

# Machine Learning

## Hidden Markov Models

Mannes Poel

Gwenn Englebienne

# Introduction

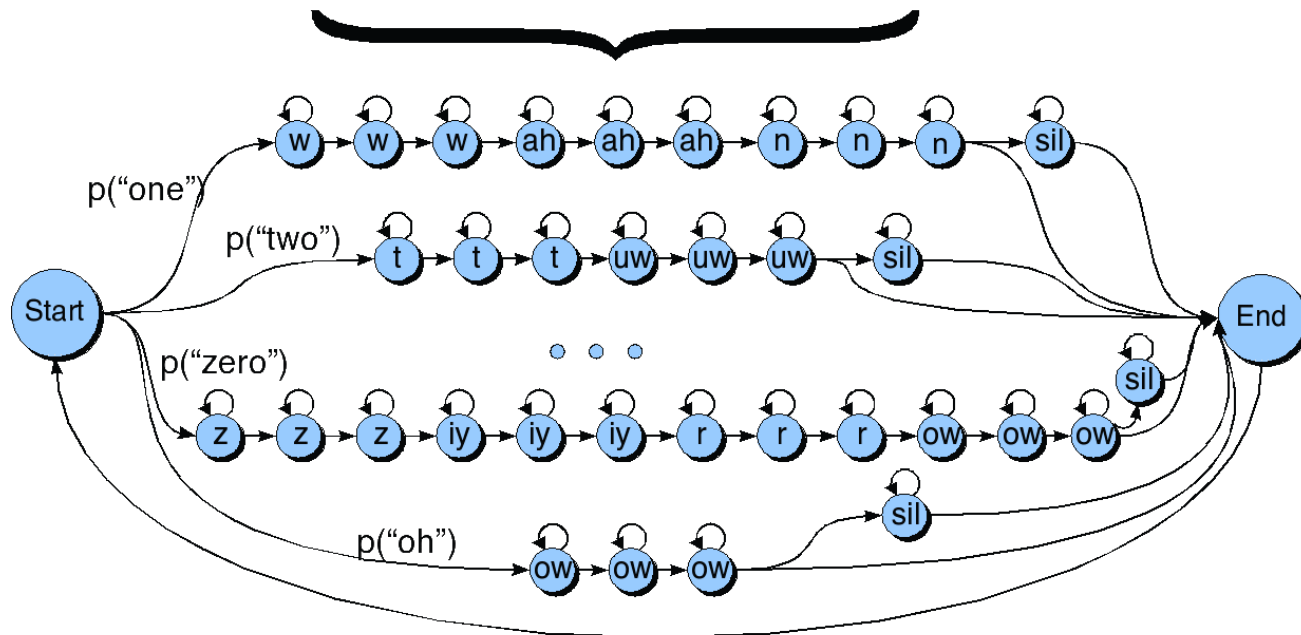
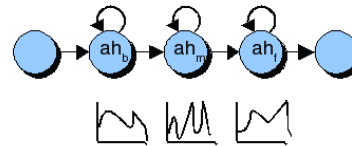
- Modeling dependencies in input; no longer iid
- Sequences:
  - ***Temporal:***
    - ***In speech; phonemes in a word (dictionary), words in a sentence (syntax, semantics of the language).***
    - ***In handwriting, pen movements, gesture recognition.***
    - ***In brain signals; temporal aspects such as a P300***
  - ***Spatial: In a DNA sequence; base pairs***

# Speech detection

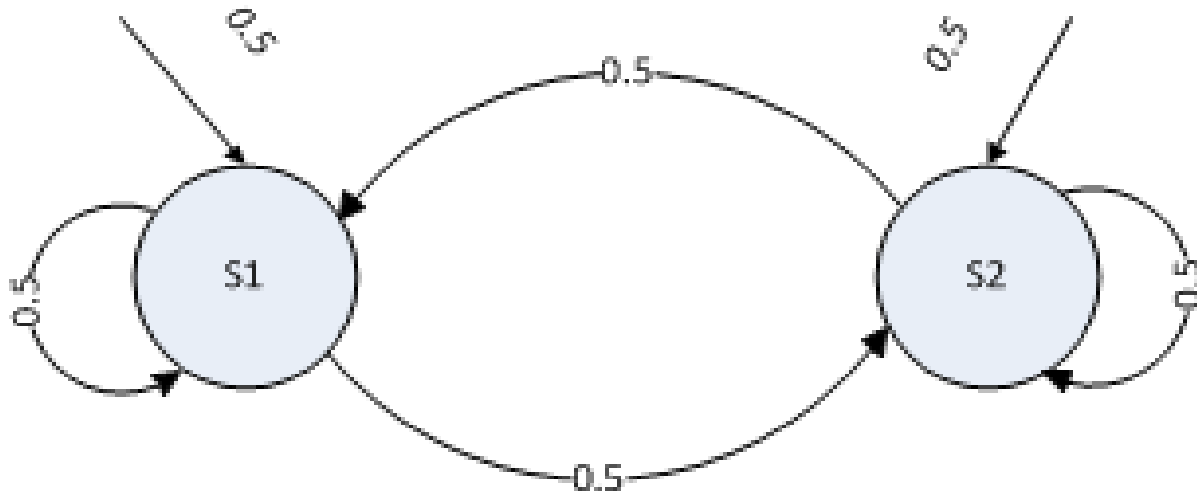
Lexicon

one	w ah n
two	t uw
three	th r iy
four	f ao r
five	f ay v
six	s ih k s
seven	s eh v ax n
eight	ey t
nine	n ay n
zero	z iy r ow
oh	ow

Phone HMM



# Probabilistic Transition Models, Stochastic Automaton, Markov Models



- Gives rise to state sequences of the form:

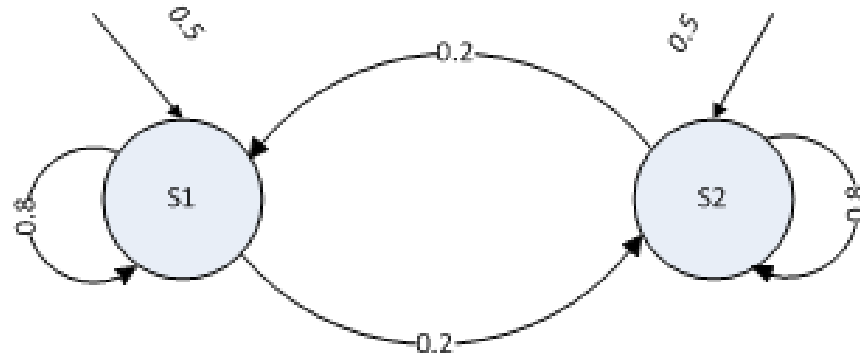
***S1S1S2S2S1S1S2S1***

# Quiz

Given the following state sequence  $\mathbf{s}$

**$\mathbf{s} = S1S2S1S2S1S2S1S2$**

and stochastic automaton  $\mathbf{M}$

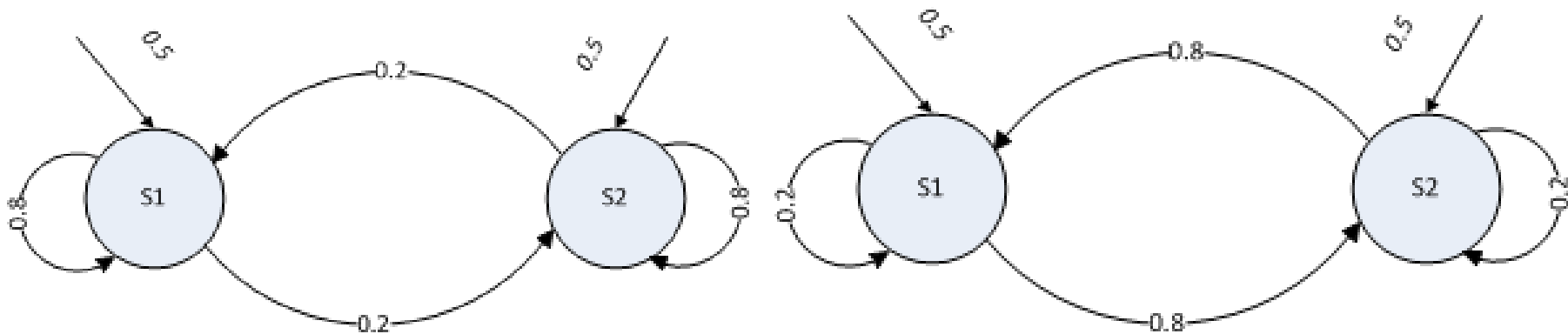


What is the likelihood that  $\mathbf{s}$  is generated by  $\mathbf{M}$ ; how to compute  $P(\mathbf{s}/\mathbf{M})$ ?

# Quiz

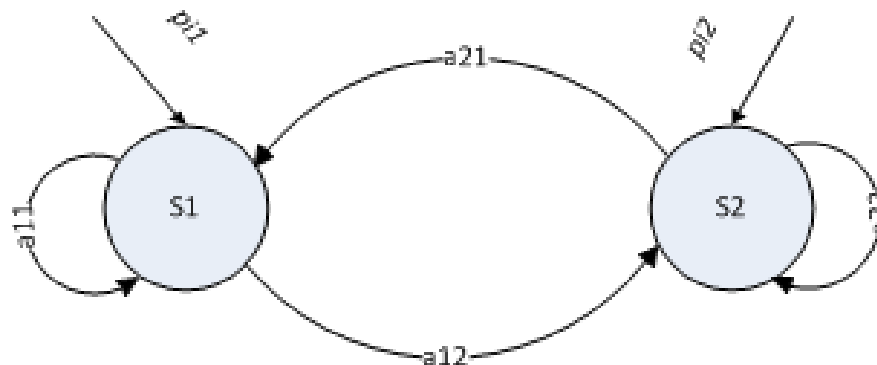
- Given the following state sequence  $\mathbf{s}$  which model is more likely to generate such a sequence; i.e.  $\max_M P(\mathbf{s}/M)$ ?

**$\mathbf{s} = S1S2S1S2S1S2S1S2$**



# Quiz

- Given the sequences:
  - S1S2S2S1S1S1S2S2S1S1S2
  - S1S1S1S2S2S1S1S2S2S1S1
  - S2S2S2S1S1S2S1S2S2S1S1
- How to estimate the probabilities?



# Discrete Markov Process

- $N$  states:  $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_N$ 
  - State at “time”  $t$ ,  $\mathbf{q}_t = \mathbf{S}_i$

- First-order Markov

$$P(\mathbf{q}_{t+1} = \mathbf{S}_j \mid \mathbf{q}_t = \mathbf{S}_i, \mathbf{q}_{t-1} = \mathbf{S}_k, \dots) = P(\mathbf{q}_{t+1} = \mathbf{S}_j \mid \mathbf{q}_t = \mathbf{S}_i)$$

- Transition probabilities

$$\mathbf{a}_{ij} \equiv P(\mathbf{q}_{t+1} = \mathbf{S}_j \mid \mathbf{q}_t = \mathbf{S}_i) \quad \mathbf{a}_{ij} \geq 0 \text{ and } \sum_{j=1}^N \mathbf{a}_{ij} = 1$$

- Initial probabilities

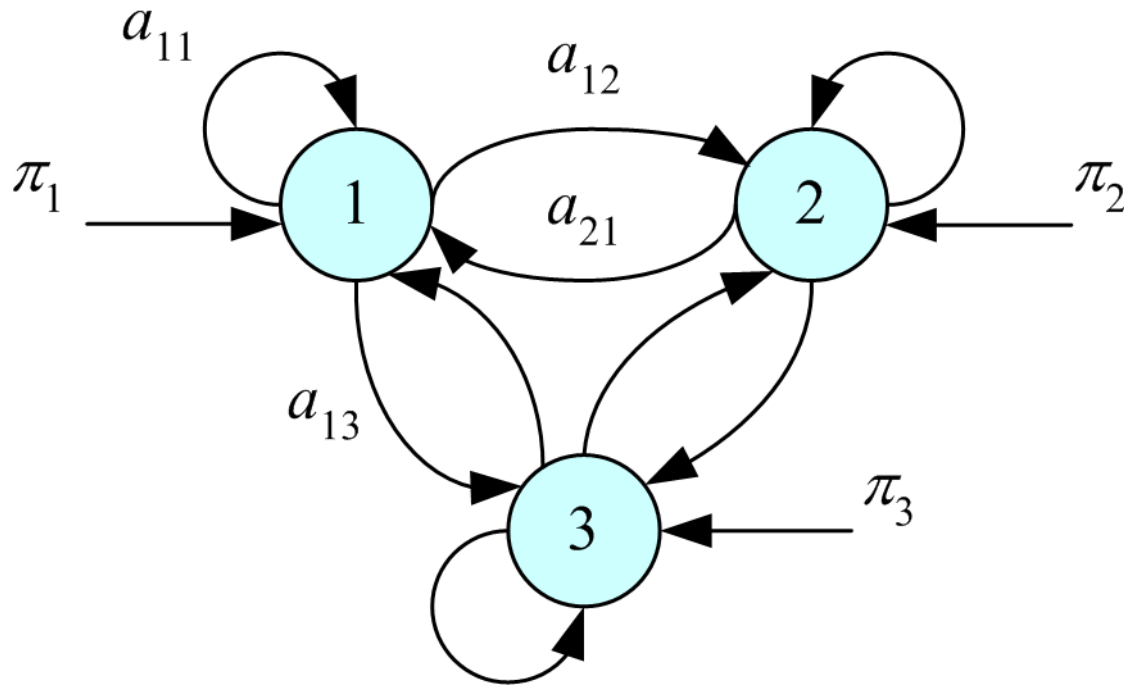
$$\pi_i \equiv P(\mathbf{q}_1 = \mathbf{S}_i) \quad \sum_{i=1}^N \pi_i = 1$$



# Example of 2 coins

- Recall the estimation of coin probabilities.  
Can be modeled as an Hidden Markov Model

# Stochastic Automaton



$$P(O = Q \mid \mathbf{A}, \Pi) = P(q_1) \prod_{t=2}^T P(q_t \mid q_{t-1}) = \pi_{q_1} a_{q_1 q_2} \cdots a_{q_{T-1} q_T}$$

# Example: Balls and Urns

- Three urns each full of balls of one color  
 $\mathbf{S}_1$ : red,  $\mathbf{S}_2$ : blue,  $\mathbf{S}_3$ : green

$$\Pi = [0.5, 0.2, 0.3]^T \quad \mathbf{A} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

$$O = \{S_1, S_1, S_3, S_3\}$$

$$P(O \mid \mathbf{A}, \Pi) = P(S_1) \cdot P(S_1 \mid S_1) \cdot P(S_3 \mid S_1) \cdot P(S_3 \mid S_3)$$

$$= \pi_1 \cdot a_{11} \cdot a_{13} \cdot a_{33}$$

$$= 0.5 \cdot 0.4 \cdot 0.3 \cdot 0.8 = 0.048$$

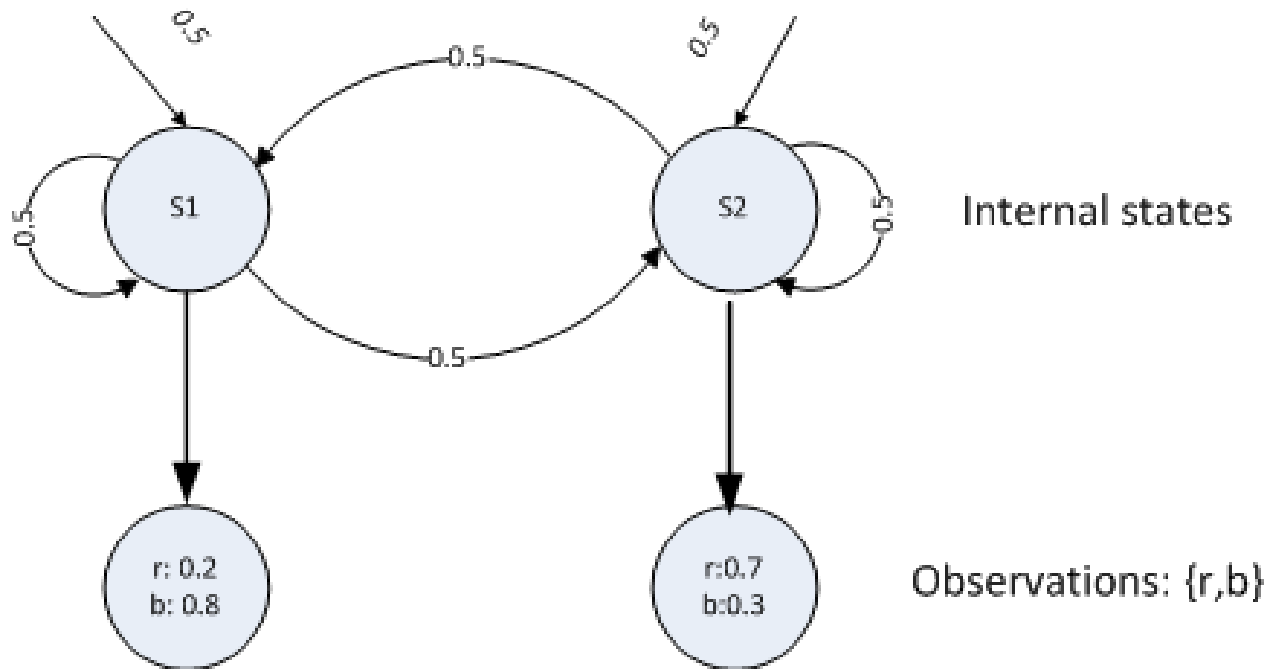
# Balls and Urns: Learning

- Given  **$K$**  example sequences of length  **$T$**

$$\hat{\pi}_i = \frac{\# \{ \text{sequences starting with } S_i \}}{\# \{ \text{sequences} \}} = \frac{\sum_k 1(q_1^k = S_i)}{K}$$

$$\begin{aligned} \hat{a}_{ij} &= \frac{\# \{ \text{transition s from } S_i \text{ to } S_j \}}{\# \{ \text{transition s from } S_i \}} \\ &= \frac{\sum_k \sum_{t=1}^{T-1} 1(q_t^k = S_i \text{ and } q_{t+1}^k = S_j)}{\sum_k \sum_{t=1}^{T-1} 1(q_t^k = S_i)} \end{aligned}$$

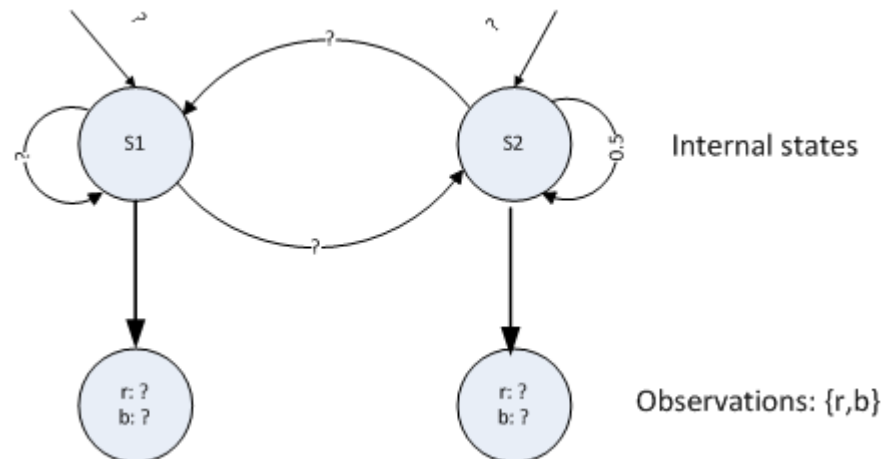
# Hidden Markov Models



- Quiz: How to compute the likelihood that the sequence ***rrrbbrbbrb*** is generated by this HMM?

# Quiz

- Given the sequences:
  - rrbbrrrbbrbrbrb
  - bbbbrrrrbbrbbbbrrbr
  - rbrbrrbbrbrbrbrbrbbr
- ***How to estimate the probabilities?***



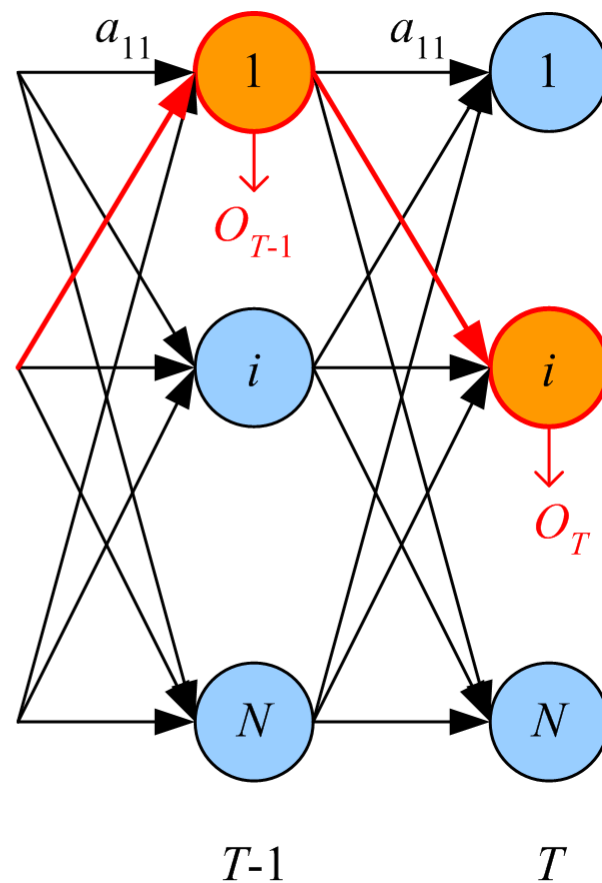
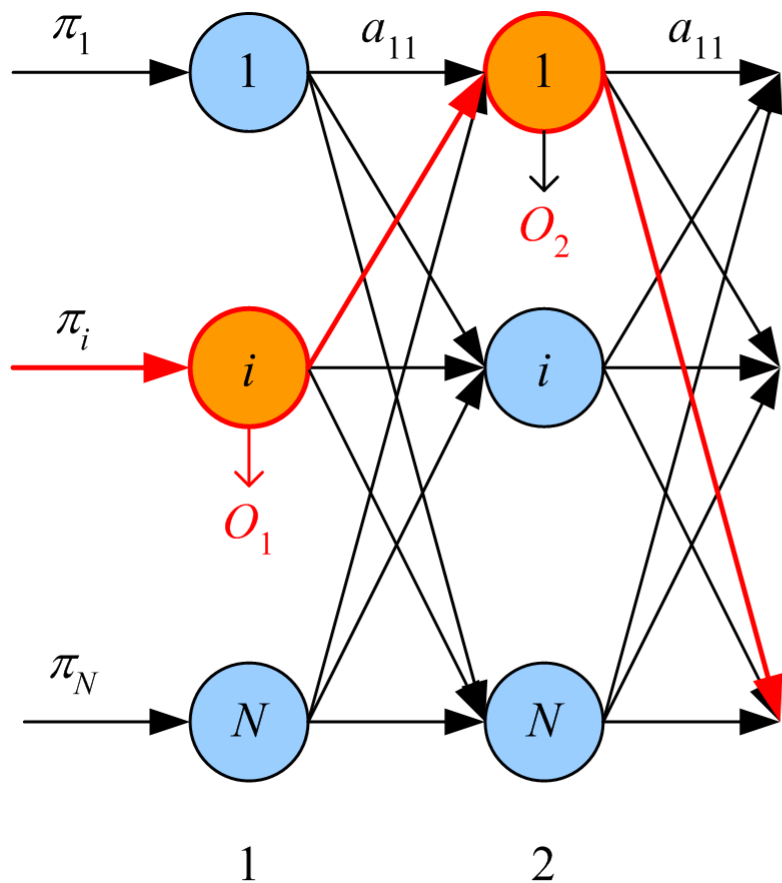
# Hidden Markov Models

- States are not observable
- Discrete observations  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_M\}$  are recorded; a probabilistic function of the state
- Emission probabilities

$$b_j(\mathbf{m}) \equiv P(\mathbf{O}_t = \mathbf{v}_m \mid \mathbf{q}_t = \mathbf{S}_j)$$

- Example: In each urn, there are balls of different colors, but with different probabilities.
- For each observation sequence, there are multiple state sequences

# HMM Unfolded in Time





# Elements of an HMM

- **$N$** : Number of states
- **$M$** : Number of observation symbols
- **$A = [a_{ij}]$** :  **$N$**  by  **$N$**  state transition probability matrix
- **$B = b_j(m)$** :  **$N$**  by  **$M$**  observation probability matrix
- **$\Pi = [\pi_i]$** :  **$N$**  by 1 initial state probability vector

**$\lambda = (A, B, \Pi)$ , parameter set of HMM**

# Three Basic Problems of HMMs

1. Evaluation: Given  $\lambda$ , and  $\mathbf{O}$ , calculate  $P(\mathbf{O} | \lambda)$
2. State sequence: Given  $\lambda$ , and  $\mathbf{O}$ , find  $\mathbf{Q}^*$  such that

$$P(\mathbf{Q}^* | \mathbf{O}, \lambda) = \max_{\mathbf{Q}} P(\mathbf{Q} | \mathbf{O}, \lambda)$$

3. *Learning: Given  $\mathcal{X}=\{\mathbf{O}^k\}_K$ , find  $\lambda^*$  such that*

$$P(\mathcal{X} | \lambda^*) = \max_{\lambda} P(\mathcal{X} | \lambda)$$

(Rabiner, 1989)

# Calculation of $P(O | \lambda)$

- Marginalize over all state sequences of length T

$$P(O | \lambda) = \sum_{all Q} P(O, Q | \lambda) = \sum_{all Q} P(O | Q, \lambda) P(Q | \lambda)$$

$$P(O | Q, \lambda) = b_{q_1}(O_1) b_{q_2}(O_2) \dots b_{q_T}(O_T)$$

$$P(Q | \lambda) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \dots a_{q_{T-1} q_T}$$

$$P(O | Q, \lambda) P(Q | \lambda) = \pi_{q_1} b_{q_1}(O_1) a_{q_1 q_2} b_{q_2}(O_2) \dots a_{q_{T-1} q_T} b_{q_T}(O_T)$$

- Complexity:  $N^T$

# Efficient calculation of $P(\mathbf{O} \mid \boldsymbol{\lambda})$

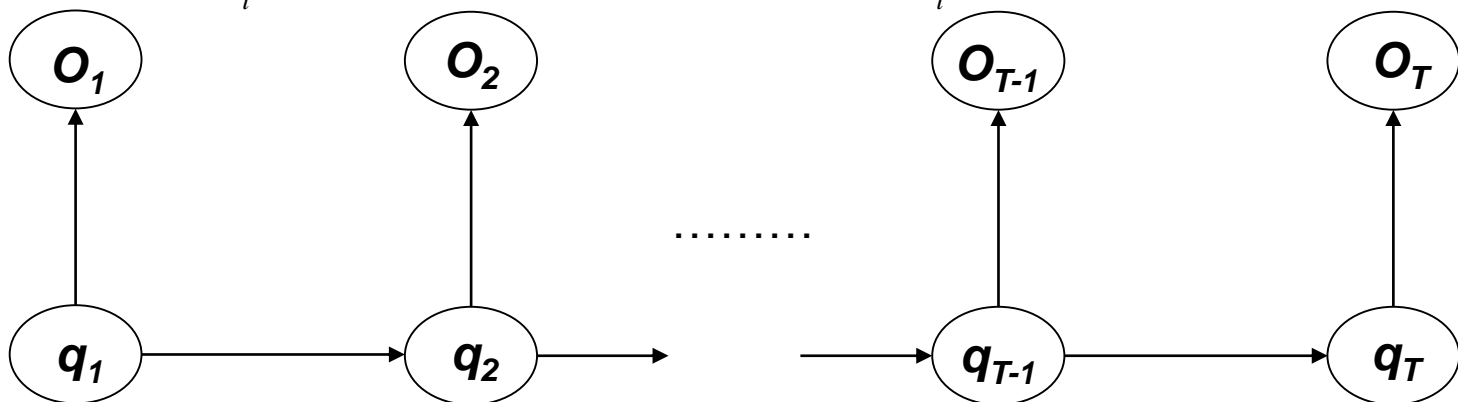
- Marginalize over last state:

$$P(\mathbf{O} \mid \boldsymbol{\lambda}) = \sum_j P(\mathbf{O}, q_T = S_j \mid \boldsymbol{\lambda}) = \sum_j \alpha_T(j)$$

- Left to compute  $\alpha_T(j)$

# Calculation of $P(\mathbf{O} | \boldsymbol{\lambda})$

$$\begin{aligned}
 \alpha_T(j) &= P(O_1 \dots O_T, q_T = S_j) = P(O_1 \dots O_T | q_T = S_j) P(q_T = S_j) \\
 &= P(O_1 \dots O_{T-1} | q_T = S_j) P(O_T | q_T = S_j) P(q_T = S_j) \\
 &= P(O_1 \dots O_{T-1}, q_T = S_j) P(O_T | q_T = S_j) \\
 &= \sum_i P(O_1 \dots O_{T-1}, q_T = S_j, q_{T-1} = S_i) b_j(O_T) \\
 &= \sum_i P(O_1 \dots O_{T-1}, q_{T-1} = S_i, q_T = S_j) b_j(O_T) \\
 &= \sum_i P(O_1 \dots O_{T-1}, q_T = S_j | q_{T-1} = S_i) P(q_{T-1} = S_i) b_j(O_T) \\
 &= \sum_i P(O_1 \dots O_{T-1} | q_{T-1} = S_i) P(q_{T-1} = S_i) P(q_T = S_j | q_{T-1} = S_i) b_j(O_T) \\
 &= \sum_i P(O_1 \dots O_{T-1}, q_{T-1} = S_i) a_{ij} b_j(O_T) = \sum_i \alpha_{T-1}(i) a_{ij} b_j(O_T)
 \end{aligned}$$



# Calculation of $P(\mathbf{O} | \lambda)$

- Forward variable:

$$\alpha_t(i) \equiv P(O_1 \cdots O_t, q_t = S_i | \lambda)$$

Initialization :

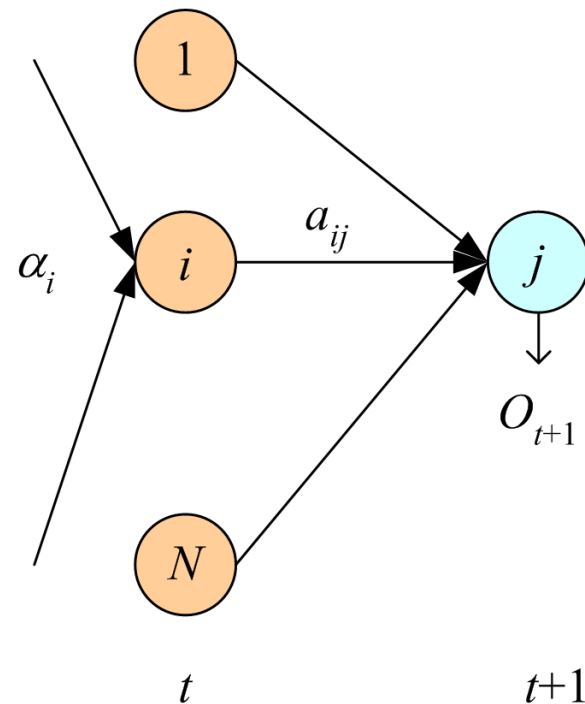
$$\alpha_1(i) = \pi_i b_i(O_1)$$

Recursion :

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

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

- Complexity?

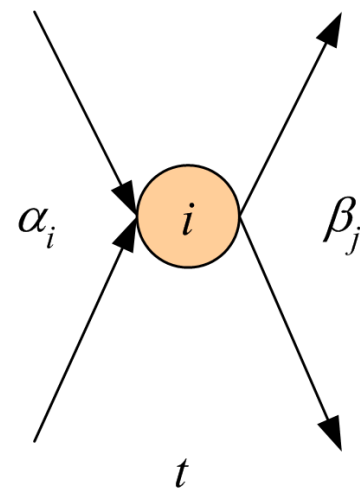


# Finding the State Sequence

$$\begin{aligned}
 \gamma_t(i) &\equiv P(q_t = S_i \mid O) = \frac{P(O \mid q_t = S_i)P(q_t = S_i)}{P(O)} \\
 &= \frac{P(O_1 \dots O_t \mid q_t = S_i)P(O_{t+1} \dots O_T \mid q_t = S_i)P(q_t = S_i)}{P(O)} \\
 &= \frac{P(O_1 \dots O_t, q_t = S_i)P(O_{t+1} \dots O_T \mid q_t = S_i)}{\sum_j P(O \mid q_t = S_j)P(S_j)} \\
 &= \frac{\alpha_t(i)\beta_t(i)}{\sum_j \alpha_t(j)\beta_t(j)}
 \end{aligned}$$

A red box highlights the term  $P(O_{t+1} \dots O_T \mid q_t = S_i)$  in the third line, with a red arrow pointing to the label  $\beta_t(i)$  below it.

A red box highlights the denominator  $\sum_j \alpha_t(j)\beta_t(j)$  in the fourth line, with an arrow pointing to a box labeled "Normalization factor".



# Finding the State Sequence

Backward variable:

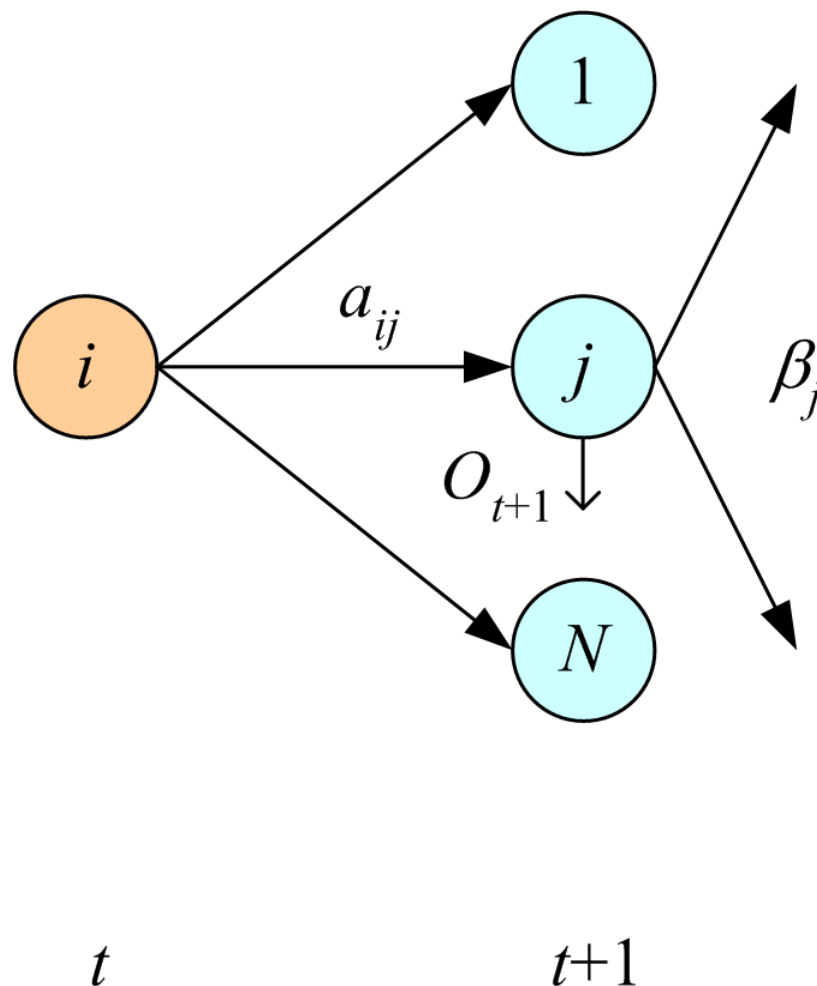
$$\beta_t(i) \equiv P(O_{t+1} \cdots O_T \mid q_t = S_i, \lambda)$$

Initializa tion :

$$\beta_T(i) = 1$$

Recursion :

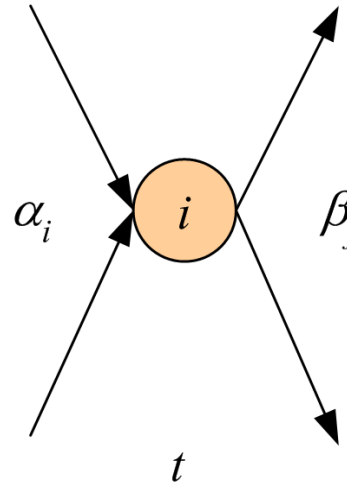
$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$





# Finding the State Sequence

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{\sum_{j=1}^N \alpha_t(j)\beta_t(j)}$$



Choose the state that has the highest probability,  
for each time step:

$$q_t^* = \arg \max_i \gamma_t(i)$$

**No!**

# Viterbi's Algorithm

$$\delta_t(\hat{i}) \equiv \max_{q_1 q_2 \dots q_{t-1}} p(q_1 q_2 \dots q_{t-1}, q_t = \mathbf{S}_{\hat{i}}, \mathbf{O}_1 \dots \mathbf{O}_t | \lambda)$$

- Initialization:

$$\delta_1(\hat{i}) = \pi_{\hat{i}} b_{\hat{i}}(\mathbf{O}_1), \psi_1(\hat{i}) = 0$$

- Recursion:

$$\delta_t(\hat{j}) = \max_i \delta_{t-1}(\hat{i}) a_{ij} b_j(\mathbf{O}_t), \psi_t(\hat{j}) = \operatorname{argmax}_i \delta_{t-1}(\hat{i}) a_{ij}$$

- Termination:

$$\mathbf{p}^* = \max_i \delta_T(\hat{i}), \mathbf{q}_T^* = \operatorname{argmax}_i \delta_T(\hat{i})$$

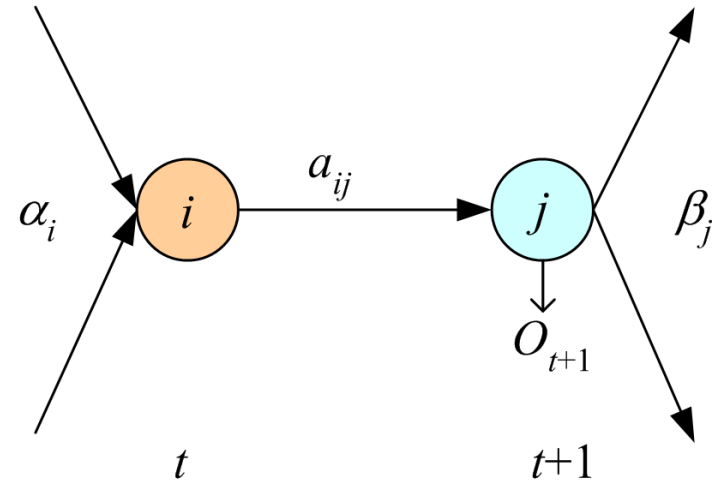
- **Path backtracking:**

$$q_{t-1}^* = \psi_t(q_t^*), t=T, T-1, \dots, 1$$

# Learning

$$\xi_t(i, j) \equiv P(q_t = S_i, q_{t+1} = S_j \mid O, \lambda)$$

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_k \sum_l \alpha_t(k) a_{kl} b_l(O_{t+1}) \beta_{t+1}(l)}$$



Baum - Welch algorithm (EM) :

$$z_i^t = \begin{cases} 1 & \text{if } q_t = S_i \\ 0 & \text{otherwise} \end{cases} \quad z_{ij}^t = \begin{cases} 1 & \text{if } q_t = S_i \text{ and } q_{t+1} = S_j \\ 0 & \text{otherwise} \end{cases}$$

# Baum-Welch (EM)

$$\text{E-step : } E[z_i^t] = \gamma_t(i) \quad E[z_{ij}^t] = \xi_t(i, j), \quad \gamma_t(i) = \sum_j \xi_t(i, j)$$

M-step :

$$\hat{\pi}_i = \frac{\sum_{k=1}^K \gamma_1^k(i)}{K} \quad \hat{a}_{ij} = \frac{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \xi_t^k(i, j)}{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \gamma_t^k(i)}$$
$$\hat{b}_j(m) = \frac{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \gamma_t^k(j) \mathbf{1}(O_t^k = v_m)}{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \gamma_t^k(i)}$$

# Continuous Observations

- Discrete:

$$P(O_t | q_t = S_j, \lambda) = \prod_{m=1}^M b_j(m)^{r_m^t} \quad r_m^t = \begin{cases} 1 & \text{if } O_t = v_m \\ 0 & \text{otherwise} \end{cases}$$

- Gaussian mixture (Discretize using **k**-means):

$$P(O_t | q_t = S_j, \lambda) = \sum_{l=1}^L P(\mathbf{G}_{jl}) p(O_t | q_t = S_j, \mathbf{G}_l, \lambda)$$

- Continuous:  $\sim \mathcal{N}(\mu_l, \Sigma_l)$

$$P(O_t | q_t = S_j, \lambda) \sim \mathcal{N}(\mu_j, \sigma_j^2)$$

Use EM to learn parameters, e.g.,  $\hat{\mu}_j = \frac{\sum_t \gamma_t(j) O_t}{\sum_t \gamma_t(j)}$

# HMM with Input

- Input-dependent observations:

$$P(O_t \mid q_t = S_j, x^t, \lambda) \sim \mathcal{N}(g_j(x^t \mid \theta_j), \sigma_j^2)$$

- Input-dependent transitions (Meila and Jordan, 1996; Bengio and Frasconi, 1996):

$$P(q_{t+1} = S_j \mid q_t = S_i, x^t)$$

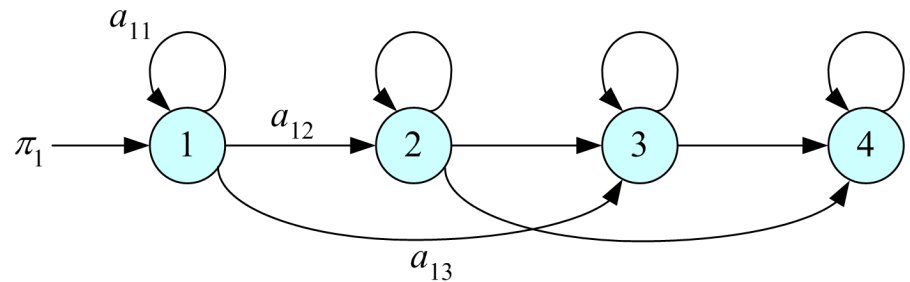
- Time-delay input:

$$x^t = f(O_{t-\tau}, \dots, O_{t-1})$$

# Model Selection in HMM

- Left-to-right HMMs:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & 0 \\ 0 & a_{22} & a_{23} & a_{24} \\ 0 & 0 & a_{33} & a_{34} \\ 0 & 0 & 0 & a_{44} \end{bmatrix}$$



- In classification, for each  $C_i$ , estimate  $\mathbf{P}(\mathbf{O} \mid \lambda_i)$  by a separate HMM and use Bayes' rule

$$P(\lambda_i \mid O) = \frac{P(O \mid \lambda_i)P(\lambda_i)}{\sum_j P(O \mid \lambda_j)P(\lambda_j)}$$

