

Tuning the knobs of neural networks using Bayesian Optimization

Niall Turbitt



About me

- Machine Learning Engineer @ Argos
- Previously
 - Data Scientist @ Lucidworks
 - Masters in Statistics @ UCD Dublin
 - Mathematics & Economics @ Trinity College Dublin

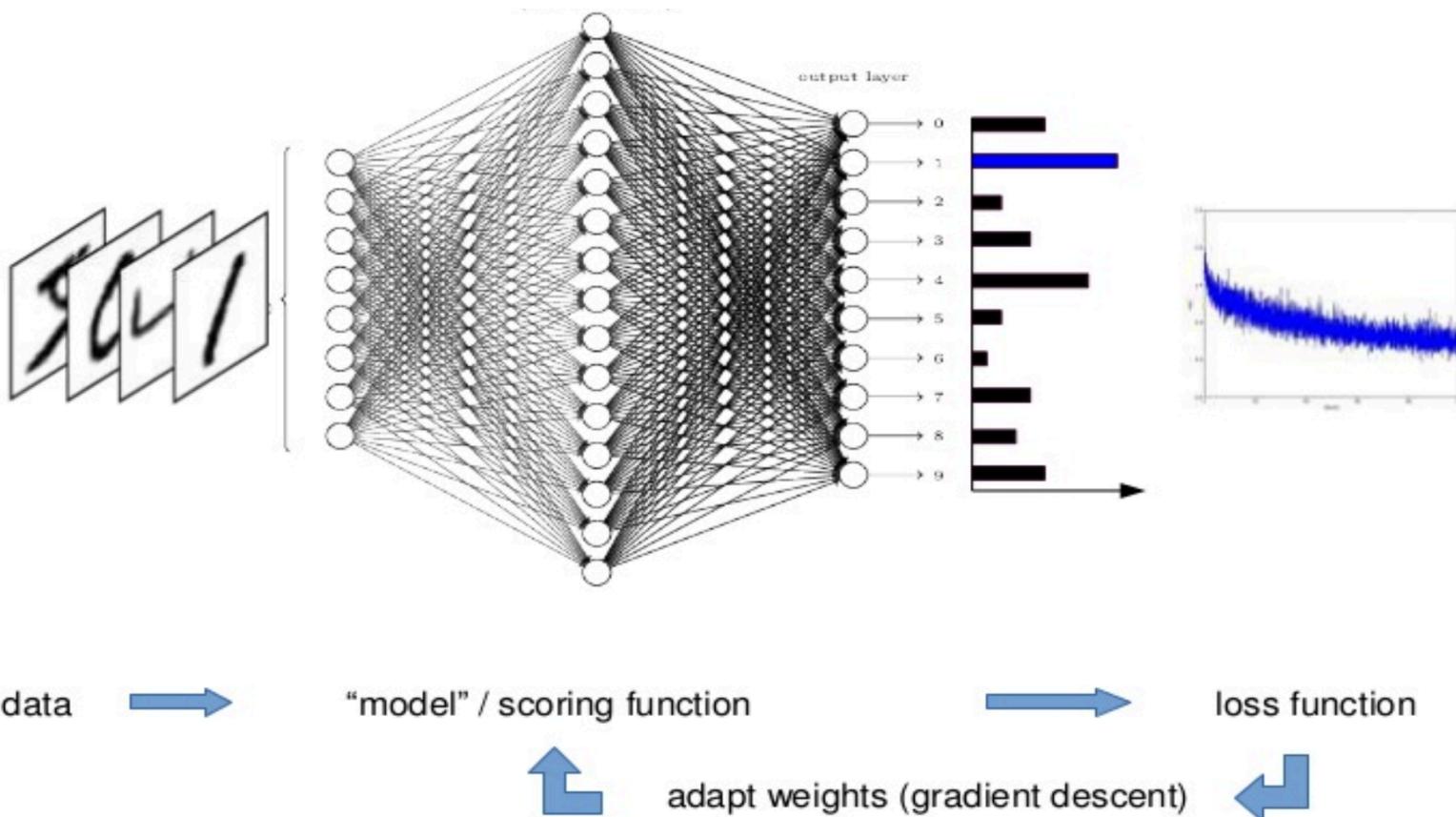
<https://github.com/niallturb/pylondinium18>

What this talk will cover

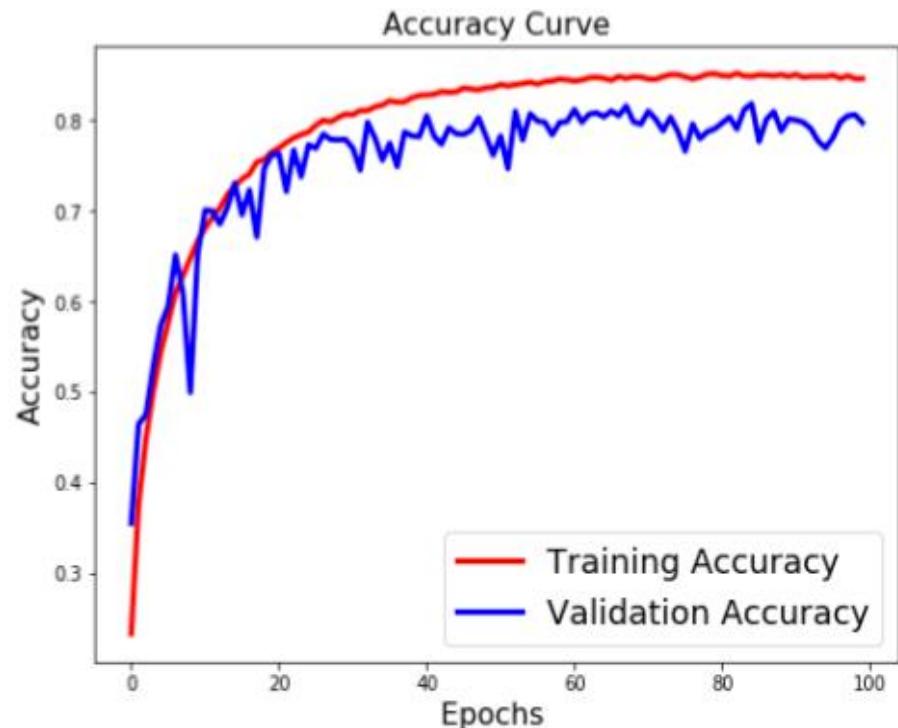
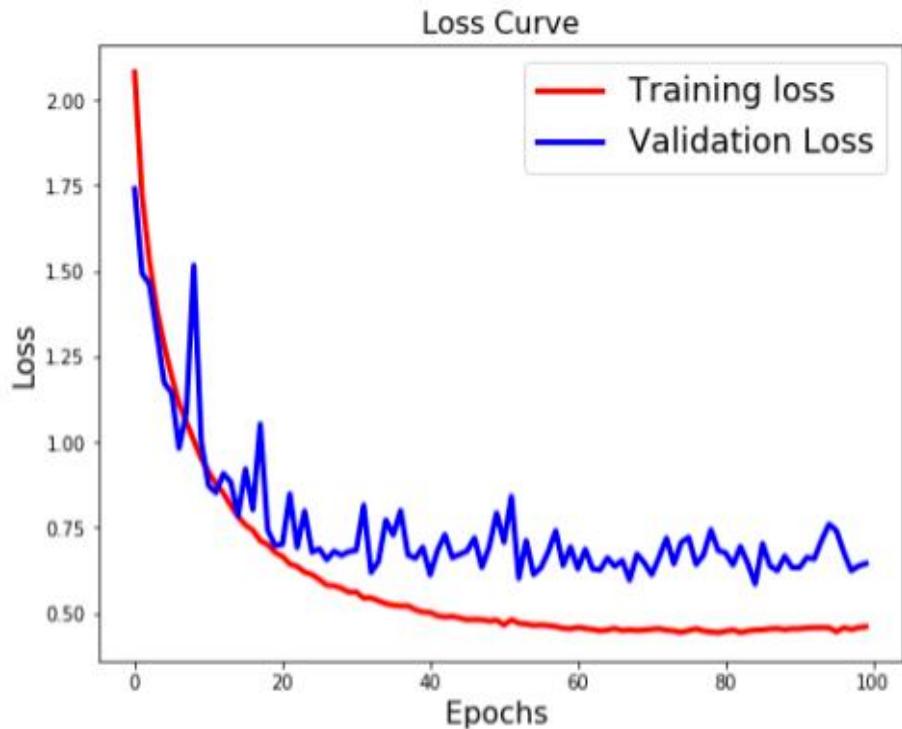
- Overview of neural networks
 - Convolutional Neural Networks
- The hyperparameter tuning problem
- Bayesian Optimization
- Fashion MNIST demo

Neural Networks – Brief Overview

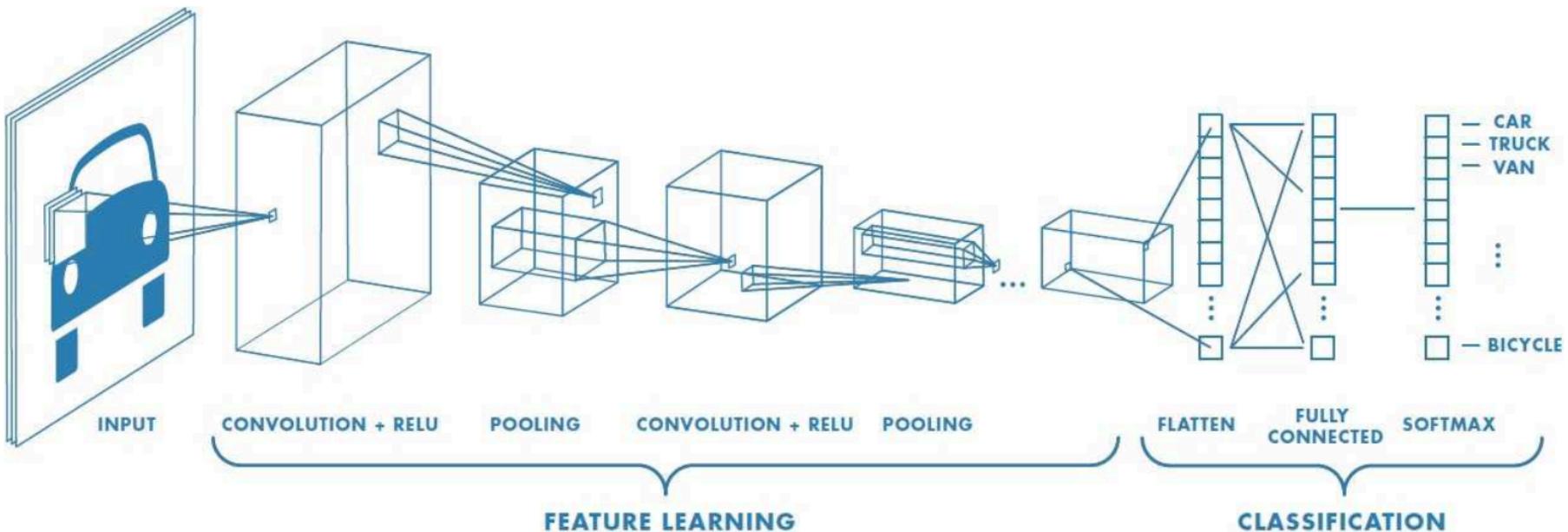
Training a Neural Network



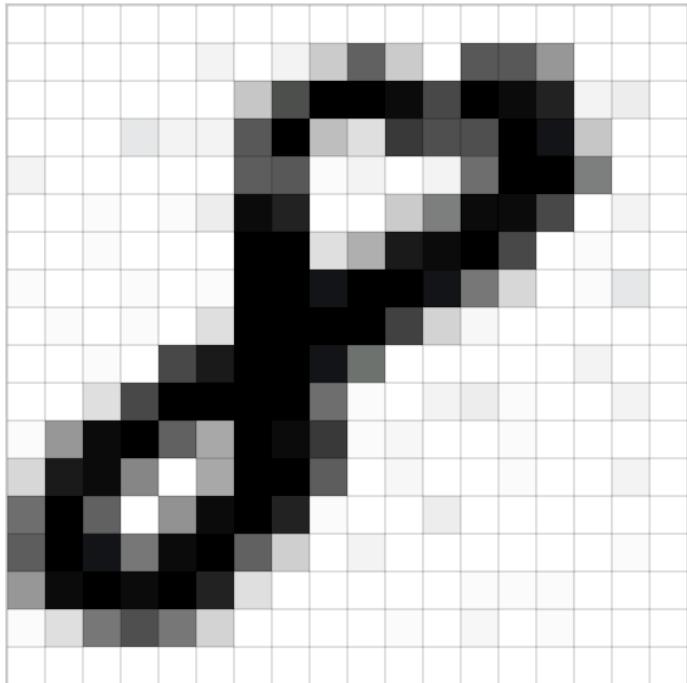
Training should converge



Convolutional Neural Networks (CNNs)



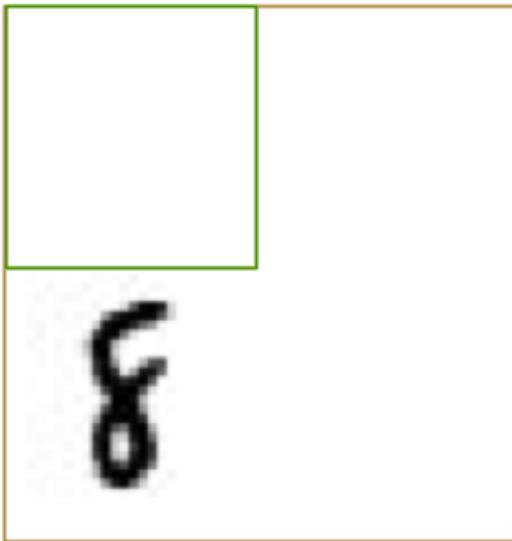
How can we input image data?



2

Tunnel Vision

Test Image



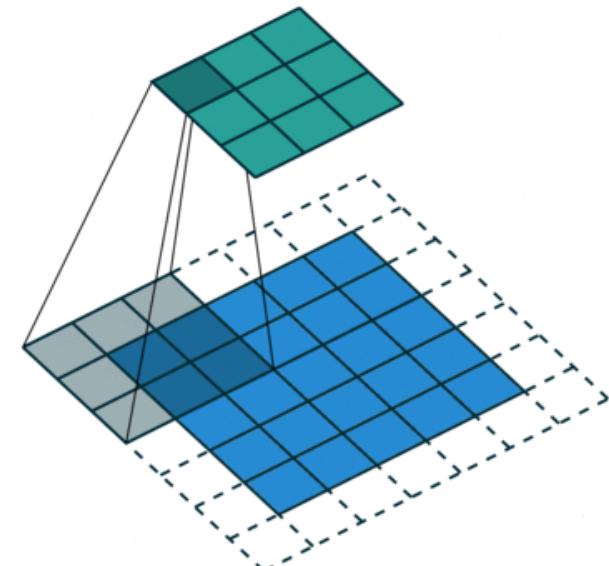
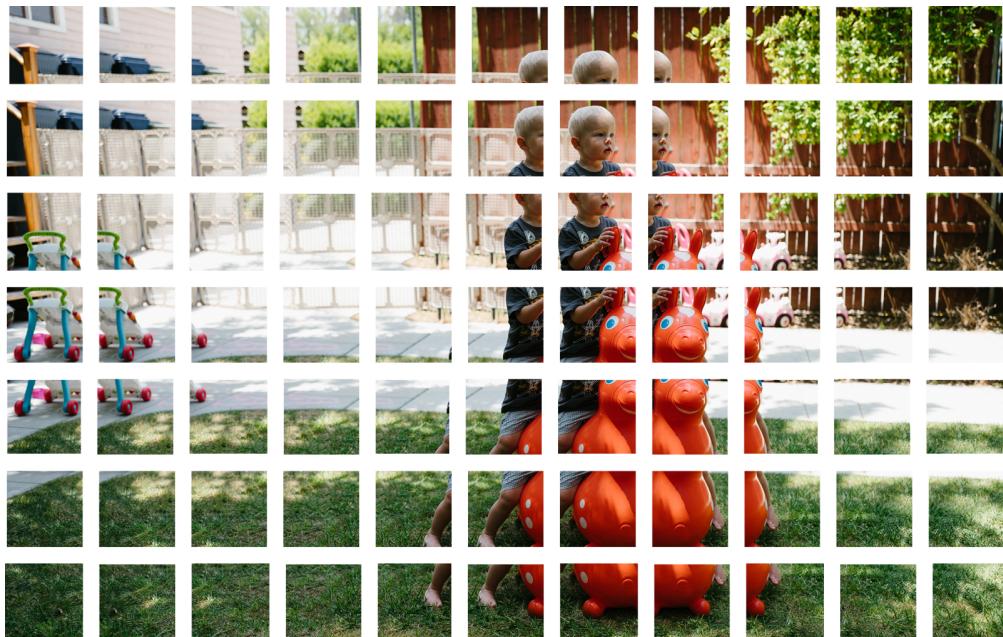
Prediction from
our network

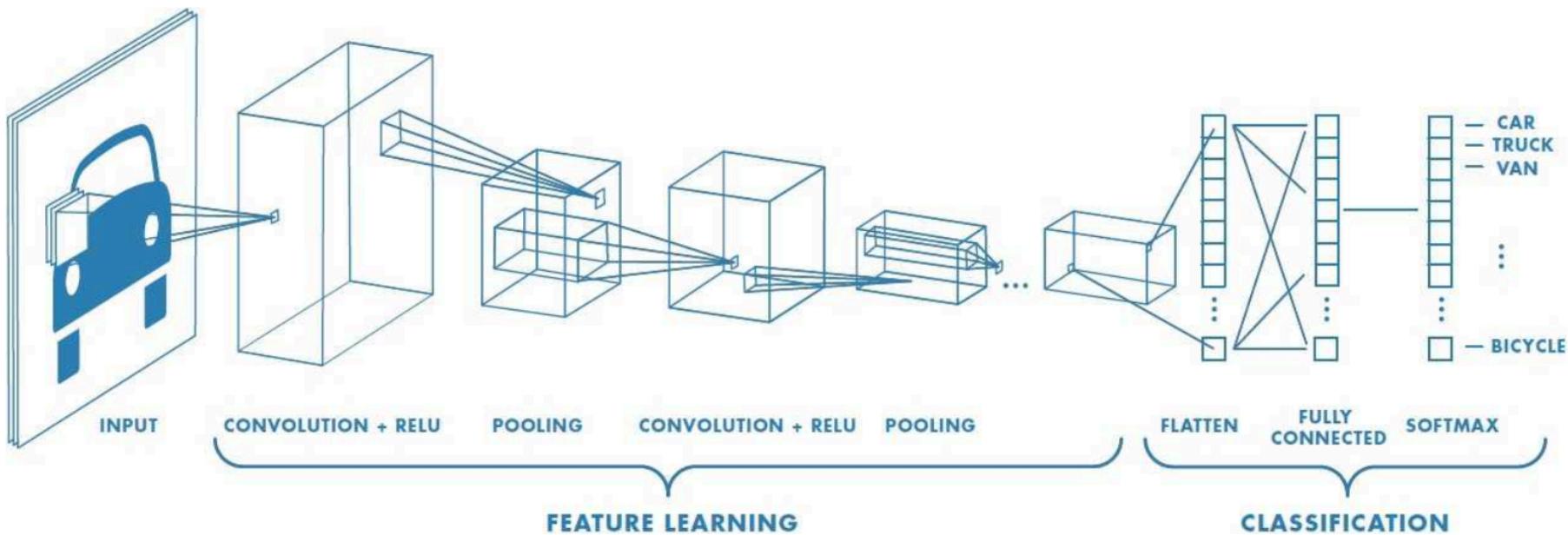
No idea!?!

Convolution is the solution



How Convolution Works





Constructing the right CNN

- Number of layers
- Number of convolution filters
- Learning rate
- Number of epochs
- Batch size
- Activation function
- Weight initialization
- Dropout rate

Hyperparameter tuning

- Grid Search
 - $O(n^k)$
- Random Search (Bergstra & Bengio, 2012)
- Bayesian Optimization

Bayesian Optimization

Bayesian Optimization

- Tool for globally optimizing objective functions which are costly or slow to evaluate
- Sequential model-based optimization (SMBO)

$$x_M = \arg \min_{x \in \mathcal{X}} f(x)$$

1. Initialize a Gaussian Process on a small set of samples from our domain, \mathcal{X} to compute a prior probability model
2. Sequentially select new locations within the domain by optimizing some *acquisition function*
3. At each iteration, we update our belief of what our objective function f looks like, having computed the loss for a new configuration of hyperparameters.

Gaussian Process

- A Gaussian Process is the generalization of a Gaussian distribution to a distribution over functions

$$f(x) \sim \mathcal{GP}(\mu(x), k(x, x'))$$

- Induces a posterior distribution over the loss function that is analytically tractable

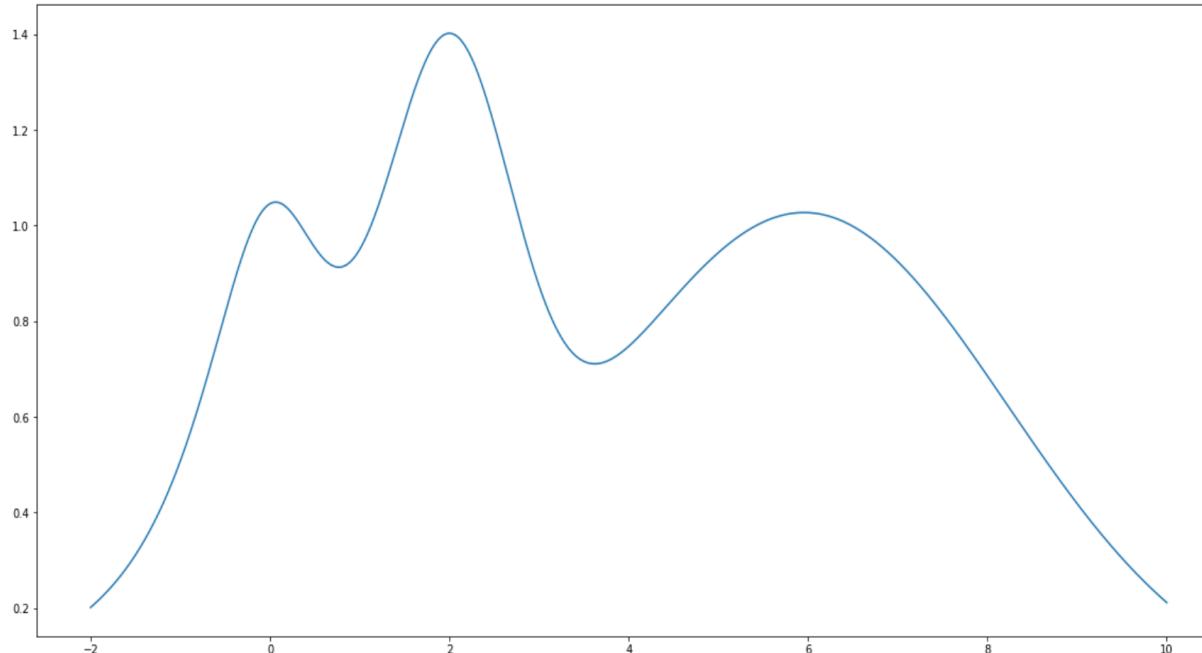
Posterior probability \propto Likelihood \times Prior probability

Acquisition Function

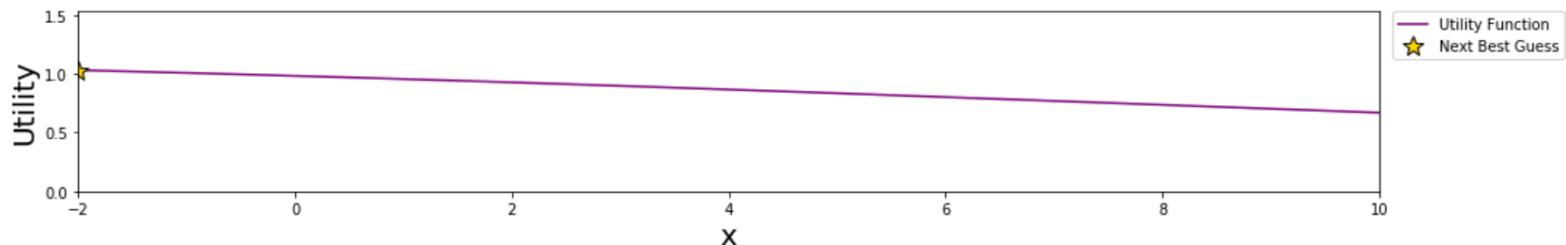
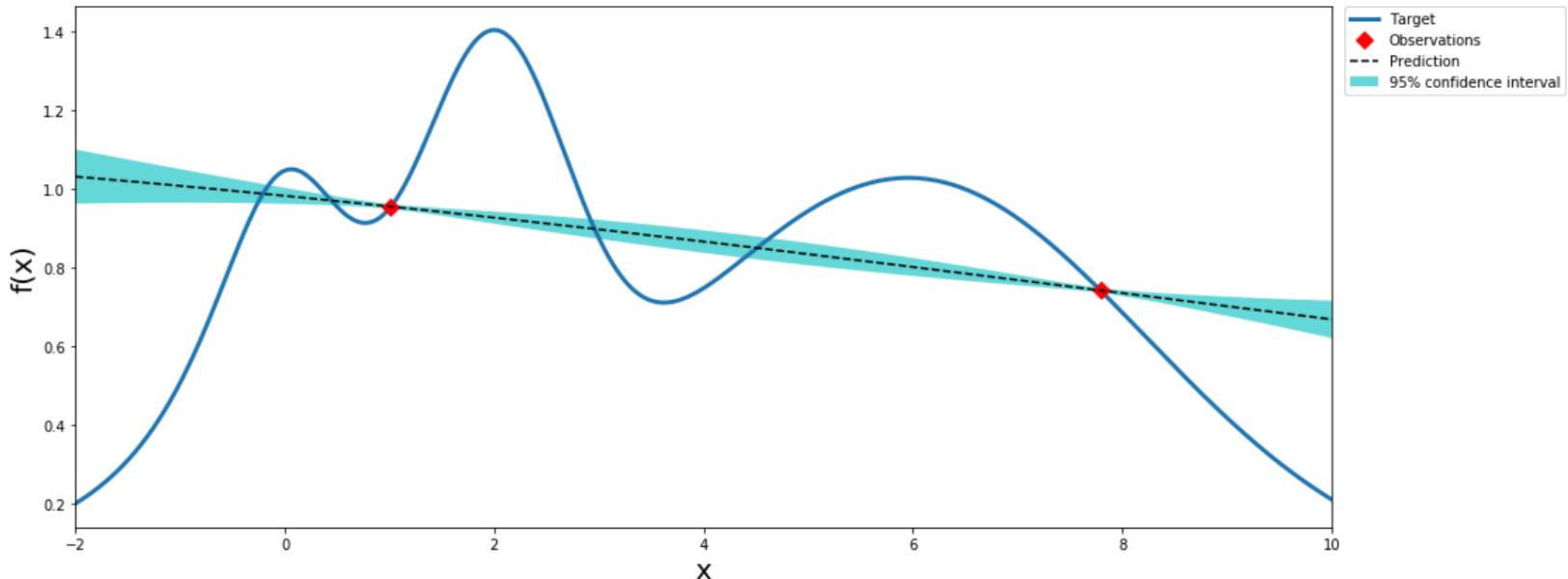
- Surrogate function that represents our beliefs of $f(x)$
 - Expected Improvement
 - Upper Confidence Bound (UCB)
 - Maximum probability of improvement (MPI)
- Exploration/Exploitation trade off

One Dimensional Bayesian Optimization Example

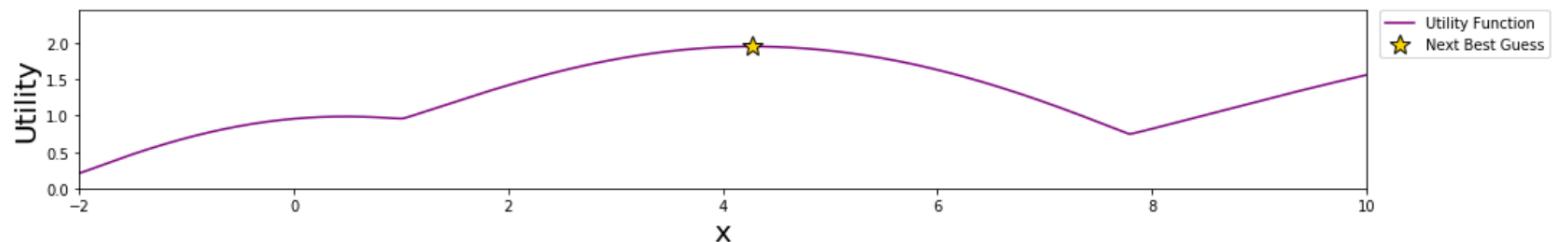
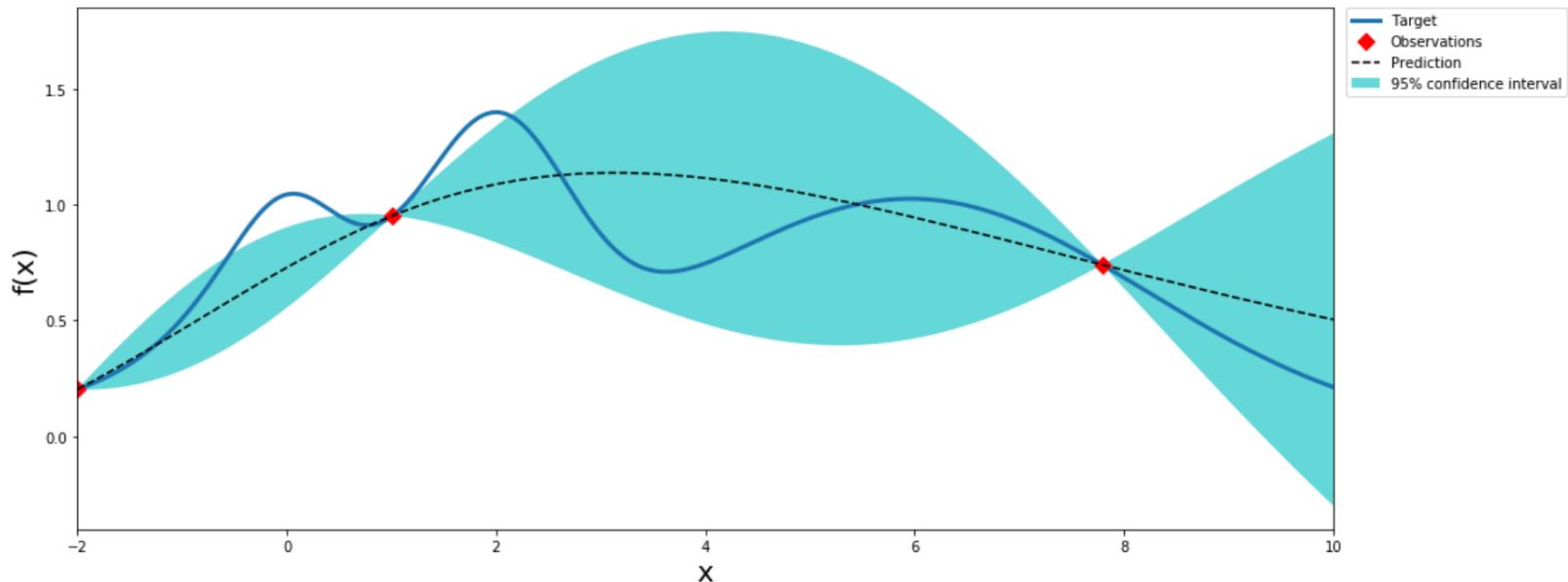
$$f(x) = e^{-(x-2)^2} + e^{-\frac{(x-6)^2}{10}} + \frac{1}{x^2 + 1}$$



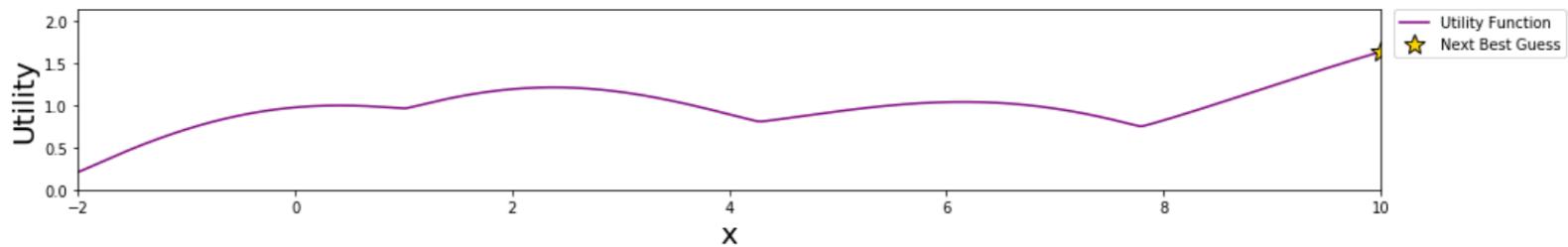
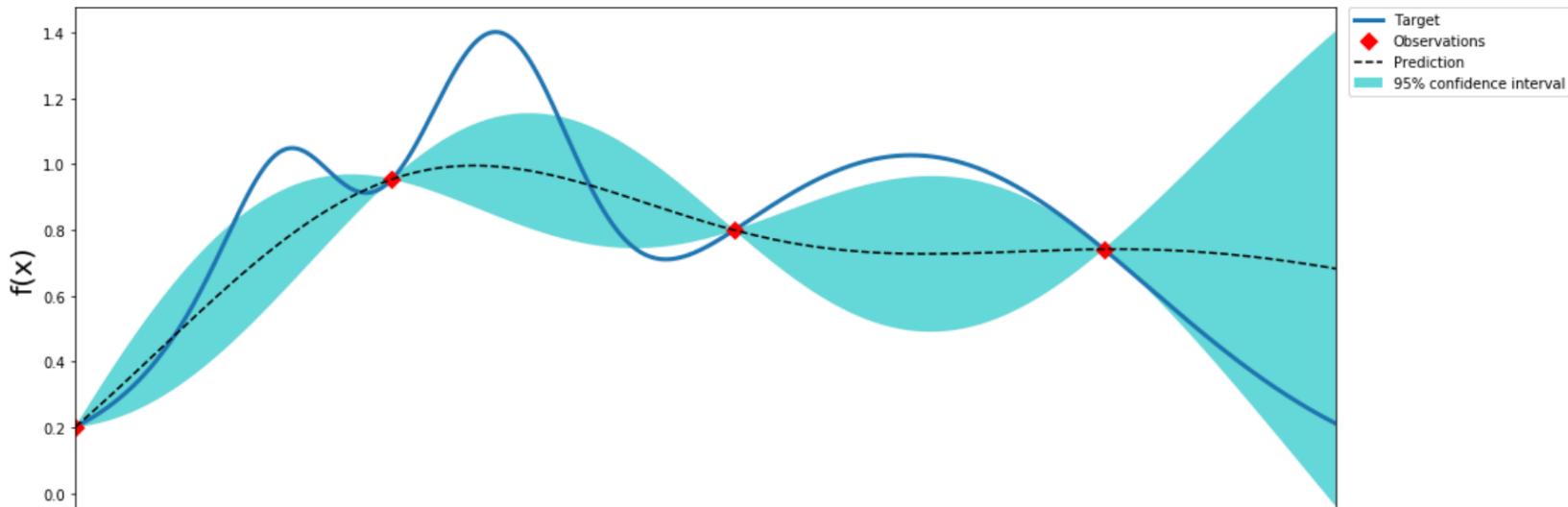
Initialize with 2 random points



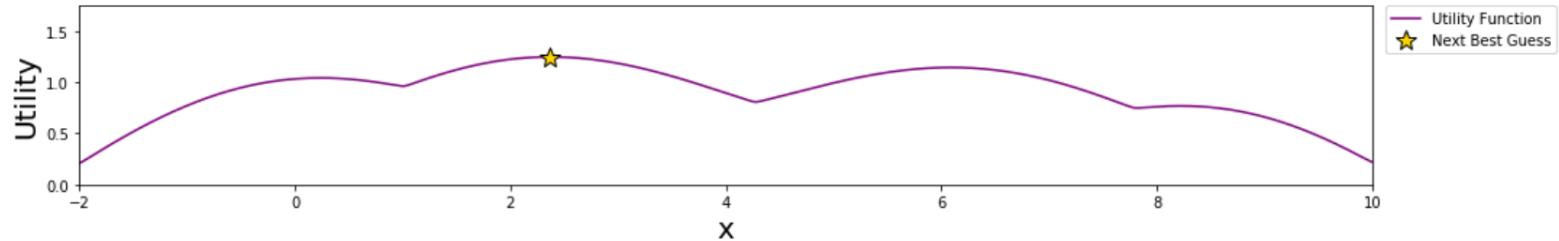
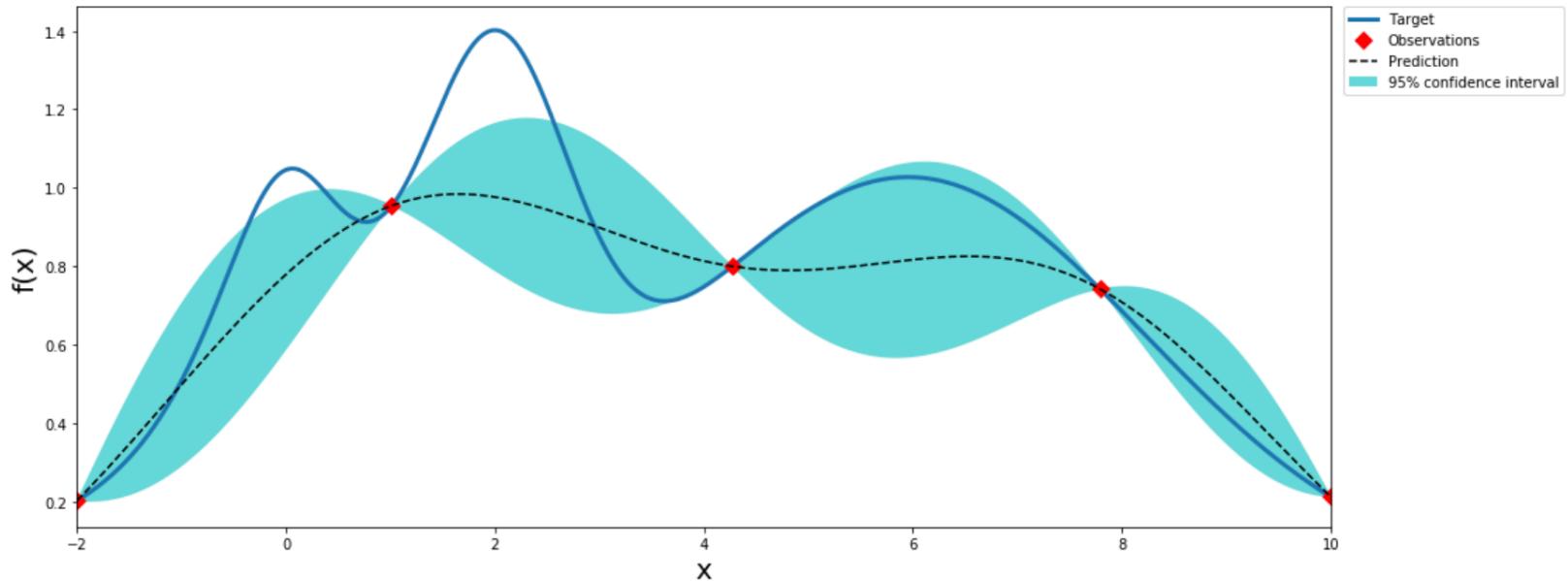
After 1 step of GP (and two random points)



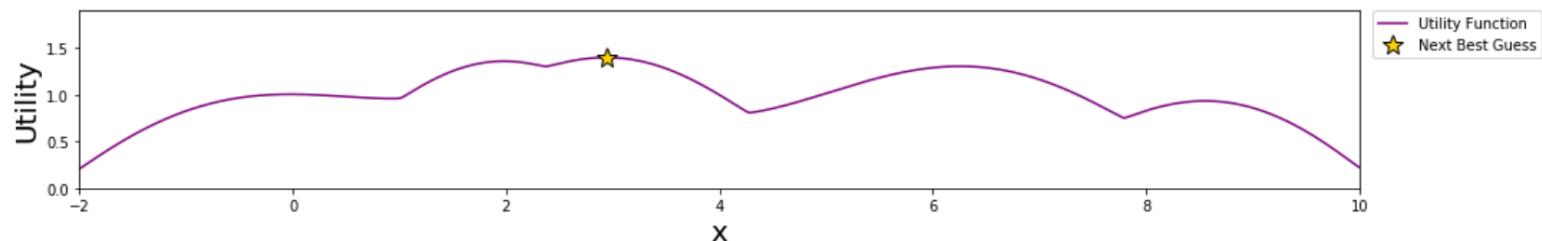
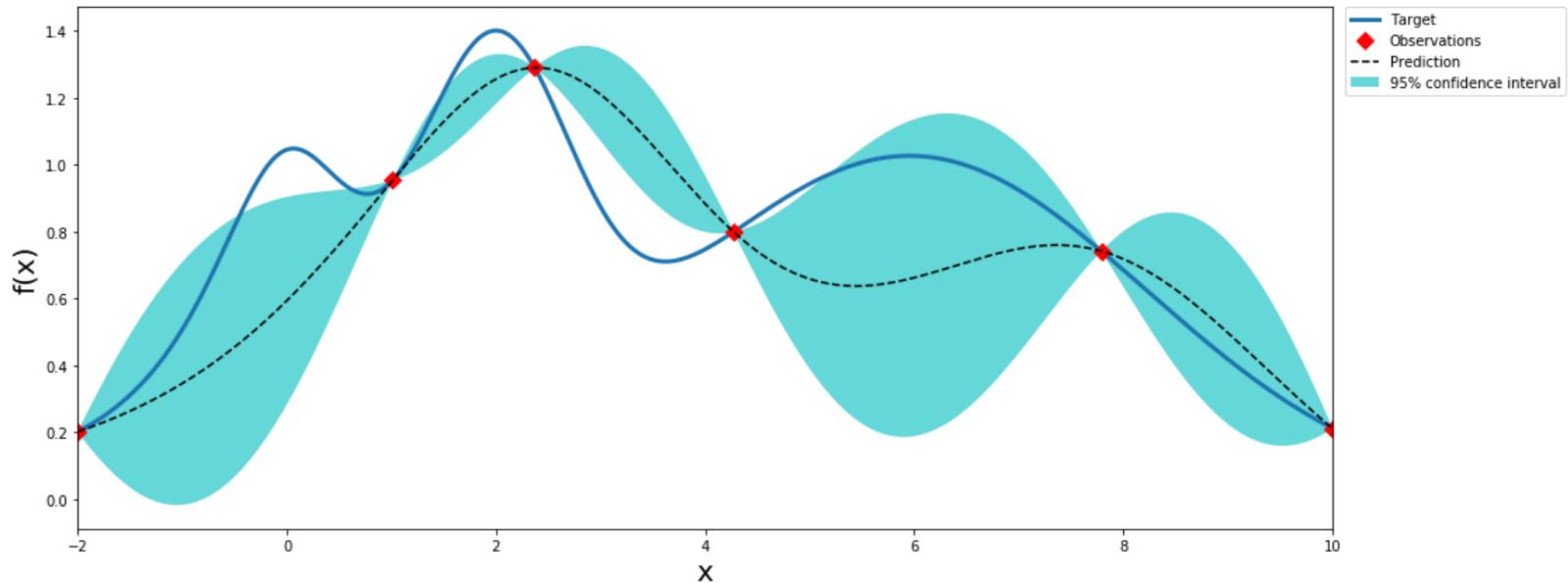
After 2 steps of GP (and two random points)



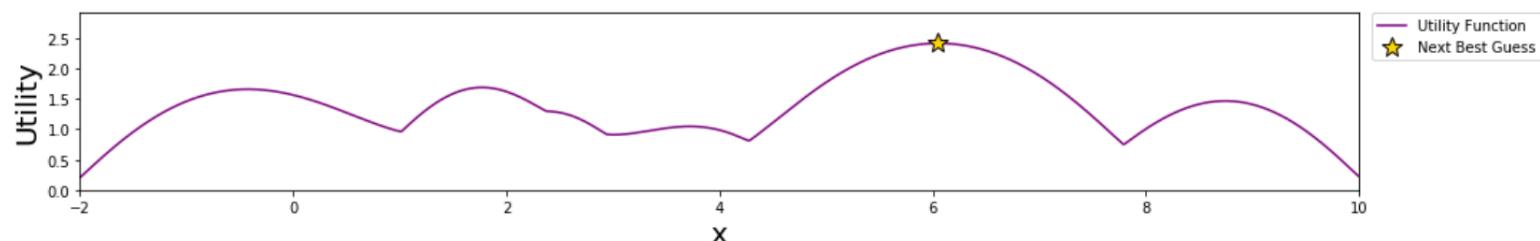
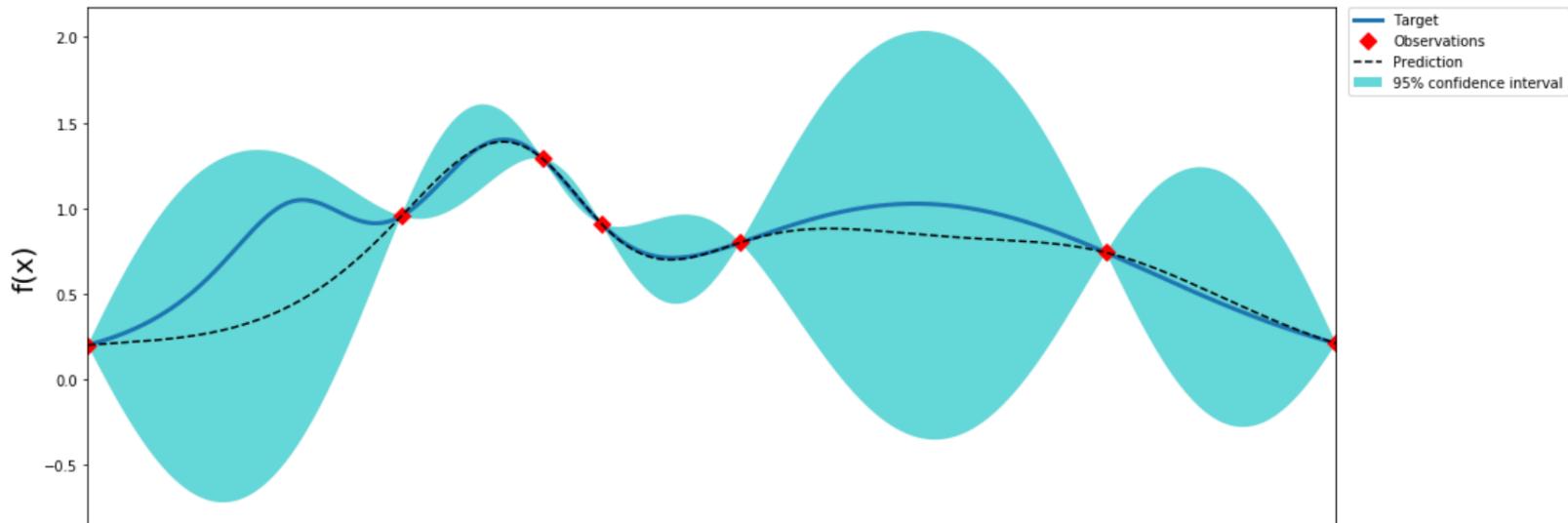
After 3 steps of GP (and two random points)



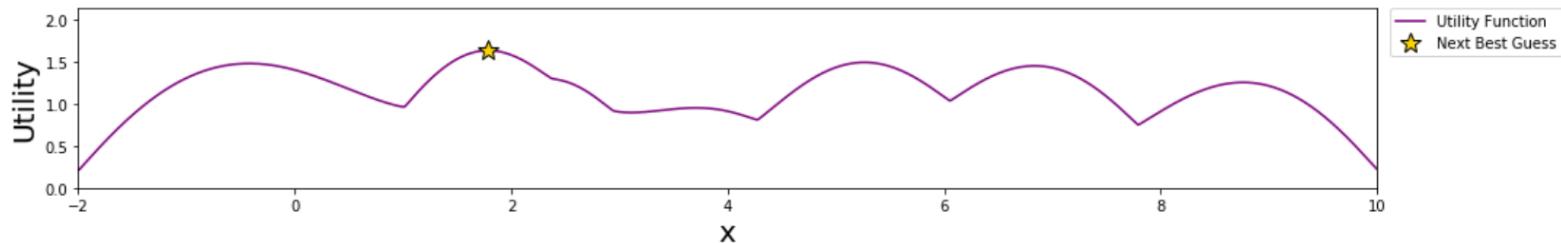
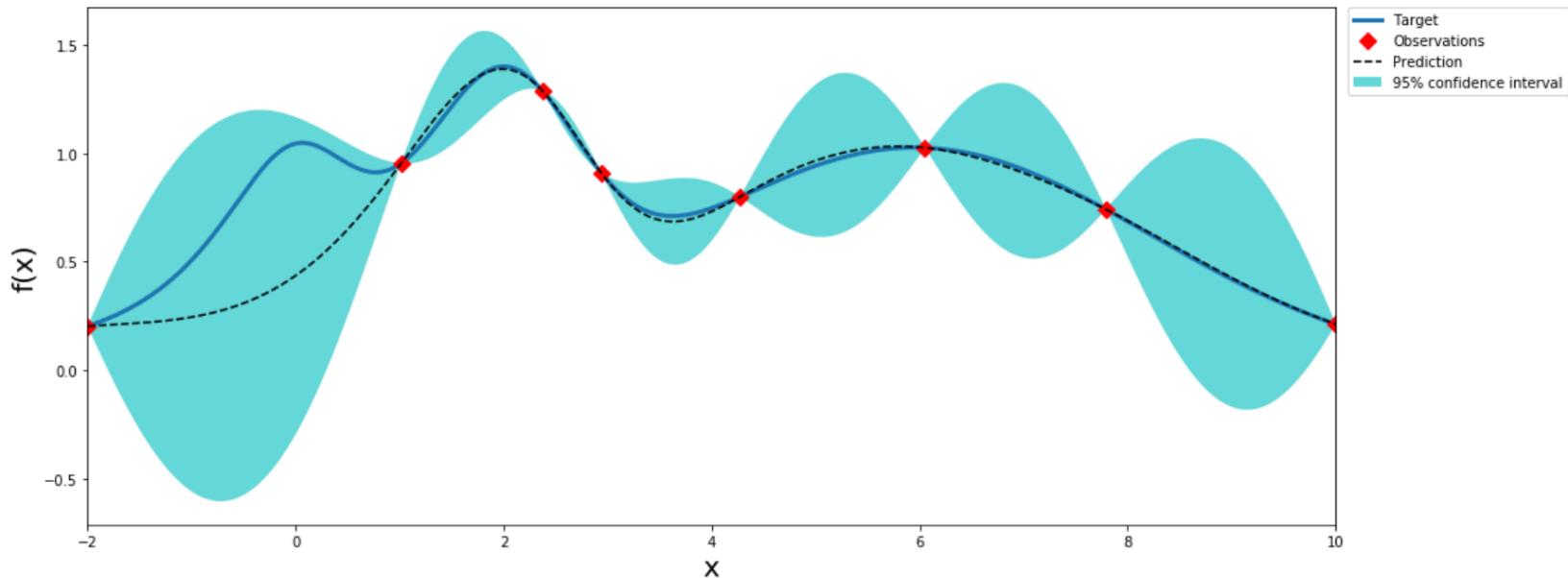
After 4 steps of GP (and two random points)



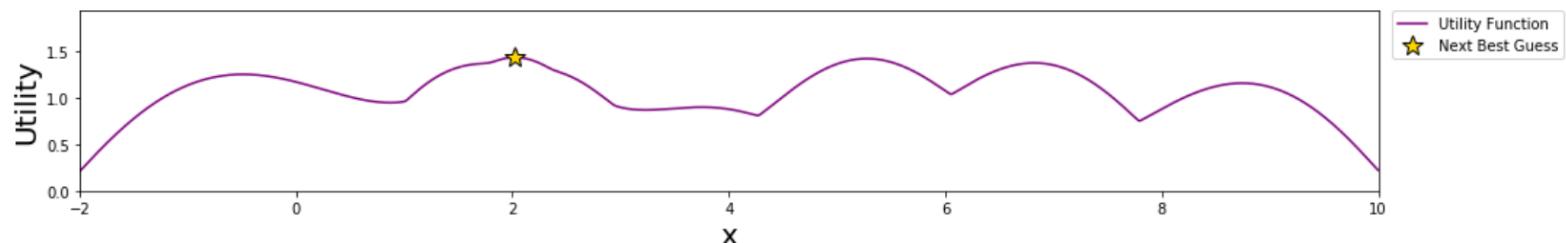
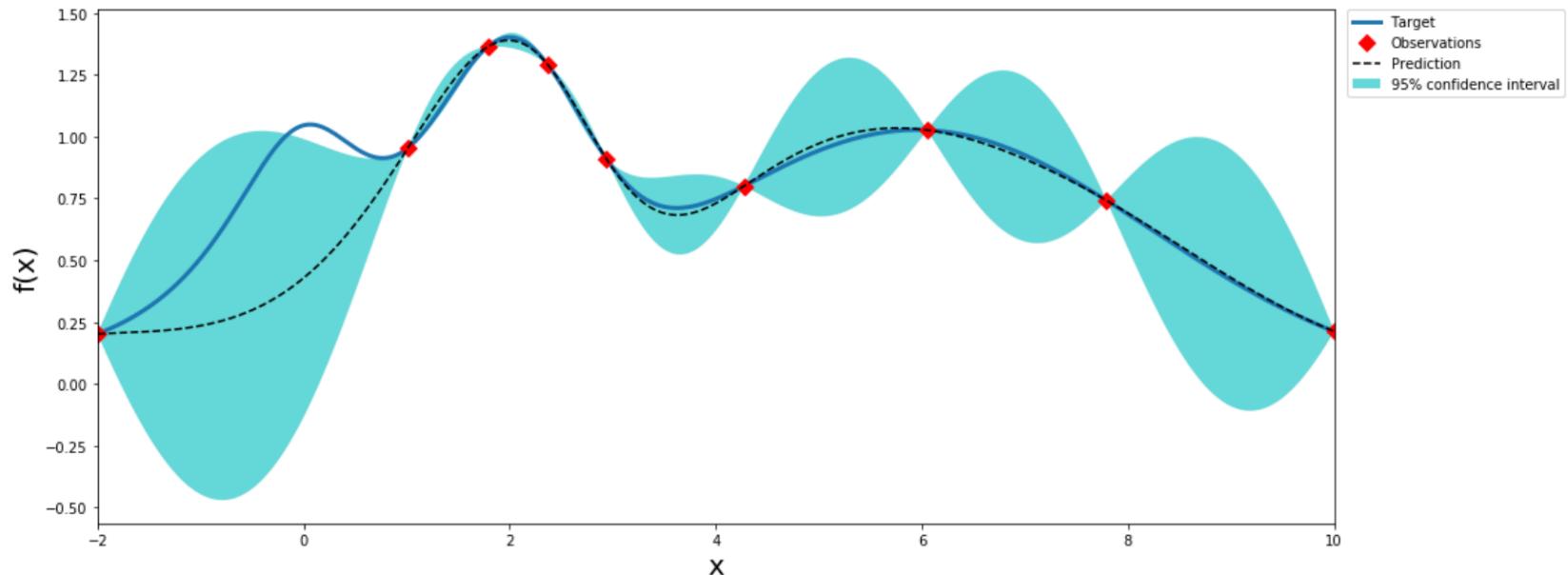
After 5 steps of GP (and two random points)



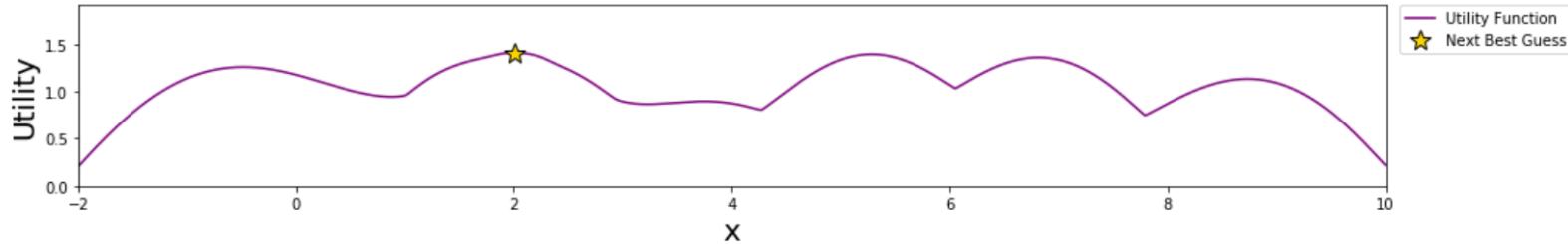
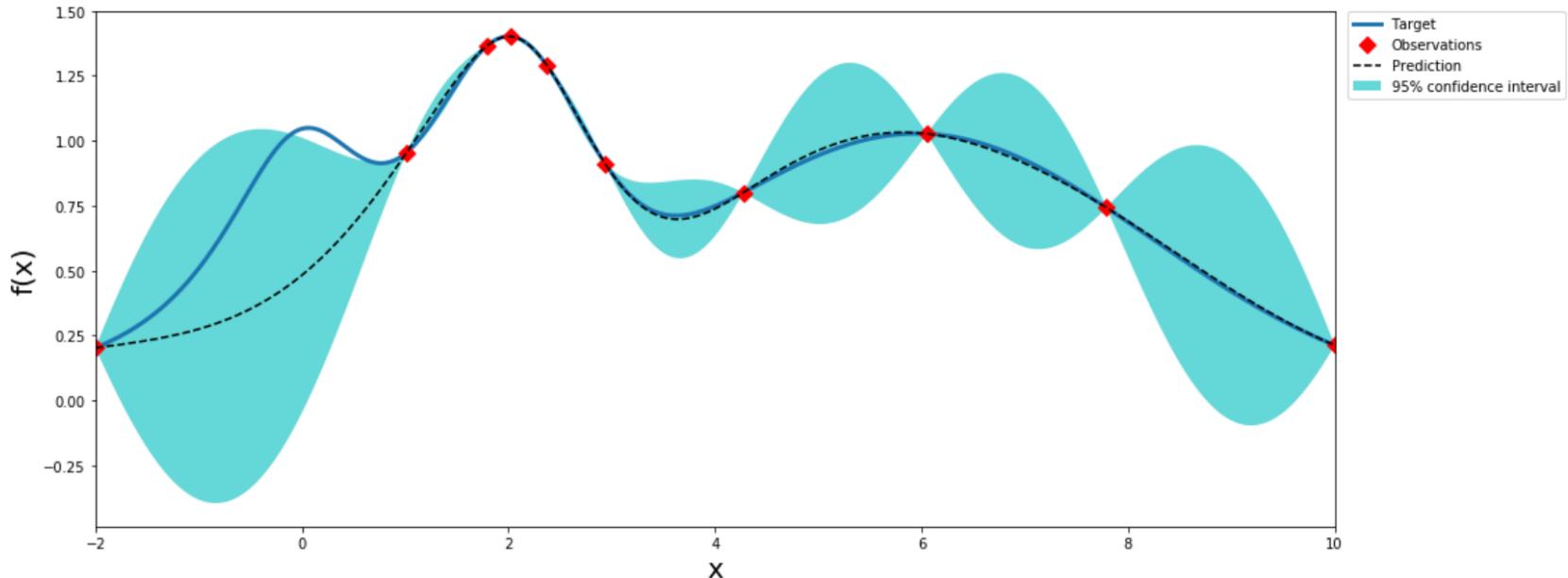
After 6 steps of GP (and two random points)



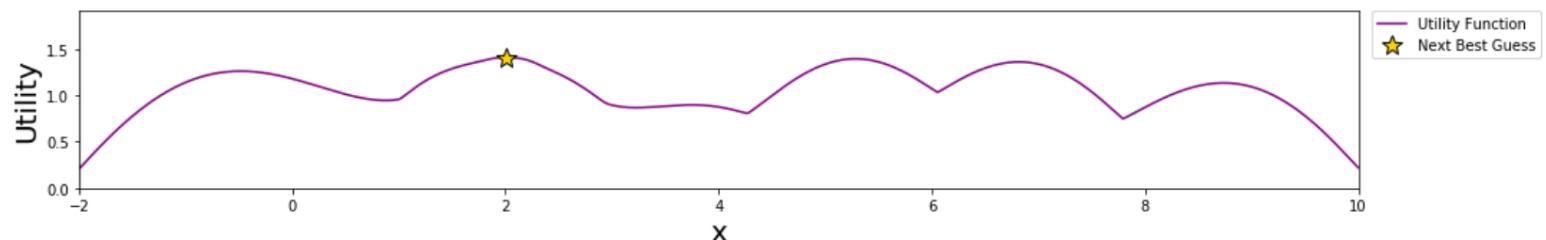
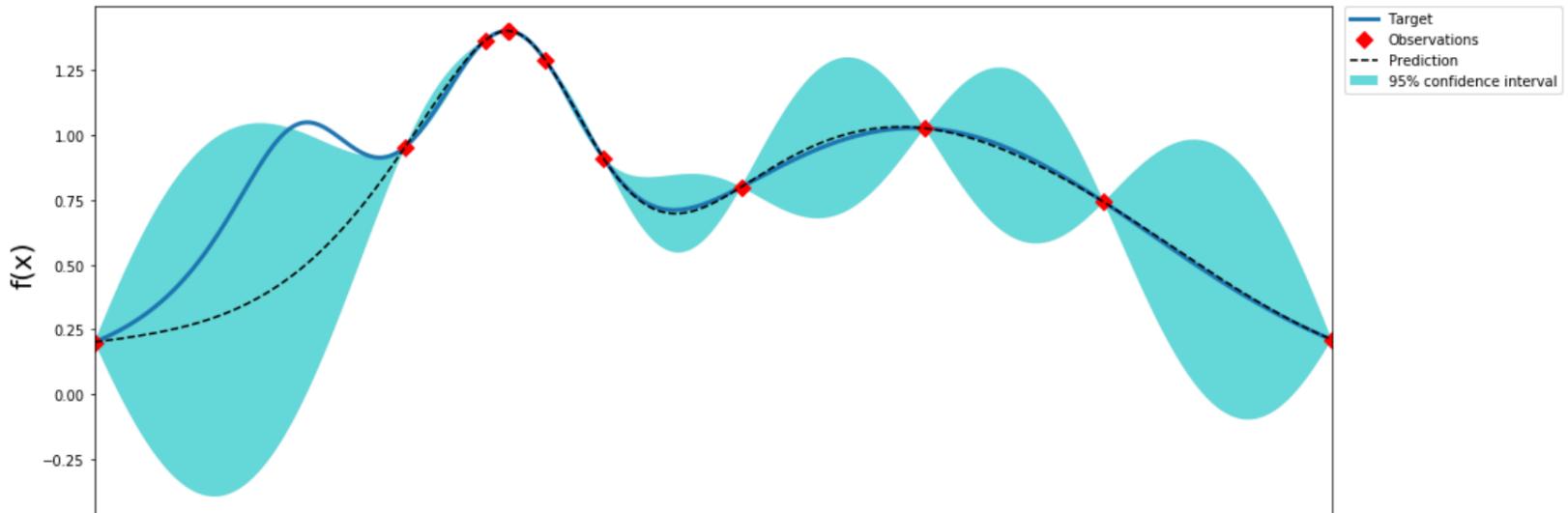
After 7 steps of GP (and two random points)



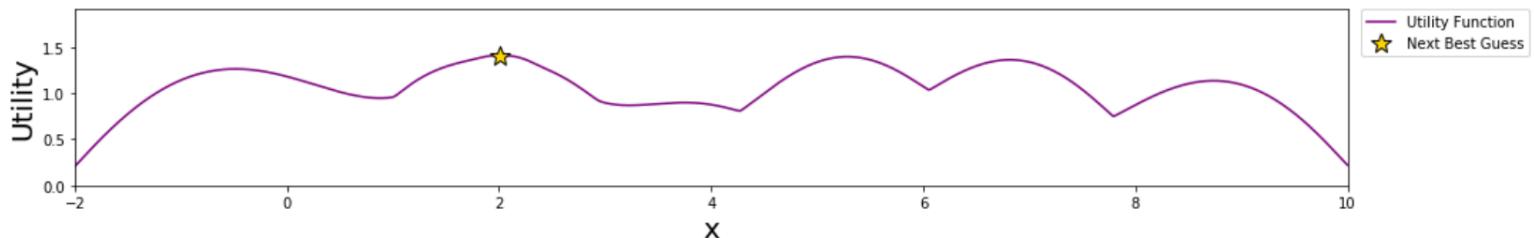
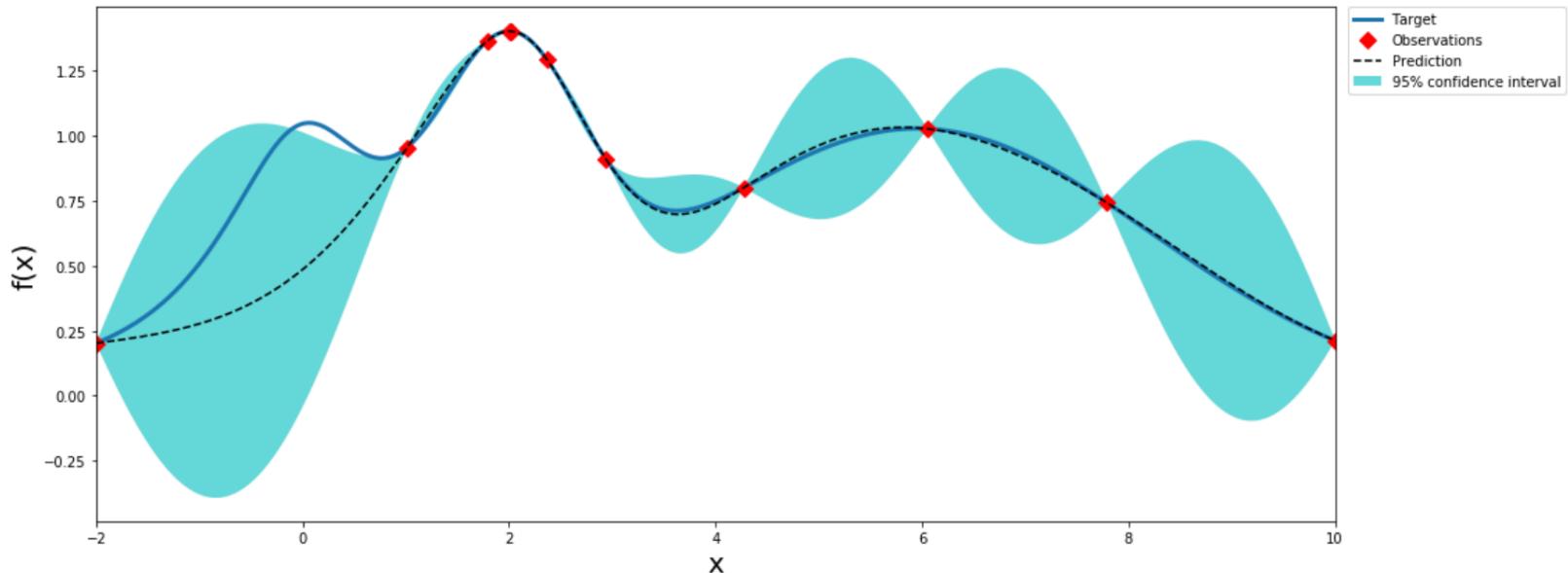
After 8 steps of GP (and two random points)



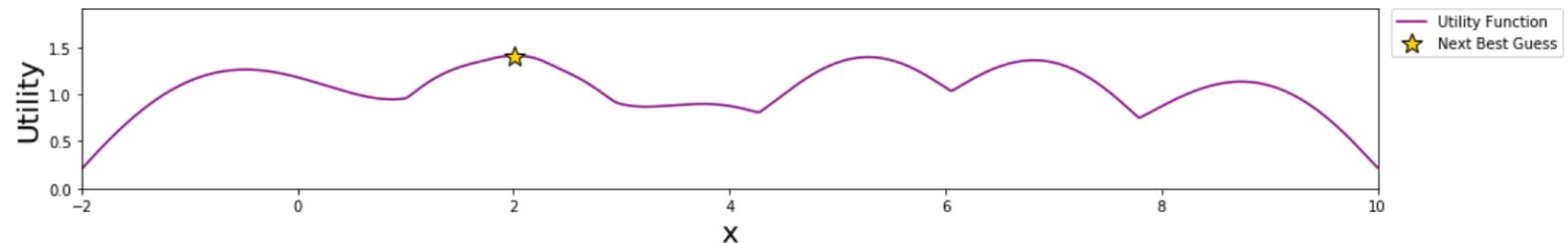
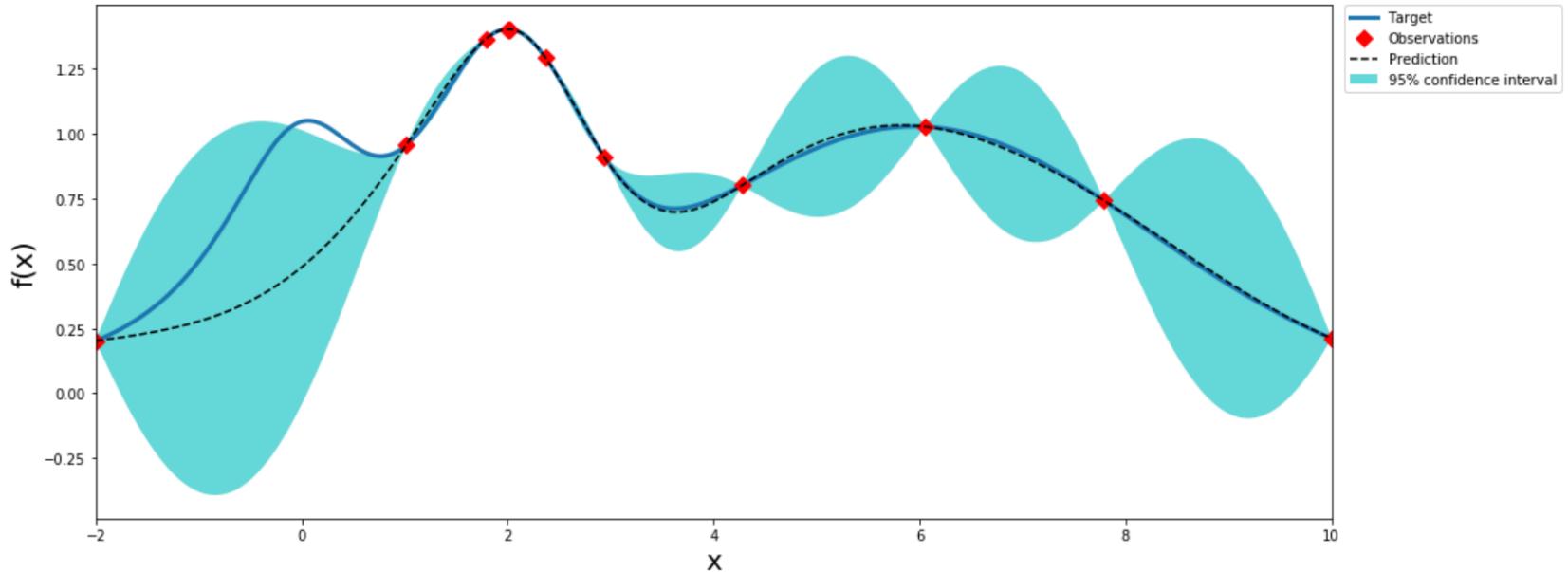
After 9 steps of GP (and two random points)



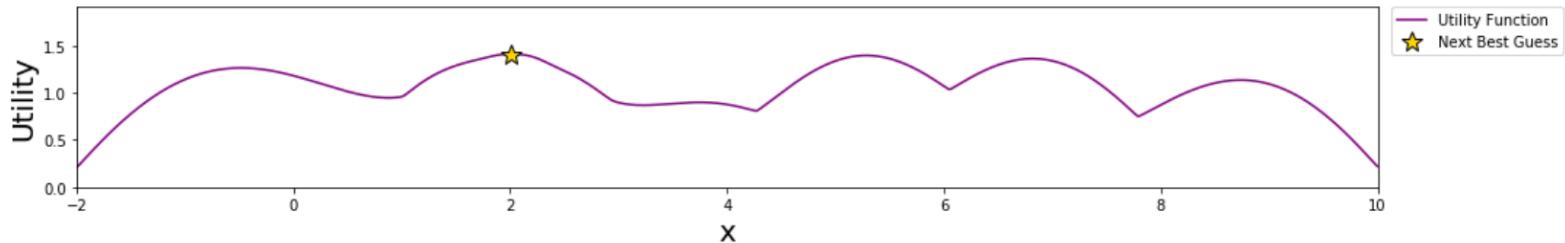
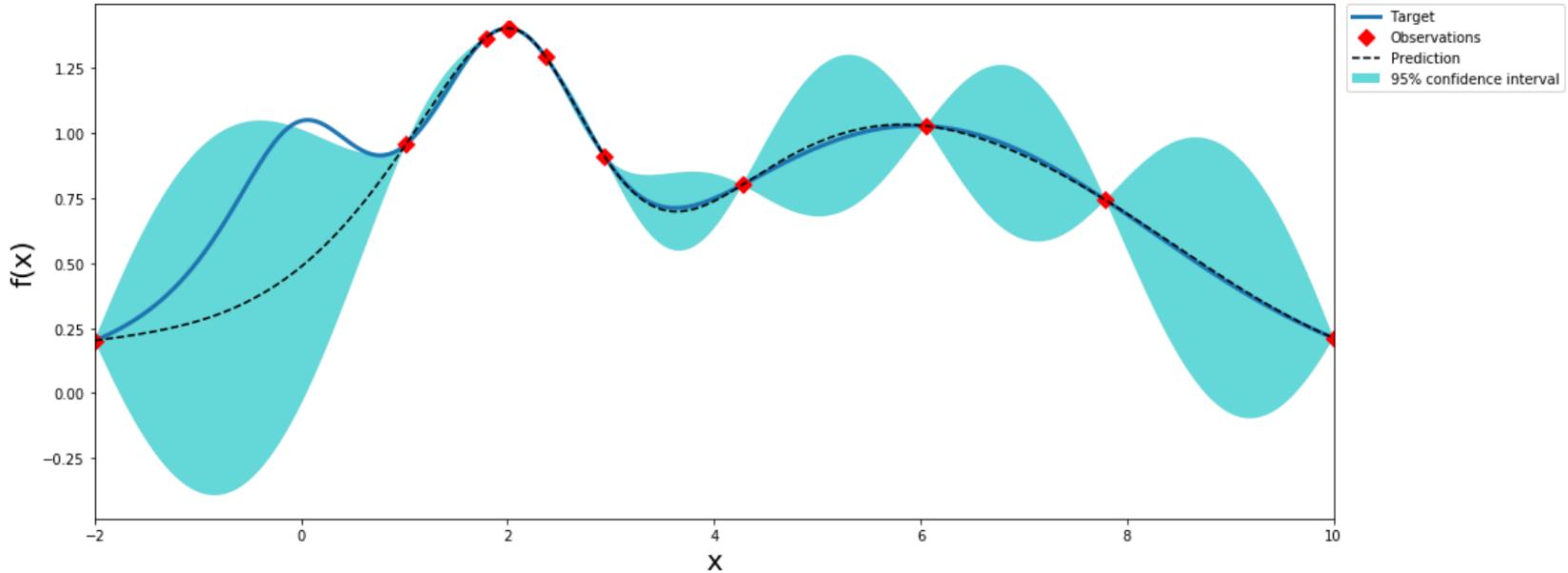
After 10 steps of GP (and two random points)



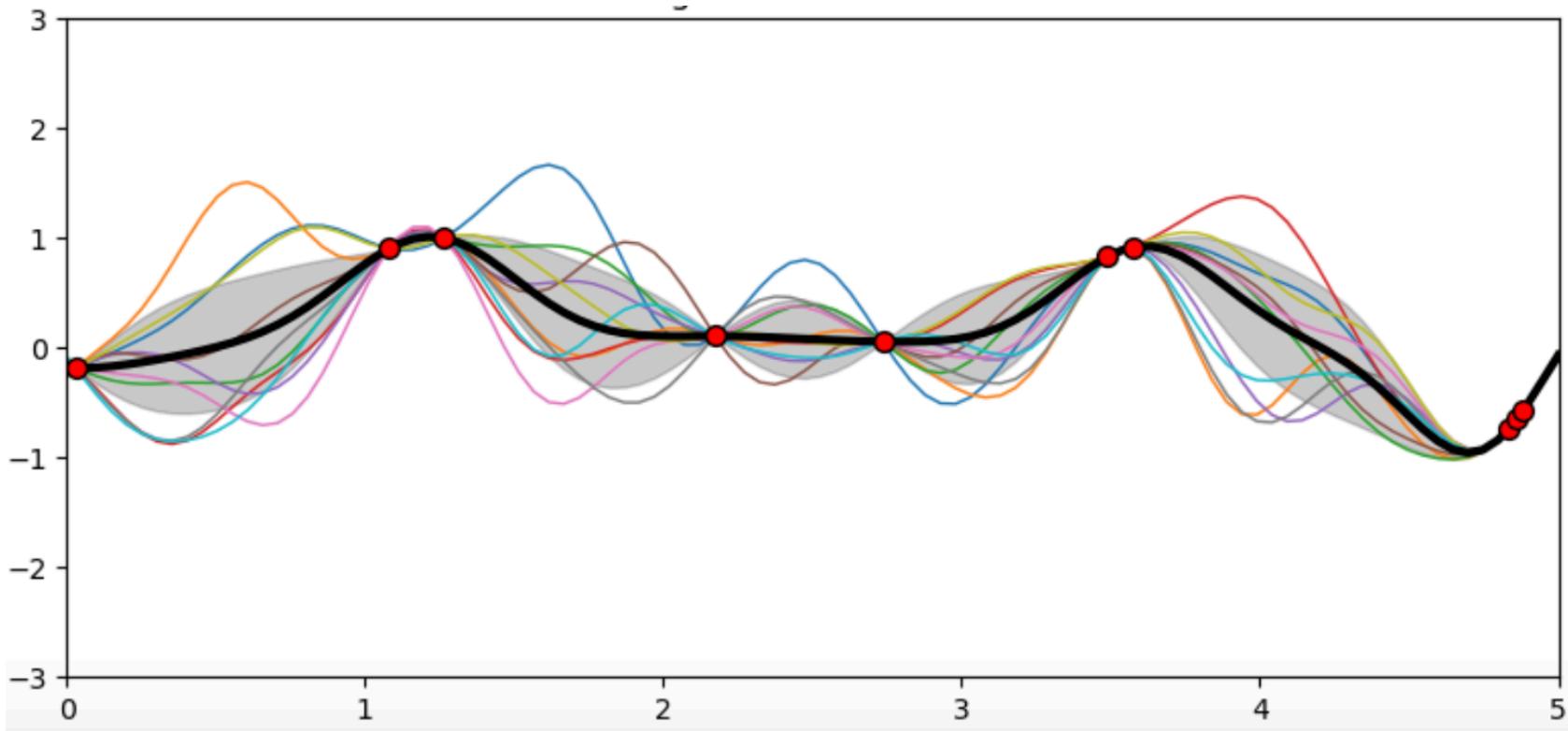
After 11 steps of GP (and two random points)



After 12 steps of GP (and two random points)



Functional Perspective



Fashion MNIST Demo

Fashion MNIST Dataset

- 60,000 training examples
- 10,000 test examples
- Each example is a 28x28 grayscale image
- 10 label classes
- Intended replacement for the original MNIST dataset



Label	Description
0	T-shirt/top
1	Trouser
2	Pullover
3	Dress
4	Coat
5	Sandal
6	Shirt
7	Sneaker
8	Bag
9	Ankle boot

<https://github.com/zalandoresearch/fashion-mnist>

Demo

https://nbviewer.jupyter.org/github/niallturb/pylondinium18/blob/master/fashion_mnist_bayes_opt.ipynb

References

Snoek et al., 2012 – Practical Bayesian Optimization of Machine Learning Algorithms

<http://papers.nips.cc/paper/4522-practical-bayesian-optimization-of-machine-learning-algorithms.pdf>

Bergstra & Bengio, 2012 – Random Search for Hyperparameter Optimization

<http://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf>

Rasmussen & Williams, 2006 – Gaussian Processes for Machine Learning

<http://www.gaussianprocess.org/gpml/>

François Chollet, 2017 – Deep Learning with Python

<https://www.manning.com/books/deep-learning-with-python>

Thank you

Questions?