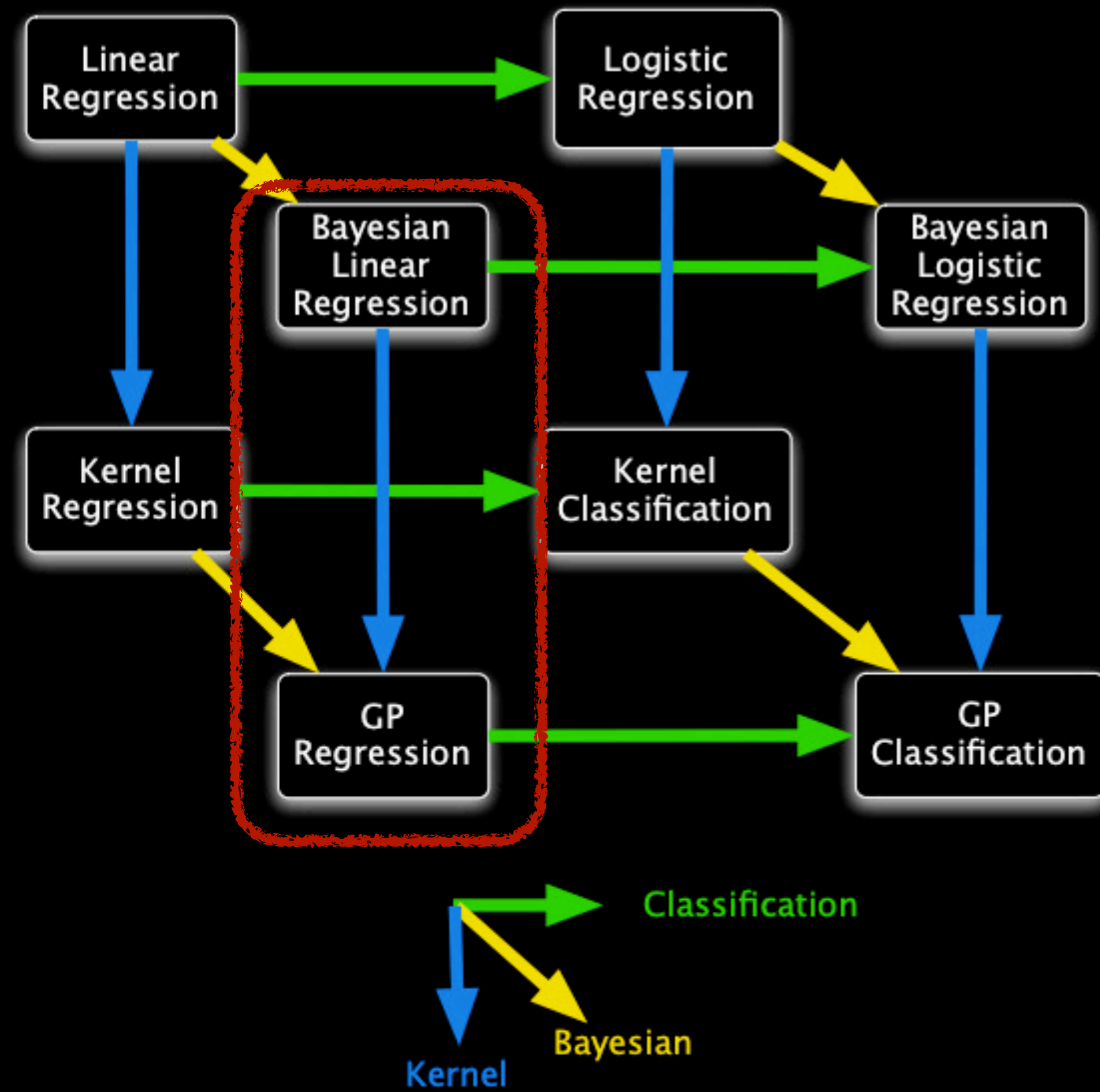


Gaussian process regression

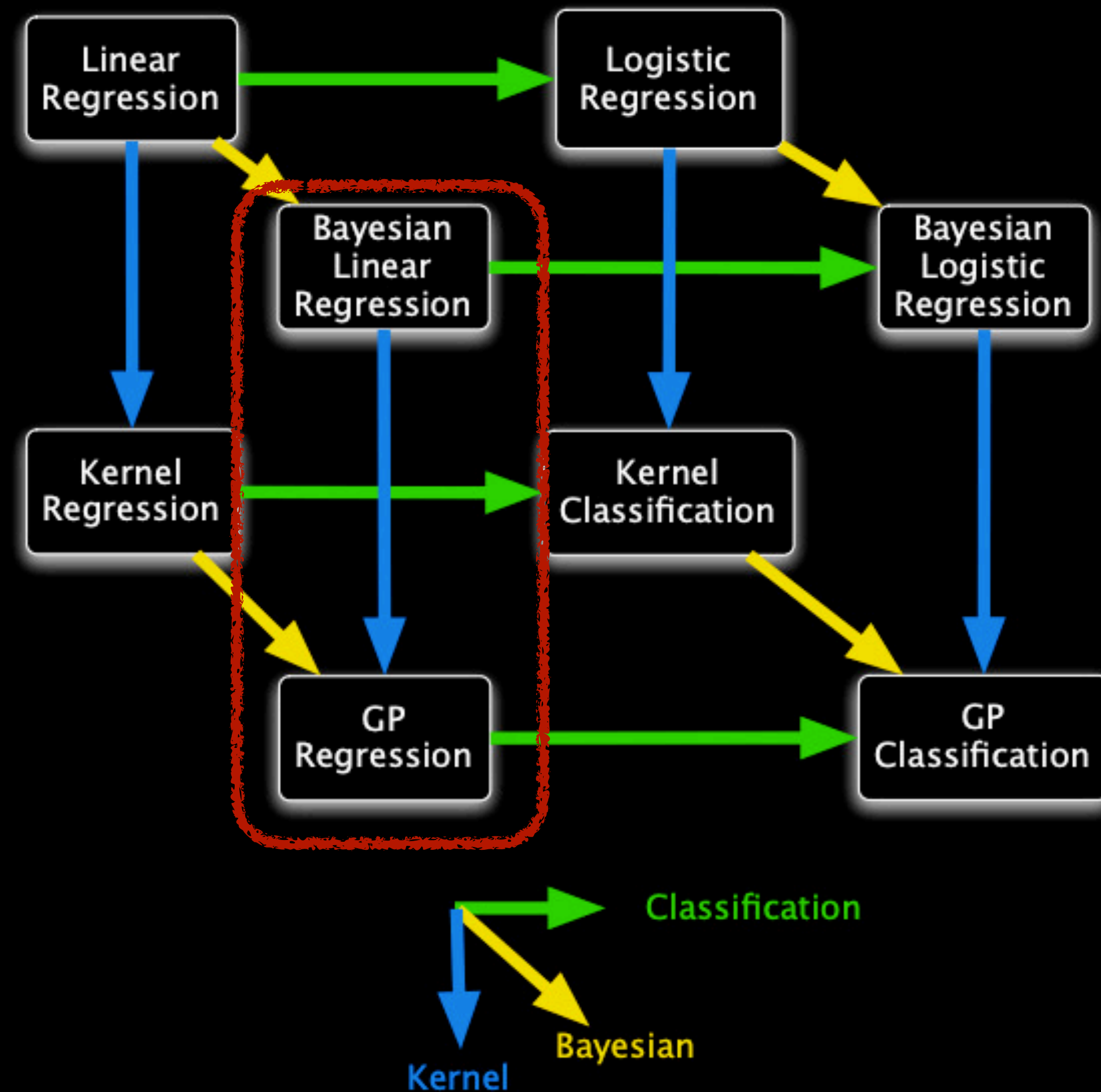
Week 4b - Statistical ML / Thang Bui / ANU / 2025 S1 - Thanks to Yuan-Sen Ting for many slides

Gaussian Process - Weight-Space Perspective



Gaussian Process - Weight-Space Perspective

Gaussian Process =
Kernelising Bayesian
Linear Regression



Recap : Linear Regression

$\Phi \in \mathbb{R}^{N \times D}$: feature matrix
 $\mathbf{y} \in \mathbb{R}^N$: target matrix

D : number of features
 N : number of training data

Recap : Linear Regression

$\Phi \in \mathbb{R}^{N \times D}$: feature matrix
 $\mathbf{y} \in \mathbb{R}^N$: target matrix

D : number of features
 N : number of training data

Week 2a

Features

$$f(\mathbf{x}) = \underbrace{\phi(\mathbf{x})^T}_{\mathbb{R}^D} \underbrace{\theta^*}_{\mathbb{R}^D}$$

Recap : Linear Regression

$\Phi \in \mathbb{R}^{N \times D}$: feature matrix
 $\mathbf{y} \in \mathbb{R}^N$: target matrix

D : number of features
 N : number of training data

Week 2a

Features

$$f(\mathbf{x}) = \underbrace{\phi(\mathbf{x})^T}_{\mathbb{R}^D} \underbrace{\theta^*}_{\mathbb{R}^D}$$

$$\theta^* = \underbrace{(\lambda \mathbf{I}_D)}_{\text{Regularisation}} + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T \mathbf{y}$$

Recap : Linear Regression

$\Phi \in \mathbb{R}^{N \times D}$: feature matrix
 $y \in \mathbb{R}^N$: target matrix

D : number of features
 N : number of training data

Week 2a

Features

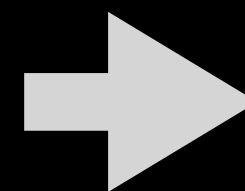
$$f(\mathbf{x}) = \underbrace{\phi(\mathbf{x})^T}_{\mathbb{R}^D} \underbrace{\theta^*}_{\mathbb{R}^D}$$

$$\theta^* = \underbrace{(\lambda \mathbf{I}_D)}_{\text{Regularisation}} + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T y$$

Week 4a

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

Kernelised



Recap : Linear Regression

$\Phi \in \mathbb{R}^{N \times D}$: feature matrix
 $y \in \mathbb{R}^N$: target matrix

D : number of features
 N : number of training data

Week 2a

Features

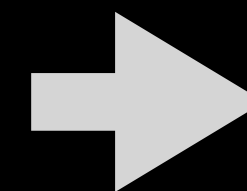
$$f(\mathbf{x}) = \underbrace{\phi(\mathbf{x})^T}_{\mathbb{R}^D} \underbrace{\theta^*}_{\mathbb{R}^D}$$

$$\theta^* = \underbrace{(\lambda \mathbf{I}_D)}_{\text{Regularisation}} + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T y$$

Week 4a

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

Kernelised



$$= k(\mathbf{x})^T (\lambda \mathbf{I}_N + \underbrace{\mathbf{K}}_{\mathbb{R}^{N \times N}})^{-1} y$$

(Bishop eq 6.9)

Gram matrix (positive semi-definite)

$$K_{mn} \equiv k(\mathbf{x}_m, \mathbf{x}_n) = \Phi \Phi^T$$

(Bishop eq 6.10, 6.11)

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

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

Recap : Bayesian Linear Regression

Features

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

$$\theta^* = (\lambda \mathbf{I}_D + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T \mathbf{y}$$

Week 2a

Recap : Bayesian Linear Regression

Features

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

$$\theta^* = (\lambda \mathbf{I}_D + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T \mathbf{y}$$

Week 2a

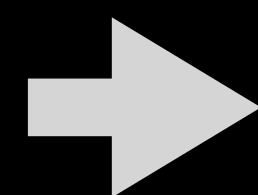
$\mathbb{R}^{D \times D}$

Week 2b

Posterior :

$$p(\theta | \underbrace{X, \mathbf{y}}_{\text{Training data}}) ?$$

Training data



Bayesian

Recap : Bayesian Linear Regression

Features

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

$$\theta^* = (\lambda \mathbf{I}_D + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T \mathbf{y}$$

Week 2a

$\mathbb{R}^{D \times D}$

Week 2b

Posterior :

$$p(\theta | \underbrace{\mathbf{X}, \mathbf{y}}_{\text{Training data}}) ?$$

Training data

Likelihood :

$$p(\mathbf{y} | \mathbf{X}, \theta) = \prod_{n=1}^N \mathcal{N}(y_n; \theta^T \phi(\mathbf{x}_n), \sigma^2)$$

(Bishop eq 3.10, GP book eq 2.3)

➔
Bayesian

Recap : Bayesian Linear Regression

Features

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

$$\theta^* = (\lambda \mathbf{I}_D + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T \mathbf{y}$$

Week 2a

$\mathbb{R}^{D \times D}$

Week 2b

Posterior :

$$p(\theta | \underbrace{\mathbf{X}, \mathbf{y}}_{\text{Training data}}) ?$$

Training data

Likelihood :

$$p(\mathbf{y} | \mathbf{X}, \theta) = \prod_{n=1}^N \mathcal{N}(y_n; \theta^T \phi(\mathbf{x}_n), \sigma^2)$$

(Bishop eq 3.10, GP book eq 2.3)

Bayesian

(Conjugate) Prior :

$$p(\theta) = \mathcal{N}(\theta; 0, \sigma_0^2 \mathbf{I}_D)$$

(Bishop eq 3.52, GP book eq 2.4)

Recap : Bayesian Linear Regression

Features

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

$$\theta^* = (\lambda \mathbf{I}_D + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T \mathbf{y}$$

Week 2a

$\mathbb{R}^{D \times D}$

Week 2b

Posterior :

$$p(\theta | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\theta; \mu, \Sigma)$$

(Bishop eq 3.49)

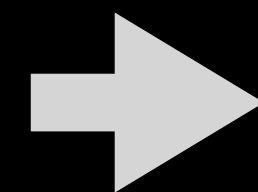
$$\Sigma^{-1} = \sigma_0^{-2} \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi$$

(Bishop eq 3.54,
GP book eq 2.7, 2.8)

$$\mu = \sigma^{-2} \Sigma \Phi^T \mathbf{y}$$

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

(Bishop eq 3.53,
GP book eq 2.7, 2.8)



Bayesian
version

Recap : Bayesian Linear Regression

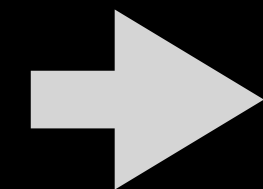
Features

$$f(\mathbf{x}) = \phi^T(\mathbf{x}) \theta^*$$

$$\theta^* = (\lambda \mathbf{I}_D + \underbrace{\Phi^T \Phi}_{\mathbb{R}^{D \times D}})^{-1} \Phi^T \mathbf{y}$$

Week 2a

$\mathbb{R}^{D \times D}$



Bayesian
version

Week 2b

Posterior :

$$p(\theta | X, \mathbf{y}) = \mathcal{N}(\theta; \mu, \Sigma)$$

(Bishop eq 3.49)

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

(Bishop eq 3.54,
GP book eq 2.7, 2.8)

$$\mu = \sigma^{-2} \Sigma \Phi^T \mathbf{y}$$

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

(Bishop eq 3.53,
GP book eq 2.7, 2.8)

Recap : Bayesian Linear Regression

Predictive distribution of y^* , given a new data \mathbf{x}^*

$$p(y^* | \mathbf{x}^*, \mathbf{X}, y) = \int p(y^* | \mathbf{x}^*, \theta) p(\theta | \mathbf{X}, y) d\theta$$

Recap : Bayesian Linear Regression

Predictive distribution of y^* , given a new data \mathbf{x}^*

$$p(y^* | \mathbf{x}^*, \mathbf{X}, y) = \int \overbrace{p(y^* | \mathbf{x}^*, \theta)}^{\text{Gaussian}} \overbrace{p(\theta | \mathbf{X}, y)}^{\text{Gaussian}} d\theta$$

Recap : Bayesian Linear Regression

Predictive distribution of y^* , given a new data \mathbf{x}^*

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{X}, y) &= \int \overbrace{p(y^* | \mathbf{x}^*, \theta)}^{\text{Gaussian}} \overbrace{p(\theta | \mathbf{X}, y)}^{\text{Gaussian}} d\theta \\ &= \mathcal{N}(y^*; \mu^T \phi(\mathbf{x}^*), \sigma^2(\mathbf{x}^*)) \end{aligned}$$

(Bishop eq 3.58, GP book eq 2.9)

$$\sigma^2(\mathbf{x}^*) = \underbrace{\sigma^2}_{\text{Observation noise}} + \phi(\mathbf{x}^*)^T \Sigma \phi(\mathbf{x}^*)$$

(Bishop eq 3.59,
GP book eq 2.9)


Observation noise

Recap : Bayesian Linear Regression

Predictive distribution of y^* , given a new data \mathbf{x}^*

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{X}, y) &= \int \overbrace{p(y^* | \mathbf{x}^*, \theta)}^{\text{Gaussian}} \overbrace{p(\theta | \mathbf{X}, y)}^{\text{Gaussian}} d\theta \\ &= \mathcal{N}(y^*; \mu^T \phi(\mathbf{x}^*), \sigma^2(\mathbf{x}^*)) \end{aligned}$$

(Bishop eq 3.58, GP book eq 2.9)

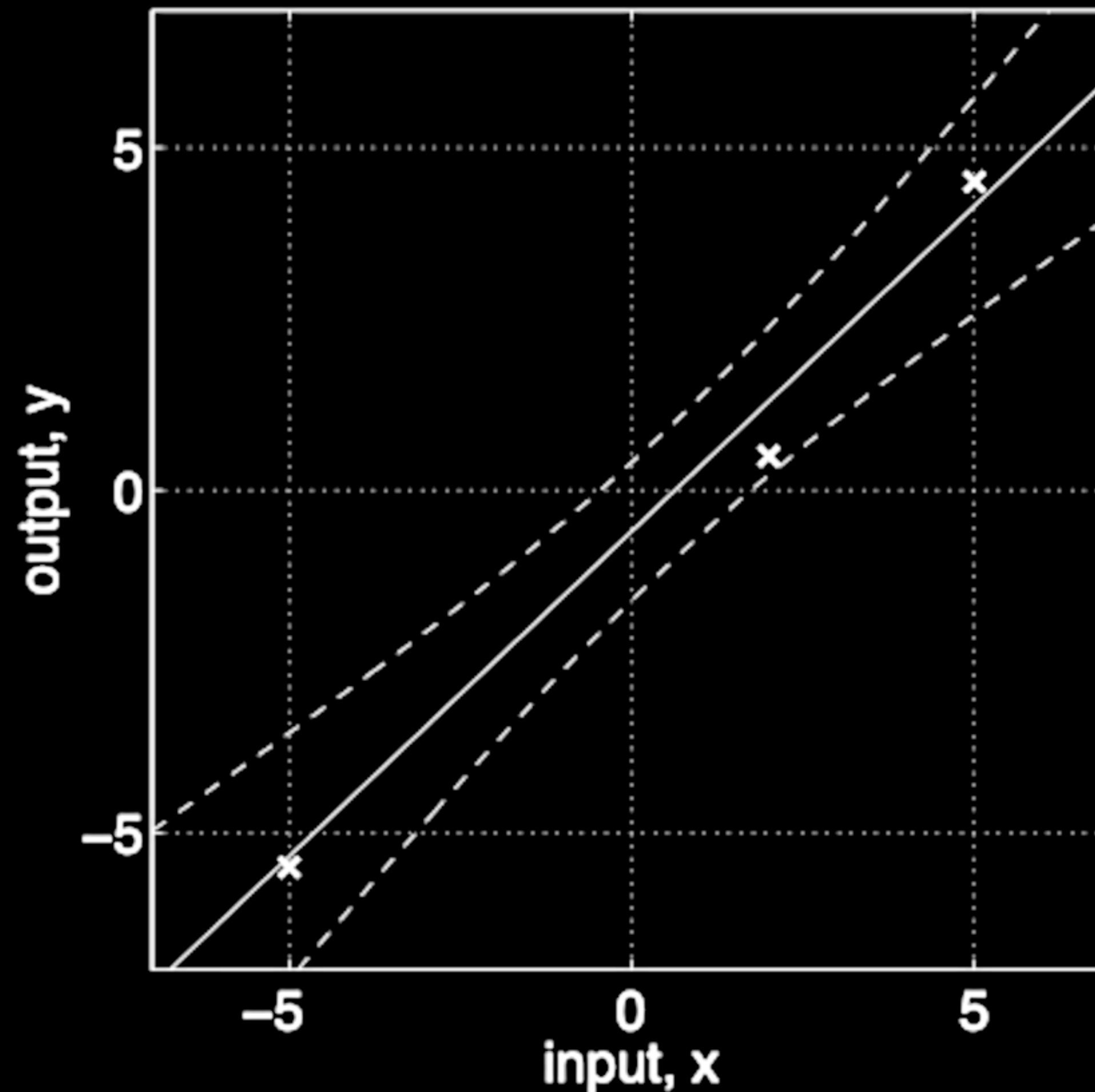
$$\sigma^2(\mathbf{x}^*) = \underbrace{\sigma^2}_{\text{Observation noise}} + \underbrace{\phi(\mathbf{x}^*)^T \Sigma \phi(\mathbf{x}^*)}_{\text{Uncertainty due to } \theta}$$

(Bishop eq 3.59,
GP book eq 2.9)

Why Full Bayesian Treatment Can Be Advantageous

Predictive distribution of y^* , given a new data \mathbf{x}^*

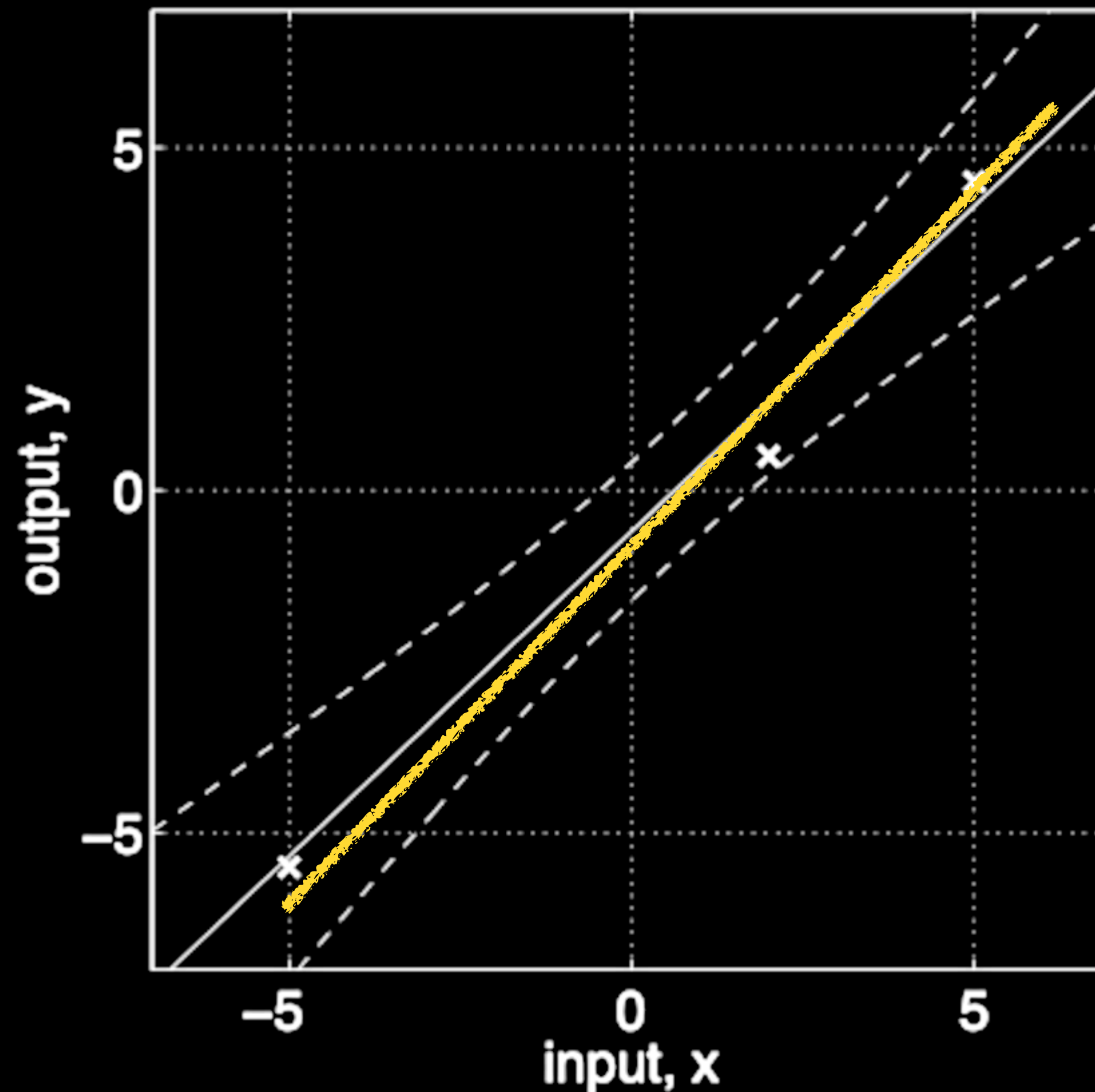
$$p(y^* | x^*, \mathbf{X}, \mathbf{y})$$



Why Full Bayesian Treatment Can Be Advantageous

Predictive distribution of y^* , given a new data \mathbf{x}^*

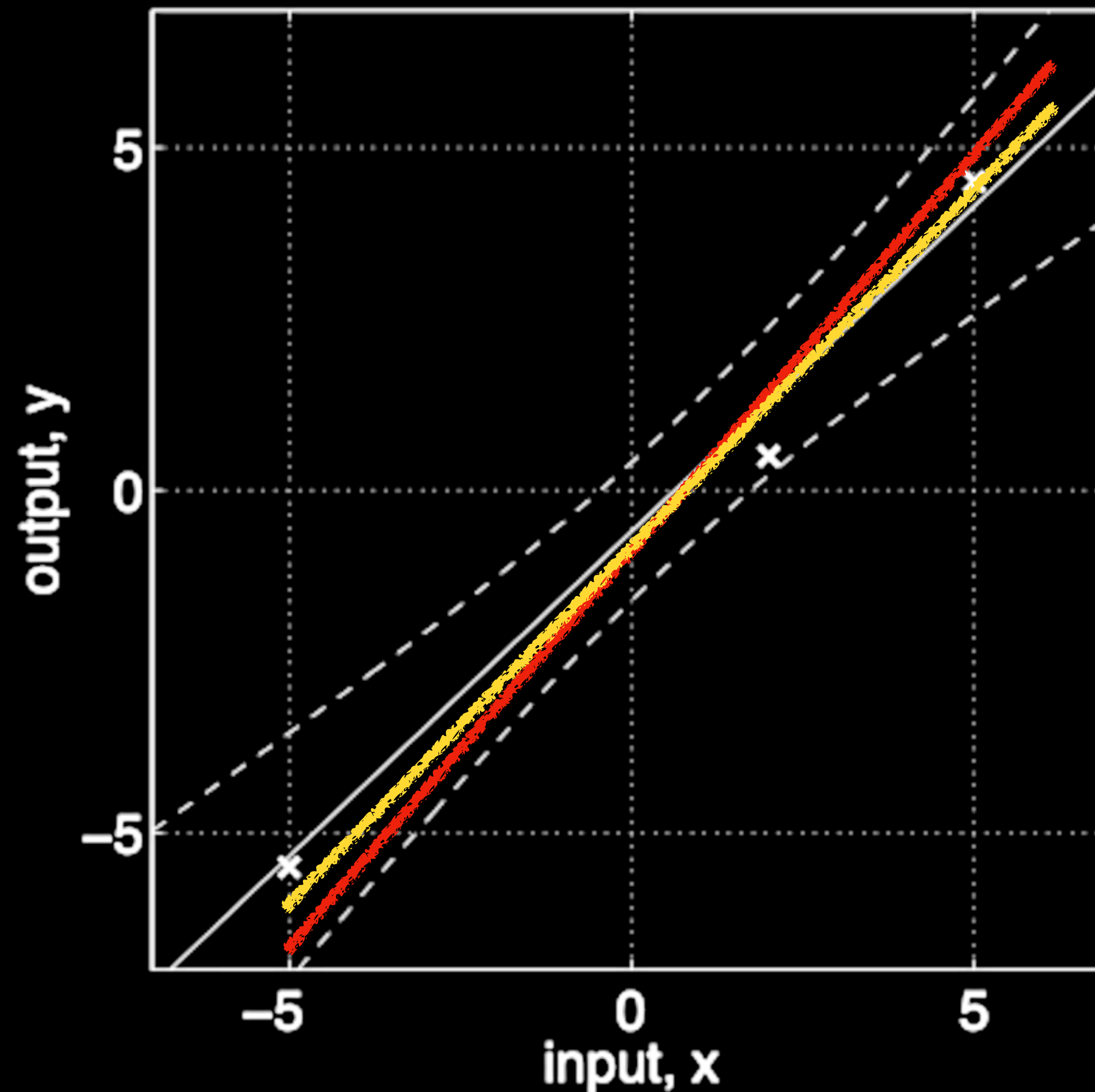
$$p(y^* | x^*, \mathbf{X}, \mathbf{y})$$



Why Full Bayesian Treatment Can Be Advantageous

Predictive distribution of y^* , given a new data \mathbf{x}^*

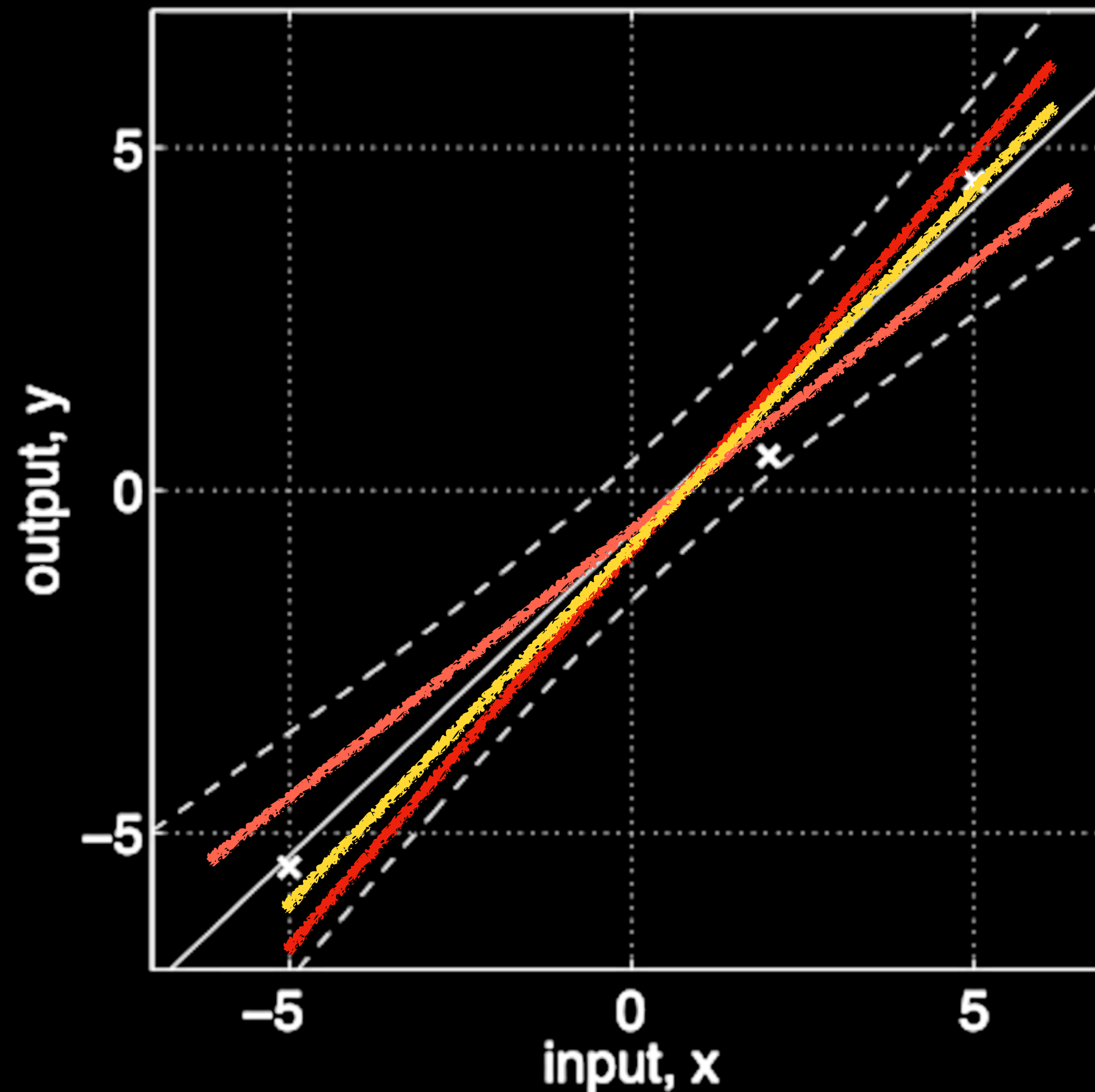
$$p(y^* | x^*, \mathbf{X}, \mathbf{y})$$



Why Full Bayesian Treatment Can Be Advantageous

Predictive distribution of y^* , given a new data \mathbf{x}^*

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

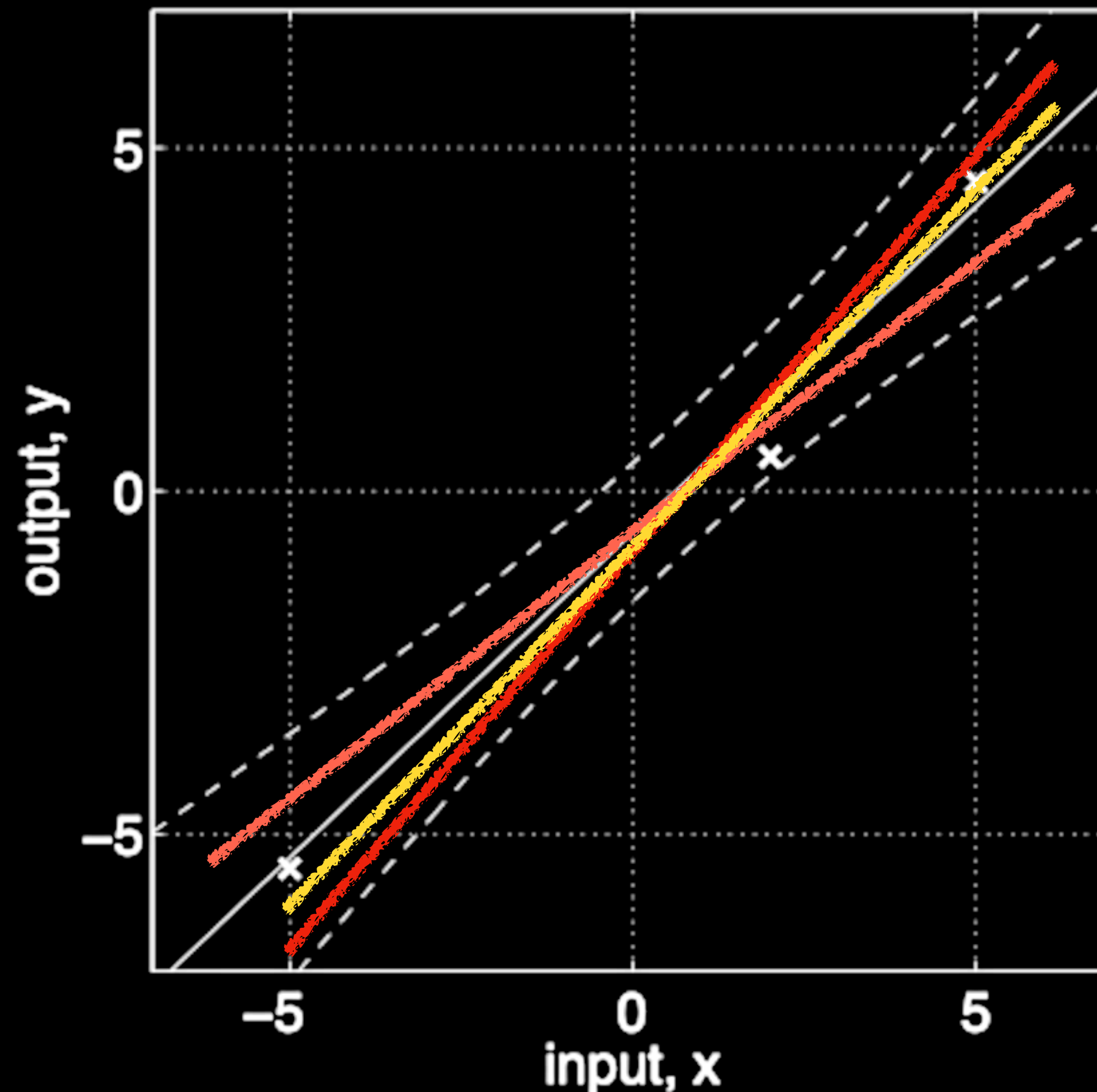


Why Full Bayesian Treatment Can Be Advantageous

Predictive distribution of y^* , given a new data x^*

Integrating all possible θ

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



Yesterday:

Primal

parameters θ

$$\text{optimal } \theta = (\Phi^T \Phi + \lambda I_D)^{-1} \Phi^T \mathbf{y}$$

$$\text{prediction } f(x_*) = \phi_*^T \theta$$

$$\text{complexity } \mathcal{O}(D^3 + D^2 N)$$

Dual

$$\text{weights } \alpha \text{ with } \theta = \sum_n \alpha_n \phi_n$$

$$\text{optimal } \alpha = (\Phi \Phi^T + \lambda I_N)^{-1} \mathbf{y}$$

$$\text{prediction } f(x_*) = \sum_n \alpha_n \phi_n^T \phi_*$$

$$\text{complexity } \mathcal{O}(N^3 + N^2 D)$$

Primal

Dual

Yesterday:

parameters θ

weights α with $\theta = \sum_n \alpha_n \phi_n$

optimal $\theta = (\Phi^\top \Phi + \lambda I_D)^{-1} \Phi^\top \mathbf{y}$

optimal $\alpha = (\Phi \Phi^\top + \lambda I_N)^{-1} \mathbf{y}$

prediction $f(x_*) = \phi_*^\top \theta$

prediction $f(x_*) = \sum_n \alpha_n \phi_n^\top \phi_*$

complexity $\mathcal{O}(D^3 + D^2 N)$

complexity $\mathcal{O}(N^3 + N^2 D)$

Can we write down

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

In the form of

$$k(\mathbf{x}^*, \mathbf{X}), k(\mathbf{X}, \mathbf{X})$$

Gaussian Process - Weight-Space Perspective

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

Gaussian Process - Weight-Space Perspective

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

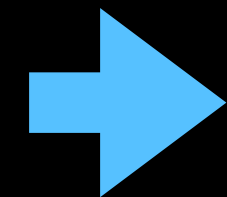
$$m(\mathbf{x}^*) = \phi^T(\mathbf{x}^*) \sigma^{-2} (\sigma_0^2 \mathbf{I}_D + \sigma^{-2} \Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

Gaussian Process - Weight-Space Perspective

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

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

Kernel trick
(see Lecture 4a)



$$m(\mathbf{x}^*) = \sigma_0^2 \phi^T(\mathbf{x}^*) \Phi^T (\sigma^2 \mathbf{I}_N + \sigma_0^2 \Phi \Phi^T)^{-1} \mathbf{y}$$

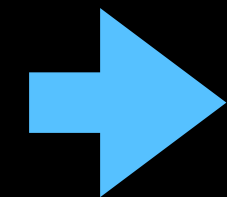
(GP book eq 2.12)

Gaussian Process - Weight-Space Perspective

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

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

Kernel trick
(see Lecture 4a)



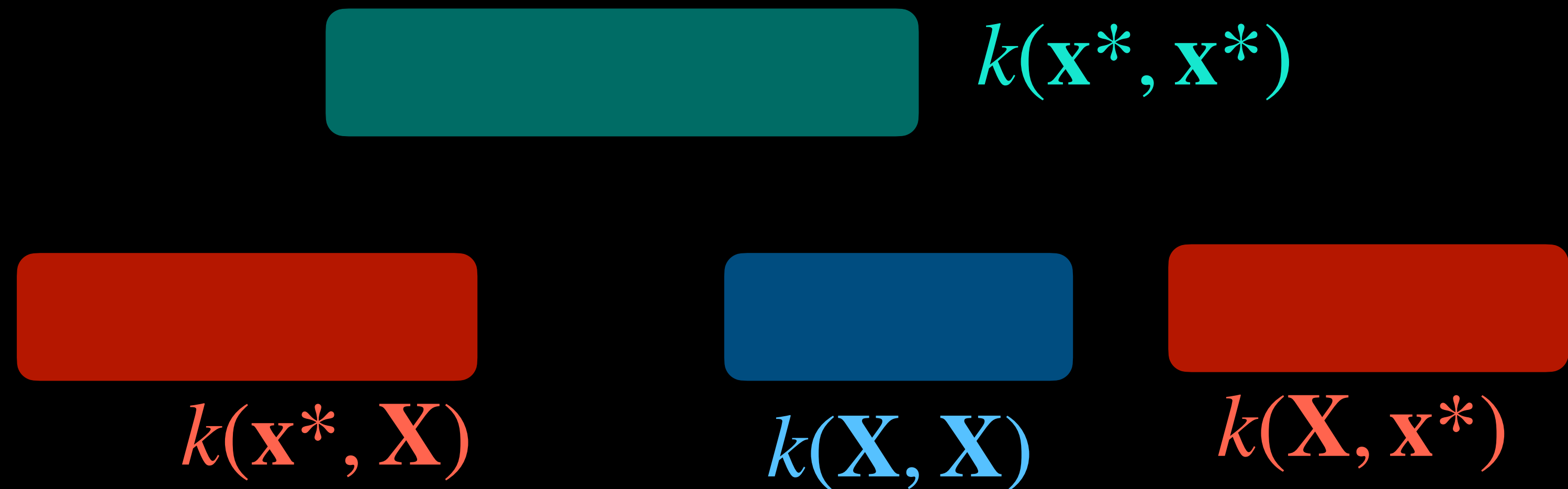
$$m(\mathbf{x}^*) = \underbrace{\sigma_0^2 \phi^T(\mathbf{x}^*) \Phi^T}_{k(\mathbf{x}^*, \mathbf{X})} (\sigma^2 \mathbf{I}_N + \underbrace{\sigma_0^2 \Phi \Phi^T}_{k(\mathbf{X}, \mathbf{X})})^{-1} y$$

(GP book eq 2.12)

Gaussian Process - Weight-Space Perspective

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

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

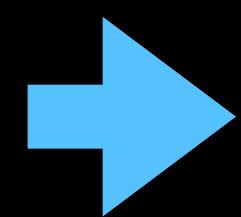


Gaussian Process - Weight-Space Perspective

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

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

Kernel trick
(see Lecture 4a)

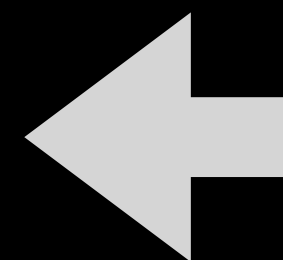


$$= \sigma^2 + \sigma_0^2 \phi(\mathbf{x}^*)^T \phi(\mathbf{x}^*) \quad k(\mathbf{x}^*, \mathbf{x}^*)$$

(GP book eq 2.12)

$$- \underbrace{\sigma_0^2 \phi(\mathbf{x}^*)^T \Phi^T}_{k(\mathbf{x}^*, \mathbf{X})} (\sigma^2 \mathbf{I}_N + \underbrace{\sigma_0^2 \Phi \Phi^T}_{k(\mathbf{X}, \mathbf{X})})^{-1} \underbrace{\sigma_0^2 \Phi \phi(\mathbf{x}^*)}_{k(\mathbf{X}, \mathbf{x}^*)}$$

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

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

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

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

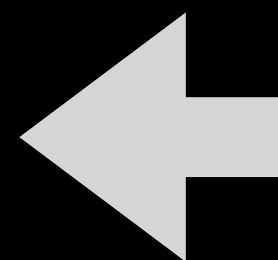
(Bishop eq 3.53)

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

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

(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

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

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

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

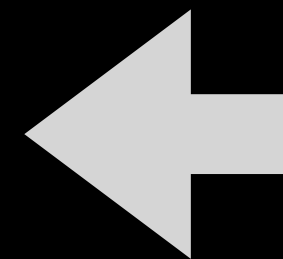
(Bishop eq 3.53)

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

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

(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)

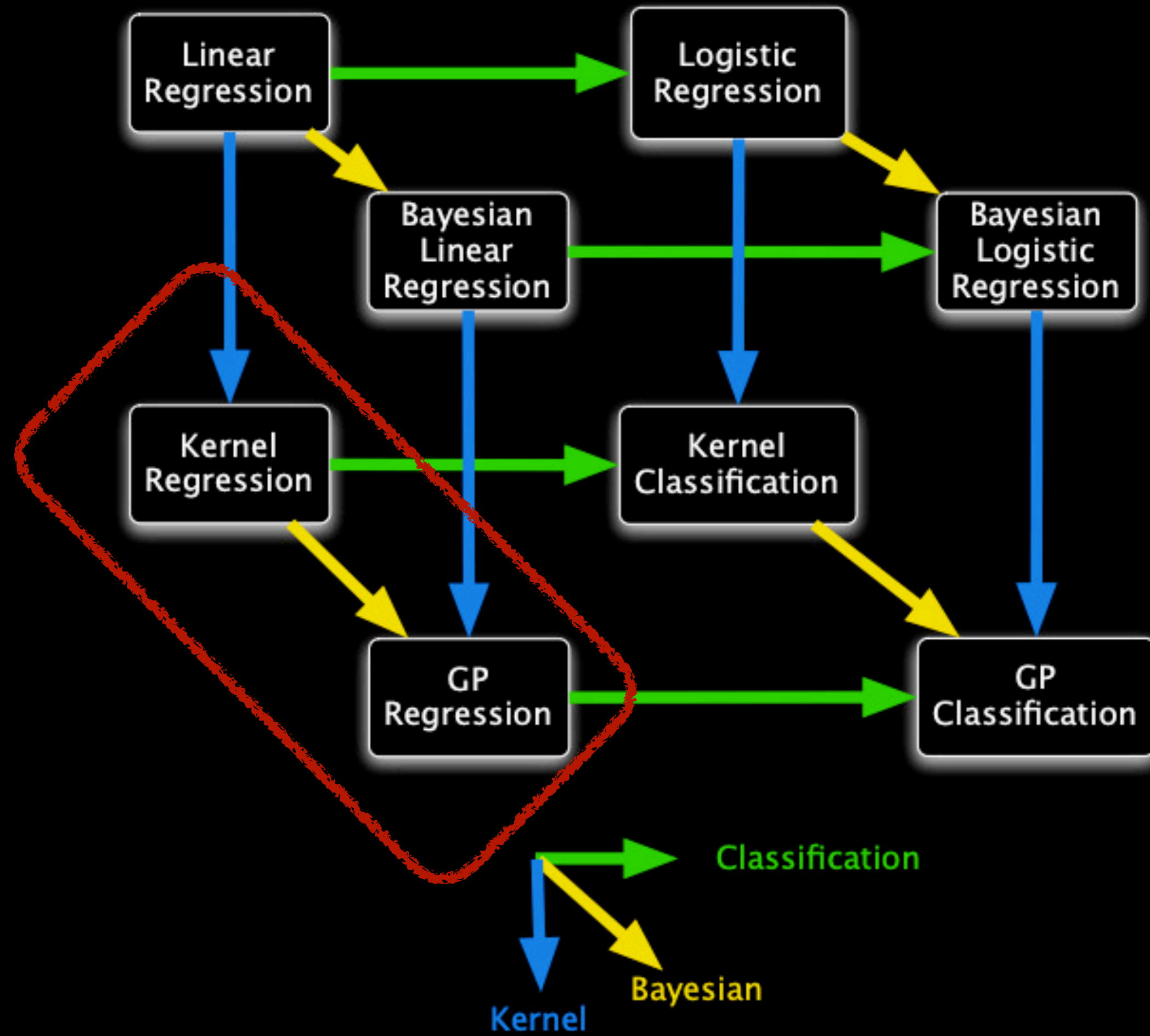
$$\sigma^2(\mathbf{x}^*) = \sigma^2 + \phi(\mathbf{x}^*)^T \Sigma \phi(\mathbf{x}^*)$$

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

(Bishop eq 3.54)

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

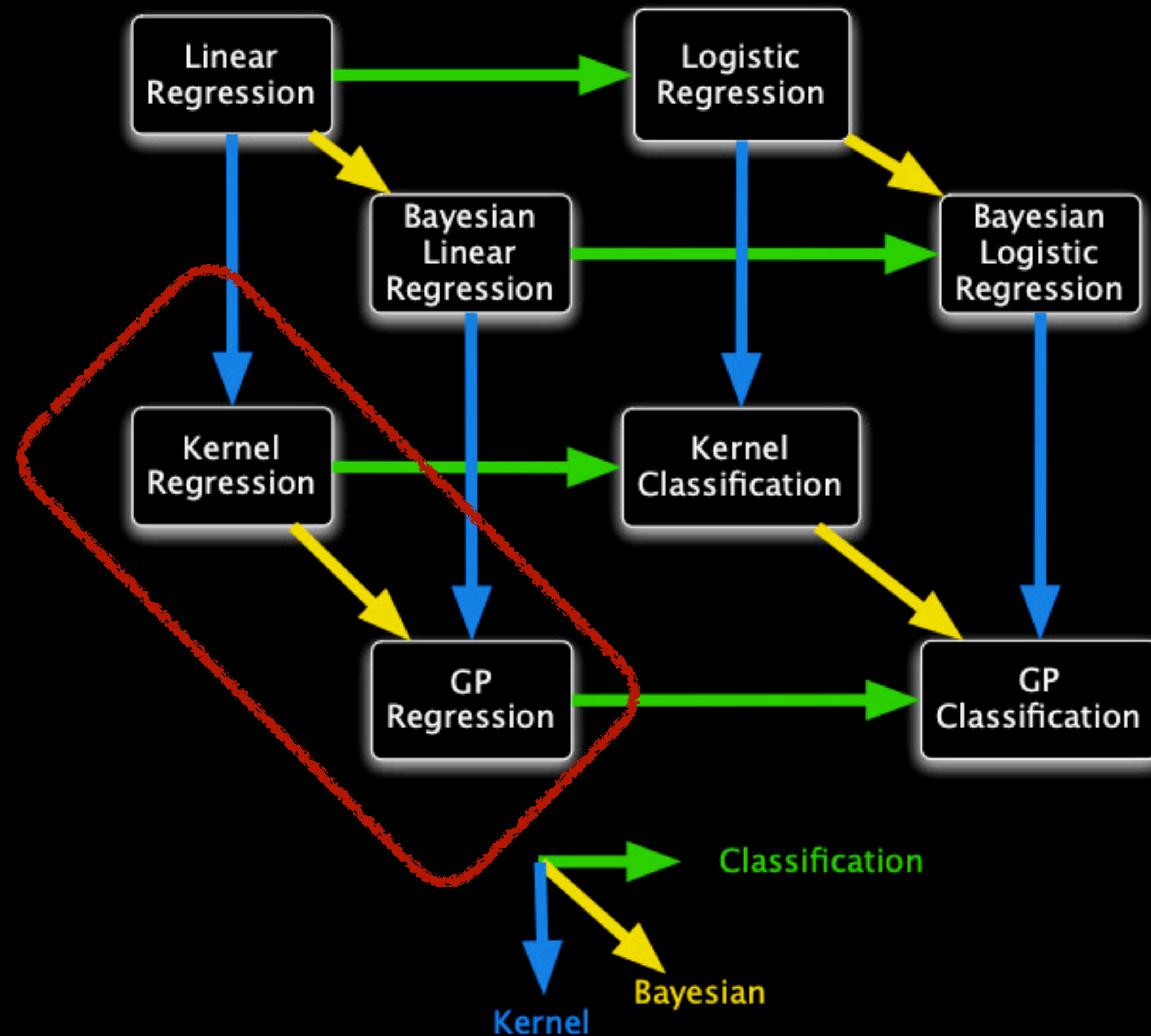
Gaussian Process - Function-Space Perspective



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**.

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

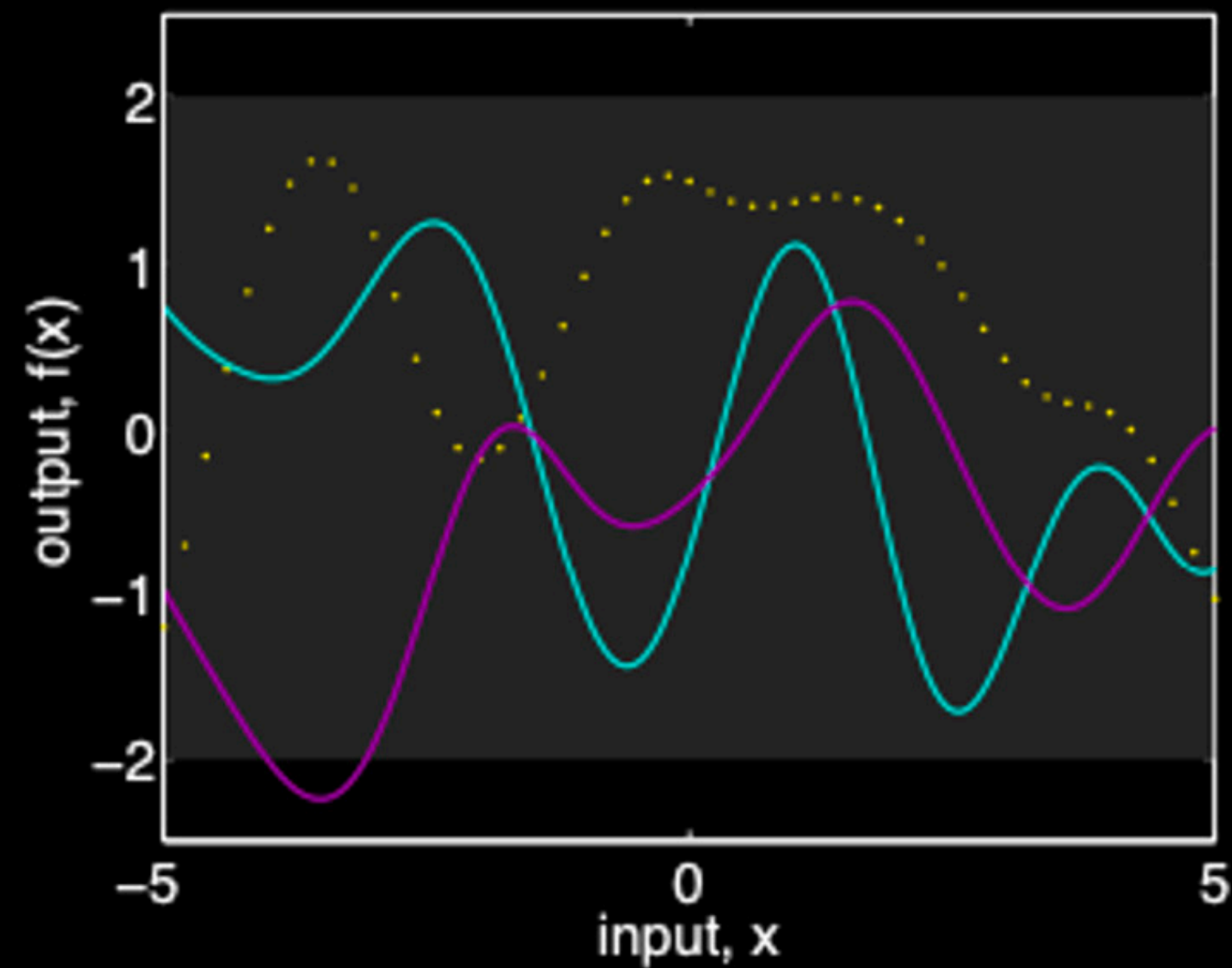
(Bishop eq 6.60, GP Book eq 2.17)

Assuming the prior “weight”

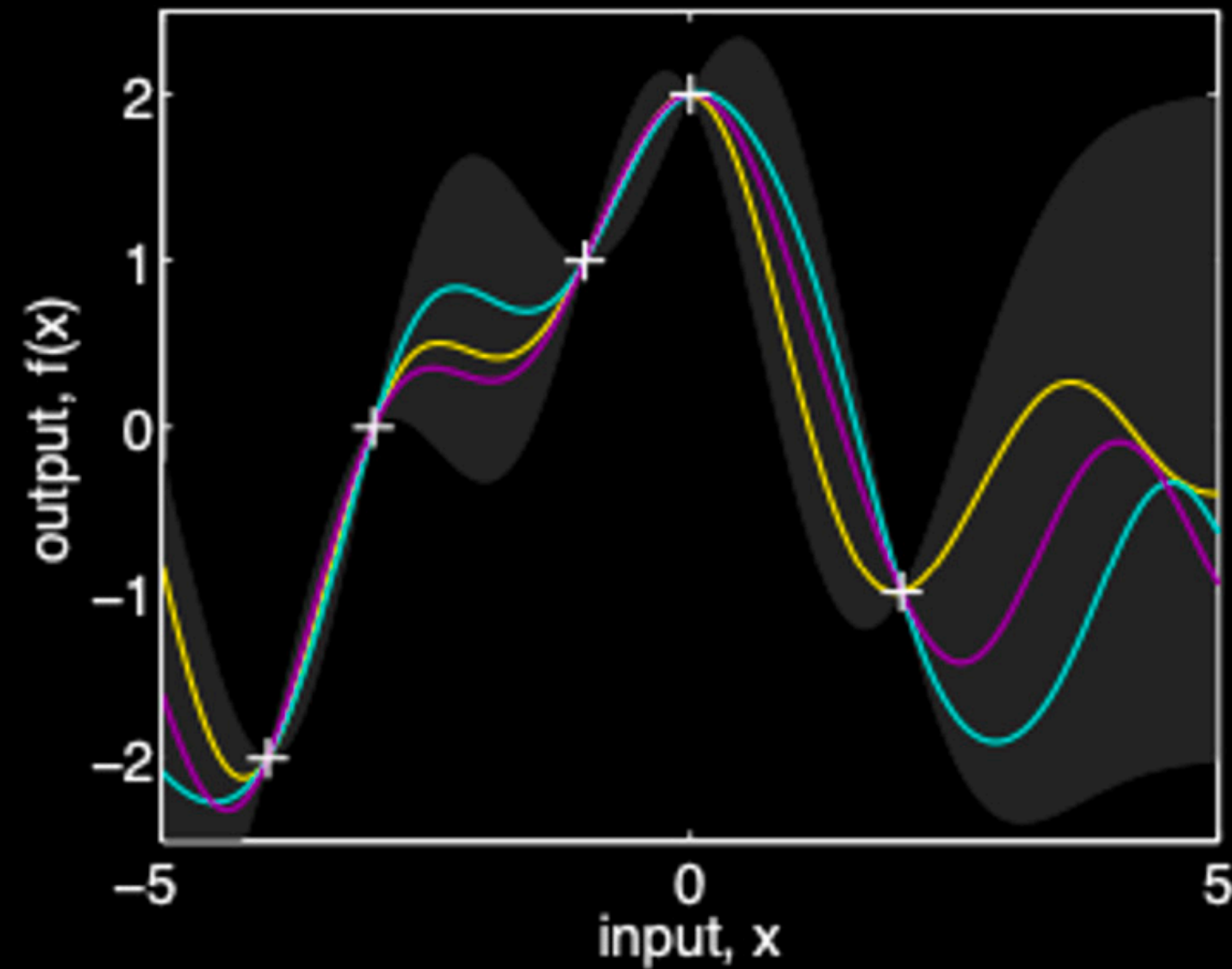
$p(w)$ to have mean zero

Function-Space Perspective

Note : prior mean of zero



(a), prior



(b), posterior

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

$$p(y, y^*) = \mathcal{N}(\mathbf{0}, \mathbf{K}_y) \quad (\text{Bishop eq 6.64})$$


 $\in \mathbb{R}^{N+1, N+1}$

$$K_{y,mn} = k(x_m, x_n) + \underbrace{\sigma^2 \delta_{mn}}_{\text{observation noise}} \quad (\text{Bishop eq 6.62})$$

Our prior belief over the kind of
functions (smoothness) we expect

$$p(y, y^*) = \mathcal{N}(\mathbf{0}, \mathbf{K}_y) \quad (\text{Bishop eq 6.64})$$

$\underbrace{\hspace{1.5cm}}_{\in \mathbb{R}^{N+1, N+1}}$

$$K_{y,mn} = k(x_m, x_n) + \underbrace{\sigma^2 \delta_{mn}}_{\text{observation noise}} \quad (\text{Bishop eq 6.62})$$

Our prior belief over the kind of functions (smoothness) we expect

$$p(y, y^*) = \mathcal{N}(\mathbf{0}, \underbrace{\mathbf{K}_y}_{\in \mathbb{R}^{N+1, N+1}}) \quad (\text{Bishop eq 6.64})$$

Also explains why we can deal with infinite-dimensional functional space in GP. We only study a “finite” distribution on data points that concerns us

$$K_{y,mn} = k(x_m, x_n) + \underbrace{\sigma^2 \delta_{mn}}_{\text{observation noise}} \quad (\text{Bishop eq 6.62})$$

Joint “prior” distribution

$$p(y, y^*) = \mathcal{N}(\mathbf{0}, \mathbf{K}_y)$$

Joint “prior” distribution

$$p(y, y^*) = \mathcal{N}(\mathbf{0}, \mathbf{K}_y)$$

Conditional distribution of Gaussians is trivial

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

Joint “prior” distribution

$$p(y, y^*) = \mathcal{N}(\mathbf{0}, \mathbf{K}_y)$$

$$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

(Bishop eq 2.81-2.82)

Conditional distribution of Gaussians is trivial

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

Joint “prior” distribution

$$p(y, y^*) = \mathcal{N}(\mathbf{0}, \mathbf{K}_y)$$

$$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

(Bishop eq 2.81-2.82)

Conditional distribution of Gaussians is trivial

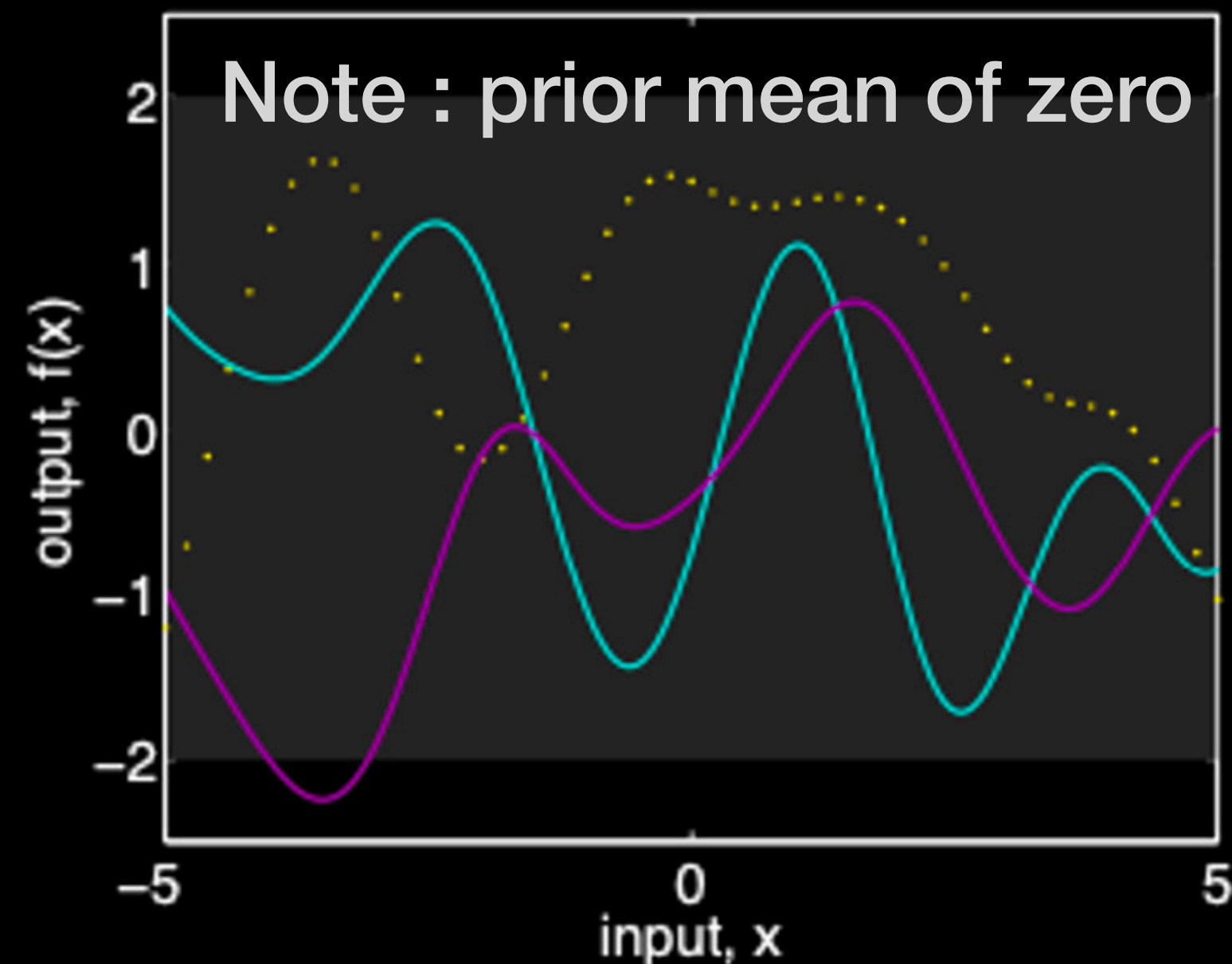
$$p(y^* | y) = p(y^* | \mathbf{x}^*, \mathbf{X}, y) = \mathcal{N}(y^*; m(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$$

$$m(x^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} y \quad (\text{Bishop eq 6.66, GP Book eq 2.19, 2.25})$$

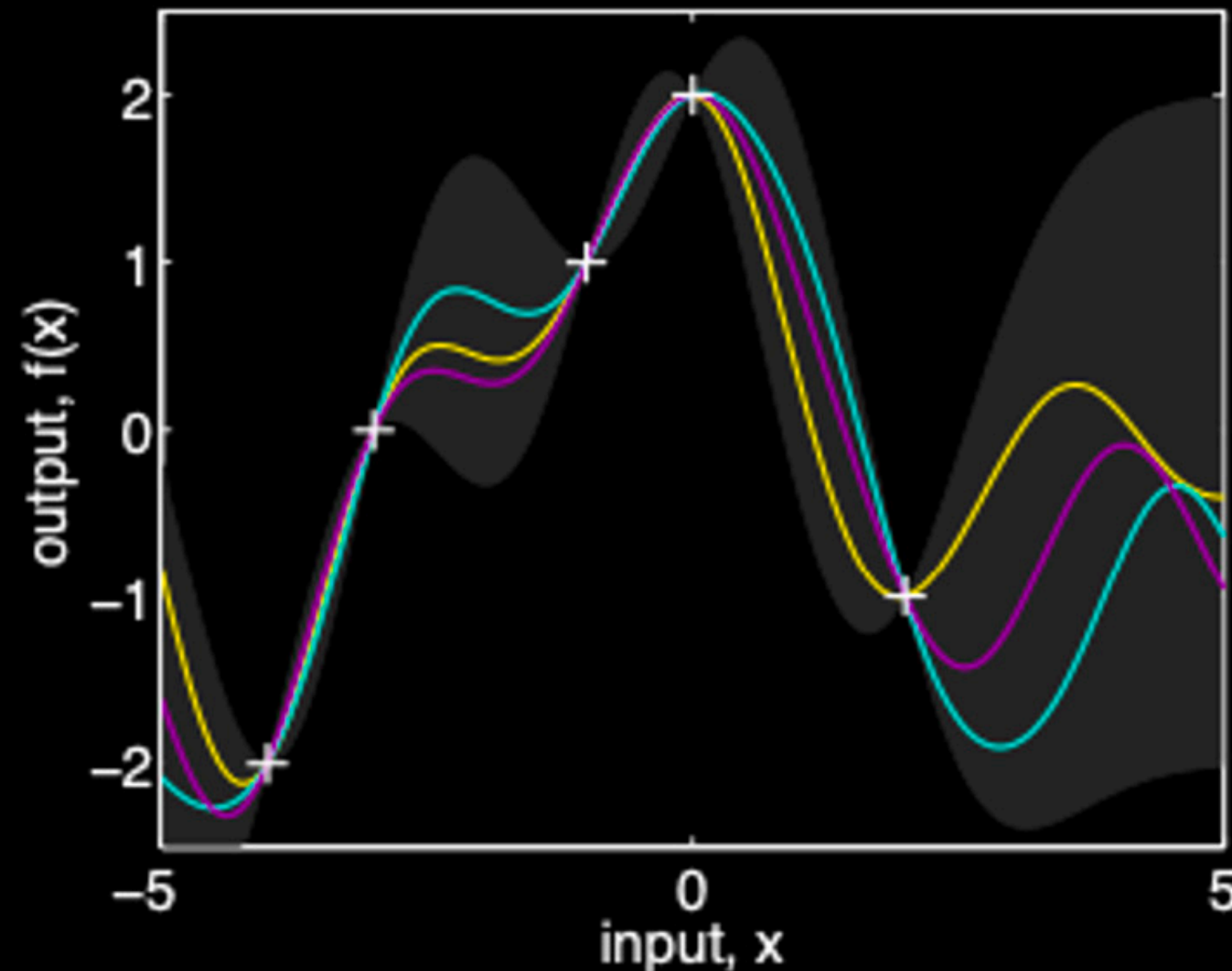
$$\sigma^2(x^*) = \sigma^2 + k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*)$$

(Bishop eq 6.67, GP Book eq 2.19, 2.26)

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



(a), prior



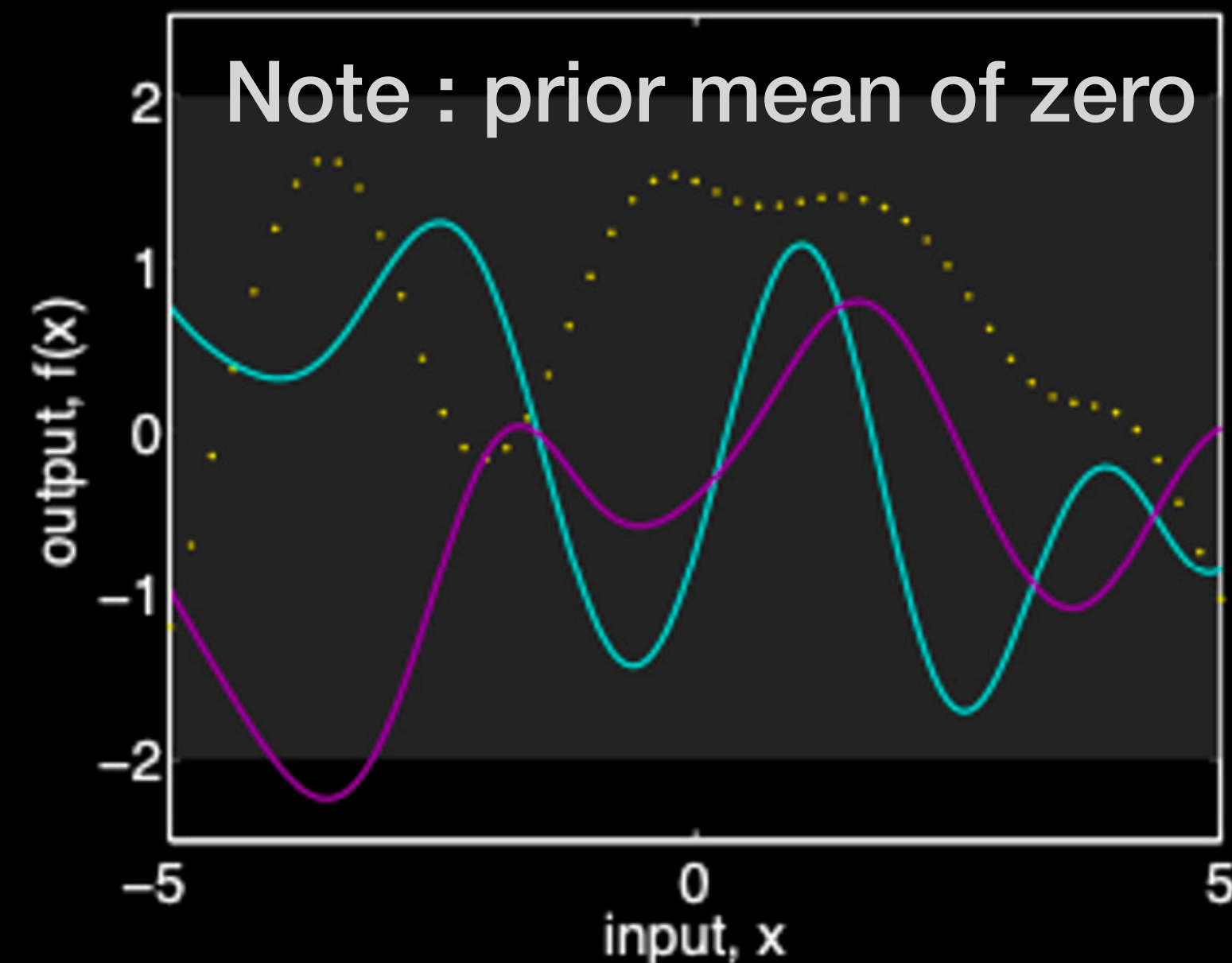
(b), posterior

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

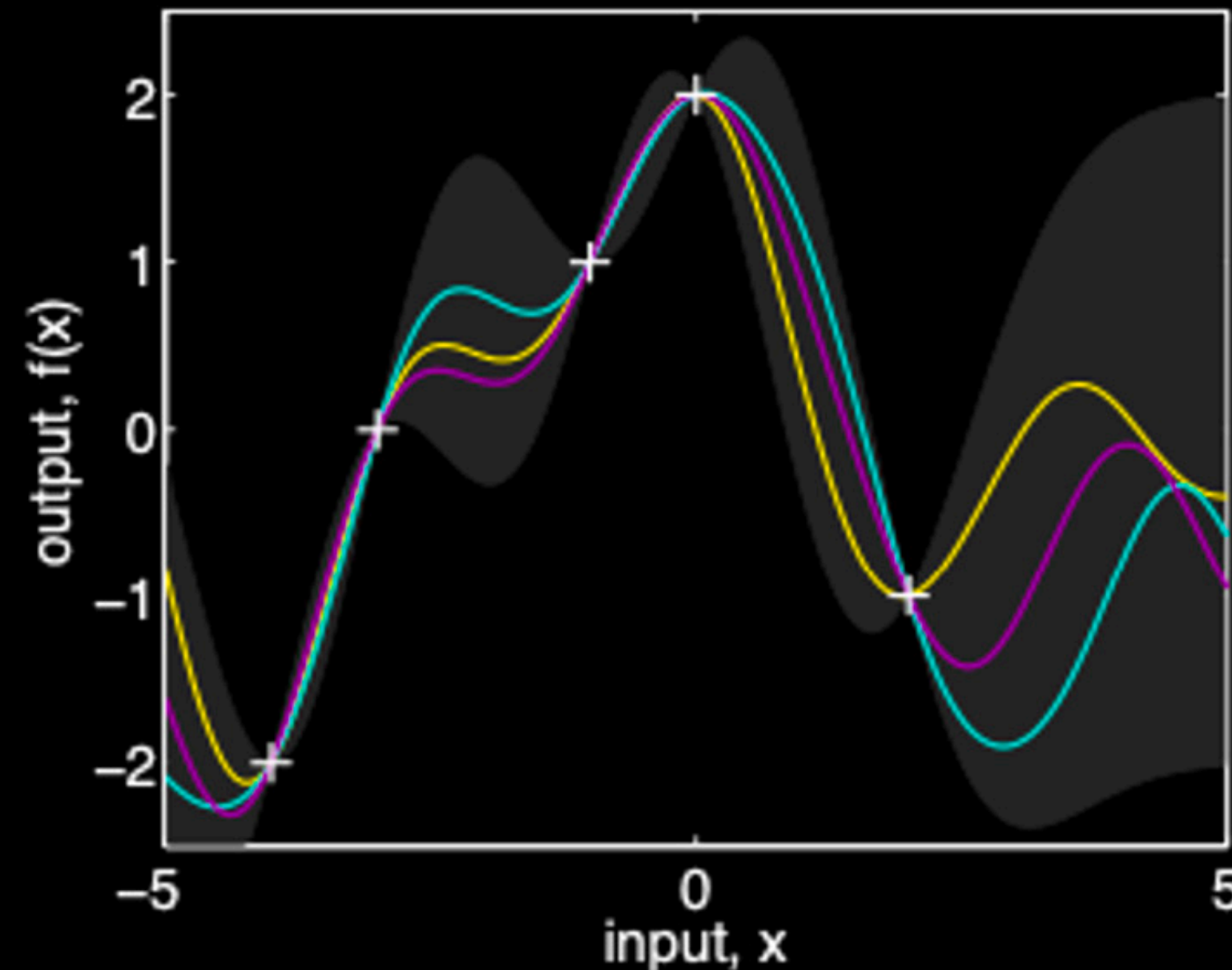
$$m(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X}) (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

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

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



(a), prior



(b), posterior

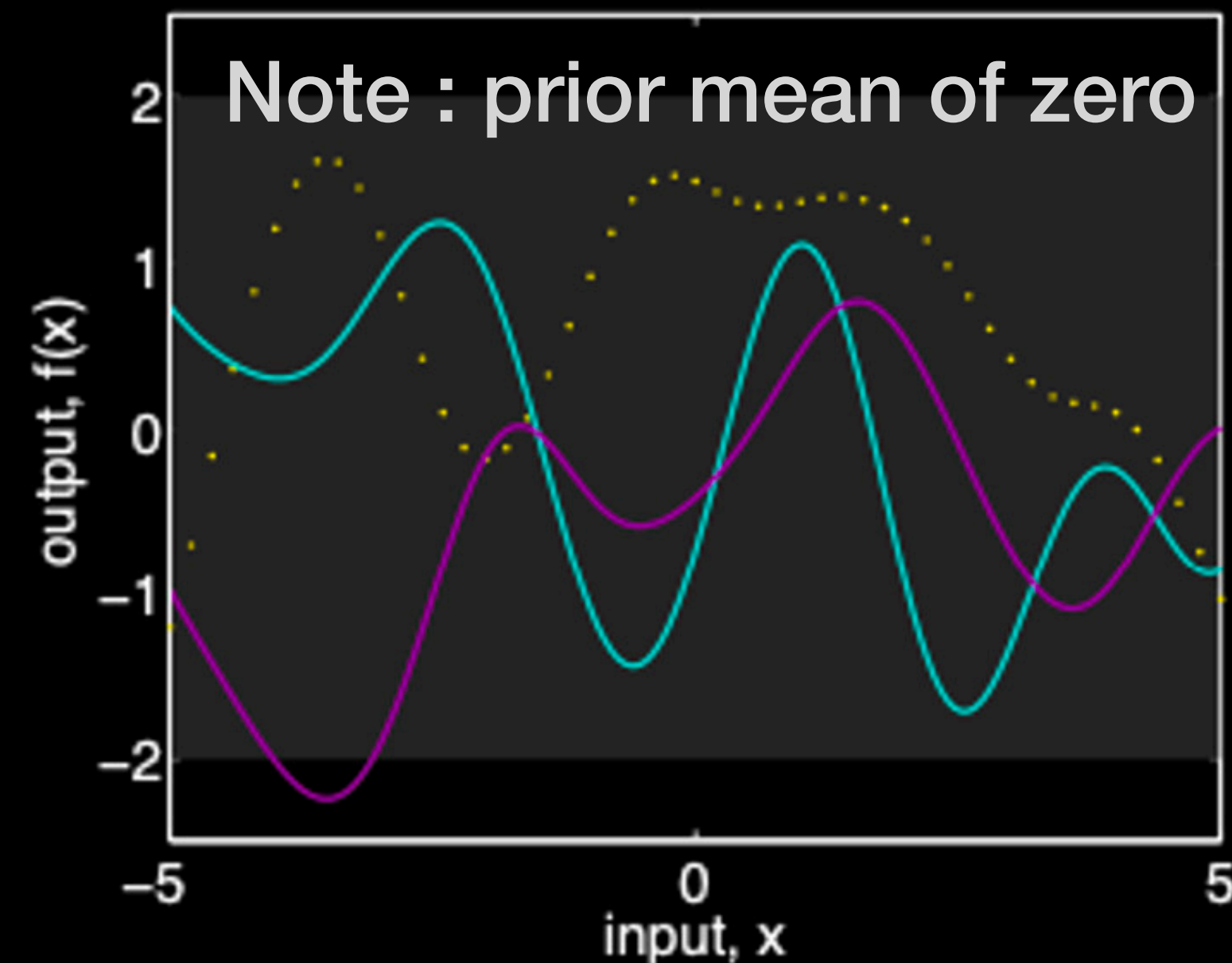
Posterior distribution of functions : “rejecting” functions that do not satisfy the constraints

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

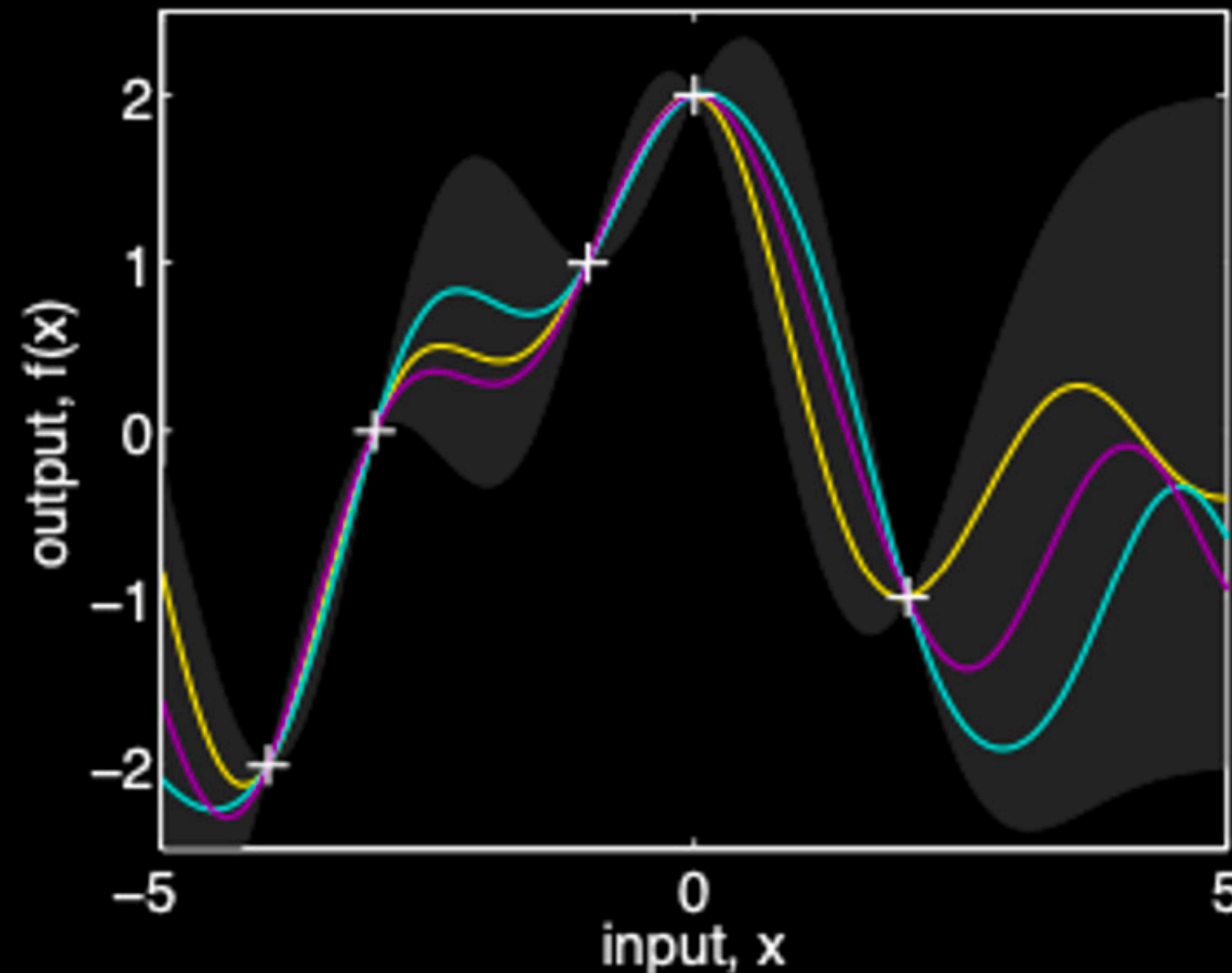
The influence depends on the proximity

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

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)

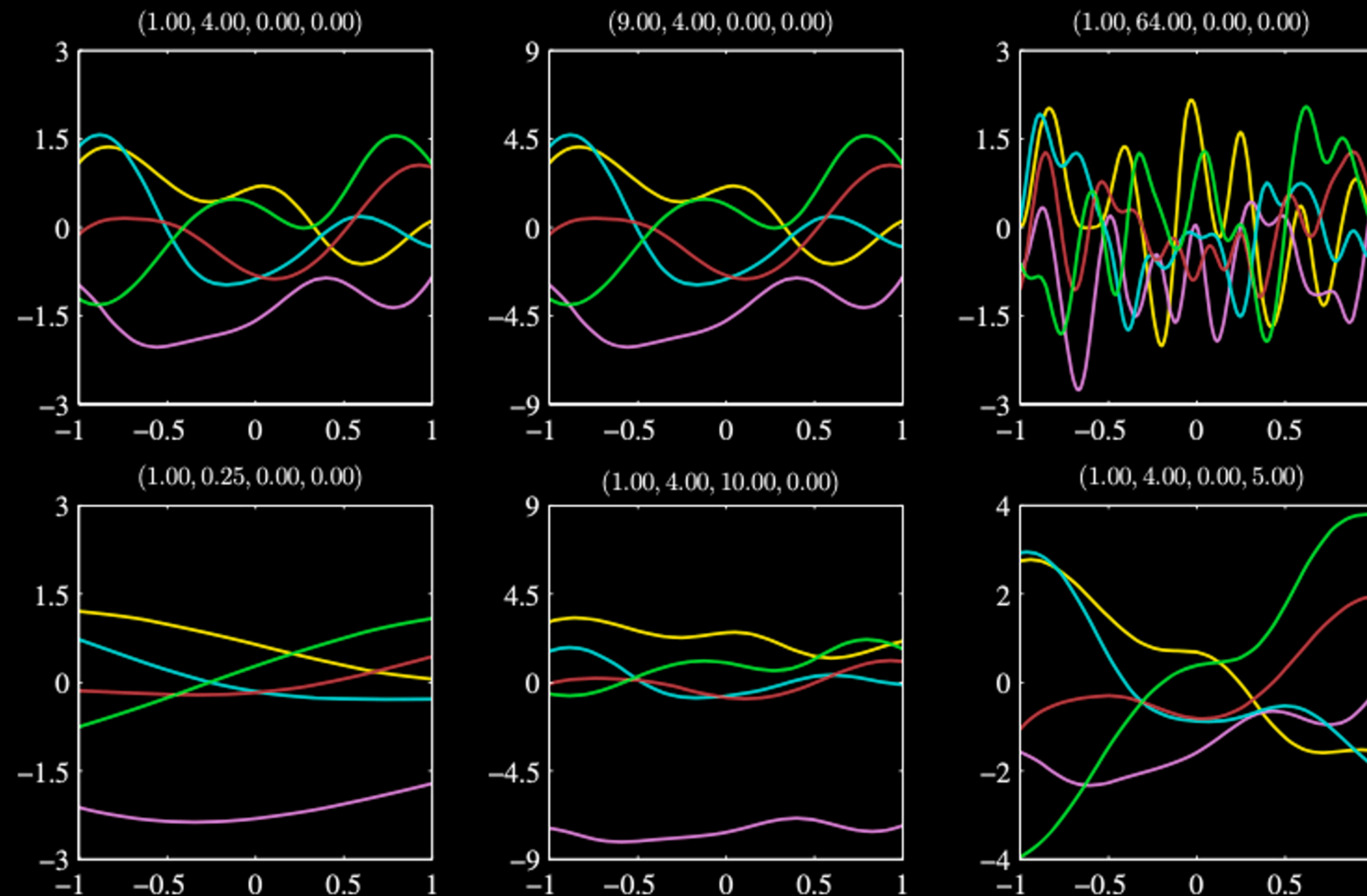


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

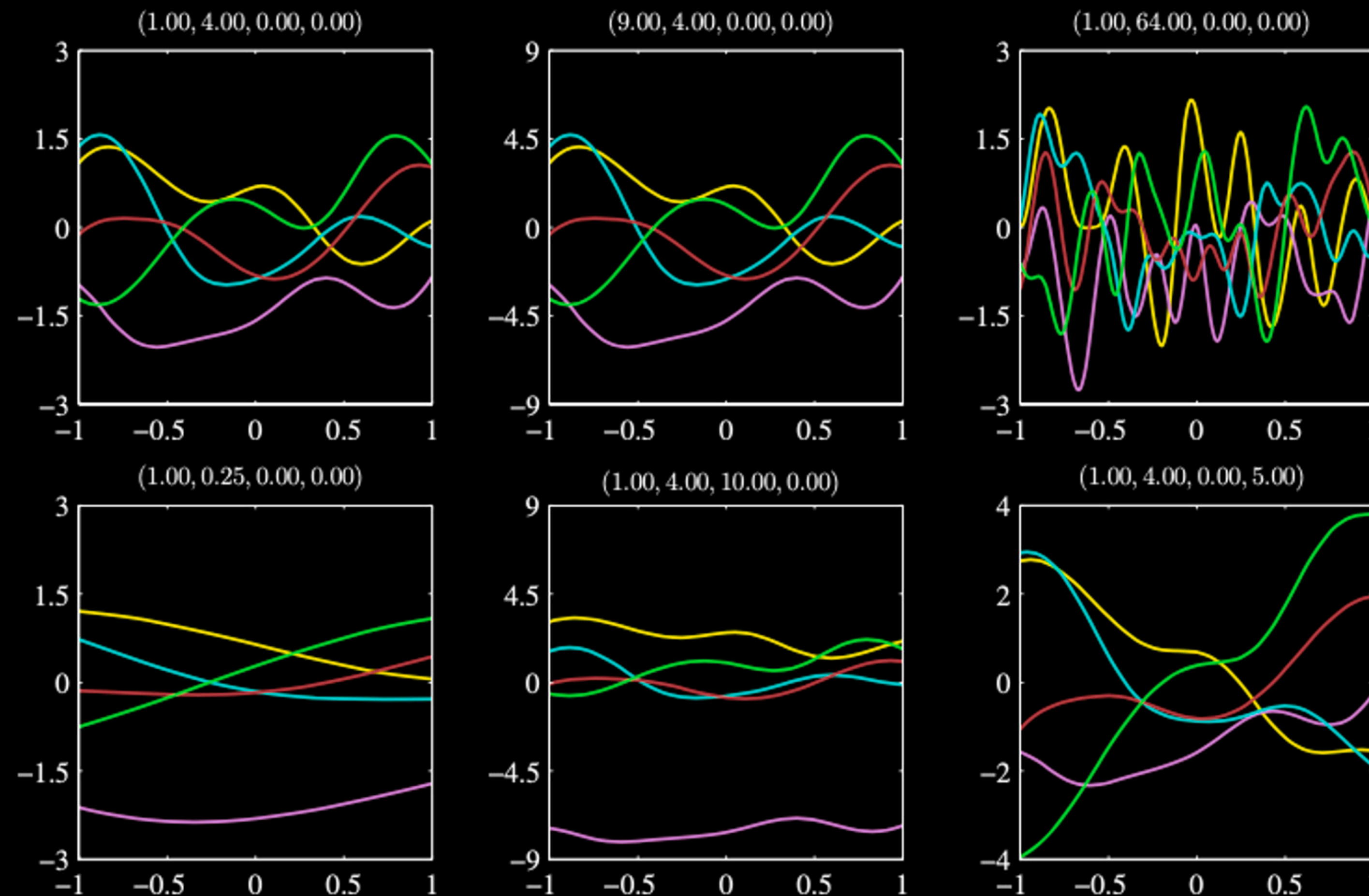


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

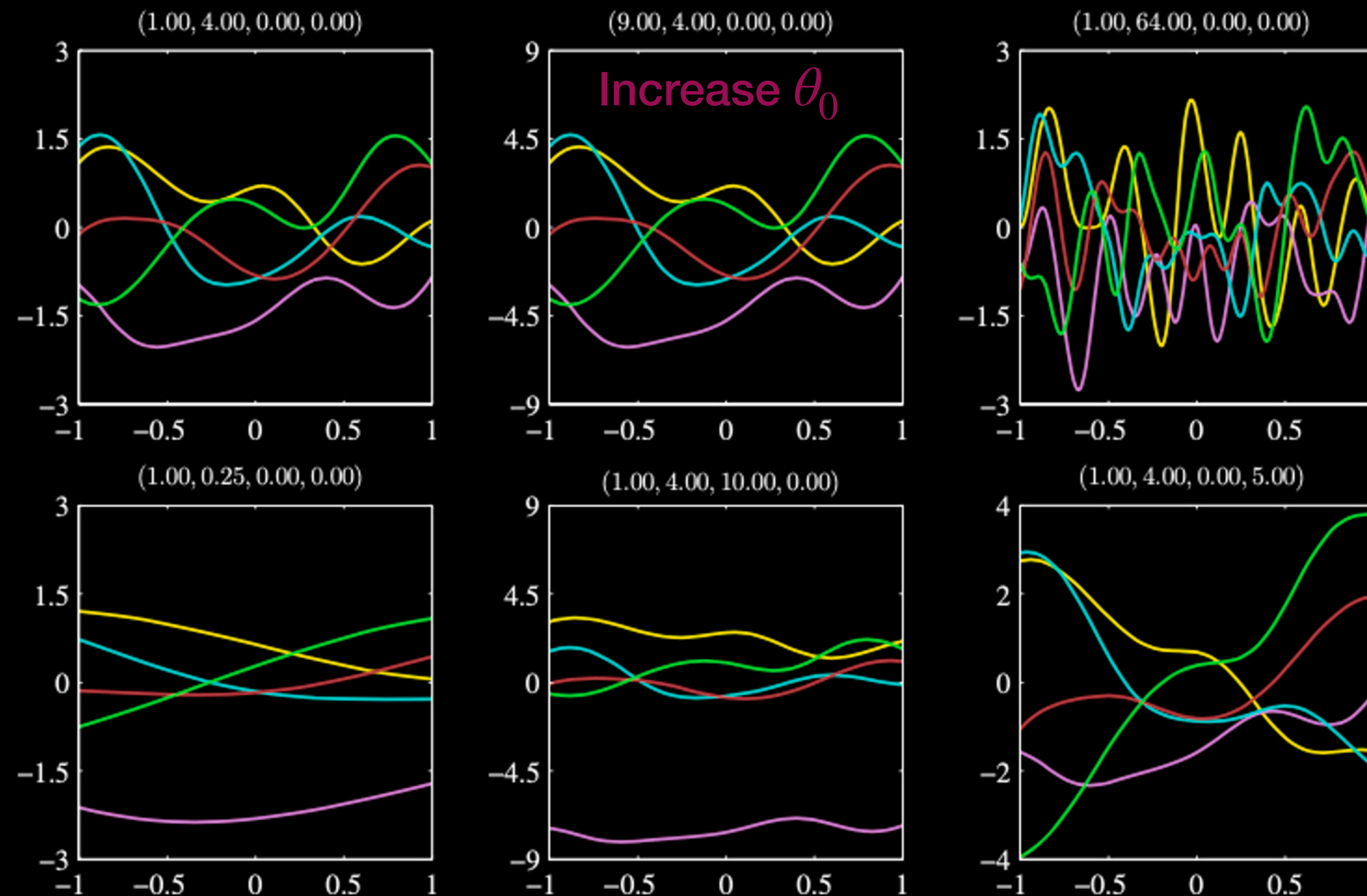


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

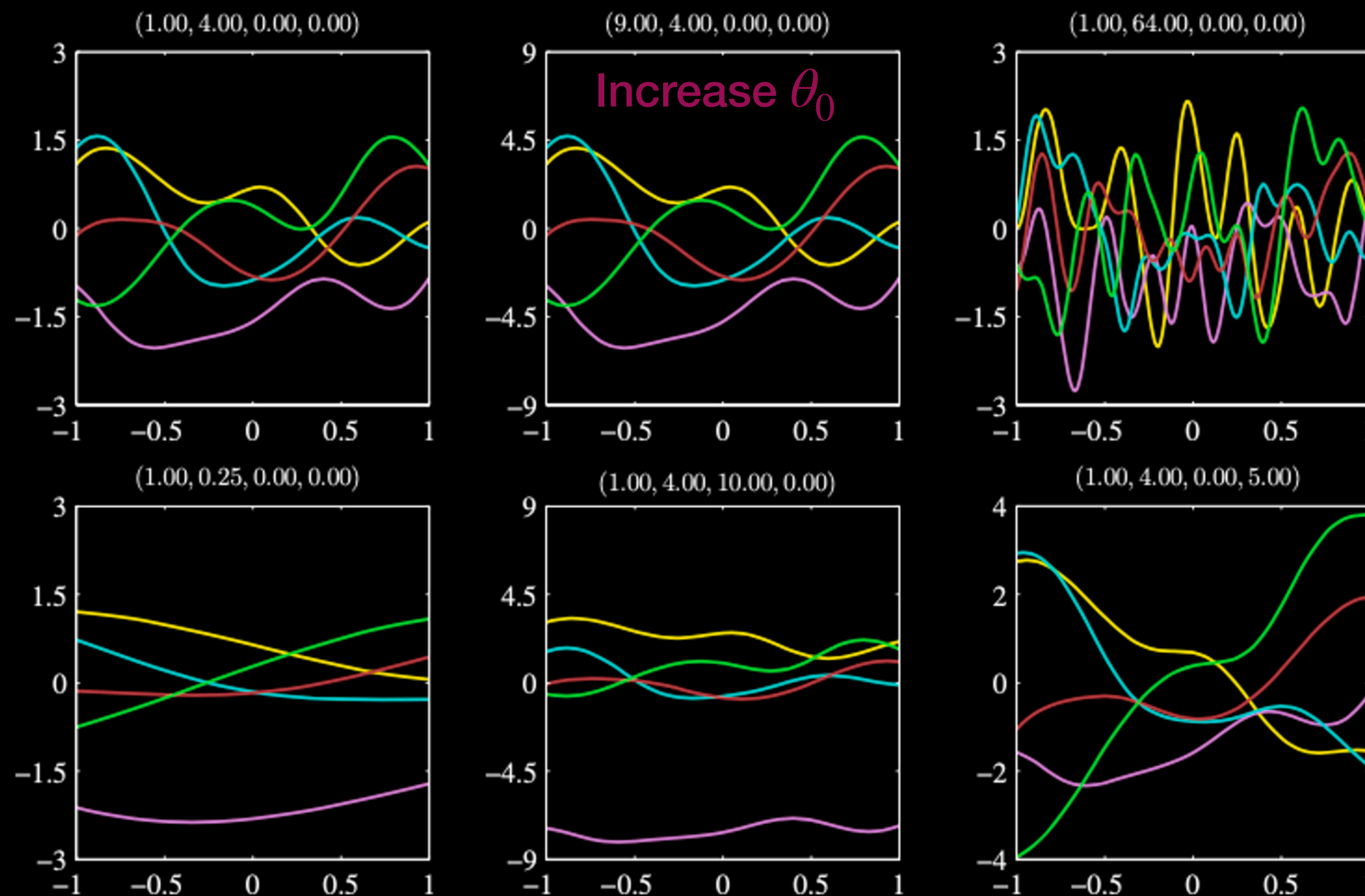


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

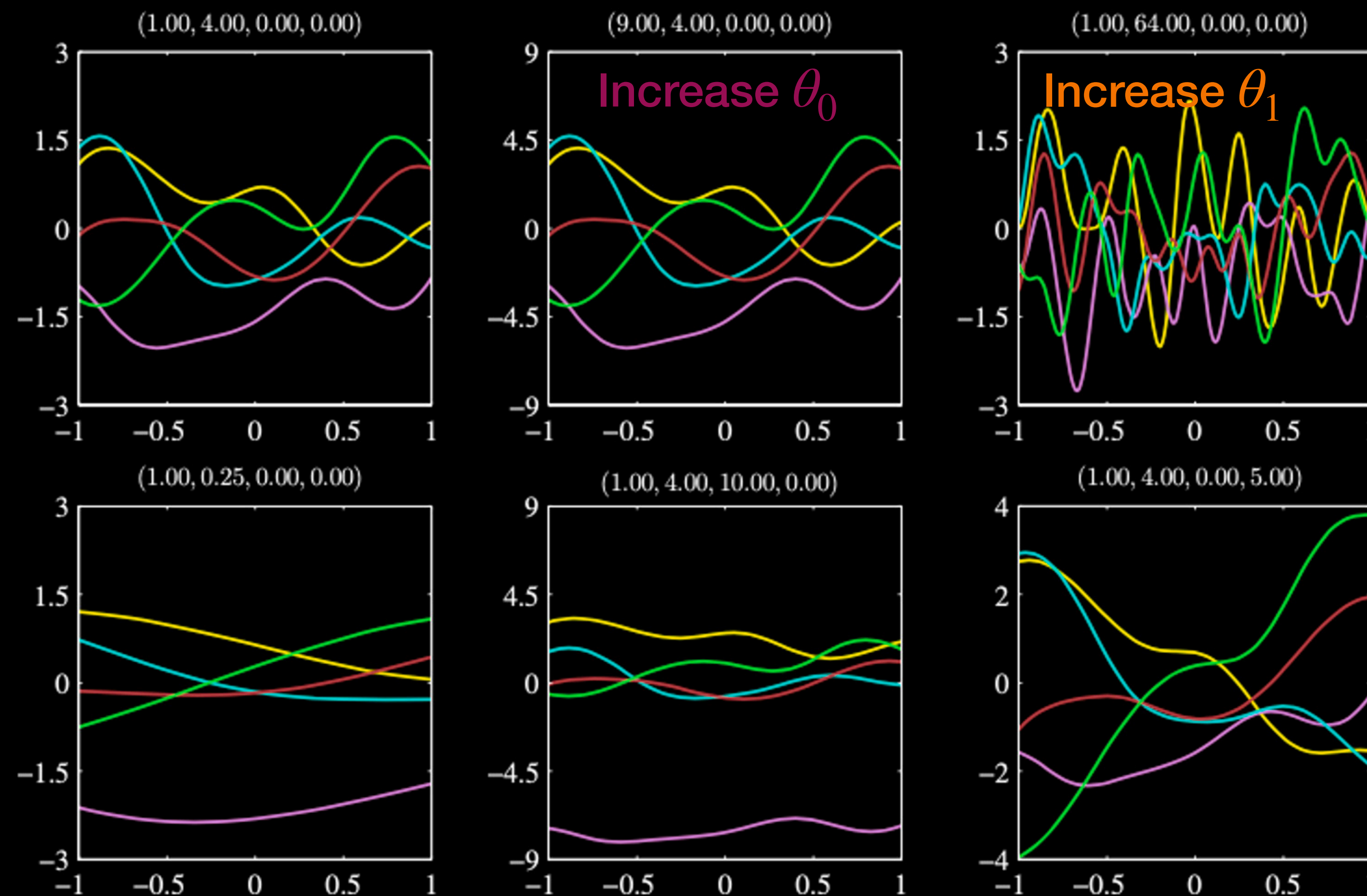


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

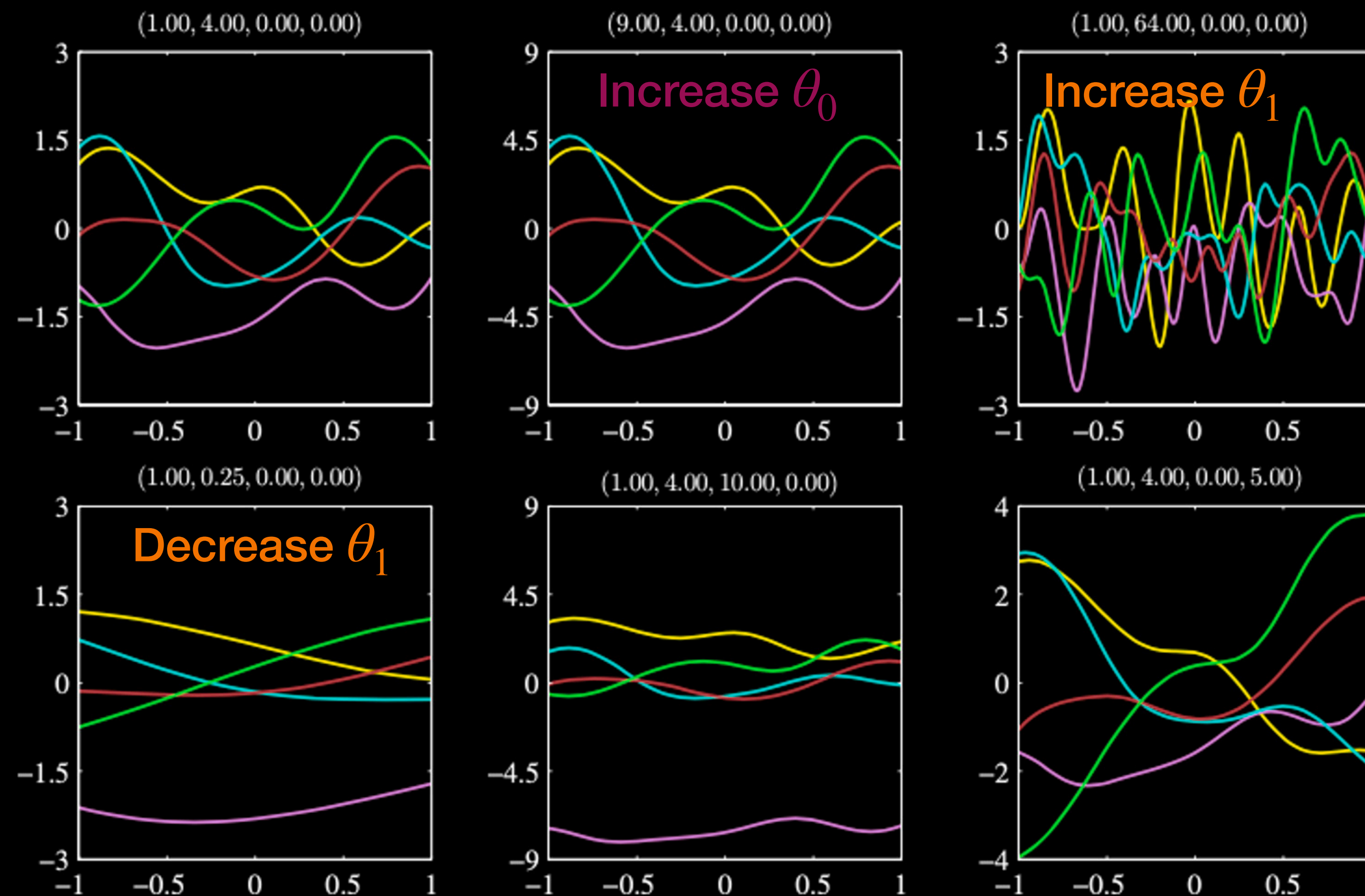


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

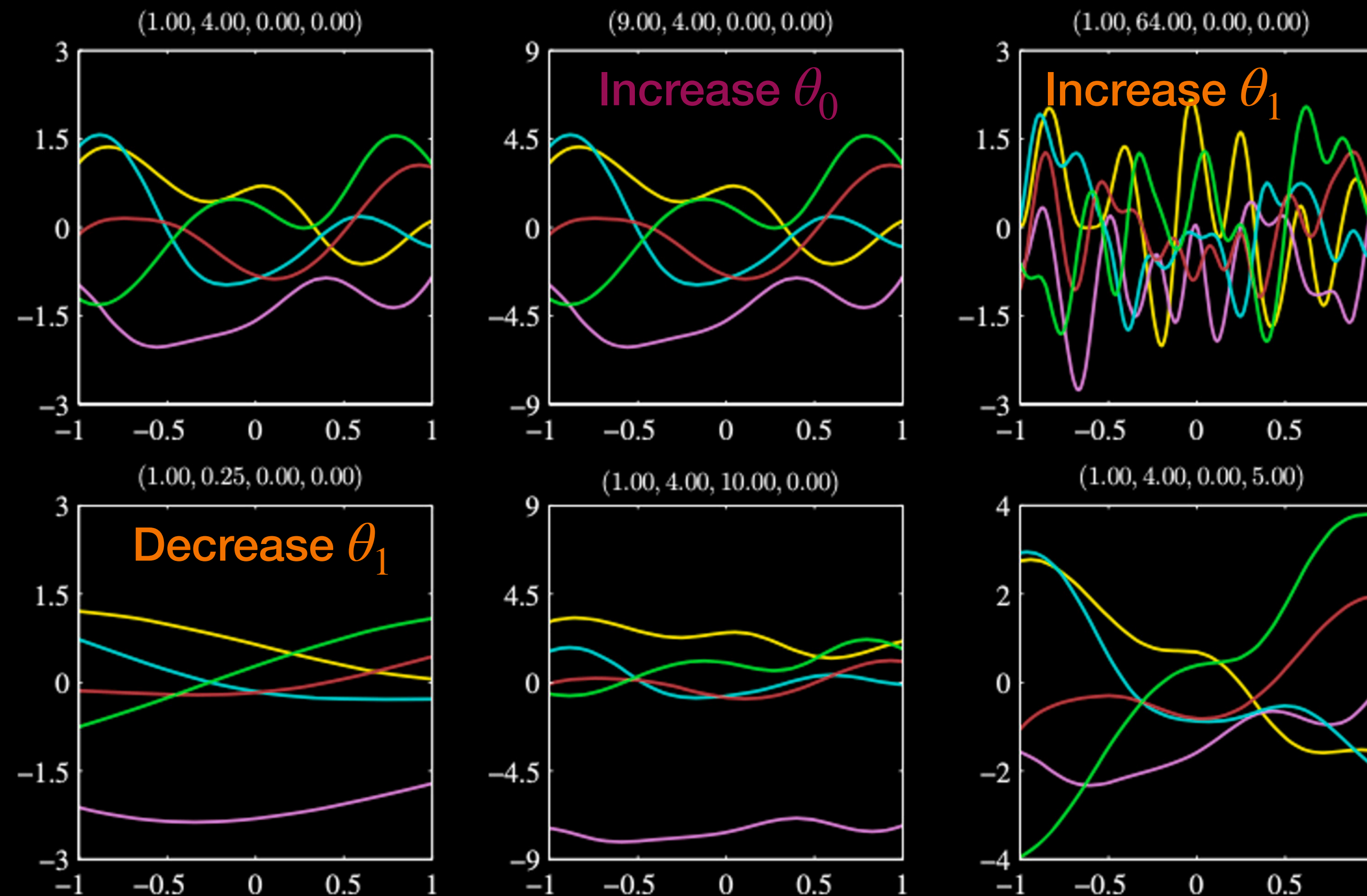


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

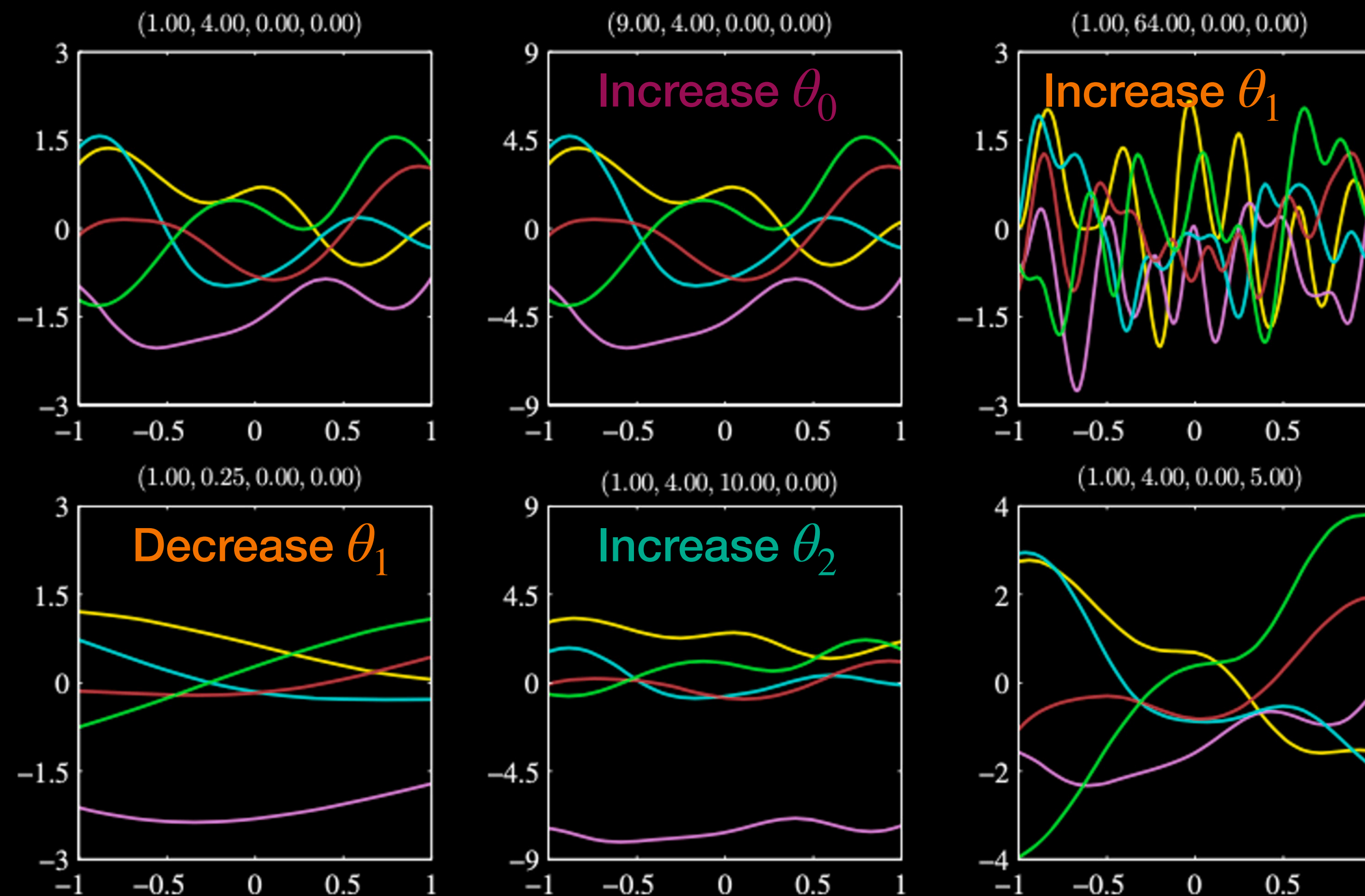


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

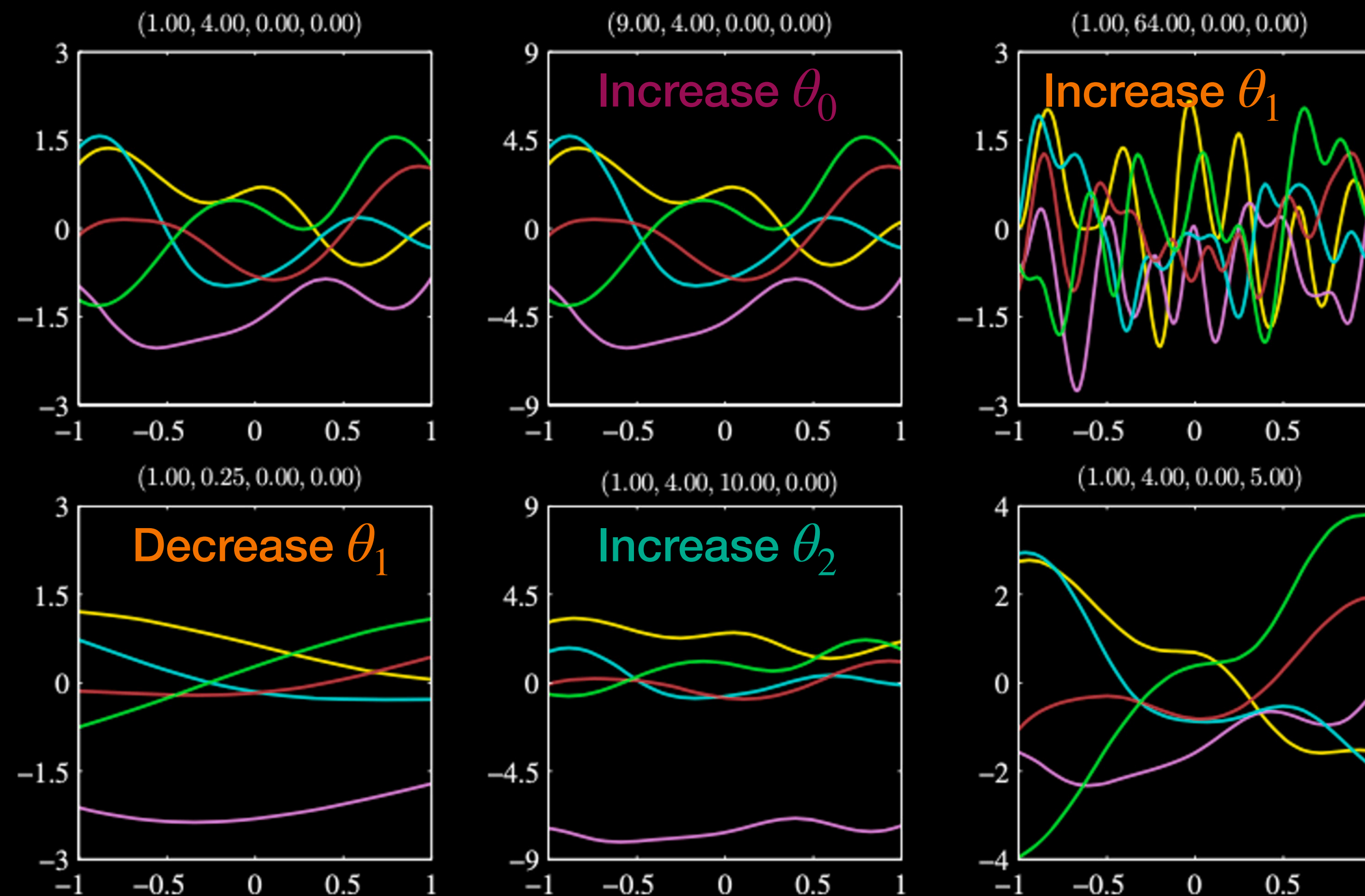


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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)

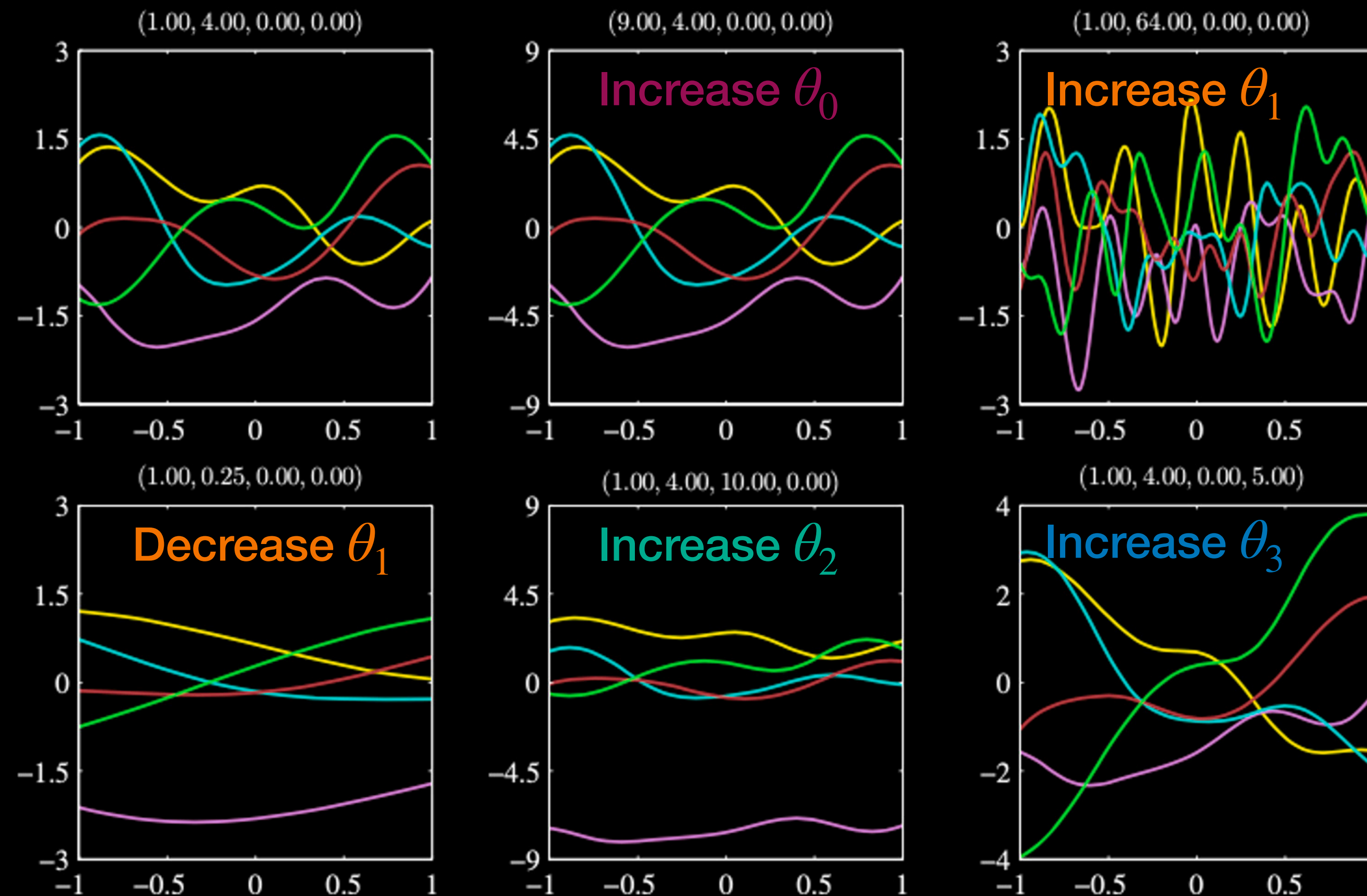


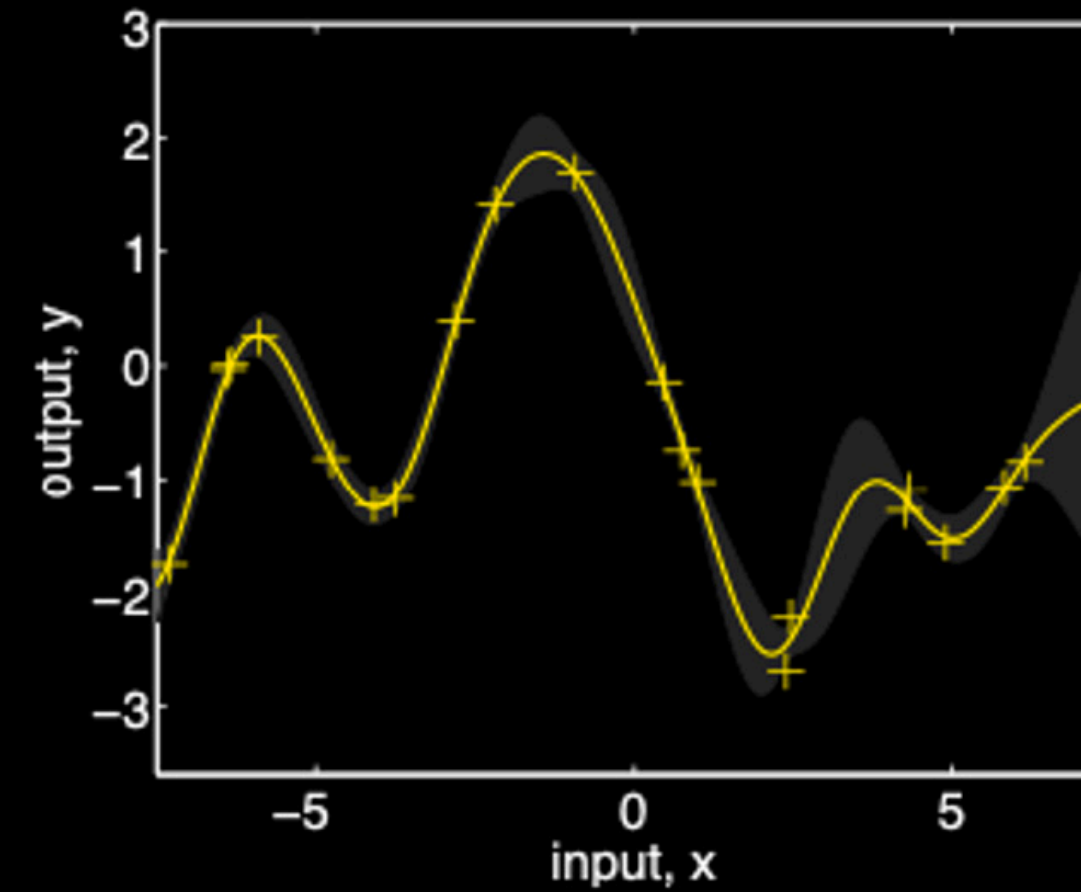
Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

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

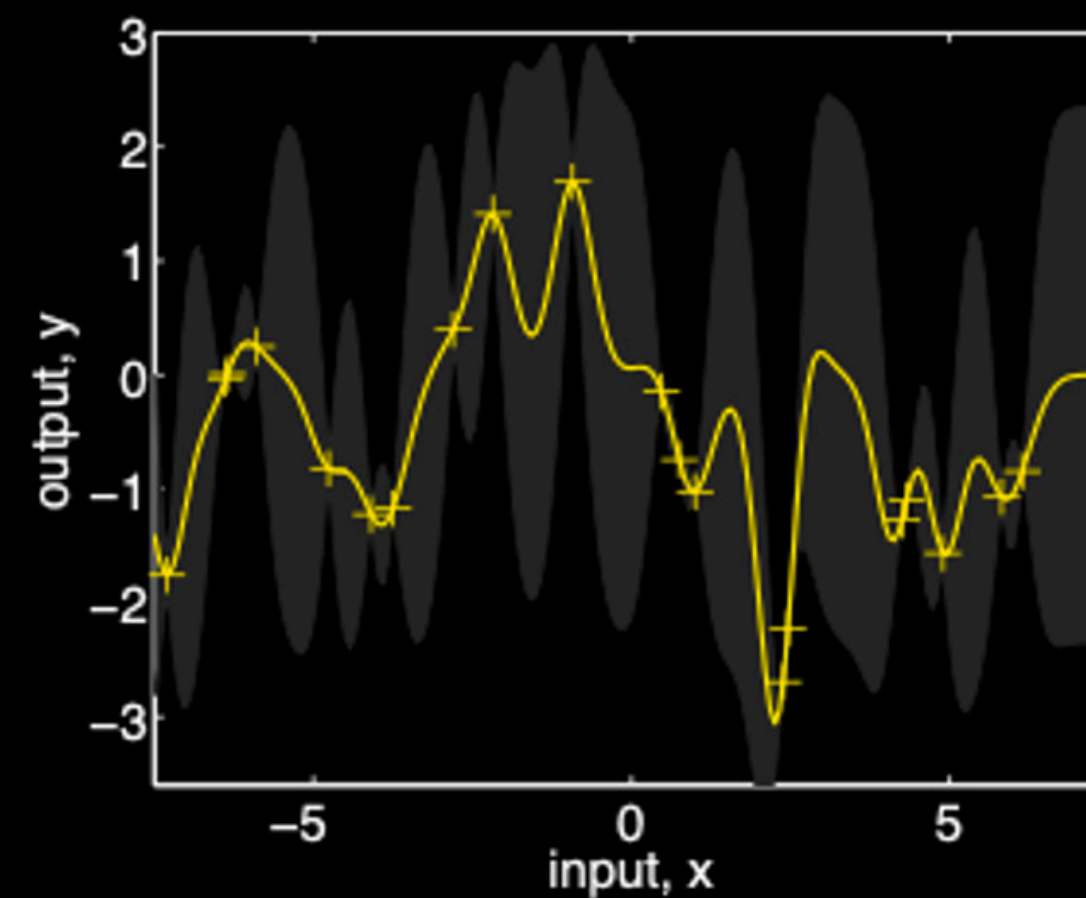
Visualising Kernels

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

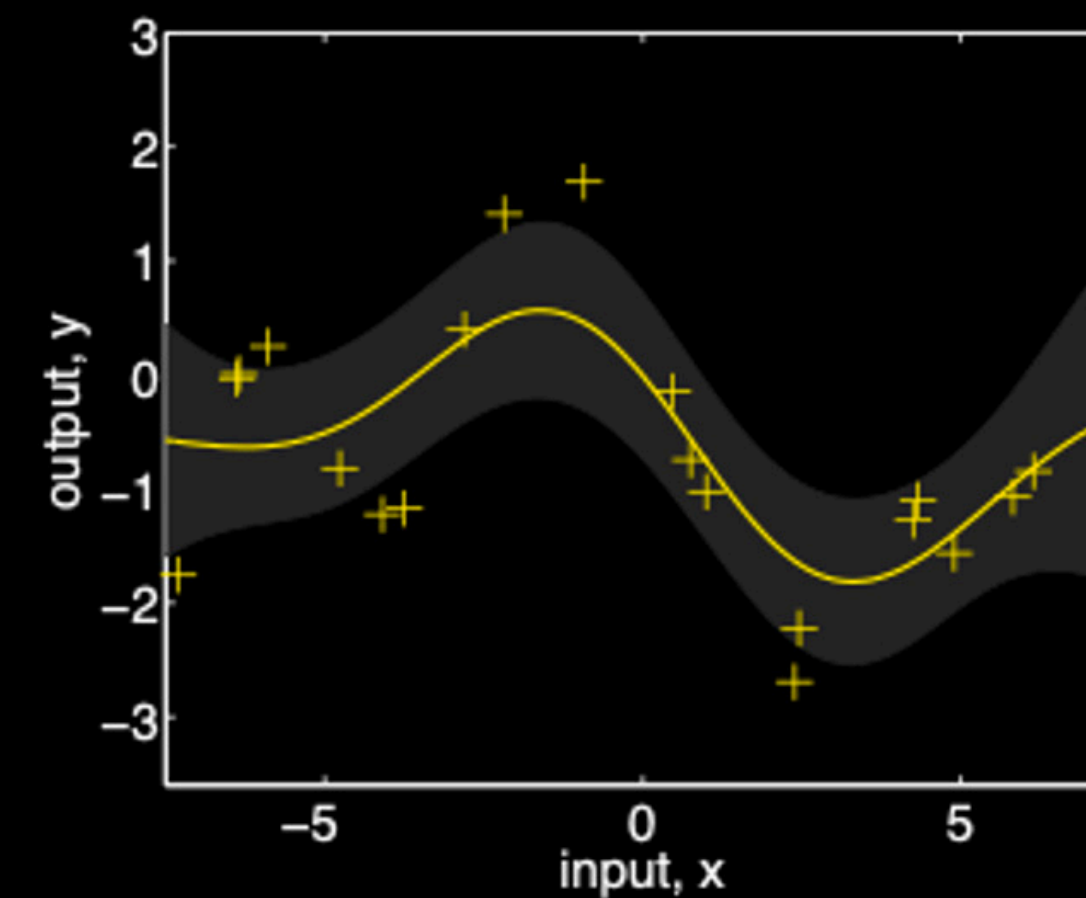
Gaussian Process,
with constraints



(a), $\ell = 1$



(b), $\ell = 0.3$



(c), $\ell = 3$

(GP book)

Bayesian Linear Regression

Gaussian Process

Bayesian Linear Regression

- When $D \ll N$, faster training and inference

Gaussian Process

Bayesian Linear Regression

- When $D \ll N$, faster training and inference
- But **restrict** to the functional class we assume

Gaussian Process

Bayesian Linear Regression

- When $D \ll N$, faster training and inference
 - But **restrict** to the functional class we assume
-

Gaussian Process

- Prior distribution assumed to be **all possible functions** (subjected to the kernel assumption and hyper parameters)

Bayesian Linear Regression

- When $D \ll N$, faster training and inference
 - But **restrict** to the functional class we assume
-

Gaussian Process

- Prior distribution assumed to be **all possible functions** (subjected to the kernel assumption and hyper parameters)
- Implicitly assume that $D \gg N$ [depend on kernels], do not need to worry about lack of degrees of freedom to fit all the data

Bayesian Linear Regression

- When $D \ll N$, faster training and inference
 - But **restrict** to the functional class we assume
-

Gaussian Process

- Prior distribution assumed to be **all possible functions** (subjected to the kernel assumption and hyper parameters)
- Implicitly assume that $D \gg N$ [depend on kernels], do not need to worry about lack of degrees of freedom to fit all the data
- Training limited N^3 **complexity** (extra reading - Chapter 8 of the GP book, GP approximations in Week 6)