

Gaussian Processes and Bayesian Optimization: An Introduction

Rama Vasudevan
Center for Nanophase Materials Sciences,

AI for atoms: How to machine learn STEM
December 9th, 2020

ORNL is managed by UT-Battelle, LLC for the US Department of Energy



U.S. DEPARTMENT OF
ENERGY

Gaussian Process Regression

- Splines, polynomials, other functional fits
 - Problem: Uncertainty depends on the model. Also, which model is the correct model?
- Bayesian approach:
 - Use Gaussian Processes
- A Gaussian process is a non-parametric Bayesian approach to regression, that finds a distribution over functions $f(x)$ that are consistent with the observed data
 - The similarity between points is defined by a covariance matrix
 - The covariance is determined by a kernel function

Gaussian Process Regression

- Why bother with multivariate normal distributions? Two reasons:
 - (1) Easy to marginalize: when you have many, if you want to limit yourself to some subset, it is simple to integrate out the variables you don't want
 - (2) Easy to condition: Can write analytical solutions for conditioning. The conditioned distribution is also normal.

Multivariate Gaussian Distribution

$$y \sim N(\mu, \Sigma)$$

Mean vector, covariance matrix

Generalizes to a Gaussian Process

$$f \sim GP(m(x), k(x, x'))$$

Mean function, covariance function

Covariance Function returns the Covariance Matrix

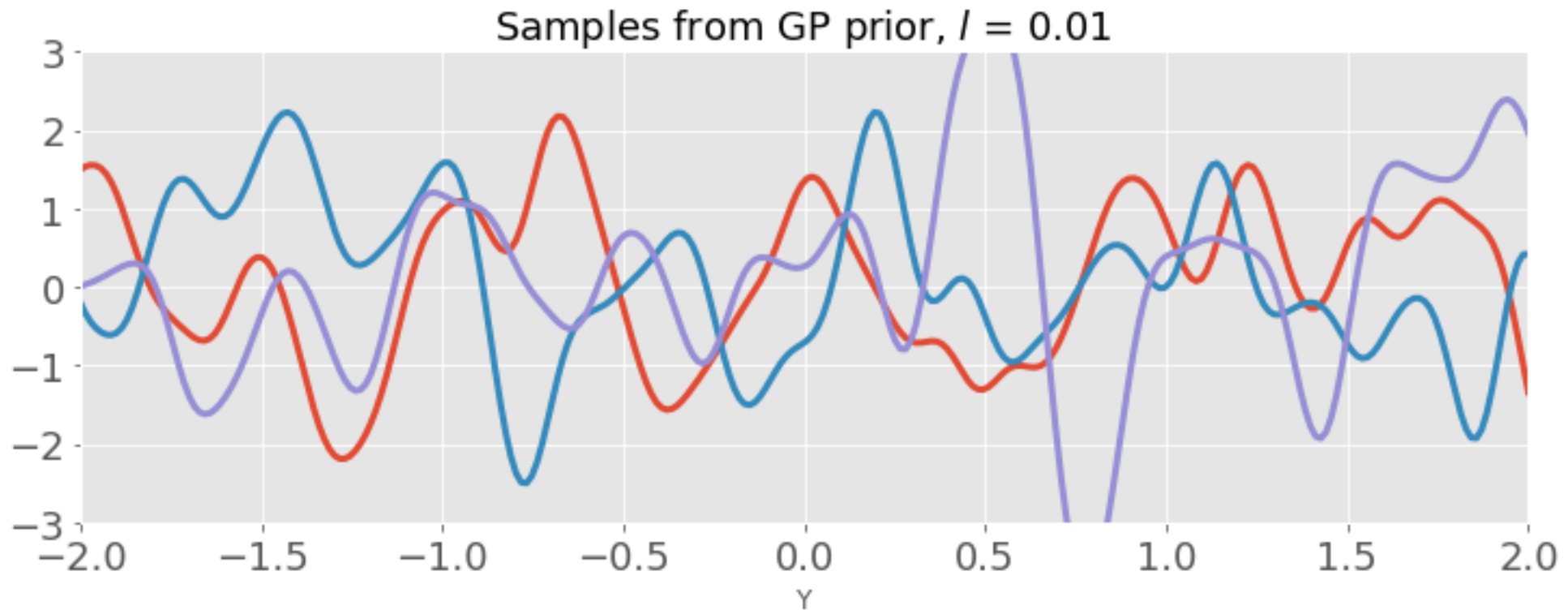
$$\Sigma \sim k(x, x' | \Phi)$$

Pass function values, return Cov. Mat.

Gaussian Process Regression

- Covariance matrix determines what type of functions we will allow.

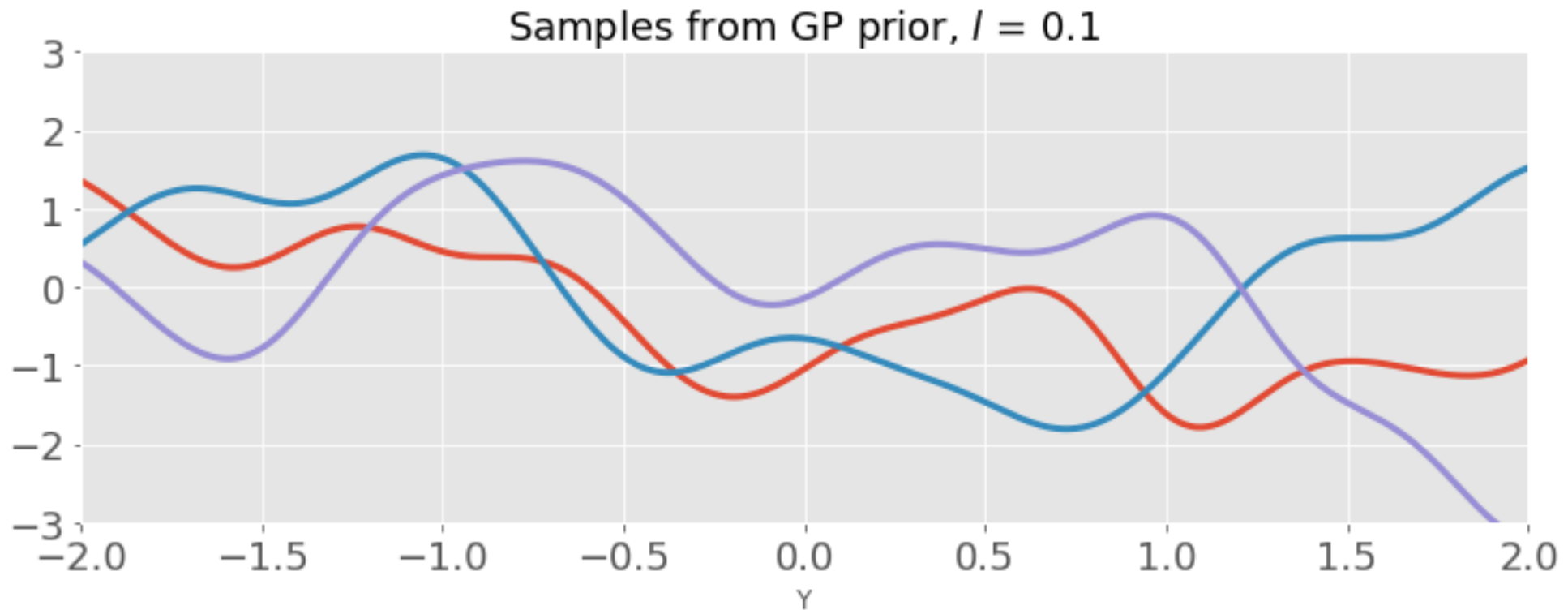
$$k(x, x') = \exp\left(-\frac{1}{2l}(x - x')^2\right)$$



Gaussian Process Regression

- Covariance matrix determines what type of functions we will allow.

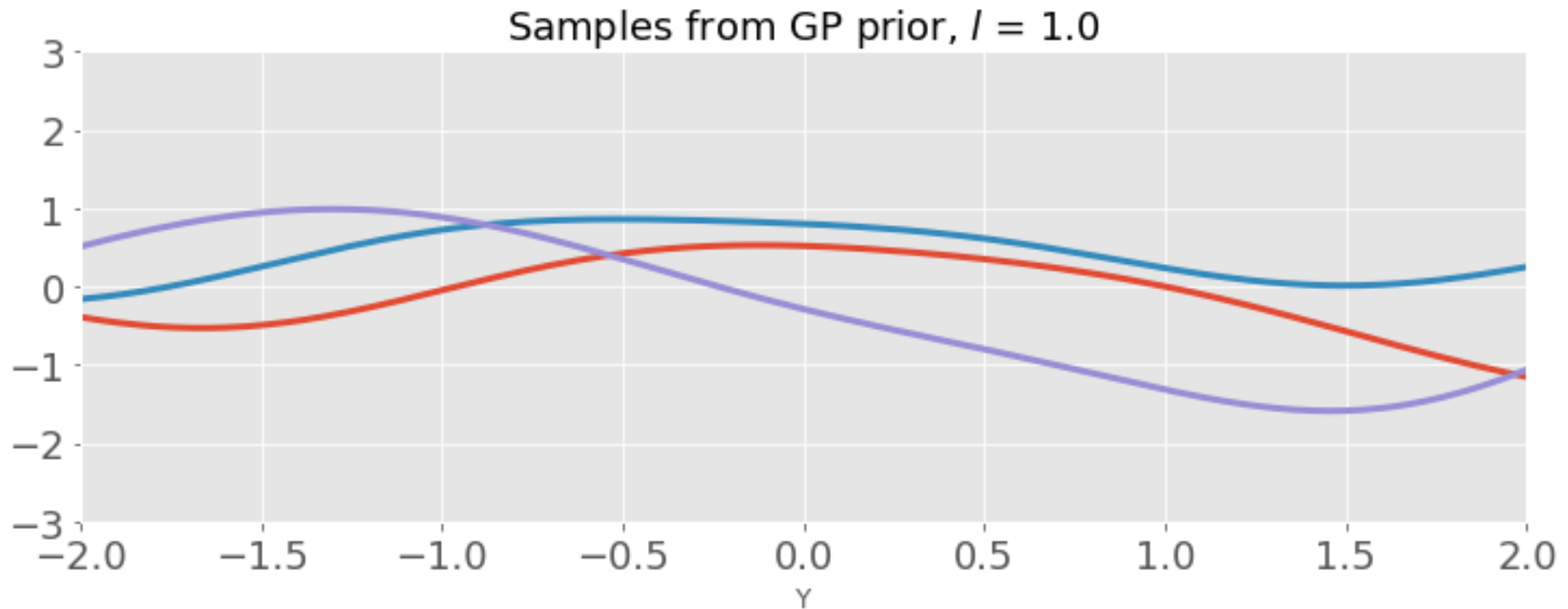
$$k(x, x') = \exp\left(-\frac{1}{2l}(x - x')^2\right)$$



Gaussian Process Regression

- Covariance matrix (kernel) determines what type of functions we will allow.

$$k(x, x') = \exp\left(-\frac{1}{2l}(x - x')^2\right)$$



l controls the lengthscale – sort of how far points should be to make them independent of each other.

Gaussian Process Regression

- Consider a (gaussian) joint probability distribution with two variables, as follows

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \right)$$

Variance in x_1 (points to σ_{11})

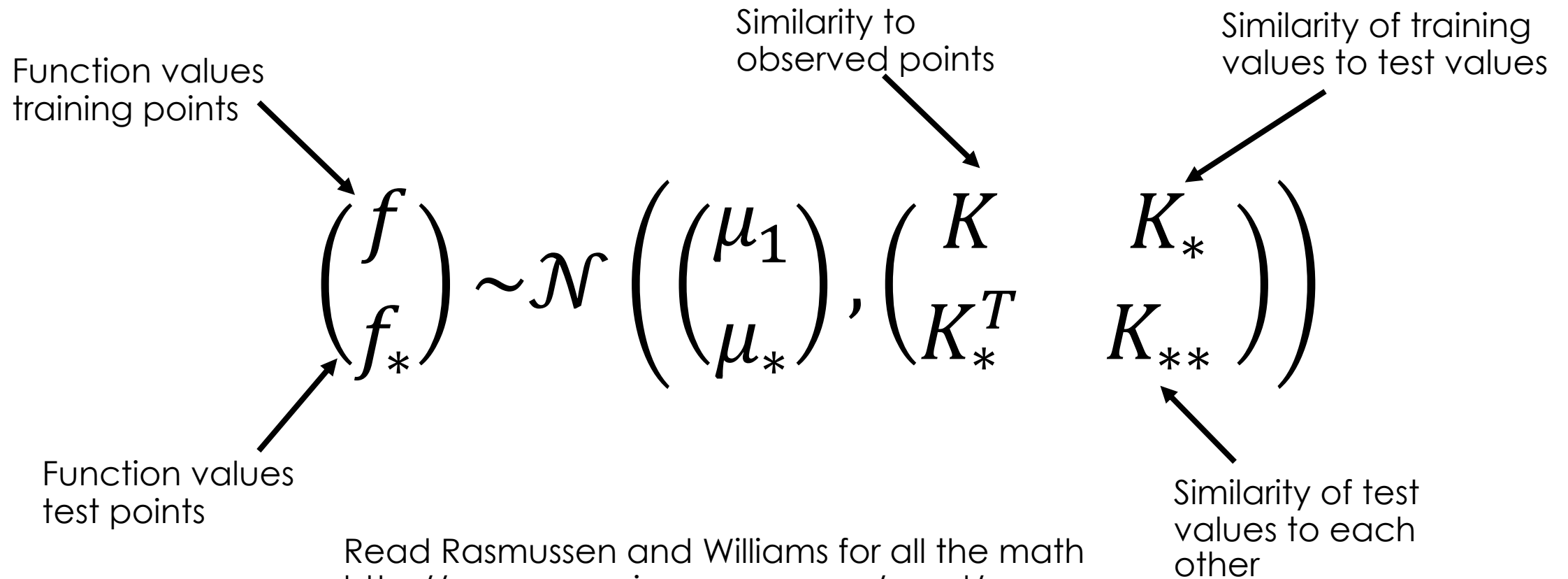
Similarity of x_1 values to x_2 values (points to σ_{12})

Variance in x_2 (points to σ_{22})

Nonzero entries off-diagonal indicate the correlations

Gaussian Process Regression

- Generalize this to the function space
- When we observe the data, we now have to calculate the posterior over the functions. This posterior is a joint distribution over function values observed and those not observed.



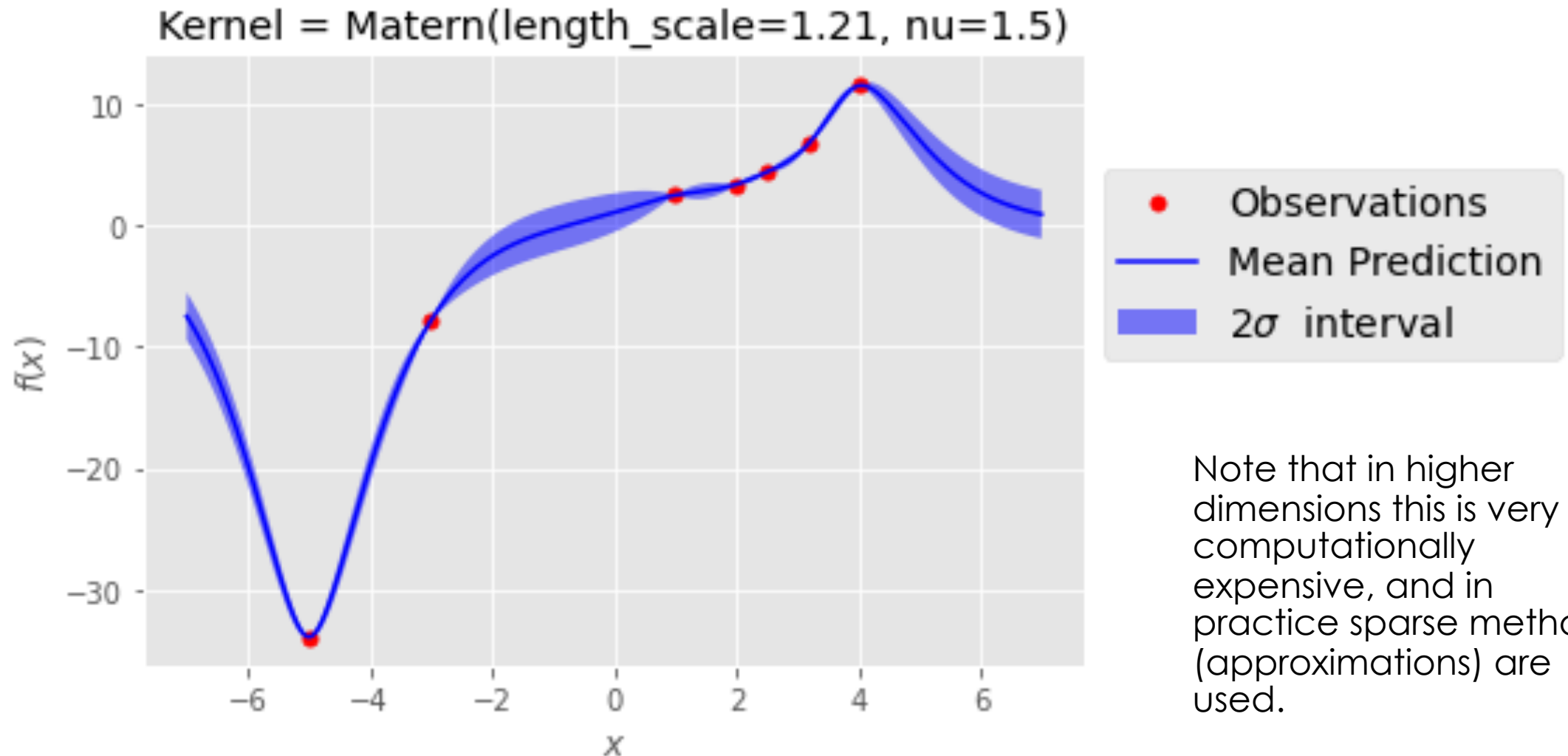
The diagram illustrates the Gaussian Process Regression equation with several annotations:

- Function values training points**: Points to the f component of the vector $\begin{pmatrix} f \\ f_* \end{pmatrix}$.
- Function values test points**: Points to the f_* component of the vector $\begin{pmatrix} f \\ f_* \end{pmatrix}$.
- Similarity to observed points**: Points to the K component of the covariance matrix $\begin{pmatrix} K & K_*^T \\ K_* & K_{**} \end{pmatrix}$.
- Similarity of training values to test values**: Points to the K_* component of the covariance matrix.
- Similarity of test values to each other**: Points to the K_{**} component of the covariance matrix.

$$\begin{pmatrix} f \\ f_* \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_1 \\ \mu_* \end{pmatrix}, \begin{pmatrix} K & K_*^T \\ K_* & K_{**} \end{pmatrix} \right)$$

Read Rasmussen and Williams for all the math
<http://www.gaussianprocess.org/gpml/>

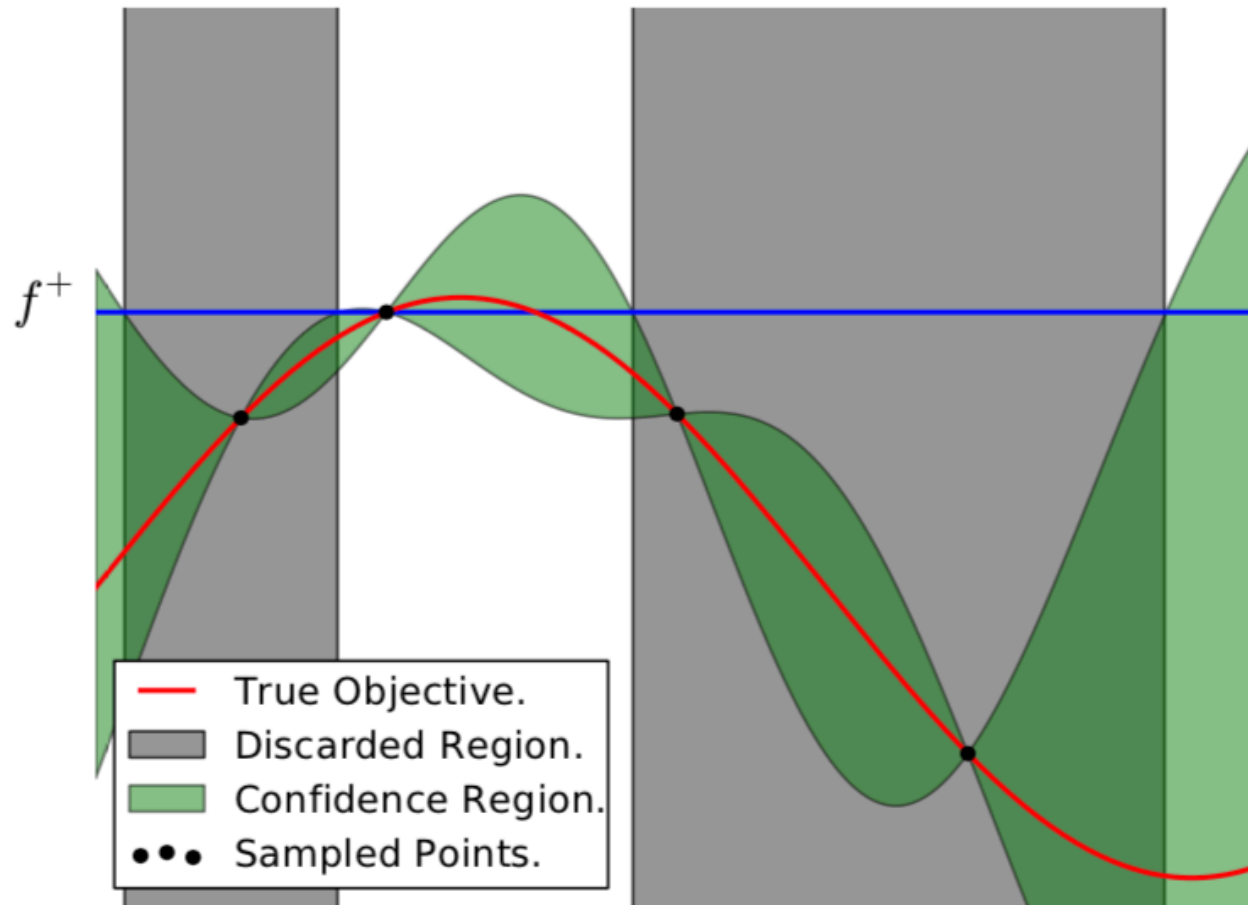
Gaussian Process Regression



Gaussian Process Regression

Open (01_Gaussian_Processes.ipynb)

Bayesian Optimization

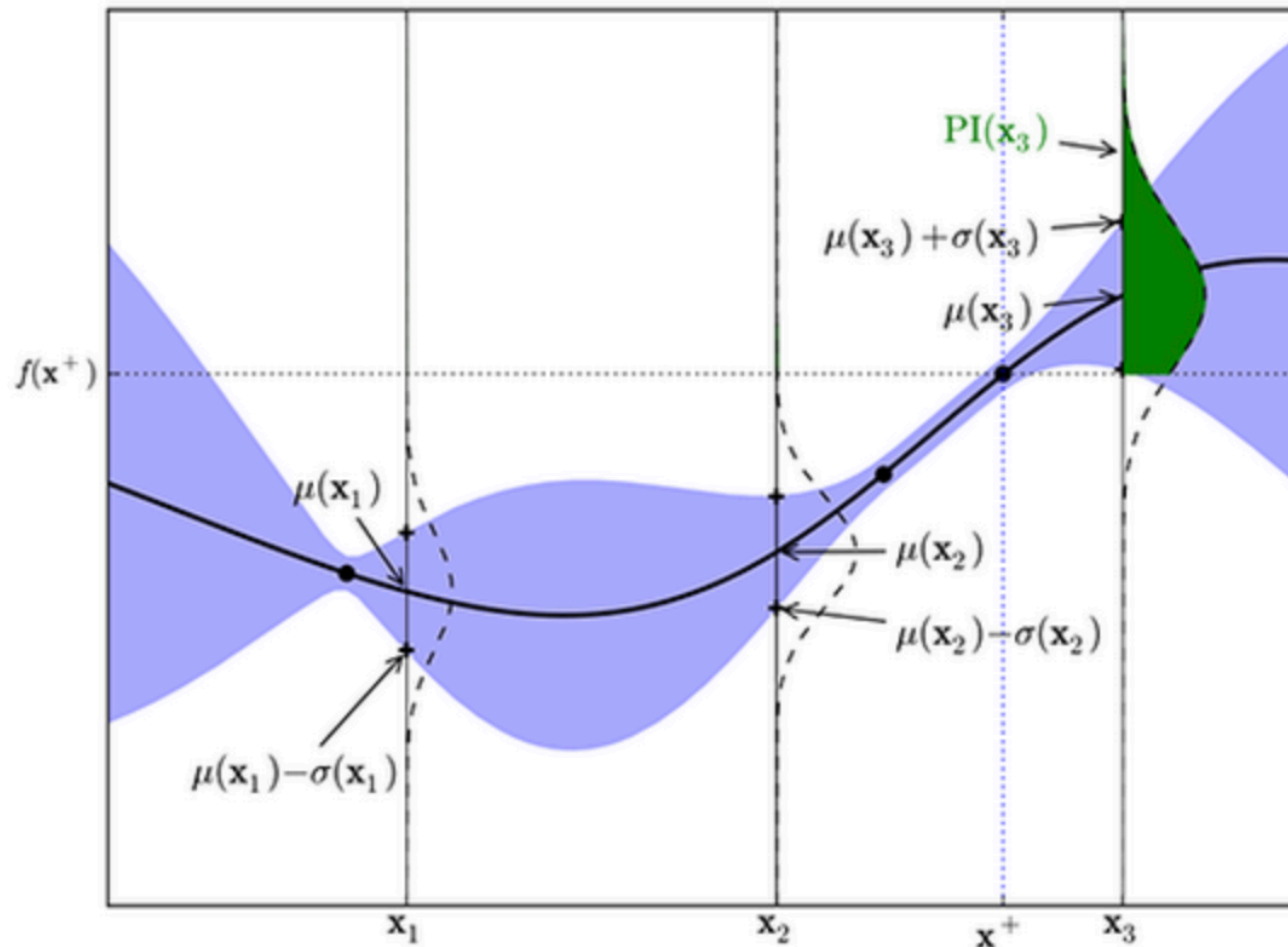


- We have some measurements in space X , and we want to maximize some property $f(X)$.
- How can we decide what point to measure next to best maximize f ?
- We need to balance the exploration of the space with exploitation of regions near we have already know

N. de Freitas et al., Taking the Human Out of the Loop: A Review of Bayesian Optimization ,
Proceedings of the IEEE **104**, 148 (2015)

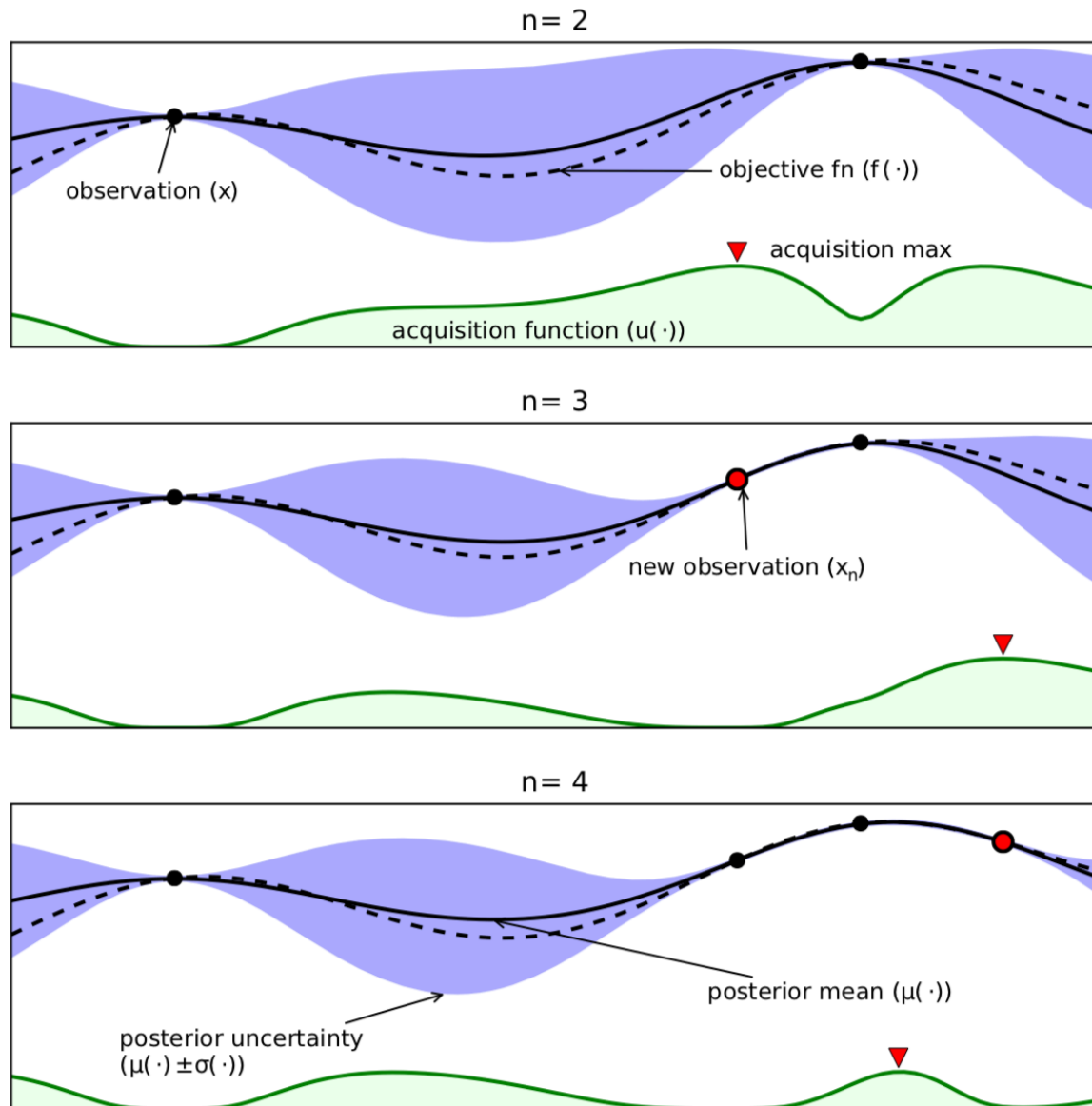
Acquisition Functions

Probability of Improvement Acquisition Function



1. **Confidence bound:** simplest possible - just take the upper confidence bound from the prediction
2. **Probability of Improvement:** Integral from current functional maximum to upper limit of distribution as test point
3. **Expected Improvement:** Instead of probability of improvement, we want to maximize the expected increase in the function value
4. **There are (always) more...**

“Active Learning”



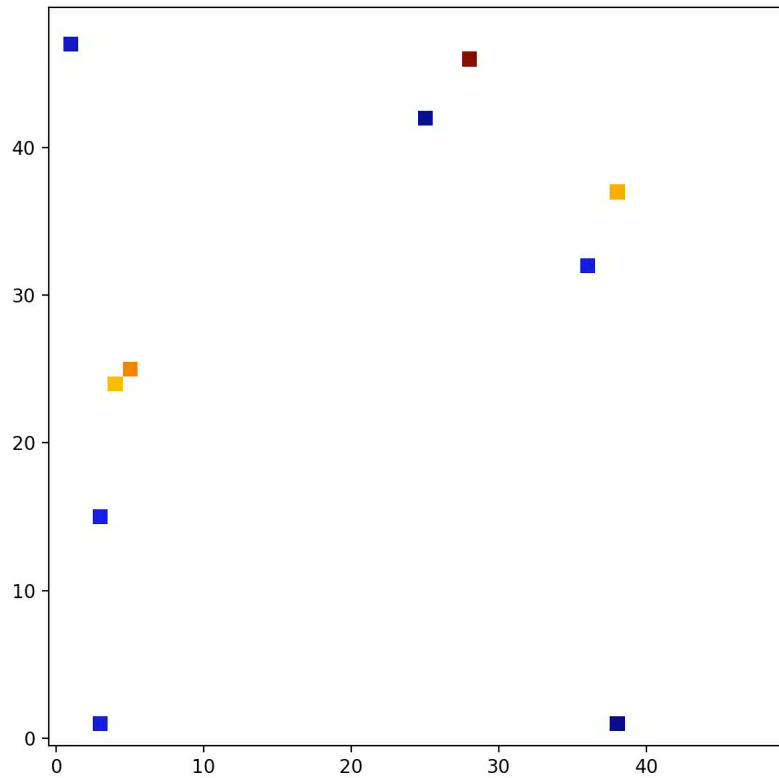
- Shown here are three iterations of Bayesian optimization
- We can do this for dictating how to collect new observations in our experiment

N. de Freitas et al., Taking the Human Out of the Loop: A Review of Bayesian Optimization, *Proceedings of the IEEE* **104**, 148 (2015)

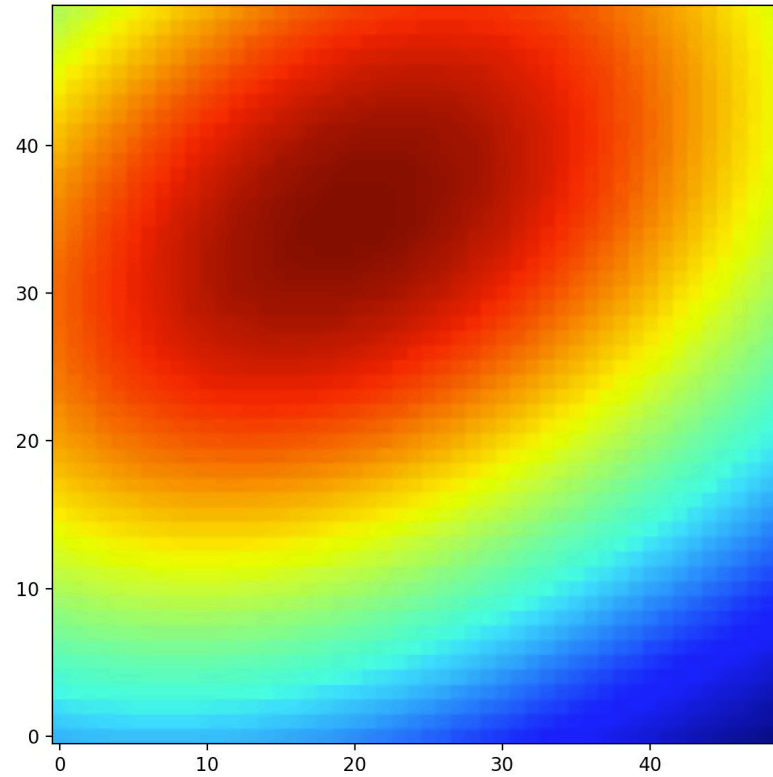
Example: Hysteresis Loops Area Maximization

Batch size = 10. 2.35s/pixel, and each iteration took ~6s on dgx. 410 points in total (40 GP iterations).
1200s in total. Total time for all pixels: 5875s (20% the time)

Measured Loop Areas



GP Prediction



GP Uncertainty

