

Fast Bayesian inference with Gaussian Processes

DSU Sydney 2022

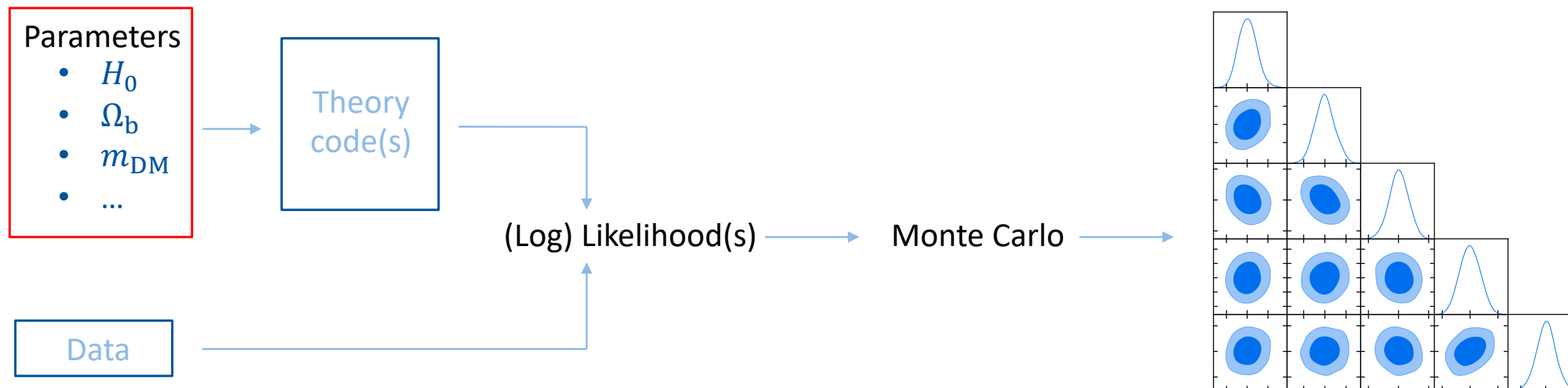
arXiv:2211.02045

<https://github.com/jonasegammal/GPry>

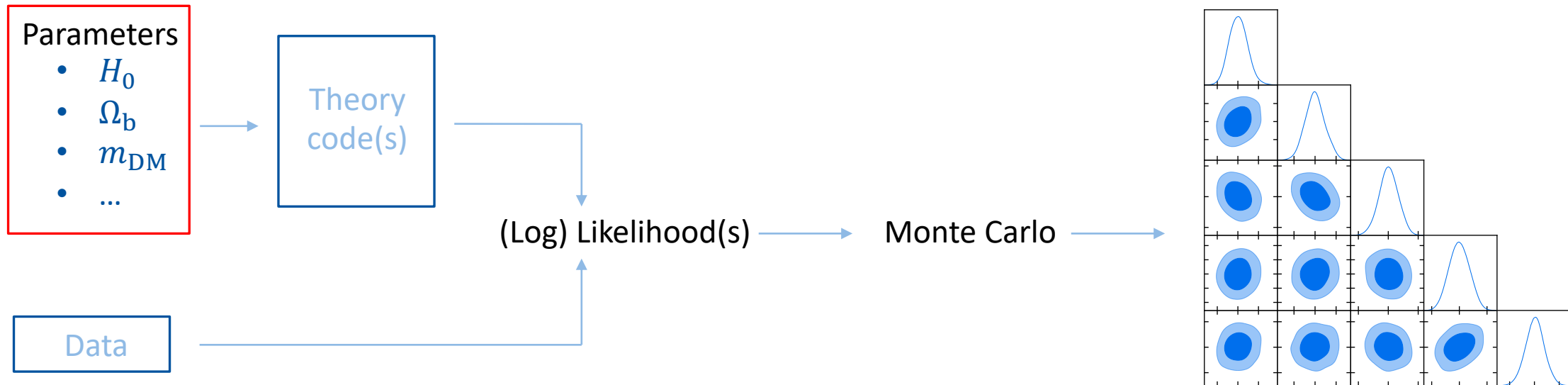
JONAS EL GAMMAL (UNIVERSITY OF STAVANGER)

WITH J. TORRADO, N. SCHÖNEBERG, C. FIDLER

1. Idea

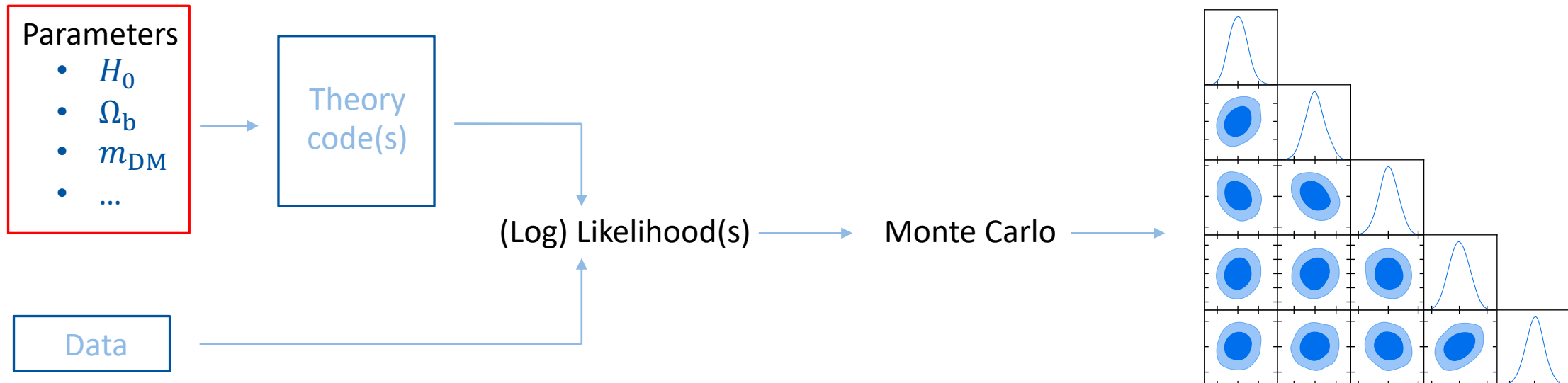


1. Idea



Example: 8d $\sim 10^5$ samples for MCMC

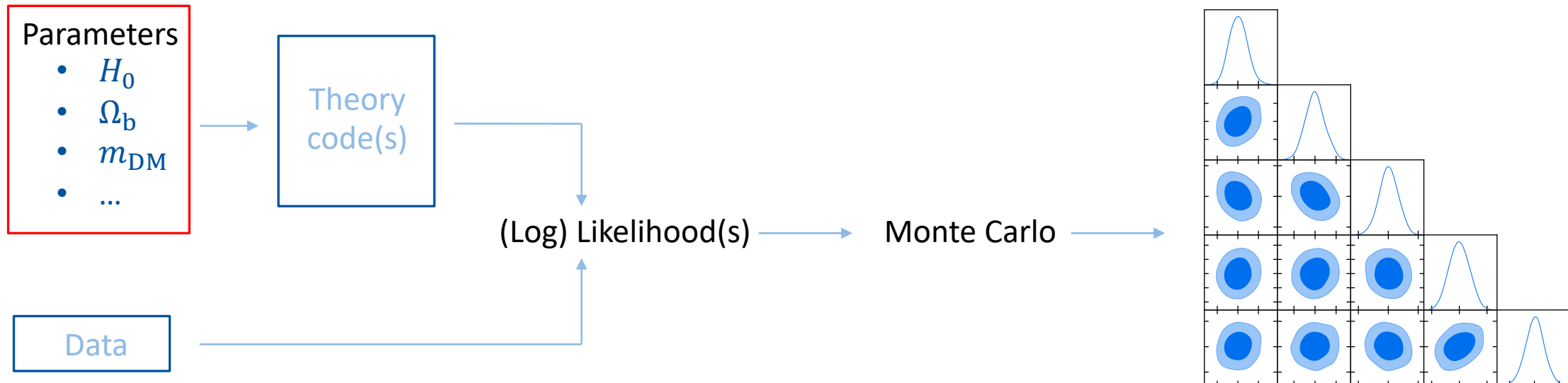
1. Idea



Example: $8\text{d} \sim 10^5$ samples for MCMC

Likelihood eval. time	Total time for inference
1 s	~ 1 day
1 min	
10 min	

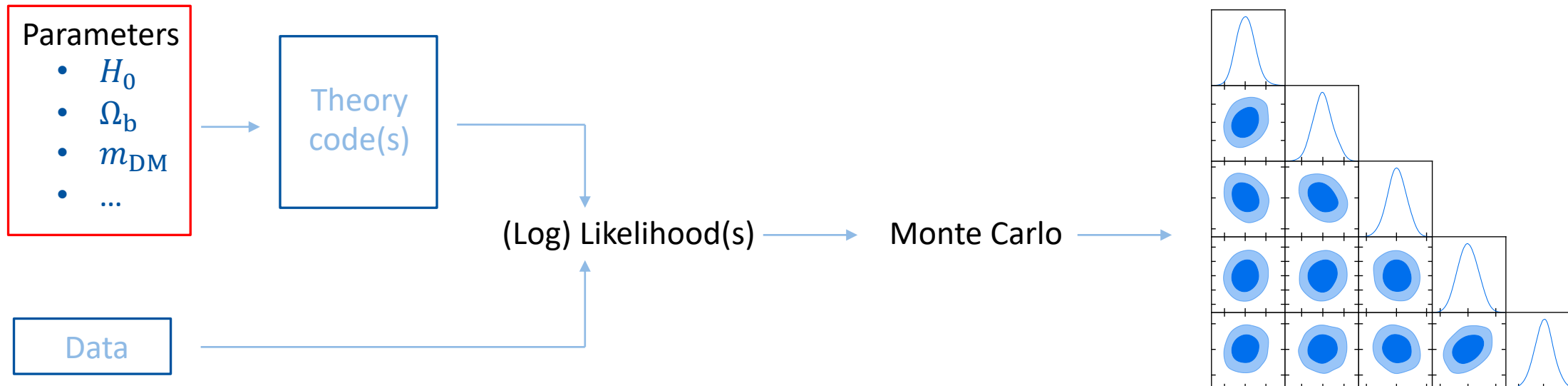
1. Idea



Example: $8\text{d} \sim 10^5$ samples for MCMC

Likelihood eval. time	Total time for inference
1 s	~ 1 day
1 min	~ 1 month
10 min	

1. Idea

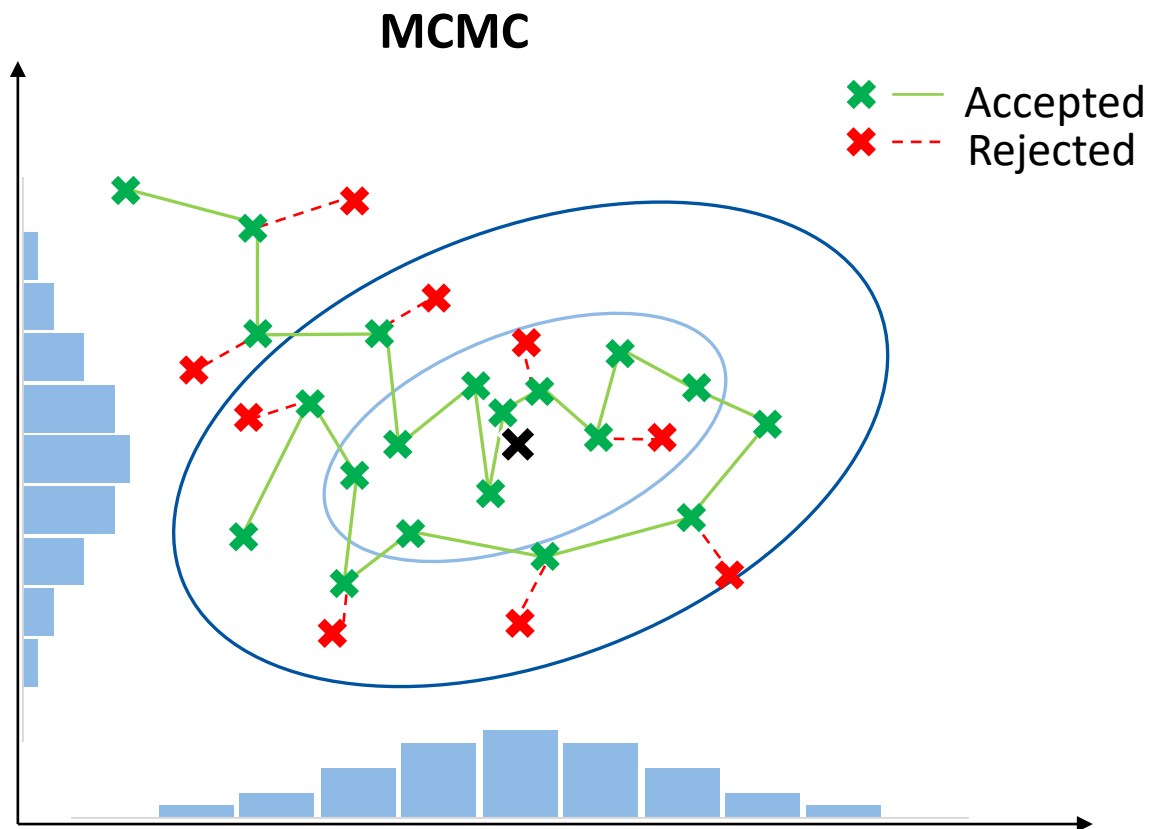


Example: $8d \sim 10^5$ samples for MCMC

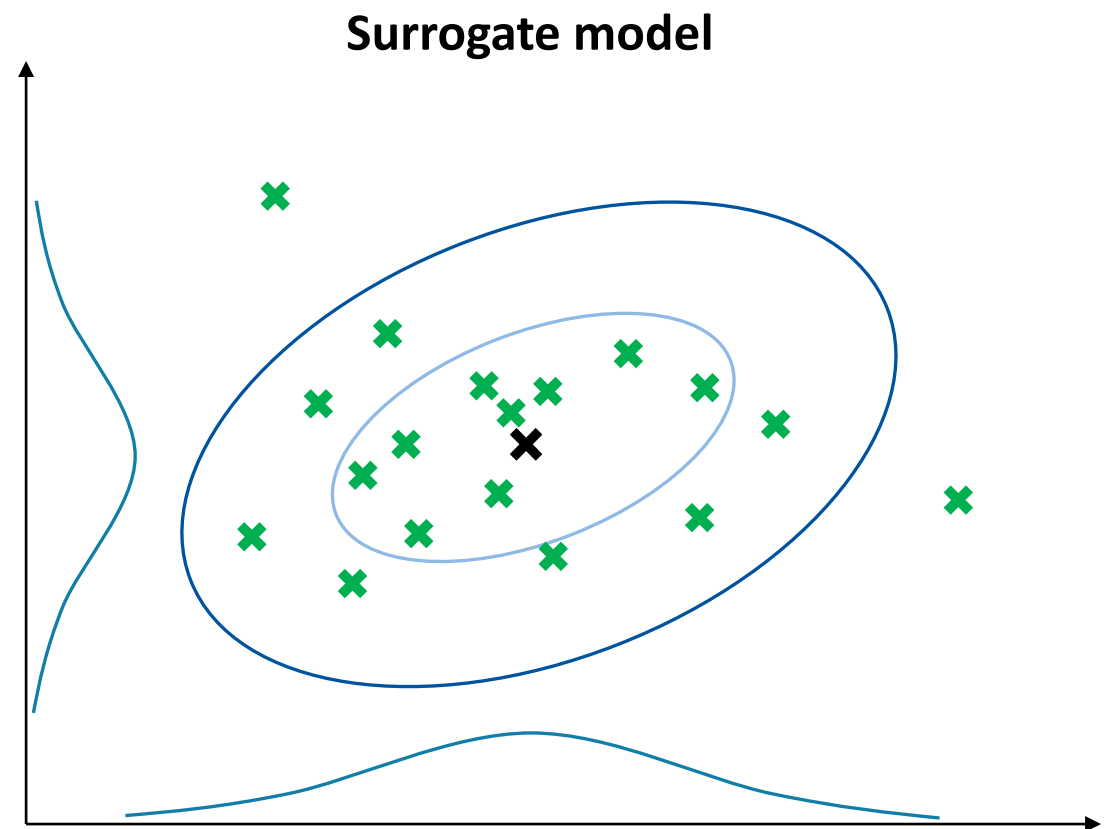
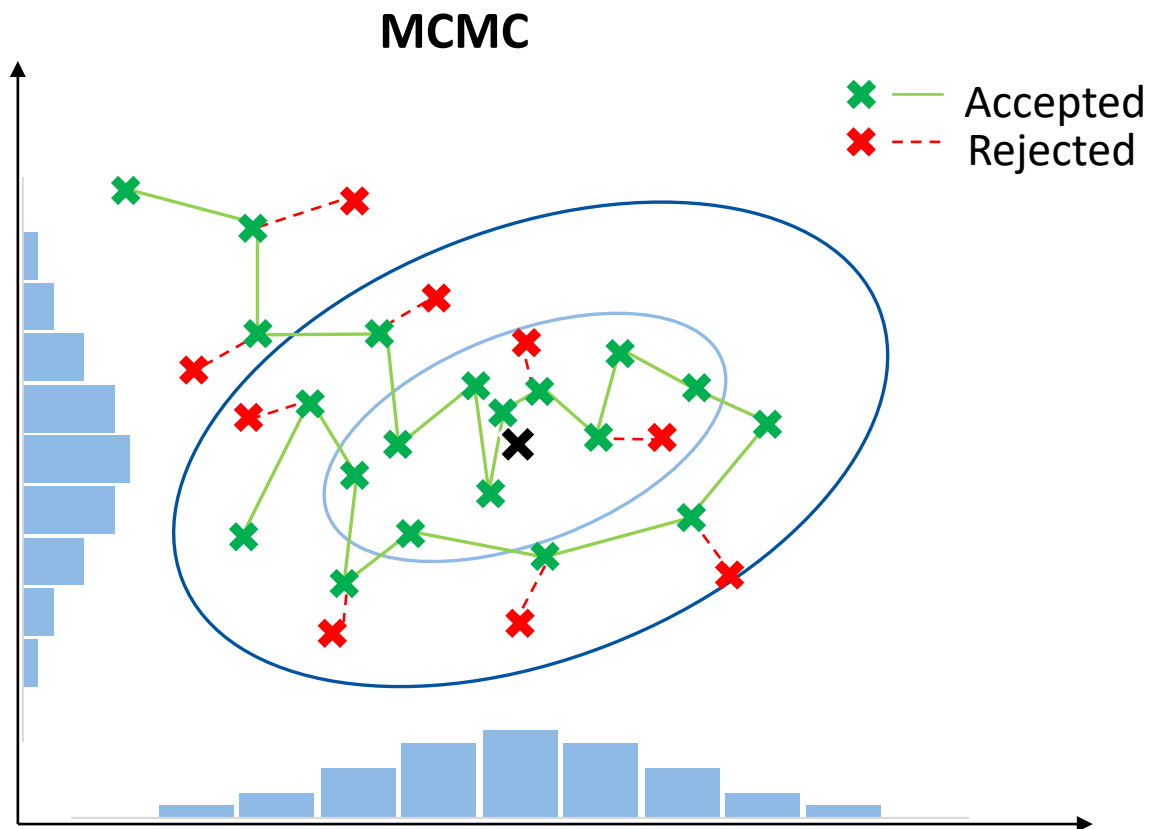
Likelihood eval. time	Total time for inference
1 s	~ 1 day
1 min	~ 1 month
10 min	~ 1 year



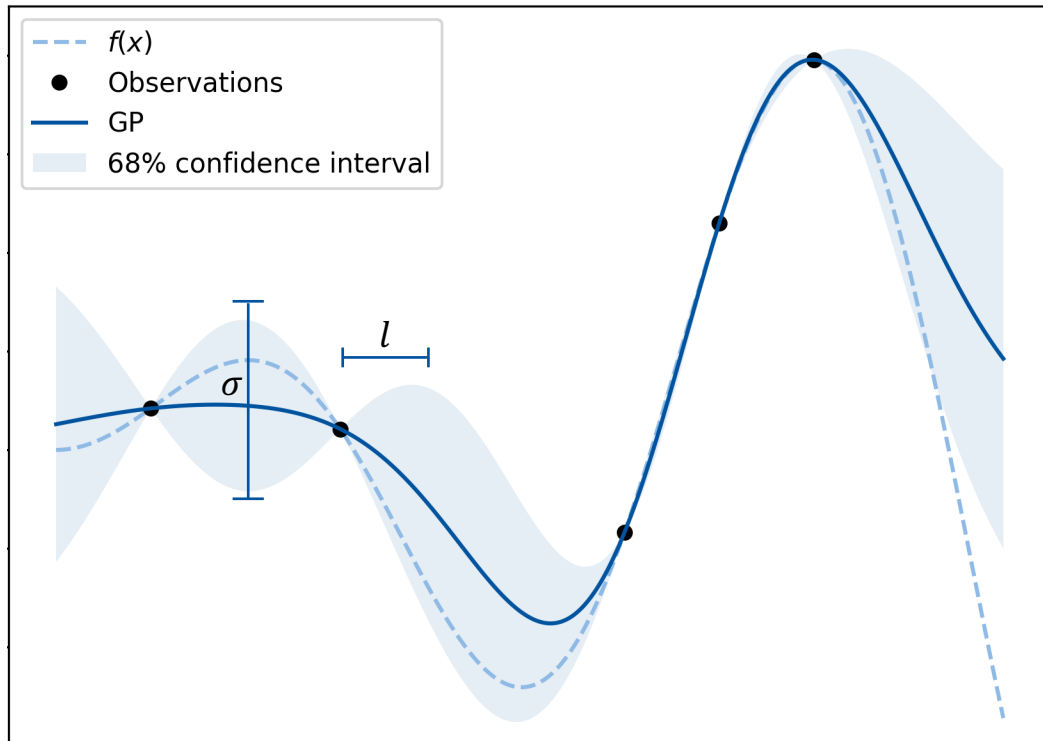
1. Idea



1. Idea



2. Gaussian Process Surrogate

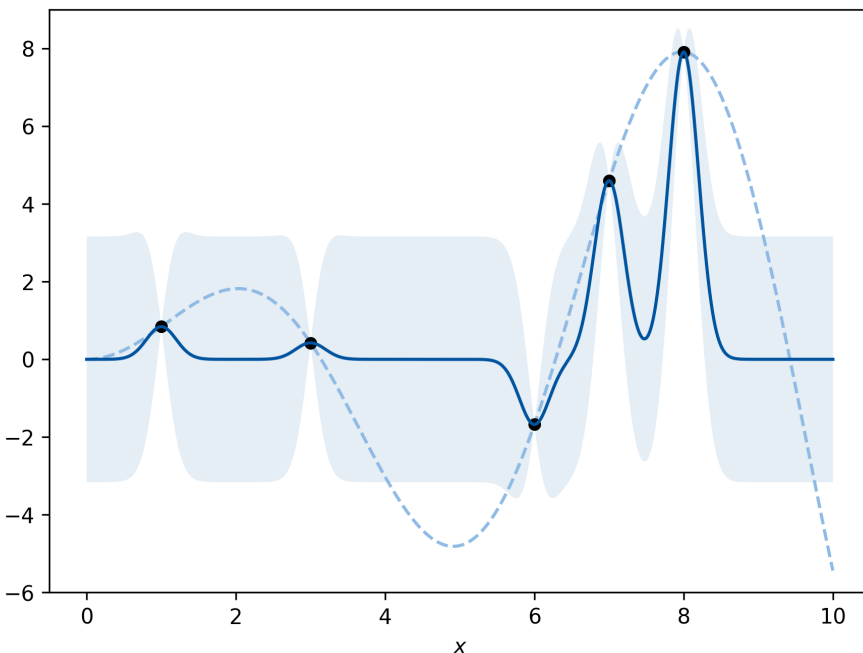


$$k(x, x') = \sigma^2 \cdot \exp\left(-\frac{(x - x')^2}{2l^2}\right)$$

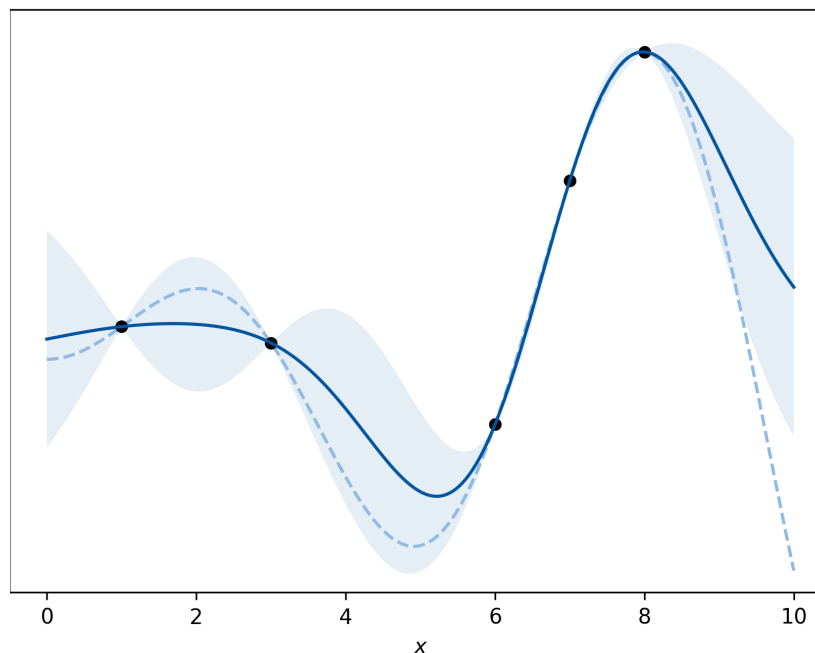
2. Gaussian Process Surrogate

$$k(x, x') = \sigma^2 \cdot \exp\left(-\frac{|x - x'|^2}{2l^2}\right)$$

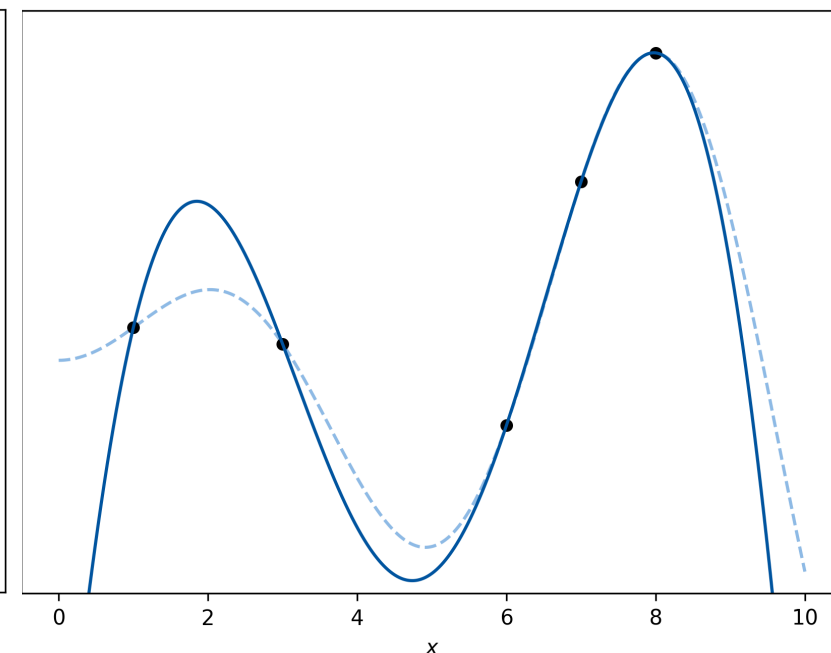
$\sigma^2 = 10, l = 0.2$



$\sigma^2 = 17.25, l = 1.25$

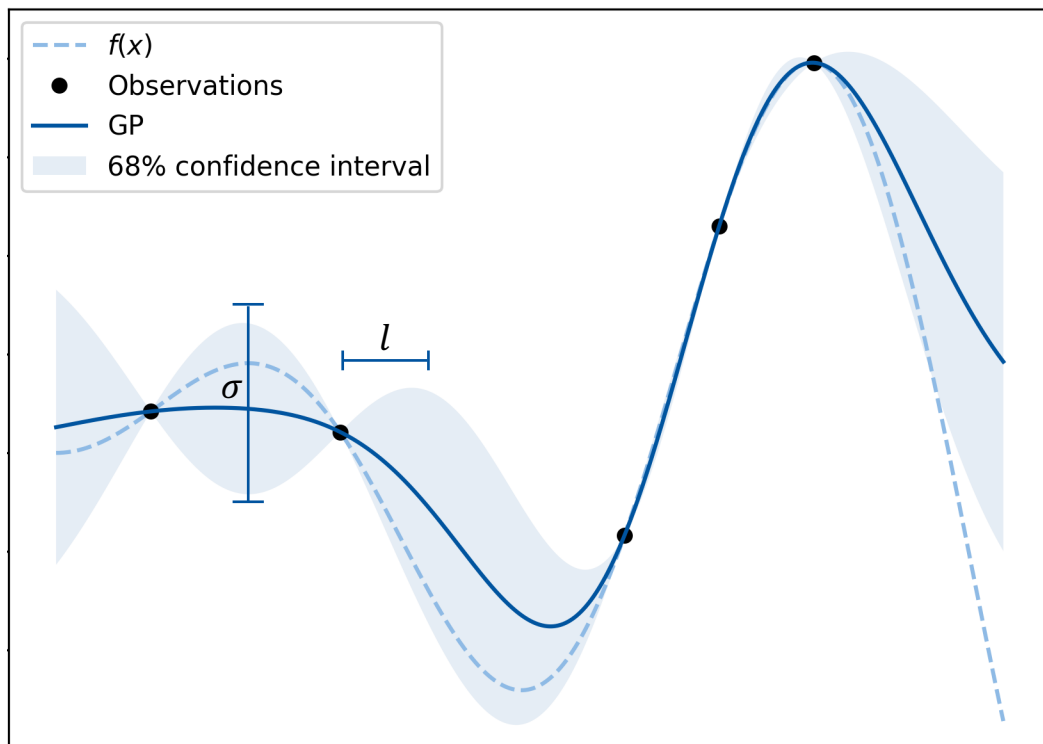


$\sigma^2 = 1, l = 5$

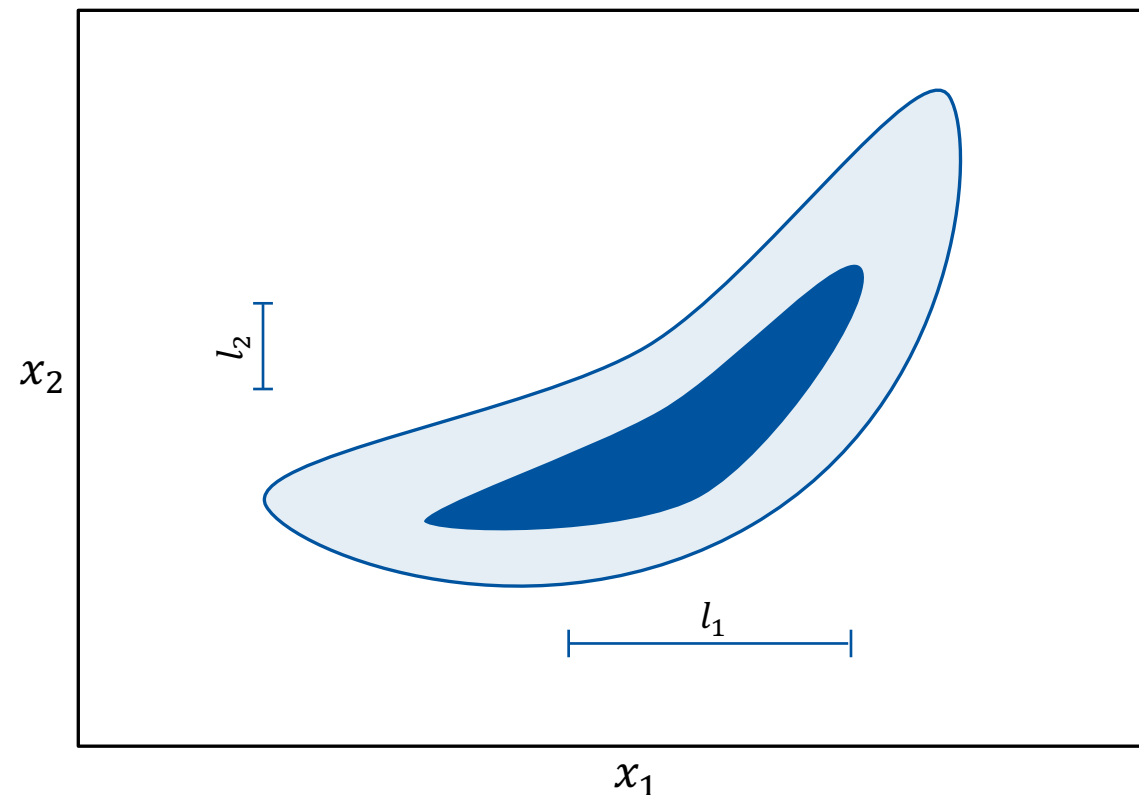


Can use MAP to get the best estimate

2. Gaussian Process Surrogate

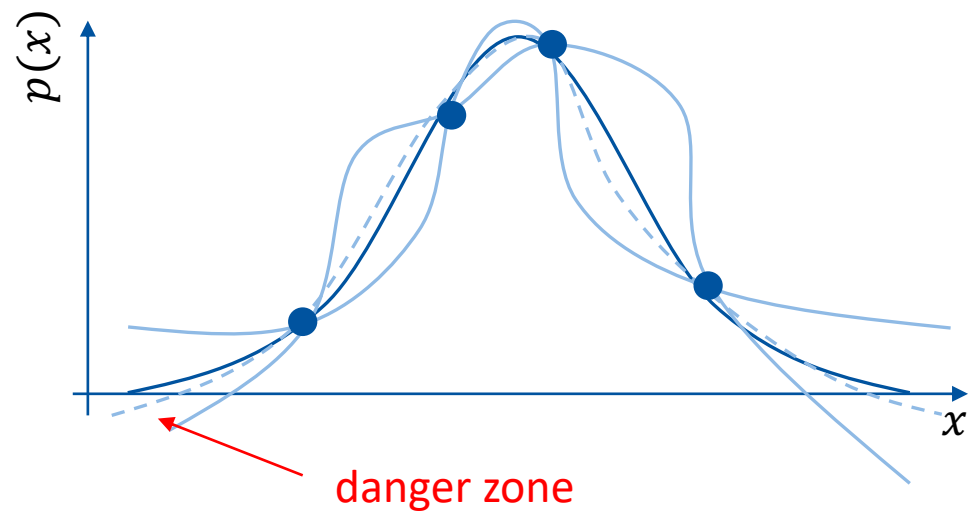


$$k(x, x') = \sigma^2 \cdot \exp\left(-\frac{(x - x')^2}{2l^2}\right)$$

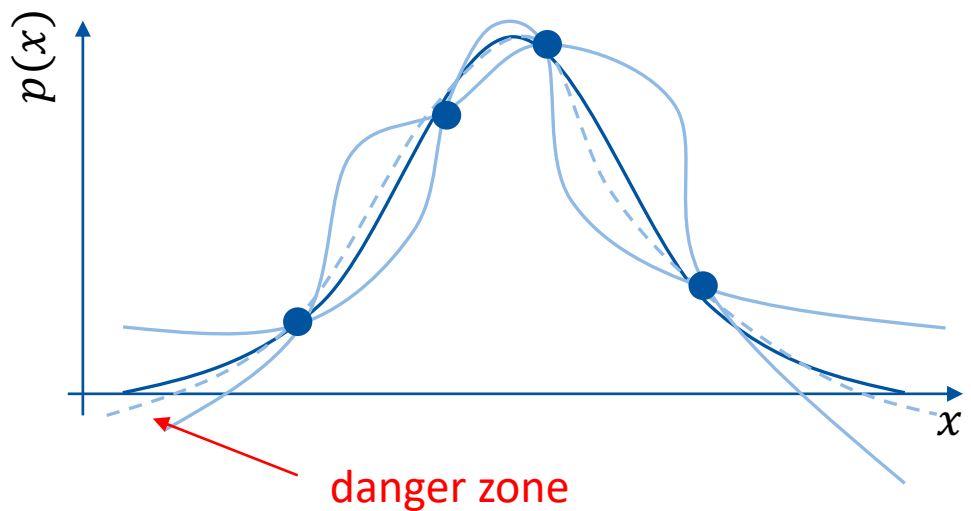


$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \cdot \exp\left(-\sum \frac{(x_i - x'_i)^2}{2l_i^2}\right)$$

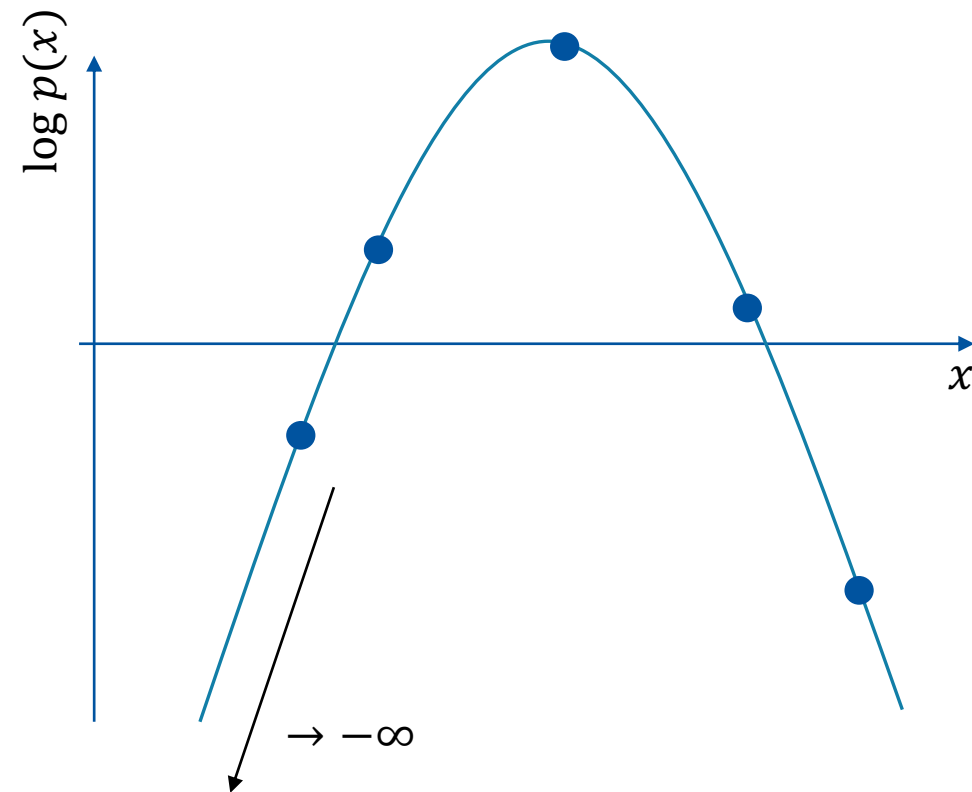
3. Region of interest



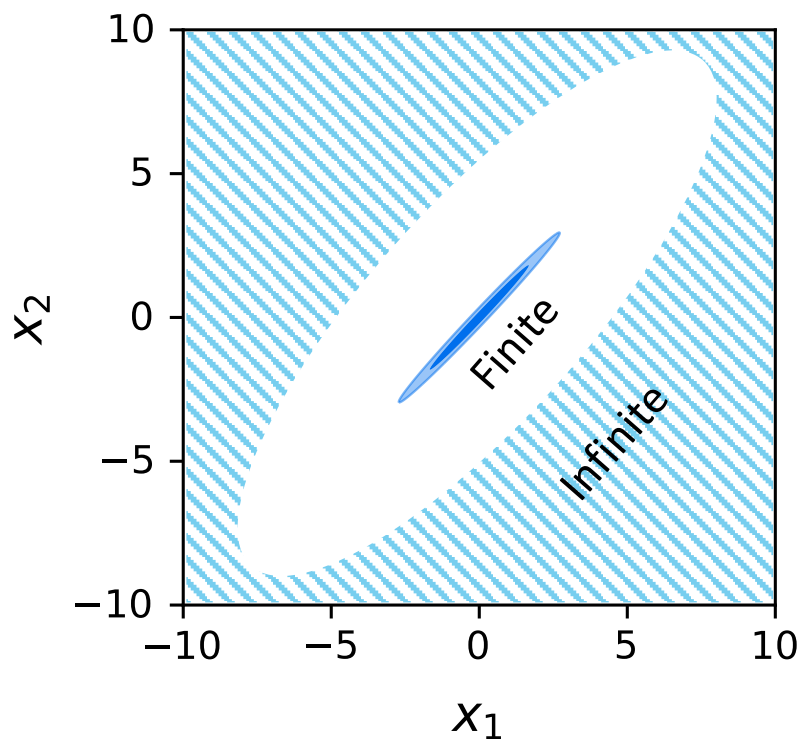
3. Region of interest



⇒ Interpolate **log-posterior** to enforce positivity

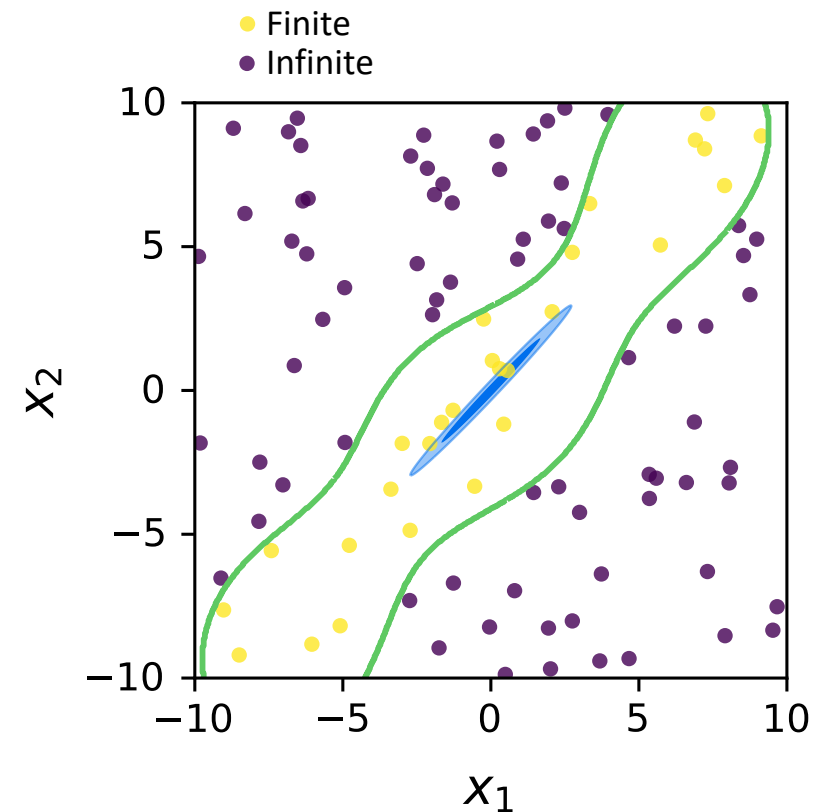


3. Region of interest



Solution: SVM Classifier

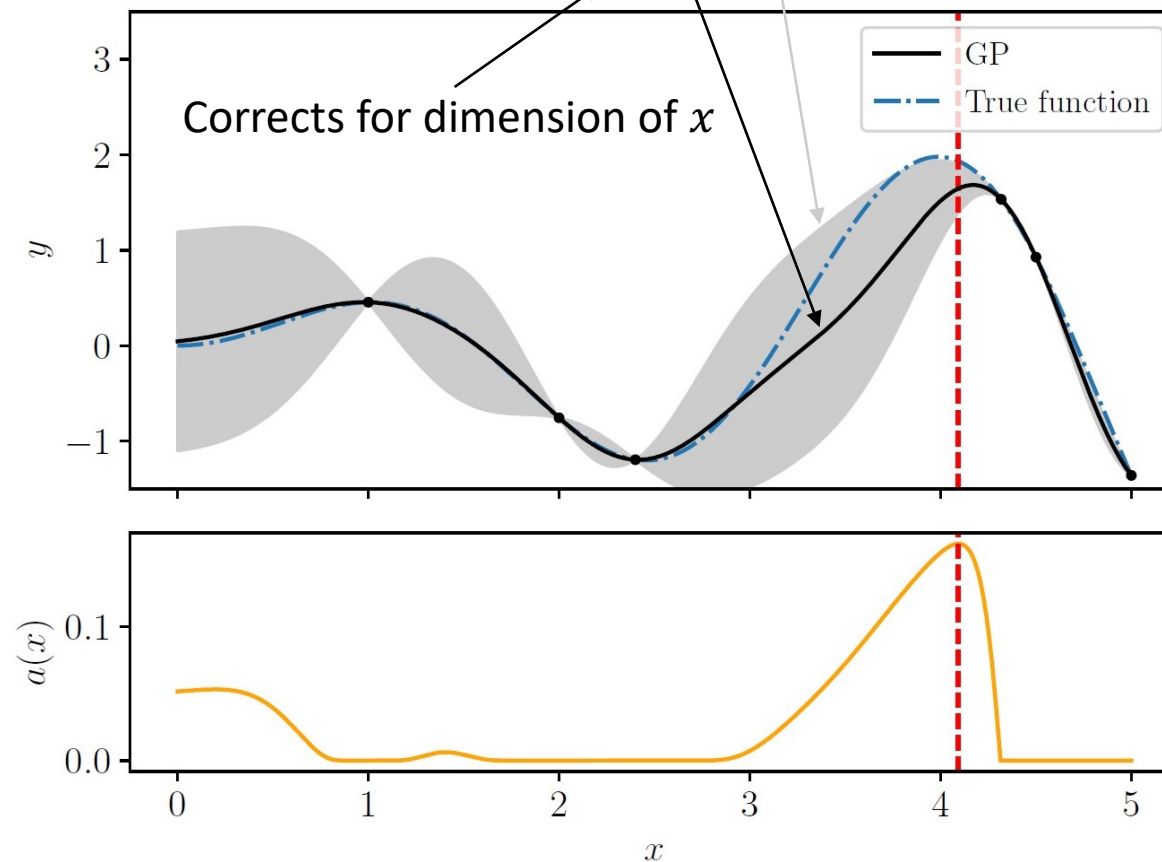
Multiply μ with $-\infty$ where SVM classifies as "infinite"



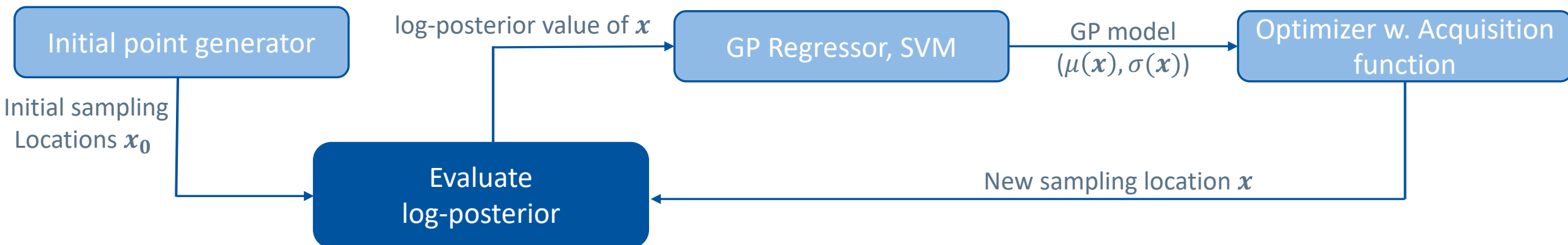
4. Active sampling

Propose samples by maximizing an **acquisition function**

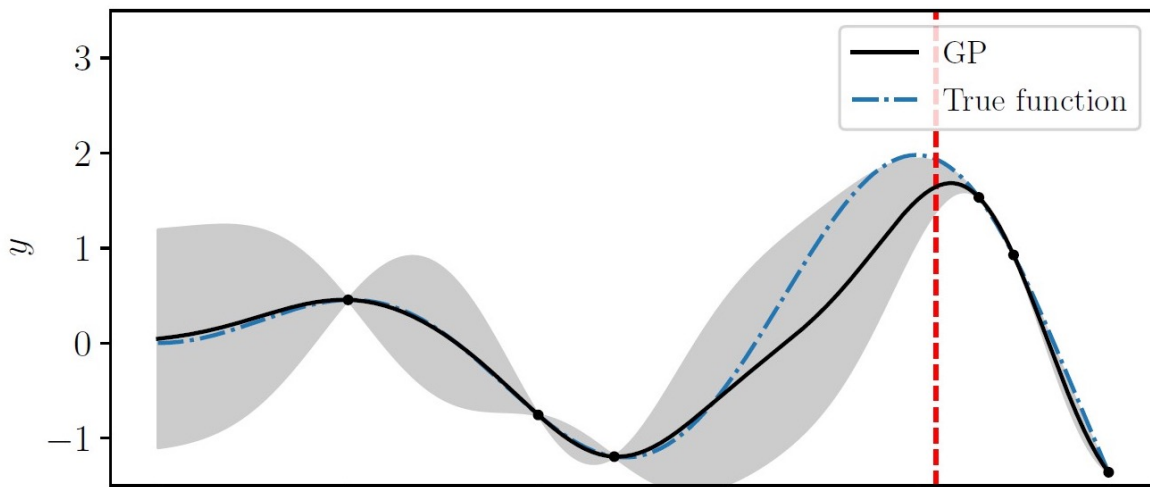
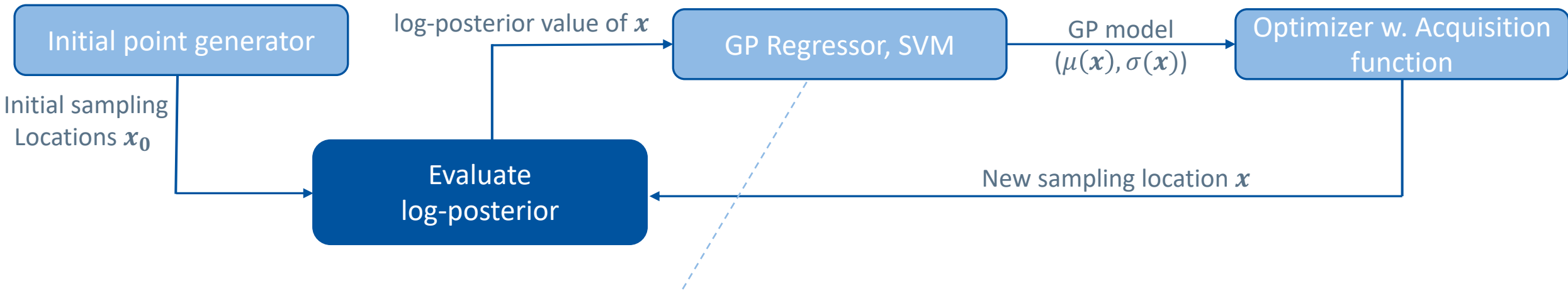
$$a(x) = \exp(2\zeta \cdot \bar{\mu}) \cdot \sigma_{\bar{\mu}}(x)$$



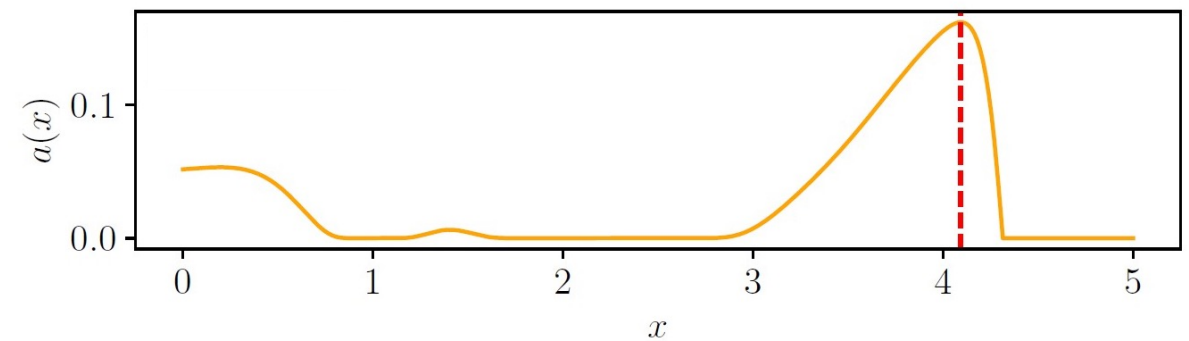
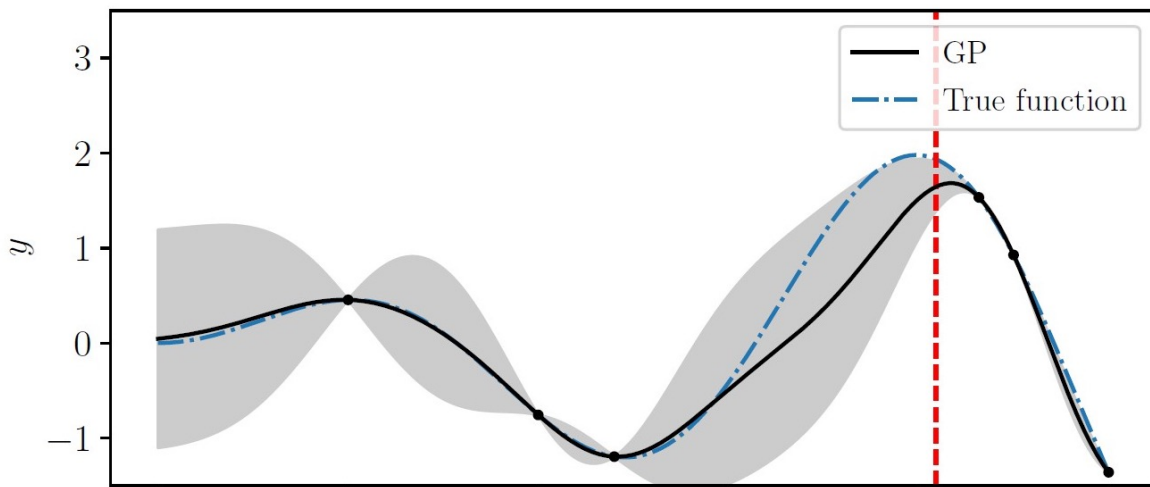
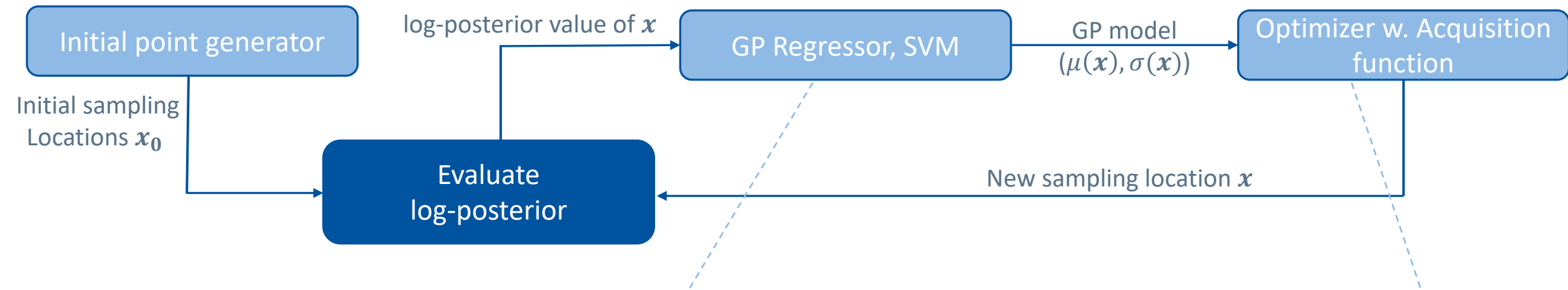
5. The Algorithm



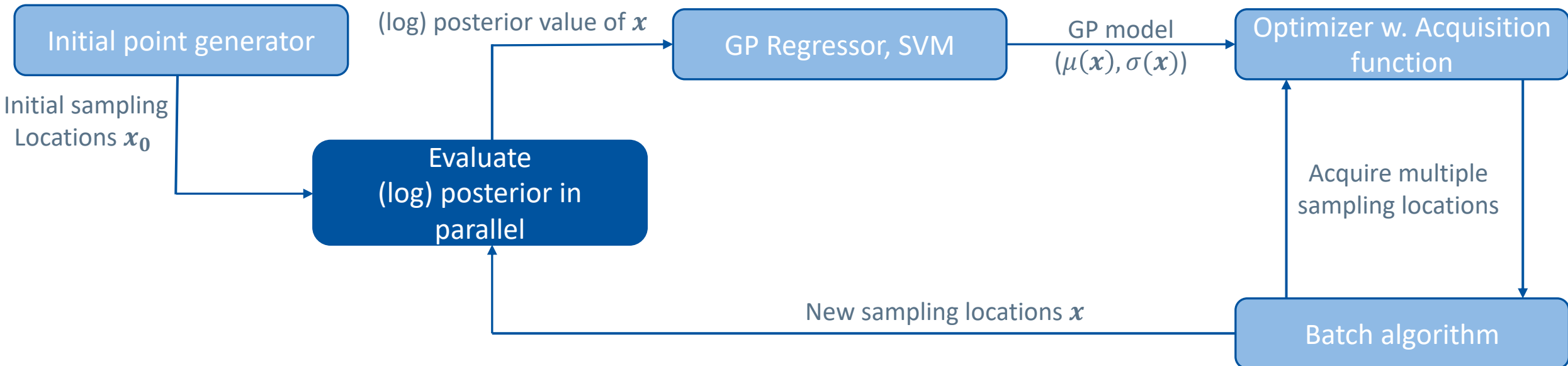
5. The Algorithm



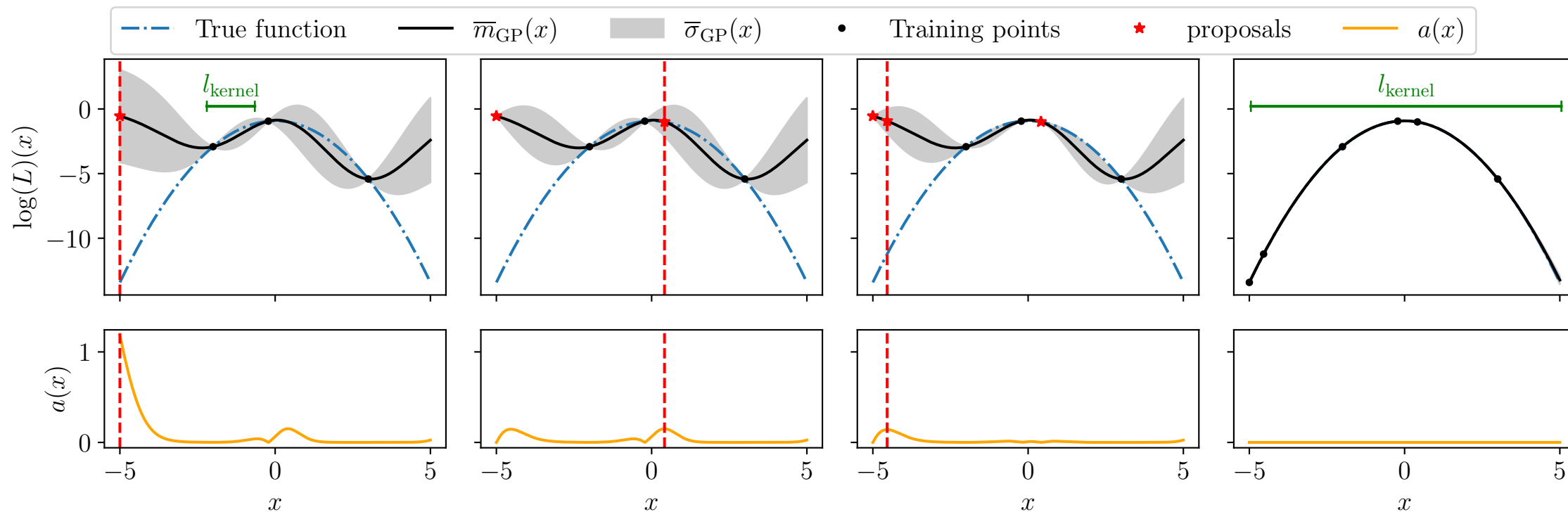
5. The Algorithm



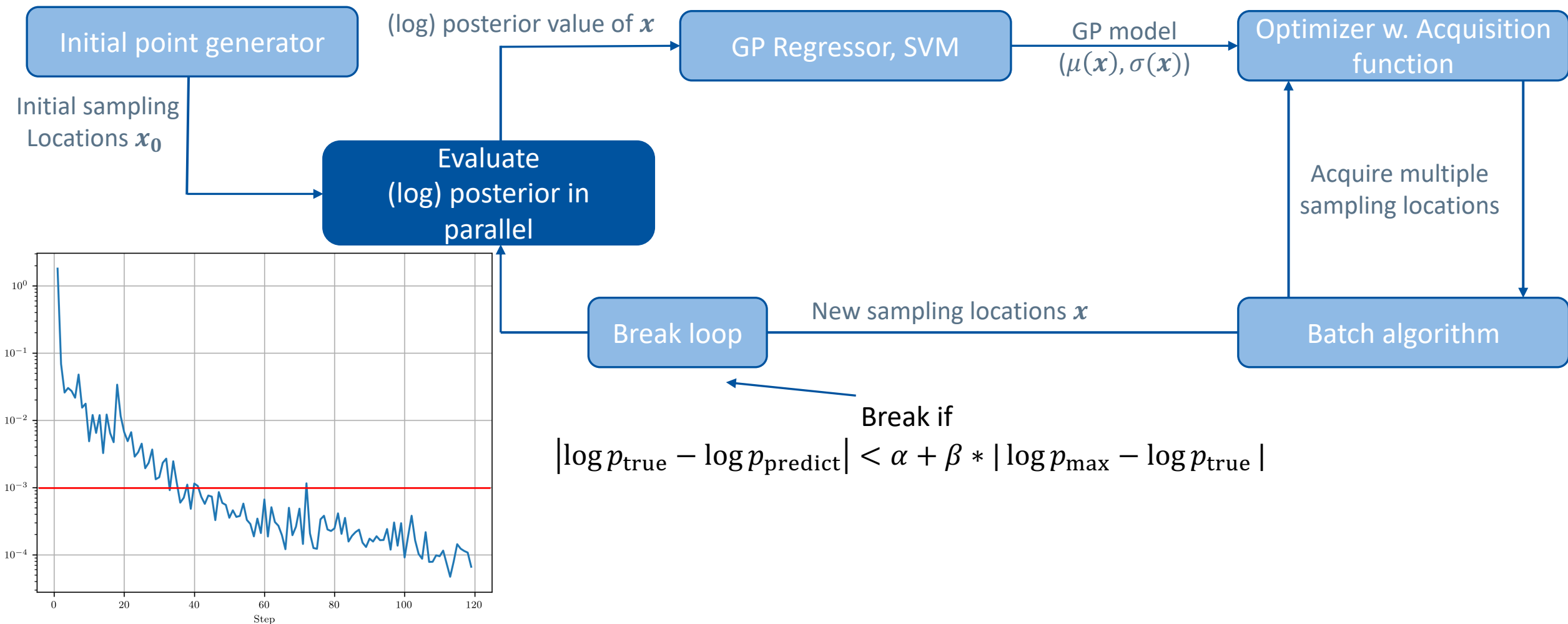
5. The Algorithm



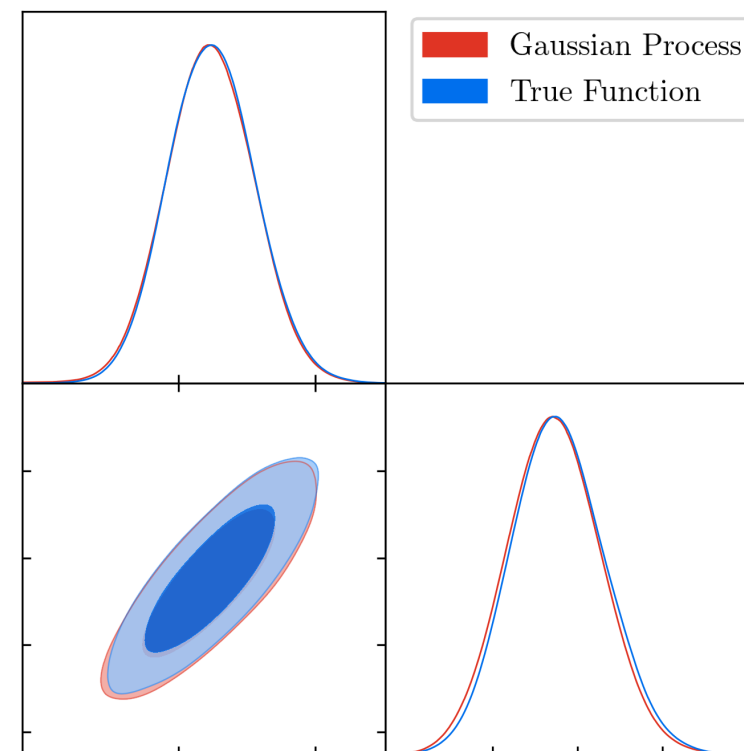
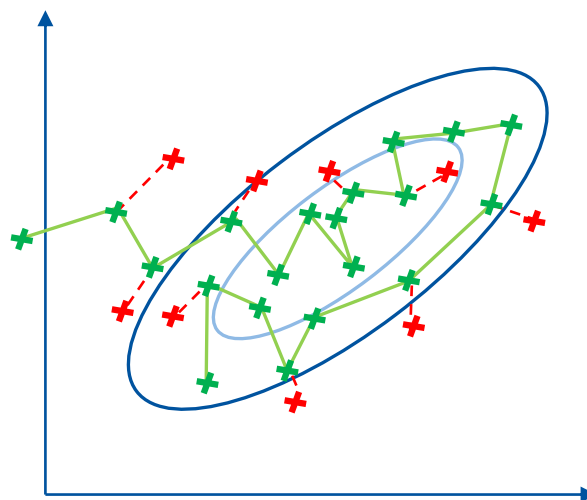
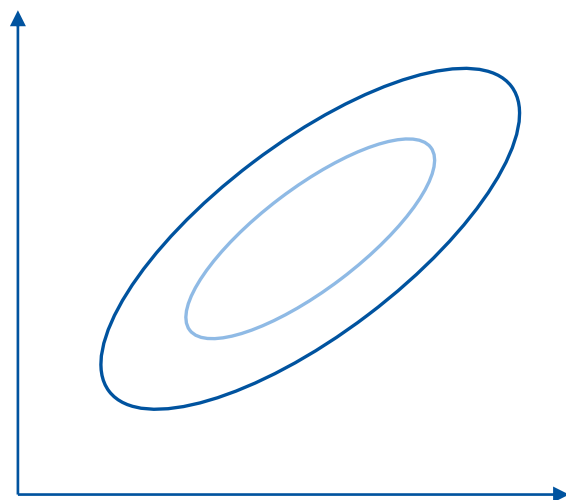
5. The Algorithm



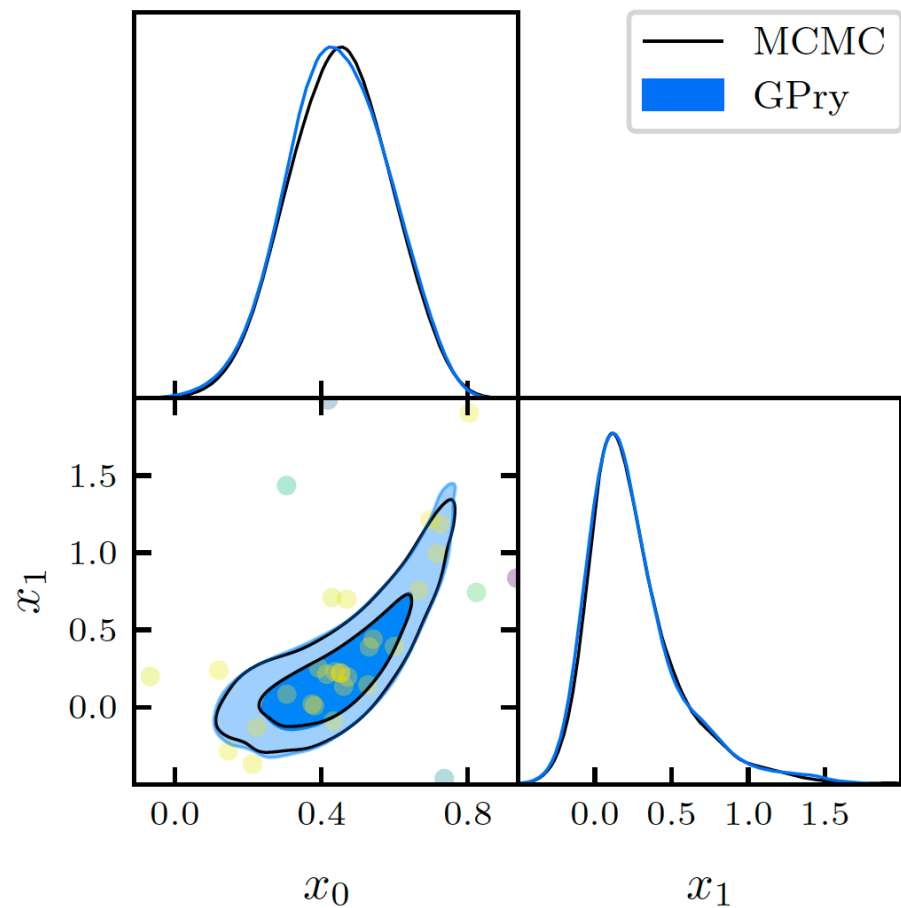
5. The Algorithm



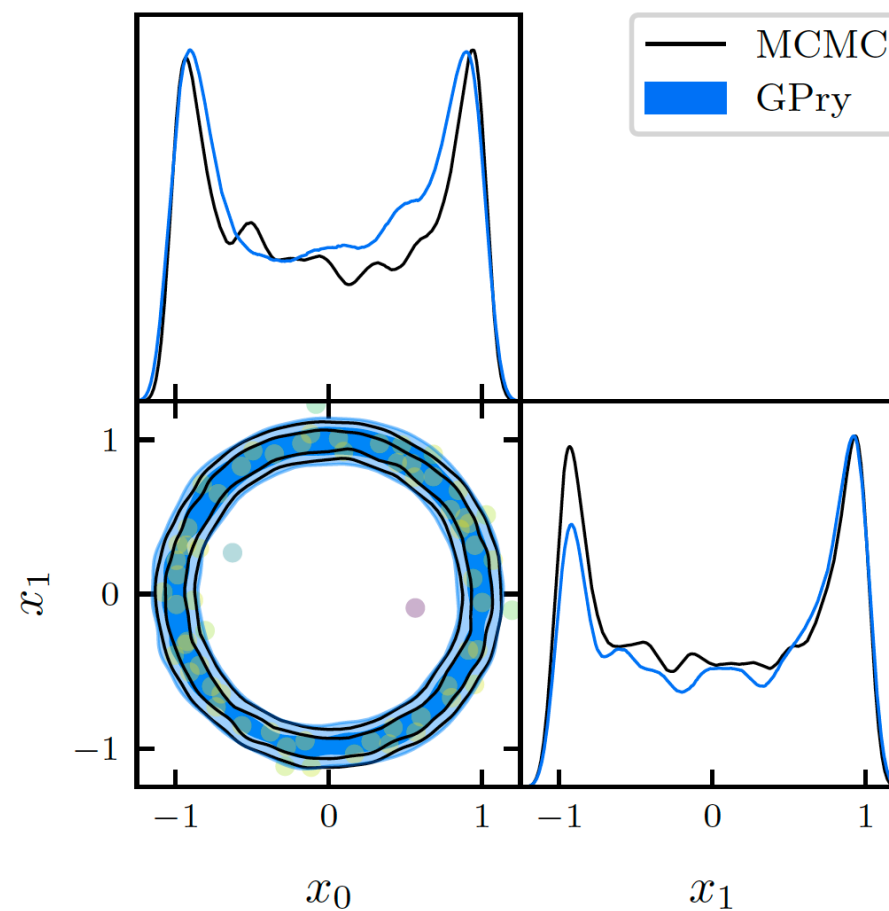
6. Marginalised quantities



7. Experiments

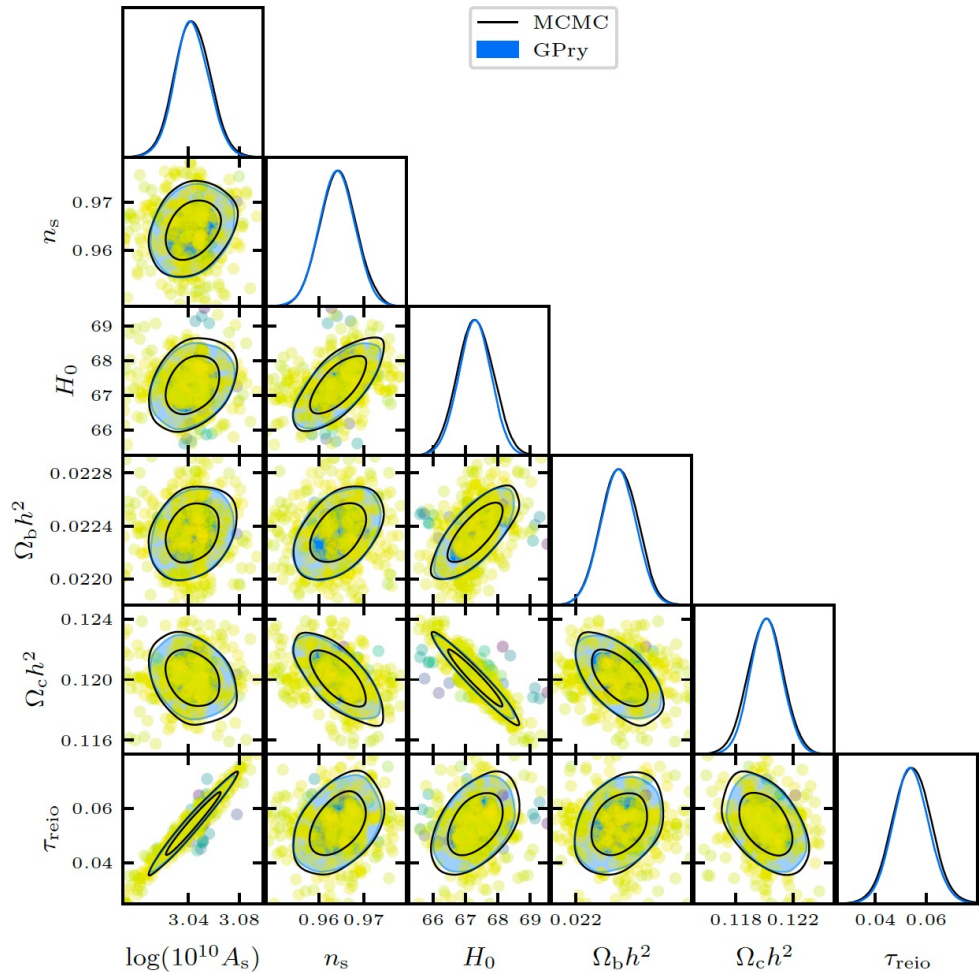


40 posterior evaluations

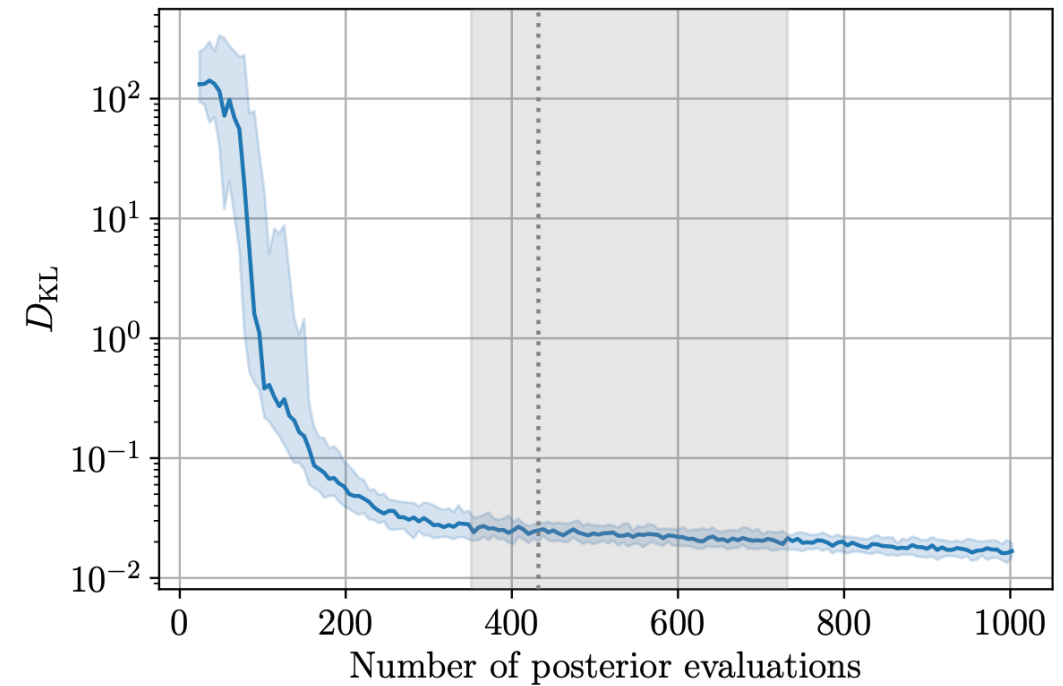


68 posterior evaluations

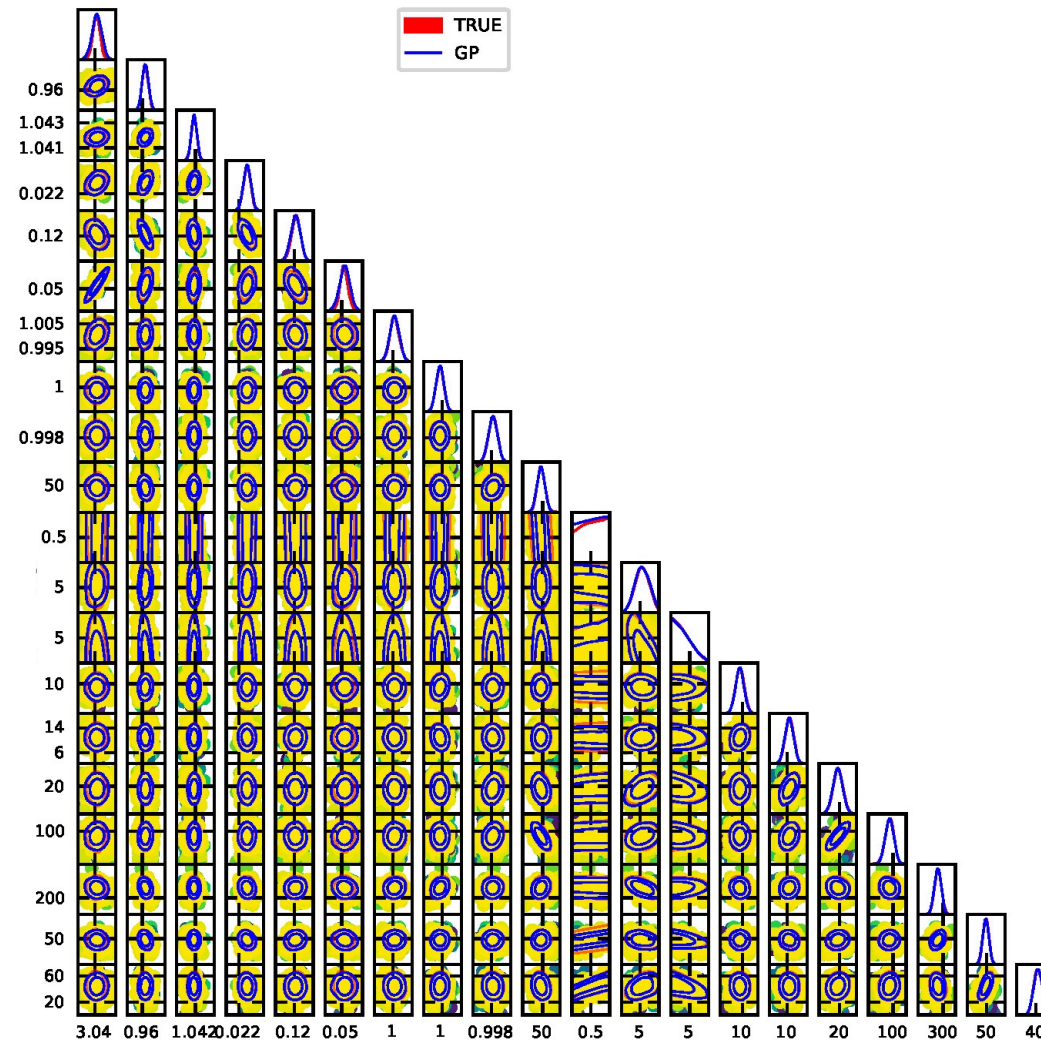
7. Experiments



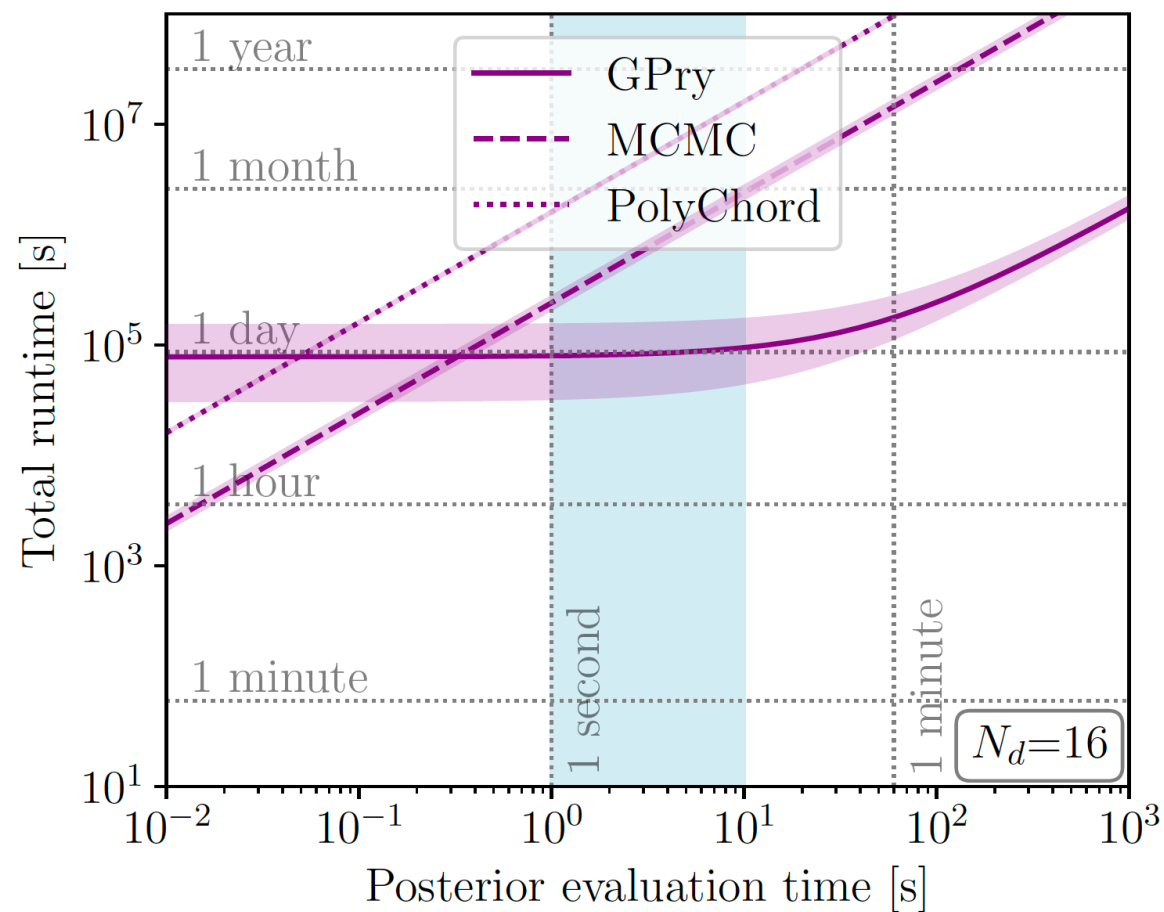
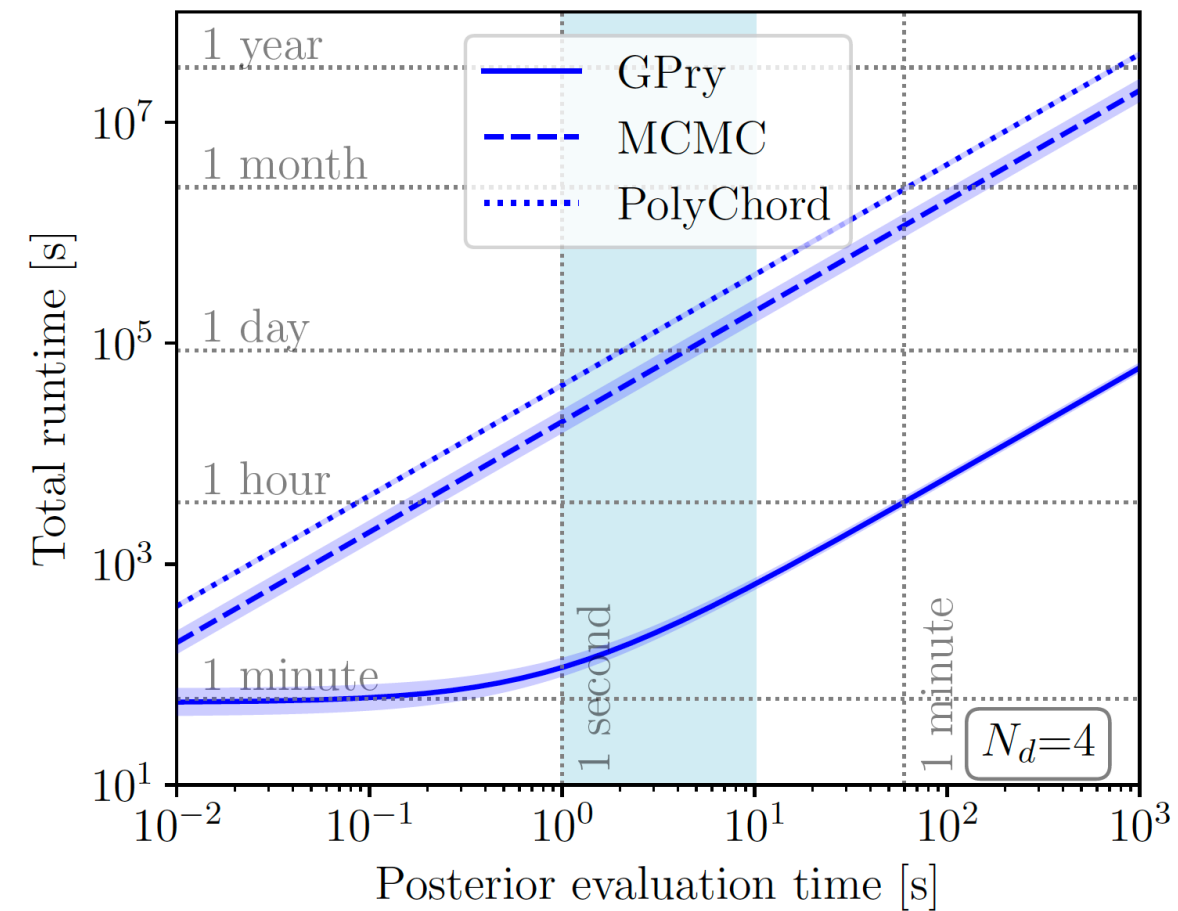
420 posterior evaluations



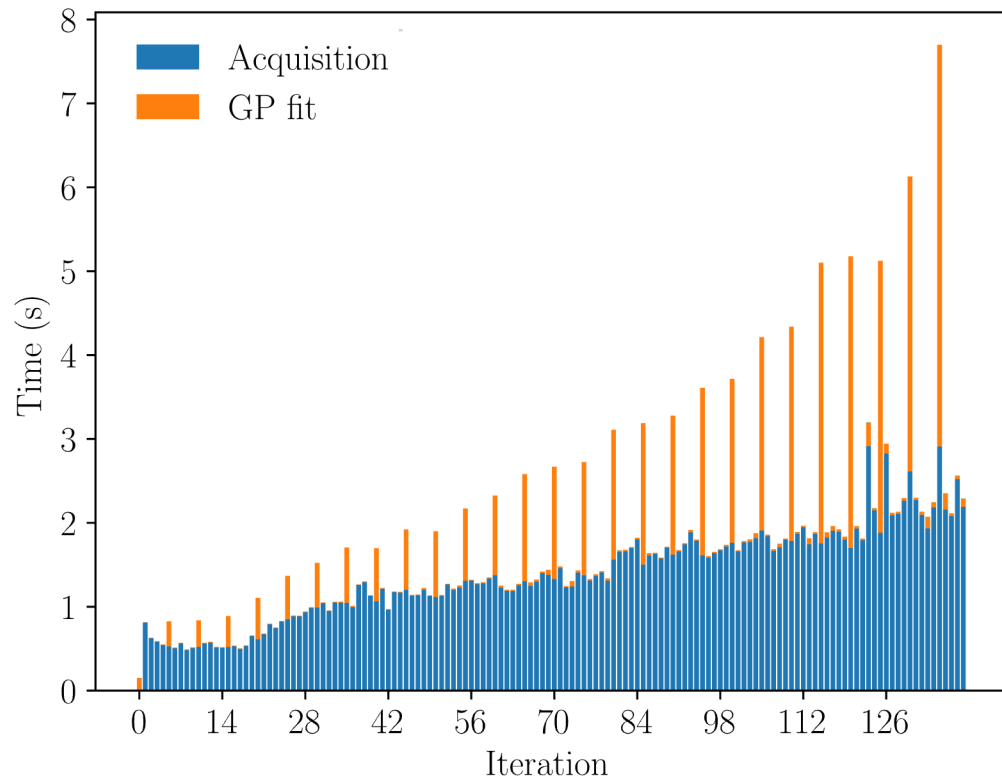
7. Experiments



8. Performance



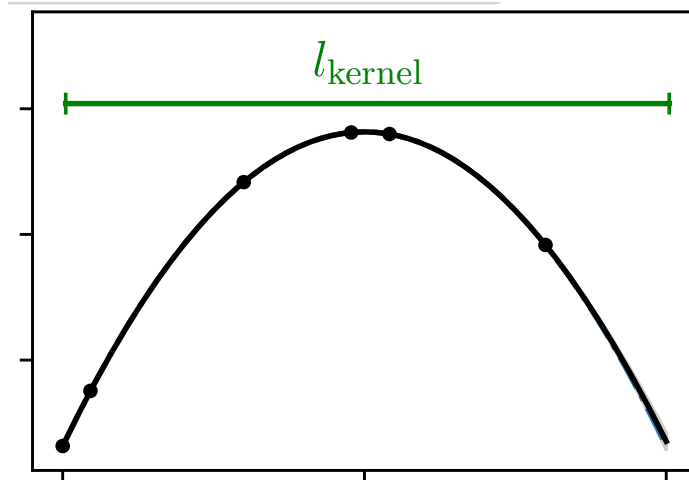
9. Limitations



Overhead

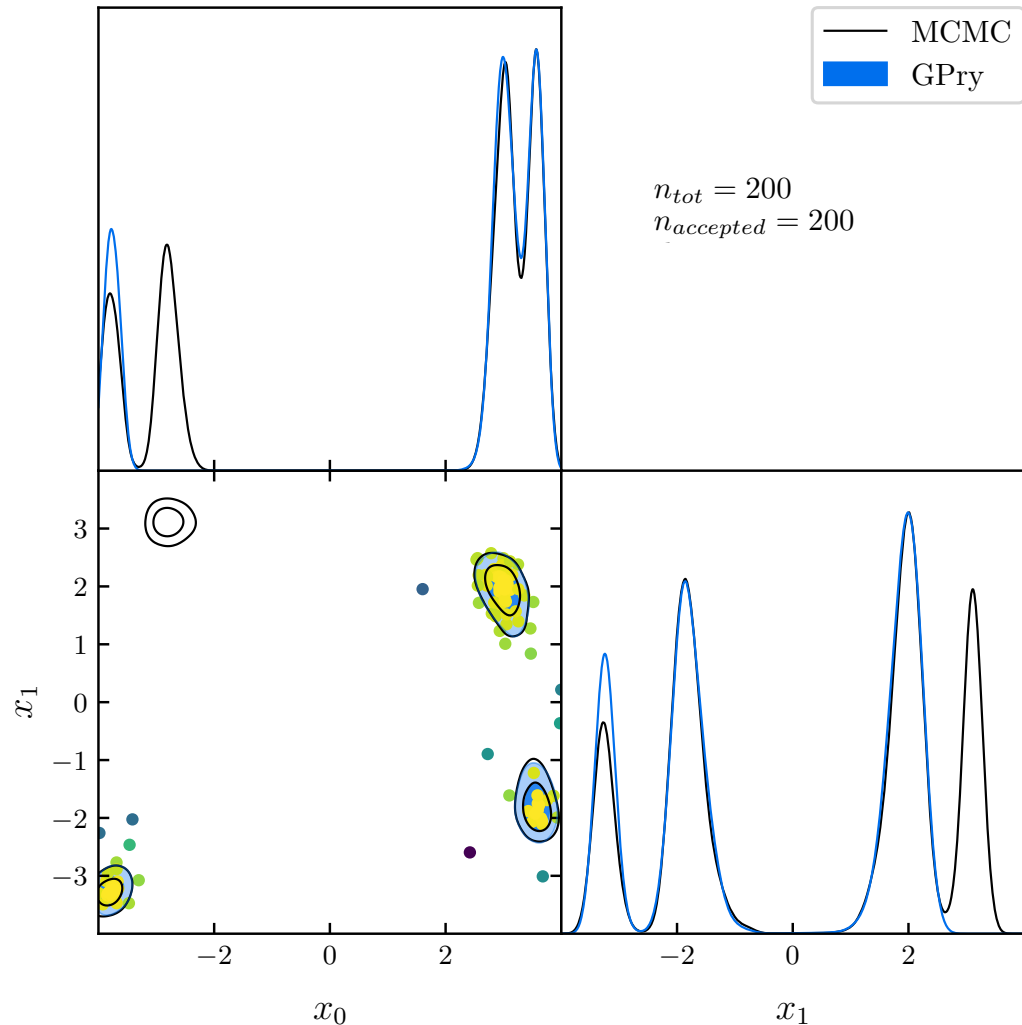
9. Limitations

Overhead



“Overfitting”

9. Limitations



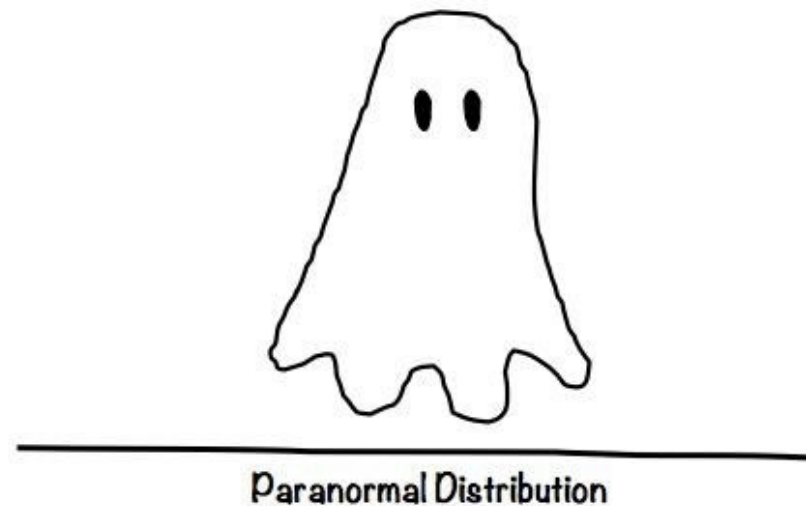
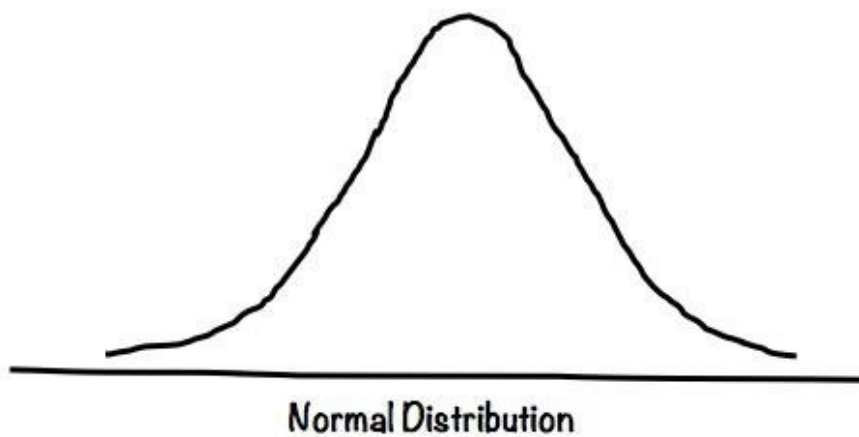
Overhead

“Overfitting”

Multimodality

We work on solving those problems...

Thank you!



<https://www.memedroid.com/memes/detail/3518248/Normal-vs-paranormal-distribution>

Backup

Gaussian Process Regression

Assume that

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

(usually $m(\mathbf{x}) = 0 \forall \mathbf{x}$)

Then $f_* = f(\mathbf{x}_*)$ drawn from the GP given the measurements $\mathbf{y} = \mathbf{f}(X)$ are distributed as

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} k(X, X) & k(X, X_*) \\ k(X_*, X) & k(X_*, X_*) \end{pmatrix}\right)$$

This means, that **test functions are distributed according to**

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*))$$

with

$$\bar{\mathbf{f}}_* = E[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X) \cdot K^{-1}(X, X) \mathbf{y}$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) \cdot K^{-1}(X, X) K(X, X_*)$$

Computational complexity
scales with n^3

Gaussian Process Regression

$$\text{Marginalize: } p(\mathbf{y}|X) = \int p(\mathbf{y}|\mathbf{f}, X) p(\mathbf{f}|X) d\mathbf{f}$$

$$= \mathcal{N}(\mathbf{f}, 0) = \mathbf{f}$$

$$= \mathcal{N}(0, K), \quad K_{ij} = k(x_i, x_j)$$

$$\Rightarrow \log(p(\mathbf{y}|X)) = -\frac{1}{2} \mathbf{y}^T K^{-1} \mathbf{y} - \frac{1}{2} \log|K| - \frac{n}{2} \log(2\pi)$$

↳ **Can use MLE to get the best estimate**

Active sampling

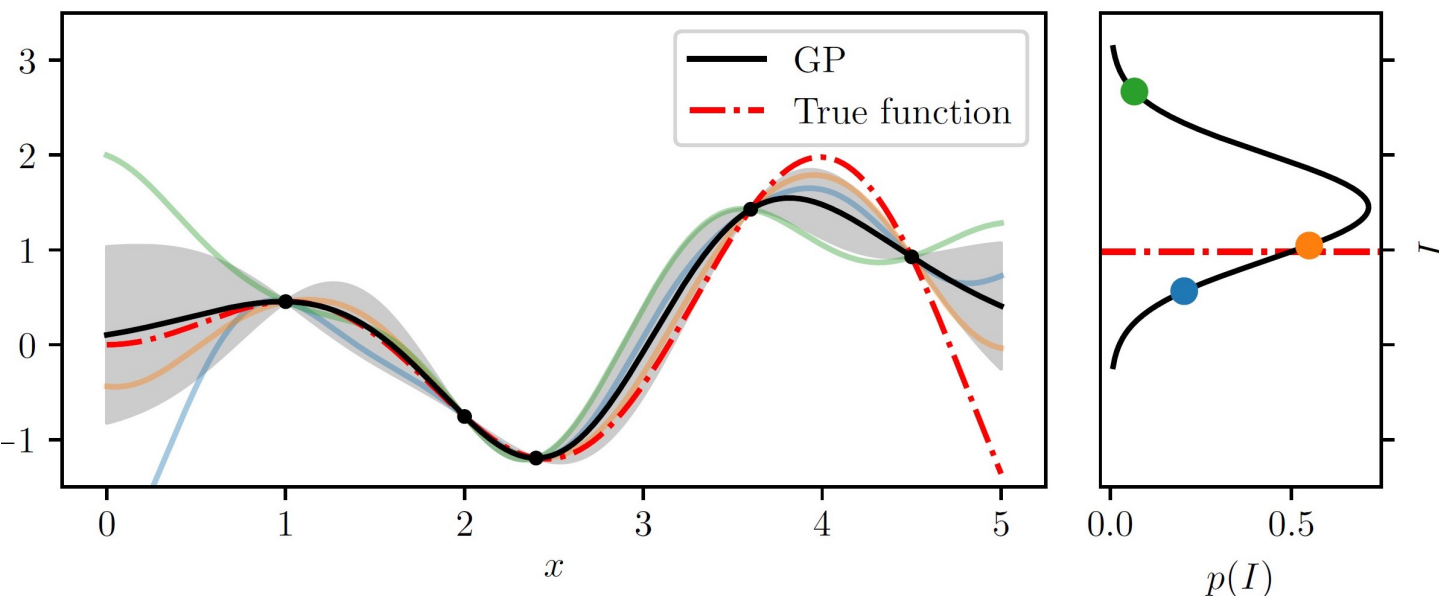
- To get marginalised quantities we want to integrate

$$\int L(x)\pi(x) dx$$

- With a GP we can get a model for $L(x)\pi(x) \sim \mathcal{GP}(0, k(x, x'))$

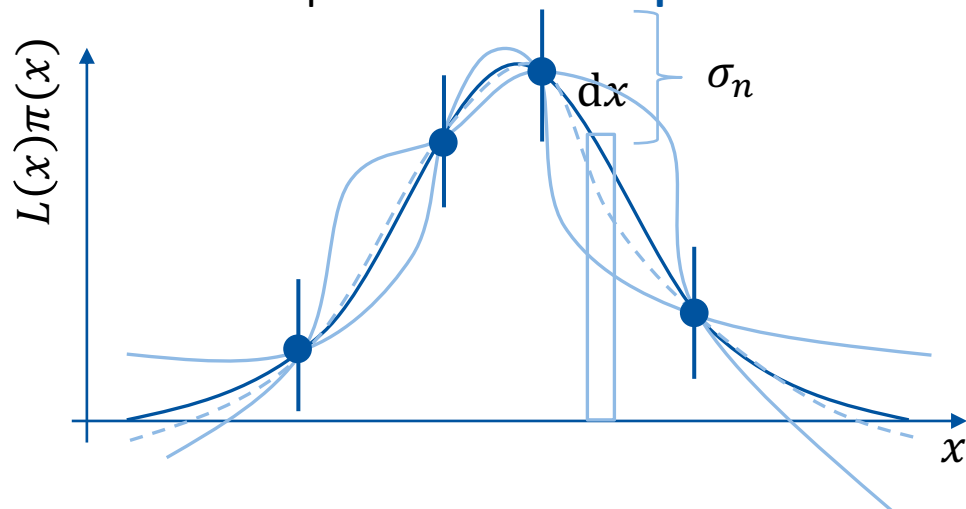
- We can integrate that model by integrating $\int \mu(x) dx = \int \bar{f}(x) dx$

- We can use $\mu(x)$ and $\sigma(x) = \sqrt{\text{cov}(f_*(x, x))}$ to find the next most informative point to sample



Active sampling

⇒ At each step maximize an **acquisition function**



$L(x)\pi(x)$ is **always positive**

$$\Rightarrow a(x) = \mu(x) \cdot \sigma(x)$$

$L(x)\pi(x)$ has **high dynamic range**

⇒ Sample log-posterior:

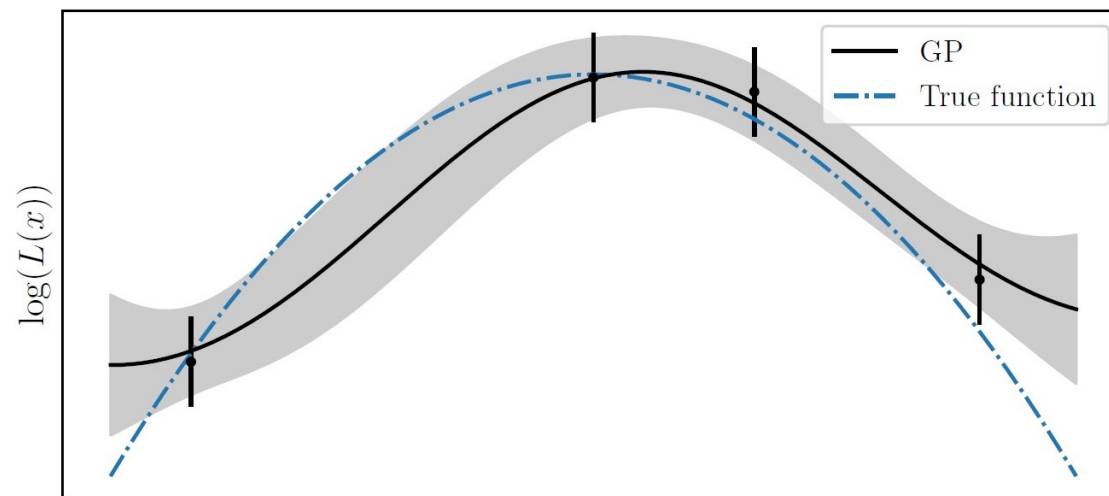
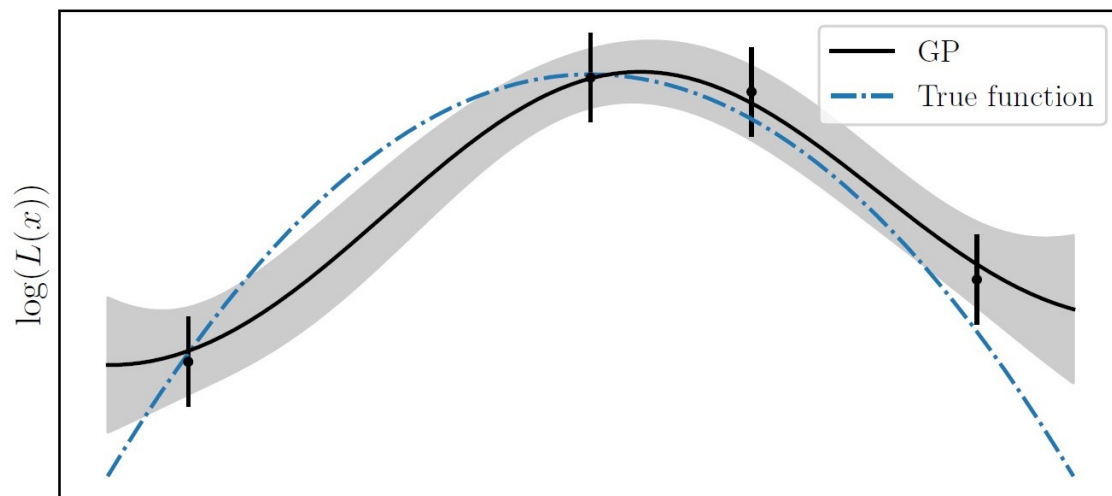
$$a(x) = \exp(2 \cdot \bar{\mu}) \cdot \sigma_{\bar{\mu}}(x)$$

$\bar{\mu}$ = Mean of GP fit to log-posterior

Correction factor ζ and statistical noise σ_n

$$a(x) = \exp(2\zeta \cdot \bar{\mu}) \cdot (\sigma_{\bar{\mu}}(x) - \sigma_n)$$

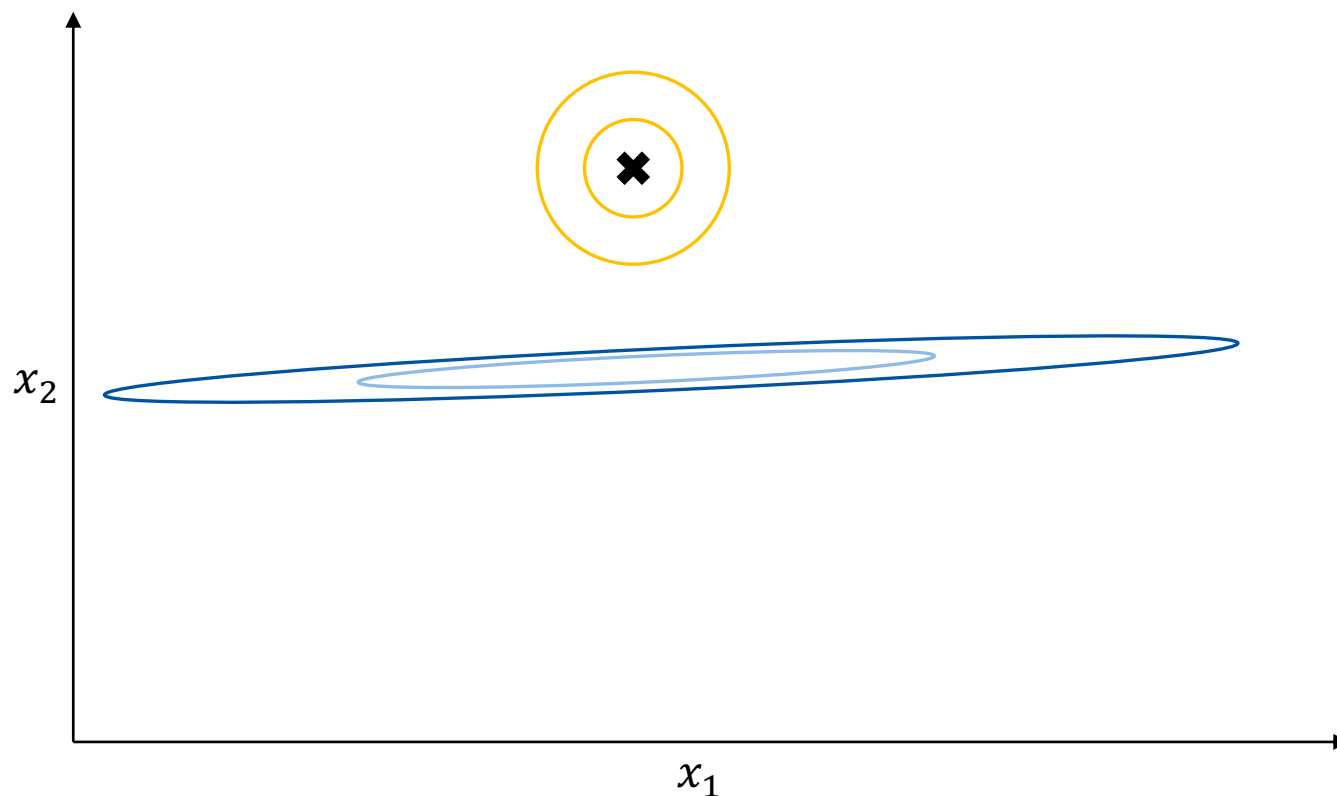
The acquisition function



Flat over large areas \Rightarrow We take the log of the acquisition function when actually optimizing it

Preprocessing

Problem 1: Different scales



Do two things:

1. Scale the priors such that they occupy the unit hypercube (every parameter is in $[0,1]$)
2. Make kernel asymmetric

$$k(x, x') = \sigma^2 \cdot \prod_{i=1}^d \exp\left(-\frac{(x_i - x'_i)^2}{2l^2}\right)$$

More hyperparameters to fit ($d + 1$) but robust!

Preprocessing

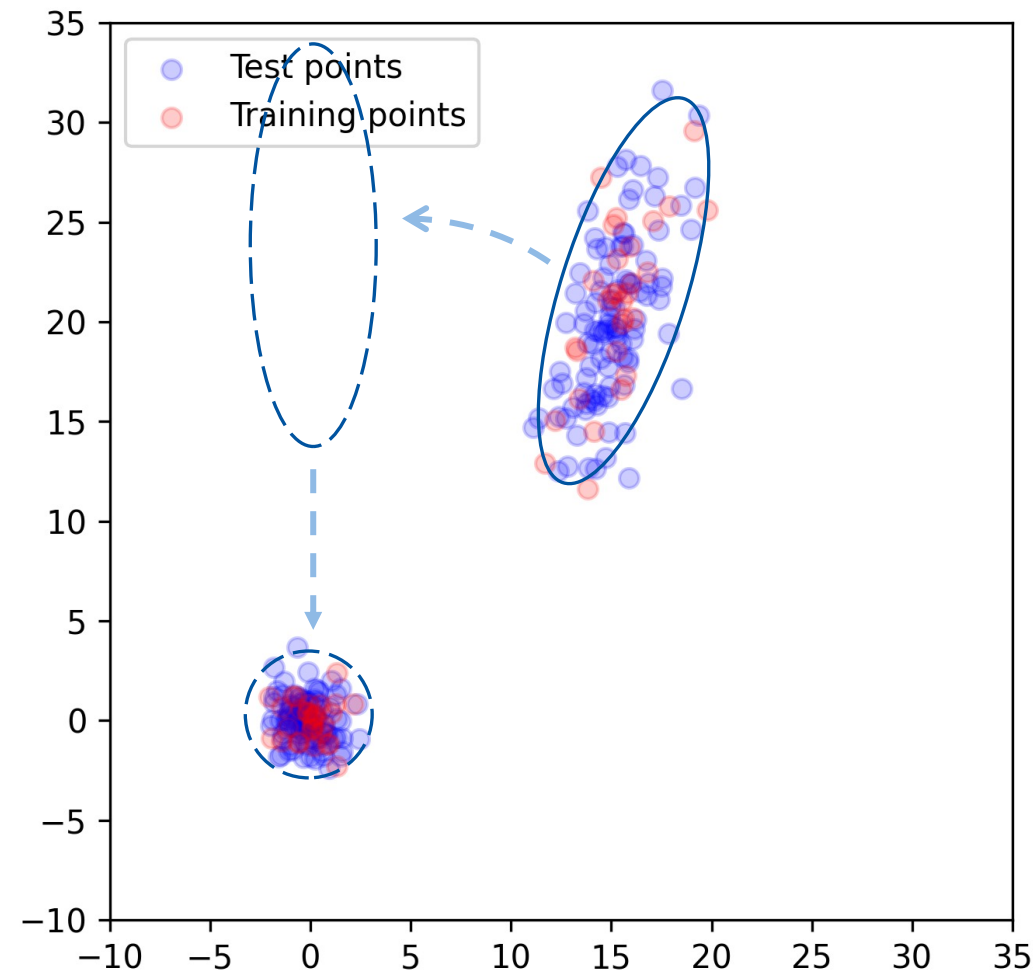
Alternative: Whitening

$$x_i \rightarrow x_i' = \frac{R_{ij}(x - \hat{\mu})_j}{\hat{\Sigma}_{ii}}$$

with

- $\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$
- $\hat{\Sigma}_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \hat{\mu}_j)(x_{ik} - \hat{\mu}_k)$
(empirical mean and covariance along each dimension)
- $\hat{\Sigma} = R\Lambda R^T$ with Λ diagonal

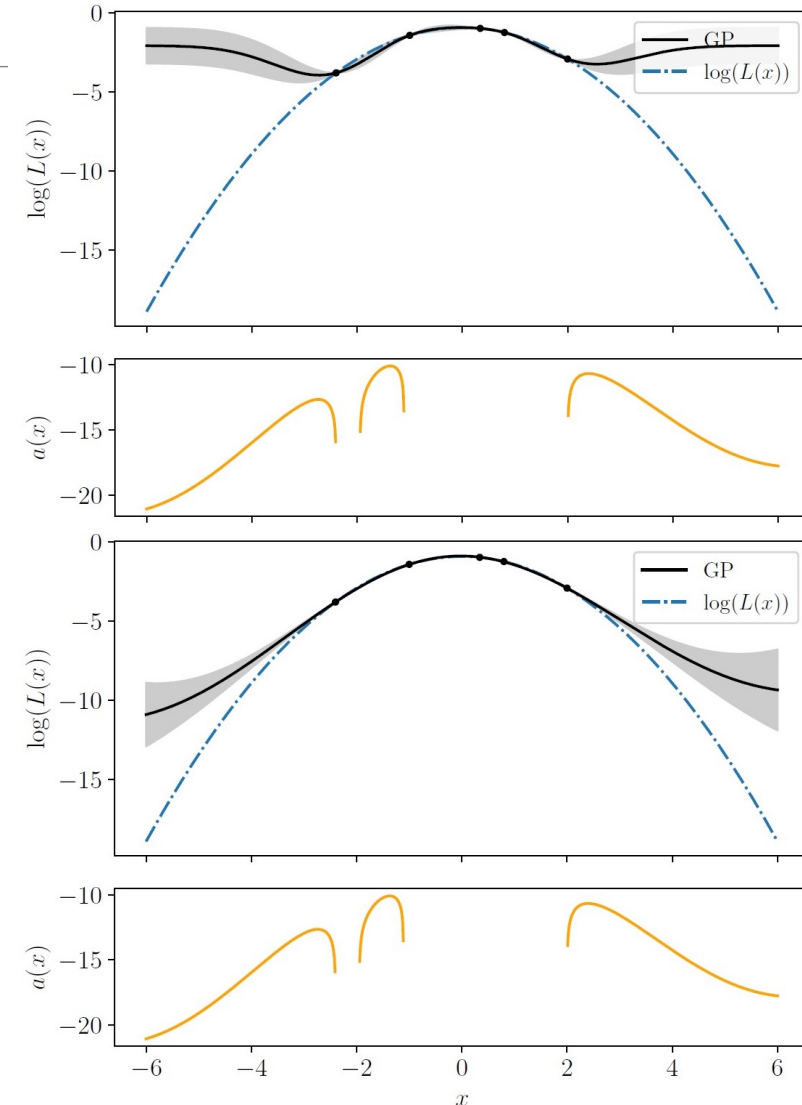
⇒ No need to make kernel asymmetric but less robust



Preprocessing

What about log-posterior values?

- Transform such that they have zero mean and unit variance
- Encourages exploration when lots of high values of the log-posterior
- Encourages exploitation when lots of low values of the log-posterior



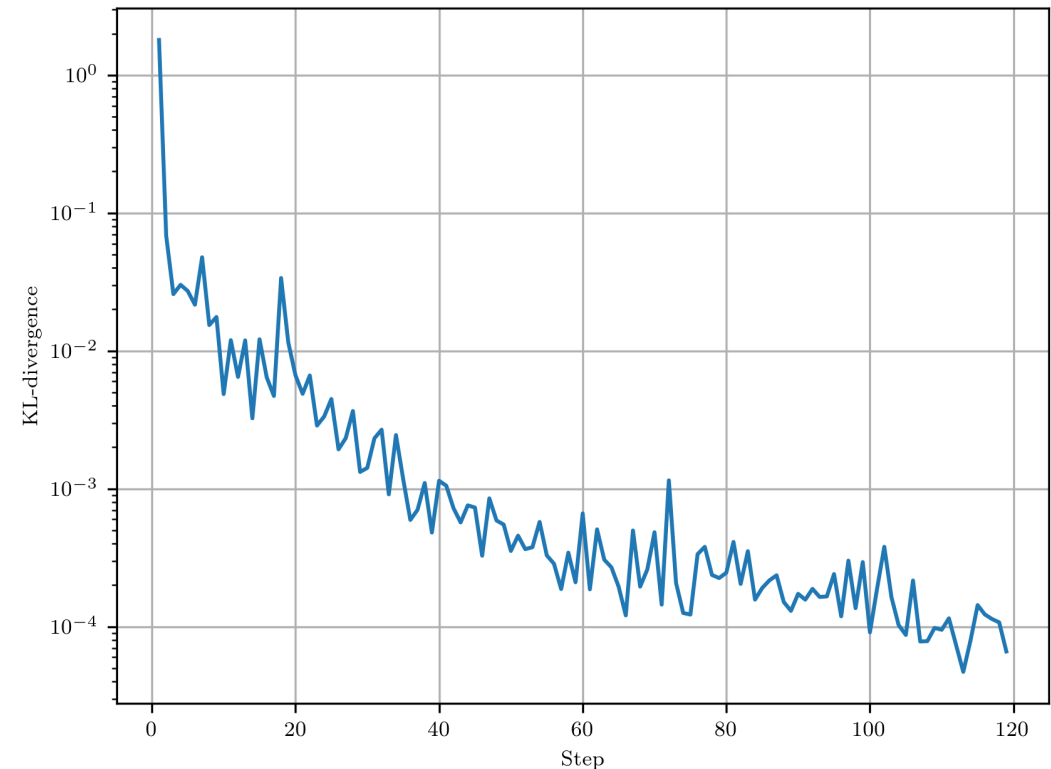
Kullback-Leibler divergence

Kullback-Leibler (KL) divergence:

$$D_{\text{KL}}(P_{n+1} || P_n) = \sum_{x \in \mathcal{X}} P_{n+1}(x) \log \left(\frac{P_{n+1}(x)}{P_n(x)} \right)$$

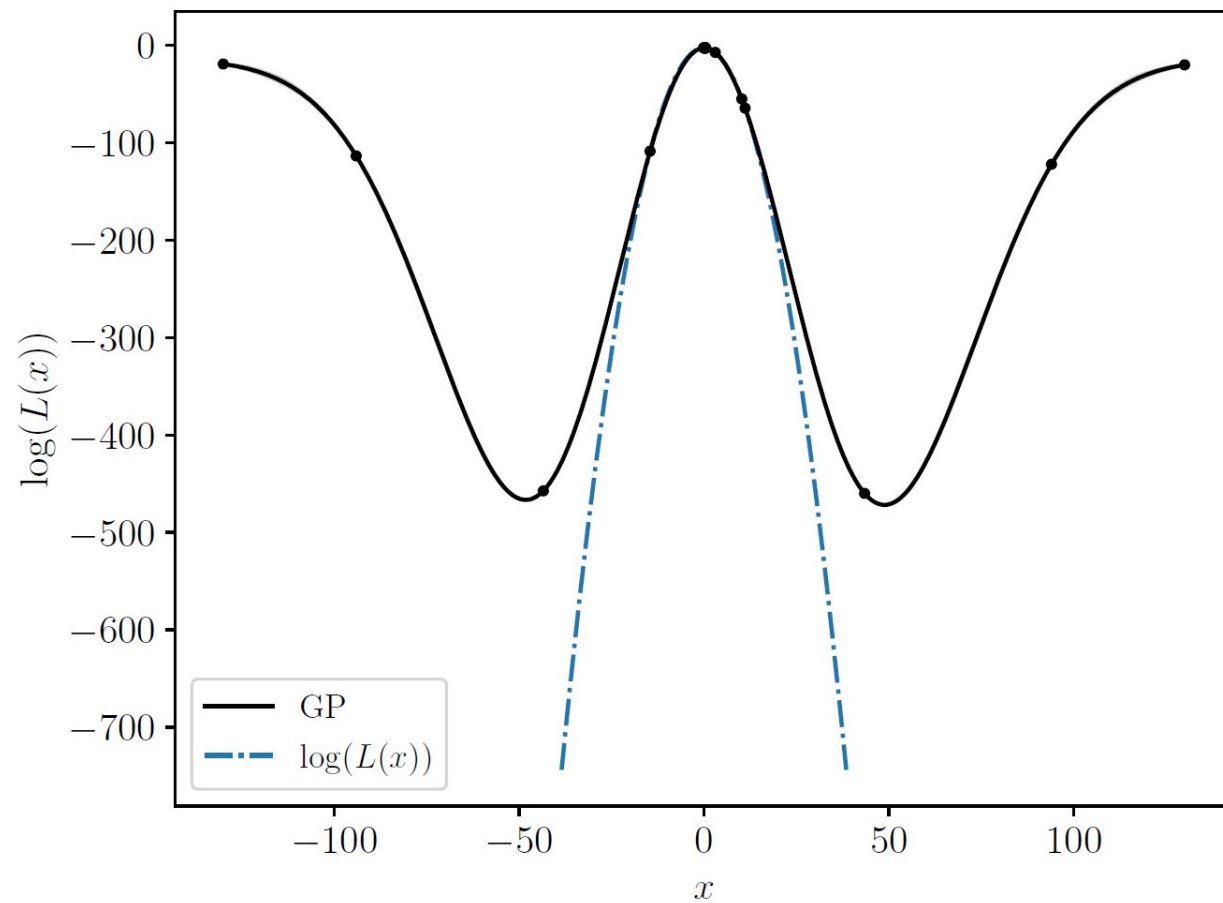
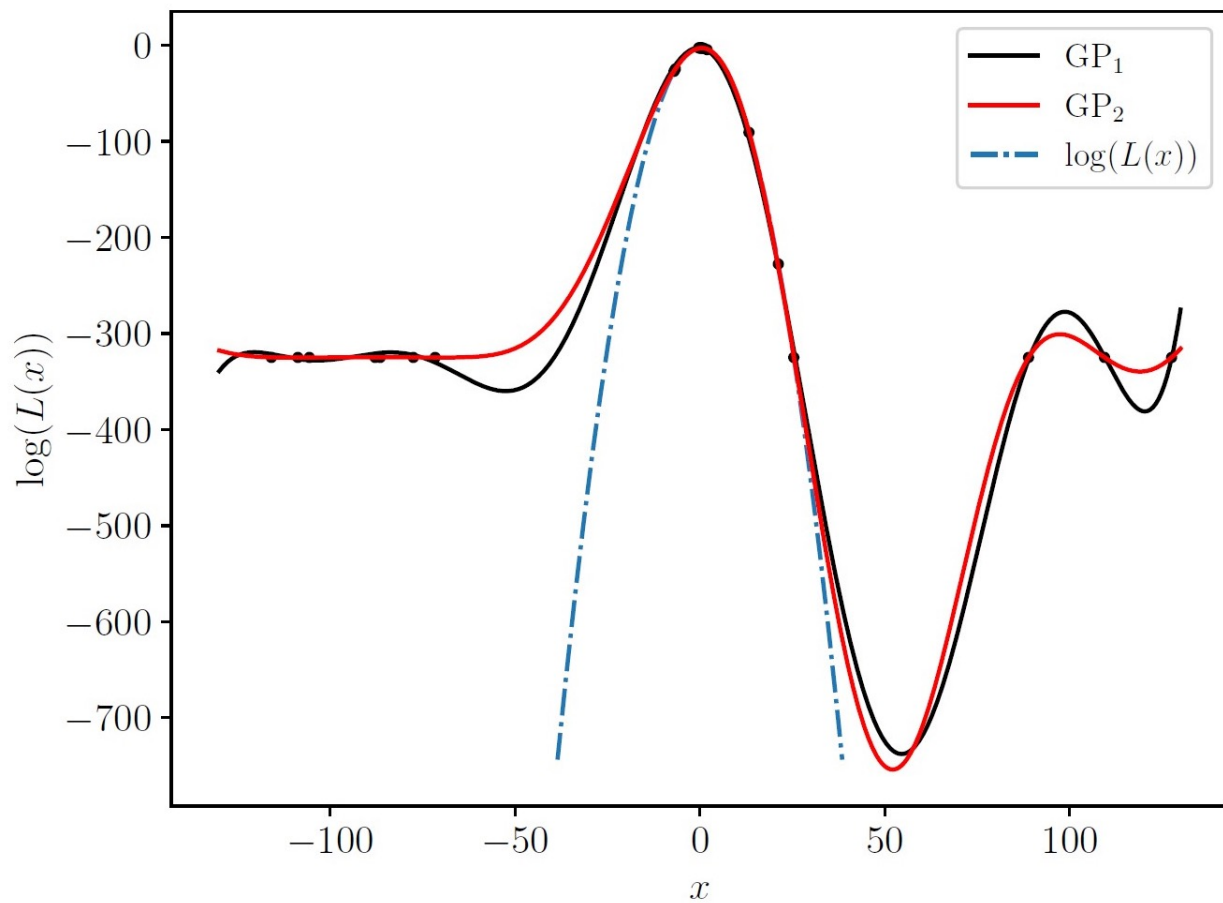
In case of a multivariate Gaussian this is just

$$D_{\text{KL}}(P || Q) = \frac{1}{2} \left[\log \frac{|\Sigma_q|}{|\Sigma_p|} - d + \text{tr} \left(\Sigma_q^{-1} \Sigma_p \right) + (\mu_q - \mu_p)^T \Sigma_q^{-1} (\mu_q - \mu_p) \right]$$



For now: Take empirical mean and covariance of the **training points**

The problem with infinity



Are we preserving Bayesianity?

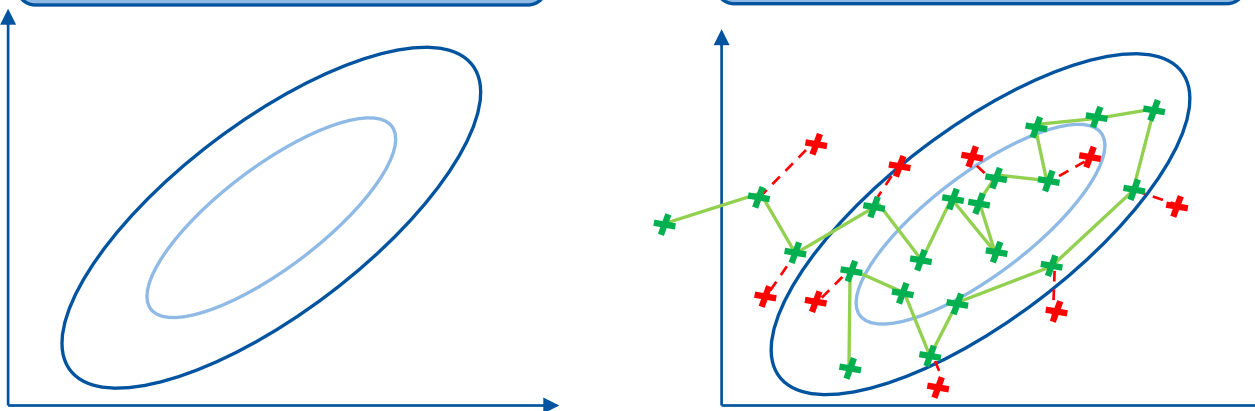
We are violating Bayesianity at **two points**:

$$\Rightarrow \log(p(y|X)) = -\frac{1}{2}y^T K^{-1}y - \frac{1}{2}\log|K| - \frac{n}{2}\log(2\pi)$$

We are maximizing this with **MLII**. Correct Bayesian way:
Sampling the posterior distribution but **very expensive!**

GP model $(\mu(x), \sigma(x))$

MCMC sampler



→ Ignoring $\sigma_{\text{GP}}(x)$

Correct:

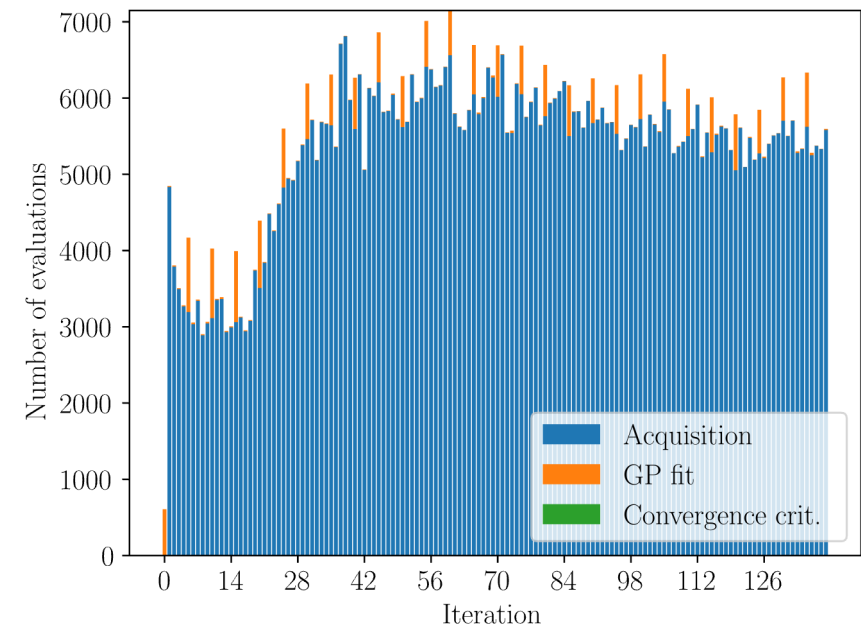
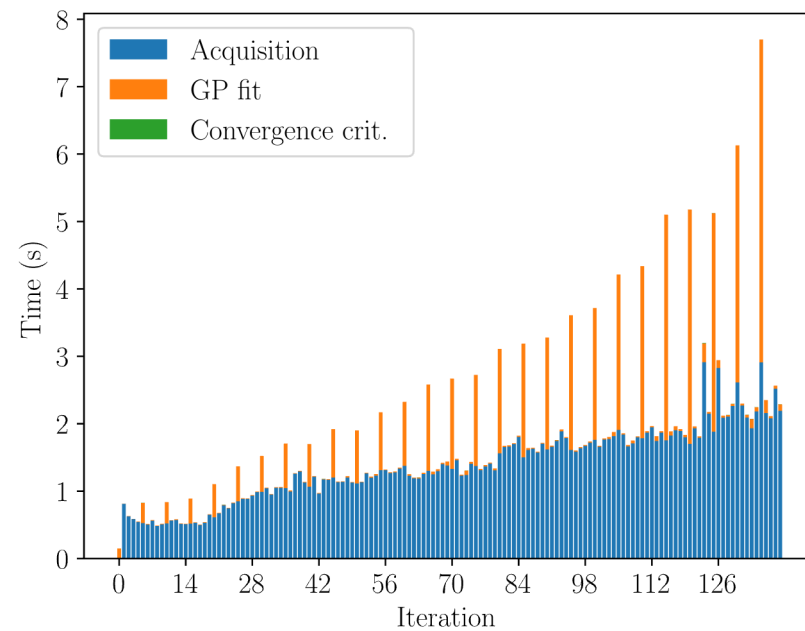
Get var(posterior) by solving

$$\iint k_{f,D}(x, x') dx dx'$$

But **super expensive!**

Overhead

8 dimensions
2 Kriging believer
steps/iteration
In total 300 accepted
samples



Refitting GP hyperparameters requires many inversions
of the kernel matrix, scales $\mathcal{O}(N_{\text{samples}}^3)$