Likelihood-free inference from galaxy surveys Prospects for Euclid

Florent Leclercq

www.florent-leclercq.eu

Imperial Centre for Inference and Cosmology Imperial College London

> Wolfgang Enzi, Alan Heavens, Jens Jasche, Guilhem Lavaux,

and the Aquila Consortium www.aquila-consortium.org

June 27th, 2019

Imperial College London

Imperial Centre for Inference & Cosmology

Vocabulary considerations I: *What is the likelihood?*



In cosmology, the (true?) likelihood should live at the level of the map of the CMB or LSS. e.g. Wiener filtering for the CMB, BORG for the LSS (a 256³-dimensional Poisson likelihood):



Jasche & Lavaux 2019, 1806.11117 - FL, Lavaux & Jasche, in prep.

Expert knowledge of the likelihood is needed to beat the curse of dimensionality: conditionals/gradients of the likelihood are required by the samplers (Gibbs/Hamiltonian).

Florent Leclercq

Vocabulary considerations II:

You may already be an LFI specialist!



Florent Leclercq

- Likelihood-free inference (LFI) techniques bypass the need for a map-level likelihood, by relying instead only on a "black-box".
- The likelihood is replaced by a measure of the distance/discrepancy △ between simulated and observed statistical summaries of the data.
- e.g. d = full galaxy survey data
 - $\Phi = {\widehat{P}(k)}$ estimated power spectrum
 - $\Delta =$ Mahalanobis distance with covariance matrix Σ

$$\Delta(\mathbf{\Phi}_{\mathbf{\theta}}, \mathbf{\Phi}_{\mathrm{O}}) = \sqrt{\sum_{k,k'} \left[\widehat{P}_{\mathbf{\theta}}(k) - \widehat{P}_{\mathrm{O}}(k) \right]^{\mathsf{T}} \Sigma_{k,k'}^{-1} \left[\widehat{P}_{\mathbf{\theta}}(k') - \widehat{P}_{\mathrm{O}}(k') \right]}$$

Note that this is what many people would call... (square root of -2 times) the log-likelihood!

 What is "primordial" depends mostly on your ambition...

Likelihood-free rejection sampling (LFRS)

- Iterate many times:
 - Sample θ from a proposal distribution $q(\theta)$
 - Simulate Φ_{θ} using the black-box
 - Compute the distance $\Delta(\Phi_{\theta}, \Phi_{O})$ between simulated and observed data
 - Retain θ if $\Delta(\Phi_{\theta}, \Phi_{O}) \leq \epsilon$, otherwise reject

 ϵ can be adaptively reduced (Population Monte Carlo)



Beyond LFRS: two scenarios

The "number of simulations" route:

- Specific cosmological models ($d \lesssim 10$), general exploration of parameter space
- Density Estimation for Likelihood-Free Inference (DELFI)

Papamakarios & Murray 2016, 1605.06376 Alsing, Feeney & Wandelt 2018, 1801.01497 Alsing, Charnock, Feeney & Wandelt 2019, 1903.00007

 Bayesian Optimisation for Likelihood-Free Inference (BOLFI)

Gutmann & Corander 2016, 1501.03291 FL 2018, 1805.07152

The "number of parameters" route:

- Model-independent theoretical parametrisation (d ≥ 100), strong existing constraints in parameter space
- Simulator Expansion for Likelihood-Free Inference (SELFI)

FL, Enzi, Jasche & Heavens 2019, 1902.10149

I thought of the name <u>after</u> developing the method!

Likelihood-free inference from galaxy surveys

The "number of simulations" route: BOLFI

Bayesian Optimisation for Likelihood-Free Inference

Florent Leclercq



The optimal acquisition function for ABC can be written down:

the Expected Integrated Variance (ExpIntVar).

Järvenpää et al. 2017, 1704.00520 (expression of ExpIntVar in the non-parametric approach) FL 2018, 1805.07152 (expression of ExpIntVar in the parametric approach)

Florent Leclercq

BOLFI: Re-analysis of the JLA supernova sample

Betoule et al. 2014, 1401.4064



- The number of required simulations is reduced by:
 - 2 orders of magnitude with respect to likelihood-free rejection sampling (for a much better approximation of the posterior)
 - 3 orders of magnitude with respect to exact Markov Chain Monte Carlo sampling

FL 2018, 1805.07152

 Cheap numerical data models (to be evaluated O(10³) times) will be required for a Euclid-like survey.

Florent Leclercq

The "number of parameters" route: SELFI

Simulator Expansion for Likelihood-Free Inference

Florent Leclercq

SELFI: Method



FL, Enzi, Jasche & Heavens 2019, 1902.10149 Florent Leclercq

- Gaussian prior + Gaussian effective likelihood
- Linearisation of the black-box around an expansion point + finite differences: h

$$\sim \mathbf{\hat{\Phi}}_{\mathbf{\theta}} \approx \mathbf{f}_0 + \nabla \mathbf{f}_0 \cdot (\mathbf{\theta} - \mathbf{\theta}_0)$$

The posterior is Gaussian and analogous to a Wiener filter:

expansion point observed summaries $\gamma \equiv \theta_0 + \Gamma (\nabla f_0)^{\intercal} C_0^{-1} (\Phi_O - f_0)$ $\Gamma \equiv \left[(\nabla f_0)^{\intercal} C_0^{-1} \nabla f_0 + S_0^{-1} \right]^{-1}$ prior covariance covariance of summaries gradient of the black-box

 f_0, C_0 and ∇f_0 can be evaluated through simulations only. The number of required simulations is fixed a priori.



SELFI + Simbelmynë: Proof-of-concept



Florent Leclercq

SELFI + Simbelmynë: Proof-of-concept



posterior

Robust inference of cosmological parameters can be easily performed a posteriori once the linearised data model is learnt

 pyselfi will be made publicly available soon

FL, Enzi, Jasche & Heavens 2019, 1902.10149

Florent Leclercq

Concluding thoughts

- Goal: developing and using algorithms for targeted questions, allowing the use of simulators including all relevant physical and observational effects.
- Bayesian analyses of galaxy surveys with fully non-linear numerical black-box models is not an impossible task!
- The "number of simulations route" (BOLFI):
 - The optimal acquisition function can be derived: the Expected Integrated Variance.
 - The number of simulations is reduced by several orders of magnitude.
- The "number of parameters route" (SELFI):
 - High-dimensional likelihood-free problems can be addressed.
 - The computational workload is fixed *a priori* and perfectly parallel.