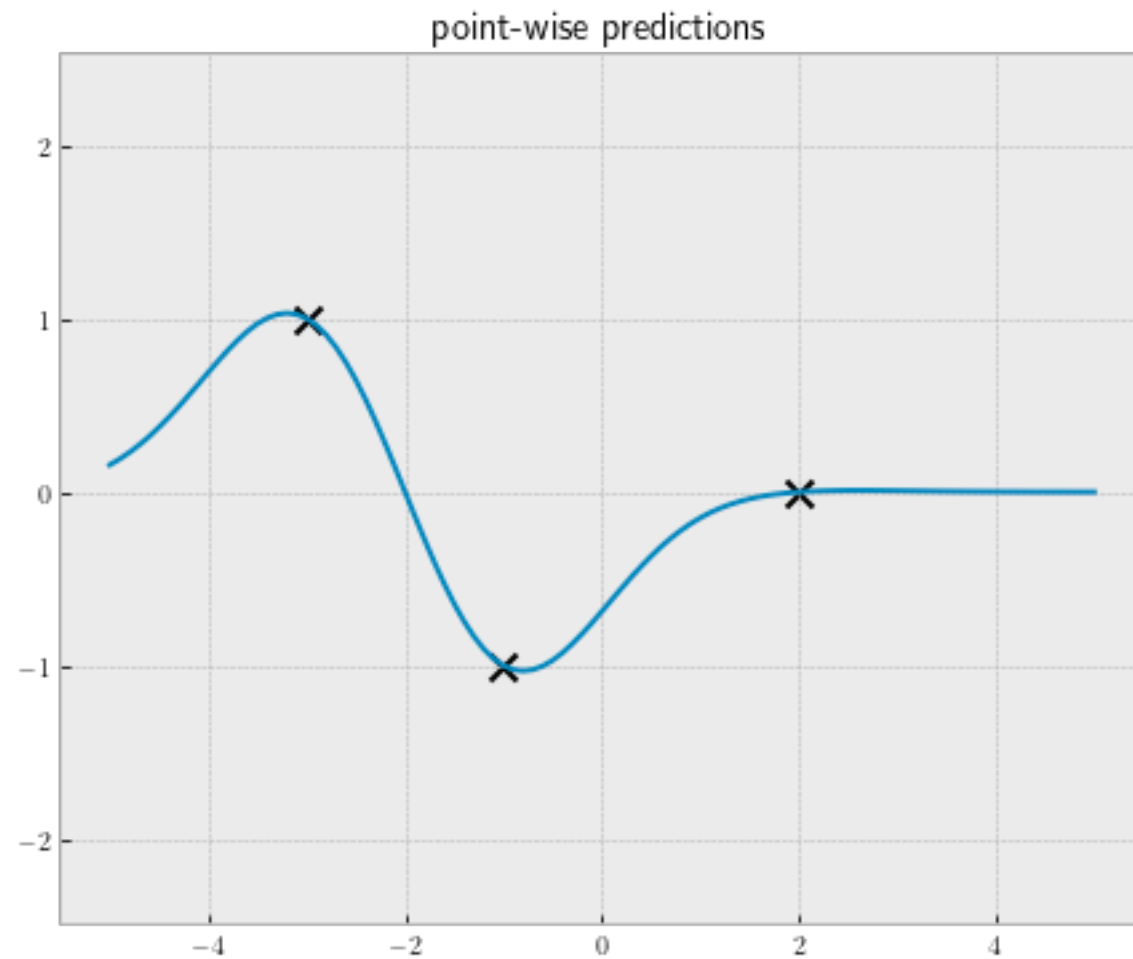


CSE 517 — MACHINE LEARNING

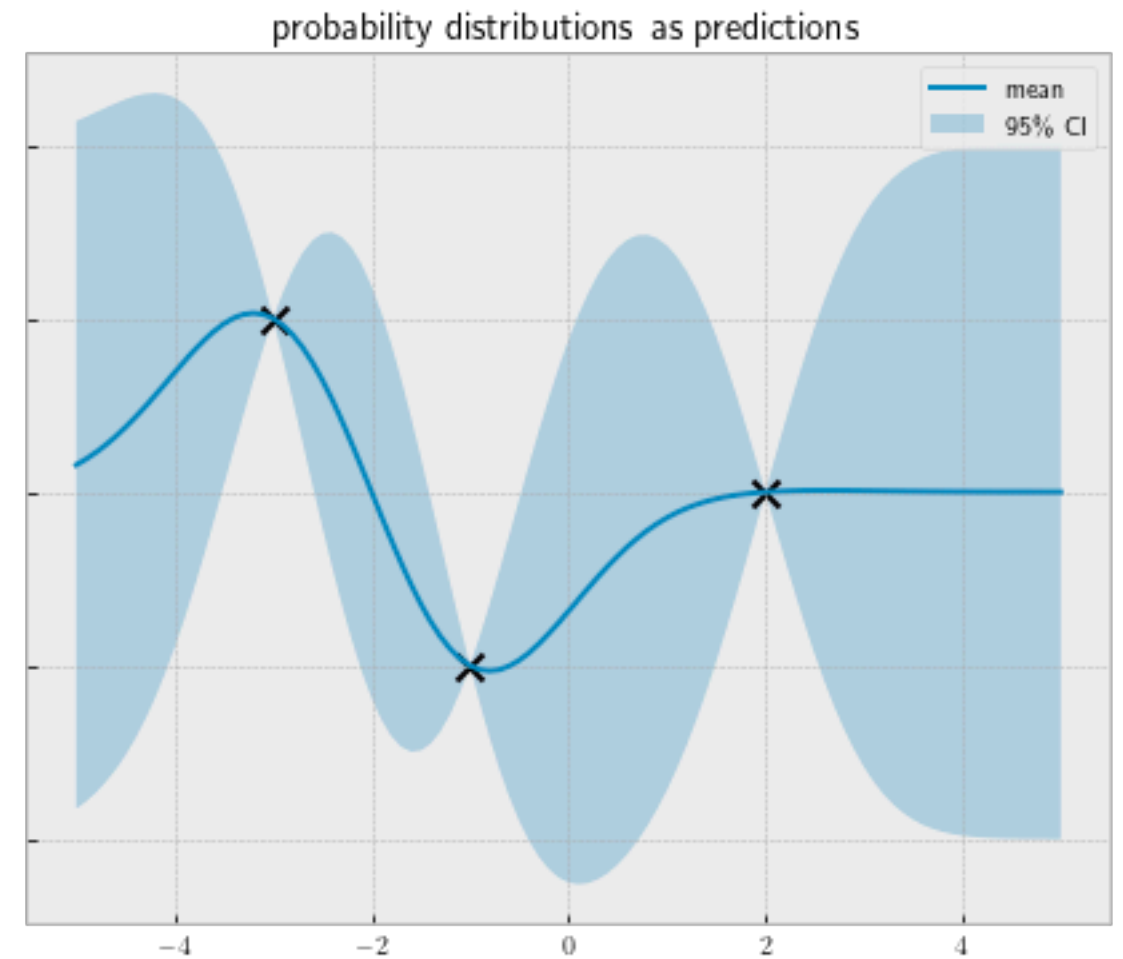
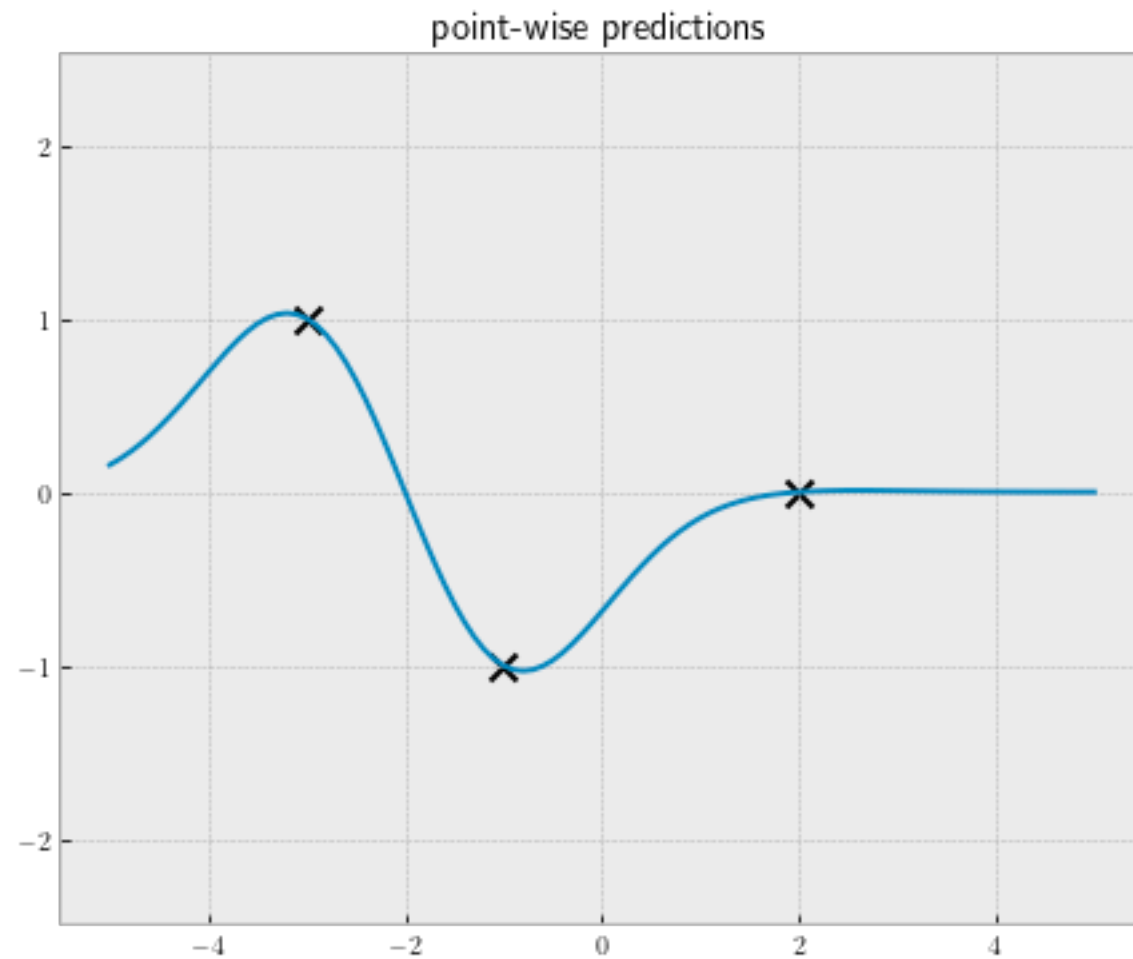
QUAN NGUYEN

GAUSSIAN PROCESSES

POINT-ESTIMATES VS. PREDICTIVE DISTRIBUTIONS



POINT-ESTIMATES VS. PREDICTIVE DISTRIBUTIONS



BEING GAUSSIAN IS VERY CONVENIENT...

BEING GAUSSIAN IS VERY CONVENIENT...

Consider multivariate normal random variable $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$ with

mean $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and covariance matrix $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$

BEING GAUSSIAN IS VERY CONVENIENT...

Consider multivariate normal random variable $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$ with

mean $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and covariance matrix $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$

$$p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Marginalization: $p(y_1) = \int_{y_2} p(y_1, y_2) \, dy_2 = N(\mu_1, \Sigma_{11})$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Marginalization: $p(y_1) = \int_{y_2} p(y_1, y_2) \, dy_2 = N(\mu_1, \Sigma_{11})$

- ▶ When we don't care about the nuisance parameter y_2

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Summation: $p(y_1) = N(\mu_1, \Sigma_{11})$

$$p(y_2) = N(\mu_2, \Sigma_{22})$$

$$p(y_1 + y_2) = N(\mu_1 + \mu_2, \Sigma_{11} + \Sigma_{22})$$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Conditioning: $p(y_1 \mid y_2) = N(\mu'_1, \Sigma'_{11})$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Conditioning: $p(y_1 \mid y_2) = N(\mu'_1, \Sigma'_{11})$

$$\mu'_1 = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (y_2 - \mu_2)$$

$$\Sigma'_{11} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

...BECAUSE GAUSSIANTY IS PRESERVED IN MANY OPERATIONS

$$p(y_1, y_2) = p \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

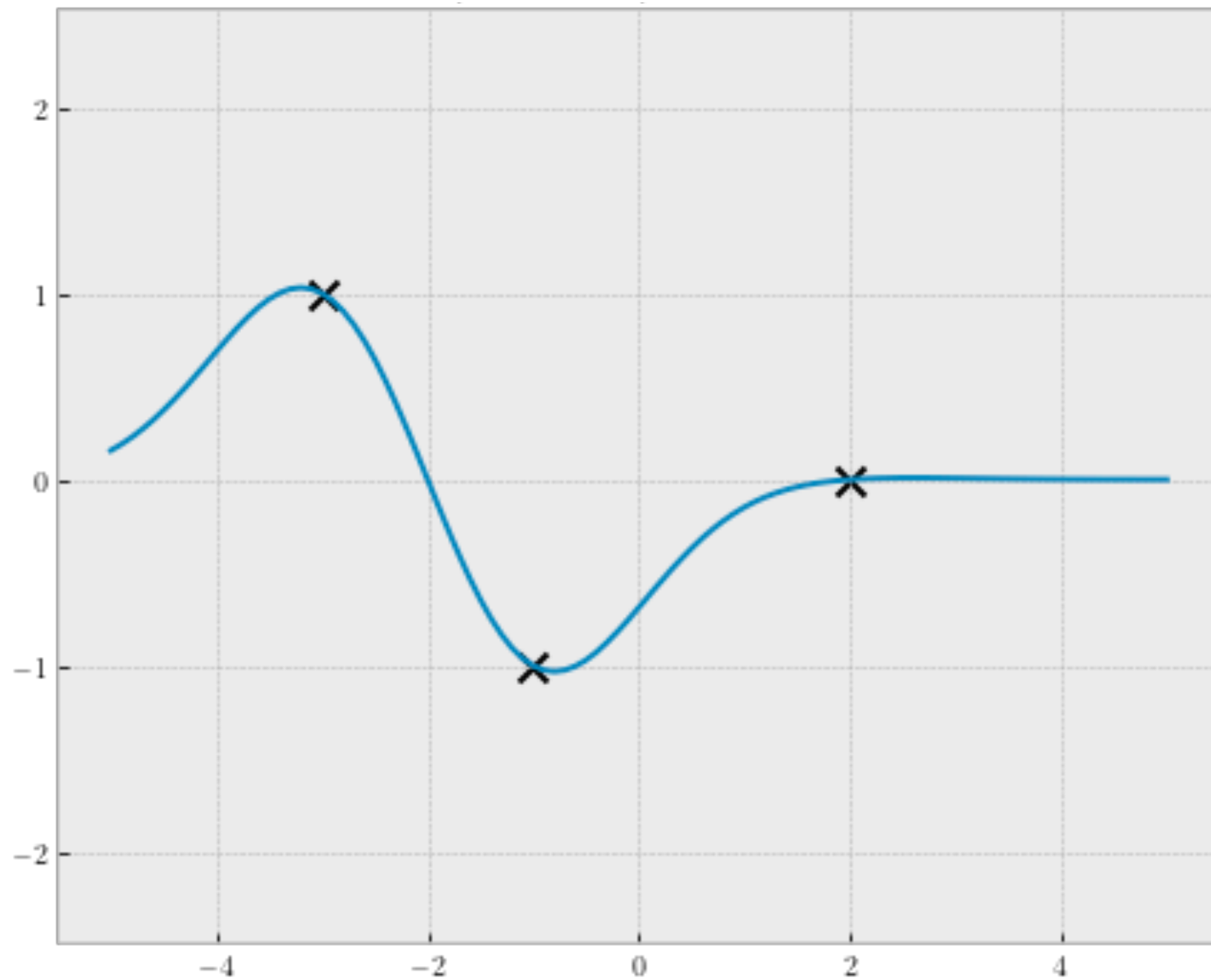
Conditioning: $p(y_1 \mid y_2) = N(\mu'_1, \Sigma'_{11})$

$$\mu'_1 = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (y_2 - \mu_2)$$

$$\Sigma'_{11} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Applies to **any** multivariate Gaussian

CONSIDER A REGRESSION TASK



THE (KERNEL) REGRESSION MODEL

THE (KERNEL) REGRESSION MODEL

$$y = f(\mathbf{x}) + \varepsilon = \mathbf{w}^\top \mathbf{x} + \varepsilon, \text{ where } \varepsilon \sim N(0, \sigma_n^2)$$

THE (KERNEL) REGRESSION MODEL

$$y = f(\mathbf{x}) + \varepsilon = \mathbf{w}^\top \mathbf{x} + \varepsilon, \text{ where } \varepsilon \sim N(0, \sigma_n^2)$$

Learning \mathbf{w} :

$$\mathbf{w}^\top \mathbf{x}_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}$$

THE (KERNEL) REGRESSION MODEL

$$y = f(\mathbf{x}) + \varepsilon = \mathbf{w}^\top \mathbf{x} + \varepsilon, \text{ where } \varepsilon \sim N(0, \sigma_n^2)$$

Learning \mathbf{w} : $\mathbf{w}^\top \mathbf{x}_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}$

What is the **prediction** for label y_* of test point \mathbf{x}_* ?

THE (KERNEL) REGRESSION MODEL

$$y = f(\mathbf{x}) + \varepsilon = \mathbf{w}^\top \mathbf{x} + \varepsilon, \text{ where } \varepsilon \sim N(0, \sigma_n^2)$$

Learning \mathbf{w} :
$$\mathbf{w}^\top \mathbf{x}_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}$$

What is the **prediction** for label y_* of test point \mathbf{x}_* ?

$$p(y_* \mid \mathbf{x}_*, \mathbf{w}) = N(\mathbf{w}^\top \mathbf{x}_*, \sigma_n^2) = N \left(K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}, \sigma_n^2 \right)$$

THE (KERNEL) REGRESSION MODEL

$$y = f(\mathbf{x}) + \varepsilon = \mathbf{w}^\top \mathbf{x} + \varepsilon, \text{ where } \varepsilon \sim N(0, \sigma_n^2)$$

Learning \mathbf{w} :
$$\mathbf{w}^\top \mathbf{x}_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}$$

What is the **prediction** for label y_* of test point \mathbf{x}_* ?

$$p(y_* \mid \mathbf{x}_*, \mathbf{w}) = N(\mathbf{w}^\top \mathbf{x}_*, \sigma_n^2) = N \left(K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}, \sigma_n^2 \right)$$

Since label y_* is what we really care about, weight vector \mathbf{w} is a **nuisance parameter**.

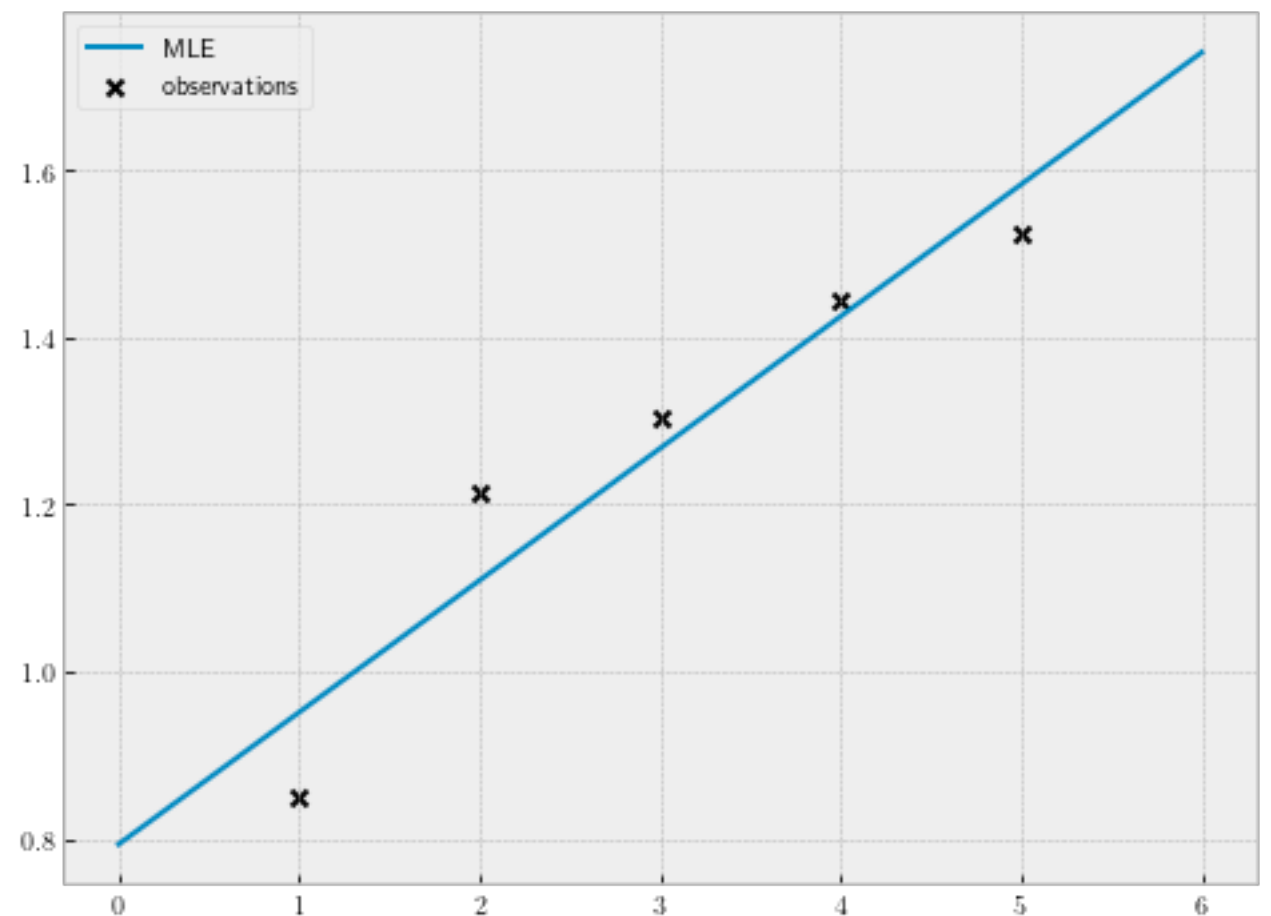
MARGINALIZING OVER WEIGHT VECTORS

MARGINALIZING OVER WEIGHT VECTORS

$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$

MARGINALIZING OVER WEIGHT VECTORS

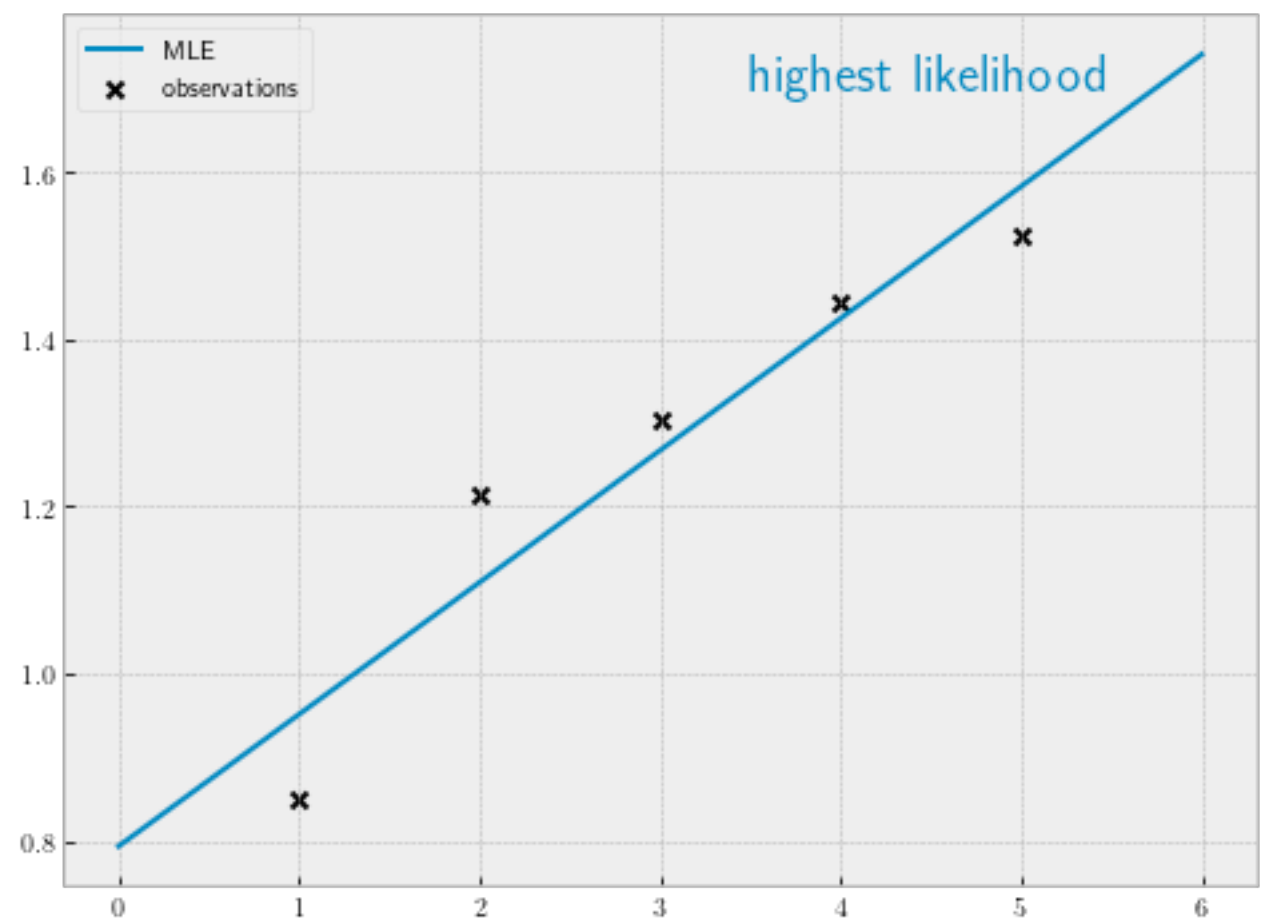
$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$



MARGINALIZING OVER WEIGHT VECTORS

$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$

MLE: $\max p(D \mid \mathbf{w})$



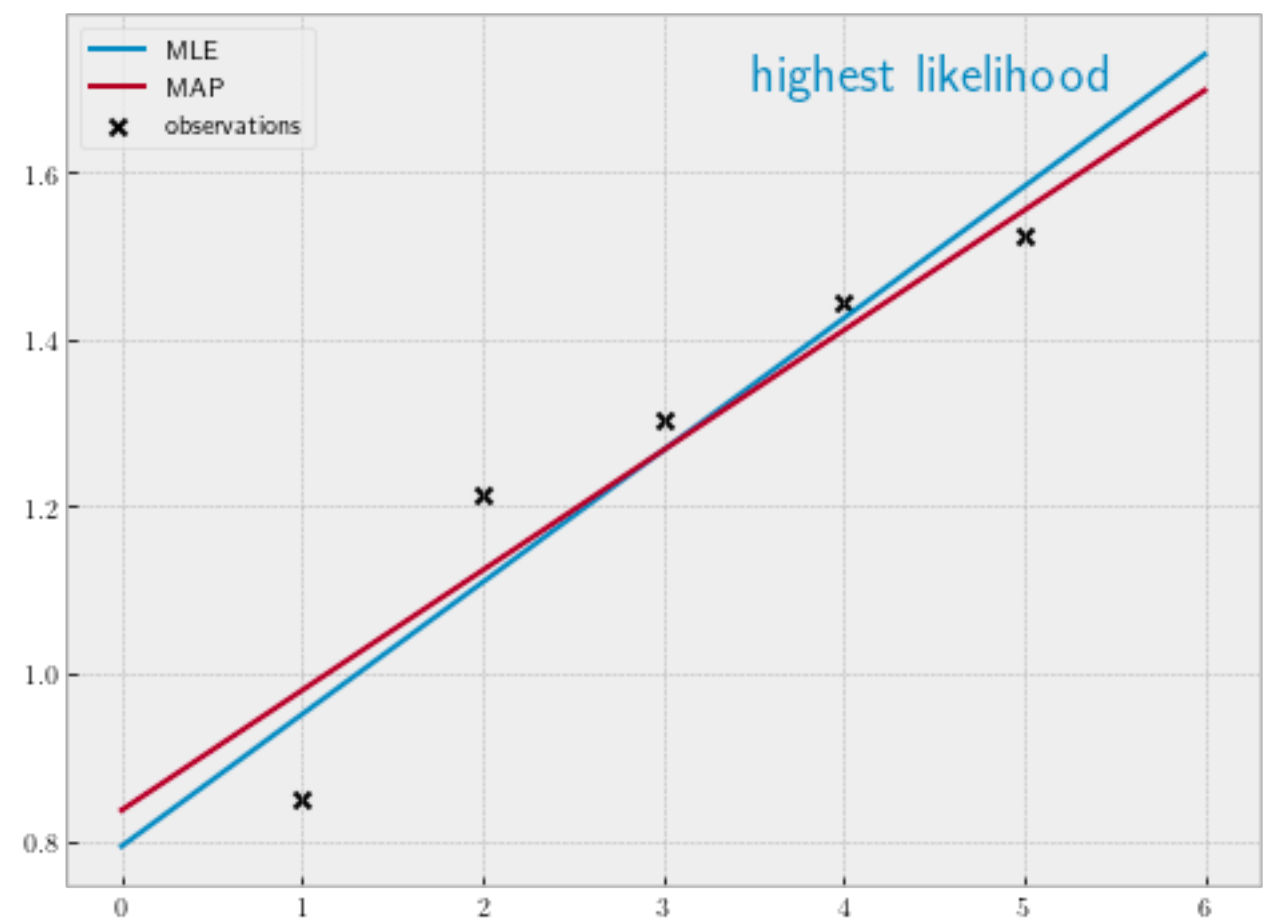
MARGINALIZING OVER WEIGHT VECTORS

$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$

MLE: $\max p(D \mid \mathbf{w})$

MAP: $\max p(\mathbf{w} \mid D)$

$$\propto p(D \mid \mathbf{w}) p(\mathbf{w})$$



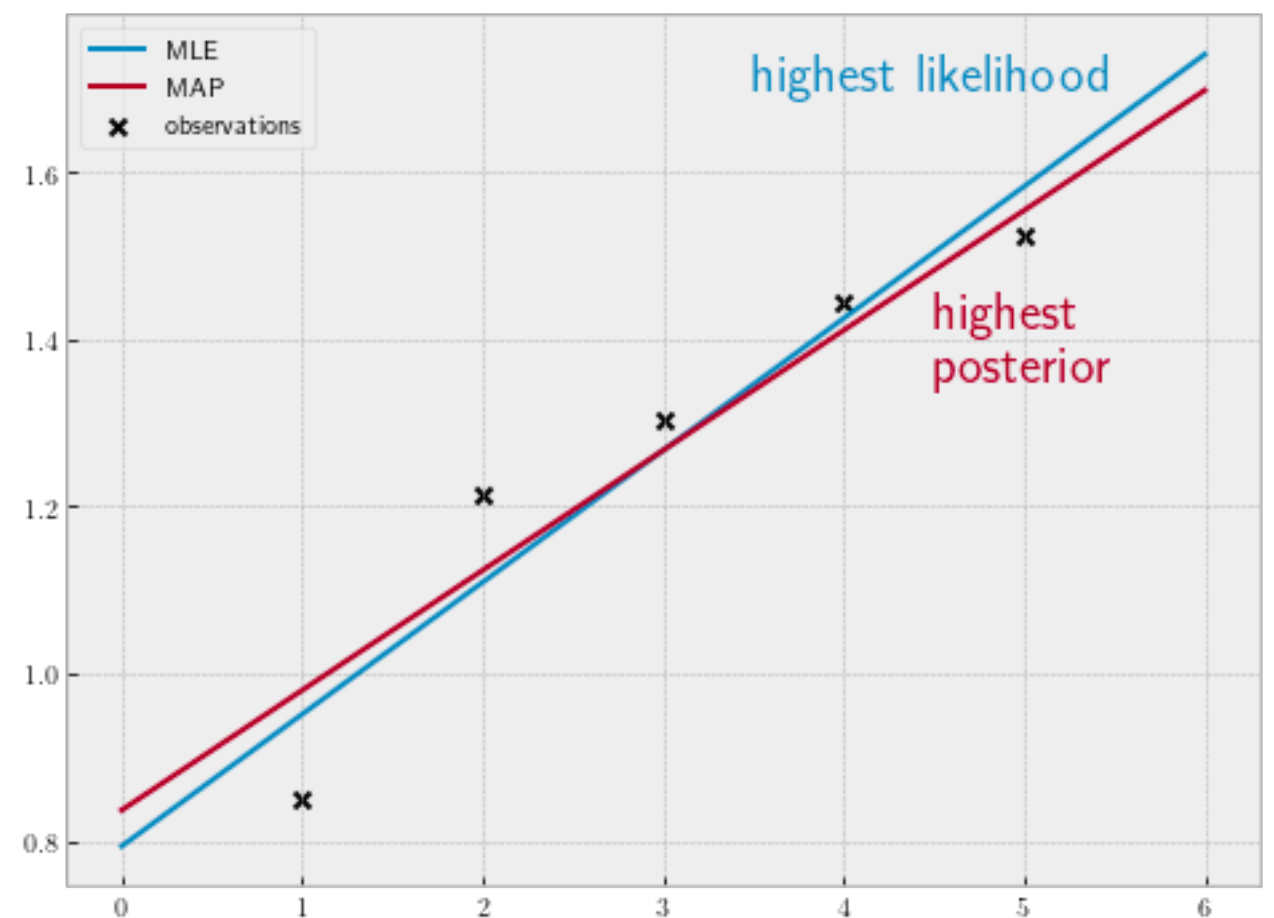
MARGINALIZING OVER WEIGHT VECTORS

$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$

MLE: $\max p(D \mid \mathbf{w})$

MAP: $\max p(\mathbf{w} \mid D)$

$$\propto p(D \mid \mathbf{w}) p(\mathbf{w})$$



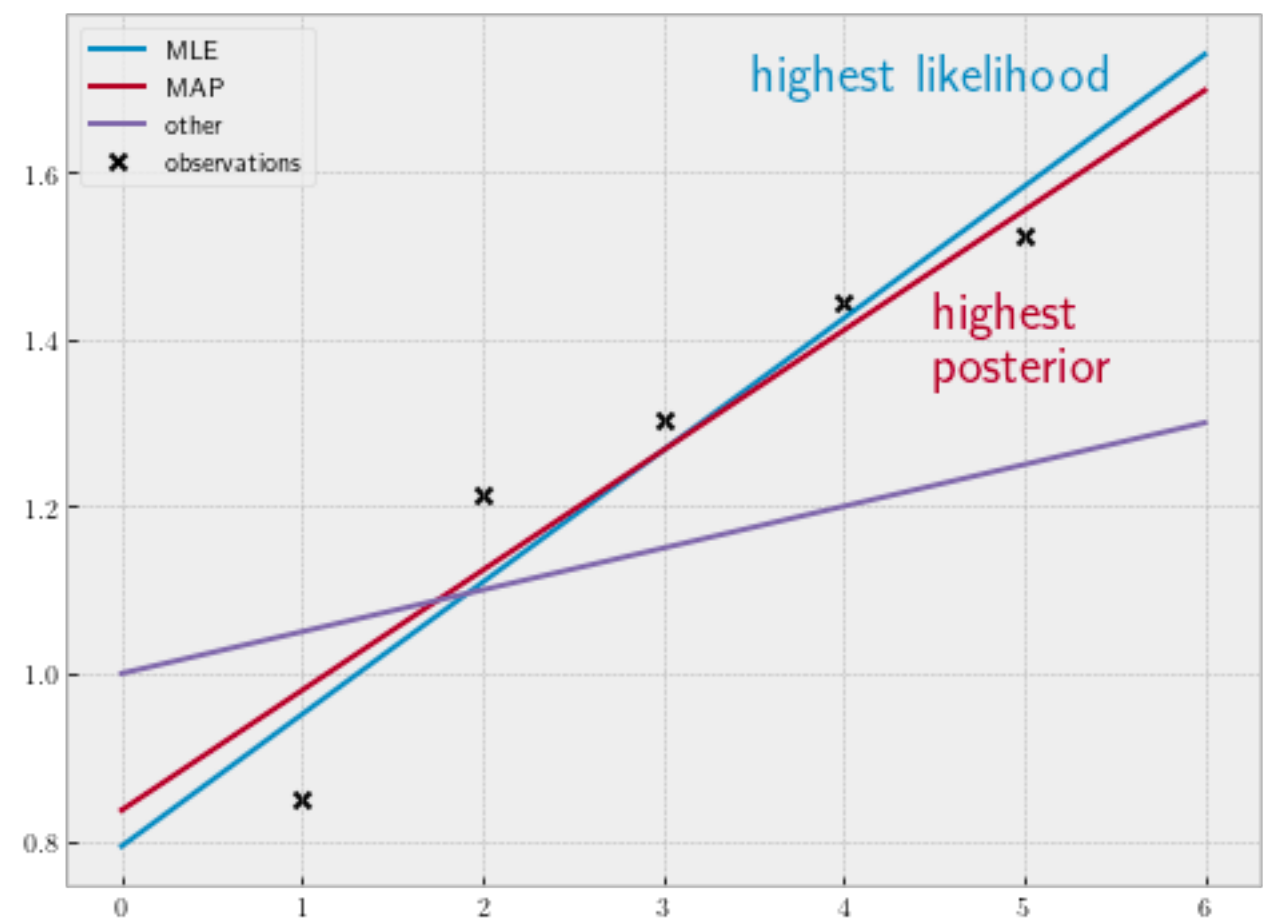
MARGINALIZING OVER WEIGHT VECTORS

$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$

MLE: $\max p(D \mid \mathbf{w})$

MAP: $\max p(\mathbf{w} \mid D)$

$$\propto p(D \mid \mathbf{w}) p(\mathbf{w})$$



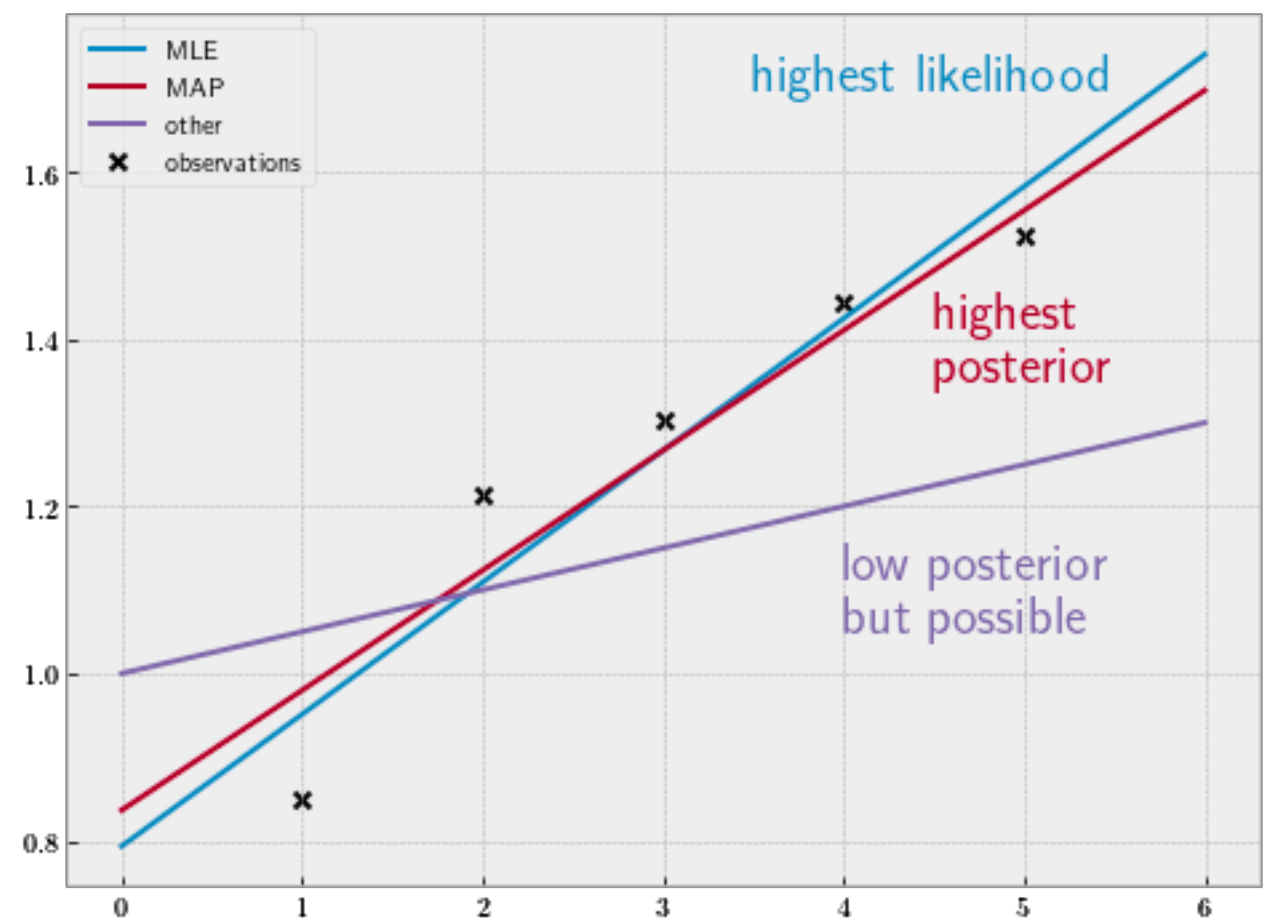
MARGINALIZING OVER WEIGHT VECTORS

$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$

MLE: $\max p(D \mid \mathbf{w})$

MAP: $\max p(\mathbf{w} \mid D)$

$$\propto p(D \mid \mathbf{w}) p(\mathbf{w})$$



MARGINALIZING OVER WEIGHT VECTORS

$$p(y \mid \mathbf{x}, D) = \int_{\mathbf{w}} p(y, \mathbf{w} \mid \mathbf{x}, D) \, d\mathbf{w} = \int_{\mathbf{w}} p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{w} \mid D) \, d\mathbf{w}$$

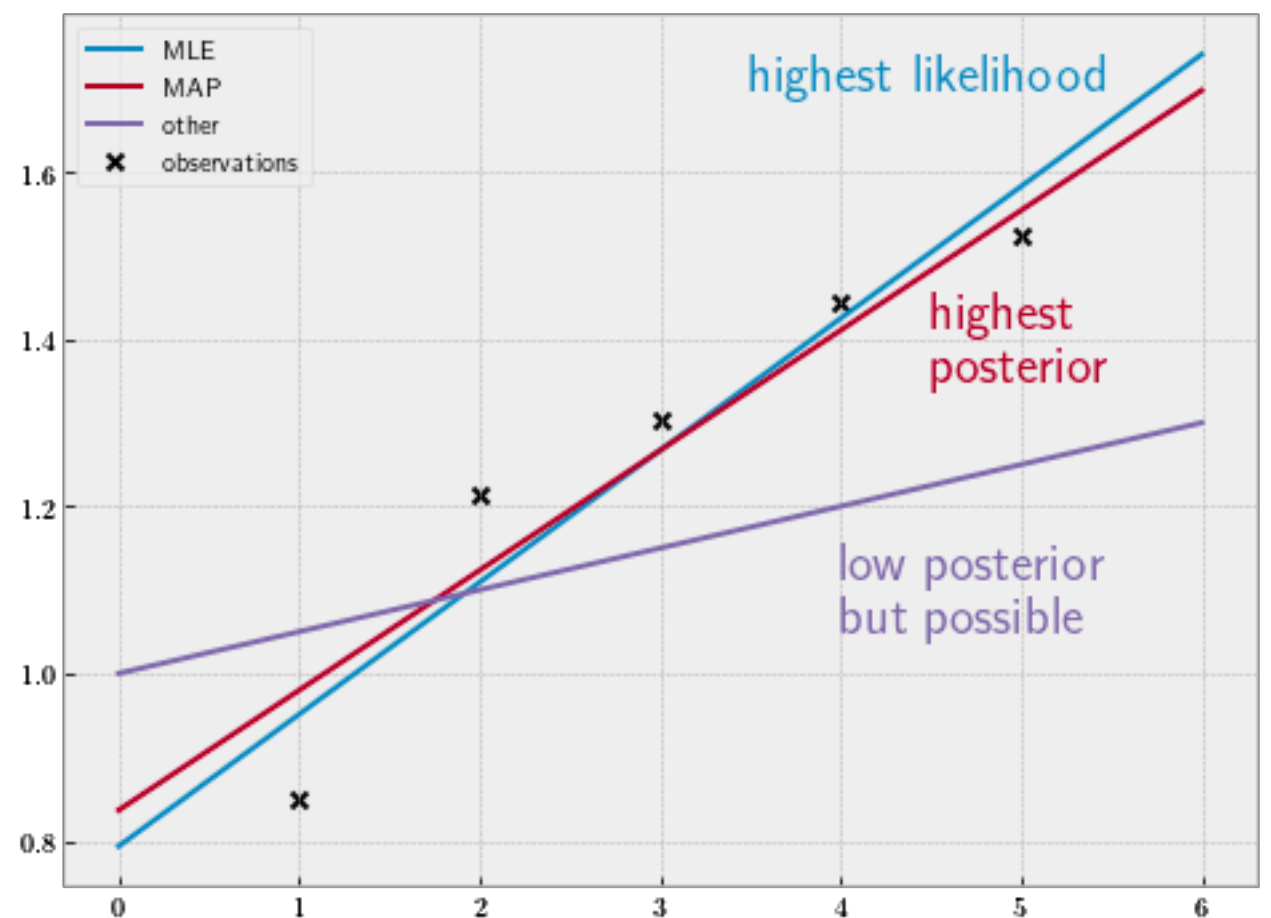
MLE: $\max p(D \mid \mathbf{w})$

MAP: $\max p(\mathbf{w} \mid D)$

$$\propto p(D \mid \mathbf{w}) p(\mathbf{w})$$

Others: low $p(D \mid \mathbf{w})$

or low $p(\mathbf{w} \mid D)$



SIDE-STEPPING THE WEIGHT VECTORS

SIDE-STEPPING THE WEIGHT VECTORS

Training data $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and test point \mathbf{x}_*

SIDE-STEPPING THE WEIGHT VECTORS

Training data $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and test point \mathbf{x}_*

$$p \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y_* \end{bmatrix} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{x}_* \right) = N(\mu, \Sigma)$$

SIDE-STEPPING THE WEIGHT VECTORS

Training data $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and test point \mathbf{x}_*

$$p \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y_* \end{bmatrix} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{x}_* \right) = N(\mu, \Sigma)$$

Conditioning on the training data D :

SIDE-STEPPING THE WEIGHT VECTORS

Training data $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and test point \mathbf{x}_*

$$p \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y_* \end{bmatrix} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{x}_* \right) = N(\mu, \Sigma)$$

Conditioning on the training data D :

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2)$$

CLOSED-FORM PREDICTIONS

CLOSED-FORM PREDICTIONS

Given training set $D = (X, y)$ and test point \mathbf{x}_* :

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

CLOSED-FORM PREDICTIONS

Given training set $D = (X, y)$ and test point \mathbf{x}_* :

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

where

$$\mu_* = K(\mathbf{x}_*, X) K^{-1}(X, X) y$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) K^{-1}(X, X) K^\top(\mathbf{x}_*, X)$$

CLOSED-FORM PREDICTIONS

Given training set $D = (X, y)$ and test point \mathbf{x}_* :

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

where

$$\mu_* = K(\mathbf{x}_*, X) K^{-1}(X, X) y$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) K^{-1}(X, X) K^\top(\mathbf{x}_*, X)$$

$K(\mathbf{x}, \mathbf{x}')$ is the covariance between \mathbf{x} and \mathbf{x}' .

- PSD.

INTERPRETING THE PREDICTIONS

INTERPRETING THE PREDICTIONS

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

$$\mu_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) K^{-1}(X, X) K^\top(\mathbf{x}_*, X)$$

INTERPRETING THE PREDICTIONS

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

$$\mu_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y} \text{ : kernel regression!!}$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) K^{-1}(X, X) K^\top(\mathbf{x}_*, X)$$

INTERPRETING THE PREDICTIONS

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

$$\mu_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}: \text{kernel regression!!}$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) K^{-1}(X, X) K^\top(\mathbf{x}_*, X)$$

- $\sigma_*^2 \leq K(\mathbf{x}_*, \mathbf{x}_*)$: uncertainty **never** increases

INTERPRETING THE PREDICTIONS

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

$$\mu_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}: \text{kernel regression!!}$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) K^{-1}(X, X) K^\top(\mathbf{x}_*, X)$$

- $\sigma_*^2 \leq K(\mathbf{x}_*, \mathbf{x}_*)$: uncertainty **never** increases
- Uncertainty quantification

INTERPRETING THE PREDICTIONS

$$p(y_* \mid \mathbf{x}_*, D) = N(\mu_*, \sigma_*^2),$$

$$\mu_* = K(\mathbf{x}_*, X) K^{-1}(X, X) \mathbf{y}: \text{kernel regression!!}$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) K^{-1}(X, X) K^\top(\mathbf{x}_*, X)$$

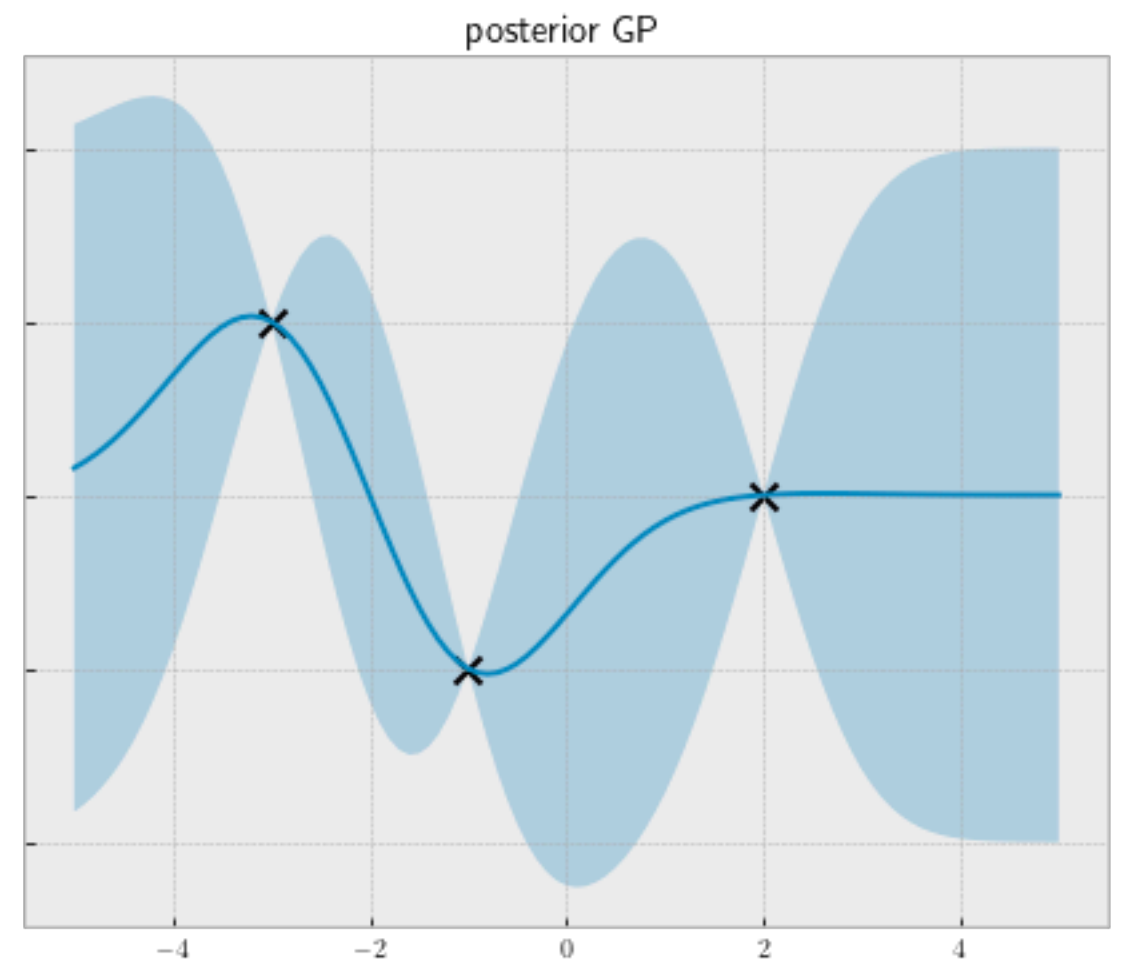
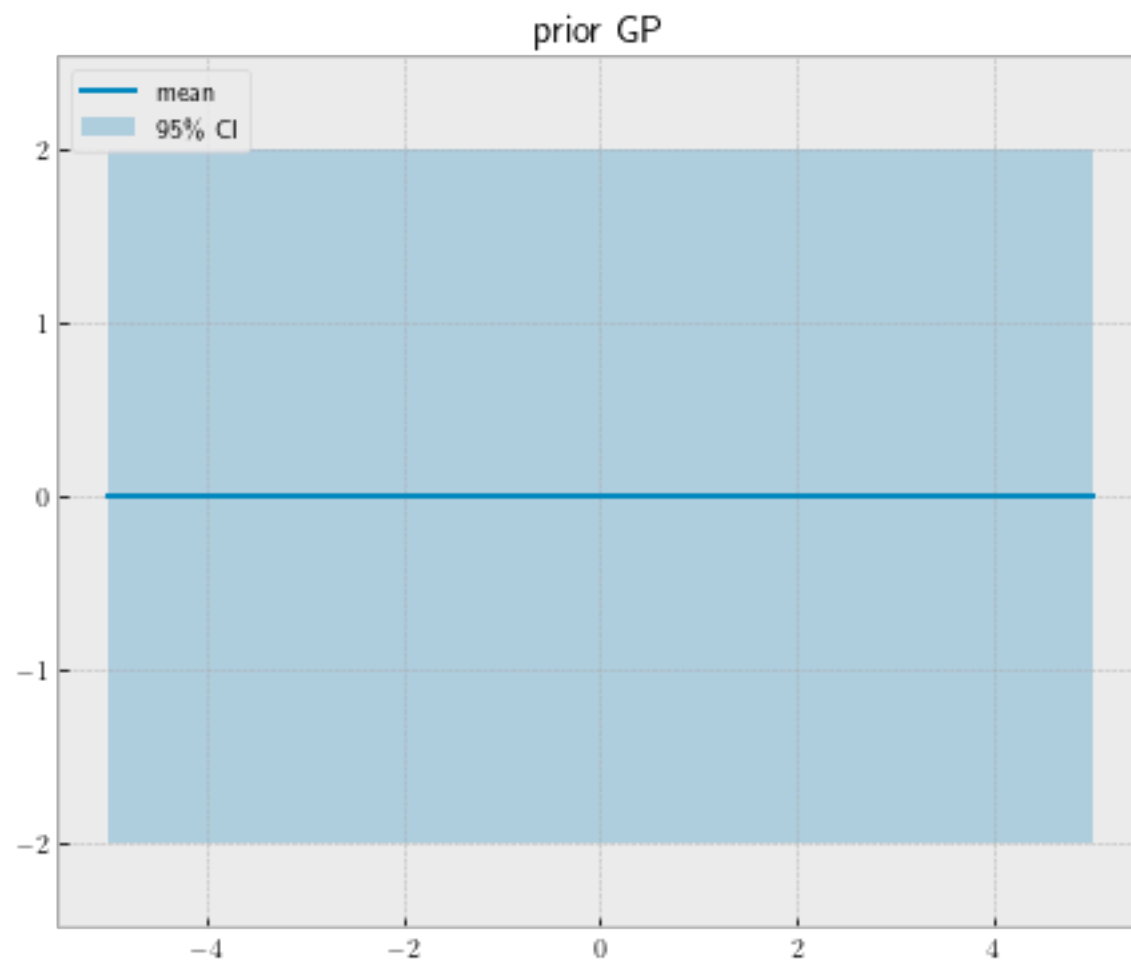
- $\sigma_*^2 \leq K(\mathbf{x}_*, \mathbf{x}_*)$: uncertainty **never** increases
- Uncertainty quantification

\mathbf{x}_* can be anything

- Distribution over **functions**

GP IN ACTION

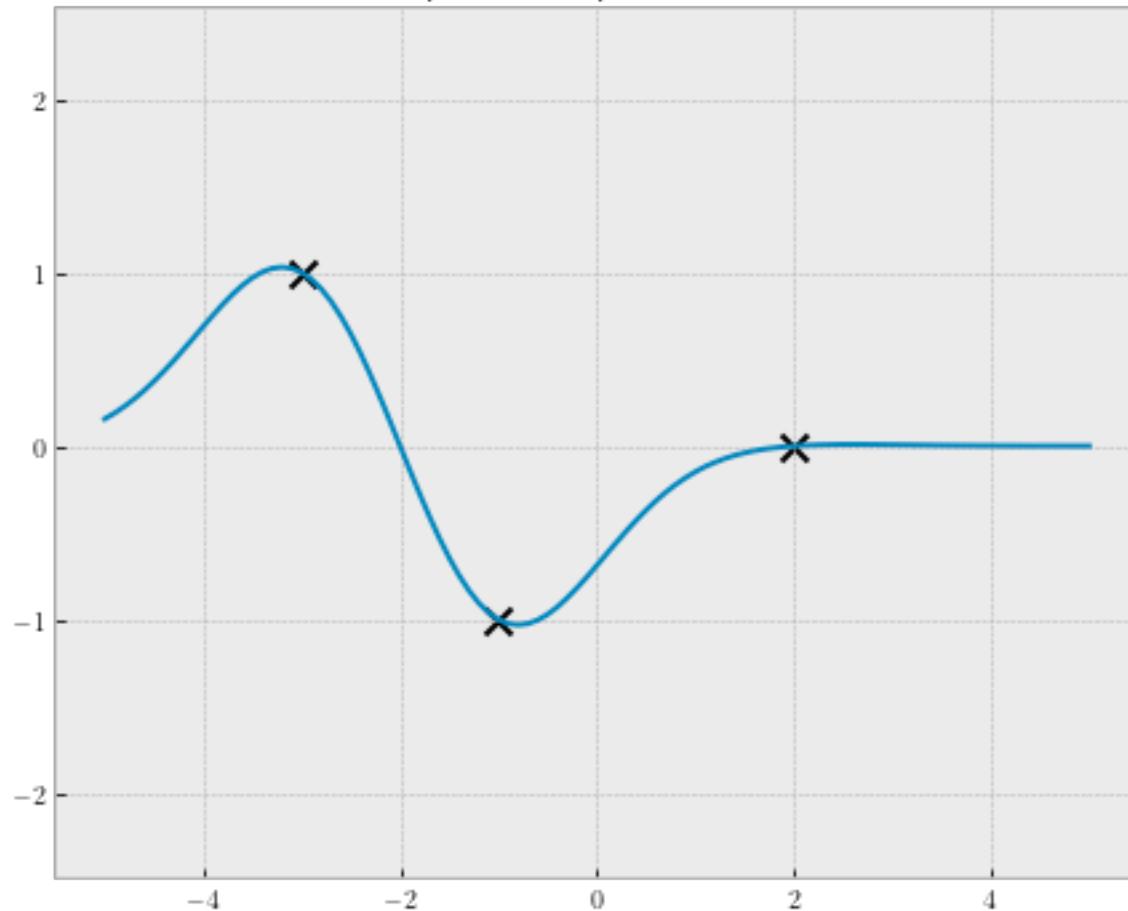
RBF kernel: $K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$



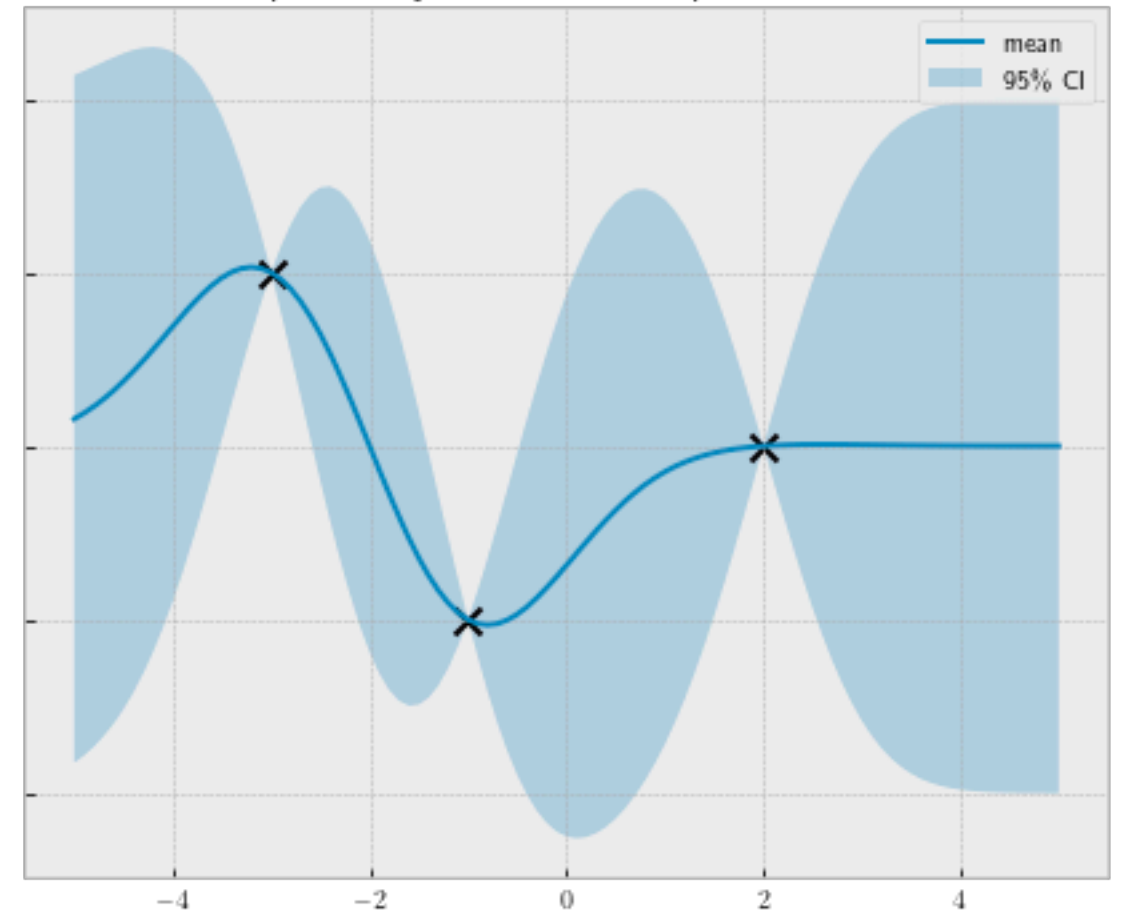
GP IN ACTION

$$\text{RBF kernel: } K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$$

point-wise predictions

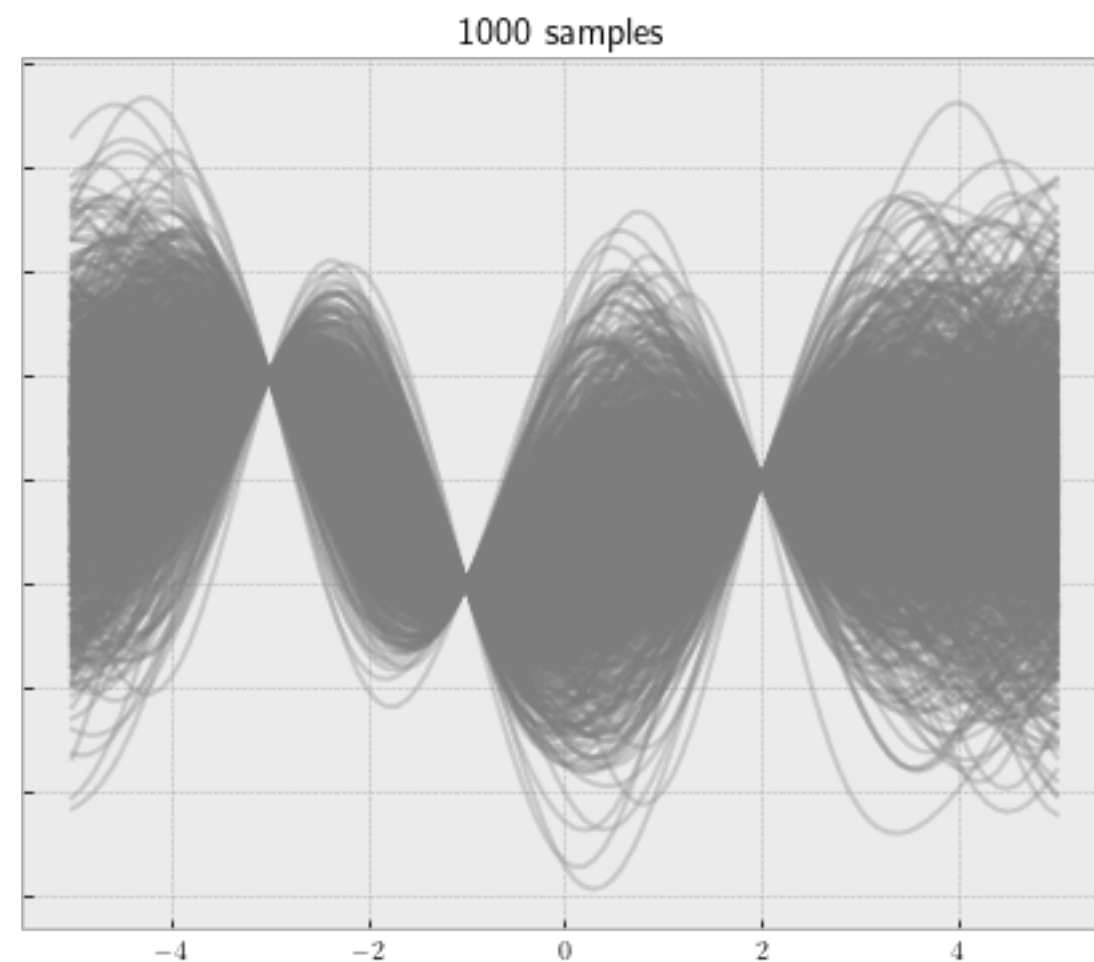
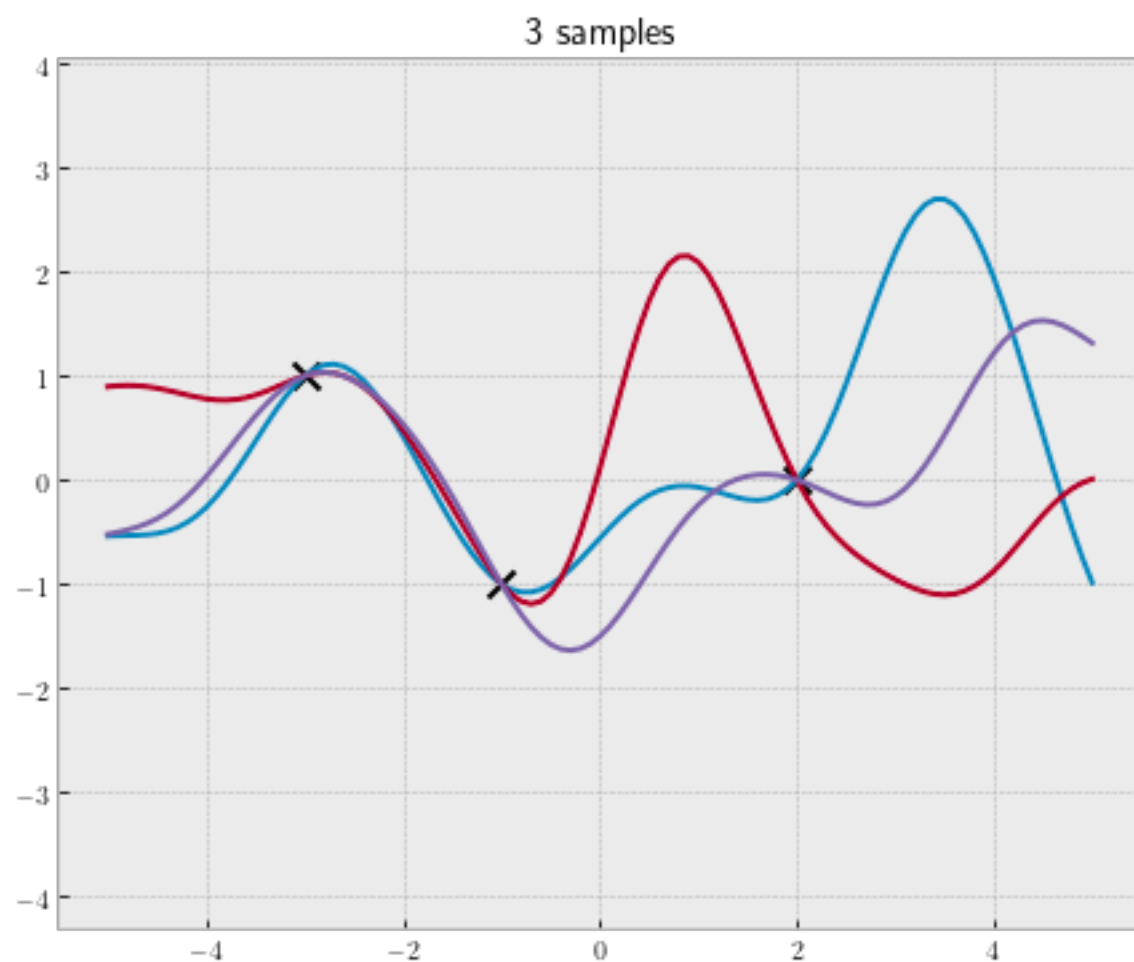


probability distributions as predictions



GP IN ACTION

RBF kernel: $K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$



NOISY OBSERVATIONS

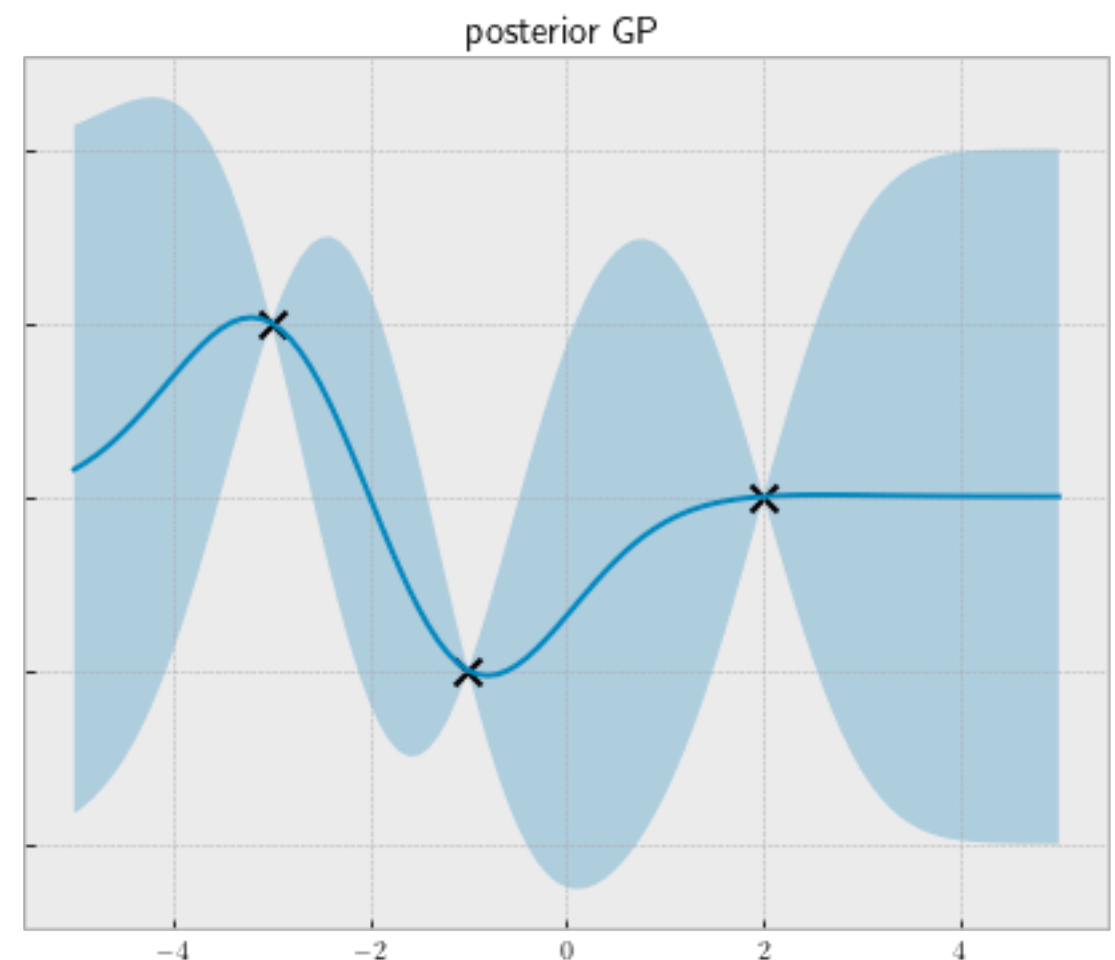
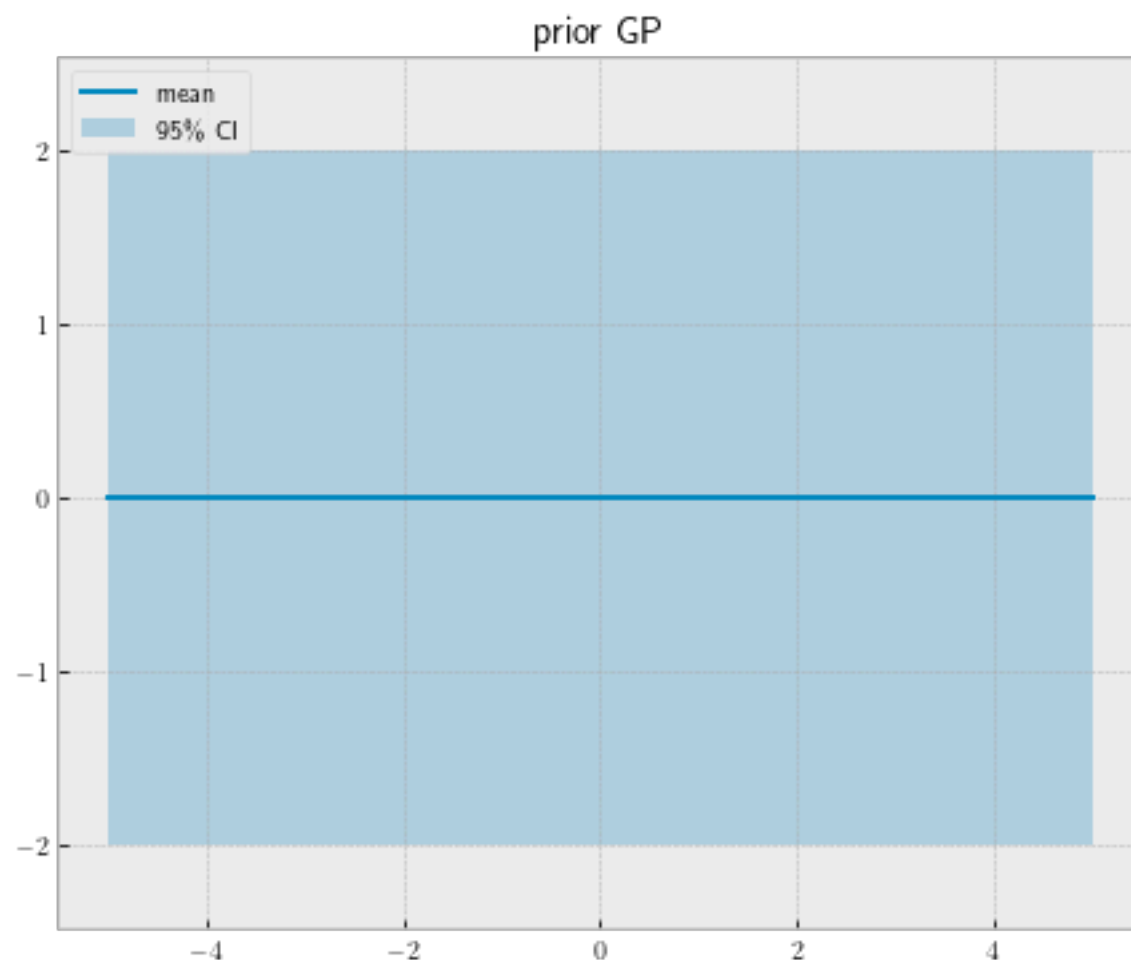
Each observation y follows a slightly different normal:

$$p(y \mid \mathbf{x}) = N(y; \mu(\mathbf{x}), K(\mathbf{x}, \mathbf{x}) + \sigma_n^2) \quad K(X, X) \rightarrow K(X, X) + \sigma_n^2 I$$

NOISY OBSERVATIONS

Each observation y follows a slightly different normal:

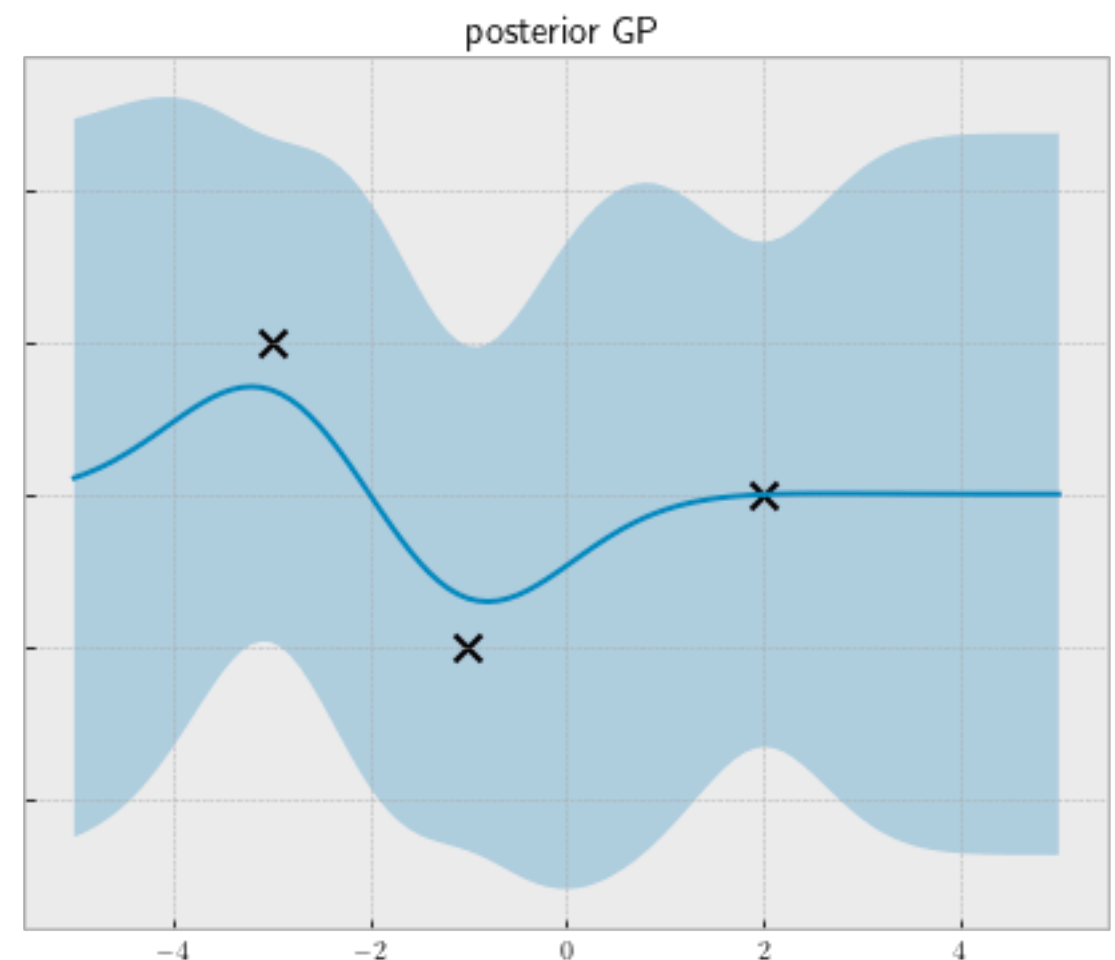
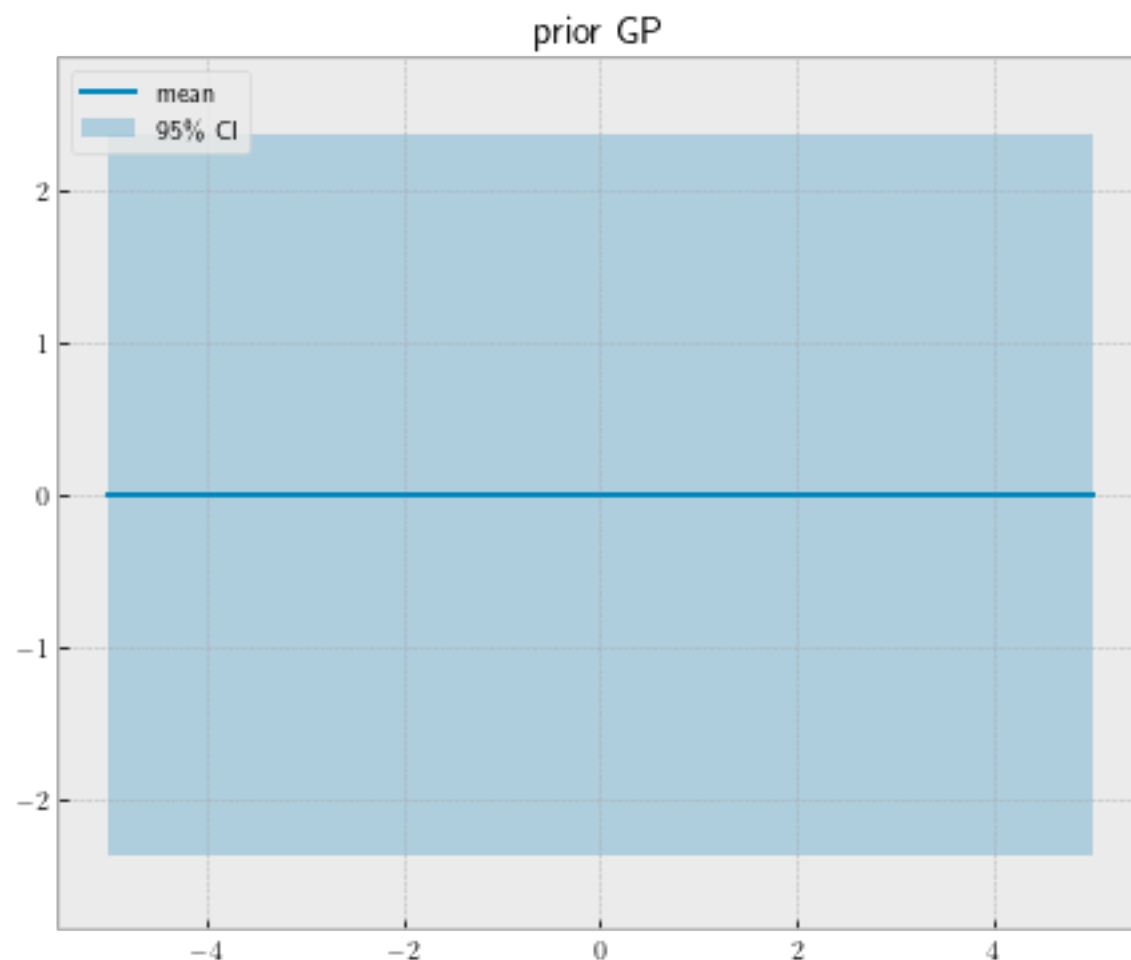
$$p(y \mid \mathbf{x}) = N(y; \mu(\mathbf{x}), K(\mathbf{x}, \mathbf{x}) + \sigma_n^2) \quad K(X, X) \rightarrow K(X, X) + \sigma_n^2 I$$



NOISY OBSERVATIONS

Each observation y follows a slightly different normal:

$$p(y \mid \mathbf{x}) = N(y; \mu(\mathbf{x}), K(\mathbf{x}, \mathbf{x}) + \sigma_n^2) \quad K(X, X) \rightarrow K(X, X) + \sigma_n^2 I$$



FORMAL DEFINITION OF A GP

FORMAL DEFINITION OF A GP

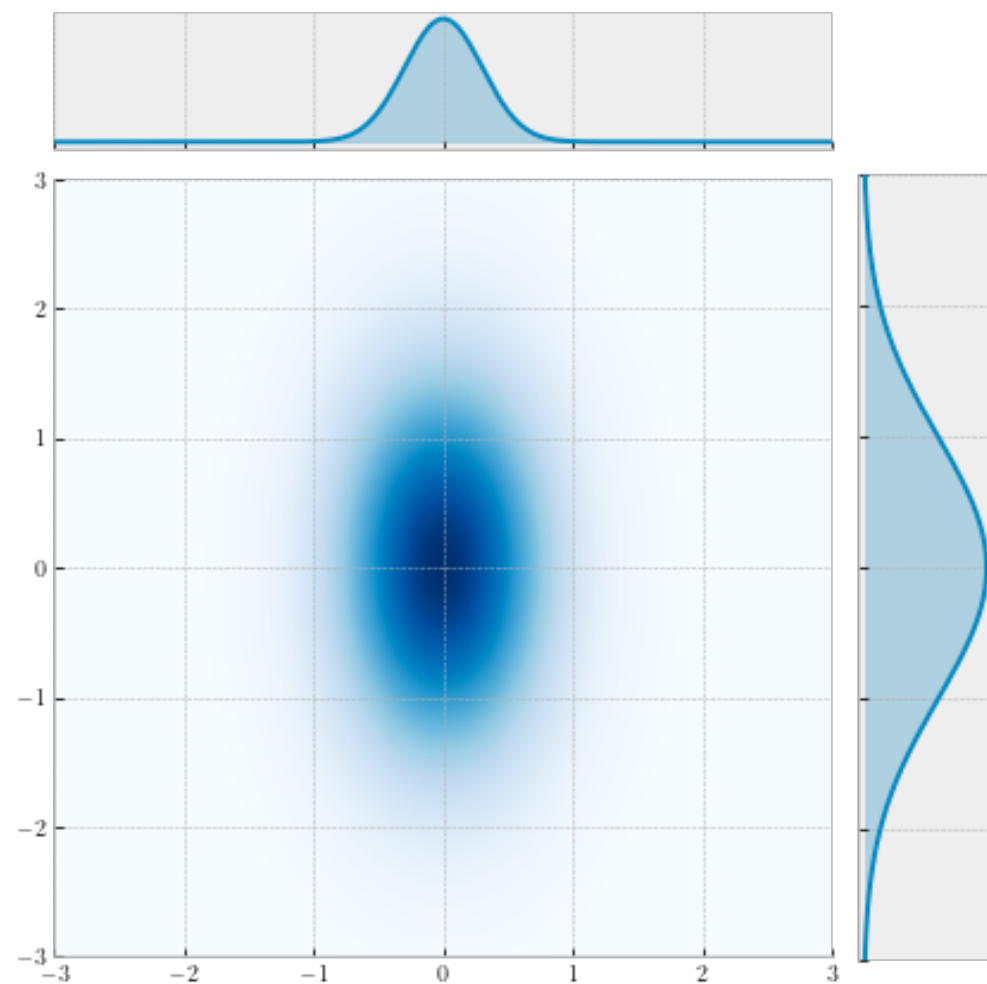
Definition: A GP is a (potentially infinite) collection of RVs such that the **joint** distribution of **every** finite subset of RVs is multivariate Gaussian.

UNCORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 0.3 & 0 \\ 0 & 1 \end{bmatrix}$$

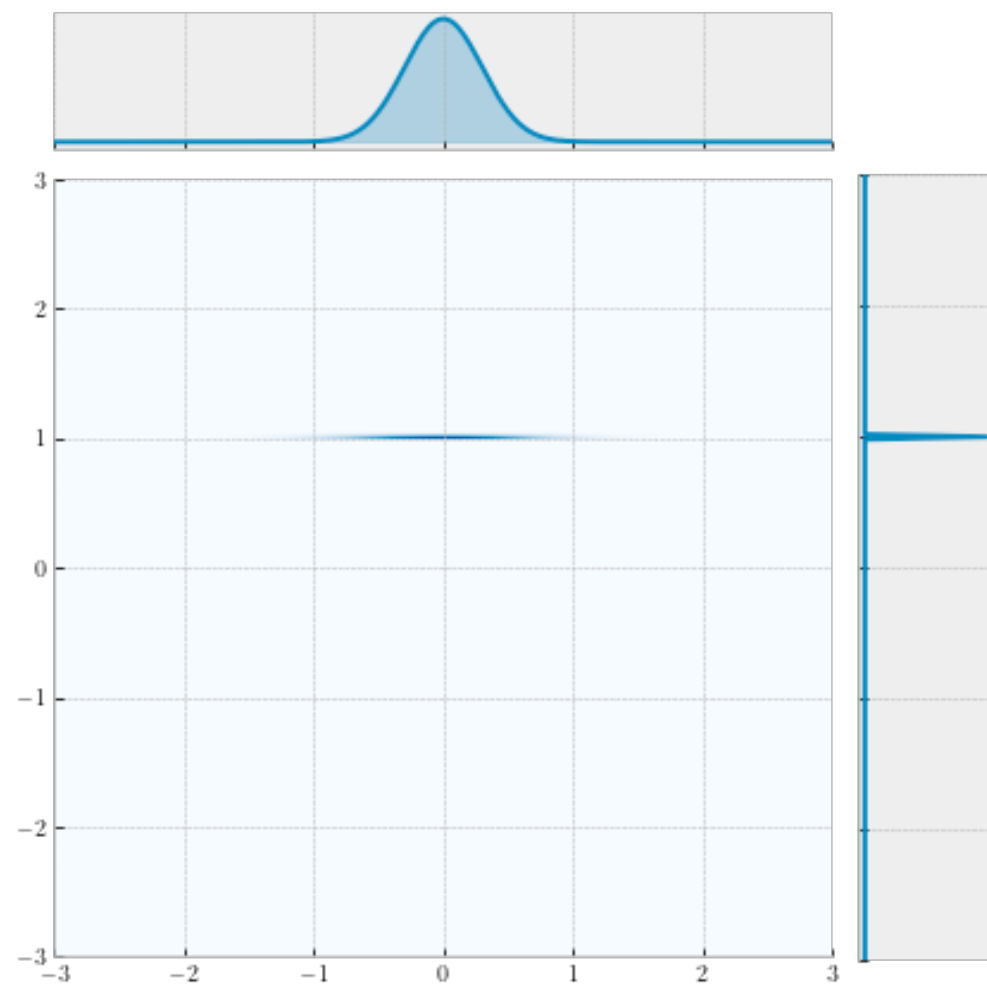
UNCORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 0.3 & 0 \\ 0 & 1 \end{bmatrix}$$



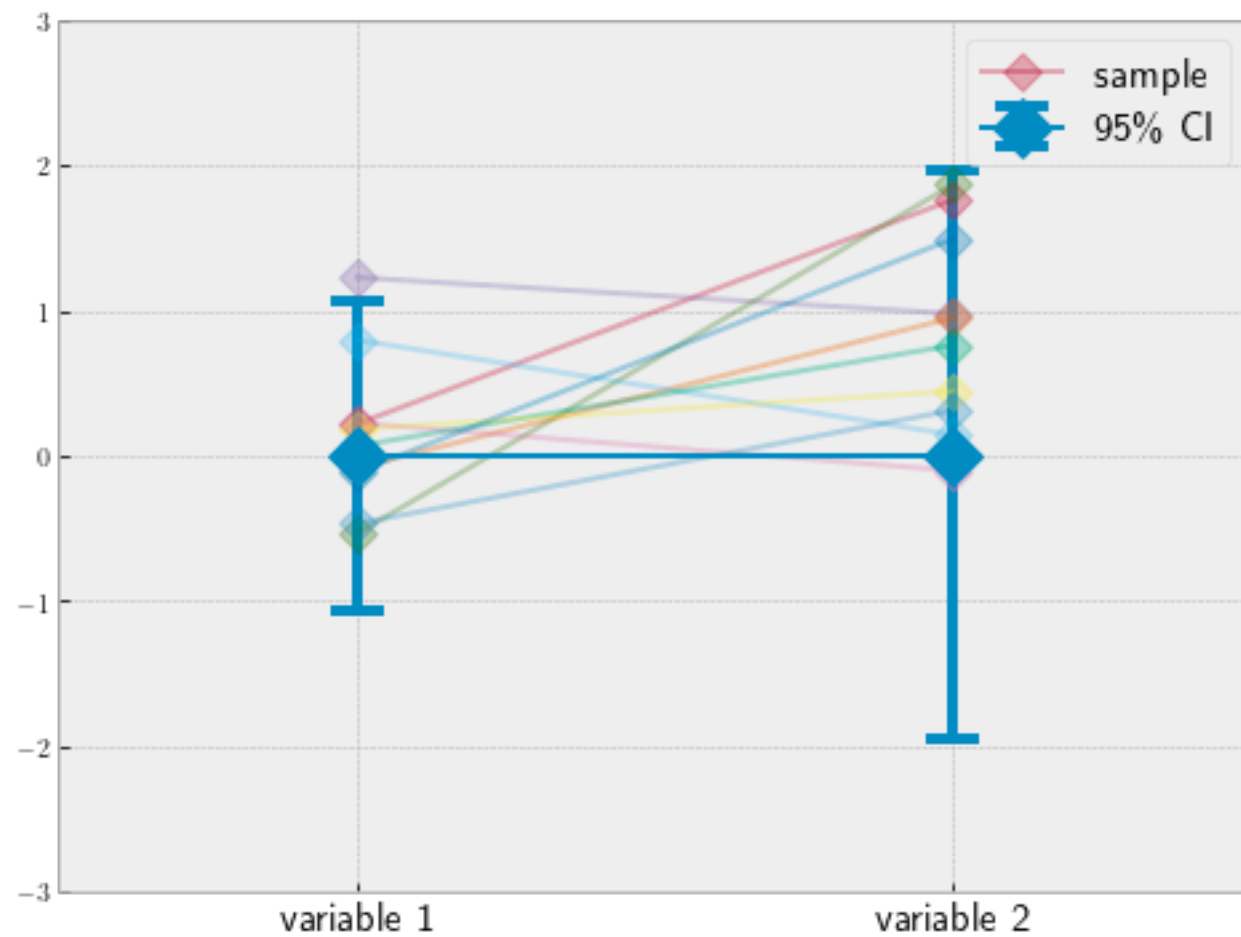
UNCORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 0.3 & 0 \\ 0 & 1 \end{bmatrix}$$



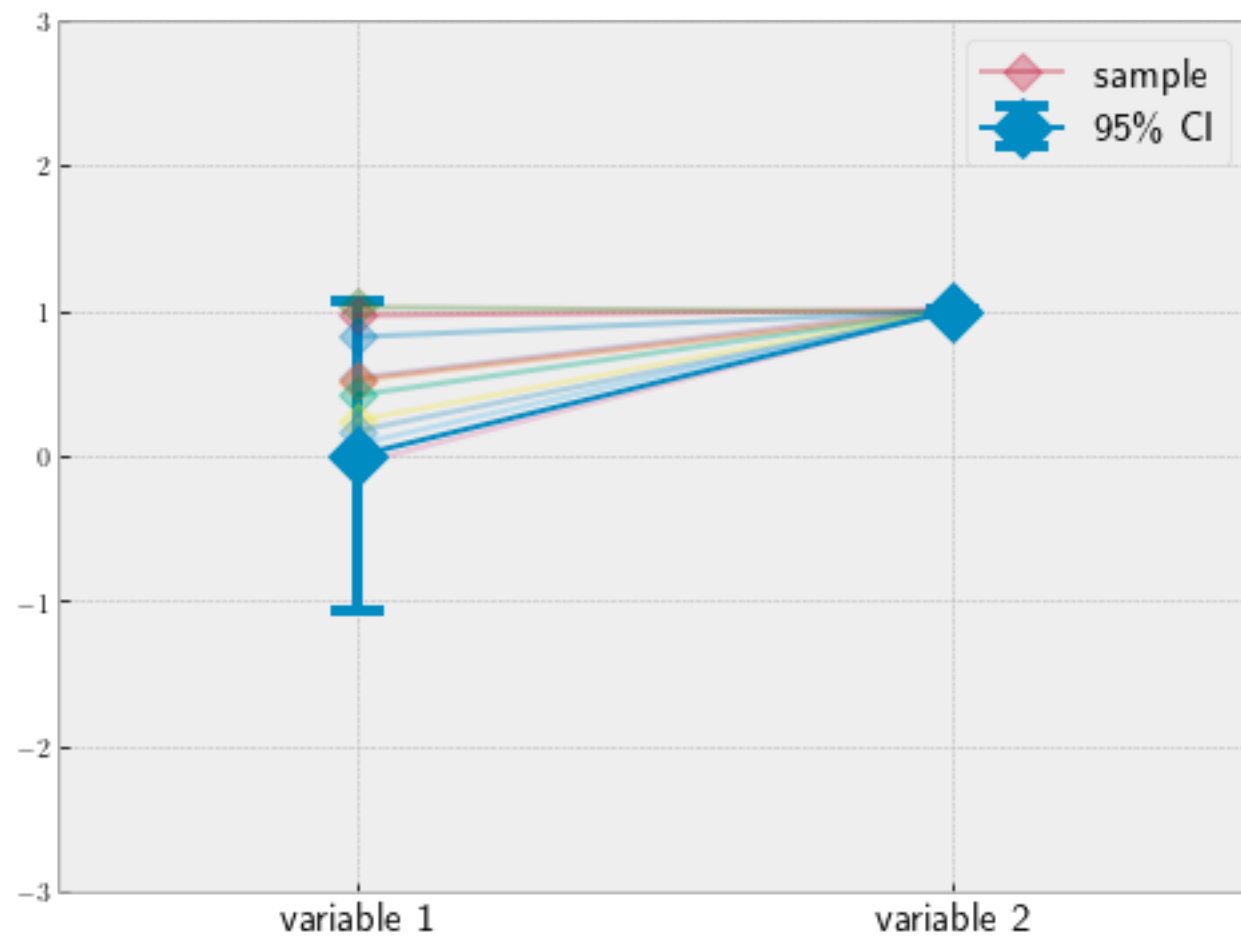
UNCORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 0.3 & 0 \\ 0 & 1 \end{bmatrix}$$



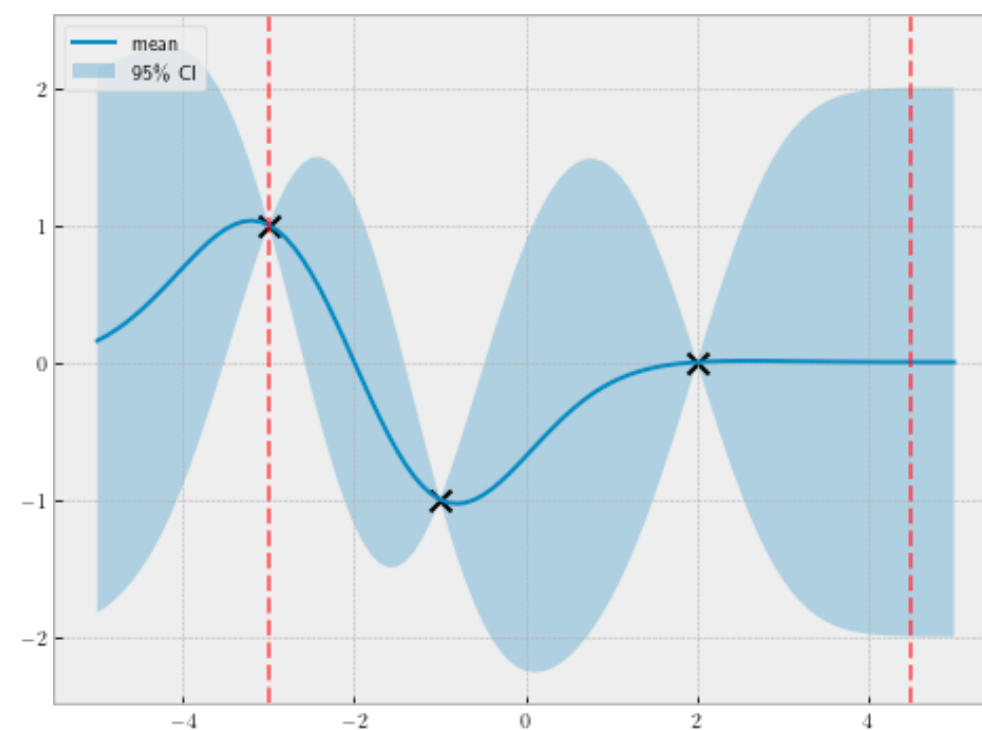
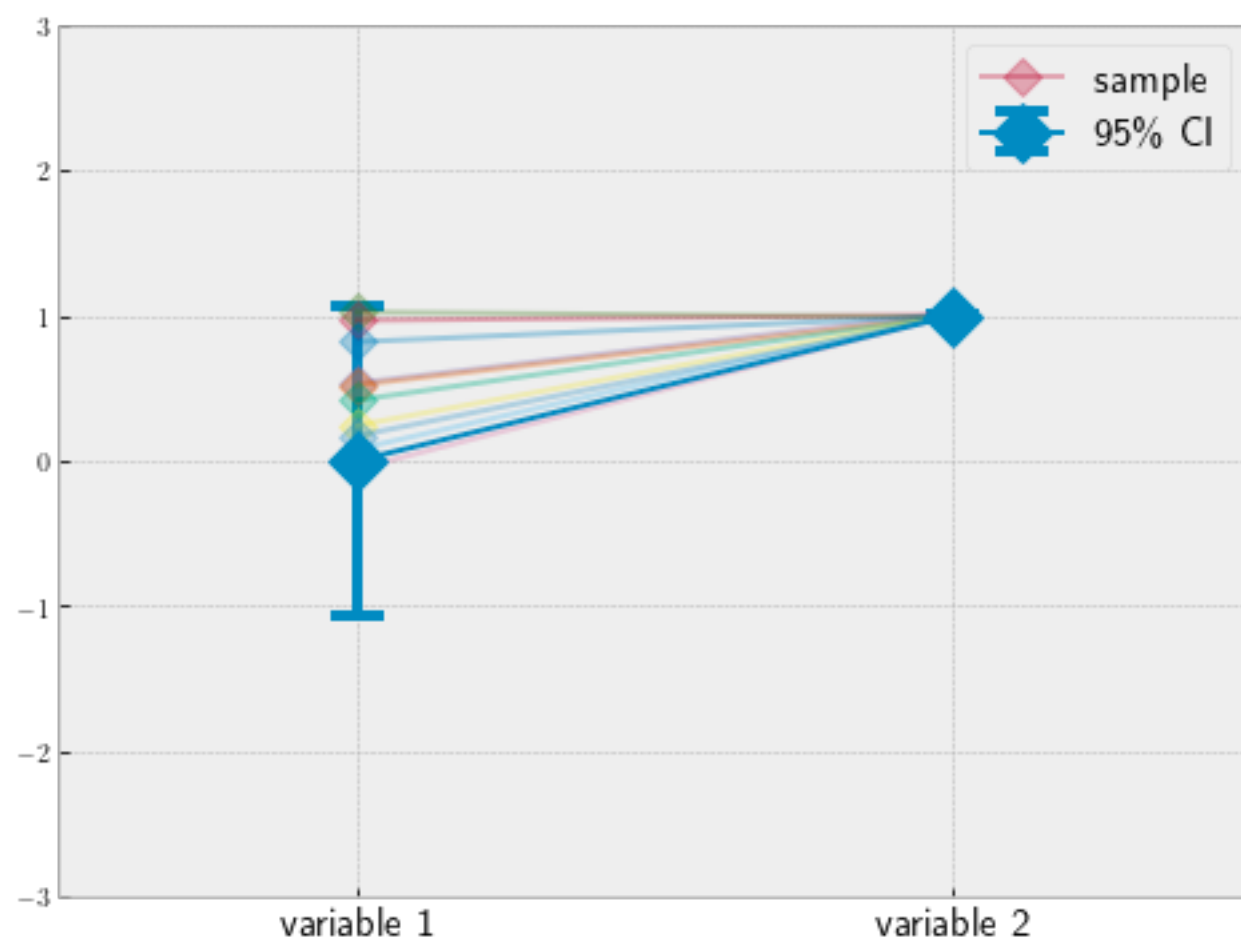
UNCORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 0.3 & 0 \\ 0 & 1 \end{bmatrix}$$



UNCORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 0.3 & 0 \\ 0 & 1 \end{bmatrix}$$

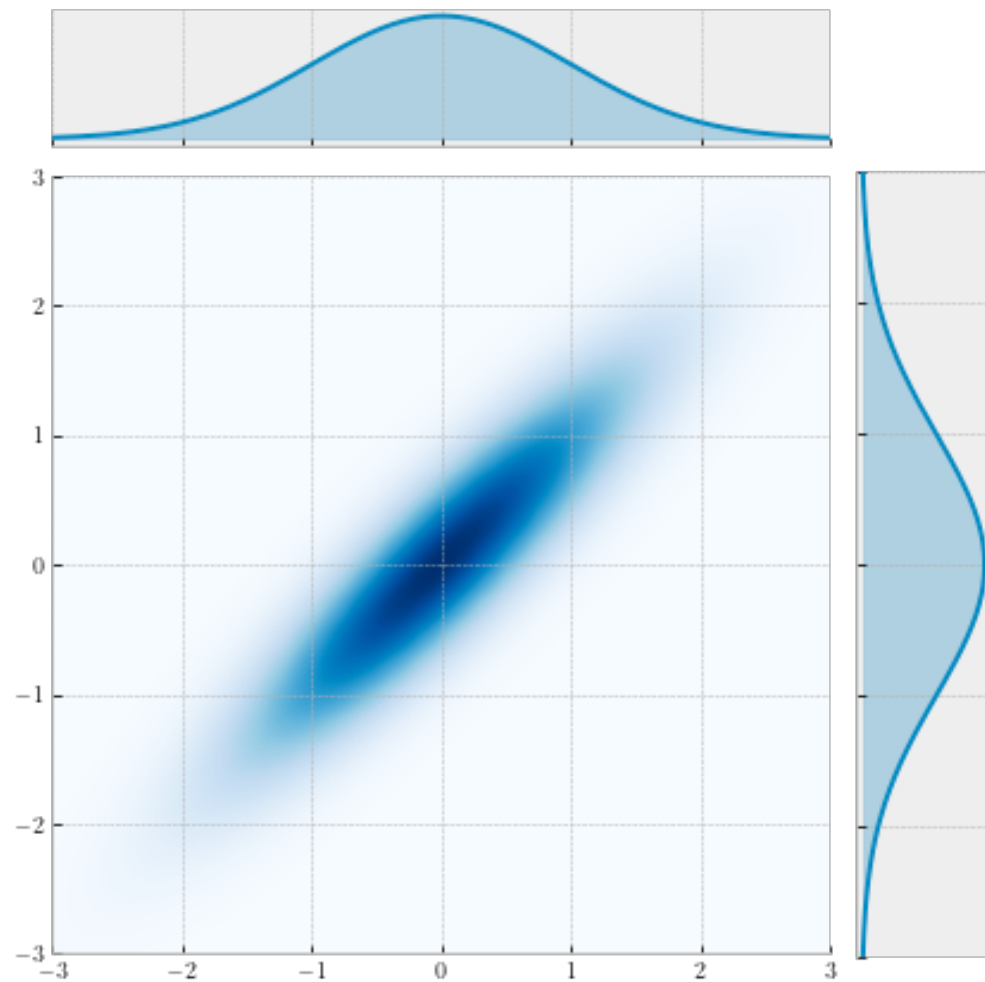


CORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$

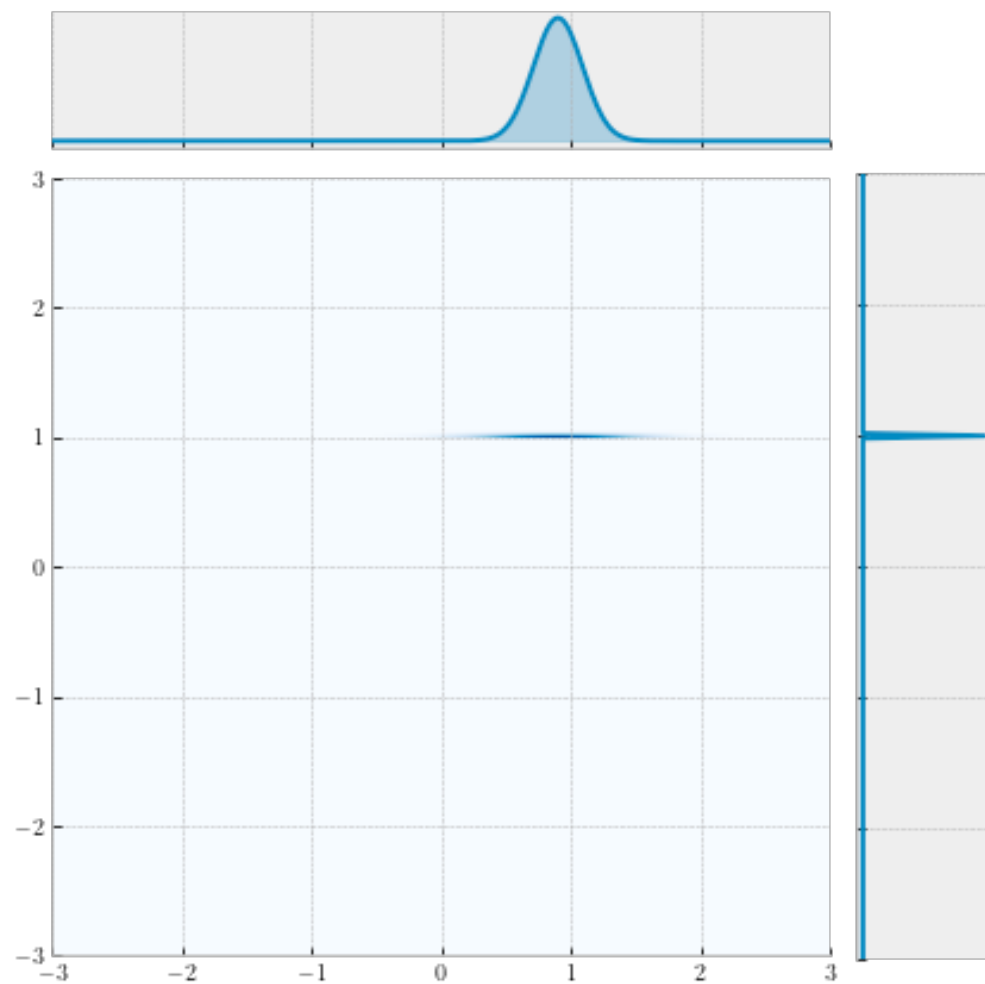
CORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$



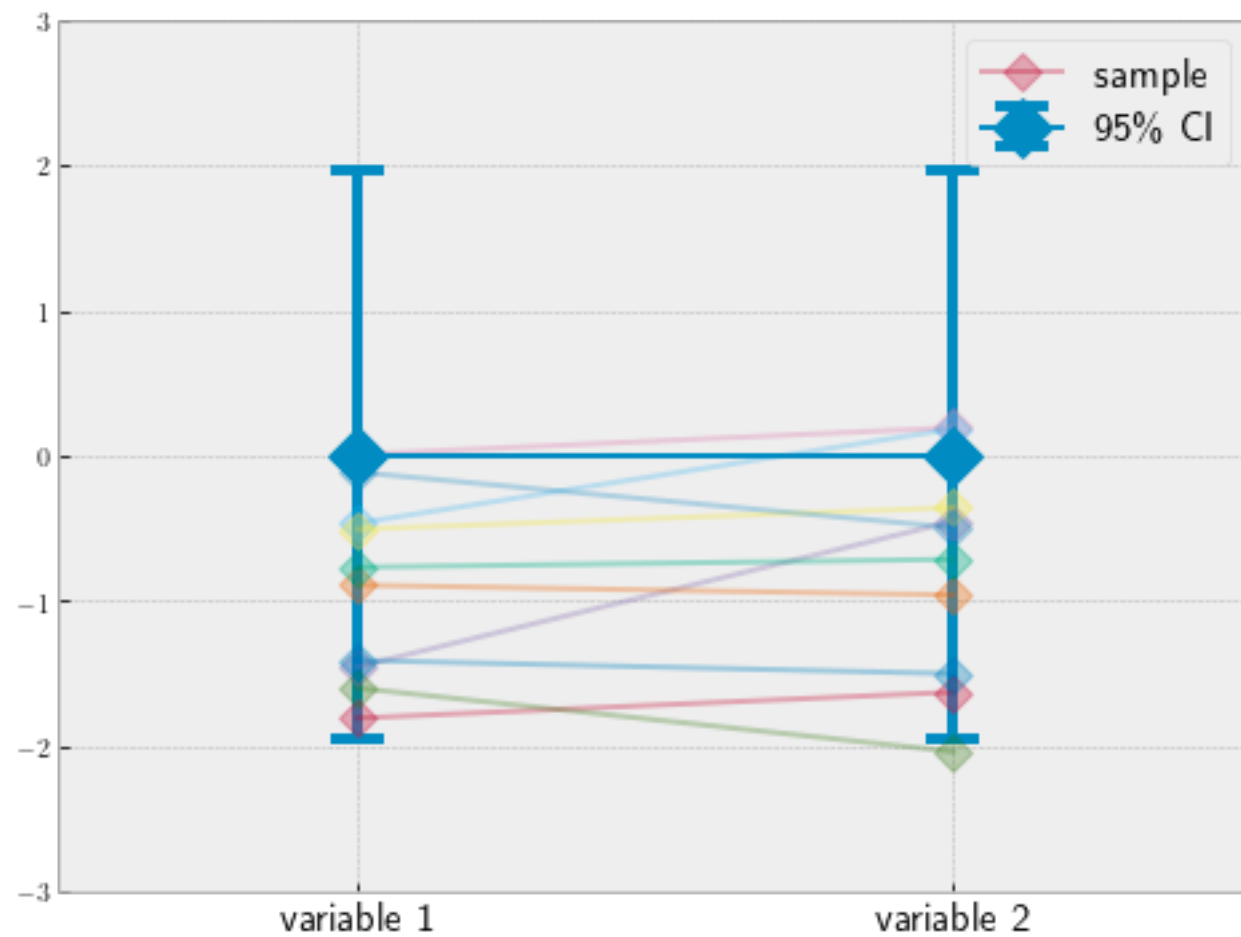
CORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$



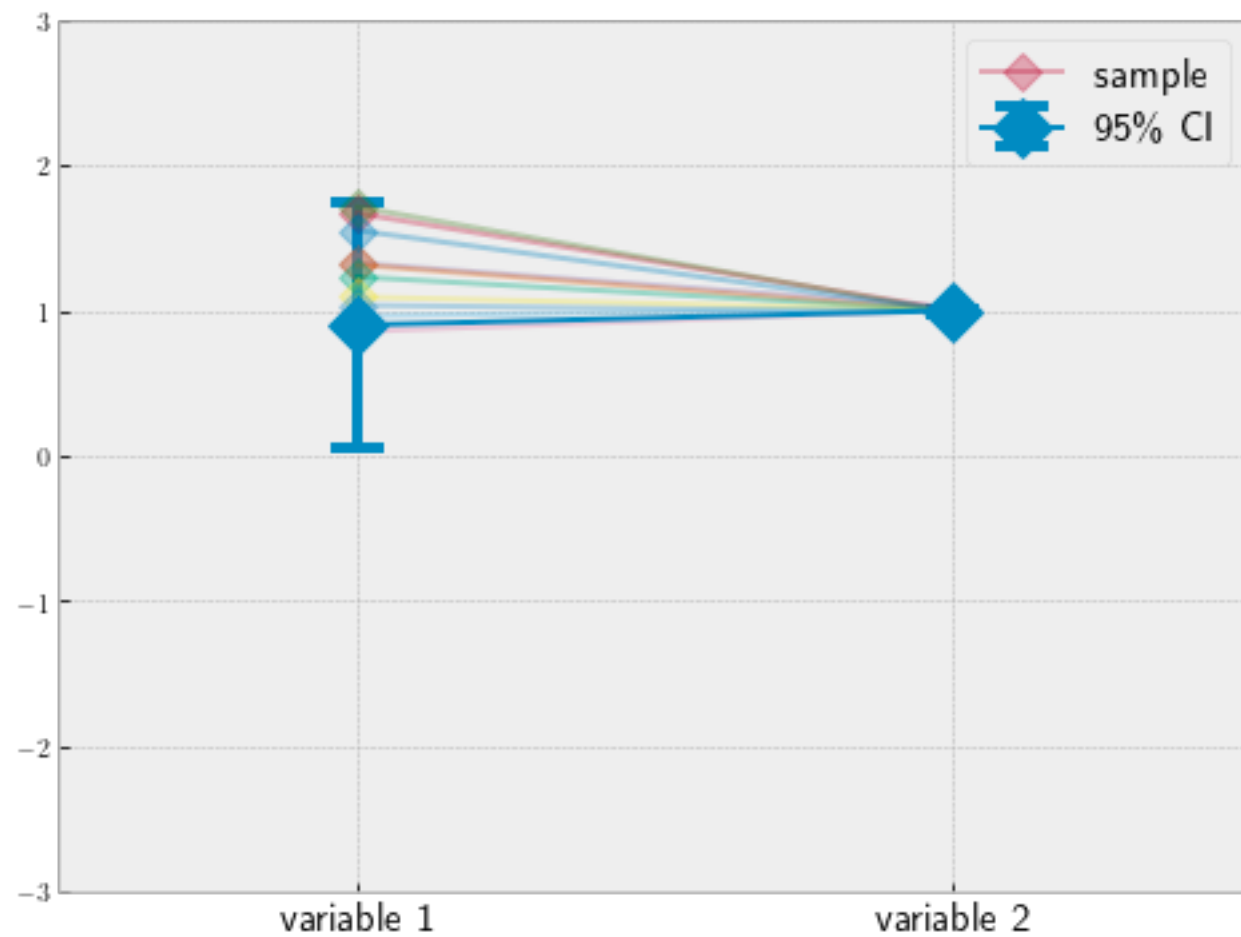
CORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$



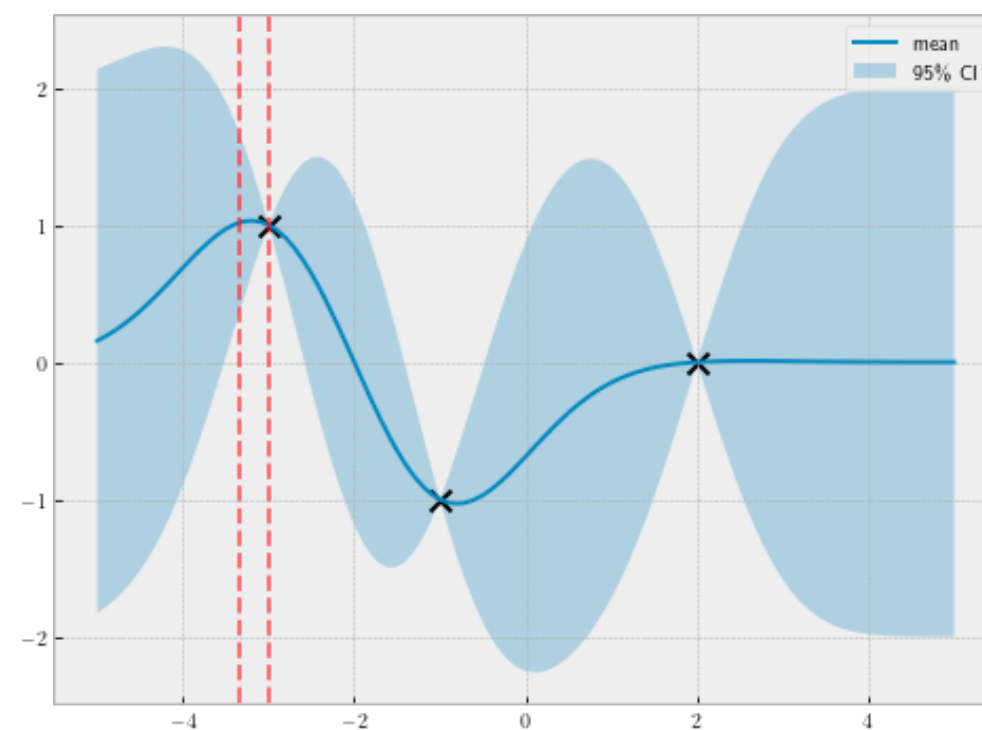
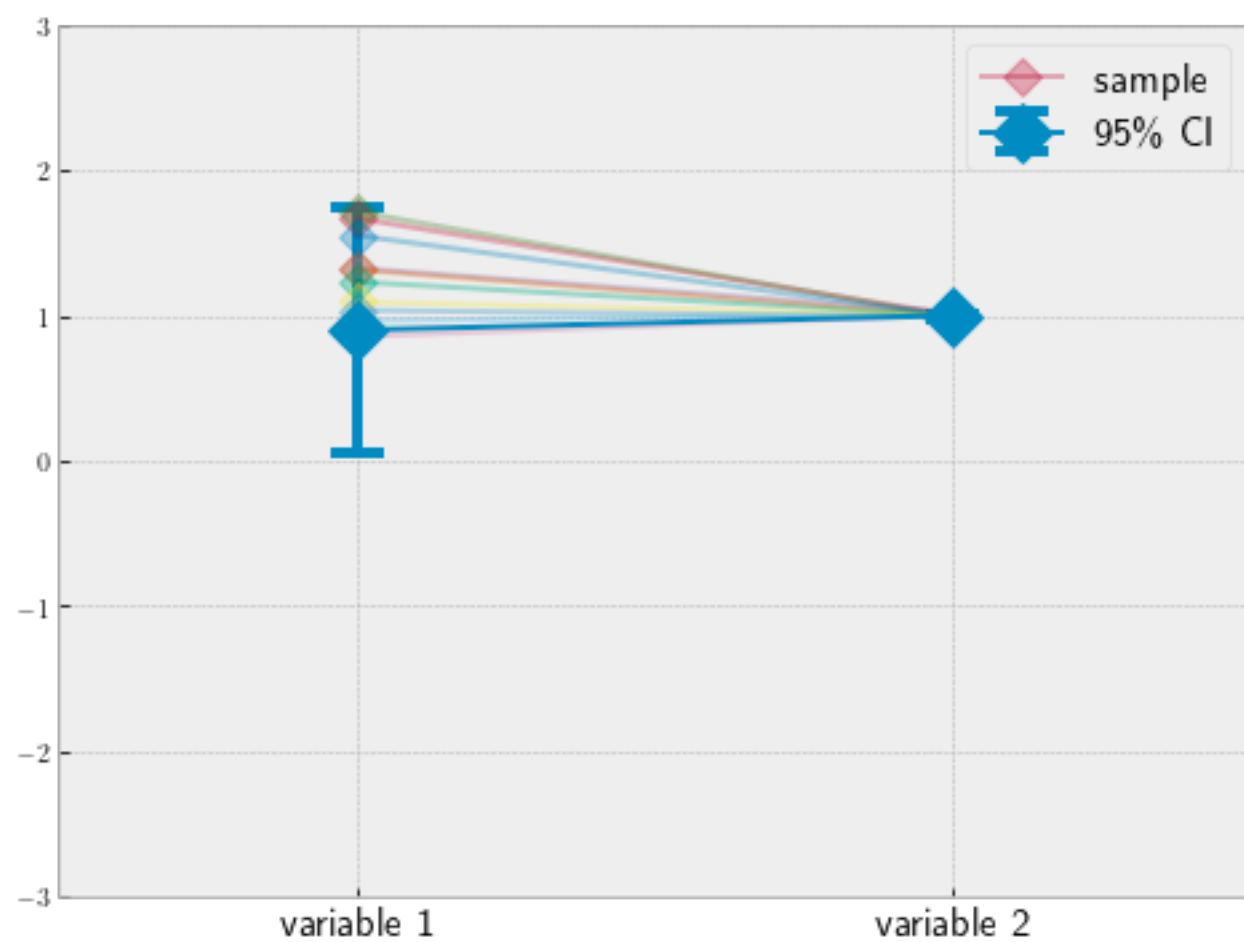
CORRELATED VARIABLES

$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$

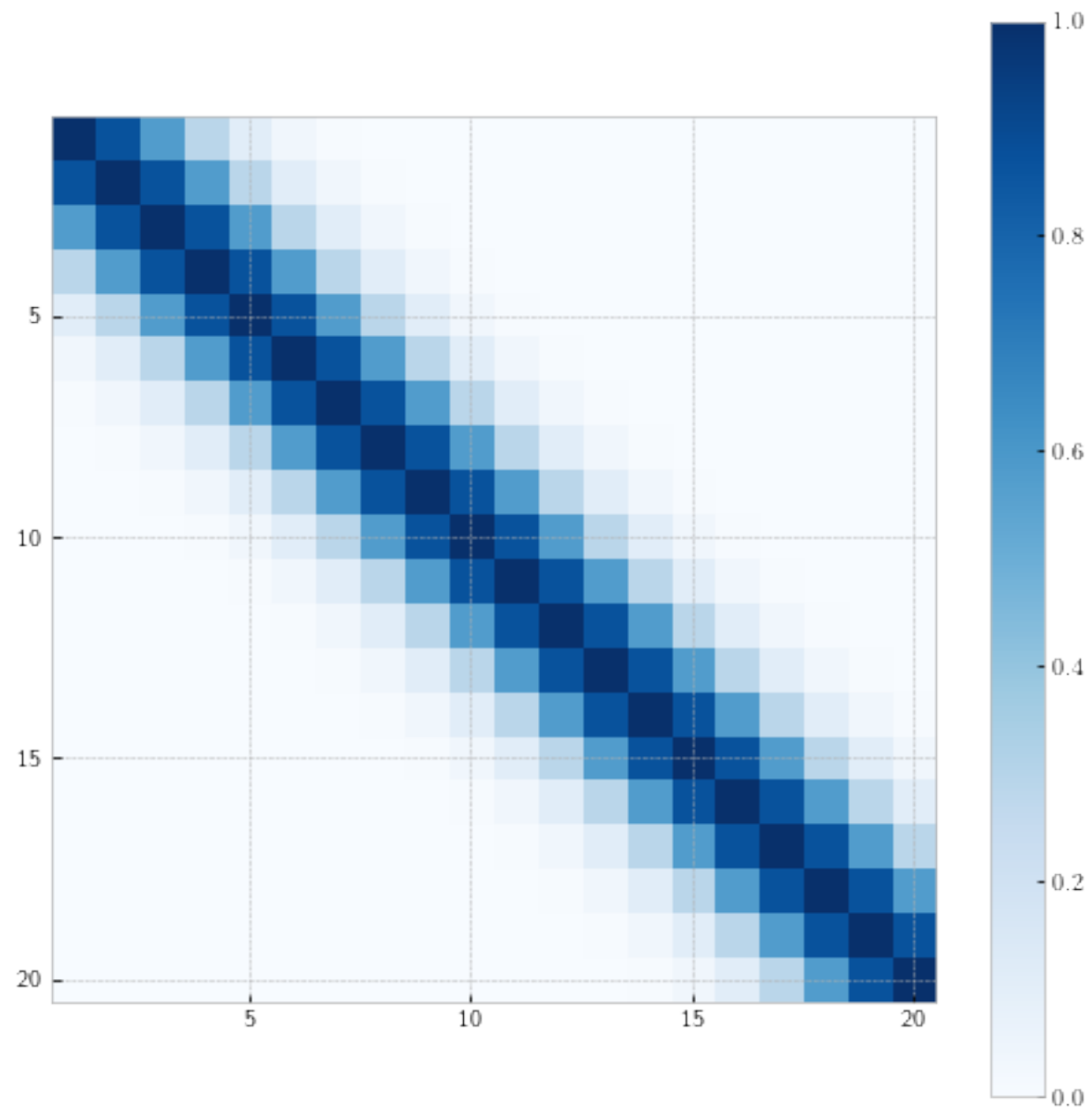


CORRELATED VARIABLES

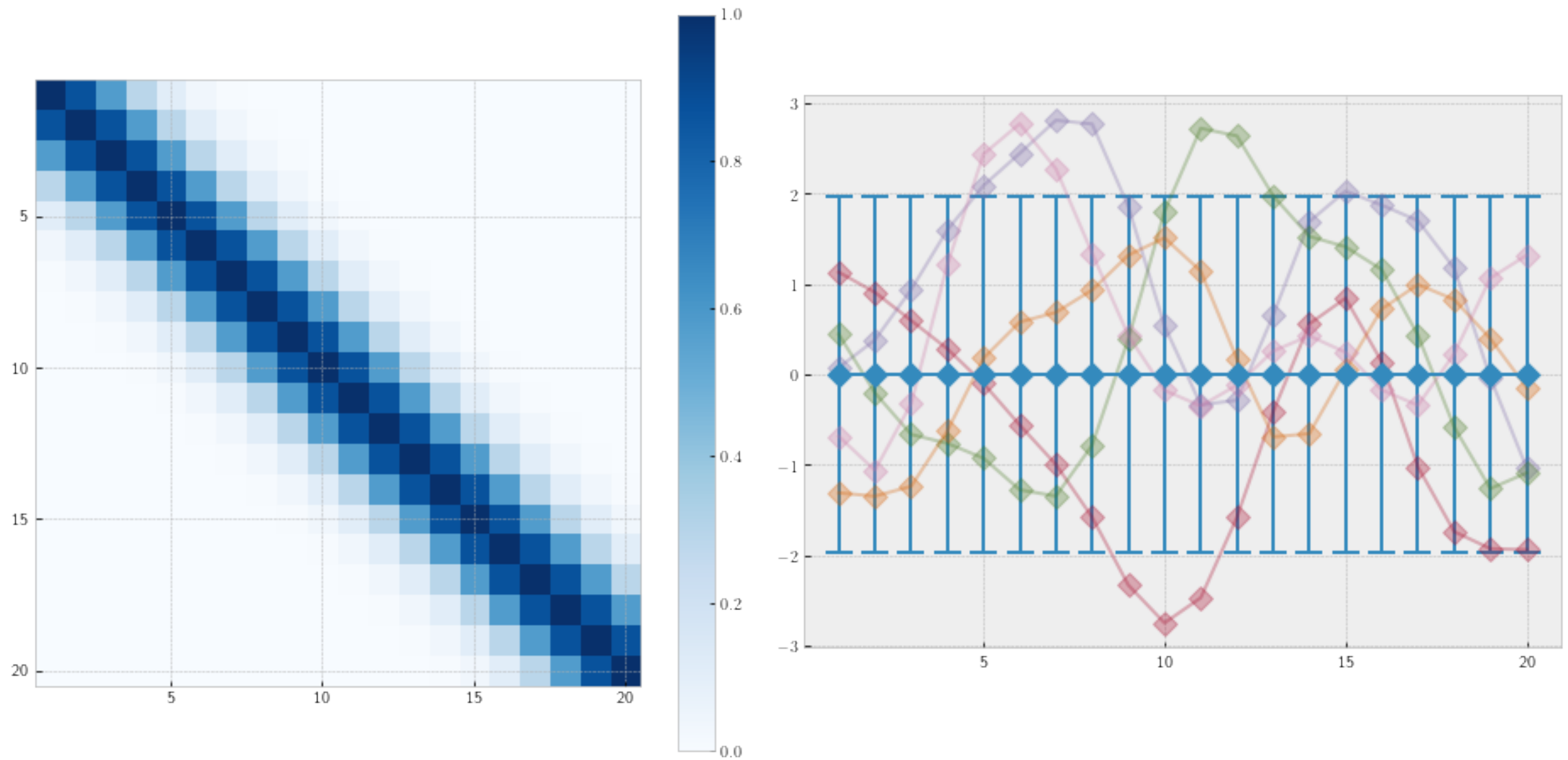
$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$



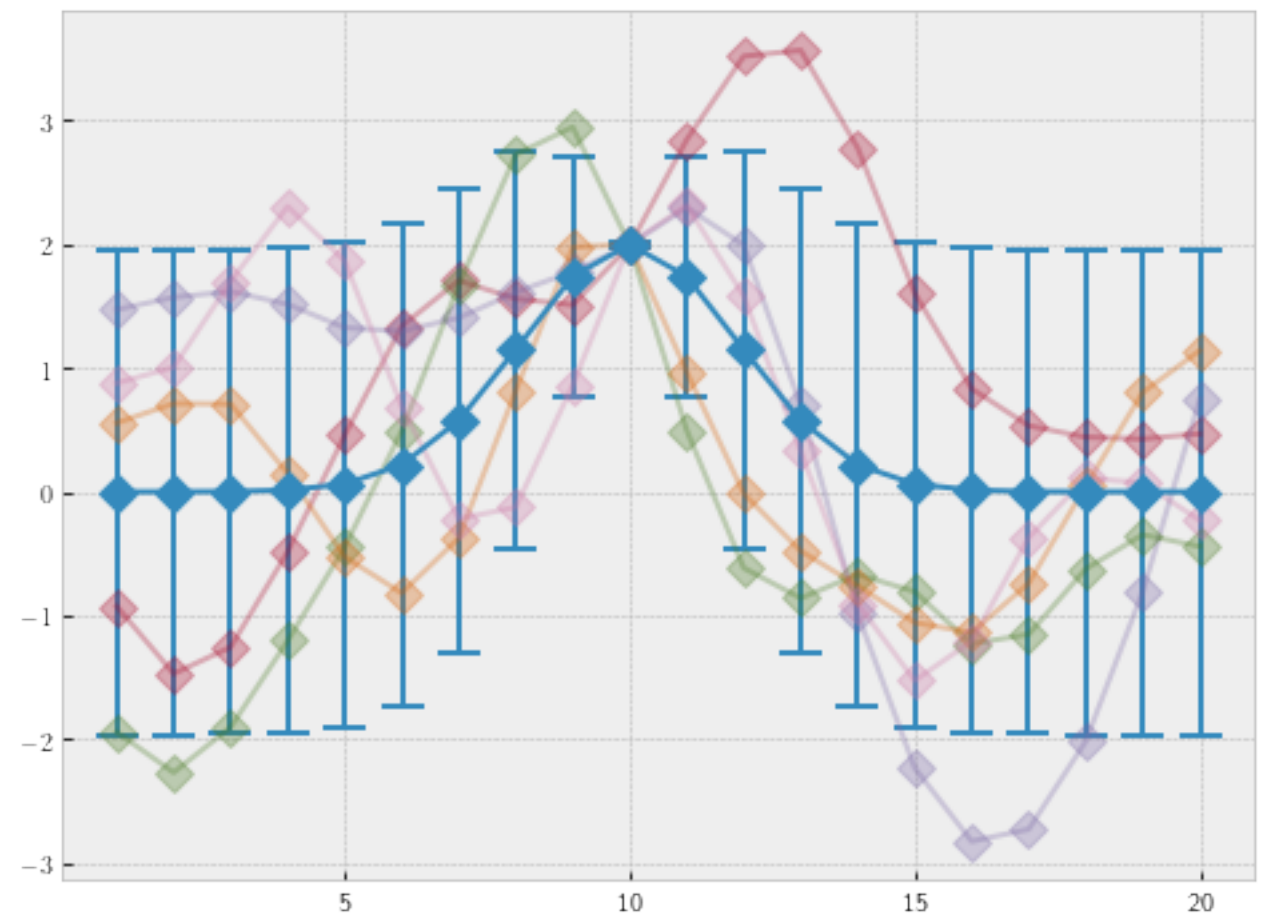
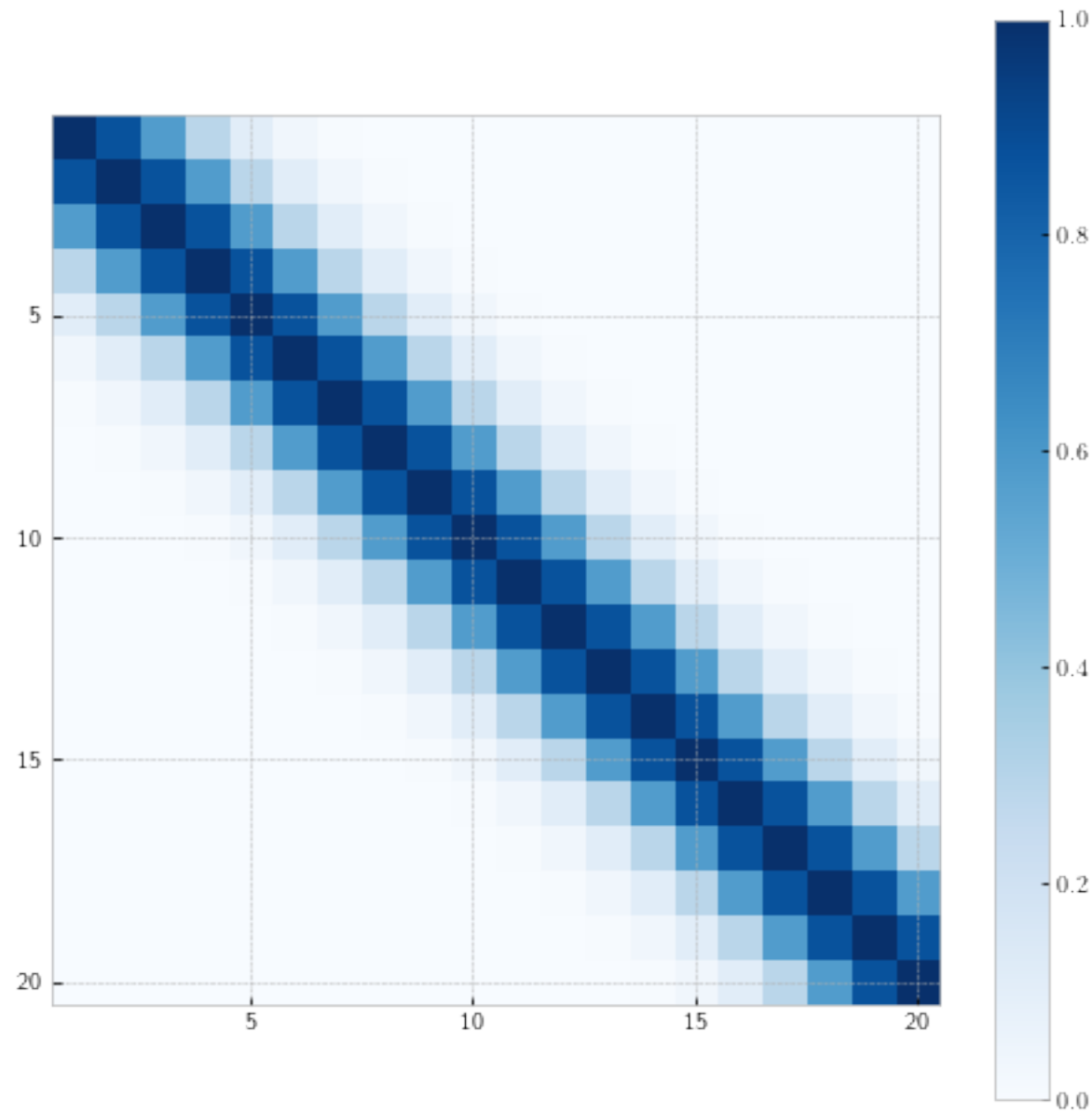
FROM 2D TO HIGHER/INFINITE DIMENSIONS



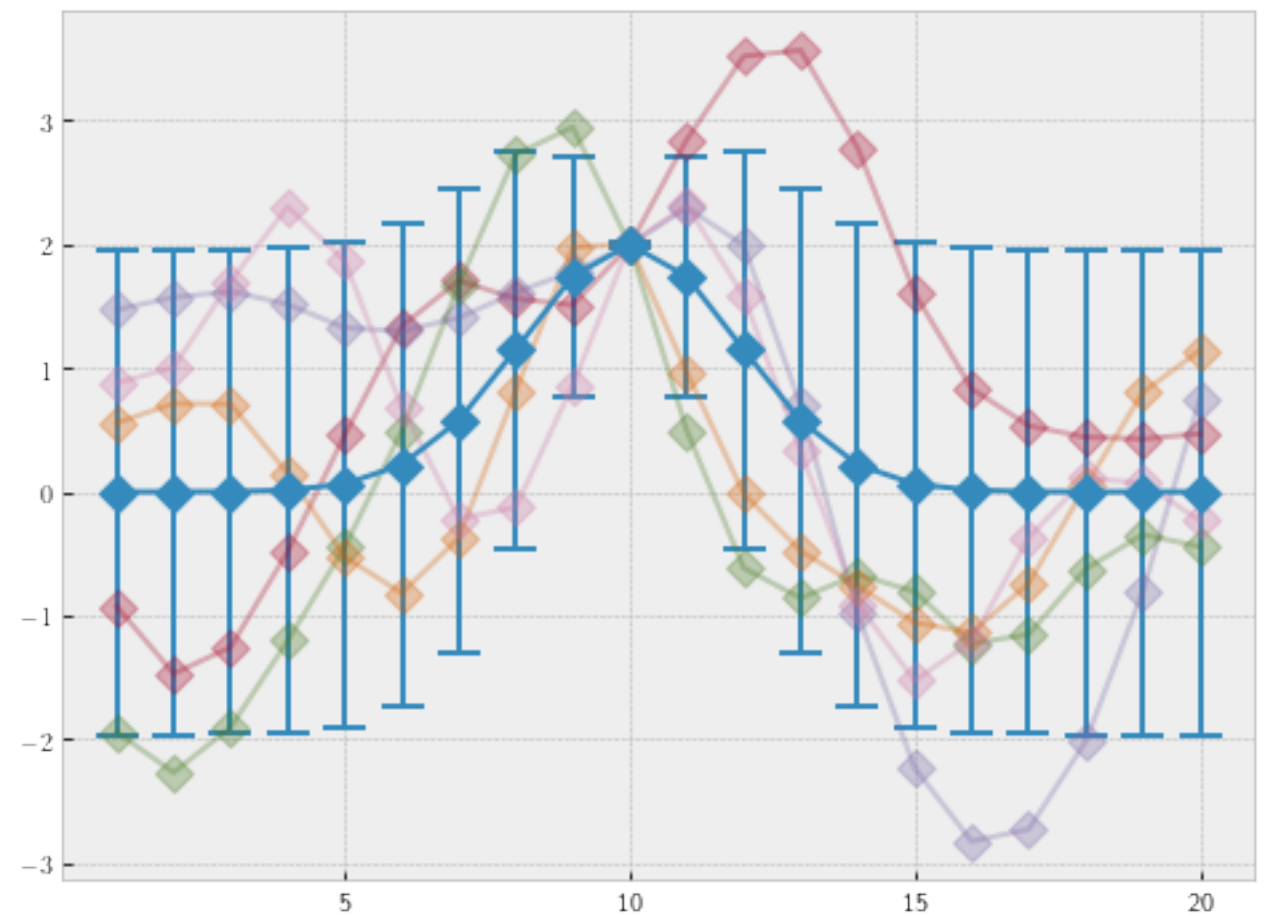
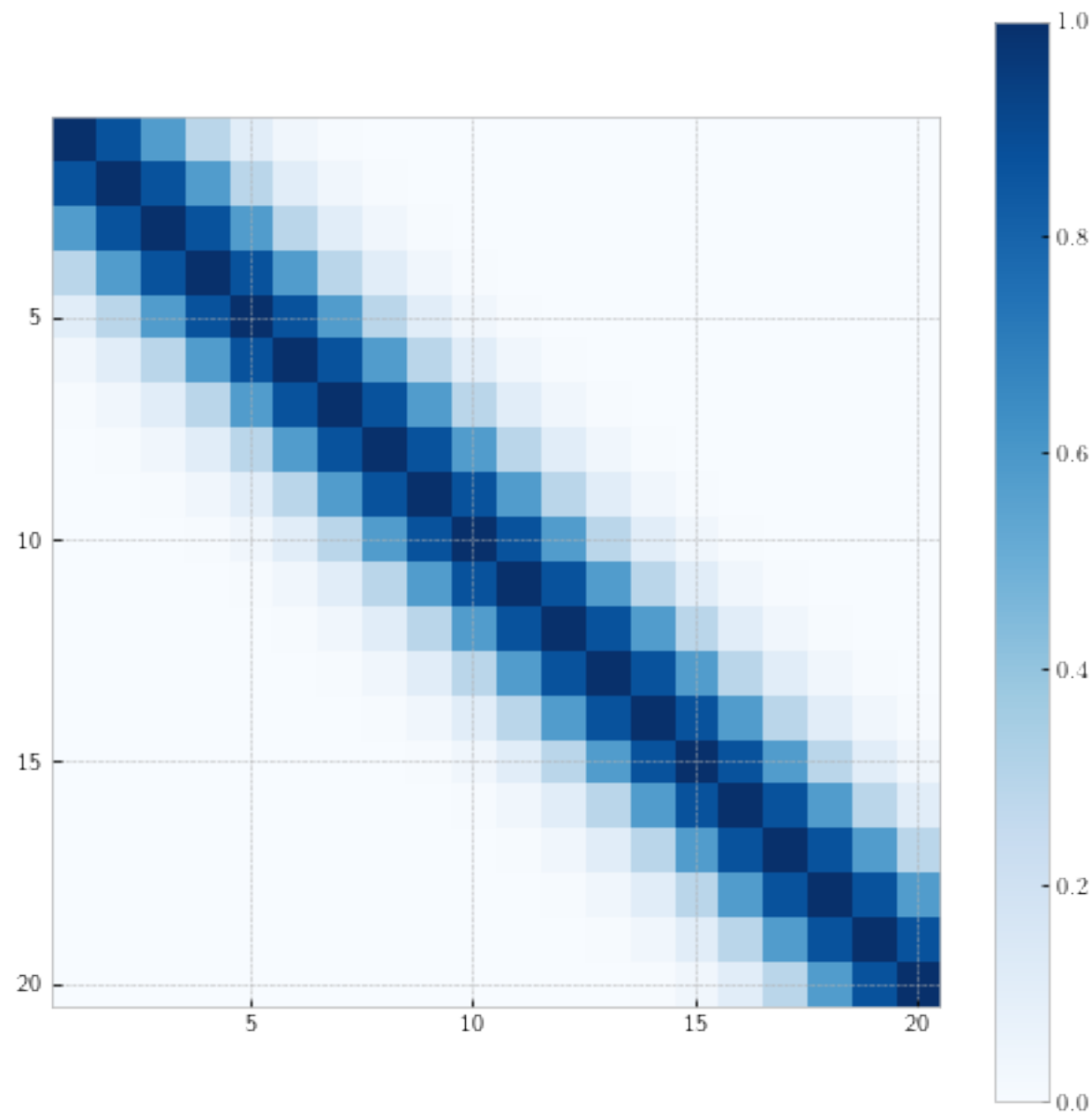
FROM 2D TO HIGHER/INFINITE DIMENSIONS



FROM 2D TO HIGHER/INFINITE DIMENSIONS



FROM 2D TO HIGHER/INFINITE DIMENSIONS



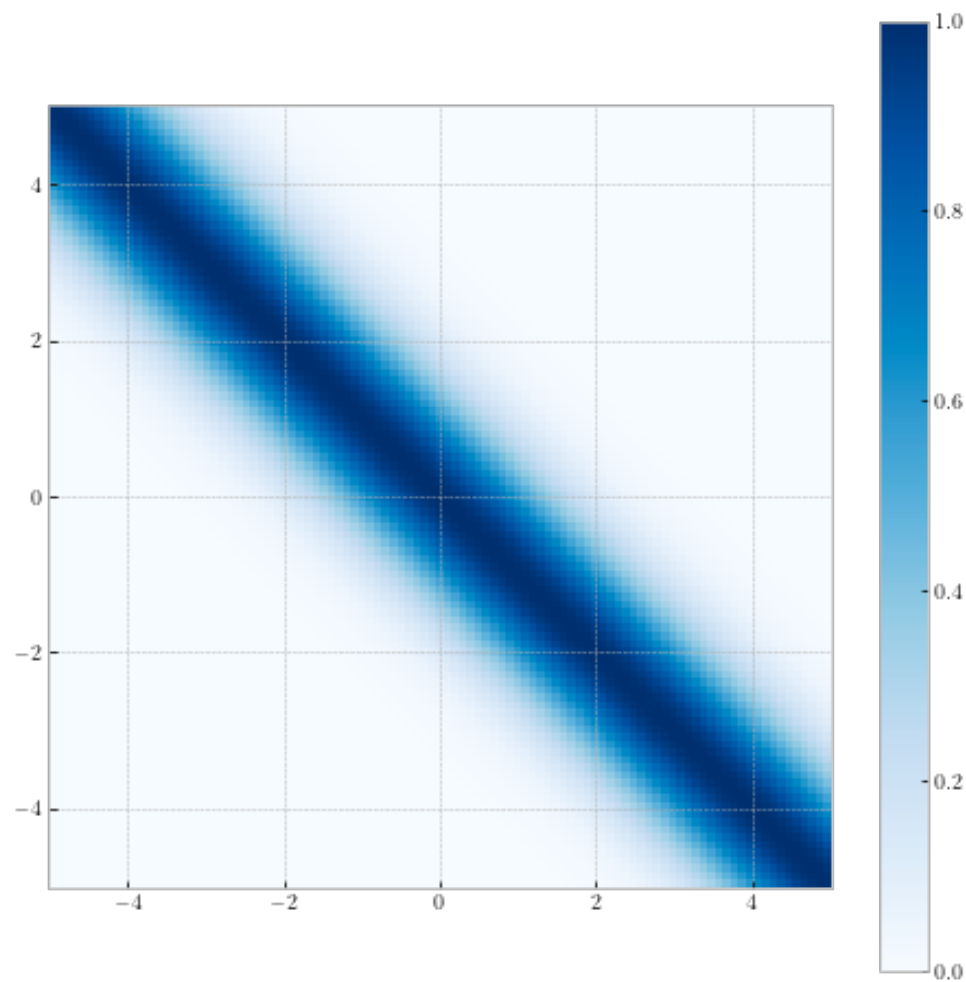
→ Tying a **knot**

RBF KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2} \right)$$

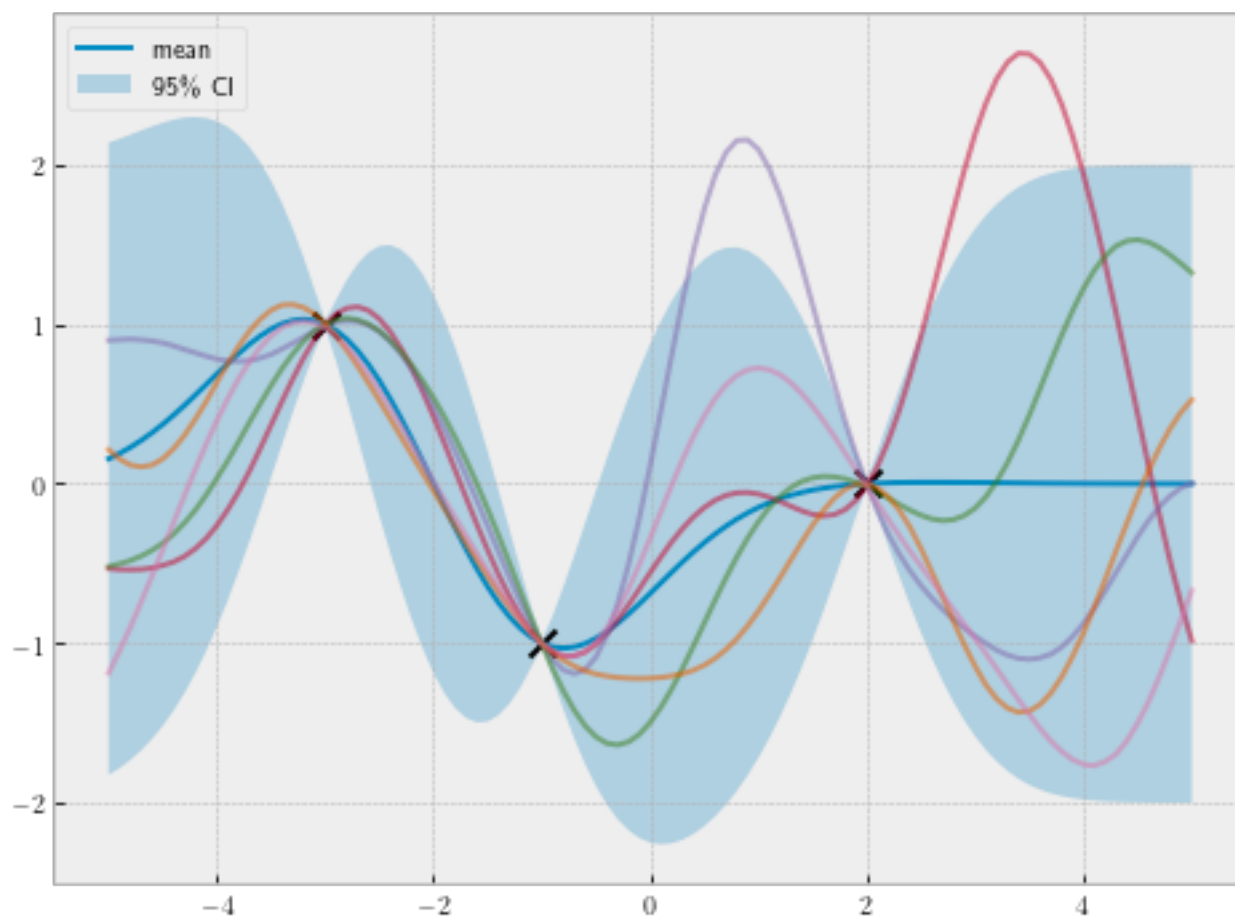
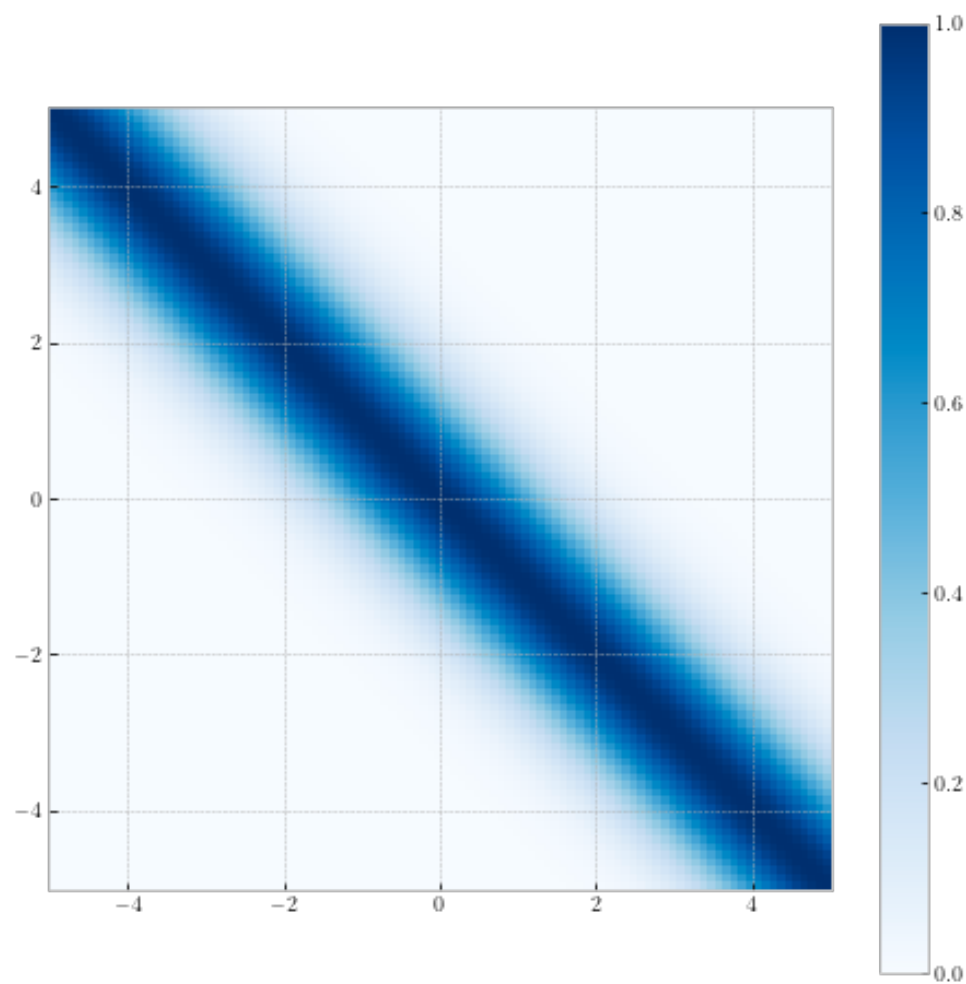
RBF KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$$



RBF KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$$

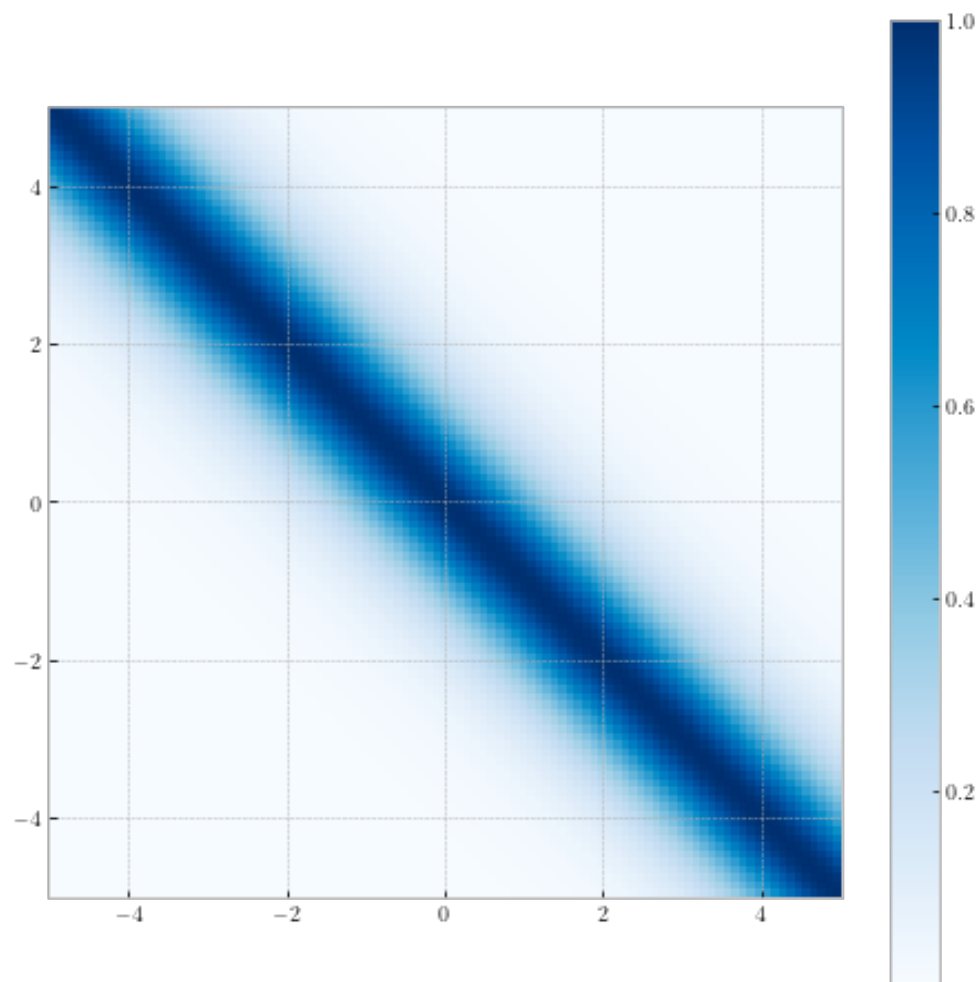


MATERN 5/2 KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{5}d + 5d^2\right) \exp\left(-\sqrt{5}d\right)$$

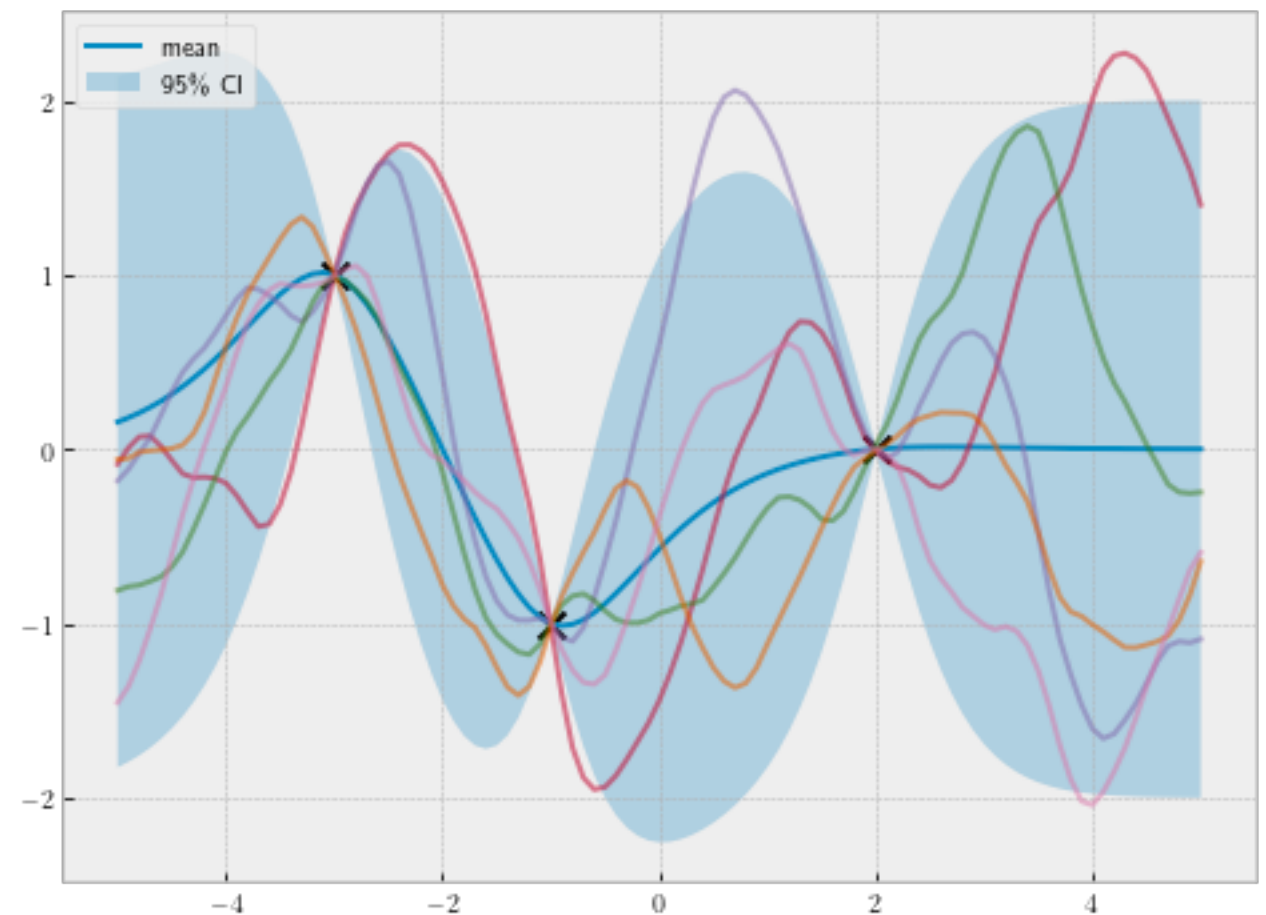
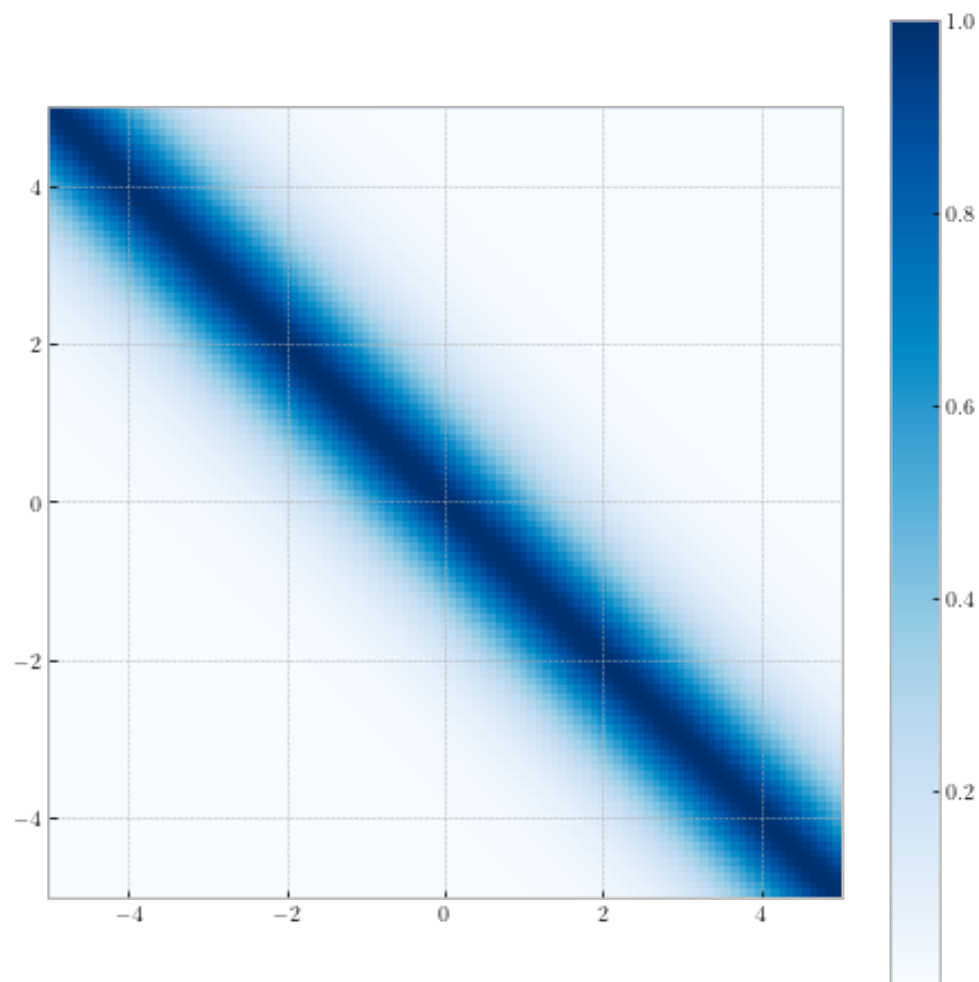
MATERN 5/2 KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{5}d + 5d^2\right) \exp\left(-\sqrt{5}d\right)$$



MATERN 5/2 KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{5}d + 5d^2\right) \exp\left(-\sqrt{5}d\right)$$

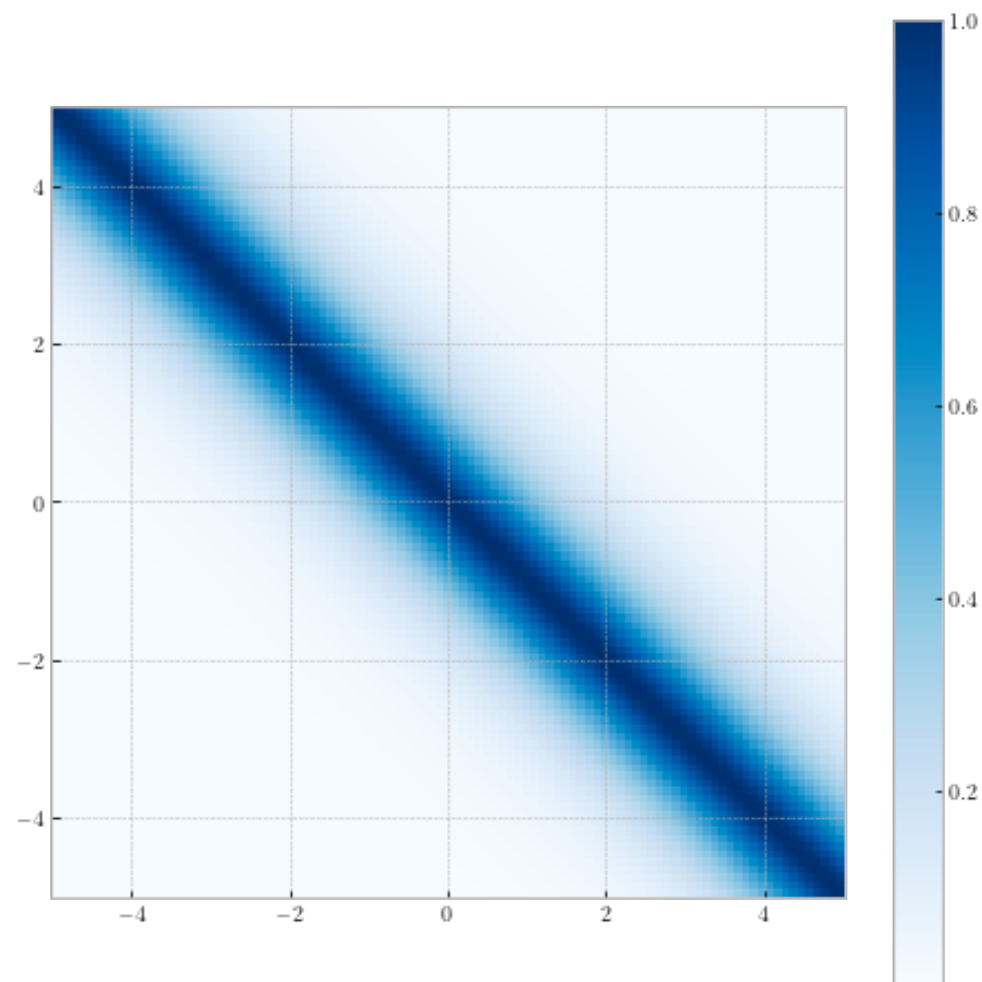


MATERN 3/2 KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{3}d\right) \exp\left(-\sqrt{3}d\right)$$

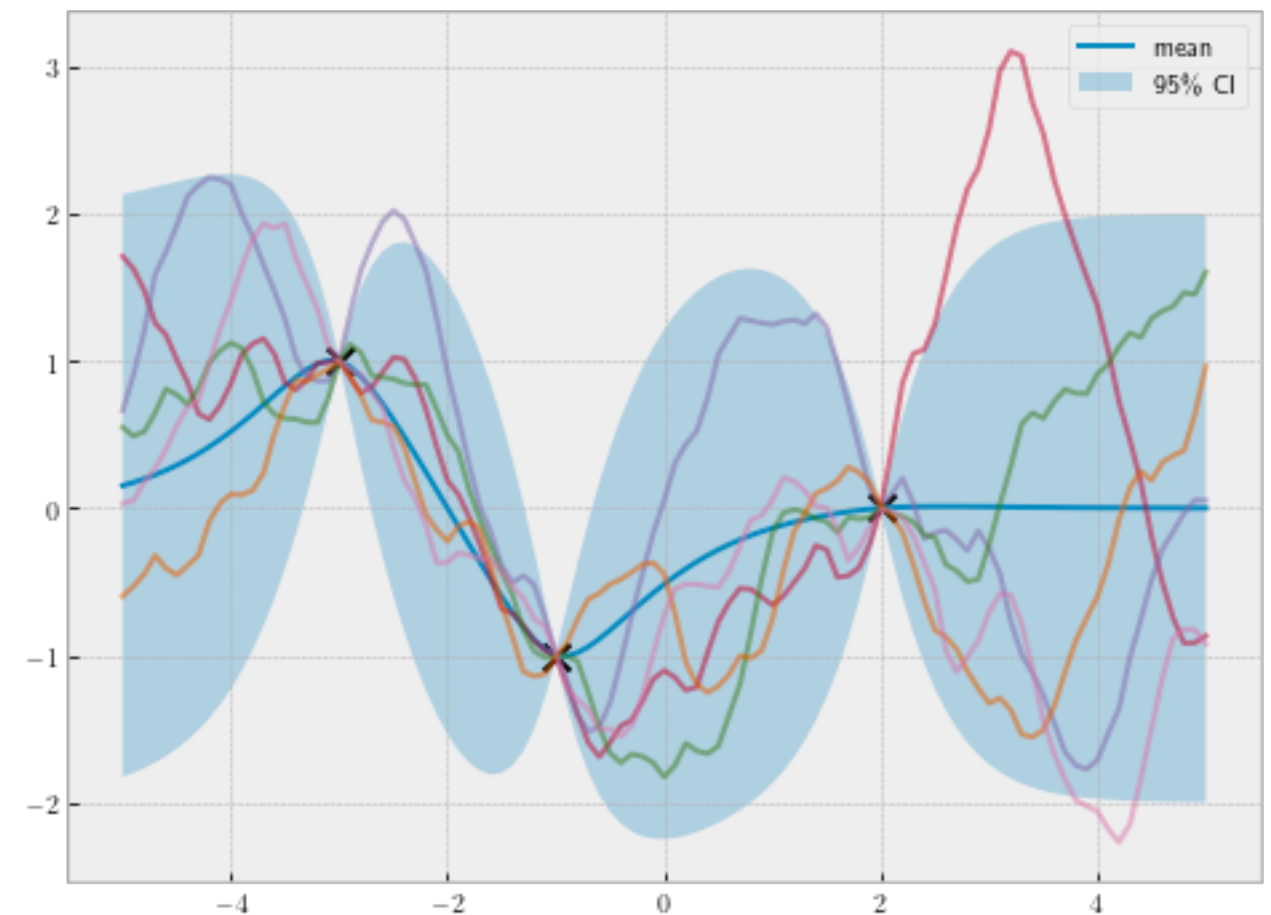
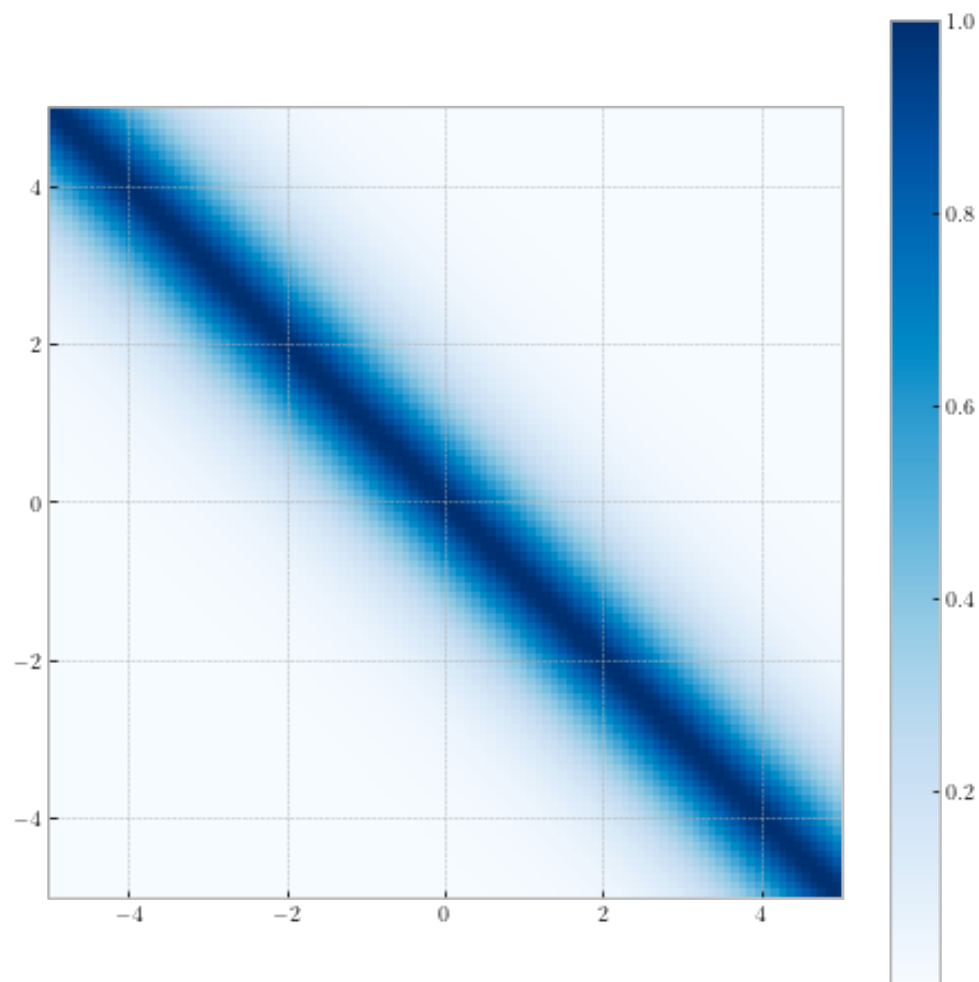
MATERN 3/2 KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{3}d\right) \exp\left(-\sqrt{3}d\right)$$



MATERN 3/2 KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{3}d\right) \exp\left(-\sqrt{3}d\right)$$

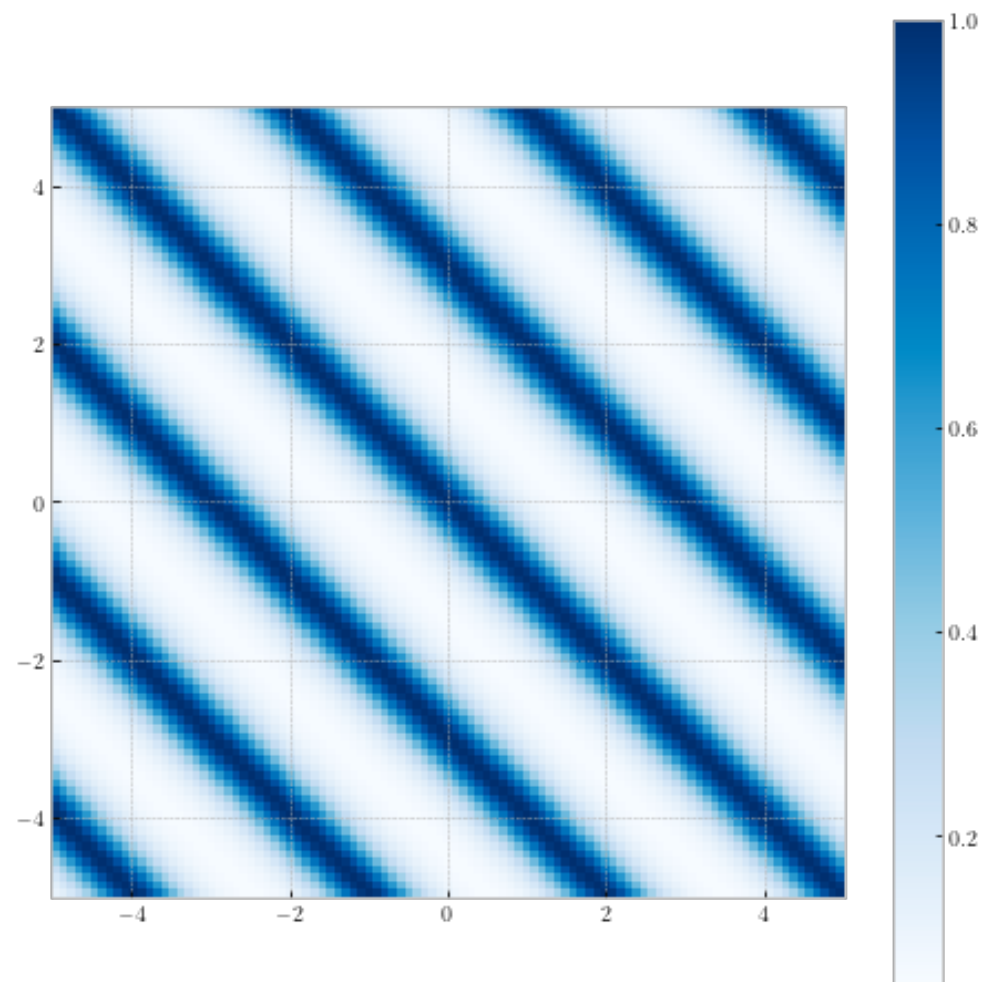


PERIODIC KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \exp \left(-2 \sum_i \sin^2(\pi(x_i - x'_i)) \right)$$

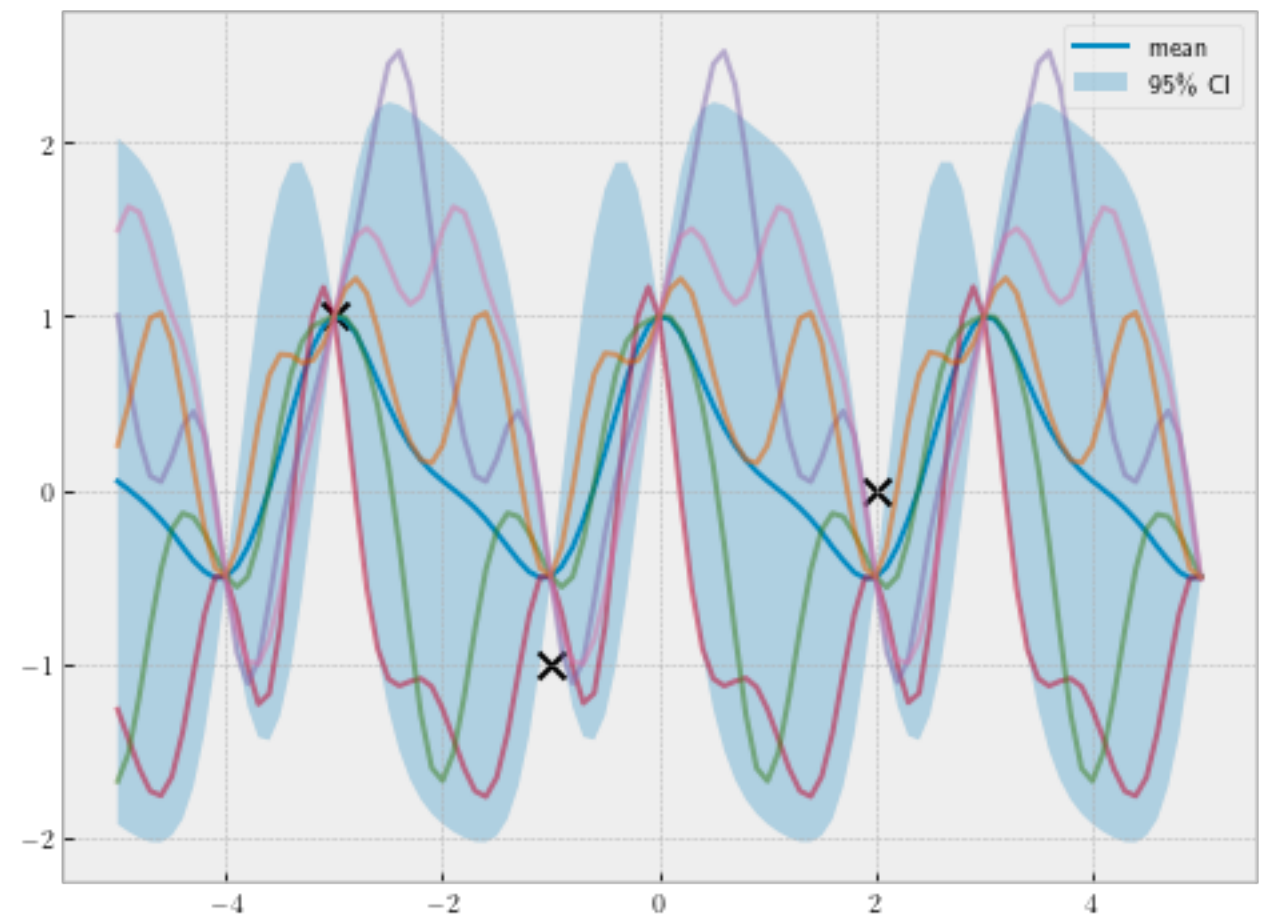
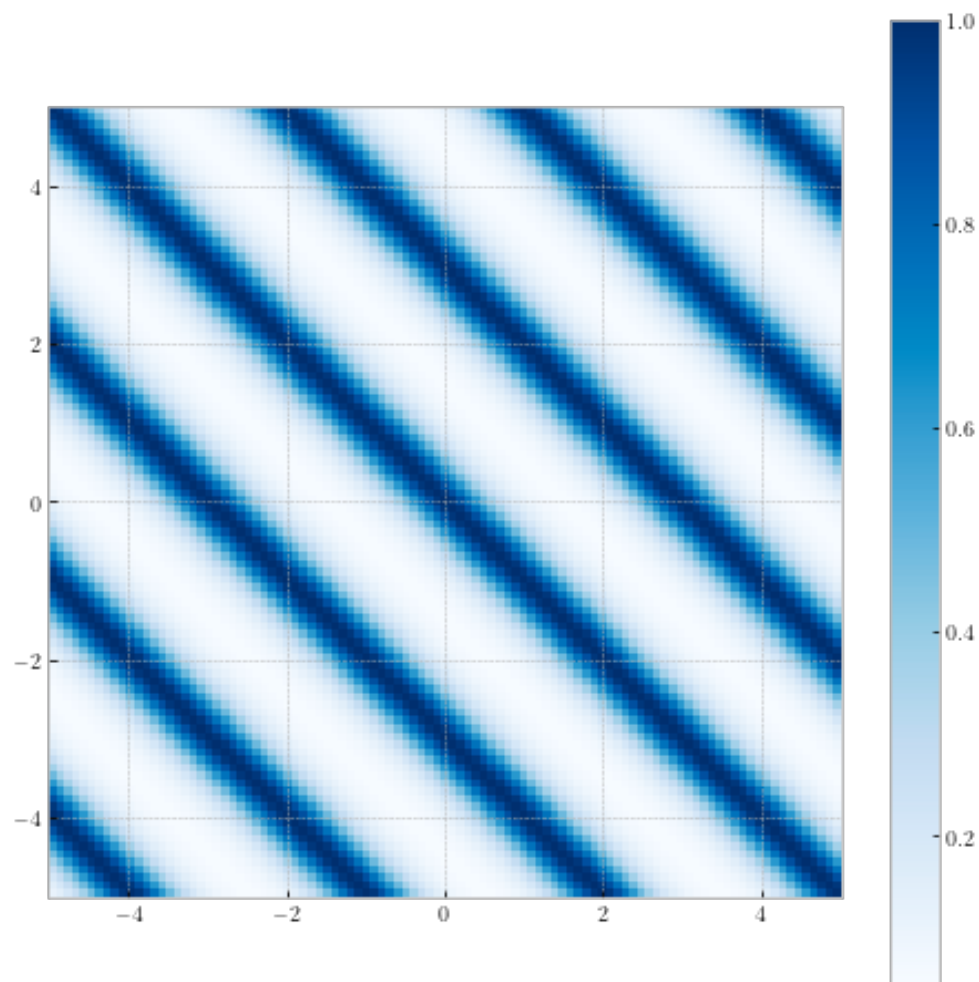
PERIODIC KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \exp \left(-2 \sum_i \sin^2(\pi(x_i - x'_i)) \right)$$



PERIODIC KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \exp \left(-2 \sum_i \sin^2 \left(\pi(x_i - x'_i) \right) \right)$$

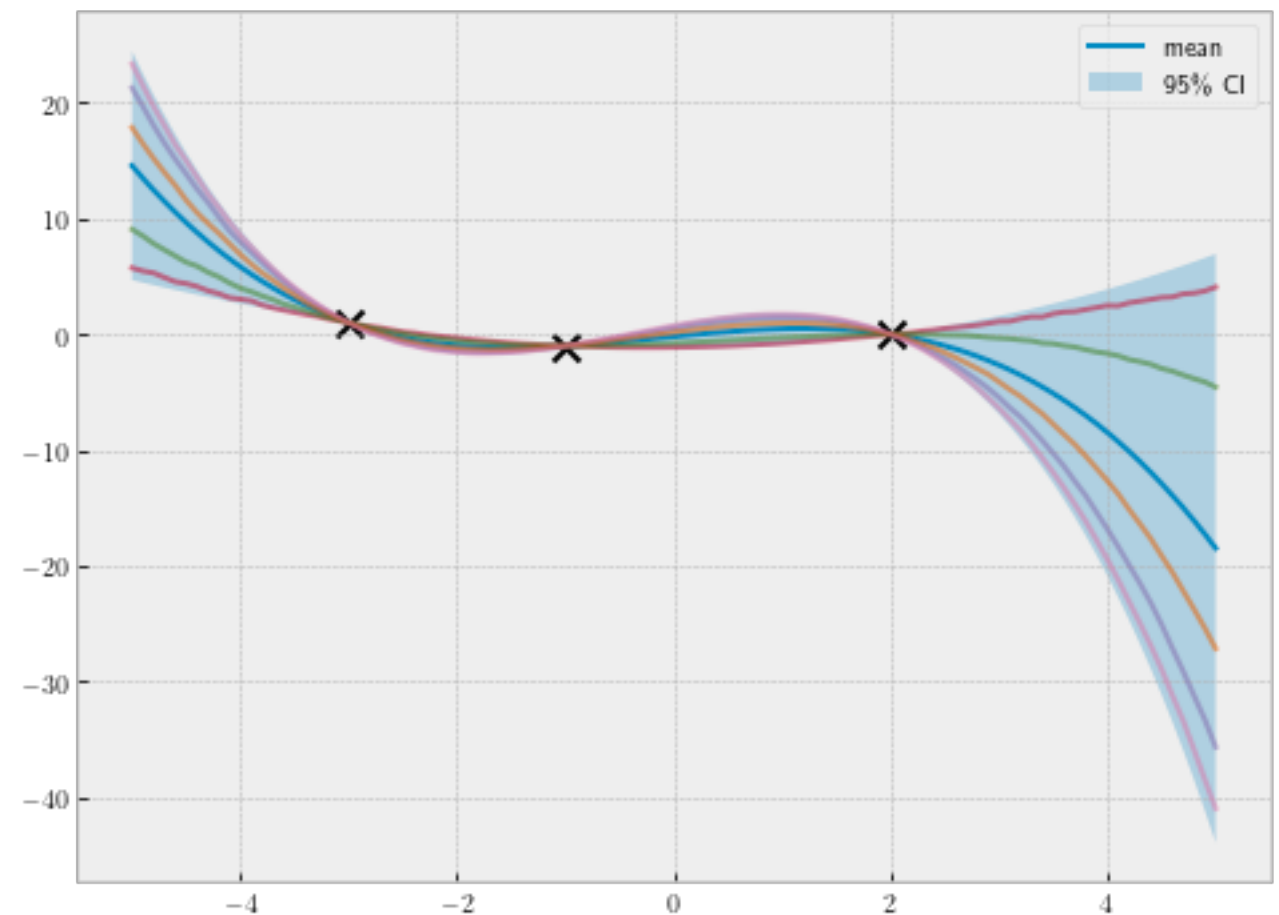
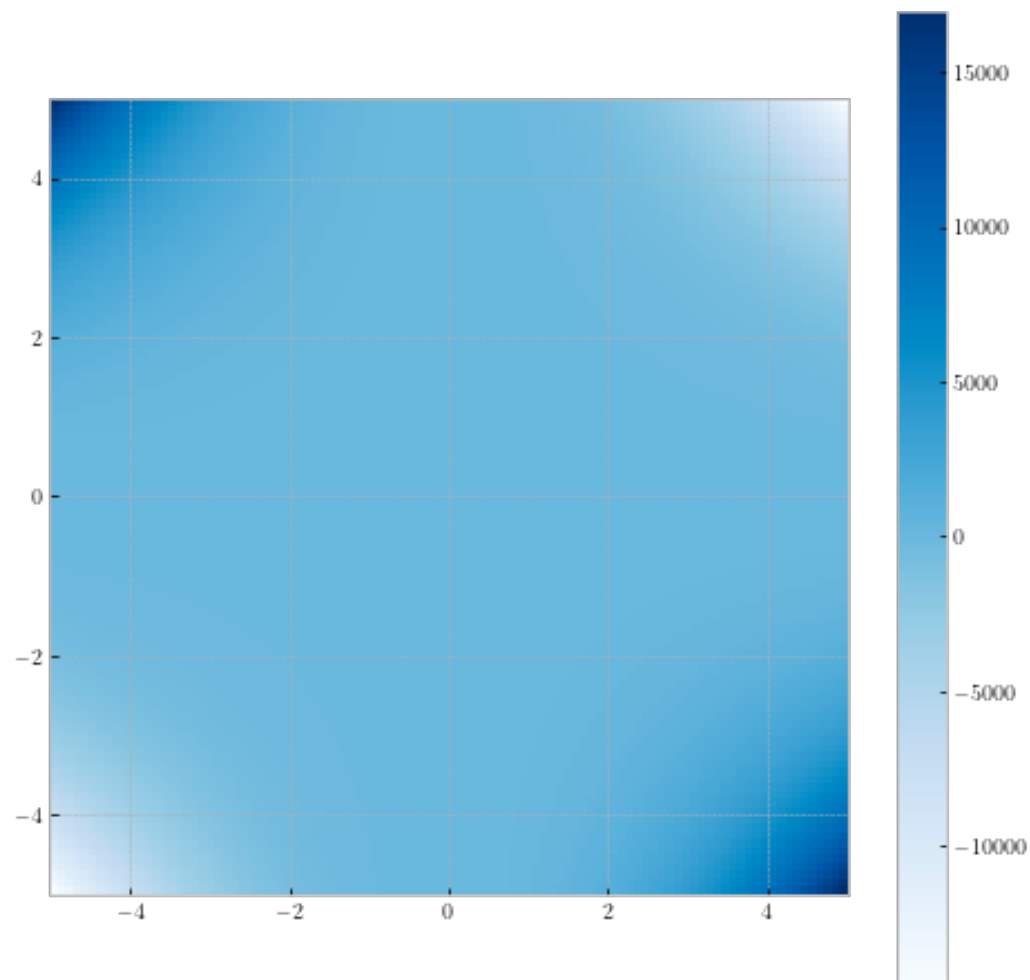


POLYNOMIAL KERNEL

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^d \text{ of degree } d$$

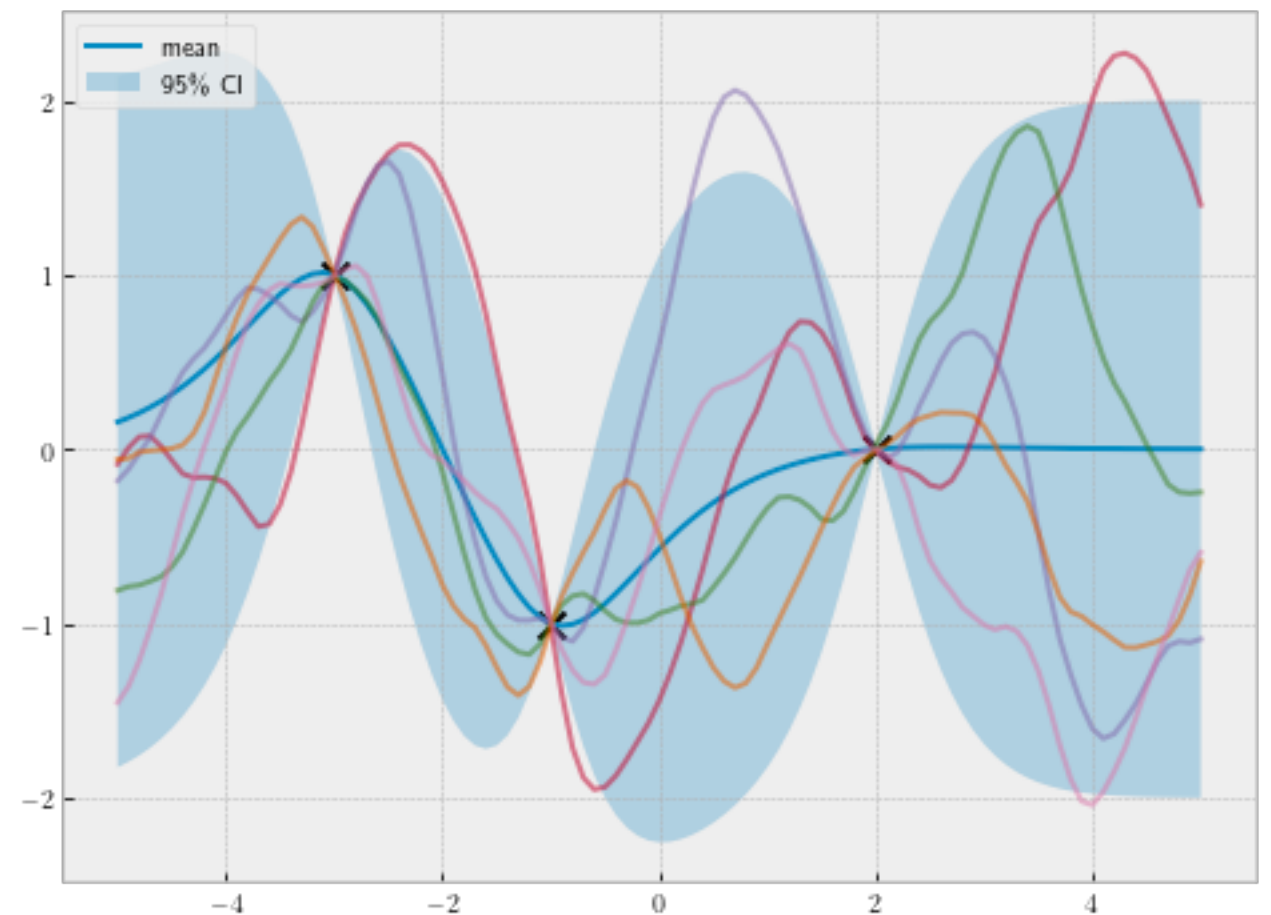
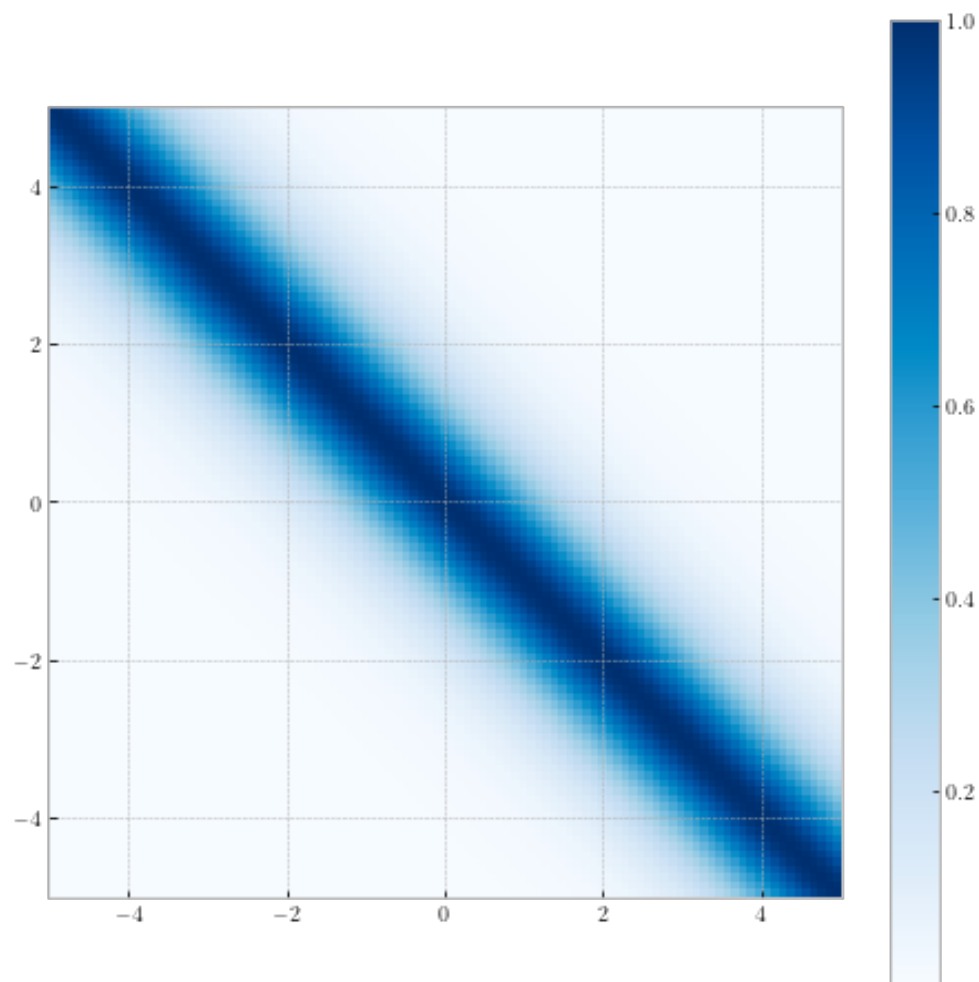
POLYNOMIAL KERNEL

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^d \text{ of degree } d$$



MATERN 5/2 KERNEL

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{5}d + 5d^2\right) \exp\left(-\sqrt{5}d\right)$$



GP HYPERPARAMETERS

GP HYPERPARAMETERS

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell} \right)$$

GP HYPERPARAMETERS

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell} \right)$$

Length scale ℓ controls how “wiggly” the function is

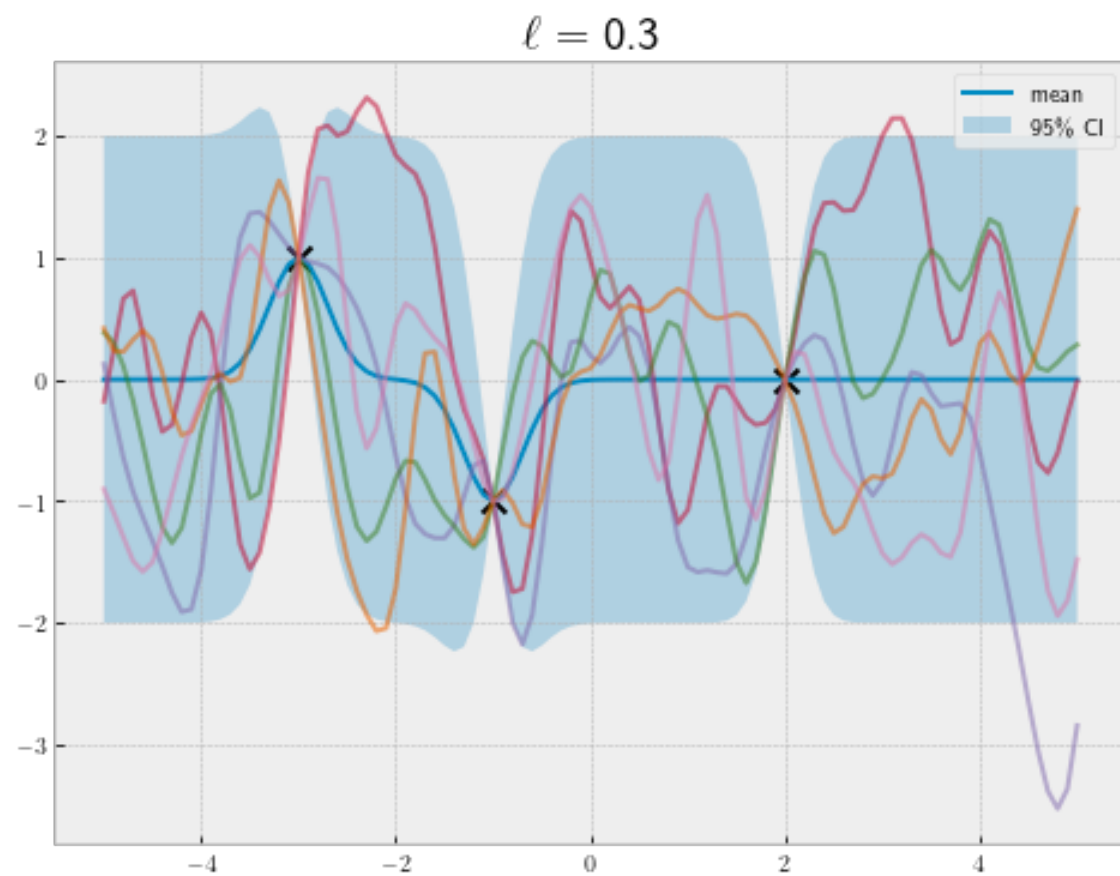
Output scale σ^2 controls the range of the function

LENGTH SCALE CONTROLS WIGGLINESS

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell} \right)$$

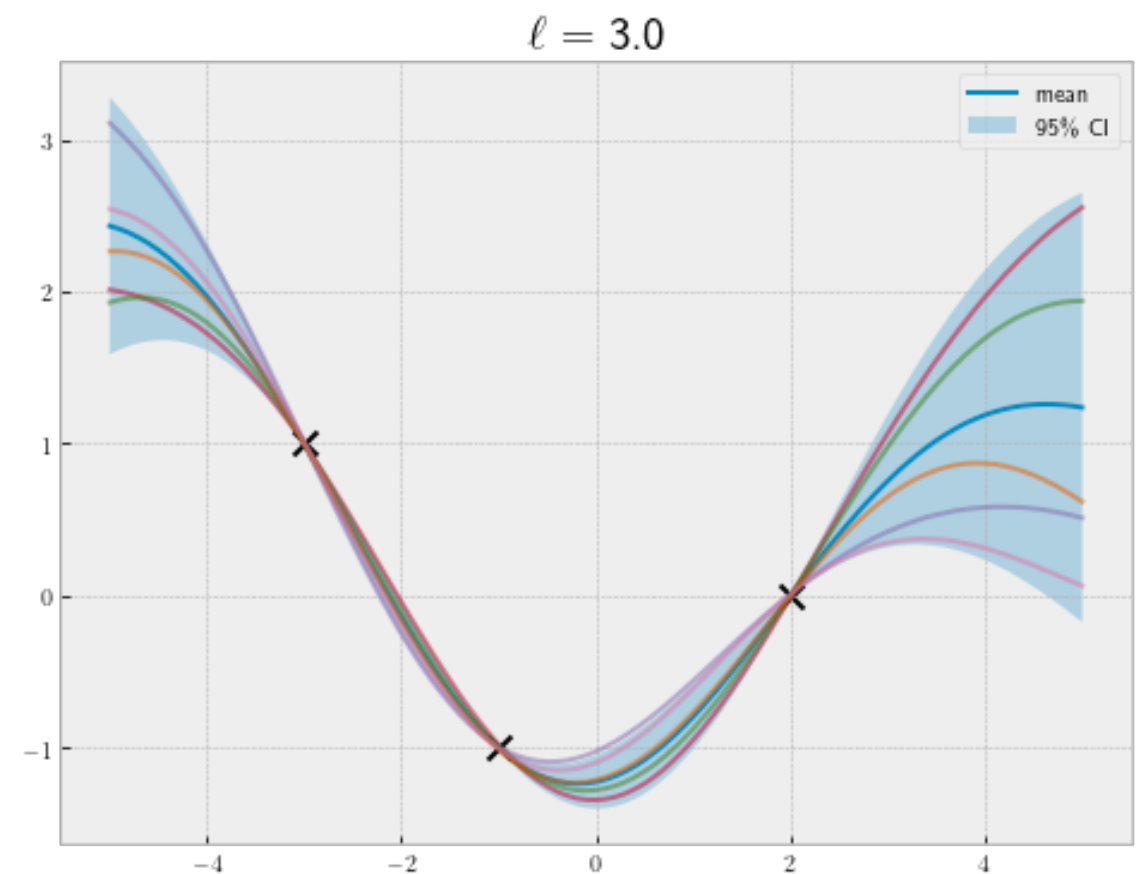
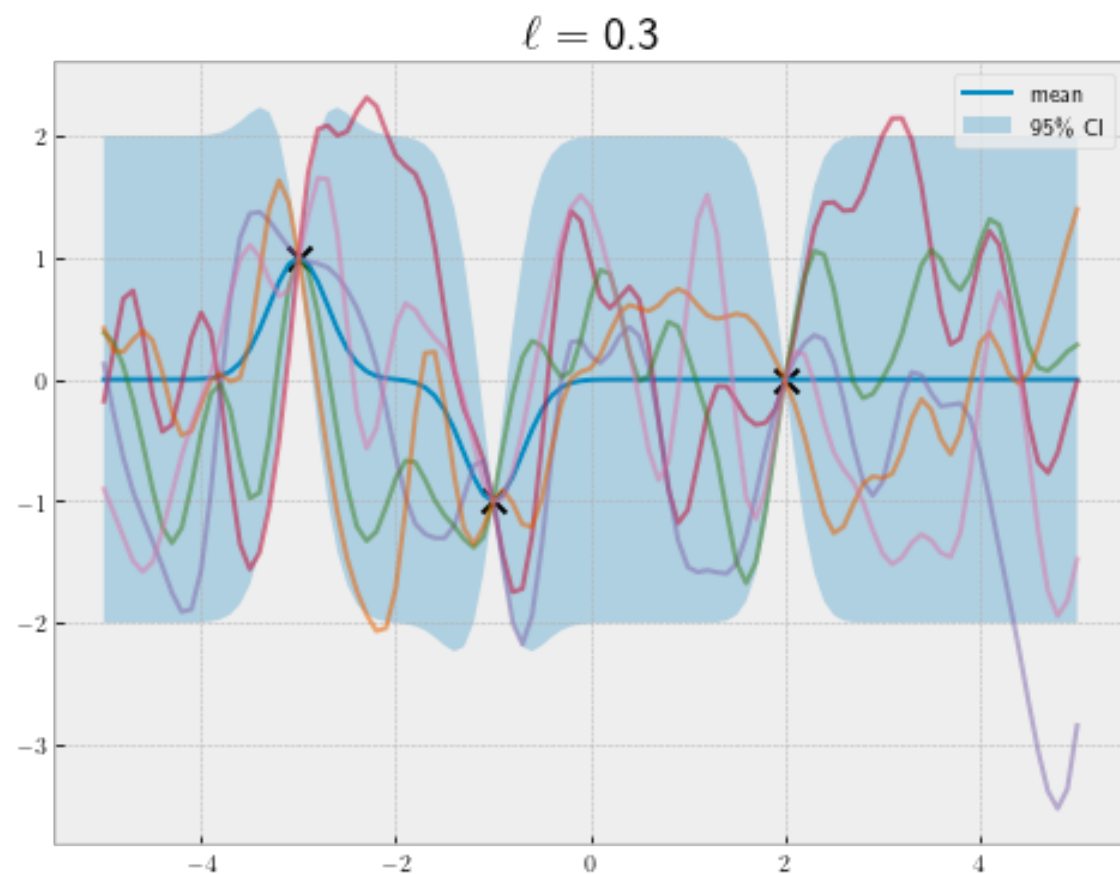
LENGTH SCALE CONTROLS WIGGLINESS

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell} \right)$$



LENGTH SCALE CONTROLS WIGGLINESS

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell}\right)$$

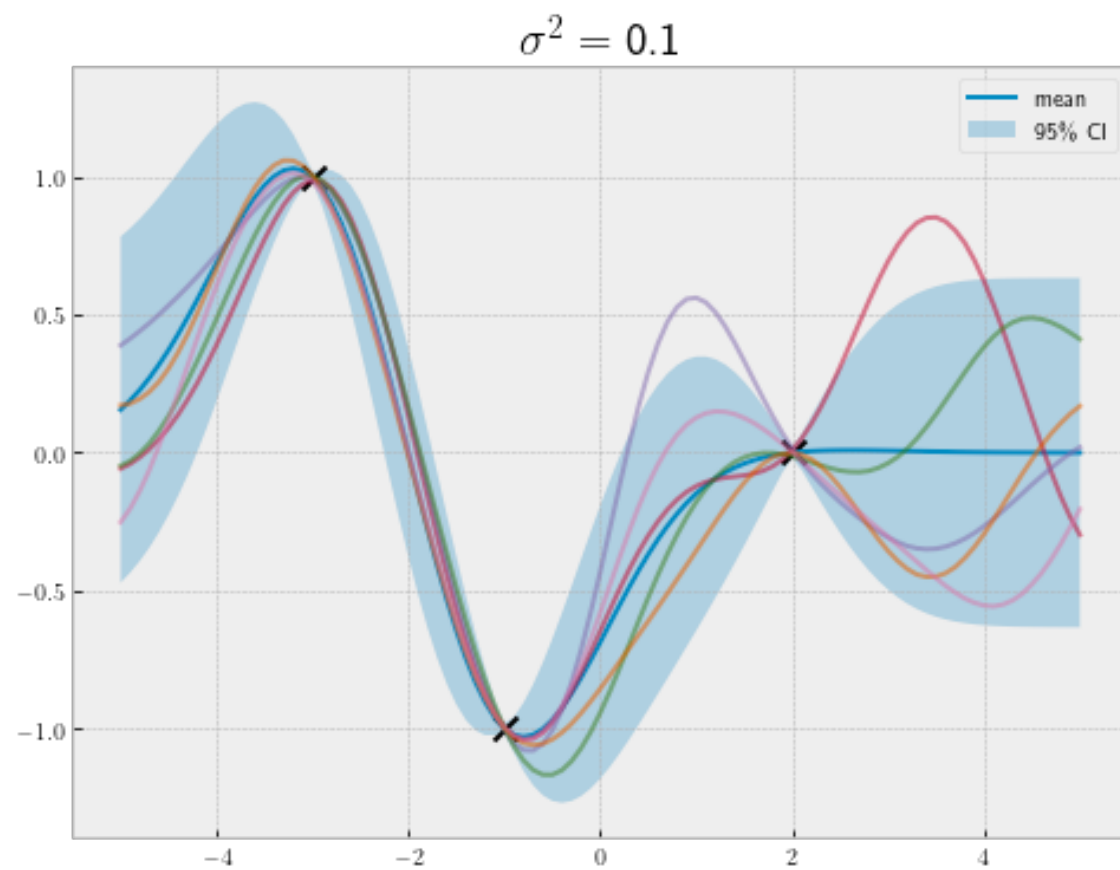


OUTPUT SCALE CONTROLS RANGE

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell} \right)$$

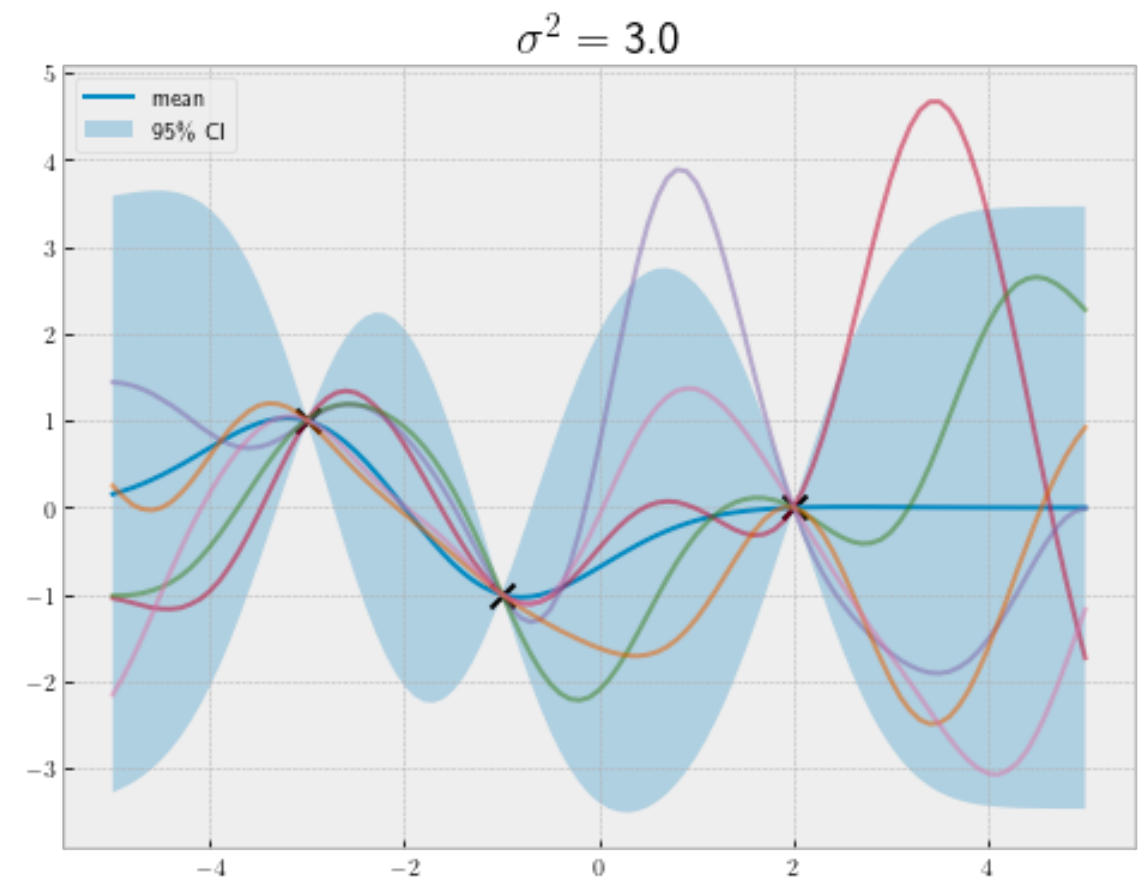
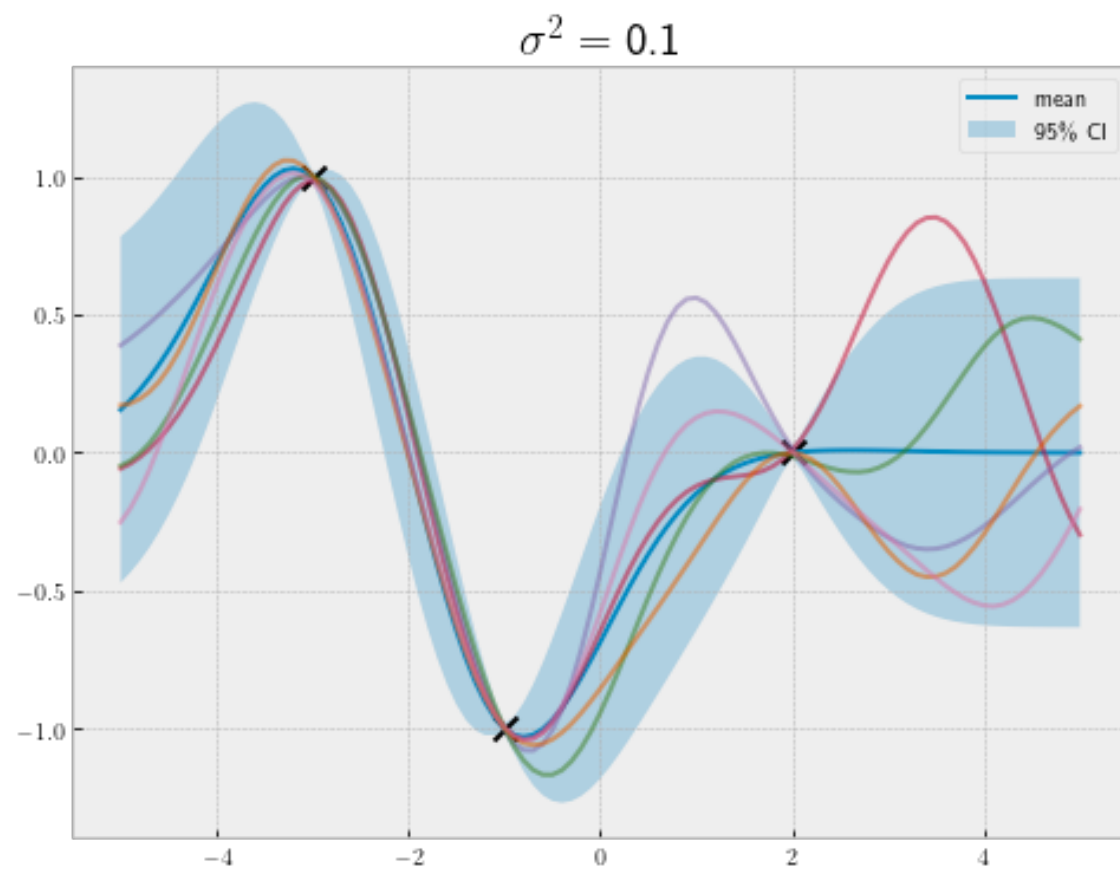
OUTPUT SCALE CONTROLS RANGE

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell}\right)$$



OUTPUT SCALE CONTROLS RANGE

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell}\right)$$



FINDING THE BEST HYPERPARAMETERS

FINDING THE BEST HYPERPARAMETERS

Given the hyperparameter vector θ , consider the **negative**

log marginal likelihood:

FINDING THE BEST HYPERPARAMETERS

Given the hyperparameter vector θ , consider the **negative**

log marginal likelihood:

$$-\log p(\mathbf{y} \mid \theta) = \frac{1}{2} \log \det K(\mathbf{X}, \mathbf{X}) + \frac{1}{2} \mathbf{y}^\top K^{-1}(\mathbf{X}, \mathbf{X}) \mathbf{y} + \frac{N}{2} \log(2\pi)$$

FINDING THE BEST HYPERPARAMETERS

Given the hyperparameter vector θ , consider the **negative**

log marginal likelihood:

$$-\log p(\mathbf{y} \mid \theta) = \frac{1}{2} \log \det K(X, X) + \frac{1}{2} \mathbf{y}^\top K^{-1}(X, X) \mathbf{y} + \frac{N}{2} \log(2\pi)$$

- ▶ Closed-form gradient w.r.t. elements in θ .

ONE LENGTH SCALE FOR EACH DIMENSION

ONE LENGTH SCALE FOR EACH DIMENSION

Automatic relevance determination (ARD) covariance function

ONE LENGTH SCALE FOR EACH DIMENSION

Automatic relevance determination (ARD) covariance function

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{1}{2} \sum_d \frac{(x_d - x'_d)^2}{\ell_d^2} \right)$$

ONE LENGTH SCALE FOR EACH DIMENSION

Automatic relevance determination (ARD) covariance function

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{1}{2} \sum_d \frac{(x_d - x'_d)^2}{\ell_d} \right)$$

Length scale ℓ_d for dimension d determines how relevant the dimension is

ONE LENGTH SCALE FOR EACH DIMENSION

Automatic relevance determination (ARD) covariance function

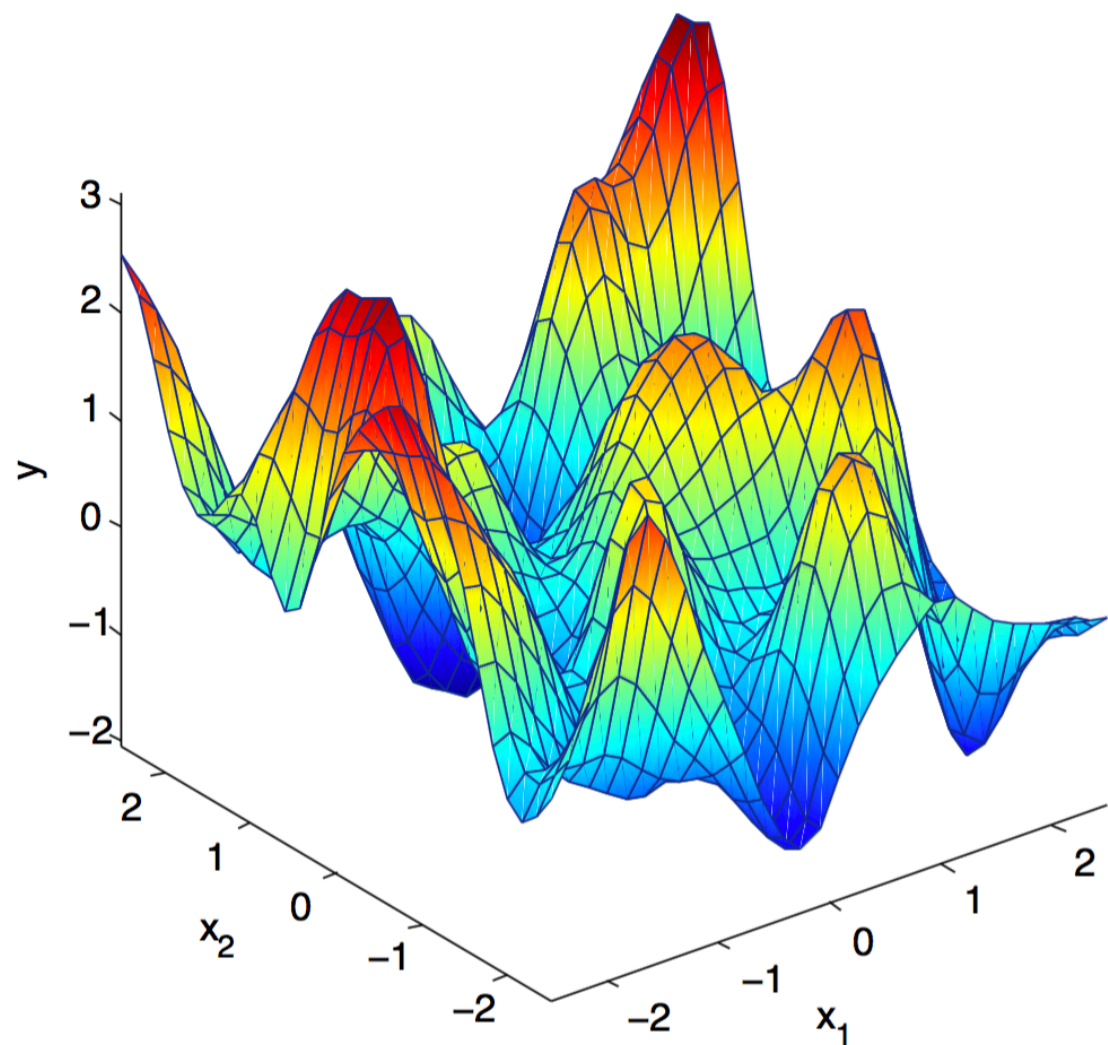
$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{1}{2} \sum_d \frac{(x_d - x'_d)^2}{\ell_d} \right)$$

Length scale ℓ_d for dimension d determines how relevant the dimension is

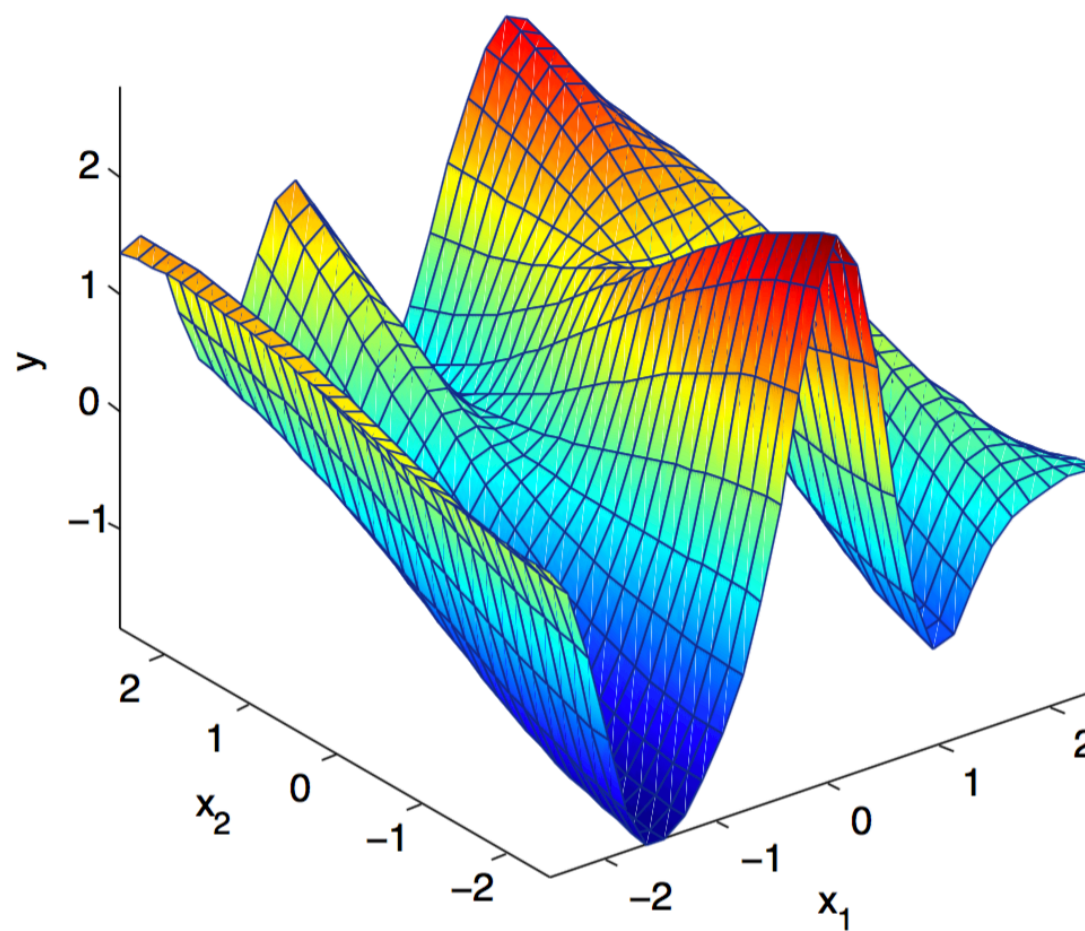
- As ℓ_d gets **larger**, dimension d becomes **less relevant**

ARD IN ACTION

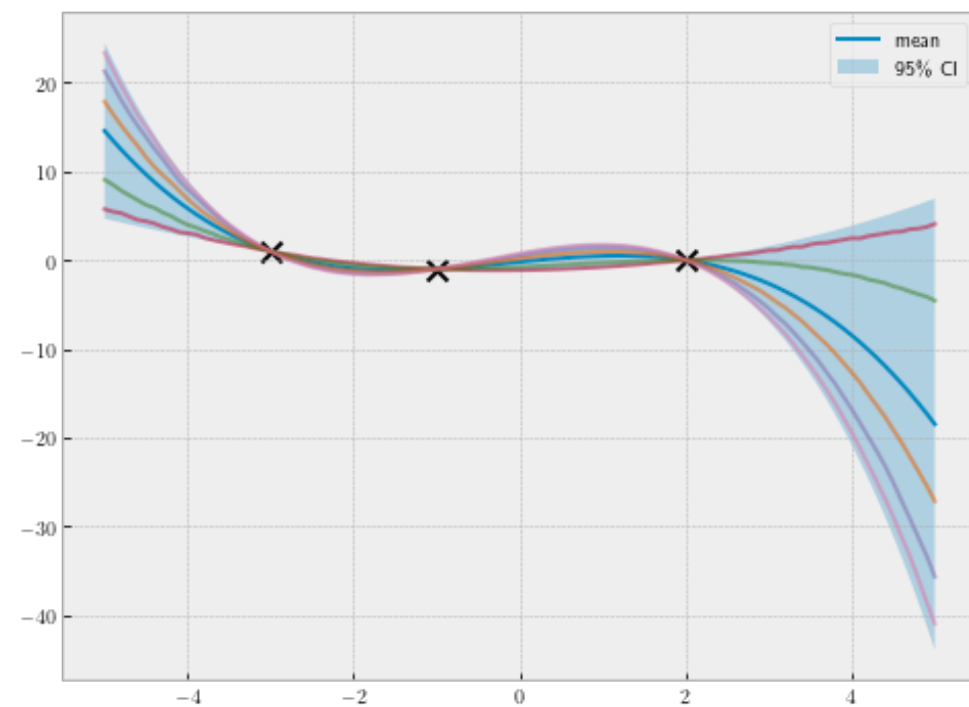
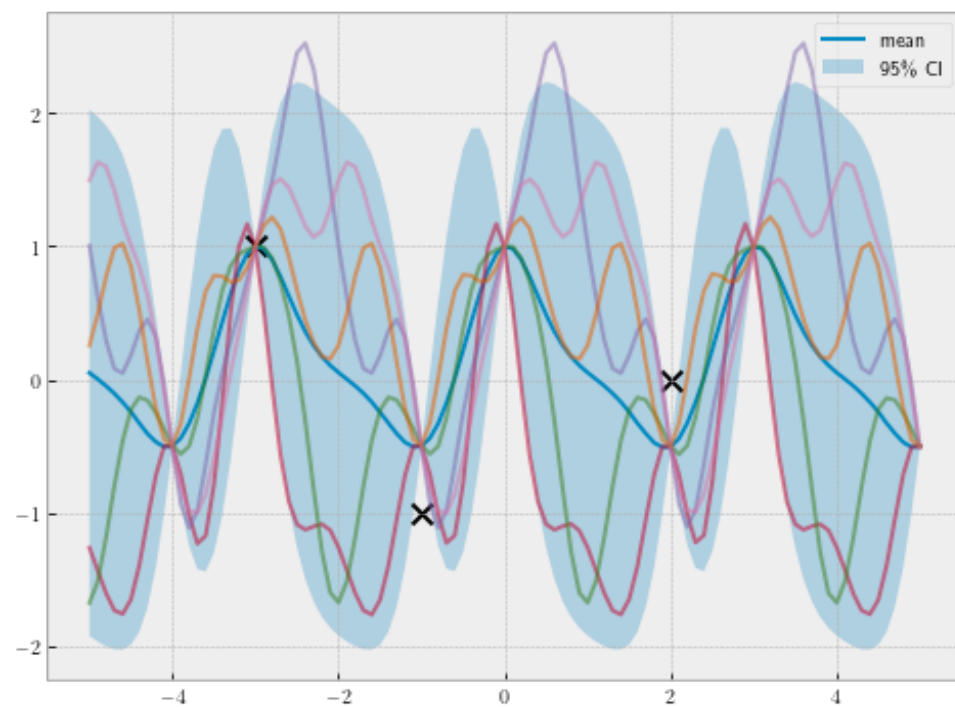
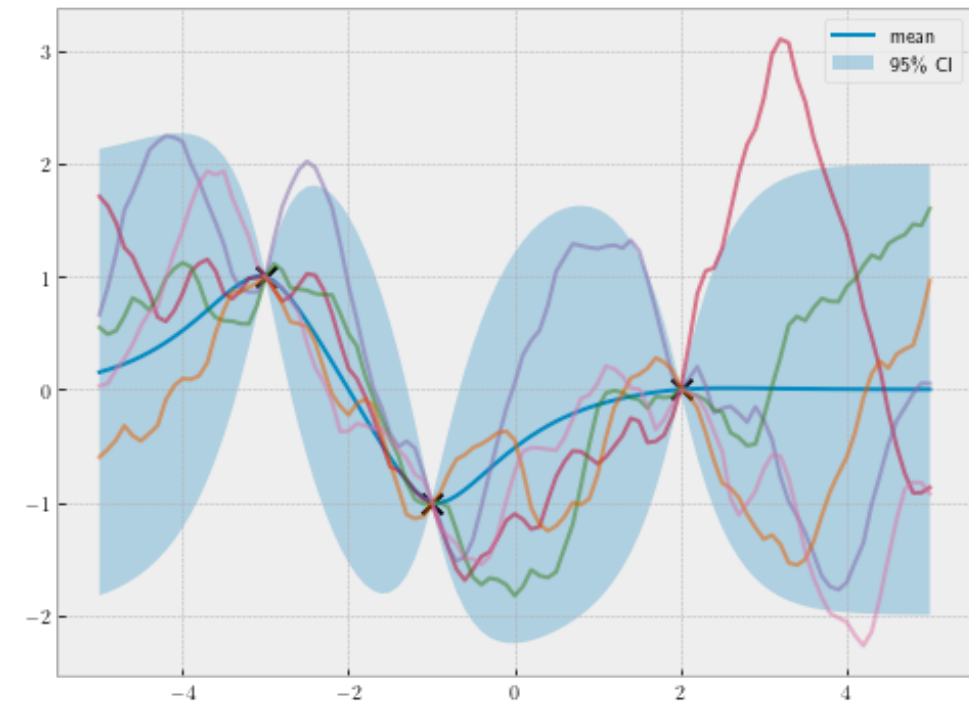
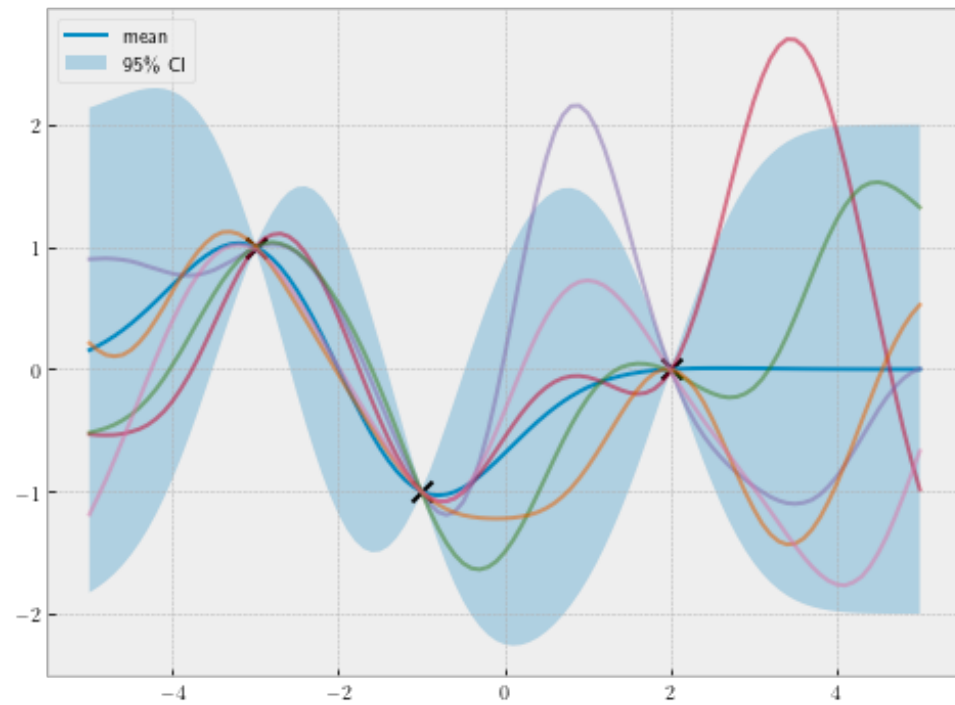
$$l_1 = 1.0, l_2 = 1.0$$



$$l_1 = 1.0, l_2 = 3.0$$



GP KERNELS



MODELING MULTIPLE TRENDS SIMULTANEOUSLY

CO₂ concentration

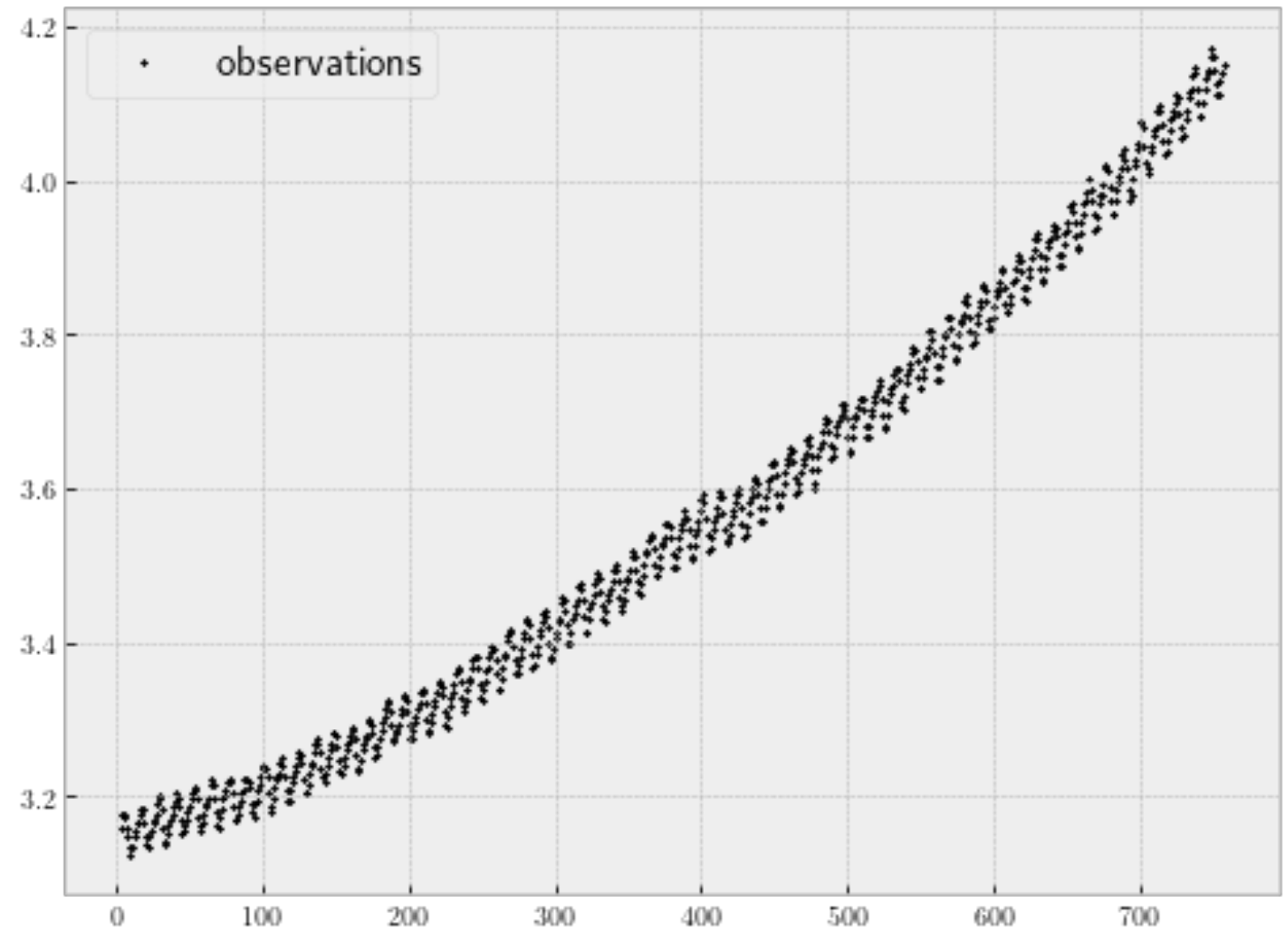
by **month**:

- ▶ Long-term rising trend
- ▶ Seasonal changes

MODELING MULTIPLE TRENDS SIMULTANEOUSLY

CO₂ concentration
by **month**:

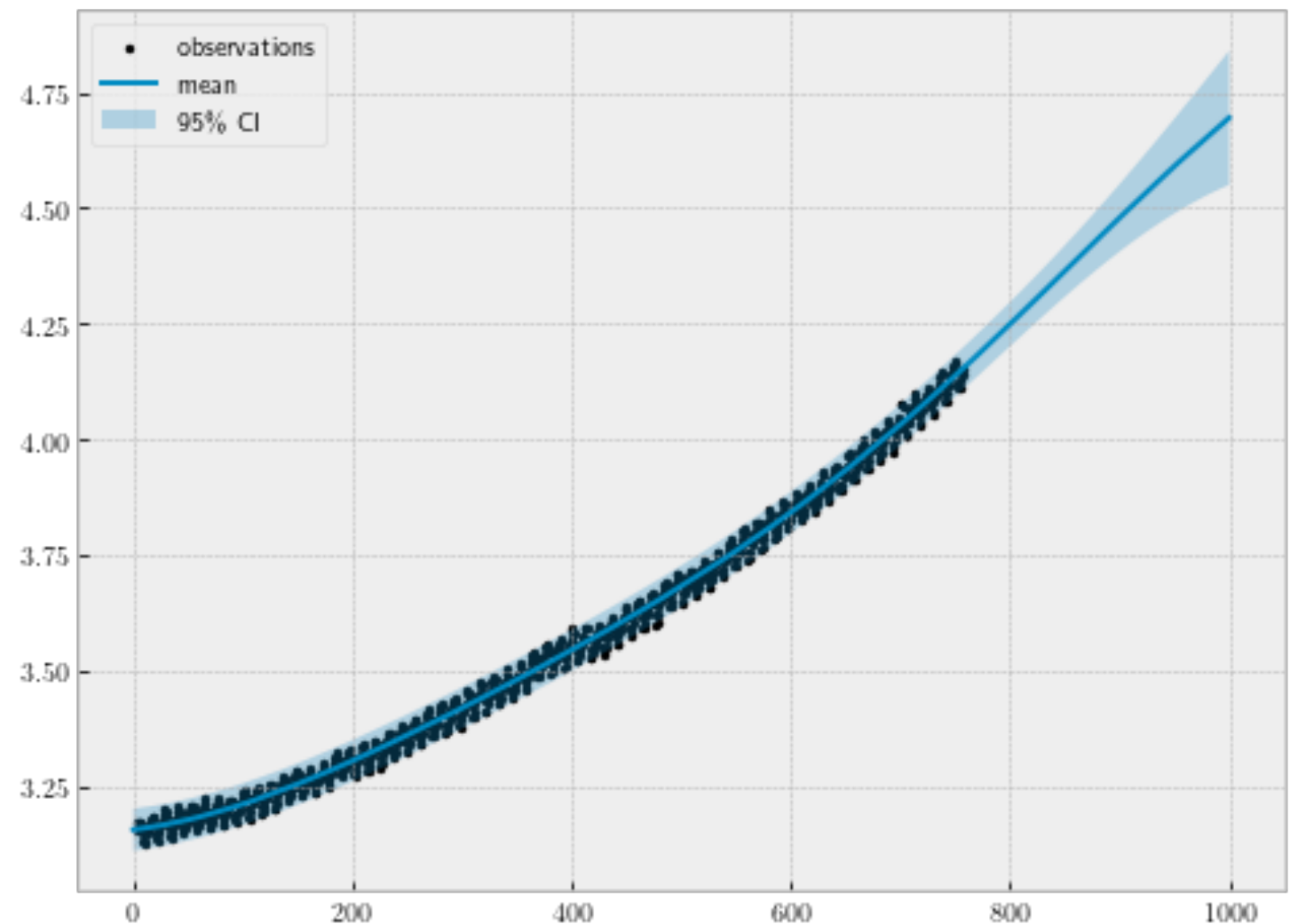
- ▶ Long-term rising trend
- ▶ Seasonal changes



MODELING MULTIPLE TRENDS SIMULTANEOUSLY

CO₂ concentration
by **month**:

- ▶ Long-term rising trend
- ▶ Seasonal changes

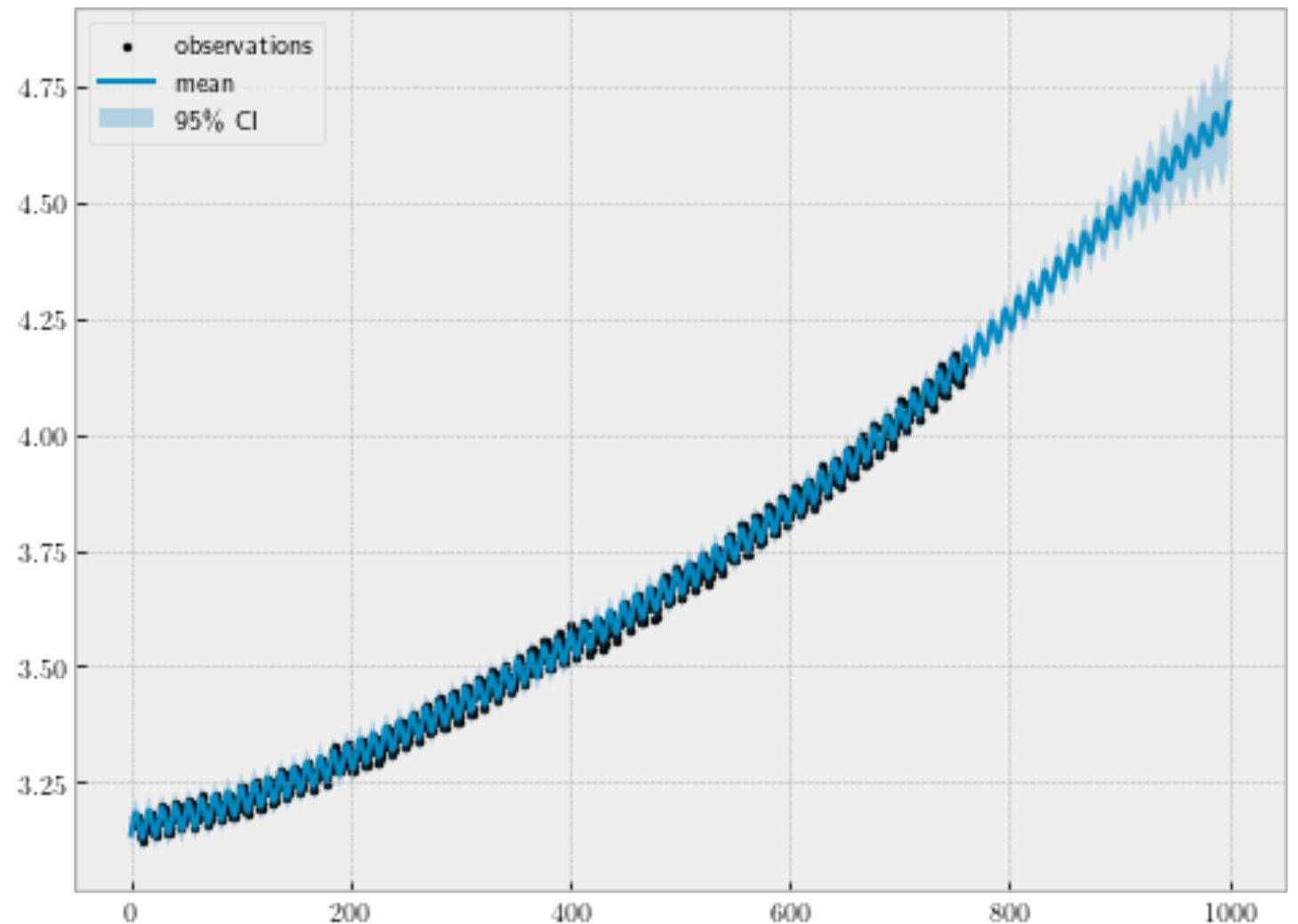


$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell} \right) + \sigma_n^2$$

MODELING MULTIPLE TRENDS SIMULTANEOUSLY

CO₂ concentration
by **month**:

- ▶ Long-term rising trend
- ▶ Seasonal changes



$$K(\mathbf{x}, \mathbf{x}') = \sigma_1^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell}\right) + \sigma_2^2 \cos\left(\pi \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{p}\right) + \sigma_n^2$$

DON'T TRY THIS IN NC!

KERNEL CONSTRUCTION

KERNEL CONSTRUCTION

- ▶ Kernel grammar

KERNEL CONSTRUCTION

- ▶ Kernel grammar
 - ▶ Adding, multiplying, exponentiating, etc.

KERNEL CONSTRUCTION

- ▶ Kernel grammar
 - ▶ Adding, multiplying, exponentiating, etc.
 - ▶ How to find the **best** “formula”?: AutoML

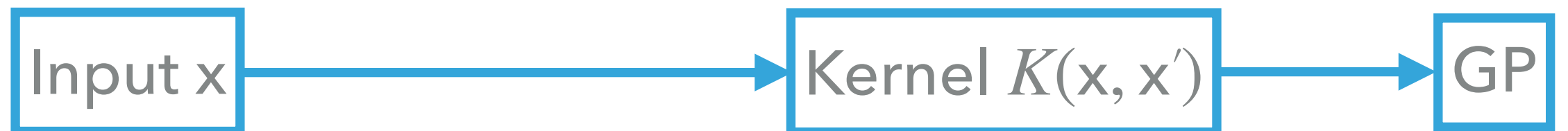
KERNEL CONSTRUCTION

- ▶ Kernel grammar
 - ▶ Adding, multiplying, exponentiating, etc.
 - ▶ How to find the **best** “formula”?: AutoML
- ▶ Deep kernel learning

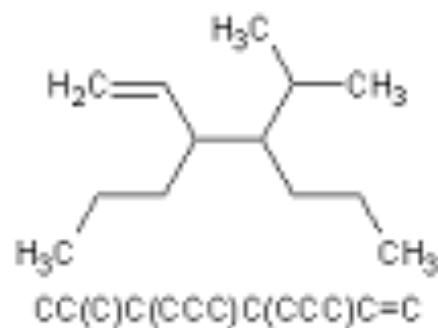
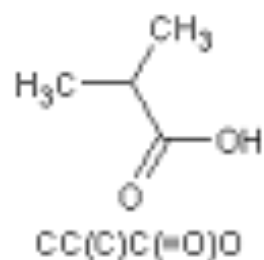
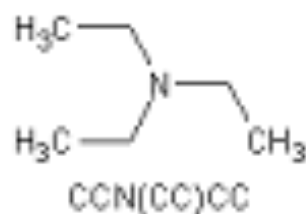
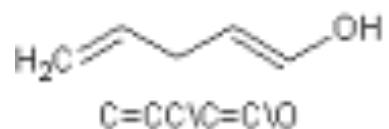
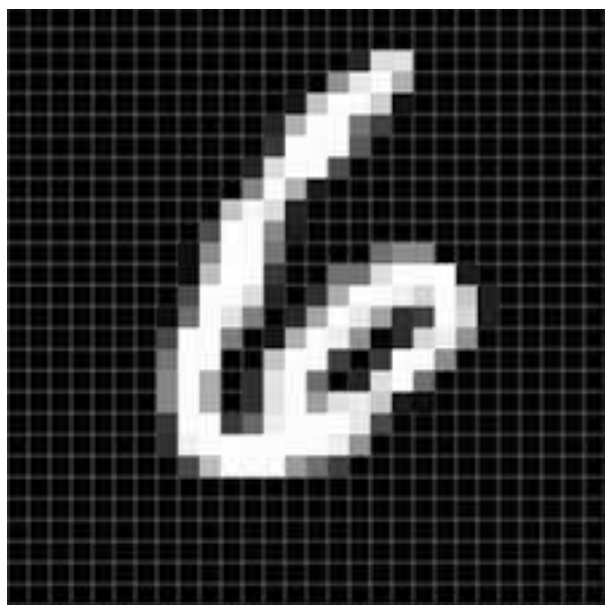
KERNEL CONSTRUCTION

- ▶ Kernel grammar
 - ▶ Adding, multiplying, exponentiating, etc.
 - ▶ How to find the **best** "formula"?: AutoML
- ▶ Deep kernel learning
 - ▶ Flexible, good for **structured** data

DEEP KERNEL LEARNING

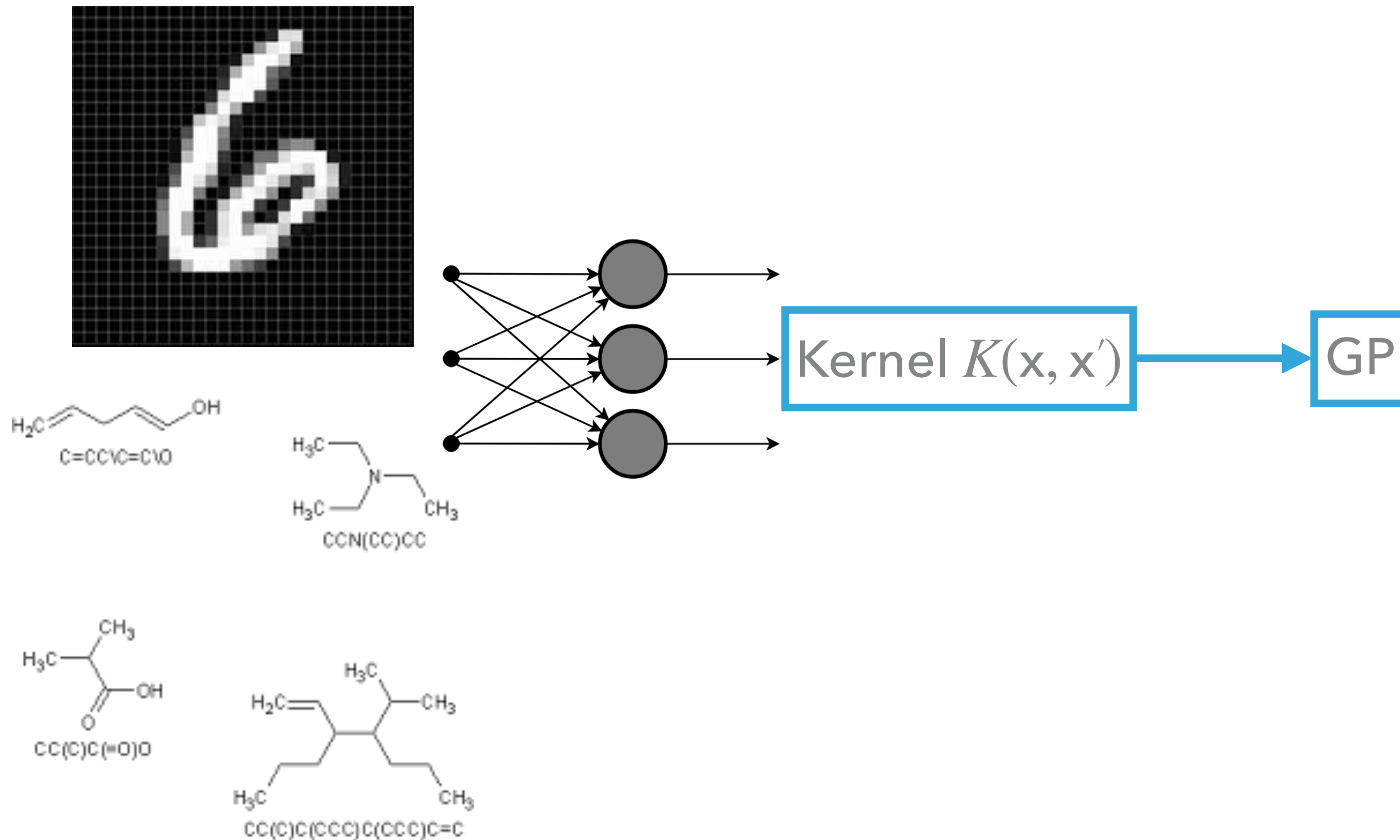


DEEP KERNEL LEARNING

Kernel $K(x, x')$

GP

DEEP KERNEL LEARNING



BAYESIAN OPTIMIZATION

BAYESIAN OPTIMIZATION

Expensive, blackbox optimization problems are common

BAYESIAN OPTIMIZATION

Expensive, blackbox optimization problems are common

- ▶ No known functional form

BAYESIAN OPTIMIZATION

Expensive, blackbox optimization problems are common

- ▶ No known functional form
- ▶ Cost associated with querying the function

BAYESIAN OPTIMIZATION

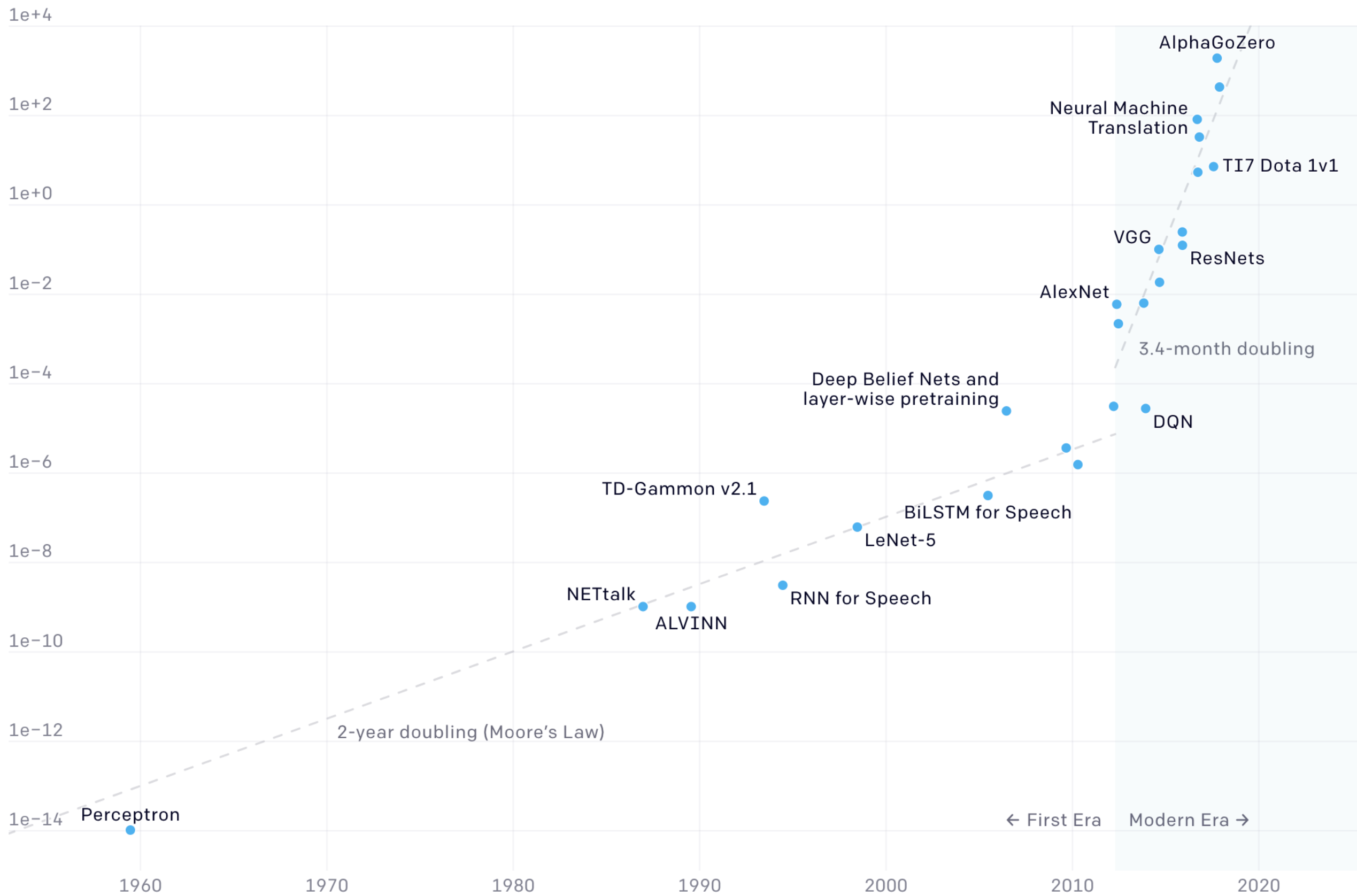
Expensive, blackbox optimization problems are common

- ▶ No known functional form
- ▶ Cost associated with querying the function
- ▶ No derivative information

BAYESIAN OPTIMIZATION EXAMPLES

Hyperparameter tuning, user's preference optimization, drug discovery, etc.

Petaflop/s-days



Bayesian Optimization for a Better Dessert

**Greg Kochanski, Daniel Golovin, John Karro, Benjamin Solnik,
Subhodeep Moitra, and D. Sculley**

{gpk, dgg, karro, bsolnik, smoitra, dsculley}@google.com; Google Brain Team

Abstract

We present a case study on applying Bayesian Optimization to a complex real-world system; our challenge was to optimize chocolate chip cookies. The process was a mixed-initiative system where both human chefs, human raters, and a machine optimizer participated in 144 experiments. This process resulted in highly rated cookies that deviated from expectations in some surprising ways – much less sugar in California, and cayenne in Pittsburgh. Our experience highlights the importance of incorporating domain expertise and the value of transfer learning approaches.

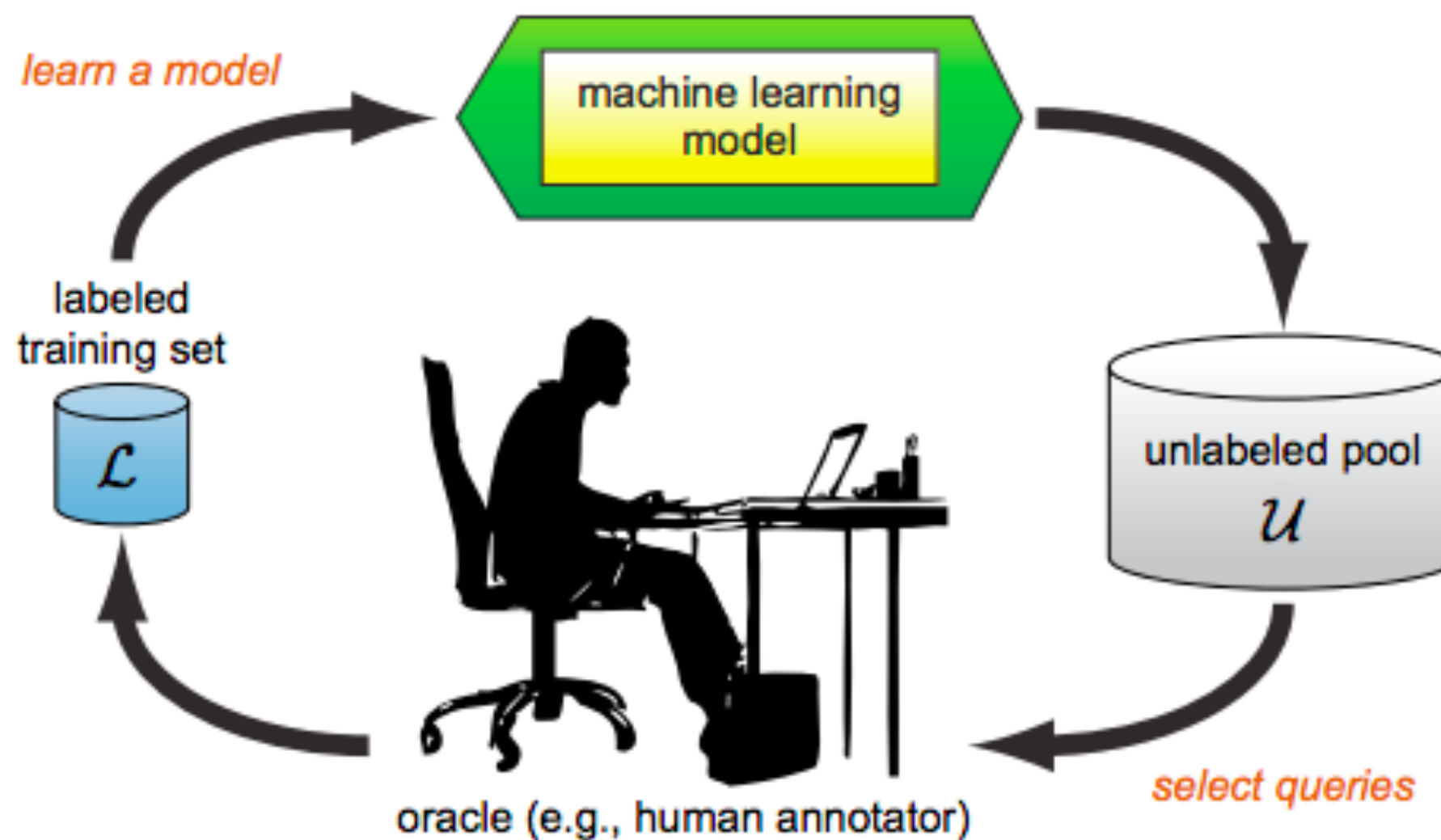
1 Introduction: The Challenge of Chocolate-Chip Cookies

Bayesian Optimization and black-box optimization are used extensively to optimize hyperparameters in machine learning (e.g. [13, 11, 3]) but less so outside that area, and even less so in fields like the culinary arts. We conjecture that the primary barrier to adoption is not technical, but rather cultural and educational. Just as it took years for application communities to frame tasks as supervised learning problems, it likewise takes time to recognize when black-box optimization can provide value. We seek to accelerate this process of cross-disciplinary adoption by creating a challenge that would help practitioners across disciplines recognize problems suitable for black-box optimization in their own settings.

The challenge was to optimize the recipe for chocolate chip cookies. This task highlights key qualities of problems well suited to Bayesian Optimization. The number of tunable parameters is relatively small (e.g. amounts of flour, sugar, etc; see e.g. Table 2). The cost of each experimental iteration is relatively high, requiring manual labor to mix and bake, then taste and report scores. And the problem is familiar enough so that people outside the optimization community can savor the complexities of the challenge and the resulting solutions.

In the remainder of this paper we will lay out a case study of our experience, including methods used, complications that arose in practice, and a summation of lessons learned.

THE ACTIVE LEARNING LOOP



TWENTY QUESTIONS