

SELFI enhanced: robustness to model misspecification and Euclid forecast



Euclid-France Theory and Likelihood workshop Institut d'Astrophysique de Paris



Florent Leclercq

www.florent-leclercq.eu

Institut d'Astrophysique de Paris CNRS & Sorbonne Université



www.aquila-consortium.org

28 November 2022



Why I decided to go "likelihood-free" for Galaxy Clustering Additional Probes (GC:AP)

Note: likelihood-free inference \approx simulation-based inference \approx implicit likelihood inference

• A question of <u>accuracy</u>: first, avoid biases.



 The HOWLS group also found that some WL additional probes have a non-Gaussian distribution.



• A question of <u>precision</u>: can numerical forward models be used to push further than $k \gtrsim 0.15 h/Mpc$? The full field contains much more information.

HOWLS-KP paper 1, Ajani et al., in prep.



SELFI enhanced: robustness to model misspecification and Euclid forecast 28/11/2022 2

The issue of model misspecification in Bayesian inference and in simulation-based inference (SBI)

- Model misspecification arises when model differs from actual data-generating process.
- Field-based inference techniques have a successful track record at handling model misspecification, e.g. automatically reporting unknown data contaminations.





Model misspecification: a major challenge particularly for approaches that marginalise over latent variables, such as simulation-based inference (SBI).

Porqueres, Ramanah, Jasche & Lavaux, 1812.05113 Lavaux, Jasche & FL, 1909.06396



Florent Leclercq



 Typical cosmological example: the galaxy power spectrum at large scales.



SELFI enhanced: robustness to model misspecification and Euclid forecast 28/11/2022

A general class of Bayesian hierarchical models (BHMs): Complex observations of a latent function controlled by top-level parameters





Key idea: a two-step SBI process that recycles simulations



- 1. Inference of the latent function θ , to check for model misspecification:
 - SELFI algorithm



Key idea: a two-step SBI process that recycles simulations



- 1. Inference of the latent function θ , to check for model misspecification:
 - SELFI algorithm
- 2. Simulation-based inference of ω :
 - Approximate Bayesian Computation (ABC), Likelihood-Free Rejection Sampling
 - Density/ratio estimation (DELFI / NRE)
 - Bayesian optimisation (BOLFI)

others...

Important: the simulations necessary for step **1**. are recycled for data compression, which is required for step **2**.



Latent function inference: the SELFI approach (Simulator Expansion for Likelihood-Free Inference)



Florent Leclercq

• We aim at inferring the latent function θ , which usually <u>contains most/all of the information</u> on ω .

(initial power spectrum in cosmology, prey/predator population functions in ecology)

- This requires doing SBI in $d = \mathcal{O}(100) \mathcal{O}(1,000)$
- If we trust the results of earlier experiments, we can Taylor-expand the black-box around an expansion point θ₀:

$$\hat{\Phi}_{\theta} \approx \mathbf{f}_0 + \nabla \mathbf{f}_0 \cdot (\theta - \theta_0) + \frac{1}{2} (\theta - \theta_0)^{\mathsf{T}} \cdot \mathbf{H} \cdot (\theta - \theta_0) + \dots$$

SELFI-2 (second order): coming soon!

 Gradients, Hessian matrix, etc. of the black-box can be evaluated via finite differences in parameter space.

Galaxy Clustering Additional Probes pipeline: diagnostics of the linearised black-box data model



- Using only here the (final, non-linearly evolved) power as summary statistics.
- Any AP can go in the data vector, since we need the simulations anyway!



Florent Leclercq

Latent function inference: the SELFI approach (Simulator Expansion for Likelihood-Free Inference)



- Linearisation of the black-box data model: $\hat{\Phi}_{\theta} \approx \mathbf{f}_0 + \nabla \mathbf{f}_0 \cdot (\mathbf{\theta} - \mathbf{\theta}_0)$
- Further assume:

Florent Leclercg

- Gaussian prior: $\mathcal{P}(\boldsymbol{\theta}) = \mathcal{G}(\boldsymbol{\theta}_0, \mathbf{S})$
- Gaussian effective likelihood: $\mathcal{P}(\mathbf{\Phi}|\mathbf{\theta}) = \mathcal{G}[\mathbf{f}(\mathbf{\theta}), \mathbf{C}_0]$

• The posterior is Gaussian and analogous to a Wiener filter:

 $\begin{array}{ll} \mbox{expansion point} & \mbox{observed summaries} \\ \mbox{mean:} \ensuremath{\boldsymbol{\gamma}} \equiv \ensuremath{\boldsymbol{\theta}}_0 + \ensuremath{\boldsymbol{\Gamma}} (\nabla \mathbf{f}_0)^\intercal \ensuremath{\mathbf{C}}_0^{-1} (\ensuremath{\boldsymbol{\Phi}}_O - \ensuremath{\mathbf{f}}_0) \\ \mbox{covariance:} \ensuremath{\boldsymbol{\Gamma}} \equiv \left[(\nabla \mathbf{f}_0)^\intercal \ensuremath{\mathbf{C}}_0^{-1} \nabla \mathbf{f}_0 + \ensuremath{\mathbf{S}}_{-1}^{-1} \right]^{-1} \\ \mbox{covariance of summaries} \\ \ensuremath{\text{gradient of the black-box}} \end{array}$

- f_0, C_0 and ∇f_0 can be evaluated through simulations only.
- The number of required simulations is fixed *α priori* (contrary to MCMC).
- The workload is perfectly parallel.



SELFI-1 Euclid forecast (cosmic variance limit)

- $V = (3780 \text{ Mpc}/h)^3$ (volume of the Euclid flagship simulation)
- Gaussian random field data model; 6,060 simulations
- 100 parameters are simultaneously inferred





0.2

0.1

1.1

0.9

0.8

 $P(k)/P_0(k)$

SELFI-1 Euclid vs BOSS

 $\boldsymbol{\theta}_0$ (prior) Numerical data models $\boldsymbol{\gamma}$ (reconstruction) allow using the galaxy $\theta_{\rm gt}$ (ground truth) 1.2 power spectrum as BOSS NGC-0.2 $\leq z < 0.5$ BOSS SGC 0.2 < z < 0.5summary statistics up to at $P_0(k)$ least $k \gtrsim 0.5 h/Mpc$ safely $N_{\rm modes} \propto k^3$: 5 times more modes are used in the P(k)1.0 analysis. θ 0.91.1 $P(k)/P_0(k)$ 0.8 0.9 0.7 10^{-1} $k \, [h/{\rm Mpc}]$ 0.1 0.2 0.30.4 0.50.6 $k \, [h/{
m Mpc}]$ Data points from Beutler et al., 1607.03149 Florent Leclercq SELFI enhanced: robustness to model misspecification and Euclid forecast 28/11/2022 11

1.3

SELFI with model misspecification: inference of latent population functions of a prey-predator model



Florent Leclercq

SELFI enhanced: robustness to model misspecification and Euclid forecast 28/11/2022 12

Check for model misspecification and data compression for SBI

 $\mathcal{P}(\boldsymbol{\omega})$

w

θ

 $\mathcal{P}(\mathbf{\Phi}|\mathbf{\theta})$

 $\mathbf{\Phi}$

ũ

- Qualitatively: the shape of the reconstructed θ is useful as a check for model misspecification (independent theoretical understanding).
- Quantitatively: we can use the Mahalanobis distance between the reconstruction γ and the prior distribution $\mathcal{P}(\boldsymbol{\theta})$:

$$d_{\mathrm{M}}(\boldsymbol{\gamma}, \boldsymbol{\theta}_0 | \mathbf{S}) \equiv \sqrt{(\boldsymbol{\gamma} - \boldsymbol{\theta}_0)^{\mathsf{T}} \mathbf{S}^{-1} (\boldsymbol{\gamma} - \boldsymbol{\theta}_0)}$$

- In the example:
 - $d_{\rm M}(\boldsymbol{\gamma}, \boldsymbol{\theta}_0 | \mathbf{S}) \approx 5.35$ for model A
 - $d_{\rm M}(\boldsymbol{\gamma}, \boldsymbol{\theta}_0 | \mathbf{S}) \approx 12.54$ for model B
 - $\langle d_{\rm M}(\mathcal{T}(\boldsymbol{\omega}), \boldsymbol{\theta}_0 | \mathbf{S}) \rangle \approx 9.43$ in fiducial simulations

• The score function
$$\nabla_{\boldsymbol{\omega}} \hat{\ell}_{\boldsymbol{\omega}0}$$
 is the gradient of the log-likelihood at fiducial point $\boldsymbol{\omega}_0$ in parameter space.

A guasi maximum-likelihood estimator for the parameters is

$$\mathcal{C}(\boldsymbol{\Phi}) = \widetilde{\boldsymbol{\omega}} \equiv \boldsymbol{\omega}_0 + \mathbf{F}_0^{-1} \left[(\nabla_{\boldsymbol{\omega}} \mathbf{f}_0)^\mathsf{T} \underline{\mathbf{C}_0^{-1}} (\boldsymbol{\Phi} - \underline{\mathbf{f}_0}) \right]$$

Fisher matrix:
$$\mathbf{F}_0 = (\nabla_{\boldsymbol{\omega}} \mathbf{f}_0)^{\mathsf{T}} \mathbf{C}_0^{-1} \nabla_{\boldsymbol{\omega}} \mathbf{f}_0$$

Already computed Cheap via finite for SELFI differences Score compression is optimal in the sense that it preserves the Fisher

information content of the data.

 $\nabla_{\boldsymbol{\omega}} \mathbf{f}_0 = |\nabla \mathbf{f}_0| \cdot |\nabla_{\boldsymbol{\omega}} \mathcal{I}$

Alsing & Wandelt, 1712.00012



Florent Leclercg

SELFI enhanced: robustness to model misspecification and Euclid forecast 28/11/2022 13

Simulation-based inference of top-level target parameters



Florent Leclercq

- Any SBI algorithm can be used to obtain the posterior $\mathcal{P}(\boldsymbol{\omega}|\widetilde{\boldsymbol{\omega}}_{\mathrm{O}})$.
- Final inference:
 - does not depend on the assumptions made to check for model misspecification,
 - is unbiased (only more conservative) in case data compression is lossy.
- Non-parametric approaches can use the Fisher-Rao distance between simulated summaries $\tilde{\omega}$ and observed summaries $\tilde{\omega}_{O}$:

$$d_{\rm FR}(\widetilde{\boldsymbol{\omega}},\widetilde{\boldsymbol{\omega}}_{\rm O}) \equiv \sqrt{(\widetilde{\boldsymbol{\omega}}-\widetilde{\boldsymbol{\omega}}_{\rm O})^{\mathsf{T}} \mathbf{F}_0(\widetilde{\boldsymbol{\omega}}-\widetilde{\boldsymbol{\omega}}_{\rm O})}$$

Prey-predator model: inference of target population parameters using Likelihood-Free Rejection Sampling



Conclusion: the statistical framework is in place for the GC:AP pipeline

- A novel <u>two-step simulation based Bayesian approach</u>, combining SELFI and SBI, to tackle the issue of model misspecification for a large class of BHMs.
- Advantages of the first step (SELFI):
 - Even if the inference is in high dimension, the simulator remains a black-box.
 - The number of simulations is fixed *a priori* by the user.
 - The computational workload is perfectly parallel.
 - The linearised data model is trained once and for all independently of the data vector (amortisation).
- Advantages of the second step (SBI):
 - SELFI quantities provide a score compressor for free.
 - General advantages of SBI with respect to likelihood-based methods are preserved.
 - Inference does not depend on the assumptions made to check for model misspecification.
- A computationally efficient and easily applicable framework to perform <u>SBI of BHMs while</u> <u>checking for model misspecification</u>.

pySELFI is publicly available at <u>https://pyselfi.florent-leclercq.eu</u>.



Additional slides



28/11/2022 17

A prey-predator model with observational effects



28/11/2022 18

A prey-predator model with observational effects



A family of priors for population functions in prey-predator systems

Assumptions:

- 1. The population functions θ are Gaussiandistributed.
- 2. They are strongly constrained to live close to $\theta_0 = T(\boldsymbol{\omega}_0)$.
- 3. x(t) and y(t) are smooth functions of time.
- 4. The uncertainty on x(t) and y(t) grows with time.
- → Gaussian prior: Overall prior uncertainty mean: θ_0 covariance: $\mathbf{S} \equiv \alpha_{norm}^2 \mathbf{K} \circ \mathbf{V}$



• The 3 free hyperparameters $\{\alpha_{norm}, t_{smooth}, t_{chaos}\}$ can be optimised using simulations.

Smoothness of the population function Chaotic behaviour of the system

$$\mathbf{K}_{z}_{ij} \equiv \begin{bmatrix} -\frac{1}{2} \begin{pmatrix} t_{i} - t_{j} \\ t_{\text{smooth}} \end{pmatrix}^{2} \end{bmatrix} \qquad \mathbf{K} \equiv \begin{pmatrix} \mathbf{K}_{x} & 0 \\ 0 & \mathbf{K}_{y} \end{pmatrix} \qquad \mathbf{V} \equiv \begin{pmatrix} x_{0} \mathbf{u} \mathbf{u}^{\mathsf{T}} & 0 \\ 0 & y_{0} \mathbf{u} \mathbf{u}^{\mathsf{T}} \end{pmatrix} \qquad (\mathbf{u})_{i} \equiv 1 + \frac{t_{i}}{t_{\text{chaos}}}$$

