

BAYESIAN INFERENCE AND RELATED COMPUTATIONAL METHODS

BRIEF OVERVIEW

L. Martino

May, 2022

- ▶ **Preamble - or a very long introduction from a signal processing point of view...**
- ▶ Bayesian Inference
- ▶ Computational methods for Bayesian Inference
- ▶ Monte Carlo sampling methods - brief overview

BASIC, STANDARD PROBLEM

- ▶ In many applications, we are interested in inferring a variable of interest,

$$\boldsymbol{\theta} = [\theta_1, \dots, \theta_{d_\theta}] \in \Theta \subseteq \mathbb{R}^{d_\theta},$$

given a set of observations $\mathbf{y} \in \mathbb{R}^{d_Y}$.

- ▶ **We want to know $\boldsymbol{\theta}$ given \mathbf{y} :**

$$\boldsymbol{\theta} \implies \mathbf{y}$$

STANDARD - TYPICAL APPROACH

- ▶ Minimizing a cost function:

$$C(\boldsymbol{\theta}) = \text{Loss}(\boldsymbol{\theta}, \mathbf{y}),$$

obtaining

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} C(\boldsymbol{\theta}).$$

REGULARIZATION:

- ▶ fighting against overfitting,
- ▶ avoiding numerical problems and increasing the numerical stability.

CLASSICAL APPROACH + REGULARIZATION

Cost function:

$$C(\theta) = \text{Loss}(\theta, \mathbf{y}) + \text{Reg}(\theta)$$

Again, we minimize it (**optimization**).

- ▶ Why pass to a probabilistic domain/approach?
- ▶ **from optimization** \implies **to sampling**, why?
- ▶ a first answer in the next two slides.

Example: nonlinear model for regression

Regression problem: Assume that we have N data $\{x_i, y_i\}$.
We consider M bases

$$\phi_m(x) : \mathbb{R}^D \longrightarrow \mathbb{R}, \quad m = 1, \dots, M.$$

- ▶ We want that the solution has the following mathematical form:

$$\hat{f}(x) = \sum_{m=1}^M \theta_m \phi_m(x).$$

- ▶ We want to find the θ_m 's, we will consider

$$M \leq N.$$

Our nonlinear model

Assuming that we have N data $\{x_i, y_i\}$,

$$y_i = \sum_{m=1}^M \theta_m \phi_m(x_i) + \text{error} \dots$$

with

$$M \leq N.$$

Can we still associate a linear system?

Example: our observation model

the model is **non-linear**...but...

it is still linear...

with respect to the coefficients $\theta_1, \theta_2, \dots, \theta_M$, —
YES is still linear ! for this reason, it can be
analytically solved

Then, we can construct some matrices and vectors...

EXAMPLE: RECTANGULAR LINEAR SYSTEM ($M \leq N$)

We can define:

$$\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_M]^\top, \quad M \times 1$$

$$\boldsymbol{\Phi} = \begin{bmatrix} \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_M(x_1) \\ \phi_1(x_2) & \phi_2(x_2) & \dots & \phi_M(x_2) \\ \vdots & \vdots & & \\ \phi_1(x_N) & \phi_2(x_N) & \dots & \phi_M(x_N) \end{bmatrix} \quad N \times M,$$

$$\mathbf{y} = [y_1, y_2, \dots, y_N]^\top, \quad N \times 1.$$

Example: rectangular linear system

The system can be written as:

$$\Phi\theta = \mathbf{y}.$$

Check the dimensions

$$[N \times M] \times [M \times 1] = N \times 1.$$

Example: rectangular linear system

Since the system is rectangular, we cannot write

$$\theta = \Phi^{-1} \mathbf{y}, \quad \text{NOOOO!!!}$$

the inverse matrix Φ^{-1} does not exist ! since Φ is rectangular !

EXAMPLE

Cost function of the Regularized Least Squares:

$$C(\boldsymbol{\theta}) = \|\mathbf{y} - \Phi\boldsymbol{\theta}\|^2 + \lambda\|\boldsymbol{\theta}\|^2.$$

we want to minimize $C(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$.

EXAMPLE

Solution:

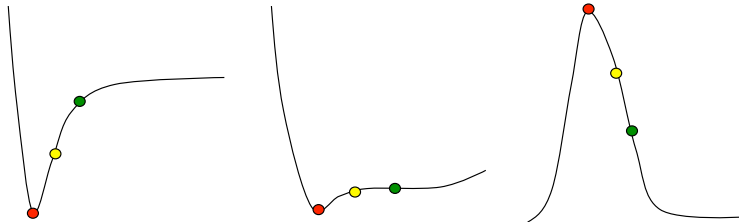
$$\hat{\boldsymbol{\theta}} = [\hat{\theta}_1, \dots, \hat{\theta}_M]^\top = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \lambda \mathbf{I}_M)^{-1} \boldsymbol{\Phi}^\top \mathbf{y}. \quad (1)$$

Recall that solution has the following mathematical form:

$$\hat{f}(x) = \sum_{m=1}^M \hat{\theta}_m \phi_m(x).$$

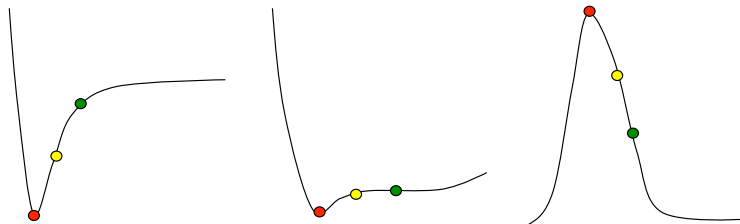
PROBABILISTIC APPROACH

(a) Differ among different cost functions, **(b)** different possible “points” for summarizing the cost function, **(c)** compute areas (e.g., for credible intervals) ...



PROBABILISTIC APPROACH

in the “probabilistic” approach: the MAP estimator, the MMSE estimator (*mean - expected value*), the median estimator are well-defined, and also the “areas” have a “meaning” and can be computed...



Probabilistic version...

Cost function:

$$C(\theta) = \underbrace{\text{Loss}(\theta, \mathbf{y})}_{\text{neg. log-likelihood}} + \underbrace{\text{Reg}(\theta)}_{\text{neg. log-prior}}$$

$$\text{likelihood} = p(\mathbf{y}|\theta) \propto \exp(-\text{Loss}(\theta, \mathbf{y})),$$

$$\text{prior} = p(\theta) \propto \exp(-\text{Reg}(\theta)).$$

Bayesian “slang”

Cost function:

$$C(\theta) = \underbrace{\text{Loss}(\theta, \mathbf{y})}_{\text{neg. log-likelihood}} + \underbrace{\text{Reg}(\theta)}_{\text{neg. log-prior}}$$

Bayesian Inference:

$$\text{posterior} \propto \exp\{-C(\theta)\} = \text{likelihood} \times \text{prior}$$

Bayesian Inference

Bayesian Inference:

posterior \propto likelihood \times prior.

where

$$\text{posterior} = p(\boldsymbol{\theta}|\mathbf{y}),$$

$$\text{likelihood} = p(\mathbf{y}|\boldsymbol{\theta}) \propto \exp(-\text{Loss}(\boldsymbol{\theta}, \mathbf{y})),$$

$$\text{prior} = p(\boldsymbol{\theta}) \propto \exp(-\text{Reg}(\boldsymbol{\theta})).$$

Prior versus Regularization

$$\text{prior} = p(\boldsymbol{\theta}) \propto \exp(-\text{Reg}(\boldsymbol{\theta})).$$

Main difference: the prior density must/should be normalized (“normalizable”)... (in some case, this condition can be also relaxed)

Normalization \implies since it represents probabilities \implies now we have more interpretability of different situations (think on different regularizations - e.g., previous figures - and different priors ...)

- ▶ Preamble - or a very long introduction from a signal processing point of view...
- ▶ **Bayesian Inference**
- ▶ Computational methods for Bayesian Inference
- ▶ Monte Carlo sampling methods - brief overview

Bayesian Inference:

1. **“Main Actors” in Bayesian inference**
2. Important considerations and consistency
3. Goals
4. Levels of inference
5. Type of Priors - choice of the priors
6. Reasons to be Bayesian

PROBLEM STATEMENT AND MAIN ACTORS

- ▶ In many applications, we are interested in inferring a variable of interest,

$$\boldsymbol{\theta} = [\theta_1, \dots, \theta_{d_\theta}] \in \Theta \subseteq \mathbb{R}^{d_\theta},$$

given a set of observations $\mathbf{y} \in \mathbb{R}^{d_Y}$.

PROBLEM STATEMENT AND MAIN ACTORS

- ▶ The **posterior probability density function (pdf)** is

$$\bar{\pi}(\theta) = p(\theta|\mathbf{y}) = \frac{\ell(\mathbf{y}|\theta)g(\theta)}{p(\mathbf{y})} \propto \ell(\mathbf{y}|\theta)g(\theta), \quad (2)$$

where

- ▶ $\ell(\mathbf{y}|\theta) = p(\mathbf{y}|\theta)$ is the likelihood function (induced by the observation model);
- ▶ $g(\theta) = p(\theta)$ is the prior pdf,
- ▶ $Z = p(\mathbf{y})$: marginal likelihood/Bayesian evidence.

(note that we have 2 conditionals $p(\theta|\mathbf{y})$, $\ell(\mathbf{y}|\theta)$, and 2 marginals, $g(\theta)$ and $p(\mathbf{y})$ - with the prior and the likelihood, we create a joint pdf of θ, \mathbf{y})

MARGINAL LIKELIHOOD - BAYESIAN EVIDENCE

- ▶ Given \mathbf{y} , the **marginal likelihood - Bayesian model evidence** is an integral (a *normalizing constant*):

$$Z = p(\mathbf{y}) = \int_{\Theta} \ell(\mathbf{y}|\boldsymbol{\theta})g(\boldsymbol{\theta})d\boldsymbol{\theta}.$$

is fixed and, in general, is **unknown**.

- ▶ $Z =$ **Weighted average of the likelihood values !!**
- ▶ Note that

$$0 \leq \min p(\mathbf{y}|\boldsymbol{\theta}) \leq Z \leq \max p(\mathbf{y}|\boldsymbol{\theta}) = p(\mathbf{y}|\hat{\boldsymbol{\theta}}_{\text{ML}}).$$

- ▶ $Z = p(\mathbf{y})$ useful for model selection purposes.

EXAMPLES OF “MODEL SELECTION”

Model selection: (some examples)

- ▶ Tuning of the parameters of the observation model (i.e., of the likelihood).
- ▶ Tuning of the parameters of prior.
- ▶ **Choose the best observation/measurement model among a set of possible models.**
- ▶ **Select the order/complexity in a model** (for instance, the order of a polynomial in a regression, or the order of an AR - FIR - filter etc.)
- ▶ **Variable selection.**
- ▶ etc.

EXAMPLES OF “MODEL SELECTION”

Model selection: - The previous examples can be divided into two main scenarios:

- ▶ Basic model selection
- ▶ Nested models

This classification is important for the possible choice of the priors.

- F. Llorente, L. Martino, E. Cuberlo, J. Lopez-Santiago, D. Delgado, "On the safe use of prior densities for Bayesian model selection", viXra:2110.0032 , 2021.

UNNORMALIZED POSTERIOR

- ▶ Since Z is generally unknown: then, in many cases, we are only able to evaluate the unnormalized pdf

$$\pi(\boldsymbol{\theta}) = \ell(\mathbf{y}|\boldsymbol{\theta})g(\boldsymbol{\theta}) \propto \bar{\pi}(\boldsymbol{\theta}).$$

Note that

$$\bar{\pi}(\boldsymbol{\theta}) = \frac{1}{Z}\pi(\boldsymbol{\theta}).$$

- ▶ Note that

$$Z = \int_{\Theta} \pi(\boldsymbol{\theta})d\boldsymbol{\theta}.$$

BREAK TO EXPLAIN: “EVALUATION VERSUS SAMPLING”

“Evaluation versus Sampling of a density” in this slides

- ▶ **Evaluation:** evaluate point-wise a function/density.
- ▶ **Sampling:** GENERATE RANDOM NUMBERS according to a density.

TYPICAL EXAMPLE OF APPLICATION: BAYESIAN INVERSION

- ▶ **Observation model - inducing likelihood:**

$$\mathbf{y} = \mathbf{G}(\boldsymbol{\theta}) + \mathbf{v},$$

where $\mathbf{G}(\boldsymbol{\theta})$ is a “physical” model (for instance) and \mathbf{v} is an independent Gaussian noise (for instance).

- ▶ **Likelihood:**

$$p(\mathbf{y}|\boldsymbol{\theta}) = \ell(\mathbf{y}|\boldsymbol{\theta}) \propto \exp\left(-\frac{\|\mathbf{y} - \mathbf{G}(\boldsymbol{\theta})\|^2}{2\sigma_v^2}\right).$$

- ▶ the “goal” is: virtually, $\boldsymbol{\theta} = \mathbf{G}^{-1}(\mathbf{y})$. For this reason, it is called “inversion”

Bayesian Inference:

1. “Main Actors” in Bayesian inference
2. **Important considerations and consistency**
3. Goals
4. Levels of inference
5. Type of Priors - choice of the priors
6. Reasons to be Bayesian

MAIN DIFFERENCES WITH FREQUENTIST APPROACH

- ▶ **The vector of data \mathbf{y} is given and fixed.**
- ▶ θ "should" be considered random, since we assume $\theta \sim g(\theta)$.
- ▶ But "practical Bayesians" and/or "Bayesians with common sense", considers/knows that it exists a (fixed) θ_{true} , that we desire to infer.
- ▶ In fact, under mild conditions, **the Bayesian estimators are consistent as the number of data grows.**
- ▶ (note that bias, variance and MSE of an estimator are more frequentist ideas/quantities since consider expectation over \mathbf{y} ... we can extend these concepts here, considering different posteriors - one for each \mathbf{y}' generated according to the model- and then make an average...)

CONSISTENCY

- ▶ Under mild conditions, **the Bayesian estimators are consistent as the number of data grows.**
- ▶ **Consistency:** As the number of data grows, the posterior becomes more tighter and tighter around θ_{true} .

Bayesian Inference:

1. "Main Actors" in Bayesian inference
2. Important considerations and consistency
3. **Goals**
4. Levels of inference
5. Type of Priors - choice of the priors
6. Reasons to be Bayesian

MAIN GOAL

- ▶ **Goal:** extract and summarize the statistical information contained in the posterior pdf $\bar{\pi}(\theta)$ and compute Z .

GOAL - QUADRATURE PROBLEMS

- ▶ More specifically, *in many cases*, our goal is to compute efficiently some integral involving π ,

$$\mathbf{I} = \frac{1}{Z} \int_{\Theta} \mathbf{f}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}, \quad (3)$$

where $\mathbf{f}(\boldsymbol{\theta}) : \Theta \rightarrow \mathbb{R}^n$, and

$$Z = \int_{\Theta} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

In most of the cases, Z is also unknown and we have to estimate it (useful for [model selection](#) purposes).

GOAL - QUADRATURE PROBLEMS - MOMENTS

- ▶ **Example:** If $f(\theta) = \theta$, the integral $\mathbf{I} = \int_{\Theta} \theta \bar{\pi}(\theta) d\theta$ represents the **MMSE estimator** - the reason of this name required more hours of course...
- ▶ More generally, **all the moments of the posterior are:**

$$\mathbf{I}_k = \int_{\Theta} \theta^k \bar{\pi}(\theta) d\theta,$$

$$k = 1, 2, 3 \dots$$

MODEL SELECTION - MARGINAL LIKELIHOOD - AGAIN QUADRATURE PROBLEM

- ▶ Marginal likelihood - Bayesian model evidence:

$$Z = p(\mathbf{y}) = \int_{\Theta} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int_{\Theta} \ell(\mathbf{y}|\boldsymbol{\theta}) g(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

- ▶ $Z =$ **Weighted average of the likelihood values !!**
- ▶ Note that

$$0 \leq \min p(\mathbf{y}|\boldsymbol{\theta}) \leq Z \leq \max p(\mathbf{y}|\boldsymbol{\theta}) = p(\mathbf{y}|\hat{\boldsymbol{\theta}}_{\text{ML}}).$$

GOAL

- ▶ **Goal:** extract and summarize the statistical information contained in the posterior pdf $\bar{\pi}(\boldsymbol{\theta})$ and compute Z .
- ▶ **The problem of extracting information from $\bar{\pi}(\boldsymbol{\theta})$ is mainly converted in a quadrature problem.**

Bayesian Inference:

1. “Main Actors” in Bayesian inference
2. Important considerations and consistency
3. Goals
4. **Levels of inference**
5. Type of Priors - choice of the priors
6. Reasons to be Bayesian

LEVELS OF INFERENCE

In the standard/basic framework $(g(\theta), \ell(\mathbf{y}|\theta), Z)$:

- ▶ **Level 1:** Inference about θ .
- ▶ **Level 2:** Learning/obtaining Z .

different levels \implies different “rules”, **in this sense different priors can be used or have different meanings...**

MORE LEVELS: HIERARCHICAL MODELING

Hierarchical modeling: $h(\nu)$ prior over ν , $g(\theta|\nu)$ prior over θ

- ▶ **Level 0:** Inference about ν .
- ▶ **Level 1:** Inference about θ .
- ▶ **Level 2:** Learning/obtaining Z ; in this case we have also several $Z|\nu$.

We can use Level 2 for learning ν or use a full-Bayesian solution. See Section of:

- F. Llorente, L. Martino, E. Cuberlo, J. Lopez-Santiago, D. Delgado, "On the safe use of prior densities for Bayesian model selection", viXra:2110.0032 , 2021.

Bayesian Inference:

1. “Main Actors” in Bayesian inference
2. Important considerations and consistency
3. Goals
4. Levels of inference
5. **Type of Priors - choice of the priors**
6. Reasons to be Bayesian

TYPE OF PRIORS !

- ▶ There are several type of priors. See as an example **Section 3.3:**

- F. Llorente, L. Martino, E. Cuberlo, J. Lopez-Santiago, D. Delgado, "On the safe use of prior densities for Bayesian model selection", viXra:2110.0032 , 2021.

PRIORS FOR MODEL SELECTION

- ▶ Marginal likelihood - Bayesian model evidence:

$$Z = p(\mathbf{y}) = \int_{\Theta} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int_{\Theta} \ell(\mathbf{y}|\boldsymbol{\theta}) g(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

- ▶ Issues when the prior is improper...

- F. Llorente, L. Martino, E. Cuberlo, J. Lopez-Santiago, D. Delgado, "On the safe use of prior densities for Bayesian model selection", viXra:2110.0032 , 2021.

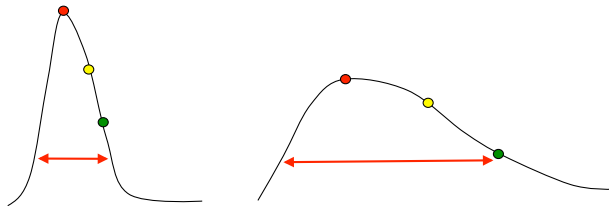
- F. Llorente, L. Martino, D. Delgado, J. Lopez-Santiago, "Marginal likelihood computation for model selection and hypothesis testing: an extensive review", (to appear) SIAM Review, 2022. arXiv:2005.08334

Bayesian Inference:

1. “Main Actors” in Bayesian inference
2. Important considerations and consistency
3. Goals
4. Levels of inference
5. Type of Priors - choice of the priors
6. **Reasons to be Bayesian**

REASONS TO USE A BAYESIAN APPROACH

- ▶ include **prior information** in our model
- ▶ **different possible point estimators** (not only maximum...)



REASONS TO USE A BAYESIAN APPROACH

- ▶ Provide **complete posterior information**:
 - ▶ **including information** by the prior density
 - ▶ **quantify uncertainties** (histograms)
 - ▶ **credible intervals** (areas)
 - ▶ **quantiles** (areas)
 - ▶ number of modes and modes
 - ▶ **correlations among parameters** (components of θ) (multi-dimensional histograms)
 - ▶ **dependence/sensibility analysis of the model with respect to the components of θ** (related to the gradient/derivatives of the model/transformation and the variance of the histograms - marginal densities)

REASONS TO USE A BAYESIAN APPROACH

- ▶ **easier application and interpretation of statistical quantities and procedures**: credible intervals (easier than confidence intervals), hypothesis testing, model selection etc... **in my opinion, this is the main benefit with respect to the frequentist approach.**
- ▶ **regularization** - numerical stability (due to the prior)
- ▶ **model selection** (marginal likelihood - Z)

EXAMPLES OF θ

Inference about θ :

- ▶ θ can be a vector of **parameters** of a model
- ▶ θ can be a vector of **hyper-parameters**
- ▶ θ can be a **model index** (model selection),
- ▶ θ can include the **number of parameters** in a model (complexity of the model)
- ▶ More specifically, θ can represent the position or the trajectory of a target, the volatility in a financial time series, velocity and direction of the wind etc.

- ▶ Preamble - or a very long introduction from a signal processing point of view...
- ▶ Bayesian Inference
- ▶ **Computational methods for Bayesian Inference**
- ▶ Monte Carlo sampling methods - brief overview

COMPUTATIONAL METHODS FOR BAYESIAN INFERENCE

- ▶ In many applications, we are not able to compute analytically \mathbf{I} and Z .
- ▶ **Numerical approximations:**
 1. Deterministic quadrature rules - Gaussian-Hermite, Cubature rules...
 2. Monte Carlo techniques
 3. Variance Reduction - Quasi Monte Carlo (negative correlation)
 4. Variational inference techniques
 5. other modern quadrature rules

COMPUTATIONAL METHODS FOR BAYESIAN INFERENCE (2)

- ▶ The methods 1-2-3 of the previous list can be considered **quadrature techniques**.

COMPUTATIONAL METHODS FOR BAYESIAN INFERENCE (2)

- ▶ The methods 1-2-3 of the previous list can be summarized with this formula:

$$\mathbf{I} = \frac{1}{Z} \int_{\Theta} \mathbf{f}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \approx \hat{\mathbf{I}}_N = \sum_{n=1}^N \bar{w}_n \mathbf{f}(\boldsymbol{\theta}_n), \quad (4)$$

- ▶ **Error bounds:**

$$\|\mathbf{I} - \hat{\mathbf{I}}_N\| \leq V(\mathbf{f}) D_N(\boldsymbol{\theta}_{1:N}), \quad (5)$$

- $V(\mathbf{f})$ depends on variation of \mathbf{f} in Θ
- $D_N(\boldsymbol{\theta}_{1:N})$ depends on choice of the nodes/samples $\boldsymbol{\theta}_{1:N}$
- ▶ **Clearly, $D_N(\boldsymbol{\theta}_{1:N}) \rightarrow 0$ when $N \rightarrow \infty$**

- ▶ Preamble - or a very long introduction from a signal processing point of view...
- ▶ Bayesian Inference
- ▶ Computational methods for Bayesian Inference
- ▶ **Monte Carlo sampling methods - brief overview**

Monte Carlo Sampling Methods

- ▶ They are **random number generators** \implies from a **generic density - given an available random source**
- ▶ that can be used for building/designing **stochastic quadrature rules in Bayesian inference.**

Important: Monte Carlo sampling methods are random number generators, they have life out of the “Bayesian world...” - but the main application is in Bayesian inference.

MORE THAN A QUADRATURE RULE AND AN OPTIMIZER...

with Monte Carlo:

- ▶ We can also optimize $\bar{\pi}(\boldsymbol{\theta})$ (or $\pi(\boldsymbol{\theta})$ is the same): **global optimization, the “true” optimization !!**
- ▶ But with optimization we just get one point - as I said before, we want to **extract and summarize the statistical information contained in $\bar{\pi}(\boldsymbol{\theta})$.**

SAMPLING >> OPTIMIZATION

- ▶ With a Monte Carlo sampling method, we have also an “optimizer”.
- ▶ Optimization can be used for obtaining better “samplers”.

APPROXIMATION OF THE “MEASURE OF THE POSTERIOR”

- ▶ We obtain a particle approximation (with N samples)

$$\hat{\pi}(\boldsymbol{\theta}) = \sum_{i=1}^N \bar{w}_i \delta(\boldsymbol{\theta} - \boldsymbol{\theta}^{(i)}).$$

STANDARD MONTE CARLO APPROXIMATION

- ▶ If we can generate N i.i.d. random vectors θ_n distributed according to $\bar{\pi}(\theta)$, $n = 1, \dots, N$, then

$$\hat{\mathbf{I}}_N = \frac{1}{N} \sum_{n=1}^N \mathbf{f}(\theta_n) \approx \mathbf{I}, \quad \theta_n \sim \bar{\pi}(\theta).$$

- ▶ However:
 - ▶ Often it is not possible to draw from $\bar{\pi}(\theta)$.
 - ▶ Even in this "ideal" case, it is not straightforward to approximate Z , i.e., to find $\hat{Z} \approx Z$.

SAMPLING METHODS

If it is not possible to draw directly from the **target** pdf $\bar{\pi}(\boldsymbol{\theta})$:

- ▶ One can draw samples from a simpler **proposal pdf**, $q(\boldsymbol{\theta})$, and then *filter properly* these samples for obtaining an approximation of I and Z .

Sampling algorithm: all the steps corresponding to this “filtering operations”.

$$\mathbf{z}_1, \dots, \mathbf{z}_M \sim q(\boldsymbol{\theta}) \longrightarrow \boxed{\text{MC sampling}} \longrightarrow \hat{\pi}(\boldsymbol{\theta}) = \sum_{m=1}^M \bar{w}_m \delta(\boldsymbol{\theta} - \mathbf{z}_m)$$

$M \geq N =$ effective number of samples

$$\sum_{m=1}^M \bar{w}_m = 1 \quad \longrightarrow \quad \begin{array}{ll} \bar{w}_m \propto \{0, 1\} & * \text{ accept/reject} \\ \bar{w}_m = \text{by repetition} & * \text{ MCMC} \\ \bar{w}_m = \text{generic} & * \text{ importance} \\ \bar{w}_m = \frac{1}{M} & * \text{ standard MC} \end{array}$$

SAMPLING METHODS: CLASSIFICATION

(STATIC SCENARIO)

4 main classes of algorithms:

1. Direct methods (based on random variable transformation).
 - ▶ **Independent samples.** (the best, almost)
 - ▶ computational effort: lowest.
 - ▶ applicability: low.
2. Rejection sampling.
 - ▶ **Independent samples.** (the best, almost)
 - ▶ computational effort: higher (depending on the acceptance rate).
 - ▶ applicability: wider of direct methods, but in general low.
3. Importance sampling (IS).
 - ▶ **Weighted samples.**
 - ▶ computational effort: low.
 - ▶ applicability: always. - **Easy approx of Z**
4. Markov chain Monte Carlo (MCMC).
 - ▶ **(positive) Correlated samples.**
 - ▶ computational effort: low.
 - ▶ applicability: always. - **Exploration of the space**

(the best scenario is: negative correlated samples \Rightarrow , e.g., including deterministic “parts” within the method)

SAMPLING VS VARIATIONAL/OTHERS APPROACHES

Benefits:

- ▶ **Applicability - flexible** (easy to apply to different problems/framework) - the unique requirement is to be able to evaluate $\pi(\boldsymbol{\theta}) \propto \bar{\pi}(\boldsymbol{\theta})$; this condition can even relax.
- ▶ Complete approximation of the posterior that **can be improved increasing the computational cost**.

Drawbacks:

- ▶ **“Slower”** - more computational demanding (depending on the specific requirements of the considered application).

SPEED UP - MONTE CARLO

To speed up:

- ▶ given the application, specific algorithm design:
 - better proposal choice
 - include more information of the posterior (e.g., gradient)
 - include determinism (when it is possible, in a proper way)

PERFORMANCE OF A SAMPLING METHOD

Black-box point of view:

- ▶ The performance strictly depends on the choice of $q(\theta)$.
- ▶ We desire $q(\theta) \approx \bar{\pi}(\theta)$.

This is the reason for employing **adaptive techniques**.

- ▶ I worked (a lot) with adaptive MCMC and adaptive IS.

STATIC - DYNAMIC PARAMETERS

$$\theta = [\mathbf{x}, \boldsymbol{\lambda}]$$

Dynamic static

(a factorization of the posterior is available)

$$\mathbf{x} = [x_1, \dots, x_{d_x}] \in \mathbb{R}^{d_x},$$

$$\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_{d_\lambda}] \in \mathbb{R}^{d_\lambda},$$

$$d_\theta = d_x + d_\lambda.$$

DIFFERENT FRAMEWORKS - APPLICATIONS

TABLE: Scenarios where we are able to evaluate $\pi(\theta|\mathbf{y})$.

Computational Scenarios	Monte Carlo approach
$\pi(\boldsymbol{\lambda} \mathbf{y})$	static-batch scenario
$\pi(\lambda_1 \mathbf{y}, \lambda_{2:d_\lambda} = \text{fixed}),$ $\pi(\lambda_2 \mathbf{y}, \lambda_1 = \text{fixed}, \lambda_{3:d_\lambda} = \text{fixed}),$ $\pi(\lambda_3 \mathbf{y}, \lambda_{1:2} = \text{fixed}, \lambda_{4:d_\lambda} = \text{fixed}), \dots$	static-batch scenario with Gibbs approach , component-wise approach
$\pi(\boldsymbol{\lambda} y_1), \pi(\boldsymbol{\lambda} y_1, y_2), \pi(\boldsymbol{\lambda} y_{1:\ell}) \dots$ $\pi(\boldsymbol{\lambda} y_{1:d_Y} = \mathbf{y}).$	data tempering or online inference (data streaming)
$\pi(\boldsymbol{\lambda} y_1, y_2, y_3), \pi(\boldsymbol{\lambda} y_4), \pi(\boldsymbol{\lambda} y_5, y_6) \dots$	parallel - distributed - diffused estimation Big Data
$\pi(x_1 \mathbf{y}), \pi(x_1, x_2 \mathbf{y}), \pi(x_1 x_2, x_3 \mathbf{y}) \dots$	the dimension of x 's increases progressively (classifier and regressor chains)
$\pi(x_1 y_1), \pi(x_1, x_2 y_1, y_2), \pi(x_{1:3} y_{1:3}) \dots$ $\pi(x_{1:d_x} = \mathbf{x} y_{1:d_Y} = \mathbf{y}).$	completely sequential scenario - HMM Kalman Filters; Particle Filters state space models
$\pi(x_1, \boldsymbol{\lambda} y_1), \pi(x_{1,2}, \boldsymbol{\lambda} y_{1,2}), \pi(x_{1:3}, \boldsymbol{\lambda} y_{1:3}) \dots$ $\pi(x_{1:d_x} = \mathbf{x}, \boldsymbol{\lambda} y_{1:d_Y} = \mathbf{y}).$	"Tracking and parameter estimation in state-space models"

DIFFERENT FRAMEWORKS - APPLICATIONS

TABLE: More “strange” scenarios; e.g., we cannot evaluate $\pi(\theta|\mathbf{y})$.

Computational Scenarios	Monte Carlo approach
$Z_X(\theta) = \int \ell(\mathbf{y} \theta) d\mathbf{y}$ unknown	methods for “double intractable” posteriors pseudo-marginal MCMC
costly likelihood, or impossible to evaluate the likelihood, or “too much” data	Approximate Bayesian Computation (ABC) , pseudo-marginal MCMC, noisy MC Monte Carlo for Big Data
unknown dimension d_θ of $\theta = [\theta_1, \dots, \theta_{d_\theta}]$	“tracking with unknown number of targets”, “change point detection”, inference also about d_θ , Reversible Jump MCMC, Particle learning
model selection	inference + choose the best model (related to the previous point)

▶ Markov Chain Monte Carlo (MCMC)

- L. Martino, V. Elvira. "Metropolis Sampling", Wiley StatsRef: Statistics Reference Online, 2017. arXiv:1704.04629
- L. Martino, "A Review of Multiple Try MCMC algorithms for Signal Processing", Digital Signal Processing, Volume 75, Pages: 134-152, 2018.

MARKOV CHAIN MONTE CARLO (MCMC)

- ▶ Markov Chain Monte Carlo (MCMC) techniques yield an ergodic Markov chain

$$\theta_1 \rightarrow \theta_2 \rightarrow \dots \theta_t \rightarrow \dots \theta_{T-1} \rightarrow \theta_T,$$

with a **stationary/invariant density**, that is exactly the **posterior** $\bar{\pi}(\theta)$.

- ▶ There exists a $t_b < \infty$ (length of the **burn-in period**), such that

$$\theta_t \sim \bar{\pi}(\theta), \quad \text{for } t \geq t_b, \quad (6)$$

i.e., the marginal pdf of θ_t is the posterior.

KERNEL OF A MCMC ALGORITHM

- ▶ An MCMC method is completely defined by the probability to obtain a new state θ_t given the previous one, θ_{t-1} .
- ▶ The corresponding conditional density $K(\theta_t|\theta_{t-1})$ is usually called **kernel**.
- ▶ $K(\theta_t|\theta_{t-1})$ summarizes all the steps of the MCMC algorithm.

INVARIANT/STATIONARY DISTRIBUTION

- ▶ Definition of invariant-stationary pdf $p_S(\mathbf{x}_t)$:

$$\int_{\Theta} K(\boldsymbol{\theta}_t|\boldsymbol{\theta}_{t-1})p_S(\boldsymbol{\theta}_{t-1})d\boldsymbol{\theta}_{t-1} = p_S(\boldsymbol{\theta}_t). \quad (7)$$

- ▶ **MCMC method:** design $K(\boldsymbol{\theta}_t|\boldsymbol{\theta}_{t-1})$ in order to have

$$p_S(\boldsymbol{\theta}) = \bar{\pi}(\boldsymbol{\theta}). \quad (8)$$

EIGENFUNCTIONS

- ▶ This problem is related to the search of **eigenvalues** and **eigenfunctions** in the equation

$$\int_{\Theta} K(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}) \phi(\boldsymbol{\theta}_{t-1}) d\boldsymbol{\theta}_{t-1} = \mu \phi(\boldsymbol{\theta}_t) \quad (9)$$

where μ is an eigenvalue and $\phi(\cdot)$ is an eigenfunction (corresponding to μ).

BALANCE CONDITION

- ▶ Use the definition of invariance is difficult, in general.
- ▶ The **balance condition**

$$K(\boldsymbol{\theta}_t|\boldsymbol{\theta}_{t-1})\bar{\pi}(\boldsymbol{\theta}_{t-1}) = K(\boldsymbol{\theta}_{t-1}|\boldsymbol{\theta}_t)\bar{\pi}(\boldsymbol{\theta}_t), \quad (10)$$

is a sufficient condition to prove the invariance.

- ▶ If a density satisfies the balance condition, then is invariant w.r.t. the kernel $K(\boldsymbol{\theta}_t|\boldsymbol{\theta}_{t-1})$.
- ▶ In this case, the chain is **reversible**.

METROPOLIS-HASTINGS (MH) ALGORITHM

- ▶ Recall $\pi(\boldsymbol{\theta}) \propto \bar{\pi}(\boldsymbol{\theta})$.
- ▶ Proposal pdf: $q(\boldsymbol{\theta}|\boldsymbol{\theta}_{t-1})$.

MH algorithm:

- Choose $\boldsymbol{\theta}_0$.
- For $t = 1, \dots, T$:
 1. Draw $\boldsymbol{\theta}'$ from $q(\boldsymbol{\theta}|\boldsymbol{\theta}_{t-1})$.
 2. Set $\boldsymbol{\theta}_t = \boldsymbol{\theta}'$ with probability

$$\alpha = \min \left[1, \frac{\pi(\boldsymbol{\theta}')q(\boldsymbol{\theta}_{t-1}|\boldsymbol{\theta}')}{\pi(\boldsymbol{\theta}_{t-1})q(\boldsymbol{\theta}'|\boldsymbol{\theta}_{t-1})} \right]. \quad (11)$$

Otherwise, set $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1}$ (with probability $1 - \alpha$).

- Output: $\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_T\}$.
-

EFFECTIVE SAMPLE SIZE (ESS) FOR MCMC

- ▶ The samples are **(positive) correlated**.
- ▶ Due to the ergodicity:

$$\tilde{\mathbf{I}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{f}(\boldsymbol{\theta}_t) \approx \mathbf{I}. \quad (12)$$

(recall that we should consider only $t \geq t_b$)

- ▶ **Effective Sample Size (ESS):**

$$T_{eff} = T \frac{\text{var}[\hat{\mathbf{I}}_T]}{\text{var}[\tilde{\mathbf{I}}_T]} \approx \frac{T}{1 + 2 \sum_{k=1}^{\infty} \rho_k}, \quad (13)$$

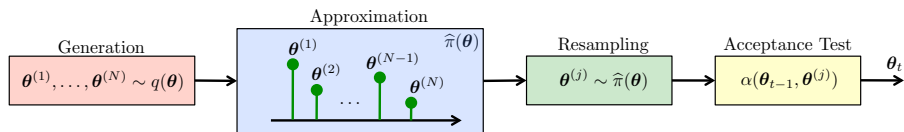
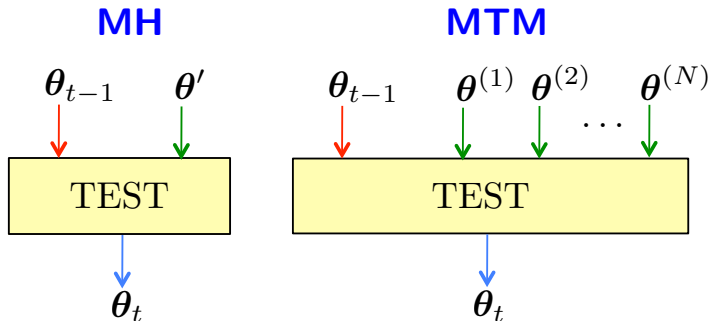
where $\rho_k = \frac{\text{COV}[f(\boldsymbol{\theta}_t), f(\boldsymbol{\theta}_{t+k})]}{\text{var}[f(\boldsymbol{\theta}_t)]}$.

IMPROVE PERFORMANCE

To reduce the correlation, speed up the convergence:

- ▶ MCMC with **and adaptive proposal density**
- ▶ Adding **gradient information** to the proposal pdf - Hamiltonian Monte Carlo
- ▶ Design more efficient algorithms: **Multiple Try Metropolis (MTM)**
- ▶ In high dimension, work component by component - **Gibbs sampling**
- ▶ Design **non-reversible MCMC methods**.

MULTIPLE TRY METROPOLIS (MTM)



► Importance Sampling (IS)

- V. Elvira, L. Martino, D. Luengo, M. F. Bugallo, "Generalized Multiple Importance Sampling", Statistical Science, Volume 34, Number 1, Pages 129-155, 2019.
- L. Martino, V. Elvira, F. Louzada, "Effective Sample Size for Importance Sampling Based on Discrepancy Measures", Signal Processing, Volume 131, Pages: 386-401, 2017

IMPORTANCE SAMPLING (IS)

- ▶ Consider the following equality:

$$\begin{aligned} \mathbf{I} = E_{\bar{\pi}}[\mathbf{f}(\boldsymbol{\theta})] &= \int_{\Theta} \mathbf{f}(\boldsymbol{\theta}) \bar{\pi}(\boldsymbol{\theta}) d\boldsymbol{\theta}, \\ &= \frac{1}{Z} \int_{\Theta} \mathbf{f}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}, \\ &= \frac{1}{Z} \int_{\Theta} \mathbf{f}(\boldsymbol{\theta}) \frac{\pi(\boldsymbol{\theta})}{q(\boldsymbol{\theta})} q(\boldsymbol{\theta}) d\boldsymbol{\theta}, \\ &= \frac{1}{Z} \int_{\Theta} \mathbf{f}(\boldsymbol{\theta}) w(\boldsymbol{\theta}) q(\boldsymbol{\theta}) d\boldsymbol{\theta}, \\ &= \frac{1}{Z} E_q[\mathbf{f}(\boldsymbol{\theta}) w(\boldsymbol{\theta})] \end{aligned}$$

where $w(\boldsymbol{\theta}) = \frac{\pi(\boldsymbol{\theta})}{q(\boldsymbol{\theta})}$.

- ▶ **Importance Sampling:** apply the standard MC procedure for approximating $E_q[\mathbf{f}(\boldsymbol{\theta}) w(\boldsymbol{\theta})]$ (when Z is known).

IS ESTIMATORS (WITH UNKNOWN Z)

1. **Sampling:** N samples from the proposal $q(\theta)$

$$\theta_n \sim q(\theta), \quad n = 1, \dots, N.$$

2. **Weighting:** Each sample is “corrected” by the importance weight

$$w_n = \frac{\pi(\theta_n)}{q(\theta_n)}, \quad n = 1, \dots, N.$$

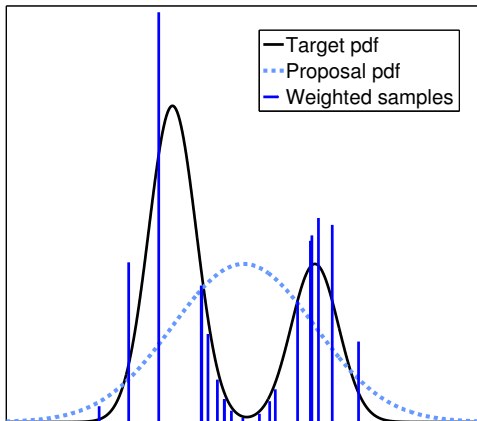
3. **Estimators:**

$$\tilde{\mathbf{I}}_N = \sum_{n=1}^N \bar{w}_n \mathbf{f}(\theta_n), \quad \hat{Z} = \frac{1}{N} \sum_{n=1}^N w_n,$$

where

$$\bar{w}_n = \frac{w_n}{\sum_{i=1}^N w_i} = \frac{w_n}{N\hat{Z}}.$$

EXAMPLE - IS



PROPER WEIGHTING

- ▶ Is the previous weighting scheme unique? No.
- ▶ Consider an extended proposal pdf $q_e(\boldsymbol{\theta}, w) = q(w|\boldsymbol{\theta})q(\boldsymbol{\theta})$.
- ▶ **Properly weighted samples** with respect to $\bar{\pi}$:

$$E_{q_e}[W(\boldsymbol{\theta}) \mathbf{f}(\boldsymbol{\theta})] = cE_{\bar{\pi}}[\mathbf{f}(\boldsymbol{\theta})], \quad (14)$$

where $c > 0$ is a constant value.

- ▶ Different possible weighting schemes.
- ▶ Easy to see when different proposal densities are used jointly.

MULTIPLE IMPORTANCE SAMPLING (MIS)

Consider N proposal densities, $q_1(\boldsymbol{\theta}), \dots, q_N(\boldsymbol{\theta})$.

- ▶ **Sampling:** $\boldsymbol{\theta}_n \sim q_n(\boldsymbol{\theta})$ with $n = 1, \dots, N$.
- ▶ **Classical Weighting (CW):**

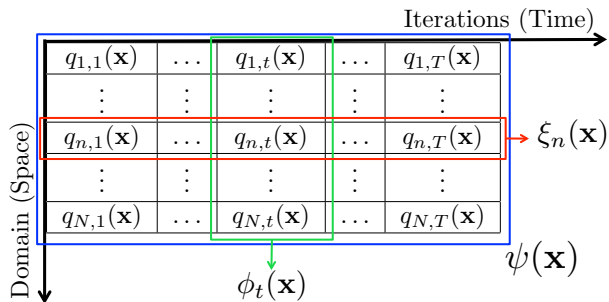
$$w_n = \frac{\pi(\boldsymbol{\theta}_n)}{q(\boldsymbol{\theta}_n)}, \quad n = 1, \dots, N.$$

- ▶ **Deterministic Mixture (DM) Weighting:**

$$w_n = \frac{\pi(\boldsymbol{\theta}_n)}{\frac{1}{N} \sum_{k=1}^N q_k(\boldsymbol{\theta}_n)}, \quad n = 1, \dots, N.$$

- ▶ The DM-IS estimators have lower variance than the CW-IS estimators (but more costly; a bit).
- ▶ There are even more possibilities: for instance, **the partial DM weights**.

ADAPTATION AND MIS



- Trade-off: complexity - performance

▶ Particle Filtering

- L. Martino, J. Read, V. Elvira, F. Louzada, "Cooperative Parallel Particle Filters for on-Line Model Selection and Applications to Urban Mobility" Digital Signal Processing Volume 60, Pages: 172-185, 2017.

PARTICLE FILTERING

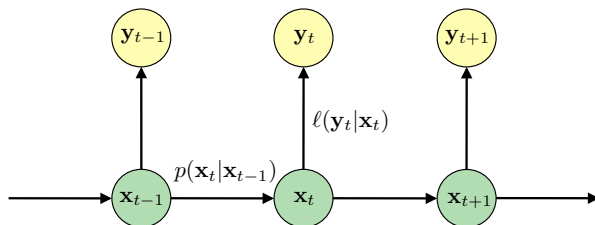
- ▶ **Particle Filtering = Sequential Importance Sampling + Resampling**
- ▶ In signal processing, mainly used in state-space models.

STATE-SPACE MODELS

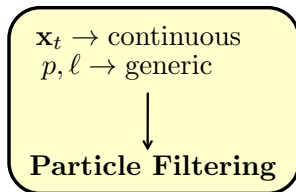
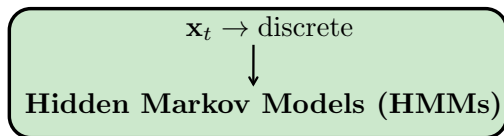
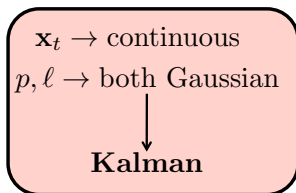
- ▶ $t \in \mathbb{N}$: discrete iteration index,
- ▶ $\mathbf{x}_t \in \mathbb{R}^{d_x}$: state variable that we want to infer,
- ▶ $\mathbf{y}_t \in \mathbb{R}^{d_y}$: observation at time t ,

$$\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}), \quad (\text{propagation}), \quad (15)$$

$$\mathbf{y}_t \sim \ell(\mathbf{y}_t | \mathbf{x}_t), \quad (\text{likelihood}). \quad (16)$$



RELATIONSHIPS WITH OTHER METHODS



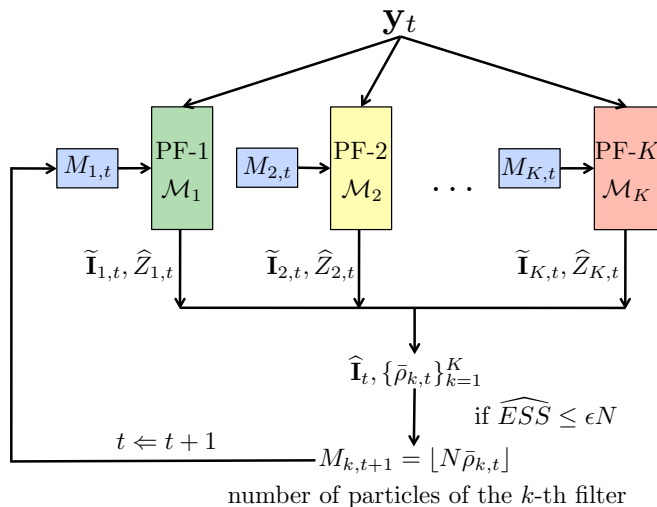
STATE-SPACE MODELS WITH UNKNOWN PARAMETER

$$\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}, \boldsymbol{\lambda}_p), \quad (\text{propagation}), \quad (17)$$

$$\mathbf{y}_t \sim \ell(\mathbf{y}_t | \mathbf{x}_t, \boldsymbol{\lambda}_\ell), \quad (\text{likelihood}). \quad (18)$$

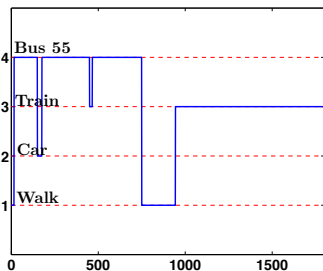
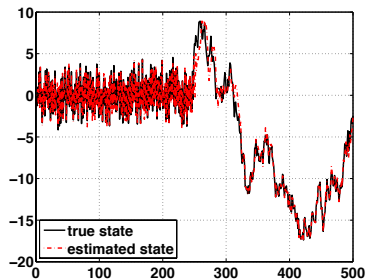
We can also select the best state-space model within a set of possible models.

COOPERATIVE PARALLEL PARTICLE FILTERS



We also adapt the number of particles in each filter (but the sum of all particles is fixed)

COOPERATIVE PARALLEL PARTICLE FILTERS



- ▶ Thank you very much!
- ▶ Any questions?