

10-701: Introduction to Machine Learning

Lecture 26 – Gaussian Processes

Henry Chai

4/22/24

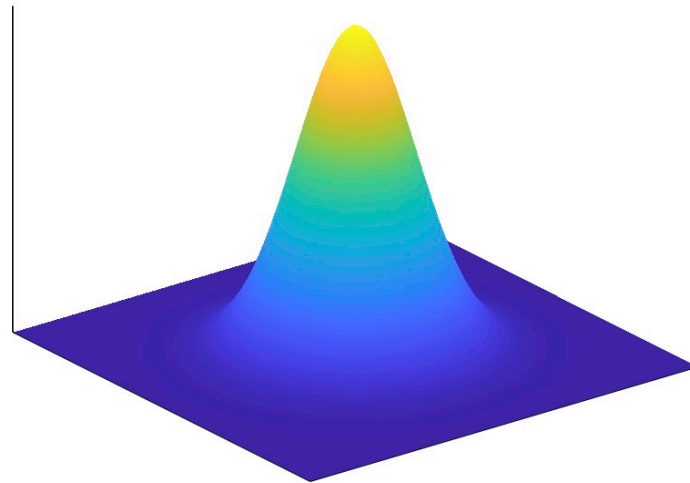
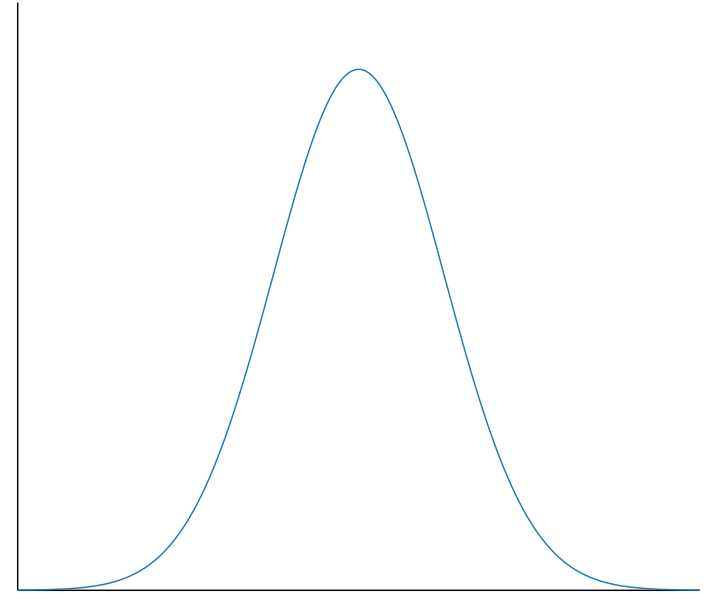
Front Matter

- Announcements:
 - Exam 2 on 5/6 **from 1 PM – 3 PM in TEP 1403**
 - You are allowed to bring one letter-/A4-size sheet of notes; you can put *whatever* you want on *both sides*
 - **Pre-midterm material may be referenced but will not be the primary focus of any question**
 - Project Final Reports due on 4/26 (Friday) at 11:59 PM
 - **No late days can be used on project deliverables**
- Recommended Readings:
 - Murphy, Chapters 15.1-15.2

Gaussians

(Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$



- Multivariate Gaussians:

$$\mathbf{x} = [x_1, \dots, x_D]^T$$

$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

Some fun facts about Gaussians

- Closure under linear transformations:
- Closure under addition
- Closure under conditioning:

Some old
friends

Gaussian process =
Bayesian linear regression + Kernels

Recall: MAP for Linear Regression

- If we assume a linear model with additive Gaussian noise

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \rightarrow \mathbf{y} \sim N(X\mathbf{w}, \sigma^2 I_N)$$

and independent identical Gaussian priors on the weights...

$$\mathbf{w} \sim N\left(\mathbf{w}_{D+1}, \frac{\sigma^2}{\lambda} I_{D+1}\right) \rightarrow p(\mathbf{w}) \propto \exp\left(-\frac{1}{2\sigma^2} (\lambda \mathbf{w}^T \mathbf{w})\right)$$

- ... then, the MAP of \mathbf{w} is the ridge regression solution!

$$\begin{aligned} \mathbf{w}_{MAP} &= \underset{\mathbf{w}}{\operatorname{argmin}} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \\ &= (X^T X + \lambda I_{D+1})^{-1} X^T \mathbf{y} \end{aligned}$$

Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\mathbf{y} \sim N(X\mathbf{0}_{D+1} + \mathbf{0}_N = \mathbf{0}_N, X\Sigma X^T + \sigma^2 I_N)$$

Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0}_{D+1} \\ \mathbf{0}_N \end{bmatrix}, \begin{bmatrix} \Sigma & ??? \\ ??? & X\Sigma X^T + \sigma^2 I_N \end{bmatrix} \right)$$

Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0}_{D+1} \\ \mathbf{0}_N \end{bmatrix}, \begin{bmatrix} \Sigma & X\Sigma \\ X\Sigma & X\Sigma X^T + \sigma^2 I_N \end{bmatrix} \right)$$

Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\mathbf{w} \mid \mathbf{y} \sim N(\boldsymbol{\mu}_{POST}, \boldsymbol{\Sigma}_{POST})$$

where

$$\begin{aligned}\boldsymbol{\mu}_{POST} &= \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y}, \\ \boldsymbol{\Sigma}_{POST} &= \Sigma - \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma\end{aligned}$$

Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then given a new test data point \mathbf{x}' , the prediction is

$$y' \mid \mathbf{y} = \mathbf{x}'^T \mathbf{w} \mid \mathbf{y} \sim N(\mathbf{x}'^T \boldsymbol{\mu}_{POST}, \mathbf{x}'^T \Sigma_{POST} \mathbf{x}')$$

where

$$\begin{aligned} \boldsymbol{\mu}_{POST} &= \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y}, \\ \Sigma_{POST} &= \Sigma - \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \end{aligned}$$

Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then given a new test data point \mathbf{x}' , the prediction is

$$y' \mid \mathbf{y} = \mathbf{x}'^T \mathbf{w} \mid \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = \mathbf{x}'^T \Sigma \mathbf{x}' - \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \mathbf{x}'$$

Some old
friends

Gaussian process =
Bayesian linear regression + Kernels

Bayesian Linear Regression...

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then given a new test data point \mathbf{x}' , the prediction is

$$y' \mid \mathbf{y} = \mathbf{x}'^T \mathbf{w} \mid \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = \mathbf{x}'^T \Sigma \mathbf{x}' - \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \mathbf{x}'$$

Bayesian Linear Regression can be kernelized!

$$\Phi = \begin{bmatrix} 1 & \phi(\mathbf{x}^{(1)})^T \\ 1 & \phi(\mathbf{x}^{(2)})^T \\ \vdots & \vdots \\ 1 & \phi(\mathbf{x}^{(N)})^T \end{bmatrix}$$

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point \mathbf{x}' , the prediction is

$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED}$$

$$= \phi(\mathbf{x}')^T \Sigma \phi(\mathbf{x}') - \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \Phi \Sigma \phi(\mathbf{x}')$$

Bayesian Linear Regression can be kernelized!

$$\Phi = \begin{bmatrix} 1 & \phi(\mathbf{x}^{(1)})^T \\ 1 & \phi(\mathbf{x}^{(2)})^T \\ \vdots & \vdots \\ 1 & \phi(\mathbf{x}^{(N)})^T \end{bmatrix}$$

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point \mathbf{x}' , the prediction is

$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \boldsymbol{\Sigma}_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\boldsymbol{\Sigma}_{PRED}$$

$$= \phi(\mathbf{x}')^T \Sigma \phi(\mathbf{x}') - \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \Phi \Sigma \phi(\mathbf{x}')$$

- Define the kernel function to be

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

Bayesian Linear Regression can be kernelized!

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point \mathbf{x}' , the prediction is

$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = K(\mathbf{x}', \mathbf{x}') - K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} K(X, \mathbf{x})$$

- Define the kernel function to be

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

Some old
friends

Gaussian process =
Bayesian linear regression + Kernels

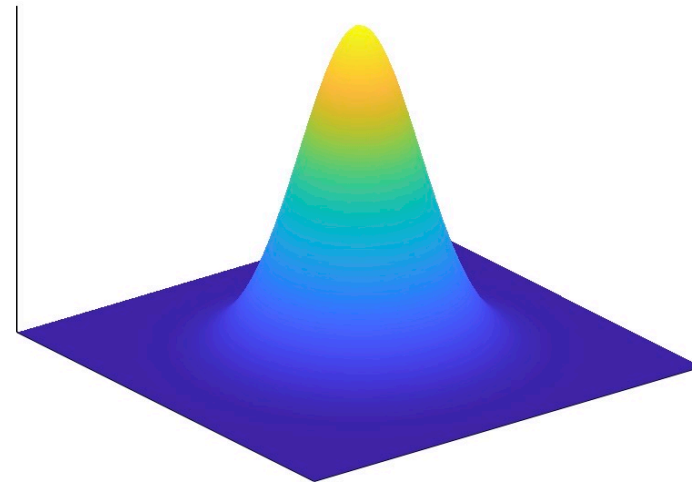
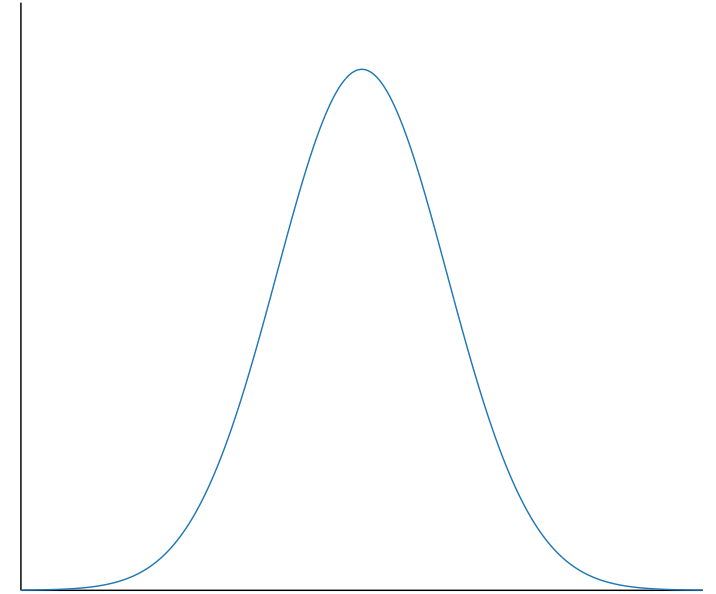
A new
perspective

Gaussian process =
The extension of a Gaussian
distribution to functions

Gaussians

- (Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$



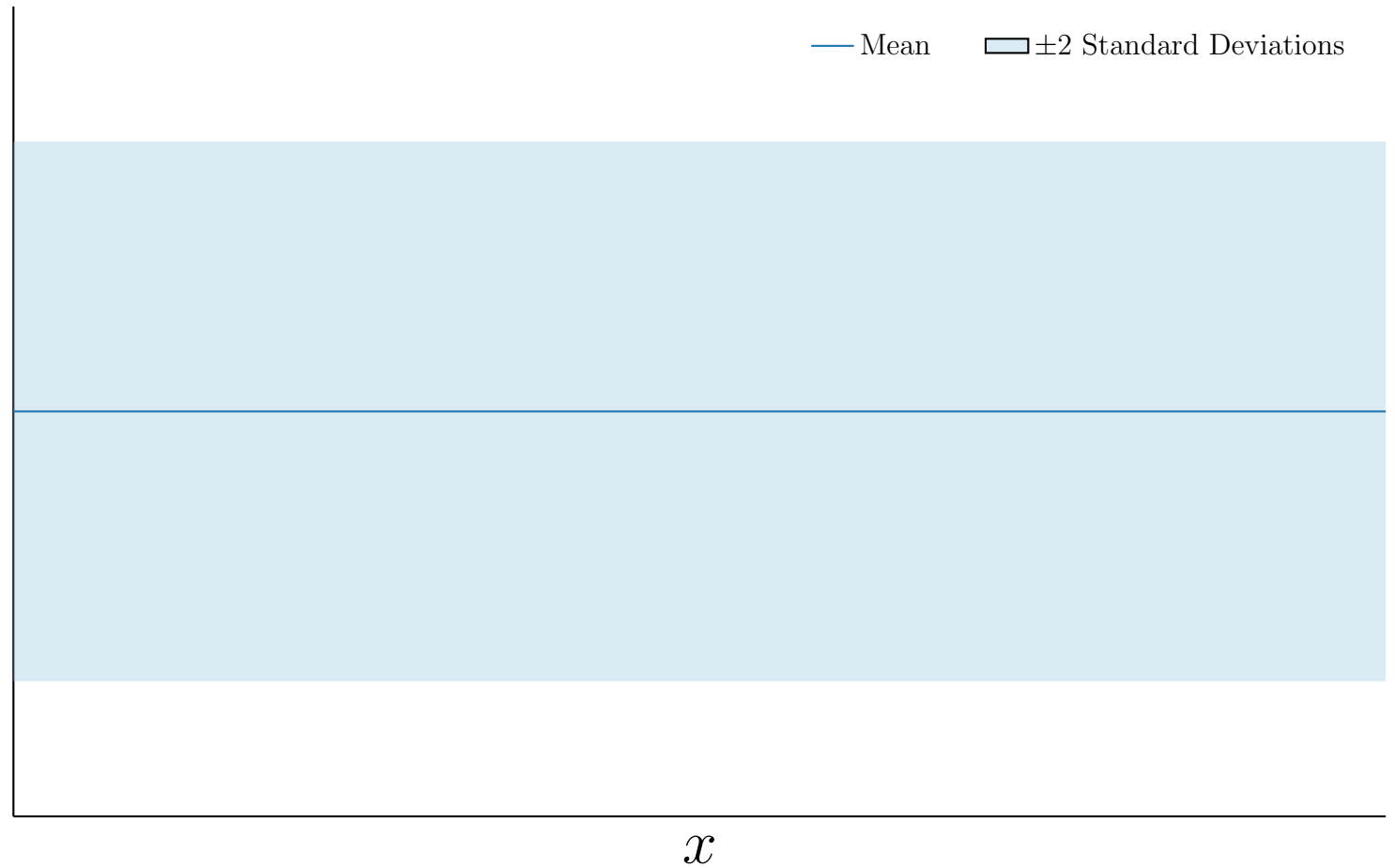
- Multivariate Gaussians:

$$\mathbf{x} = [x_1, \dots, x_D]^T$$

$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

Gaussian Process (GP)

$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x), \Sigma(x, x'))$$

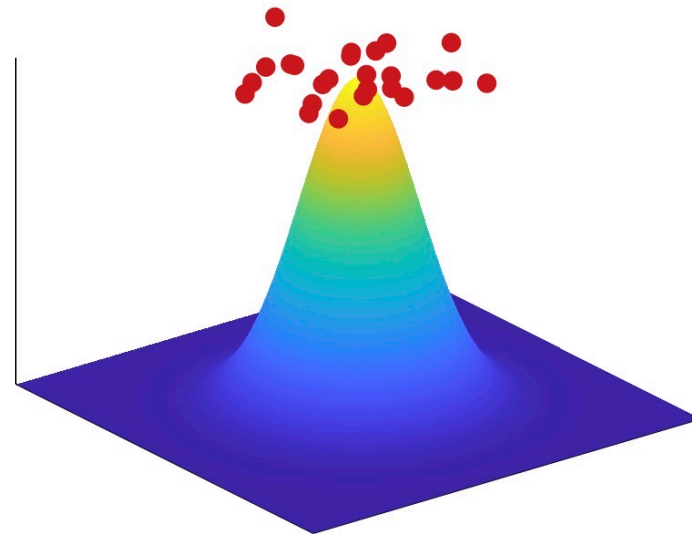
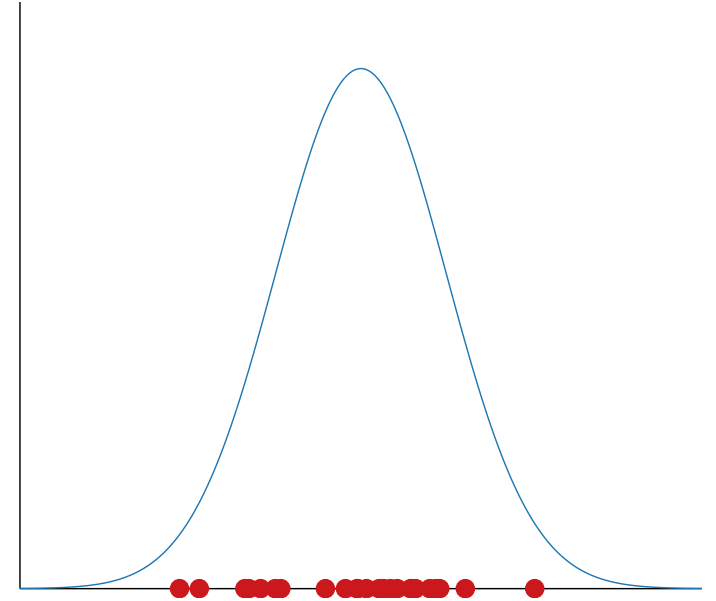


$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

Gaussians

- (Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$



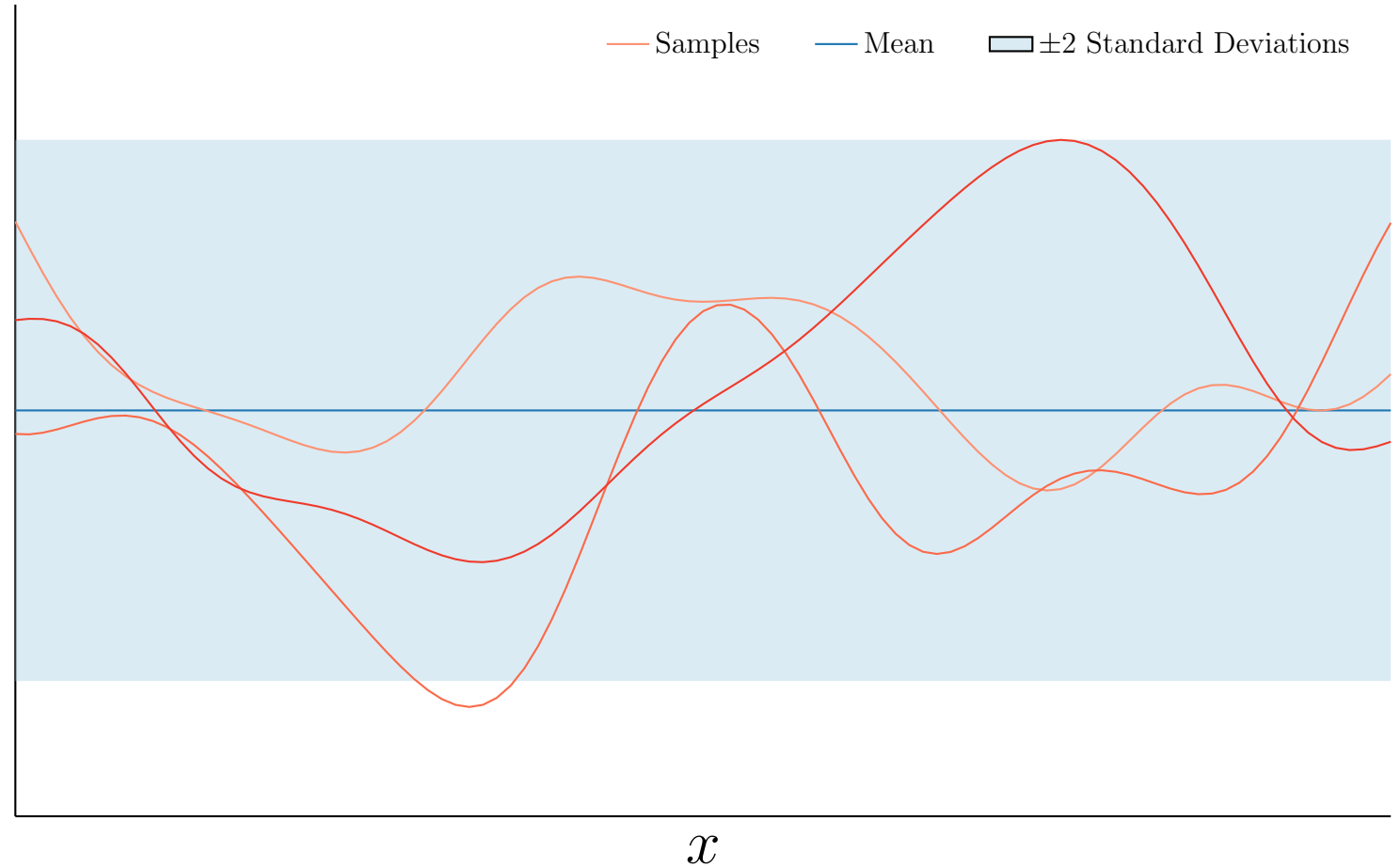
- Multivariate Gaussians:

$$\mathbf{x} = [x_1, \dots, x_D]^T$$

$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

Gaussian Process (GP)

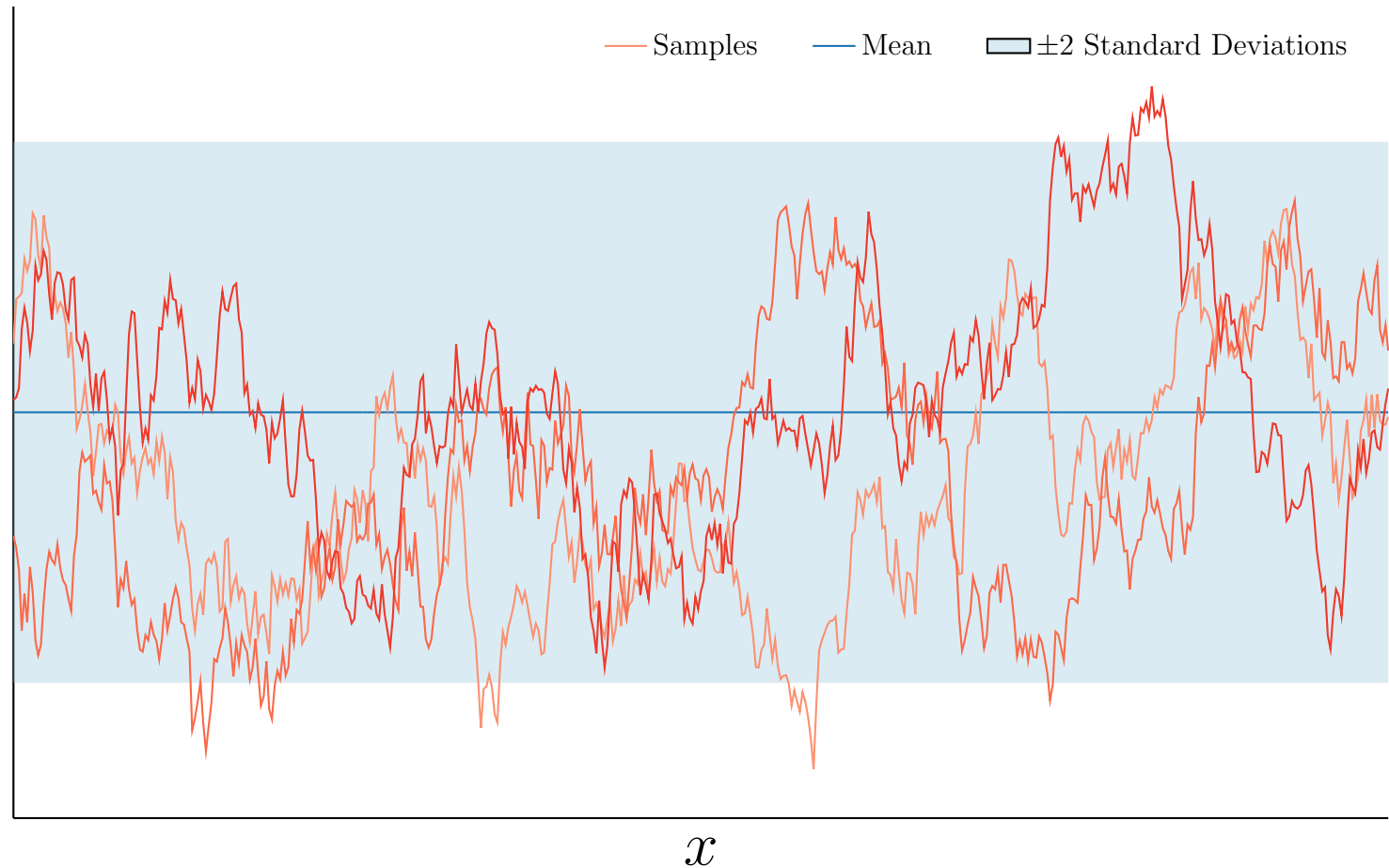
$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-(x - x')^2))$$



$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

Gaussian Process (GP)

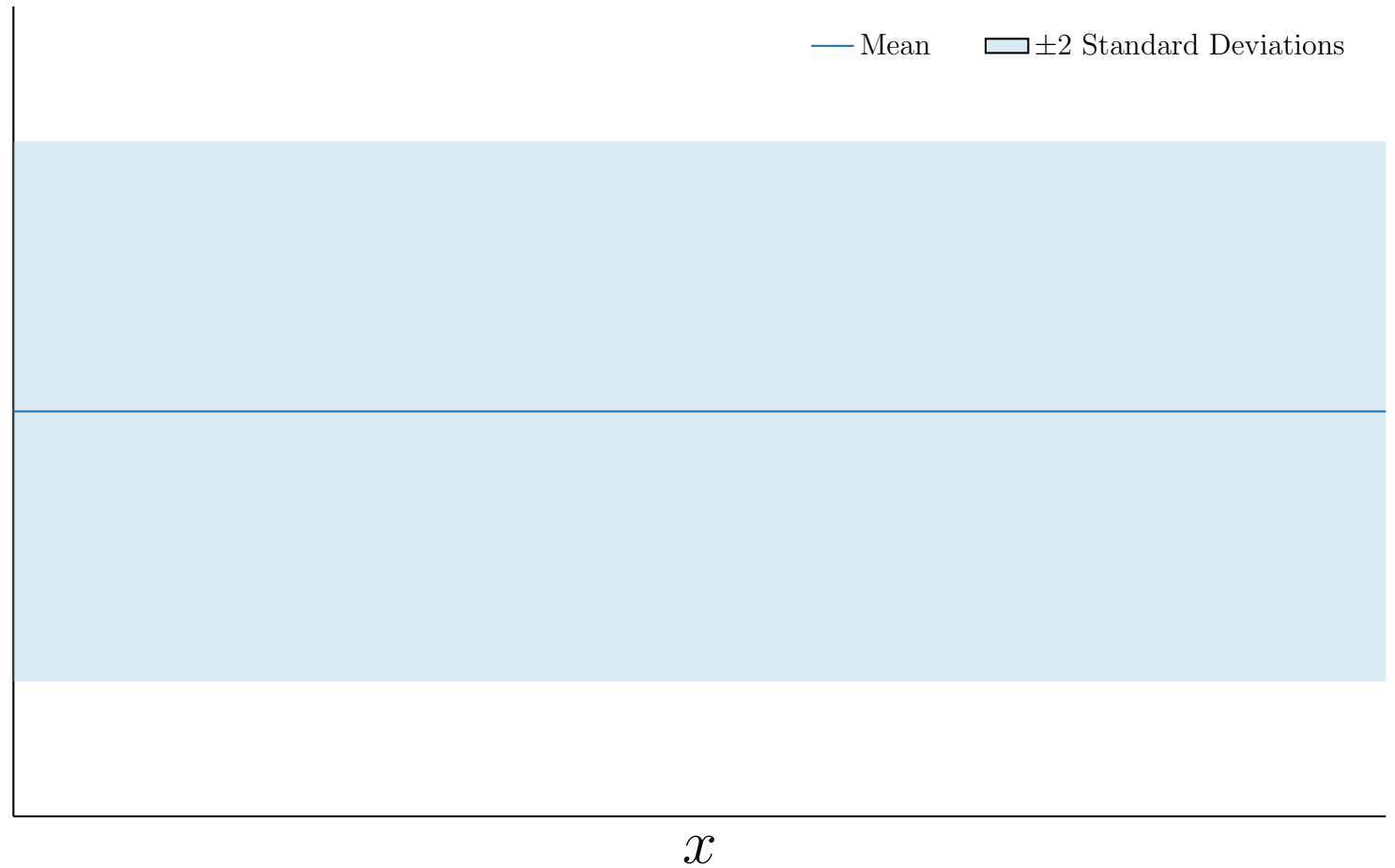
$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-|x - x'|))$$



$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

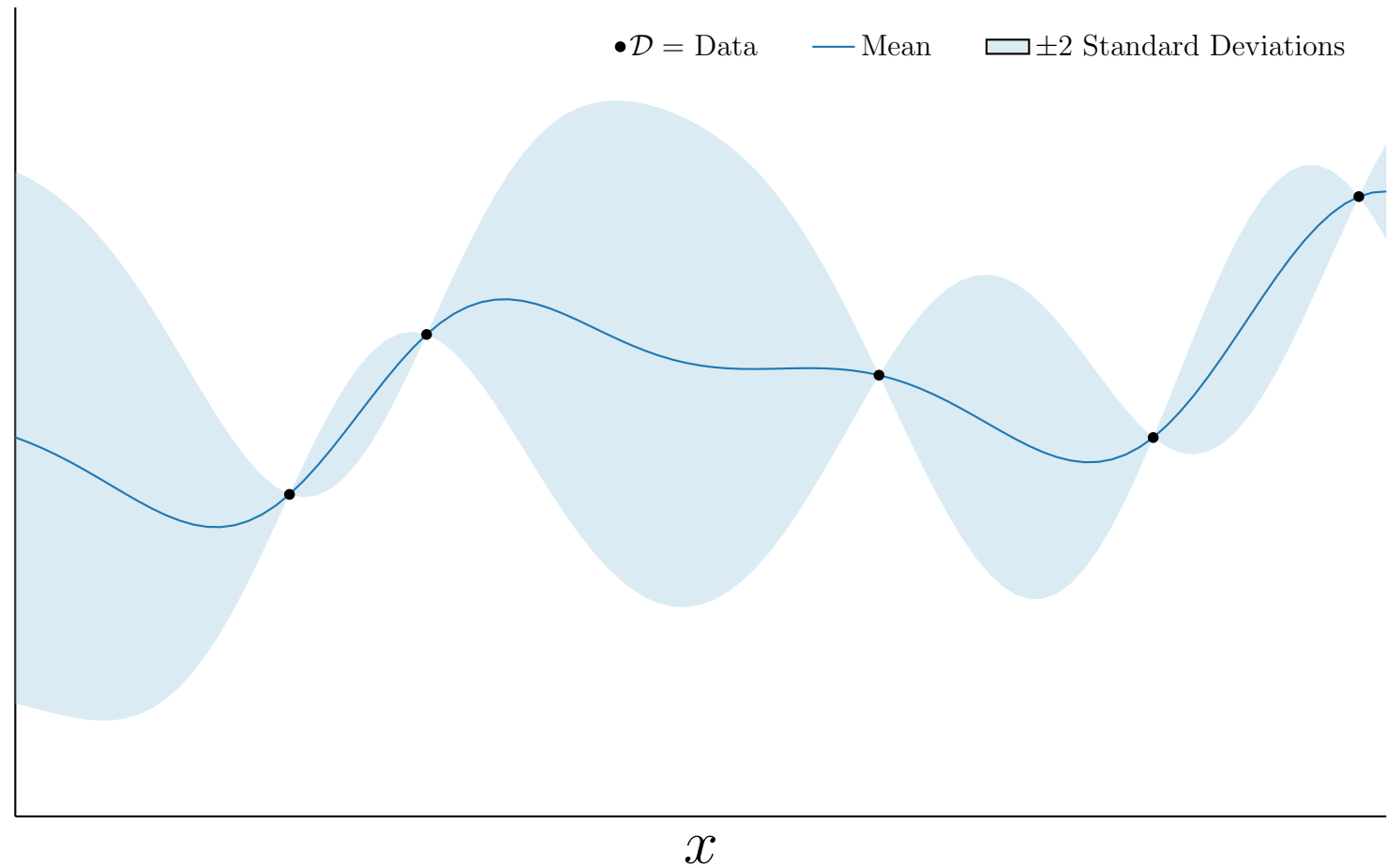
GP Prior

$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-(x - x')^2))$$



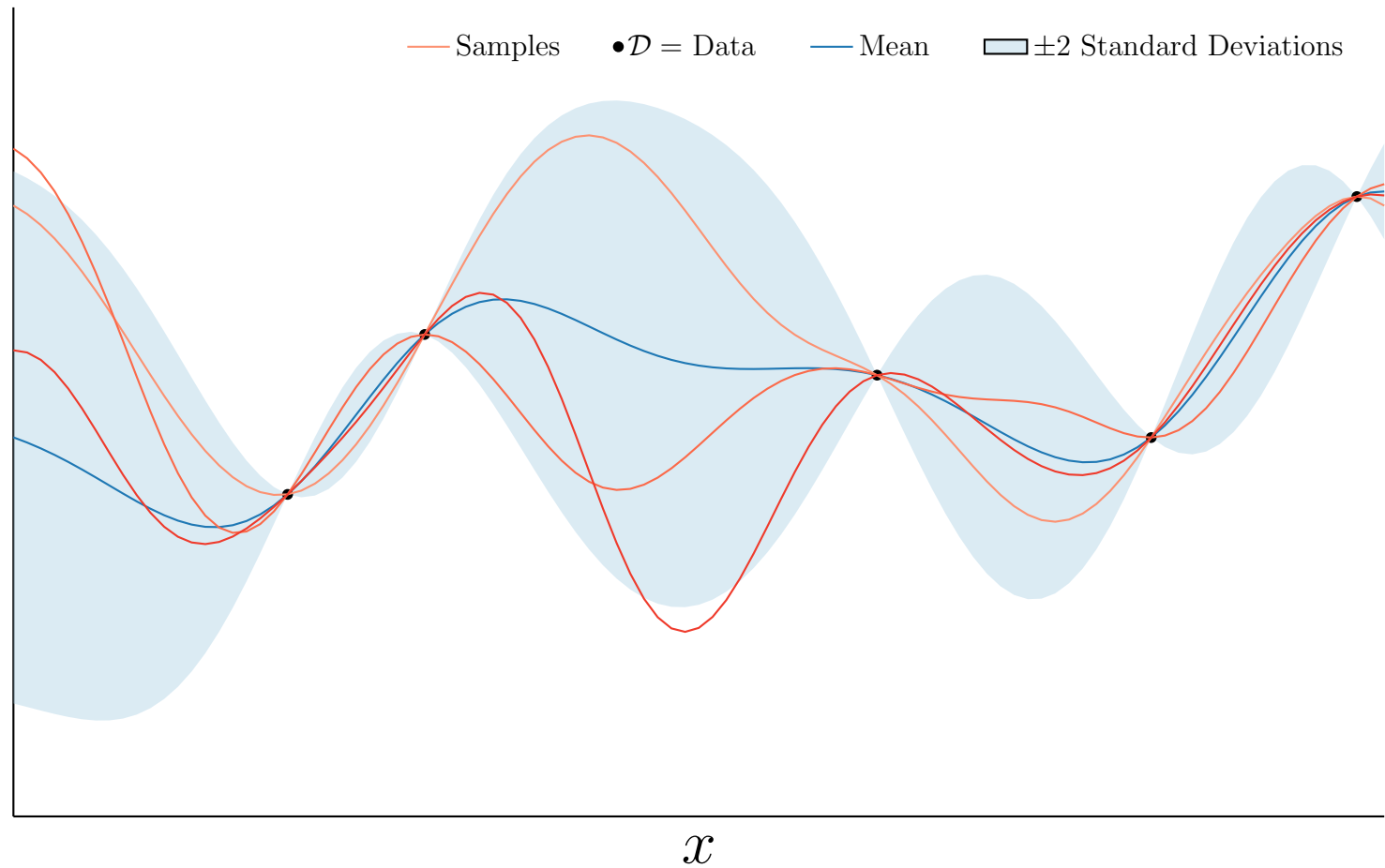
GP Posterior

$$f | \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$$



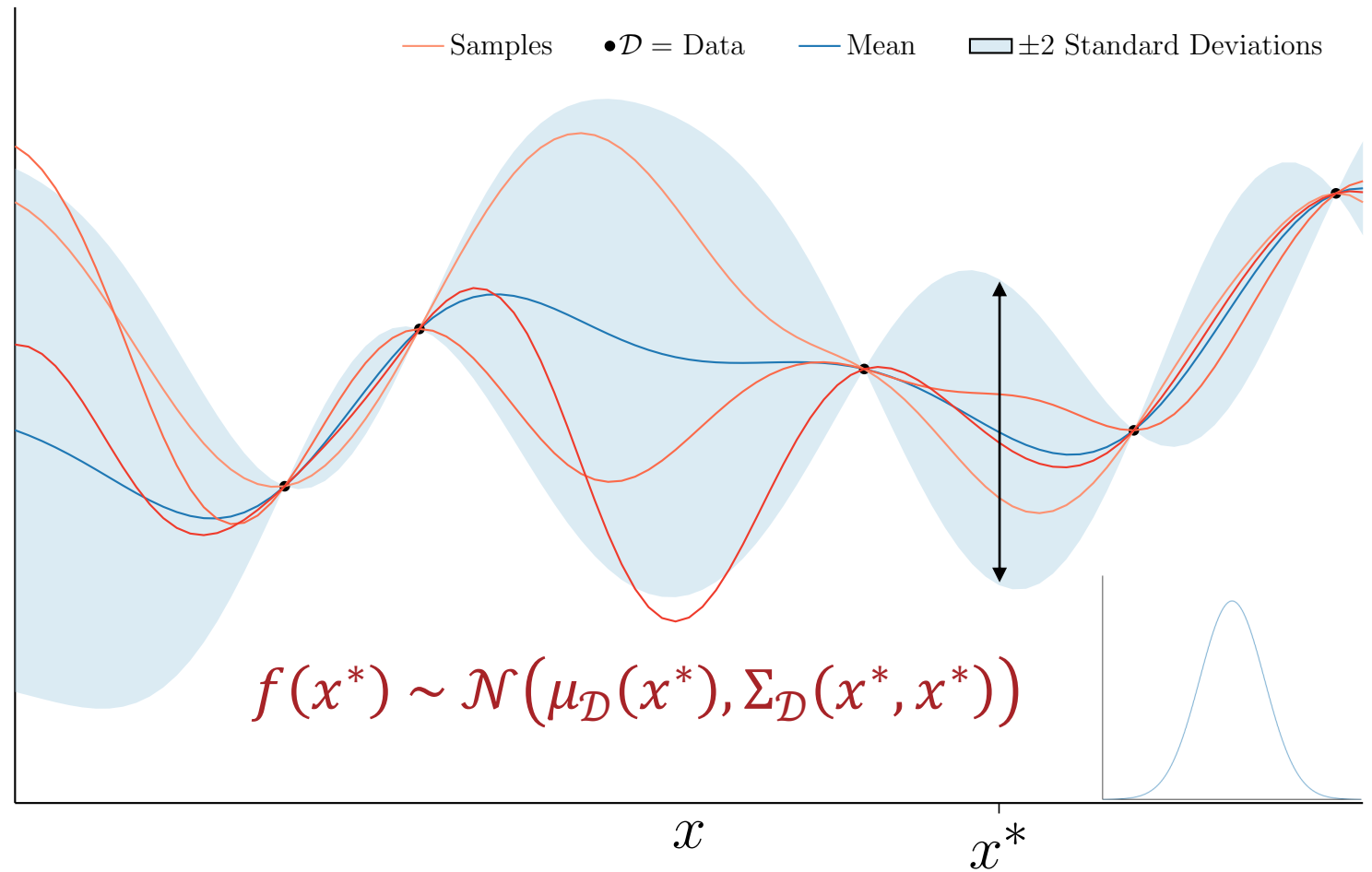
GP Posterior

$$f | \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$$

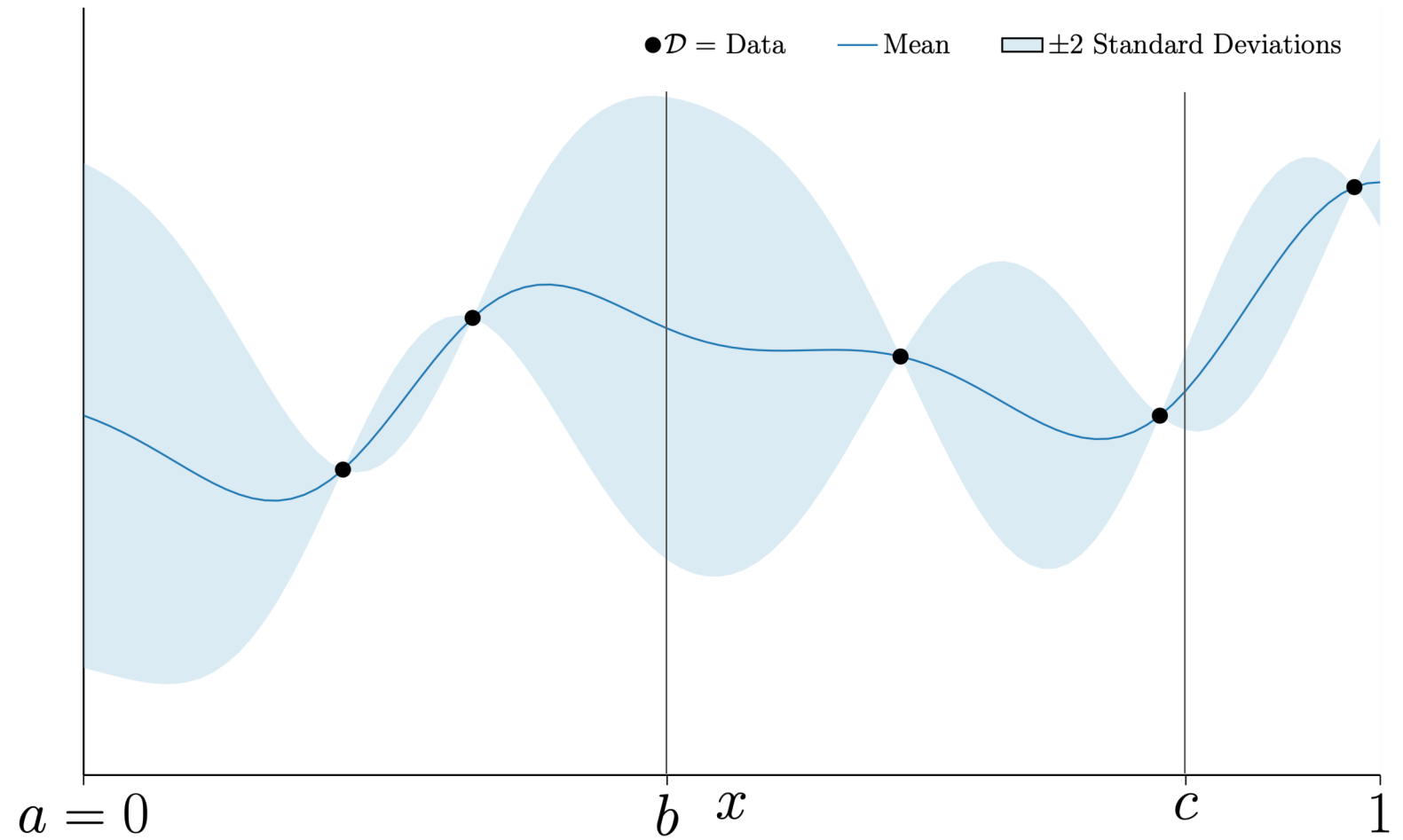


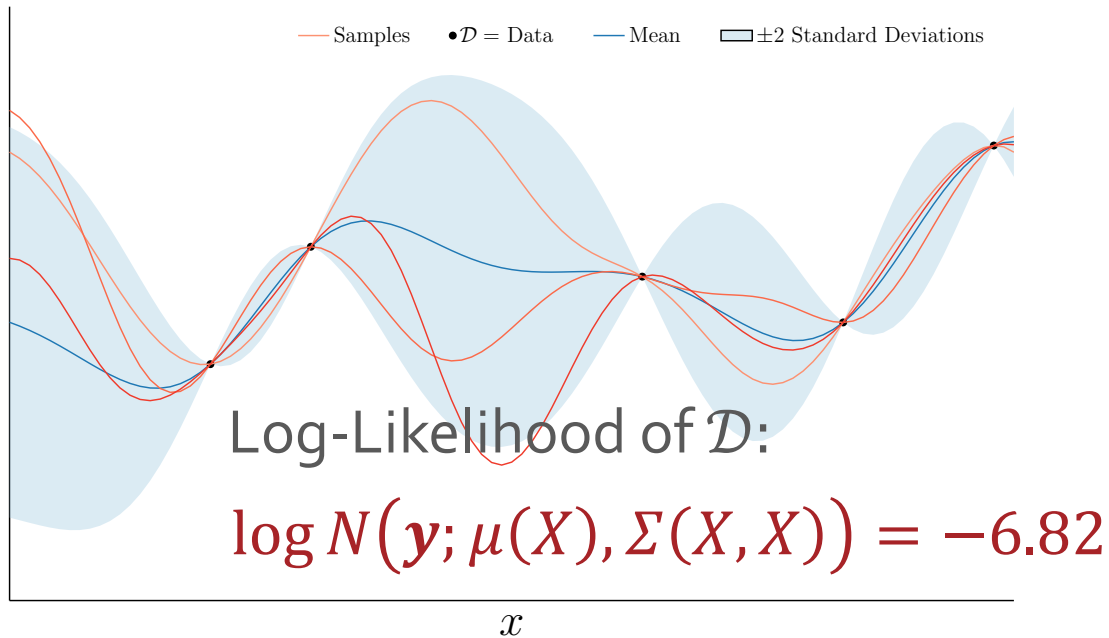
GP Posterior

$$f | \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$$

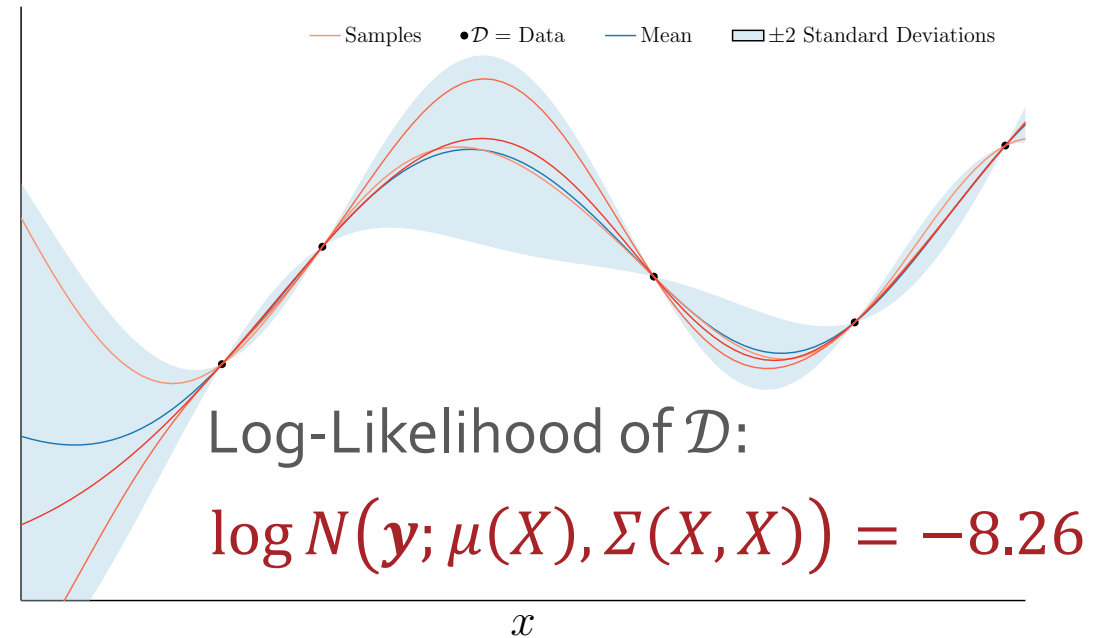


Active Learning





$$f \sim \mathcal{GP} \left(f; 0, (1^2) \exp \left(-\frac{(x - x')^2}{1^2} \right) \right)$$



$$f \sim \mathcal{GP} \left(f; 0, (2^2) \exp \left(-\frac{(x - x')^2}{2^2} \right) \right)$$

Kernel Hyperparameters

- Can be set via MLE
- As long as μ and Σ are differentiable, the log-likelihood is differentiable with respect to the kernel hyperparameters

Noise

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point \mathbf{x}' , the prediction is

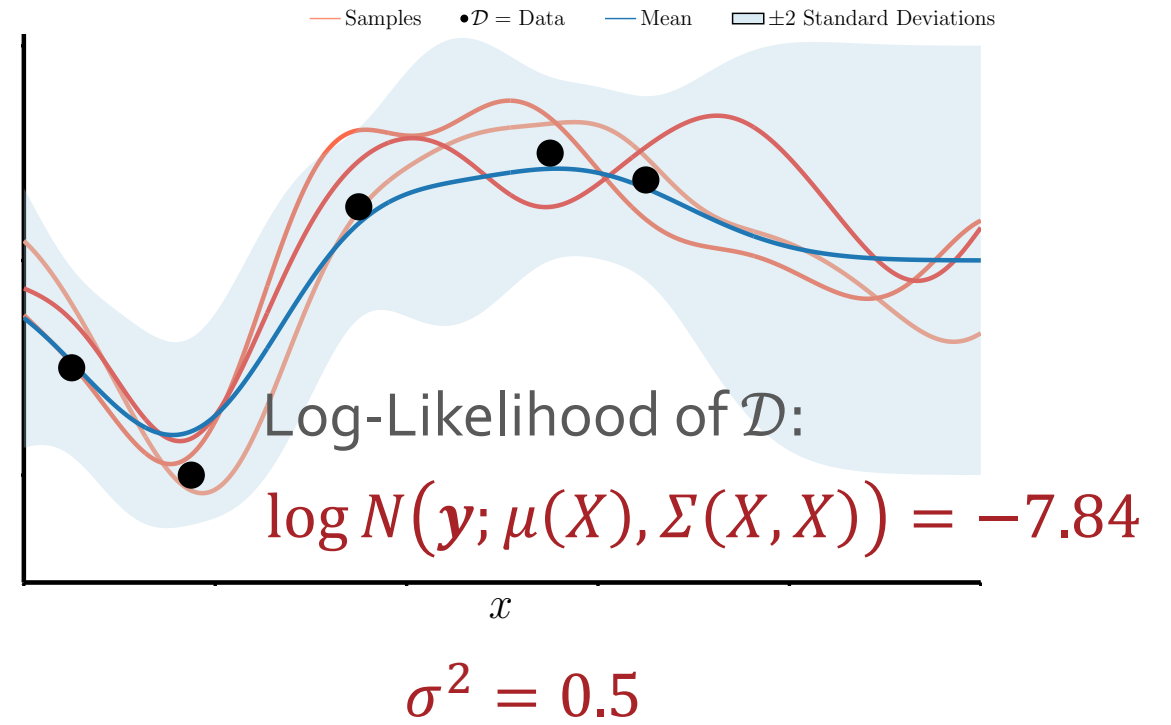
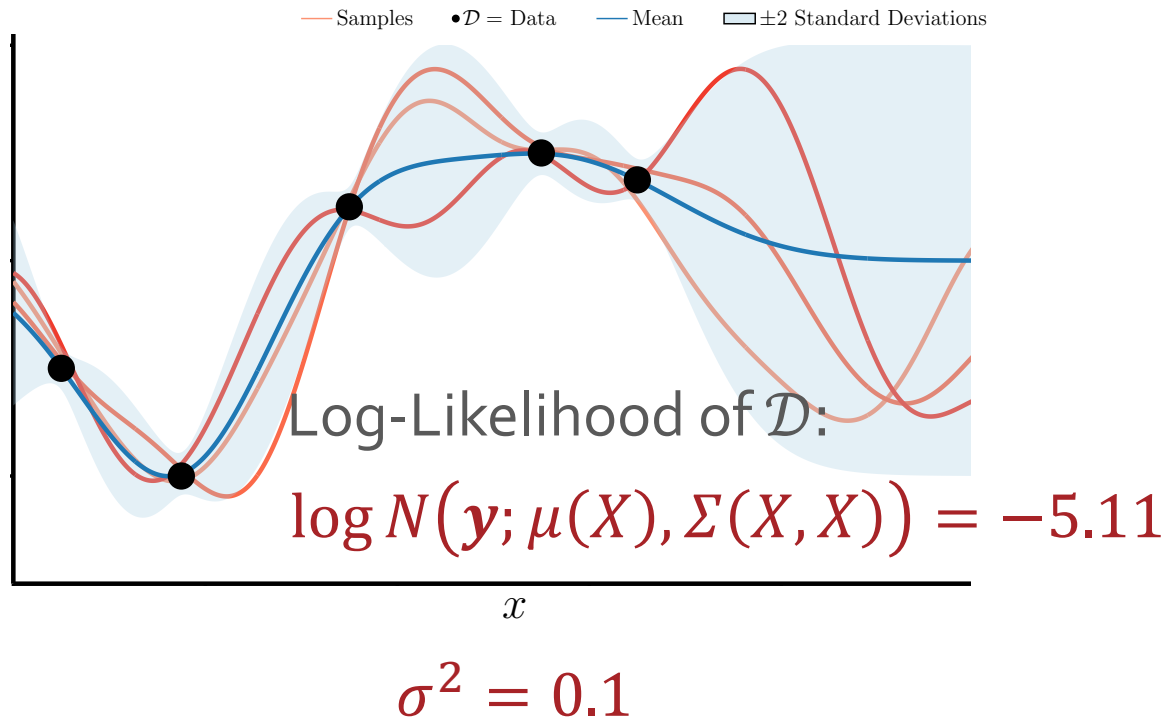
$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = K(\mathbf{x}', \mathbf{x}') - K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} K(X, \mathbf{x})$$

- σ^2 is another hyperparameter we can tune
 - $\sigma^2 = 0$ is a noiseless fit: the mean will always pass through the observations exactly; $\sigma^2 > 0$ allows for deviations



Noise