

CSCI 699 - ProbGen

Probabilistic and Generative Models

Willie Neiswanger

Lecture 13 - Bayesian Optimization

Today

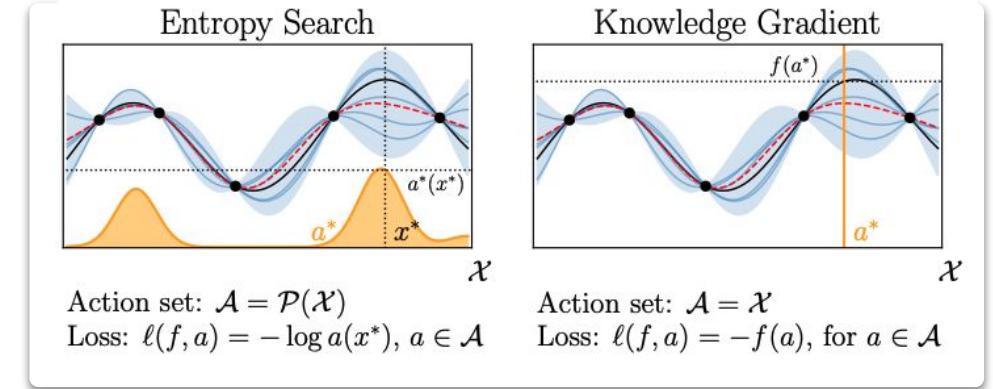
Today

Lecture: Review sequential decision making and active learning, then finish the lecture on Bayesian optimization.

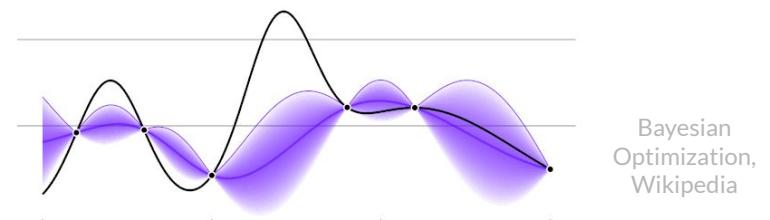
Today

Lecture: Review sequential decision making and active learning, then finish the lecture on Bayesian optimization.

- Decision making under uncertainty.
- Active learning.
- Bayesian Optimization
- E.g., UCB, PI, EI, KG, ES, etc.



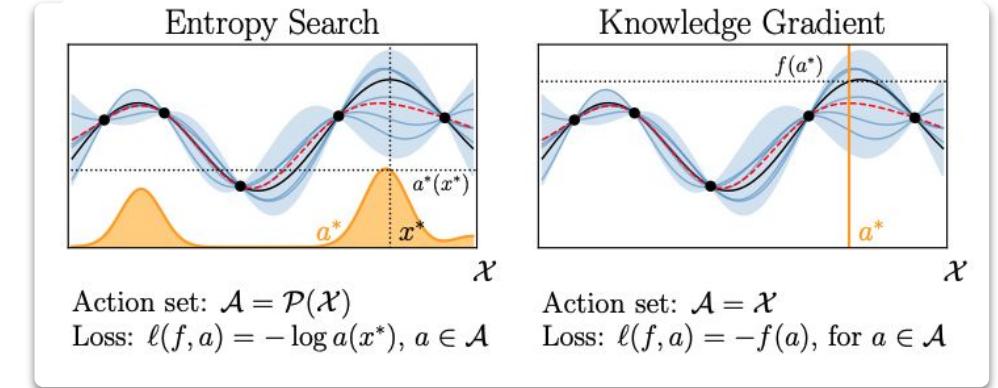
"Generalizing Bayesian Optimization with Decision-theoretic Entropies", 2022



Today

Lecture: Review sequential decision making and active learning, then finish the lecture on Bayesian optimization.

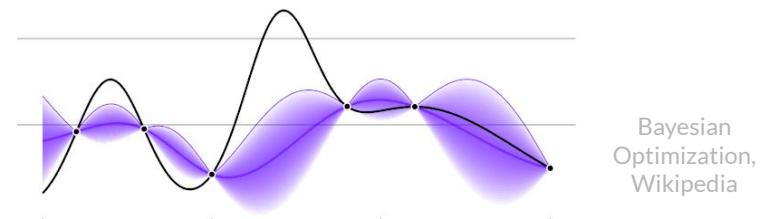
- Decision making under uncertainty.
- Active learning.
- Bayesian Optimization
- E.g., UCB, PI, EI, KG, ES, etc.



"Generalizing Bayesian Optimization with Decision-theoretic Entropies", 2022

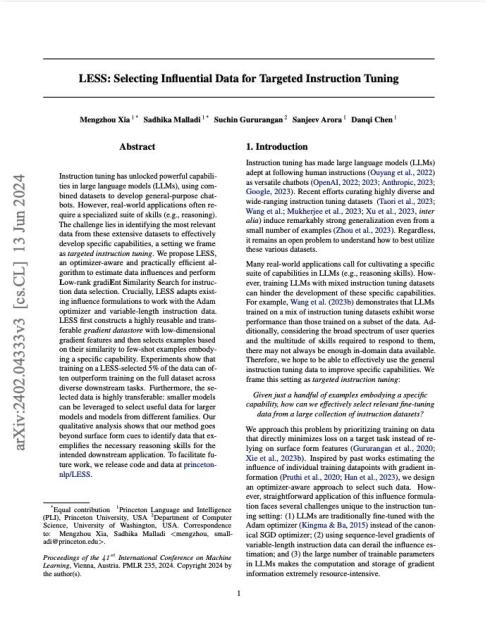
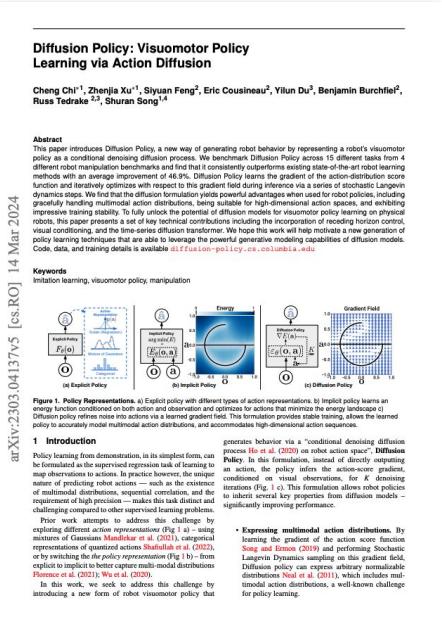
After:

- Paper presentations from students.



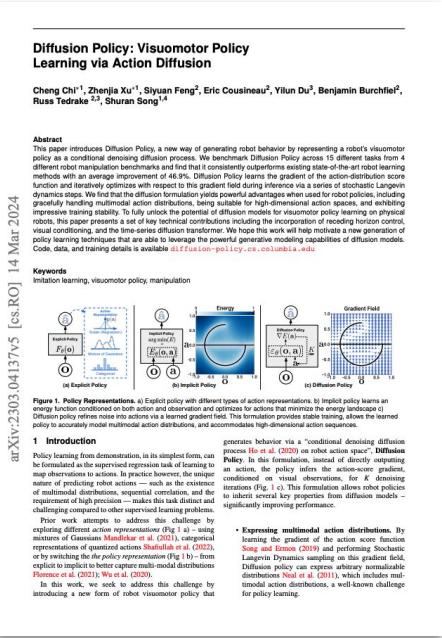
Today

After: Next batch of paper presentations. — 5 today!

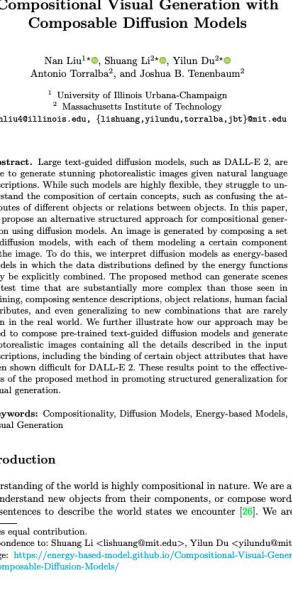


Today

After: Next batch of paper presentations.



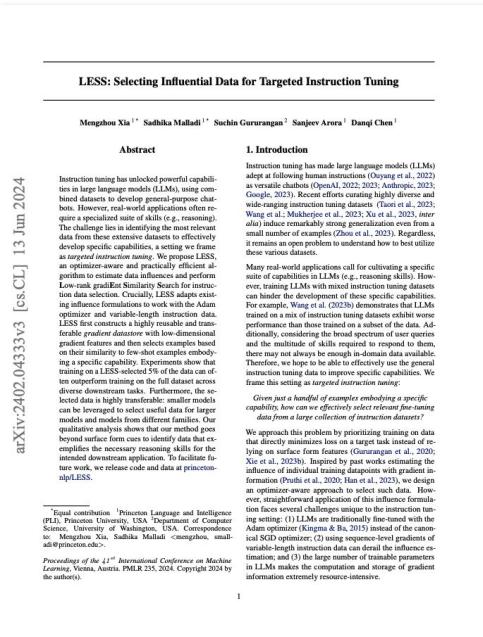
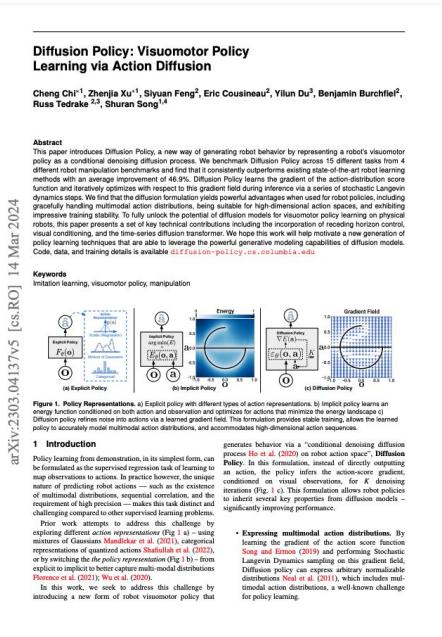
arXiv:2206.01714v6 [cs.CV] 17 Jan 2023



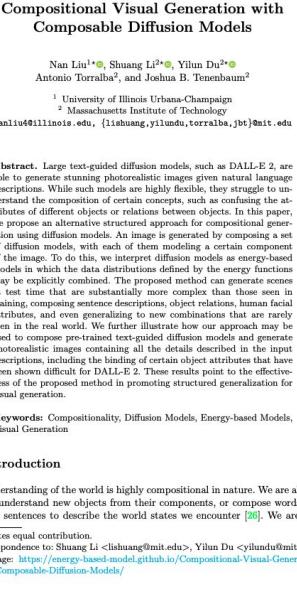
Diffusion Policy – denoising diffusion for robotics

Today

After: Next batch of paper presentations.



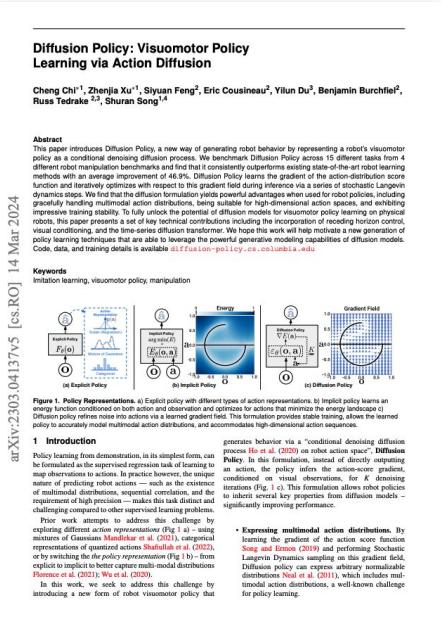
arXiv:2206.01714v6 [cs.CV] 17 Jan 2023



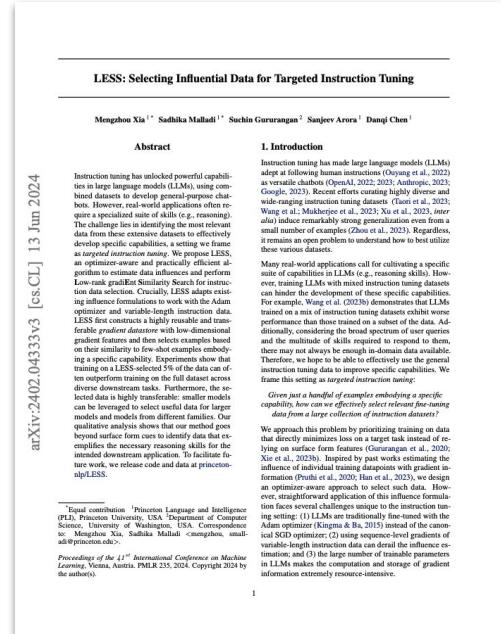
LESS – gradient-based instruction data selection

Today

After: Next batch of paper presentations.



arXiv:2303.04137v5 [cs.RO] 14 Mar 2024



arXiv:2402.04333v3 [cs.CE] 13 Jun 2024



arXiv:2010.02502v4 [cs.LG] 5 Oct 2022



arXiv:2210.01776v2 [q-bio.BM] 11 Feb 2023

arXiv:2206.01714v6 [cs.CV] 17 Jan 2023

Compositional Visual Generation with Composable Diffusion Models

Nan Liu^{1*}, Shuang Li^{2*}, Yilun Du^{2*}, Antonio Torralba², and Joshua B. Tenenbaum²

¹ University of Illinois Urbana-Champaign

² Massachusetts Institute of Technology

nanliu4@illinois.edu, {liushuang,yilundu,torralba,jbt}@mit.edu

Abstract. Large text-guided diffusion models, such as DALL-E 2, are able to generate stunning photorealistic images given natural language descriptions. While search engines highly flexibly, they struggle with maintaining the precision of certain objects such as ligands, the attributes of different objects or relations between objects. In this paper, we propose an alternative structured approach for compositional generation using diffusion models. An image is generated by composing a set of diffusion models, with each of them modeling a certain component of the ligand. In the domain of molecular docking, we find that energy-based models may be explicitly combined. The proposed method can generate scenes that are substantially more complex than those seen in training, composing sentences, descriptions, object relations, human facial attributes, etc., generating images from the perspective of the target object. The paper also discusses the limitations of the scoring function used in the proposed method. We further illustrate how our approach can be used to compose pre-trained text-guided diffusion models and generate photorealistic images containing all the details described in the input descriptions, including the binding of certain object attributes that have been shown difficult for DALL-E 2. These results point to the effectiveness of the proposed method in promoting structured generalization for visual generation.

Keywords: Compositional, Diffusion Models, Energy-based Models, Visual Generation

1 Introduction

The biological functions of proteins can be modulated by small molecule ligands (such as drugs) binding to them. This is a crucial task in computational drug design is *molecular docking*—predicting the position and orientation of a ligand when it binds to a target protein—from which often the ligand’s any might be inferred. Traditionally approaches such as (Korpi & Ochoa, 2010; Halgren et al., 2004) rely on scoring-functions that estimate the correctness of a proposed structure pose, and an optimization algorithm to find the global minimum of the scoring function. If the scoring function is not and the landscape of the scoring function rugged, these methods tend to be slow and inaccurate, especially for high-throughput workflows.

Recent works (Stärk et al., 2022; Lu et al., 2022) have demonstrated the ability to produce samples comparable to GANs, without having to perform adversarial training. To achieve this, many denoising autoregressive models (DDPMs) (van den Oord et al., 2016b) and normalizing flows (Rezende & Mohamed, 2015; Dinh et al., 2016). However, GANs require very specific choices in optimization and architectures that are hard to tune. Normalizing flows are also slower than DDPMs (Ho et al., 2020; Ho et al., 2021), and could fail to cover modes of the data distribution (Zhai et al., 2021).

Recent work on iterative generative models (Ho et al., 2021a), such as denoising diffusion probabilities (DDIMs) (Ho et al., 2021b) have demonstrated the ability to produce samples comparable to GANs, without having to perform adversarial training. To achieve this, many denoising autoregressive models (DDPMs) (van den Oord et al., 2016b) and normalizing flows (Rezende & Mohamed, 2015; Dinh et al., 2016). However, GANs require very specific choices in optimization and architectures that are hard to tune. Normalizing flows are also slower than DDPMs (Ho et al., 2020; Ho et al., 2021), and could fail to cover modes of the data distribution (Zhai et al., 2021).

A critical drawback of these models is that they require many iterations to produce a high quality image. For example, the DDPMs (Ho et al., 2021b) need to run 1000 steps to produce a high quality image. This is a significant challenge for this type of model, as it requires many iterations to produce a high quality image. By contrast, the proposed DiFFDOCK only needs to run 10 steps to produce a high quality image. This is a significant improvement over the DDPMs (Ho et al., 2021b).

To close this efficiency gap between DDPMs and GANs, we present denoising diffusion implicit models (DDIMs). DDIMs are implicit probabilistic models (Mohamed & Lakshminarayanan, 2016) and are closely related to DDPMs, in the sense that they are trained with the same objective function.

*Equal contribution. Correspondence to {gcorvo, hstark, bjing}@mit.edu

**indicates equal contribution.

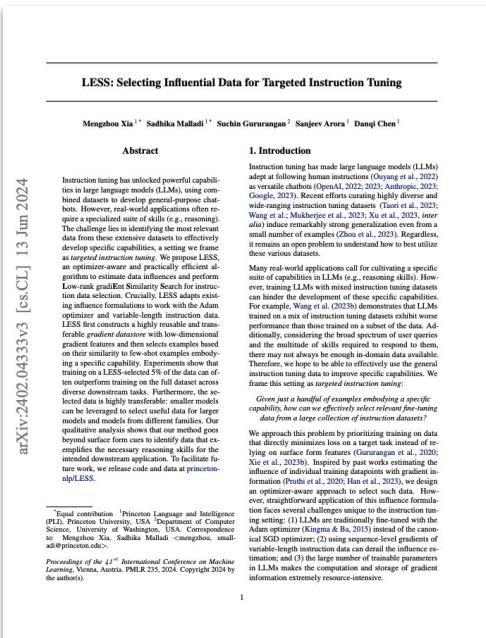
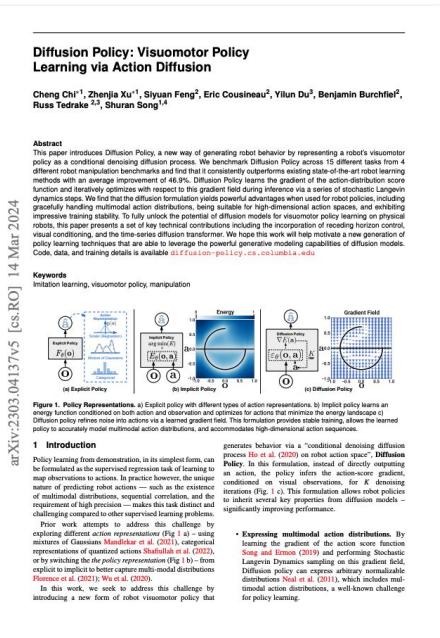
Correspondence to: Shuang Li <liushuang@mit.edu>, Yilun Du <yilundu@mit.edu>

Webpage: <https://energy-based-model.github.io/Compositional-Visual-Generation-with-Composable-Diffusion-Models/>

DDIM – classic method for efficient generation

Today

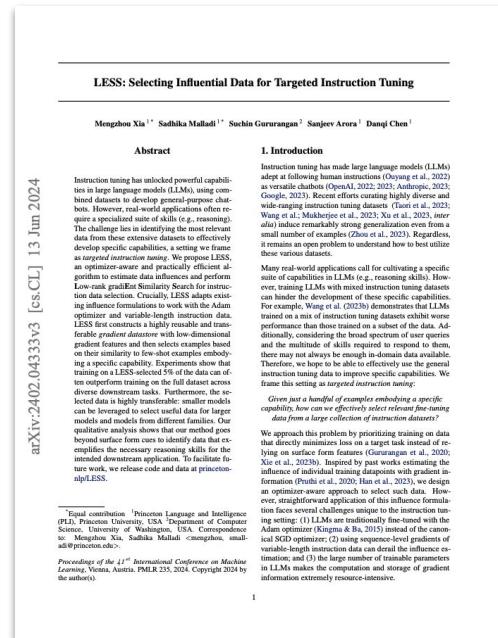
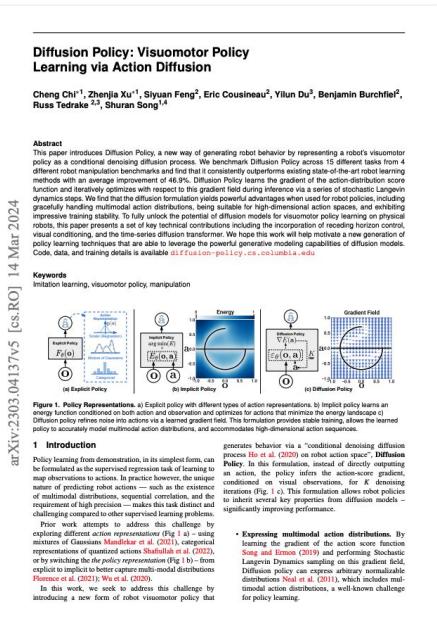
After: Next batch of paper presentations.



DiffDock – diffusion for molecular docking

Today

After: Next batch of paper presentations.



arXiv:2206.01714v6 [cs.CV] 17 Jan 2023

Compositional Visual Generation with Composable Diffusion Models

Nan Liu^{1*}, Shuang Li^{2*}, Yilun Du^{2*}, Antonio Torralba², and Joshua B. Tenenbaum²
¹ University of Illinois Urbana-Champaign
² Massachusetts Institute of Technology
nanliu4@illinois.edu, {liushuang,yilundu,torralba,jbt}@mit.edu

Abstract. Large text-guided diffusion models, such as DALL-E 2, are able to generate stunning photorealistic images given natural language descriptions. While such models are highly flexible, they struggle under-composed inputs that require reasoning about multiple objects and their attributes. In this paper, we propose an alternative structured approach for compositional generation using diffusion models. An image is generated by composing a set of diffusion models, with each of them modeling a certain component of the scene. In the domains in which the data distributions defined by the encoding functions may be explicitly combined, the proposed method can generate scenes at test time that are substantially more complex than those seen in training, composing sentences, descriptions, object relations, human facial attributes, and even generating images from scratch. Evaluation approaches are also provided, showing that our approach can be used to compose pre-trained text-guided diffusion models and generate photorealistic images containing all the details described in the input descriptions, including the binding of certain object attributes that have been shown difficult for DALL-E 2. These results point to the effectiveness of the proposed method in promoting structured generalization for visual generation.

Keywords: Compositional, Diffusion Models, Energy-based Models, Visual Generation

1 Introduction

The biological functions of proteins can be modulated by small molecule ligands (such as drugs) binding to them. This is a crucial task in computational drug design—predicting the position, orientation, and binding of a ligand when it binds to a target protein—from which often the ligand’s 3D shape may be inferred. Traditionally, approaches such as molecular docking (Trott & Olson, 2010; Halgren et al., 2004) rely on scoring functions that estimate the correctness of a proposed structure pose, and an optimization algorithm is used for the minimization of the scoring function. If the scoring function is robust and the landscape of the scoring function is smooth, these methods tend to be too slow and inaccurate, especially for high-throughput workflows.

Recent works (Stark et al., 2022; Lu et al., 2022) have developed deep learning models to predict the binding of a ligand to a target protein. Compared to search-based methods, these have yet to demonstrate significant improvements in accuracy. We argue this may be because the regression-based paradigm corresponds implicitly to a “forward” diffusion process, while the scoring function corresponds to a “backward” one. Energy metrics restrict the likelihood of the data under the predictive model rather than a regression loss. That is, the forward diffusion process—given a ligand and target protein structure, we learn a distribution over ligand poses.

To this end, we develop DifffDock, a diffusion generative model (DGM) over the space of ligand poses for molecular docking. We define a diffusion process over the degrees of freedom involved in docking. The forward diffusion process is a sequence of random walks over the degrees of freedom in the pocket, and the torsion angles describing the conformation. DifffDock samples poses by running the learned inverse diffusion process, which iteratively transforms an initial noisy pose distribution of ligands into a learned model (see Figure 1). Interestingly, this process can be viewed as the progressive refinement of random poses via updates of their translations, rotations, and torsion angles.

*Equal contribution. Correspondence to: {gcorso, hstark, bjing}@mit.edu

[†] indicates equal contribution.

Correspondence to: Shuang Li <liushuang@mit.edu>, Yilun Du <yilundu@mit.edu>

Webpage: <https://energy-based-model.github.io/Compositional-Visual-Generation-with-Composable-Diffusion-Models/>

Combining diffusion models for compositional generation

Course Projects

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

- **Introduction to the topic and motivation for your project:** the problem you are trying to solve, why it's useful in the world/downstream applications, why it's difficult to solve.

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

- **Introduction to the topic and motivation for your project:** the problem you are trying to solve, why it's useful in the world/downstream applications, why it's difficult to solve.
- **Prior work:** previous methods that are related to your work, high-level overview on any prior material needed to understand your work.

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

- **Introduction to the topic and motivation for your project:** the problem you are trying to solve, why it's useful in the world/downstream applications, why it's difficult to solve.
- **Prior work:** previous methods that are related to your work, high-level overview on any prior material needed to understand your work.
- **Technical contribution:** your main *contributions, methods, or technical work* in this project.

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

- **Introduction to the topic and motivation for your project:** the problem you are trying to solve, why it's useful in the world/downstream applications, why it's difficult to solve.
- **Prior work:** previous methods that are related to your work, high-level overview on any prior material needed to understand your work.
- **Technical contribution:** your main *contributions, methods, or technical work* in this project.
- **Experimental results:** the experiments you carried out and results you achieved, tying it back to your original project goals and motivation.

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

- **Introduction to the topic and motivation for your project:** the problem you are trying to solve, why it's useful in the world/downstream applications, why it's difficult to solve.
- **Prior work:** previous methods that are related to your work, high-level overview on any prior material needed to understand your work.
- **Technical contribution:** your main *contributions, methods, or technical work* in this project.
- **Experimental results:** the experiments you carried out and results you achieved, tying it back to your original project goals and motivation.
- **Discussion:** Thoughts/debrief on your project – how did the idea work out, and did you accomplish your goal? Any next steps for the project? Any new perspectives on promising directions in this subfield?

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

- First presentation day: **April 25**
 - Groups: 1, 3, 9 → Thank you for presenting first!

Course Projects

Final Presentations – **25 to 30 minute** presentation! Cover the full project scope, including:

- First presentation day: **April 25**
 - Groups: 1, 3, 9 → Thank you for presenting first!

- First presentation day: **May 2**
 - Groups: 2, 4, 5, 6, 7, 8 → Thank you also :-).

We will stay on schedule to get through everyone!

Course Projects

Final Reports

Course Projects

Final Reports – 8-10 pages long. Similar scope as final presentation, including:

Course Projects

Final Reports – 8-10 pages long. Similar scope as final presentation, including:

- **Introduction**
- **Related Work**
- **Methods and Technical Contribution**
- **Experiments and Results**
- **Discussion**

Course Projects

Final Reports – 8-10 pages long. Similar scope as final presentation, including:

- **Introduction**
- **Related Work**
- **Methods and Technical Contribution**
- **Experiments and Results**
- **Discussion**

Due May 9th!

Course Projects

Final Reports – 8-10 pages long. Similar scope as final presentation, including:

- **Introduction**
- **Related Work**
- **Methods and Technical Contribution**
- **Experiments and Results**
- **Discussion**

Due May 9th!

Probabilistic and Generative Models: Final Report Template

Author List Test*
Department List
University of Southern California
Los Angeles, CA 90007
email@usc.edu

1 Introduction
This is a template for the course project final report for *CSCI 699: Probabilistic and Generative Models*. You are welcome to adjust the following section titles, though these give a guideline for what content is expected. Note that the final report should aim to be around eight-to-ten pages long in total.
Due date: Each group should have one team member upload the midway report as a PDF on Brightspace by May 9, 2025, end-of-day.
In this first section (*Introduction*), you should include an introduction for your project, including the motivation for your project, the problem you are trying to solve, why it's useful in the world/downstream applications, why it's difficult to solve.

2 Related Work
Include a literature review of prior related work here.

3 Methods and Technical Contribution
You are welcome to choose a better name for this section :-). This section should include your main contributions, methods, and/or technical work in this project. Ideally, you should also provide an overview of any background technical material needed to understand your work.

4 Experiments and Results
This section should include the experiments you carried out and results you achieved, tying it back to your original goals/motivation for this project. Make sure to describe details of any evaluation metrics you use, and baseline methods you compare against.

5 Discussion
You should also include a discussion section, covering things like: how did the idea work out, and did you accomplish your goal? Any next steps/future plans for the project? Any new perspectives on promising directions in this subfield?

*More information can go here.

Have LaTeX template to share

Active Learning and Bayesian Optimization

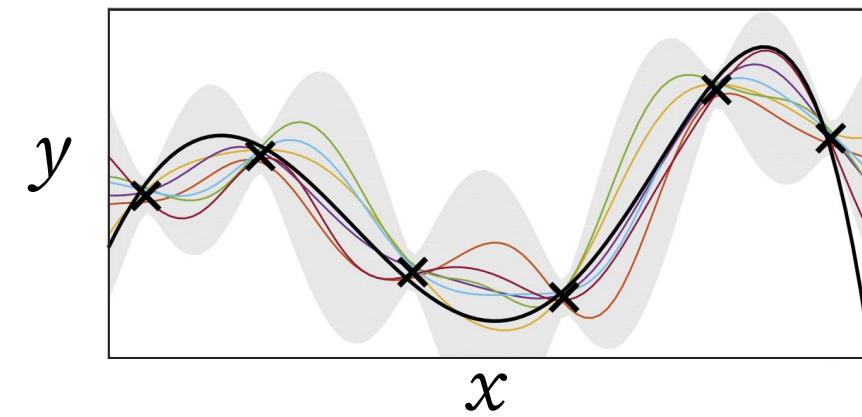
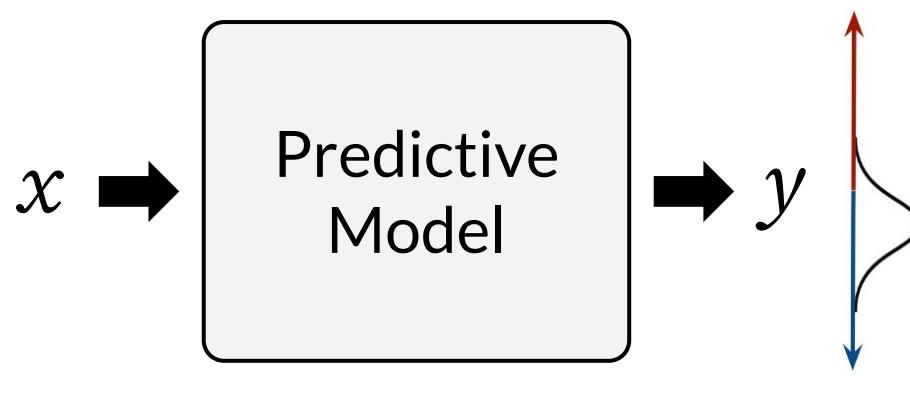
Bayesian Optimization – Review

Bayesian Optimization – Review

What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.

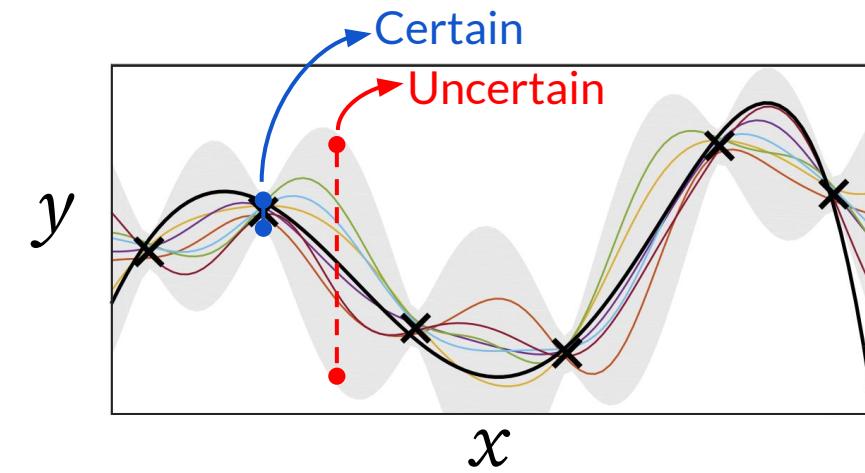
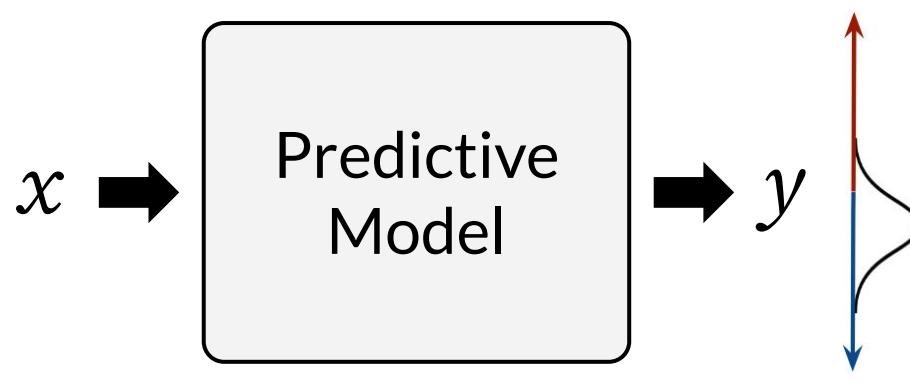
Bayesian Optimization – Review

What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.



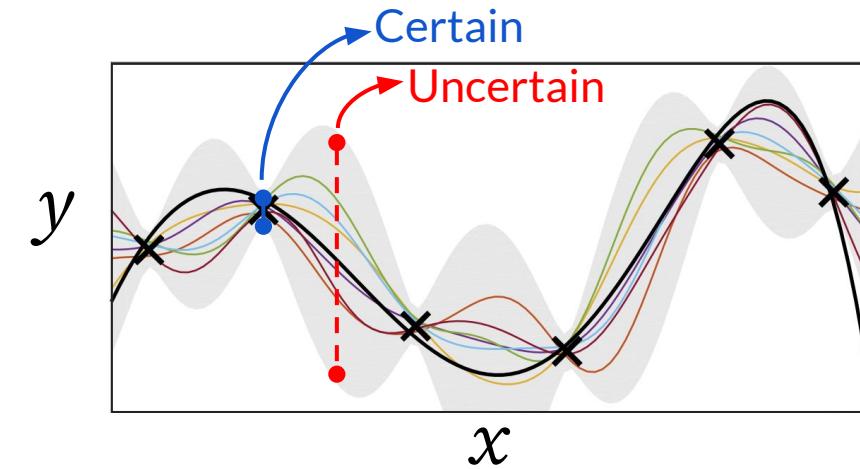
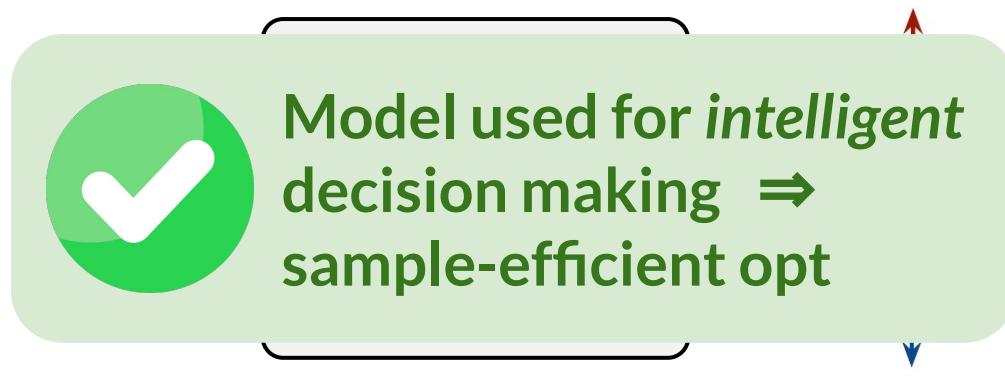
Bayesian Optimization – Review

What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.



Bayesian Optimization – Review

What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.



Bayesian Optimization – Algorithm Overview

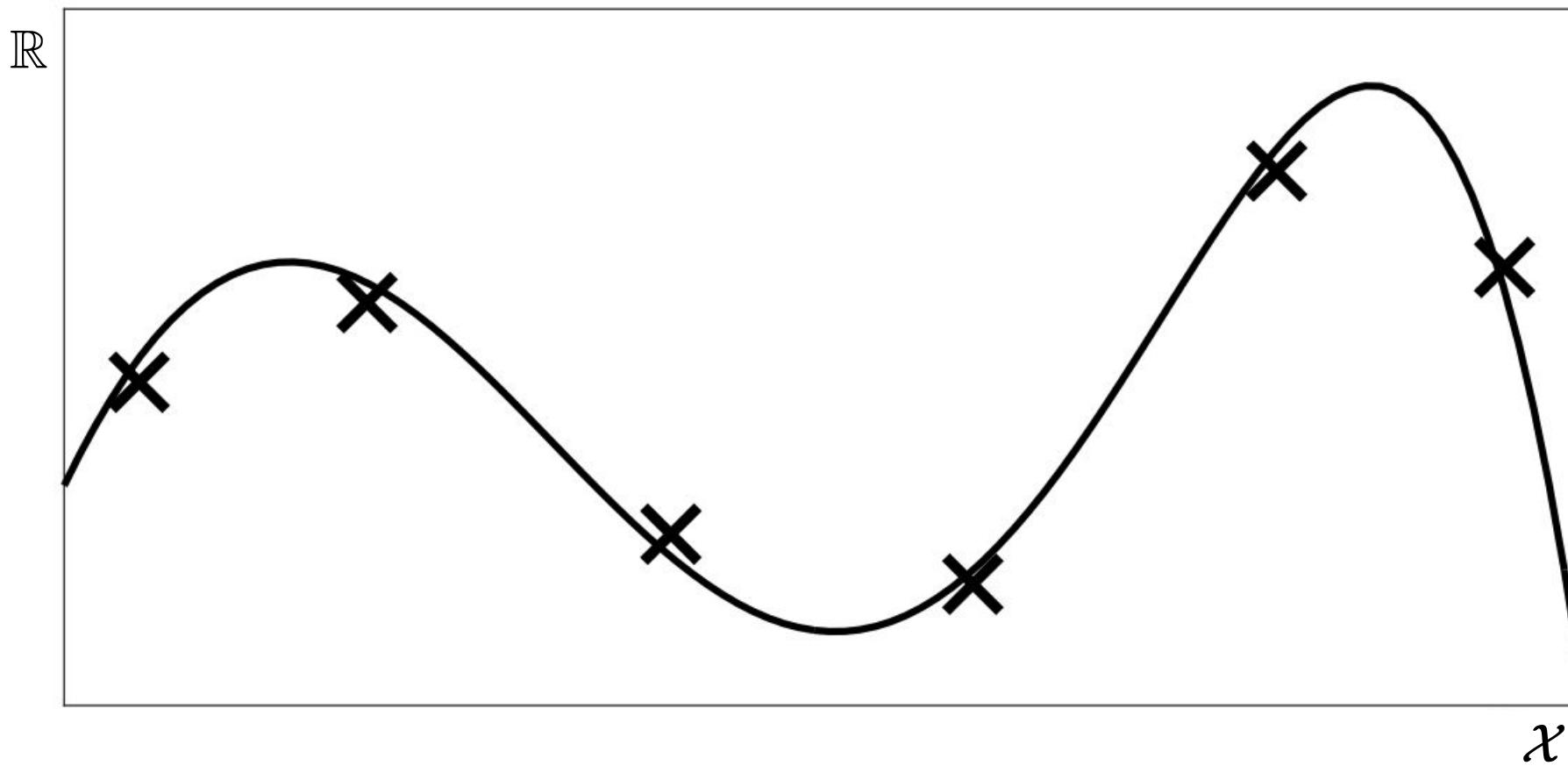
At each iteration of BO, three steps:

Bayesian Optimization – Algorithm Overview

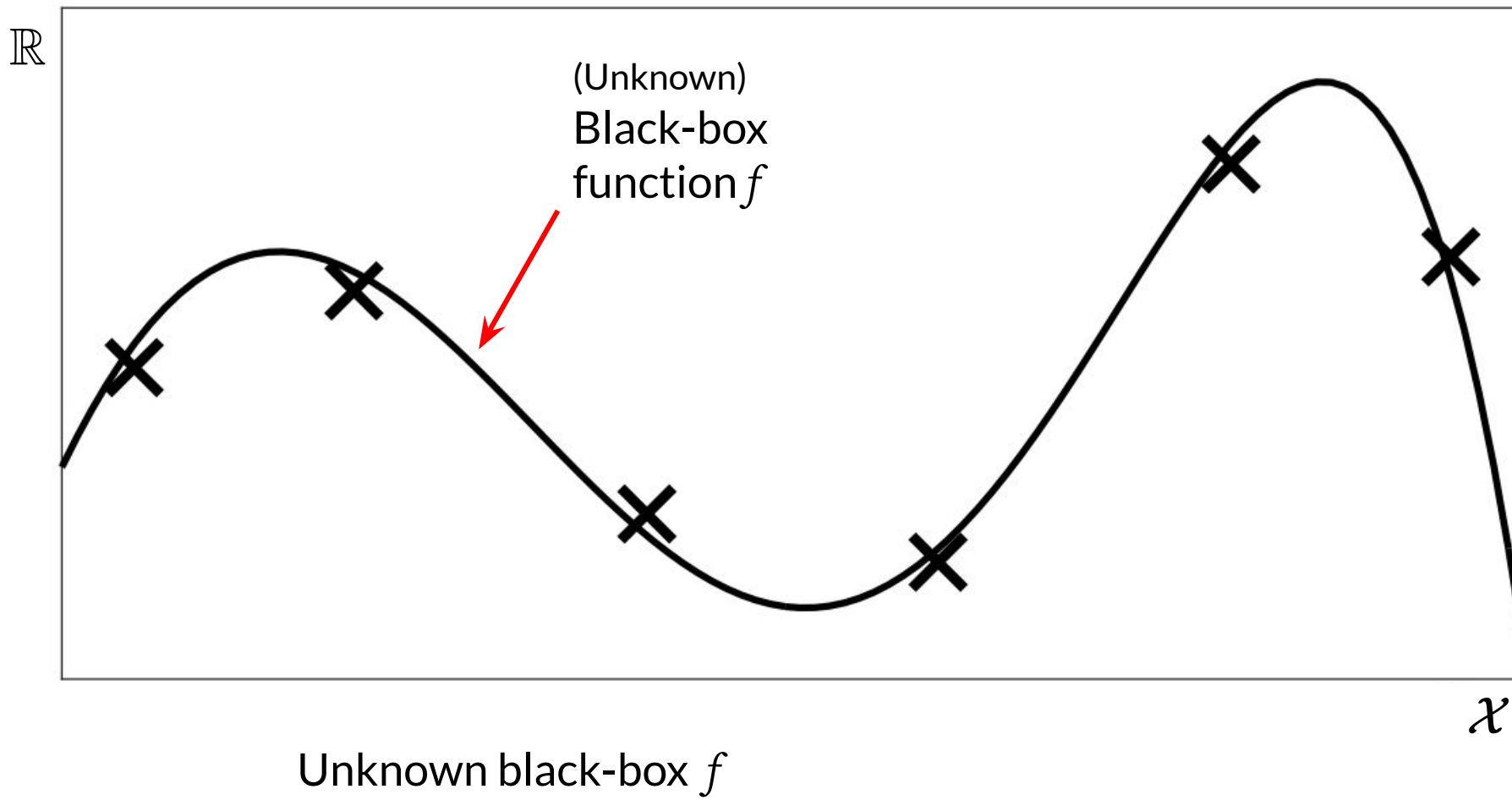
At each iteration of BO, three steps:

1. Fit model to dataset.
2. Optimize an *acquisition function* to select next point to query.
3. Query black-box function.

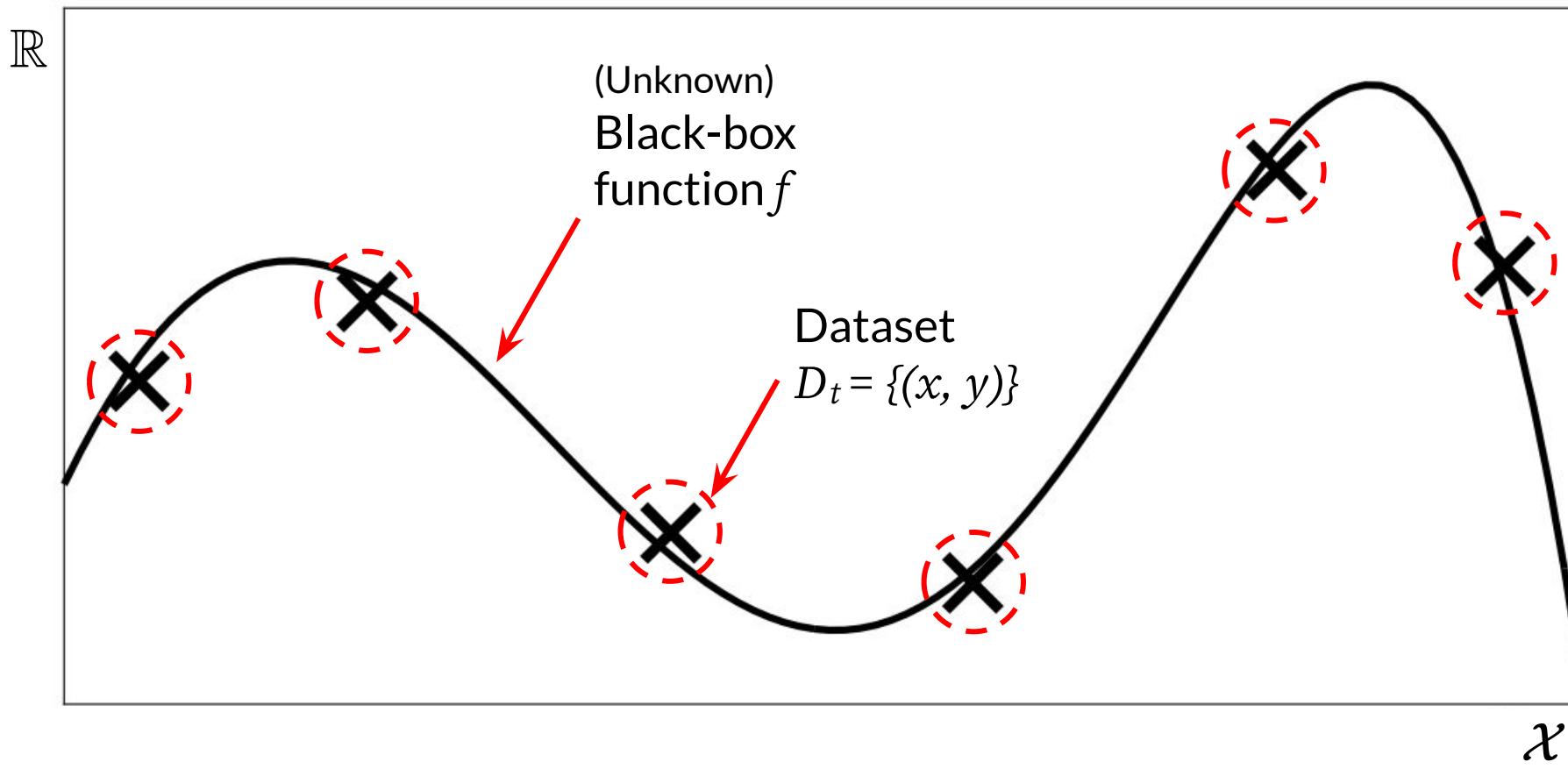
Bayesian Optimization – Visualization



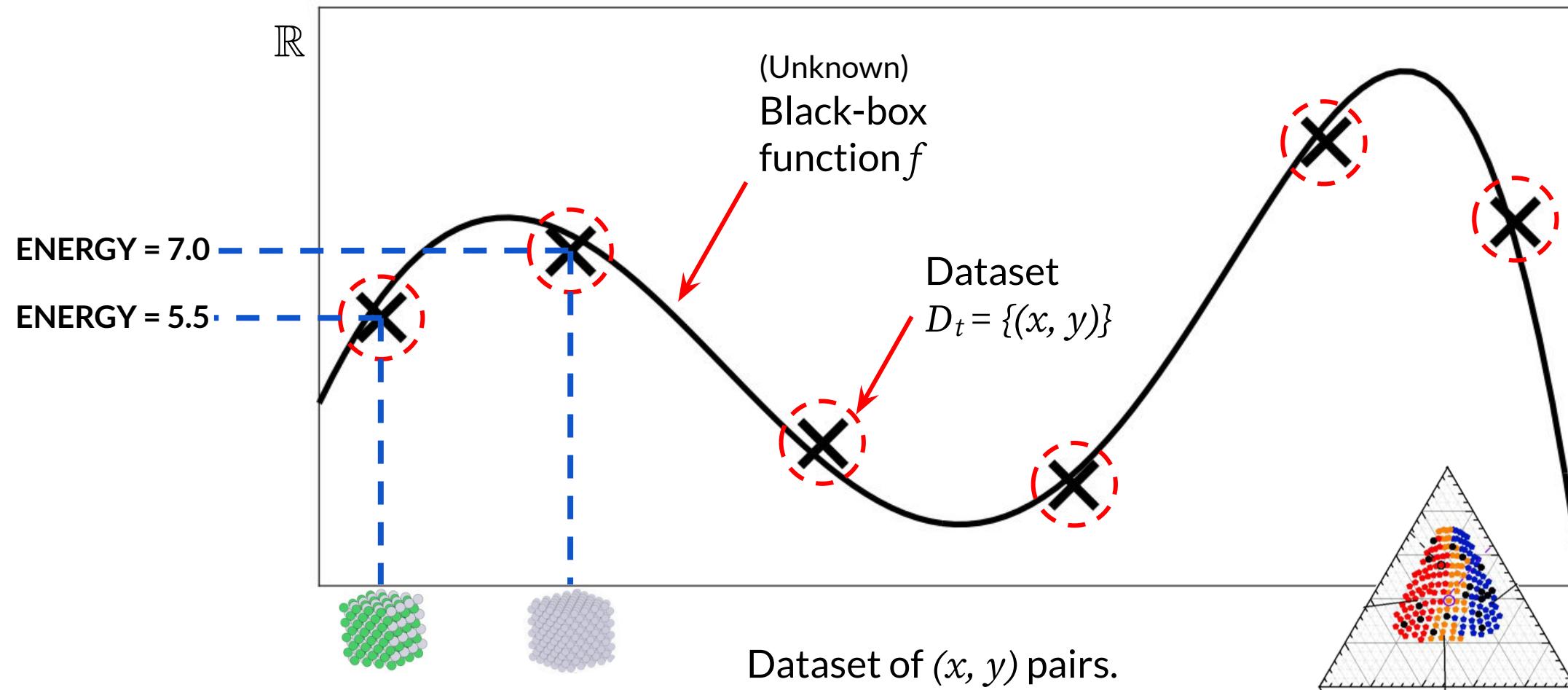
Bayesian Optimization – Visualization



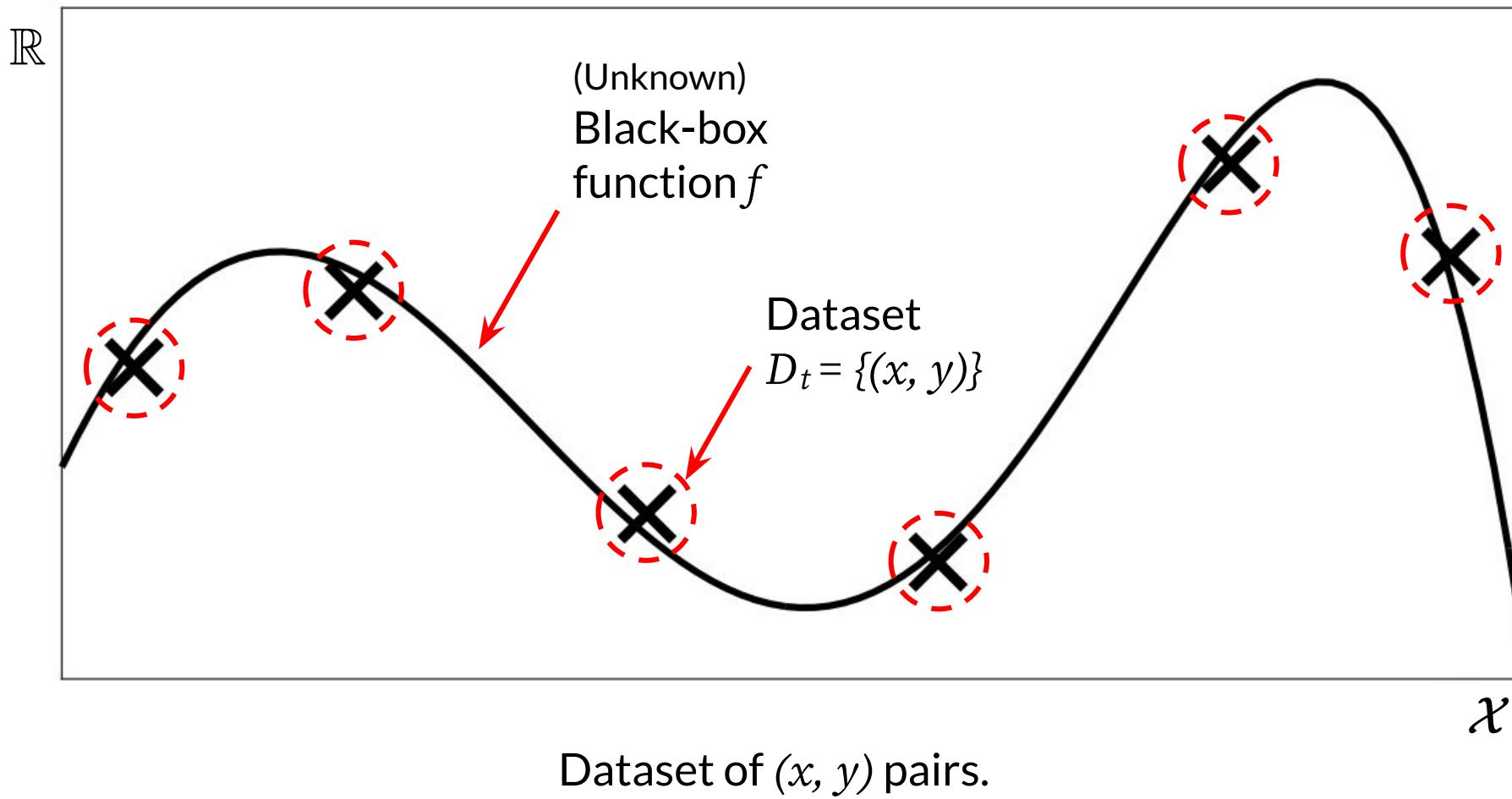
Bayesian Optimization – Visualization



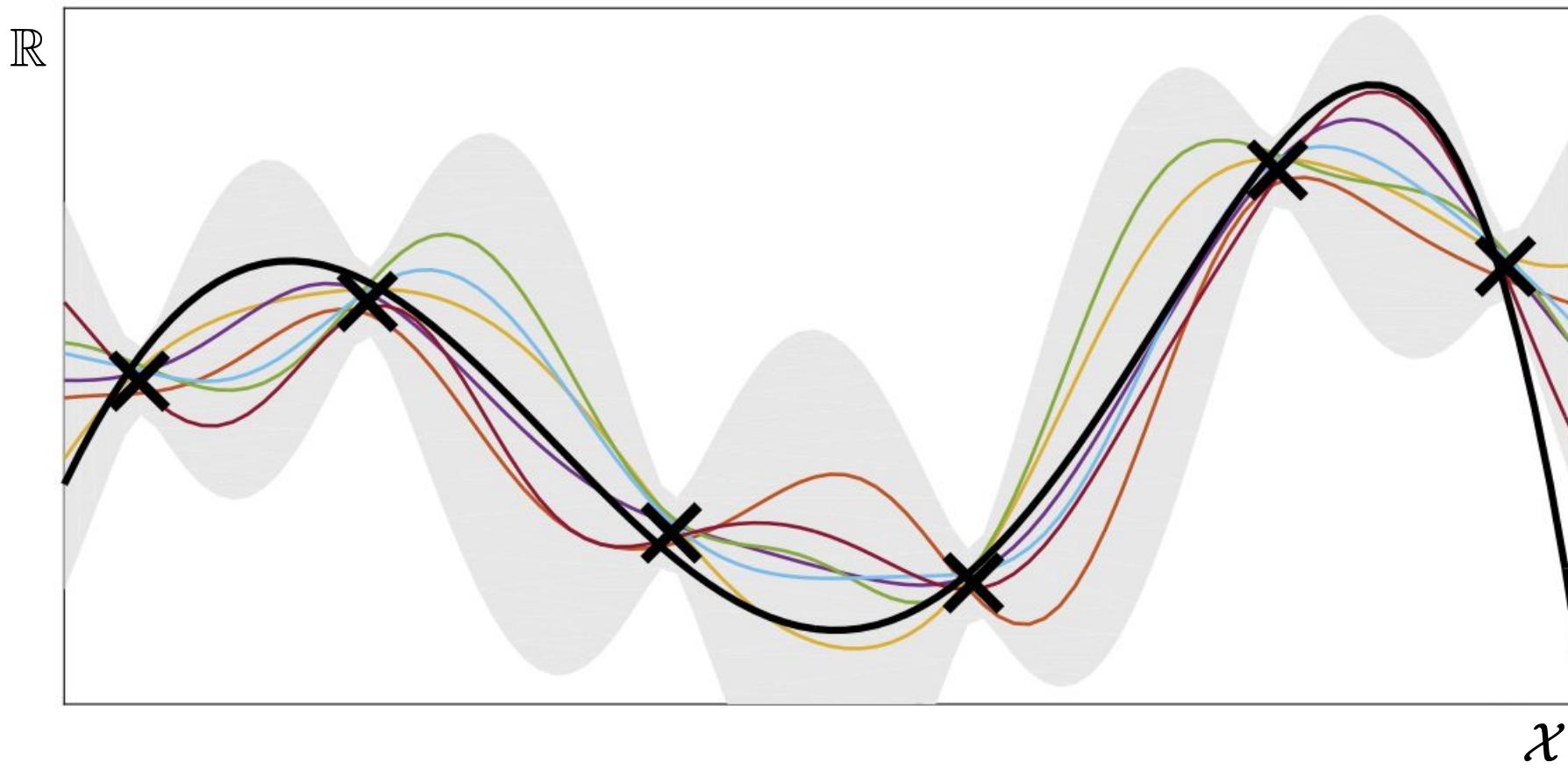
Bayesian Optimization – Visualization



Bayesian Optimization – Visualization

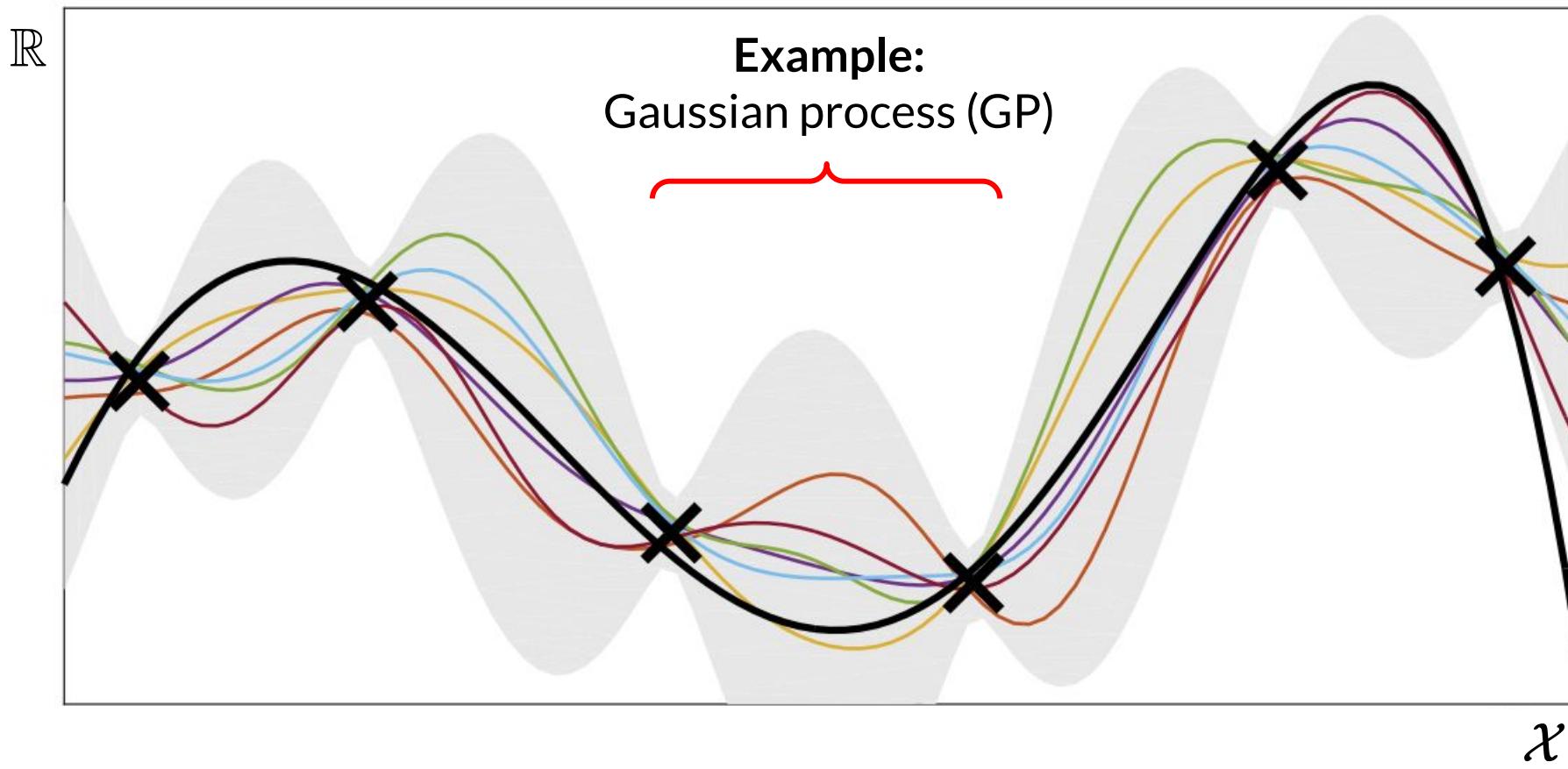


Bayesian Optimization – Visualization



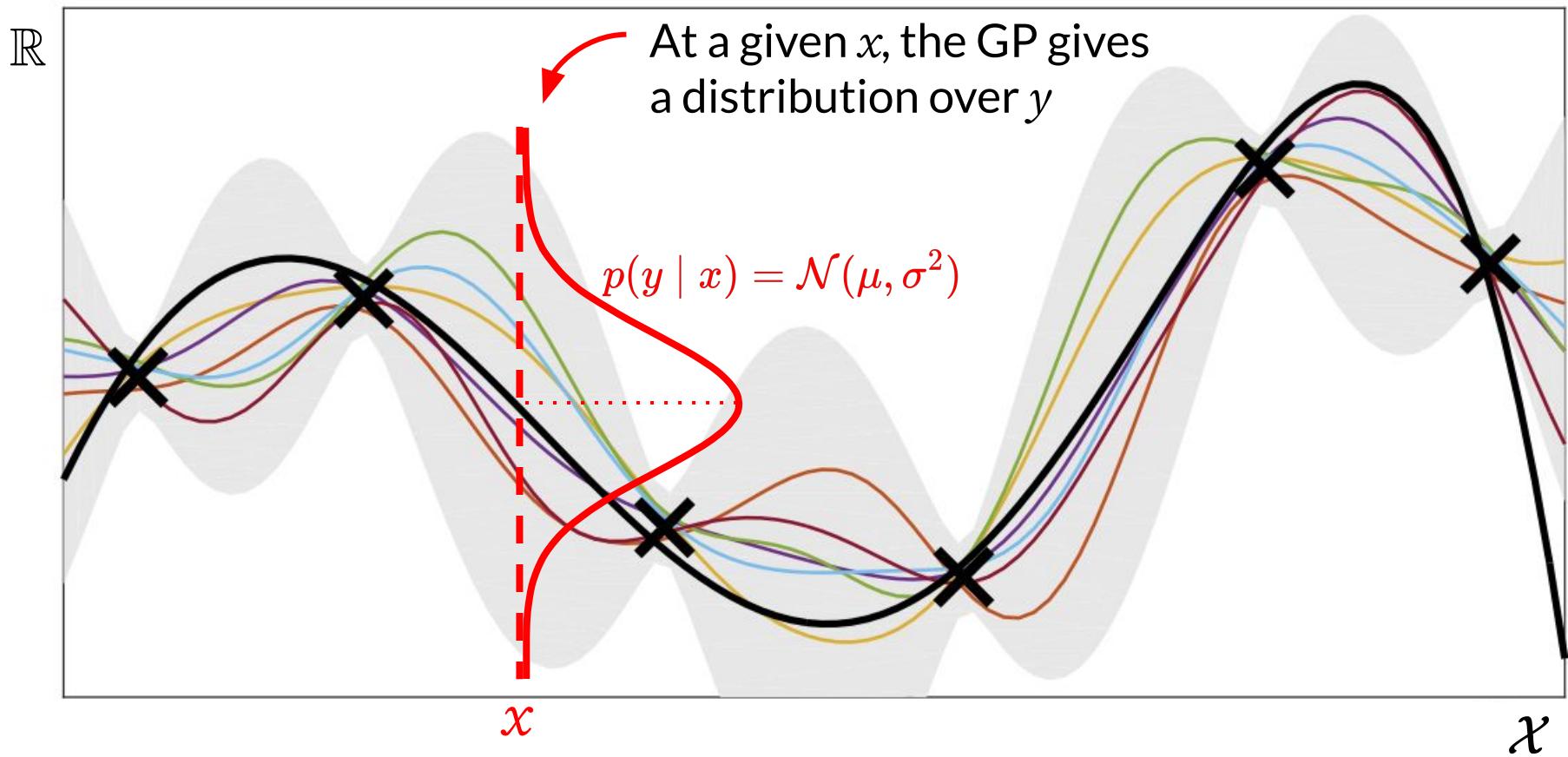
Given this dataset, can fit probabilistic model.

Bayesian Optimization – Visualization



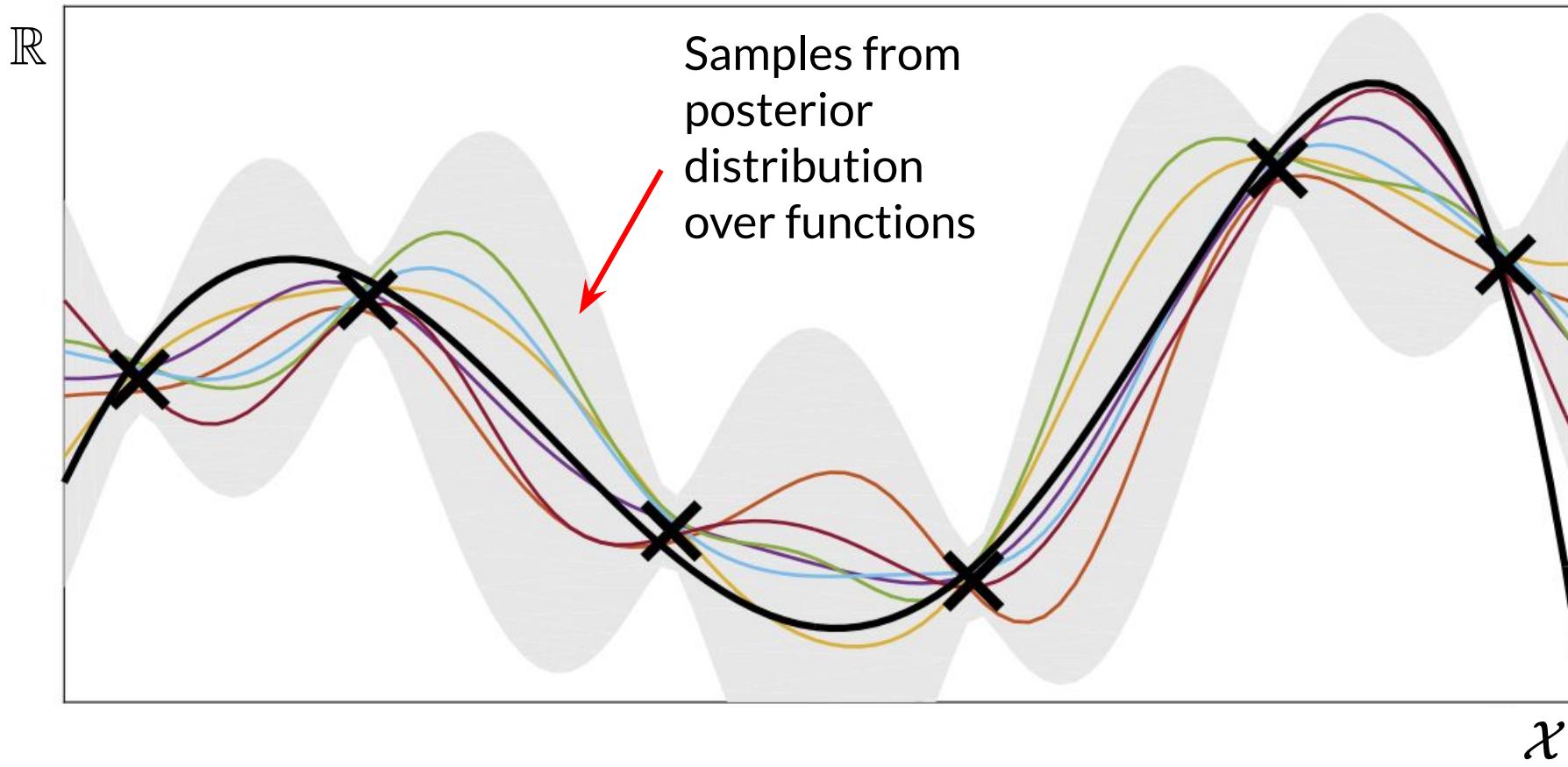
A popular choice for this model is a **Gaussian Process**.

Bayesian Optimization – Visualization



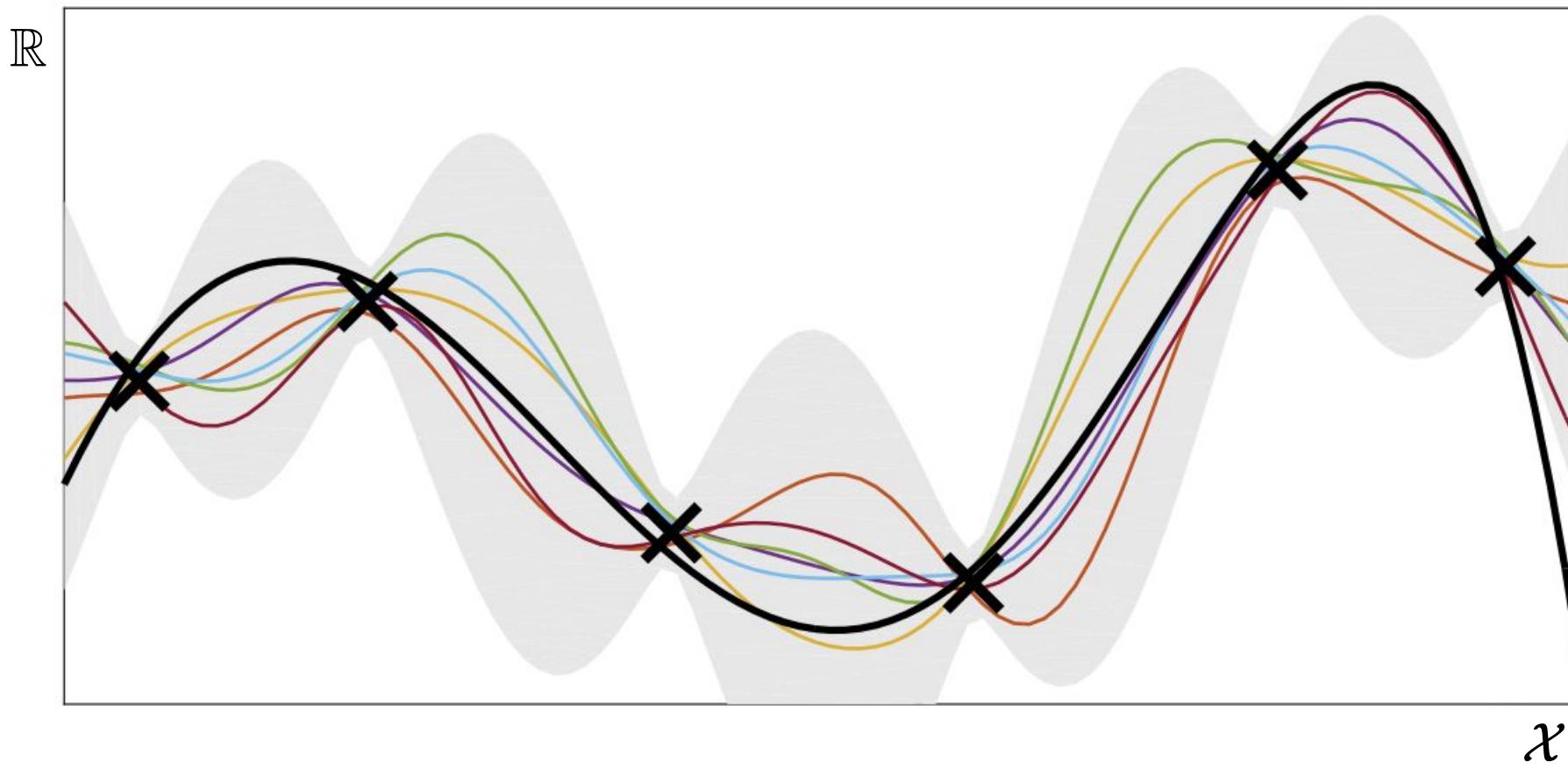
A popular choice for this model is a **Gaussian Process**.

Bayesian Optimization – Visualization



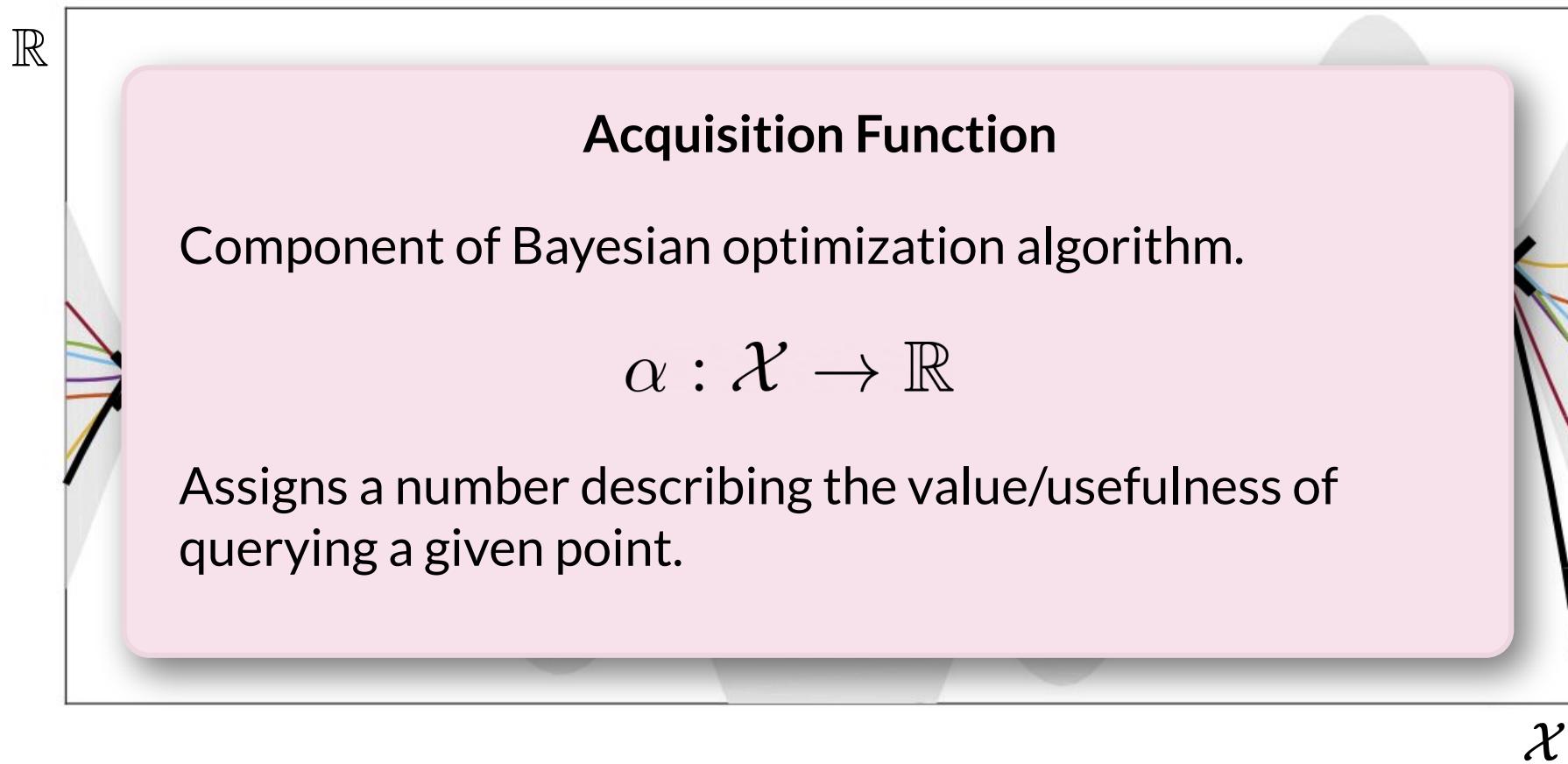
A popular choice for this model is a **Gaussian Process**.

Bayesian Optimization – Visualization

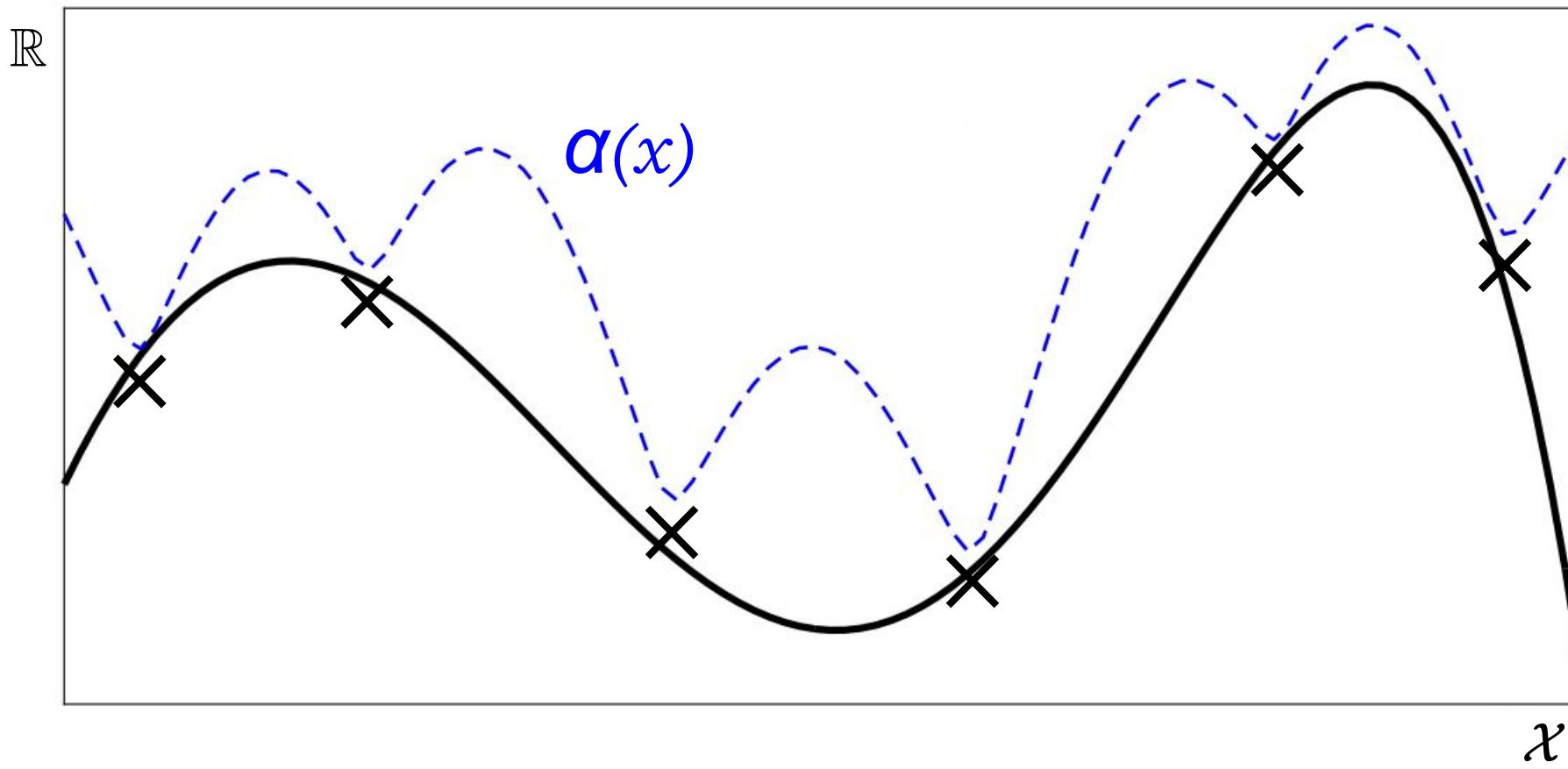


Define acquisition function using model.

Bayesian Optimization – Visualization

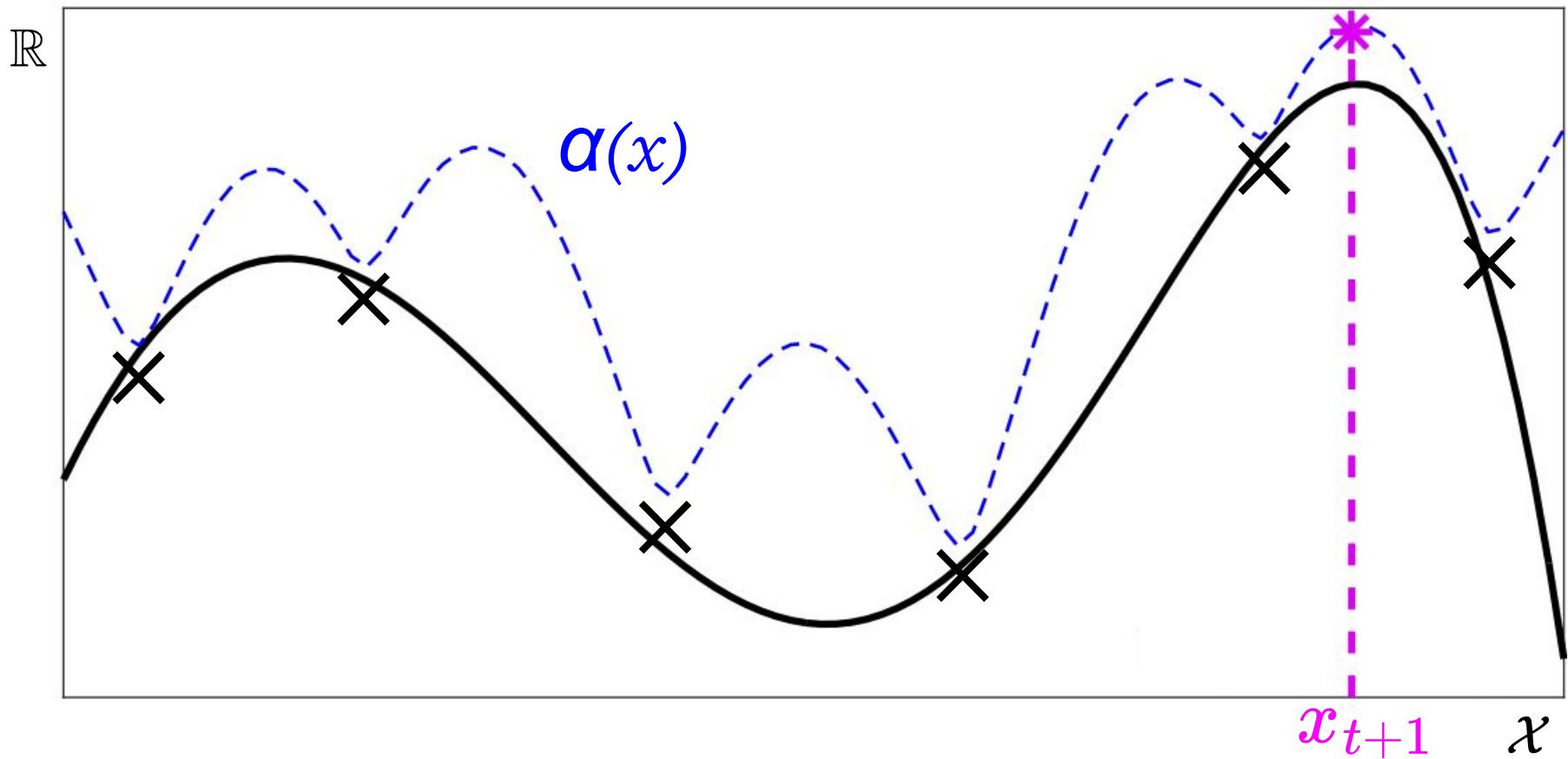


Bayesian Optimization – Visualization



Define acquisition function using model.

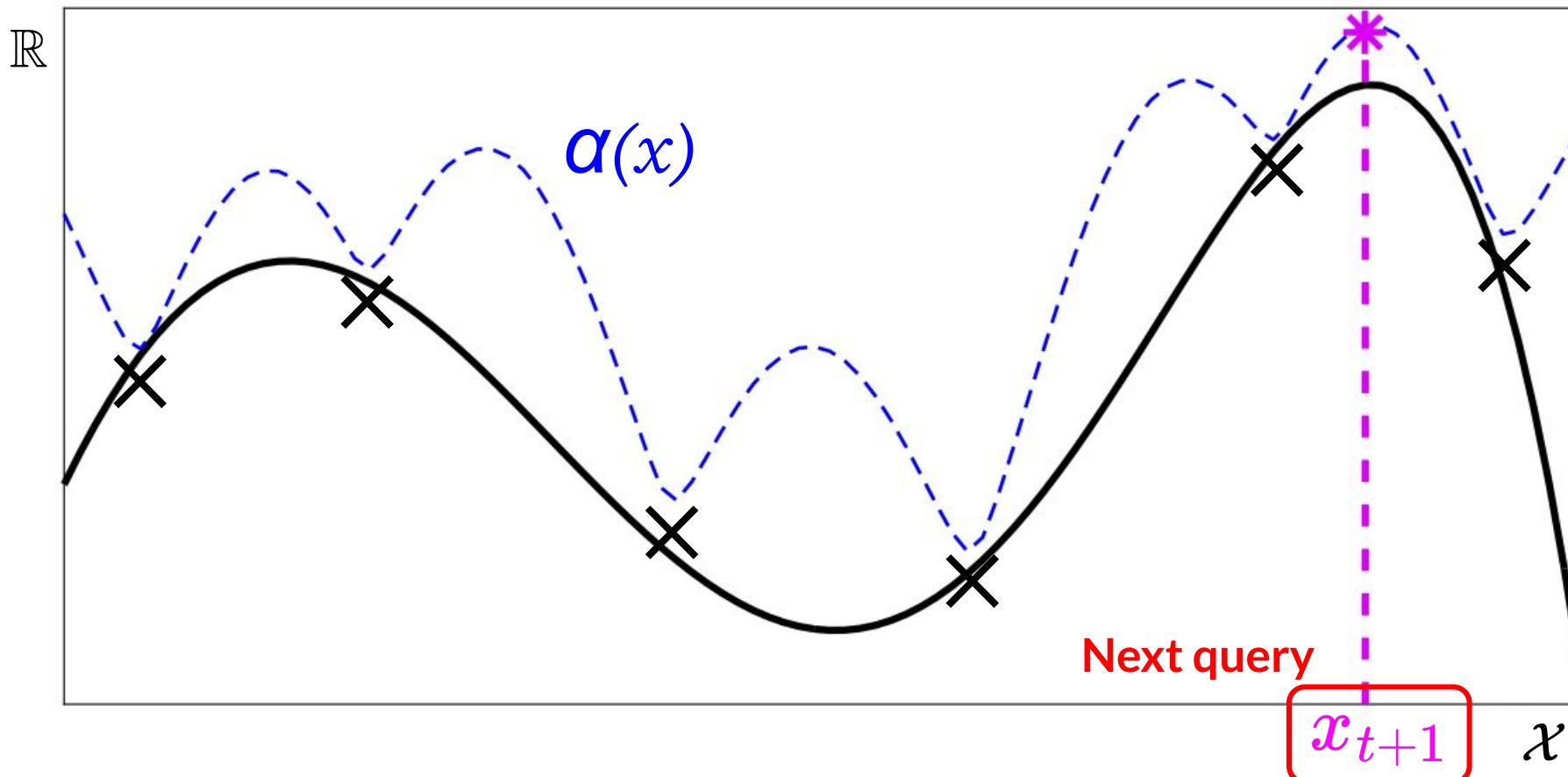
Bayesian Optimization – Visualization



Optimize acquisition function \Rightarrow yields next point to query.

$$x_{t+1} = \arg \max_{x \in \mathcal{X}} \alpha(x)$$

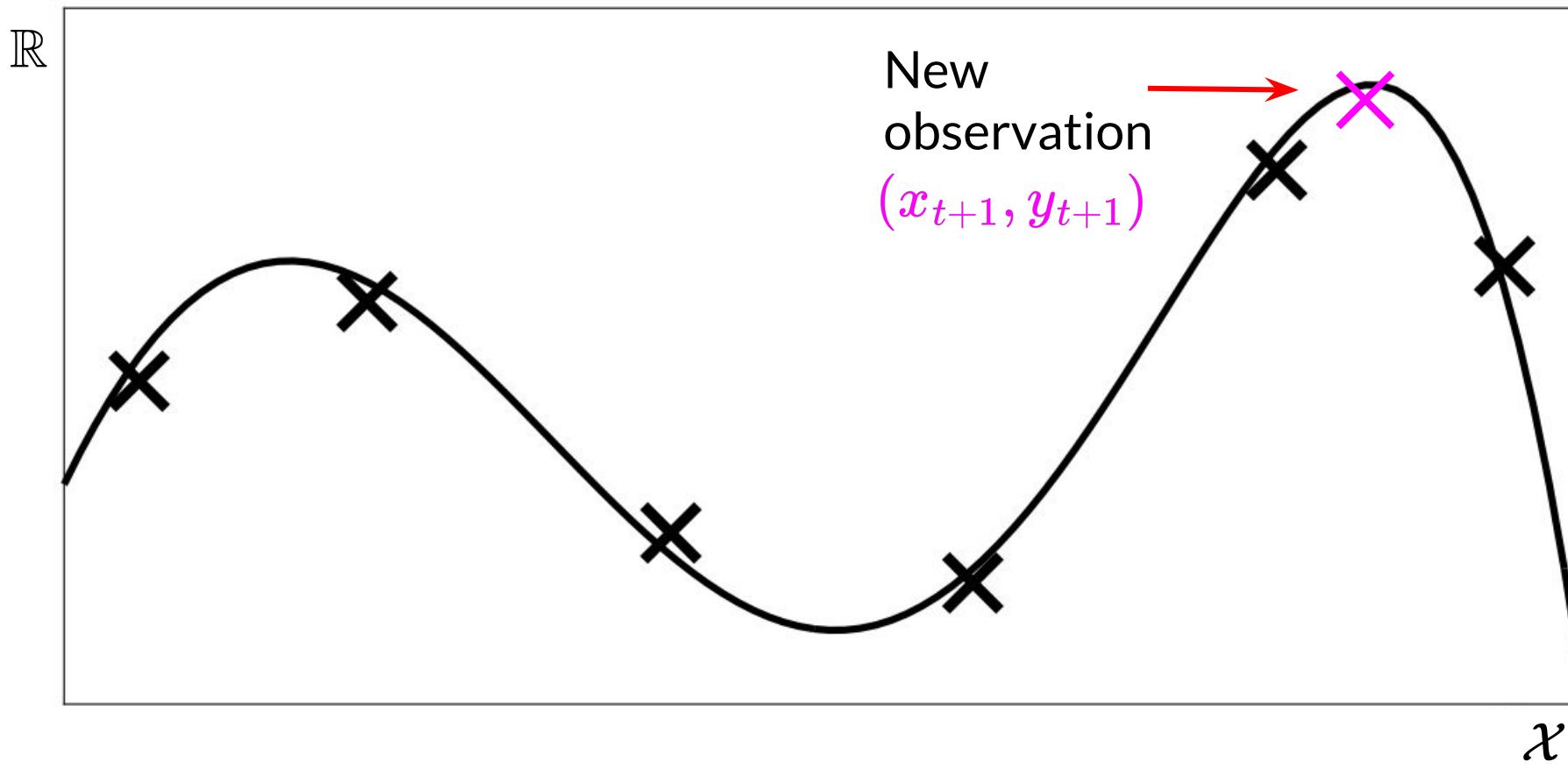
Bayesian Optimization – Visualization



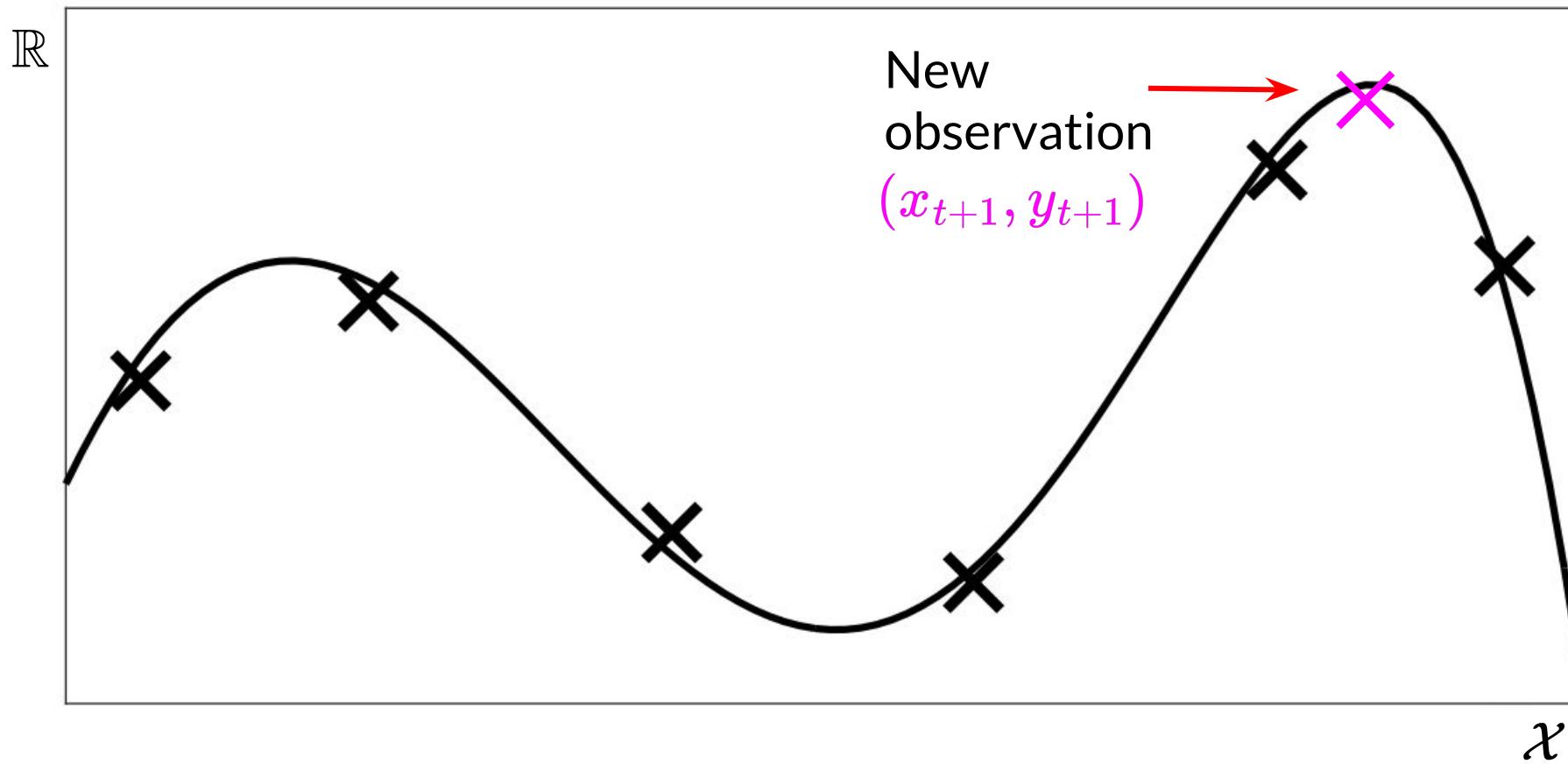
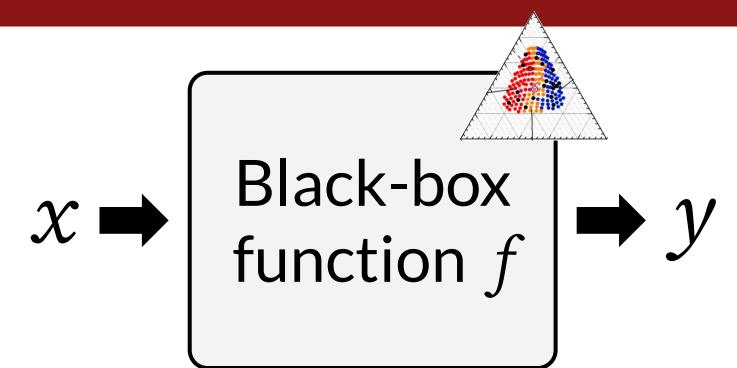
Optimize acquisition function \Rightarrow yields next point to query.

$$x_{t+1} = \arg \max_{x \in \mathcal{X}} \alpha(x)$$

Bayesian Optimization – Visualization



Bayesian Optimization – Visualization



Query black-box f at x , observe y , and update dataset.

Acquisition Function – Key Component of BO

Acquisition Function – Key Component of BO

Acquisition functions are, effectively, (one of) the main things that change between different BO procedures.

⇒ they dictate the BO algorithm, and choose the point queried at each iteration.

Acquisition Function – Key Component of BO

Acquisition functions are, effectively, (one of) the main things that change between different BO procedures.

⇒ they dictate the BO algorithm, and choose the point queried at each iteration.

Goal when designing acquisition functions:

- Want to be “optimal” in some sense; to choose the point that is *most informative*.
- And/or: most-likely to accomplish goals of algorithm, e.g., carry out optimization.

Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

UCB (*upper confidence bound*) – which we saw visualized before.

Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

UCB (*upper confidence bound*) – which we saw visualized before.

PI (*probability of improvement*) – one of the simplest acquisition functions.

Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

UCB (*upper confidence bound*) – which we saw visualized before.

PI (*probability of improvement*) – one of the simplest acquisition functions.

EI (*expected improvement*) – perhaps the most popular of the acquisition functions.

Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

UCB (*upper confidence bound*) – which we saw visualized before.

PI (*probability of improvement*) – one of the simplest acquisition functions.

EI (*expected improvement*) – perhaps the most popular of the acquisition functions.

KG (*knowledge gradient*) – bit more complex, but performs well and is well-justified.

Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

UCB (*upper confidence bound*) – which we saw visualized before.

PI (*probability of improvement*) – one of the simplest acquisition functions.

EI (*expected improvement*) – perhaps the most popular of the acquisition functions.

KG (*knowledge gradient*) – bit more complex, but performs well and is well-justified.

ES (*entropy search*) – rooted in information theory and *Bayesian optimal experimental design (BOED)*.

Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

UCB (*upper confidence bound*) – which we saw visualized before.

PI (*probability of improvement*) – one of the simplest acquisition functions.

EI (*expected improvement*) – perhaps the most popular of the acquisition functions.

KG (*knowledge gradient*) – bit more complex, but performs well and is well-justified.

ES (*entropy search*) – rooted in information theory and *Bayesian optimal experimental design (BOED)*.

But first, a quick discussion on terminology!

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

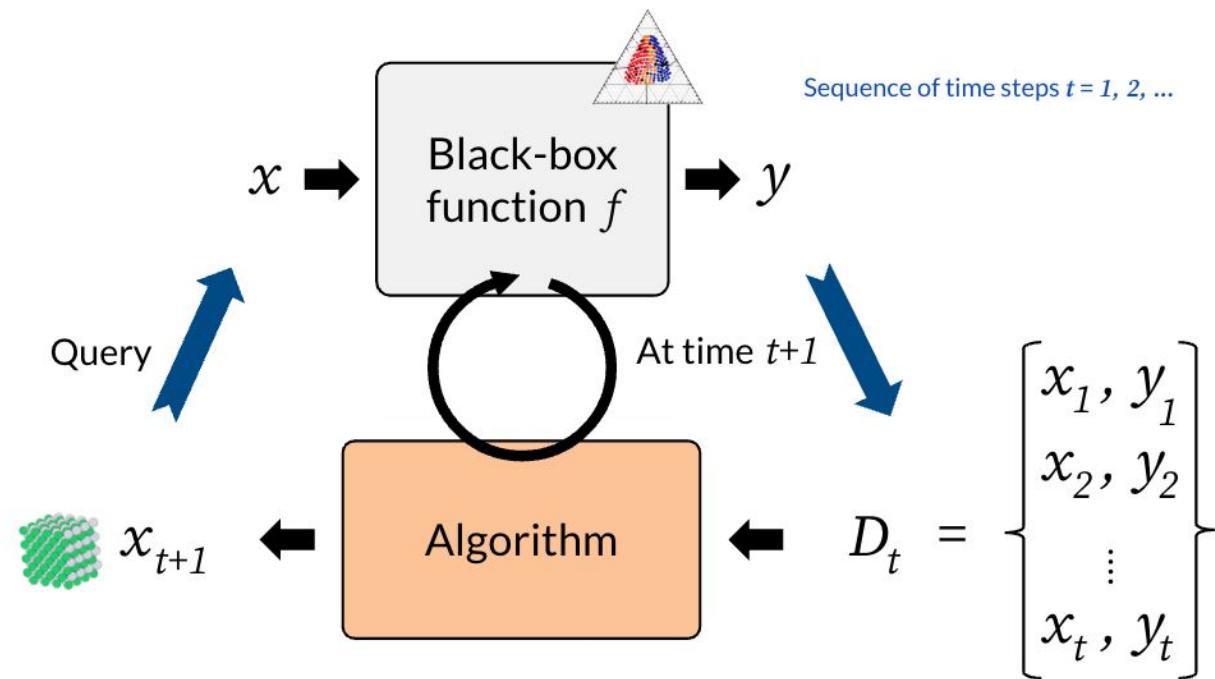
These terms all refer to related procedures, but there are slight differences between them...

Similarities: these all refer to sequential decision making (data acquisition) procedures.

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

Similarities: these all refer to sequential decision making (data acquisition) procedures.



Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

Differences:

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

Differences:

Active Learning

Originally from CS/ML.

Classically: iteratively choosing points to label from a set.

More recently: any learning method where you choose a data point to observe at each iteration .

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

Differences:

Active Learning

Originally from CS/ML.

Classically: iteratively choosing points to label from a set.

More recently: any learning method where you choose a data point to observe at each iteration .

Bayesian Optimization

Originally from operations research (among other fields).

Can view as: a special case of active learning, aiming to carry out the task of optimization.

(Which could potentially use an experimental design strategy, e.g., in *entropy search*).

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

Differences:

Active Learning

Originally from CS/ML.

Classically: iteratively choosing points to label from a set.

More recently: any learning method where you choose a data point to observe at each iteration .

Bayesian Optimization

Originally from operations research (among other fields).

Can view as: a special case of active learning, aiming to carry out the task of optimization.

(Which could potentially use an experimental design strategy, e.g., in entropy search).

Experimental Design

Originally from statistics.

Classically: goal is to efficiently estimate a model parameter (often by reducing posterior entropy).

Can view parameter estimation as one strategy for active learning/Bayesian optimization.

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

A bit more on experimental design:

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

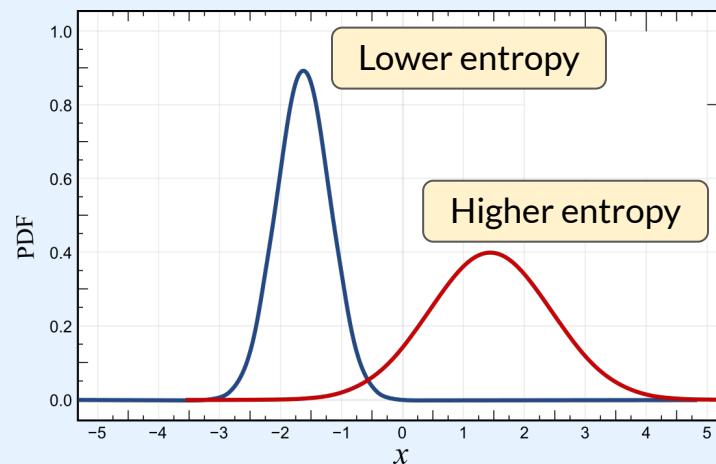
A bit more on experimental design: often referred to as *information-based experimental design* (if it involves reducing posterior entropy), or *Bayesian optimal experimental design (BOED)*.

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

A bit more on experimental design: often referred to as *information-based experimental design* (if it involves reducing posterior entropy), or *Bayesian optimal experimental design (BOED)*.

BOED: have model with an (unknown) parameter of interest.

- Choose experiments that most reduce uncertainty about parameter.
- **Uncertainty:** **entropy** of posterior distribution over parameter.



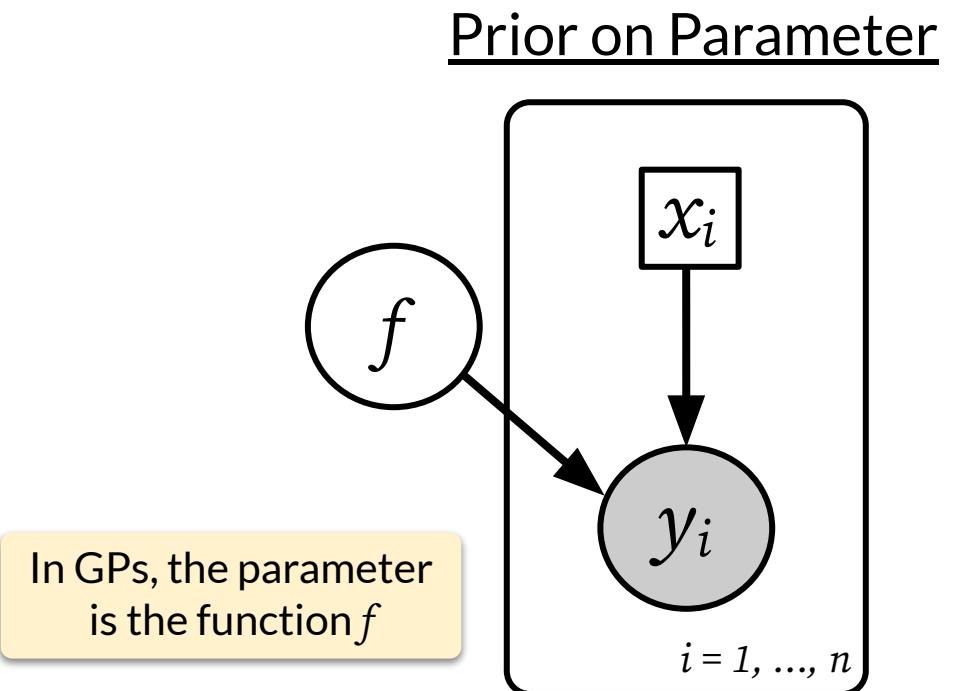
Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

How would experimental design look for GPs?

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

How would experimental design look for GPs?

Recall the PGM for GPs:

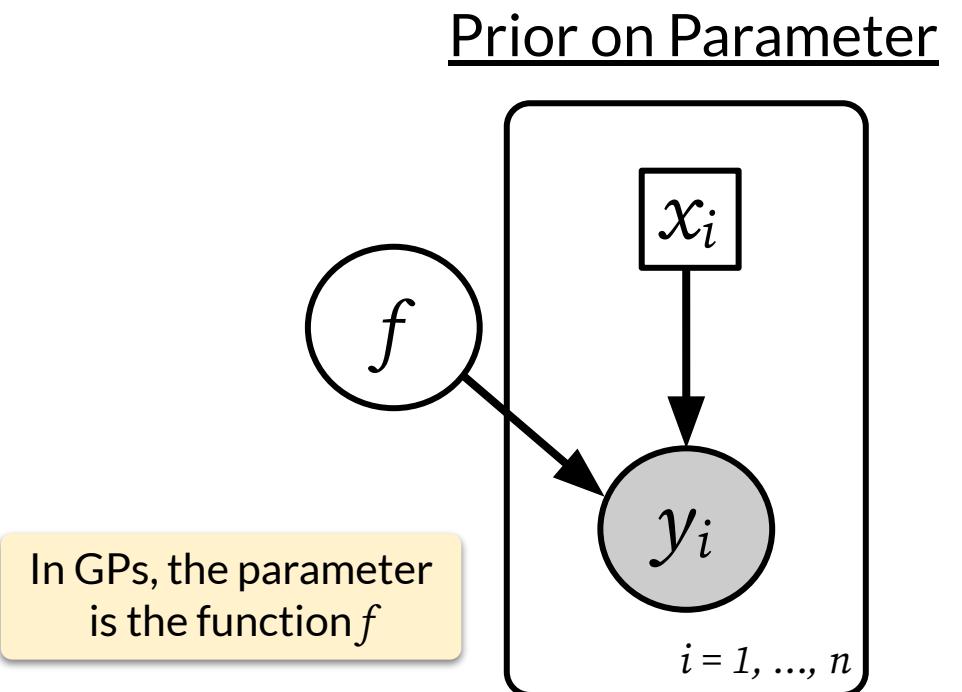


Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

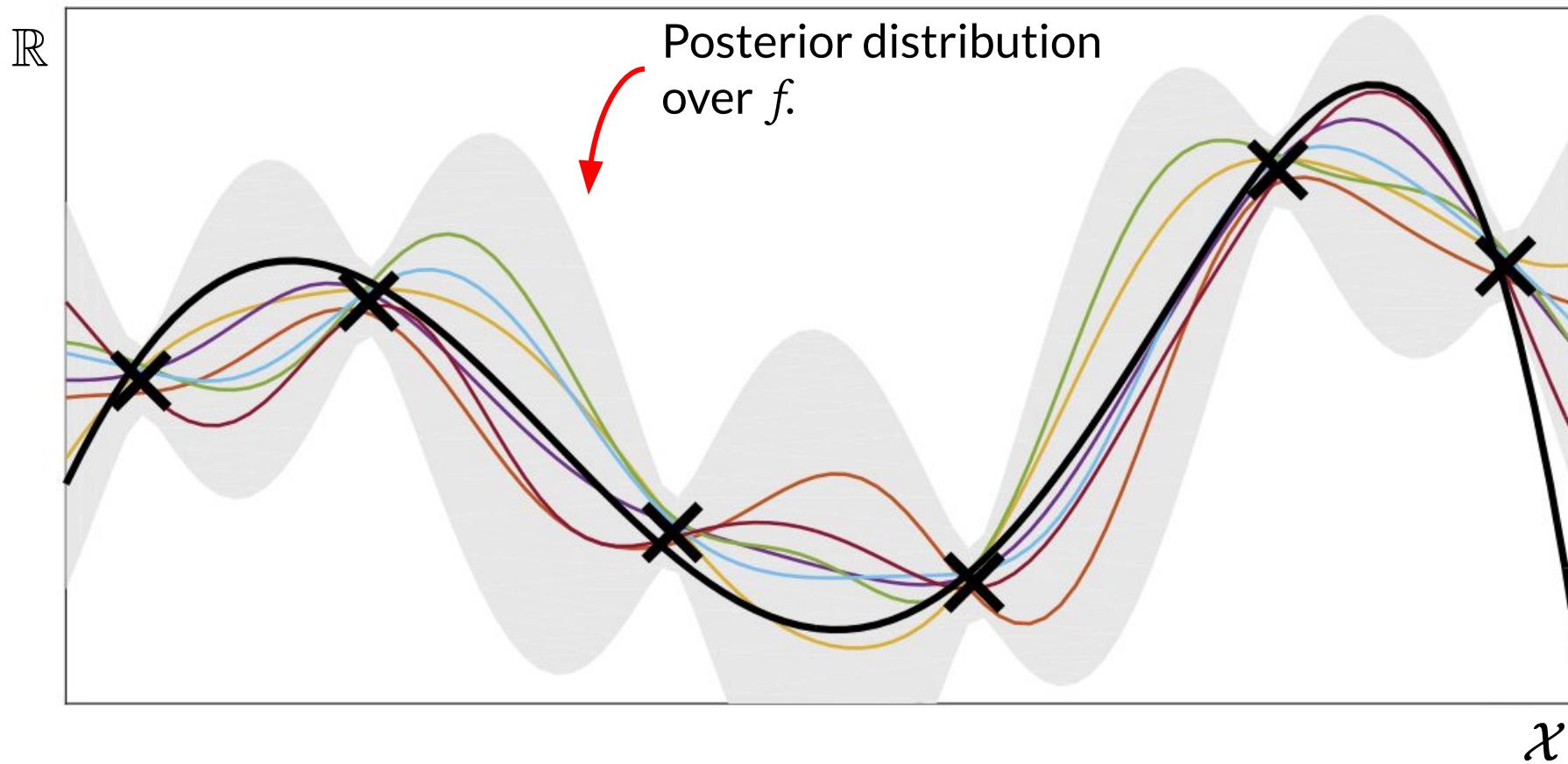
How would experimental design look for GPs?

Recall the PGM for GPs:

- We have a prior on the model parameter.
- Where the model parameter can be viewed as the function f .
- And via Bayesian inference we compute a posterior distribution over f .

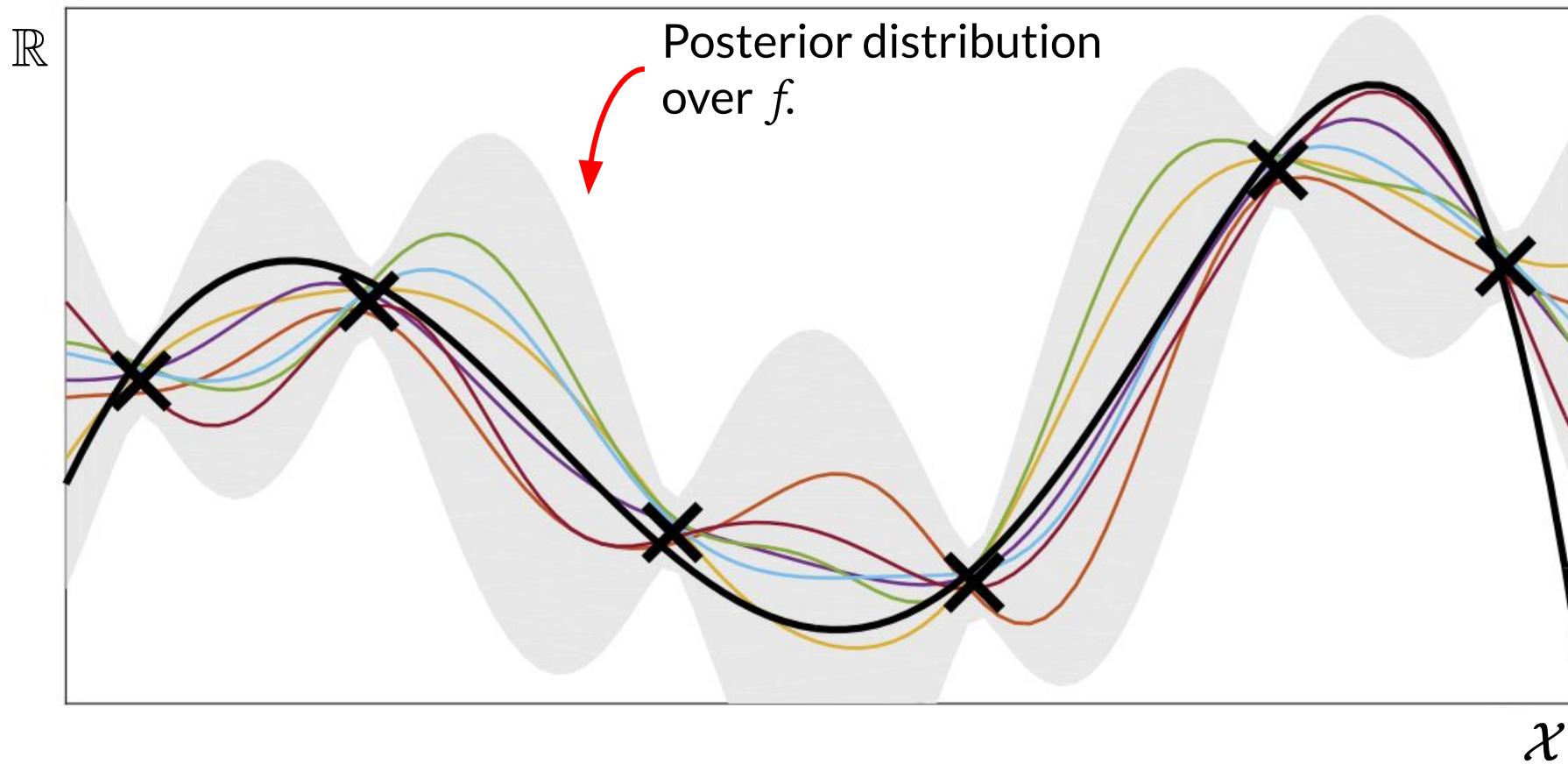


Terminology – Active Learning, Bayes Opt, Experimental Design, etc.



Every time we get an observation, we reduce the posterior uncertainty over f .

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.



Experimental design \Rightarrow choose a point to maximally shrink the uncertainty (entropy).

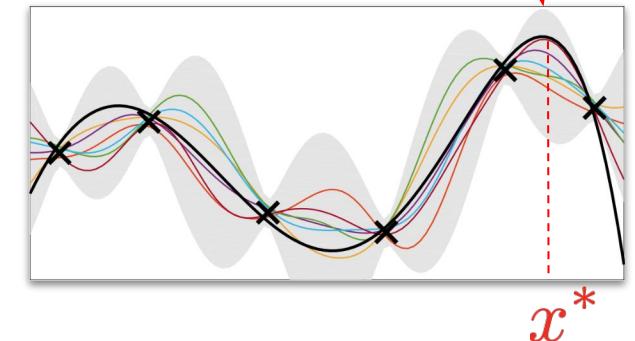
Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

As mentioned, this can be used for BO – but we don't care about reducing uncertainty over full function, just over the region near the optima.

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

As mentioned, this can be used for BO – but we don't care about reducing uncertainty over full function, just over the region near the optima.

Only care about reducing uncertainty around here



Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

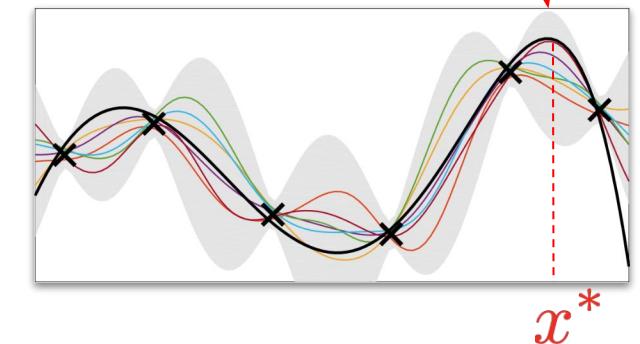
As mentioned, this can be used for BO – but we don't care about reducing uncertainty over full function, just over the region near the optima.

This leads to *information-based Bayesian optimization*:

Info-based BO: entropy search (ES) , predictive ES , max-value ES .

- All rooted in **Bayesian optimal experimental design (BOED)**.

Only care about reducing uncertainty around here



Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

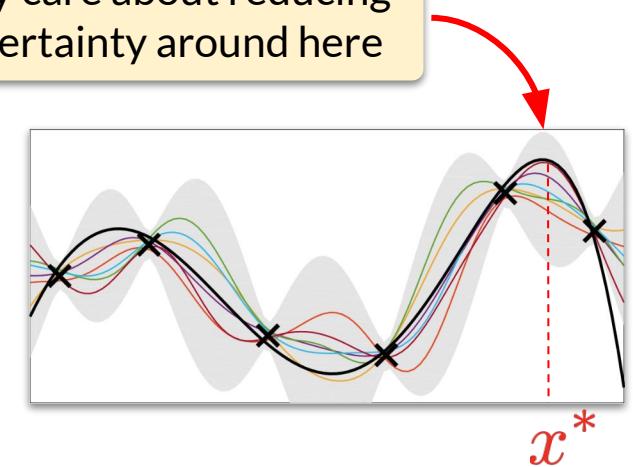
As mentioned, this can be used for BO – but we don't care about reducing uncertainty over full function, just over the region near the optima.

This leads to *information-based Bayesian optimization*:

Info-based BO: entropy search (ES) , predictive ES , max-value ES .

- All rooted in **Bayesian optimal experimental design (BOED)**.

Only care about reducing uncertainty around here

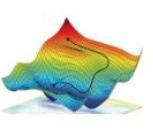
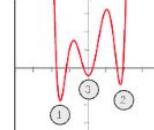
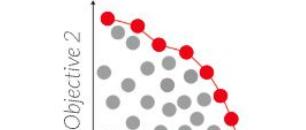
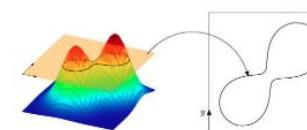


Which we will talk about a little later...

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

Side note: there are many other potential special cases of active learning beyond global optimization, as in BO...

Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

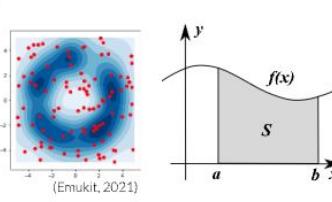
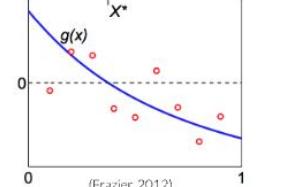
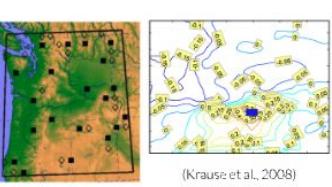
Optimization Variants (global, local, top-k optima)	Multi-objective Opt. (Pareto frontiers)	Level Set Estimation (superlevel set, sublevel set)	Search (subset matching criteria)	Phase Identification (boundaries, partitions)
 Local optima	 Top-k optima		 (Mororó and van der Meer, 2020)	 (molecules)
Applications	Applications	Applications	Applications	Applications
- ML model training - HPO / NAS - Systems tuning	- Process optimization - Portfolio optimization - Laboratory equipment / machines	- Catalyst design - Active learning / weak supervision - Environmental monitoring	- Drug discovery - Fraud detection - Targeted opinion polling	- High throughput materials design/discovery

Quadrature, Integration
(integrals, expectations)

Graph Algorithms
(shortest paths)

Root Finding, Bisection
(roots)

Sensor Placement
(function value at locations)

Quadrature, Integration (integrals, expectations)	Graph Algorithms (shortest paths)	Root Finding, Bisection (roots)	Sensor Placement (function value at locations)
 (Emukiti, 2021)	 Geometric shortest path problem (GSP)	 (Frazier, 2012)	 (Krause et al., 2008)
Applications	Applications	Applications	Applications
- Probabilistic modeling (marginal distributions, normalization constants) - Estimating center of mass (and centroids)	- Transportation networks - Shipping networks - Social networks	- Target localization (airborne radar) - Edge detection (computer vision) - Biology / microbiology (phase shifts)	- Water distribution systems - Outbreak detection in networks - Weather monitoring

Acquisition Functions – Back to Acquisition Functions

Acquisition Functions – Back to Acquisition Functions

Acquisition Function

Component of Bayesian optimization algorithm.

$$\alpha : \mathcal{X} \rightarrow \mathbb{R}$$

Assigns a number describing the value/usefulness of querying a given point.

Acquisition Functions – Upper Confidence Bound (UCB)

Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

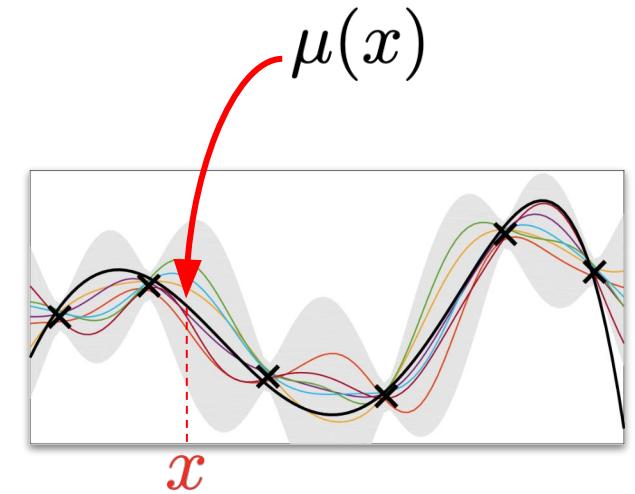
Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

where:

$\mu(x)$ - Mean of *marginal posterior* at input x .



Acquisition Functions – Upper Confidence Bound (UCB)

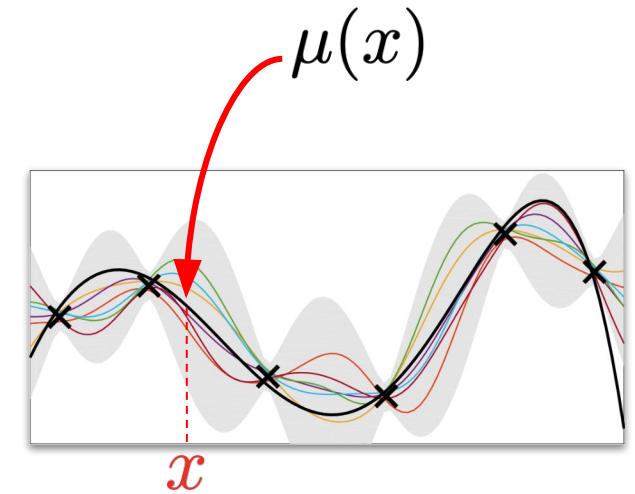
The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

where:

$$p(f(x) | \mathcal{D}_t)$$

$\mu(x)$ - Mean of *marginal posterior* at input x .



Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

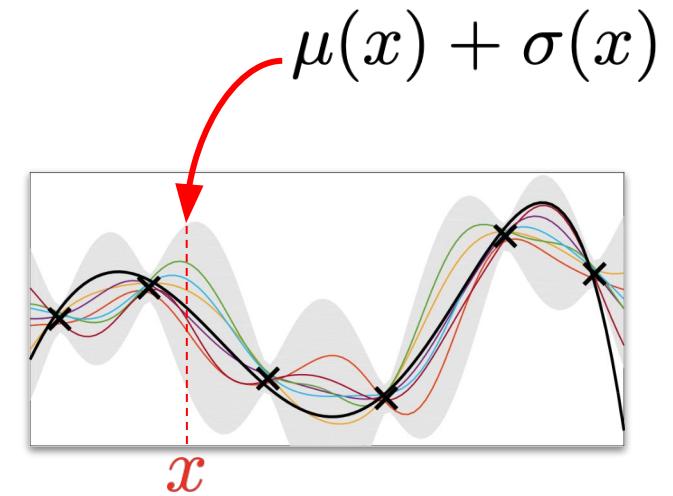
$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

where:

$$p(f(x) | \mathcal{D}_t)$$

$\mu(x)$ - Mean of *marginal posterior* at input x .

$\sigma(x)$ - Standard deviation of *marginal posterior* at input x .



Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

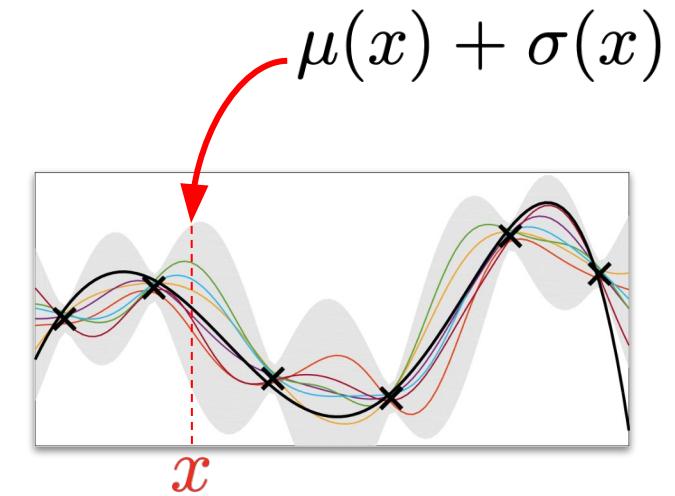
where:

$$p(f(x) | \mathcal{D}_t)$$

$\mu(x)$ - Mean of *marginal posterior* at input x .

$\sigma(x)$ - Standard deviation of *marginal posterior* at input x .

β_t - Hyperparameter controlling tradeoff between exploration and exploitation.

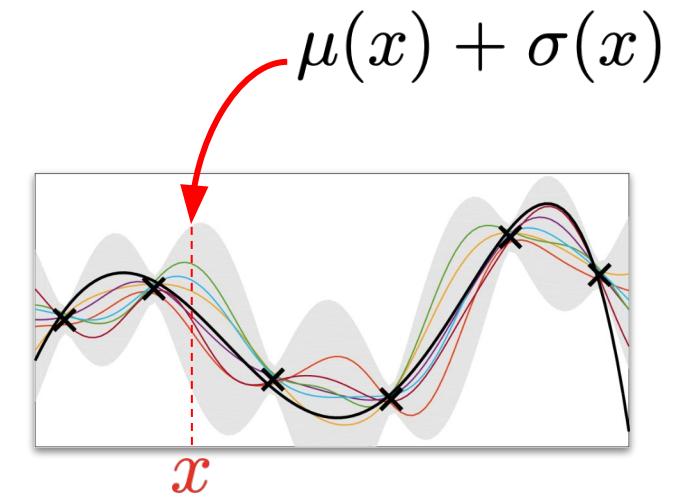


Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:



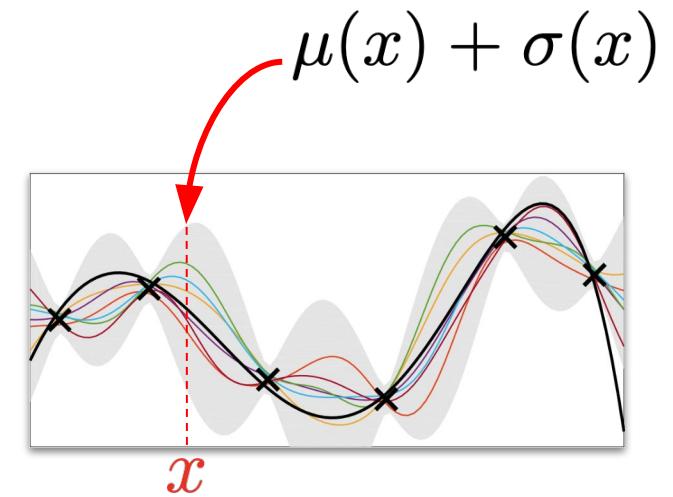
Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:

Maximizing std-dev \rightarrow queries most *unknown* part of function
 \Rightarrow most “informative”.



Acquisition Functions – Upper Confidence Bound (UCB)

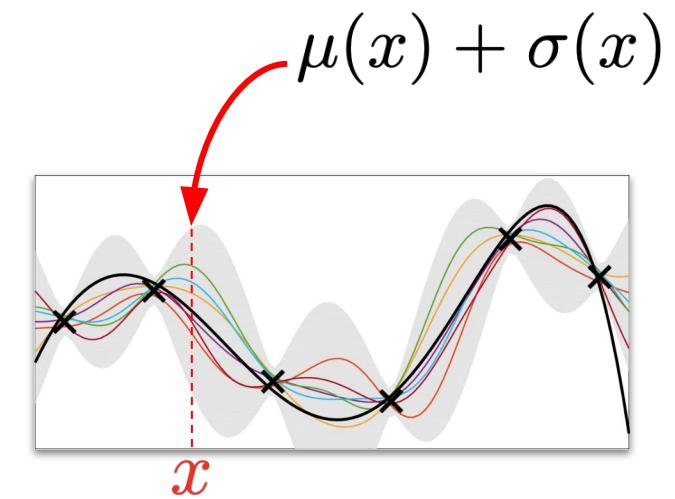
The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:

Maximizing std-dev \rightarrow queries most *unknown* part of function
 \Rightarrow most “informative”.

Maximizing mean \rightarrow queries *best guess of highest value*, based on current knowledge of function.



Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

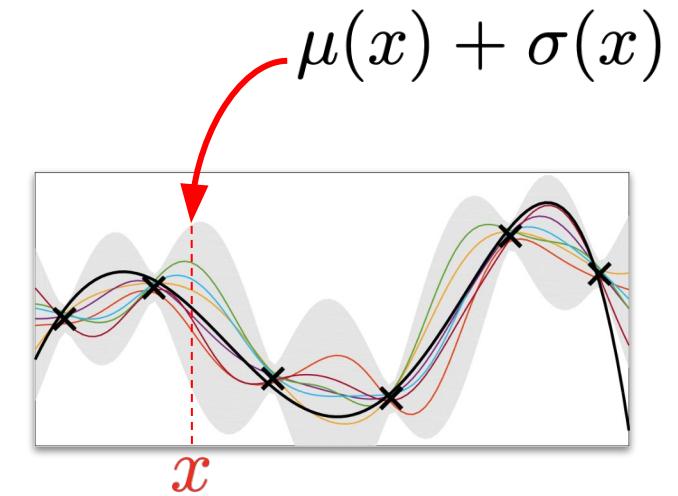
$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:

Maximizing std-dev \rightarrow queries most *unknown* part of function
 \Rightarrow most “informative”.

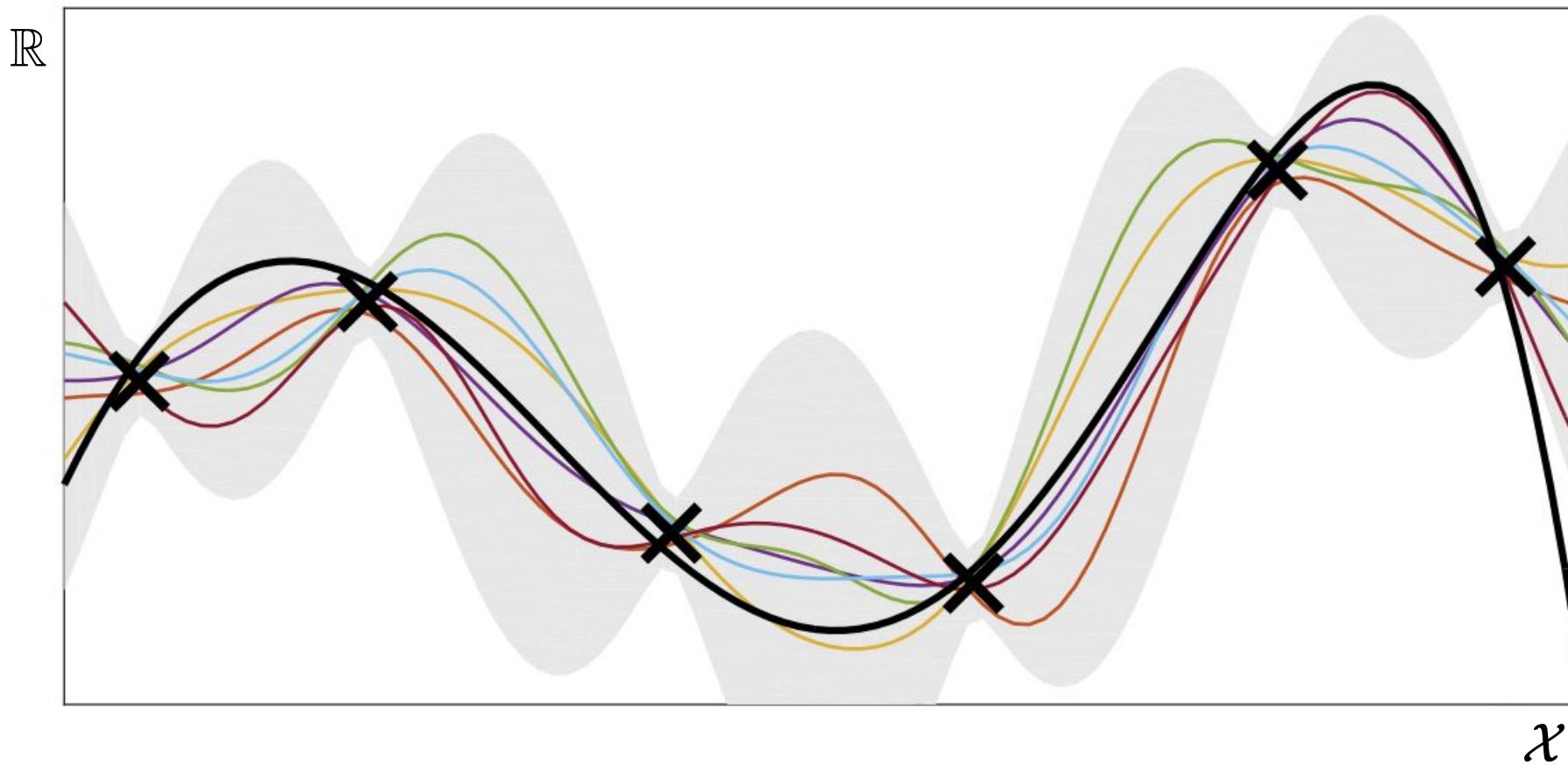
Maximizing mean \rightarrow queries *best guess of highest value*, based on current knowledge of function.

We want to trade off between these two quantities.



Acquisition Functions – Upper Confidence Bound (UCB)

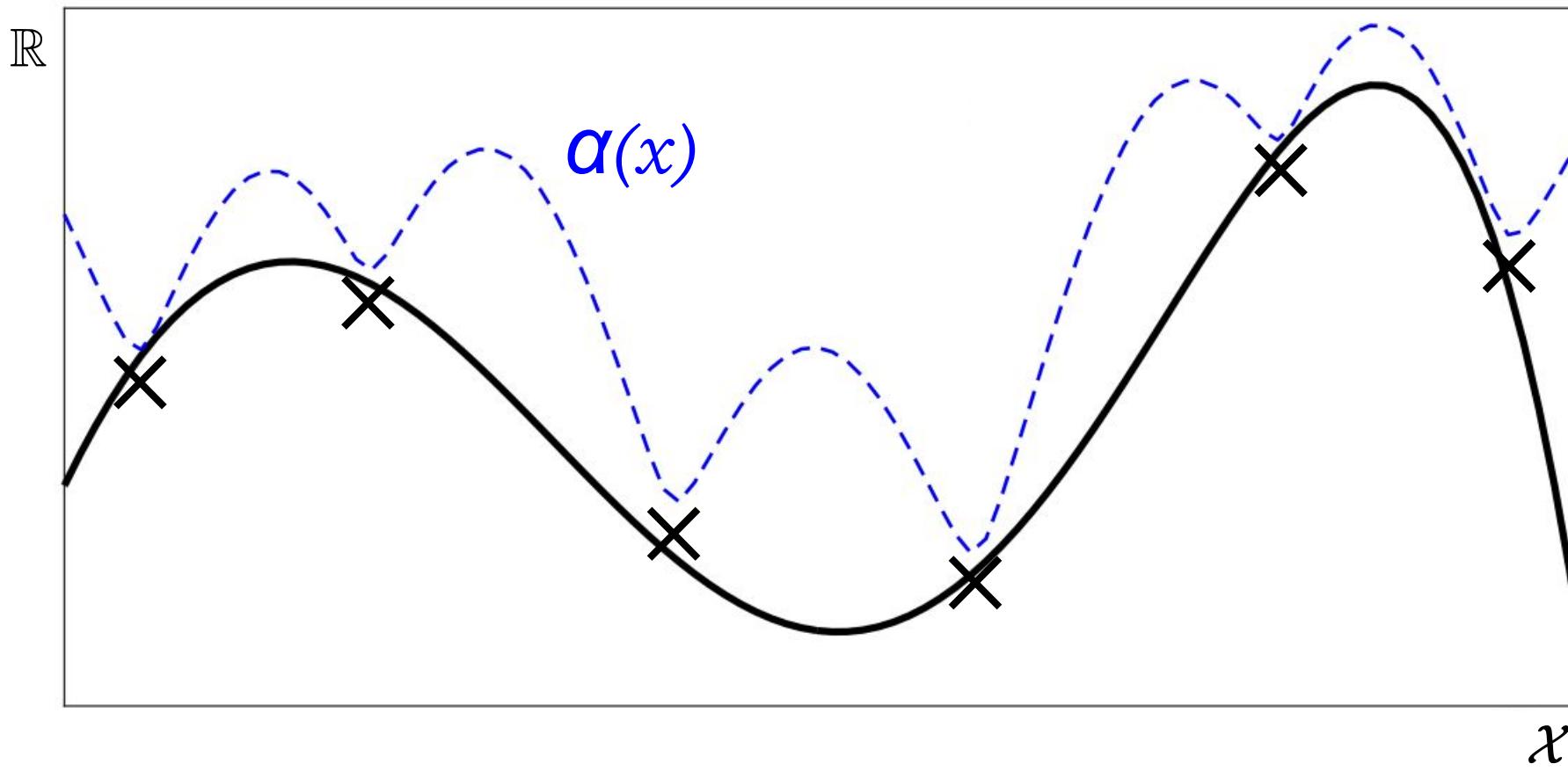
Visualizing UCB



Given our probabilistic model...

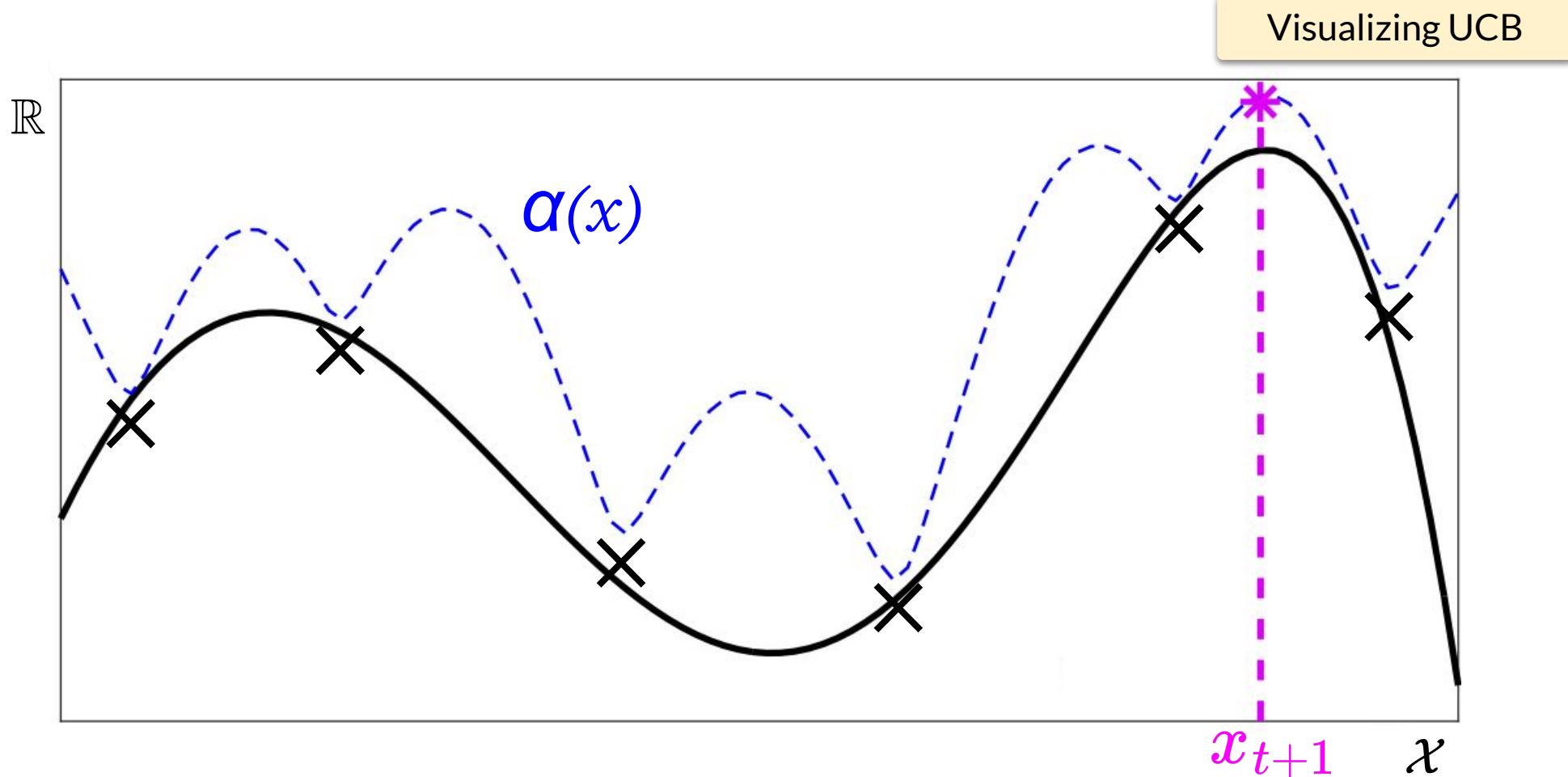
Acquisition Functions – Upper Confidence Bound (UCB)

Visualizing UCB



... define the **UCB** acquisition function.

Acquisition Functions – Upper Confidence Bound (UCB)



Acquisition Functions – Experimental Design (EIG about f)

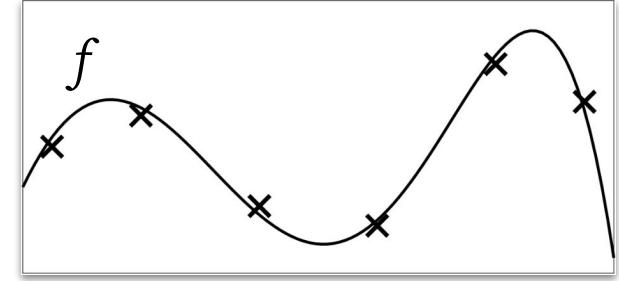
Quick aside on experimental design, to give some perspective on UCB:

Acquisition Functions – Experimental Design (EIG about f)

Quick aside on experimental design, to give some perspective on UCB:

Suppose you don't want to do optimization, but instead want to learn the full function landscape (i.e., learn "all of f ").

⇒ technically speaking, this is not *optimization* anymore.

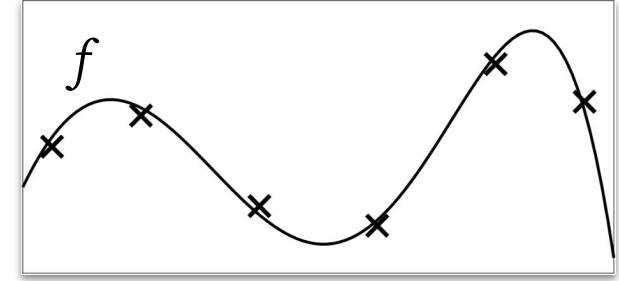


Acquisition Functions – Experimental Design (EIG about f)

Quick aside on experimental design, to give some perspective on UCB:

Suppose you don't want to do optimization, but instead want to learn the full function landscape (i.e., learn "all of f ").

⇒ technically speaking, this is not *optimization* anymore.



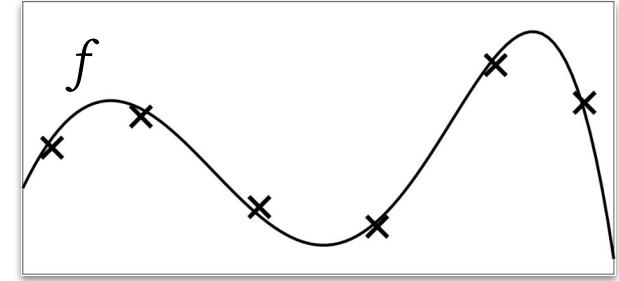
An *experimental design approach* would be to query the point that **maximally reduces the posterior entropy about f** .

Acquisition Functions – Experimental Design (EIG about f)

Quick aside on experimental design, to give some perspective on UCB:

Suppose you don't want to do optimization, but instead want to learn the full function landscape (i.e., learn "all of f ").

⇒ technically speaking, this is not *optimization* anymore.



An *experimental design approach* would be to query the point that **maximally reduces the posterior entropy about f** .

(This is equivalently also commonly called: *expected information gain*, or *EIG, about f* .)

Acquisition Functions – Experimental Design (EIG about f)

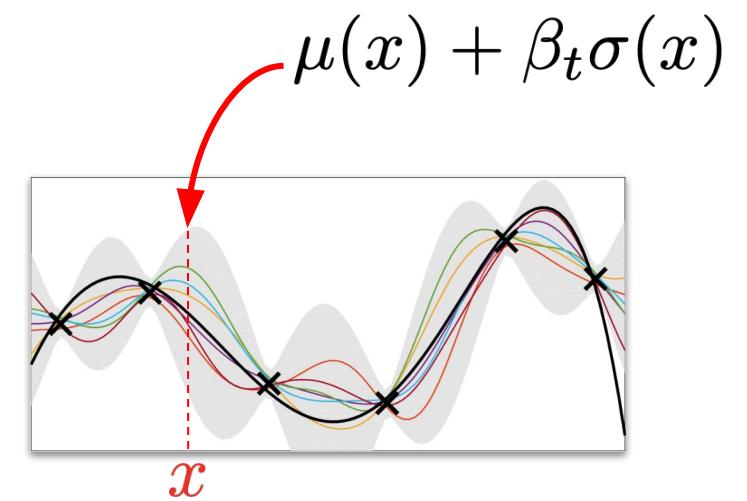
You can prove (at least for GP models) that maximally reducing posterior entropy is equivalent to querying point with maximum standard deviation (of posterior marginal).

Acquisition Functions – Experimental Design (EIG about f)

You can prove (at least for GP models) that maximally reducing posterior entropy is equivalent to querying point with maximum standard deviation (of posterior marginal).

So in UCB \Rightarrow

The second term (standard deviation term) is in fact balancing exploration (EIG about f) with exploitation (best expected input x).



Acquisition Functions – Probability of Improvement (PI)

Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\alpha_t(x) = p(f(x) > f(x^*) \mid \mathcal{D}_t)$$

Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\alpha_t(x) = p(f(x) > f(x^*) \mid \mathcal{D}_t)$$

where $f(x^*)$ is the function value of the best point observed so far in \mathcal{D}_t .

I.e., the posterior probability that the $f(x)$ is greater than the best point observed so far.

Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\alpha_t(x) = p(f(x) > f(x^*) \mid \mathcal{D}_t)$$

where $f(x^*)$ is the function value of the best point observed so far in \mathcal{D}_t .

I.e., the posterior probability that the $f(x)$ is greater than the best point observed so far.

You can prove this $\alpha_t(x)$ is equal to ...

Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\begin{aligned}\alpha_t(x) &= p(f(x) > f(x^*) \mid \mathcal{D}_t) \\ &= \Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

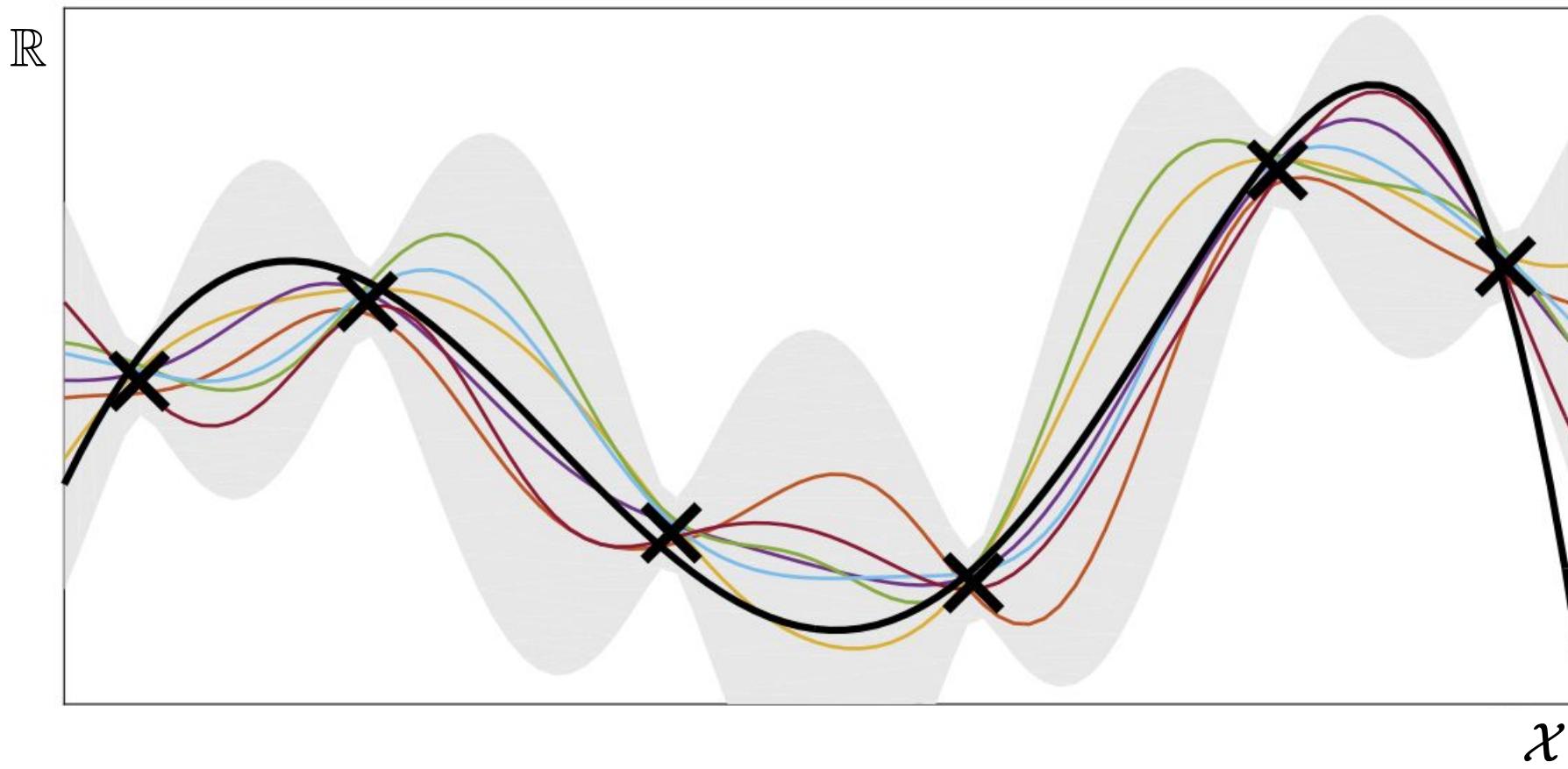
$$\begin{aligned}\alpha_t(x) &= p(f(x) > f(x^*) \mid \mathcal{D}_t) \\ &= \Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

where $\Phi(\cdot)$ denotes the CDF of the standard normal distribution.

(And, as before, $\mu(x)$ and $\sigma(x)$ are the mean and SD of marginal posterior at x .)

Acquisition Functions – Probability of Improvement (PI)

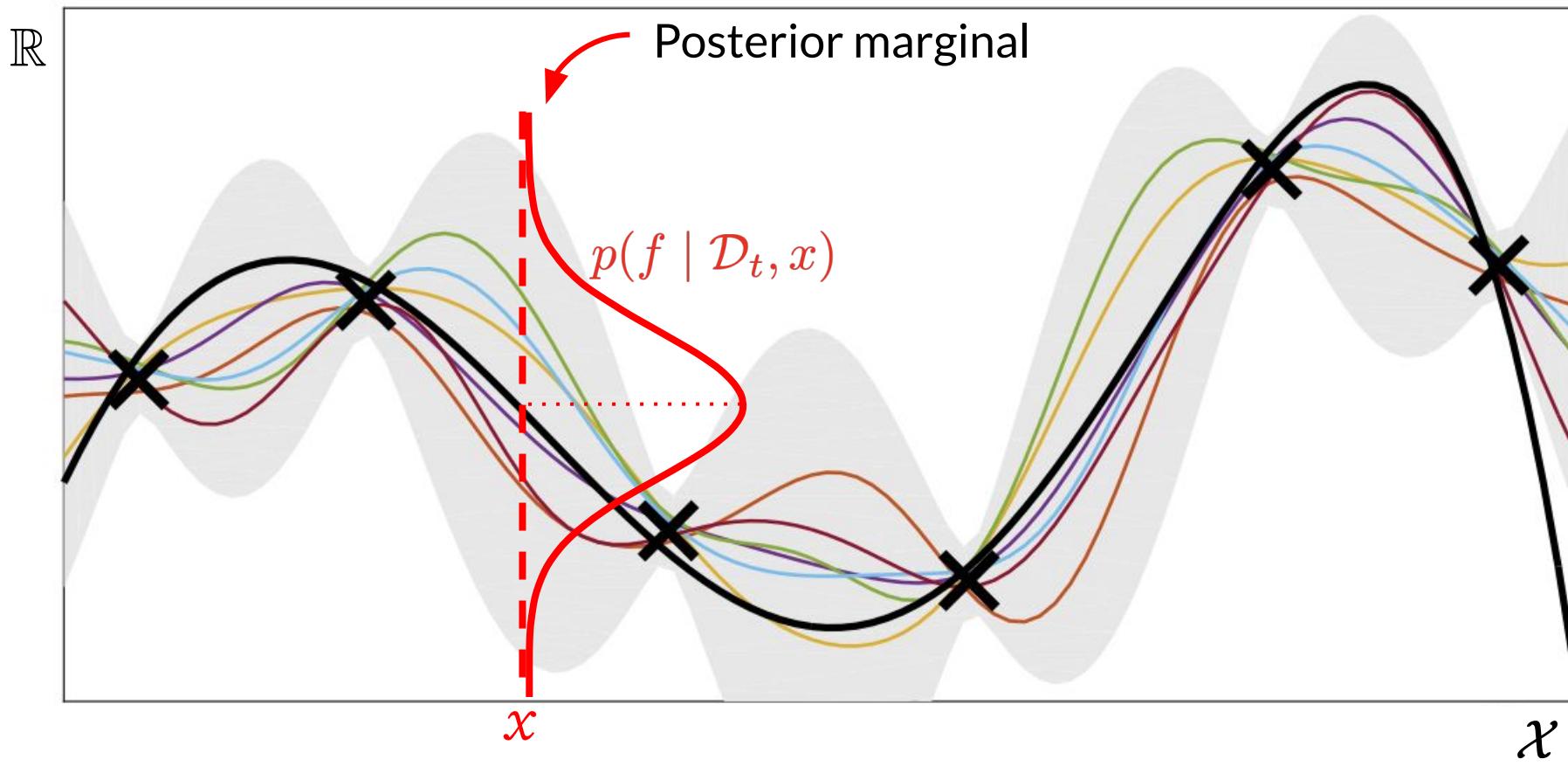
Visualizing PI



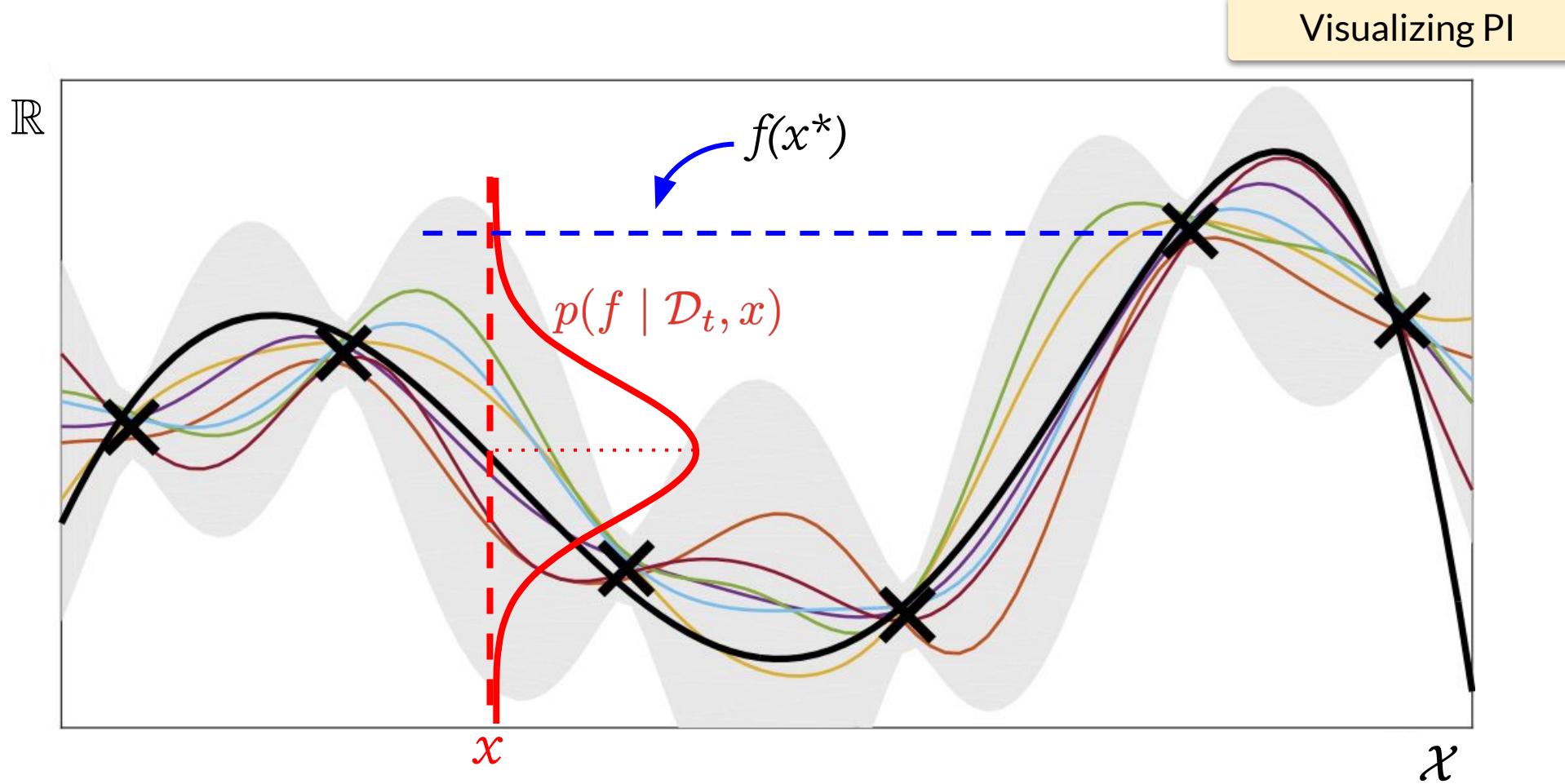
Given our probabilistic model...

Acquisition Functions – Probability of Improvement (PI)

Visualizing PI

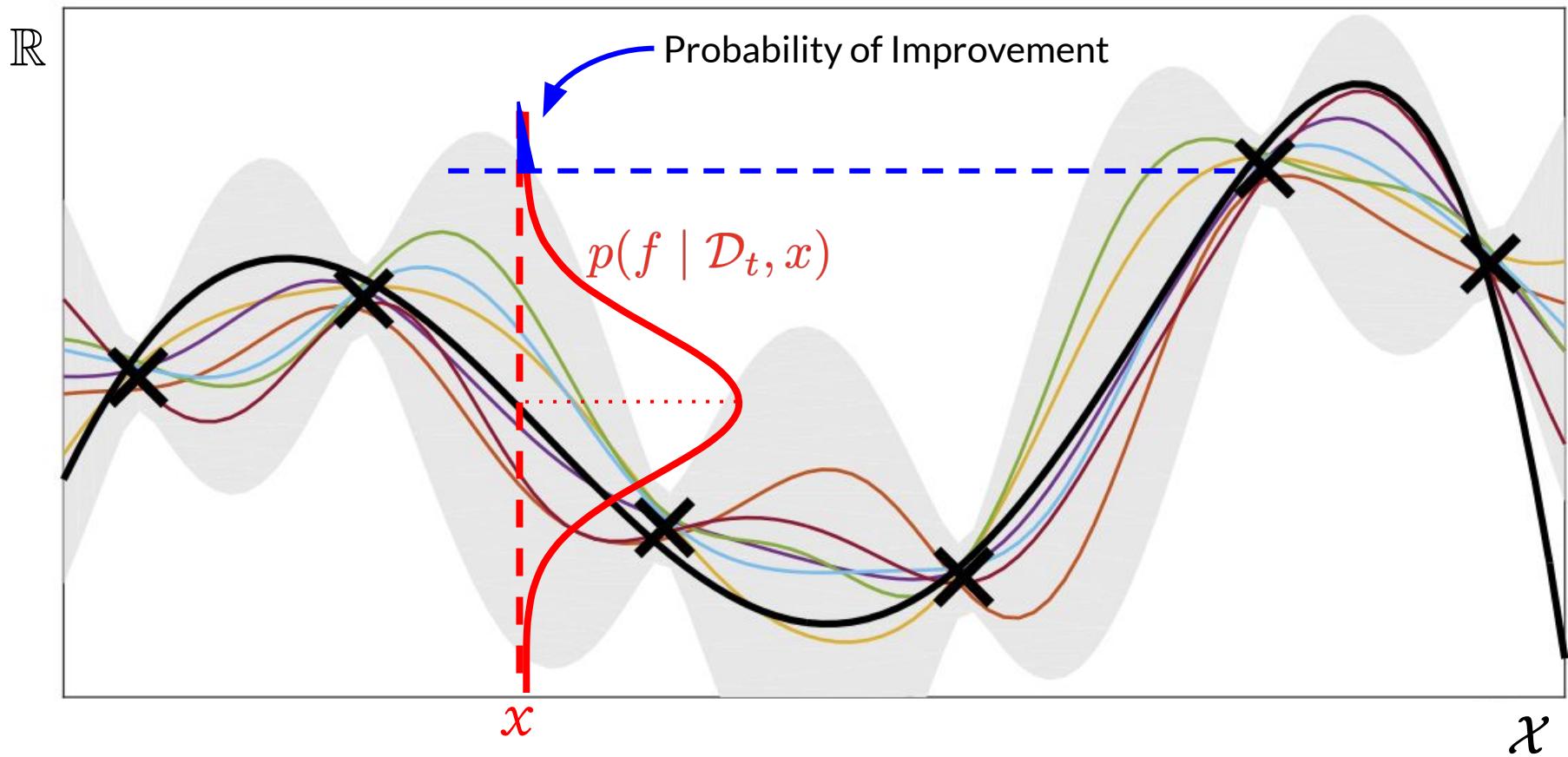


Acquisition Functions – Probability of Improvement (PI)



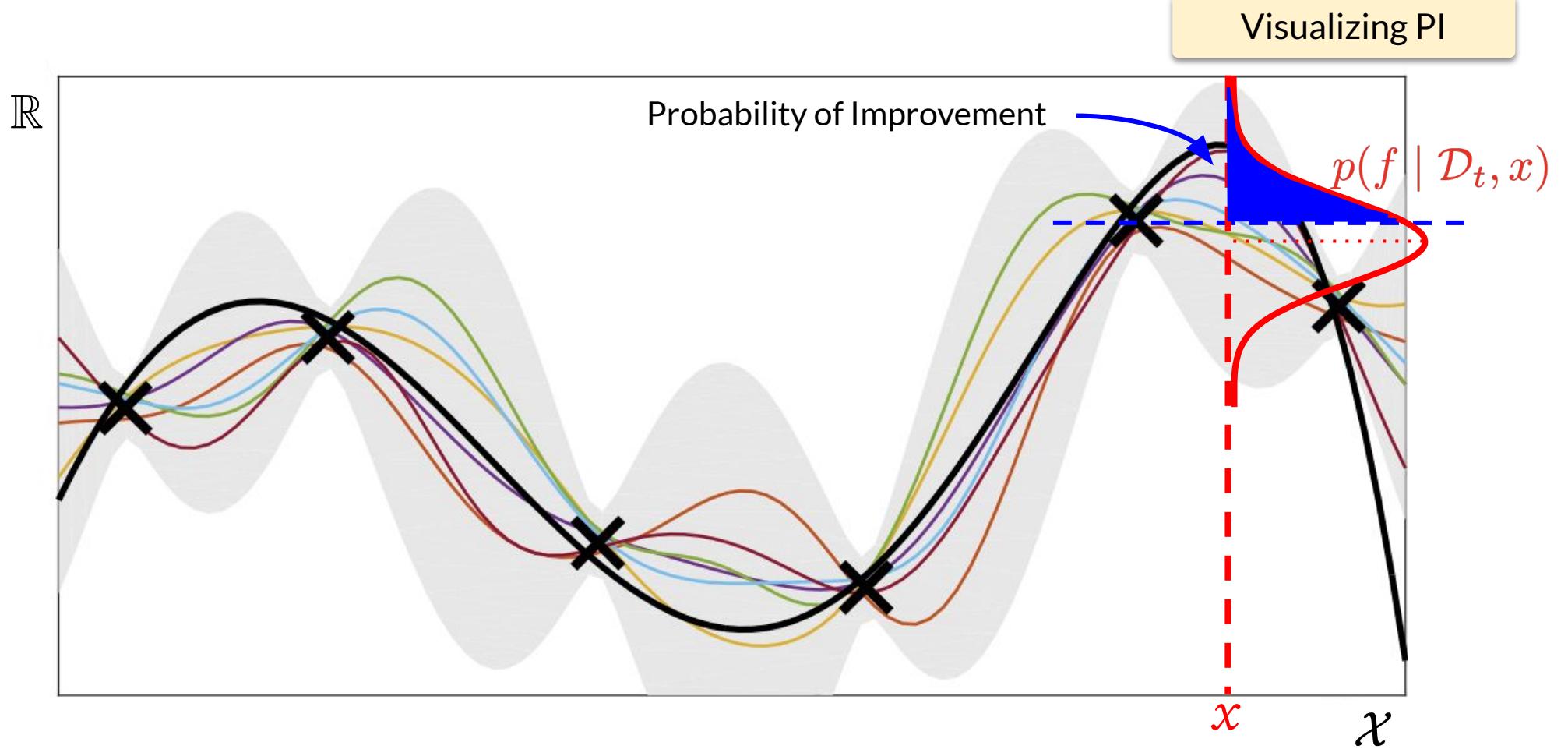
Acquisition Functions – Probability of Improvement (PI)

Visualizing PI



... define the PI acquisition function.

Acquisition Functions – Probability of Improvement (PI)



Acquisition Functions – Probability of Improvement (PI)

However, note that:

Acquisition Functions – Probability of Improvement (PI)

However, note that:

PI maximizes the **probability** that $f(x)$ is greater than then best point so far.

Does ***not*** try to maximize the **amount** that $f(x)$ is greater than the best point so far.

Acquisition Functions – Probability of Improvement (PI)

However, note that:

PI maximizes the **probability** that $f(x)$ is greater than then best point so far.

Does *not* try to maximize the **amount** that $f(x)$ is greater than the best point so far.

(i.e., it might choose an x that is only *slightly better* than best point so far, if it is more confident about it – rather than a point that is *significantly greater in expectation*, but it is less confident about it).

Acquisition Functions – Probability of Improvement (PI)

However, note that:

PI maximizes the **probability** that $f(x)$ is greater than then best point so far.

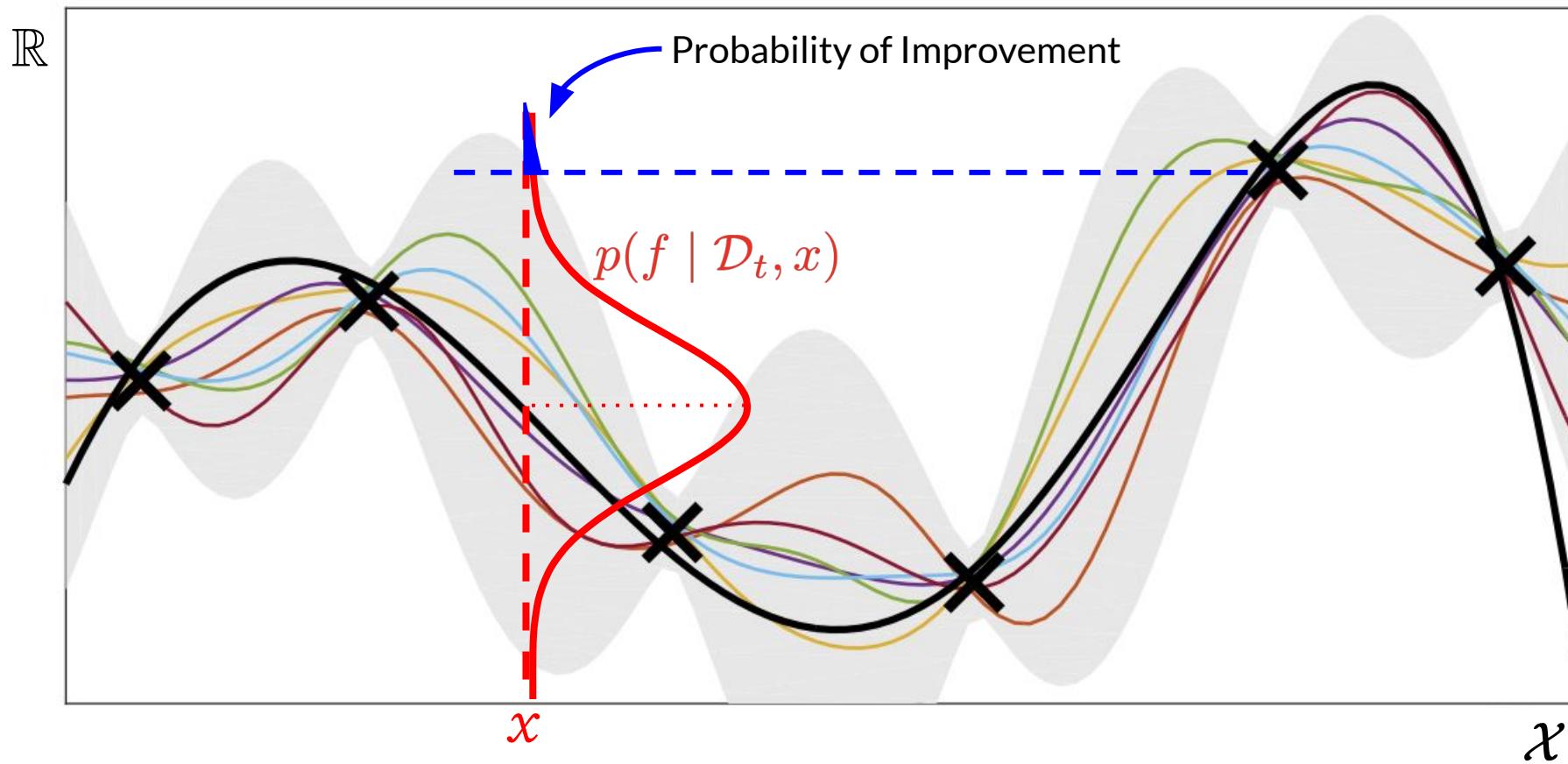
Does *not* try to maximize the **amount** that $f(x)$ is greater than the best point so far.

(i.e., it might choose an x that is only *slightly better* than best point so far, if it is more confident about it – rather than a point that is *significantly greater in expectation*, but it is less confident about it).

Visualizing this...

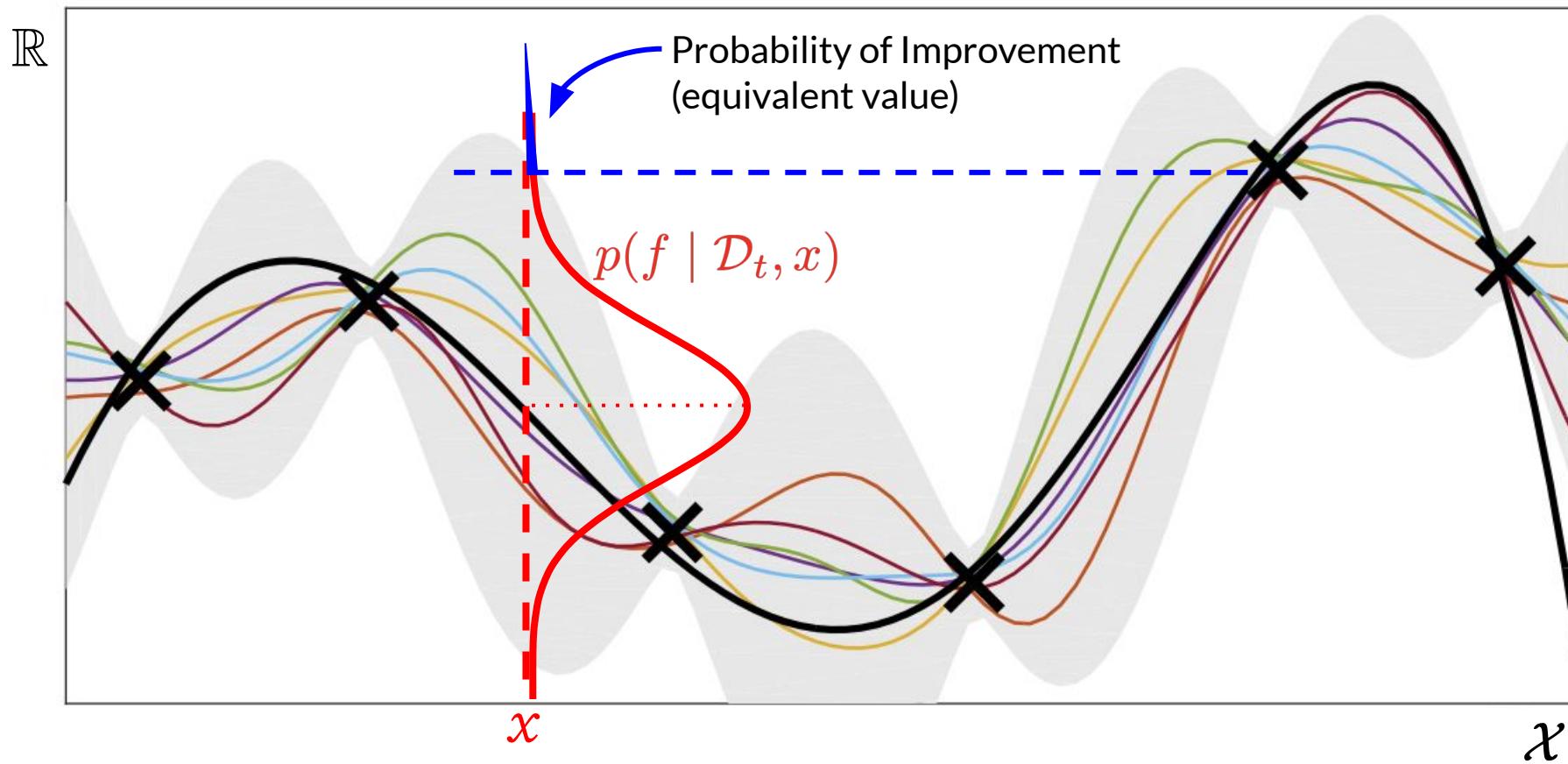
Acquisition Functions – Probability of Improvement (PI)

Visualizing PI

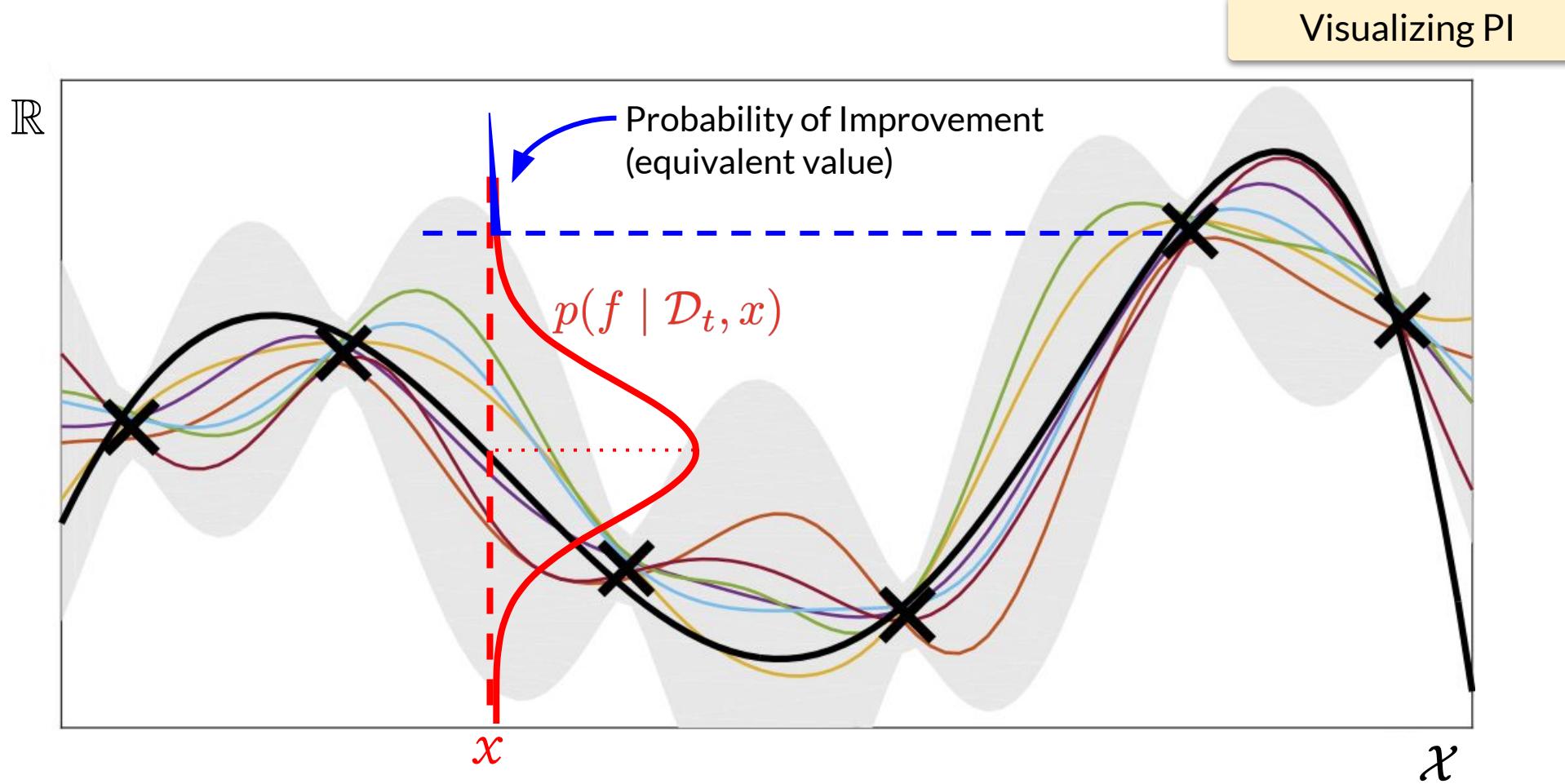


Acquisition Functions – Probability of Improvement (PI)

Visualizing PI



Acquisition Functions – Probability of Improvement (PI)



In contrast, we could try to maximize the *expected amount of improvement* over the best point so far...

Acquisition Functions – Expected Improvement (EI)

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

where $f(x^*)$ is the function value of the best point observed so far in \mathcal{D}_t .

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

where $f(x^*)$ is the function value of the best point observed so far in \mathcal{D}_t .

I.e., the expected amount that $f(x)$ is greater than the best point observed so far (while ignoring the magnitude if $f(x)$ is less than the best point so far).

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

where $f(x^*)$ is the function value of the best point observed so far in \mathcal{D}_t .

I.e., the expected amount that $f(x)$ is greater than the best point observed so far (while ignoring the magnitude if $f(x)$ is less than the best point so far).

You can prove this is equal to ...

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\begin{aligned}\alpha_t(x) &= \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}] \\ &= (\mu(x) - f(x^*))\Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right) + \sigma(x)\phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\begin{aligned}\alpha_t(x) &= \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}] \\ &= (\mu(x) - f(x^*))\Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right) + \sigma(x)\phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

where $\Phi(\cdot)$ denotes the CDF of the standard normal distribution and $\phi(\cdot)$ denotes the PDF of the standard normal distribution.

(And, as before, $\mu(x)$ and $\sigma(x)$ are the mean and SD of marginal posterior at x .)

Acquisition Functions – Expected Improvement (EI)

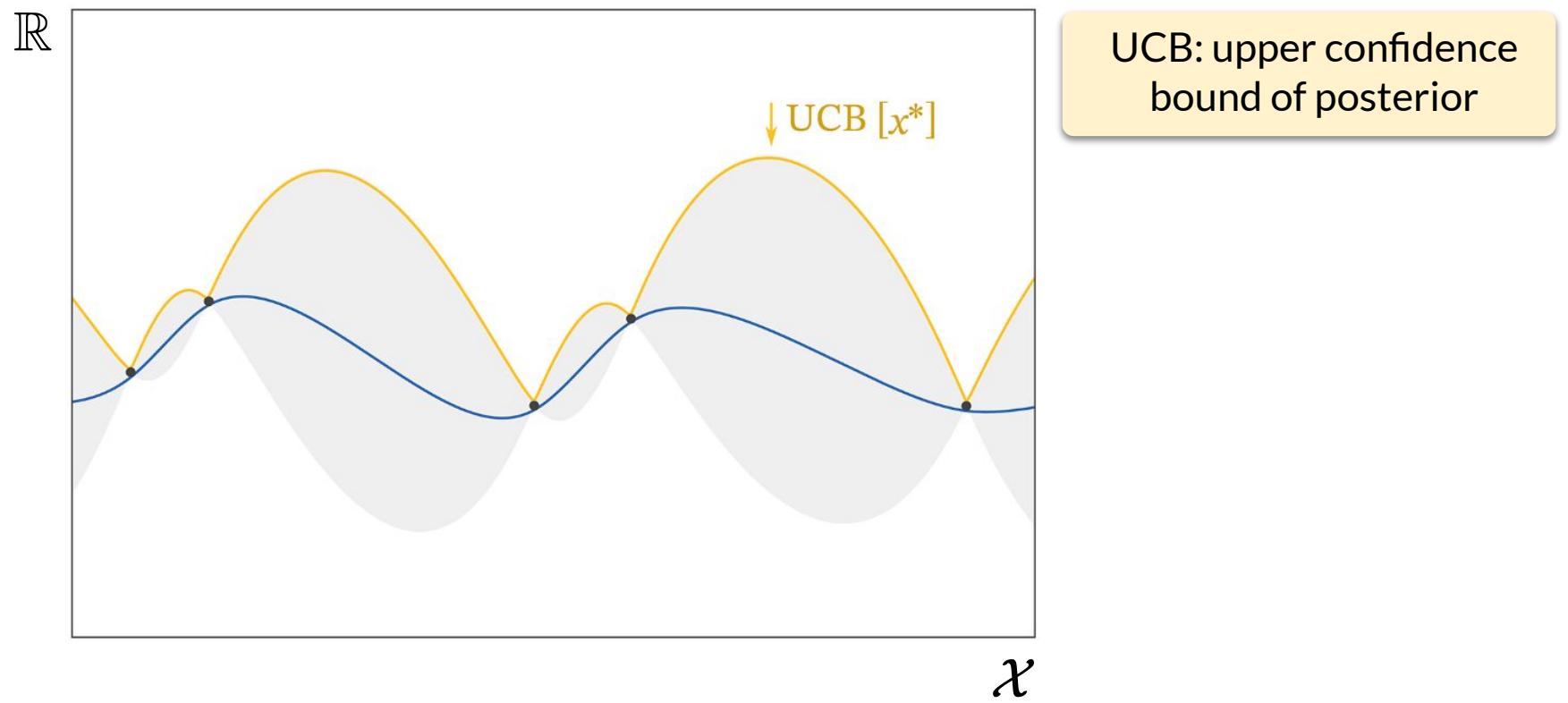
Nice visual comparing UCB vs. PI vs. EI:

Source: M. O. Ahmed, S. Prince, "Bayesian optimization"

Acquisition Functions – Expected Improvement (EI)

Nice visual comparing UCB vs. PI vs. EI:

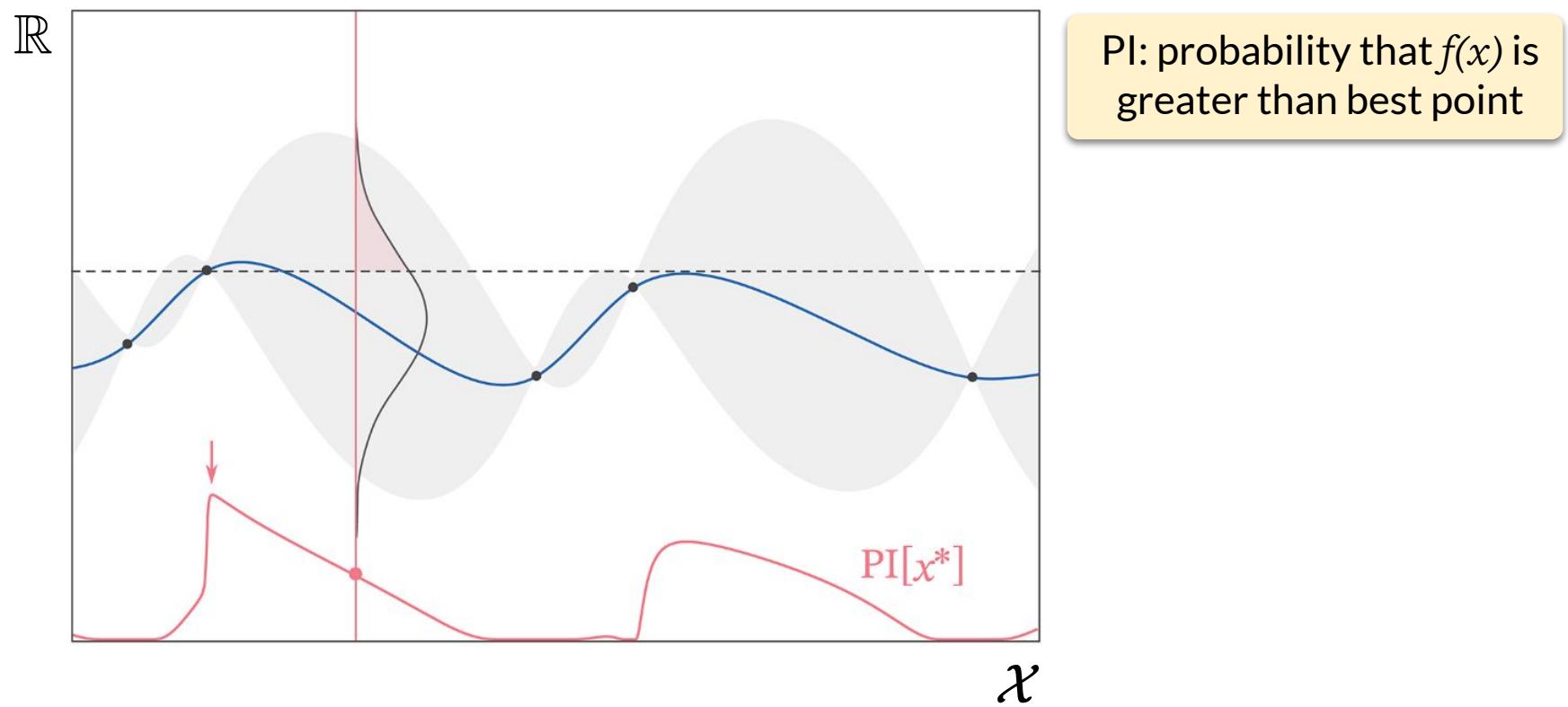
Source: M. O. Ahmed, S. Prince, "Bayesian optimization"



Acquisition Functions – Expected Improvement (EI)

Nice visual comparing UCB vs. PI vs. EI:

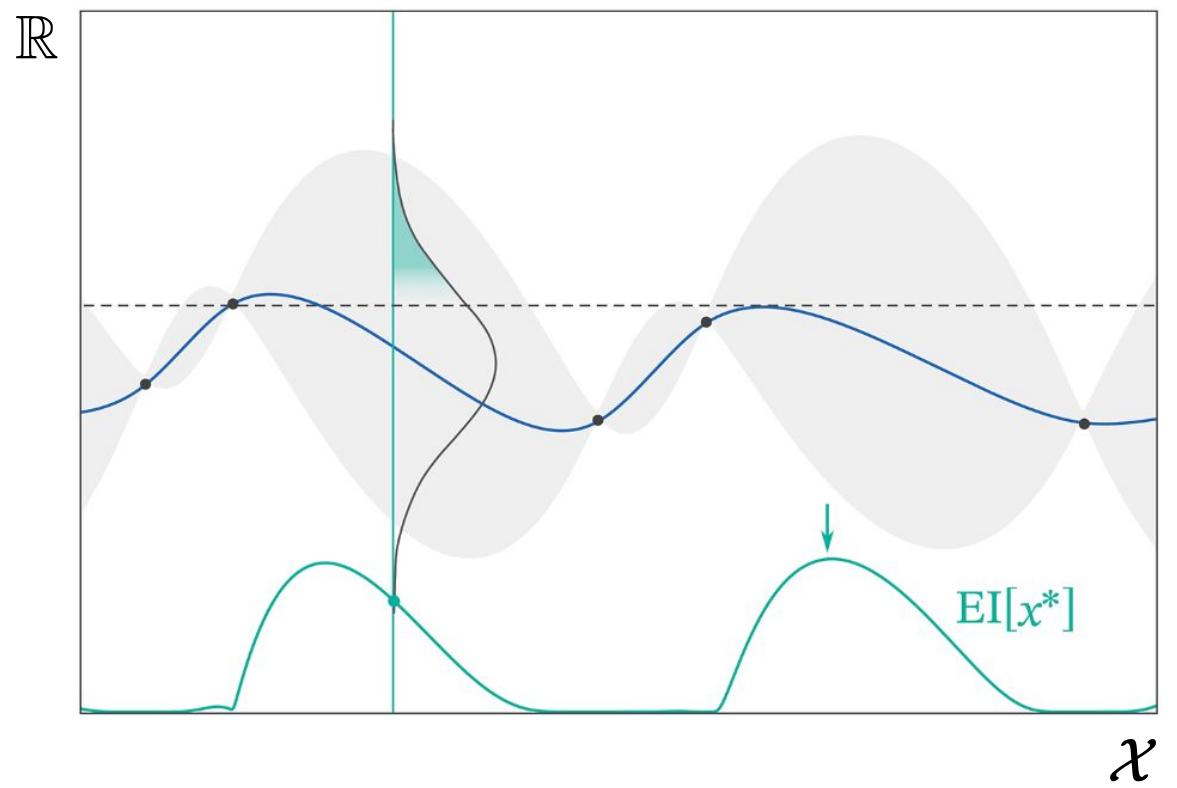
Source: M. O. Ahmed, S. Prince, "Bayesian optimization"



Acquisition Functions – Expected Improvement (EI)

Nice visual comparing UCB vs. PI vs. EI:

Source: M. O. Ahmed, S. Prince, "Bayesian optimization"



EI: expected amount $f(x)$ is greater than best point
(and ignoring amount that it's less than best point)

Acquisition Functions – Entropy Search (ES)

Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

Specifically, this means:

Reducing posterior uncertainty (entropy) over quantity of interest – in this case, the location x^* of the function maximizer – rather than, e.g., the full landscape of f .

Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

Specifically, this means:

Reducing posterior uncertainty (entropy) over quantity of interest – in this case, the location x^* of the function maximizer – rather than, e.g., the full landscape of f .

As mentioned, this is an *info-based* BO acquisition function called **entropy search**.

Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

Specifically, this means:

Reducing posterior uncertainty (entropy) over quantity of interest – in this case, the location x^* of the function maximizer – rather than, e.g., the full landscape of f .

As mentioned, this is an *info-based* BO acquisition function called **entropy search**.

Why would you want to use this acquisition function?

⇒ Intuitively: if you get rewarded for *how good your guess of x^* is* – rather than, the value of the maximizer, $f(x^*)$ – ES should be more optimal.

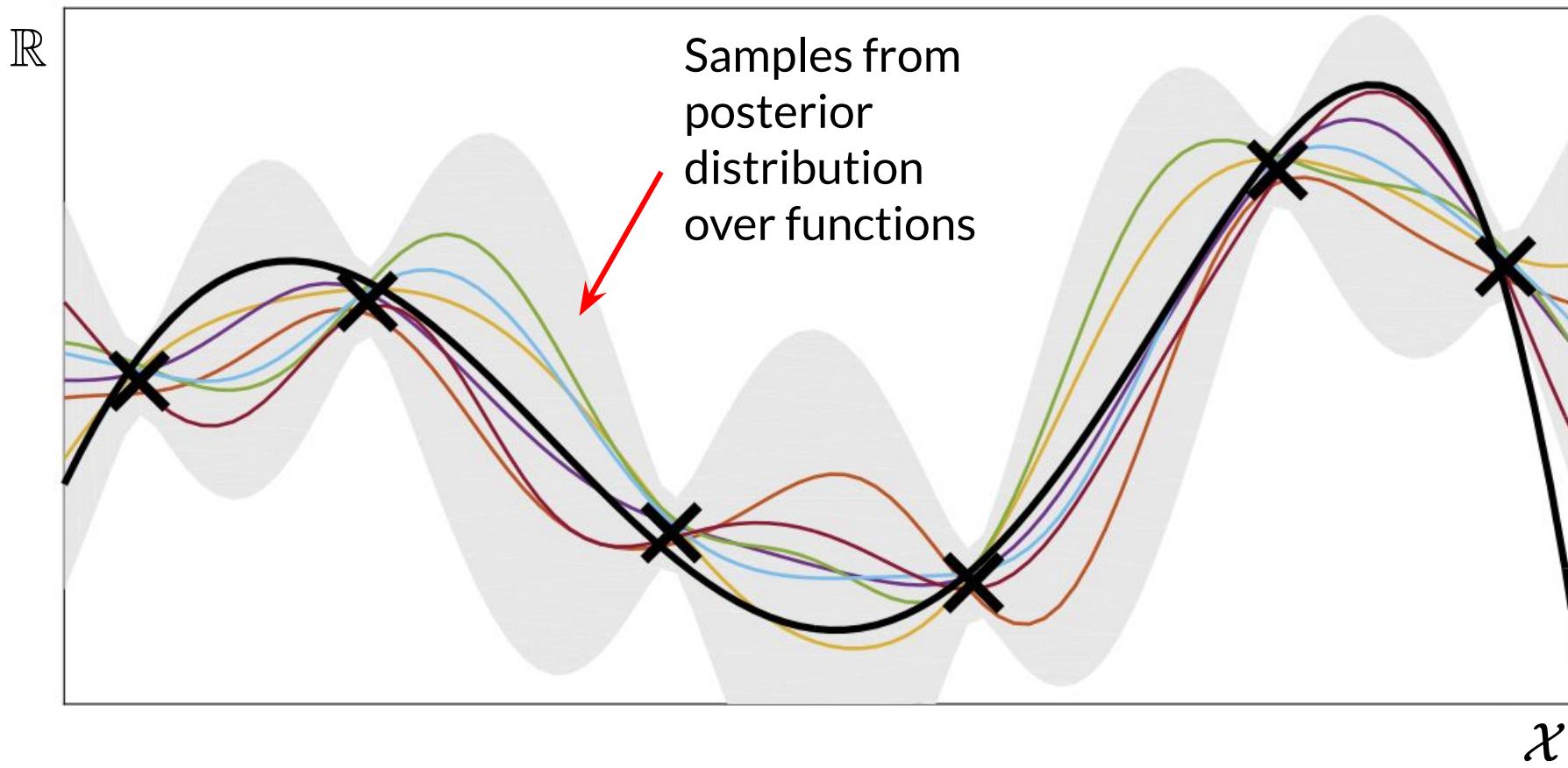
Acquisition Functions – Entropy Search

Visualizing ES

To explain this acquisition function, let's visualize...

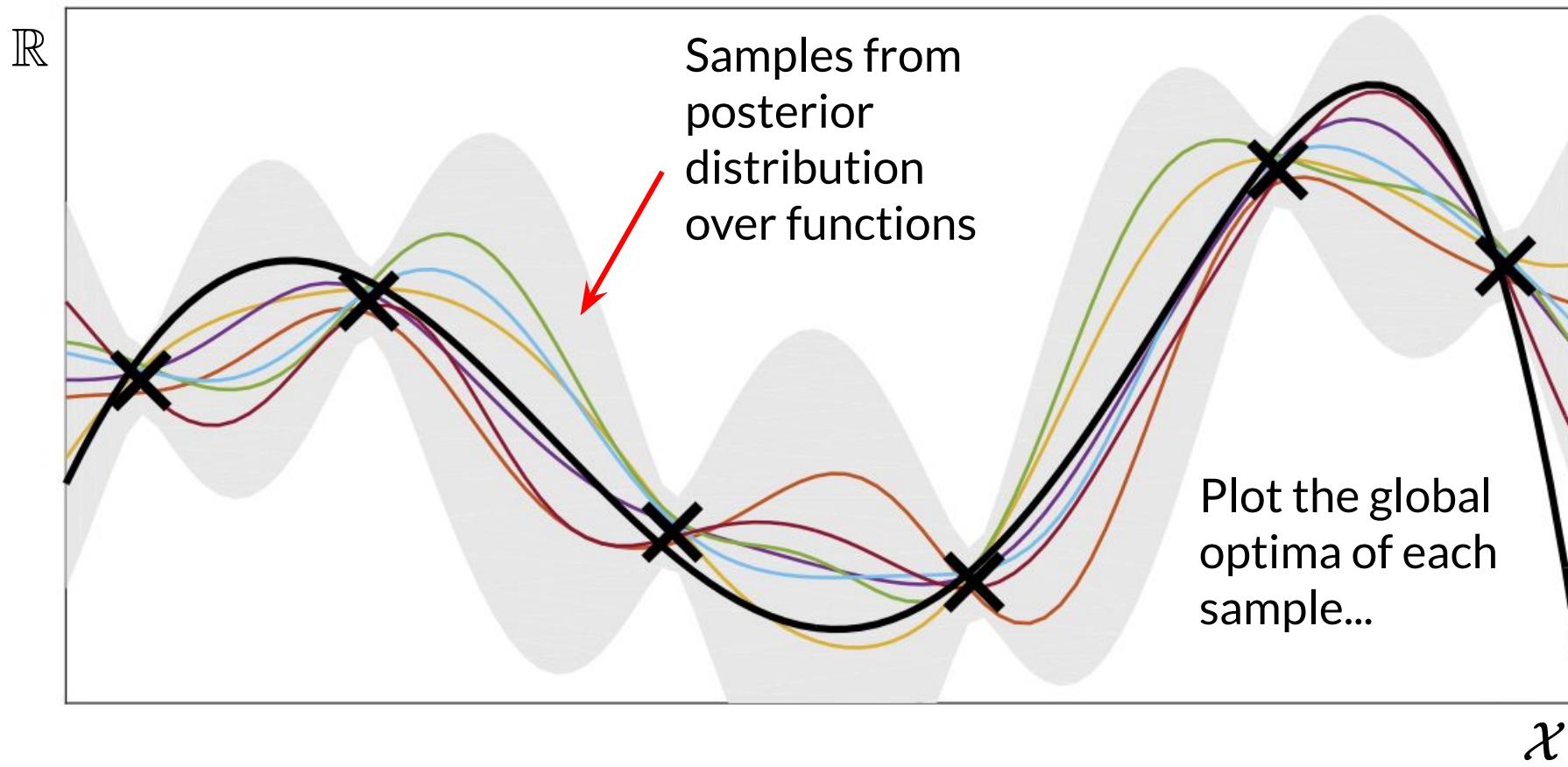
Acquisition Functions – Entropy Search

Visualizing ES



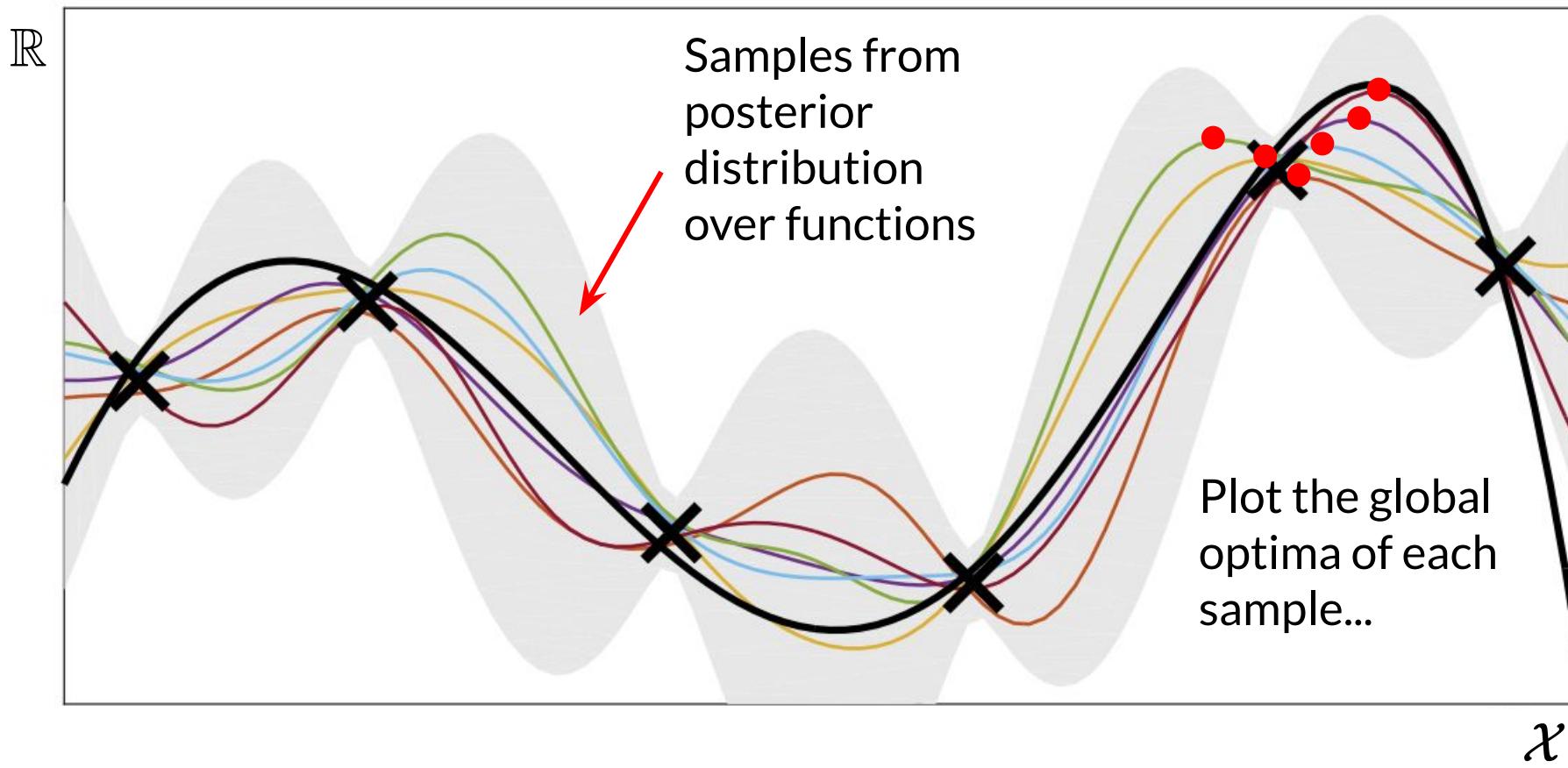
Acquisition Functions – Entropy Search

Visualizing ES



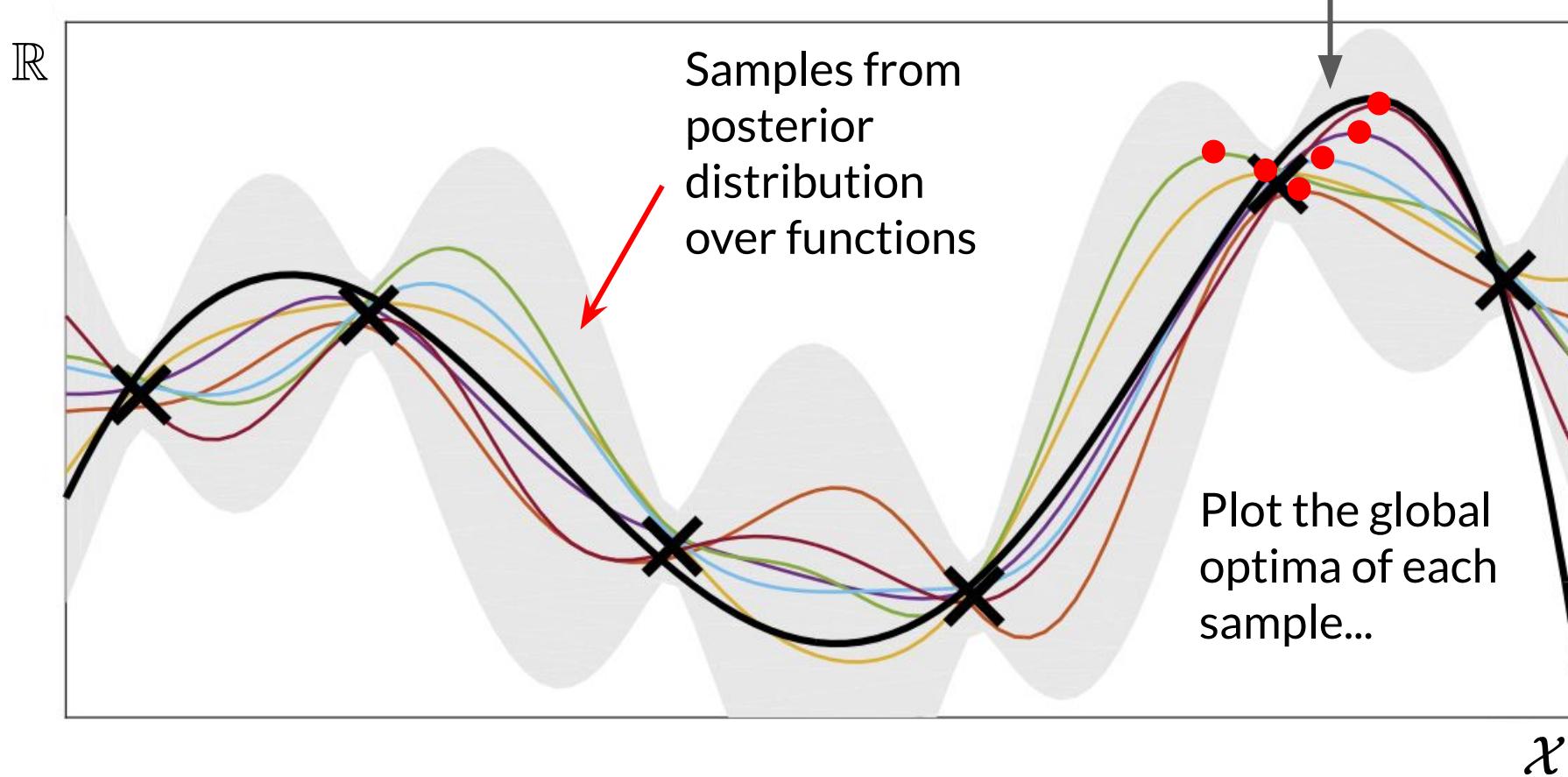
Acquisition Functions – Entropy Search

Visualizing ES



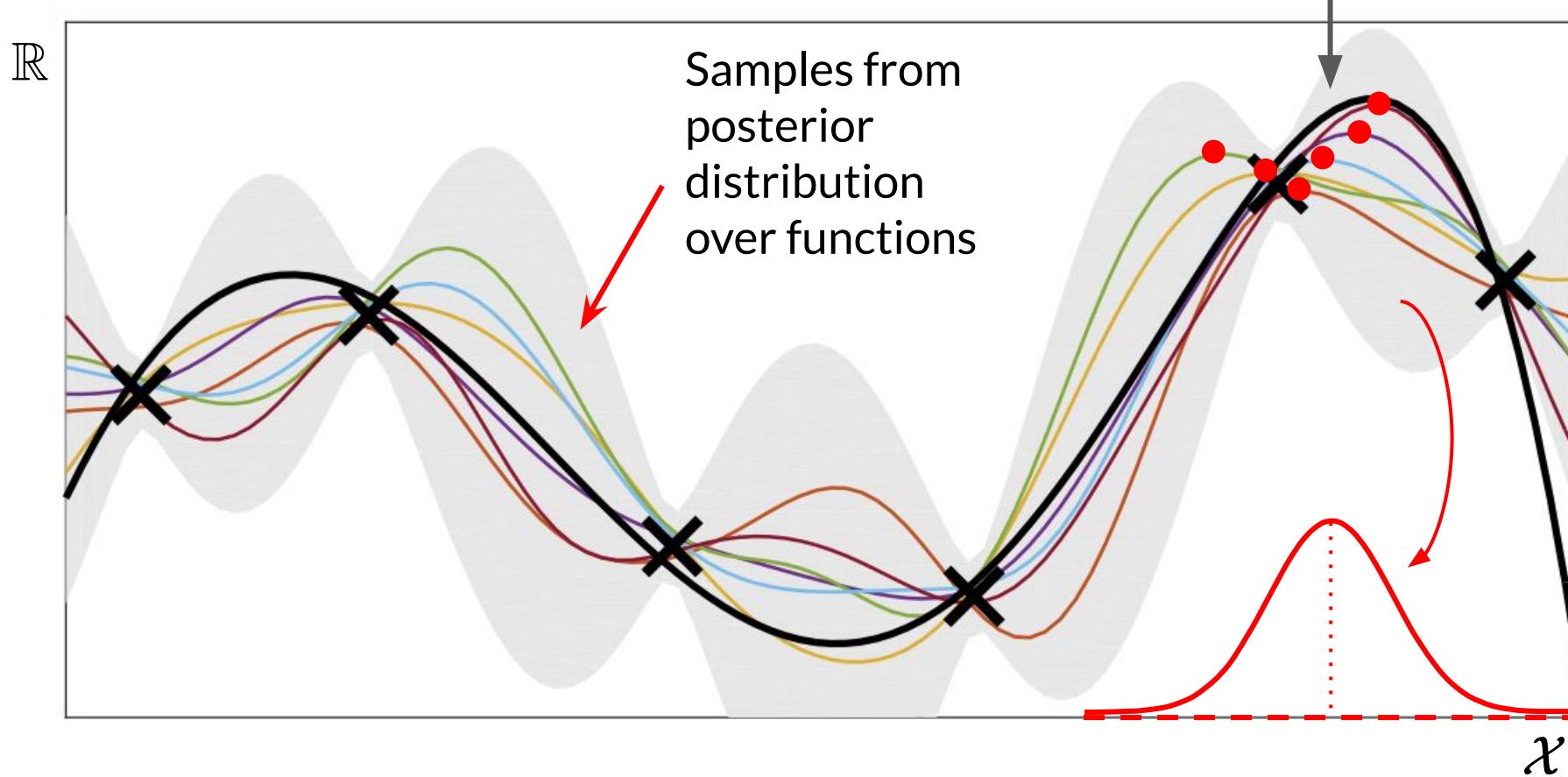
Acquisition Functions – Entropy Search

There is a posterior distribution over global optima induced by probabilistic model: $p(x^* | \mathcal{D}_t)$



Acquisition Functions – Entropy Search

There is a posterior distribution over global optima induced by probabilistic model: $p(x^* | \mathcal{D}_t)$



Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* | \mathcal{D}_t)] - \mathbb{E}_{p(y_x | \mathcal{D}_t)}[H[p(x^* | \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over
global optima x^*

minus

expected entropy of posterior distribution
over global optima x^* ,
... if we were to make a query at x

Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* | \mathcal{D}_t)] - \mathbb{E}_{p(y_x | \mathcal{D}_t)}[H[p(x^* | \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over
global optima x^*

minus

expected entropy of posterior distribution
over global optima x^* ,
... if we were to make a query at x

Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over global optima x^*

minus

expected entropy of posterior distribution over global optima x^* ,
... if we were to make a query at x

```
graph LR; A["\alpha_t(x) = H[p(x^* | \mathcal{D}_t)] - \mathbb{E}_{p(y_x | \mathcal{D}_t)}[H[p(x^* | \mathcal{D}_t \cup \{(x, y_x)\})]]"] --> B["entropy of posterior distribution over global optima x*"]; A --> C["expected entropy of posterior distribution over global optima x*, ... if we were to make a query at x"]; A --- D["minus"];
```

Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over global optima x^*

minus

expected entropy of posterior distribution over global optima x^* ,
... if we were to make a query at x

Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over global optima x^*

minus

expected entropy of posterior distribution over global optima x^* ,
... if we were to make a query at x

Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over global optima x^*

minus

expected entropy of posterior distribution over global optima x^* ,
... if we were to make a query at x

Expected information gain (EIG) – expected decrease in entropy if we were to query f at x .

Acquisition Functions – Additional

Acquisition Functions – Additional

A few acquisition functions that are useful for further study:

Acquisition Functions – Additional

A few acquisition functions that are useful for further study:

- **Knowledge gradient** – somewhat complex, but one of the *most-optimal* (both theoretically and in practice) BO acquisition functions.

Acquisition Functions – Additional

A few acquisition functions that are useful for further study:

- **Knowledge gradient** – somewhat complex, but one of the *most-optimal* (both theoretically and in practice) BO acquisition functions.
- **Thompson sampling** – a stochastic acquisition function (*i.e.*, draw one posterior, and use it within an acquisition function)

Acquisition Functions – Additional

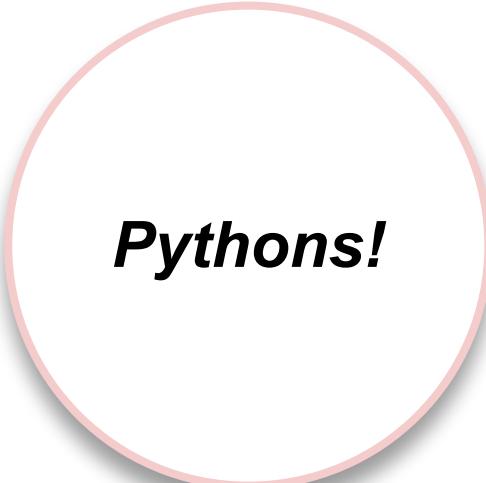
A few acquisition functions that are useful for further study:

- **Knowledge gradient** – somewhat complex, but one of the *most-optimal* (both theoretically and in practice) BO acquisition functions.
- **Thompson sampling** – a stochastic acquisition function (*i.e.*, draw one posterior, and use it within an acquisition function)
- **Non-BO acq fns.** – acquisition functions for non-optimization tasks, such as level sets, quadrature, phase boundaries, and more-general tasks.

Next Time

Next Time

Final presentations! Groups 1, 3, 9.



Pythons!



Geoguessrs!



GenIMG!

