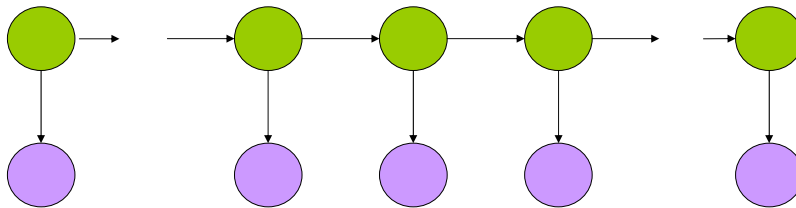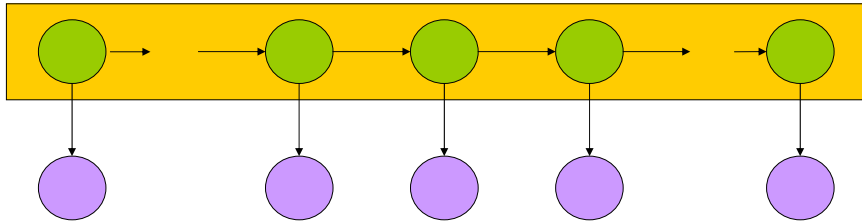# Hidden Markov Models

**David Meir Blei**
**November 1, 1999**
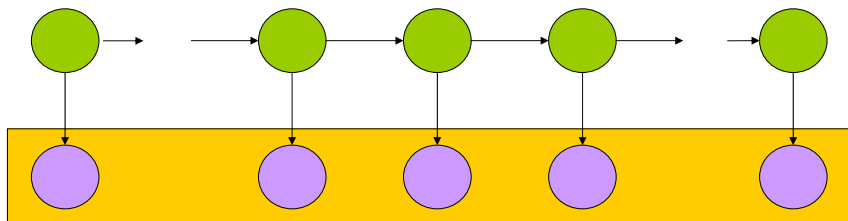
---

# What is an HMM?



- Graphical Model
- Circles indicate states
- Arrows indicate probabilistic dependencies between states
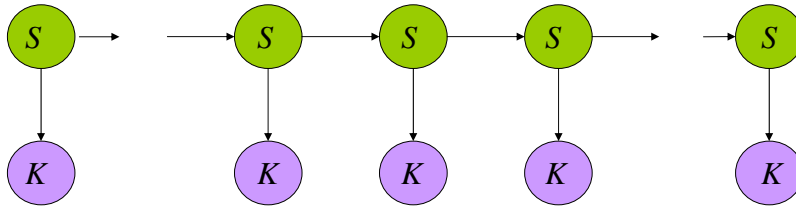
# What is an HMM?

- Green circles are *hidden states*
- Dependent only on the previous state
- "The past is independent of the future given the present."
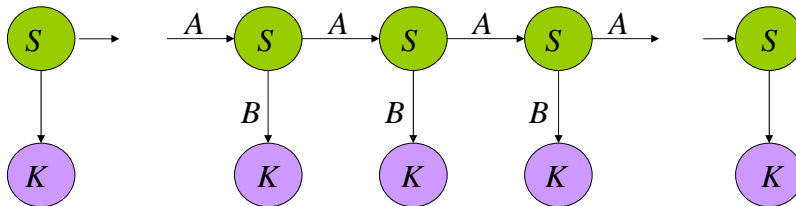
# What is an HMM?

- Purple nodes are *observed states*
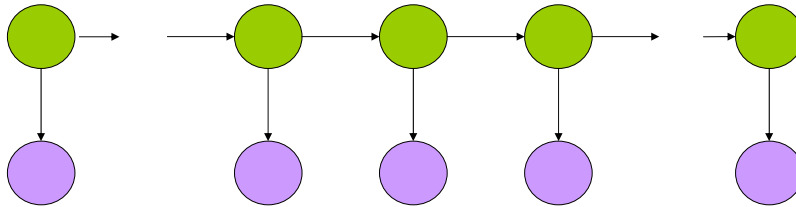- Dependent only on their corresponding hidden state

# HMM Formalism



- $\{S, K, \Pi, A, B\}$
- $S : \{s_1 \ldots s_N\}$ are the values for the hidden states
- $K : \{k_1 \ldots k_M\}$ are the values for the observations
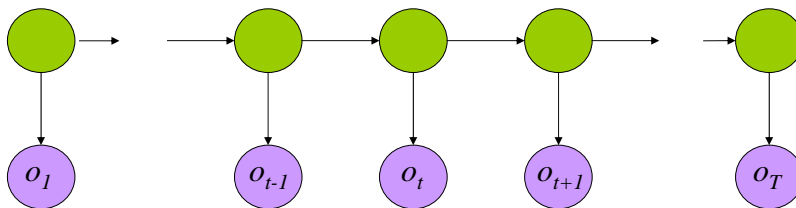
# HMM Formalism



- $\{S, K, \Pi, A, B\}$
- $\Pi = \{\pi_\iota\}$ are the initial state probabilities
- $A = \{a_{ij}\}$ are the state transition probabilities
- $B = \{b_{ik}\}$ are the observation state probabilities

# Inference in an HMM



- Compute the probability of a given observation sequence
- Given an observation sequence, compute the most likely hidden state sequence
- Given an observation sequence and set of possible models, which model most closely fits the data?
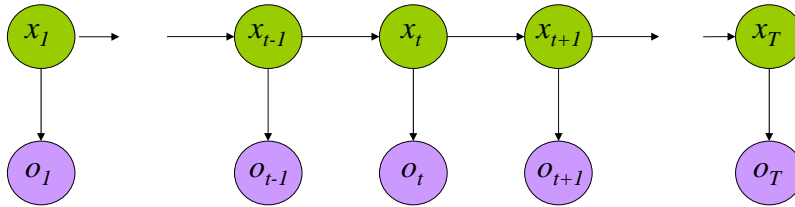
# Decoding



Given an observation sequence and a model, compute the probability of the observation sequence
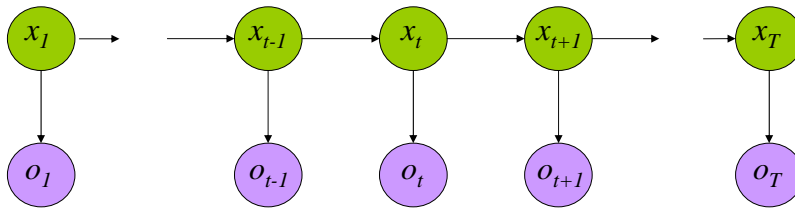
$$O = (o_1...o_T), \mu = (A, B, \Pi)$$

Compute $P(O \mid \mu)$

# Decoding



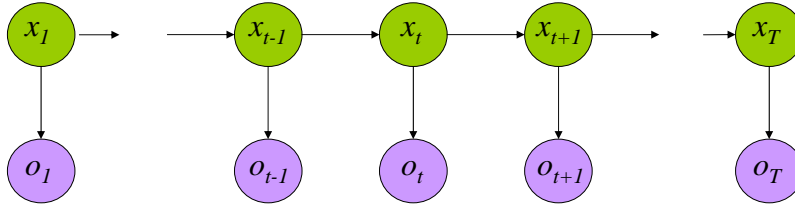$$P(O \mid X, \mu) = b_{x_1 o_1} b_{x_2 o_2} ... b_{x_T o_T}$$

# Decoding



$$P(O \mid X, \mu) = b_{x_1 o_1} b_{x_2 o_2} ... b_{x_T o_T}$$

$$P(X \mid \mu) = \pi_{x_1} a_{x_1 x_2} a_{x_2 x_3} ... a_{x_{T-1} x_T}$$

# Decoding



$$P(O \mid X, \mu) = b_{x_1 o_1} b_{x_2 o_2} ... b_{x_T o_T}$$

$$P(X \mid \mu) = \pi_{x_1} a_{x_1 x_2} a_{x_2 x_3} ... a_{x_{T-1} x_T}$$

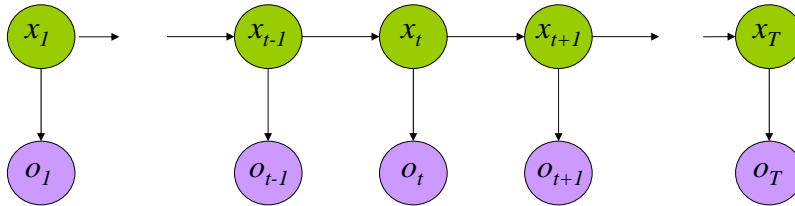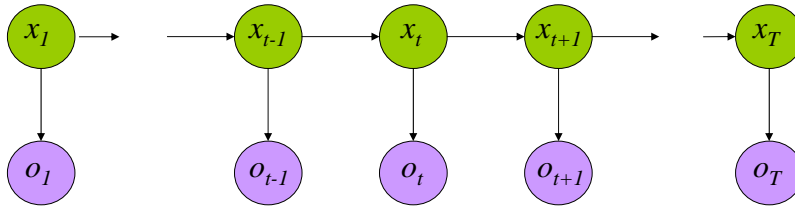$$P(O, X \mid \mu) = P(O \mid X, \mu) P(X \mid \mu)$$

# Decoding



$$P(O \mid X, \mu) = b_{x_1 o_1} b_{x_2 o_2} ... b_{x_T o_T}$$

$$P(X \mid \mu) = \pi_{x_1} a_{x_1 x_2} a_{x_2 x_3} ... a_{x_{T-1} x_T}$$
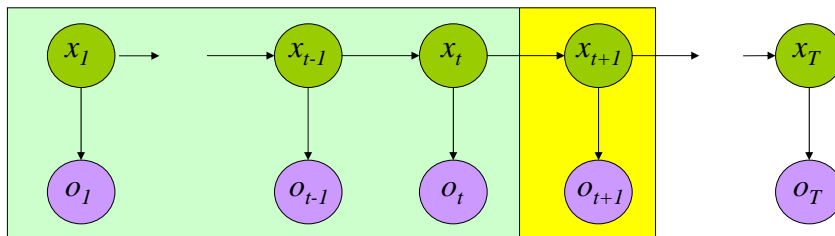
$$P(O, X \mid \mu) = P(O \mid X, \mu) P(X \mid \mu)$$

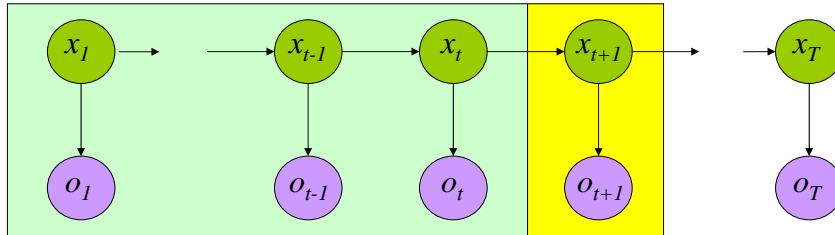$$P(O \mid \mu) = \sum_X P(O \mid X, \mu) P(X \mid \mu)$$

6

# Decoding



$$P(O \mid \mu) = \sum_{\{x_1...x_T\}} \pi_{x_1} b_{x_1 o_1} \prod_{t=1}^{T-1} a_{x_t x_{t+1}} b_{x_{t+1} o_{t+1}}$$

# Forward Procedure



- Special structure gives us an efficient solution using *dynamic programming*.
- **Intuition**: Probability of the first *t* observations is the same for all possible *t*+1 length state sequences.
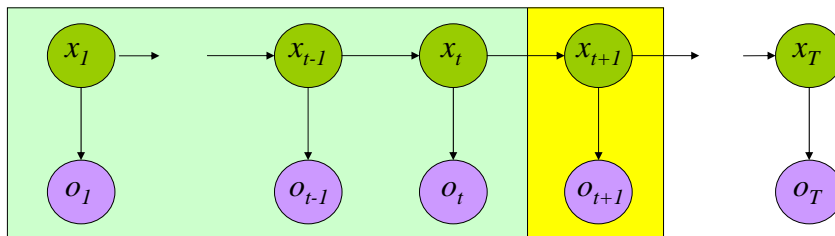- **Define:** $\alpha_i(t) = P(o_1...o_t, x_t = i \mid \mu)$

# Forward Procedure



$\alpha_j(t+1)$

$$= P(o_1...o_{t+1}, x_{t+1} = j)$$
$$= P(o_1...o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
$$= P(o_1...o_t \mid x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
$$= P(o_1...o_t, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

# Forward Procedure



$\alpha_j(t+1)$

$$= P(o_1...o_{t+1}, x_{t+1} = j)$$
$$= P(o_1...o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
$$= P(o_1...o_t \mid x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
$$= P(o_1...o_t, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

# Forward Procedure



$\alpha_j(t+1)$

$$= P(o_1...o_{t+1}, x_{t+1} = j)$$
$$= P(o_1...o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
$$= P(o_1...o_t \mid x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
$$= P(o_1...o_t, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

# Forward Procedure



$\alpha_j(t+1)$

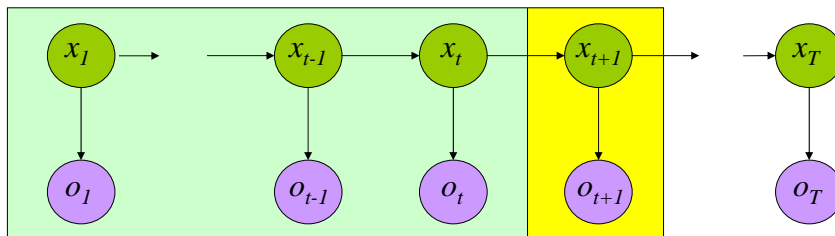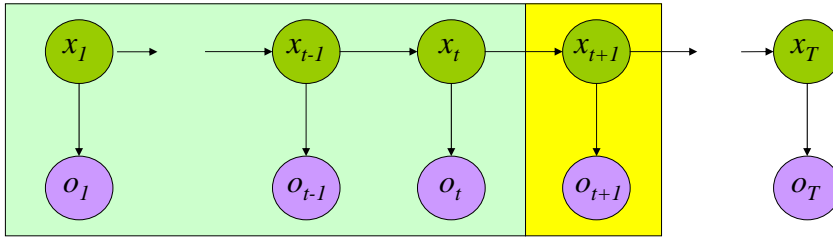$$= P(o_1...o_{t+1}, x_{t+1} = j)$$
$$= P(o_1...o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
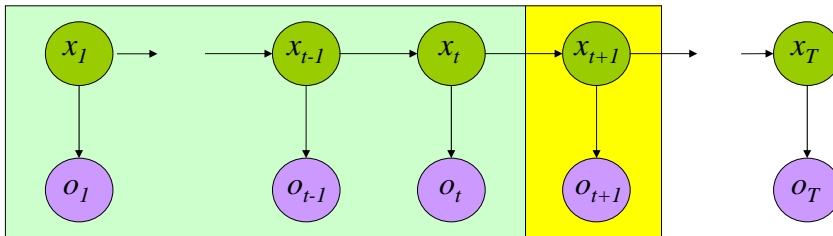$$= P(o_1...o_t \mid x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j)$$
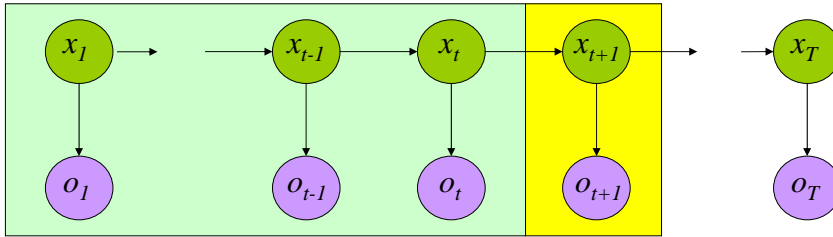$$= P(o_1...o_t, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

# Forward Procedure



$$= \sum_{i=1...N} P(o_1...o_t, x_t = i, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_{t+1} = j \mid x_t = i)P(x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_t = i)P(x_{t+1} = j \mid x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} \alpha_i(t) a_{ij} b_{jo_{t+1}}$$

# Forward Procedure



$$= \sum_{i=1...N} P(o_1...o_t, x_t = i, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_{t+1} = j \mid x_t = i)P(x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_t = i)P(x_{t+1} = j \mid x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} \alpha_i(t) a_{ij} b_{jo_{t+1}}$$

# Forward Procedure



$$= \sum_{i=1...N} P(o_1...o_t, x_t = i, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_{t+1} = j \mid x_t = i)P(x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_t = i)P(x_{t+1} = j \mid x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} \alpha_i(t)a_{ij}b_{jo_{t+1}}$$

# Forward Procedure



$$= \sum_{i=1...N} P(o_1...o_t, x_t = i, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_{t+1} = j \mid x_t = i)P(x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

$$= \sum_{i=1...N} P(o_1...o_t, x_t = i)P(x_{t+1} = j \mid x_t = i)P(o_{t+1} \mid x_{t+1} = j)$$

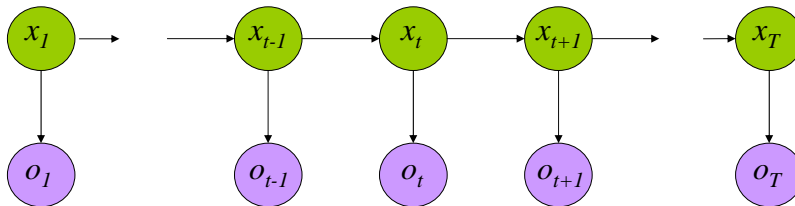$$= \sum_{i=1...N} \alpha_i(t)a_{ij}b_{jo_{t+1}}$$

# Backward Procedure



$$\beta_i(T+1) = 1$$

$$\beta_i(t) = P(o_t...o_T \mid x_t = i)$$

$$\beta_i(t) = \sum_{j=1...N} a_{ij} b_{io_t} \beta_j(t+1)$$

Probability of the rest of the states given the first state

# Decoding Solution
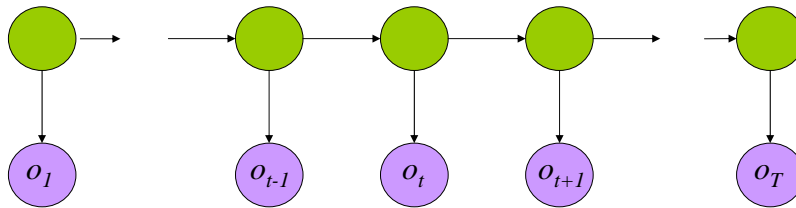


$$P(O \mid \mu) = \sum_{i=1}^{N} \alpha_i(T) \qquad \text{Forward Procedure}$$

$$P(O \mid \mu) = \sum_{i=1}^{N} \pi_i \beta_i(1) \qquad \text{Backward Procedure}$$

$$P(O \mid \mu) = \sum_{i=1}^{N} \alpha_i(t) \beta_i(t) \qquad \text{Combination}$$
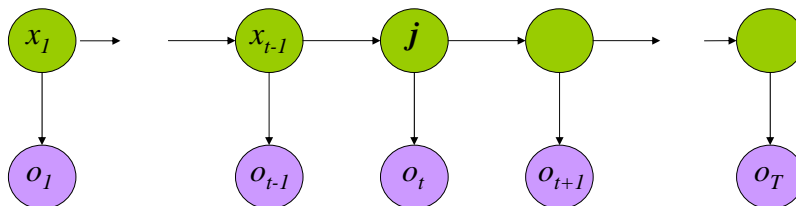
# Best State Sequence



- Find the state sequence that best explains the observations

- **Viterbi** algorithm

- $\arg\max\limits_{X} P(X \mid O)$
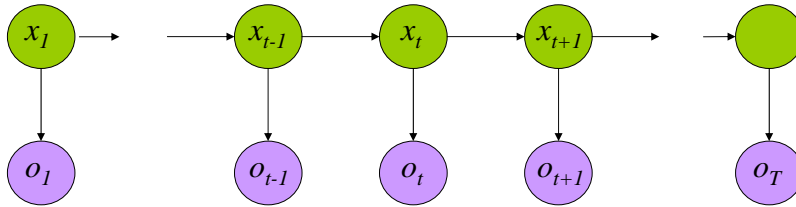
# Viterbi Algorithm



$$\delta_j(t) = \max_{x_1...x_{t-1}} P(x_1...x_{t-1}, o_1...o_{t-1}, x_t = j, o_t)$$

The state sequence which maximizes the probability of seeing the observations to time t-1, landing in state j, and seeing the observation at time t
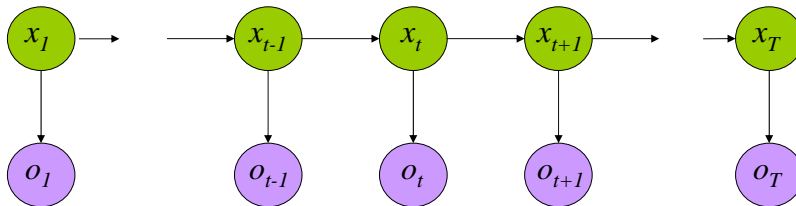
# Viterbi Algorithm



$$\delta_j(t) = \max_{x_1...x_{t-1}} P(x_1...x_{t-1}, o_1...o_{t-1}, x_t = j, o_t)$$

$$\delta_j(t+1) = \max_i \delta_i(t) a_{ij} b_{jo_{t+1}}$$

$$\psi_j(t+1) = \arg\max_i \delta_i(t) a_{ij} b_{jo_{t+1}}$$
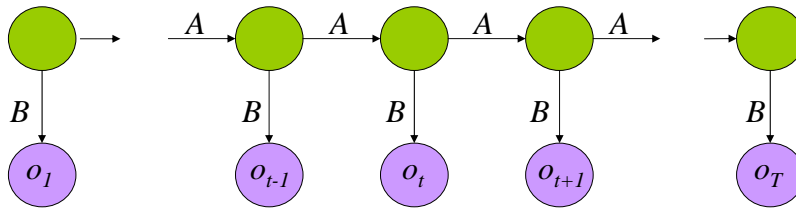
Recursive Computation

# Viterbi Algorithm



$$\hat{X}_T = \arg\max_i \delta_i(T)$$

$$\hat{X}_t = \psi_{\hat{X}_{t+1}}(t+1)$$

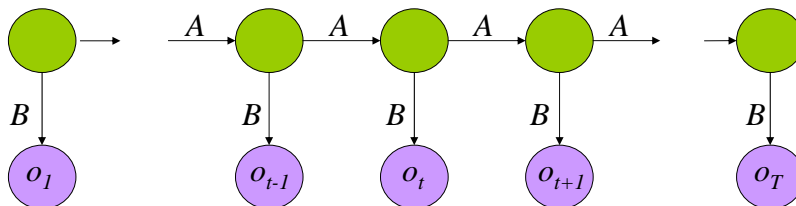$$P(\hat{X}) = \arg\max_i \delta_i(T)$$

Compute the most likely state sequence by working backwards

# Parameter Estimation



- Given an observation sequence, find the model that is most likely to produce that sequence.
- No analytic method
- Given a model and observation sequence, update the model parameters to better fit the observations.
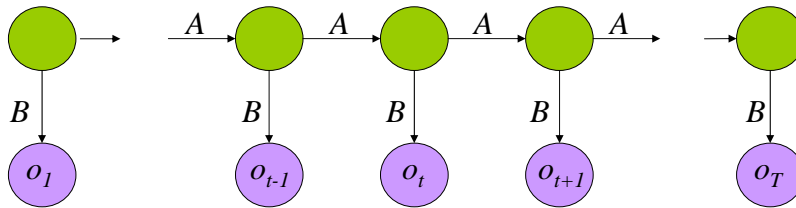
# Parameter Estimation



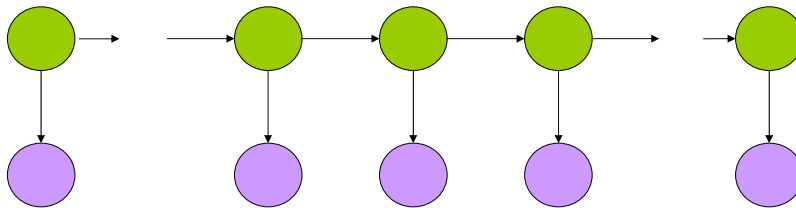| | |
|---|---|
| $$p_t(i, j) = \dfrac{\alpha_i(t) a_{ij} b_{j o_{t+1}} \beta_j(t+1)}{\sum\limits_{m=1...N} \alpha_m(t) \beta_m(t)}$$ | Probability of traversing an arc |
| $$\gamma_i(t) = \sum_{j=1...N} p_t(i, j)$$ | Probability of being in state $i$ |

15

# Parameter Estimation

$$\hat{\pi}_i = \gamma_i(1)$$

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T} p_t(i,j)}{\sum_{t=1}^{T} \gamma_i(t)}$$

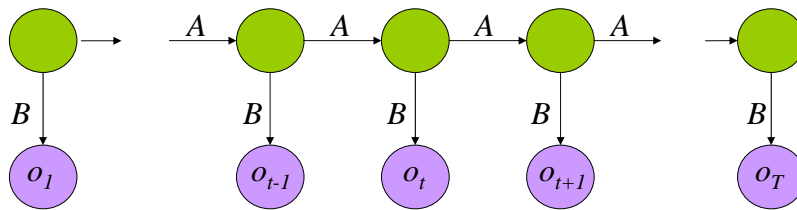$$\hat{b}_{ik} = \frac{\sum_{\{t:o_t=k\}} \gamma_t(i)}{\sum_{t=1}^{T} \gamma_i(t)}$$

Now we can compute the new estimates of the model parameters.

# HMM Applications

- Generating parameters for n-gram models
- Tagging speech
- Speech recognition

# The Most Important Thing



We can use the special structure of this model to do a lot of neat math and solve problems that are otherwise not solvable.