

# Bayesian algorithms for practical accelerator control

Ryan Roussel

11/1/2022

*rroussel@slac.stanford.edu*



U.S. DEPARTMENT OF  
**ENERGY**

Stanford  
University

**SLAC** NATIONAL  
ACCELERATOR  
LABORATORY

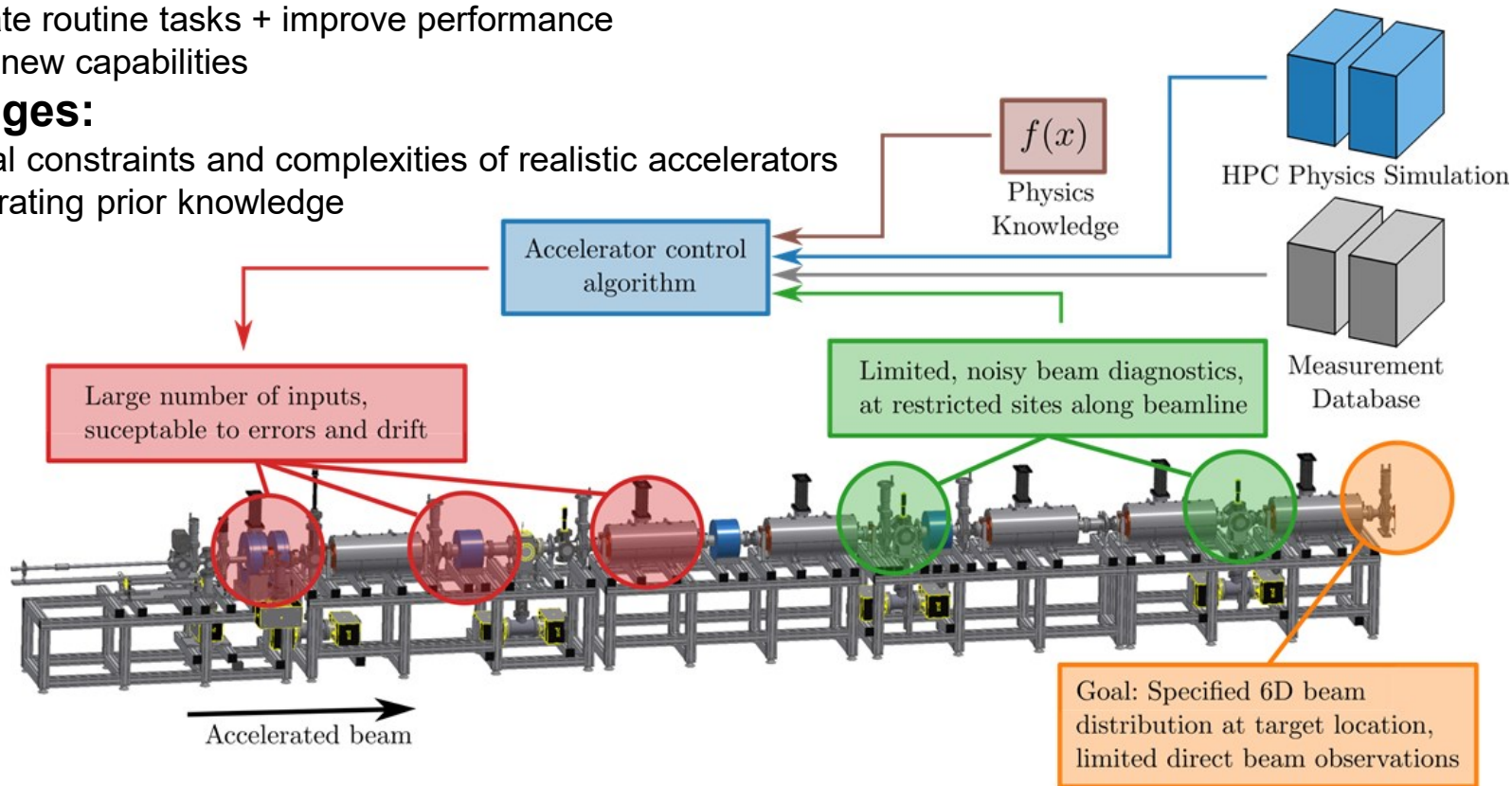
# Machine Learning Based Accelerator Control

## Goals:

- Automate routine tasks + improve performance
- Enable new capabilities

## Challenges:

- Practical constraints and complexities of realistic accelerators
- Incorporating prior knowledge
- Scaling



Can you calculate  
gradients easily?

- Analytical models

Yes

Gradient descent  
(SGD, Adam etc.)

Scales to >10k parameters  
(see ML training)

No

- Simulations  
- Experiments

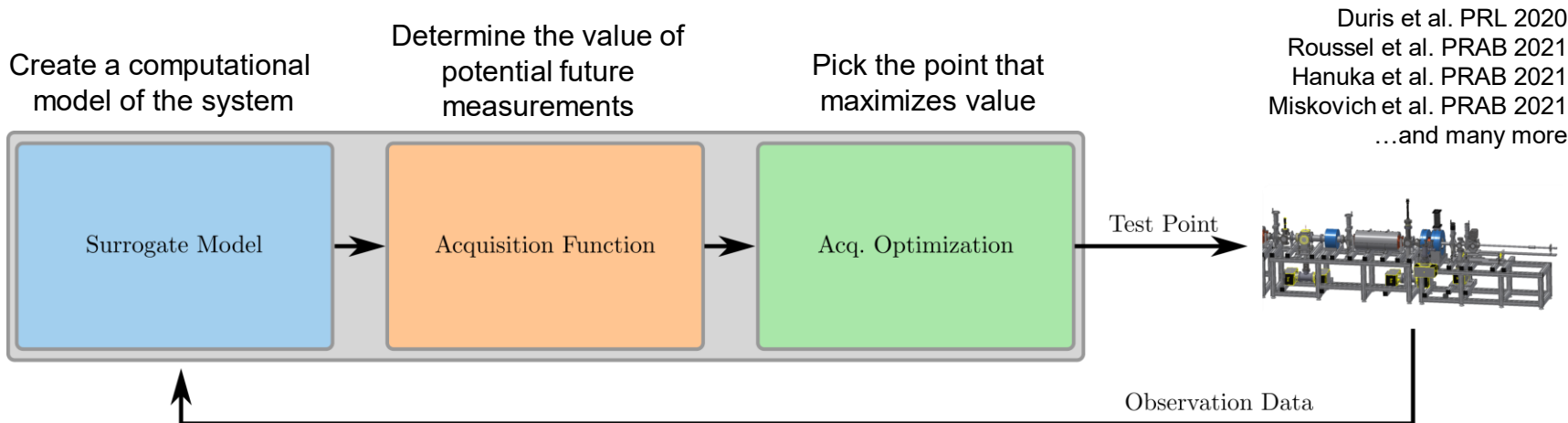
Black box  
optimization  
algorithms

Scales poorly with  
input dimension if not unimodal

Go ahead, try it with simplex...

# Model Based Optimization of Accelerators

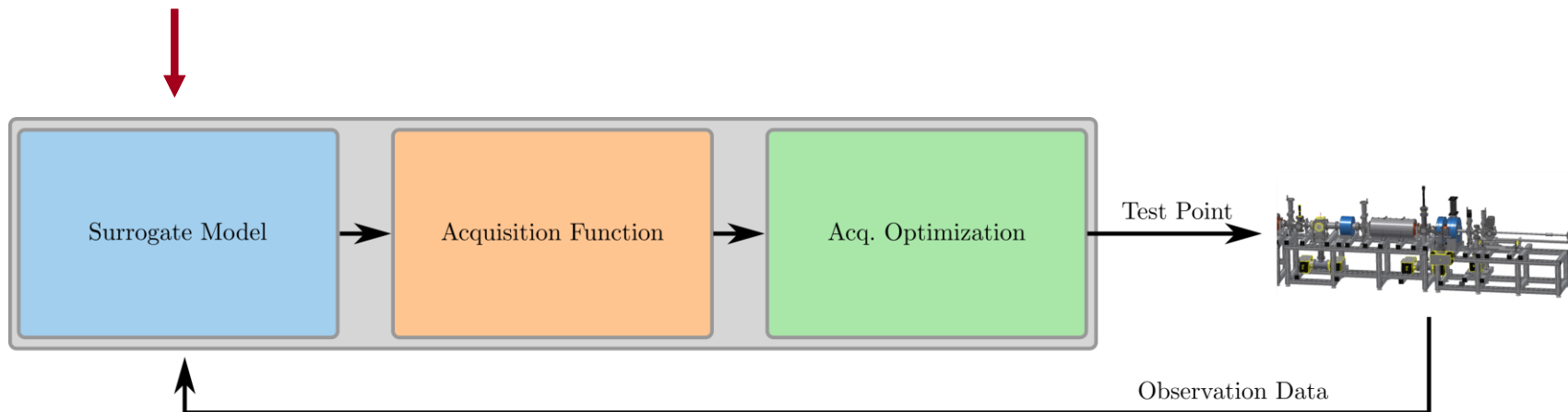
Can we improve black box optimization with surrogate models? **YES**



... **but** we must consider the cost of creating models with enough information

# Model Based Optimization of Accelerators

Gaussian processes (GPs)



Why?

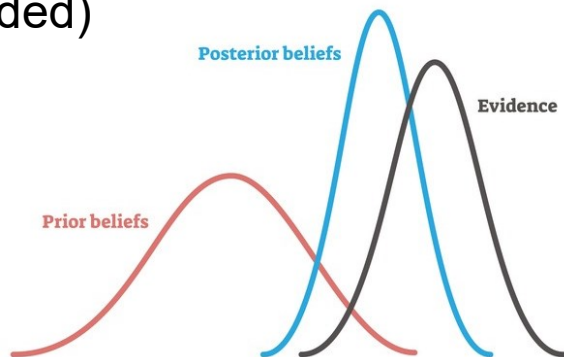
- Extracts a lot of information from a small number of data points → efficient
- Produces explicit uncertainty estimations -> advantageous for global optimization

# Bayesian Statistics: A Practical Example

$$p(B|A) = \frac{\text{likelihood } p(A|B) \text{ prior } p(B)}{\text{marginal likelihood } p(A)}$$

**A:** our measurement  
(we measure that the sun  
has/has not exploded)

**B:** our prediction  
(the sun has/has  
not exploded)



DID THE SUN JUST EXPLODE?  
(IT'S NIGHT, SO WE'RE NOT SURE.)



FREQUENTIST STATISTICIAN:

THE PROBABILITY OF THIS RESULT HAPPENING BY CHANCE IS  $\frac{1}{36} = 0.027$ . SINCE  $p < 0.05$ , I CONCLUDE THAT THE SUN HAS EXPLODED.



BAYESIAN STATISTICIAN:

BET YOU \$50 IT HASN'T.



# Bayesian Statistics: A (More) Practical Example

Model:

$$y = w_0 + w_1 x + \epsilon$$

Gaussian noise

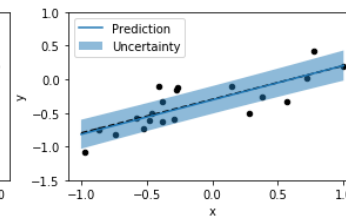
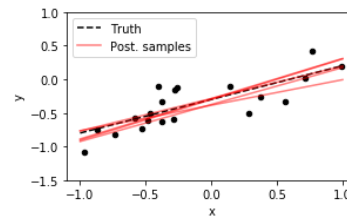
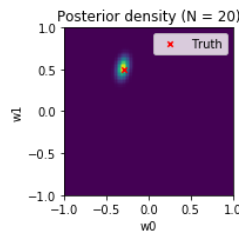
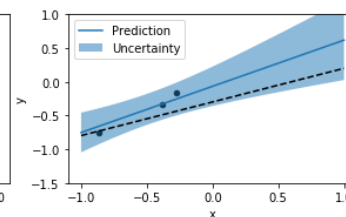
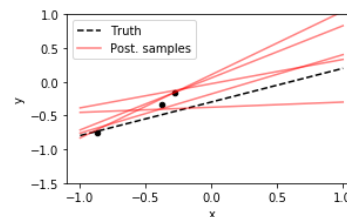
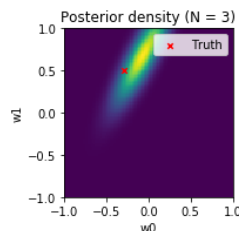
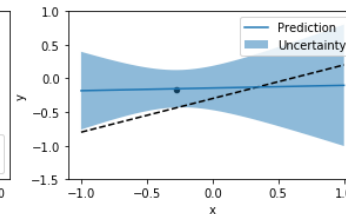
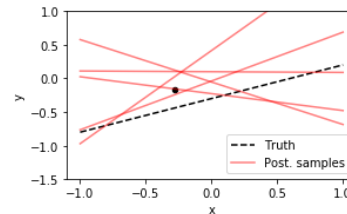
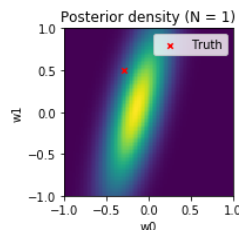
Bayesian regression  
(determining the weights):

$$p(w_0, w_1 | \mathbf{y}_N, \mathbf{x}_N) \propto p(\mathbf{y}_N | w_0, w_1, \mathbf{x}_N) p(w_0, w_1)$$

Making predictions:

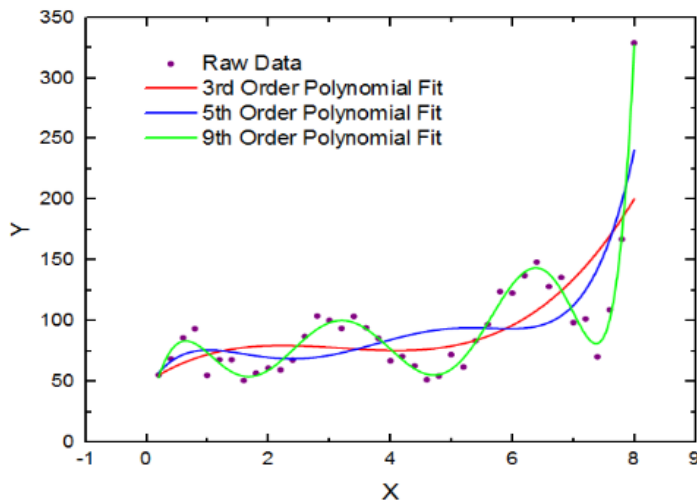
$$p(\mathbf{y}_M | w_0, w_1, \mathbf{x}_M)$$

Note: least squares fitting is equivalent to using a uniform prior and Gaussian likelihood



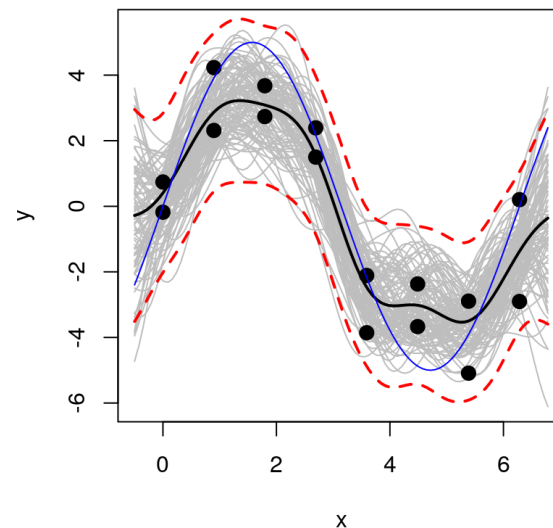
# Parametric vs. Non-Parametric Modeling

## Parametric modeling



$$f(x) = f(x; \theta)$$

## Non-Parametric modeling

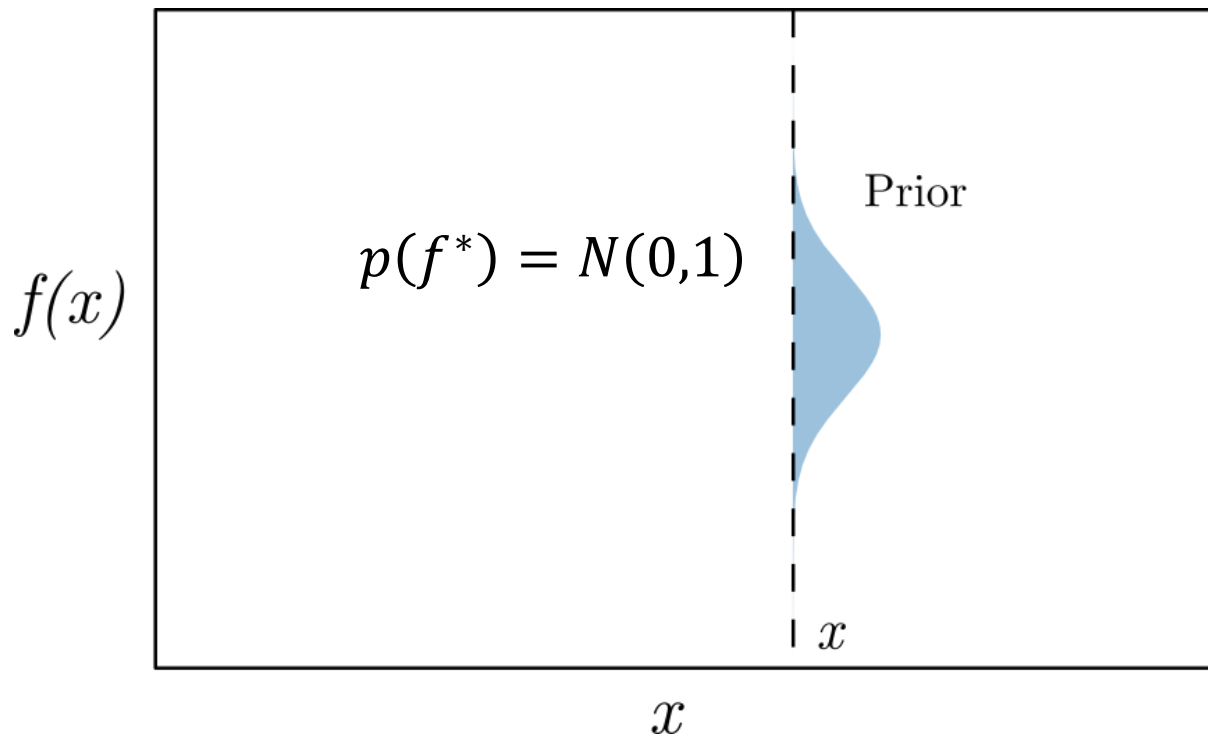


$$f = \{\theta_0(x_0), \theta_1(x_1), \dots, \theta_N(x_N)\}$$



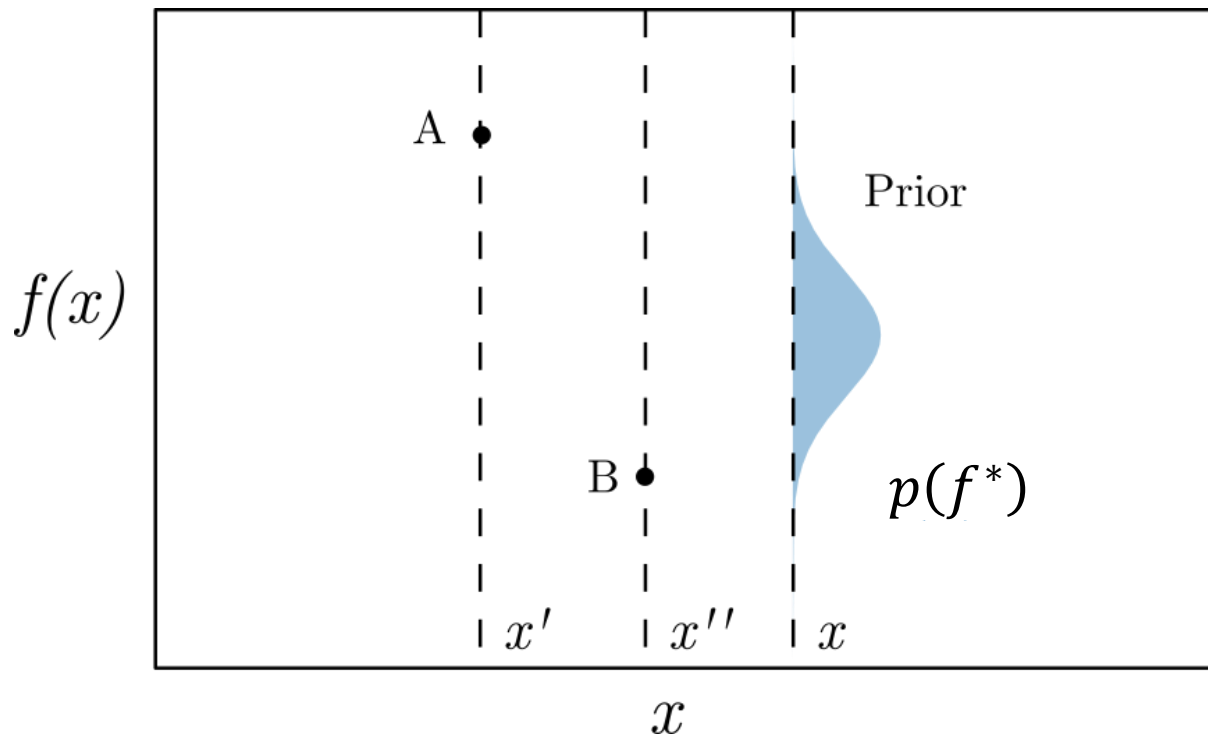
# Some intuition...

Let's predict the function value  $f^*$  at the point  $x$



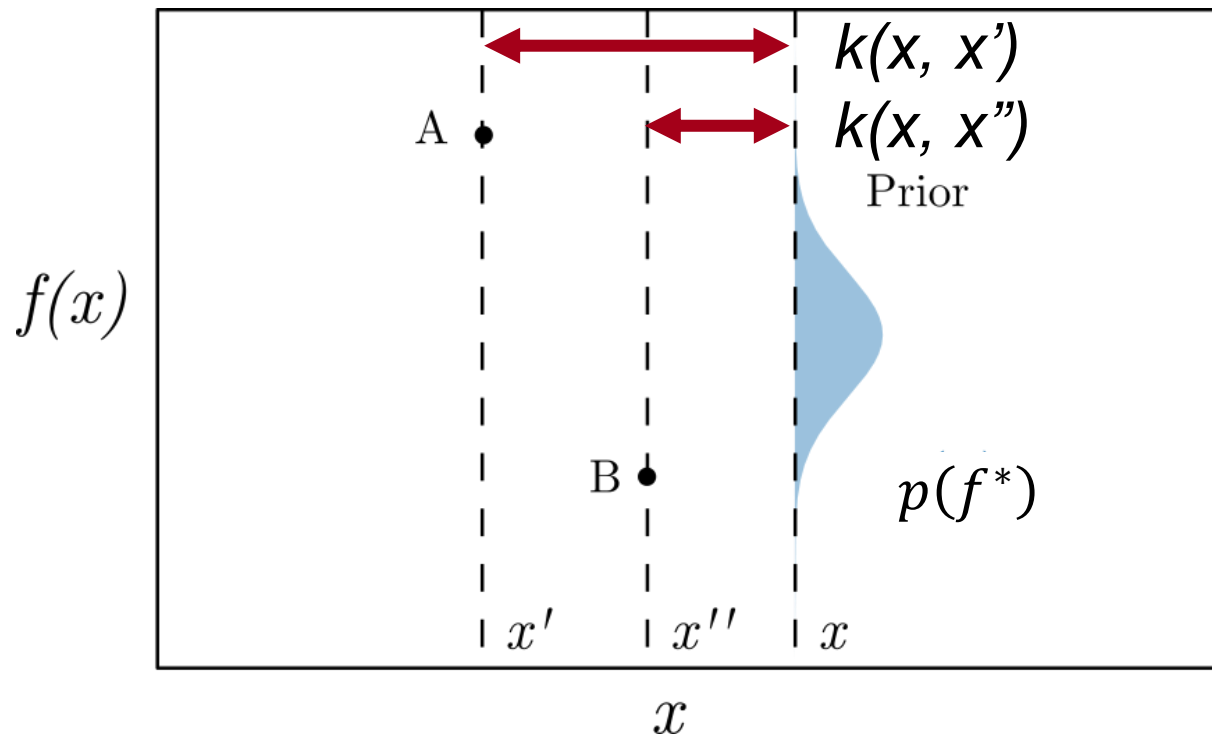
## Some intuition...

Which observation will have a larger impact on changing  $p(f)$ ?



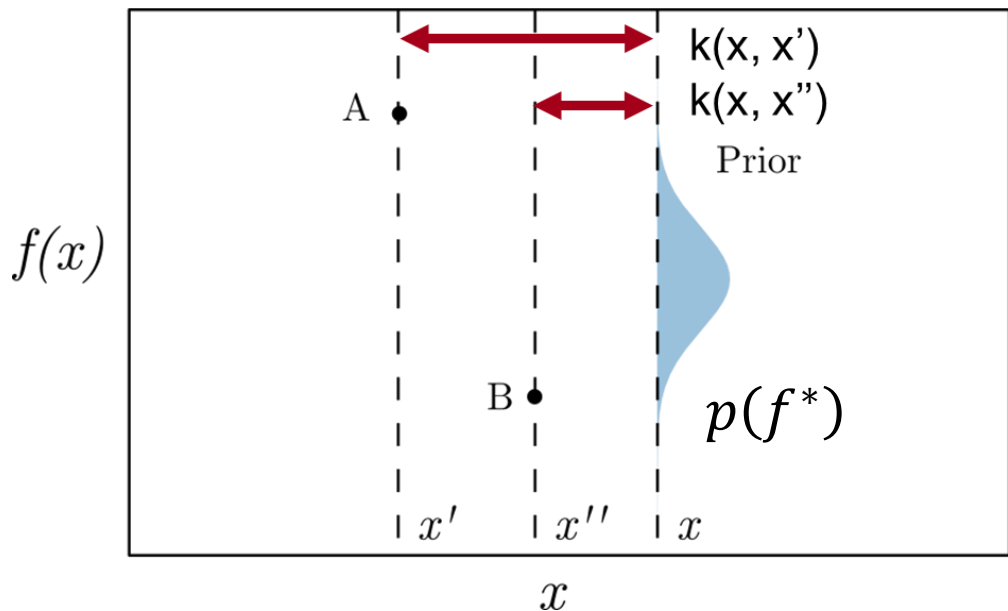
# Adding some math

Which observation will have a larger impact on changing  $p(f)$ ?



$$k(x, x') < k(x, x'')$$

# Adding some math



$$p(f_A, f_B, f^*) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

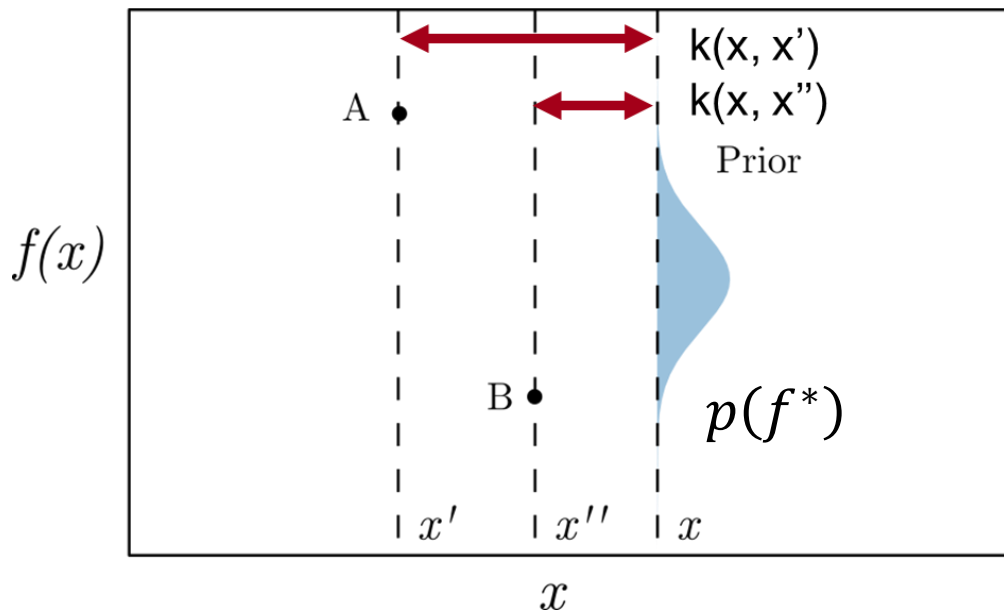
$$\boldsymbol{\Sigma} = \begin{pmatrix} k(x', x') & k(x', x'') & k(x', x) \\ k(x', x'') & k(x'', x'') & k(x'', x) \\ k(x', x) & k(x'', x) & k(x, x) \end{pmatrix}$$



$$p(f^* | f_A, f_B) = \frac{p(f_A, f_B | f^*) p(f^*)}{p(f_A, f_B)} = \frac{p(f_A, f_B, f^*)}{p(f_A, f_B)}$$

Bayes rule

# Adding some math



$$p(f^* | f_A, f_B) = N(\boldsymbol{\mu}^*, \boldsymbol{\sigma}^*)$$

$$\boldsymbol{\mu}^* = \boldsymbol{\mu} + K^* K^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

$$\boldsymbol{\sigma}^* = K^{**} - K^{*T} K^{-1} K^*$$

$$p(f^* | f_A, f_B) = \frac{p(f_A, f_B | f^*) p(f^*)}{p(f_A, f_B)} = \frac{p(f_A, f_B, f^*)}{p(f_A, f_B)}$$

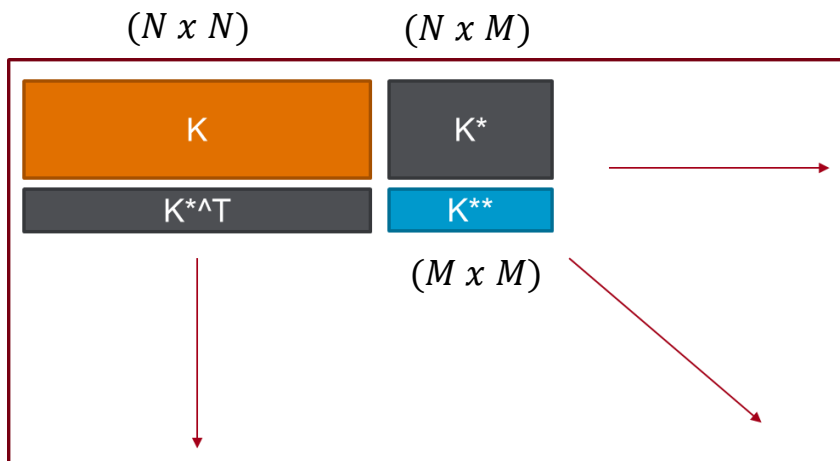


$K^{-1} \sim \mathcal{O}(N^3)!$

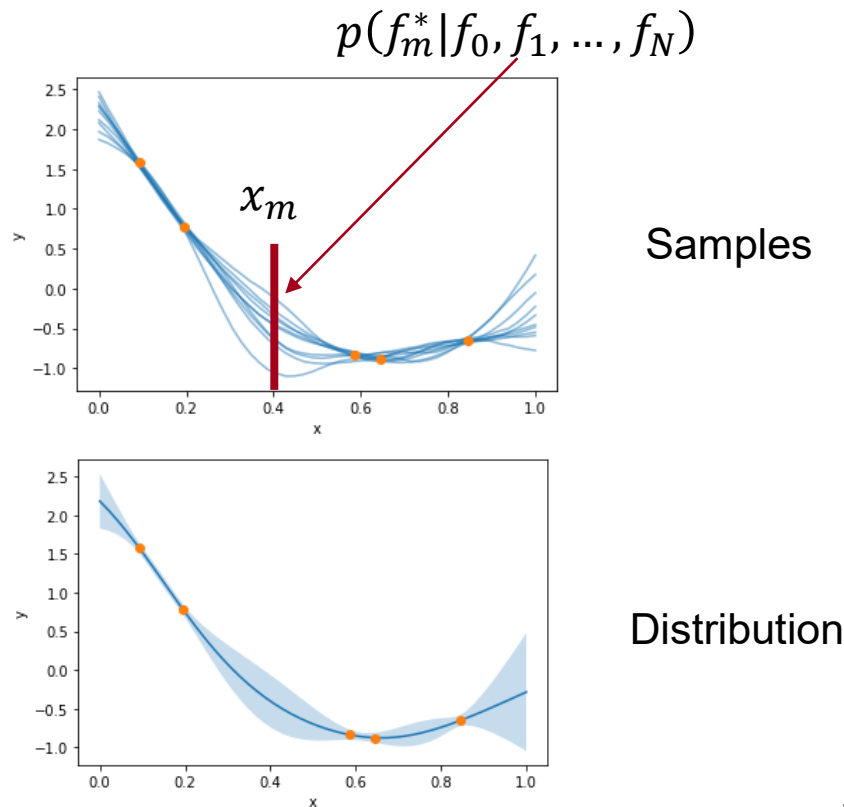
# Making predictions with GP's

What about multiple predictions?

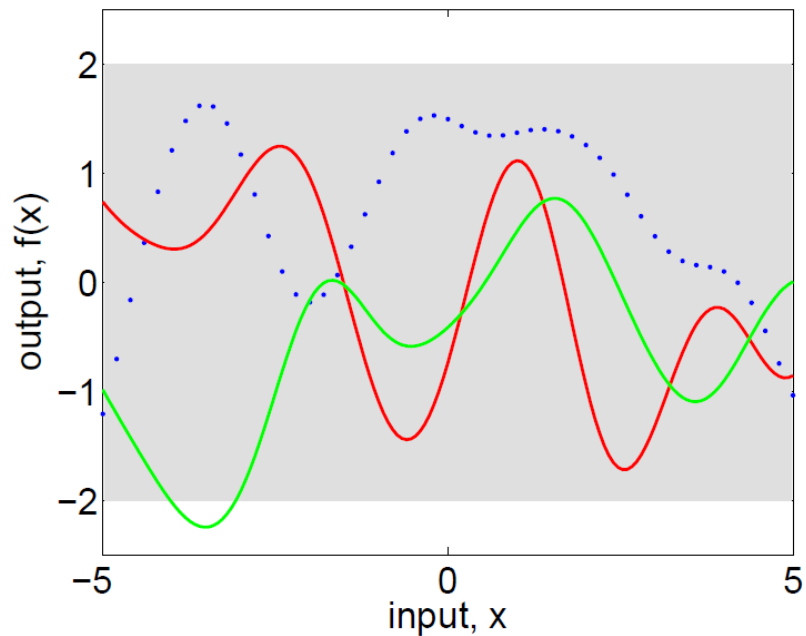
$$p(f_0^*, f_1^*, \dots, f_M^* | f_0, f_1, \dots, f_N) = N(\mu^*, \sigma^*)$$



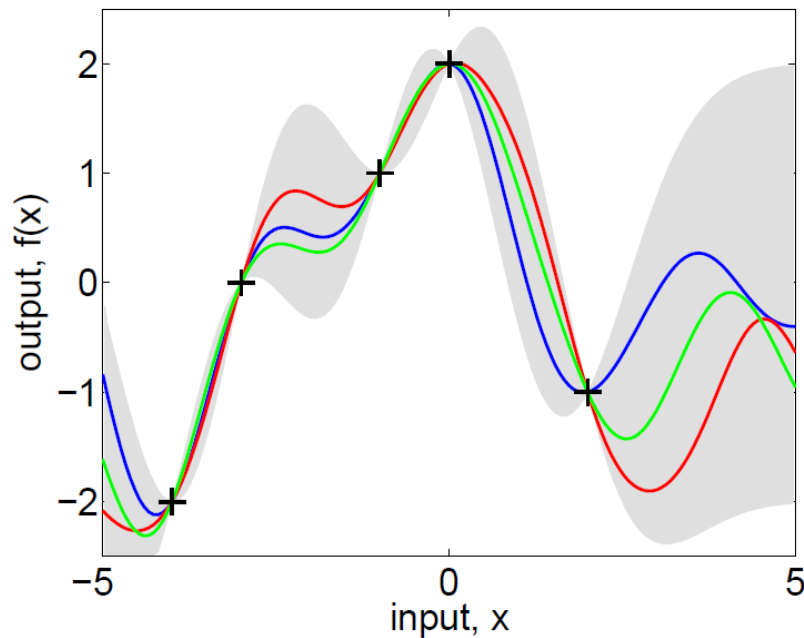
Draw function samples? Sample from the joint posterior distribution at requested points



# Another perspective



(a), prior



(b), posterior

# Fitting Gaussian Processes to Data

Need to determine the form and hyperparameters of the kernel

Each kernel has hyperparameters that control the overall function behavior.

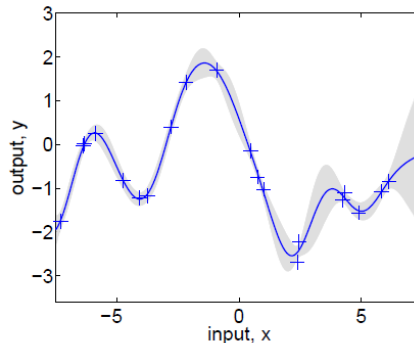
Radial Basis Function:

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

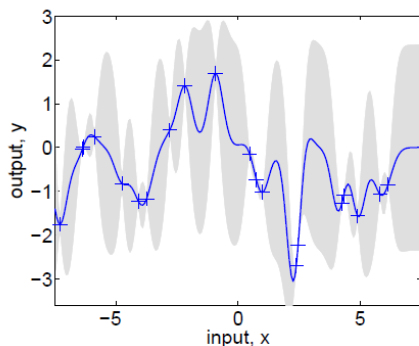
Kernel amplitude

Kernel length scale

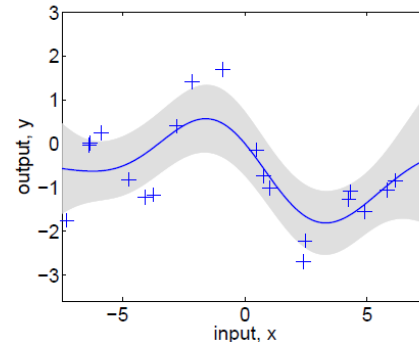
Noise



(a),  $\ell = 1$



(b),  $\ell = 0.3$



(c),  $\ell = 3$  R&W



# Fitting Gaussian Processes to Data

Need to determine the form and hyperparameters of the kernel

We fit kernel parameters to data by using Maximum Likelihood Estimation (MLE or MLE-II).

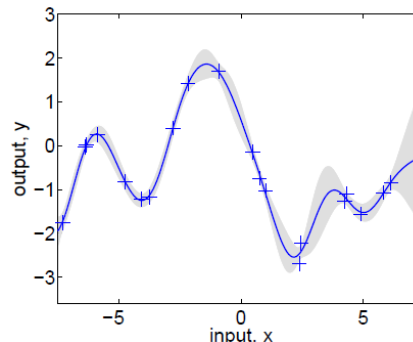
$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \log p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{X})$$

With a Gaussian likelihood the log likelihood can be calculated analytically:

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

Predictive accuracy

Model complexity



(a),  $\ell = 1$

# More Expressive Kernels

We can also encode low dimensional structure into the kernel.

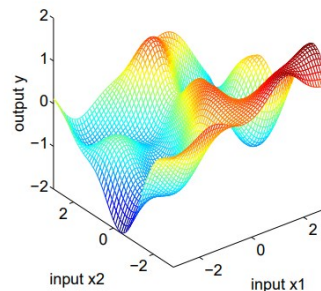
$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \boldsymbol{\Sigma}(\mathbf{x} - \mathbf{x}')\right)$$

Automatic relevance determination

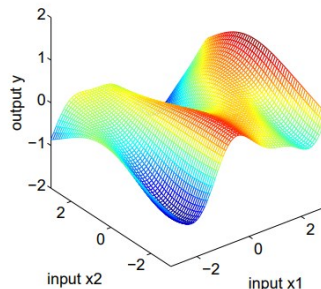
$$\boldsymbol{\Sigma} = \text{diag}(\mathbf{l})^{-2}$$

Factor analysis distance

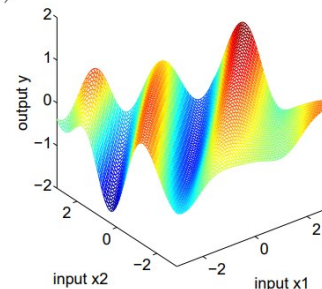
$$\boldsymbol{\Sigma} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^T + \text{diag}(\mathbf{l})^{-2}$$



(a)



(b)



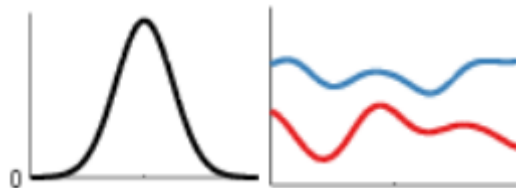
(c)

# Fitting Gaussian Processes to Data

Need to determine the form and hyperparameters of the kernel

Radial Basis Function:

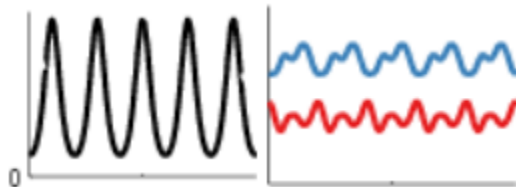
$$k_{\text{SE}}(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$$



Stationary

Periodic:

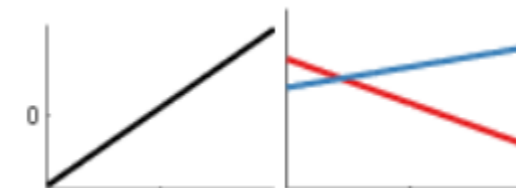
$$k_{\text{Per}}(x, x') = \sigma^2 \exp\left(-\frac{2 \sin^2(\pi|x-x'|/p)}{\ell^2}\right)$$



Stationary

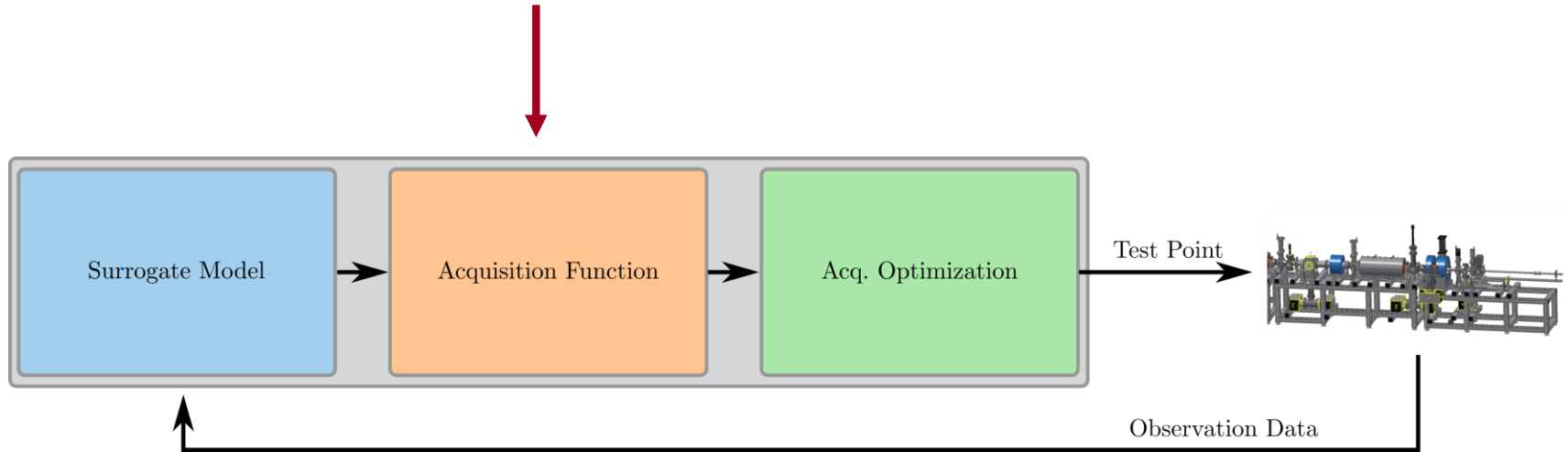
Linear:

$$k_{\text{Lin}}(x, x') = \sigma_b^2 + \sigma_v^2(x - c)(x' - c)$$



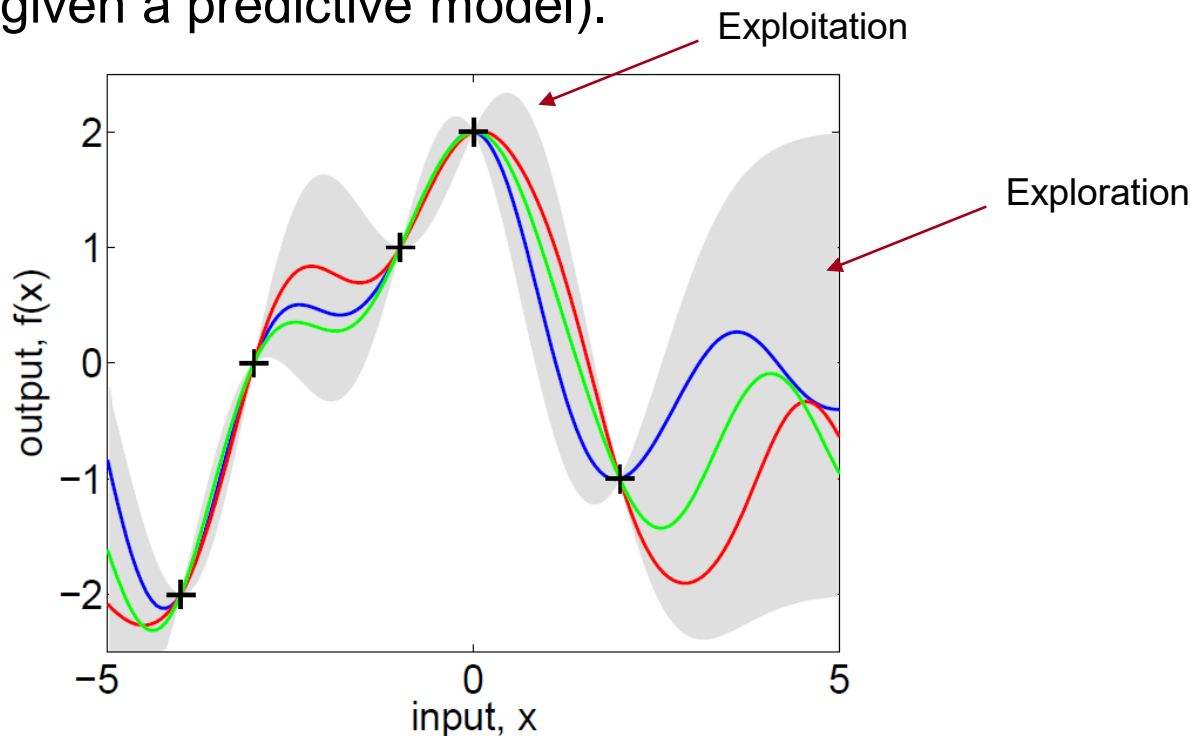
Non-stationary

# Model Based Optimization of Accelerators

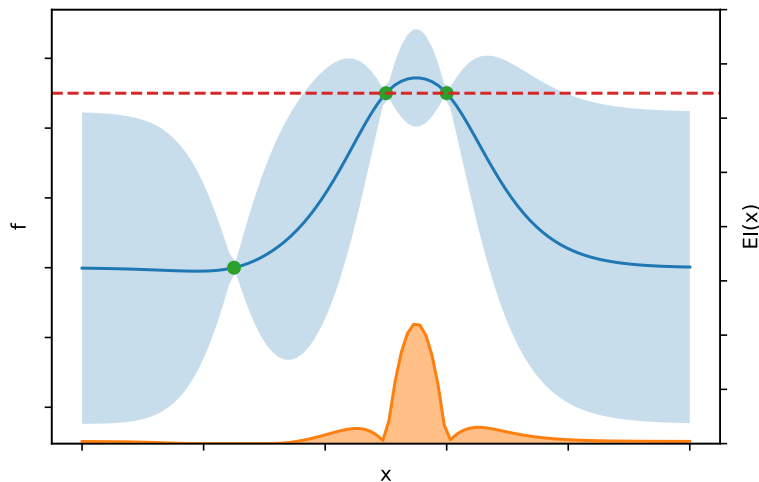


# The Acquisition Function

Define a function that characterizes the value of making a potential measurement (given a predictive model).

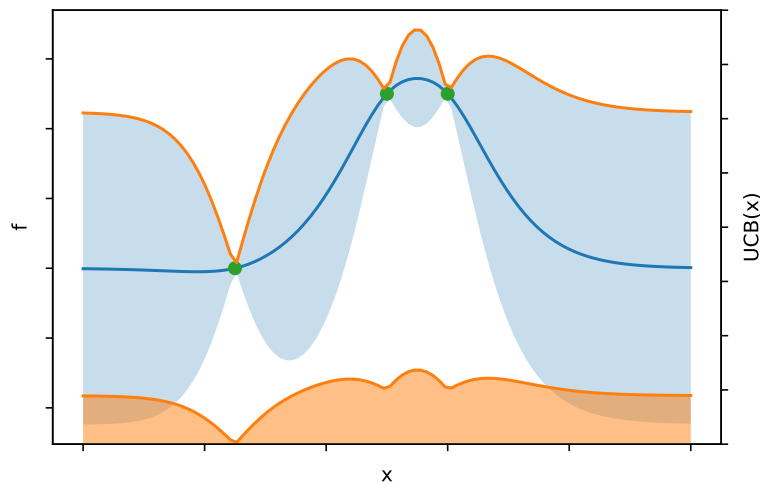


## Expected Improvement



$$EI(x) = \mathbb{E}[\max(f(x) - f^*)]$$

## Upper Confidence Bound



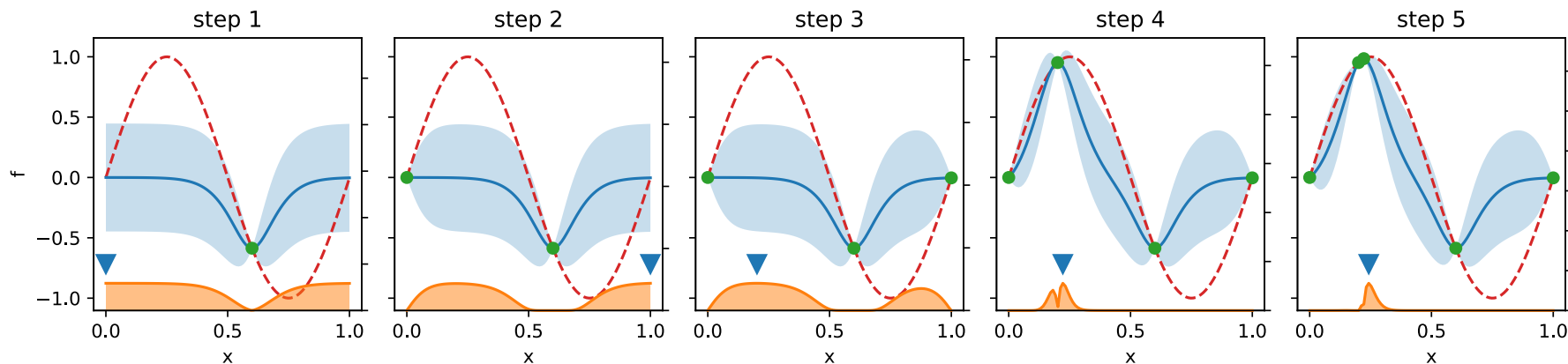
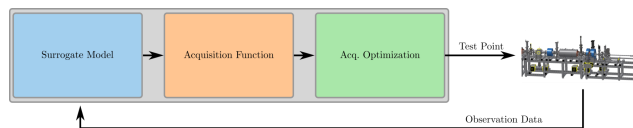
$$UCB(x) = \mu(x) + \sqrt{\beta}\sigma(x)$$

Note: most implementations of BO assume maximization

# Single Objective Optimization

$$EI(\mathbf{x}) = \mathbb{E}[\max(f(\mathbf{x}) - f^*)]$$

$$\mathbf{x}_{t+1} = \operatorname{argmax}_{\mathbf{x}} EI(\mathbf{x})$$

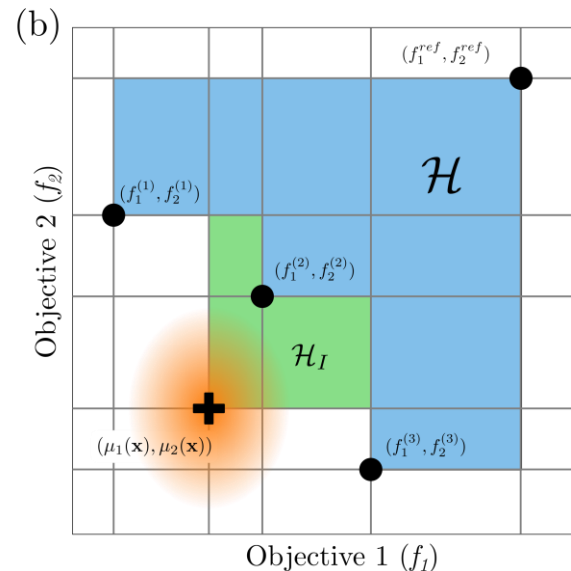
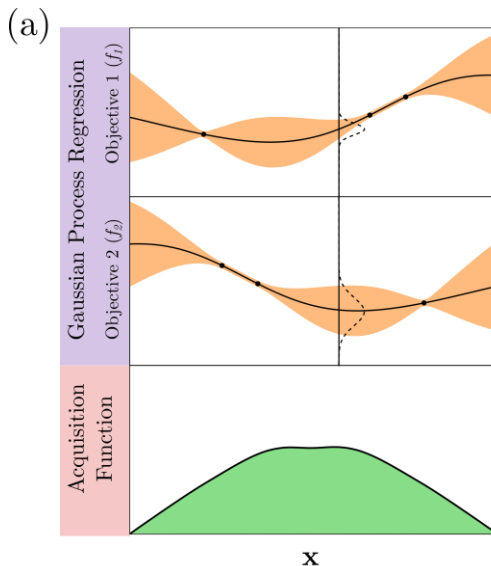
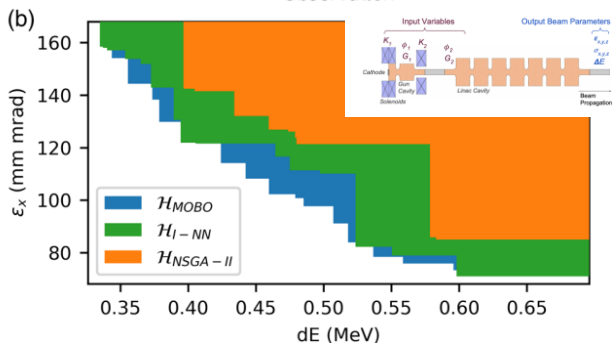
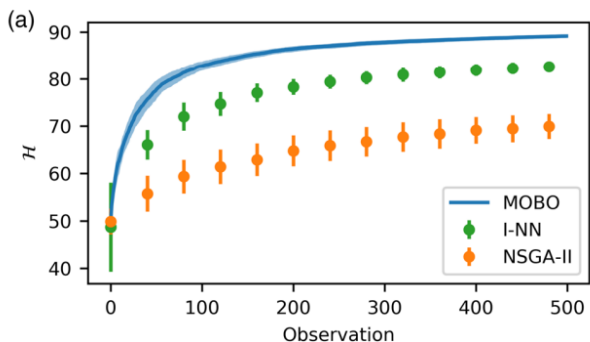


Some notes:

- The model accuracy improves in the region of interest!
- Initially the model uncertainty is maximized at the domain boundaries, BO likes to sample those
- Helpful if the acquisition function is differentiable → use gradient descent to optimize

# Multi-Objective Optimization

Determine the optimal trade-off between objectives -> the Pareto front



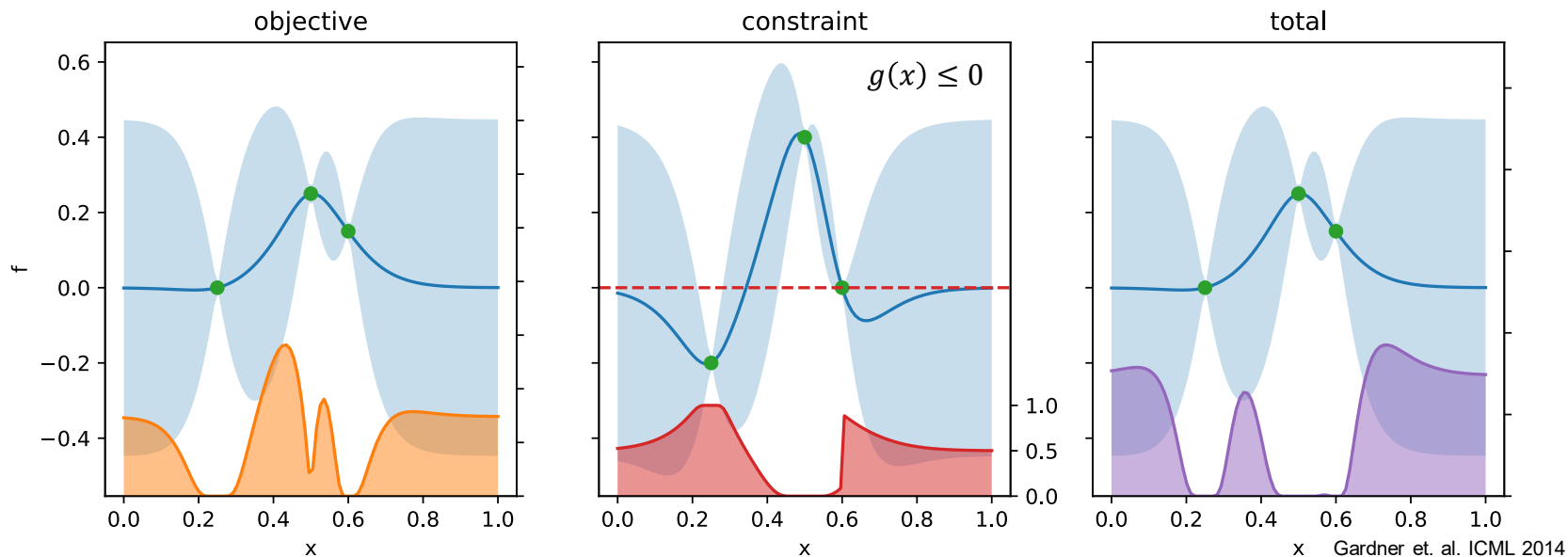
$$\alpha_{EHVI}(\mu, \sigma, \mathcal{P}, \mathbf{r}) := \int_{\mathbb{R}^P} \text{HVI}(\mathcal{P}, \mathbf{y}, \mathbf{r}) \cdot \xi_{\mu, \sigma}(\mathbf{y}) d\mathbf{y}$$



# Incorporating Constraints

Weight the acquisition function by the probability that constraints are satisfied

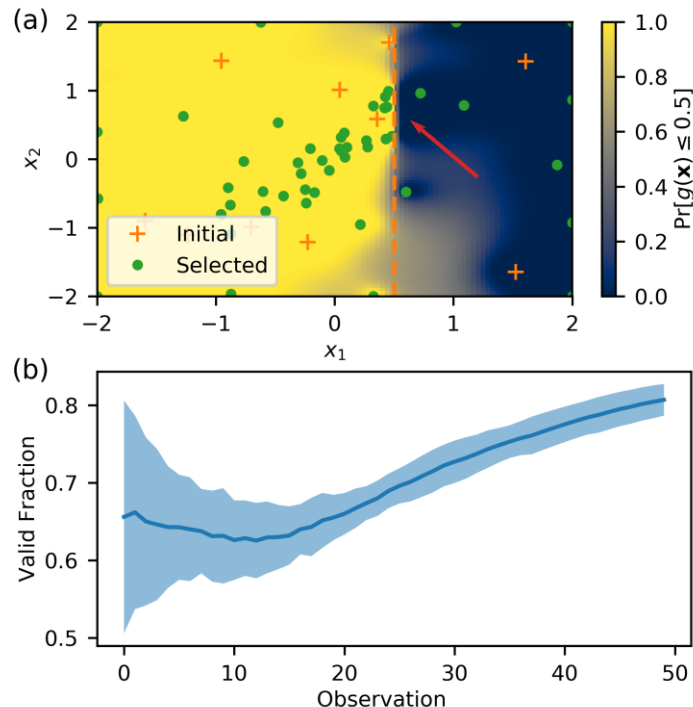
$$\hat{\alpha}(x) \rightarrow \alpha(x) \prod_i p[g_i(x) \leq h_i] \quad \text{Warning: Requires } \alpha(x) \geq 0$$



# Incorporating Constraints

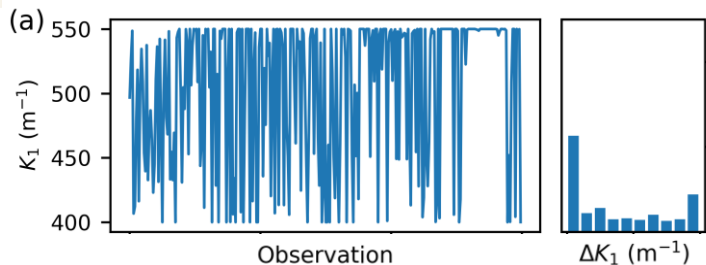
Weight the acquisition function by the probability that constraints are satisfied

$$\alpha(x) \rightarrow \alpha(x) \prod_i p[g_i(x) \leq h_i]$$



# Proximal Biasing

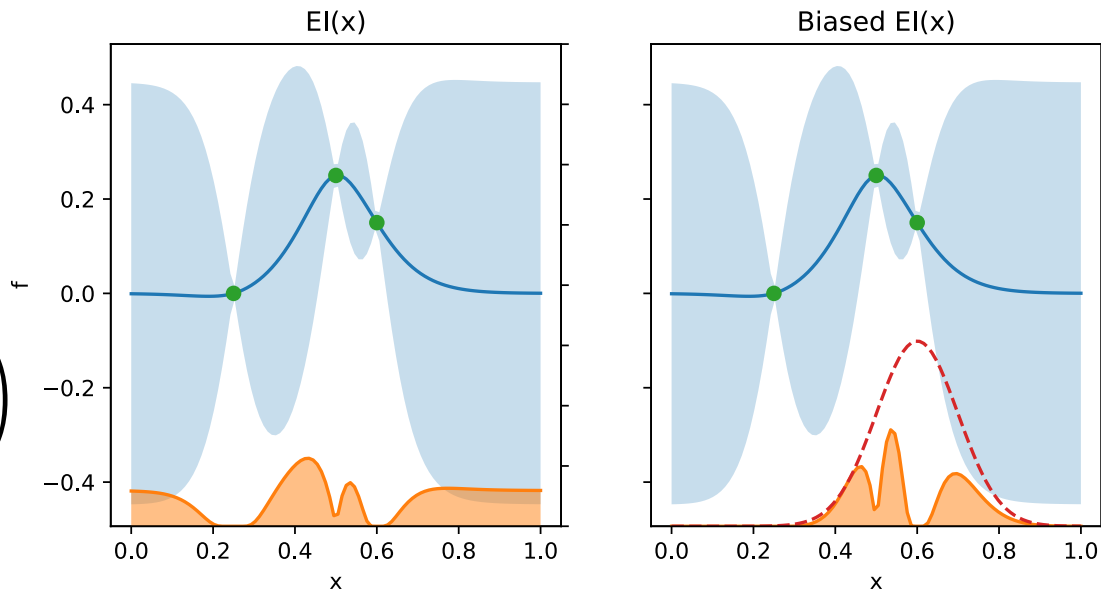
Poor optimization behavior for experimental beamlines



Weight the acquisition function by travel distance  $\rightarrow$  better than hard limits

$$\hat{\alpha}(x) \rightarrow \alpha(x) \exp\left(-\frac{(x - x_0)^2}{2\sigma^2}\right)$$

Warning: Requires  $\alpha(x) \geq 0$



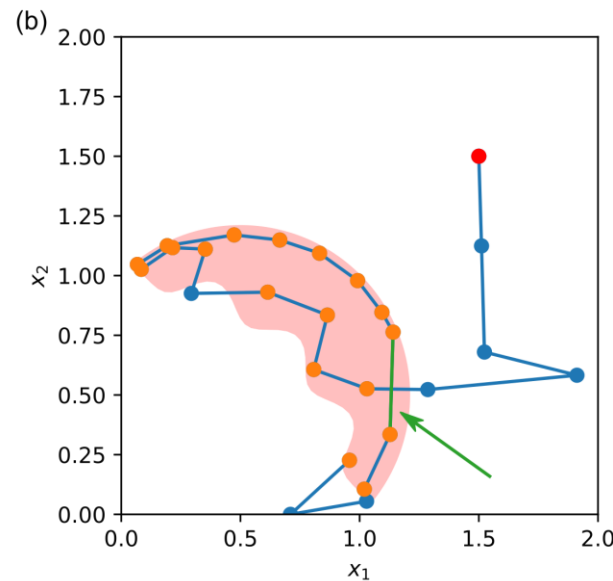
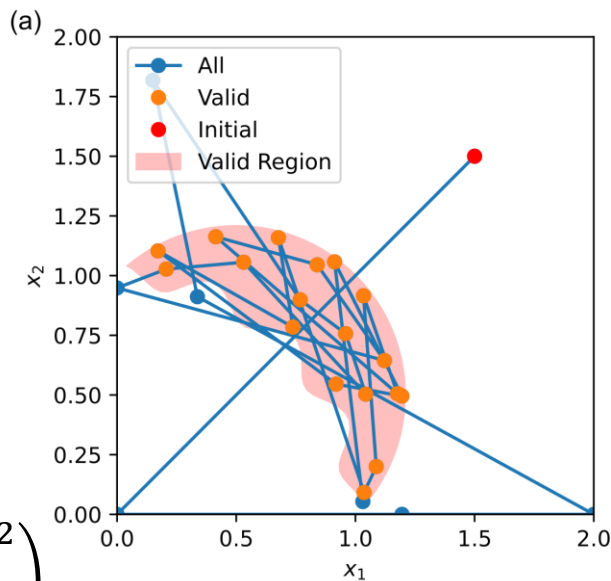
# Proximal Biasing

Poor optimization behavior for experimental beamlines

Weight the acquisition function by travel distance  $\rightarrow$  better than hard limits

$$\hat{\alpha}(x) \rightarrow \alpha(x) \exp\left(-\frac{(x - x_0)^2}{2\sigma^2}\right)$$

Warning: Requires  $\alpha(x) \geq 0$



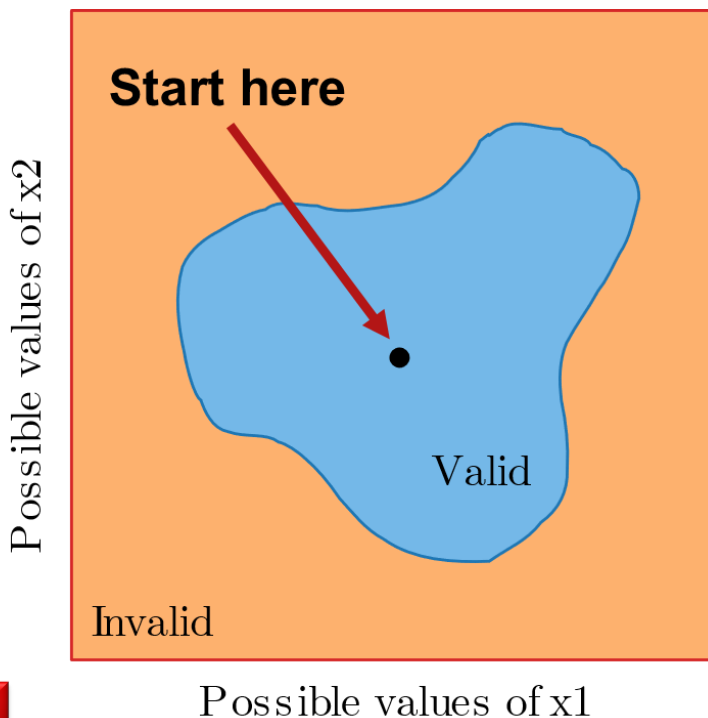
# Case study: Bayesian Exploration

Autonomous characterization of novel systems.

- Starts with a single valid observation AND **no prior information except for hardware limitations**
- Samples points in a quasi-uniform grid where the grid spacing is learned automatically based on the beam response
- Learns where the valid region is w/limited invalid observations
- Considers costs associated with changing parameters
- Natively applicable to multi-dimensional exploration
- General purpose that works in almost any case

Click button →

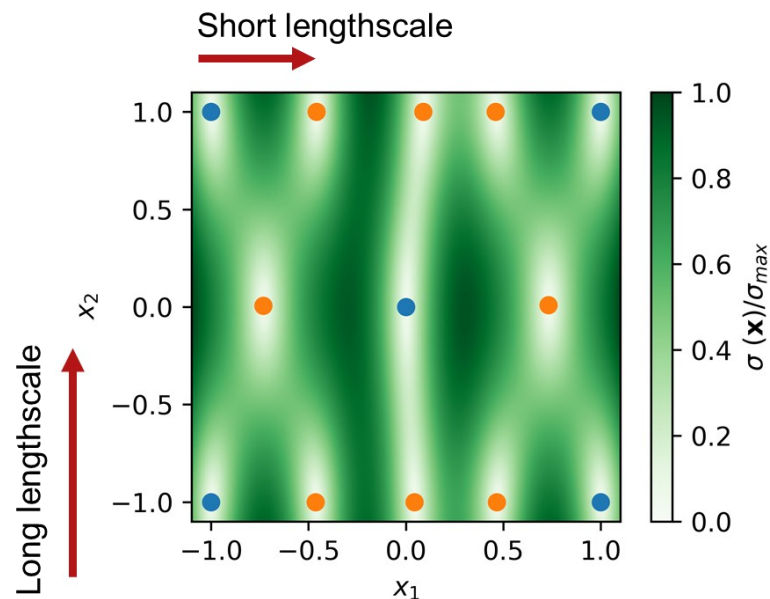
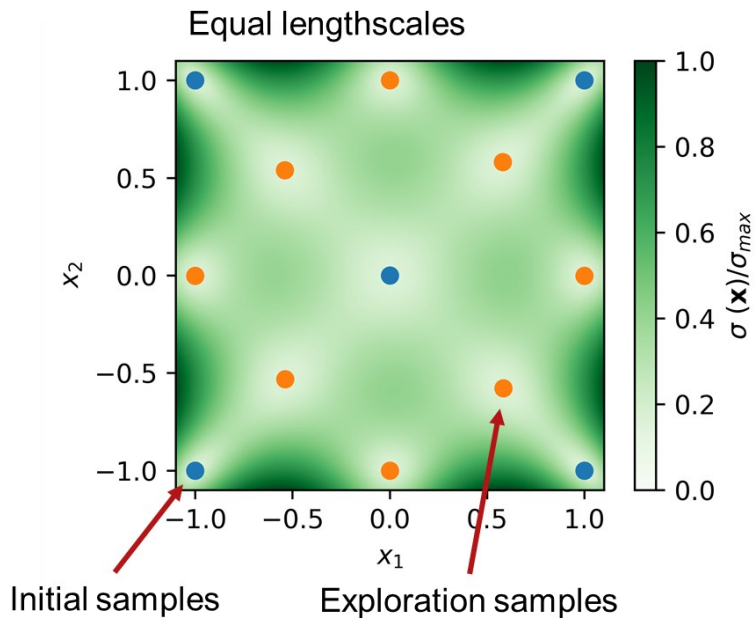
**GO!**



# Uncertainty Sampling

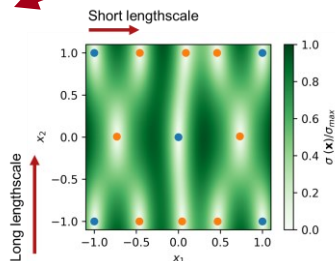
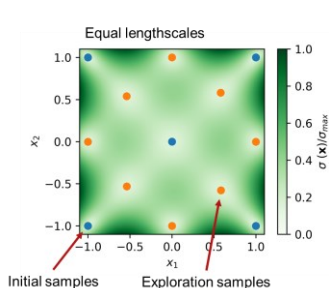
If the function changes more rapidly along one axis, sample more points along that axis!

$$\alpha(\mathbf{x}) = \sigma(\mathbf{x})$$

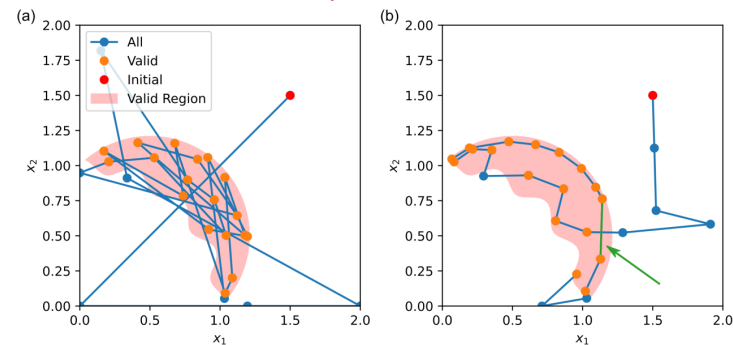


$$\alpha(\mathbf{x}) = \sigma(\mathbf{x}) \prod_{i=1}^N p(g_i(\mathbf{x}) \geq h_i) \Psi(\mathbf{x}, \mathbf{x}_0)$$

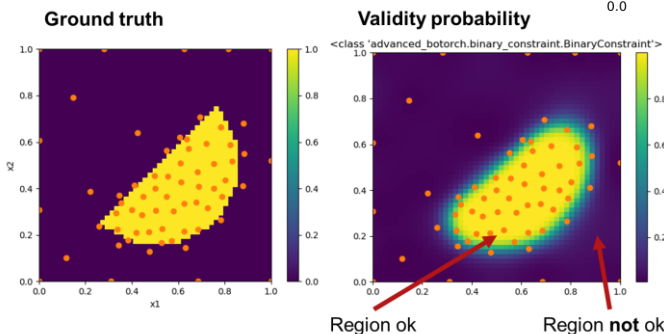
Adaptive sampling



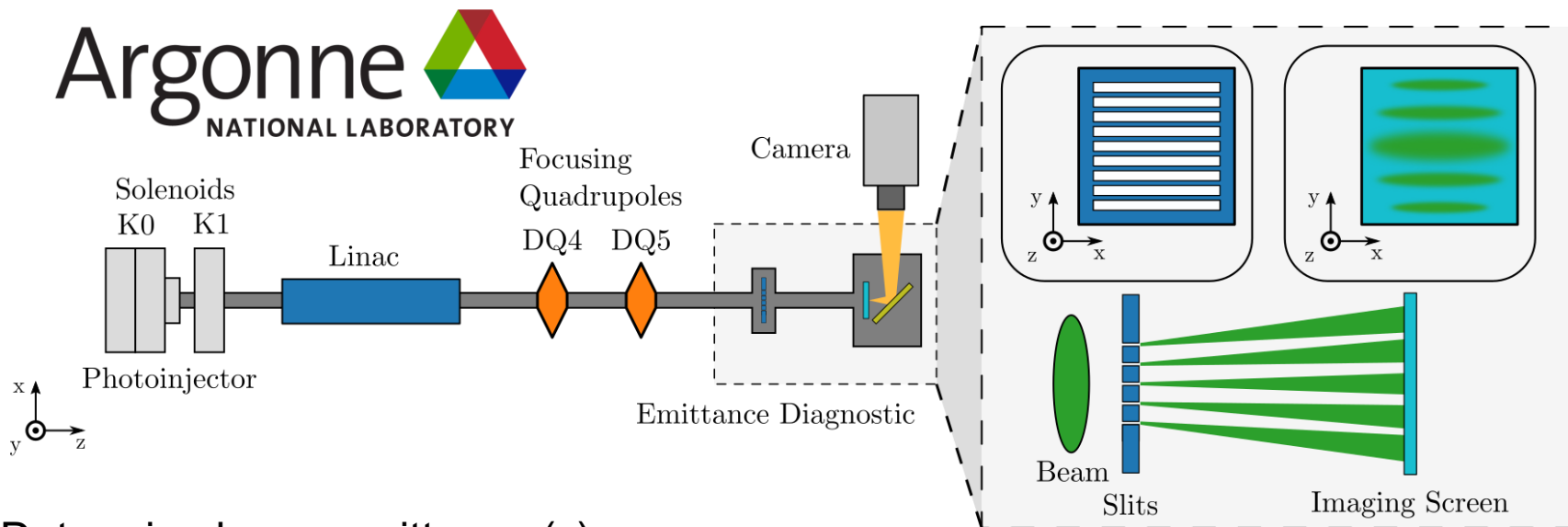
Proximal biasing



Unknown constraints



# Characterizing Photoinjector Emittance at AWA



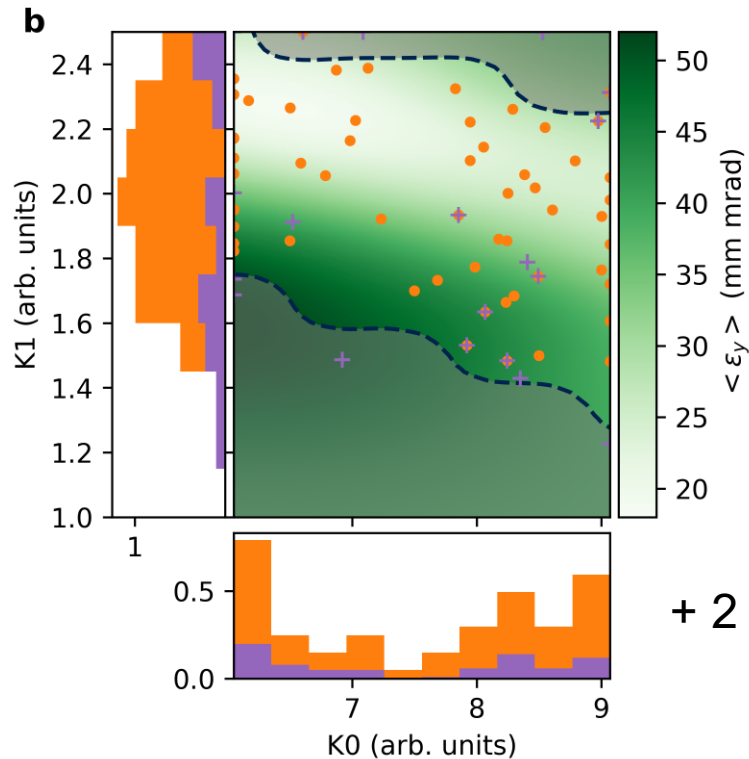
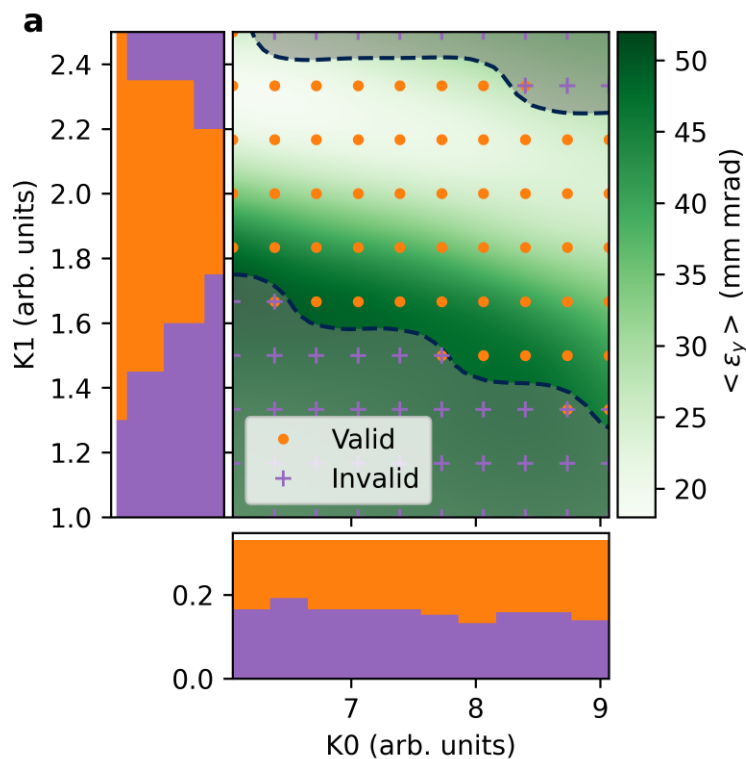
Determine beam emittance ( $\epsilon$ )

as a function of:

- 2 solenoids
- 2 quadrupoles

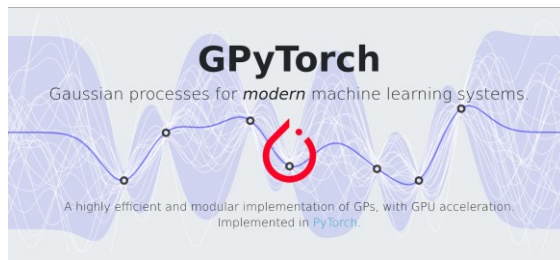


# Characterizing Photoinjector Emittance at AWA



+ 2 quads!

# Implementing Bayesian Optimization



## Use Botorch

- Single/multi-objective Bayesian optimization (serial and parallel)
- Constraints
- Proximal biasing
- MC Acquisition function evaluation
- Flexible incorporation of pyTorch Modules
- GPU support

See Examples: [https://github.com/slaclab/bo\\_tutorial](https://github.com/slaclab/bo_tutorial)



See afternoon talk/poster!

- Bayesian optimization is at its best when evaluating objective functions is **expensive** (incl. simulations)!
- Need to **balance** the costs of creating/evaluating the model vs. your application
- Take advantage of **prior information** about the objective function to speed up optimization (see NN prior poster, Nikita talk)
- Be **creative** with your acquisition functions to suit your optimization needs

- Gaussian Processes for Machine Learning (R & W) - <https://gaussianprocess.org/gpml/>
- Greenhill, Stewart, et al. "Bayesian optimization for adaptive experimental design: a review." IEEE access 8 (2020): 13937-13948.
- Balandat, Maximilian, et al. "BoTorch: a framework for efficient Monte-Carlo Bayesian optimization." Advances in neural information processing systems 33 (2020): 21524-21538.
- Tutorial Examples: [https://github.com/slaclab/bo\\_tutorial](https://github.com/slaclab/bo_tutorial)