

(1) Active Learning

(2) Bayesian Optimization

CS772A: Probabilistic Machine Learning

Piyush Rai

Plan for today

- Two important problems where model/predictive uncertainty is used
 - Active Learning
 - Bayesian Optimization

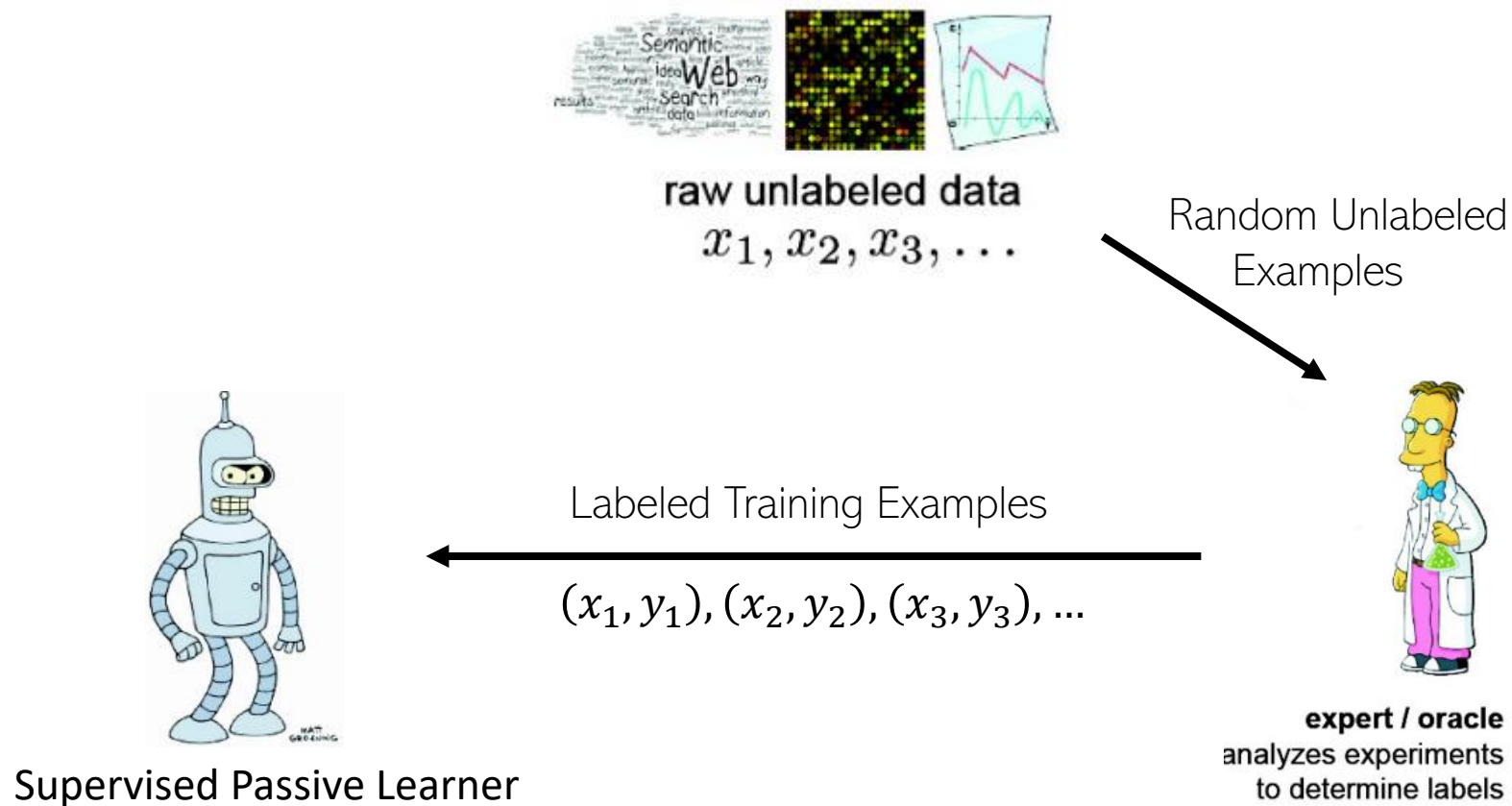


Active Learning



Passive Learning

- Standard supervised learning is passive
 - Learner has no control over what labelled training examples it gets to learn from



Active Learning

It is therefore also a sequential learning strategy (training data is not given all at once)

- In Active Learning, the learner can request specific labelled examples as it trains
 - In particular, examples **that the learner thinks will be most useful** to learn the underlying function

Will soon see what are some common notions of usefulness

Using the current model, identify the most useful example(s) from the unlabeled data pool



raw unlabeled data
 x_1, x_2, x_3, \dots

Although one example in each iteration is more common, labels of one or more than one examples can be queried ("batch" active learning)



Assumes some small amount of initial labelled training examples $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$ are available to learn an initial model

Query the expert for the true label of the selected unlabeled example, say x_1

$\langle x_1, ? \rangle$

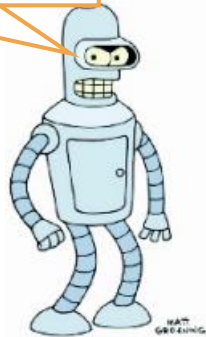
$\langle x_1, y_1 \rangle$

$\langle x_2, ? \rangle$

$\langle x_2, y_2 \rangle$



expert / oracle
analyzes experiments
to determine labels



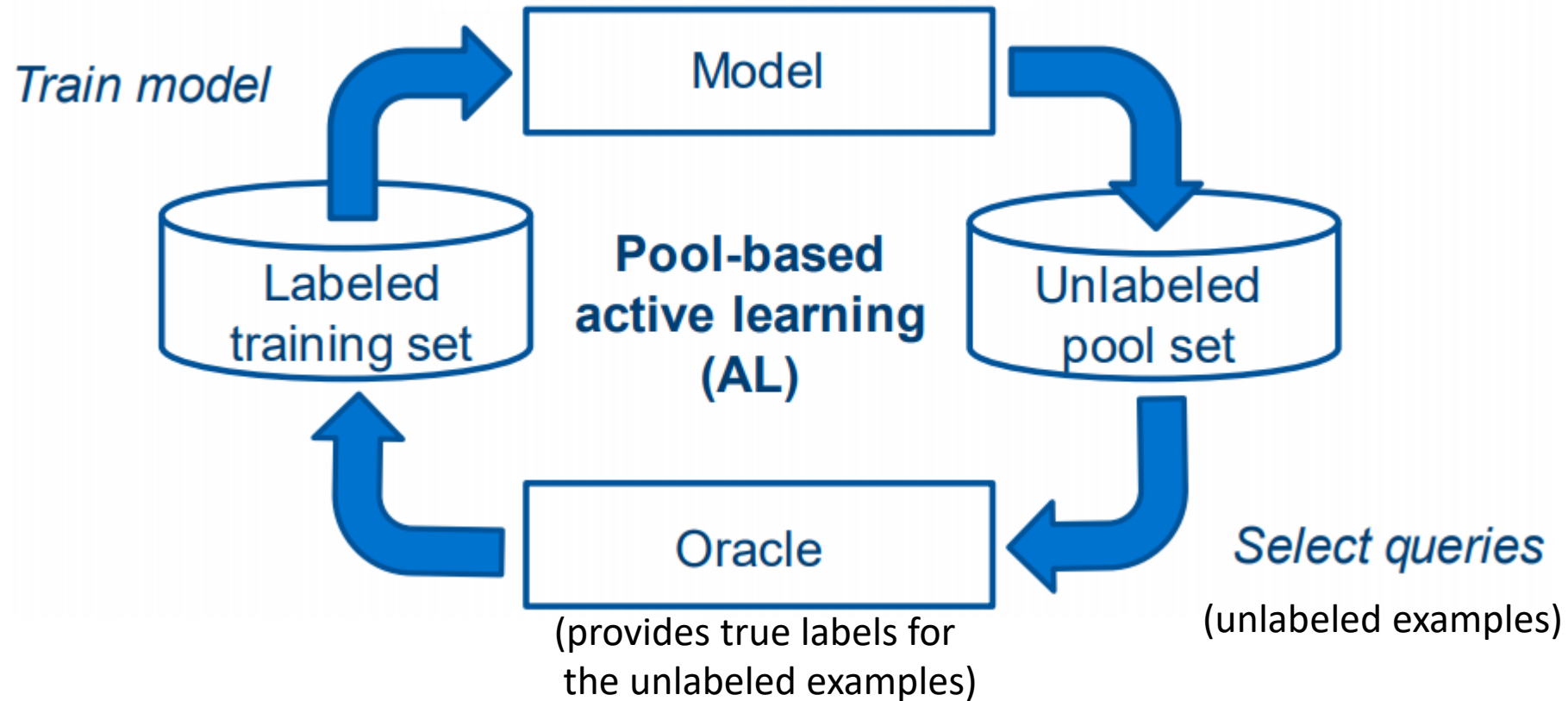
Supervised Active Learner

More labelled training examples will be acquired "actively" to improve the initial model by retraining it using the updated training data $\mathcal{D} = \{\mathcal{D} \cup (x_*, y_*)\}$ and repeating the process until we get the desired accuracy or our budget exhausts



Active Learning

- The figure below is another illustration of AL



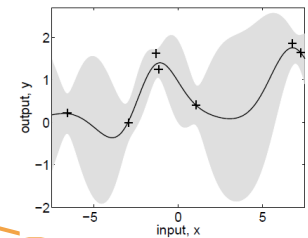
Measuring Usefulness in AL

Given the acquisition function, the most useful example can be selected from the pool as

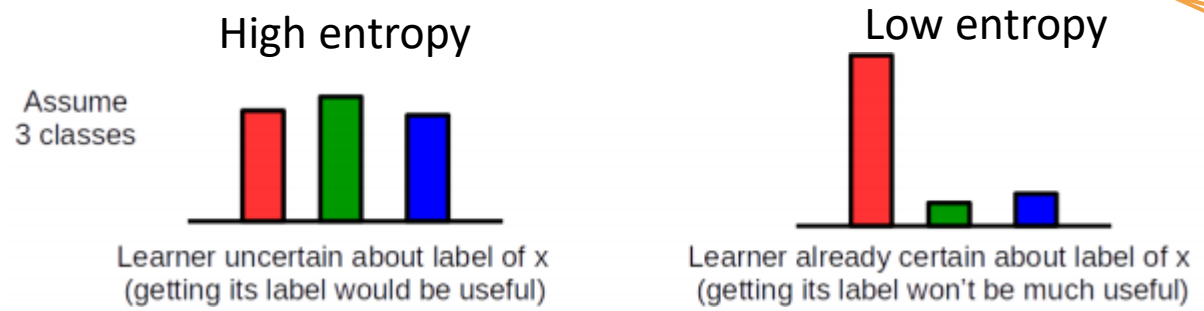
$$\hat{x}_* = \operatorname{argmax}_{x_* \in \mathcal{X}_{pool}} A(x_* | p(\theta | \mathcal{D}))$$


Note: We will use shorthand $A(x_*)$

- Various ways to measure the usefulness of an unlabeled example x_*
 - Defined by an “acquisition function” $A(x_*)$ (high value for most useful unlabeled examples)
- Approach 1: For x_* , look at uncertainty in output y_* predicted by the current model
 - Can use variance in the posterior predictive distribution: $A(x_*) = \operatorname{var}(y_*)$
 - More generally, can use entropy of the PPD: $A(x_*) = \mathbb{H}[p(y_* | x_*, \mathcal{D})]$



Note that this is the “marginal entropy” of the output distribution since the posterior predictive of the output is obtained by marginalizing over the posterior



- Approach 2: Look at how much our model will improve if we add this unlabeled example with its true label, to our training set, and retrain the model

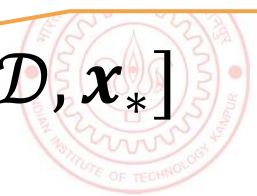
Mutual Information

$$A(x_*) = \mathbb{H}[p(\theta | \mathcal{D})] - \mathbb{E}_{p(y_* | x_*, \mathcal{D})} \mathbb{H}[p(\theta | \mathcal{D} \cup (x_*, y_*))] = \mathbb{I}[\theta; y_* | \mathcal{D}, x_*]$$

Entropy of the current posterior

Need to use expectation here since y_* is not known

Entropy of the new posterior after including the new example in our training set



Batch Active Learning

- Approaches we saw work by querying and adding one example at a time
- Expensive in practice since we have to retrain every time after including a new example
 - Especially true for deep learning models which are computationally expensive to train
- In practice, we want to use AL to jointly query the labels of $B > 1$ examples

$$(\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_B) = \operatorname{argmax}_{(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_B) \in \mathcal{X}_{pool}} A(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_B | p(\theta | \mathcal{D}))$$

- Difficult to construct such joint acquisition function and maximize them
- A greedy scheme is to simply select the B highest scoring points

$$A(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_B | p(\theta | \mathcal{D})) = \sum_{b=1}^B \mathbb{I}[\theta; y_b | \mathcal{D}, \mathbf{x}_b]$$

- The above however is myopic and ignores correlations among the selected points
 - Some recent works have addresses this issue¹

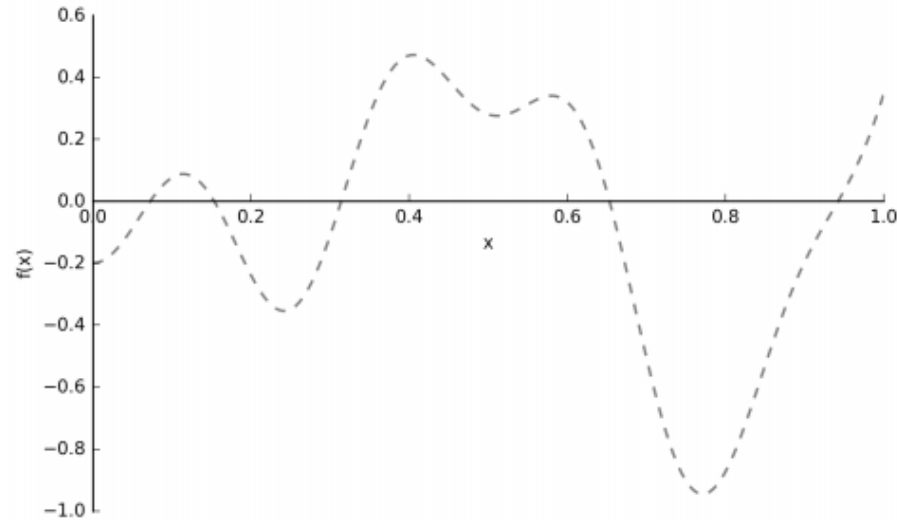


Bayesian Optimization



Bayesian Optimization: The Basic Formulation

- Consider finding the optima \mathbf{x}_* (say minima) of a function $f(\mathbf{x})$



- Caveat: We don't know the form of the function; can't get its gradient, Hessian, etc
- Can only query the **function's values** at certain points (i.e., only “black-box” access)
 - The values may or may not be noisy (i.e., we may be given $f(\mathbf{x})$ or $f(\mathbf{x}) + \epsilon$)



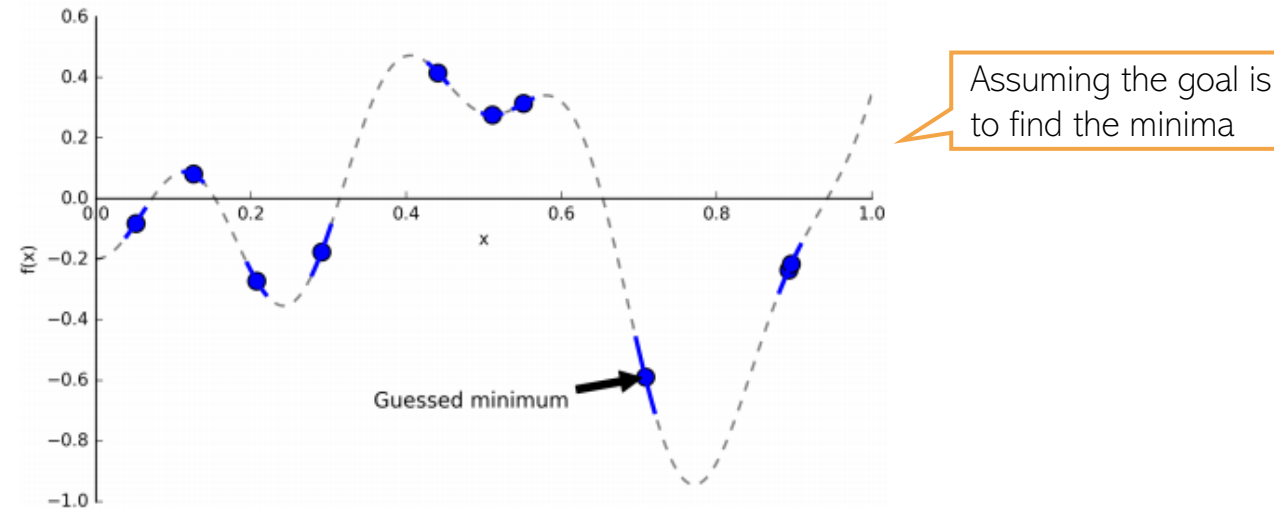
Bayesian Optimization: Some Applications

- Drug Design: Want to find the optimal chemical composition for a drug
 - Optimal composition will be the one that has the best efficacy
 - But we don't know the efficacy function
 - Can only know the efficacy via doing clinical trials
 - Each trial is expensive; can't do too many trials
- Hyperparameter Optimization: Want to find the optimal hyperparameters for a model
 - Optimal hyperparam values will be those that give the lowest test error
 - Don't know the true "test error" function
 - Need to train the model each time with different h.p. values and compute test error
 - Training every time will be expensive (e.g., for deep nets)
 - Note: Hyperparams here can even refer to the structure of a deep net (depth, width, etc)
- Many other applications: Website design via A/B testing, material design, optimizing physics based models (e.g., aircraft design), etc



Bayesian Optimization

- Can use BO to find maxima or minima
- Would like to locate the optima by querying the function's values (say, from an oracle)



- We would like to do so using as few queries as possible
- Reason: The function's evaluation may be time-consuming or costly



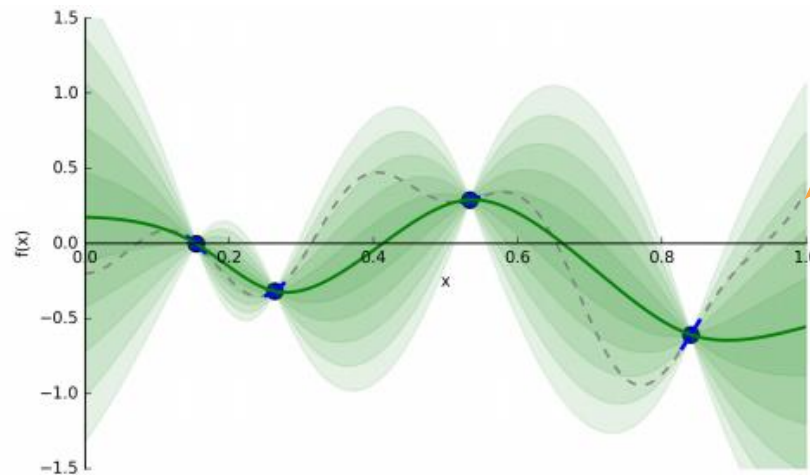
Bayesian Optimization

- Suppose we are allowed to make the queries sequentially
- This information will be available to us in form of query-function value pairs
- Queries so far can help us estimate the function

Note: Function values can be noisy too, e.g., $f(x_n) + \epsilon_n$

$\{(x_n, f(x_n))\}_{n=1}^N$

By solving a regression problem



Dotted curve: True function
Green curve: Current estimate ("surrogate") of the function
Shaded region: Uncertainty in the function's estimate

- BO uses **past queries** + **function's estimate+uncertainty** to decide **where to query next**
- Similar to Active Learning but the goal is to learn **f** as well as finds its optima



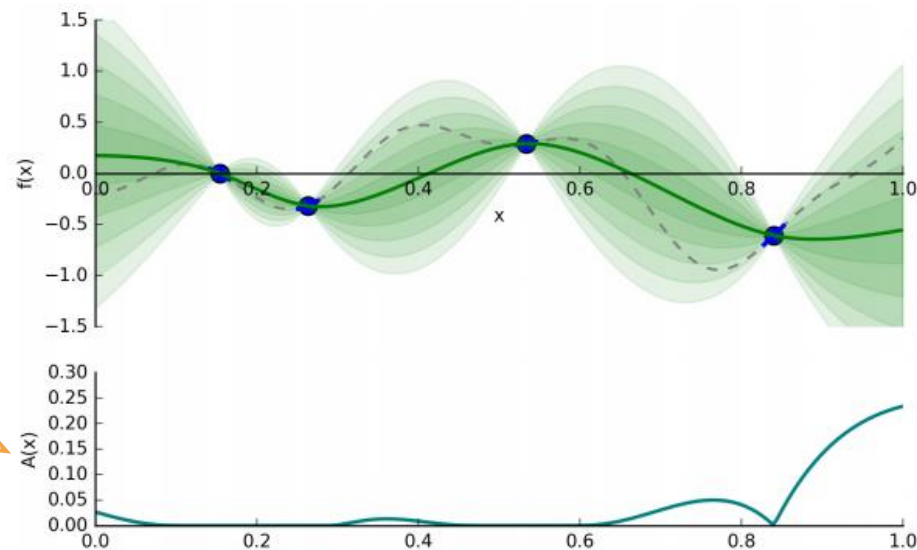
Bayesian Optimization

- BO requires two ingredients
 - A regression model to learn a surrogate of $f(x)$ given previous queries $\{(x_n, f(x_n))\}_{n=1}^N$
 - An acquisition function $A(x)$ to tell us where to query next

Note: Function values can be noisy too, e.g., $f(x_n) + \epsilon_n$

Assumption: $A(x)$ should be easier to optimize than $f(x)$

A typical example of what $A(x)$ might look like, assuming that the goal is to find the **maxima** of $f(x)$



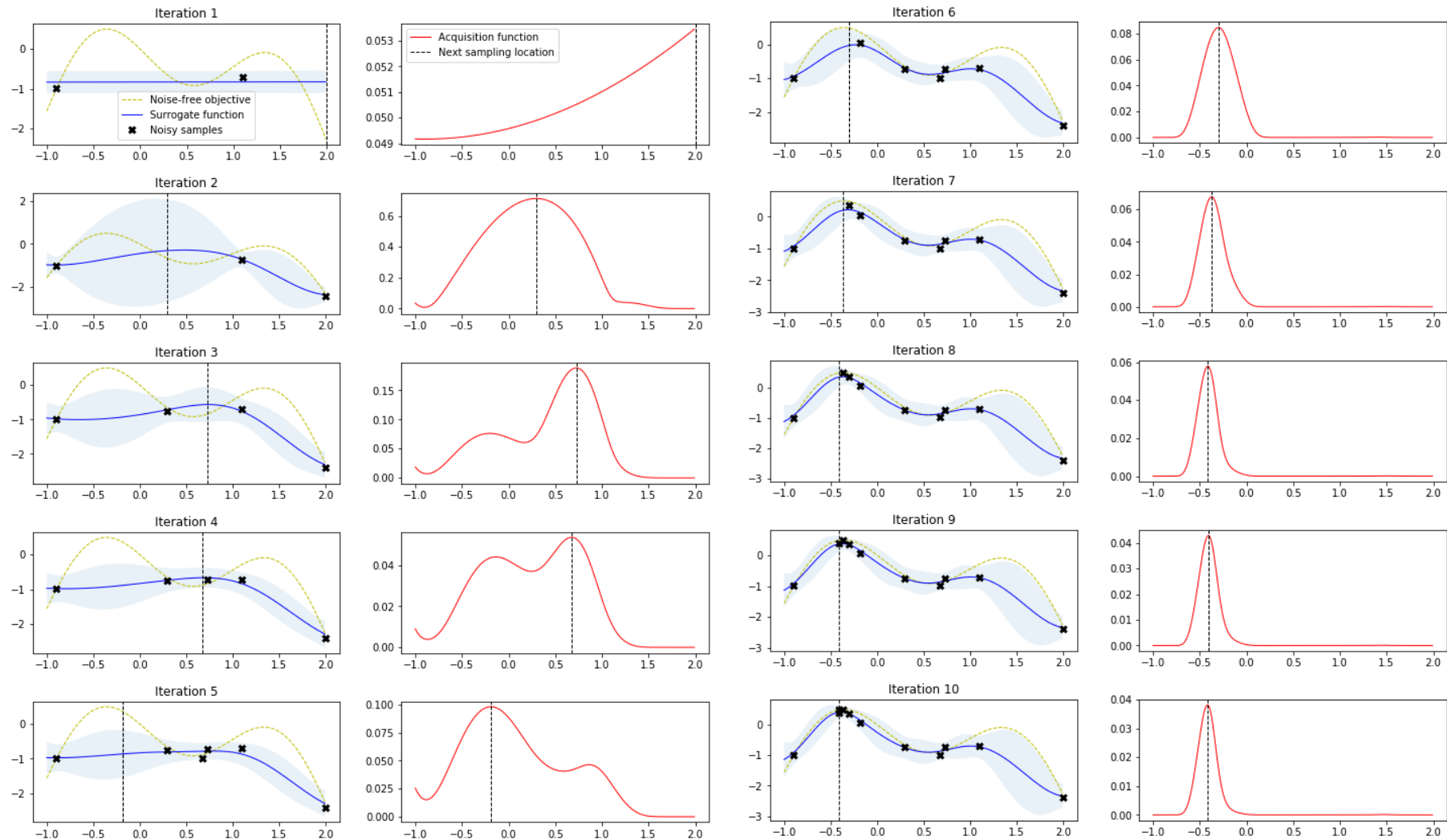
Dotted curve: True function
Green curve: Current surrogate of the function
Shaded region: Uncertainty in the function's estimate

- Note: The regression model must also have estimate of function's uncertainty
 - Bayesian nonlinear regression, such as GP, Bayesian Neural network, etc would be ideal



Bayesian Optimization: An Illustration

- Suppose our goal is to find the maxima of $f(x)$ using BO



Some Basic Acquisition Functions for BO

(assuming we are finding the minima)



Acquisition Functions: Probability of Improvement¹⁷

- Assume past queries $\mathcal{D}_N = (\mathbf{X}, \mathbf{f}) = (x_n, f(x_n))_{n=1}^N$ and suppose $f_{min} = \min \mathbf{f}$
- Suppose f_{new} denotes the function's value at the next query point x_{new}
- We have an improvement if $f_{new} < f_{min}$ (recall we are doing minimization)
- Assuming the function is real-valued, suppose the posterior predictive for x_{new} is

$$p(f_{new} | x_{new}, \mathcal{D}_N) = \mathcal{N}(f_{new} | \mu(x_{new}), \sigma^2(x_{new}))$$

- We can define a probability of improvement based acquisition function

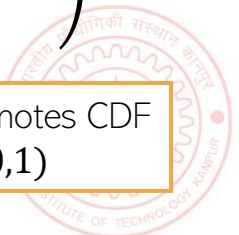
$$A_{PI}(x_{new}) = p(f_{new} \leq f_{min}) = \int_{-\infty}^{f_{min}} \mathcal{N}(f_{new} | \mu(x_{new}), \sigma^2(x_{new})) df_{new} = \Phi\left(\frac{f_{min} - \mu(x_{new})}{\sigma(x_{new})}\right)$$

- The optimal query point will be one that maximizes $A_{PI}(x_{new})$

$$x_* = \operatorname{argmax}_{x_{new}} A_{PI}(x_{new})$$

Exercise: Verify

$\Phi()$ denotes CDF of $\mathcal{N}(0,1)$



Acquisition Functions: Expected Improvement

- PI doesn't take into account the amount of improvement
- Expected Improvement (EI) takes this into account and is defined as

$$A_{EI}(x_{new}) = \mathbb{E}[f_{min} - f_{new}] = \int_{-\infty}^{f_{min}} (f_{min} - f_{new}) \mathcal{N}(f_{new} | \mu(x_{new}), \sigma^2(x_{new})) df_{new}$$

Exercise: Prove this result

$$= (f_{min} - \mu(x_{new})) \Phi\left(\frac{f_{min} - \mu(x_{new})}{\sigma(x_{new})}\right) + \sigma(x_{new}) \mathcal{N}\left(\frac{f_{min} - \mu(x_{new})}{\sigma(x_{new})}; 0, 1\right)$$

- The optimal query point will be one that maximizes $A_{EI}(x_{new})$

$$x_* = \operatorname{argmax}_{x_{new}} A_{EI}(x_{new})$$

Focus on points where the function has small values (since we are looking for its minima)

Focus on points where the function has high uncertainty (so that including them improves our estimate of the function)

- Note that the above acquisition function trades off exploitation vs exploration
 - Will prefer points with small predictive mean $\mu(x_{new})$: Exploitation
 - Will prefer points with large predictive variance $\sigma(x_{new})$: Exploration



Acquisition Functions: Lower Confidence Bound

- Lower Confidence Bound (LCB) also takes into account exploitation vs exploration

- Used when the regression model is a Gaussian Process (GP)

When using BO for maximization, we use Upper Confidence Bound (UCB) defined as $A_{UCB}(x_{new}) = \mu(x_{new}) + \kappa \sigma(x_{new})$ and $x_* = \operatorname{argmax}_{x_{new}} A_{UCB}(x_{new})$

- Assume the posterior predictive for a new point to be

$$p(f_{new}|x_{new}, \mathcal{D}_N) = \mathcal{N}(f_{new}|\mu(x_{new}), \sigma^2(x_{new}))$$

- The LCB based acquisition function is defined as

$$A_{LCB}(x_{new}) = \mu(x_{new}) - \kappa \sigma(x_{new})$$

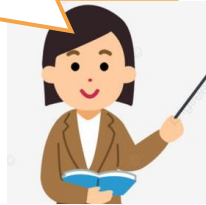
Thus prefer points at which the function has low mean but high variance

- Point with the smallest LCB is selected as the next query point

$$x_* = \operatorname{argmin}_{x_{new}} A_{LCB}(x_{new})$$

- κ is a parameter to trade-off exploitation (low mean) and exploration (high variance)

- Under certain conditions, the iterative application of this acquisition function will converge to the true global optima of f (Srinivas et al. 2010)



Bayesian Optimization: The Overall Algo

- Initialize $\mathcal{D} = \{\}$
- For $n = 1, 2, \dots, N$ (or until the budget doesn't exhaust)
 - Select the next query point \mathbf{x}_n by optimizing the acquisition function

$$\mathbf{x}_n = \operatorname{argopt}_x A(\mathbf{x})$$

- Get function's value from the black-box oracle: $f_n = f(\mathbf{x}_n)$

- $\mathcal{D} = \{\mathcal{D} \cup (\mathbf{x}_n, f_n)\}$

Can get the function's minima from this set of function's values

- Update the regression model for f using data \mathcal{D}



BO: Some Challenges/Open Problems

- Learning the regression model for the function
 - GPs are flexible but can be expensive as N grows
 - Bayesian neural networks can be an more efficient alternative to GPs (Snoek et al, 2015)
 - Hyperparams of the regression model itself (e.g., GP cov. function, Bayesian NN hyperparam)
- High-dimensional Bayesian Optimization (optimizing functions of many variables)
 - Most existing methods work well only for a moderate-dimensional x
 - Number of function evaluations required would be quite large in high dimensions
 - Lot of recent work on this (e.g., based on dimensionality reduction)
- Multitask Bayesian Optimization (joint BO for several related functions)
 - Basic idea: If two functions are similar their optima would also be nearby



BO: Some Further Resources

- Some survey papers:
 - A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning (Brochu et al., 2010)
 - Taking the Human Out of the Loop: A Review of Bayesian Optimization (Shahriari et al., 2015)
- Some open source software libraries
 - BoTorch: Bayesian Optimization in PyTorch
 - GPflowOpt: Bayesian Optimization in Tensorflow (and using GP for modeling the function)
 - Also available in scikit-optimize

