

## Welcome to Sequential Monte Carlo methods!!

### Lecture 1 – Introduction and probabilistic modelling

Johan Alenlöv, Linköping University 2025-02-03

**Aim:** To provide an introduction to the theory and application of sequential Monte Carlo (SMC) methods.

**Aim:** To provide an introduction to the theory and application of sequential Monte Carlo (SMC) methods.

After the course you should be able to derive your own SMC-based algorithms allowing you to do inference in nonlinear models.

**Aim:** To provide an introduction to the theory and application of sequential Monte Carlo (SMC) methods.

After the course you should be able to derive your own SMC-based algorithms allowing you to do inference in nonlinear models.

Day 1-3: Focus on state space models (SSMs). How to learning them from data and how to estimate their hidden states.

Day 4-5: Using SMC for inference in general probabilistic models.

**Aim:** Compute the **position** of a person moving around indoors using variations in the ambient magnetic field and the motion of the person (acceleration and angular velocities). All of this observed using sensors in a standard smartphone.

**Aim:** Compute the **position** of a person moving around indoors using variations in the ambient magnetic field and the motion of the person (acceleration and angular velocities). All of this observed using sensors in a standard smartphone.



Key ingredients of the solution:

- 1. The particle filter for computing the position
- 2. The Gaussian process for building and representing the map of the ambient magnetic field
- 3. Inertial sensor signal processing

Key ingredients of the solution:

- 1. The particle filter for computing the position
- 2. The Gaussian process for building and representing the map of the ambient magnetic field
- 3. Inertial sensor signal processing

Movie - map making: www.youtube.com/watch?v=enlMiUqPVJo

Movie – indoor positioning result

Key ingredients of the solution:

- 1. The particle filter for computing the position
- 2. The Gaussian process for building and representing the map of the ambient magnetic field
- 3. Inertial sensor signal processing

Movie - map making: www.youtube.com/watch?v=enlMiUqPVJo

Movie – indoor positioning result

Arno Solin, Simo Särkkä, Juho Kannala and Esa Rahtu. Terrain navigation in the magnetic landscape: Particle filtering for indoor positioning. Proc. of the European Navigation Conf. (ENC), Helsinki, Finland, June, 2016.



Arno Solin, Manon Kok, Niklas Wahlström, Thomas B. Schön and Simo Särkkä. Modeling and interpolation of the ambient magnetic field by Gaussian processes. *IEEE Trans. on Robotics*, 34(4):1112-1127, 2018.



Carl Jidling, Niklas Wahlström, Adrian Wills and Thomas B. Schön. Linearly constrained Gaussian processes. Advances in Neural Information Processing Systems (NIPS), Long Beach, CA, USA, December, 2017.

**Aim:** To learn a model explaining the seasonal influenza epidemics and then make use of this model to compute predictions.

**Aim:** To learn a model explaining the seasonal influenza epidemics and then make use of this model to compute predictions.

Susceptible/infected/recovered (SIR) model:

$$S_{t+dt} = S_t + \mu \mathcal{P}dt - \mu S_t dt - (1 + Fv_t)\beta_t S_t \mathcal{P}^{-1} I_t dt,$$
  

$$I_{t+dt} = I_t - (\gamma + \mu)I_t dt + (1 + Fv_t)\beta_t S_t \mathcal{P}^{-1} I_t dt,$$
  

$$R_{t+dt} = R_t + \gamma I_t dt - \mu R_t dt,$$
  

$$\beta_t = R_0 (\gamma + \mu)(1 + \alpha \sin(2\pi t/12)),$$

Measurements:

$$y_k = \rho \operatorname{logit}(\overline{l}_k/\mathcal{P}) + e_k, \qquad e_k \sim \mathcal{N}(0, \sigma^2).$$

**Aim:** To learn a model explaining the seasonal influenza epidemics and then make use of this model to compute predictions.

Susceptible/infected/recovered (SIR) model:

$$S_{t+dt} = S_t + \mu \mathcal{P}dt - \mu S_t dt - (1 + Fv_t)\beta_t S_t \mathcal{P}^{-1}I_t dt,$$
  

$$I_{t+dt} = I_t - (\gamma + \mu)I_t dt + (1 + Fv_t)\beta_t S_t \mathcal{P}^{-1}I_t dt,$$
  

$$R_{t+dt} = R_t + \gamma I_t dt - \mu R_t dt,$$
  

$$\beta_t = R_0(\gamma + \mu)(1 + \alpha \sin(2\pi t/12)),$$

Measurements:

$$y_k = \rho \text{logit}(\overline{l}_k/\mathcal{P}) + e_k, \qquad e_k \sim \mathcal{N}(0, \sigma^2).$$

Information about the unknown parameters  $\theta = (\gamma, R_0, \alpha, F, \rho, \sigma)$  and states  $x_t = (S_t, I_t, R_t)$  has to be learned from measurements.

Compute  $p(\theta, x_{1:T} | y_{1:T})$ , where  $y_{1:T} = (y_1, y_2, \dots, y_T)$  and use it to compute the predictive distribution.



Disease activity (number of infected individuals  $I_t$ ) over an eight year period.

Fredrik Lindsten, Michael I. Jordan and Thomas B. Schön. Particle Gibbs with ancestor sampling. Journal of Machine Learning Research (JMLR), 15:2145-2184, June 2014.

A **probabilistic program** encodes a **probabilistic model** according to the semantics of a particular probabilistic programming language, giving rise to a **programmatic model**.

A **probabilistic program** encodes a **probabilistic model** according to the semantics of a particular probabilistic programming language, giving rise to a **programmatic model**.

The memory state of a running probabilistic program evolves dynamically and stochastically in time and so is a **stochastic process**.

A **probabilistic program** encodes a **probabilistic model** according to the semantics of a particular probabilistic programming language, giving rise to a **programmatic model**.

The memory state of a running probabilistic program evolves dynamically and stochastically in time and so is a **stochastic process**.

**SMC** is a common inference method for programmatic model.

Creates a clear separation between the model and the inference methods. Opens up for the automation of inference! x ~ Bernoulli(p); assume(x)
if (x) { value(x)
y ~ N(0,1); assume(y)
} else {
y <- 0;</pre>

A **probabilistic program** encodes a **probabilistic model** according to the semantics of a particular probabilistic programming language, giving rise to a **programmatic model**.

The memory state of a running probabilistic program evolves dynamically and stochastically in time and so is a **stochastic process**.

**SMC** is a common inference method for programmatic model.

Creates a clear separation between the model and the inference methods. Opens up for the automation of inference! x - Bernoulli(p); assume(x)
if (x) { value(x)
 y - N(0,1); assume(y)
} else {
 y <- 0;</pre>

More during lecture 17.

### Course structure – overview

- 17 lectures (45 min. each)
- Credits offered: 6ECTS (Swedish system)
- Practicals (solve exercises and hand-in assignments, discuss and ask questions)
- Discussions (discuss concepts)
- Hand-in assignments. You can collaborate, but the reports with the solutions are individual.
  - One set to be done between day 2 and 3
  - $\cdot$  One set to be done after day 5
- Complete course information (including lecture slides) is available from the course website: https:

//www.ida.liu.se/divisions/stima/fokurser/smc2025/

Feel free to ask questions at any time!

The only way to really learn something is by implementing it on your own.

### Outline – 5 days

Day 1 – Probabilistic modelling and particle filtering basics

- a) Probabilistic modelling of dynamical systems and filtering
- b) Introduce Monte Carlo and derive the bootstrap particle filter
- Day 2 Particle filtering
  - a) Auxiliary particle filter, full adaptation and practicalities
- Day 3 Parameter learning
  - a) Maximum likelihood parameter learning, convergence
- Day 4 Bayesian parameter learning
  - a) Particle Metropolis Hastings
  - b) Particle Gibbs
- Day 5 Beyond state space models (outlooks)
  - a) General target sequences and SMC samplers
  - b) Diffusion models

**Aim:** Introduce the course and provide background on probabilistic modelling.

### Outline:

- 1. Course introduction and practicalities
- 2. Probabilistic modelling
- 3. Key probabilistic objects
- 4. Ex. probabilistic autoregressive modelling (if there is time)

# Probabilistic modelling

### Modelling

Mathematical model: A compact representation—set of assumptions—of the data that in precise mathematical form captures the key properties of the underlying situation.

### Modelling

Mathematical model: A compact representation—set of assumptions—of the data that in precise mathematical form captures the key properties of the underlying situation.

Most of the course (day 1-3) is concerned with dynamical phenomena. The methods are more general than that and during the last day we will broaden the scope significantly.

Dynamical phenomena produce temporal measurements (data) arriving as a **sequence** 

$$y_{1:t} = (y_1, y_2, \ldots, y_t).$$

### Modelling

Mathematical model: A compact representation—set of assumptions—of the data that in precise mathematical form captures the key properties of the underlying situation.

Most of the course (day 1-3) is concerned with dynamical phenomena. The methods are more general than that and during the last day we will broaden the scope significantly.

Dynamical phenomena produce temporal measurements (data) arriving as a **sequence** 

$$y_{1:t} = (y_1, y_2, \ldots, y_t).$$

#### Nice introduction to probabilistic modelling in Machine Learning

Ghahramani, Z. Probabilistic machine learning and artificial intelligence. Nature 521:452-459, 2015.

# It is important to maintain a solid representation of uncertainty in all mathematical objects and throughout all calculations.



### The two basic rules from probability theory

Let X and Y be continuous random variables<sup>1</sup>. Let  $p(\cdot)$  denote a general probability density function.

- 1. Marginalization (integrate out a variable):  $p(x) = \int p(x, y) dy$ .
- 2. Conditional probability: p(x, y) = p(x | y)p(y).

<sup>&</sup>lt;sup>1</sup>Notation: Upper-case letters for random variables (r.v.) X when we talk about models. Lower-case letters for realizations of the r.v., X = x and in algorithms. We will not use bold face for vectors.

### The two basic rules from probability theory

Let X and Y be continuous random variables<sup>1</sup>. Let  $p(\cdot)$  denote a general probability density function.

- 1. Marginalization (integrate out a variable):  $p(x) = \int p(x, y) dy$ .
- 2. Conditional probability: p(x, y) = p(x | y)p(y).

Combine them into Bayes' rule:

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{p(\mathbf{y})} = \frac{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{\int p(\mathbf{y} | \mathbf{x})p(\mathbf{x})d\mathbf{x}}$$

<sup>&</sup>lt;sup>1</sup>Notation: Upper-case letters for random variables (r.v.) X when we talk about models. Lower-case letters for realizations of the r.v., X = x and in algorithms. We will not use bold face for vectors.

Measurements  $y_{1:T} = (y_1, y_2, ..., y_T)$ : The measured data somehow obtained from the phenomenon we are interested in.

Measurements  $y_{1:T} = (y_1, y_2, ..., y_T)$ : The measured data somehow obtained from the phenomenon we are interested in.

**Unknown (static) model parameters**  $\theta$ : Describes the model, but unknown (or not known well enough) to the user.

**Measurements**  $y_{1:T} = (y_1, y_2, \dots, y_T)$ : The measured data somehow obtained from the phenomenon we are interested in.

**Unknown (static) model parameters**  $\theta$ : Describes the model, but unknown (or not known well enough) to the user.

Unknown model variables  $x_t$  (changing over time): Describes the state of the phenomenon at time t (in the indoor positioning example above  $x_t$  includes the unknown position).

**Measurements**  $y_{1:T} = (y_1, y_2, ..., y_T)$ : The measured data somehow obtained from the phenomenon we are interested in.

**Unknown (static) model parameters**  $\theta$ : Describes the model, but unknown (or not known well enough) to the user.

Unknown model variables  $x_t$  (changing over time): Describes the state of the phenomenon at time t (in the indoor positioning example above  $x_t$  includes the unknown position).

**Explanatory variables** *u*: Known variables that we do not bother to model as stochastic.

**Measurements**  $y_{1:T} = (y_1, y_2, ..., y_T)$ : The measured data somehow obtained from the phenomenon we are interested in.

**Unknown (static) model parameters**  $\theta$ : Describes the model, but unknown (or not known well enough) to the user.

Unknown model variables  $x_t$  (changing over time): Describes the state of the phenomenon at time t (in the indoor positioning example above  $x_t$  includes the unknown position).

**Explanatory variables** *u*: Known variables that we do not bother to model as stochastic.

A key task is often to learn  $\theta$  and/or  $x_t$  based on the available measurements  $y_{1:T}$ .

Full probabilistic model: The joint model of the data and parameters.

 $p(\theta, y) = p(y | \theta)p(\theta)$ 

**Full probabilistic model:** The joint model of the data and parameters.  $p(\theta, y) = p(y | \theta)p(\theta)$ 

**Likelihood function:**  $\mathcal{L}(\theta; y) = p(Y = y | \theta)$ 

**Full probabilistic model:** The joint model of the data and parameters.  $p(\theta, y) = p(y | \theta)p(\theta)$ 

**Likelihood function:**  $\mathcal{L}(\theta; y) = p(Y = y | \theta)$ 

**Posterior distribution:** Condition on *y* instead of  $\theta$ .

$$p(\theta | y) = \frac{p(\theta, y)}{p(y)} = \frac{p(y | \theta)}{\int p(y | \theta) p(\theta) d\theta}$$

**Full probabilistic model:** The joint model of the data and parameters.  $p(\theta, y) = p(y | \theta)p(\theta)$ 

**Likelihood function:**  $\mathcal{L}(\theta; y) = p(Y = y | \theta)$ 

**Posterior distribution:** Condition on y instead of  $\theta$ .

$$p(\theta | y) = \frac{p(\theta, y)}{p(y)} = \frac{p(y | \theta)}{\int p(y | \theta) p(\theta) d\theta}$$

**Prediction:** A statement about a future event  $\bar{y}$  that has not been observed.

14/22

### Computational problems

The problem of learning a model based on data leads to computational challenges, both

 Integration: e.g. the high-dimensional integrals arising during marg. (averaging over all possible parameter values θ):

$$p(y_{1:T}) = \int p(y_{1:T} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{\theta}.$$

• **Optimization:** e.g. when extracting point estimates, for example by maximizing the posterior or the likelihood

$$\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{arg\,max}} p(\boldsymbol{y}_{1:T} \mid \boldsymbol{\theta})$$

### Computational problems

The problem of learning a model based on data leads to computational challenges, both

• Integration: e.g. the high-dimensional integrals arising during marg. (averaging over all possible parameter values  $\theta$ ):

$$p(y_{1:T}) = \int p(y_{1:T} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{\theta}.$$

• **Optimization:** e.g. when extracting point estimates, for example by maximizing the posterior or the likelihood

$$\widehat{\theta} = \operatorname*{arg\,max}_{\theta} p(y_{1:T} \mid \theta)$$

Typically impossible to compute exactly, use approximate methods

- Monte Carlo (MC), Markov chain MC (MCMC), and sequential MC (SMC).
- Variational inference (VI).

An autoregressive model of order *n* is given by

 $Y_t = A_1 Y_{t-1} + A_2 Y_{t-2} + \dots + A_n Y_{t-n} + E_t, \quad E_t \sim \mathcal{N}(\mu, \tau^{-1})$ 

where  $\mu$  and  $\tau$  are known explanatory variables ( $\mu = 0, \tau \neq 0$ ).

An autoregressive model of order *n* is given by

$$Y_t = A_1 Y_{t-1} + A_2 Y_{t-2} + \dots + A_n Y_{t-n} + E_t, \quad E_t \sim \mathcal{N}(\mu, \tau^{-1})$$

where  $\mu$  and  $\tau$  are known explanatory variables ( $\mu = 0, \tau \neq 0$ ). The unknown model variables are collected as

 $\boldsymbol{\theta} = (A_1, A_2, \ldots, A_n)^{\mathsf{T}}$ 

with the prior

 $\theta \sim \mathcal{N}(0, \rho^{-1} I_n),$  where  $\rho$  assumed to be known.

An autoregressive model of order *n* is given by

$$Y_t = A_1 Y_{t-1} + A_2 Y_{t-2} + \dots + A_n Y_{t-n} + E_t, \quad E_t \sim \mathcal{N}(\mu, \tau^{-1})$$

where  $\mu$  and  $\tau$  are known explanatory variables ( $\mu = 0, \tau \neq 0$ ). The unknown model variables are collected as

 $\boldsymbol{\theta} = (A_1, A_2, \ldots, A_n)^{\mathrm{T}}$ 

with the prior

 $\theta \sim \mathcal{N}(0, \rho^{-1} I_n),$  where  $\rho$  assumed to be known.

**Task:** Compute the posterior  $p(\theta | y_{1:T})$ .

Full probabilistic model  $p(\theta, y_{1:T}) = p(y_{1:T} | \theta)p(\theta)$ , where the data distribution is given by

$$p(y_{1:T} | \boldsymbol{\theta}) = p(y_T | y_{1:T-1}, \boldsymbol{\theta}) p(y_{1:T-1} | \boldsymbol{\theta}) = \cdots = \prod_{t=1}^{l} p(y_t | y_{1:t-1}, \boldsymbol{\theta}).$$

Full probabilistic model  $p(\theta, y_{1:T}) = p(y_{1:T} | \theta)p(\theta)$ , where the data distribution is given by

$$p(y_{1:T} | \boldsymbol{\theta}) = p(y_T | y_{1:T-1}, \boldsymbol{\theta}) p(y_{1:T-1} | \boldsymbol{\theta}) = \cdots = \prod_{t=1}^{l} p(y_t | y_{1:t-1}, \boldsymbol{\theta}).$$

From the model we have that

$$p(y_t | y_{1:t-1}, \boldsymbol{\theta}) = \mathcal{N}(y_t | \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{z}_t, \tau^{-1}),$$
  
where  $\boldsymbol{Z}_t = (\boldsymbol{Y}_{t-1}, \boldsymbol{Y}_{t-2}, \dots, \boldsymbol{Y}_{t-n})^{\mathsf{T}}.$ 

Full probabilistic model  $p(\theta, y_{1:T}) = p(y_{1:T} | \theta)p(\theta)$ , where the data distribution is given by

$$p(y_{1:T} | \boldsymbol{\theta}) = p(y_T | y_{1:T-1}, \boldsymbol{\theta}) p(y_{1:T-1} | \boldsymbol{\theta}) = \cdots = \prod_{t=1}^T p(y_t | y_{1:t-1}, \boldsymbol{\theta}).$$

From the model we have that

$$p(y_t \mid y_{1:t-1}, \boldsymbol{\theta}) = \mathcal{N}(y_t \mid \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{z}_t, \tau^{-1}),$$
  
where  $Z_t = (Y_{t-1}, Y_{t-2}, \dots, Y_{t-n})^{\mathsf{T}}$ . Hence,  
$$p(y_{1:T} \mid \boldsymbol{\theta}) = \prod_{t=1}^{\mathsf{T}} \mathcal{N}(y_t \mid \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{z}_t, \tau^{-1}) = \mathcal{N}(\boldsymbol{y} \mid \boldsymbol{z}\boldsymbol{\theta}, \tau^{-1}\boldsymbol{l}_{\mathsf{T}}),$$

where we have made use of  $\mathbf{Y} = (Y_1, Y_2, ..., Y_T)^T$  and  $\mathbf{Z} = (Z_1, Z_2, ..., Z_T)^T$ .

$$p(\boldsymbol{\theta}, \boldsymbol{y}) = \underbrace{\mathcal{N}(\boldsymbol{y} \mid \boldsymbol{z}\boldsymbol{\theta}, \tau^{-1} \boldsymbol{l}_{T})}_{p(\boldsymbol{y} \mid \boldsymbol{\theta})} \underbrace{\mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{0}, \rho^{-1} \boldsymbol{l}_{n})}_{p(\boldsymbol{\theta})}$$
$$= \mathcal{N}\left(\begin{pmatrix} \boldsymbol{\theta} \\ \boldsymbol{y} \end{pmatrix} \mid \begin{pmatrix} \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix}, \begin{pmatrix} \rho^{-1} \boldsymbol{l}_{2} & \rho^{-1} \boldsymbol{z}^{T} \\ \rho^{-1} \boldsymbol{z} & \tau^{-1} \boldsymbol{l}_{T} + \rho^{-1} \boldsymbol{z} \boldsymbol{z}^{T} \end{pmatrix} \right).$$

$$p(\boldsymbol{\theta}, \boldsymbol{y}) = \underbrace{\mathcal{N}(\boldsymbol{y} \mid \boldsymbol{z}\boldsymbol{\theta}, \tau^{-1} \boldsymbol{l}_{T})}_{p(\boldsymbol{y} \mid \boldsymbol{\theta})} \underbrace{\mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{0}, \rho^{-1} \boldsymbol{l}_{n})}_{p(\boldsymbol{\theta})}$$
$$= \mathcal{N}\left(\begin{pmatrix} \boldsymbol{\theta} \\ \boldsymbol{y} \end{pmatrix} \mid \begin{pmatrix} \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix}, \begin{pmatrix} \rho^{-1} \boldsymbol{l}_{2} & \rho^{-1} \boldsymbol{z}^{T} \\ \rho^{-1} \boldsymbol{z} & \tau^{-1} \boldsymbol{l}_{T} + \rho^{-1} \boldsymbol{z} \boldsymbol{z}^{T} \end{pmatrix} \right).$$

The posterior is given by

$$p(\boldsymbol{\theta} \mid \boldsymbol{y}) = \mathcal{N}(\boldsymbol{\theta} \mid m_T, S_T),$$

where

$$m_{T} = \tau S_{T} \mathbf{z}^{\mathsf{T}} \mathbf{y},$$
  
$$S_{T} = \left(\rho^{-1} l_{2} + \sigma \mathbf{z}^{\mathsf{T}} \mathbf{z}\right)^{\mathsf{T}}.$$

### Ex) Situation before any data is used

$$Y_t = A_1 Y_{t-1} + A_2 Y_{t-2} + E_t, \qquad E_t \sim \mathcal{N}(0, 0.2).$$



Prior

7 samples from the prior

White dot – true value for  $\theta = (0.6, 0.2)$ .

### Ex) Situation after $y_1$ is obtained



Likelihood

Posterior

7 samples from the posterior

### Ex) Situation after $y_{1:2}$ and $y_{1:20}$



posterior

Mathematical model: A compact representation—set of assumptions—of some phenomenon of interest.

**Probabilistic modelling:** Provides the capability to represent and manipulate **uncertainty** in data, models, decisions and predictions.

**Full probabilistic model:** The joint distribution of all observed (here  $y_{1:T}$ ) and unobserved (here  $\theta$ ) variables.

**Data distribution/likelihood:** Distribution describing the observed data conditioned on unobserved variables.

**Prior distribution:** Encodes initial assumptions on the unobserved variables.

**Posterior distribution:** Conditional distribution of the unobserved variables given the observed variables.