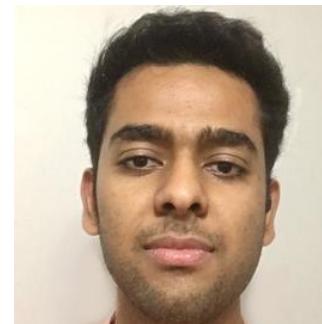
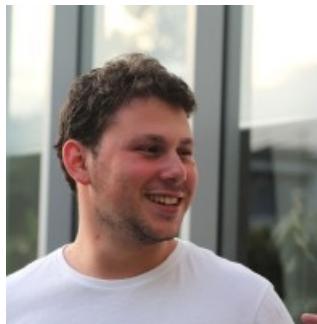


Efficient model selection with Bayesian optimisation

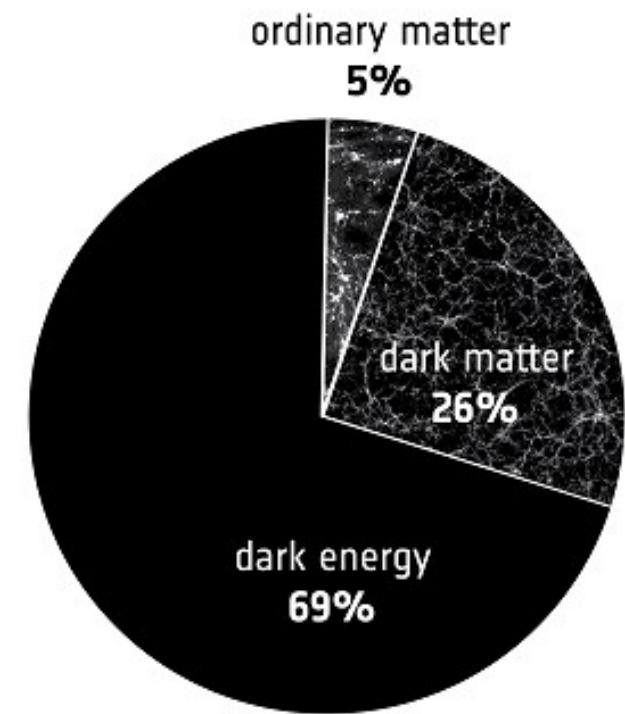
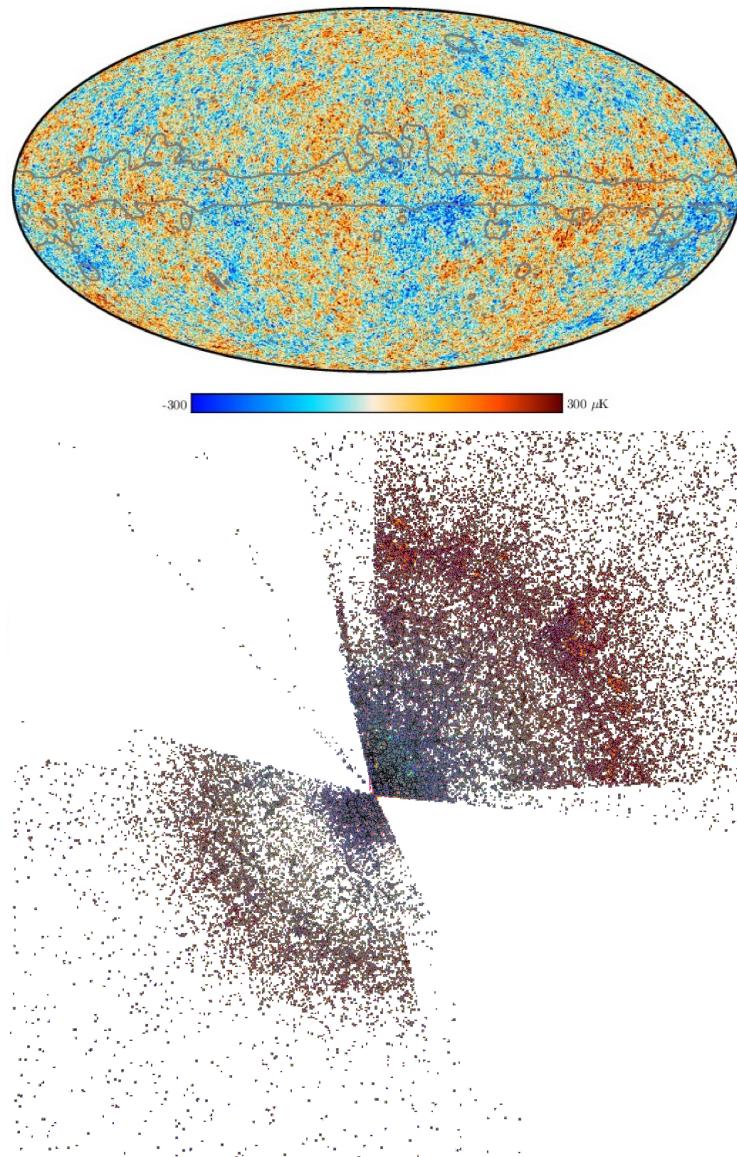
Jan Hamann

based on [JCAP 03 \(2022\) 03, 036 \[arXiv:2112.08571\]](#) with [Julius Wons](#)
and work in progress with [Nathan Cohen](#) and [Ameek Malhotra](#)



UNSW
SYDNEY

GGI Neutrino Frontiers focus week
July 2024



Parameter inference/optimisation

Cosmological model \mathcal{M}

Parameters θ

For instance:

Standard LCDM

$$\theta = (\omega_b, \omega_{\text{cdm}}, H_0, \tau, A_s, n_s)$$

or

LCDM + Neff

$$\theta = (\omega_b, \omega_{\text{cdm}}, H_0, \tau, A_s, n_s, N_{\text{eff}})$$

etc.

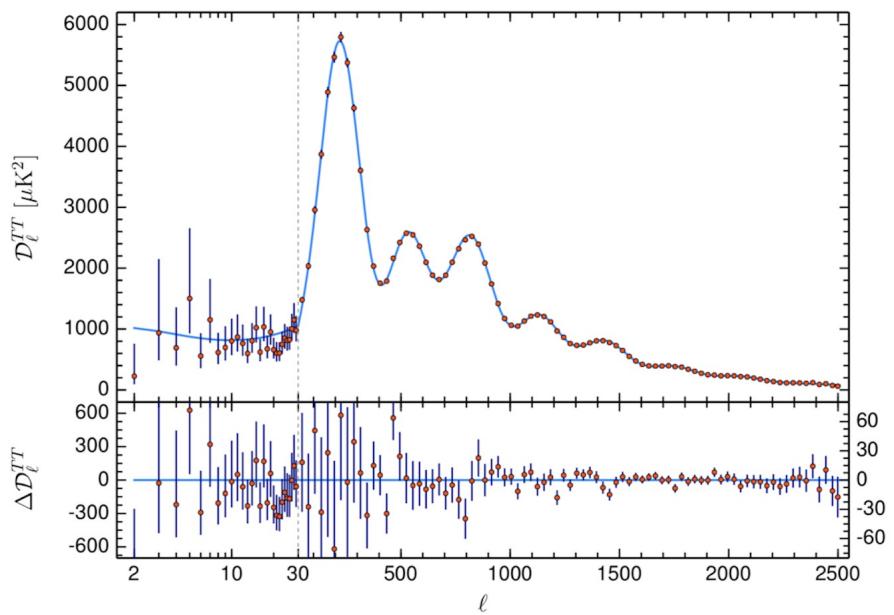
Parameter inference/optimisation

Cosmological model \mathcal{M}
Parameters θ

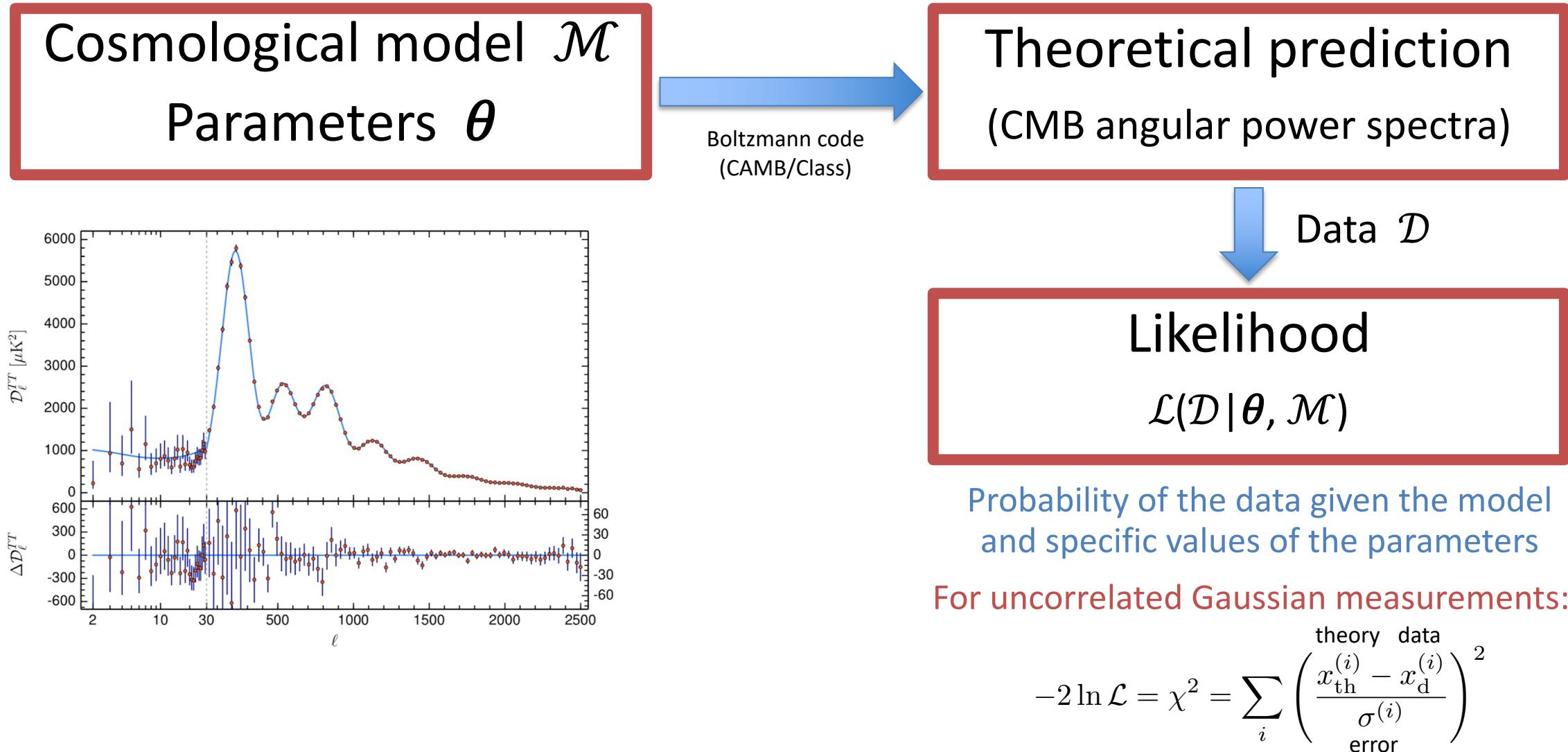


Boltzmann code
(CAMB/Class)

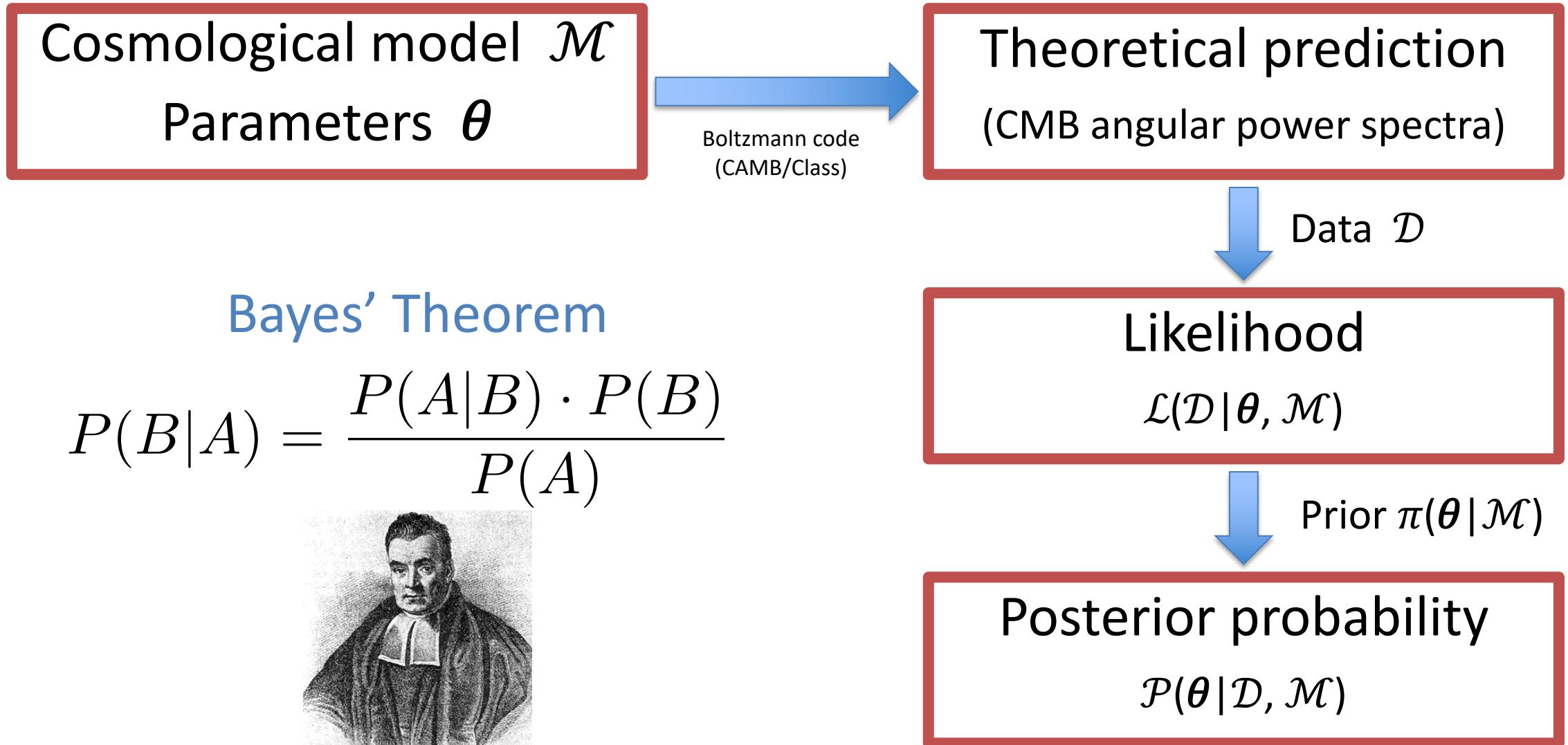
Theoretical prediction
(CMB angular power spectra)



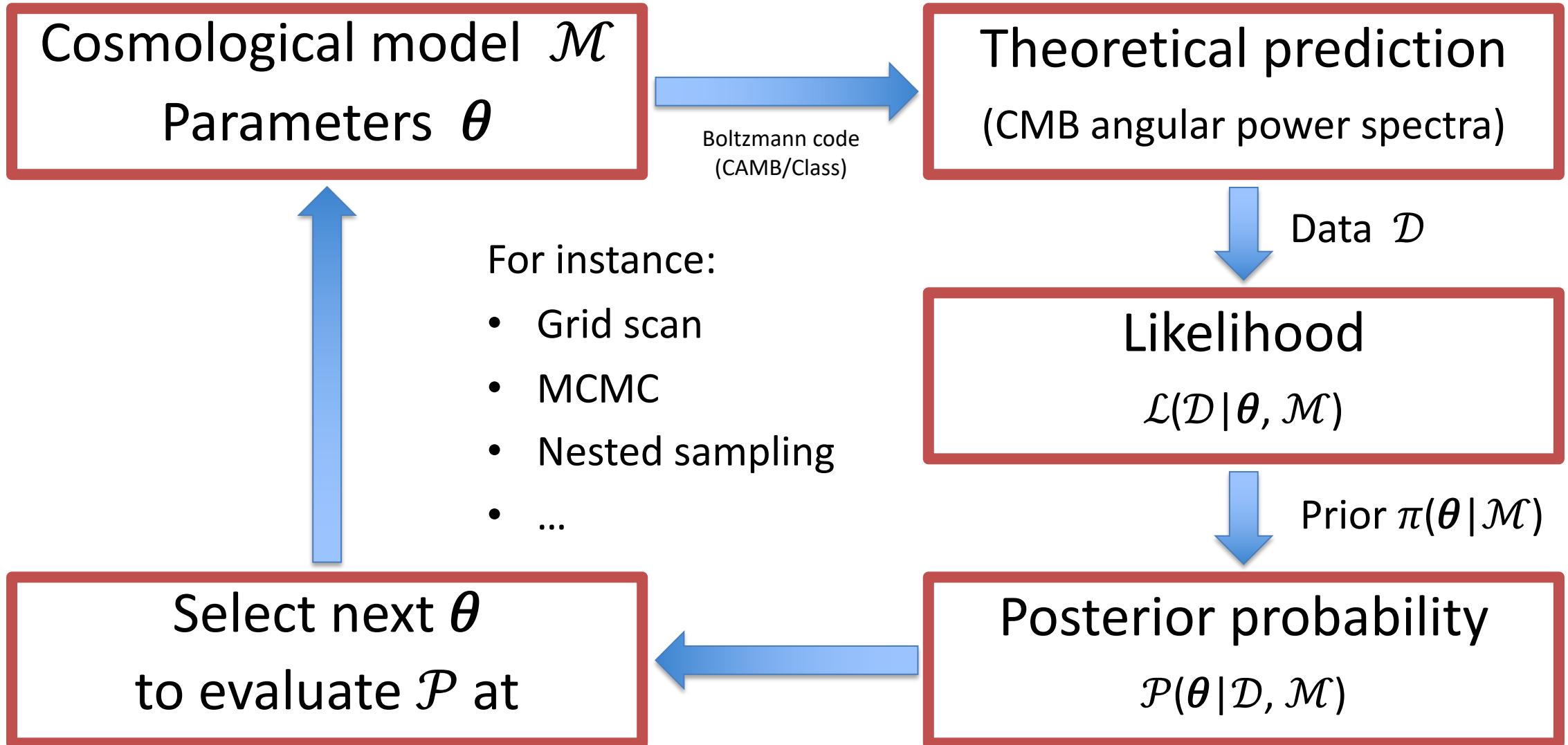
Parameter inference/optimisation



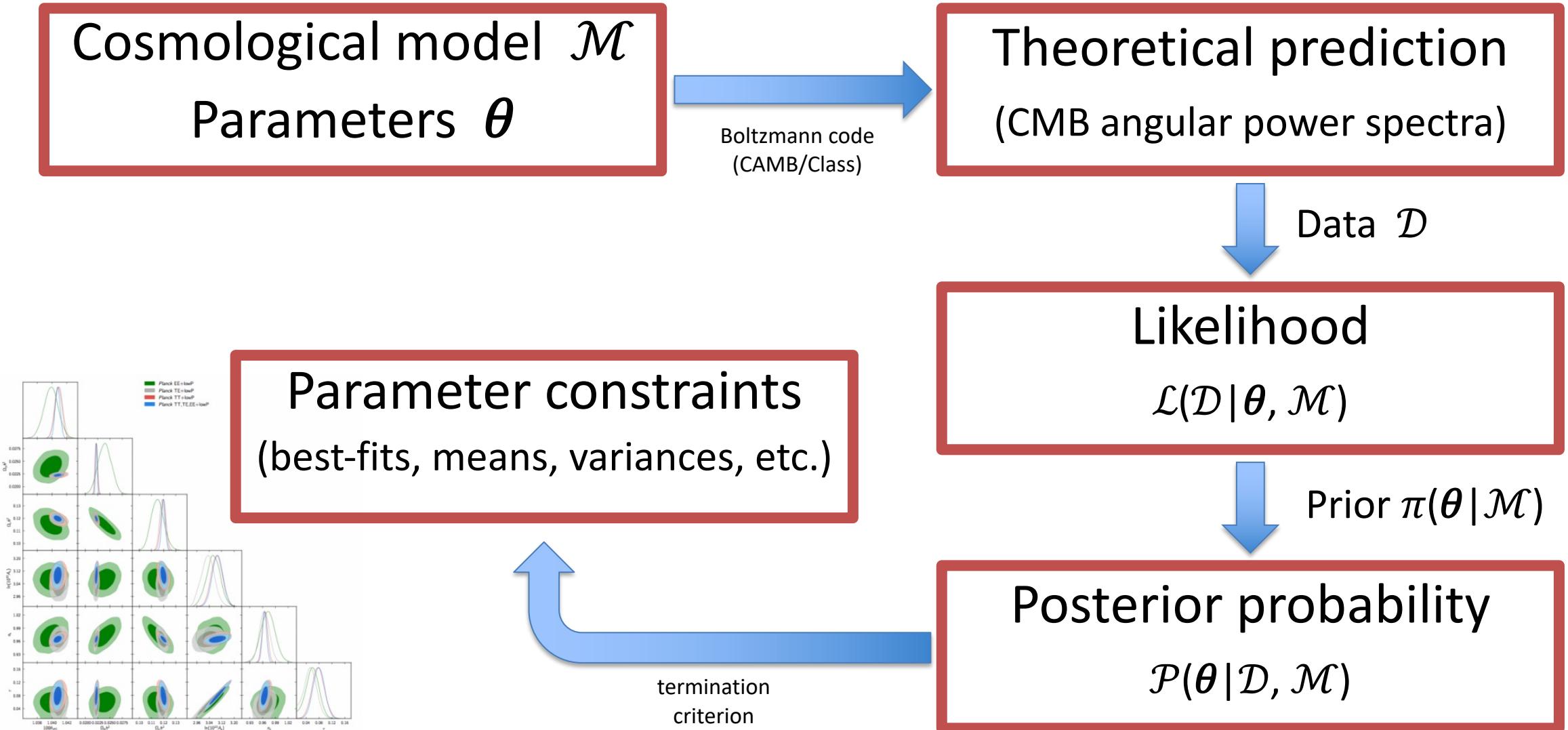
Parameter inference/optimisation



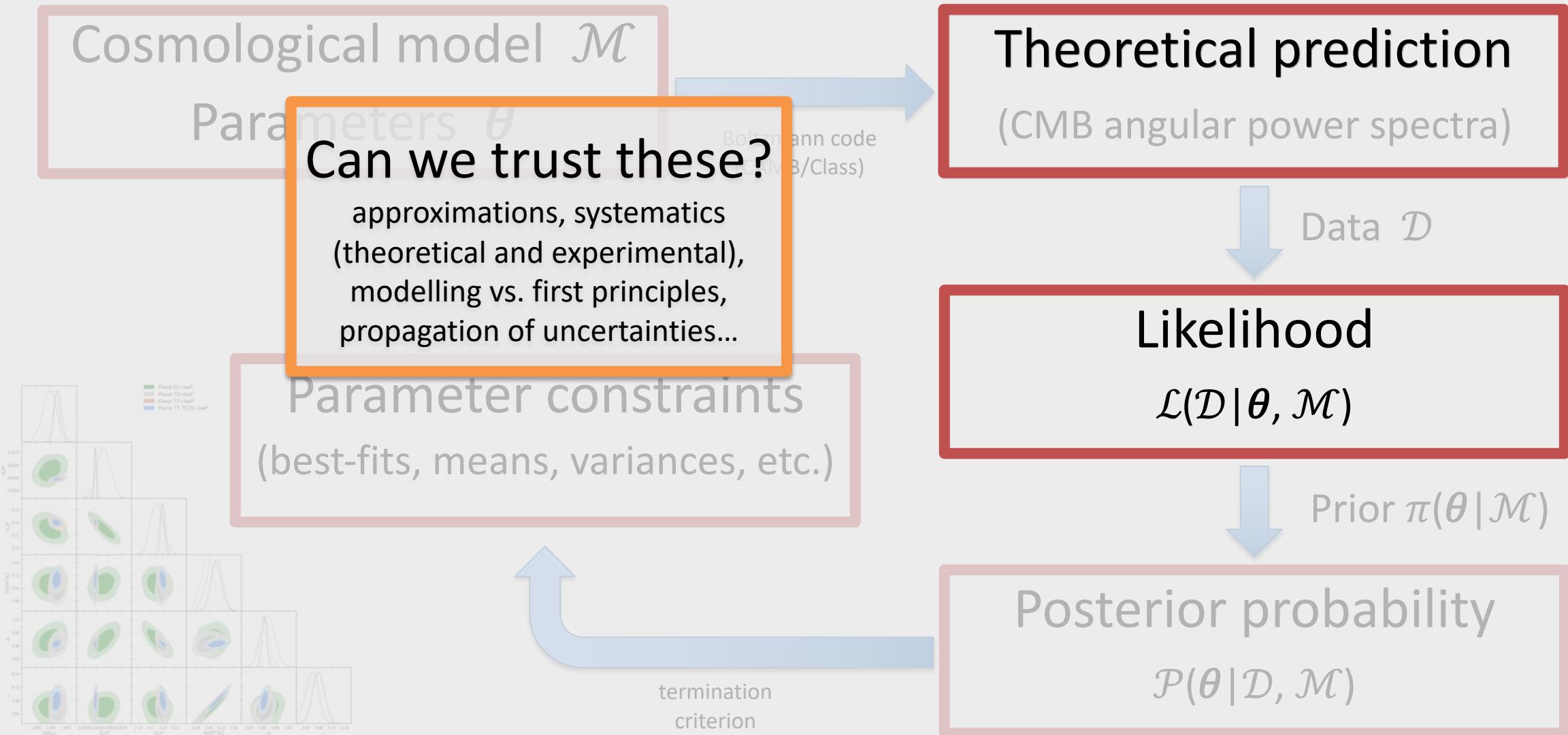
Parameter inference/optimisation



Parameter inference/optimisation

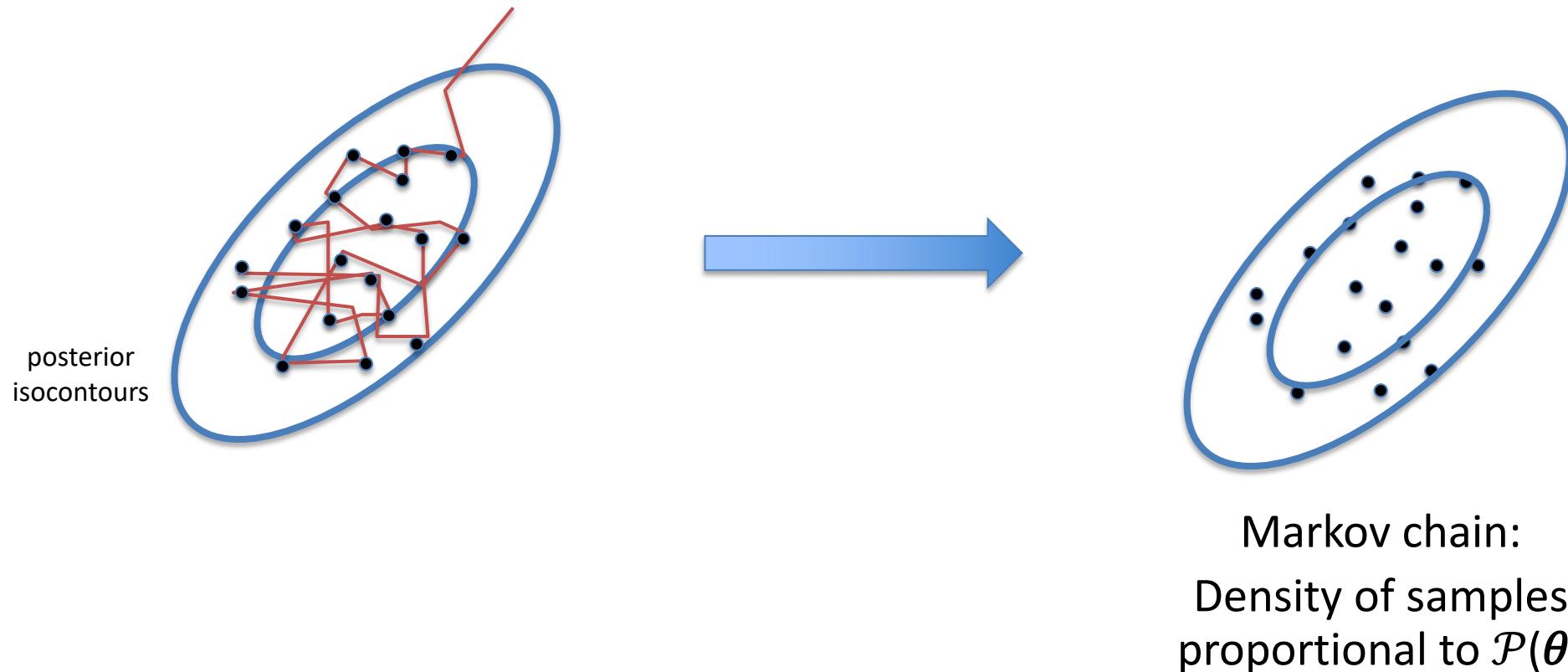


Parameter inference/optimisation



The usual approach: Markov chain Monte Carlo

- Basic idea: random walk in parameter space that explores $\mathcal{P}(\theta)$



The usual approach: Markov chain Monte Carlo

Metropolis-Hastings algorithm:

[Metropolis et al. (1953)]

1. Start at point θ in parameter space
2. Save θ to Markov chain
3. Propose a step to a new point θ'
4. Decide whether to accept the proposal and take the step:
 - If $\mathcal{P}(\theta') \geq \mathcal{P}(\theta)$, accept the proposal
 - If $\mathcal{P}(\theta') < \mathcal{P}(\theta)$, accept the proposal with a probability $p = \mathcal{P}(\theta')/\mathcal{P}(\theta)$, otherwise reject
5. If step was accepted set $\theta' = \theta$
6. Go to 2.

Animated illustration:

<http://chi-feng.github.io/mcmc-demo/app.html?algorithm=RandomWalkMH&target=standard>

[Feng et al., Github]

Pros and cons of MCMC

- + easily implemented
- + easily parallelisable
- + essentially zero overhead
- + mild scaling of number of required samples with dimension N of parameter space (power law $\sim N^\alpha$ rather than exponential)
- + works great for near-Gaussian posteriors (most of cosmology)
 - o not very precise at finding the maximum
 - o typically requires $\mathcal{O}(10^4)$ function evaluations for $N = \mathcal{O}(10)$
- struggles with complicated (multi-modal, non-Gaussian, non-linearly correlated, etc.) posteriors
- not very smart: most of the information is ignored!

Bayesian optimisation

Step 1: Regression

Guess the shape of the function based on known function values (“data”)

Step 2: Selection

Decide at which point to evaluate the next function value

Goal: find global maximum of function
(and learn general shape in the process)

Bayesian optimisation

Step 1: Regression

Gaussian

Process
Regression

Guess the shape of the function based on known function values ("data")

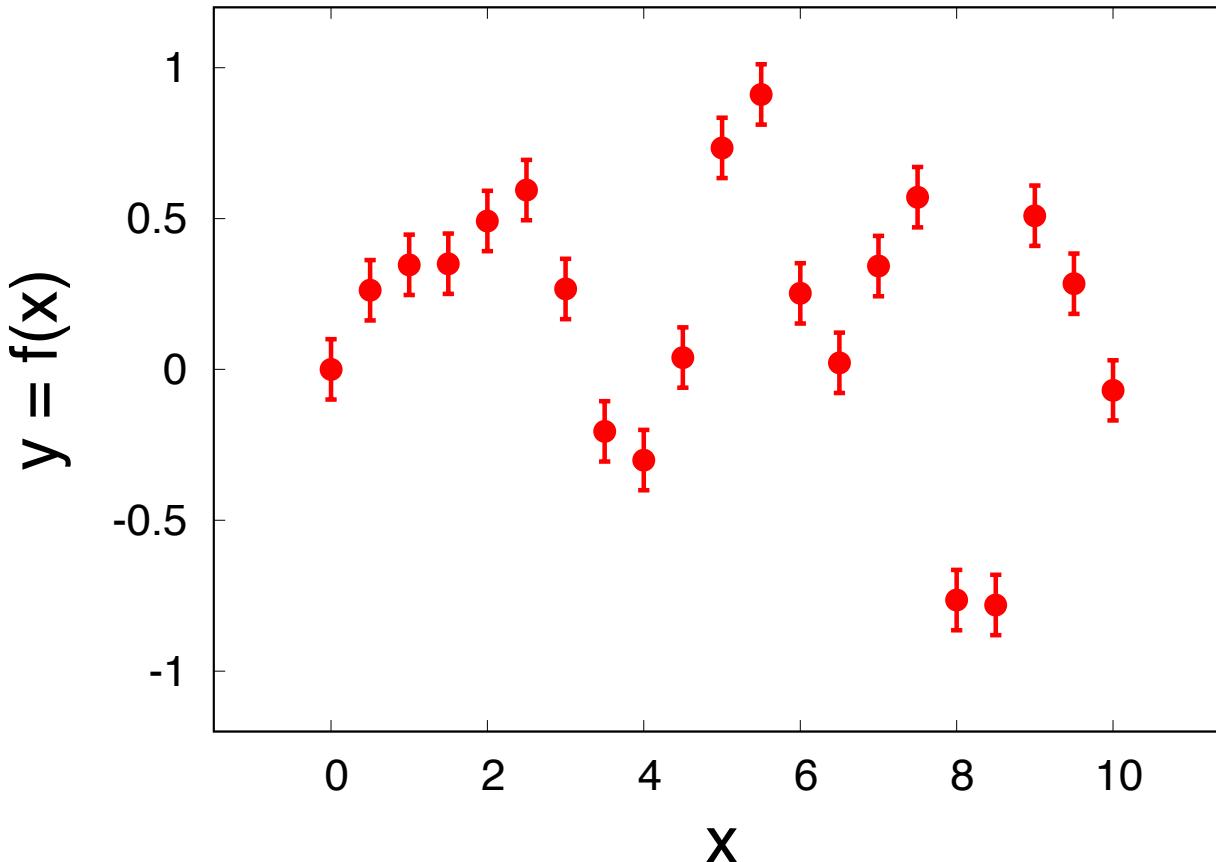
Step 2: Selection

Acquisition function:
Expected Improvement

Decide at which point to evaluate the next function value

Goal: find global maximum of function
(and learn general shape in the process)

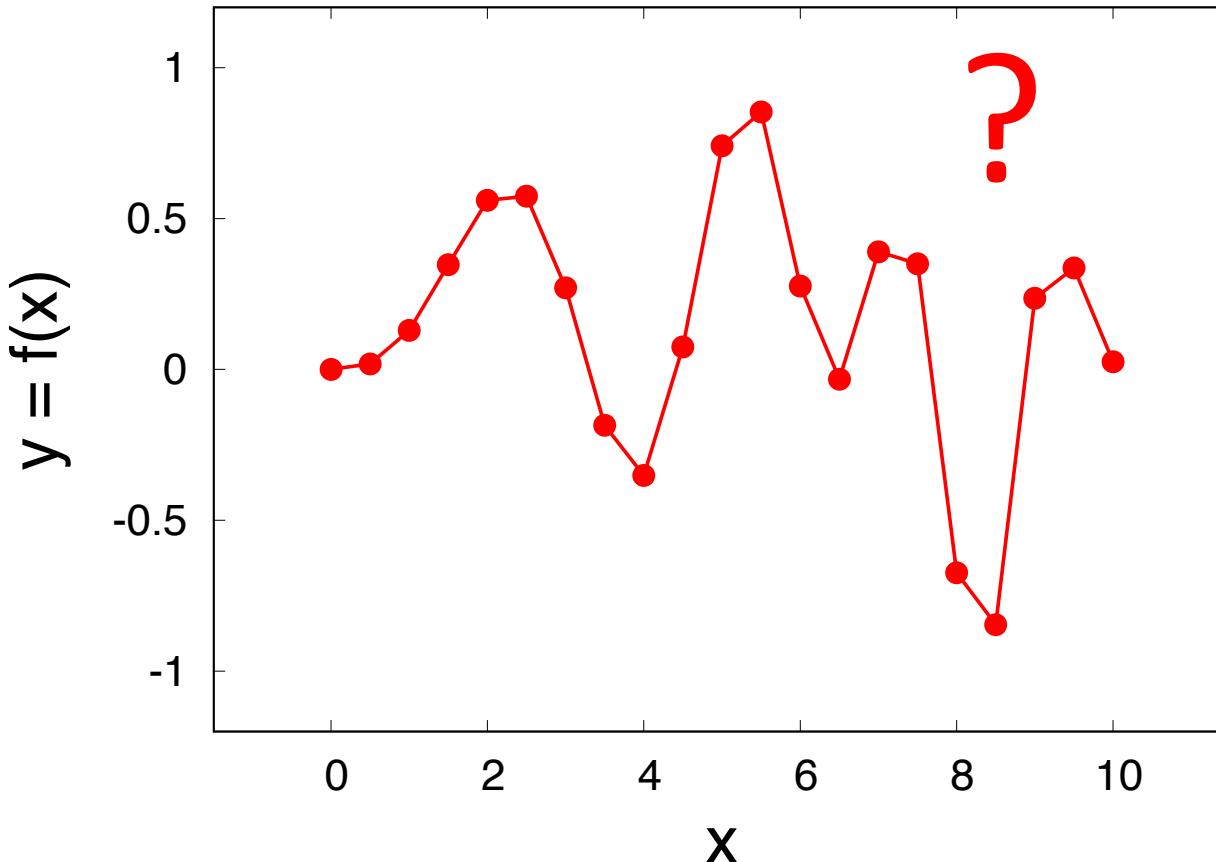
Gaussian Process Regression (GPR)



Data: (x_i, y_i)

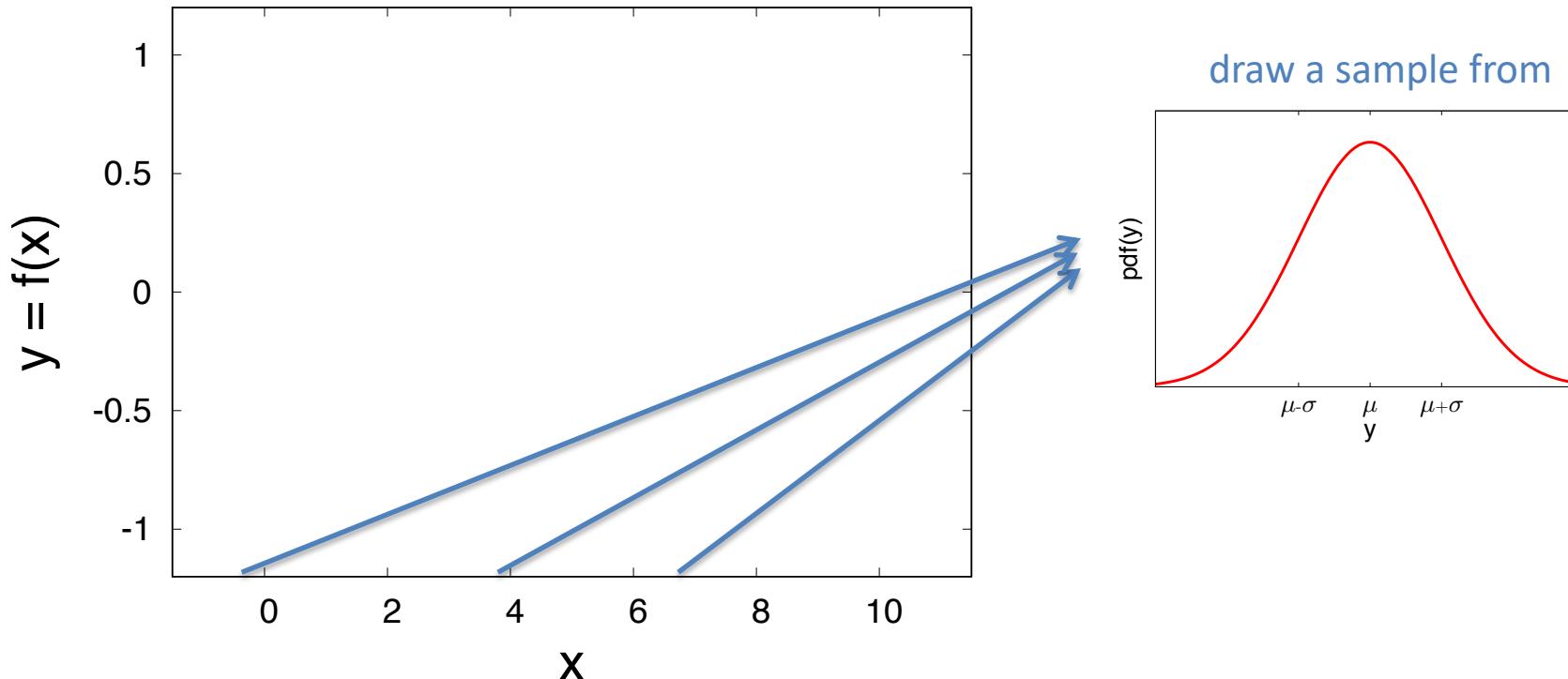
Covariance of the data: Σ_{ij}

Gaussian Process Regression (GPR)



Data: (x_i, y_i)

Gaussian Process



- Imagine for each x , $f(x)$ is a random Gaussian variate drawn from $\mathcal{N}(\mu, \sigma^2)$
- \uparrow \uparrow
mean variance

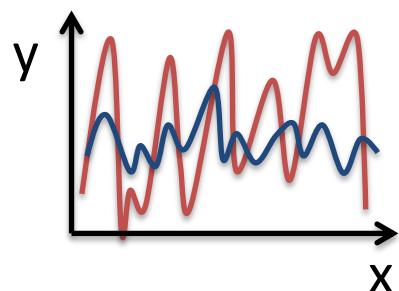
Gaussian Process

- Gaussian process is specified by

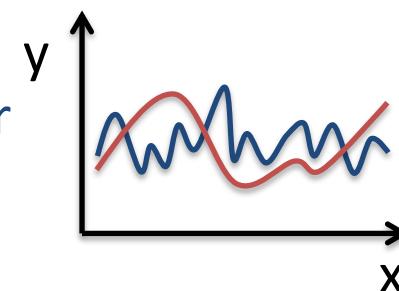
$$\mathcal{N}(\mu(x), K(x, x'))$$


covariance function

- Example covariance function:



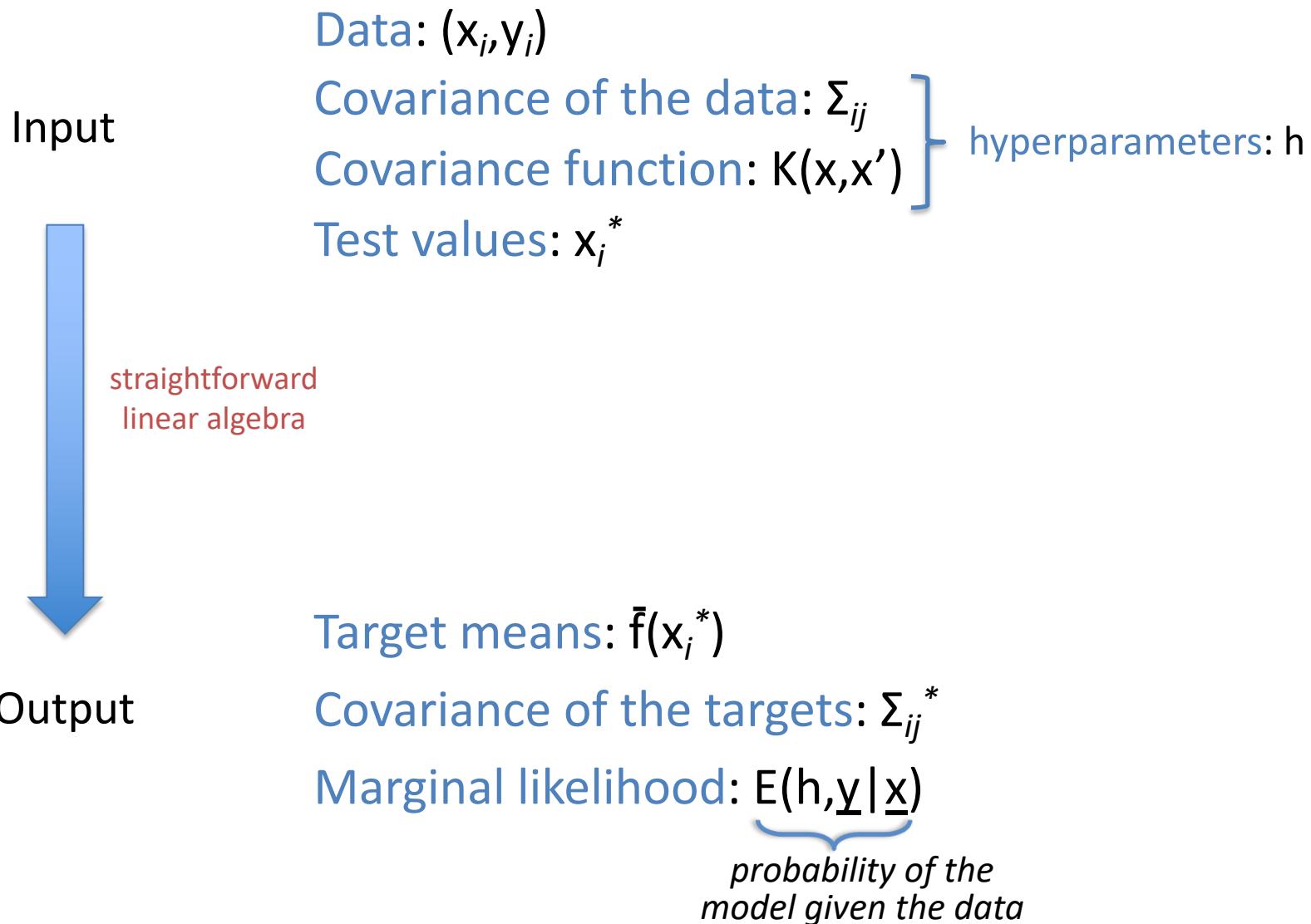
larger/smaller
A and *L*



Gaussian Process Regression (GPR)

- Take function space generated by Gaussian process GP as *prior probability distribution* for f
- Impose data (x_i, y_i) as condition on GP
- $GP | (x_i, y_i)$ is the *posterior probability distribution* for f (this is still an infinite space of functions, but one can now evaluate, e.g., expectation value and covariance)

Gaussian Process Regression



Gaussian Process Regression

Input



straightforward
linear algebra

Output

Data: (x_i, y_i)

Covariance of the data: Σ_{ij}

Covariance function: $K(x_i, x_j)$

Test values: x_i^*

Target means: $f(x_i^*)$

Covariance of the targets: Σ_{ij}

Marginal likelihood: $E(h, y | x)$

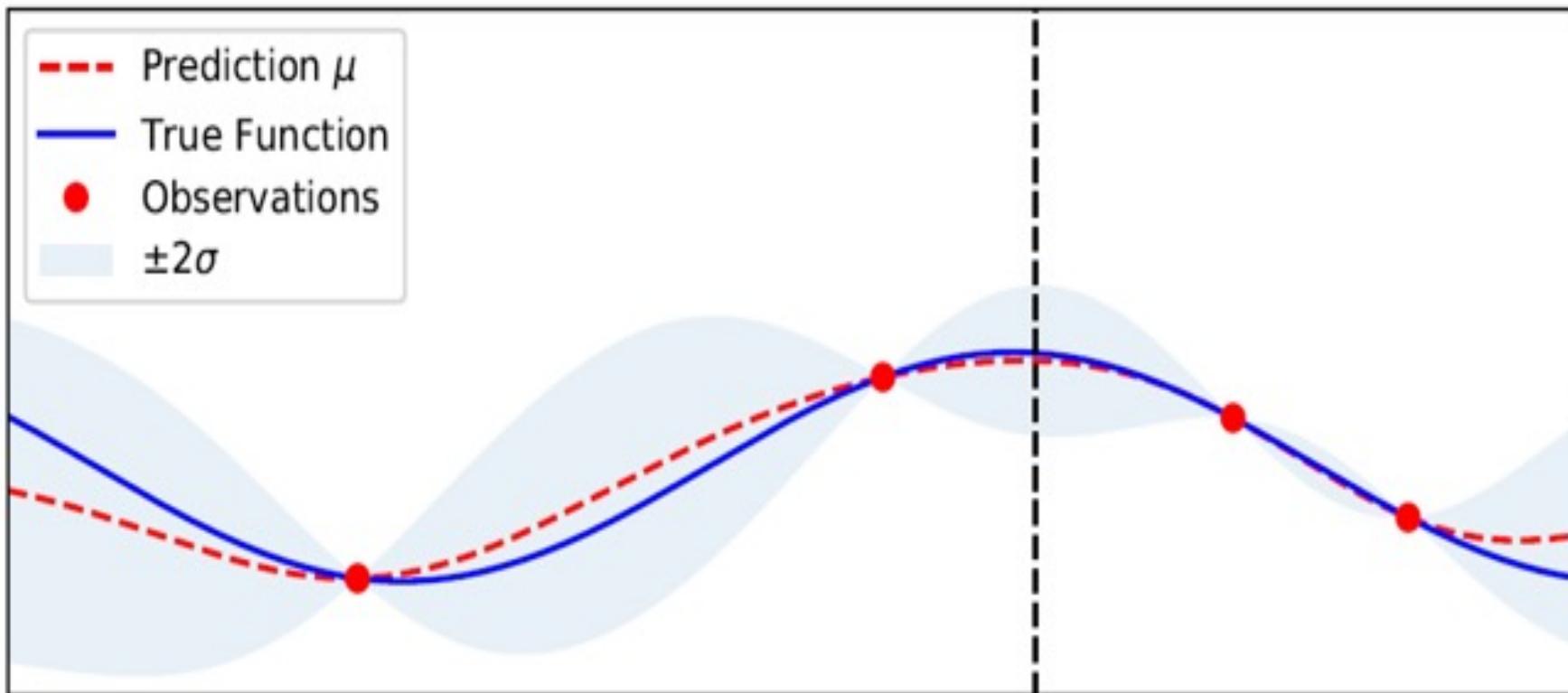
probability of the
model given the data

**Maximise marginal
likelihood as function of
hyperparameters**

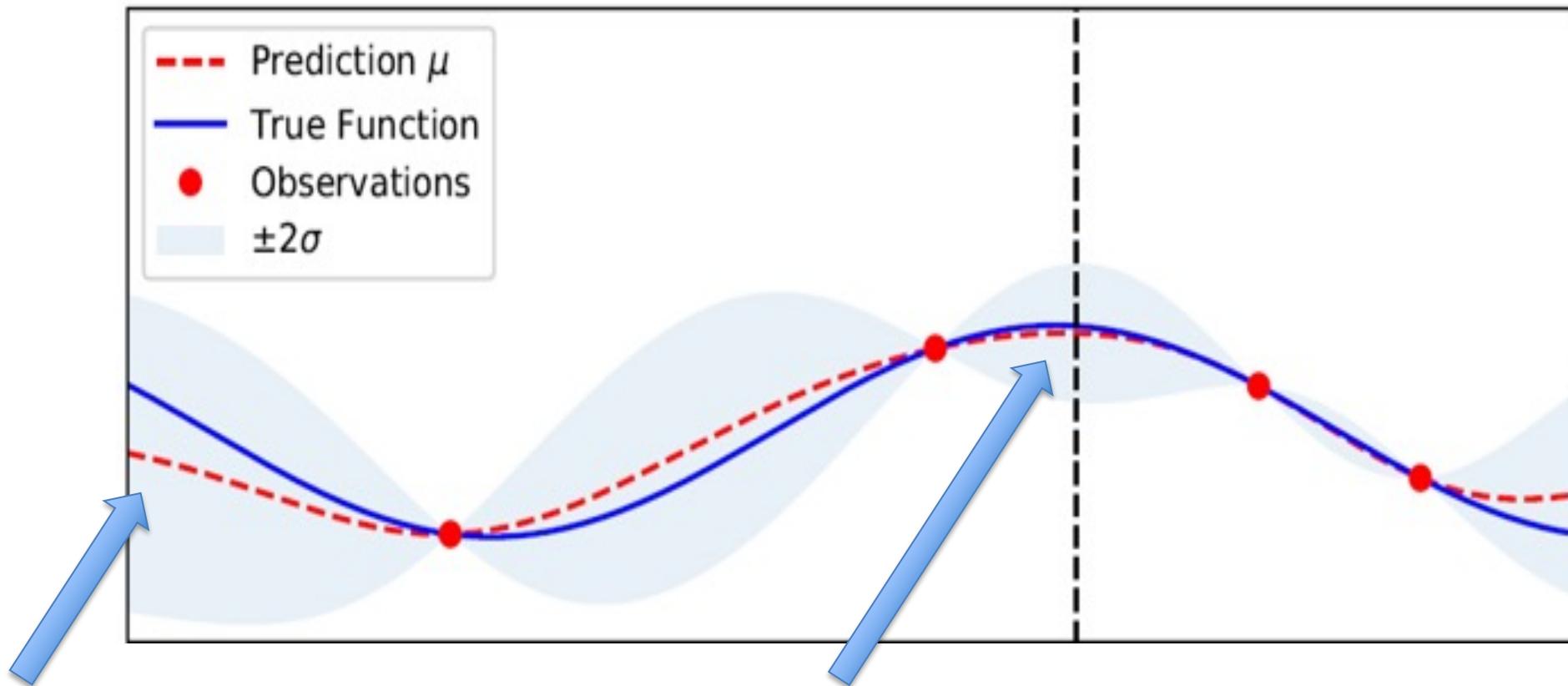
=

*Let the data decide on
the most appropriate
Gaussian process!*

Gaussian Process Regression



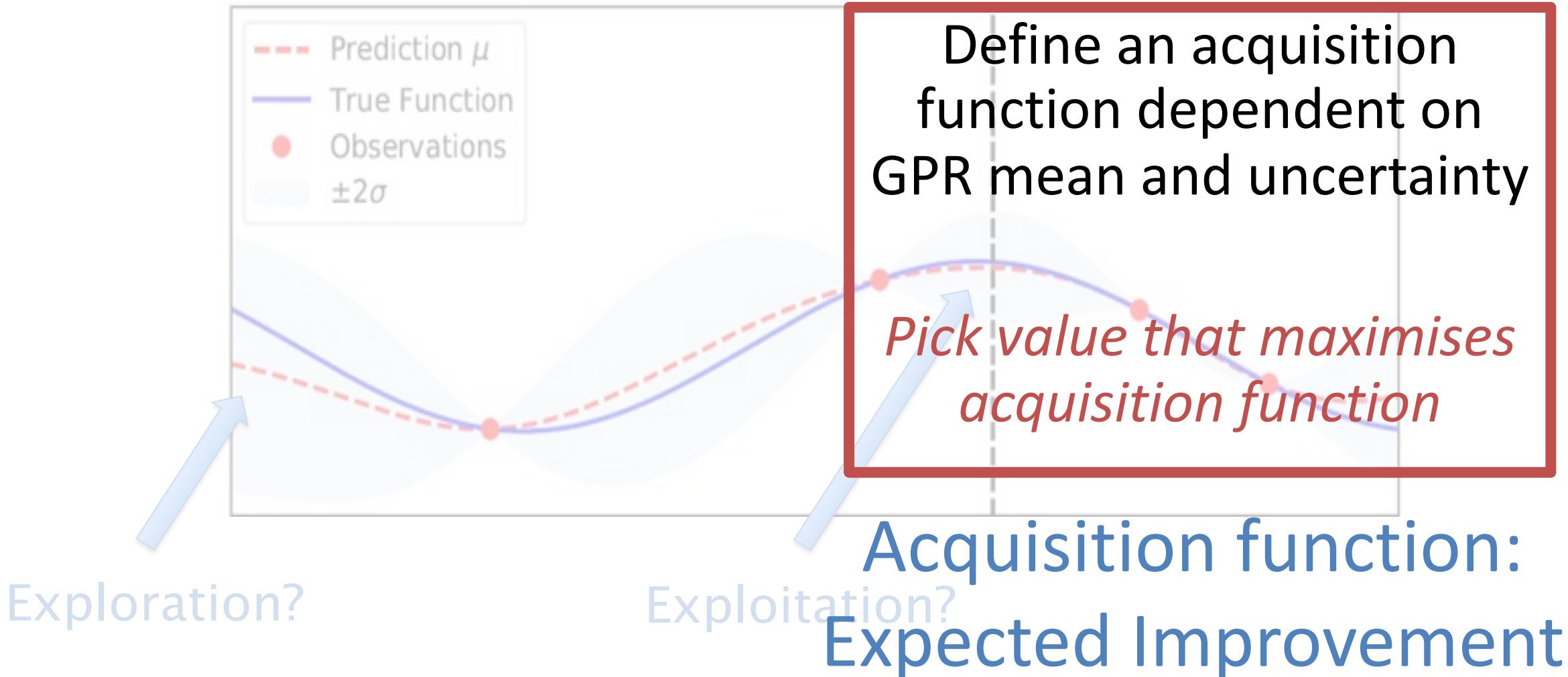
Where to draw the next sample?



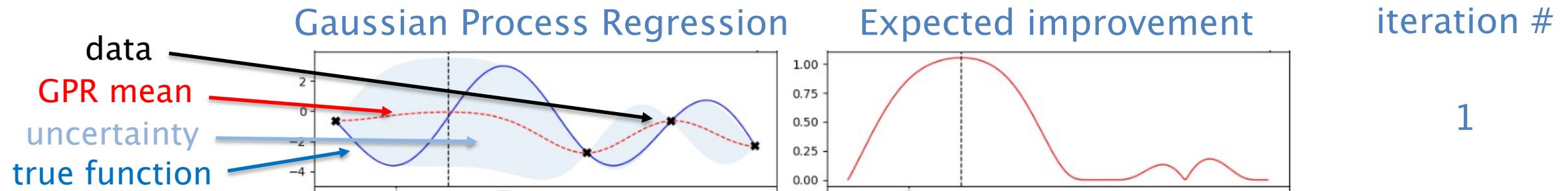
Exploration?

Exploitation?

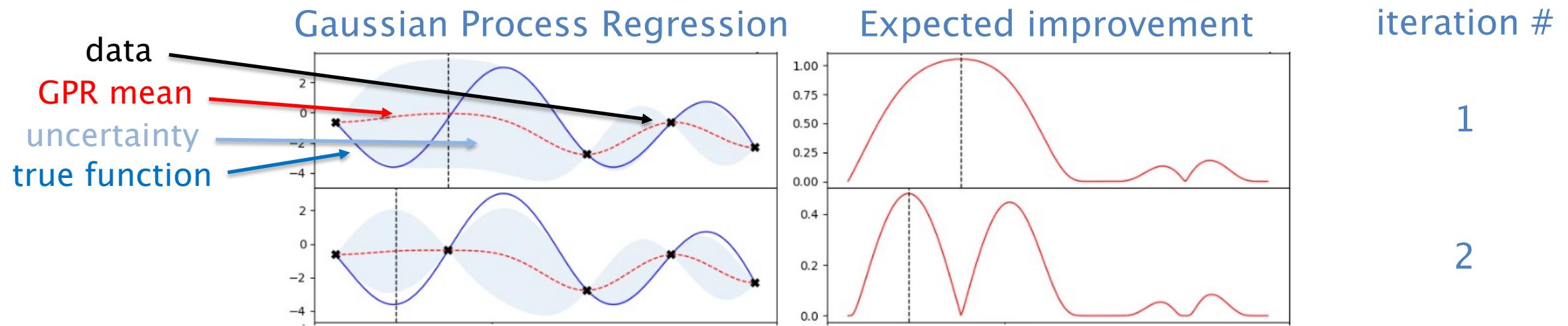
Where to draw the next sample?



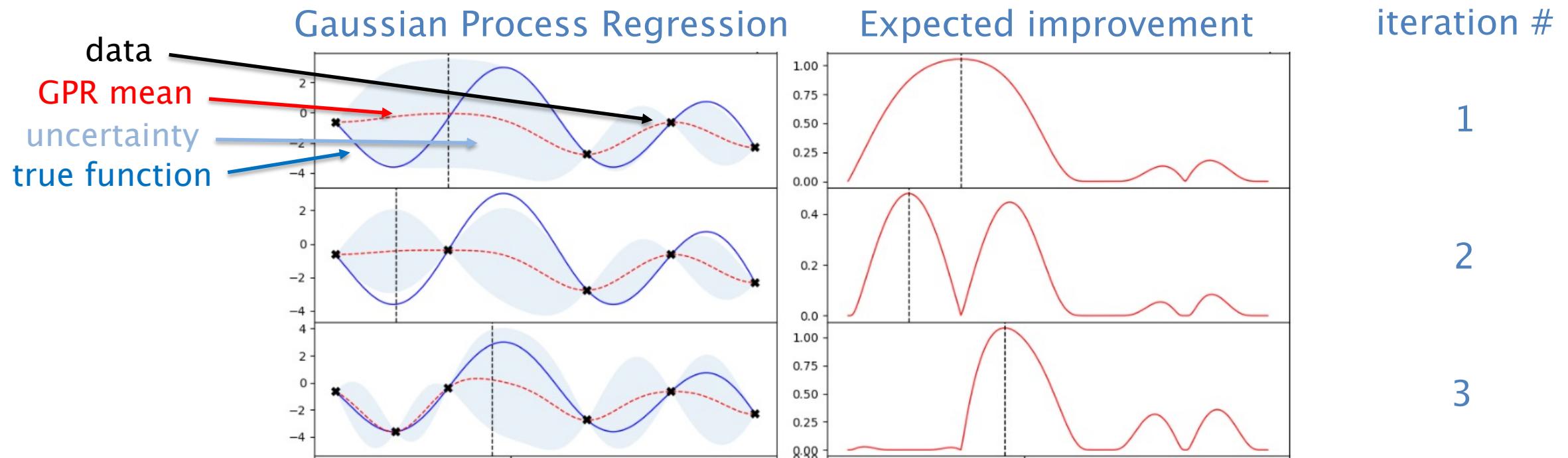
Bayesian optimisation



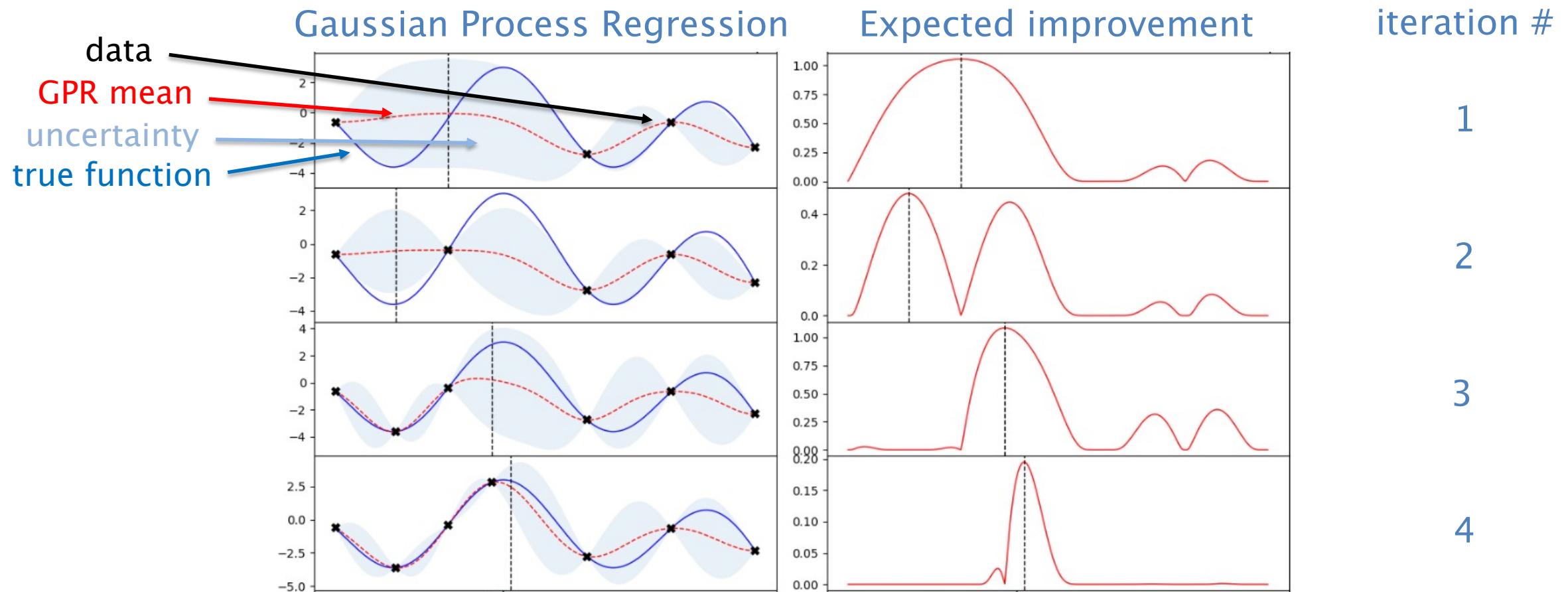
Bayesian optimisation



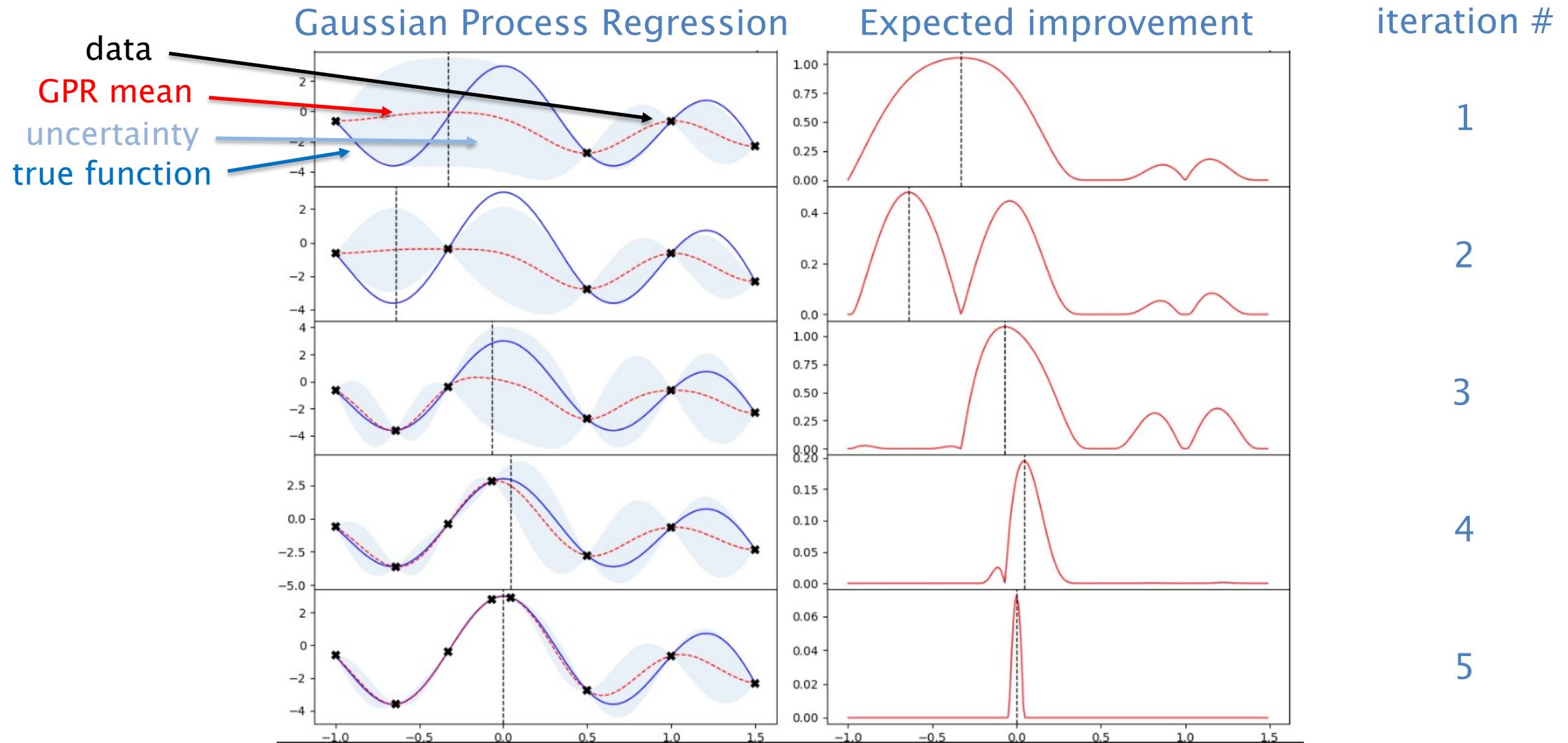
Bayesian optimisation



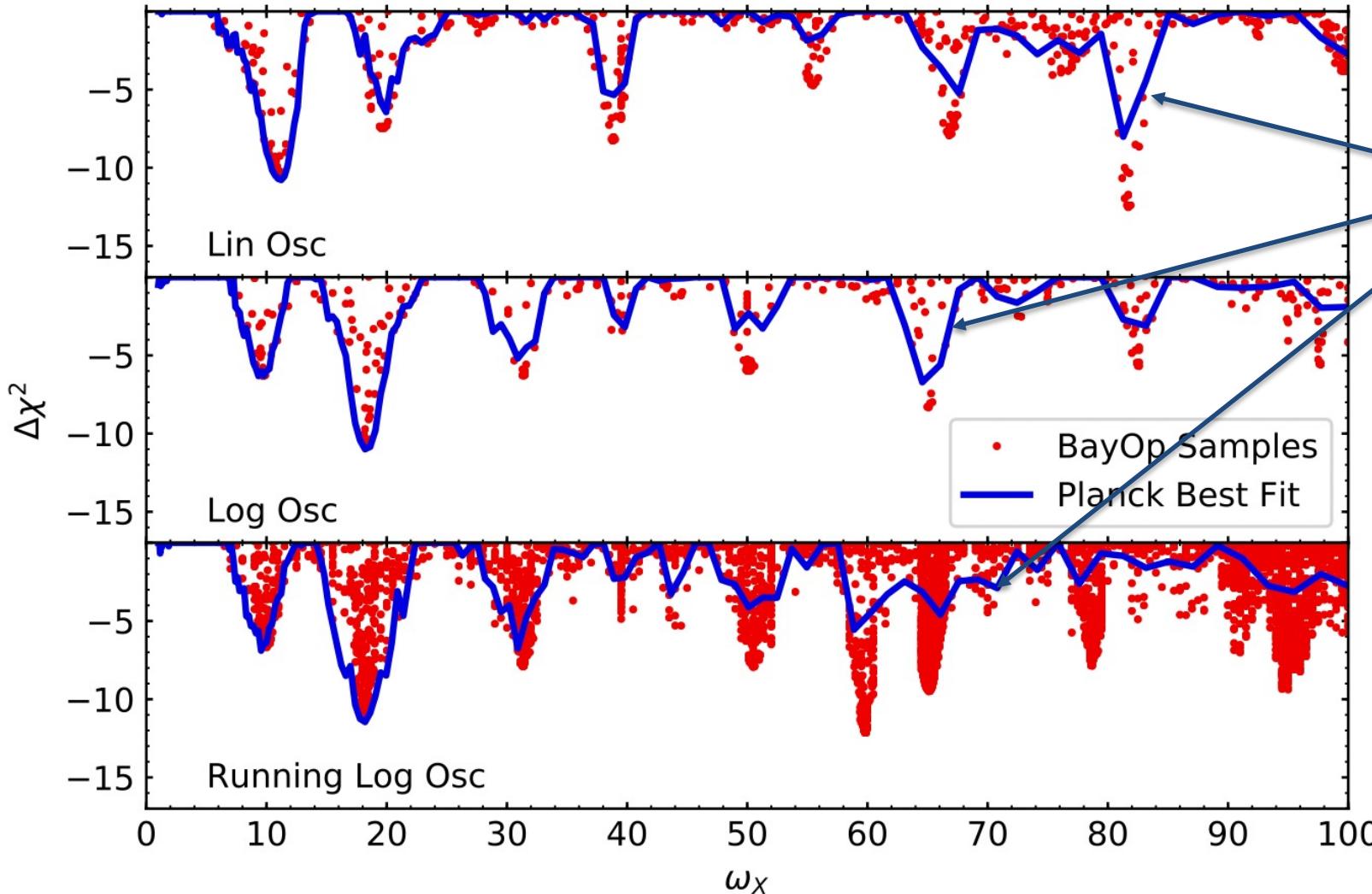
Bayesian optimisation



Bayesian optimisation



An example application: inflation models with modulated primordial power spectra



[*Planck inflation 2018*]

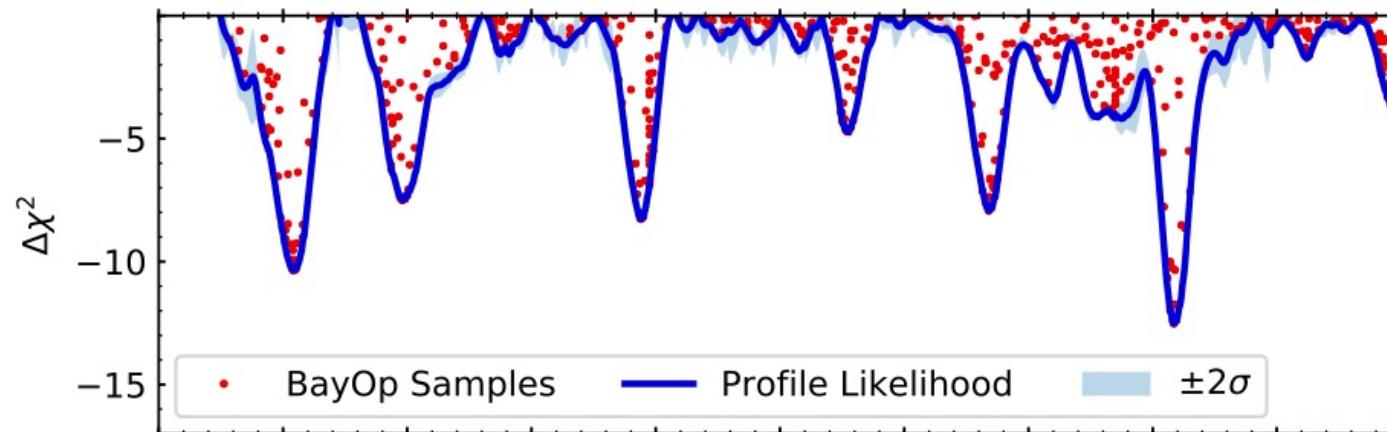
using nested sampling,
 $O(10^5)$ samples

red dots: our results with BO

- Two orders of magnitude fewer function evaluations
- Much better at finding global and local extrema

[JH & Wons, 2021]

BayOp – not only good for optimisation



1700 samples
8 frequency bins

... it also learns the global shape of the function

Pros and cons of Bayesian Optimisation

- + high efficiency
- + excellent at finding global maximum
- + very good at determining overall shape, profiles of functions
- + works even for very nasty (non-Gaussian, multimodal, etc.) functions
- + does not require user input or fine-tuning of settings to work
- may struggle with higher-dimensional problems ($D \gtrsim 10$)
- non-trivial computational overhead (CPU time, memory)

Bayesian optimisation for parameter inference

- Learn shape of posterior probability density
- Replace (potentially expensive) calculation of theoretical prediction and likelihood evaluation with (cheap!) GPR emulation
- Implemented in a Python package: [GPy](#) [El Gammal et al., 2022]

But this assumes we know the right model...

Model selection: Bayesian method

$$P(\mathcal{M}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathcal{M}) \cdot P(\mathcal{M})}{P(\mathcal{D})}$$

Probability of model \mathcal{M}
given the data \mathcal{D}

Bayesian evidence

$$P(\mathcal{D}|\mathcal{M}) = \int d\theta \mathcal{L}(\mathcal{D}|\theta, \mathcal{M}) \pi(\theta|\mathcal{M})$$

Comparing two models:
Bayes factor B_{12}

$$B_{12} = \frac{P(\mathcal{D}|\mathcal{M}_1)}{P(\mathcal{D}|\mathcal{M}_2)}$$

“Model \mathcal{M}_1 is B_{12} times more
probable than \mathcal{M}_2 ”

Model selection: Bayesian method

Bayesian evidence

$$P(\mathcal{D}|\mathcal{M}) = \int d\theta \mathcal{L}(\mathcal{D}|\theta, \mathcal{M}) \pi(\theta|\mathcal{M})$$

- Integral over entire parameter space
- Rewards models that make *risky* predictions and *get it right* over generic models that can *fit anything*
- Natural implementation of Occam's razor:

Numquam ponenda est pluralitas sine necessitate

Plurality must never be posited without necessity



Bayesian model selection

- Multi-dimensional integration is a challenging task
- Standard approach: Nested sampling algorithm

[Skilling 2004, Feroz et al. 2013, Handley et al. 2015]

- typically requires $\mathcal{O}(10^5\text{-}10^6)$ function evaluations for features models

This is even harder than parameter inference
Can Bayesian Optimisation help?

Evidence calculation with Bayesian optimisation

- Goal is to select next function value to be evaluated in such a way that it maximises the expected reduction in uncertainty of the integral
- Use a different acquisition function: Integrated Mean Square Prediction Error (IMSPE)

$$\text{IMSPE}(\theta) = \int d\theta' \sigma_{\widehat{\text{GP}}(\theta)}(\theta')$$

GPR uncertainty

Pretend to take a sample at θ ,
then do a new GPR

Very convenient:
gives estimate of the uncertainty of the evidence integral

Evidence calculation with Bayesian optimisation

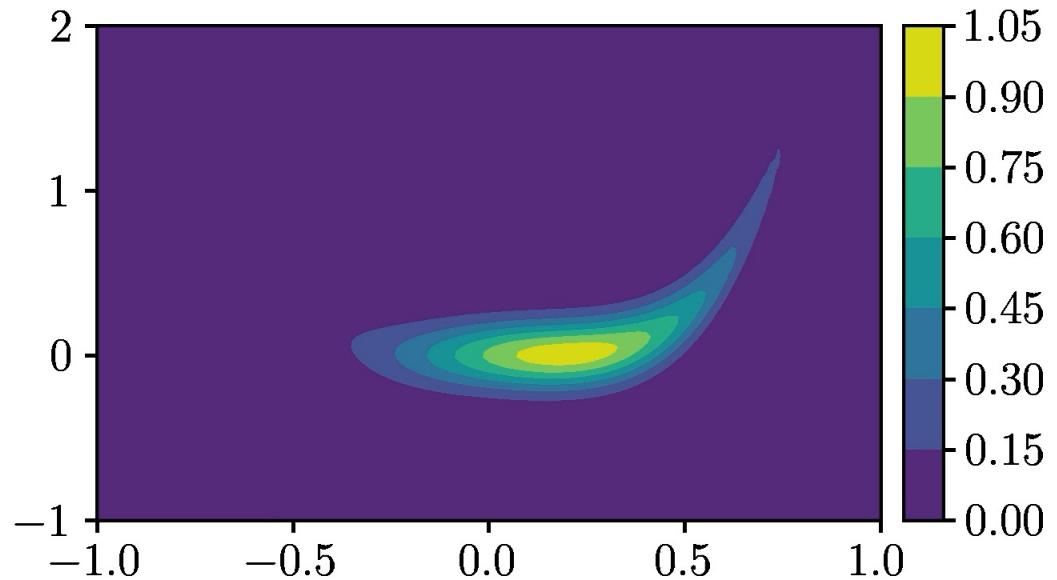
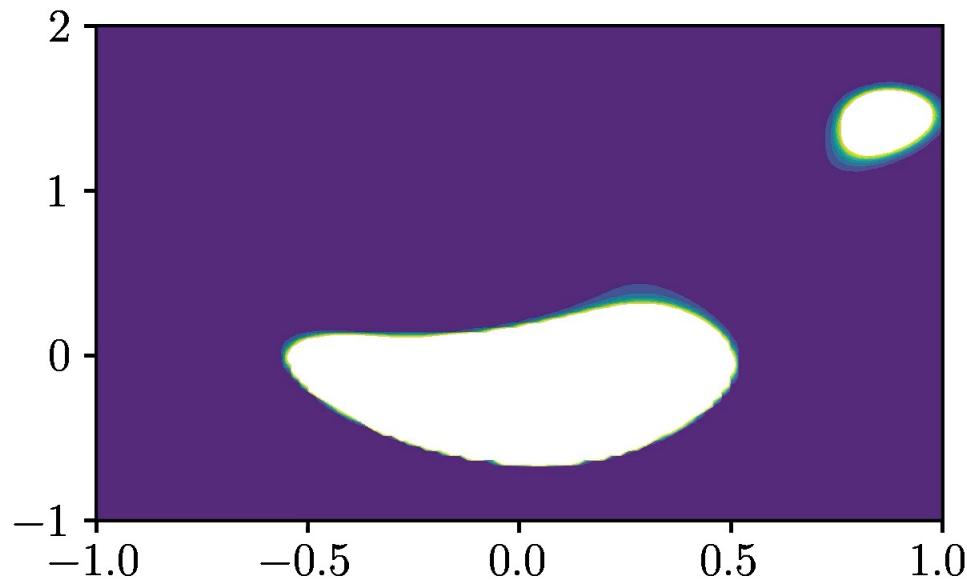
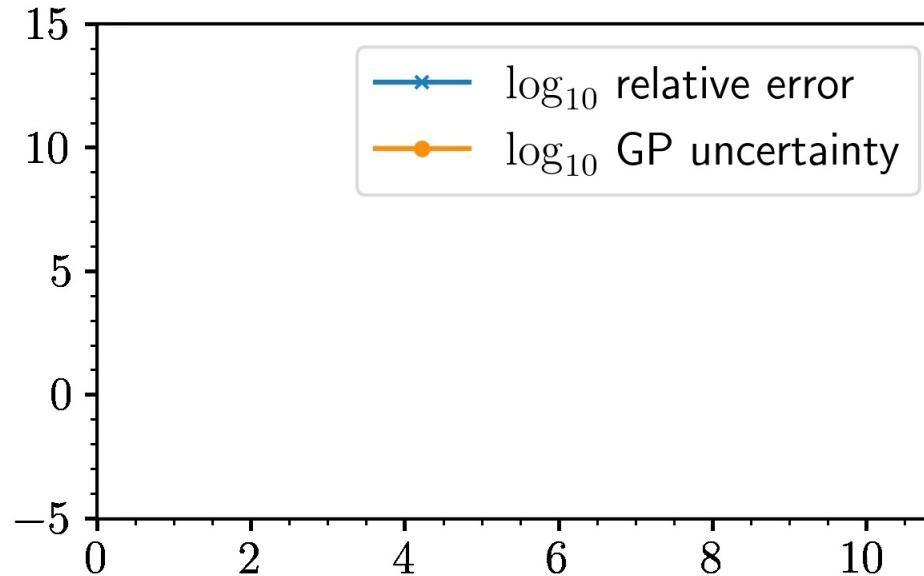
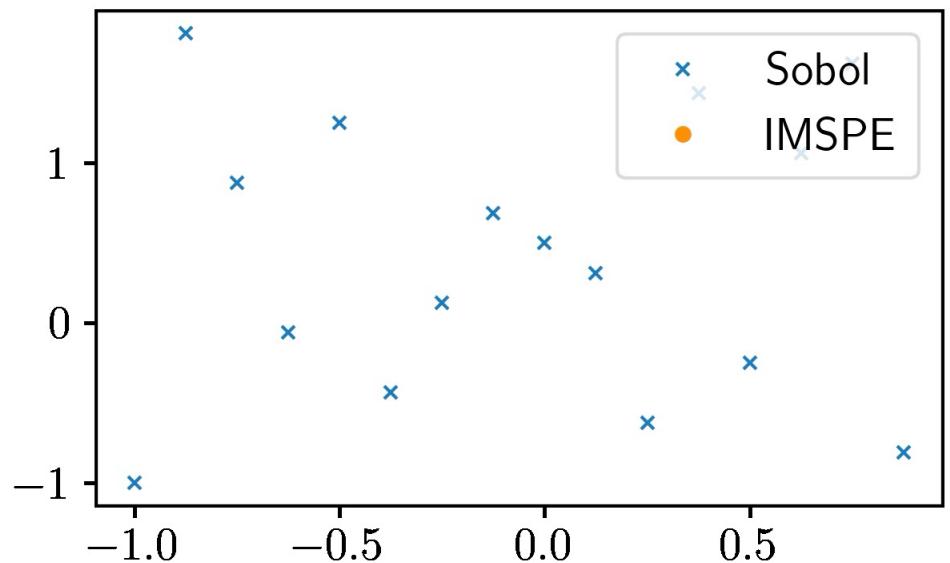
- Our code still in development...
- Code based largely on existing Python frameworks ([BoTorch](#))

[Balandat et al. (2019)]

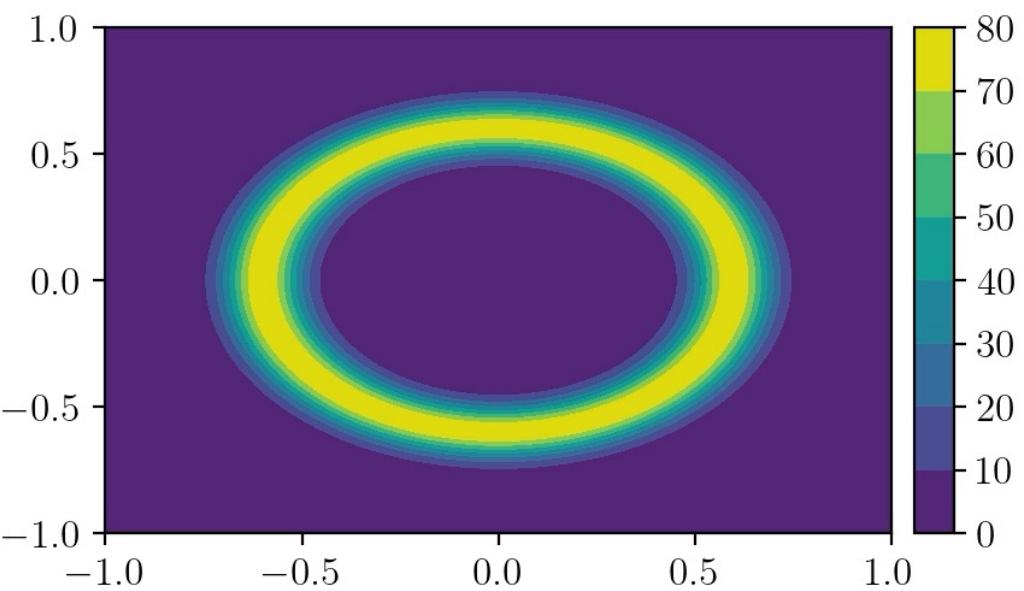
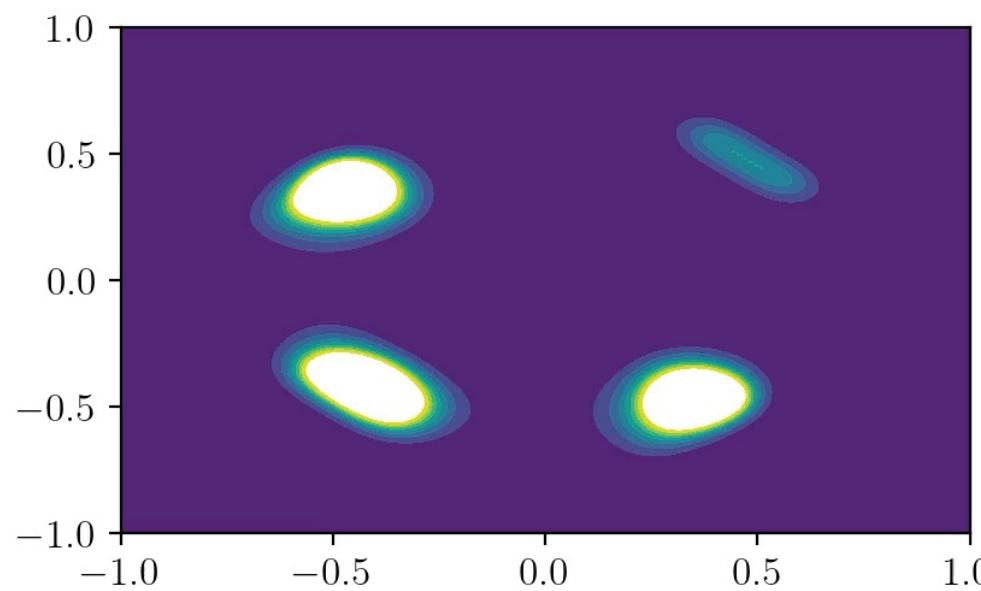
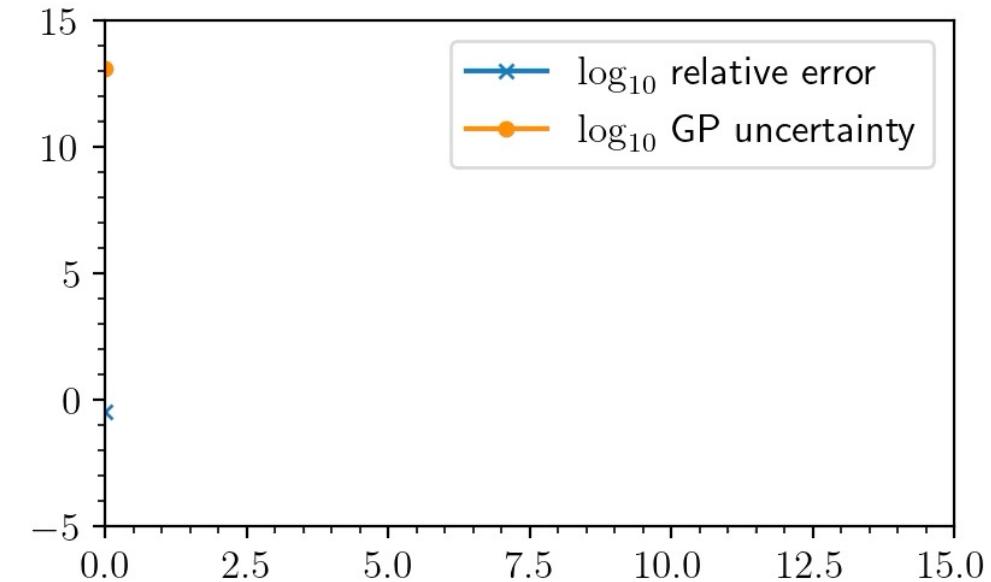
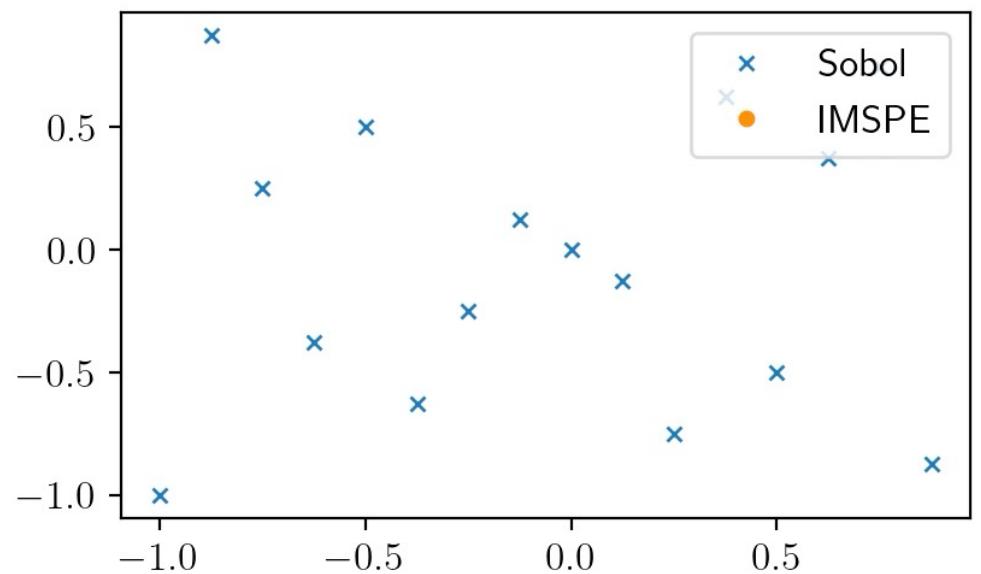
- Uses clever method for dealing with hyperparameters and acquisition function maximization ([Sparse Axis-Aligned Subspace Bayesian Optimisation \(SAASBO\)](#))
- Sampling from hyperparameter space PDF instead of maximizing (overengineering? – but more Bayesian in spirit)

[Eriksson & Jankowiak (2021)]

Step 0



Step 0



Conclusions

- Bayesian optimisation is a very efficient machine-learning technique for extremising unknown functions
- It can also be applied to cosmological parameter estimation and Bayesian model comparison
- Most useful for expensive-to-calculate likelihoods and complicated posterior distributions
- Paper and code for Bayesian evidence calculation out soon!