Steven L. Scott

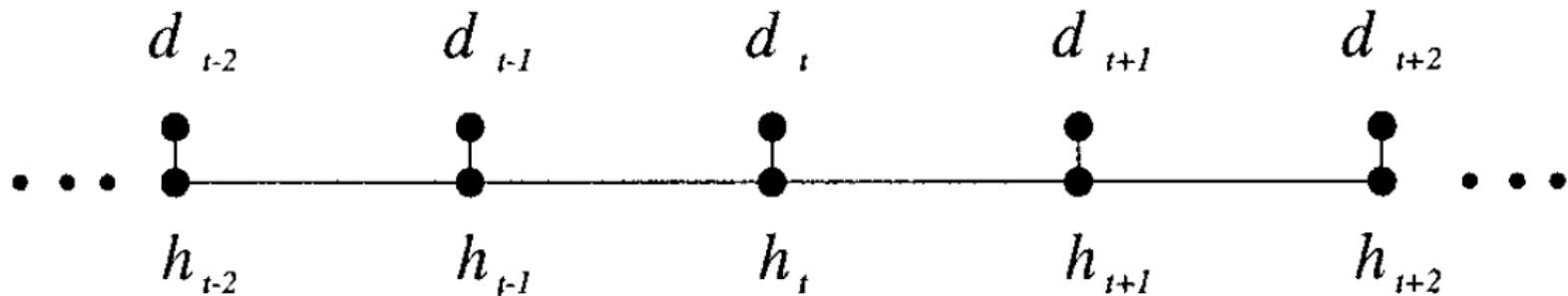# Bayesian Methods for Hidden Markov Models: Recursive Computing in the 21st Century

Presented by Ahmet Engin Ural

# Outline

- Overview of HMM
  - Evaluating likelihoods
    - The Likelihood Recursion
    - The Forward-Backward Recursion
- Sampling HMM
  - DG and FB samplers
  - Autocovariance of samplers
  - Some issues with samplers (in general)
- Estimation
  - Marginal
  - MAP
  - Size of the state space

# Hidden Markov Models



| | |
|---|---|
| h | $(h_1 \dots h_n)$ |
| **Q** | $Q_{ij} = q\,(h_i\,,\,h_j\,)$ (stationary) |
| **$\pi_0$** | Initial state |
| **θ** | $P_0 \dots P_{s-1}$ |

$$p(d_t \mid d_{-t}, \mathbf{h}, \theta, \mathbf{Q}, \pi_0) = P_{h_t}(d_t \mid \theta), \qquad (2)$$

# Calculating the likelihood

$$p(d_1^n \mid \theta) = \sum_{\mathbf{h} \in \mathcal{S}^n} \pi_0(h_1) P_{h_1}(d_1 \mid \theta) \prod_{t=2}^{n} q(h_{t-1}, h_t) P_{h_t}(d_t \mid \theta).$$
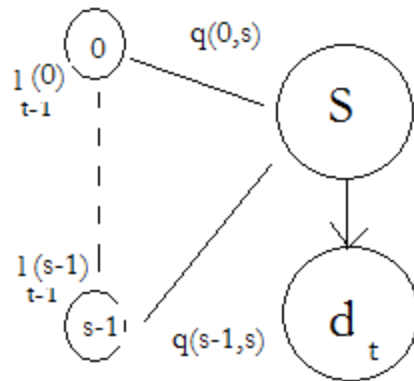
$$(3)$$

Sum over all possible hidden state sequences, the probability of the observed generated by that hidden state sequence

# Calculating the likelihood

$$p(d_1^n \mid \theta) = \sum_{\mathbf{h} \in \mathcal{S}^n} \pi_0(h_1) P_{h_1}(d_1 \mid \theta) \prod_{t=2}^{n} q(h_{t-1}, h_t) P_{h_t}(d_t \mid \theta).$$

$$(3)$$

Instead, likelihood recursion $O(S^2\, n)$ steps

Forward variable: $\quad \ell_t(s) \; = P_s(d_t \mid \theta) \sum_{r=0}^{S-1} q(r, s) \ell_{t-1}(r).$

# Forward Backward Recursions

- Forward recursion is as likelihood recursion
- Backward variable: $\pi_t(s \mid \theta) = \ell_t(s)/\ell_t^*$

  where $\ell_t^* = \sum_{s=0}^{S-1} \ell_t(s)$
- Transition probabilities, *p( r - > s at time t | we observed until t)*

$$p_{trs} \propto p(h_{t-1} = r, h_t = s, d_t \mid d_1^{t-1}, \theta)$$

$$= \pi_{t-1}(r \mid \theta) q(r, s) P_s(d_t \mid \theta),$$

- Backward recursion

$$p_{trs}' = p(h_{t-1} = r \mid h_t = s, d_1^n, \theta) p(h_t = s \mid d_1^n, \theta)$$

# Forward Backward Recursions

- Forward recursion is as likelihood recursion
- Backward variable: $\pi_t(s \mid \theta) = \ell_t(s)/\ell_t^*$
  where $\ell_t^* = \sum_{s=0}^{S-1} \ell_t(s)$
- Transition probabilities, *p( r -> s at time t | we observed until t)*

$$p_{trs} \propto p(h_{t-1} = r, h_t = s, d_t \mid d_1^{t-1}, \theta)$$

$$= \pi_{t-1}(r \mid \theta) q(r,s) P_s(d_t \mid \theta),$$

- Backward recursion

$$p'_{trs} = p(h_{t-1} = r \mid h_t = s, d_1^n, \theta) p(h_t = s \mid d_1^n, \theta)$$

$$= p(h_{t-1} = r \mid h_t = s, d_1^t, \theta) \pi_t'(s \mid \theta)$$

$$\pi_t'(s \mid \theta) \equiv \Pr(h_t = s \mid d_1^n, \theta)$$

# Forward Backward Recursions

- Forward recursion is as likelihood recursion
- Backward variable: $\pi_t(s \mid \theta) = \ell_t(s)/\ell_t^*$

  where $\ell_t^* = \sum_{s=0}^{S-1} \ell_t(s)$
- Transition probabilities, *p( r - > s at time t | we observed until t)*

  $$p_{trs} \propto p(h_{t-1} = r, h_t = s, d_t \mid d_1^{t-1}, \theta)$$

  $$= \pi_{t-1}(r \mid \theta) q(r, s) P_s(d_t \mid \theta),$$
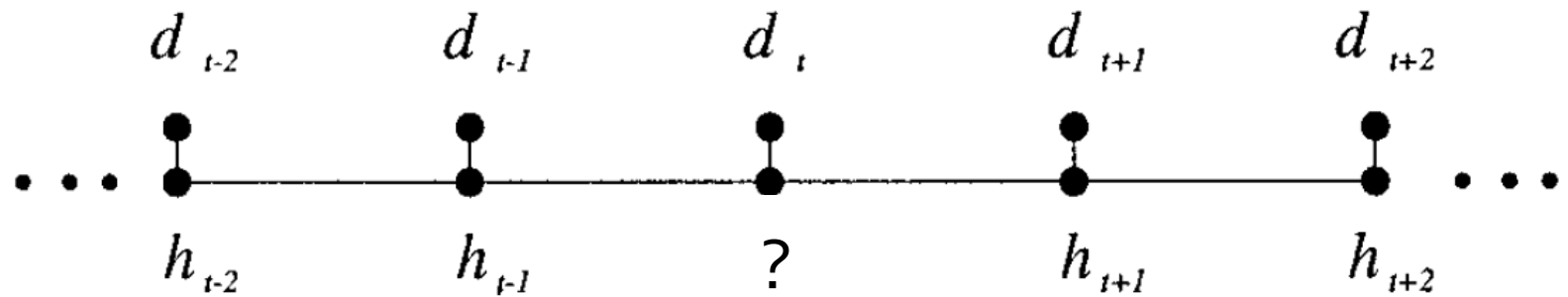- Backward recursion

  $$p_{trs}' = p(h_{t-1} = r \mid h_t = s, d_1^n, \theta) p(h_t = s \mid d_1^n, \theta)$$

  $$= p(h_{t-1} = r \mid h_t = s, d_1^t, \theta) \pi_t'(s \mid \theta)$$

  $$= p_{trs} \frac{\pi_t'(s \mid \theta)}{\pi_t(s \mid \theta)},$$

# Sampling

- Direct Gibbs Sampling

$$d_{t-2} \qquad d_{t-1} \qquad d_{t} \qquad d_{t+1} \qquad d_{t+2}$$

$$\cdots \qquad \qquad \qquad \qquad \qquad \cdots$$

$$h_{t-2} \qquad h_{t-1} \qquad ? \qquad h_{t+1} \qquad h_{t+2}$$
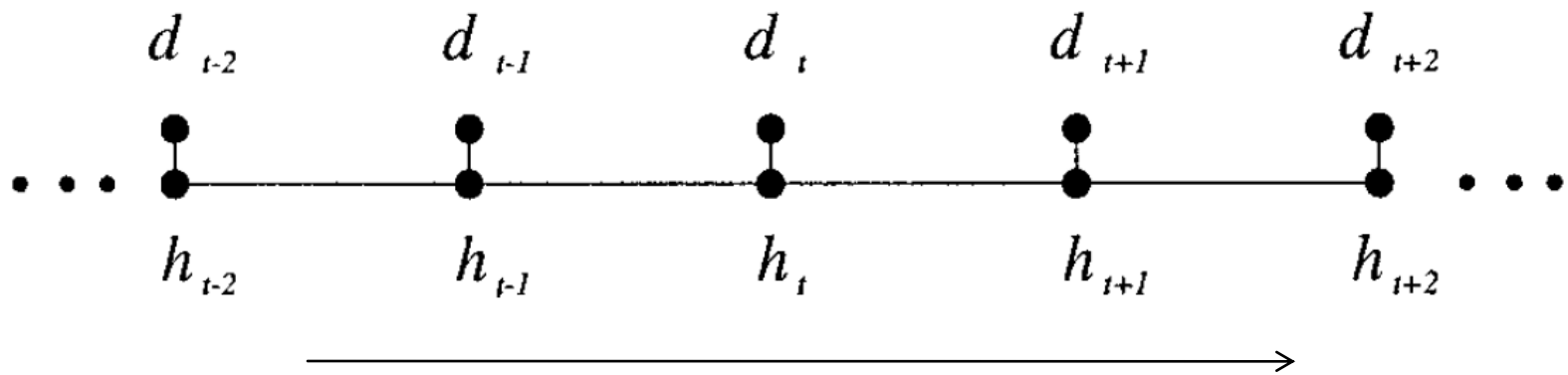
$$p(h_t = s \mid h_{-t}, d_1^n, \theta) \propto q(h_{t-1}, s) q(s, h_{t+1}) P_s(d_t \mid \theta),$$
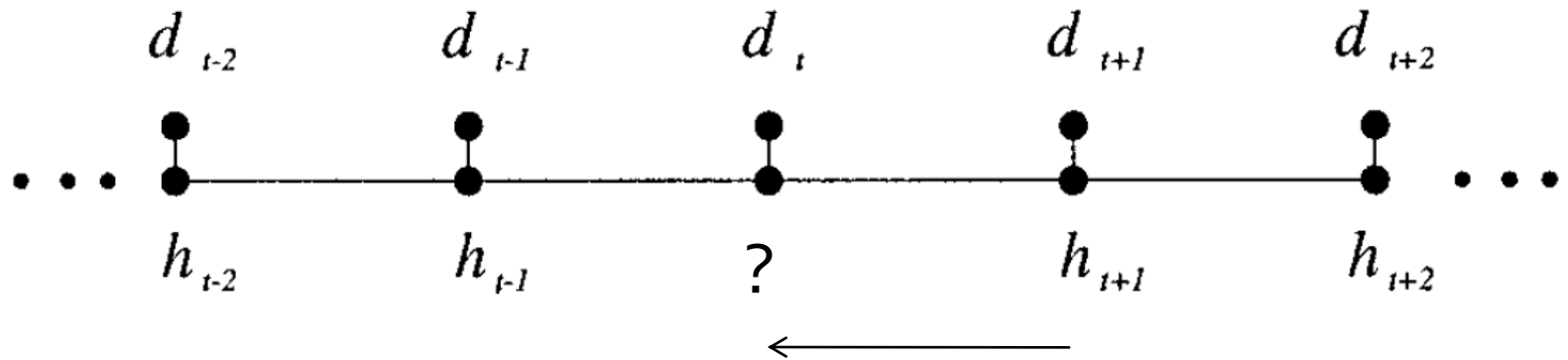
# Sampling

- Forward Backward Recursion sampling

$$d_{t-2} \qquad d_{t-1} \qquad d_t \qquad d_{t+1} \qquad d_{t+2}$$



$$h_{t-2} \qquad h_{t-1} \qquad h_t \qquad h_{t+1} \qquad h_{t+2}$$

At the forward step, the transition matrices, (P) are produced;

$$p_{trs} \propto p(h_{t-1} = r, h_t = s, d_t \mid d_1^{t-1}, \theta)$$

$$= \pi_{t-1}(r \mid \theta) q(r, s) P_s(d_t \mid \theta),$$

# Sampling

- Forward Backward Recursion sampling



At the backward step, the state is sampled by

$$p(h_{n-t} = r \mid h_{n-t+1}^n, d_1^n, \theta) \quad \propto \quad P_{n-t+1, r, h_{t+1}}.$$

# Autocovariance

- *T* is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration $\tau$. ($\mathbf{T}_1$ is for time 1)

.

$$\mathbf{T}_1 = \begin{matrix} \mathsf{I}(h_1, h_1) \\ \mathsf{I}(h_1, h_2) \\ . \\ . \\ . \\ \mathsf{I}(h_s, h_s) \end{matrix}$$
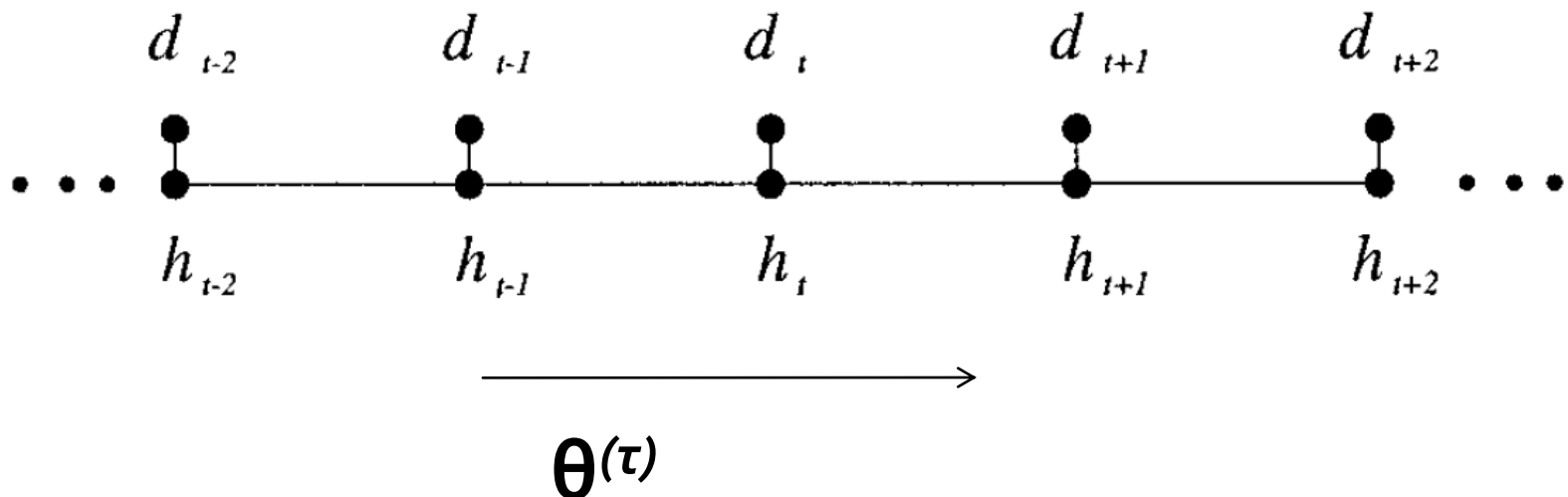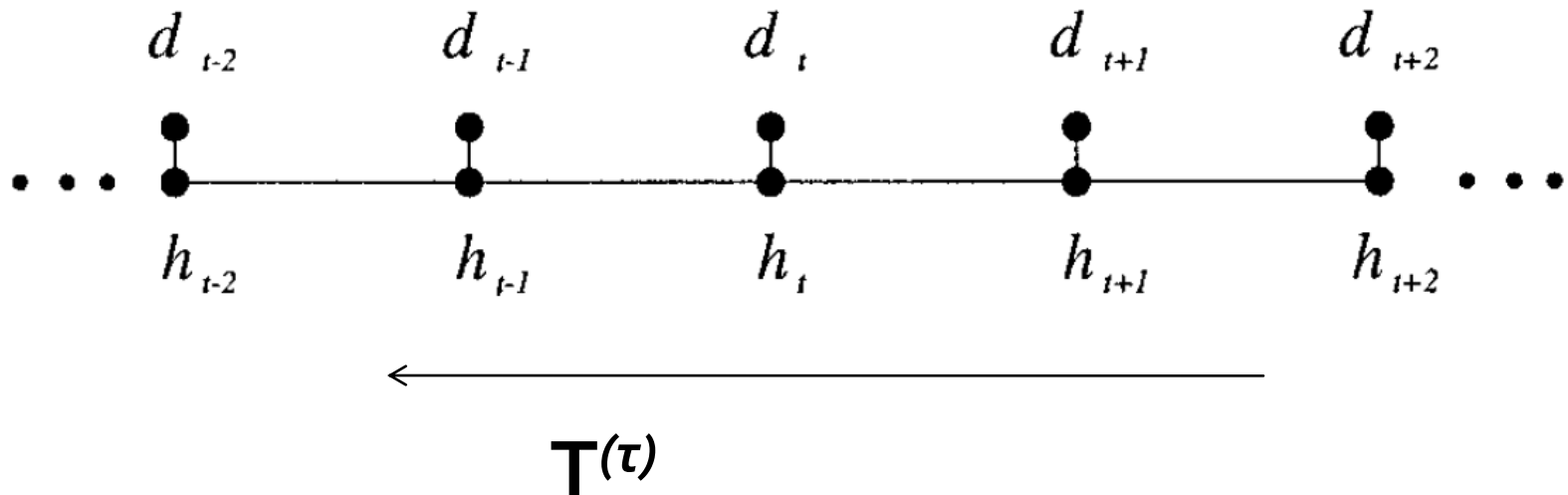
$$T = \sum_{j=1}^{m} T_j.$$

# Autocovariance

- $T$ is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration $\tau$. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emissions.

# Autocovariance

- *T* is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration τ. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.

# Autocovariance

- *T* is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration τ. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
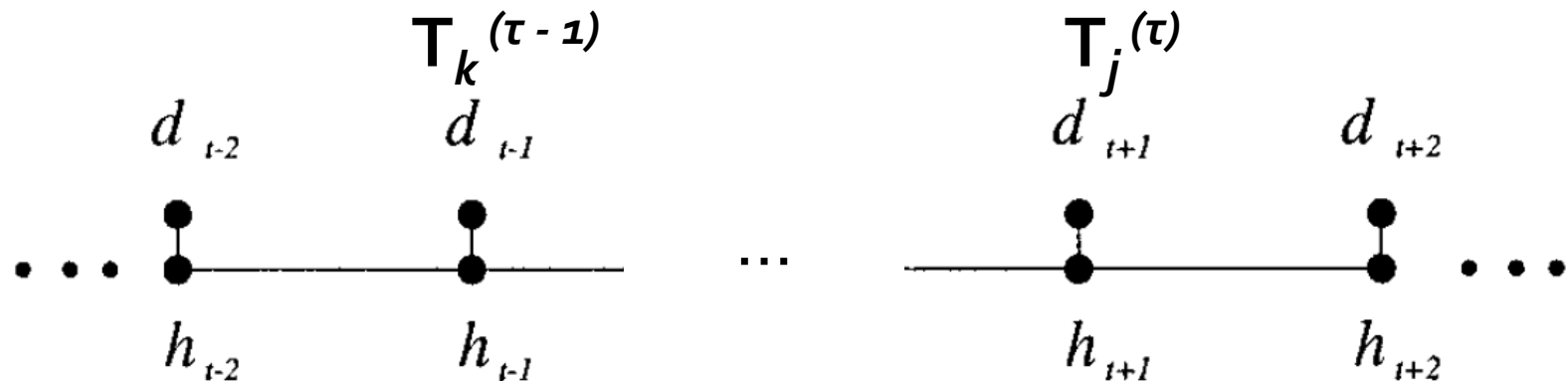


$$T^{(\tau)}$$

# Autocovariance

- $T$ is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration $\tau$. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
- DG: $T_j^{(\tau-1)}$ is conditionally independent of $T_k^{(\tau)}$ given $\theta^{(\tau)}$.

# Autocovariance

- $T$ is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration $\tau$. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
- DG: $T_j^{(\tau-1)}$ is conditionally independent of $T_k^{(\tau)}$ given $\theta^{(\tau)}$.

$$T_k^{(\tau-1)} \qquad\qquad\qquad T_j^{(\tau)}$$

$$d_{t-2} \qquad d_{t-1} \qquad\qquad d_{t+1} \qquad d_{t+2}$$

$$\ldots \qquad\qquad \ldots \qquad\qquad \ldots$$

$$h_{t-2} \qquad h_{t-1} \qquad\qquad h_{t+1} \qquad h_{t+2}$$

$T_k^{(\tau-1)}$ cond indep $T_j^{(\tau)}$ , if $\theta^{(\tau)}$ , $T_p^{(\tau-1)}$ and $T_l^{(\tau)}$ are given.
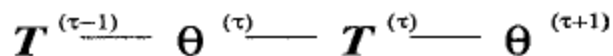
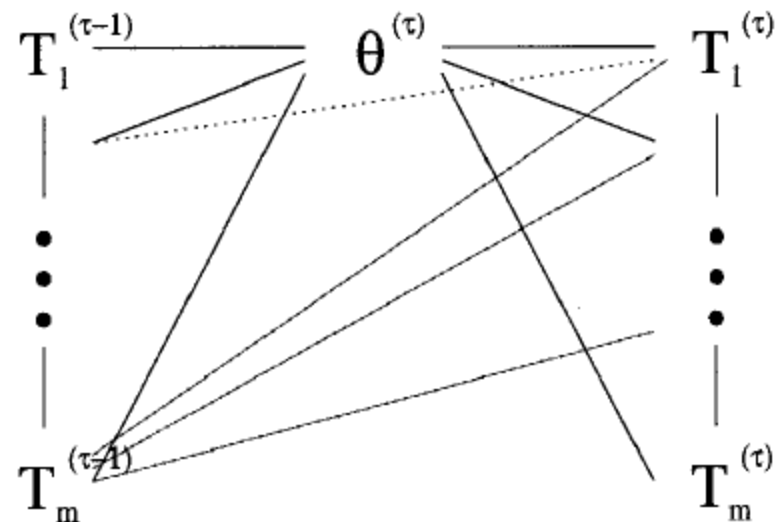$\{p > j\}$ $\qquad\qquad$ $\{l < j\}$

# Autocovariance

- *T* is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration $\tau$. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
- DG : $T_j^{(\tau-1)}$ is conditionally independent of $T_k^{(\tau)}$ given $\theta^{(\tau)}$.



(a) FB

(b) DG

# Autocovariance

- $T$ is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration τ. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
- DG : $T_j^{(\tau-1)}$ is conditionally independent of $T_k^{(\tau)}$ given $\theta^{(\tau)}$.

$$\mathrm{cov}(T^{(\tau-1)}, T^{(\tau)}) = E\{\mathrm{cov}(T^{(\tau-1)}, T^{(\tau)} \mid \theta^{(\tau)})\}$$

$$+ \mathrm{cov}\{E(T^{(\tau-1)} \mid \theta^{(\tau)}), E(T^{(\tau)} \mid \theta^{(\tau)})\}$$

$$= E\{\mathrm{cov}(T^{(\tau-1)}, T^{(\tau)} \mid \theta^{(\tau)})\}$$

$$+ \mathrm{var}\{E(T^{(\tau)} \mid \theta^{(\tau)})\}, \tag{11}$$

# Autocovariance

- *T* is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration $\tau$. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
- DG : $T_j^{(\tau-1)}$ is conditionally independent of $T_k^{(\tau)}$ given $\theta^{(\tau)}$.

$$\text{cov}(T^{(\tau-1)}, T^{(\tau)}) = E\{\text{cov}(T^{(\tau-1)}, T^{(\tau)} \mid \theta^{(\tau)})\}$$

$$+ \text{cov}\{E(T^{(\tau-1)} \mid \theta^{(\tau)}), E(T^{(\tau)} \mid \theta^{(\tau)})\}$$

$$= \boxed{E\{\text{cov}(T^{(\tau-1)}, T^{(\tau)} \mid \theta^{(\tau)})\}} \quad \text{= 0 for FB}$$

$$+ \text{var}\{E(T^{(\tau)} \mid \theta^{(\tau)})\}, \tag{11}$$

# Autocovariance

- *T* is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration τ. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
- DG : $T_j^{(\tau-1)}$ is conditionally independent of $T_k^{(\tau)}$ given $\theta^{(\tau)}$.

$$\mathrm{cov}(T^{(\tau-1)}, T^{(\tau)}) = E\{\mathrm{cov}(T^{(\tau-1)}, T^{(\tau)} \mid \theta^{(\tau)})\}$$

$$+ \mathrm{cov}\{E(T^{(\tau-1)} \mid \theta^{(\tau)}), E(T^{(\tau)} \mid \theta^{(\tau)})\}$$

$$= E\{\mathrm{cov}(T^{(\tau-1)}, T^{(\tau)} \mid \theta^{(\tau)})\}$$

same for FB and DG <=
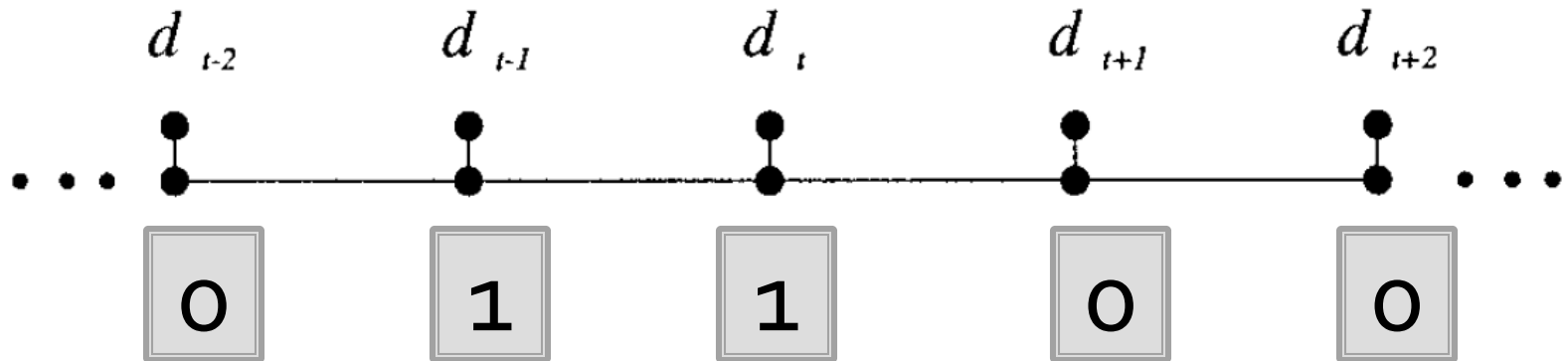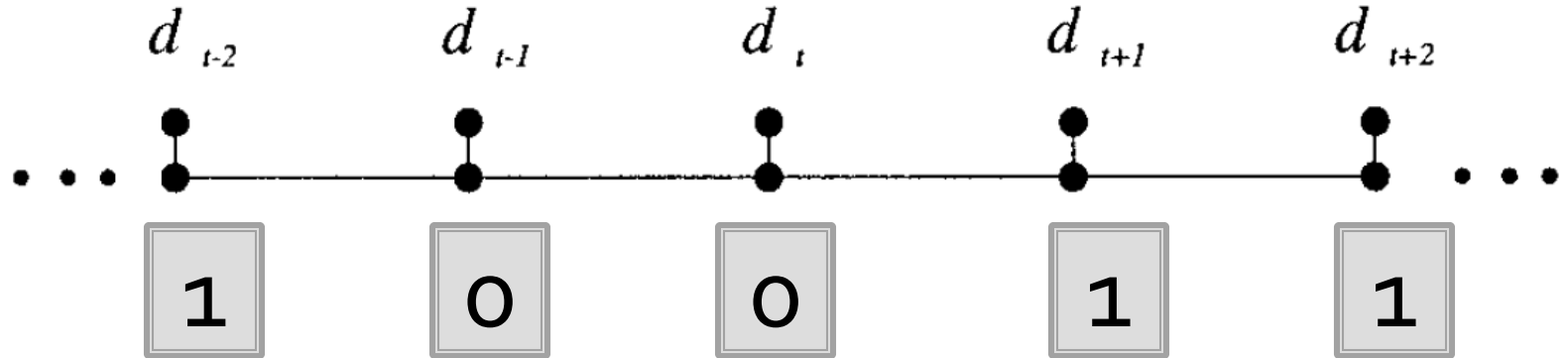$$+ \mathrm{var}\{E(T^{(\tau)} \mid \theta^{(\tau)})\}, \tag{11}$$

# Autocovariance

- $T$ is a vector that has sufficient statistics for state transitions. $T^{(\tau)}$ is the set of all such vectors iteration τ. ($T_1$ is for time 1)
- *Let* $\theta^{(\tau)}$ be the sufficient statistics for emission probabilities.
- FB: $T^{(\tau)}$ is conditionally independent of $T^{(\tau+1)}$ given $\theta^{(\tau)}$.
- DG : $T_j^{(\tau-1)}$ is conditionally independent of $T_k^{(\tau)}$ given $\theta^{(\tau)}$.

$$\mathrm{cov}_{\mathrm{DG}}\left(T^{(\tau-1)}, T^{(\tau)}\right) = \mathrm{cov}_{\mathrm{FB}}\left(T^{(\tau-1)}, T^{(\tau)}\right)$$

$$+ E_{\mathrm{DG}}\left\{\mathrm{cov}_{\mathrm{DG}}\left(T^{(\tau-1)}, T^{(\tau)} \mid \theta^{(\tau)}\right)\right\}.$$
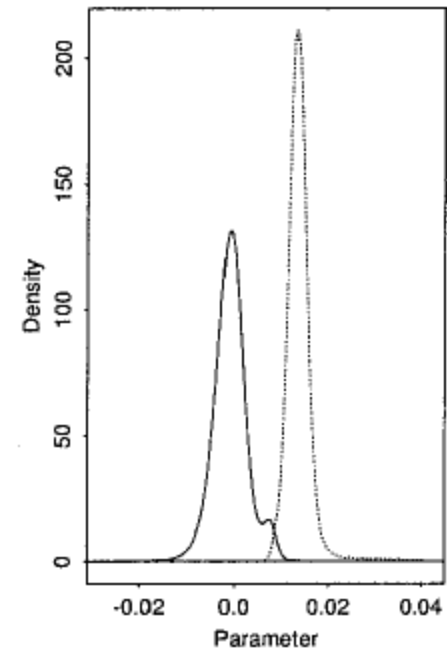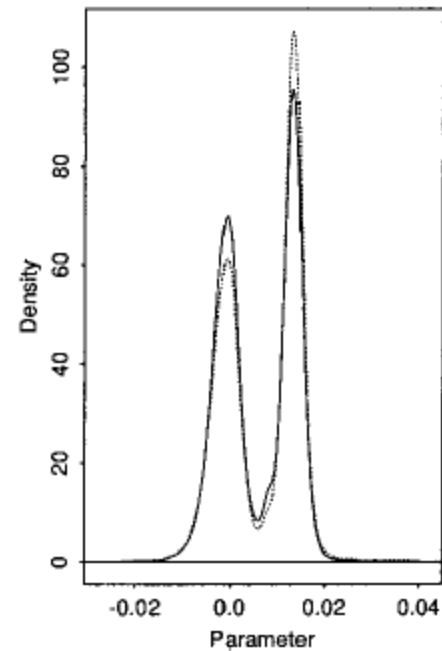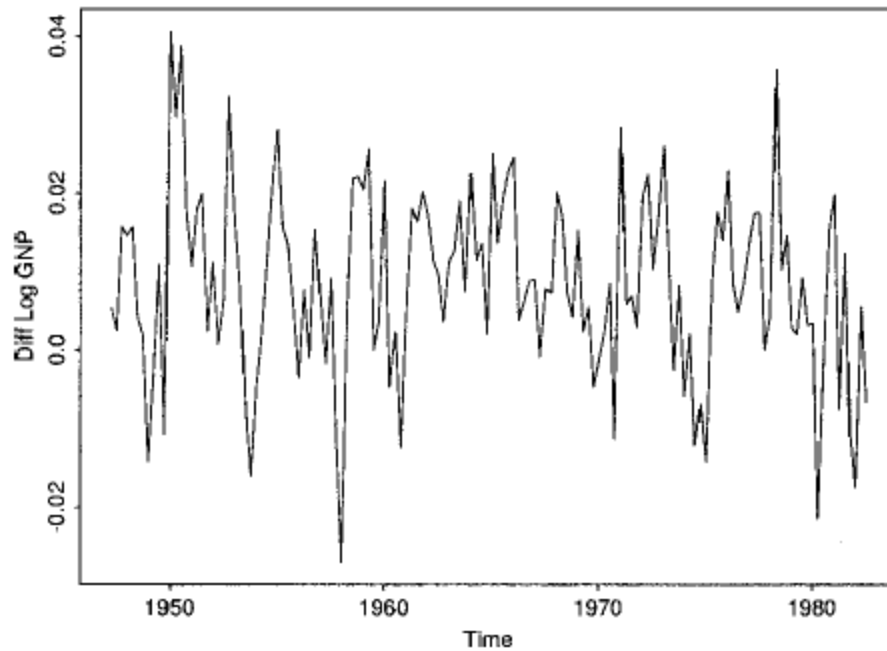
# Some issues

- Label switching

# Some issues

- Label switching
  - Implications
  - Solution: constraints

# Some issues

- Label switching
- Collapsed states
  - May be evidence for over parameterizations
  - Priors

# Estimating the hidden states

- Stating $\pi_t(s)$ is usually sufficient.

# Estimating the hidden states

- Stating $\pi_t(s)$ is usually sufficient.
- Overall configuration may be needed;
  - Marginal Distributions
  - MAP estimates

# Estimating the hidden states

- Stating $\pi_t(s)$ is usually sufficient.
- Overall configuration may be needed;
  - Marginal Distributions
    - Averaging over all runs $(1 - m)$ with indicator function:

$$\tilde{\pi}'_t(s) = 1/m \sum_{j=1}^{m} I(h_t^{(j)} = s).$$

# Estimating the hidden states

- Stating $\pi_t(s)$ is usually sufficient.
- Overall configuration may be needed;
  - Marginal Distributions
    - Averaging over all runs (1 – m) with indicator function
    - Averaging over all runs (1 – m) probabilities
      (Rao-Blackwellized estimate)

$$\hat{\pi}_t'(s) = 1/m \sum_{j=1}^{m} \pi_t'(s \mid \theta^{(j)})$$

# Estimating the hidden states

- Stating $\pi_t(s)$ is usually sufficient.
- Overall configuration may be needed;
  - Marginal Distributions
    - Averaging over all runs (1 – m) with indicator function
    - Averaging over all runs (1 – m) probabilities
  - MAP estimate ( L = max p ( h, d | θ )

$$L_1(s) = \pi_0(s)P_s(d_1 \mid \theta)$$

$$L_t(s) = \max_r[L_{t-1}(r)q(r, s)]P_s(d_t \mid \theta).$$

# Estimating the hidden states

- Stating $\pi_t(s)$ is usually sufficient.
- Overall configuration may be needed;
  - Marginal Distributions
    - Averaging over all runs $(1 - m)$ with indicator function
    - Averaging over all runs $(1 - m)$ probabilities
  - MAP estimate: to find $\hat{h}$

$$\hat{h}_t = \arg\max_{r \in \mathcal{S}} L_t(r) q(r, \hat{h}_{t+1})$$

converges when it is same for all s in $h_{t+1}$.

# Size of the state space

- Calculating p(S|D)

# Size of the state space

- Calculating p(S|D)

$$p(S \mid d_1^n) = \int p(S \mid d_1^n, \theta) p(\theta \mid d_1^n) \, d\theta$$

$$\approx 1/m \sum_{j=1}^{m} p(S \mid d_1^n, \theta^{(j)}),$$

$$p(S \mid d_1^n, \theta^{(j)}) \propto p(d_1^n \mid \theta_S^{(j)}, S) p(S)$$

# Size of the state space

- Calculating p(S | D)
- Schwartz criterion C(S):

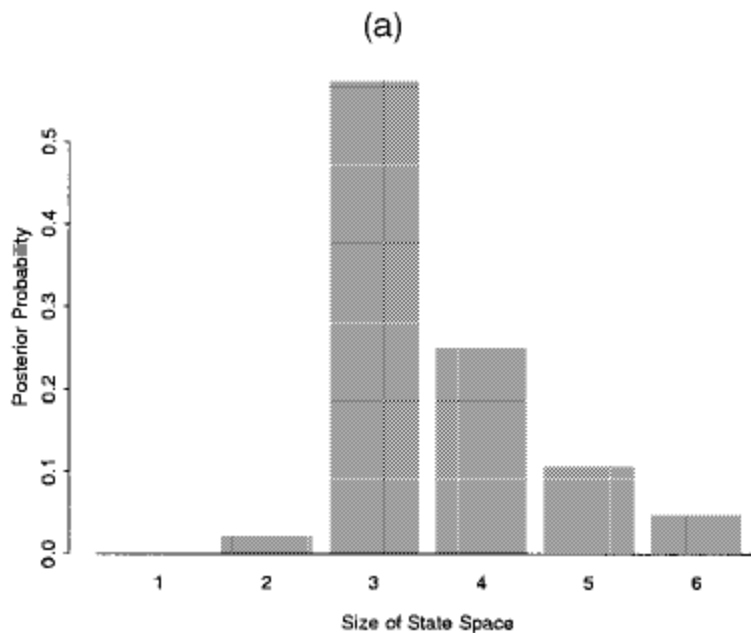$$C(S) = \log \ell - k_S \log(n)/2,$$

# Size of the state space

- Calculating p(S | D)
- Schwartz criterion C(S):
- Bayesian Information Criterion BIC:
  - p(S | D) – 2 C(S)

# Size of the state space

- Calculating p(S|D)
- Schwartz criterion C(S):
- Bayesian Information Criterion BIC

(a)

Posterior Probability

Size of State Space

(b)

| S | maximized log-posterior | $k_S$ | $C(S)$ | $BIC$ |
|---|---|---|---|---|
| 1 | -174.3 | 1 | -177.0 | 354.0 |
| 2 | -150.7 | 4 | -161.6 | 323.2 |
| 3 | -140.7 | 9 | -165.3 | 330.6 |
| 4 | -139.2 | 16 | -183.1 | 366.2 |
| 5 | -139.5 | 25 | -208.0 | 416.0 |
| 6 | -139.8 | 36 | -238.4 | 476.8 |

# Thank you