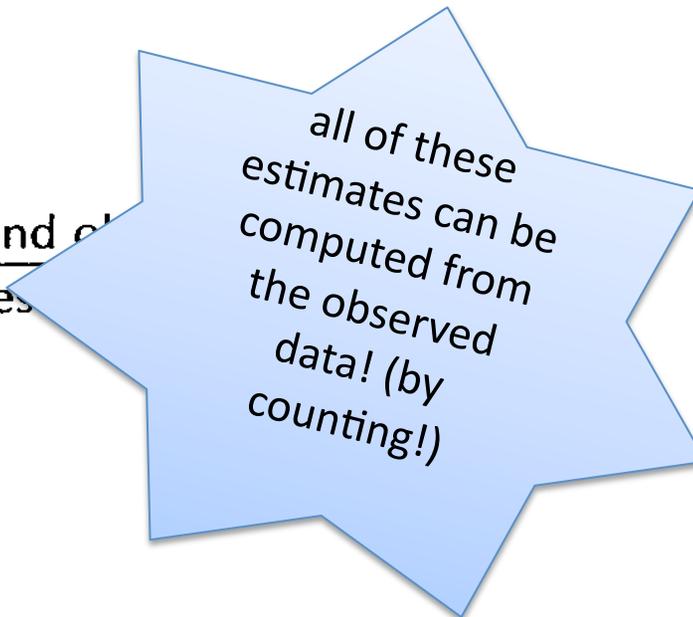
$\bar{\pi}_i$ = expected frequency (number of times) in state S_i at time $(t = 1) = \gamma_1(i)$

\bar{a}_{ij} = $\frac{\text{expected number of transitions from state } S_i \text{ to state } S_j}{\text{expected number of transitions from state } S_i}$

$$= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$\bar{b}_j(k)$ = $\frac{\text{expected number of times in state } j \text{ and observed } v_k}{\text{expected number of times in state } j}$

$$= \frac{\sum_{t=1}^T \gamma_t(j) \cdot \mathbb{1}_{\{O_t = v_k\}}}{\sum_{t=1}^T \gamma_t(j)}$$



all of these estimates can be computed from the observed data! (by counting!)

EM for λ

at a given time t we have

$$\bar{\lambda} = (\bar{A}, \bar{B}, \bar{\pi})$$

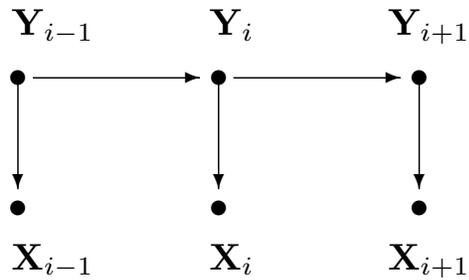
E-step calculate the likelihood of our observations \bullet
using the current estimates

M-step re-estimate the parameters (left-hand sides of
equations on previous slide) using current estimates

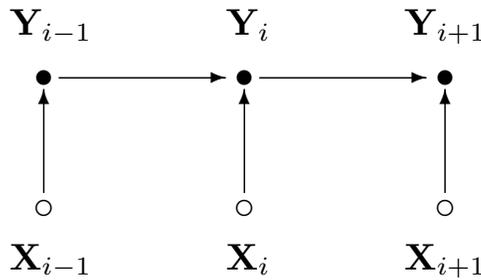
shortcomings of HMMs

- HMMs model the *joint* probability of the observations (\mathbf{x}) and (\mathbf{y})
 - ie., it's a generative model
- but what we *really* care about is the *conditional probability* of a label sequence \mathbf{y} given \mathbf{x}
- enter *conditional markov models*

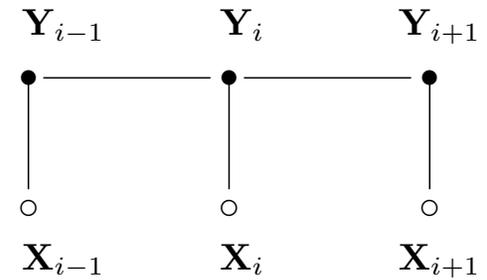
conditional models



HMM



MEMM

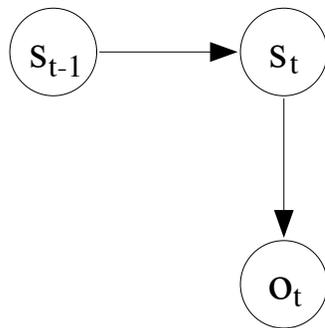


CRF

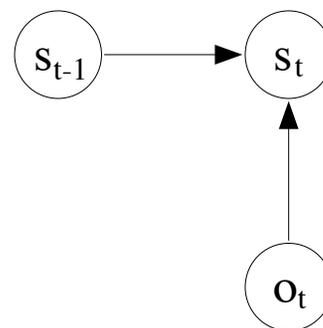
conditional models don't waste time modeling the observed data

MEMMs (2000, McCallum et al)

- MEMM: Maximum Entropy Markov Model
- output probability vector of transitioning from one state to all other states, given an input observation \mathbf{x}
 - *conditional* in that we estimate $p(\text{state})$ at time t given an observation and the preceding state



HMM



MEMM

MEMMs v. HMM

- HMM prob. of emitting \mathbf{o} and being in state s at time t

$$\alpha_{t+1}(s) = \sum_{s' \in S} \alpha_t(s') P(s|s') P(o_{t+1}|s)$$

- MEMM

$$\alpha_{t+1}(s) = \sum_{s' \in S} \alpha_t(s') P_{s'}(s|o_{t+1})$$

CRFs (2001, Lafferty et al)

- CRF: Conditional Random Field
- single model for the joint probability of the sequence of labels given the observations
- mitigates the *label bias* problem in which states with low-entropy (almost certain) transition probability vectors effectively ignore the observation

