

Short Course

State Space Models, Generalized Dynamic Systems
and
Sequential Monte Carlo Methods,
and
their applications
in Engineering, Bioinformatics and Finance

Rong Chen

Rutgers University

Peking University

Part Three: Advanced Sequential Monte Carlo

3.1 Mixture Kalman Filter

3.1.1 Conditional Dynamic Linear Models

3.1.2 Mixture Kalman Filters

3.1.3 Partial Conditional Dynamic Linear Models

3.1.4 Extend Mixture Kalman Filters

3.1.5 Future Directions

3.2 Constrained SMC

3.3 Parameter Estimation in SMC

3.4 Look-Ahead Strategies

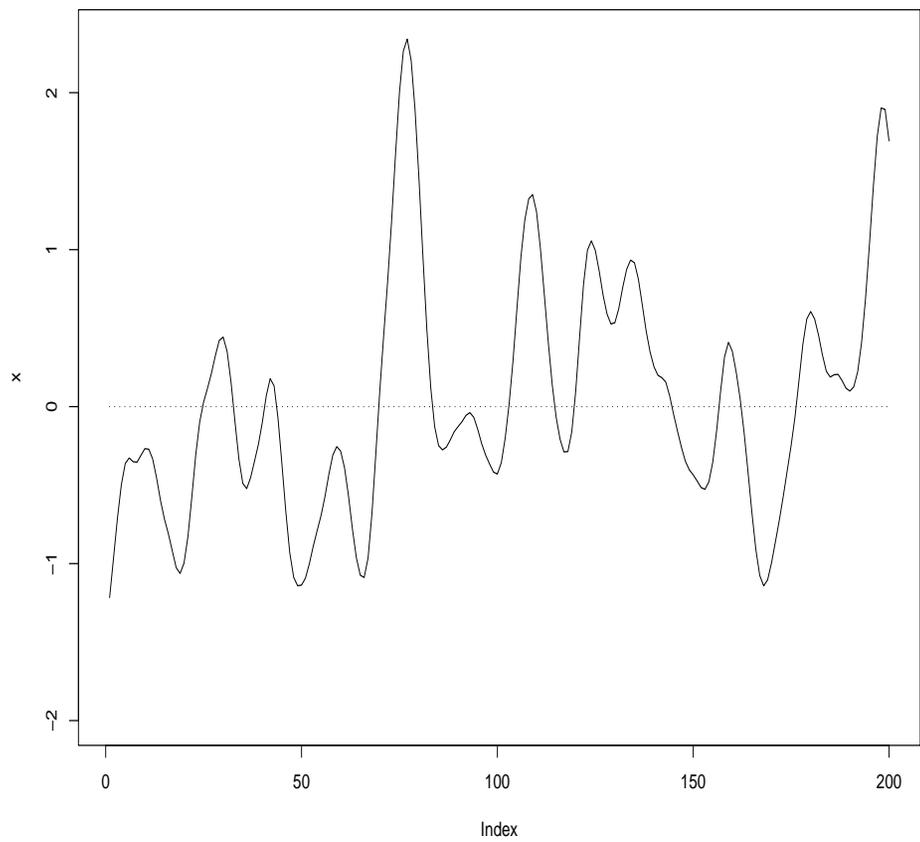
3.4.1. The principle of lookahead

- Dynamic systems often process strong 'memory'
- Future observations can reveal substantial information on the current state
- Slight delay is tolerable

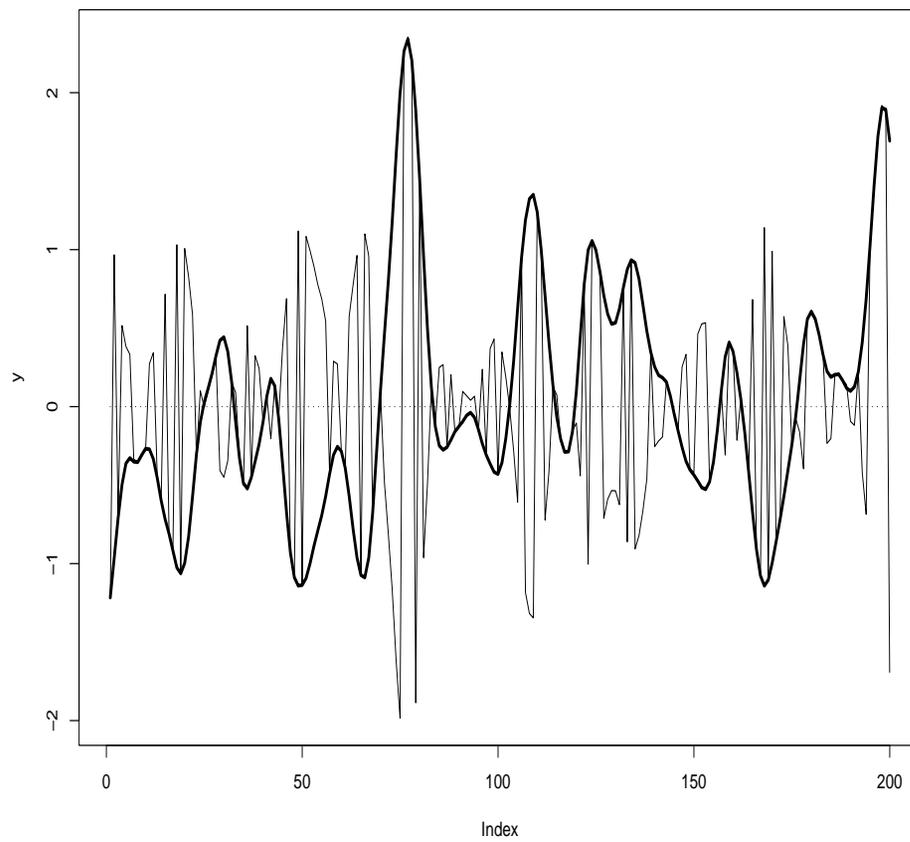
Make inference on the state x_t at time $t+d$, based on observations

$$y_1, \dots, y_t, y_{t+1}, \dots, y_{t+d}.$$

alpha



y observation



If \hat{h}_{t+d} is a consistent MC estimator of $E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d})$, then

$$E \left[\hat{h}_{t+d} - h(\mathbf{x}_t) \right]^2 = E \left[\hat{h}_{t+d} - E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}) \right]^2 + E \left[E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}) - h(\mathbf{x}_t) \right]^2$$

- The first term goes to zero with as MC sample size increases
- [*Proposition*] The second term decreases as d increases
- When MC sample size is sufficiently large, the first term is negligible comparing to the second term, then longer lookahead (larger d) always improves efficiency
- With limited sample size and limited computational time, lookahead may not always be more efficient.

3.4.2. Lookahead algorithms

(1) Exact lookahead weighting

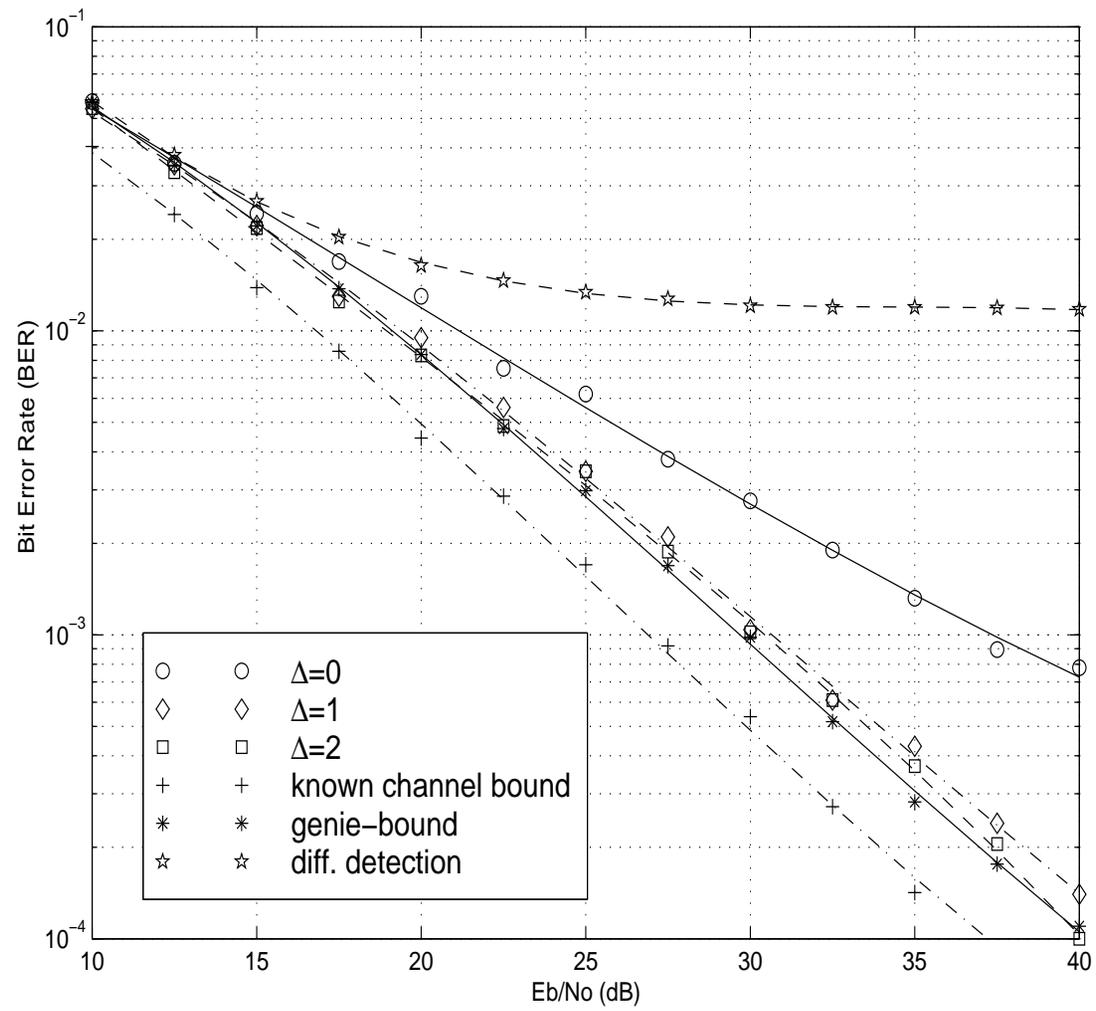
If: $(\mathbf{x}_{t+d}^{(j)}, w_{t+d}^{(j)})$ is properly weighted w.r.t. $p(\mathbf{x}_{t+d} \mid \mathbf{y}_{t+d})$

then: $(x_t^{(j)}, w_{t+d}^{(j)})$ is properly weighted w.r.t. $p(x_t \mid \mathbf{y}_{t+d})$.

Hence, inference on x_t can be made using $(x_t^{(j)}, w_{t+d}^{(j)})$, with concurrent SMC.

- Sample of x_t is drawn at time t , based on \mathbf{y}_t , with weight w_t
- Inference on x_t is made at time $t + d$, with weight w_{t+d} ;
—— w_{t+d} is based on \mathbf{y}_{t+d} and samples of $(\mathbf{x}_t, x_{t+1}, \dots, x_{t+d})$.

$$E(h(x_t) \mid \mathbf{y}_{t+d}) \approx \frac{\sum_{j=1}^m h(x_t^{(j)}) w_{t+d}^{(j)}}{\sum_{j=1}^m w_{t+d}^{(j)}}$$



In fact:

$$w_{t+d}^{(j)} = w_{t-1}^{(j)} \frac{\pi_{t+d}(\mathbf{x}_{t+d}^{(j)})}{\pi_{t-1}(\mathbf{x}_{t-1}^{(j)}) \prod_{s=t}^{t+d} g_s(x_s^{(j)} \mid \mathbf{x}_{s-1}^{(j)}, \mathbf{y}_s)}$$

An improved version (if practical) is to use

$$\tilde{w}_{t+d}^{(j)} = w_{t-1}^{(j)} \frac{\pi_{t+d}(\mathbf{x}_t^{(j)})}{\pi_{t-1}(\mathbf{x}_{t-1}^{(j)}) g_t(x_t^{(j)} \mid \mathbf{x}_{t-1}^{(j)}, \mathbf{y}_t)}$$

where $\pi_{t+d}(\mathbf{x}_t^{(j)}) = \int \pi_{t+d}(\mathbf{x}_t^{(j)}, x_{t+1}, \dots, x_{t+d}) dx_{t+1} \dots dx_{t+d}$

[Proposition]:

$$\text{Var} [w_{t+d} \mid \mathbf{y}_{t+d}] \geq \text{Var} [\tilde{w}_{t+d} \mid \mathbf{y}_{t+d}]$$

and

$$\text{Var} [w_{t+d} h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}] \geq \text{Var} [\tilde{w}_{t+d} h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}]$$

(2) Exact Lookahead Sampling

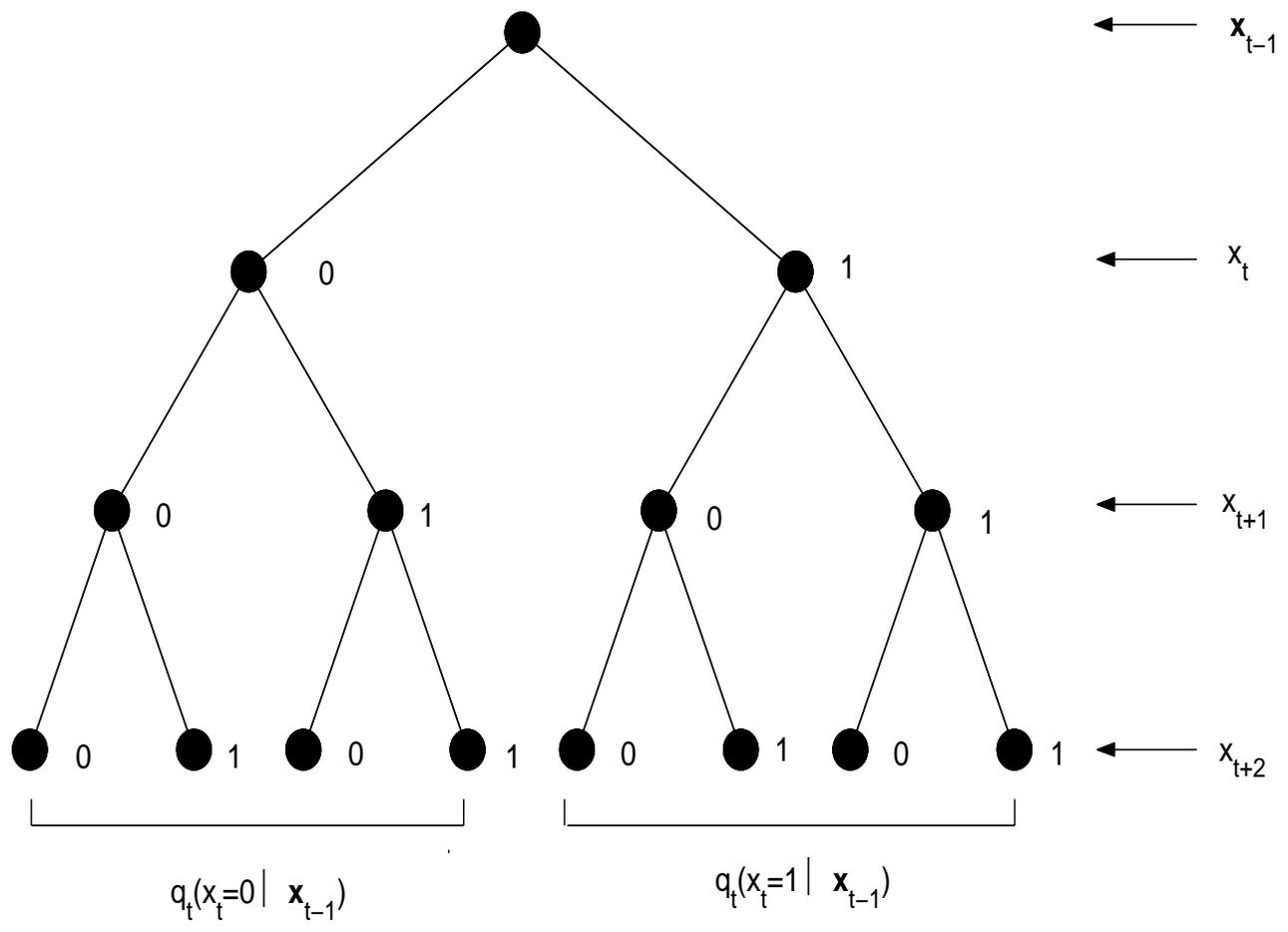
- **New target distribution:** $\pi_t^*(\mathbf{x}_t) = p(\mathbf{x}_t \mid \mathbf{y}_{t+d})$
- — recall: concurrent target distribution: $\pi_t(\mathbf{x}_t) = p(\mathbf{x}_t \mid \mathbf{y}_t)$
- Sample x_t based on a trial distribution that uses the full information \mathbf{y}_{t+d}

In particular, we can use

$$\begin{aligned} g_t(x_t \mid \mathbf{x}_{t-1}^{(j)}) &= p(x_t \mid \mathbf{x}_{t-1}^{(j)}, \mathbf{y}_{t+d}) \\ &= \int p(x_t, x_{t+1}, \dots, x_{t+d} \mid \mathbf{x}_{t-1}^{(j)}, \mathbf{y}_{t+d}) dx_{t+1} \dots x_{t+d} \end{aligned}$$

Then

$$w_t^{*(j)} \propto w_{t-1}^{*(j)} \frac{\pi_t^*(\mathbf{x}_t)}{\pi_{t-1}^*(\mathbf{x}_{t-1}) g_t(x_t \mid \mathbf{x}_{t-1})} = w_{t-1}^{*(j)} \frac{p(\mathbf{x}_{t-1} \mid \mathbf{y}_{t+d})}{p(\mathbf{x}_{t-1} \mid \mathbf{y}_{t+d-1})}$$



We compare the (improved) exact lookahead weighting with the exact lookahead sampling methods:

Suppose at time t , $(\mathbf{x}_t^{(j)}, w_t^{(j)})$ properly weighted w.r.t. $\pi_t(\mathbf{x}_t) = p(\mathbf{x}_t \mid \mathbf{y}_t)$.

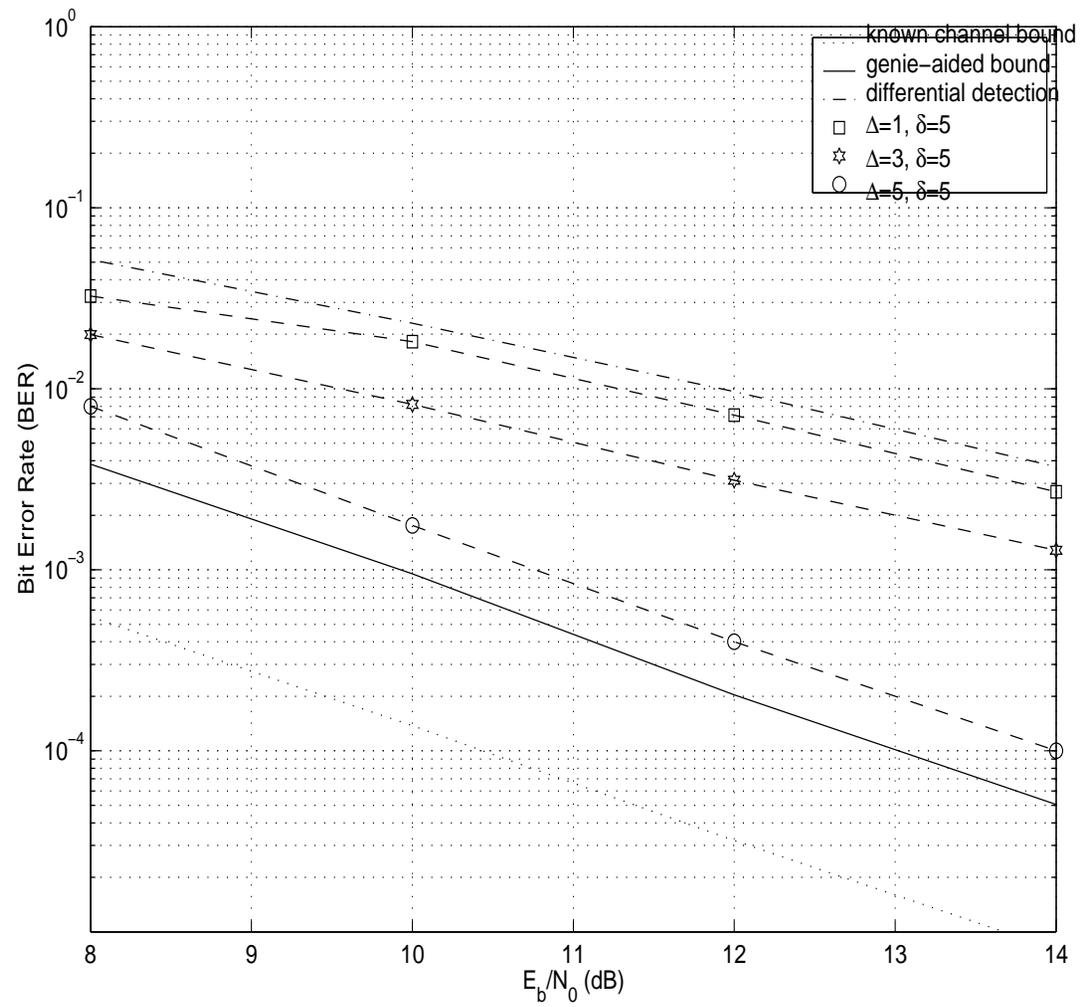
[*Proposition*]

$$\text{Var}[\tilde{w}_{t+d} \mid \mathbf{y}_{t+d}] \geq \text{Var}[w_t^* \mid \mathbf{y}_{t+d}]$$

and

$$\text{Var}[\tilde{w}_{t+d}h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}] \geq \text{Var}[w_t^*h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}]$$

However: Excessive computing cost

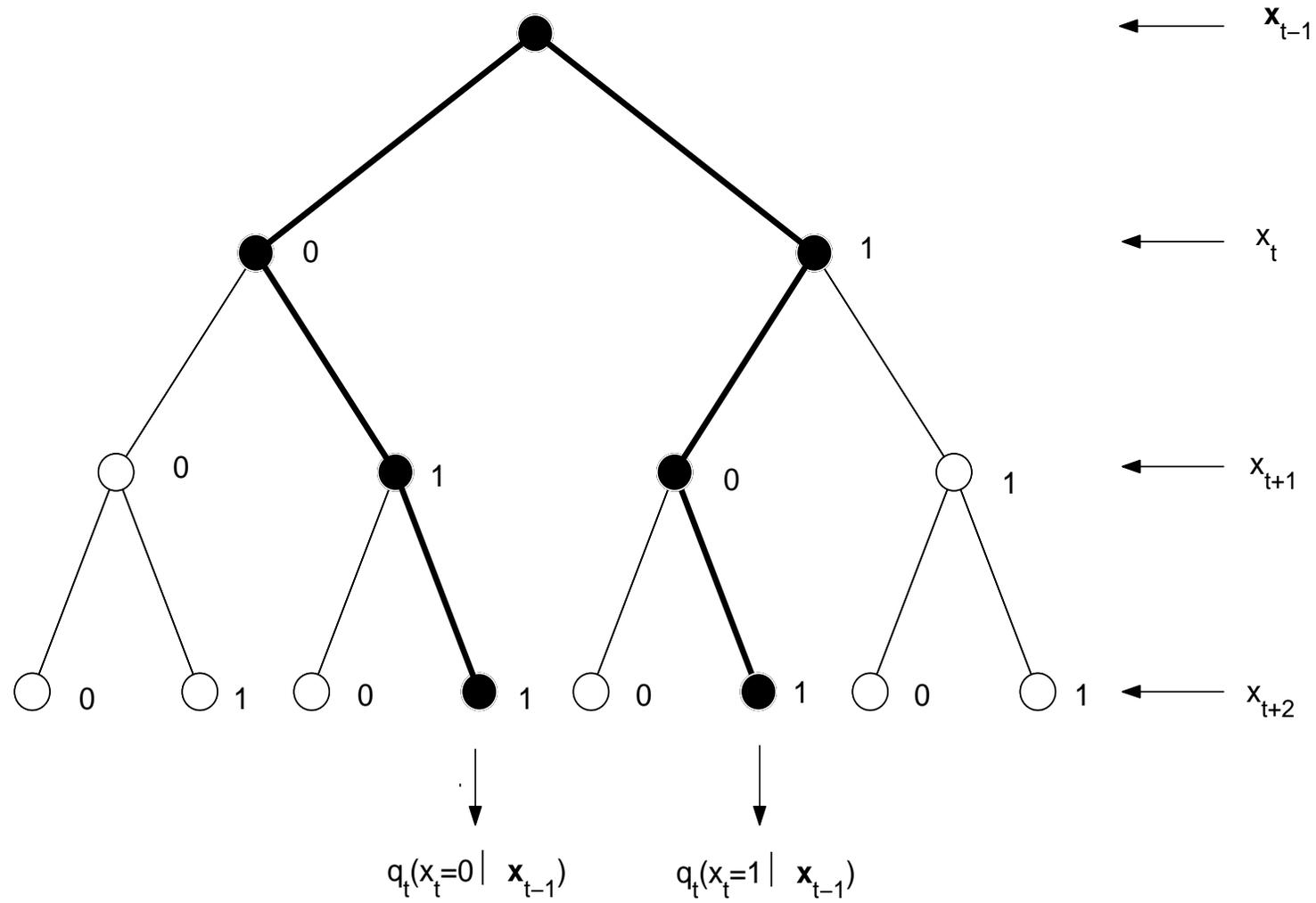


(3) Pilot lookahead sampling:

- Full exploration of future space is too expensive
- Pilots can be used to partially explore the future space
- Limited numbers of pilots are able to obtain useful future information with low computational cost

Specifically,

- Suppose x_t takes J possible values $\{a_1, \dots, a_J\}$
- Starting with each possible $x_t = a_i$, propagate to x_{t+d} with concurrent SMC with optimal sampling distribution.
- Obtain pilot incremental weight for each pilot
- Sample x_t from $\{a_1, \dots, a_J\}$ according to pilot incremental weight
- Update weight



More specifically

- For each $\mathbf{x}_{t-1}^{(j)}$ and each a_i ,
- generate x_{t+1}, \dots, x_{t+d} from

$$\prod_{s=t+1}^{t+d} \pi_s(x_s \mid \mathbf{x}_{t-1}^{(j)}, a_i, x_{t+1}, \dots, x_{s-1})$$

- Obtain the pilot incremental weight

$$U_t^{(i,j)} = \frac{\pi_{t+d}(\mathbf{x}_{t-1}^{(j)}, x_t = a_i, x_{t+1}^{(i,j)}, \dots, x_{t+d}^{(i,j)})}{\pi_{t-1}(\mathbf{x}_{t-1}^{(j)}) \prod_{s=t+1}^{t+d} \pi_s(x_s^{(i,j)} \mid \mathbf{x}_{t-1}^{(j)}, x_t = a_i, x_{t+1}^{(i,j)}, \dots, x_{s-1}^{(i,j)})}$$

- sample $x_t^{(j)}$ from $\{a_1, \dots, a_J\}$ with probability

$$g_t(x_t = a_i \mid \mathbf{x}_{t-1}^{(j)}) = \frac{U_t^{(i,j)}}{\sum_{k=1}^J U_t^{(k,j)}}$$

- New weight

$$w_t^{(j)} = w_{t-1}^{(j)} \frac{\pi_t(\mathbf{x}_t^{(j)})}{\pi_{t-1}(\mathbf{x}_{t-1}^{(j)}) g_t(x_t = x_t^{(j)} \mid \mathbf{x}_{t-1})}$$

Note:

- (x_t, w_t) is properly weighted w.r.t $p(x_t | \mathbf{y}_t)$.
- even though sampling is done with future information
- Inference using (x_t, w_t) is not efficient

$$\frac{\sum h(x_t^{(j)})w_t^{(j)}}{\sum w_t^{(j)}} \approx E(h(x_t) | \mathbf{y}_t)$$

Remedies:

- To make inference on x_t at time $t+d$, calculate

$$w_{t+d}^{**(j)} = w_{t-1}^{(j)} \sum_{i=1}^J U_t^{(i,j)}$$

- Then $(x_t^{(j)}, w_{t+d}^{**(j)})$ is properly weighted w.r.t. $p(x_t | \mathbf{y}_{t+d})$.
- To make inference on x_t at time $t+d$, use $(x_t^{(j)}, w_{t+d}^{**(j)})$

Further improvement: Multi-pilot lookahead sampling: using multiple pilots per a_i .

Comparison with exact lookahead sampling:

[*Proposition*]

$$0 \leq \text{Var}(w_t^{**}) - \text{Var}(w_t^*) \sim O(1/K)$$

where K is number of pilots per a_i and w_t^* is the weight of exact lookahead sampling.

Recall:

$$\text{Var}(w_t^*) \geq \text{Var}(\tilde{w}_t)$$

where \tilde{w}_t is the weight of improved delay weighting algorithm.

Deterministic pilot lookahead

- Pilots need not be random
- A better pilot might be the path that maximize

$$\pi_{t+d}(x_{t+1}, \dots, x_{t+d} \mid \mathbf{x}_{t-1}^{(j)}, x_t = a_i)$$

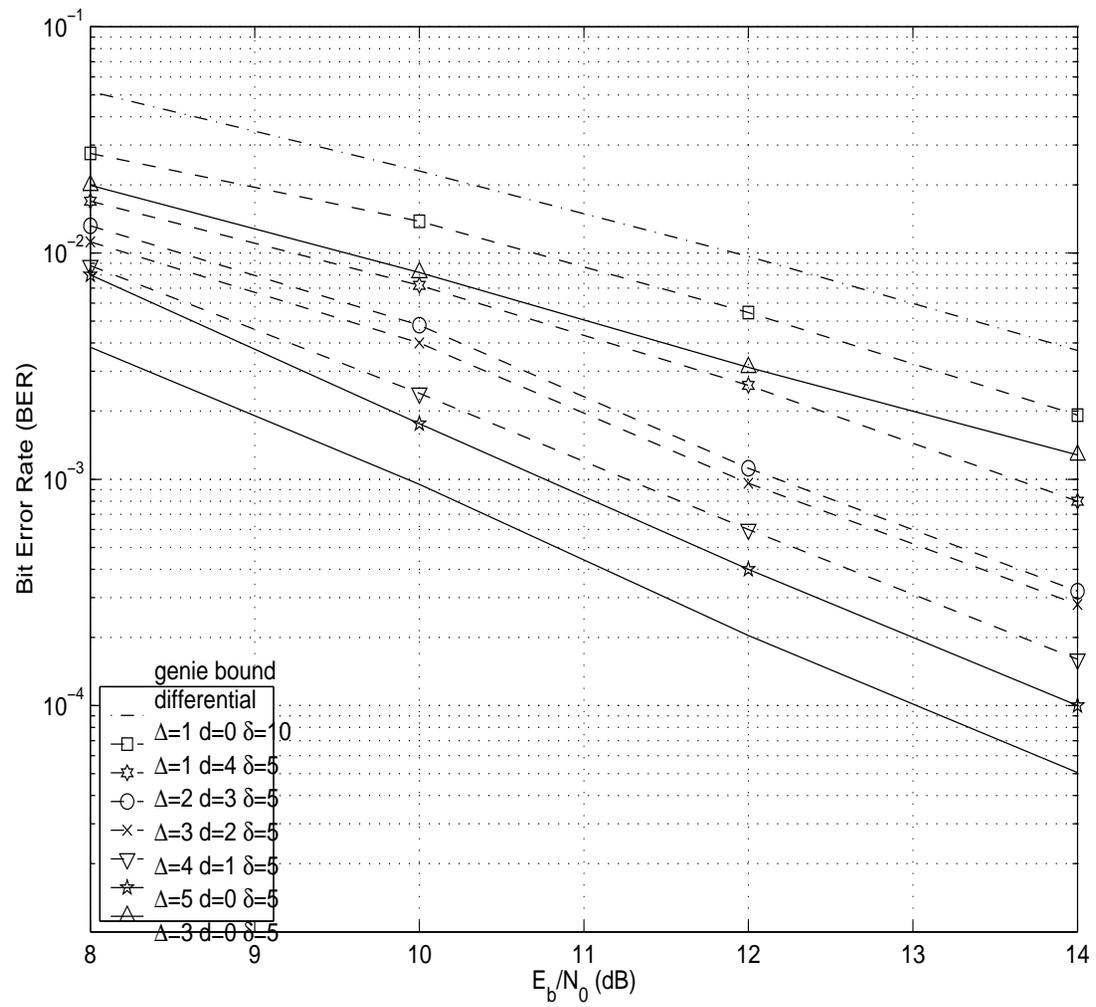
- The true maximum is too expensive to get
- A greedy sequential search: for $s = t + 1, \dots, t + d$

$$x_s^{(i,j)} = \arg \max_{x_s} \pi_s(x_s \mid \mathbf{x}_{t-1}^{(j)}, x_t = a_i, x_{t+1}^{(i,j)}, \dots, x_{s-1}^{(i,j)})$$

●

$$g_t(x_t = a_i \mid \mathbf{x}_{t-1}^{(j)}) \propto \frac{\pi_{t+d}(\mathbf{x}_{t-1}^{(j)}, x_t = a_i, x_{t+1}^{(i,j)}, \dots, x_{t+d}^{(i,j)})}{\pi_{t-1}(\mathbf{x}_{t-1}^{(j)})}$$

- Often, the samples are better than (random) one-pilot lookahead sampling.



(4) Adaptive Sampling

- Sample from a simple trial distribution when information is strong.
- Sample from a better trial distribution (e.g. lookahead) when information is weak.

Recall:

$$E \left[\hat{h}_{t+d} - h(\mathbf{x}_t) \right]^2 = E \left[\hat{h}_{t+d} - E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}) \right]^2 + E \left[E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}) - h(\mathbf{x}_t) \right]^2$$

With finite MC samples, the first term may increase as d increases.

Exact lookahead weighting: comparing $t + d - 1$ and $t + d$

$$\hat{h} = \frac{1}{m} \sum_{j=1}^m w_{t+d-1}^{(j)} h(\mathbf{x}_t^{(j)}) \rightarrow E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d-1})$$

$$\tilde{h} = \frac{1}{m} \sum_{j=1}^m w_{t+d}^{(j)} h(\mathbf{x}_t^{(j)}) \rightarrow E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d})$$

[*Proposition*]: **when**

$$E \left[\text{Var} \left(\tilde{h} \mid \mathbf{x}_{t+d-1}^{(1:m)}, \mathbf{y}_{t+d-1} \right) \mid \mathbf{y}_{t+d-1} \right] \geq 2 \text{Var} \left[E(h(\mathbf{x}_t) \mid \mathbf{y}_{t+d}) \mid \mathbf{y}_{t+d-1} \right]$$

we have

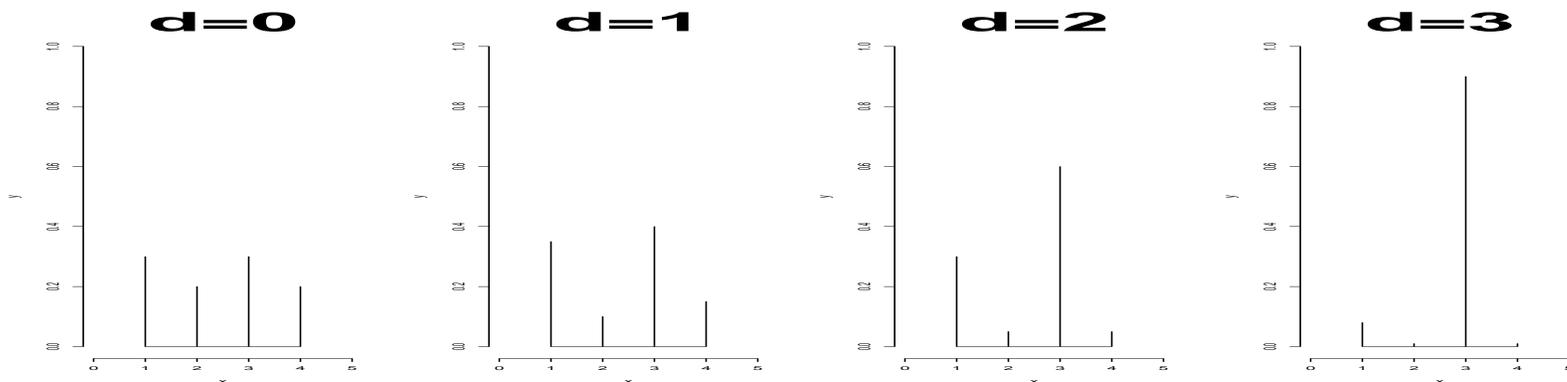
$$E \left[\left(\hat{h} - h(\mathbf{x}_t) \right)^2 \mid \mathbf{y}_{t+d-1} \right] \geq E \left[\left(\tilde{h} - h(\mathbf{x}_t) \right)^2 \mid \mathbf{y}_{t+d-1} \right].$$

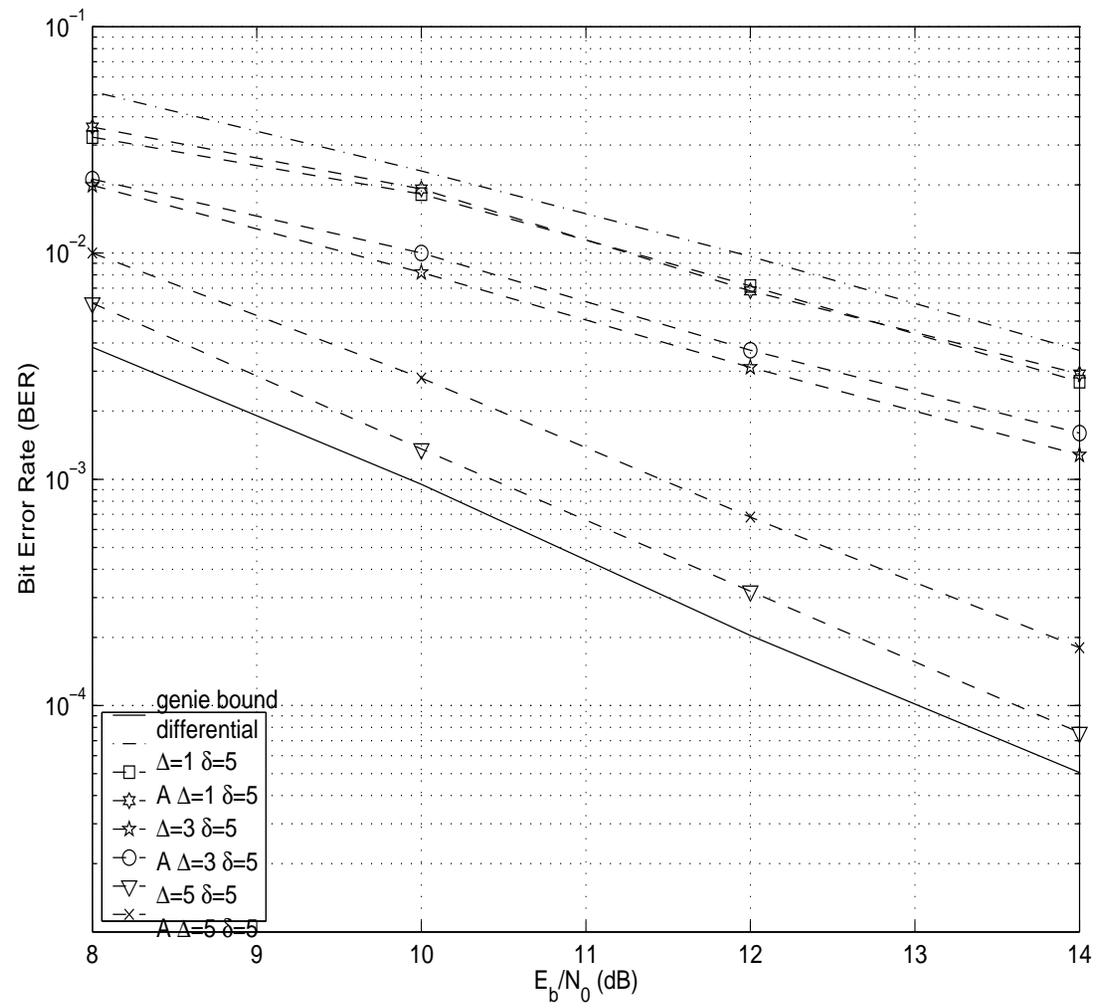
Remarks:

- When $p(\mathbf{x}_t \mid \mathbf{y}_{t+d-1}) = p(\mathbf{x}_t \mid \mathbf{y}_{t+d})$, the result holds.
- In general, the condition is difficult to check
- Instead, we check if the information is strong:
 - Specifically, iteratively try $d = 1, 2, \dots, d_{max}$. Stop when

$$\max_{i_0} \left\{ \hat{\pi}_{t+d}(x_t = a_{i_0}) \right\} \doteq \max_{i_0} \left\{ \frac{\sum_j w_{t-1}^{(j)} U_t^{(i_0,j)}}{\sum_{i,j} w_{t-1}^{(j)} U_t^{(i,j)}} \right\} > p_0,$$

for $p_0 > 0$ but close to 1. (discrete state space)





Other applications:

- Multi-target tracking in clutter
- Self-avoiding walks modelling protein structure
- Generating samples of diffusion bridges
- Signal processing in more complex fading channels