

On Graphical Models for Communications and Machine Learning: Algorithms, Bounds, and Analog Implementation

PhD Defense

Justin Dauwels

www.dauwels.com

Signal and Information Processing Laboratory
ETH Zurich

December 1

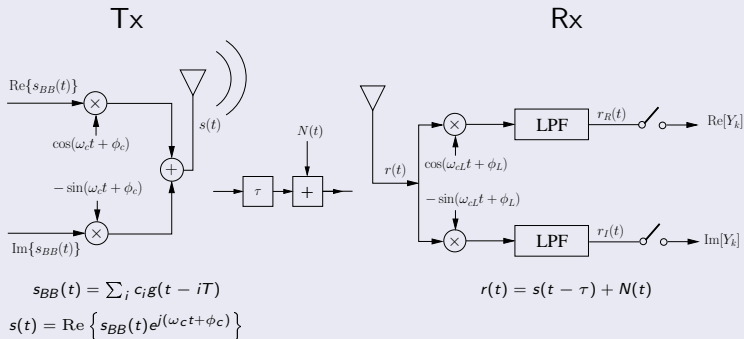


Two threads

- 1 **Particular** problem of carrier-phase synchronization in single-carrier communications systems.
- 2 **Message-passing algorithms** for various applications.

Single-Carrier Communications System

Block diagram



τ : channel delay

T : symbol length

$s_{BB}(t)$: baseband signal

ω_c : carrier frequency

$g(t)$: pulse shape

$s(t)$: passband signal

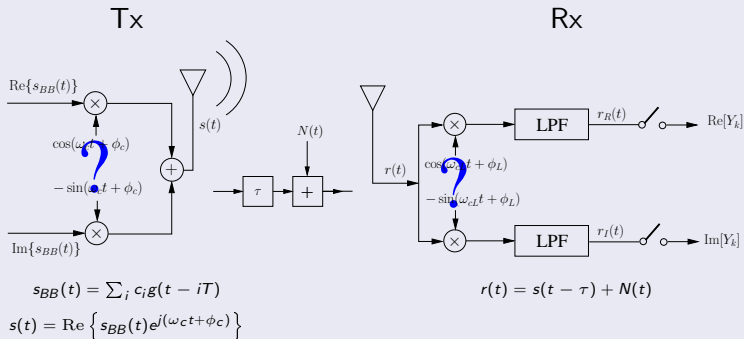
ϕ_c : carrier phase

c_i : data symbols

$N(t)$: AWGN

Single-Carrier Communications System

Block diagram



τ : channel delay

T : symbol length

$s_{BB}(t)$: baseband signal

ω_c : carrier frequency

$g(t)$: pulse shape

$s(t)$: passband signal

ϕ_c : carrier phase

c_i : data symbols

$N(t)$: AWGN

Questions

Modeling

Which **physical mechanisms** are responsible for (phase) noise?
How can (phase) noise be **modeled**?

Algorithms

How can phase-estimation algorithms **systematically** be derived?

Performance limits

How well can the (noisy) carrier phase be **estimated**?
How much does the **information rate** decrease due to phase noise?

Contributions

Modeling

Which **physical mechanisms** are responsible for (phase) noise?
How can (phase) noise be **modeled**?

Simple intuitive model for phase noise.

Algorithms

How can phase-estimation algorithms **systematically** be derived?

As message passing on factor graph of the system a hand.

Performance limits

How well can the (noisy) carrier phase be **estimated**?
How much does the **information rate** decrease due to phase noise?

Computation of Cramér-Rao bounds/information rates/capacities.

Channel model

Model for single-carrier system with slowly varying phase offset

$$Y_k = X_k e^{j\Theta_k} + W_k, \quad W_k \sim \mathcal{CN}_{0, \sigma_N^2}.$$

Constant-phase model

$$\Theta_k = \Theta \in [0, 2\pi).$$

Random-walk phase model

$$\Theta_k = (\Theta_{k-1} + N_k) \bmod 2\pi, \quad N_k \sim \mathcal{N}_{0, \sigma_N^2}.$$

σ_N^2 and σ_W^2 are assumed to be **known**.

The input symbols X_k are protected by an **error-correcting code**.

Contributions

Modeling

Which physical mechanisms are responsible for (phase) noise?
How can (phase) noise be modeled?

Simple intuitive model for phase noise.

Algorithms

How can phase-estimation algorithms **systematically** be derived?

As message passing on factor graph of the system a hand.

Performance limits

How well can the (noisy) carrier phase be estimated?
How much does the information rate decrease due to phase noise?

Computation of Cramér-Rao bounds/information rates/capacities.

Algorithms for joint decoding and phase estimation

Estimation task

Given a block of **observations** $Y \triangleq (Y_1, Y_2, \dots, Y_N)$, infer:

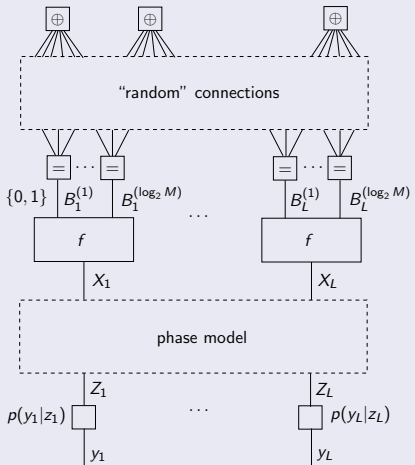
- the **coded symbols** $X \triangleq (X_1, X_2, \dots, X_N)$
- the **phase** $\Theta \triangleq (\Theta_1, \Theta_2, \dots, \Theta_N)$.

Derivation of message-passing estimation algorithms [Wiberg, 1996]

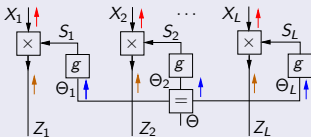
- 1 Draw factor graph of joint pdf $p(x, y, \theta)$.
- 2 Apply sum-product rule at each node.
- 3 If sum-product rule is **infeasible** at a certain node, then apply an **approximation** = choose appropriate **message types**.
- 4 Choose an update schedule.

Factor graphs

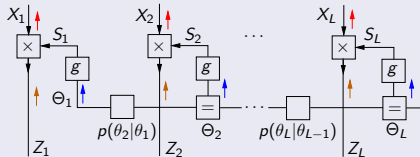
Factor graph of $p(x, y, \theta)$



Constant-phase model



Random-walk phase model

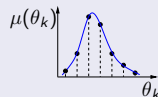


Message Types

$$\mu(x_k) \propto \int_0^{2\pi} \mu(\theta_k) e^{-|x_k e^{j\theta_k} - y_k|^2 / 2\sigma_W^2} d\theta_k$$

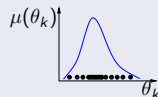
Numerical integration

$$\mu(x_k) \propto \sum_i \mu(\hat{\theta}_k^{(i)}) e^{-|x_k e^{j\hat{\theta}_k^{(i)}} - y_k|^2 / 2\sigma_W^2}$$



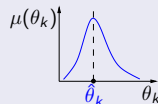
Particle method

$$\mu(x_k) \propto \sum_i e^{-|x_k e^{j\hat{\theta}_k^{(i)}} - y_k|^2 / 2\sigma_W^2}$$



Decision based

$$\mu(x_k) \approx e^{-|x_k e^{j\hat{\theta}_k} - y_k|^2 / 2\sigma_W^2}$$



Message-passing view

Unification

- Existing algorithms as **message passing on factor graphs**.
- Examples:
Particle Filtering, Markov-Chain Monte Carlo, Iterative Conditional Modes, Gradient Methods, Stochastic Approximation, Expectation maximization, SAGE, Natural-Gradient Methods, Backpropagation Algorithm, etc.
- Generic local **message-update rules** for each approach.

Motivation

- Divide and conquer: **local** computations.
- **Systematic** exploration of novel algorithms.
- Plug and play: **tabulate** message-update rules.

Joint work with Sascha Kori

Contributions

Modeling

Which physical mechanisms are responsible for (phase) noise?
How can (phase) noise be modeled?

Simple intuitive model for phase noise.

Algorithms

How can phase-estimation algorithms systematically be derived?

As message passing on factor graph of the system a hand.

Performance limits

How well can the (noisy) carrier phase be estimated?
How much does the information rate decrease due to phase noise?

Computation of Cramér-Rao bounds/information rates/capacities.

Cramér-Rao-type bounds

What?

Lower bounds on the mean-squared estimation error (MSE)

e.g., $\text{MSE}(\theta) = E_{Y;\theta} \left[(\theta - \hat{\theta}(Y))^2 \right]$.

Motivation

Assessment of **practical** estimators (in terms of MSE).

Three different types

- **Standard** Cramér-Rao bounds: parameters
- **Bayesian** Cramér-Rao bounds: random variables
- **Hybrid** Cramér-Rao bounds: parameters and random variables

Bayesian Cramér-Rao bound

Given: joint pdf $p(x, y)$ of $x \triangleq (x_1, \dots, x_M)$ and $y \triangleq (y_1, \dots, y_N)$.

Lower bound for the MSE $E_{X_k Y} [(X_k - \hat{x}_k(Y))^2]$?

From marginal

$$E_{X_k Y} [(X_k - \hat{x}_k(Y))^2] \geq J_k^{-1},$$

with $J_k \triangleq -E_{X_k Y} \left[\frac{\partial^2}{\partial^2 x_k} \log p(x_k, y) \right]$.

From joint pdf

$$E_{X_k Y} [(X_k - \hat{x}_k(Y))^2] \geq [J^{-1}]_{kk},$$

with $J_{ij} \triangleq -E_{XY} \left[\frac{\partial^2}{\partial x_i \partial x_j} \log p(x, y) \right]$.

BCRB from marginal is **tighter** than from joint pdf, but more difficult to compute.

Algorithms

Overview

We propose efficient and simple message-passing algorithms:

- for computing standard, Bayesian, hybrid Cramér-Rao bounds
- following both strategies.

Other applications

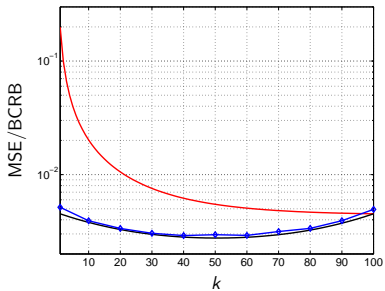
- Other types of bounds, e.g., Weiss-Weinstein (discrete variables), Bhattacharyya, etc.
- Information geometry: natural-gradient based algorithms
- Machine learning: computation of Fisher kernels.

Results: MSE and BCRB for Θ

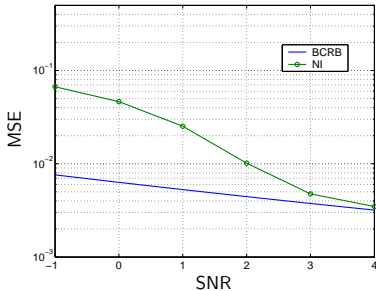
Infer $X = X_1, \dots, X_L$ and $\Theta = \Theta_1, \dots, \Theta_L$ from $Y = Y_1, \dots, Y_L$

Random walk phase model with $\sigma_W^2 = 10^{-4} \text{rad}^2$ with $L = 100$

MSE/BCRB for Θ_k (SNR = 4dB)



Average MSE/BCRB



Contributions

Modeling

Which physical mechanisms are responsible for (phase) noise?
How can (phase) noise be modeled?

Simple intuitive model for phase noise.

Algorithms

How can phase-estimation algorithms systematically be derived?

As message passing on factor graph of the system a hand.

Performance limits

How well can the (noisy) carrier phase be estimated?
How much does the information rate decrease due to phase noise?

Computation of Cramér-Rao bounds/information rates/capacities.

Information rate: introduction

Objective

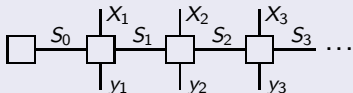
Information rate $I(X; Y) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} I(X_1, \dots, X_n; Y_1, \dots, Y_n)$ between input process $X = (X_1, X_2, \dots)$ and output process $Y = (Y_1, Y_2, \dots)$ of **time-invariant** discrete-time channel with **memory**.

State-space representation

An ergodic stochastic process $S = (S_0, S_1, S_2, \dots)$ such that

$$p(x^n, y^n, s_0^n) = p(s_0) \prod_{k=1}^n p(x_k, y_k, s_k | s_{k-1})$$

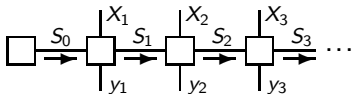
for all $n > 0$ and with $p(x_k, y_k, s_k | s_{k-1})$ not depending on k .



Basic principle

Algorithm

- 1 Sample two “very long” sequences x^n and y^n .
- 2 Compute $\log p(x^n)$, $\log p(y^n)$, and $\log p(x^n, y^n)$.
- 3 $\hat{I}(X; Y) \triangleq \frac{1}{n} \log p(x^n, y^n) - \frac{1}{n} \log p(x^n) - \frac{1}{n} \log p(y^n)$.



Discrete input space \mathcal{X} and state-space \mathcal{S} [e.g., Arnold et al.]

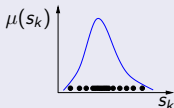
Forward sum-product sweep = forward BCJR-recursion

Continuous input space \mathcal{X} and state-space \mathcal{S}

Forward sum-product sweep by **particle filtering**.

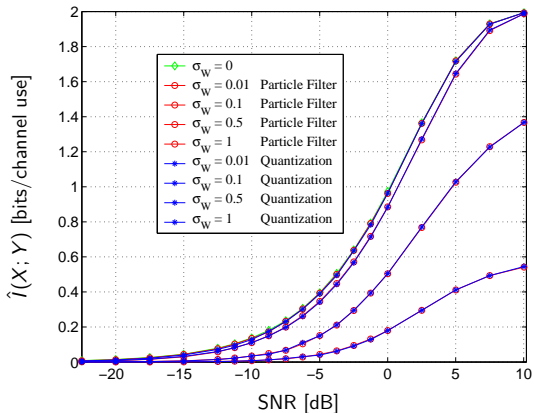
Expression $p(x_k, s_k | s_{k-1})$ not required!

E.g., stochastic differential/difference equation.



Numerical results

Random-walk phase model with i.u.d. 4-PSK input symbols X



Contributions

Modeling

Which physical mechanisms are responsible for (phase) noise?
How can (phase) noise be modeled?

Simple intuitive model for phase noise.

Algorithms

How can phase-estimation algorithms systematically be derived?

As message passing on factor graph of the system a hand.

Performance limits

How well can the (noisy) carrier phase be estimated?
How much does the information rate decrease due to phase noise?

Computation of Cramér-Rao bounds/information rates/capacities.

Capacity of continuous memoryless channel

Definition

Given: memoryless channel with law $p(y|x)$

$$C(X; Y) \triangleq \sup_{p(x)} \int_x \int_y p(x)p(y|x) \log \frac{p(y|x)}{p(y)} dx dy \triangleq \sup_{p(x)} I(X; Y)$$

Discrete input alphabet \mathcal{X}

Blahut-Arimoto algorithm.

Continuous input alphabet \mathcal{X}

Particle-based approach: $p(x) \approx \{(\hat{x}_1, w_1), (\hat{x}_2, w_2), \dots, (\hat{x}_N, w_N)\}$.

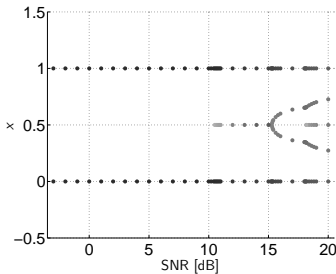
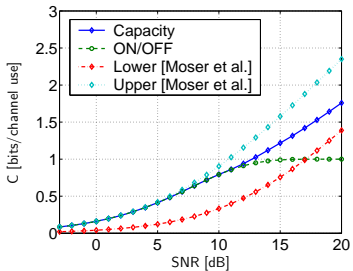
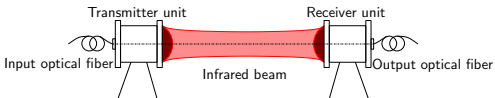
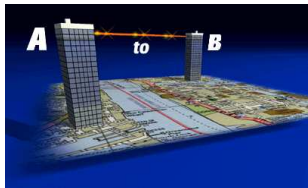
$$w^{(k)} \triangleq \operatorname{argmax}_w I(\hat{x}^{(k-1)}, w) \quad (\text{Blahut-Arimoto})$$

$$\hat{x}^{(k)} \triangleq \operatorname{argmax}_{\hat{x}} I(\hat{x}, w^{(k)}) \quad (\text{gradient method})$$

Method for channels with memory currently in development.

Results: Gaussian channel

$$Y_k = X_k + N_k \text{ with } N_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_N^2) \text{ and } \Pr[0 \leq X_k \leq 1] = 1$$



Analog circuit for PN-synchronization

Objective

Analog circuit that locks unto **pseudo-random sequences**

Method

Discrete-time **message-passing algorithm** for synchronization to LFSR-sequences converted into continuous-time.

Practical result

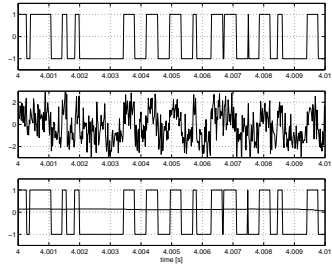
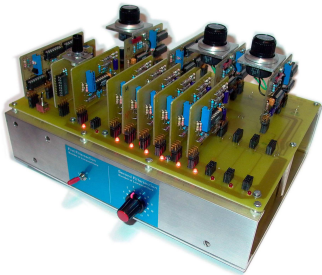
Practical circuit built and tested: it works!

Theoretical result

Connection between entrainment and ideas from estimation theory (“message passing”).

Joint work with M. Frey, N. Gershenfeld, T. Koch, P. Merkli, B. Vigoda.

Results



SNR = 0dB

*Photograph and measurements by M. Frey and P. Merkli.
Hardware built by T. Schaerer.*

Contributions

- Simple model for phase noise.
- Framework for deriving inference algorithms (joint work):
 - Factor graph = graphical representation of system
 - Algorithm = updating messages on factor graph.
- Message-passing algorithms for computing:
 - information rates
 - channel capacities
 - Cramér-Rao-type bounds.
- Analog circuit for PN-synchronization (joint work)
= message-passing algorithm as dynamical system.

Outlook

- Lower bounds on the MSE for **discrete** variables.
- **Dynamical** systems for detection/estimation:
 - Analog electronic circuits
 - Opto-electronic systems
 - Quantum-computing systems.
- Novel **applications** of message-passing methods
 - Information geometry
 - Kernel methods

On Graphical Models for Communications and Machine Learning: Algorithms, Bounds, and Analog Implementation

PhD Defense

Justin Dauwels

www.dauwels.com

Signal and Information Processing Laboratory
ETH Zurich

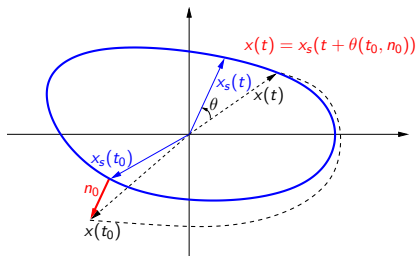
December 1



Phase Noise in Free-Running Clocks

Perturbed autonomous system

$$\frac{dx}{dt} = f(x) + N(t), \quad x \in \mathbb{R}^n, \quad N(t) \text{ is "noise".}$$



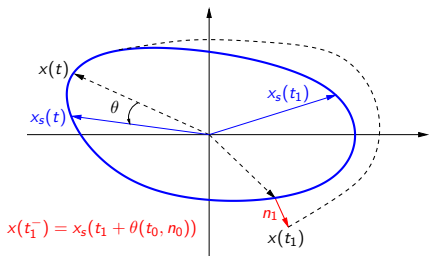
Phase offset due to "small" perturbation at $t = t_0$

$$\theta(t_0, n_0) \approx \gamma(x_s(t_0)) \cdot n_0.$$

Phase Noise in Free-Running Clocks

Perturbed autonomous system

$$\frac{dx}{dt} = f(x) + N(t), \quad x \in \mathbb{R}^n, \quad N(t) \text{ is "noise".}$$



Phase offset due to "small" perturbation at $t = t_1 \gg t_0$

$$\theta(t_0, n_0, t_1, n_1) \approx \gamma(x_s(t_0)) \cdot n_0 + \gamma(x_s(t_1 + \theta(x_s(0), n_0))) \cdot n_1.$$

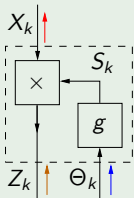
Sum-product rule

Example

Sum-product rule

$$\mu(x_k) \propto \int_0^{2\pi} \mu(\theta_k) \mu(x_k e^{j\theta_k}) d\theta_k,$$

$$\propto \int_0^{2\pi} \mu(\theta_k) e^{-|x_k e^{j\theta_k} - y_k|^2 / 2\sigma_w^2} d\theta_k$$



Intractable integral!

Unification

Particle methods

Importance Sampling, Markov-Chain Monte Carlo, Metropolis-Hastings Algorithm, Gibbs Sampling, Simulated Annealing, Particle Filtering

Decision based

Iterative Conditional Modes, Gradient Methods, Stochastic Approximation, Expectation maximization, SAGE, Gradient EM, Natural-Gradient Methods, Backpropagation Algorithm

Combinations

Monte-Carlo EM, Stochastic EM

Interpretation as message passing on factor graphs

Identified generic local message-update rules for each approach.

Joint work with Sascha Korf

Why we care . . .

Divide and conquer

Global estimation/detection problem accomplished by **simple local** computations. Complicated mathematical derivations avoided.

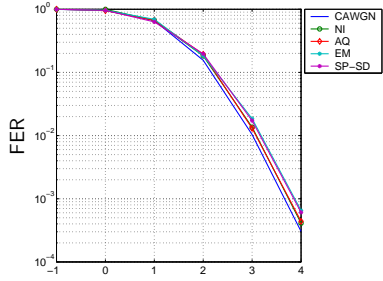
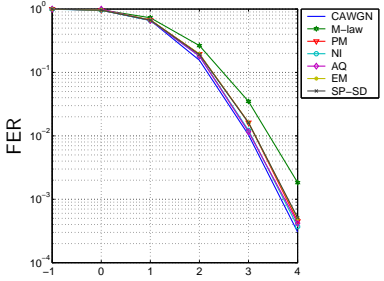
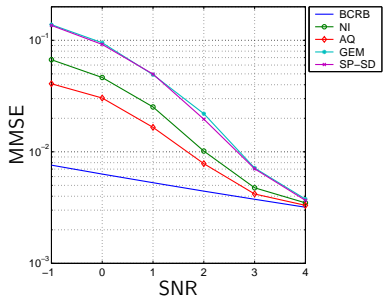
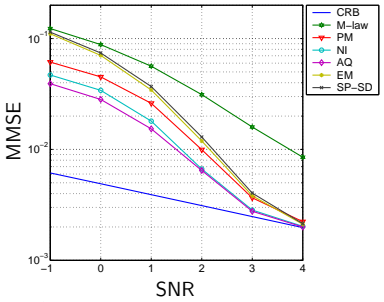
Disciplined approach

Deriving novel algorithms **systematically** by listing possible message update rules at each node in the graph.

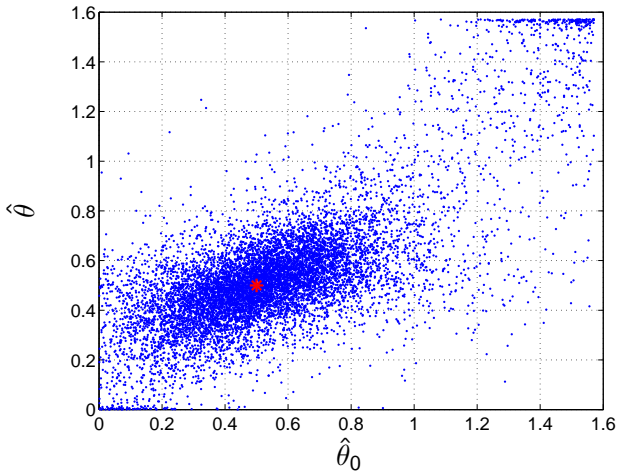
Plug and play

Deriving novel algorithms by combining **tabulated** message update rules. Efficient use of earlier work.

Phase estimation: results



EM: initial estimate vs. final estimate



Gibbs sampling

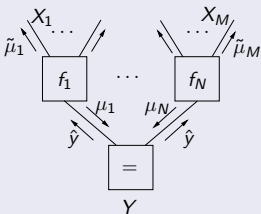
Message passing view

- 1 Select variable (EQU constraint node) Y in FG of f .
- 2 EQU node Y generates message \hat{y} by sampling from:

$$f(y) \triangleq \frac{\mu_1(y) \dots \cdot \mu_N(y)}{\sum_y \mu_1(y) \dots \cdot \mu_N(y)},$$

and broadcasts \hat{y} to its neighboring nodes f_k ($k = 1, \dots, M$).

- 3 Nodes f_k update messages $\tilde{\mu}_k$ by applying sum-product rule with as incoming messages the samples \hat{y} and \hat{x}_ℓ ($\ell = 1, \dots, M$).

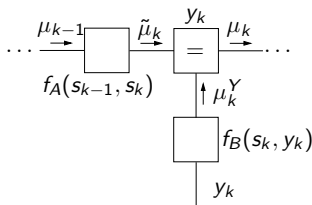


Particle filtering

Particle filtering (or “sequential Monte-Carlo integration”)
= forward-only message passing in a state-space model:

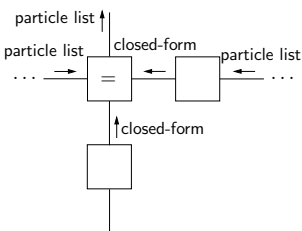
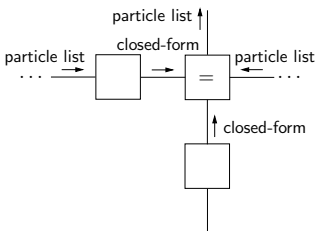
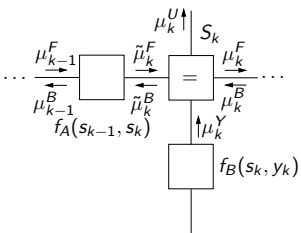
$$f(s_0, s_2, \dots, s_N, y_1, y_2, \dots, y_N) \triangleq f_A(s_0) \prod_{k=1}^N f_A(s_{k-1}, s_k) f_B(s_k, y_k),$$

where messages are represented by lists of samples.



$\tilde{\mu}_k$ is obtained from μ_{k-1} by **weighted** or **unweighted sampling**.
 μ_k is generated from $\tilde{\mu}_k$ by **importance sampling**.

Smoothing



MCMC

Algorithm

- 1 Choose an initial value \hat{x} .
- 2 Sample \hat{y} from $q(y|\hat{x})$.
- 3 Set

$$\hat{x} \triangleq \hat{y} \quad \text{with probability } p$$

where

$$p \triangleq \min \left\{ \frac{f(\hat{y})}{f(\hat{x})}, 1 \right\}$$

- 4 Iterate 2–3 a sufficient number of times.

MCMC

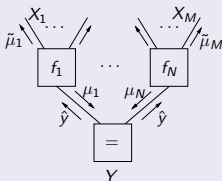
Message passing view

- 1 Select variable (equality constraint node) Y in FG of f .
- 2 Edge Y generates the message \hat{y}^{new} by sampling from $q(y|\hat{y})$.
- 3 Set $\hat{y} \triangleq \hat{y}^{\text{new}}$ with probability p where

$$p \triangleq \min \left\{ \frac{f(\hat{y}^{\text{new}})}{f(\hat{y})}, 1 \right\} \quad \text{with} \quad f(y) \triangleq \frac{\mu_1(y) \dots \mu_N(y)}{\sum_y \mu_1(y) \dots \mu_N(y)}.$$

The message \hat{y} is broadcast to the neighboring nodes f_k .

- 4 Nodes f_k update outgoing messages $\tilde{\mu}$ by applying the SP rule with as incoming messages the samples \hat{y} and \hat{x}_ℓ .



Simulated Annealing

Objective

- to sample from a multivariate function $f(x_1, \dots, x_N)$,
- to find the mode of the function f .

Algorithm

- 1 Choose an initial value $(\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N)$.
- 2 Choose an initial value α (e.g., $\alpha = 0.1$).
- 3 Sample a new value \hat{y} from $q(y|\hat{x})$.
- 4 Set $\hat{x} \stackrel{\Delta}{=} \hat{y}$ with probability p , where

$$p \stackrel{\Delta}{=} \min \left\{ \left(\frac{f(\hat{y})}{f(\hat{x})} \right)^\alpha, 1 \right\}$$

- 5 Iterate 3–4 a “large” number of times.
- 6 Increase α according to some schedule.
- 7 Iterate 5–6 until convergence or until available time is over.

Single value approximation

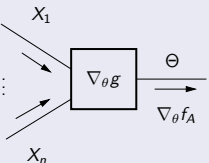
Integral-product rule evaluated by means of hard decision

$$\mu_{f \rightarrow Y}(y) \propto f(y, \hat{x}_1, \dots, \hat{x}_N),$$

where \hat{x}_k is a hard estimate of X_k , representing the message $\mu_{X_k \rightarrow f}$.

Gradient descent / sum-product

$$\nabla_{\theta} f_A(\theta) \propto \sum_{x_1, \dots, x_n} \nabla_{\theta} g(x_1, \dots, x_n, \theta) \cdot \prod_{\ell=1}^n \mu_{X_{\ell} \rightarrow g}(x_{\ell}).$$



Expectation Maximization: General problem

$$\theta_{\max} \triangleq \operatorname{argmax}_{\theta} f(\theta)$$

θ takes values in \mathbb{R} or \mathbb{R}^n

$$f(\theta) \triangleq \int_x f(x, \theta) dx$$

$\int_x g(x) dx$ stands for summation or integration.

In principle

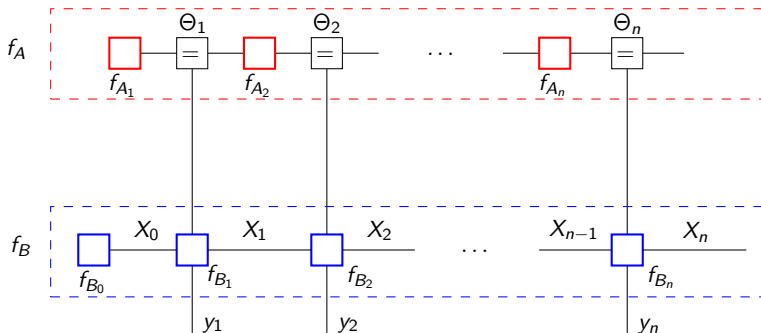
- 1 Determine $f(\theta)$ by sum-product message passing
- 2 $\theta_{\max} \triangleq \operatorname{argmax}_{\theta} f(\theta)$ by max-product message-passing

Often infeasible, since

- Sum-product rule may lead to intractable integrals
- Maximization step may be infeasible.

Parameter estimation in state-space model

$$f(x, \theta) = f_A(\theta_1) \prod_{k=1}^{n-1} f_A(\theta_k, \theta_{k+1}) \cdot f_B(x_0) \prod_{k=1}^n f_B(x_{k-1}, x_k, \theta_k, y_k)$$



Expectation Maximization

- 1 Make initial guess $\theta^{(0)}$
- 2 **Expectation** step

$$f^{(\ell)}(\theta) \triangleq \int_{\mathbf{x}} f(\mathbf{x}, \hat{\theta}^{(\ell)}) \log f(\mathbf{x}, \theta) d\mathbf{x}$$

- 3 **Maximization** step

$$\theta^{(\ell+1)} \triangleq \underset{\theta}{\operatorname{argmax}} f^{(\ell)}(\theta)$$

- 4 **Repeat** 2–3 until convergence.

EM as message passing (2)

Expectation step

$$f^{(\ell)}(\theta) \triangleq \int_{\mathbf{x}} f(\mathbf{x}, \hat{\theta}^{(\ell)}) \log f(\mathbf{x}, \theta) d\mathbf{x}$$

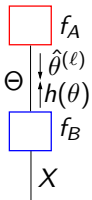
Maximization step

$$\theta^{(\ell+1)} \triangleq \operatorname{argmax}_{\theta} f^{(\ell)}(\theta)$$

$$\begin{aligned} \hat{\theta}^{(\ell+1)} &= \operatorname{argmax}_{\theta} \int_{\mathbf{x}} f(\mathbf{x}, \hat{\theta}^{(\ell)}) \log f(\mathbf{x}, \theta) d\mathbf{x} \\ &= \operatorname{argmax}_{\theta} \int_{\mathbf{x}} f_A(\hat{\theta}^{(\ell)}) f_B(\mathbf{x}, \hat{\theta}^{(\ell)}) \log (f_A(\theta) f_B(\mathbf{x}, \theta)) d\mathbf{x} \\ &= \operatorname{argmax}_{\theta} \int_{\mathbf{x}} f_B(\mathbf{x}, \hat{\theta}^{(\ell)}) (\log f_A(\theta) + \log f_B(\mathbf{x}, \theta)) d\mathbf{x} \\ &= \operatorname{argmax}_{\theta} \left(\log f_A(\theta) + \frac{\int_{\mathbf{x}} f_B(\mathbf{x}, \hat{\theta}^{(\ell)}) \log f_B(\mathbf{x}, \theta) d\mathbf{x}}{\int_{\mathbf{x}} f_B(\mathbf{x}, \hat{\theta}^{(\ell)}) d\mathbf{x}} \right) \\ &= \operatorname{argmax}_{\theta} (\log f_A(\theta) + h(\theta)) \end{aligned}$$

EM as message passing (3)

$$f(x, \theta) \triangleq f_A(\theta) f_B(x, \theta)$$



Upwards message $h(\theta)$

$$h(\theta) = \frac{\int_x f_B(x, \hat{\theta}^{(\ell)}) \log f_B(x, \theta) dx}{\int_x f_B(x, \hat{\theta}^{(\ell)}) dx}$$

$$= E_{p_B}[\log f_B(x, \theta)]$$

$$p_B(x|\hat{\theta}^{(\ell)}) \triangleq \frac{f_B(x, \hat{\theta}^{(\ell)})}{\int_x f_B(x, \hat{\theta}^{(\ell)}) dx}$$

Downwards message $\hat{\theta}^{(\ell+1)}$

$$\hat{\theta}^{(\ell+1)} = \underset{\theta}{\operatorname{argmax}} (\log f_A(\theta) + h(\theta))$$

EM as message passing (4)

Remarks

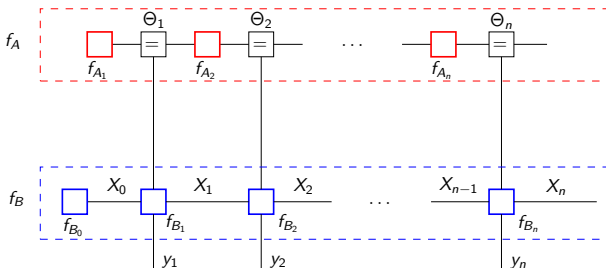
- If $f_A(\theta)$ is **constant**, then normalization may be omitted.
- Message $h(\theta)$ is **not sum-product** message!
- f_A and f_B often have a **“nice” structure**.

EM as message passing (5)

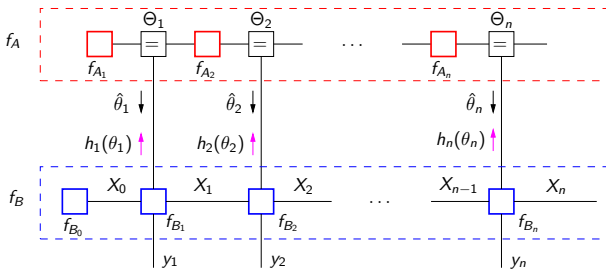
Trellis and state space models

$$f_A(\theta) \triangleq f_{A_1}(\theta_1) f_{A_2}(\theta_1, \theta_2) \dots f_{A_n}(\theta_{n-1}, \theta_n)$$

$$f_B(x, \theta) \triangleq f_{B_0}(x_0) f_{B_1}(x_0, x_1, y_1, \theta_1) f_{B_2}(x_1, x_2, y_2, \theta_2) \dots f_{B_n}(x_{n-1}, x_n, y_n, \theta_n)$$



EM as message passing (6)



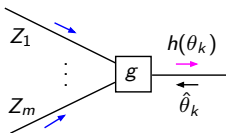
$$h(\theta) = \sum_{\ell=1}^n h_{\ell}(\theta_{\ell}) = \sum_{\ell=1}^n \int_{x_{\ell-1}} \int_{x_{\ell}} p_B(x_{\ell-1}, x_{\ell}, |y, \hat{\theta}) \log f_{B_{\ell}}(x_{\ell-1}, x_{\ell}, y, \theta_{\ell}) dx_{\ell-1} dx_{\ell}$$

$$p_B(x_{\ell-1}, x_{\ell}, |y, \hat{\theta}) = \frac{f_{B_{\ell}}(x_{\ell-1}, x_{\ell}, y, \theta_{\ell}) \mu_{X_{\ell} \rightarrow f_{B_{\ell}}}(x_{\ell}) \mu_{X_{\ell-1} \rightarrow f_{B_{\ell}}}(x_{\ell-1})}{\int_{x_{\ell-1}} \int_{x_{\ell}} f_{B_{\ell}}(x_{\ell-1}, x_{\ell}, y, \theta_k) \mu_{X_{\ell} \rightarrow f_{B_{\ell}}}(x_{\ell}) \mu_{X_{\ell-1} \rightarrow f_{B_{\ell}}}(x_{\ell-1}) dx_{\ell-1} dx_{\ell}}$$

$$(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n)^T = \operatorname{argmax}_{\theta_1, \theta_2, \dots, \theta_n} \left[\log f_{A_1}(\theta_1) + \sum_{\ell=2}^n \log f_{A_{\ell}}(\theta_{\ell-1}, \theta_{\ell}) + \sum_{\ell=1}^n h_{\ell}(\theta_{\ell}) \right]$$

EM as message passing (7)

h -messages



$$\begin{aligned} h(\theta_k) &= \gamma^{-1} \int_{\mathbf{z}} g(\mathbf{z}_1, \dots, \mathbf{z}_m, \hat{\theta}_k) \mu(\mathbf{z}_1) \dots \mu(\mathbf{z}_m) \log g(\mathbf{z}_1, \dots, \mathbf{z}_m, \theta_k) d\mathbf{z} \\ &= \int_{\mathbf{z}} p(\mathbf{z}_1, \dots, \mathbf{z}_m | \hat{\theta}_k) \log g(\mathbf{z}_1, \dots, \mathbf{z}_m, \theta_k) d\mathbf{z} \\ &= \mathbb{E}_{p(\mathbf{z}_1, \dots, \mathbf{z}_m | \hat{\theta}_k)} [\log g(\mathbf{z}_1, \dots, \mathbf{z}_m, \theta_k)] \end{aligned}$$

$$\begin{aligned} p(\mathbf{z}_1, \dots, \mathbf{z}_m | \hat{\theta}_k) &= \gamma^{-1} g(\mathbf{z}_1, \dots, \mathbf{z}_m, \hat{\theta}_k) \mu(\mathbf{z}_1) \dots \mu(\mathbf{z}_m) \\ \gamma &= \int_{\mathbf{z}} g(\mathbf{z}_1, \dots, \mathbf{z}_m, \hat{\theta}_k) \mu(\mathbf{z}_1) \dots \mu(\mathbf{z}_m) d\mathbf{z} \end{aligned}$$

$\mu(\mathbf{z}_k)$ are sum-product messages

Expectation Maximization: Properties

Theorem (Main property)

$$f(\hat{\theta}^{(k+1)}) \geq f(\hat{\theta}^{(k)}).$$

Corollary

The global maximum θ^{\max} of $f(\theta)$ is a fixed point of EM.

Theorem

The fixed points of EM are stationary points of $f(\theta)$.

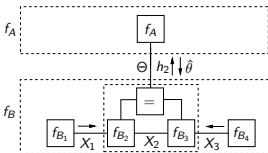
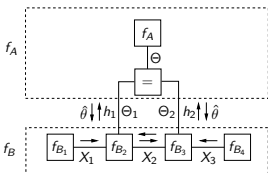
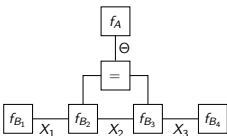
Theorem

A stationary point $\hat{\theta}^{\text{stat}}$ of f is a fixed point of EM, if $\bar{f}(\theta, \hat{\theta}^{\text{stat}})$ with

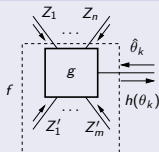
$$\bar{f}(\theta, \theta') \triangleq \sum_x f(x, \theta') \log f(x, \theta),$$

is concave in θ .

EM and compound nodes



Hybrid EM



$$\begin{aligned}
 h(\theta_k) &= \sum_{z_1} \dots \sum_{z_n} p(z_1, \dots, z_n | \hat{\theta}_k) \log f(z_1, \dots, z_n, \theta_k), \\
 &= \gamma^{-1} \sum_{z_1} \dots \sum_{z_n} f(z_1, \dots, z_n, \hat{\theta}_k) \mu(z_1) \cdots \mu(z_n) \\
 &\quad \cdot \log f(z_1, \dots, z_n, \theta_k),
 \end{aligned}$$

with

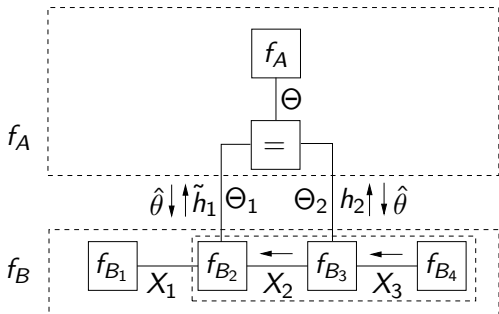
$$\gamma \triangleq \sum_{z_1} \dots \sum_{z_n} f(z_1, \dots, z_n, \hat{\theta}_k) \mu(z_1) \cdots \mu(z_n),$$

and

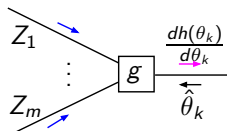
$$\begin{aligned}
 f(z_1, \dots, z_n, \theta_k) &\propto \sum_{z'_1} \dots \sum_{z'_m} g(z_1, \dots, z_n, z'_1, \dots, z'_m, \theta_k) \\
 &\quad \cdot \mu(z'_1) \cdots \mu(z'_m),
 \end{aligned}$$

where $\mu(z_1), \dots, \mu(z_n), \mu(z'_1), \dots, \mu(z'_m)$ are standard sum-product messages.

Example



Gradient EM



$$\begin{aligned}
 \frac{dh(\theta_k)}{d\theta_k} &= \gamma^{-1} \sum_{\mathbf{z}} g(z_1, \dots, z_m, \hat{\theta}_k) \mu(z_1) \dots \mu(z_m) \frac{d \log g(z_1, \dots, z_m, \theta_k)}{d\theta_k} dz, \\
 &= \sum_{\mathbf{z}} p(z_1, \dots, z_m | \hat{\theta}_k) \frac{d \log g(z_1, \dots, z_m, \theta_k)}{d\theta_k}, \\
 &= \mathbb{E}_{p(z_1, \dots, z_m | \hat{\theta}_k)} \left[\frac{d \log g(z_1, \dots, z_m, \theta_k)}{d\theta_k} \right].
 \end{aligned}$$

Gradient EM: Properties

Theorem (Cycle-free $f_B(x, \theta)$)

Assume that a factor graph of a global function $f(x, \theta) \triangleq f_A(\theta)f_B(x, \theta)$ is available whose subgraph $f_B(x, \theta)$ is **cycle-free**. The fixed points of gradient EM applied on the graph of $f(x, \theta)$ are the **stationary points** of $f(\theta)$.

Theorem (Cyclic $f_B(x, \theta)$)

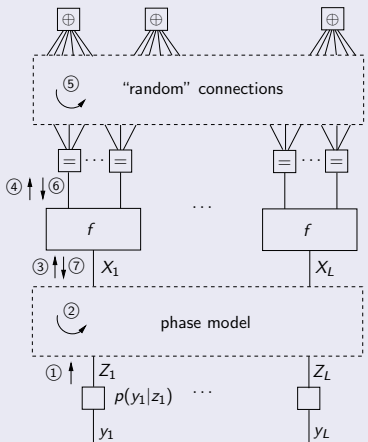
Assume that a factor graph of a global function $f(x, \theta) \triangleq f_A(\theta)f_B(x, \theta)$ is available (whose subgraph $f_B(x, \theta)$ may be **cycle-free or cyclic**). The fixed points of a gradient EM algorithm applied on that factor graph are the **stationary points** of the function $\hat{f}(\theta)$, defined as:

$$\log \hat{f}(\theta) \triangleq \log f_A(\theta) + \int_{-\infty}^{\theta} E_{b(x|\tilde{\theta})} \left[\nabla_{\theta} \log f_B(x, \tilde{\theta}) \right] d\tilde{\theta},$$

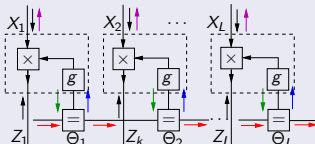
where the beliefs $b(\cdot|\theta)$ are computed by means of the sum-product messages available at convergence of the sum-product algorithm.

Scheduling

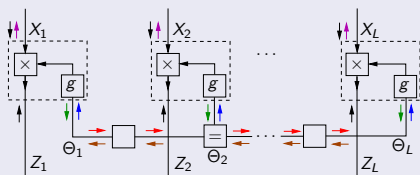
Factor graph of $p(x, y, \theta)$



Constant-phase model



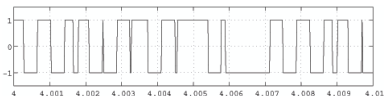
Random-walk phase model



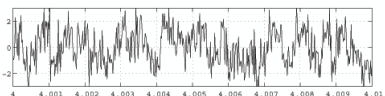
Problem statement

- Pseudo-noise signal X is transmitted over noisy channel, resulting in the noisy signal Y .
- The analog circuit estimates the signal X from the noisy signal Y .

Pseudo-noise signal X



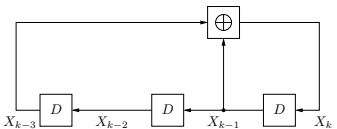
Received noisy signal Y



Applications

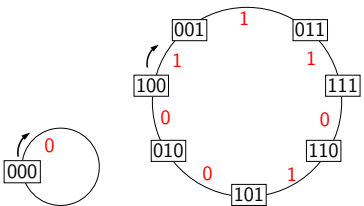
- Spread spectrum communication systems (CDMA, UWB)
- Positioning systems (GPS)

Pseudo-random sequence X generated by LFSR

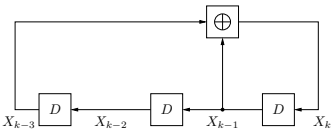


$X = [\dots, X_{k-1}, X_k, X_{k+1}, \dots]$ with $X_k = X_{k-1} \oplus X_{k-3}$

State diagram

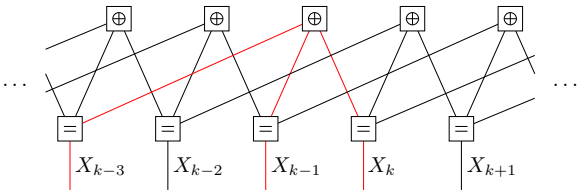


Pseudo-random sequence X generated by LFSR



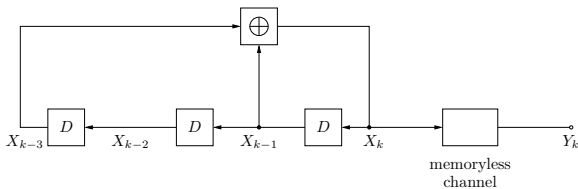
$$X = [\dots, X_{k-1}, X_k, X_{k+1}, \dots] \text{ with } X_k = X_{k-1} \oplus X_{k-3}$$

Representation as **factor graph**.



Synchronization task

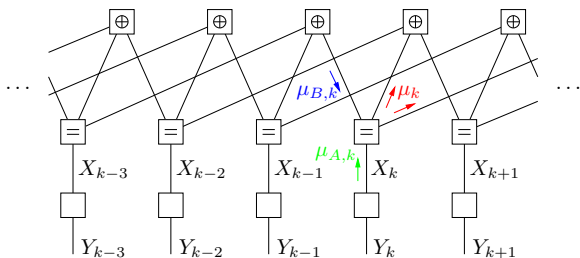
Based on the **noisy observation** Y of the sequence X , **estimate** the actual **state** of the source.



Approach:

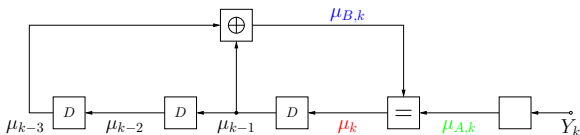
Use the **factor graph** to define a message-passing algorithm.

Forward-only message passing on the factor graph



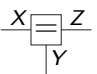
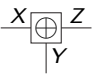
Interpretation:

Filtering of the sequence Y with a **soft** version of the LFSR.



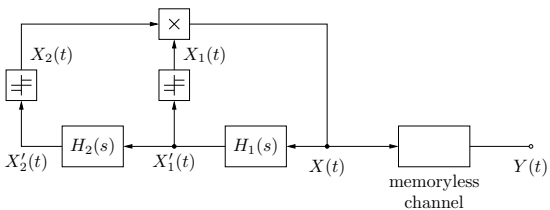
Reminder: SP for EQU and XOR-node

$$L \triangleq \log \frac{\mu(0)}{\mu(1)} \quad \Delta = \frac{\mu(0) - \mu(1)}{\mu(0) + \mu(1)}$$

 <p>$\delta[x - y]\delta[x - z]$</p>	$\begin{pmatrix} \mu_Z(0) \\ \mu_Z(1) \end{pmatrix} = \begin{pmatrix} \mu_X(0)\mu_Y(0) \\ \mu_X(0)\mu_X(1) \end{pmatrix}$ $L_Z = L_X + L_Y$ $\Delta_Z = \frac{\Delta_X + \Delta_Y}{1 + \Delta_X \Delta_Y}$
 <p>$\delta[x \oplus y \oplus z]$</p>	$\begin{pmatrix} \mu_Z(0) \\ \mu_Z(1) \end{pmatrix} = \begin{pmatrix} \mu_X(0)\mu_Y(0) + \mu_X(1)\mu_Y(1) \\ \mu_X(0)\mu_X(1) + \mu_X(1)\mu_X(0) \end{pmatrix}$ $\tanh(L_Z/2) = \tanh(L_X/2) \cdot \tanh(L_Y/2)$ $\Delta_Z = \Delta_X \Delta_Y$

Reminder: SP for EQU and XOR-node

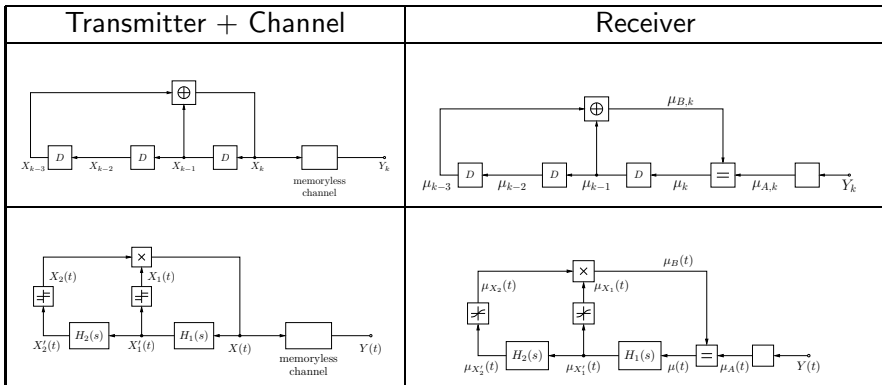
Signal Source



- **Delay elements** replaced by **linear filters**.
- Output of the filters $X_1'(t)$ and $X_2'(t) \in \mathbb{R}$.
- Introduction of **threshold functions** ($X_1(t), X_2(t)$ and $X(t) \in \{-1, +1\}$).
- **Multiplication** corresponds to **addition modulo 2**.

From Discrete-Time to Continuous-Time (3)

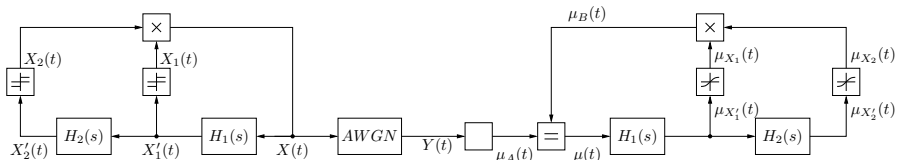
Overview



Demonstration System

Signals in the receiver are **pseudo probability functions** of the corresponding signals in the source

- **Discrete** variables represented as “**pseudo-means**” (Δ -representation)
e.g., $\tilde{E}[X(t)] = \mu_x(t) = \Pr[X(t) = +1] - \Pr[X(t) = -1]$
- **Continuous** variables ($X'_{1,2}$)
The pdf for $X'_{1,2}(t)$ is **assumed** to be **Gaussian** $\mathcal{N}(\mu_{X'_{1,2}}(t), \sigma_{1,2}^2)$.
 - Means $\mu_{X'_{1,2}}$ are **computed**
 - Variances $\sigma_{1,2}^2$ are fixed and set **manually**



Demonstration System (2)

Filters

- The **means** μ_X and $\mu_{X'_1}$ are filtered by $H_1(s)$ and $H_2(s)$.
Indeed, let $y(t) = [h \star x](t)$, then $E[y(t)] = [h \star E[x]](t)$.
- Remark: for computing the mean, the **variance is not needed!**

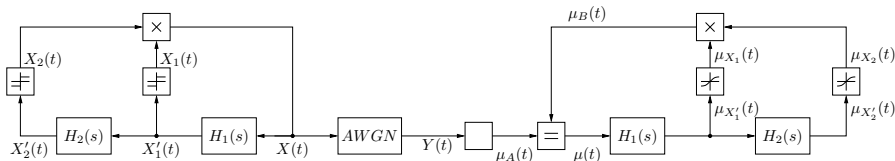
Soft-Thresholds:

$$\Pr[X_{1,2}(t) = +1] = \Pr[X'_{1,2}(t) \geq 0]$$

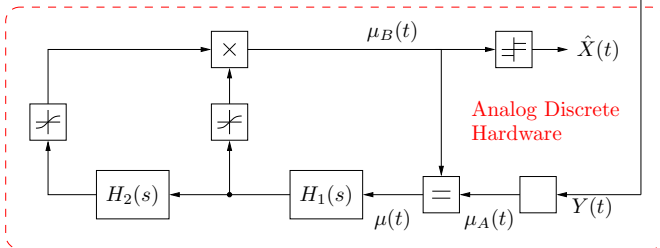
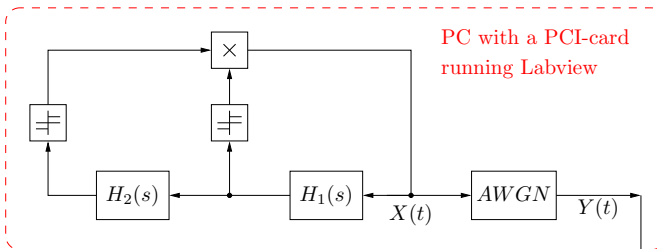
$$\Pr[X_{1,2}(t) = -1] = \Pr[X'_{1,2}(t) < 0]$$

$$\mu_{X_{1,2}}(t) = \Pr[X_{1,2}(t) = +1] - \Pr[X_{1,2}(t) = -1] = \operatorname{erf}\left(\frac{\mu_{X'_{1,2}}(t)}{\sqrt{2}\sigma_{1,2}}\right)$$

$$\mu_{X_{1,2}}(t) \approx \tanh(C \mu_{X'_{1,2}}(t))$$



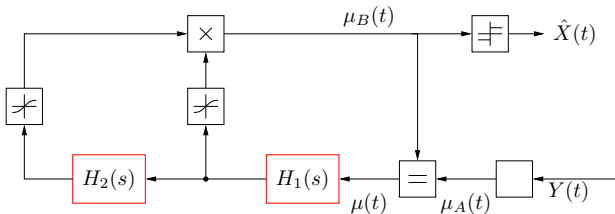
Demonstration System (3)



Demonstration System (4)

The Filters

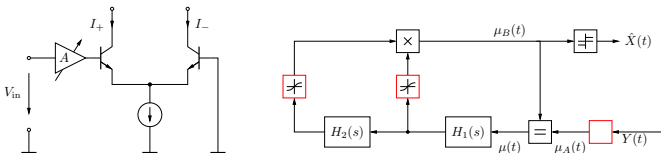
- $H_1(s)$: Butterworth Lowpass Filter, 5th order, $f_c = 1.6$ kHz
- $H_2(s)$: $4 \times H_1(s)$ in series (or $6 \times H_1(s)$ in series)



Demonstration System (5)

The Soft-Threshold Function, AWGN-Channel Estimation

Differential pair with the gain A as an adjustable parameter

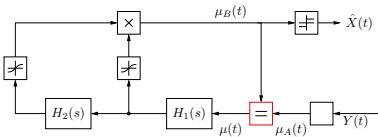
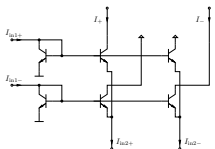


$$\log \frac{\mu_A(0)(t)}{\mu_A(1)(t)} = \log \frac{e^{-\frac{(Y(t)-1)^2}{2\sigma^2}}}{e^{-\frac{(Y(t)+1)^2}{2\sigma^2}}} = \frac{2Y(t)}{\sigma^2}$$

Demonstration System (6)

The Equality Constraint Gate

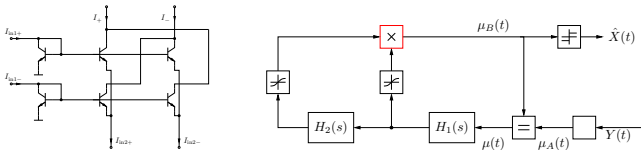
Forward-only EQU-Softgate



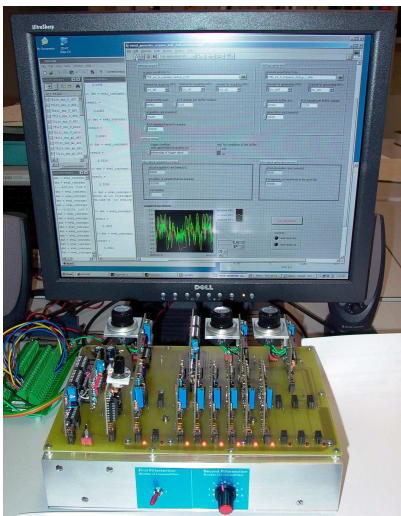
Demonstration System (7)

The Soft-XOR Gate

Forward-only XOR-Softgate corresponds to a “Gilbert multiplier”



Demonstration System (8)

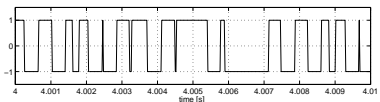
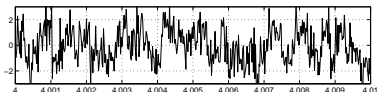
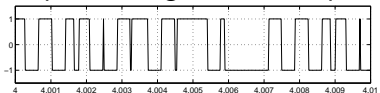


Example of a sequence at $\text{SNR} = 0$ dB

Sampling rate: 50 kHz, 500 samples shown

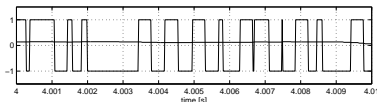
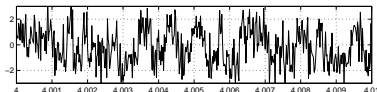
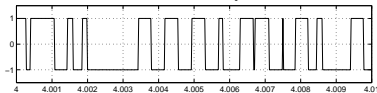
$$H_2(s): 4 \times H_1(s)$$

sequence-length: 361 samples



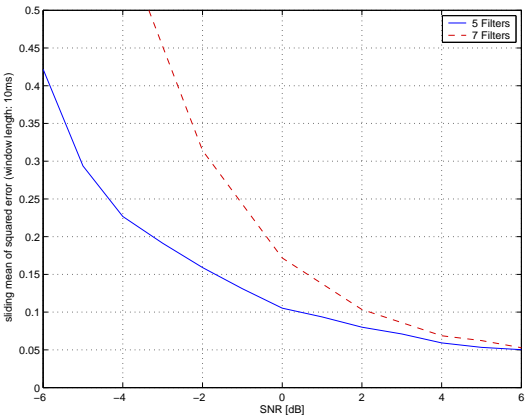
$$H_2(s): 6 \times H_1(s)$$

1'707 samples



Measurement results

MSE vs. SNR



Cramér-Rao-type bounds (2)

Example (AR model)

Let X_1, X_2, \dots be a real random process defined by:

$$X_k = a_1 X_{k-1} + a_2 X_{k-2} + \dots + a_M X_{k-M} + U_k, \quad U_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_U^2)$$

and let the process $Y = Y_1, Y_2, \dots$ be defined as:

$$Y_k = X_k + W_k, \quad U_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_W^2).$$

Task

Estimate $a_1, \dots, a_M, \sigma_U^2, \sigma_W^2$, and X from observation Y .

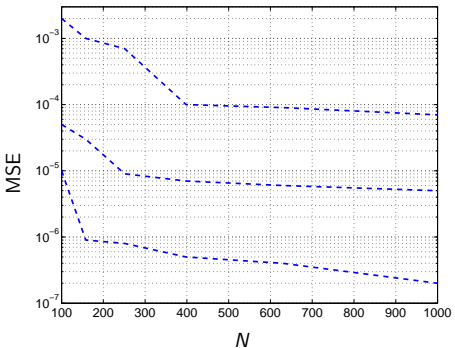
Observation

ML/MAP/MMSE-estimators are **infeasible!**

Neither can their MSE be determined \Rightarrow **lower bounds** on MMSE.

Cramér-Rao-type bounds (2)

Results for σ_W^2 ($\sigma_U^2 = 0.1$; $\sigma_W^2 = 0.001, 0.01, 0.1$)



Estimation algorithm by Sascha Korl.

Does the algorithm perform well?

Cramér-Rao-type bounds (3)

What is the CRB more precisely?

Inverse of **information matrix**.

Three different types

... as there are three different types of information matrices

- **Standard** Cramér-Rao bounds: parameters
- **Bayesian** Cramér-Rao bounds: random variables
- **Hybrid** Cramér-Rao bounds: parameters and random variables

Bayesian Cramér-Rao bound

Theorem (Bayesian Cramér-Rao bound)

Let $p(x, y)$ be the joint pdf of $x \in \mathbb{R}$ and $y \triangleq (y_1, \dots, y_N)$.
 If $p(x)$ is **zero** at boundary of its support, then for any **regular** $\hat{x}(y)$:

$$E_{XY} [(x - \hat{x}(y))^2] \geq J^{-1},$$

where the Bayesian information matrix J is defined as:

$$J \triangleq E_{XY} \left[\left(\frac{\partial}{\partial x} \log p(x, y) \right)^2 \right].$$

Properties

- **MAP-estimator** achieves bound as SNR or $N \rightarrow \infty$.
- BCRB holds for **any** regular $\hat{x}(y)$ as SNR or $N \rightarrow \infty$.

Bayesian Cramér-Rao bound: simple example

Example (Mean of a Gaussian random variable)

$Y = X + Z$ with $Z \sim \mathcal{N}(0, \sigma^2)$ with σ^2 **known** and $X \in \mathbb{R}$ **unknown**.

Estimate X from observations y_1, y_2, \dots, y_N with prior $p(X)$ for X .

$$p(x, y_1, y_2, \dots, y_N) = p(x) \prod_{k=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-y_k)^2/2\sigma^2}$$

$$\begin{aligned} \mathbf{J} &= -E_{XY} \left[\frac{d^2}{dx^2} \log p(X, Y) \right] \\ &= \frac{N}{\sigma^2} - E_X \left[\frac{d^2}{dx^2} \log p(X) \right] \end{aligned}$$

$$E_{XY}[(\hat{x}(X) - X)^2] \geq \mathbf{J}^{-1} = \left(\frac{N}{\sigma^2} - E_X \left[\frac{d^2}{dx^2} \log p(X) \right] \right)^{-1}$$

If $p(x)$ is Gaussian, then BCRB = **minimum** achievable MSE!

Vector case

Bayesian Cramér-Rao bound for component X_k

Given: joint pdf $p(x, y)$ of $x \triangleq (x_1, \dots, x_M)$ and $y \triangleq (y_1, \dots, y_N)$.

Lower bound for the MSE $E_{X_k Y} [(X_k - \hat{x}_k(Y))^2]$?

From marginal

$$E_{X_k Y} [(X_k - \hat{x}_k(Y))^2] \geq J_k^{-1},$$

with $J_k \triangleq E_{X_k Y} \left[\left(\frac{\partial}{\partial x_k} \log p(x_k, y) \right)^2 \right]$.

From joint pdf

$$E_{X_k Y} [(X_k - \hat{x}_k(Y))^2] \geq [J^{-1}]_{kk},$$

with $J_{ij} \triangleq E_{X Y} \left[\left(\frac{\partial}{\partial x_i} \log p(x, y) \right)^T \frac{\partial}{\partial x_j} \log p(x, y) \right]$.

BCRB from marginal is **tighter** than from joint pdf, but more difficult to compute.

Algorithms

From joint pdf

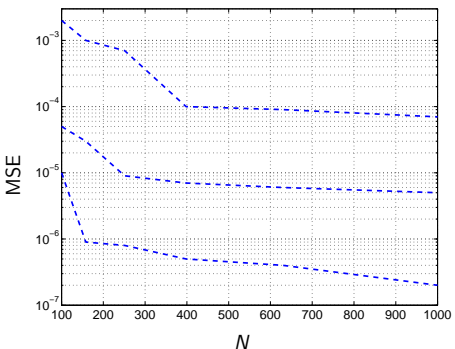
- J is often **sparse**.
- Only **diagonal** elements of inverse required.
- **Local** computations of “small” matrices (message passing).

From marginal

- J_k is usually **dense**.
- Key to J_k : $\frac{\partial}{\partial x_k} \log p(x_k, y) = E_{X_{\sim k} | x_k, Y} \left[\frac{\partial}{\partial x_k} \log p(X, Y) \right]$
- **Expectation** computed by belief propagation (or “sum-product algorithm” or “probability propagation”).

Results

Results for σ_W^2 ($\sigma_U^2 = 0.1$; $\sigma_W^2 = 0.001, 0.01, 0.1$)

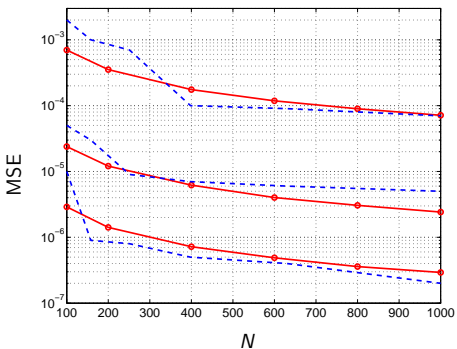


Estimation algorithm by Sascha Korl.

Does the algorithm perform well?

Results

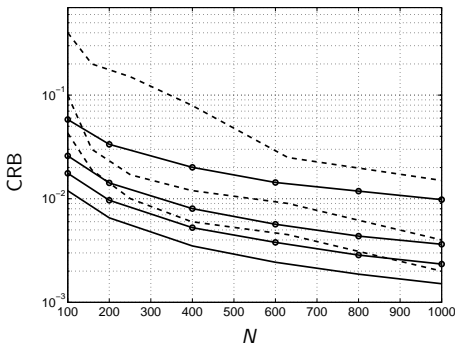
Results for σ_W^2 ($\sigma_U^2 = 0.1$; $\sigma_W^2 = 0.001, 0.01, 0.1$)



Estimation algorithm by Sascha Korl.

Results

Results for a ($\sigma_U^2 = 0.1$);



Standard CRB with **unknown** σ_U^2 and σ_W^2 (solid) for $\sigma_W^2 = 0.1/0.01/0.001$;
 MSE of algorithm by S. Korl (dashed);
 Standard CRB for a with **known** $\sigma_W^2 = 0$.

A closer look

Example (AR model)

Let X_1, X_2, \dots be a real random process defined by:

$$X_k = a_1 X_{k-1} + a_2 X_{k-2} + \dots + a_M X_{k-M} + U_k, \quad U_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_U^2)$$

and let the process $Y = Y_1, Y_2, \dots$ be defined as:

$$Y_k = X_k + W_k, \quad U_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_W^2).$$

Reminder

$$F(\theta) \triangleq \mathbf{E}_{Y; \theta} \left[\nabla_{\theta} \log p(Y; \theta) \nabla_{\theta}^T \log p(Y; \theta) \right] \text{ with } p(y; \theta) \triangleq \int_x p(x, y; \theta)$$

$$\theta = (\mathbf{a}, \sigma_U^2, \sigma_W^2)$$

$$\nabla_{\theta} \log p(y; \theta) = \mathbf{E}_{X|y} [\nabla_{\theta} \log p(X, y; \theta)].$$

Expectations $\nabla_{\theta} \log p(y; \theta) = E_{X|\theta y} [\nabla_{\theta} \log p(X, y; \theta)]$

$$p(x, y | \mathbf{a}, \sigma_W^2, \sigma_U^2) = \prod_k \underbrace{\mathcal{N}\left(x_k - \sum_{n=1}^M a_n x_{k-n} | 0, \sigma_U^2\right)}_{f_1(x_k, \dots, x_{k-M}, \mathbf{a}, \sigma_U^2)} \underbrace{\mathcal{N}(y_k - x_k | 0, \sigma_W^2)}_{f_2(x_k, \sigma_W^2, y_k)}.$$

As a consequence:

$$\nabla_{\theta} \log p(x, y | \mathbf{a}, \sigma_W^2, \sigma_U^2) = \sum_k \nabla_{\theta} \log f_1(x_k, \dots, x_{k-M}, \mathbf{a}, \sigma_U^2) + \sum_k \nabla_{\theta} \log f_2(x_k, \sigma_W^2, y_k).$$

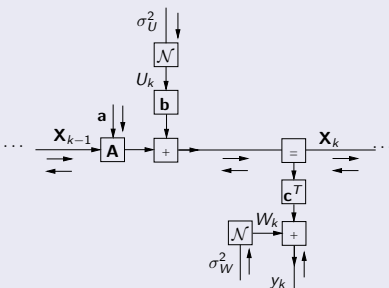
Expectations $\nabla_{\theta} \log p(y; \theta) = \mathbb{E}_{X|y} [\nabla_{\theta} \log p(X, y; \theta)]$ (2)

$$\begin{aligned} \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [\nabla_{a_i} \log f_1(X_k, \dots, X_{k-M}, \mathbf{a}, \sigma_U^2)] \\ = \frac{1}{\sigma_U^2} (\mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [X_{k-i} X_k] - \sum_{\ell=1}^M a_{\ell} \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [X_{k-i} X_{k-\ell}]) \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [\nabla_{\sigma_U^2} \log f_1(X_k, \dots, X_{k-M}, \mathbf{a}, \sigma_U^2)] \\ = -\frac{1}{2\sigma_U^2} + \frac{1}{2\sigma_U^4} \left(\mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [X_k^2] - 2 \sum_{\ell=1}^M a_{\ell} \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [X_k X_{k-\ell}] \right. \\ \left. + \sum_{\ell=1}^M \sum_{m=1}^M a_{\ell} a_m \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [X_{k-\ell} X_{k-m}] \right) \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [\nabla_{\sigma_W^2} \log f_2(X_k, \sigma_W^2, y_k)] \\ = -\frac{1}{2\sigma_W^2} + \frac{1}{2\sigma_W^4} \left(y_k^2 - 2y_k \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [X_k] + \mathbb{E}_{X|a, \sigma_W^2, \sigma_U^2, Y} [X_k^2] \right). \end{aligned}$$

$$\text{Expectations } \nabla_{\theta} \log p(y; \theta) = E_{X|\theta y} [\nabla_{\theta} \log p(X, y; \theta)] \quad (3)$$



$E_{X|a \sigma_W^2 \sigma_U^2 Y}$ computed by forward/backward **Kalman recursions**.

(= instance of belief propagation)

Computing the Fisher information matrix of AR model

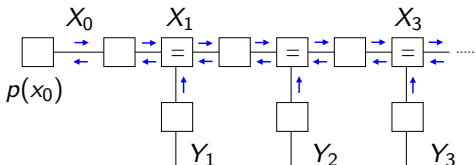
- 1 Generate a list of samples $\{\hat{y}^{(j)}\}_{j=1}^N$ from $p(y|\theta)$.
- 2 For $j = 1, \dots, N$:
 - Forward and backward Kalman recursion with $y = \hat{y}^{(j)}$.
 - Evaluate the expression:

$$E_{X|\theta, \hat{y}^{(j)}} \left[\nabla_{\theta} \log p(X, \hat{y}^{(j)}; \theta) \right].$$

- 3 Compute the estimate $\hat{\mathbf{F}}(\theta)$ for $\mathbf{F}(\theta)$:

$$\hat{\mathbf{F}}(\theta) \triangleq \frac{1}{N} \sum_{j=1}^N \left[E_{X|\theta, \hat{y}^{(j)}} \left[\nabla_{\theta} \log p(X, \hat{y}^{(j)}; \theta) \right] \right. \\ \left. E_{X|\theta, \hat{y}^{(j)}}^T \left[\nabla_{\theta} \log p(X, \hat{y}^{(j)}; \theta) \right] \right].$$

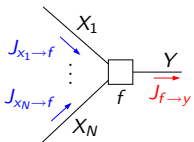
BCRB from information matrix of joint pdf



- The BCRBs of estimation in cycle-free graphical models can be computed efficiently by **message passing**.
- Messages are **matrices**.
- Messages are updated at each node according to specific **update rules**.
- The BCRBs are computed by **combining** those messages.

BCRB from information matrix of joint pdf (2)

Differentiable node function



$$J_{f \rightarrow y}^{-1}(Y) = \left(\begin{bmatrix} J_{X_1 \rightarrow f}(X_1) + E[-\Delta_{X_1}^{X_1} \log f] & \dots & E[-\Delta_{X_1}^{X_N} \log f] & E[-\Delta_{X_1}^Y \log f] \\ \vdots & \dots & \dots & \vdots \\ E[-\Delta_{X_N}^{X_1} \log f] & \dots & J_{X_N \rightarrow f}(X_N) + E[-\Delta_{X_N}^{X_N} \log f] & E[-\Delta_{X_N}^Y \log f] \\ E[-\Delta_{X_N}^{X_1} \log f] & \dots & E[-\Delta_{X_N}^Y \log f] & E[-\Delta_{X_N}^Y \log f] \end{bmatrix}^{-1} \right)_{N+1, N+1}$$

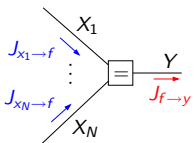
with $\Delta_{X_i}^{X_j} \triangleq \nabla_{X_i} \nabla_{X_j}^T$

Remarks

- Expectations $E[\Delta_{X_i}^{X_j} \log f]$ supposed to be [well-defined](#).
- They can [easily](#) be computed [numerically](#).
- Rows and corresponding columns can be [exchanged](#).

BCRB from information matrix of joint pdf (3)

Equality constraint node



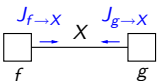
$$J_{f \to y} = \sum_{i=1}^N J_{X_i \to f}$$

Terminal node



$$J_{f \to X} = -E[\Delta_X^f \log f]$$

PCRB



$$J_{\text{tot}} = J_{f \to X} + J_{g \to X}$$

Natural gradient-based estimation

Reminder

Information matrix computed by belief propagation.
(or “sum-product algorithm” or “probability propagation”).

Approximations

Approximate information matrices by **iterative** message passing.
Approximate but **efficient** inversion by imposing **structure**.

Applications

Novel algorithms for estimation (e.g., AR model).

Framework

Conveniently derived as **message passing** on factor graphs.

Kernels from probability measures (3)

Product kernel [Jebara et al., 2004]

The product-kernel is computed as follows:

$$\kappa(\hat{y}_i, \hat{y}_j) \triangleq \sum_y p(y|\hat{y}_i)p(y|\hat{y}_j),$$

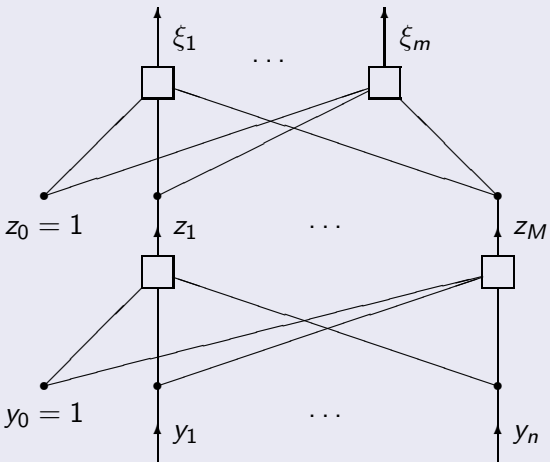
with

$$p(y|\hat{y}_i) \triangleq \sum_x p(y|x, \hat{\theta})p(x|\hat{\theta}, \hat{y}_i),$$

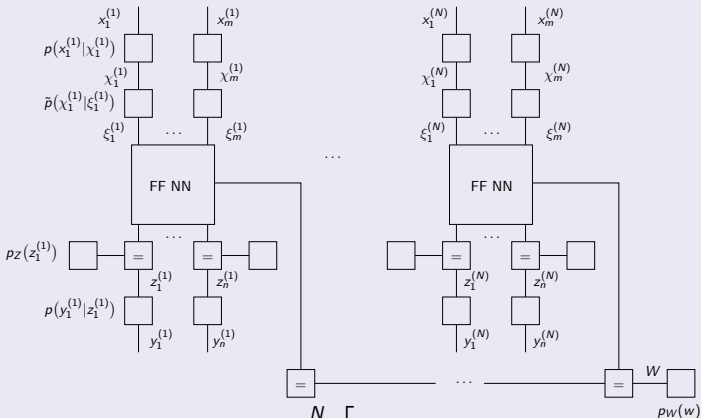
where the parameters $\hat{\theta}$ is estimated by means of the sample \hat{y}_i , e.g., by ML estimation:

$$\begin{aligned}\hat{\theta}^{\text{ML}} &\triangleq \underset{\theta}{\operatorname{argmax}} p(\hat{y}_i|\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_x p(\hat{y}_i, x|\theta).\end{aligned}$$

Feed-forward neural network

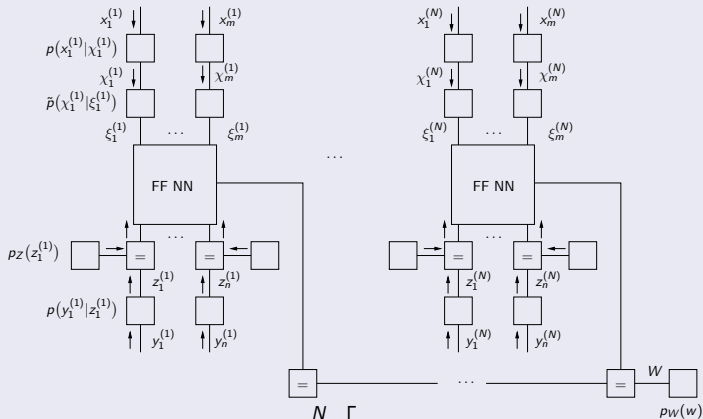


Feed-forward neural network: additional nodes



$$p(x, \xi, z, y, w) \triangleq p_Z(z) p_W(w) \prod_{i=1}^N \left[\delta(\xi^{(i)} - \xi(z^{(i)}, w)) \right. \\ \left. \cdot \left(\prod_{j=1}^m \tilde{p}(x_j^{(i)} | \chi_j^{(i)}) p(\chi_j^{(i)} | \xi_j^{(i)}) \right) \left(\prod_{j=1}^n p(y_j^{(i)} | z_j^{(i)}) \right) \right].$$

Feed-forward neural network: pre-processing



$$p(x, \xi, z, y, w) \triangleq p_Z(z) p_W(w) \prod_{i=1}^N \left[\delta(\xi^{(i)} - \xi(z^{(i)}, w)) \cdot \left(\prod_{j=1}^m \tilde{p}(x_j^{(i)} | \chi_j^{(i)}) p(\chi_j^{(i)} | \xi_j^{(i)}) \right) \left(\prod_{j=1}^n p(y_j^{(i)} | z_j^{(i)}) \right) \right].$$

Information rate: introduction

Objective

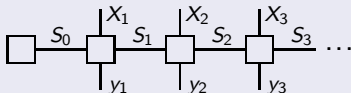
Information rate $I(X; Y) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} I(X_1, \dots, X_n; Y_1, \dots, Y_n)$
 between input process $X = (X_1, X_2, \dots)$ and output process
 $Y = (Y_1, Y_2, \dots)$ of **time-invariant** discrete-time channel with
memory.

State-space representation

An ergodic stochastic process $S = (S_0, S_1, S_2, \dots)$ such that

$$p(x^n, y^n, s_0^n) = p(s_0) \prod_{k=1}^n p(x_k, y_k, s_k | s_{k-1})$$

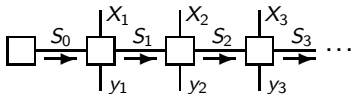
for all $n > 0$ and with $p(x_k, y_k, s_k | s_{k-1})$ not depending on k .



Basic principle

Algorithm

- 1 Sample two "very long" sequences x^n and y^n .
- 2 Compute $\log p(x^n)$, $\log p(y^n)$, and $\log p(x^n, y^n)$.
- 3 $\hat{I}(X; Y) \triangleq \frac{1}{n} \log p(x^n, y^n) - \frac{1}{n} \log p(x^n) - \frac{1}{n} \log p(y^n)$.

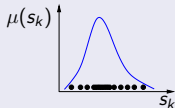


Discrete input space \mathcal{X} and state-space \mathcal{S} [e.g., Arnold et al.]

Forward sum-product sweep = forward BCJR-recursion

Continuous input space \mathcal{X} and state-space \mathcal{S}

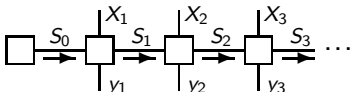
Forward sum-product sweep by **particle filtering**.
 Expression $p(x_k, s_k | s_{k-1})$ not required!
 E.g., stochastic differential/difference equation.



Forward sum-product sweep

Computation of

$$p(y^n) = \int_{x^n} \int_{s_0^n} p(x^n, y^n, s_0^n).$$



$$\mu_k(s_k) = \int_{x_k} \int_{s_{k-1}} \mu_{k-1}(s_{k-1}) p(x_k, y_k, s_k | s_{k-1}) dx_k ds_{k-1}$$

$$= \int_{x^k} \int_{s_0^{k-1}} p(x^k, y^k, s^k) dx^k ds_0^{k-1}$$

$$p(y^n) = \int_{s_n} \mu_n(s_n),$$

Results: Gaussian channel with $0 \leq X \leq 1$

Model for **free-space optical** communications channel

- Transmitter: light emitting diode (LED) or laser diode (LD)
- Signal modulated on optical intensity (ON/OFF keying)
- Direct line-of-sight path is dominant
- Noise source = ambient light
- Peak power constraint due to eye safety and potential thermal skin damage

