# Bayesian Optimization

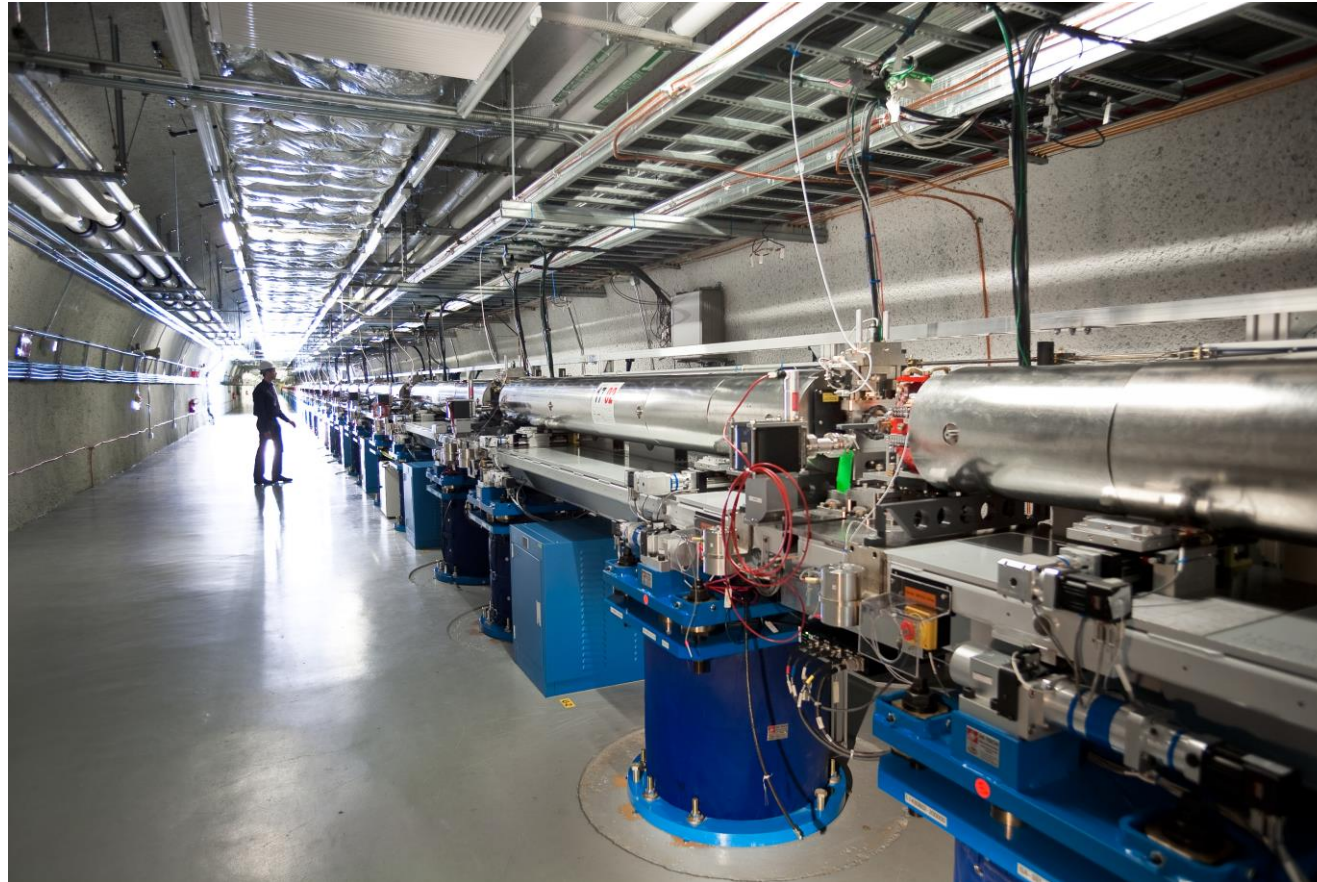**Presenter:** Adi Hanuka

**Day 5**

- Motivation

- Model-based vs model-free optimizers

- Bayesian Optimization (BO)
  - Overview
  - Acquisition functions
  - Accounting for constraints
  - Proximal optimization
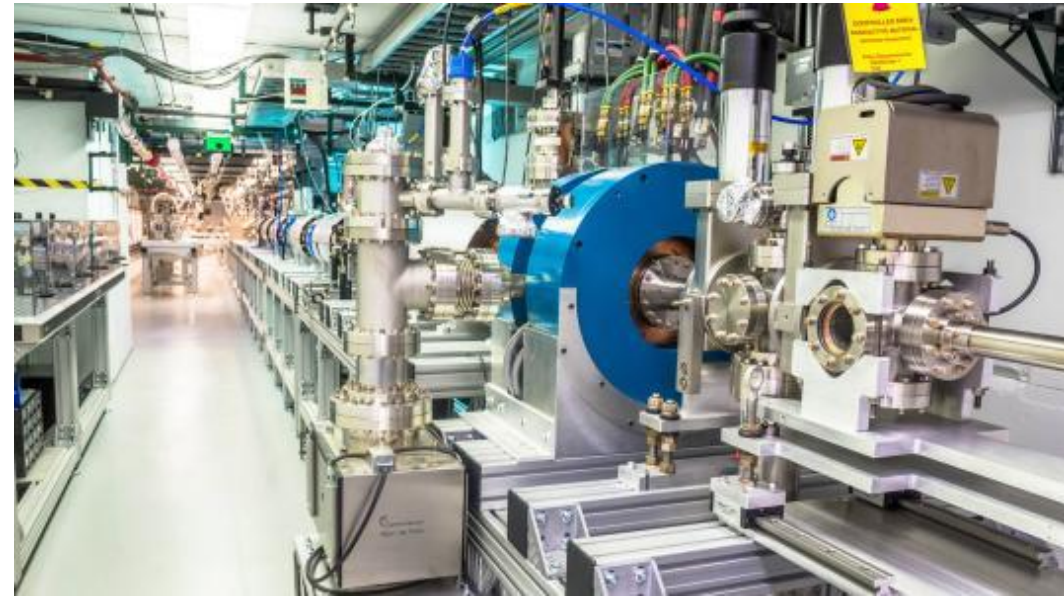
- Applications

- Summary of optimization methods

Optimize operations: maximize X-ray energy, minimize emittance, ….

Optimize design parameters

**Setup time [min]**



15 min

2-5 times/day

- Config change
- Tune to find FEL
- Quads tuning
- Undulator tuning
- Pointing / focusing

**Search space [n-D]**



~24D



Quadrupoles provide focusing
→ maintain small beam size
→ Higher X-ray pulse energy!

- Motivation

- **Model-based vs model-free optimizers**

- Bayesian Optimization (BO)
  - Overview
  - Acquisition functions
  - Accounting for constraints
  - Proximal optimization

- Applications
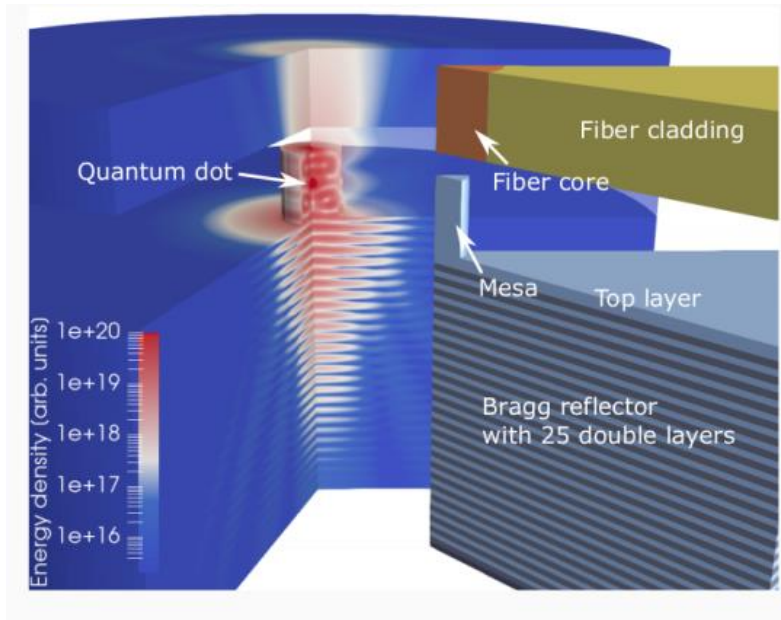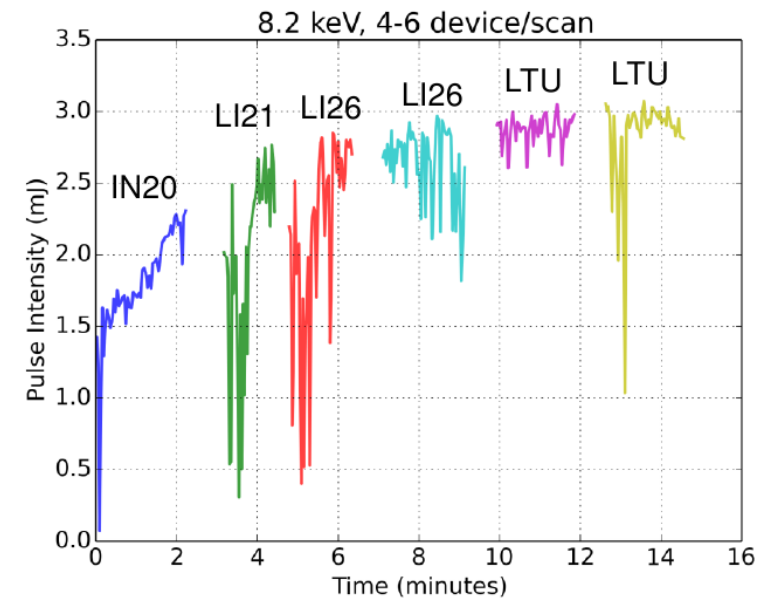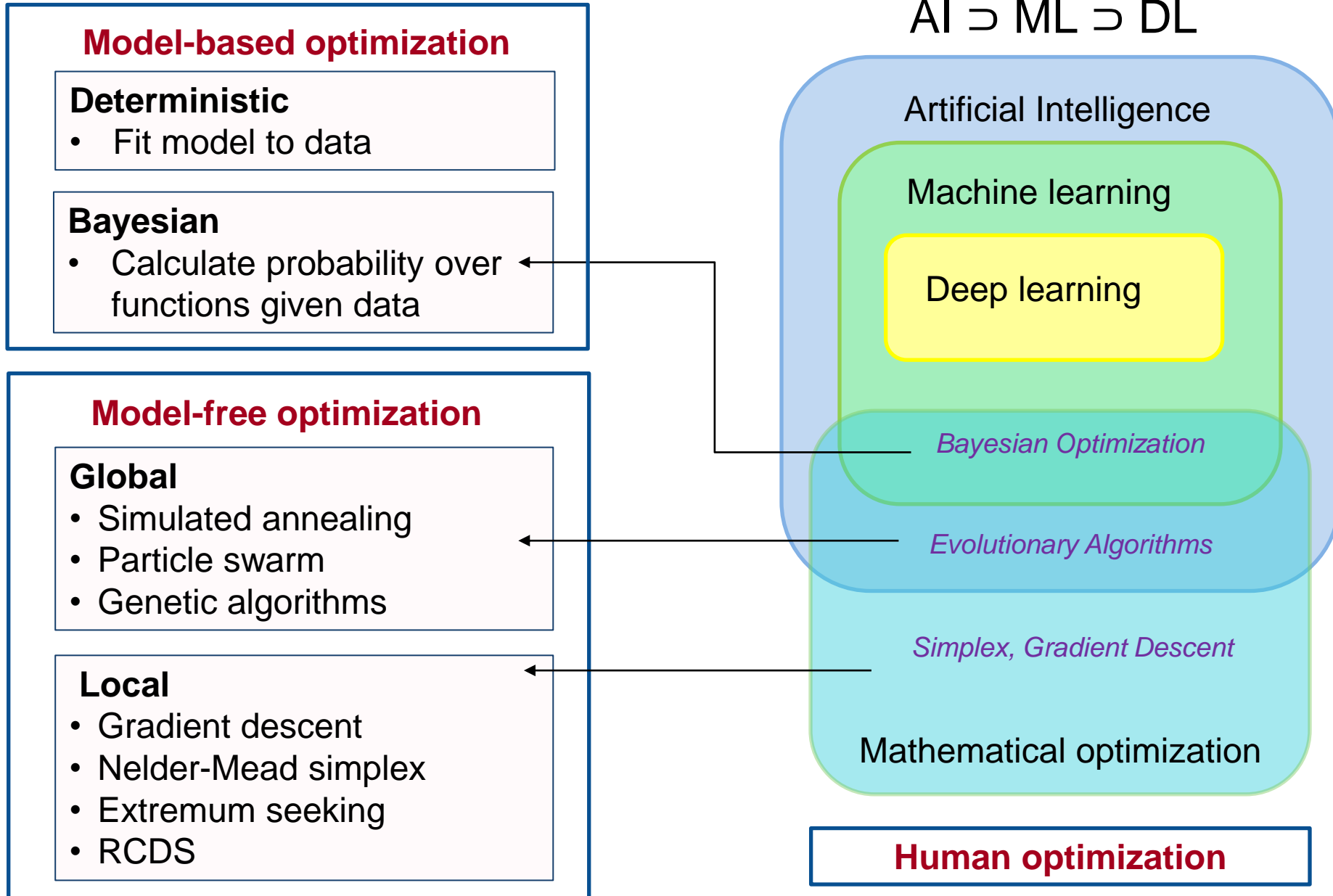
- Summary of optimization methods

# Comparison of Optimizers

## Model-based optimization

**Deterministic**
- Fit model to data

**Bayesian**
- Calculate probability over functions given data

## Model-free optimization

**Global**
- Simulated annealing
- Particle swarm
- Genetic algorithms

**Local**
- Gradient descent
- Nelder-Mead simplex
- Extremum seeking
- RCDS

AI ⊃ ML ⊃ DL

Artificial Intelligence

Machine learning

Deep learning

*Bayesian Optimization*

*Evolutionary Algorithms*

*Simplex, Gradient Descent*

Mathematical optimization

**Human optimization**

## Human optimization ≠ Numerical optimization

**Human** optimization

- Life-long learning
- Experience
- Mental modes
- (relatively) Slow decisions
- Limited working memory

**Numerical** optimization

- Bulk learning
- Cannot estimate uncertainty
- Juggle many things at once
- Fast decisions

Model-based Bayesian optimization
combines the complementary strengths of both approaches

"A good regulator of the system is a good model of that system."

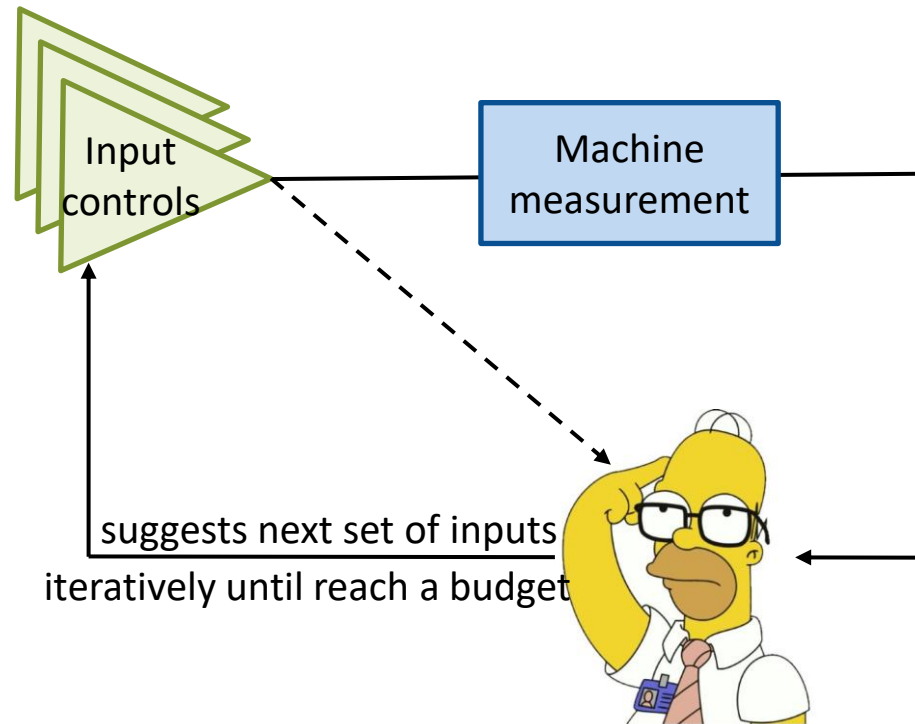ROGER C. CONANT & W. ROSS ASHBY (1970)  Science, 1:2, 89-97

- Motivation

- Model-based vs model-free optimizers

- **Bayesian Optimization (BO)**
  - Overview
  - Acquisition functions
  - Accounting for constraints
  - Proximal optimization

- Applications

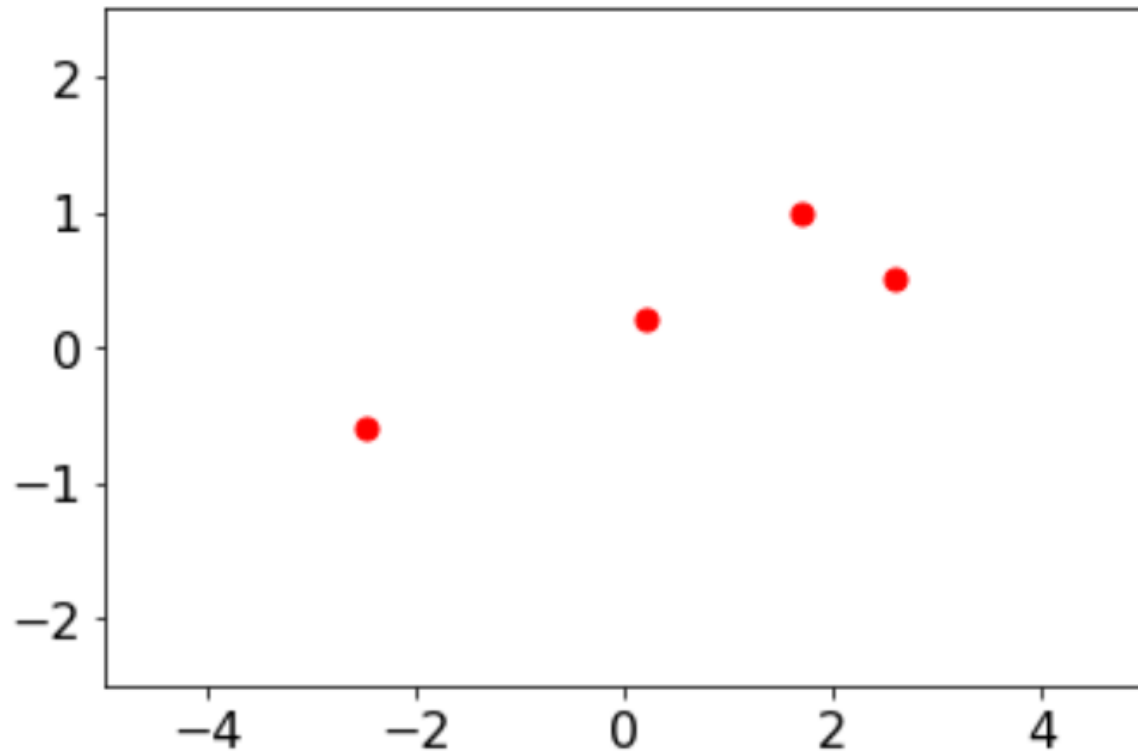- Summary of optimization methods

- Gradient-free
- Learns by experience



Input controls

Machine measurement

suggests next set of inputs iteratively until reach a budget

**Acquisition function (utility function)** that tells us where to query the system next.

Let's get some intuition...  Where is the maximum of $f$ ?

**Question**: Where should we take the next evaluation?



**Probabilistic surrogate model** for the values our function takes on unseen points.

quad   X-ray energy

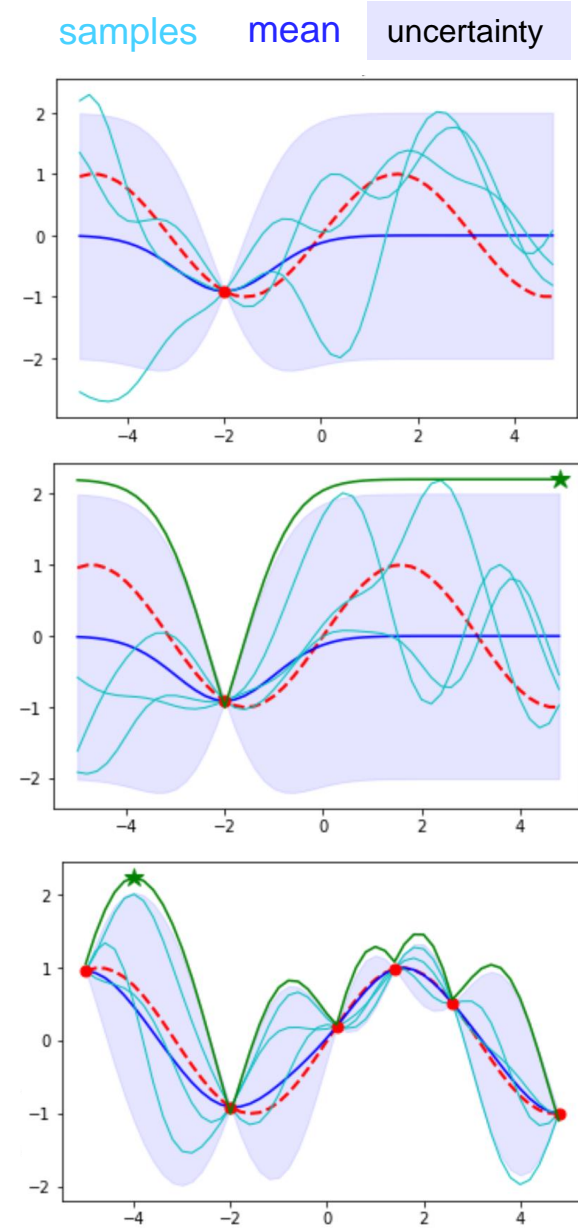unknown `objective function` $f(x);\ [x_0, f(x_0)]$

```
for each step t = 1, 2, 3, . . . . ,T:
    1. Build probabilistic model
```
$\rightarrow \hat{f}_{t-1}(x)$     Gaussian process

```
    2. Choose next point to simultaneously increase
       objective & decrease model uncertainty
```
$\rightarrow x_t = \textbf{argmax}\left(\textbf{UCB}(x|\hat{f}_{t-1})\right)$     $\textbf{UCB}(x) = \mathbf{E}[\hat{f}(x)] + \beta\hat{\sigma}(x)$

```
    3. Sample new (noisy) point
```
$\rightarrow f(x_t)$

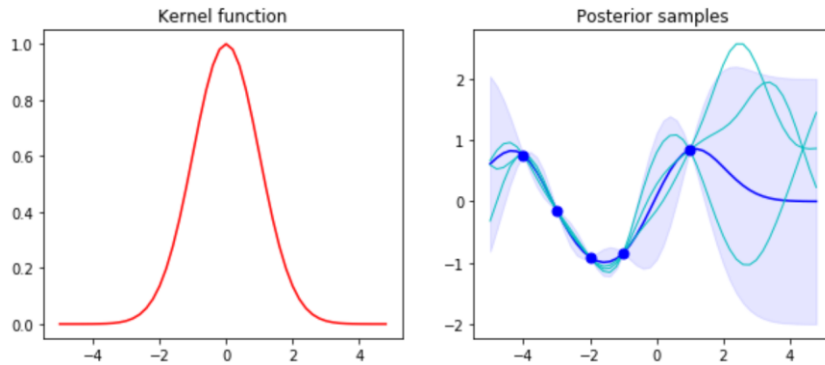samples    mean    uncertainty
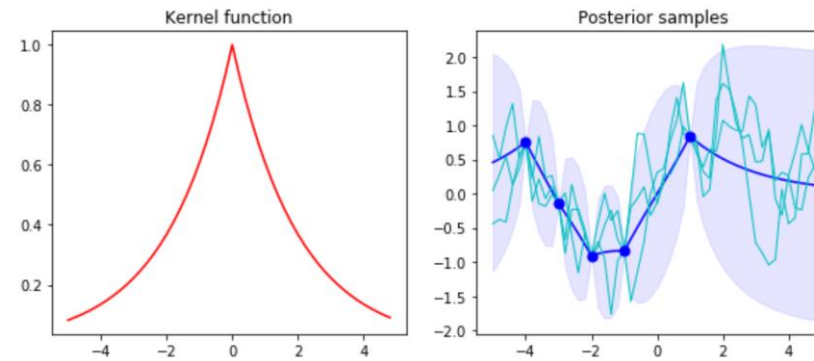
# Surrogate model: Gaussian process



- Give a reliable estimate of their own uncertainty

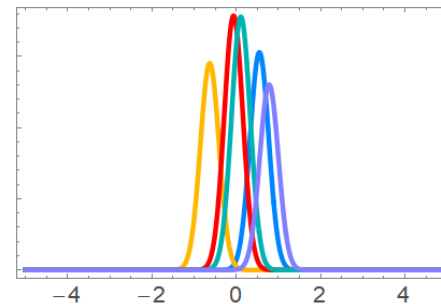- Shape our prior belief via the choice of kernel $k(x, x')$

$$k_{\text{RBF}}(x, x') = \exp\left(-\frac{(x - x')^2}{l^2}\right)$$

$$k_{\text{exponential}}(x, x') = \exp\left(-\frac{\|x - x'\|}{l^2}\right)$$





- Latent variables changing day to day
  → **optimum moves**
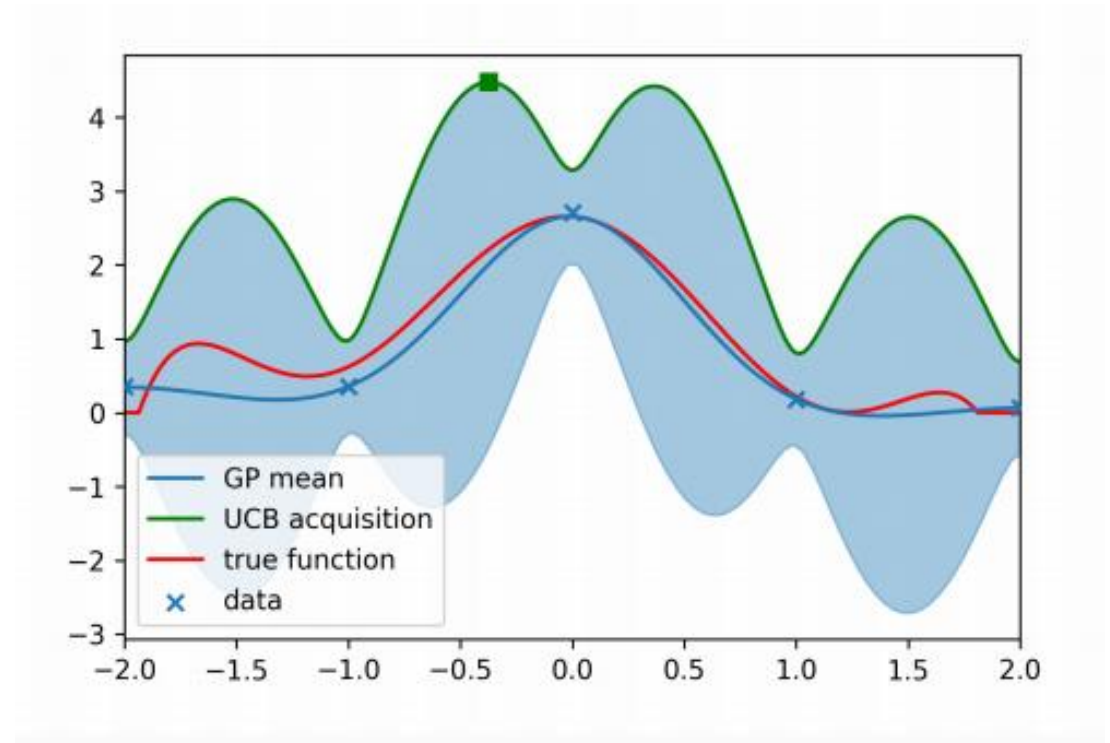  → Kernel captures **shape**

$$\text{UCB}_t(x) = \mu_t + \beta_t \sigma_t(x)$$

- $\mu_t$ - posterior mean after seeing $t$ points.
- $\sigma_t$ - posterior standard deviation after seeing $t$ points.
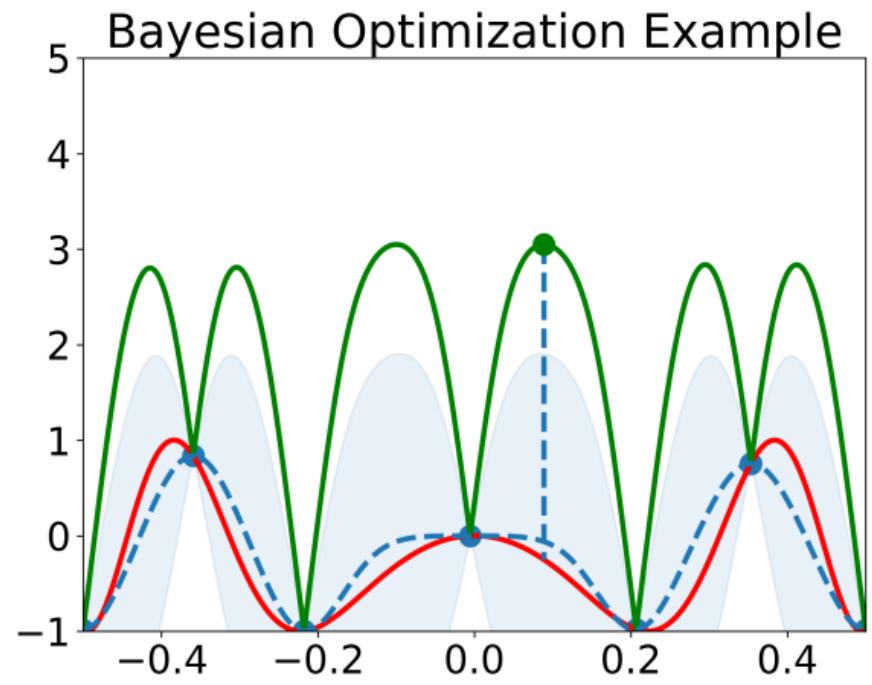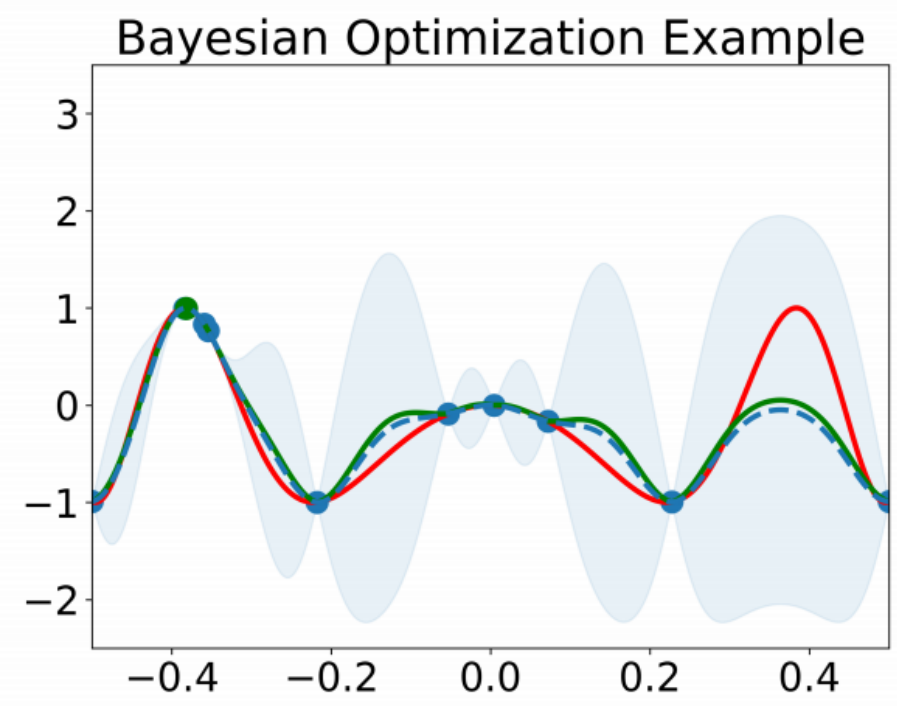
What is $\beta_t$?
- trades exploration and exploitation.
  - too small $\Rightarrow$ gets stuck/hill climbing.
  - too high $\Rightarrow$ incremental grid search.
- Common heuristic approach: $\beta \approx 2$.
- $\beta_t$ may increase with time to trade exploration as the optimization progresses.

**Question**: Which of the examples below is a better optimization process?



*Adapted from the 2nd ICFA workshop

$\beta$ small - hill climbing



*Adapted from the 2nd ICFA workshop

$\beta$ high - incremental grid search



Bayesian Optimization Example

Legend:
- evaluations
- true function
- acquisition function
- mean prediction
- 2x std. dev.

*Adapted from the 2nd ICFA workshop

$\beta$ small - hill climbing

$\beta$ high - incremental grid search





*Adapted from the 2nd ICFA workshop

$$\text{EI}_t(\boldsymbol{x}) = \mathbf{E}[\max(0, \text{f}(\boldsymbol{x}) - \text{f}(\boldsymbol{x}^+)]$$

- Analytical solution: $(\mu_t(x) - \mu(x^+))\Phi(Z) + \sigma(x)\varphi(Z)$

  where $Z = \dfrac{\mu_t - \mu(x^+)}{\sigma_t(x)}$ and $\Phi, \varphi$ are cdf and pdf of standard normal.



Chen, Yutian, et al. "Bayesian optimization in alphago." *arXiv preprint arXiv:1812.06855* (2018).

Other types of acquisition functions, each results in a different optimization process.

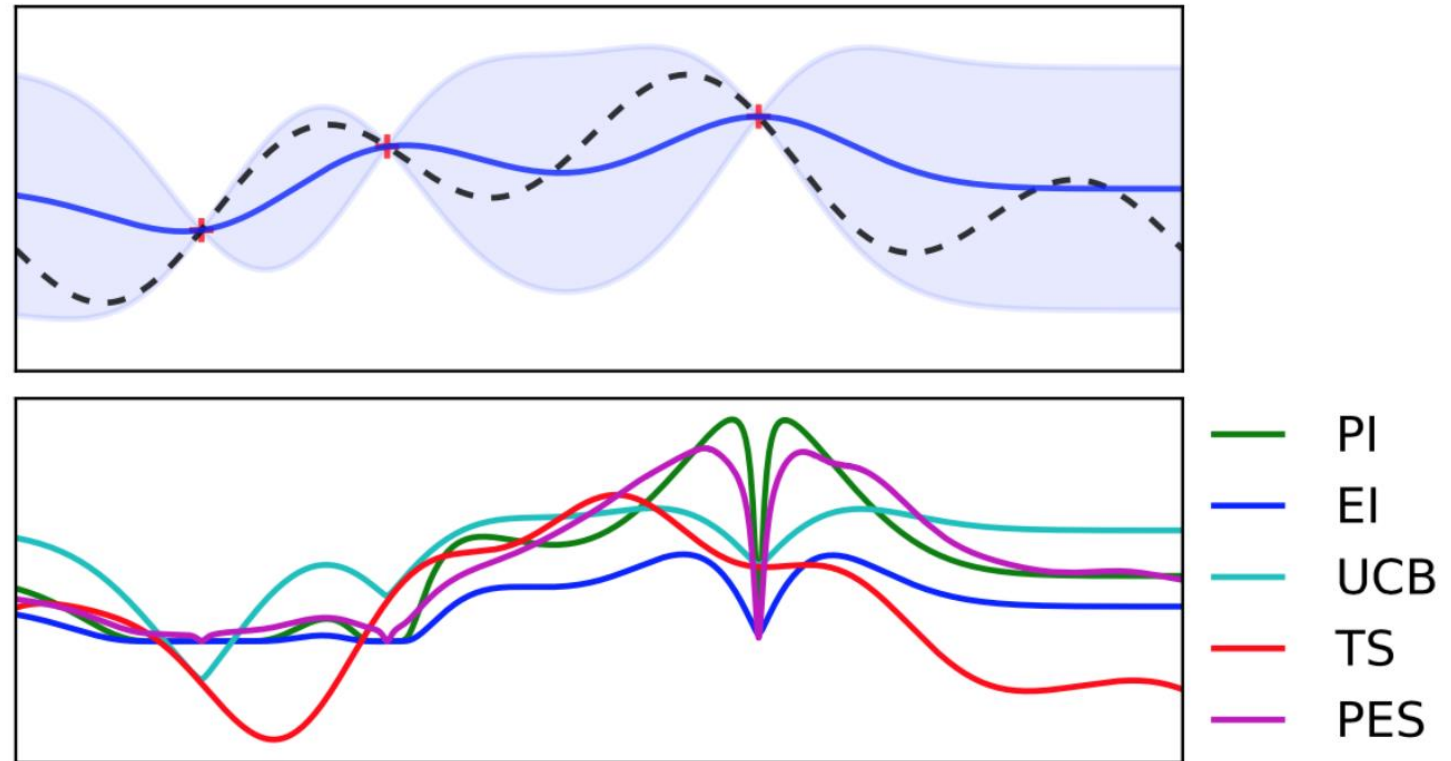- PI: Probability of improvement
- TS: Thompson sampling
- PES: Predictive entropy search



https://towardsdatascience.com/shallow-understanding-on-bayesian-optimization-324b6c1f7083

**Unknown objective**
$f(x)$

**Acquisition function**
$\text{UCB}(x) = \mathbf{E}\left[\hat{f}(x)\right] + \beta\hat{\sigma}(x)$

$-$ $\hat{f}(x)$ **esitmate**

$\hat{f}(x)$ **uncertainty**

X  **Evaluation points**

- Motivation

- Model-based vs model-free optimizers

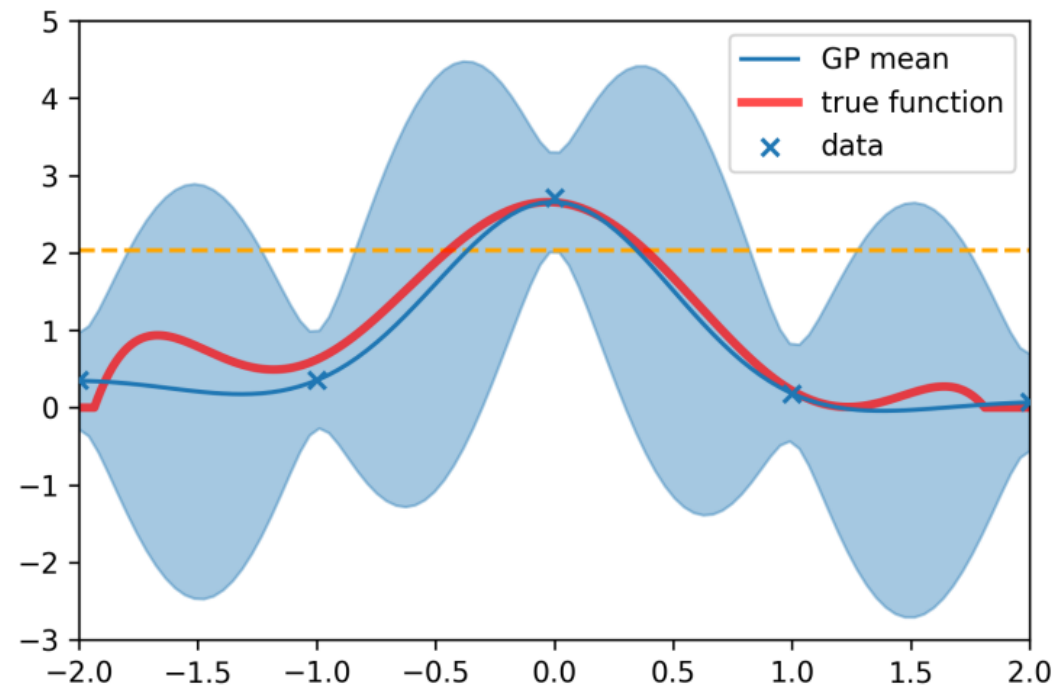- **Bayesian Optimization (BO)**
  - Overview
  - Acquisition functions
  - **Accounting for constraints**
  - Proximal optimization

- Applications

- Summary of optimization methods

Constrain the acquisition function search:
- Avoid unnecessary evaluations.
- Safe BO – not to harm the system.

*Adapted from the 2nd ICFA workshop
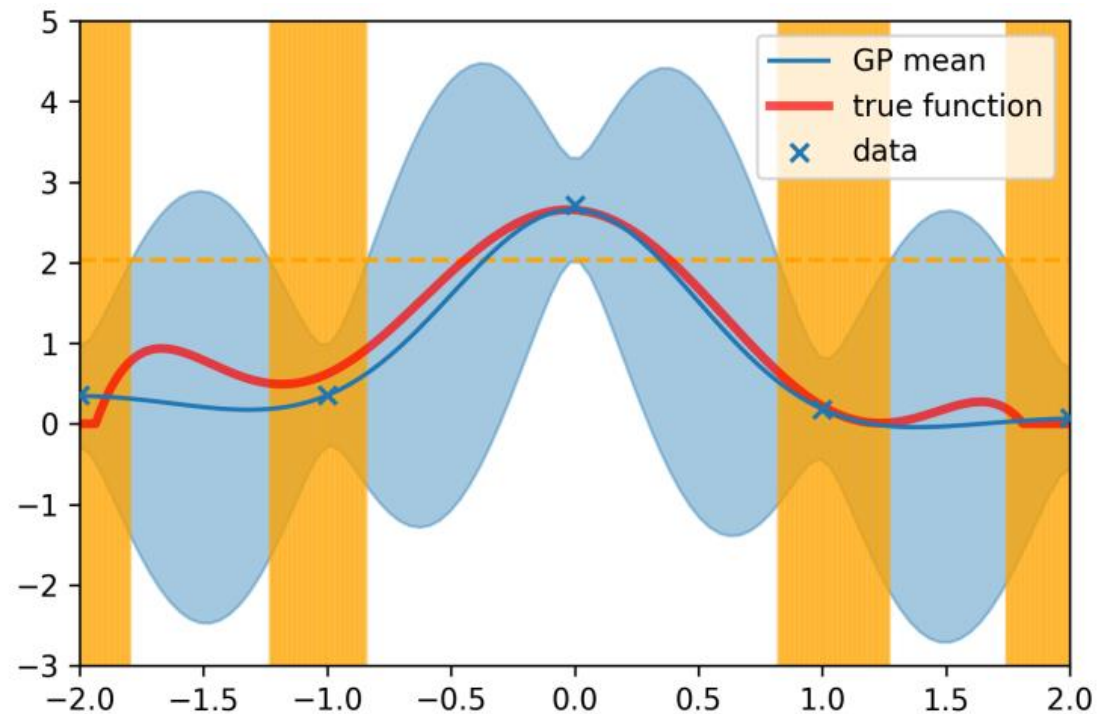
Constrain the acquisition function search:
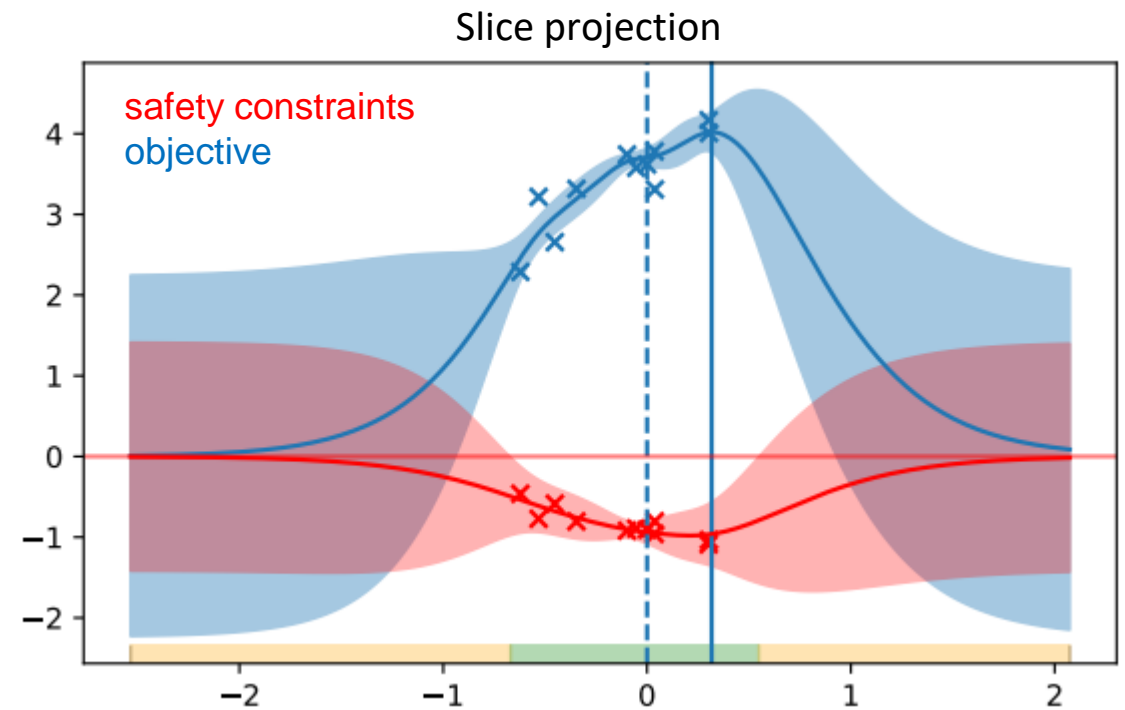- Avoid unnecessary evaluations.
- Safe BO – not to harm the system.

Orange
regions to
be avoided

Maximize FEL energy at SwissFEL using 24 parameters with constrains (lower bound on intensity).

**Problem:** making large changes in machine input parameters (magnetic field strengths, cavity phases) frequently is undesirable or infeasible.

**Solution**: Prioritize points in input space that are near the current or most recently observed parameter setting.

Done by penalizing the acquisition function (i.e. by multiplying a multivariate Gaussian distribution).

**Surrogate model:**

GP regression - $O(n^3)$

- <u>Speed</u>: Sparse GP.
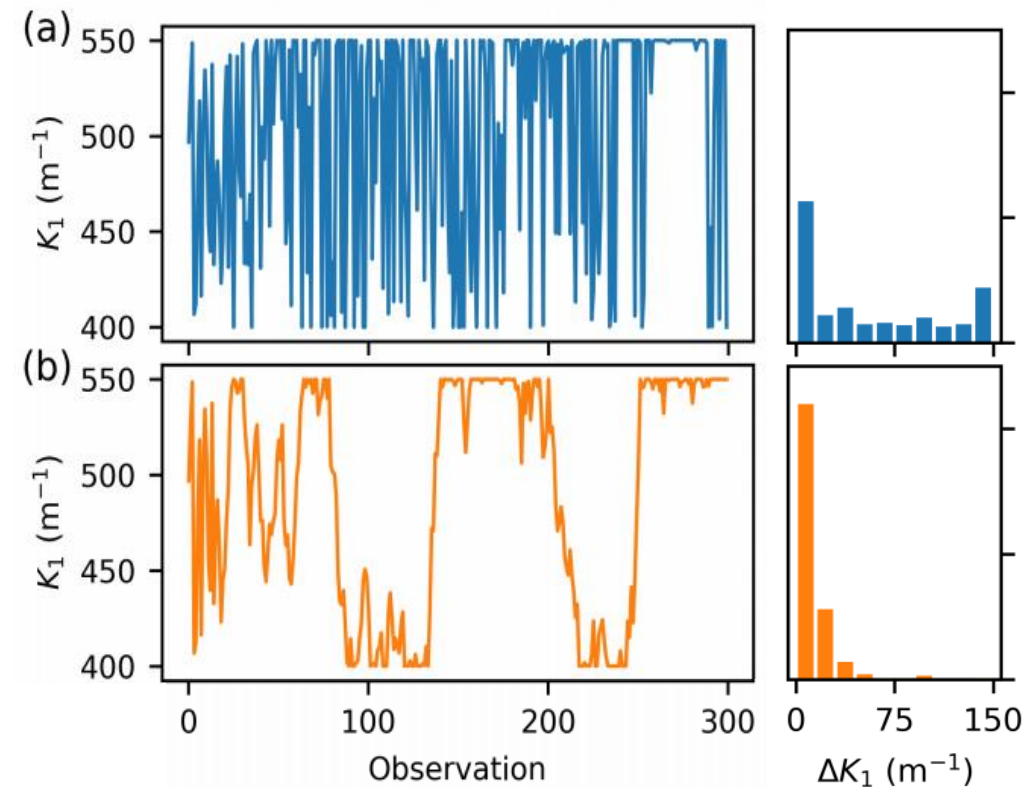- <u>Accuracy</u>: correlated kernels, non-zero prior.

**Acquisition function optimization:**

Also called "BO's inner optimization problem"; wealth of diverse methods were proposed.

- <u>Speed</u>: local optimization (LineBO), parallelism, constrains
- <u>Safety</u>: constrains.

- Motivation

- Model-based vs model-free optimizers

- Bayesian Optimization (BO)
  - Overview
  - Acquisition functions
  - Accounting for constraints
  - Proximal optimization

- **Applications**

- Summary of optimization methods

Maximize X-ray pulse energy simultaneously on 12 quadrupoles with diagonal kernel.

- GP reaches higher optimum
- GP is 4 times faster





RF Gun

Linacs (L0,L1)
220 GeV

Linac (L2)
5 GeV

Linac (L3)
14 GeV

Undulator

X-rays

Learn correlations based on physics/ historical data.

$$k_{\mathrm{RBF}}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_f^2 \exp(-(\boldsymbol{x} - \boldsymbol{x}')^T \boldsymbol{\Sigma}(\boldsymbol{x} - \boldsymbol{x}'))$$

$$\Sigma = \begin{bmatrix} L & 0 \\ 0 & L \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} L_{11} & L_{12} \\ L_{21} & L_{22} \end{bmatrix}$$

$$\Sigma = -\frac{H_{ij}}{2}$$



(a) Ground truth

(b) Isotropic kernel

(c) Correlated kernel

(d) Convergence tests

J. Duris, PRL, 2020
A. Hanuka, PRAB, 2021

Minimize vertical emittance (= maximize beam loss rate)
with 13 skew quadrupole magnets

- GP 10x speedup.





2-3 sec / step
RCDS: 6 sec / step

A. Hanuka, NeurIPS 2019
A. Hanuka, PRAB, 2021

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

mean function





GP with prior mean $m(x) = 0$ (dashed) converged slower to a lower optimum.

A. Hanuka, NeurIPS 2019
A. Hanuka, PRAB, 2021

MOBO - Find the set of Pareto-optimal points in objective space.

Simultaneously minimize transverse emittance & longitudinal bunch length in the AWA photoinjector.

**Advantages:**
- Noise robust.
- Data efficient (statistical model).
- Global guarantees.
- Can handle safety constraints.

**Caveats:**
- Computational efficiency : Maximizing the acquisition function, GP regression.
- Curse of dimensionality.
- Practical: Hyperparameters, difficult to evaluate model fit.

- Motivation
- Model-based vs model-free optimizers
- Bayesian Optimization (BO)
  - Overview
  - Acquisition functions
  - Accounting for constraints
  - Proximal optimization
- Applications
- **Summary of optimization methods**

# Summary of optimization methods

**Instructions:**
- We're going to split to breakout room.
- Each breakout room will fill in the table of comparison for one algorithm (room 1 → algo 1, etc)
- <u>Table of comparison</u>
- Optional answers – Low/Medium/High or Yes/No.
- Choose one presenter to present the table in the main room.

- Sample efficiency
- Computational cost of picking the next point
- Multi-objective
- Sensitivity to local minima
- Sensitivity to noise
- Requires to compute or estimate derivatives of $f$
- Evaluations of $f$ inherently done in parallel
- Hyper-parameters

1. Nelder-Mead
2. Gradient descent
3. Powell / RCDS
4. L-BFGS
5. Genetic algorithm
6. Bayesian optimization

# Summary of optimization methods

| | Nelder-Mead | Gradient descent | Powell / RCDS | L-BFGS | Genetic algorithm | Bayesian optimization |
|---|---|---|---|---|---|---|
| Sample efficiency | | | | | | |
| Computational cost of picking the next point | | | | | | |
| Multi-objective | | | | | | |
| Sensitivity to local minima | | | | | | |
| Sensitivity to noise | | | | | | |

# Summary of optimization methods

| | Nelder-Mead | Gradient descent | Powell / RCDS | L-BFGS | Genetic algorithm | Bayesian optimization |
|---|---|---|---|---|---|---|
| Sample efficiency | Medium | Medium | Medium/high | Medium/high | Low | High |
| Computational cost of picking the next point | Low/Medium | Low | Low | Low | Medium (e.g. sorting) | High (esp. in high dimensions) |
| Multi-objective | No | No | No | No | Yes | Yes |
| | (but can use scalarization) | | | | | |
| Sensitivity to local minima | High | High | High | High | Low | Low (builds a **global** model of $f$) |
| | (but can use multi-start) | | | | | |
| Sensitivity to noise | High | High | High (Powell) Low (RCDS) | High | Medium | Low (can model noise itself) |

# Summary of optimization methods

| | **Nelder -Mead** | **Gradient descent** | **Powell / RCDS** | **L-BFGS** | **Genetic algorithm** | **Bayesian optimization** |
|---|---|---|---|---|---|---|
| Requires to compute or estimate derivatives of $f$ | No | Yes | No | Yes | No | No |
| Evaluations of $f$ *inherently* done in parallel | No | No | No | No | Yes | No |
| Hyper-parameters | Initial simplex | Step size: $\alpha$ (+momentum: $\beta$) | # fit points<br><br>Noise level | Accuracy of hessian estimate | • Population size<br>• Mutation rate<br>• Cross-over rate<br>• Number of generations | • Kernel function<br>• Kernel length scales, amplitude<br>• Noise level<br>• Acquisition function |

1 For the weekend!

2 Lectures only! We still have lab afternoon