



DATA, ESTIMATION & INFERENCE

GAUSSIAN PROCESSES & BAYESIAN DEEP LEARNING

MICHAELMAS TERM, 2020
University of Oxford

Tim G. J. Rudner

tim.rudner@cs.ox.ac.uk

Syllabus: <https://www.notion.so/Data-Estimation-and-Inference-Part-2-d177bdb310e74f9aaf034c35652e4f3>

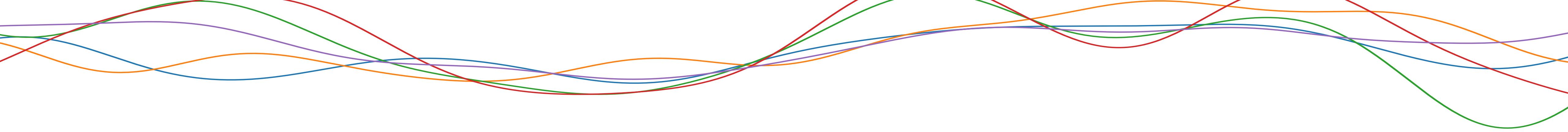
OVERVIEW

Syllabus:

<https://www.notion.so/Data-Estimation-and-Inference-Part-2-d177bdb310e74f9aafd034c35652e4f3>

Today:

- Recap
- Hyperparameter optimization for GPs
- Uncertainty quantification



RECAP: GAUSSIAN PROCESS REGRESSION 101

RECAP: GAUSSIAN PROCESS REGRESSION 101

Gaussian Processes:

Definition 2.1 A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution. \square

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$
$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})],$$
$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

Bayesian linear model:

$$f(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^\top \mathbf{w}$$

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$$

$$\mathbb{E}[f(\mathbf{x})] = \boldsymbol{\phi}(\mathbf{x})^\top \mathbb{E}[\mathbf{w}] = 0,$$

$$\mathbb{E}[f(\mathbf{x})f(\mathbf{x}')] = \boldsymbol{\phi}(\mathbf{x})^\top \mathbb{E}[\mathbf{w}\mathbf{w}^\top] \boldsymbol{\phi}(\mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^\top \Sigma_p \boldsymbol{\phi}(\mathbf{x}')$$

$$|\mathbf{x} \times \mathbf{P} \quad \mathbf{P} \times \mathbf{P} \quad \mathbf{P} \times \mathbf{x}|$$

$$\Leftrightarrow K(\mathbf{x}, \mathbf{x}')$$

$$\boldsymbol{\phi}: \mathcal{X} \rightarrow \mathcal{Z}$$

Example:

$$\boldsymbol{\phi}(\mathbf{x}) \triangleq [\mathbf{x}_1 \mathbf{x}^2]$$

RECAP: GAUSSIAN PROCESS REGRESSION 101

Gaussian Process prior:

$$\mathbf{f}_* \sim \mathcal{N}(\mathbf{0}, K(X_*, X_*))$$

Gaussian Process joint:

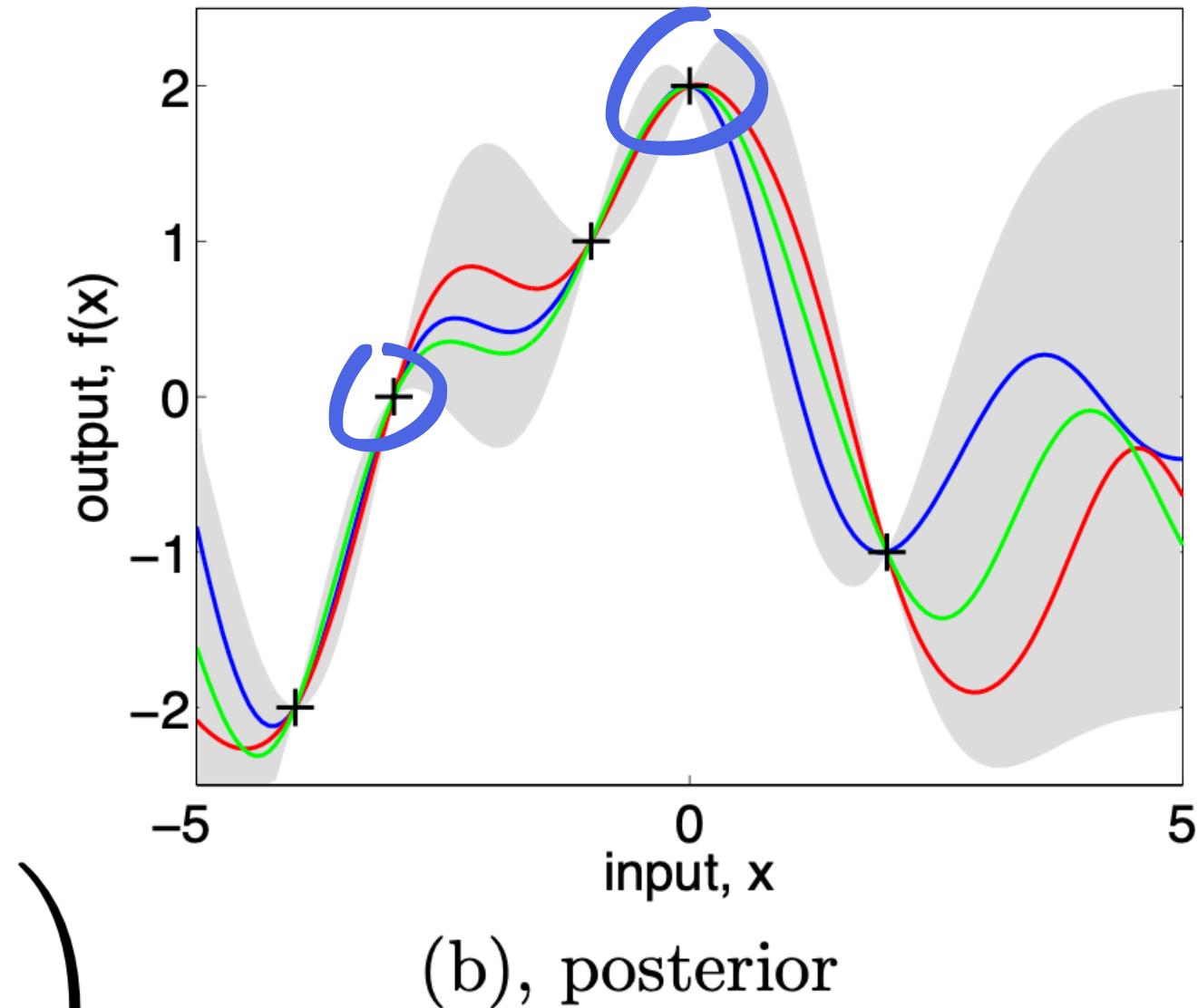
$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right)$$

Gaussian Process posterior:

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*) \in \mathbb{R}^{2 \times 2}$$



RECAP: GAUSSIAN PROCESS REGRESSION 101

Gaussian Process posterior predictive distribution:

$$\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] = K(\mathbf{X}_*, \mathbf{X})[K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X})[K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I]^{-1} \underbrace{K(\mathbf{X}, \mathbf{X}_*)}_{= ?}$$

Bayesian linear regression posterior predictive distribution:

$$f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\boldsymbol{\phi}_*^\top \boldsymbol{\Sigma}_p \boldsymbol{\Phi} (K + \sigma_n^2 I)^{-1} \mathbf{y},$$

$$\boldsymbol{\phi}_*^\top \boldsymbol{\Sigma}_p \boldsymbol{\phi}_* - \boldsymbol{\phi}_*^\top \boldsymbol{\Sigma}_p \boldsymbol{\Phi} (K + \sigma_n^2 I)^{-1} \underbrace{\boldsymbol{\Phi}^\top \boldsymbol{\Sigma}_p \boldsymbol{\phi}_*}_{=?})$$

$$\boldsymbol{\phi}(\mathbf{x}_*) = \boldsymbol{\phi}_* \quad K = \boldsymbol{\Phi}^\top \boldsymbol{\Sigma}_p \boldsymbol{\Phi}$$

Move on
this later

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Gaussian Process prior:

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

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})],$$

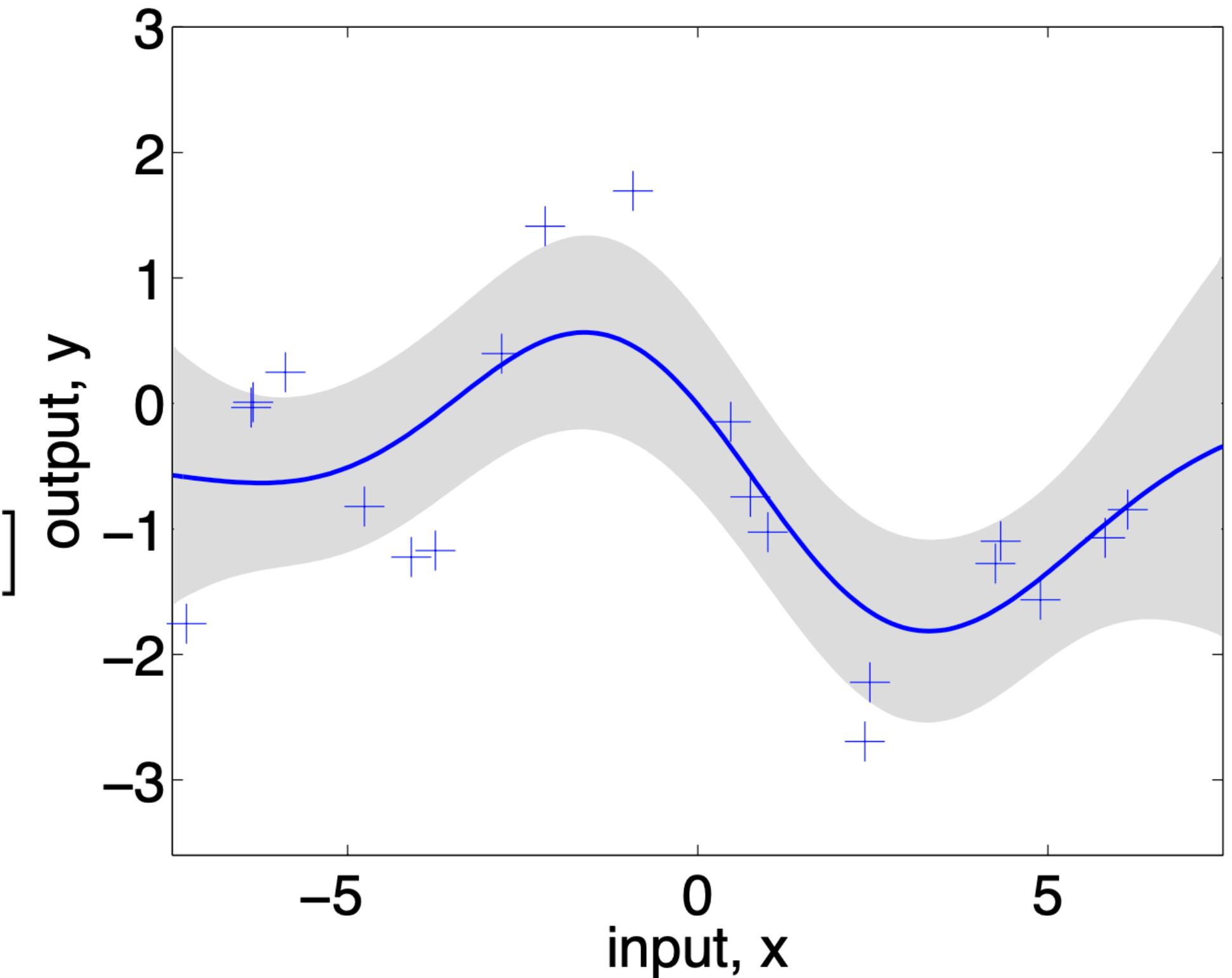
$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

Gaussian Process posterior:

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$



IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Marginal likelihood:

$$p(\mathbf{y}|X) = \int p(\mathbf{y}|\mathbf{f}, X)p(\mathbf{f}|X) d\mathbf{f}$$

$$\log p(\mathbf{y}|X) = -\frac{1}{2}\mathbf{y}^\top(K + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\log|K + \sigma_n^2 I| - \frac{n}{2}\log 2\pi \quad *$$

input: X (inputs), \mathbf{y} (targets), k (covariance function), σ_n^2 (noise level),
 \mathbf{x}_* (test input)

- 2: $L := \text{cholesky}(K + \sigma_n^2 I)$
 - 3: $\boldsymbol{\alpha} := L^\top \backslash (L \backslash \mathbf{y})$
 - 4: $\bar{f}_* := \mathbf{k}_*^\top \boldsymbol{\alpha}$
 - 5: $\mathbf{v} := L \backslash \mathbf{k}_*$
 - 6: $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^\top \mathbf{v}$
 - 7: $\log p(\mathbf{y}|X) := -\frac{1}{2}\mathbf{y}^\top \boldsymbol{\alpha} - \sum_i \log L_{ii} - \frac{n}{2}\log 2\pi$
 - 8: **return:** \bar{f}_* (mean), $\mathbb{V}[f_*]$ (variance), $\log p(\mathbf{y}|X)$ (log marginal likelihood)
- } predictive mean eq. (2.25)
} predictive variance eq. (2.26)
} eq. (2.30)
- *

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Implementation checklist:

1) GP prior, including kernel function

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

2) Posterior predictive distribution:

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$

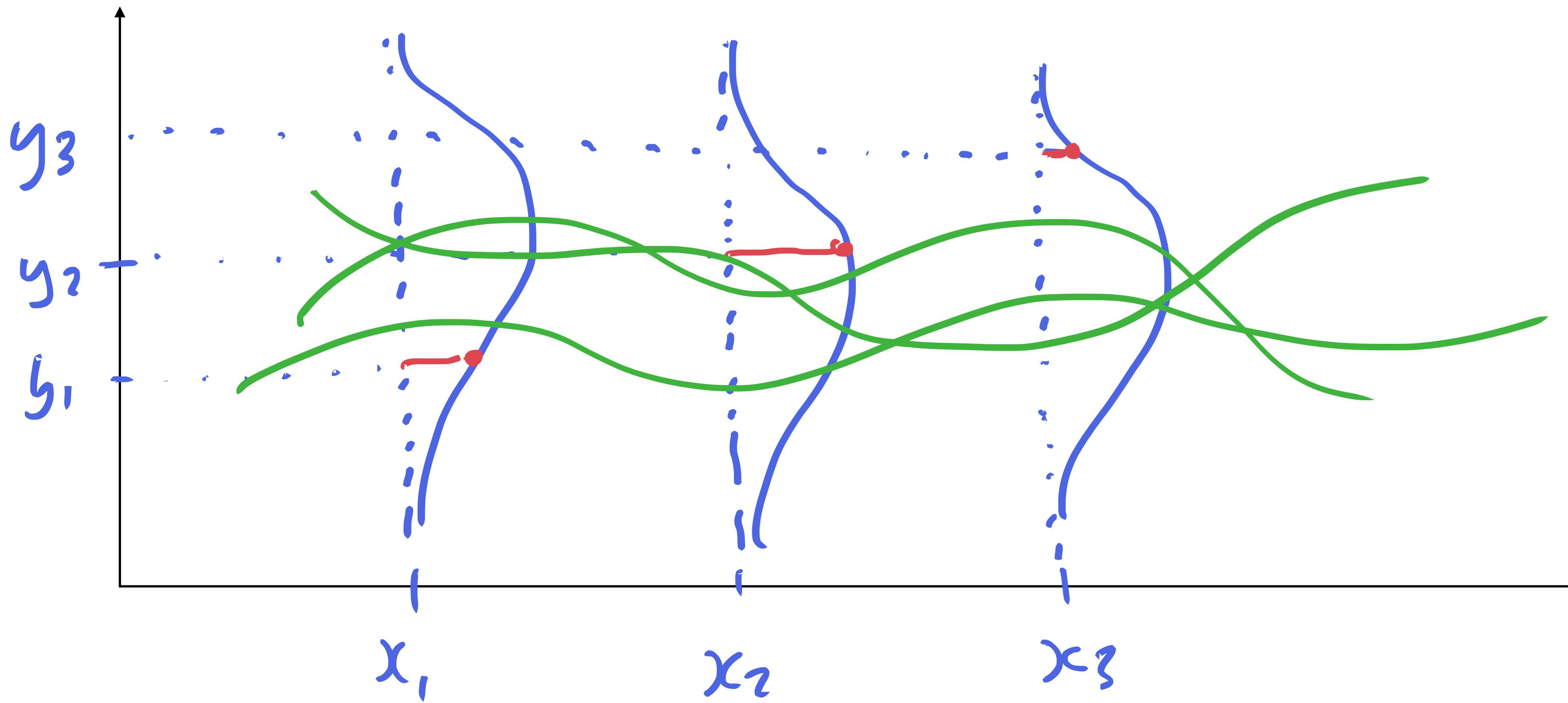
3) Marginal likelihood:

$$\log p(\mathbf{y} | X) = -\frac{1}{2} \mathbf{y}^\top (K + \sigma_n^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{n}{2} \log 2\pi$$

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Understanding the marginal likelihood

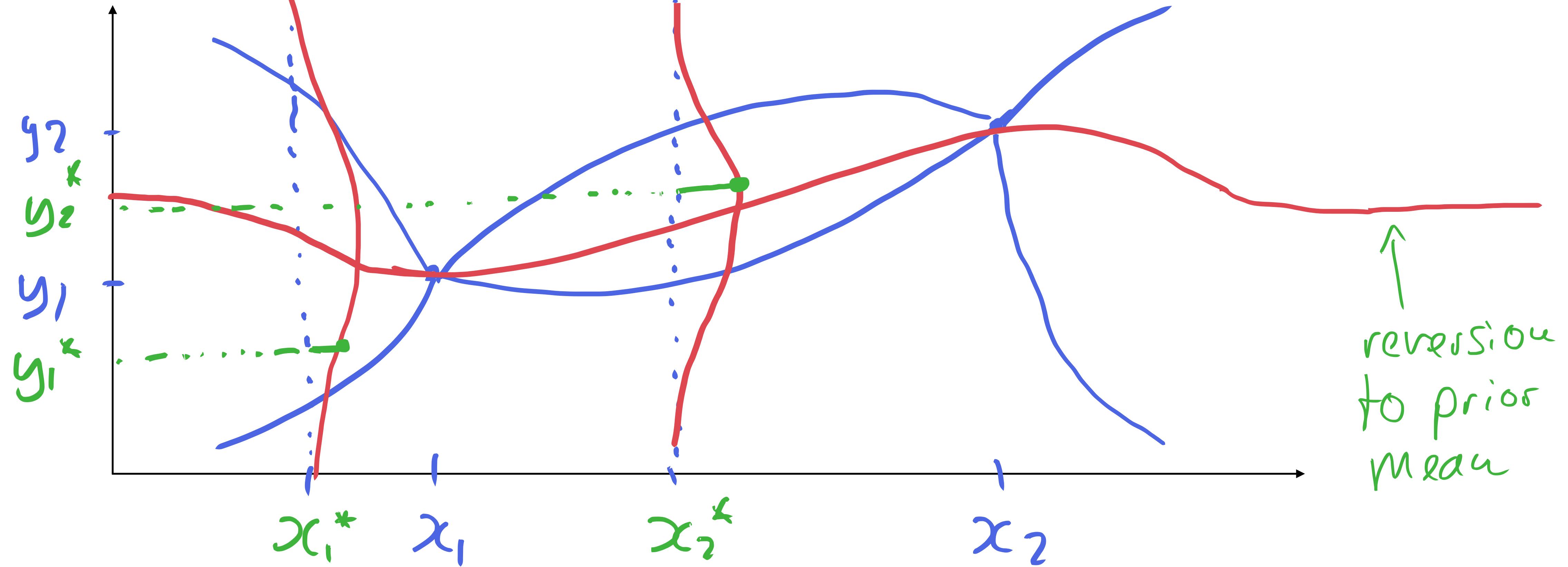
$$p(\mathbf{y}|X) = \int p(\mathbf{y}|\mathbf{f}, X)p(\mathbf{f}|X) d\mathbf{f}$$



IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Assessing performance

$$-\log p(y_* | \mathcal{D}, \mathbf{x}_*) = \frac{1}{2} \log(2\pi\sigma_*^2) + \frac{(y_* - \bar{f}(\mathbf{x}_*))^2}{2\sigma_*^2}$$



IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Covariance (~kernel) functions

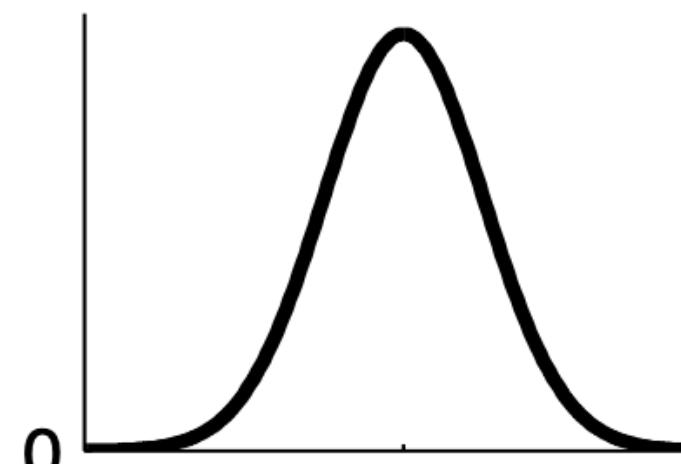
Kernel name:

$$k(x, x') =$$

Plot of $k(x, x')$:

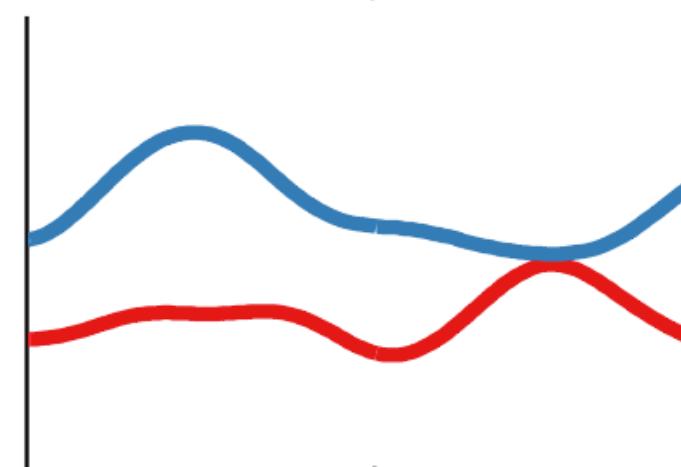
Squared-exp (SE)

$$\sigma_f^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$$



$x - x'$
↓

Functions $f(x)$
sampled from
GP prior:



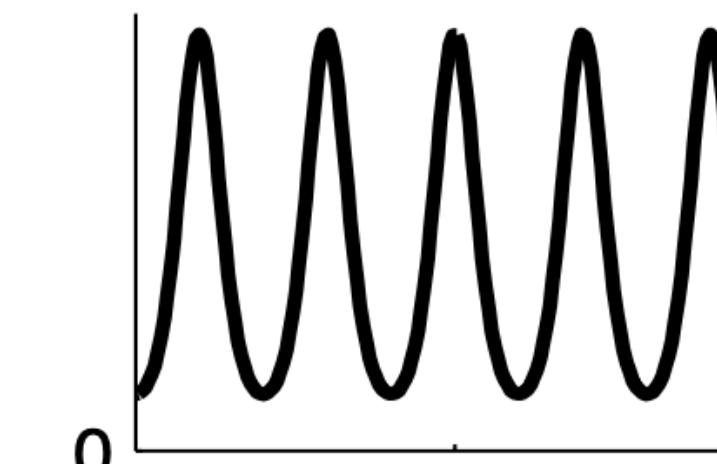
x

Type of structure:

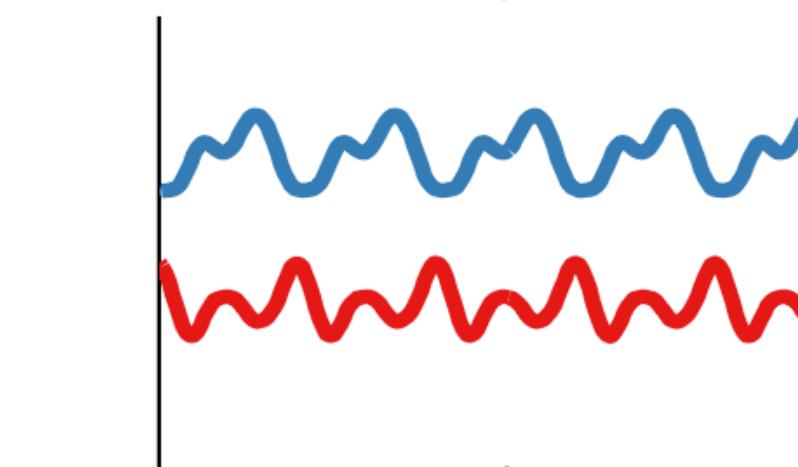
local variation

Periodic (Per)

$$\sigma_f^2 \exp\left(-\frac{2}{\ell^2} \sin^2\left(\pi \frac{x-x'}{p}\right)\right)$$



$x - x'$
↓

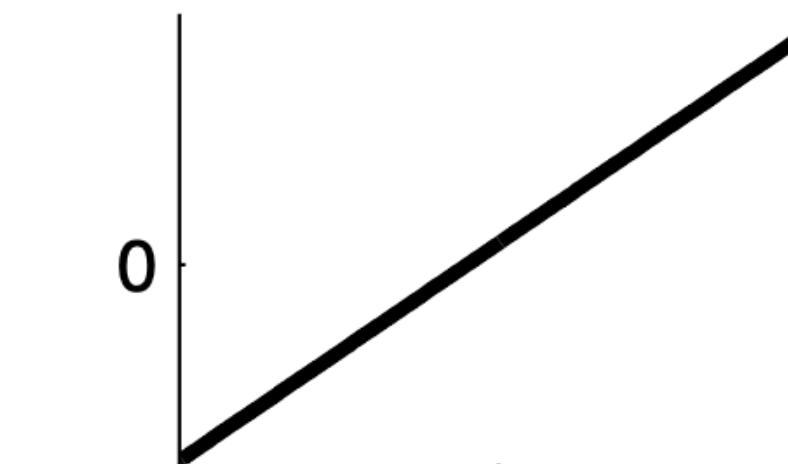


x

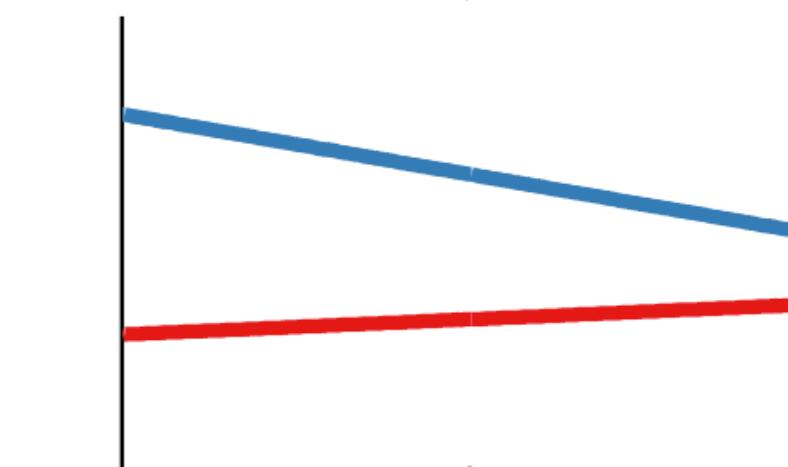
repeating structure

Linear (Lin)

$$\sigma_f^2(x - c)(x' - c)$$



x (with $x' = 1$)
↓



x

linear functions

$X^T X$

Bayesian linear regression

$\phi(x)^T \phi(x)$

↑ ↑

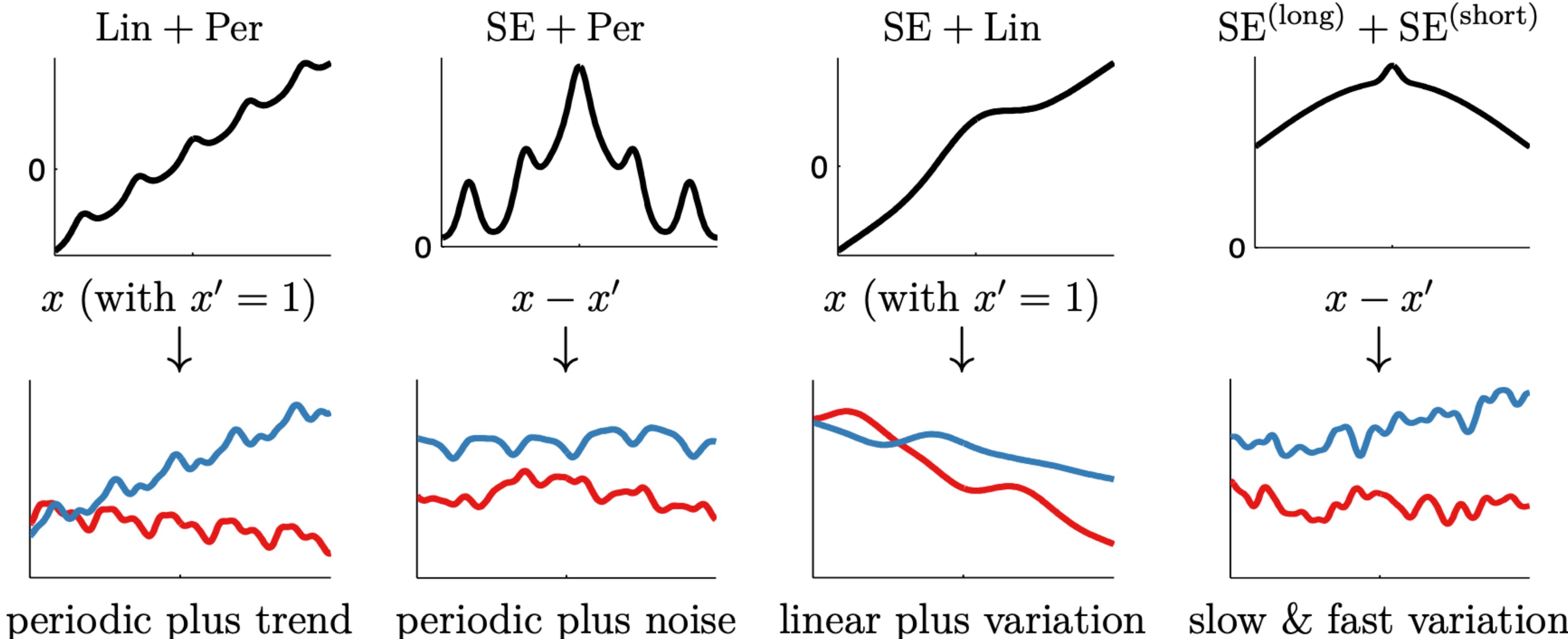
$n \times P$ $P \times 1$

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Combining covariance functions through addition

$$k_a + k_b = k_a(\mathbf{x}, \mathbf{x}') + k_b(\mathbf{x}, \mathbf{x}')$$

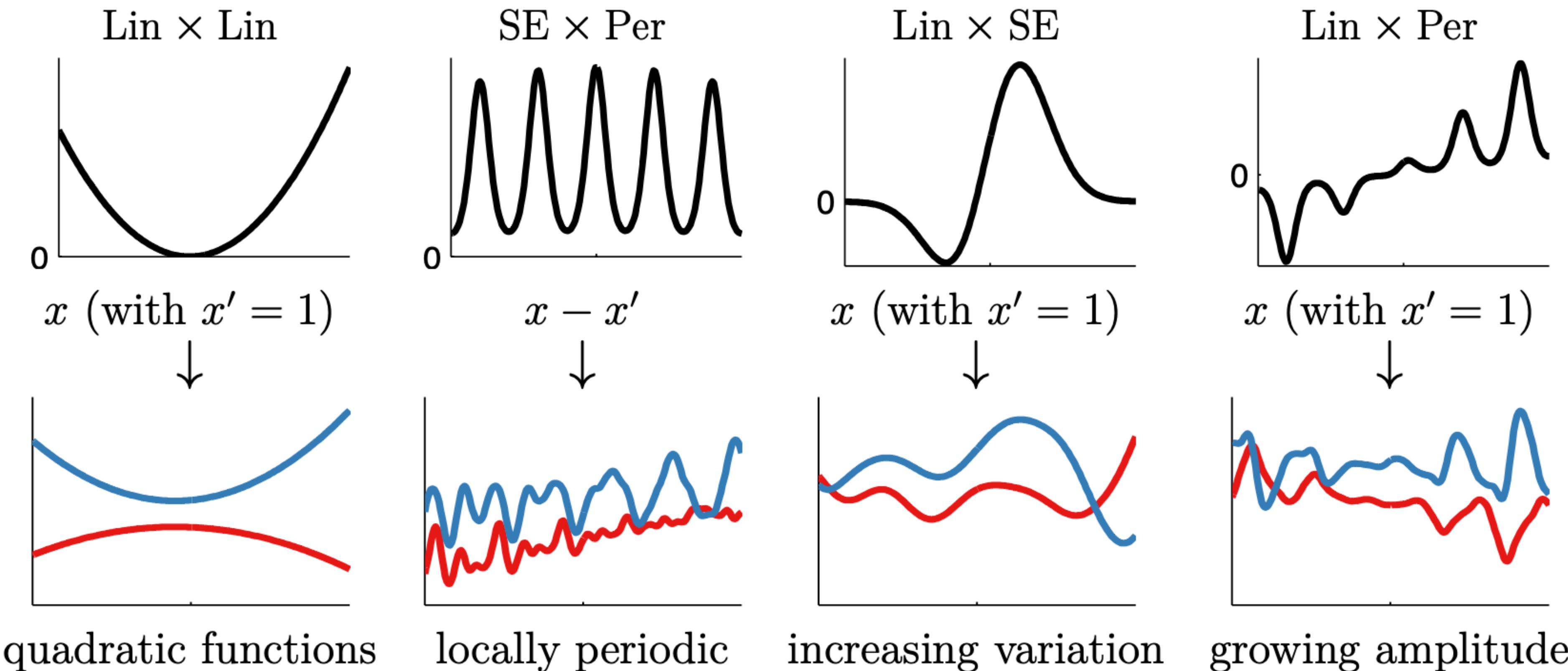
[Link to visualization](#)



IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Combining covariance functions through multiplication

$$k_a \times k_b = k_a(\mathbf{x}, \mathbf{x}') \times k_b(\mathbf{x}, \mathbf{x}')$$

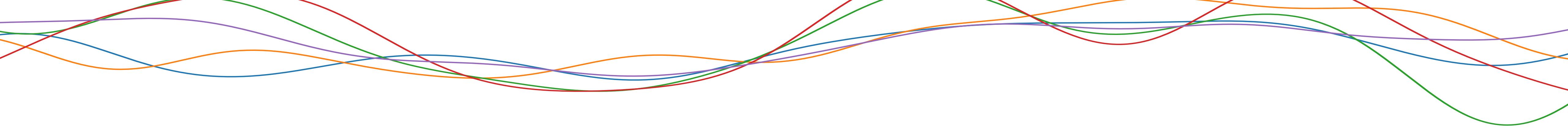


UNCERTAINTY QUANTIFICATION

UNCERTAINTY QUANTIFICATION

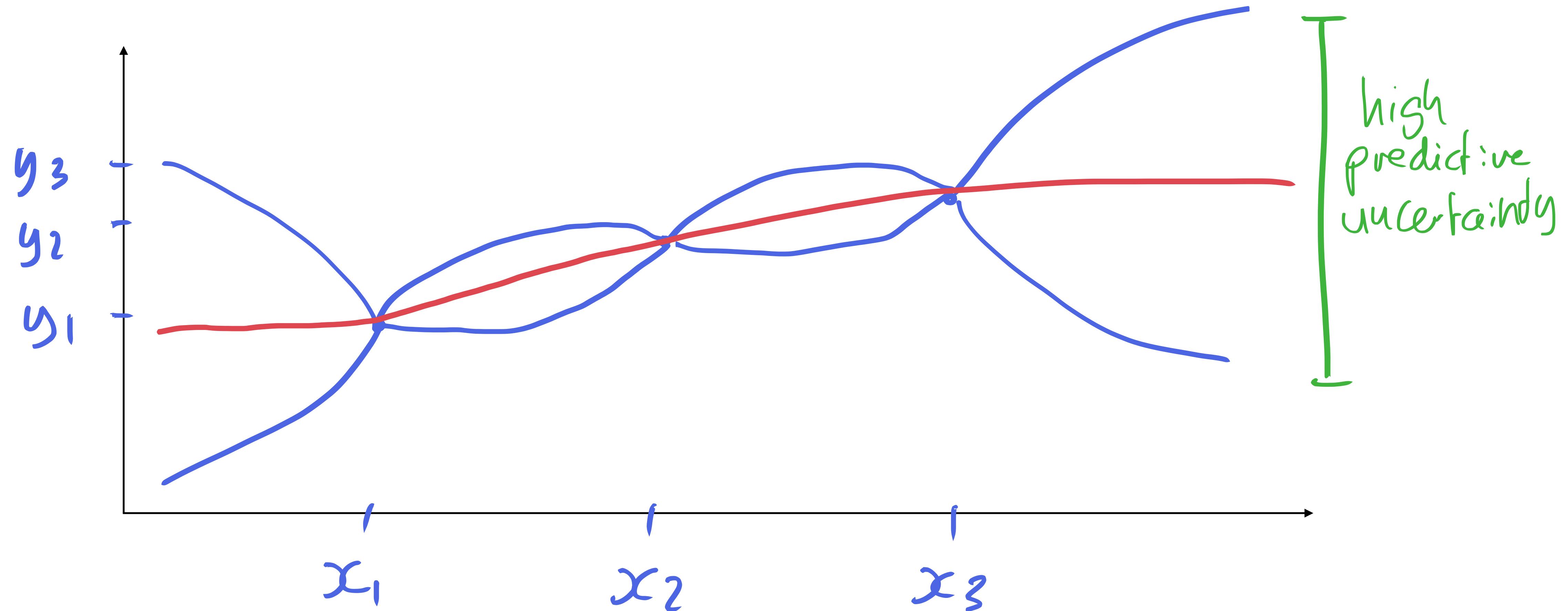
Why Bayesian uncertainty quantification?

- Know what you don't know
- AI safety
- Efficient exploration
- Mathematically-principled foundations



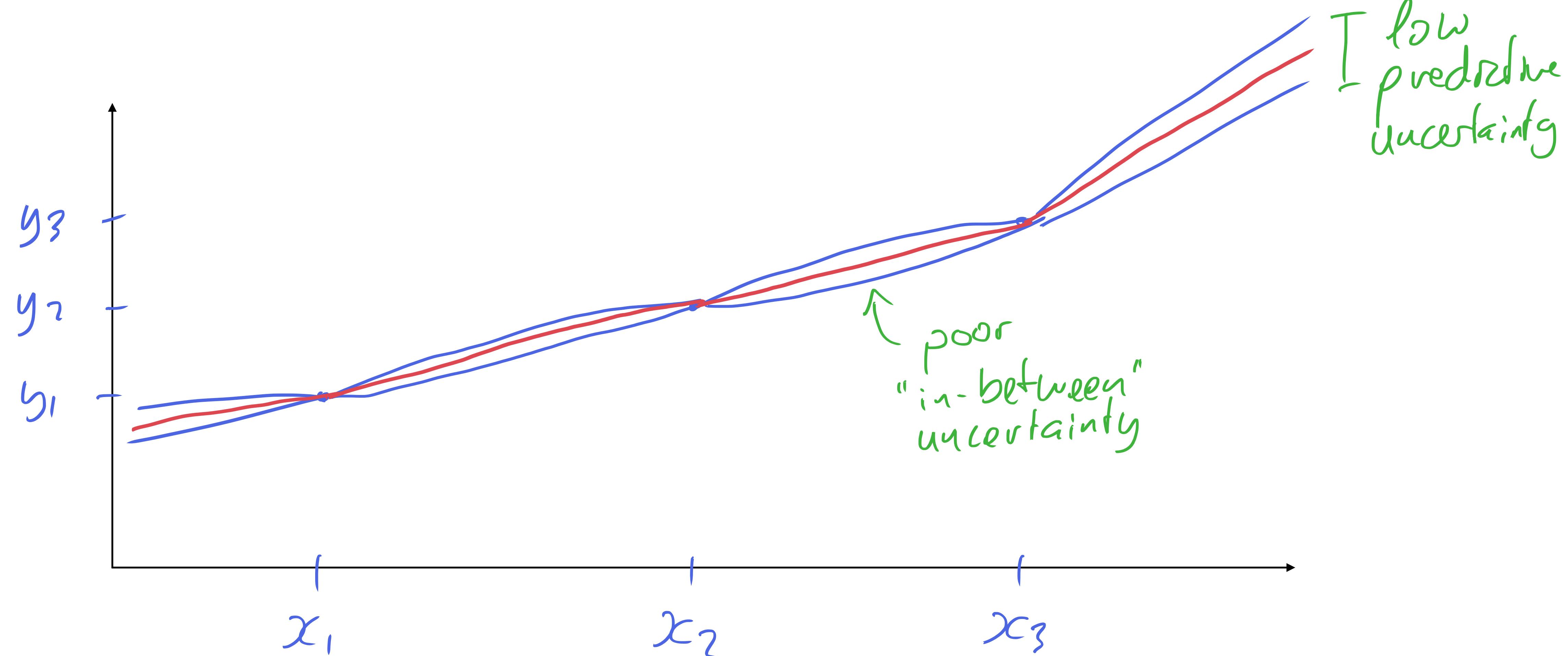
UNCERTAINTY QUANTIFICATION

What predictive uncertainty is good predictive uncertainty?



UNCERTAINTY QUANTIFICATION

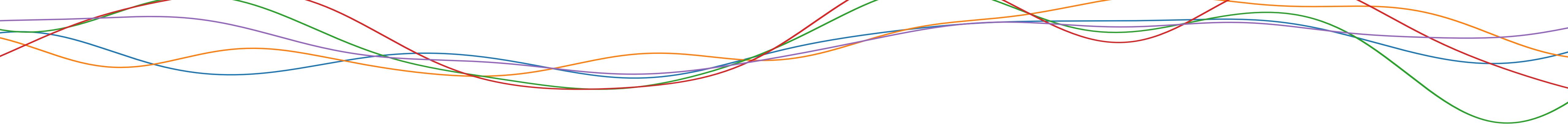
What predictive uncertainty is bad predictive uncertainty?



UNCERTAINTY QUANTIFICATION

Desirable properties in predictive uncertainty estimates

- Low predictive uncertainty close to observed data
- Reversion to the prior away from observed data
- A suitable prior predictive distribution
- Avoid confident but wrong predictions

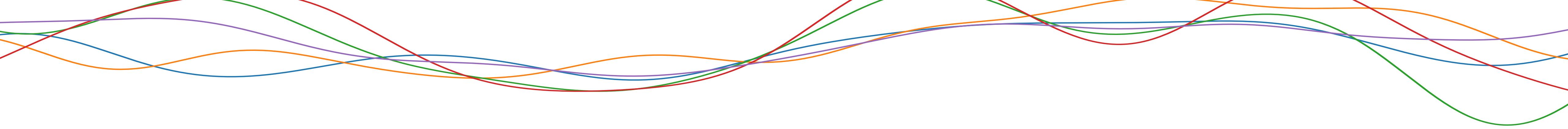


UNCERTAINTY QUANTIFICATION

What models reliably provide good predictive uncertainty estimates?

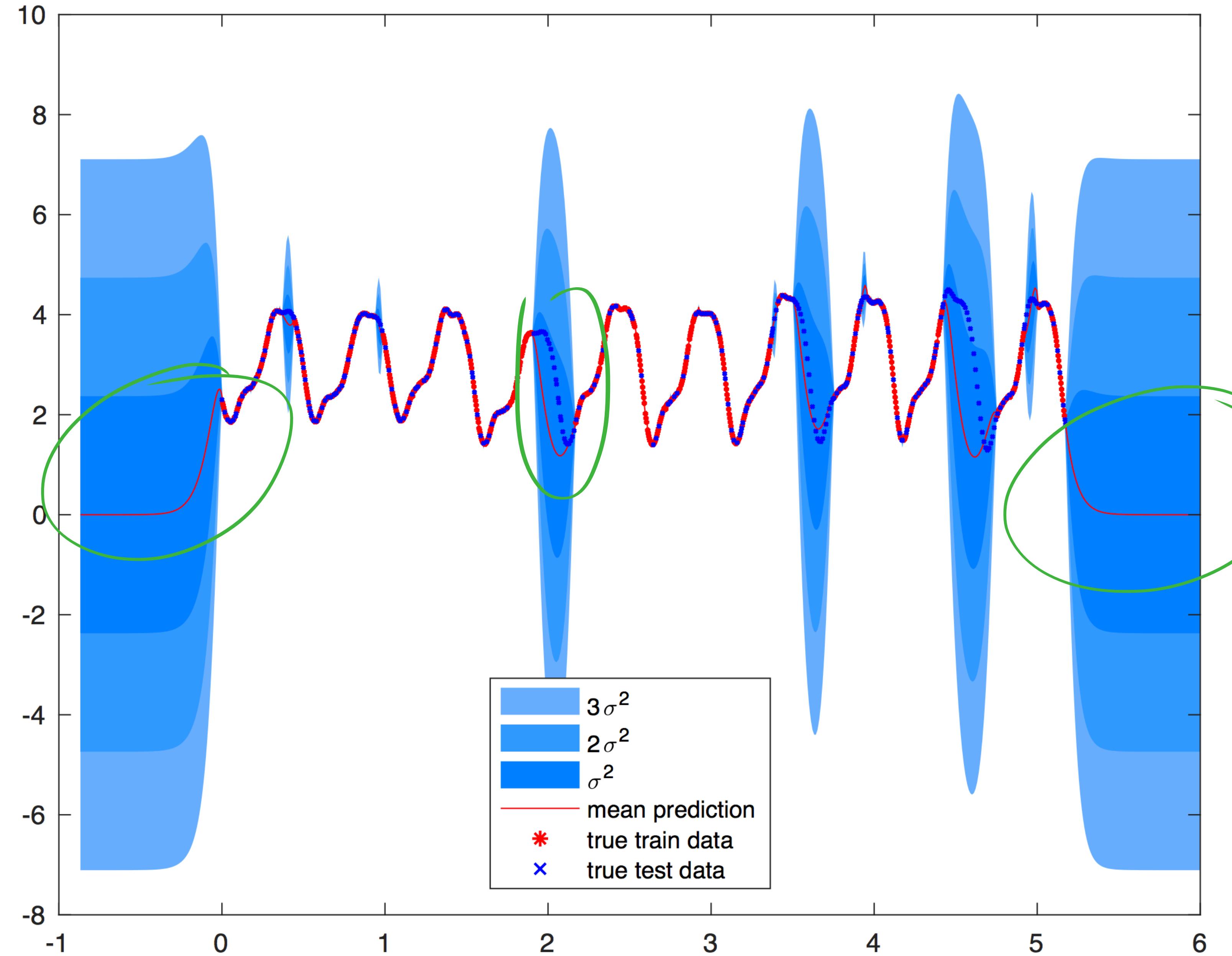
- Gaussian Processes?
- Bayesian linear models?
- Bayesian Neural Networks?

It depends!



LAB ASSIGNMENT

LAB ASSIGNMENT



LAB ASSIGNMENT

