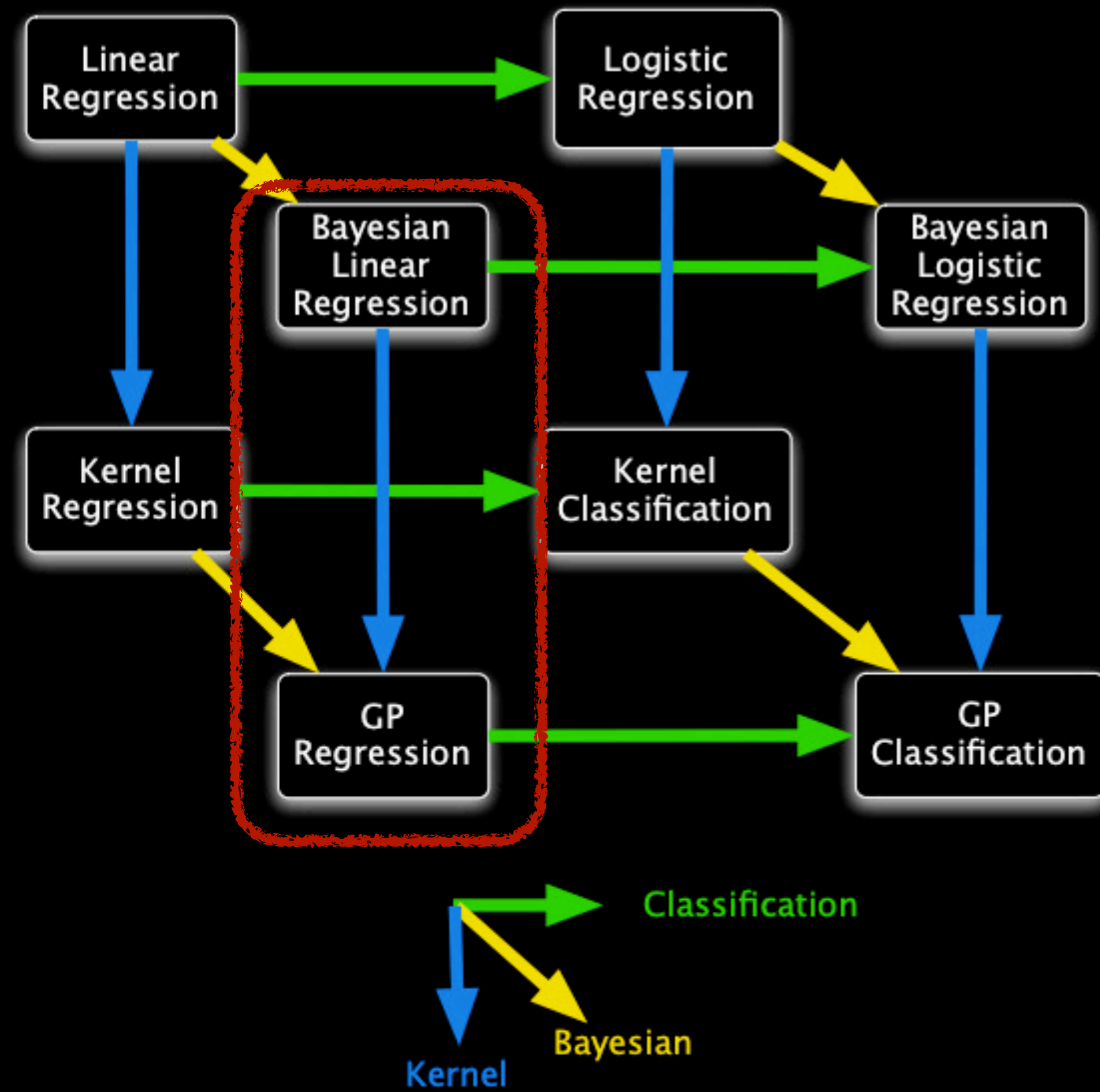


Gaussian process regression

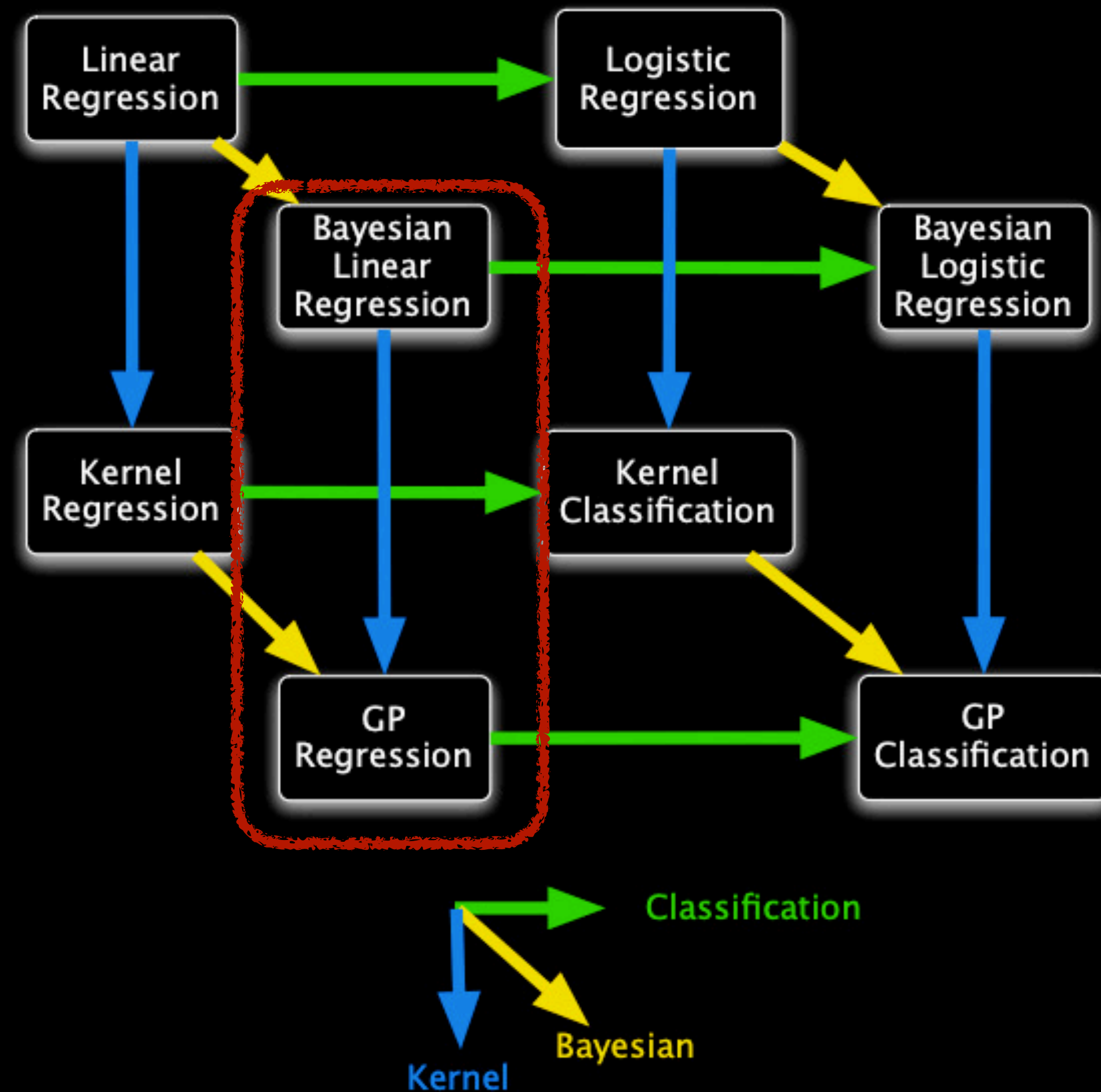
(part 2)

Gaussian Process - Weight-Space Perspective



Gaussian Process - Weight-Space Perspective

Gaussian Process =
Kernelising Bayesian
Linear Regression



Why Kernel ?

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

Why Kernel ?

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^T \mathbf{x}_n$

Why Kernel ?

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^T \mathbf{x}_n \quad \Rightarrow \quad \phi(\mathbf{x}) = \mathbf{x}$

Why Kernel ?

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^T \mathbf{x}_n \quad \Rightarrow \quad \phi(\mathbf{x}) = \mathbf{x}$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2$

Why Kernel ?

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^T \mathbf{x}_n \quad \Rightarrow \quad \phi(\mathbf{x}) = \mathbf{x}$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2 \quad \Rightarrow \quad \begin{array}{l} \text{In 2D} \\ \phi(\mathbf{x}) = \left(x_1^2, \sqrt{2}x_1x_2, x_2^2 \right) \\ \text{(Bishop eq 6.12)} \end{array}$

Why Kernel ?

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^T \mathbf{x}_n \quad \Rightarrow \quad \phi(\mathbf{x}) = \mathbf{x}$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2 \quad \Rightarrow \quad \begin{matrix} \text{In 2D} \\ \phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \end{matrix}$
(Bishop eq 6.12)

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|_2^2}{2\sigma^2}\right)$

(Bishop eq 6.23, GP Book eq 2.16)

Why Kernel ?

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^T \mathbf{x}_n \quad \Rightarrow \quad \phi(\mathbf{x}) = \mathbf{x}$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2 \quad \Rightarrow \quad \phi(\mathbf{x}) = \begin{pmatrix} x_1^2, \sqrt{2}x_1x_2, x_2^2 \end{pmatrix}$
(Bishop eq 6.12)

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|_2^2}{2\sigma^2}\right)$

(Bishop eq 6.23, GP Book eq 2.16)

\Rightarrow Infinite dimensional features

Why Kernel ?

Simplifying the process of coming up with “features”

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n)$$

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^T \mathbf{x}_n \quad \Rightarrow \quad \phi(\mathbf{x}) = \mathbf{x}$

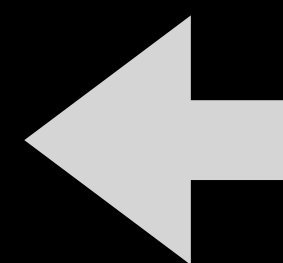
e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2 \quad \Rightarrow \quad \begin{matrix} \text{In 2D} \\ \phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \end{matrix}$
(Bishop eq 6.12)

e.g., $k(\mathbf{x}_m, \mathbf{x}_n) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|_2^2}{2\sigma^2}\right)$

(Bishop eq 6.23, GP Book eq 2.16)

\Rightarrow Infinite dimensional features

Gaussian Process



$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y})$$

Kernelised

$$= \mathcal{N}(y^*; m(x^*), \sigma^2(x^*))$$

$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\sigma^2(x^*) = \sigma^2 + k(\mathbf{x}^*, \mathbf{x}^*) \quad (\text{Bishop eq 6.66})$$

$$-k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

(Bishop eq 6.67)

Bayesian Linear Regression

$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y})$$

$$= \mathcal{N}(y^*; \mu^T \phi(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$$

$$\mu = \sigma^{-2} (\sigma_0^{-2} \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

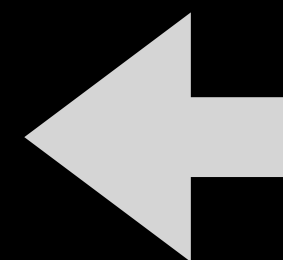
(Bishop eq 3.53)

$$\sigma^2(\mathbf{x}^*) = \sigma^2 + \phi(\mathbf{x}^*)^T \Sigma \phi(\mathbf{x}^*)$$

$$\Sigma = (\sigma_0^{-2} \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi)^{-1}$$

(Bishop eq 3.54)

Gaussian Process



$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y})$$

Kernelised

$$= \mathcal{N}(y^*; m(x^*), \sigma^2(x^*))$$

$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$
$$\sigma^2(x^*) = \sigma^2 + k(\mathbf{x}^*, \mathbf{x}^*) \quad (\text{Bishop eq 6.66})$$

$$-k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

(Bishop eq 6.67)

Inverse of an $\mathbb{R}^{N \times N}$ matrix

Bayesian Linear Regression

$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y})$$

$$= \mathcal{N}(y^*; \mu^T \phi(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$$

$$\mu = \sigma^{-2} (\sigma_0^{-2} \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

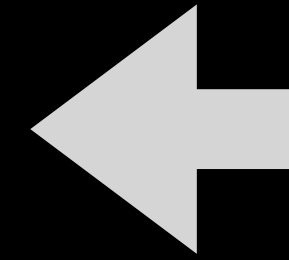
(Bishop eq 3.53)

$$\sigma^2(\mathbf{x}^*) = \sigma^2 + \phi(\mathbf{x}^*)^T \Sigma \phi(\mathbf{x}^*)$$

$$\Sigma = (\sigma_0^{-2} \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi)^{-1}$$

(Bishop eq 3.54)

Gaussian Process



$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y})$$

Kernelised

$$= \mathcal{N}(y^*; m(x^*), \sigma^2(x^*))$$

$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$
$$\sigma^2(x^*) = \sigma^2 + k(\mathbf{x}^*, \mathbf{x}^*) \quad (\text{Bishop eq 6.66})$$

$$-k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

(Bishop eq 6.67)

Inverse of an $\mathbb{R}^{N \times N}$ matrix

Bayesian Linear Regression

$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y})$$

$$= \mathcal{N}(y^*; \mu^T \phi(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$$

$$\mu = \sigma^{-2} (\sigma_0^{-2} \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

(Bishop eq 3.53)

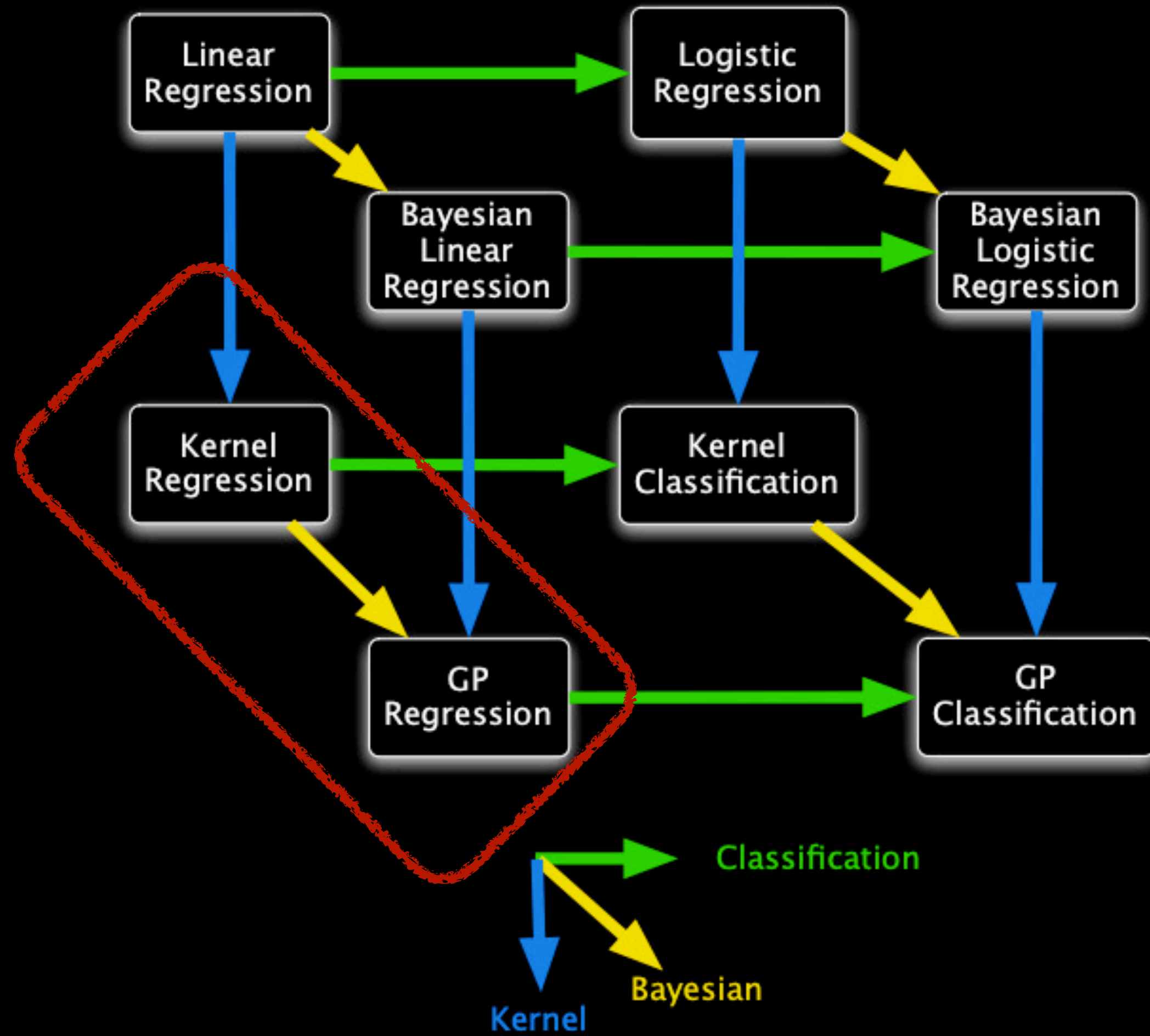
$$\sigma^2(\mathbf{x}^*) = \sigma^2 + \phi(\mathbf{x}^*)^T \Sigma \phi(\mathbf{x}^*)$$

$$\Sigma = (\sigma_0^{-2} \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi)^{-1}$$

(Bishop eq 3.54)

Inverse of an $\mathbb{R}^{D \times D}$ matrix

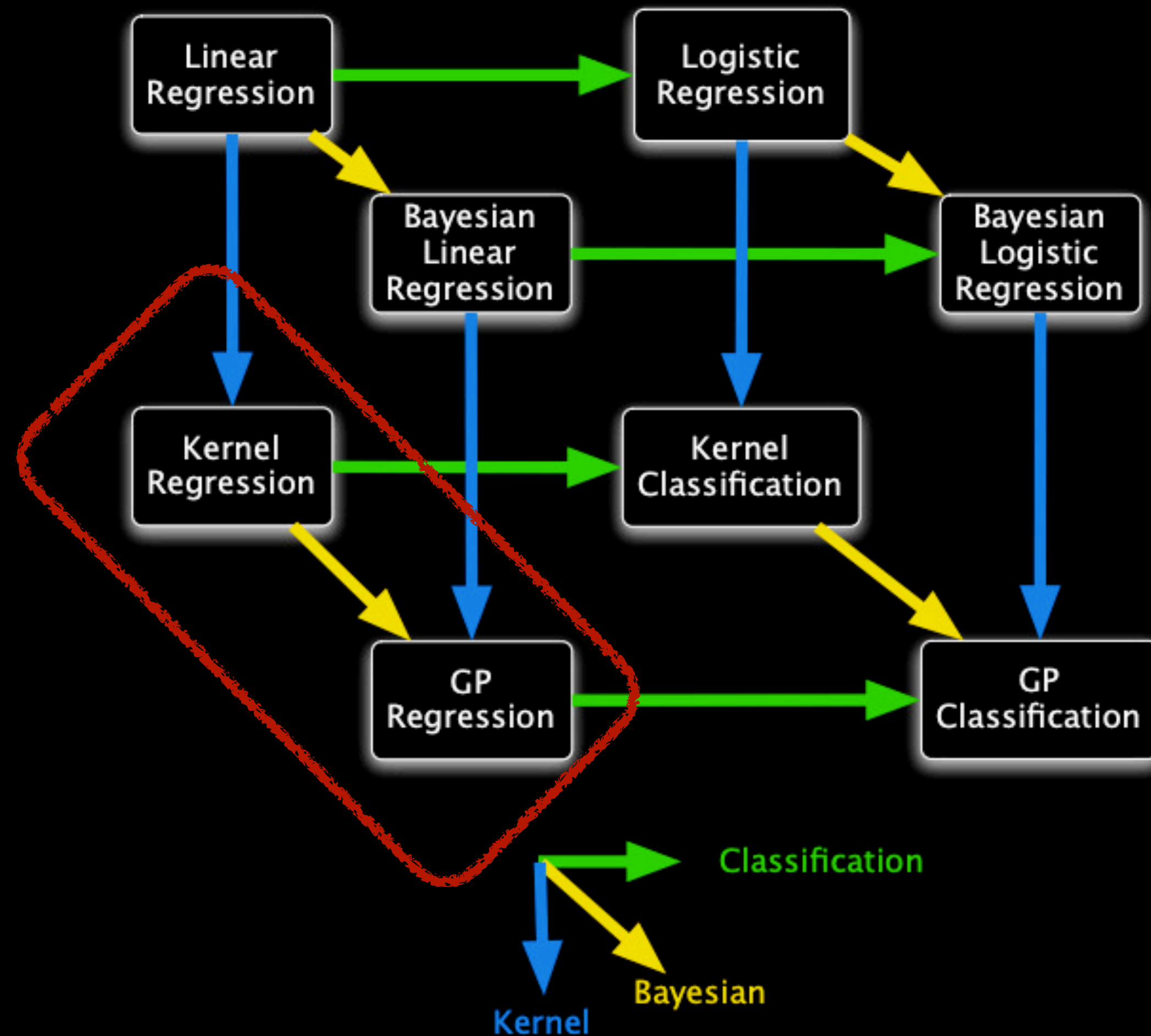
Gaussian Process - Function-Space Perspective



Gaussian Process - Function-Space Perspective

Gaussian Process =

Making Kernel
Regression “Bayesian”



Gaussian Process - Function-Space Perspective

Gaussian Process - Function-Space Perspective

- A Gaussian Process is a probability distribution over functions $f(\mathbf{x})$

such that $\forall (\mathbf{x}_1, \dots, \mathbf{x}_n), p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) = \mathcal{N}$

Gaussian Process - Function-Space Perspective

- A **Gaussian Process** is a probability distribution over functions $f(\mathbf{x})$
such that $\forall(\mathbf{x}_1, \dots, \mathbf{x}_n), p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) = \mathcal{N}$
- Commonly set mean to zero, due to no prior knowledge of $f(x)$

Gaussian Process - Function-Space Perspective

- A **Gaussian Process** is a probability distribution over functions $f(\mathbf{x})$
such that $\forall(\mathbf{x}_1, \dots, \mathbf{x}_n), p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) = \mathcal{N}$
- Commonly set mean to zero, due to no prior knowledge of $f(x)$
- Fully specified by their **covariance function** or **kernel**.

Gaussian Process - Function-Space Perspective

- A **Gaussian Process** is a probability distribution over functions $f(\mathbf{x})$

such that $\forall(\mathbf{x}_1, \dots, \mathbf{x}_n), p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) = \mathcal{N}$

- Commonly set mean to zero, due to no prior knowledge of $f(x)$
- Fully specified by their **covariance function** or **kernel**.

$$\Rightarrow p(\mathbf{f}(\mathbf{x})) = \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K})$$

(Bishop eq 6.60, GP Book eq 2.17)

Gaussian Process - Function-Space Perspective

- A **Gaussian Process** is a probability distribution over functions $f(\mathbf{x})$

such that $\forall(\mathbf{x}_1, \dots, \mathbf{x}_n), p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) = \mathcal{N}$

- Commonly set mean to zero, due to no prior knowledge of $f(x)$
- Fully specified by their **covariance function** or **kernel**.

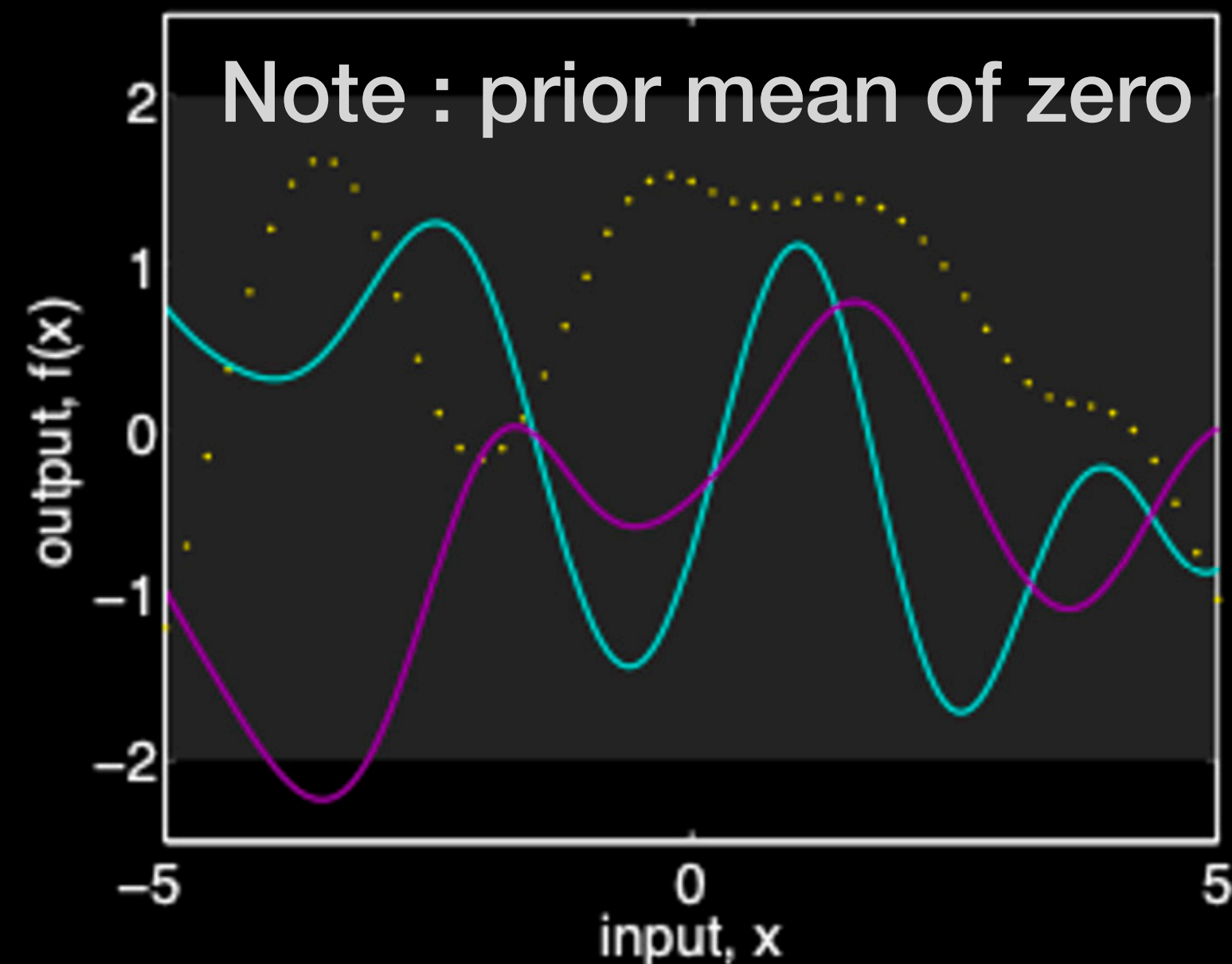
$$\Rightarrow p(\mathbf{f}(\mathbf{x})) = \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K})$$

(Bishop eq 6.60, GP Book eq 2.17)

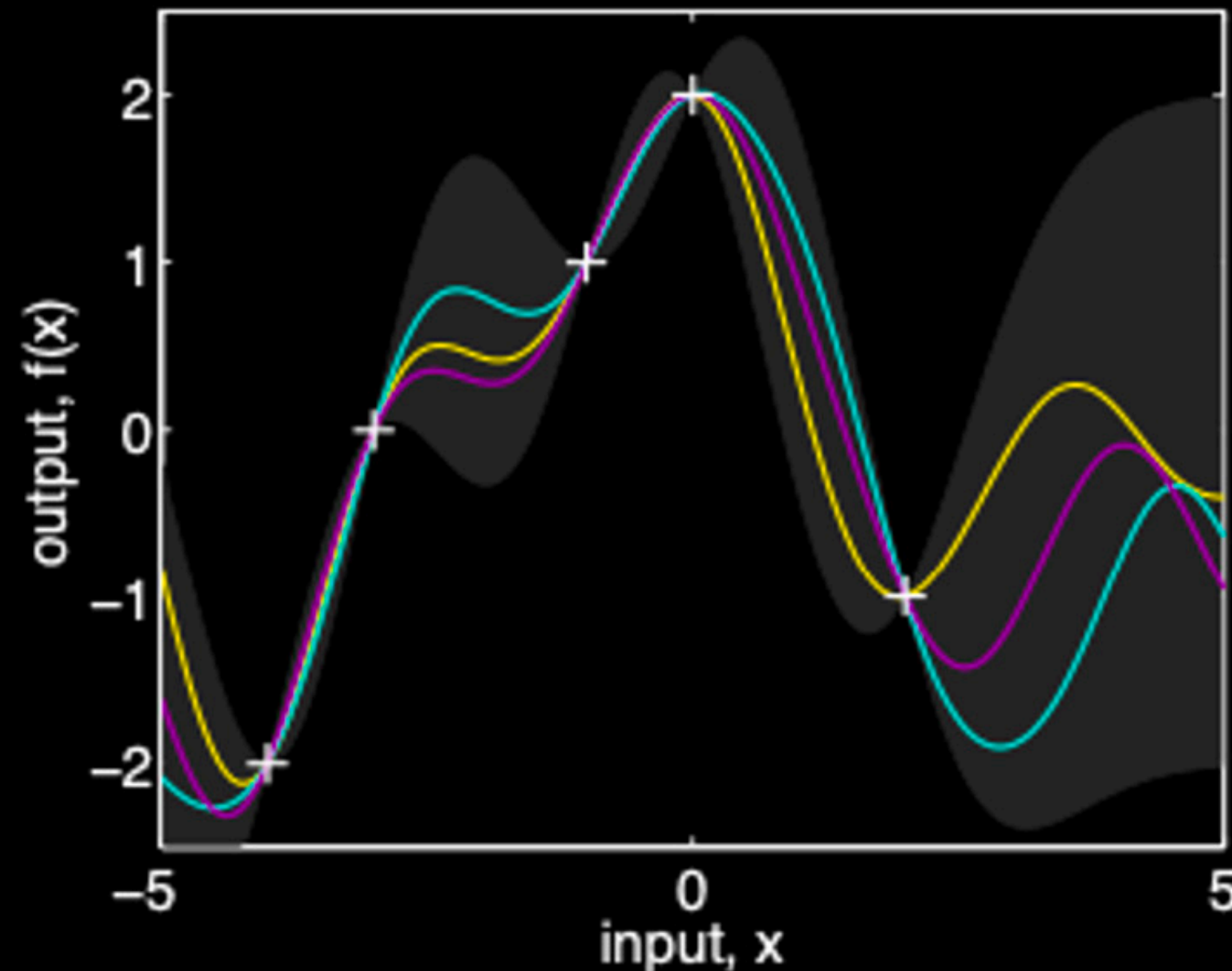
Assuming the prior “weight”

$p(w)$ to have mean zero

Function-Space Perspective $p(y^* | y) = \mathcal{N}(m(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$



(a), prior



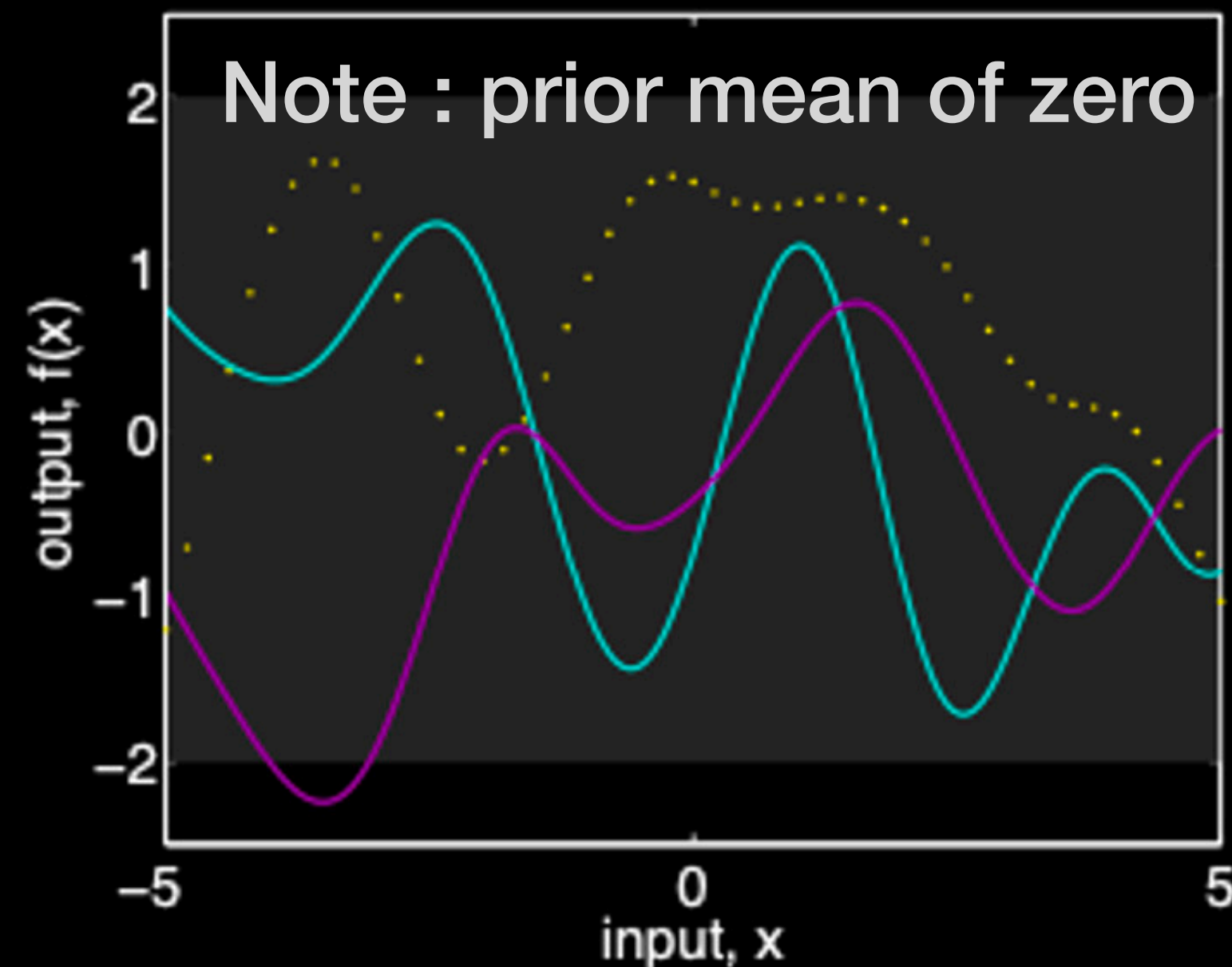
(b), posterior

Posterior distribution of functions : “rejecting” functions that do not satisfy the constraints

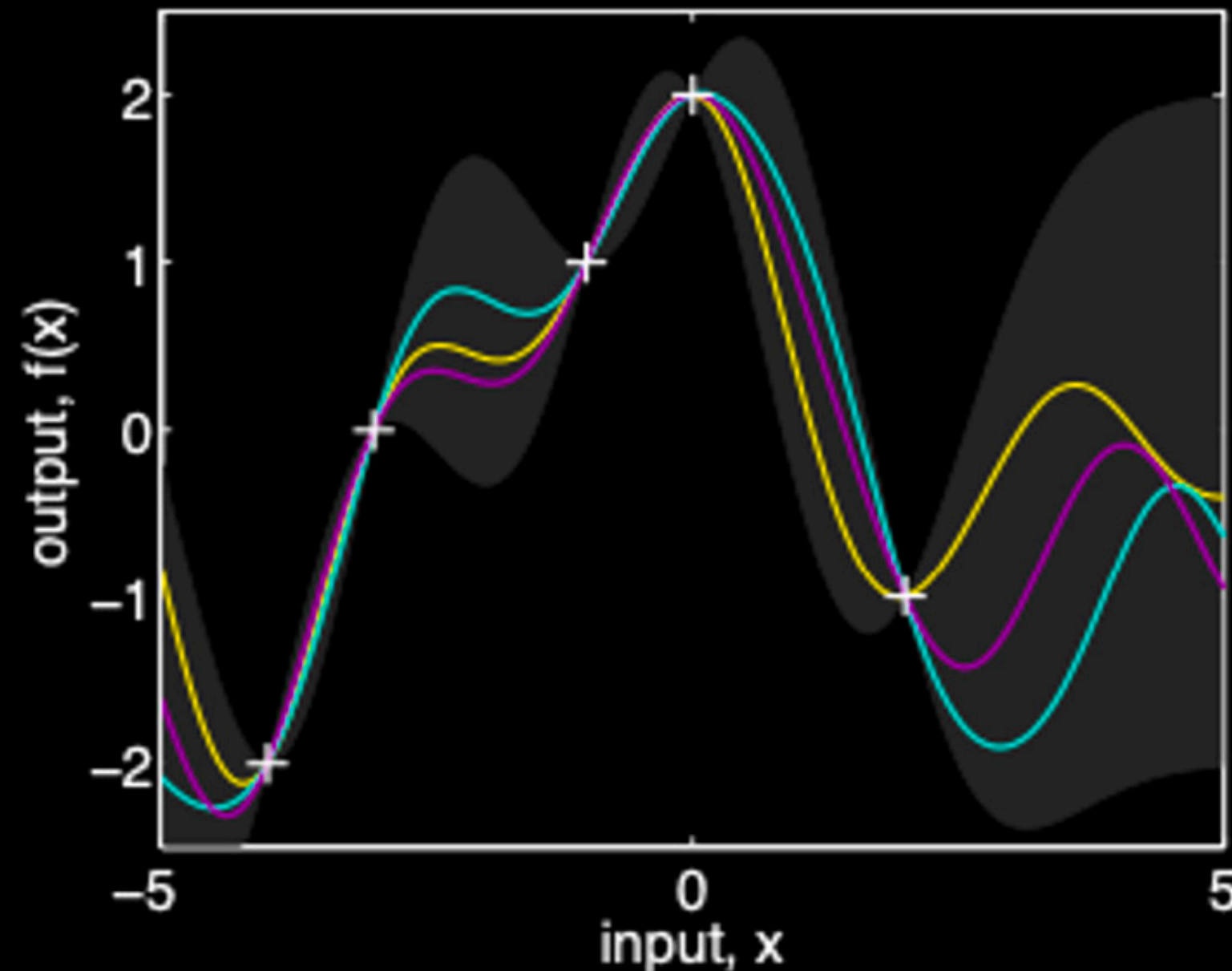
$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\sigma^2(\mathbf{x}^*) = \sigma^2 + k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

Function-Space Perspective $p(y^* | y) = \mathcal{N}(m(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$



(a), prior



(b), posterior

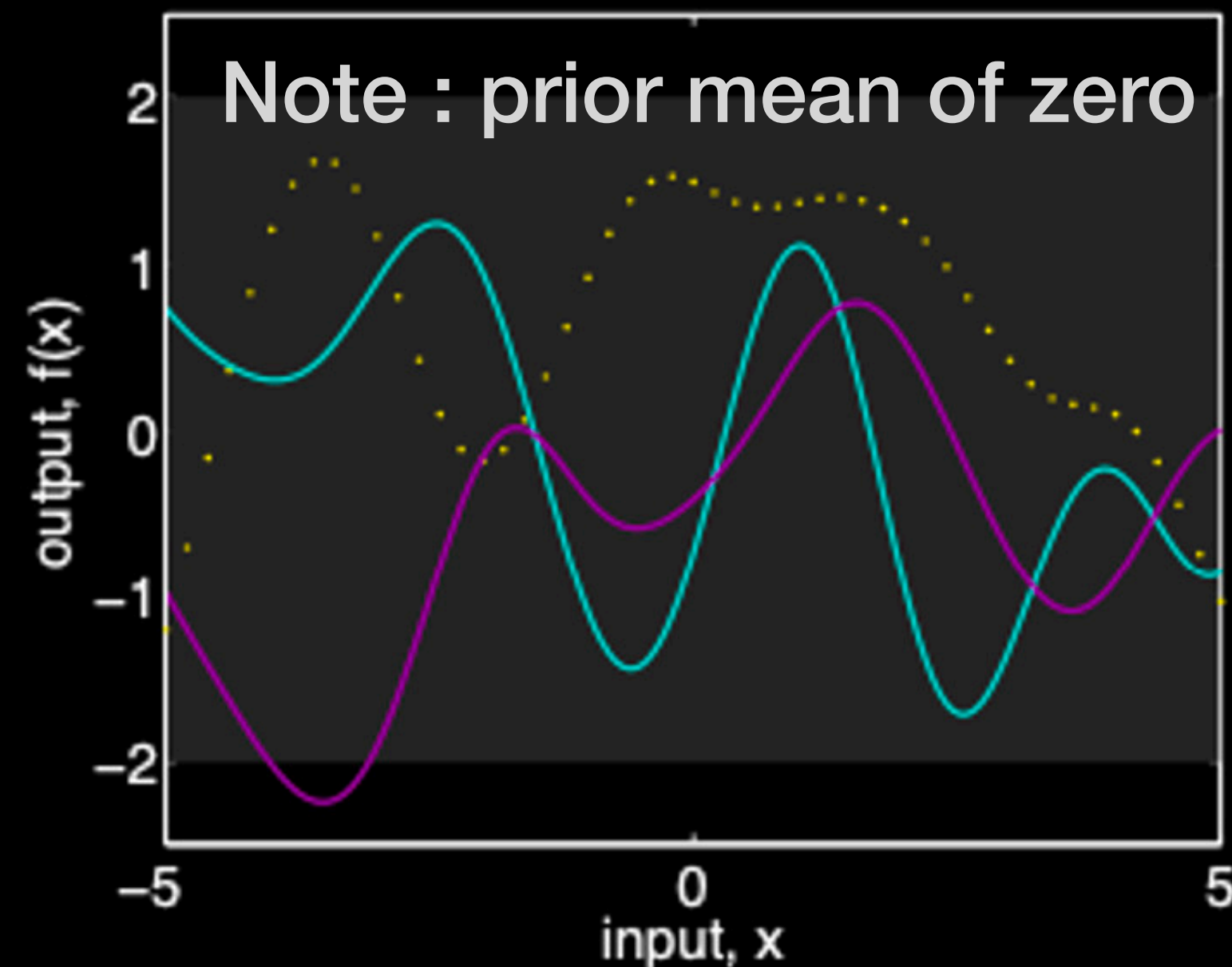
Posterior distribution of functions : “rejecting” functions that do not satisfy the constraints

$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

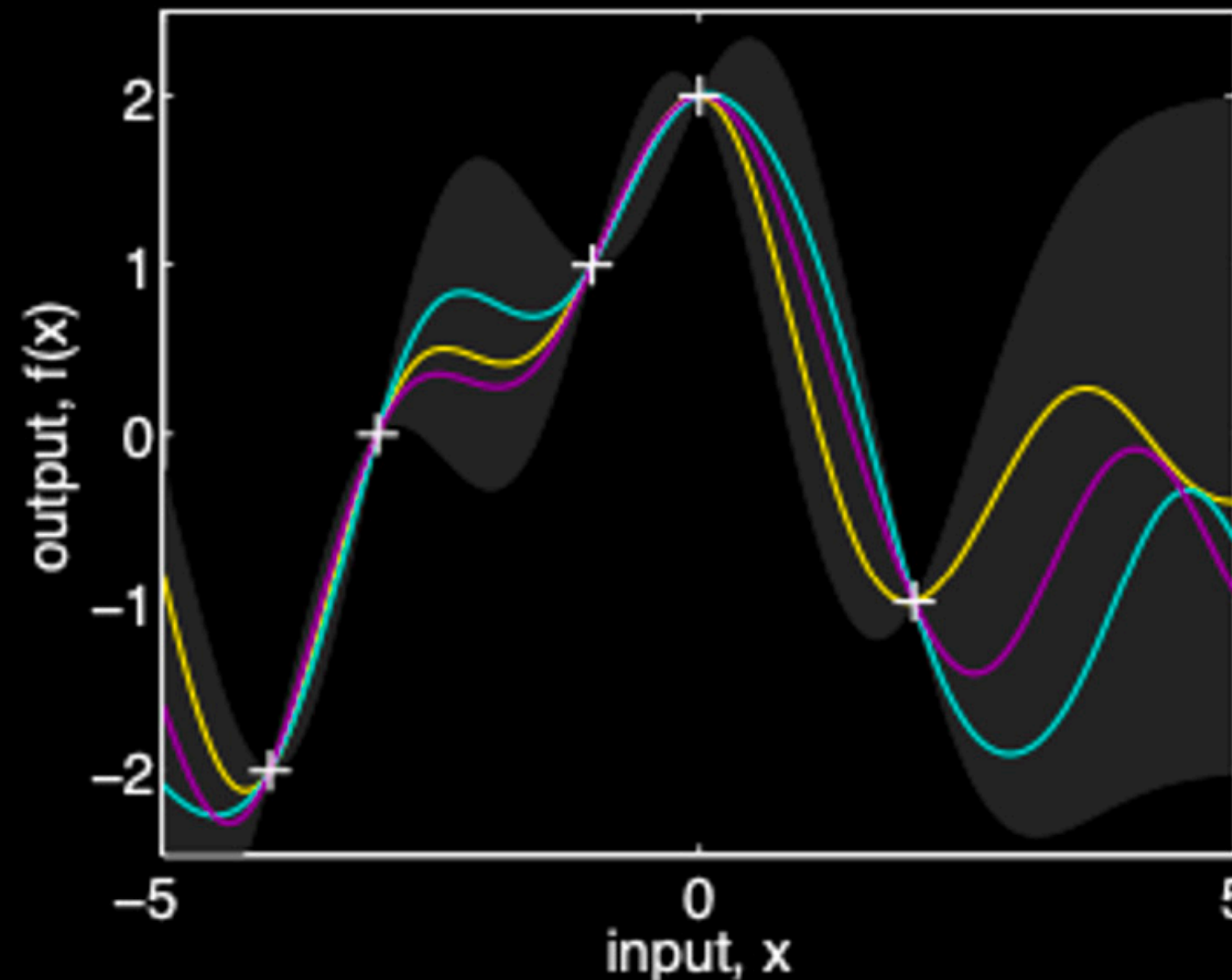
The influence depends on the proximity

$$\sigma^2(\mathbf{x}^*) = \sigma^2 + k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

Function-Space Perspective $p(y^* | y) = \mathcal{N}(m(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$



(a), prior



(b), posterior

Posterior distribution of functions : “rejecting” functions that do not satisfy the constraints

$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

The influence depends on the proximity

$$\sigma^2(\mathbf{x}^*) = \sigma^2 + k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

Uncertainty reduction given the training data, does not depend on y !

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

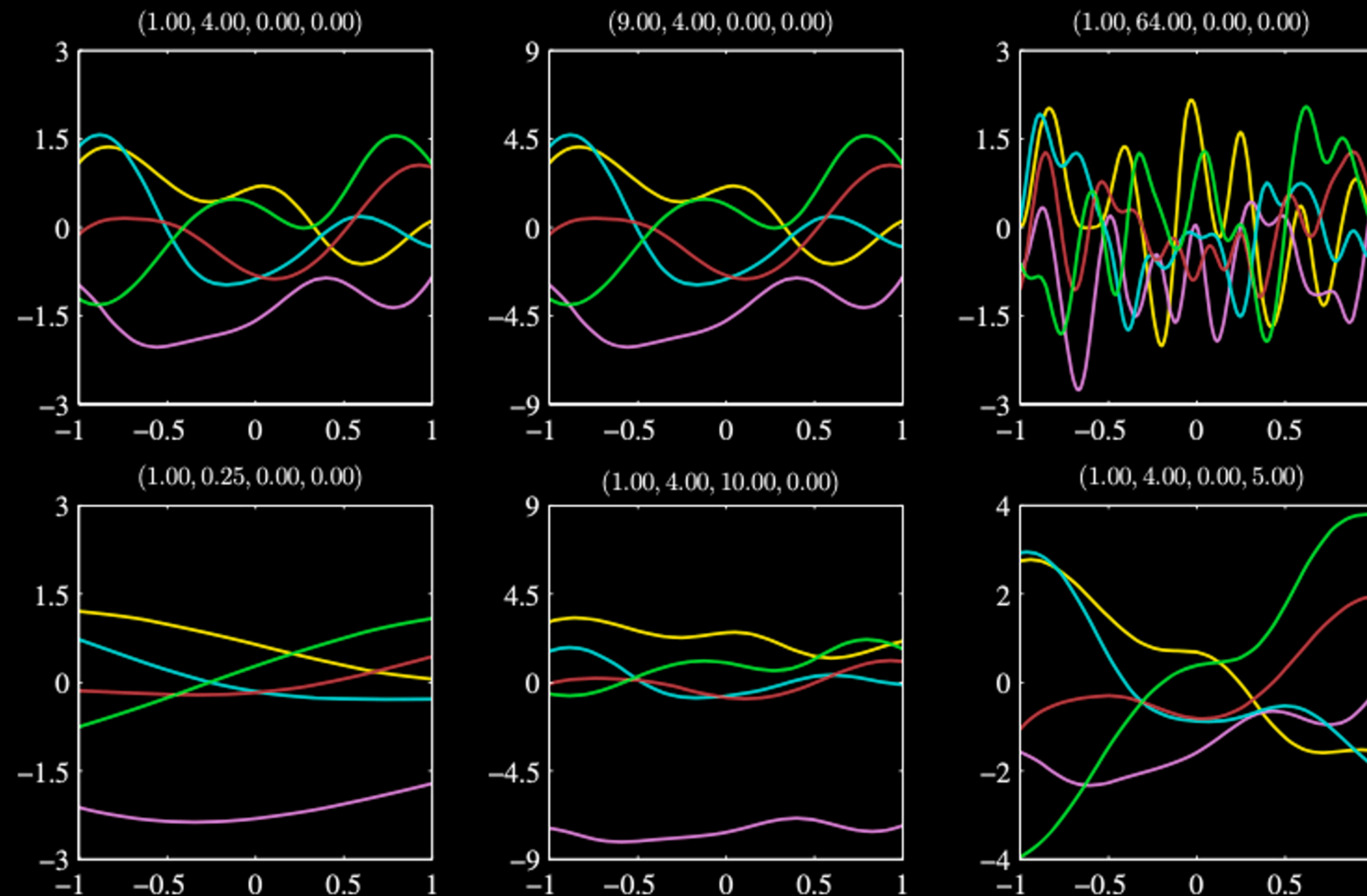


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

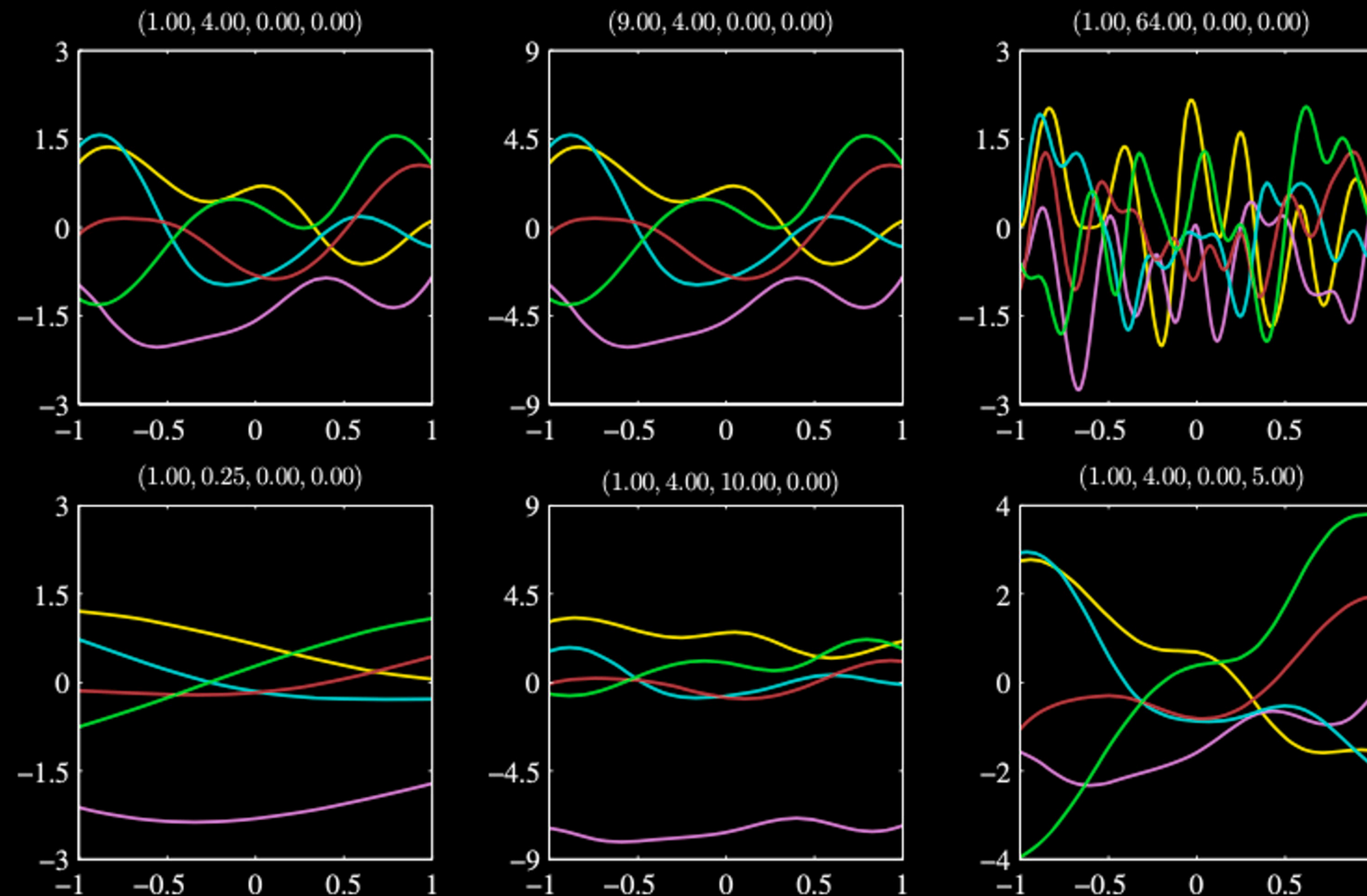


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

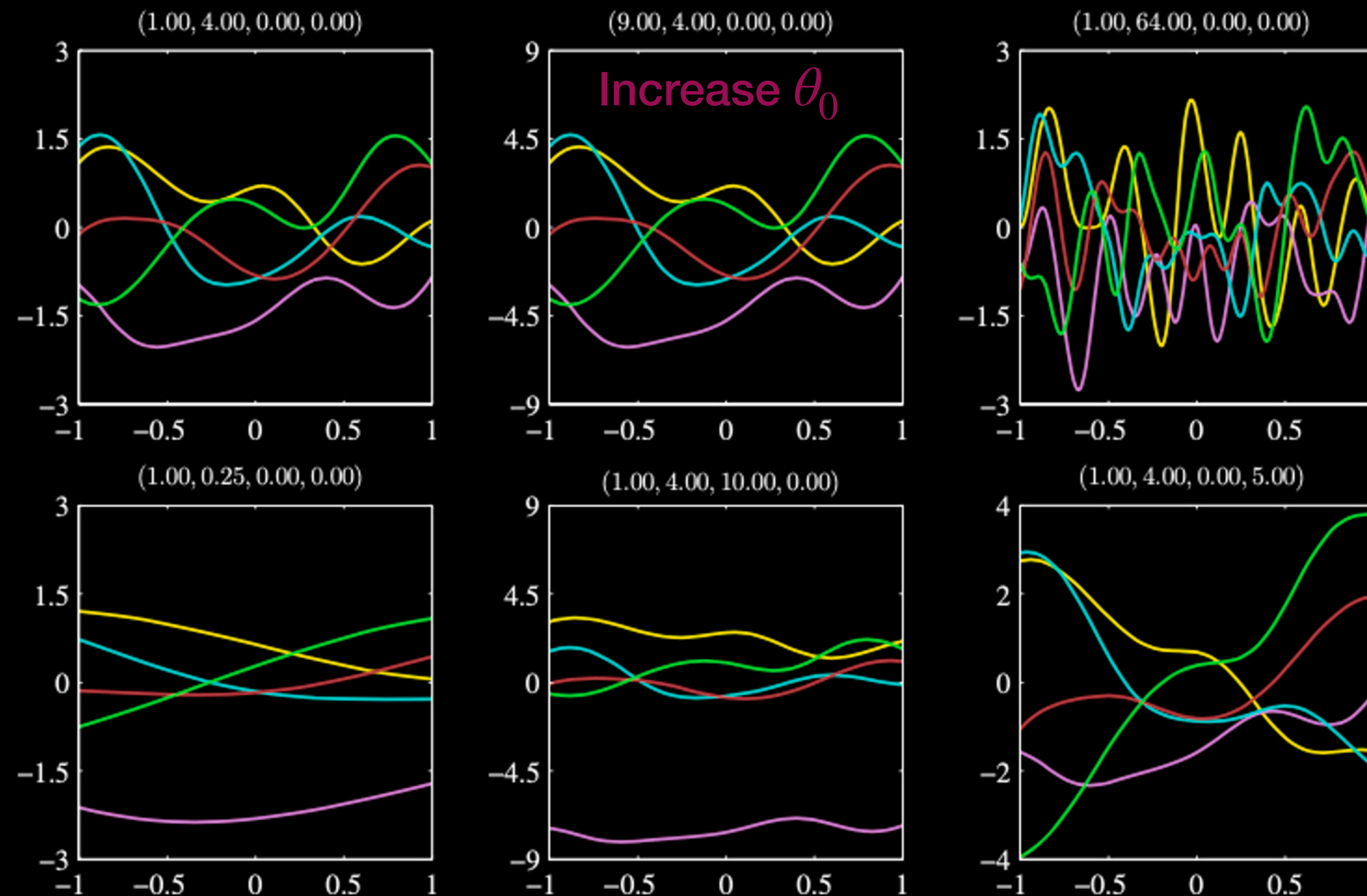


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

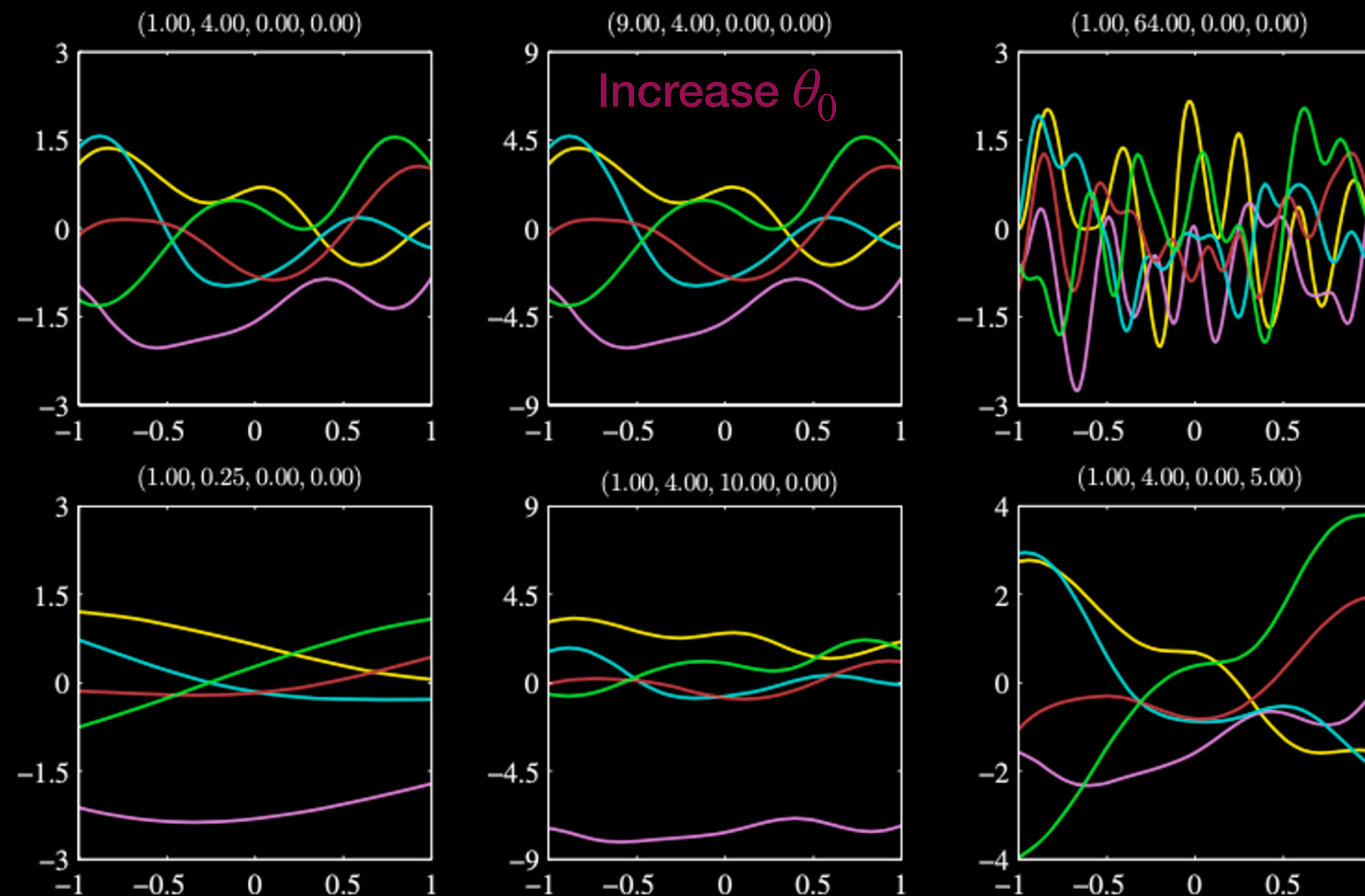


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

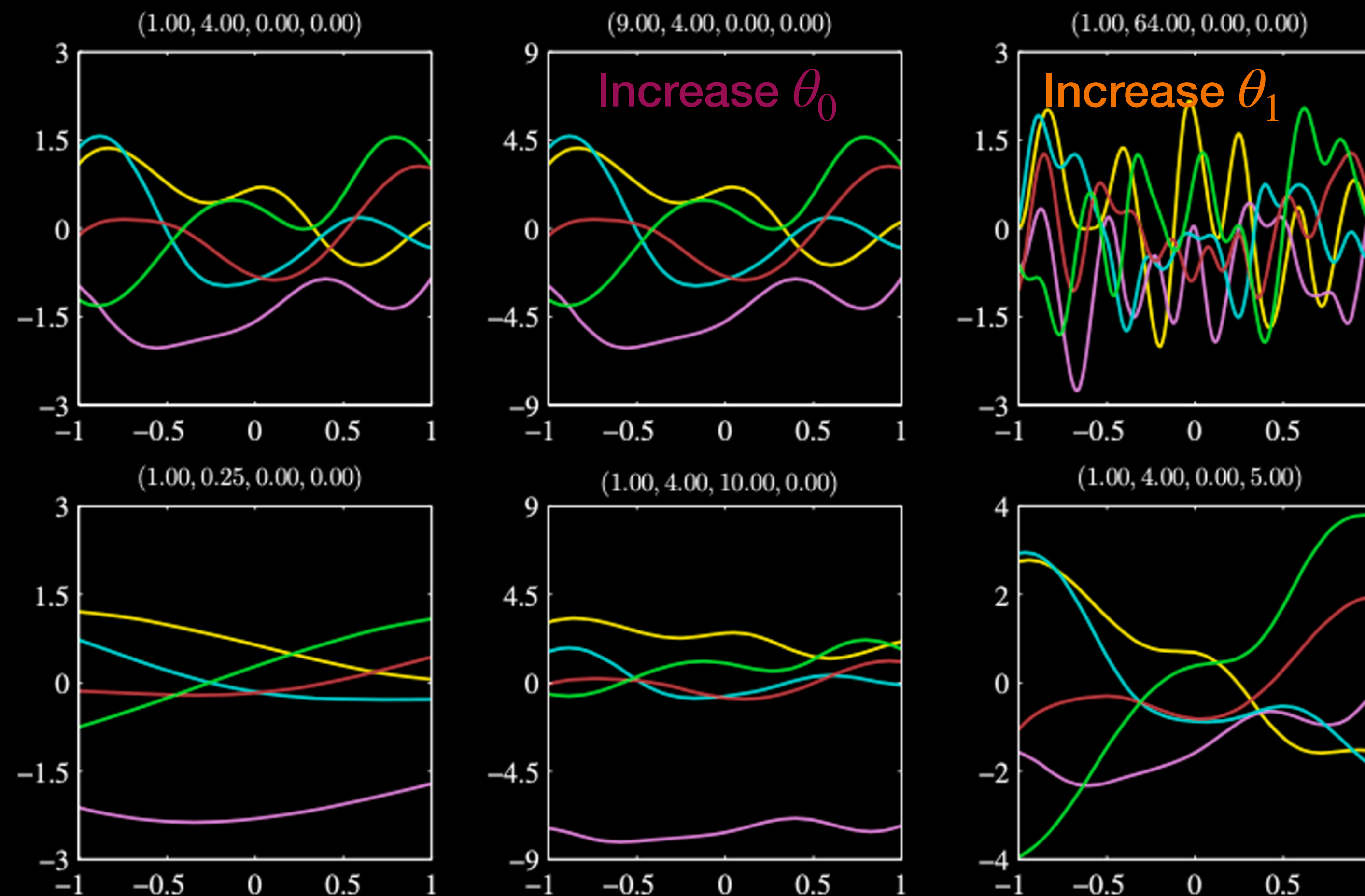


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

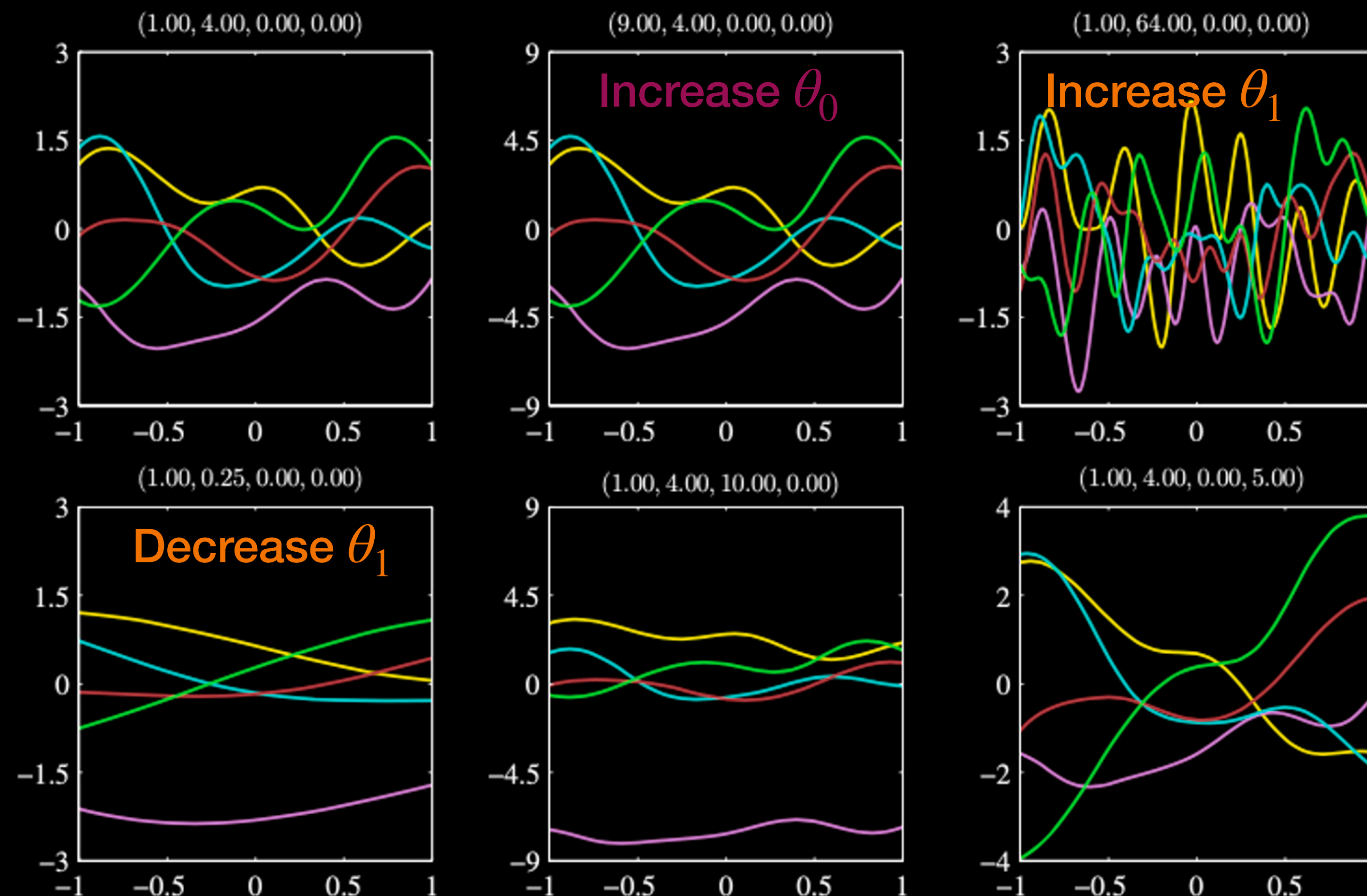


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

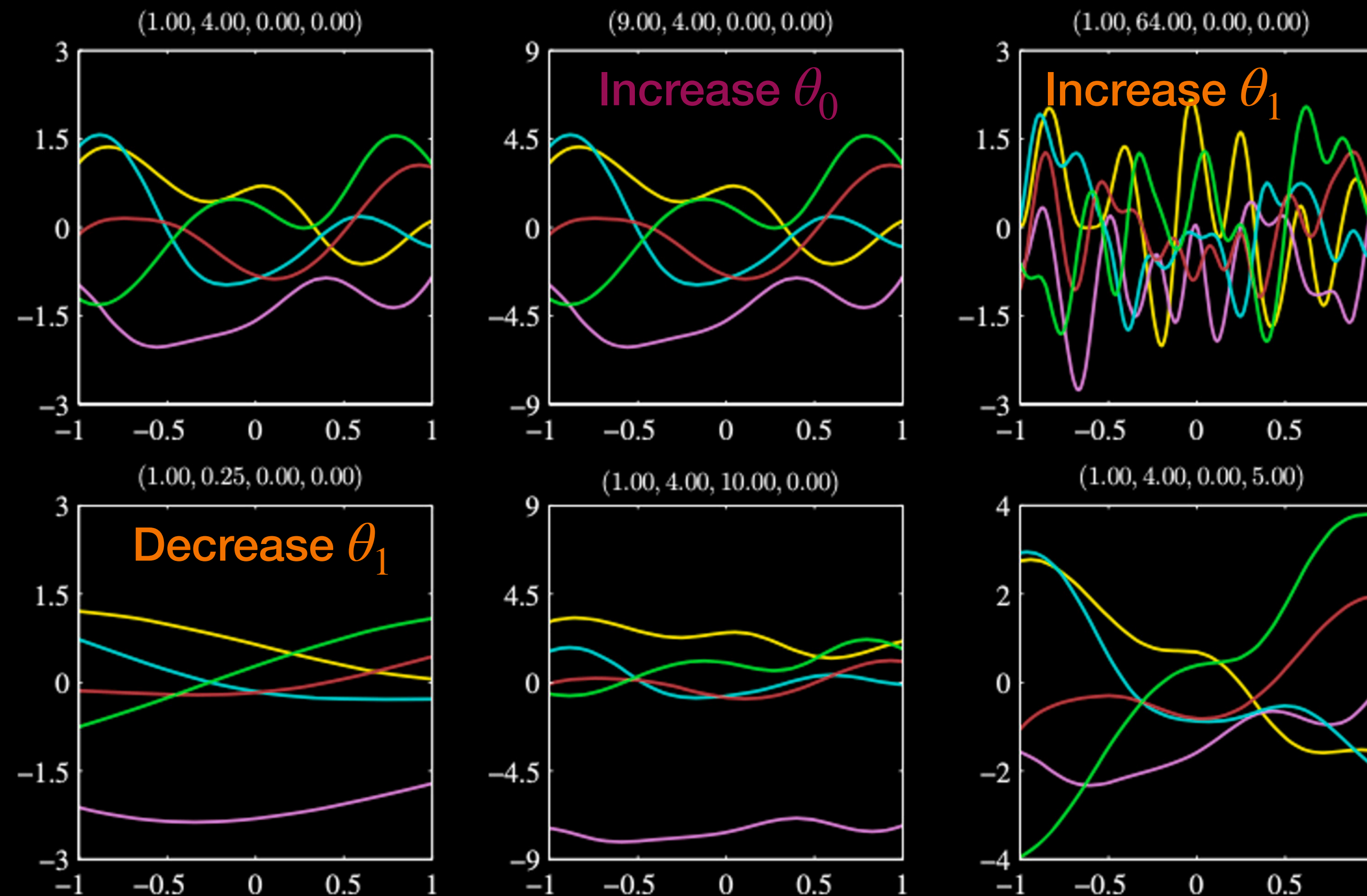


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

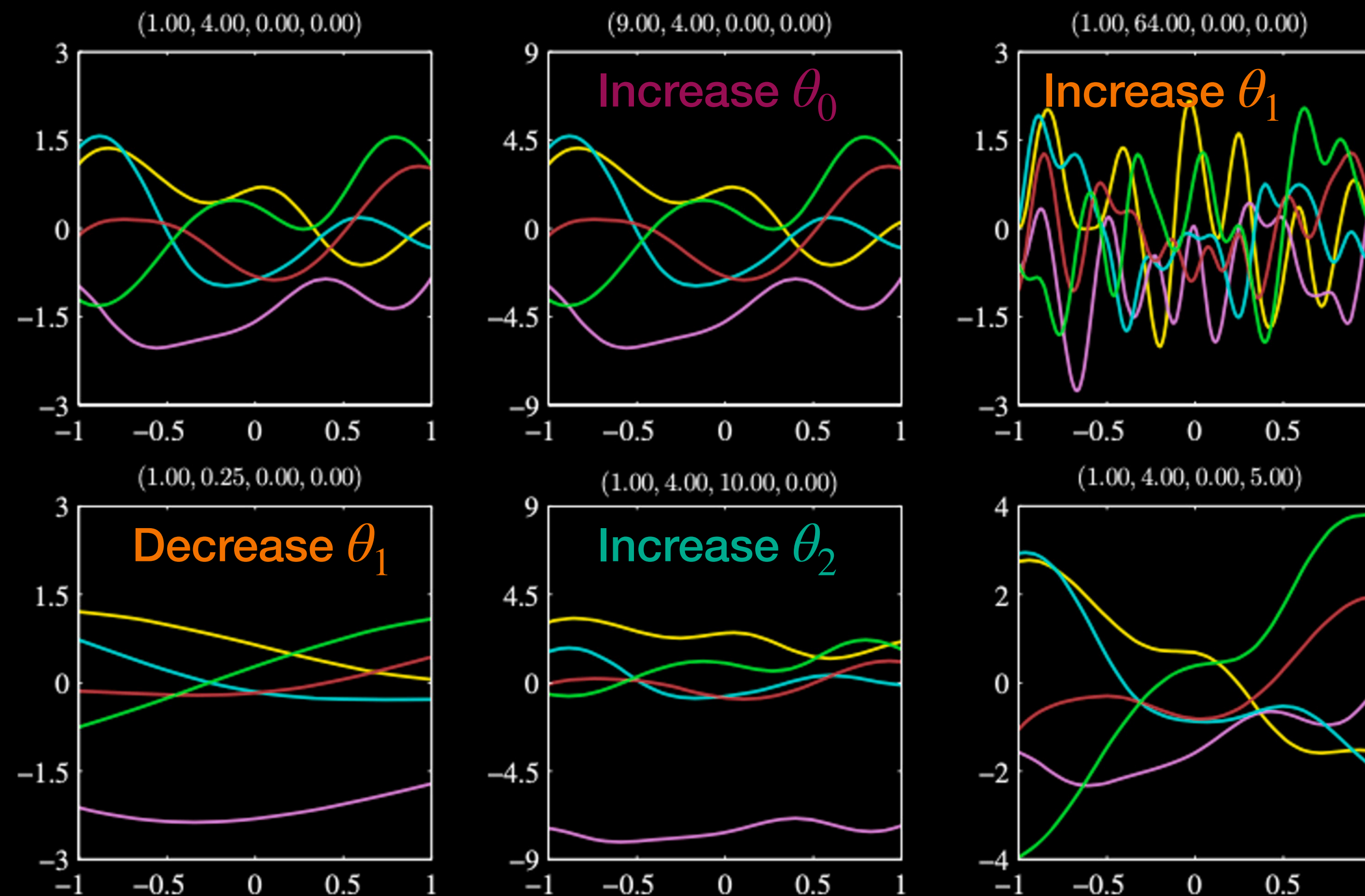


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

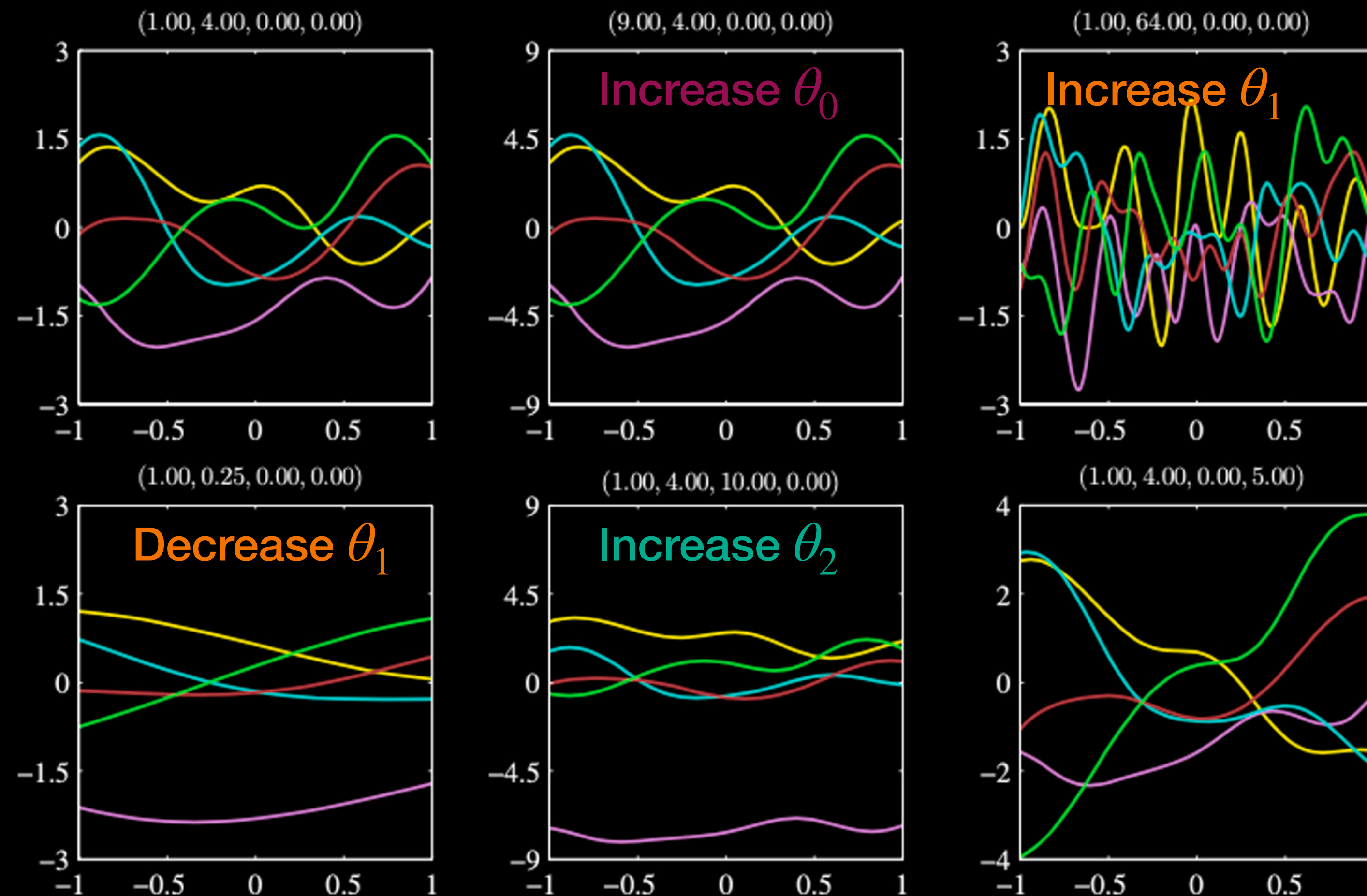


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Visualising Kernels

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Randomly sampling
from the Gaussian
Process prior
distribution

(Bishop textbook)

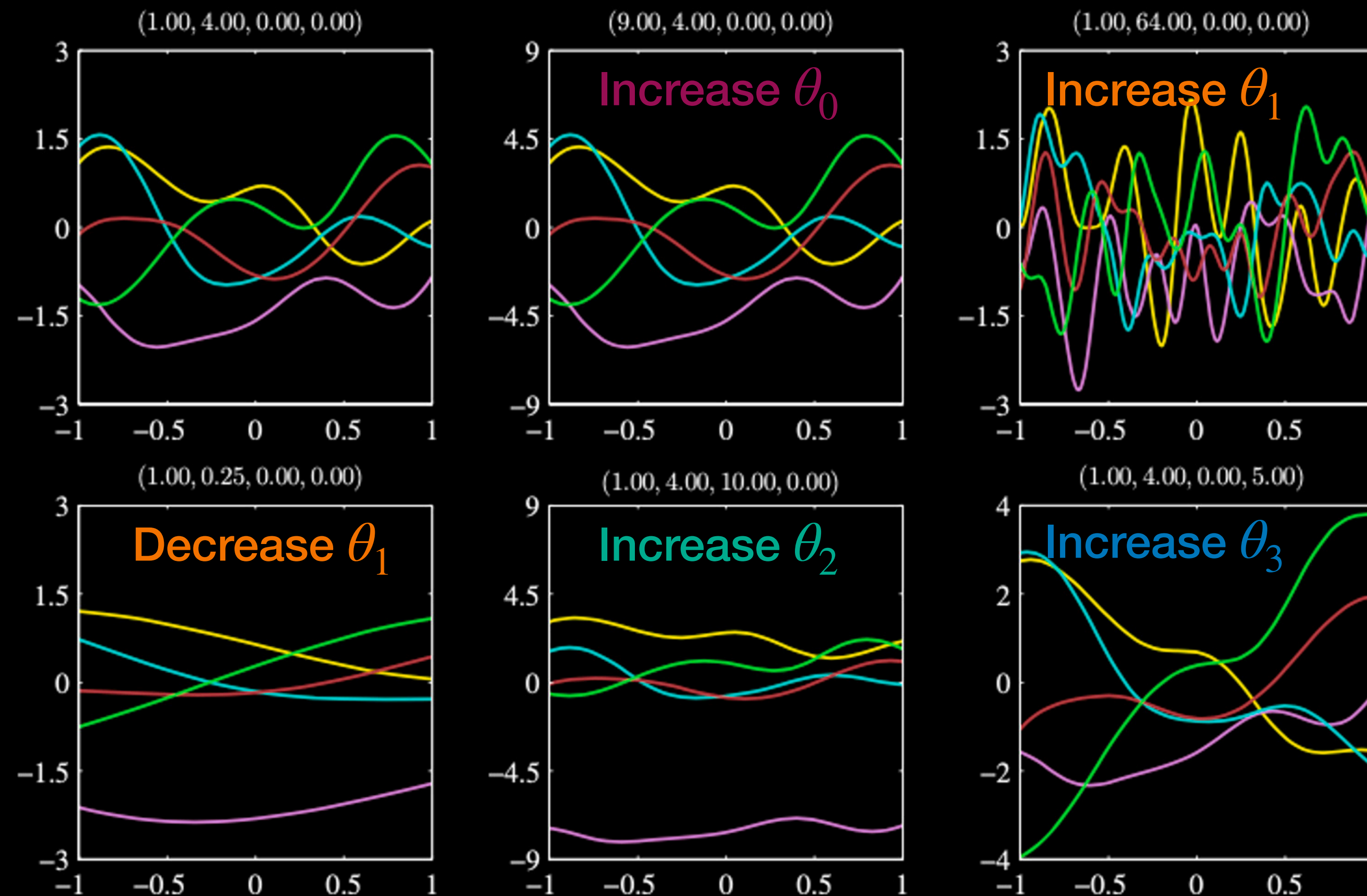


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

More examples: <https://smlbook.org/GP/>

Algorithm

$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y^*; m(x^*), \sigma^2(x^*))$$

$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y}$$

$$\sigma^2(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

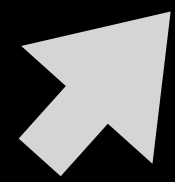
Algorithm

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}(\mathbf{y}^*; \mathbf{m}(\mathbf{x}^*), \sigma^2(\mathbf{x}^*)) \\ m(\mathbf{x}^*) &= k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ \sigma^2(\mathbf{x}^*) &= k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)^{-1} k(\mathbf{X}, \mathbf{x}^*) \end{aligned} \quad \mathcal{O}(N^3)$$

- $L = \text{Cholesky}(k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)$ $k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N = \mathbf{L} \cdot \mathbf{L}^T$
 \mathbf{L} lower-triangular

Algorithm

$N^3/3$ operations



- $L = \text{Cholesky}(k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)$

$$k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N = \mathbf{L} \cdot \mathbf{L}^T$$

\mathbf{L} lower-triangular

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}(\mathbf{y}^*; \mathbf{m}(\mathbf{x}^*), \sigma^2(\mathbf{x}^*)) \\ m(\mathbf{x}^*) &= k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ \sigma^2(\mathbf{x}^*) &= k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)^{-1} k(\mathbf{X}, \mathbf{x}^*) \end{aligned} \quad \mathcal{O}(N^3)$$

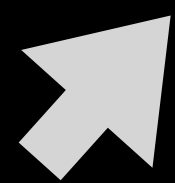
Why Cholesky?

Cost of some matrix factorizations and decompositions. A is $n \times n$, except for QR or VD , where it is $m \times n$ ($m \geq n$).

Factorization/decomposition	Number of flops
LU factorization with partial pivoting ($PA = LU$)	$2n^3/3$
LU factorization with partial pivoting of upper Hessenberg matrix ($PA = LU$)	n^2
Cholesky factorization ($A = R^*R$)	$n^3/3$
Householder QR factorization ($A = QR$)	$2n^2(m - n/3)$ for R ; $4(m^2n - mn^2 + n^3/3)$ for $m \times m$ Q ; $2n^2(m - n/3)$ for $m \times n$ Q ; $2np(2m - n)$ for QB with $m \times p$ B and Q held in factored form.
SVD ^a ($A = P\Sigma Q^*$)	$14mn^2 + 8n^3$ ($P(:, 1:n)$, Σ , and Q) ^b $6mn^2 + 20n^3$ ($P(:, 1:n)$, Σ , and Q) ^c
Hessenberg decomposition ($A = QHQ^*$)	$14n^3/3$ (Q and H), $10n^3/3$ (H only)
Schur decomposition ^a ($A = QTQ^*$)	$25n^3$ (Q and T), $10n^3$ (T only)
For Hermitian A :	
Tridiagonal reduction ($A = QTQ^*$)	$8n^3/3$ (Q and T), $4n^3/3$ (T only)
Spectral decomposition ($A = QDQ^*$)	$9n^3$ (Q and D), $4n^3/3$ (D only)

Algorithm

$N^3/3$ operations



$$p(y^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y^*; m(x^*), \sigma^2(x^*))$$

$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \quad \mathcal{O}(N^3)$$

$$\sigma^2(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

- $L = \text{Cholesky}(k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})$

$$k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} = \mathbf{L} \cdot \mathbf{L}^T$$

\mathbf{L} lower-triangular

- $\alpha = \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y})$

$$\mathbf{A}\mathbf{x} = \mathbf{b} \Rightarrow \mathbf{x} = \mathbf{A} \setminus \mathbf{b}$$

- $m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) \cdot \alpha$

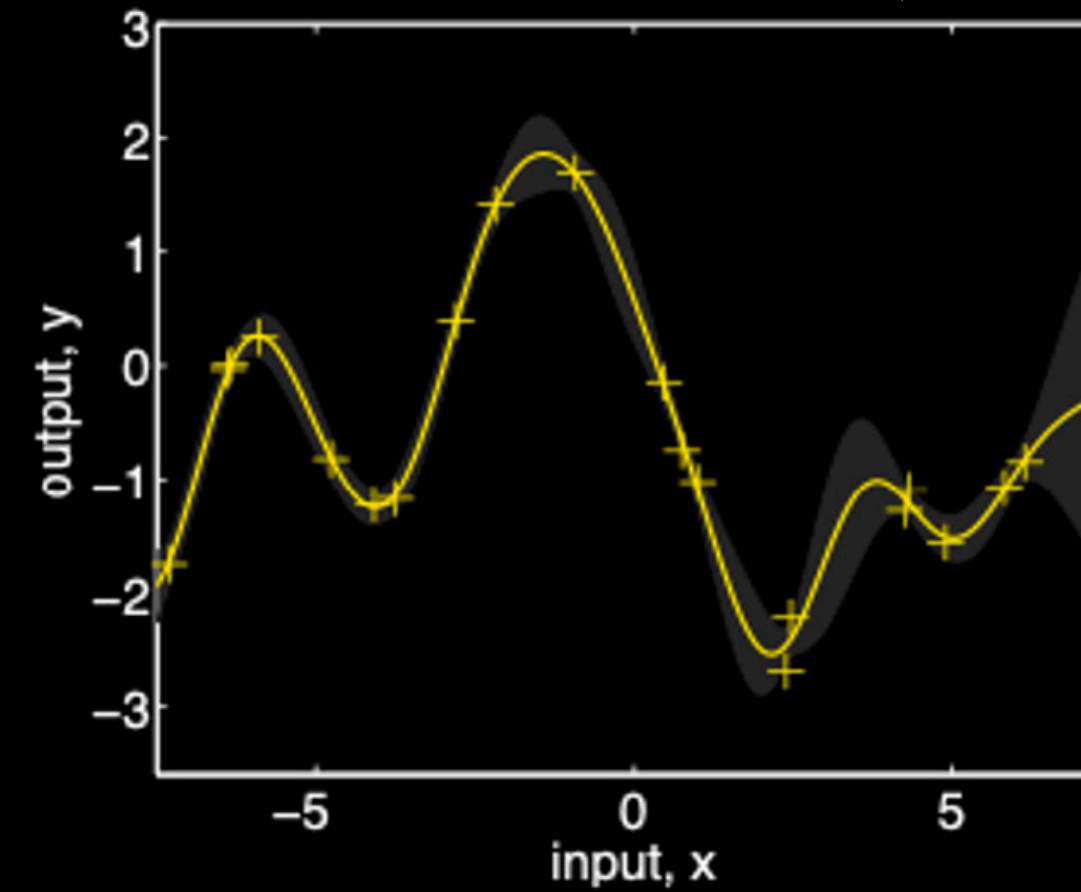
- $\mathbf{v} = \mathbf{L} \setminus k(\mathbf{X}, \mathbf{x}^*)$

- $\sigma^2(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{v}^T \mathbf{v}$

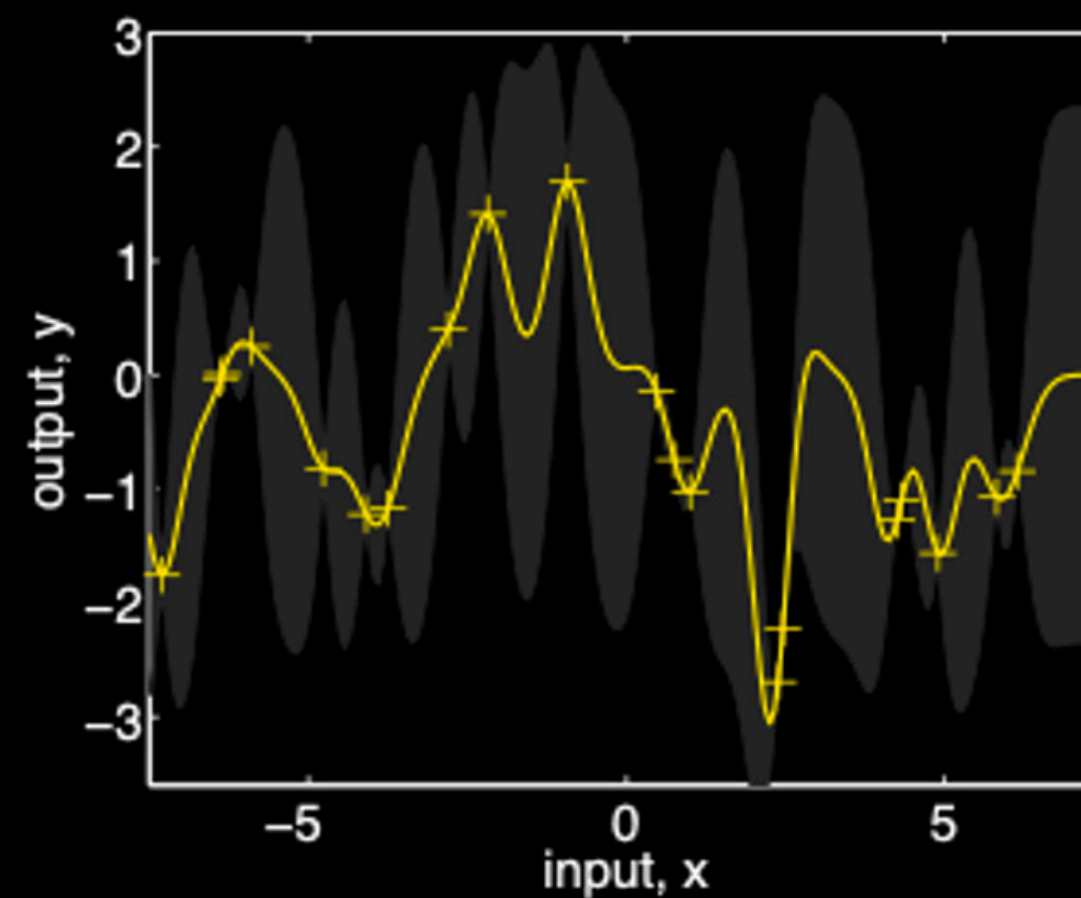
Learning Hyperparameters

$$k(\mathbf{x}_n, \mathbf{x}_m) = \sigma_f^2 \exp \left(-\frac{1}{2\ell^2} ||\mathbf{x}_n - \mathbf{x}_m||^2 \right) + \sigma_n^2 \delta_{nm}$$

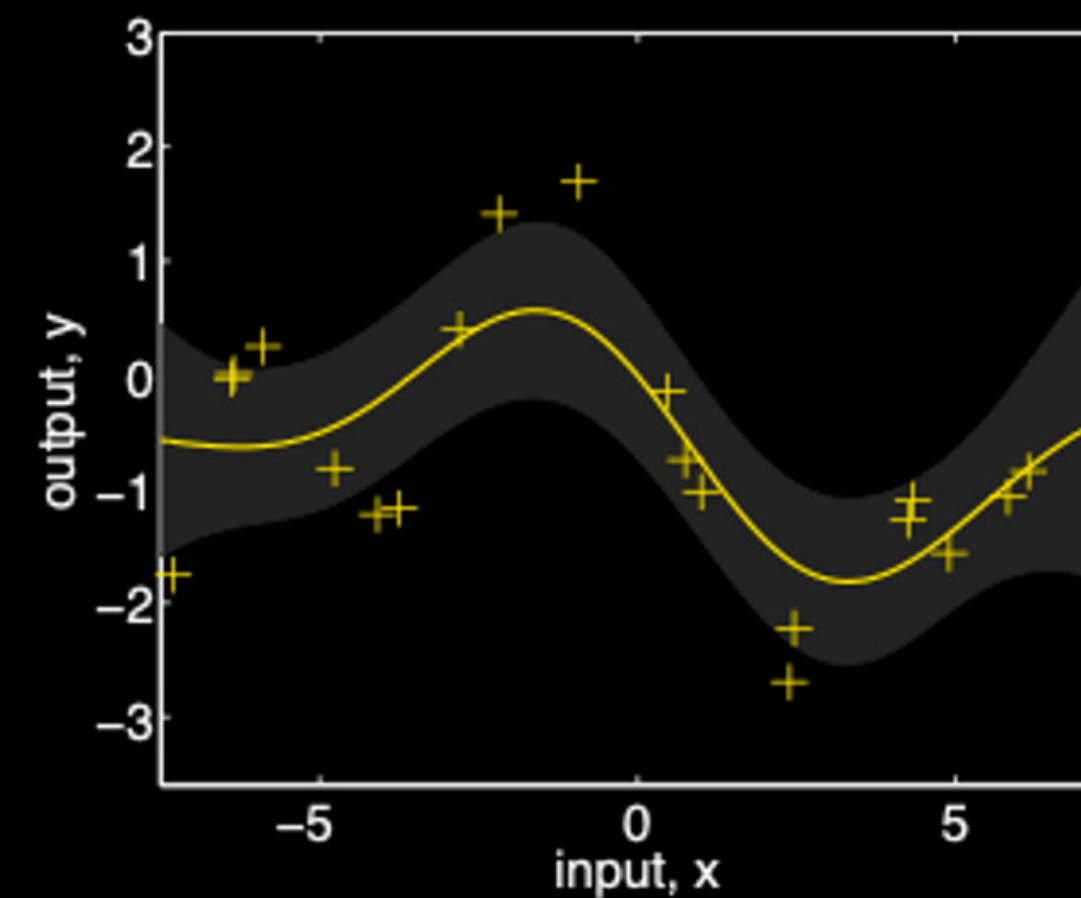
How to choose the
characteristic length /
smoothness ?



(a), $\ell = 1$



(b), $\ell = 0.3$

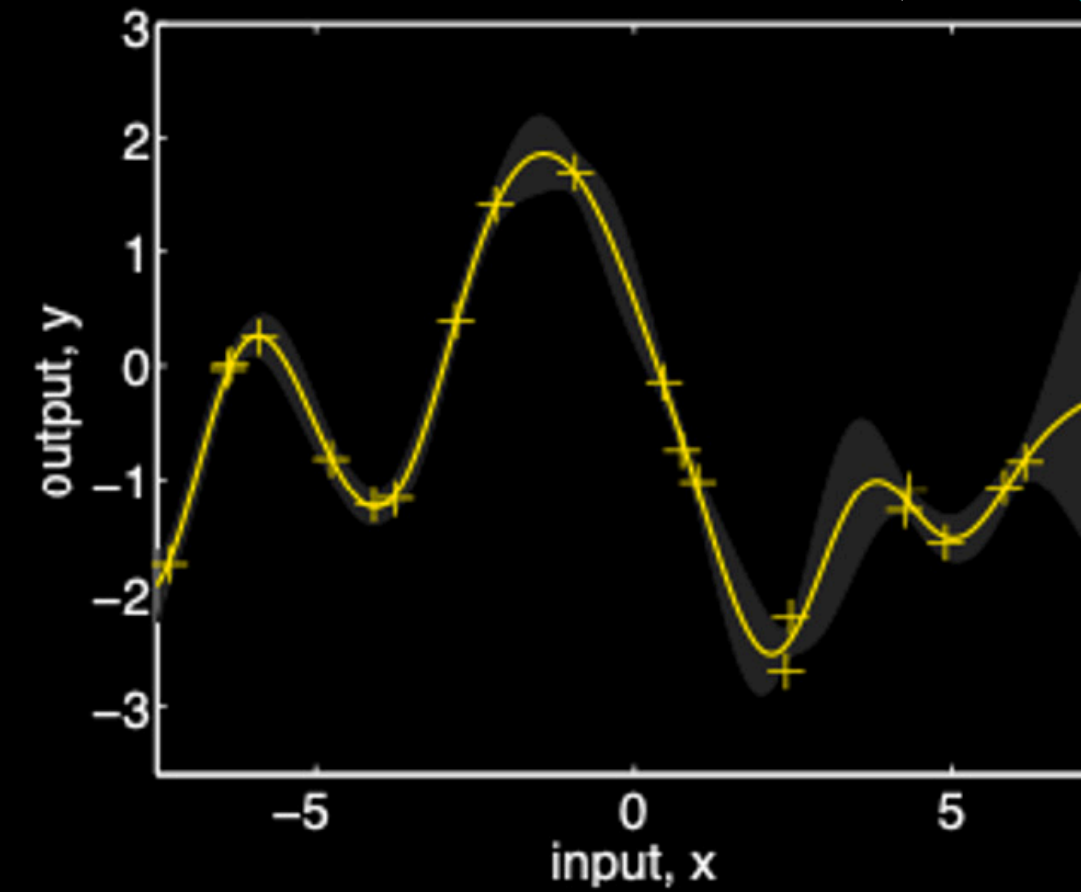


(c), $\ell = 3$

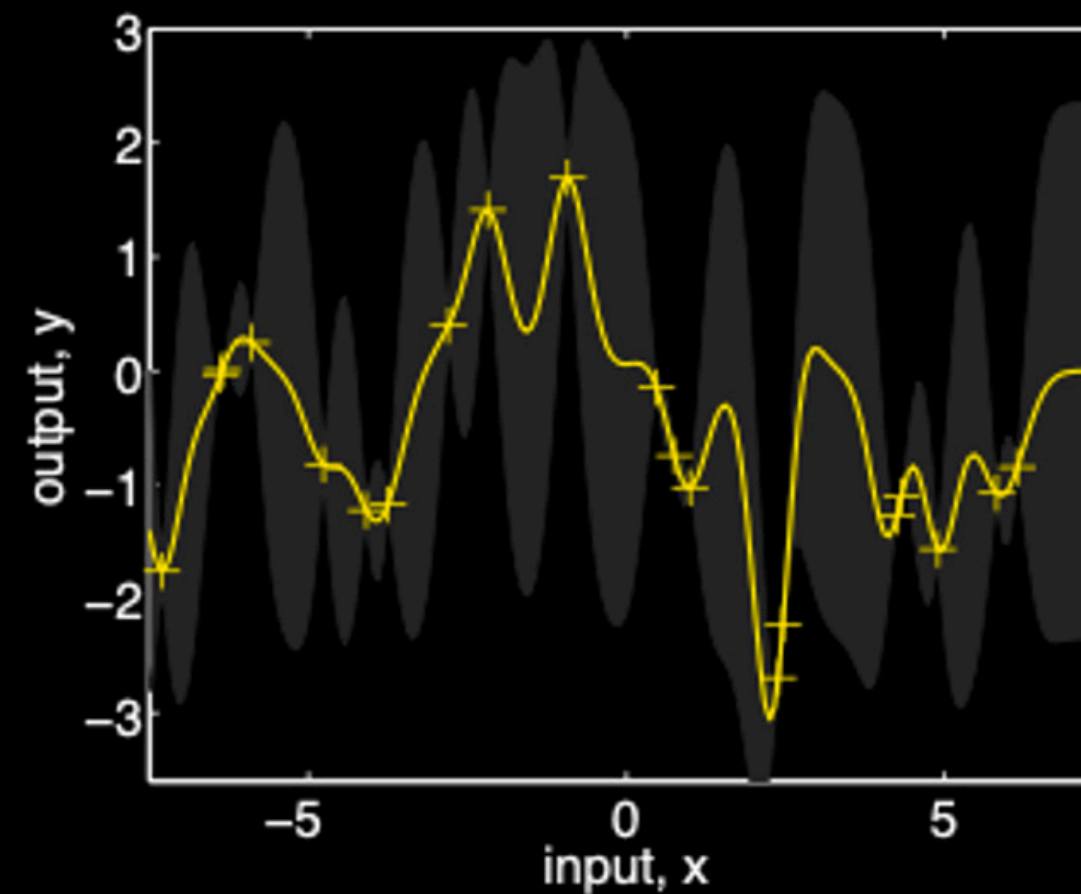
Learning Hyperparameters

$$k(\mathbf{x}_n, \mathbf{x}_m) = \sigma_f^2 \exp \left(-\frac{1}{2\ell^2} ||\mathbf{x}_n - \mathbf{x}_m||^2 \right) + \sigma_n^2 \delta_{nm}$$

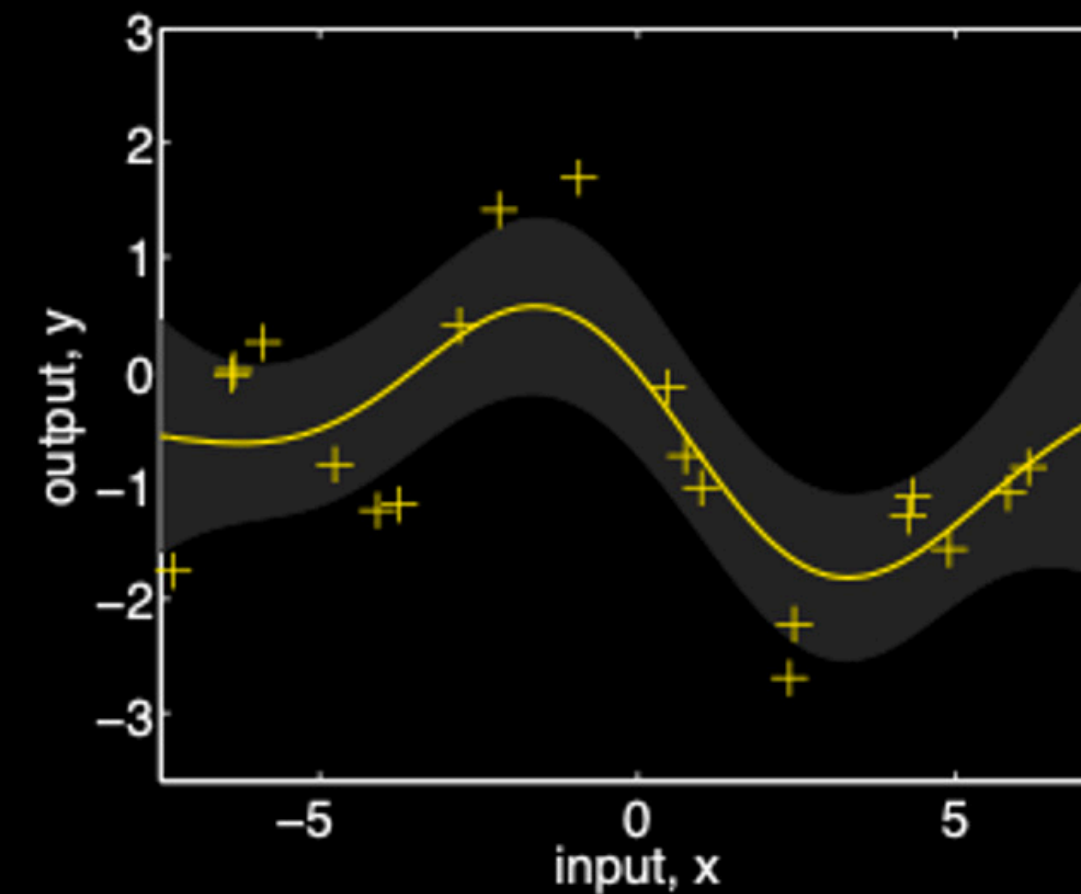
How to choose the
characteristic length /
smoothness ?



(a), $\ell = 1$



(b), $\ell = 0.3$

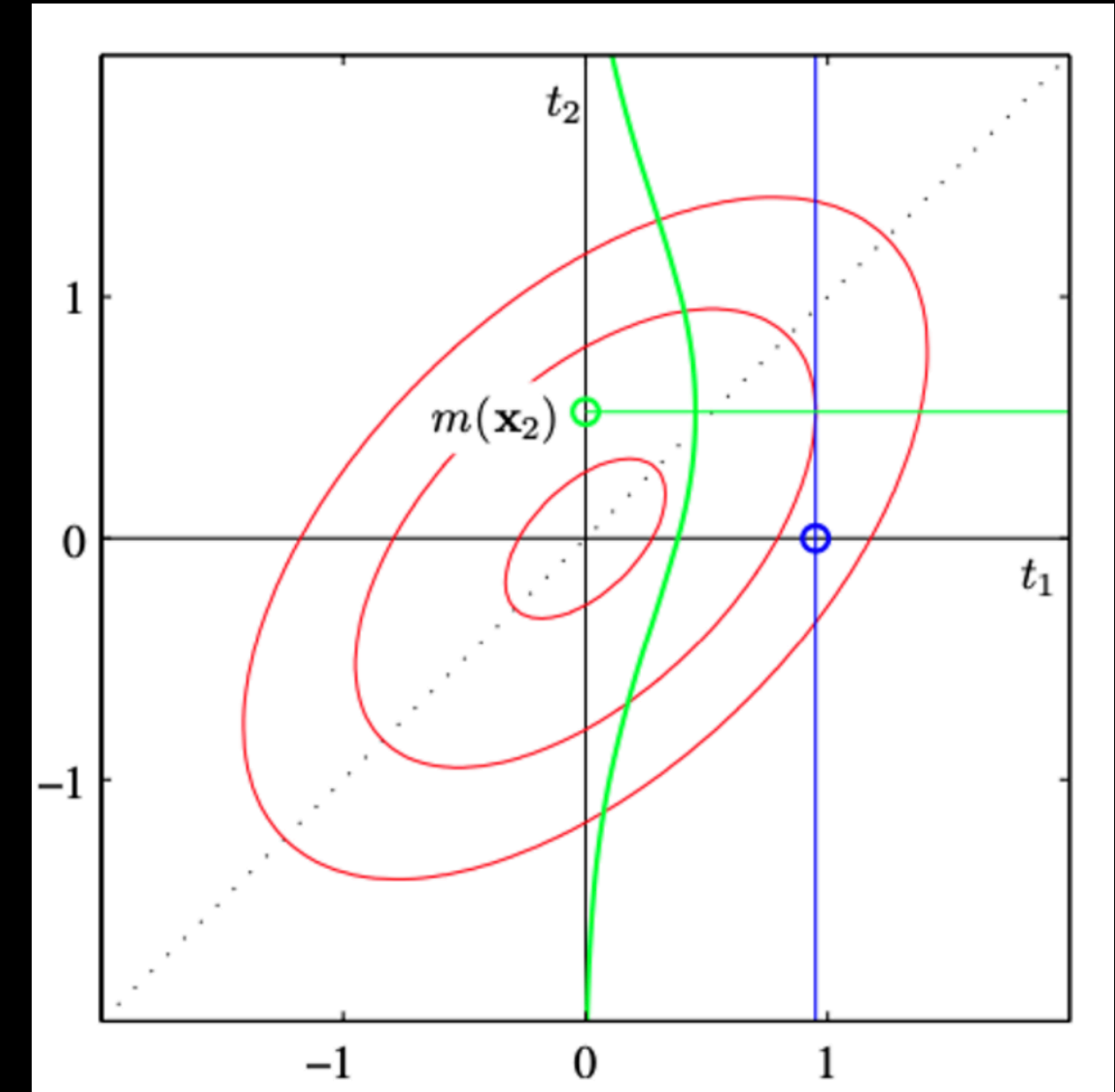


(c), $\ell = 3$

Learning Hyperparameters

Kernel hyperparameter

$$\underbrace{p(\mathbf{y} | \mathbf{X}, \theta)}_{\text{the marginal likelihood}} = \mathcal{N}(\mathbf{y}; \mathbf{0}, \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}}_{\hat{\mathbf{K}}})$$



$$\ln p(\mathbf{y} | \mathbf{X}, \theta) = - \sum_i \ln \mathbf{L}_{ii} - \frac{1}{2} \mathbf{y}^T \boldsymbol{\alpha} - \frac{N}{2} \ln(2\pi)$$

$$\begin{aligned} L &= \text{Cholesky } \hat{\mathbf{K}} \\ \boldsymbol{\alpha} &= \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y}) \end{aligned}$$

Learning Hyperparameters

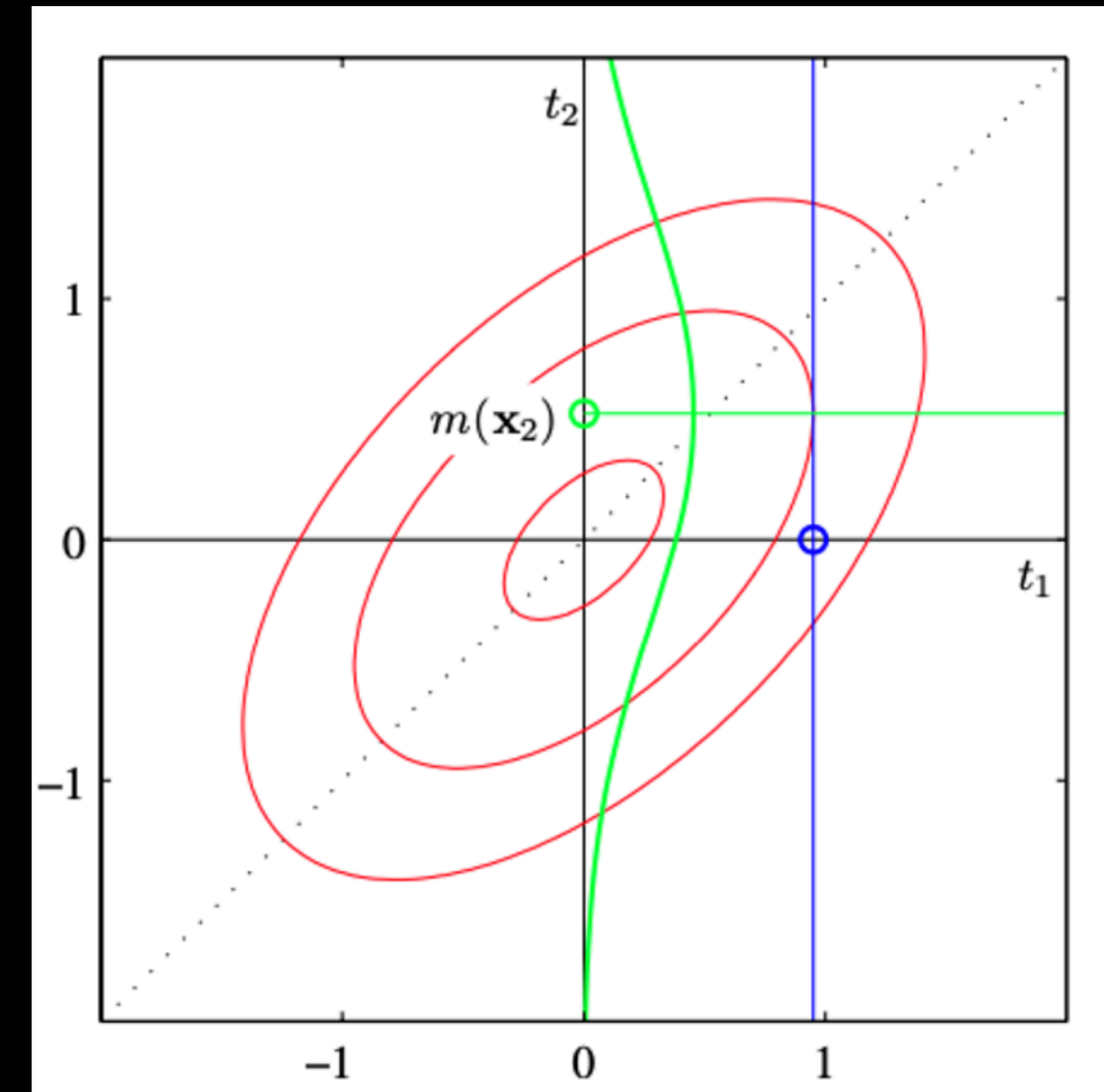
Kernel hyperparameter

$$\underbrace{p(\mathbf{y} | \mathbf{X}, \theta)}_{\text{the marginal likelihood}} = \mathcal{N}(\mathbf{y}; \mathbf{0}, \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}}_{\hat{\mathbf{K}}})$$

the marginal likelihood

$$\ln p(\mathbf{y} | \mathbf{X}, \theta) = -\frac{1}{2} \ln |\hat{\mathbf{K}}| - \frac{1}{2} \mathbf{y}^T \hat{\mathbf{K}}^{-1} \mathbf{y} - \frac{N}{2} \ln(2\pi)$$

(Bishop eq 6.69)



$$\ln p(\mathbf{y} | \mathbf{X}, \theta) = -\sum_i \ln \mathbf{L}_{ii} - \frac{1}{2} \mathbf{y}^T \boldsymbol{\alpha} - \frac{N}{2} \ln(2\pi)$$

$$\begin{aligned} L &= \text{Cholesky } \hat{\mathbf{K}} \\ \boldsymbol{\alpha} &= \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y}) \end{aligned}$$

Learning Hyperparameters

Kernel hyperparameter

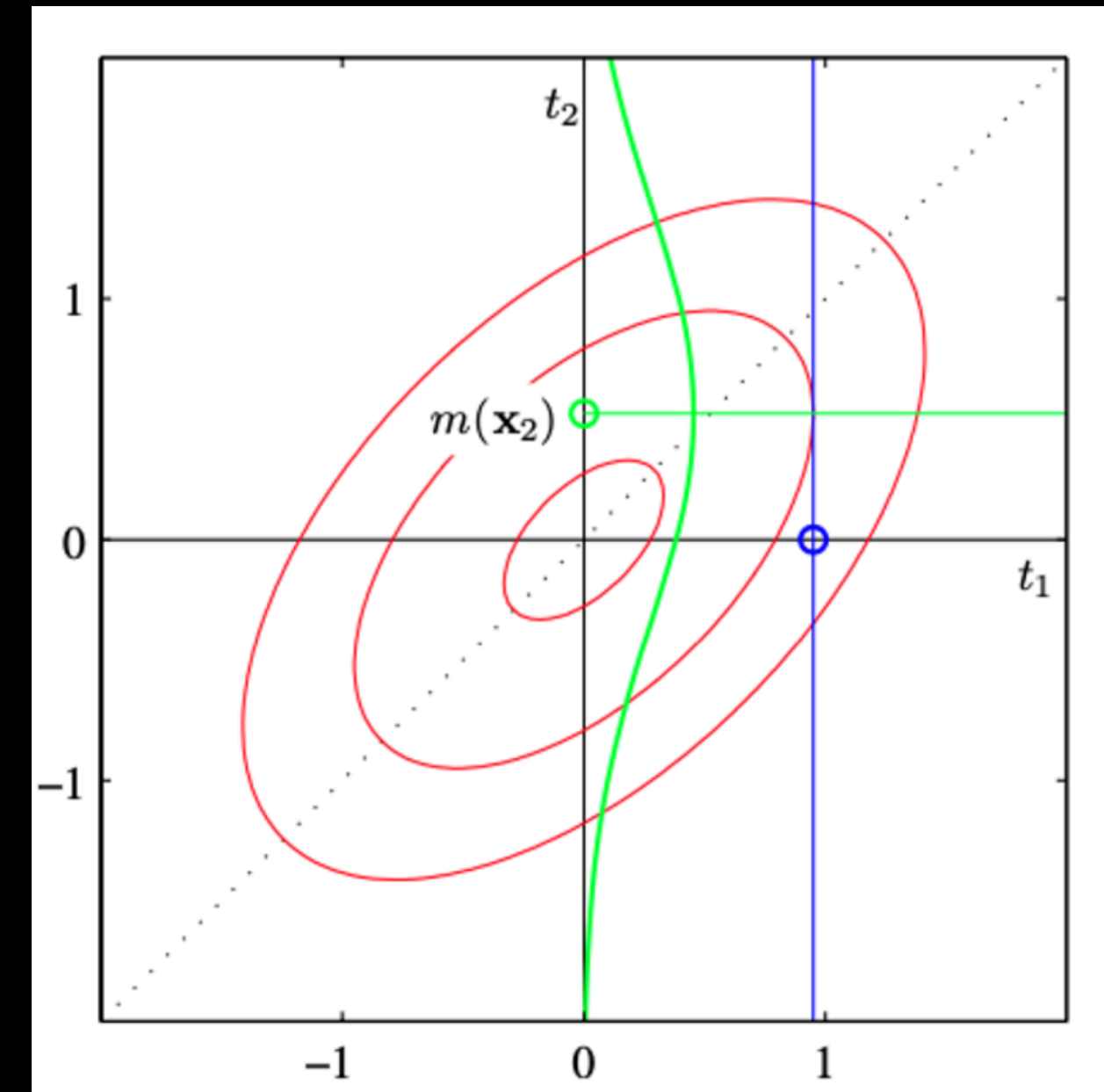
$$\underbrace{p(\mathbf{y} | \mathbf{X}, \theta)}_{\text{the marginal likelihood}} = \mathcal{N}(\mathbf{y}; \mathbf{0}, \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}}_{\hat{\mathbf{K}}})$$

$$\ln p(\mathbf{y} | \mathbf{X}, \theta) = -\frac{1}{2} \ln |\hat{\mathbf{K}}| - \frac{1}{2} \mathbf{y}^T \hat{\mathbf{K}}^{-1} \mathbf{y} - \frac{N}{2} \ln(2\pi)$$

(Bishop eq 6.69)

Learning through gradient descent

$$\ln p(\mathbf{y} | \mathbf{X}, \theta) = -\sum_i \ln L_{ii} - \frac{1}{2} \mathbf{y}^T \boldsymbol{\alpha} - \frac{N}{2} \ln(2\pi)$$



$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$

$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$$

(Bishop eq C.21, C.22)

$$L = \text{Cholesky } \hat{\mathbf{K}}$$

$$\boldsymbol{\alpha} = \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y})$$

Learning Hyperparameters

Kernel hyperparameter

$$\underbrace{p(\mathbf{y} | \mathbf{X}, \theta)}_{\text{the marginal likelihood}} = \mathcal{N}(\mathbf{y}; \mathbf{0}, \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}}_{\hat{\mathbf{K}}})$$

$$\ln p(\mathbf{y} | \mathbf{X}, \theta) = -\frac{1}{2} \ln |\hat{\mathbf{K}}| - \frac{1}{2} \mathbf{y}^T \hat{\mathbf{K}}^{-1} \mathbf{y} - \frac{N}{2} \ln(2\pi)$$

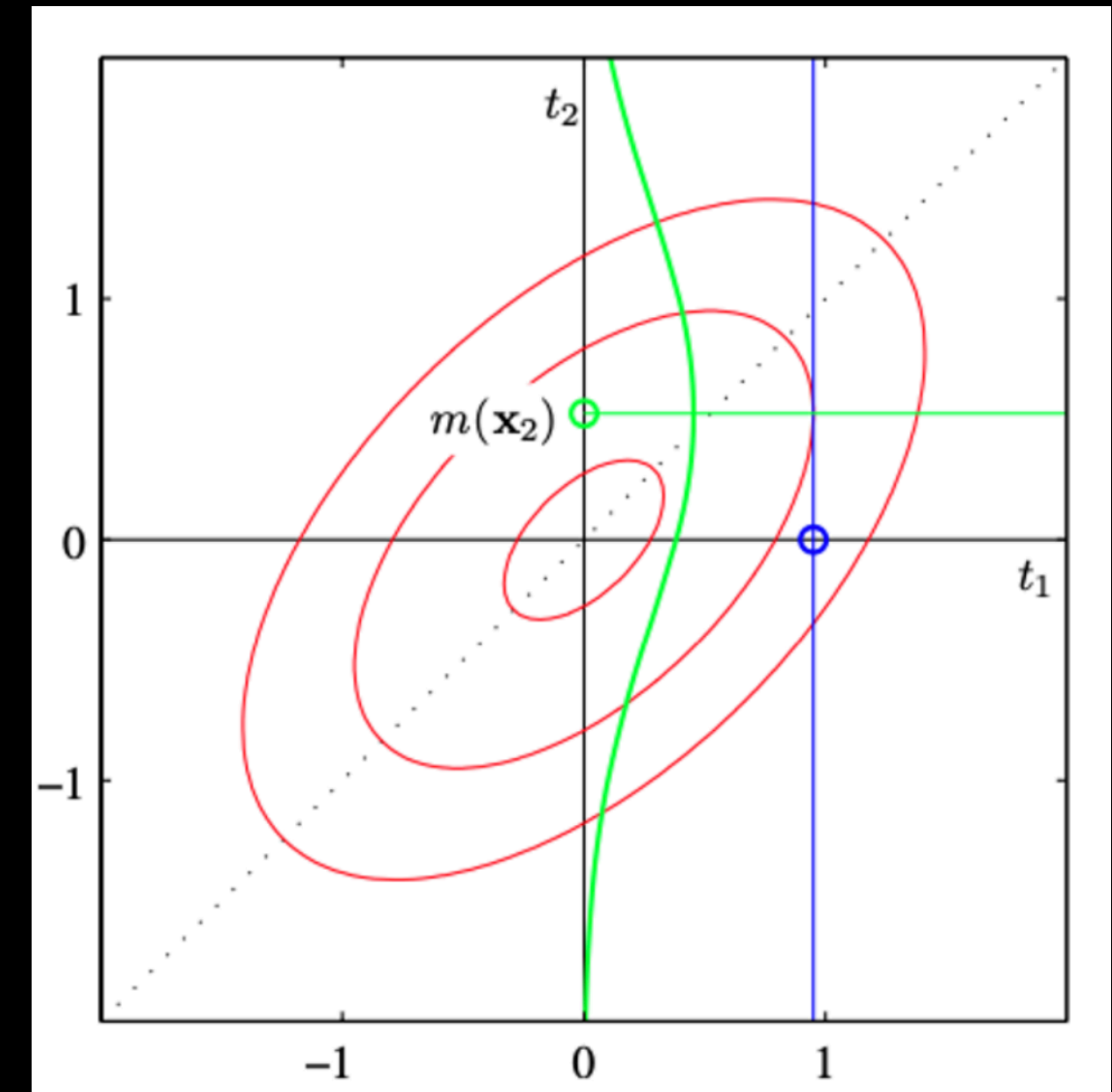
(Bishop eq 6.69)

Learning through gradient descent

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{y} | \mathbf{X}, \theta) = -\frac{1}{2} \text{Tr} \left(\hat{\mathbf{K}}^{-1} \frac{\partial \hat{\mathbf{K}}}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{y}^T \hat{\mathbf{K}}^{-1} \frac{\partial \hat{\mathbf{K}}}{\partial \theta_i} \hat{\mathbf{K}}^{-1} \mathbf{y}$$

(Bishop eq 6.70)

$$\ln p(\mathbf{y} | \mathbf{X}, \theta) = -\sum_i \ln L_{ii} - \frac{1}{2} \mathbf{y}^T \boldsymbol{\alpha} - \frac{N}{2} \ln(2\pi)$$



$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$

$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$$

(Bishop eq C.21, C.22)

$$L = \text{Cholesky } \hat{\mathbf{K}}$$

$$\boldsymbol{\alpha} = \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y})$$