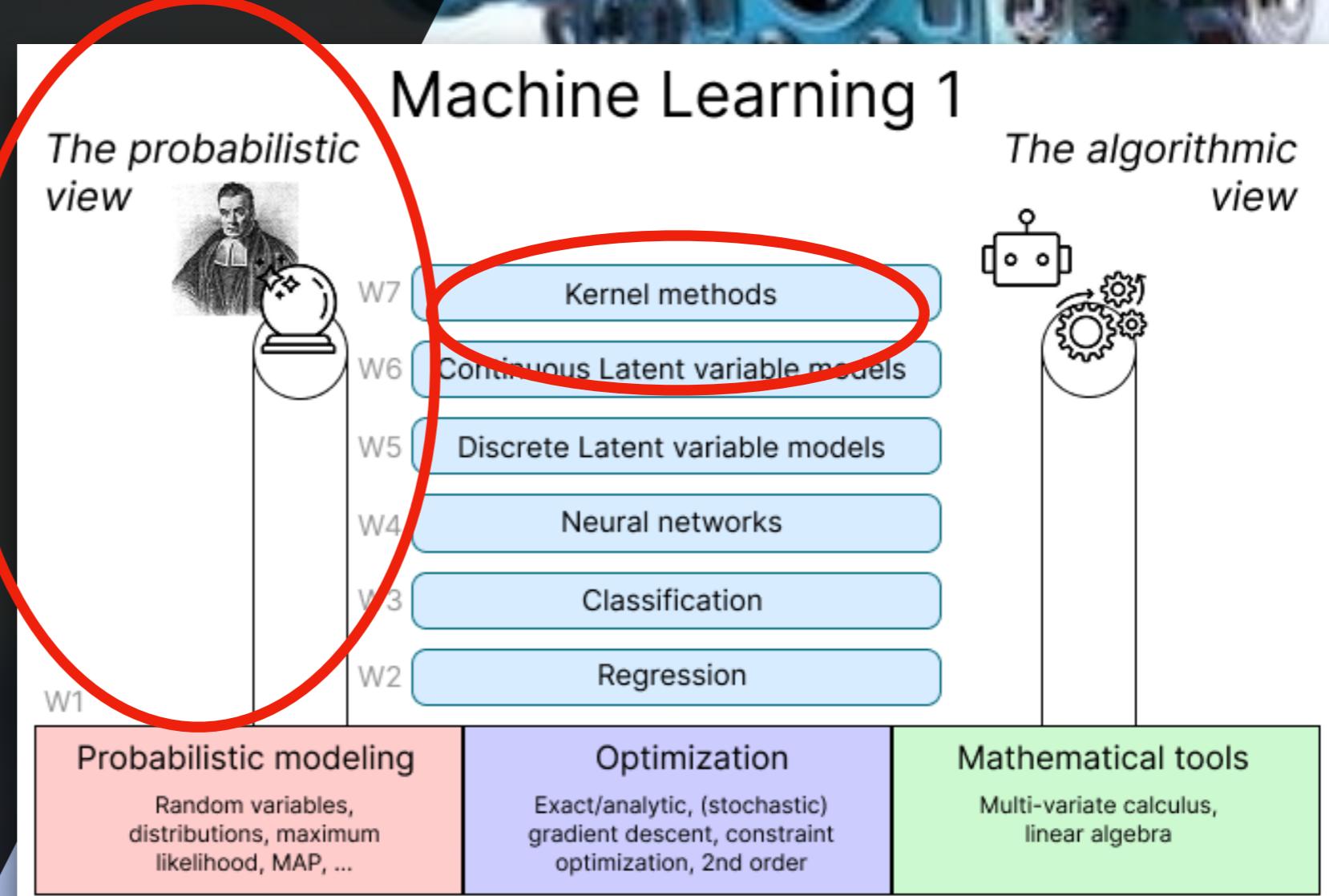


Machine Learning 1

Lecture 12 - Kernel Methods
Gaussian Processes

Erik Bekkers



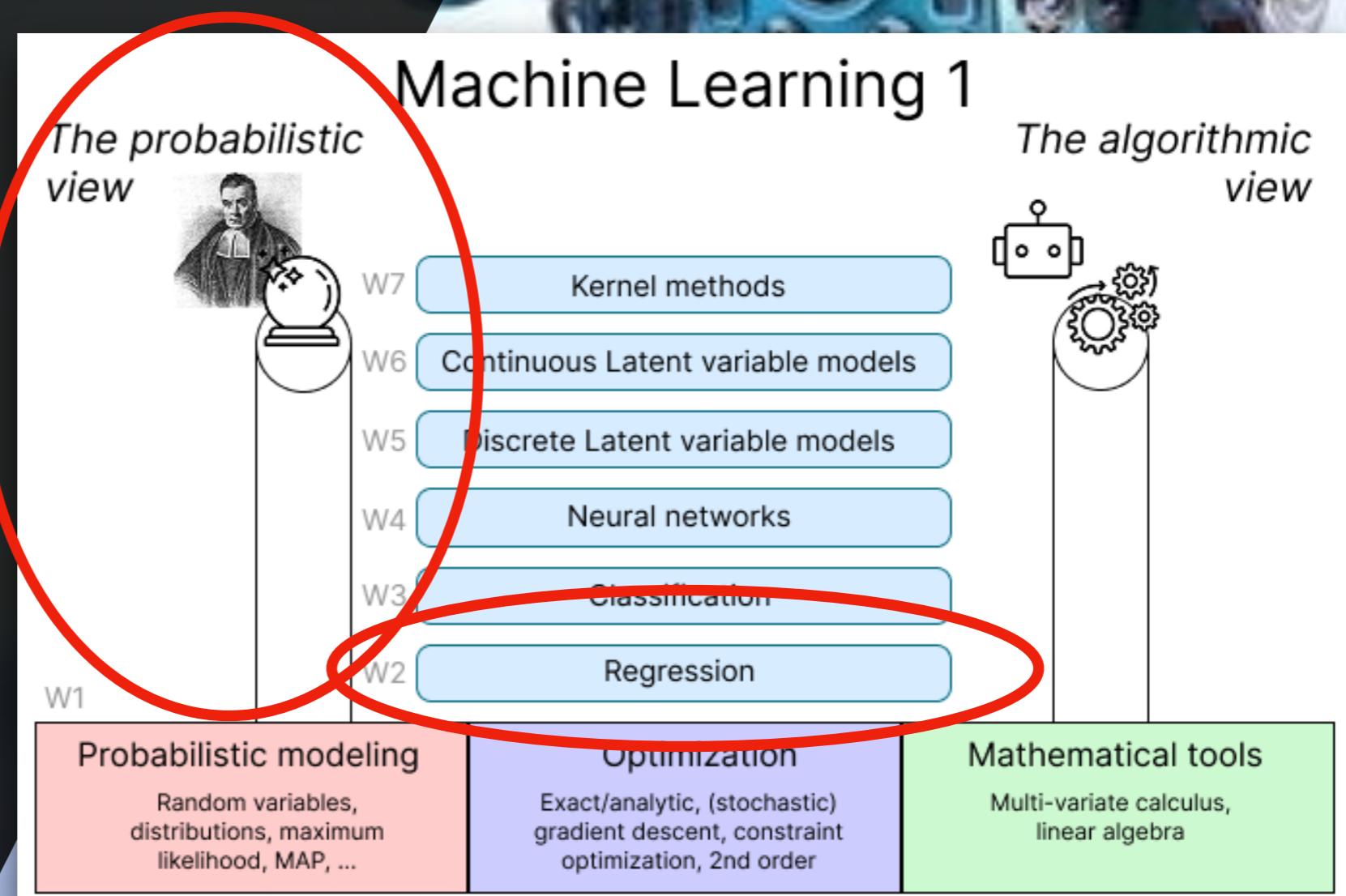
Machine Learning 1

Lecture 5.1 - Supervised Learning

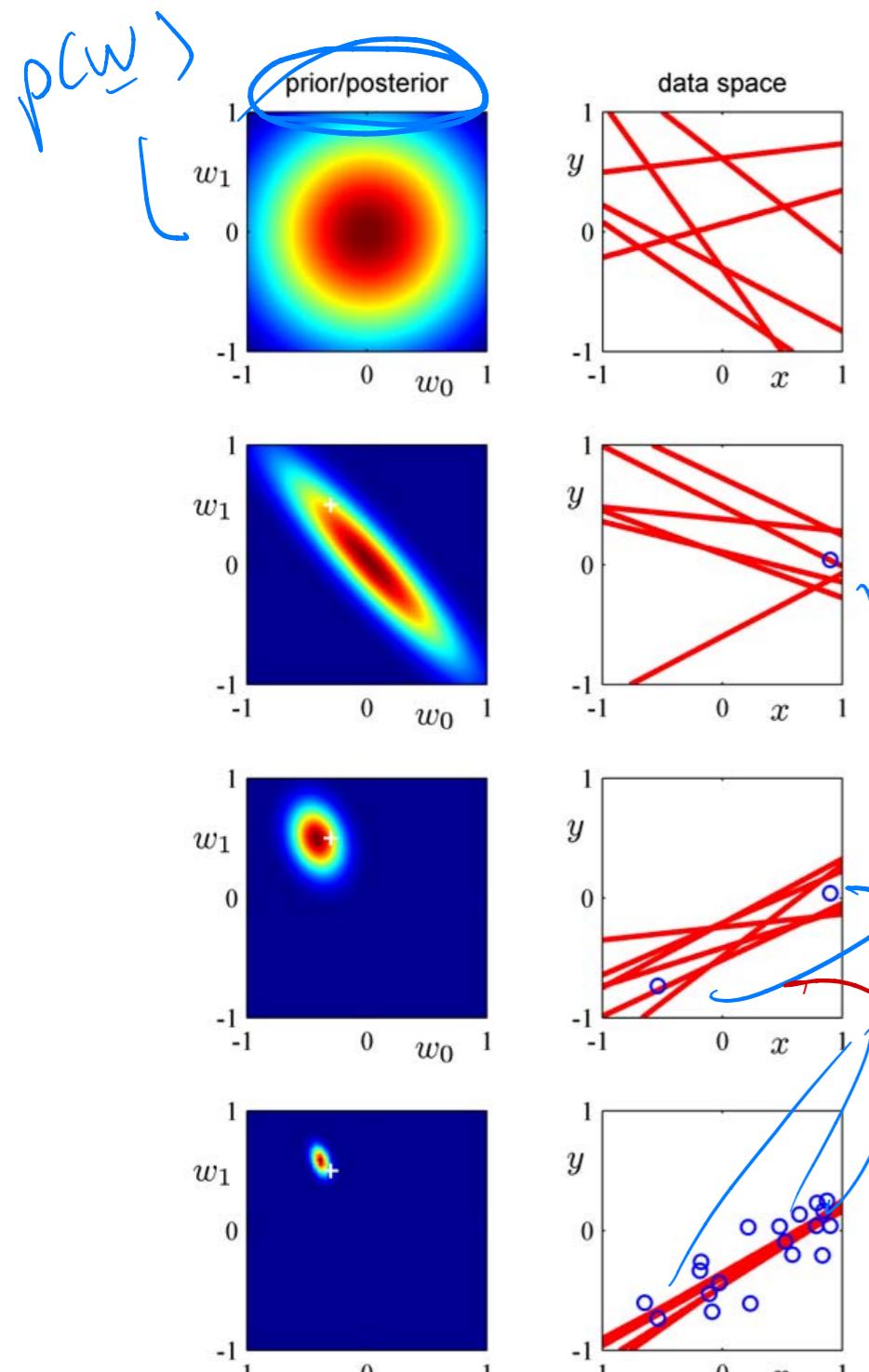
Bayesian Linear Regression - **The Equivalent Kernel**

Erik Bekkers

(Bishop 3.3.3)

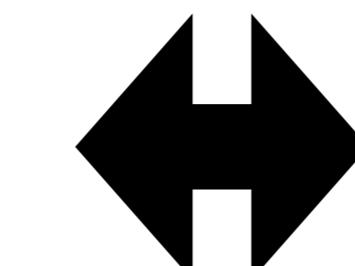
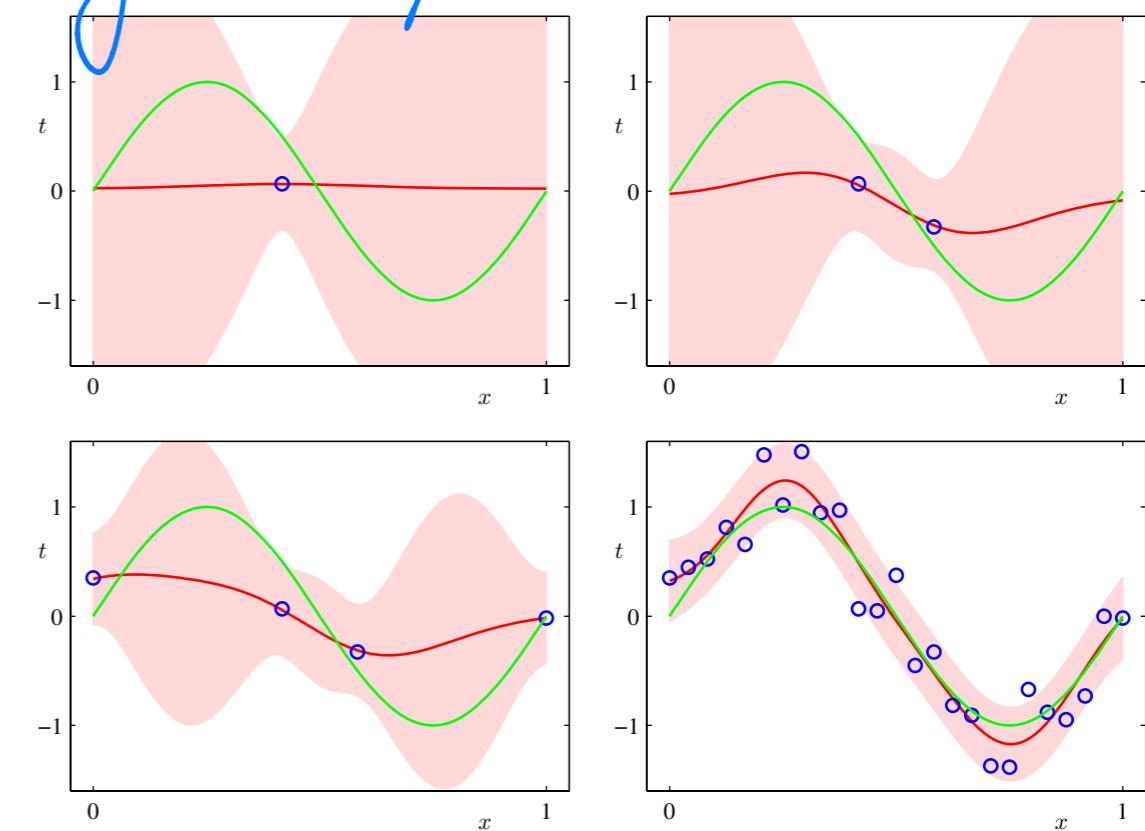


Bayesian Linear Regression



Continuous Random Functions (Gaussian Processes)

Bayesian pred. dist.



Data D

$p(w|D)$

$$1 \quad \underline{w} \sim p(w|D)$$

$$2. \text{ Plot } \underline{f}_w(x) = \underline{w}^\top \phi(x)$$

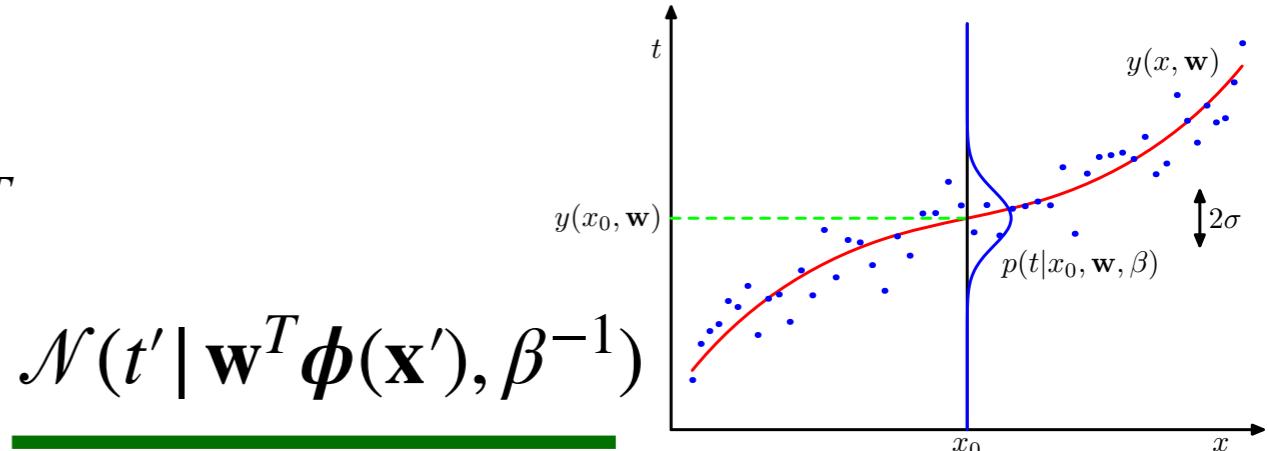
$$p(t|x, \underline{w}) = N(t | \underline{f}_w(x), \beta^{-1})$$

Bayesian Linear Regression

- Regression problem with:

- Data: $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$

- Predictive distribution** $p(t' | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t' | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}'), \beta^{-1})$



- Probabilistic model with **Gaussians**:

- Likelihood: $p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \mathcal{N}(\mathbf{t} | \Phi \mathbf{w}, \beta \mathbf{I})$

- Conjugate prior: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$

Gaussian times Gaussian is a Gaussian

- Posterior:** $p(\mathbf{w} | \mathbf{t}, \mathbf{X}) = \frac{p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w})}{p(\mathbf{t} | \mathbf{X}, \beta)} = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$

- Maximum A Posteriori estimate:

- $\mathbf{w}_{\text{MAP}} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w} | \mathbf{t}, \mathbf{X}) = \mathbf{m}_N$

Bishop Ch 2.3, Eq. 2.116

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

$$\mathbf{m}_N = \mathbf{S}_N^{-1} (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t})$$

Predictive Distribution

- Observed dataset with inputs $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ and targets $\mathbf{t} = (t_1, \dots, t_N)^T$
- Gaussian Posterior** distribution (from Gaussian prior and Gaussian likelihood)

$$\rightarrow p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \alpha, \beta) = \underline{\mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)}$$

with $\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

- Parametrized Gaussian **predictive distribution**:

$$\rightarrow p(t' | \mathbf{x}', \mathbf{w}, \beta) = \underline{\mathcal{N}(t' | \phi(\mathbf{x}')^T \mathbf{w}, \beta^{-1})}$$

- Analytic Gaussian **Bayesian predictive distribution** for new input

$$\rightarrow p(t' | \mathbf{x}', \mathbf{X}, \mathbf{t}, \alpha, \beta) = \int \underline{\mathcal{N}(t' | \phi(\mathbf{x}')^T \mathbf{w}, \beta^{-1})} \underline{\mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)} d\mathbf{w}$$

$$= \mathcal{N}(t' | \mathbf{x}', \phi(\mathbf{x}')^T \mathbf{m}_N, \sigma_N^2(\mathbf{x}'))$$

- With $\sigma_N^2(\mathbf{x}') = \beta^{-1} + \phi(\mathbf{x}')^T \mathbf{S}_N \phi(\mathbf{x}')$

$p(t', \underline{\mathbf{w}} | \mathbf{x}', \mathbf{X}, \mathbf{t}, \alpha, \beta)$
 marginalization
 "weighted sum
 of pred.-distr."

Bishop Eq.
2.115

Integral / Convolution
 with two Gaussian
 is again a Gaussian

Predictive Distribution

- ▶ Datasets:
 - ▶ $t = \sin(2\pi x) + \epsilon$
 - ▶ $\epsilon \sim \mathcal{N}(0, \beta^{-1})$
- ▶ Dataset sizes:
 - ▶ $N = 1, 2, 4, 25$
- ▶ Model:
 - ▶ $y(x, \mathbf{w}) = \boldsymbol{\phi}(x)^T \mathbf{w}$
 - ▶ $\boldsymbol{\phi}_j(x)$: Gaussian basis functions

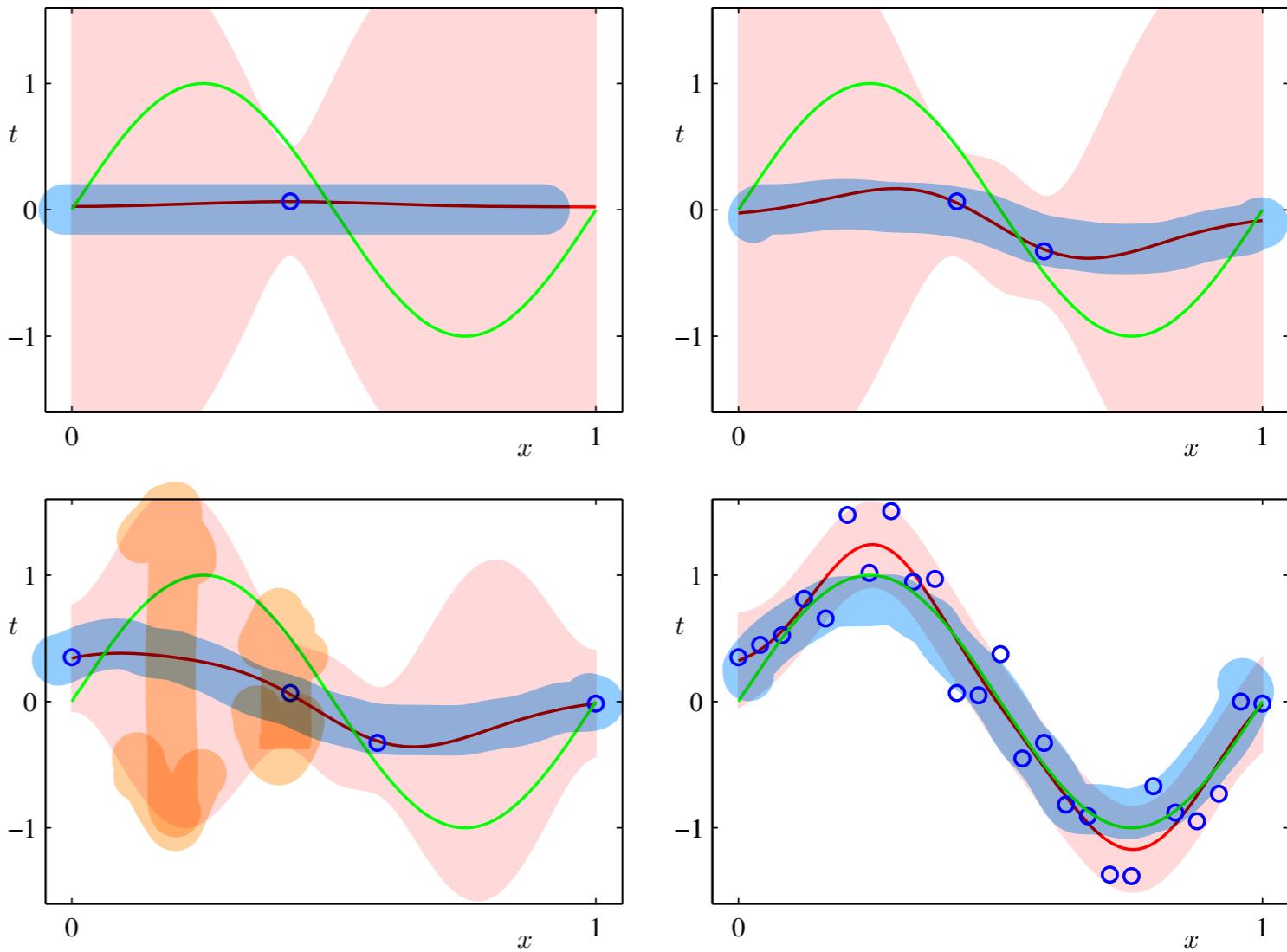
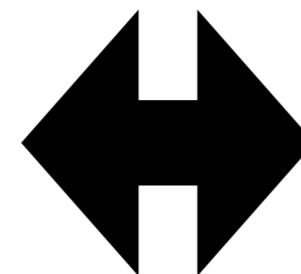
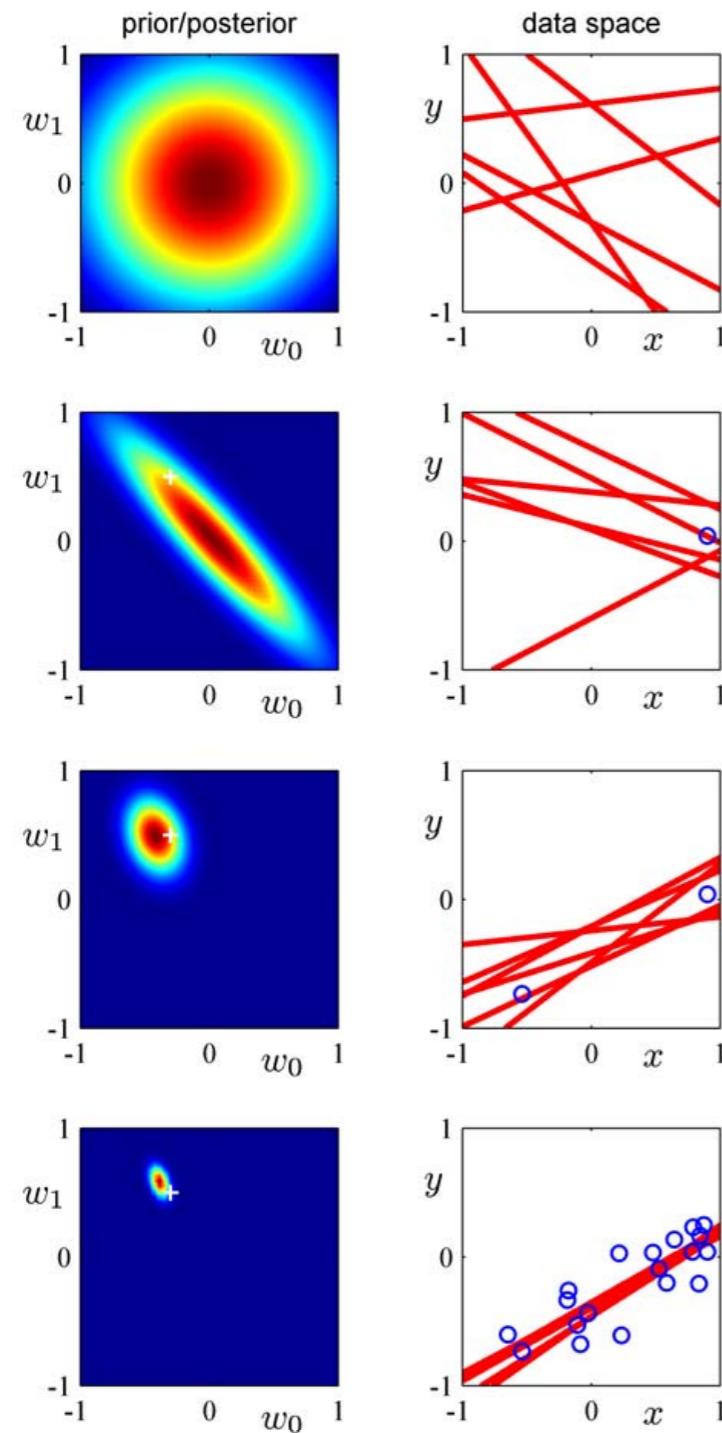


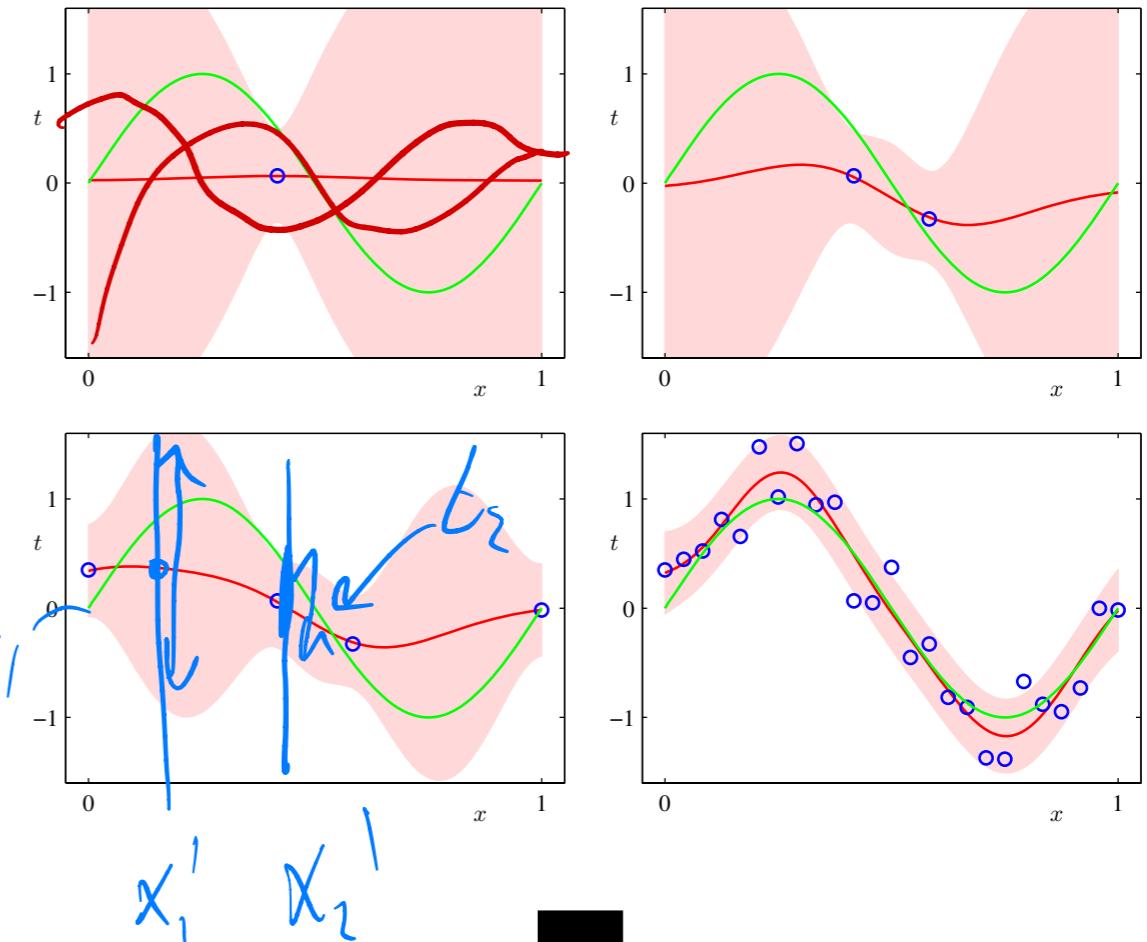
Figure: Predictive distribution (Bishop 3.8)

- ▶ **Bayesian predictive distribution:**
 - ▶ $p(t' | x', \mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t' | \mathbf{x}'^T \boldsymbol{\phi}(\mathbf{x}') \mathbf{m}_N, \sigma_N^2(\mathbf{x}'))$
- ▶ $\sigma_N^2(\mathbf{x}') = \beta^{-1} + \boldsymbol{\phi}(\mathbf{x}')^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}')$, $\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$, $\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$

Bayesian Linear Regression



Continuous Random Functions / Bayesian predictive distributions / Gaussian Processes



$cov[t_1', t_2' | x_1, x_2']$

Kernel method!

(Predictions based on
proximity to training data)

Equivalent Kernel Formulation

- ▶ Bayesian predictive distribution

$$\begin{aligned} p(t'|x', \mathbf{X}, \mathbf{t}, \alpha, \beta) &= \int p(t'|x', \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t' | \mathbf{m}_N^T \boldsymbol{\phi}(x'), \sigma_N^2(x')) \end{aligned}$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} \quad \sigma_N^2(x') = \frac{1}{\beta} + \boldsymbol{\phi}(x')^T \mathbf{S}_N \boldsymbol{\phi}(x') \quad \mathbf{S}_N^{-1} = \alpha \mathbf{1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

- ▶ predictive mean:

$$y(x', \mathbf{m}_N) = \boldsymbol{\phi}(x')^T \mathbf{m}_N = \beta \boldsymbol{\phi}(x')^T \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} = \beta \boldsymbol{\phi}(x')^T \mathbf{S}_N \sum_{n=1}^N (\boldsymbol{\Phi}^T)_{:,n} t_n$$

$$= \sum_{n=1}^N \beta \boldsymbol{\phi}(x')^T \mathbf{S}_N \boldsymbol{\phi}(x_n) t_n$$

$$= \sum_{n=1}^N k(x', x_n) t_n$$

parametric "primed" form
chosen parameters

non-parametric form

"dual" form

Equivalent kernel for Gaussian Basis Functions

Figure: Equivalent kernel $k(x', x)$ (Bishop 3.10)

- Localized kernel

$$k(x', x) = \beta \phi(x')^T \mathbf{S}_N \phi(x)$$

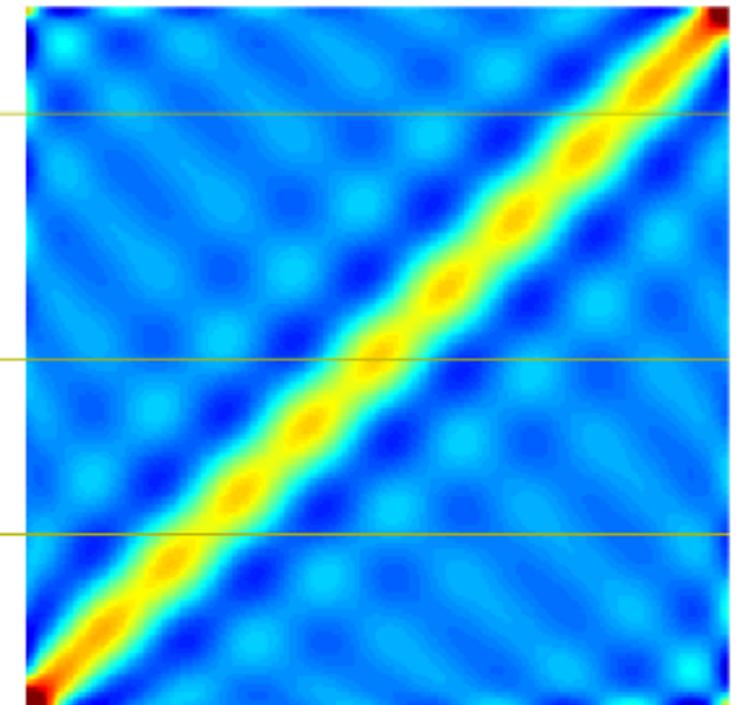
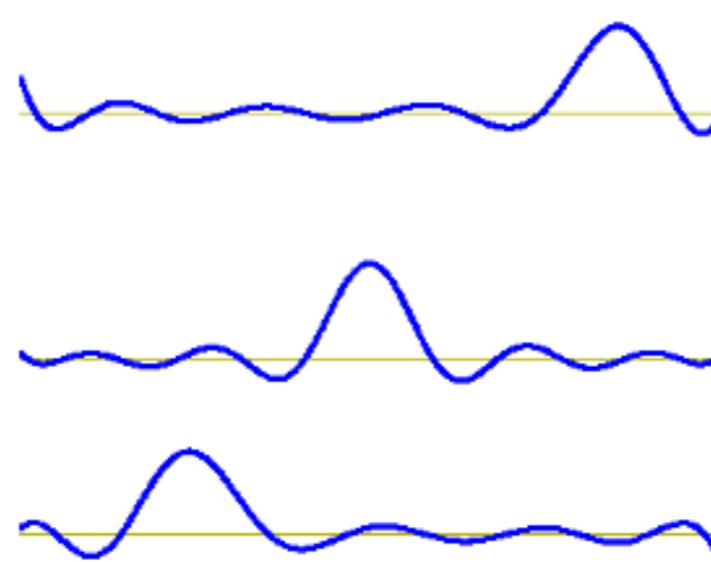
- predictive mean

$$y(x', \mathbf{m}_N) = \sum_{n=1}^N k(x', x_n) t_n$$

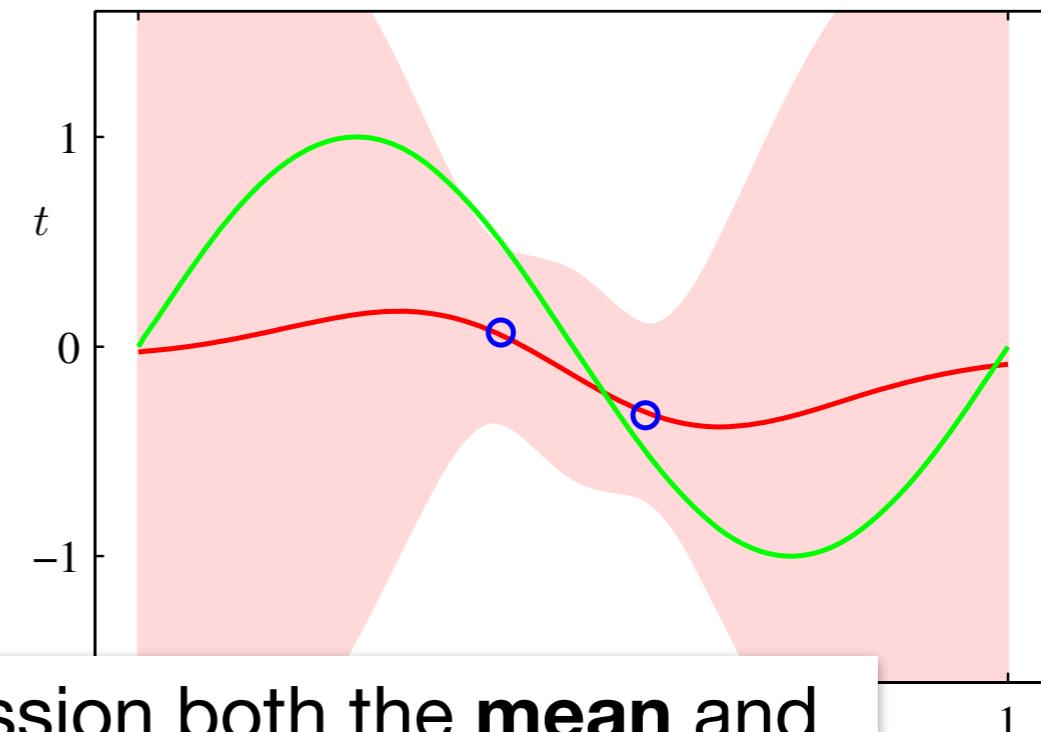
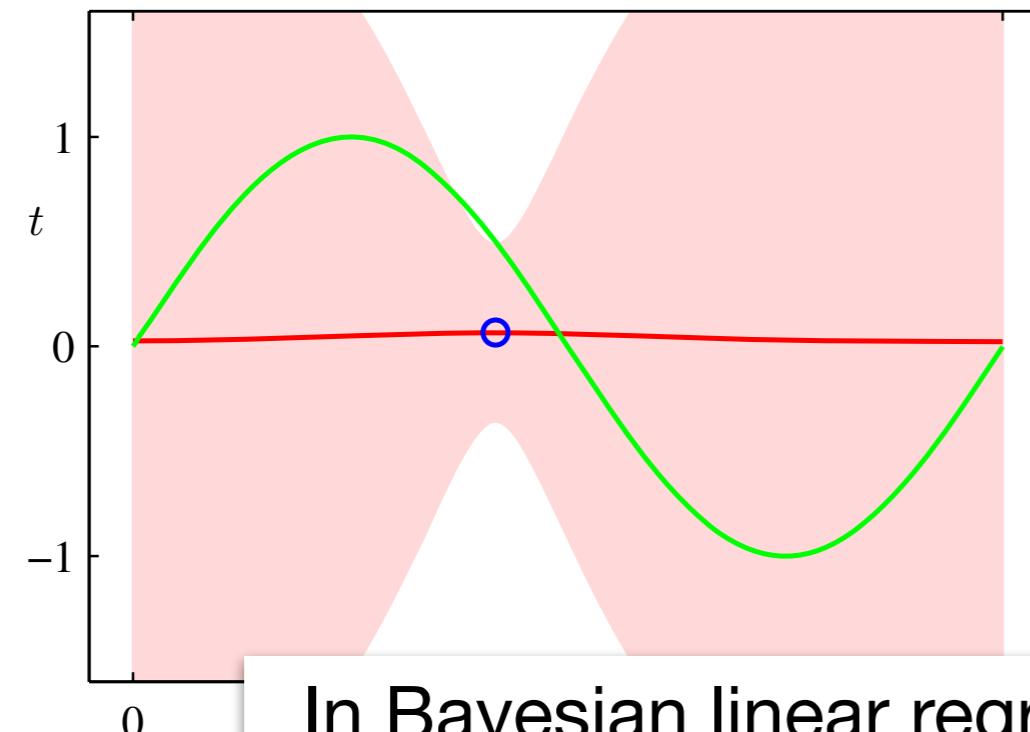
- Training points x_n close to x' contribute more!

- Covariance of between predictions:

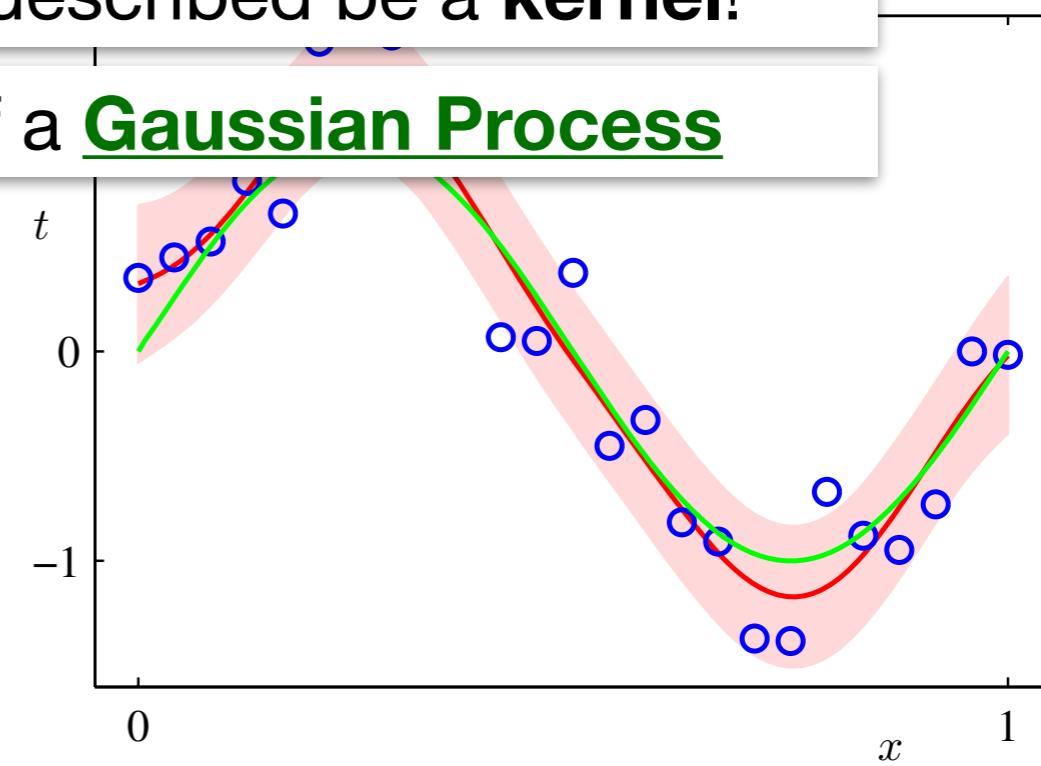
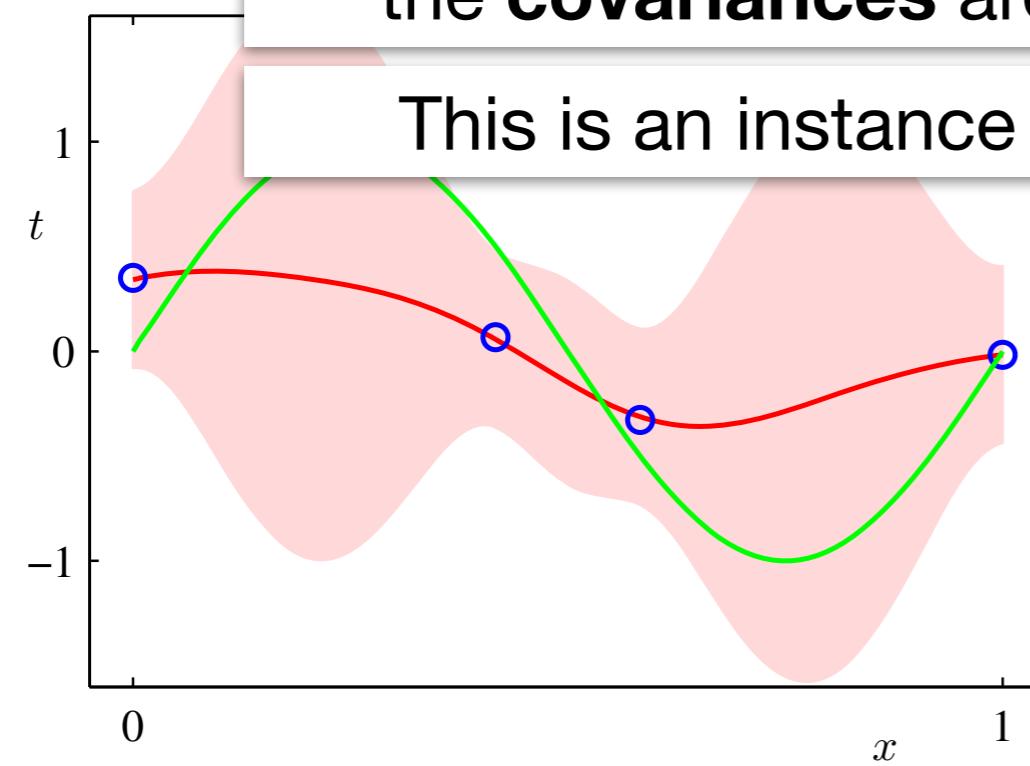
$$\begin{aligned} \text{cov}[t_1, t_2 | x_1, x_2] &= \text{cov}_{\mathbf{w}}[y(x_1, \mathbf{w}), y(x_2, \mathbf{w})] = \text{cov}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}, \mathbf{w}^T \phi(x_2)] \\ &= \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w} \mathbf{w}^T \phi(x_2)] - \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}] \mathbb{E}_{\mathbf{w}}[\mathbf{w}^T \phi(x_2)] \\ &\text{Kernel prediction at two locations} \\ &= \phi(x_1)^T (\mathbb{E}_{\mathbf{w}}[\mathbf{w} \mathbf{w}^T] - \mathbb{E}_{\mathbf{w}}[\mathbf{w}] \mathbb{E}_{\mathbf{w}}[\mathbf{w}^T]) \phi(x_2) = \phi(x_1)^T \text{cov}[\mathbf{w}, \mathbf{w}] \phi(x_2) = \phi(x_1)^T \mathbf{S}_N \phi(x_2) \\ &\quad \beta^{-1} k(x_1, x_2) \end{aligned}$$



Revisit Bayesian linear regression



In Bayesian linear regression both the **mean** and the **covariances** are described by a **kernel**!



This is an instance of a **Gaussian Process**

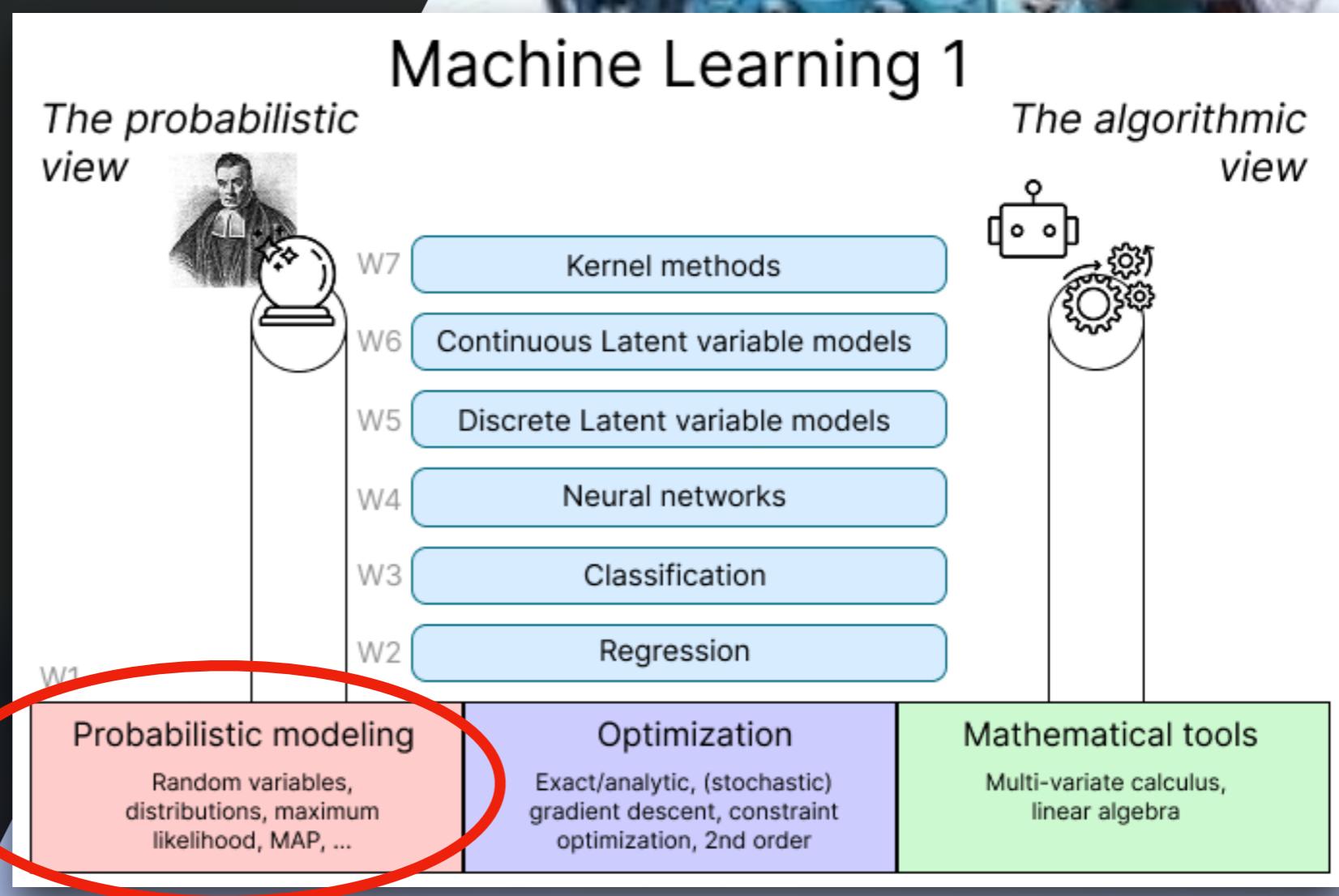
Machine Learning 1

Lecture 12.1 - Kernel Methods

Gaussian Processes - Properties of Gaussian
Random Variables

Erik Bekkers

(Bishop 2.3.1, 2.3.2)

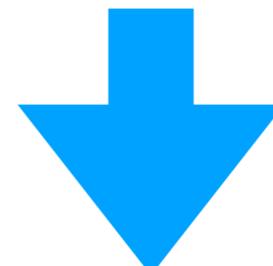


Gaussian Processes are Gaussian distributions for continuous functions

Instead of sampling (finite dimensional) vectors

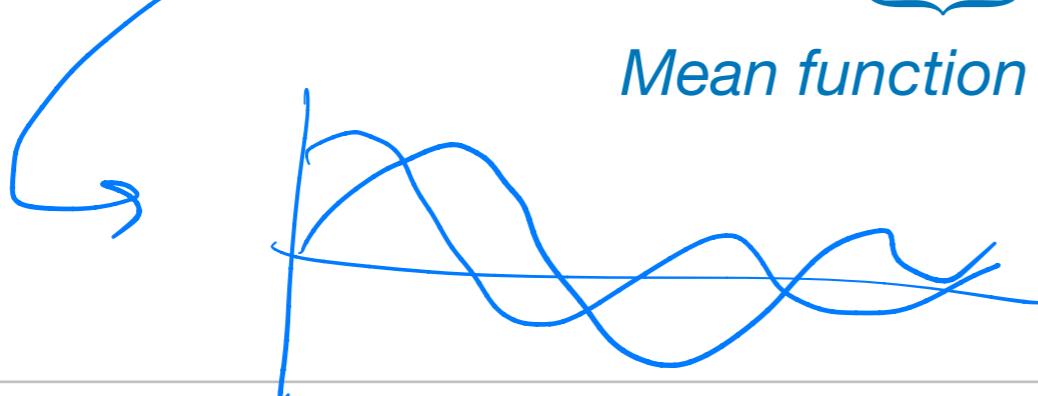
$$\mathbf{x} = \begin{pmatrix} i \\ j \end{pmatrix} \quad \mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$x = i$$



With GPs we sample (infinite dimensional) functions

$$f(\cdot) \sim GP(\underbrace{\boldsymbol{\mu}(\cdot)}_{\text{Mean function}}, \underbrace{k(\cdot, \cdot)}_{\text{Kernel}})$$



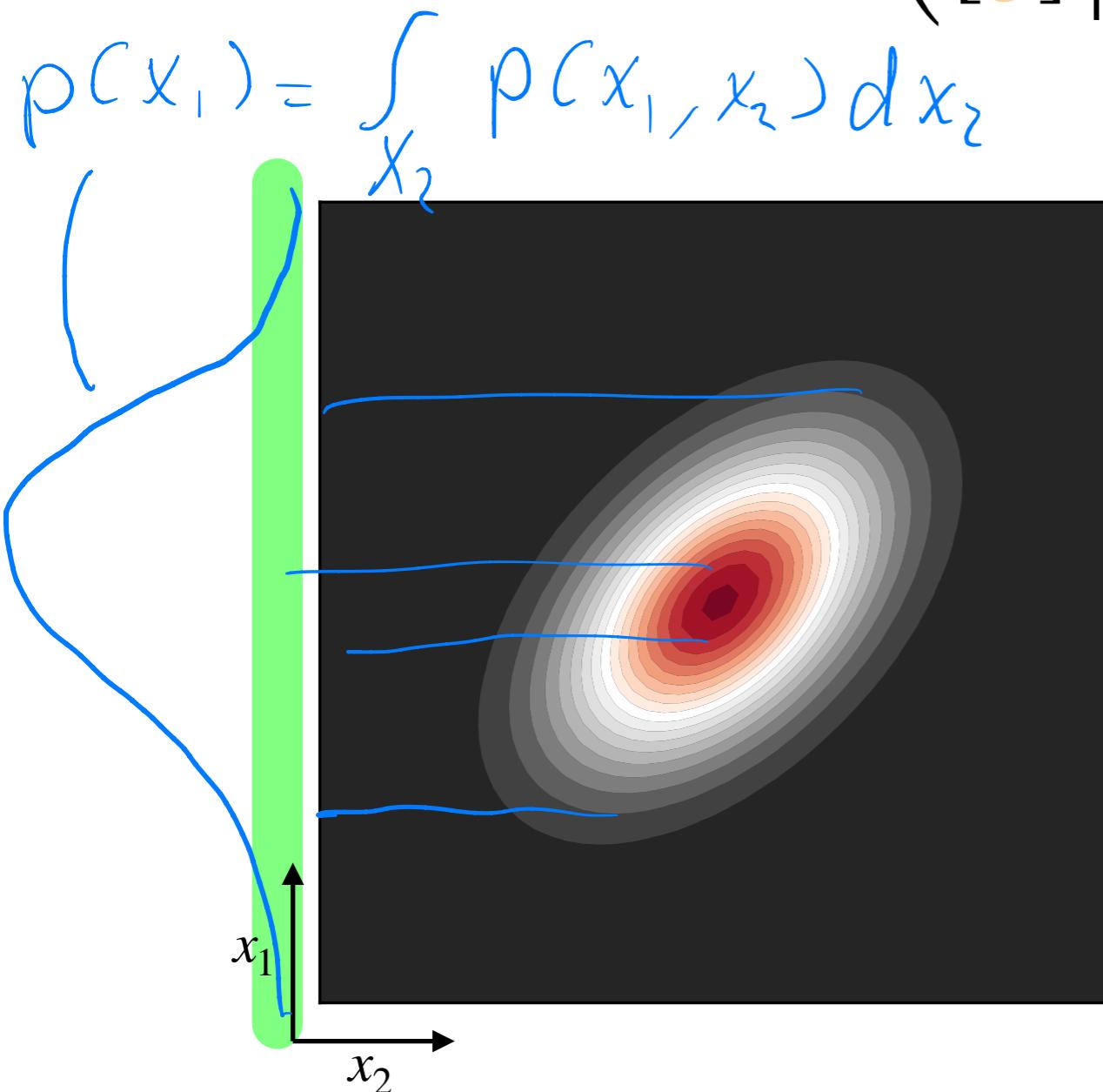
Gaussian Processes are Gaussian distributions for continuous functions

*In order to **sample from GPs** and in order to **fit them** we need to revisit some **basic properties of multi-variate Gaussians***

Gaussians: Marginalization property

- Take two random variables x_1 and x_2 , that are jointly Gaussian distributed:

$$p(x_1, x_2) = \mathcal{N} \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$



- Then the marginals are given by

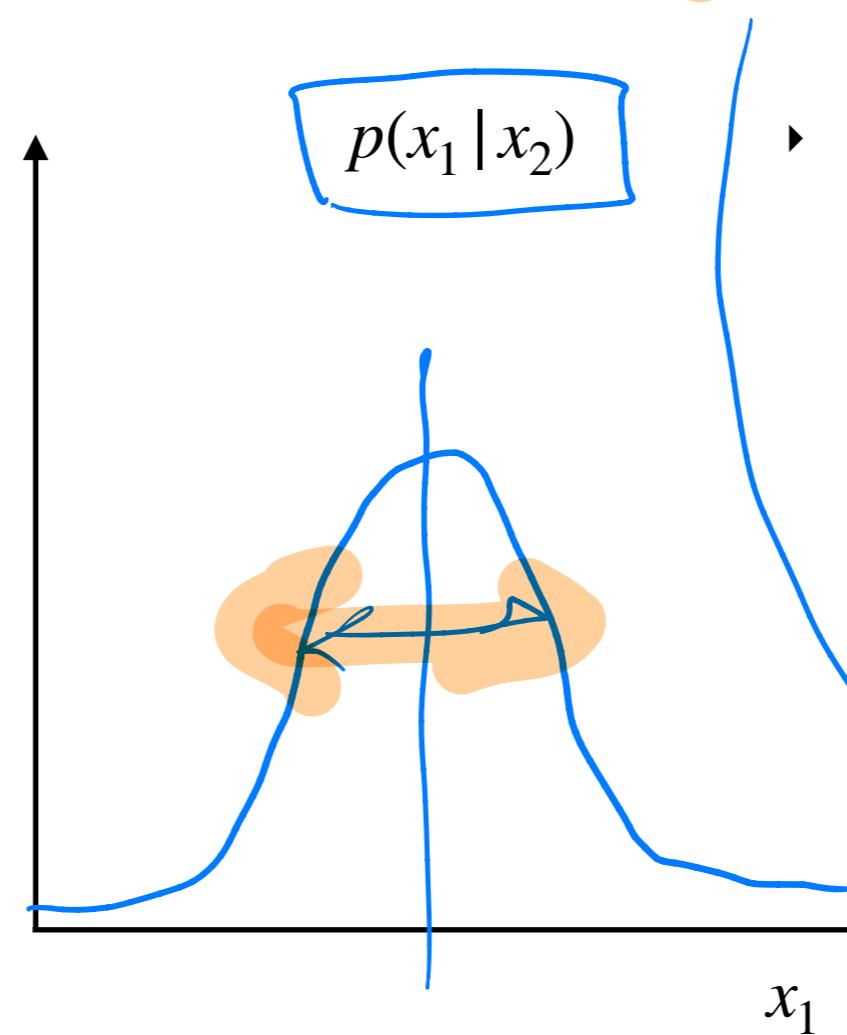
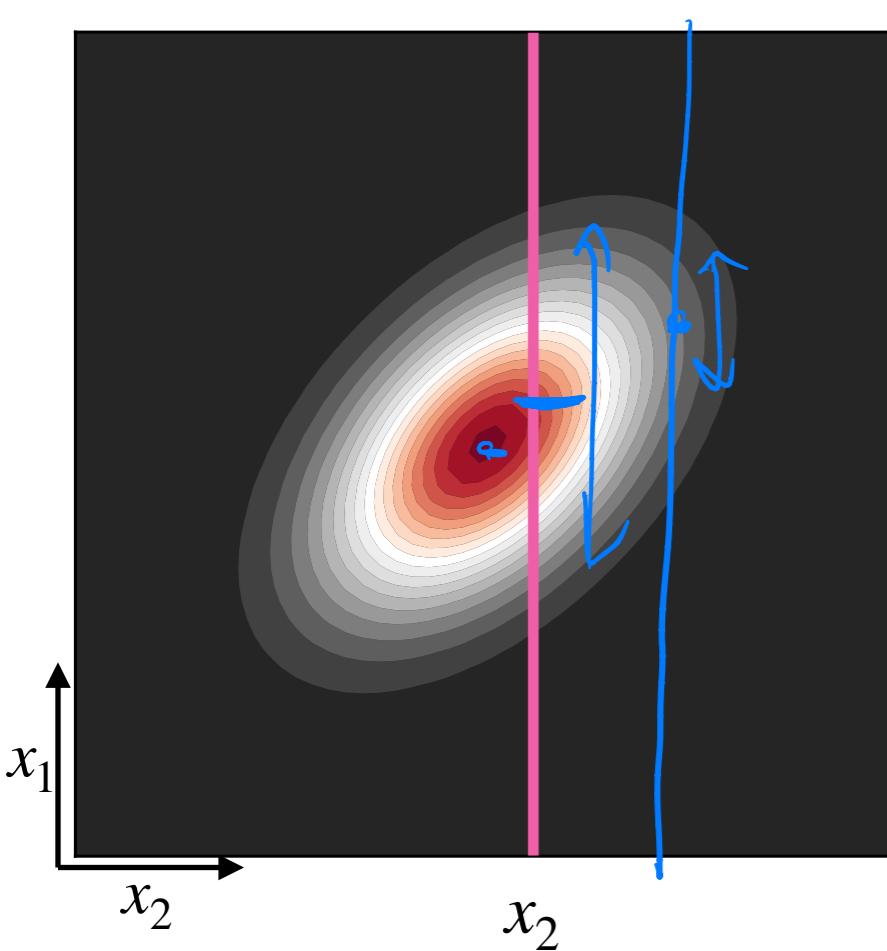
$$p(x_1) = \mathcal{N}(x_1 | \mu_1, \Sigma_{11})$$

$$p(x_2) = \mathcal{N}(x_2 | \mu_2, \Sigma_{22})$$

Gaussians: Conditioning Property

- Take two random variables x_1 and x_2 , that are jointly Gaussian distributed:

$$p(x_1, x_2) = \mathcal{N} \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$



- Then the conditional is:

$$p(x_1 | x_2) = \mathcal{N}(\mu_{1|2}, \Sigma_{1|2})$$

with

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Summing Random Variables

- The sum of two independent Gaussian random variables is also a Gaussian random variable:

- If

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

$$y \sim \mathcal{N}(\mu', \Sigma')$$

- Then

$$z = x + y \rightarrow$$

$$\text{Cov}[z] = \mathbb{E}[z^2] - \mathbb{E}[z]^2$$

$= \dots$



$$z \sim \mathcal{N}(\underline{\mu + \mu'}, \underline{\Sigma + \Sigma'})$$

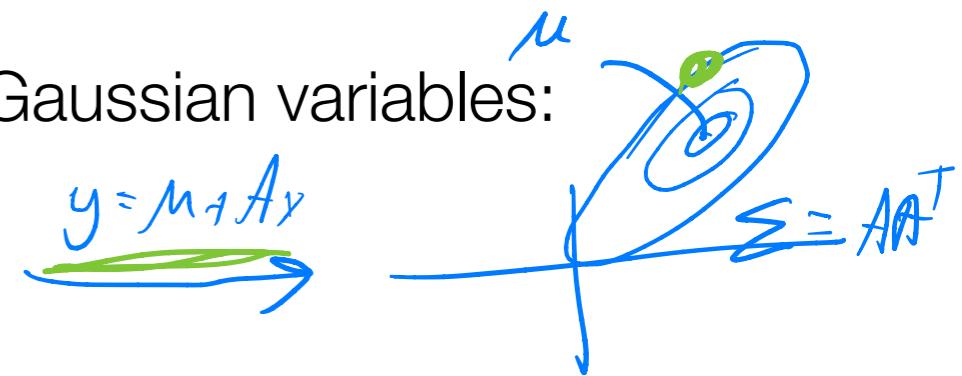
$$\mathbb{E}[z] = \mathbb{E}[x + y] = \mathbb{E}[x] + \mathbb{E}[y] = \mu + \mu'$$

Sampling Correlated Gaussian Variables

- If we have sampled a vector \mathbf{x} of uncorrelated Gaussian variables:

Recall
PPCA

$$\underline{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$



If \mathbf{y} is obtained from \mathbf{x} via $\mathbf{y} = \mu + \mathbf{Ax}$ then

$$\mathbf{y} \sim \mathcal{N}(\mu, \Sigma) \quad \text{with} \quad \Sigma = \mathbf{AA}^T$$

- **Reparametrization trick:** You want to sample correlated samples

$\mathbf{y} \sim \mathcal{N}(\mu, \Sigma)$ from a Gaussian with some covariance Σ :

1. Determine \mathbf{A} by e.g. eigendecomposition s.t. $\Sigma = \mathbf{AA}^T$
2. Simply sample from a Gaussian $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})$
3. Obtain \mathbf{y} via $\mathbf{Ax} + \mu$

- *Tips for computing \mathbf{A}*

- *For a given Σ , you can compute $\Sigma = \mathbf{AA}^T$ with a Cholesky decomposition such that \mathbf{A} is lower triangular*
- *Or you compute the eigendecomposition $\Sigma = \mathbf{U}\Lambda\mathbf{U}^T$ and take $\mathbf{A} = \mathbf{U}\Lambda^{1/2}$*

multi-variate case \mathbb{R}^2 Summary

- Consider $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$ and $p(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right)$
- Then the marginals are given by

$$p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1 \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \quad \text{and} \quad p(\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_2 \mid \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$$

- And the conditional is given by
- $$p(\mathbf{x}_1 \mid \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 \mid \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2}) \quad \text{with} \quad \boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2)$$

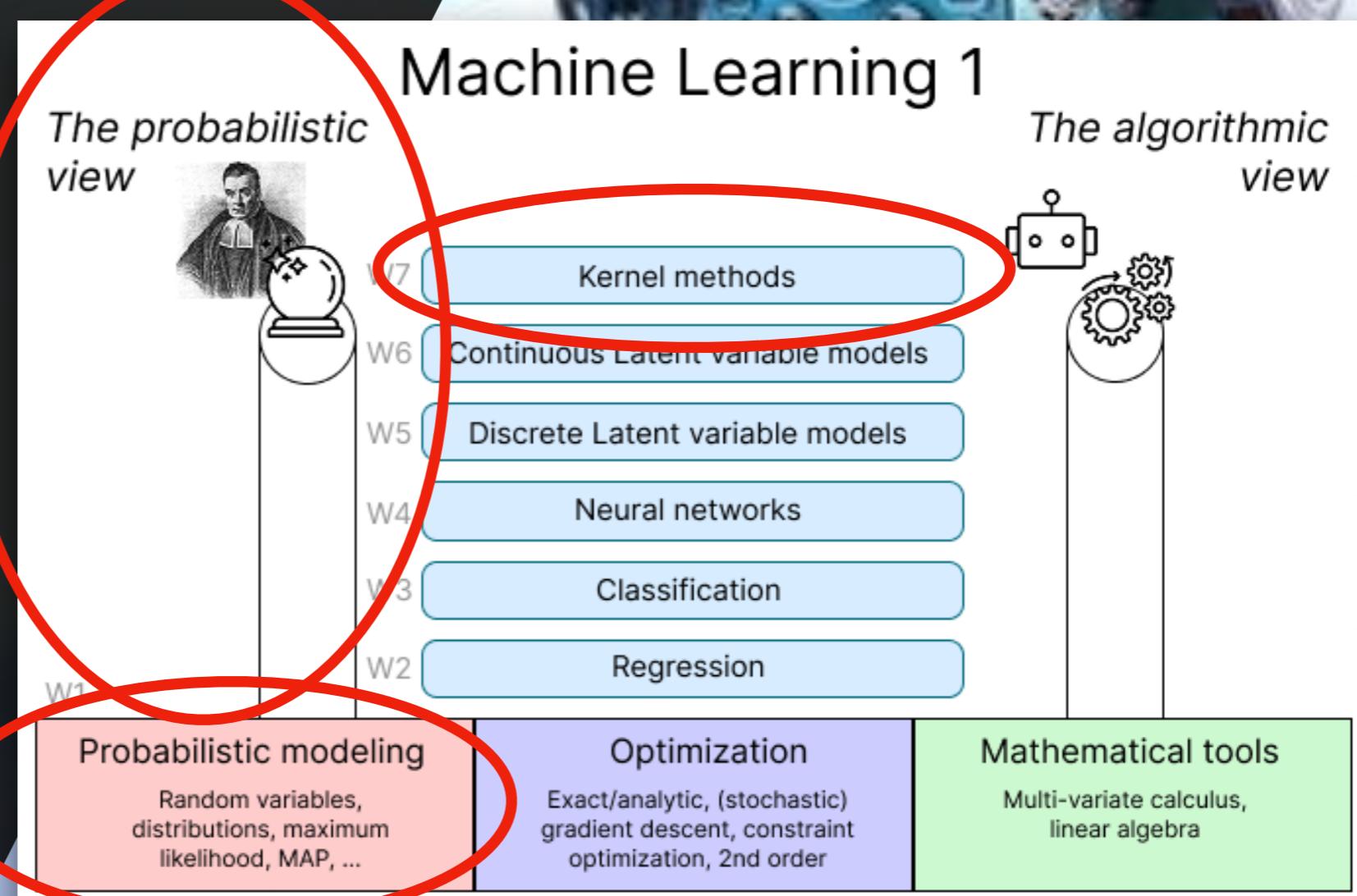
$$\boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}$$

- If \mathbf{x} is an uncorrelated Gaussian random variable (i.e., $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$)
then $\mathbf{y} = \boldsymbol{\mu} + \mathbf{A}\mathbf{x}$ is correlated $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{A}\mathbf{A}^T)$ with $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^T$

Machine Learning 1

Lecture 12.3 - Kernel Methods
Gaussian Processes - Definition

Erik Bekkers
(Bishop 6.4.1)

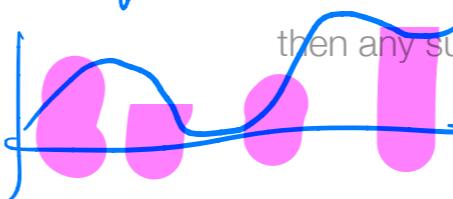


Gaussian Processes

- ▶ **Definition (Gaussian Process):**

Recall property of Gaussians: If

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} \sim N(\mu_{1-5}, \Sigma_{1-5}),$$

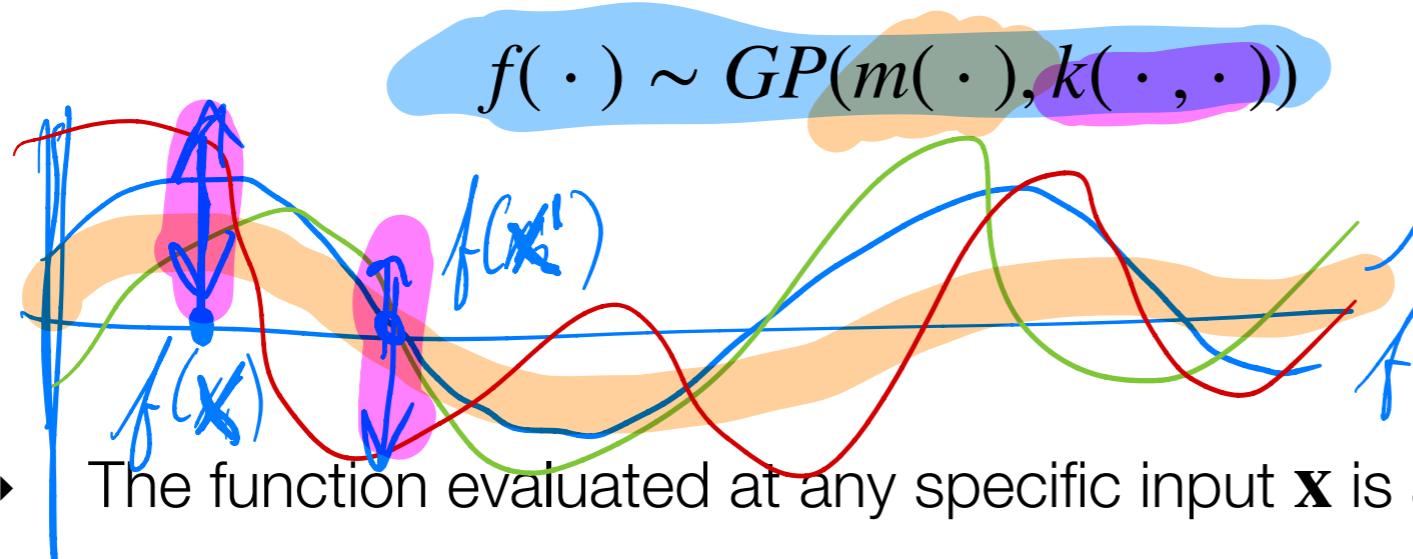


then any sub-vector is also Gaussian. E.g. $\begin{pmatrix} x_2 \\ x_4 \\ x_6 \end{pmatrix} \sim N(\mu_{2,4,6}, \Sigma_{2,4,6})$

A Gaussian process is a collection of random variables indexed with time or space, any finite number of which is jointly Gaussian distributed

essentially functions consistency property

A Gaussian process is a distribution for random functions



mean function m(·)

not x is random
but f(x) is random

- ▶ The function evaluated at any specific input \mathbf{x} is a random variable $f(\mathbf{x})$, with

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

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

Weather forecasting
temperature f as a
function of location $x \in S^2$
of location & time $x \in S^2 \times \mathbb{R}$

$$f: S^2 \times \mathbb{R} \rightarrow \mathbb{R}$$

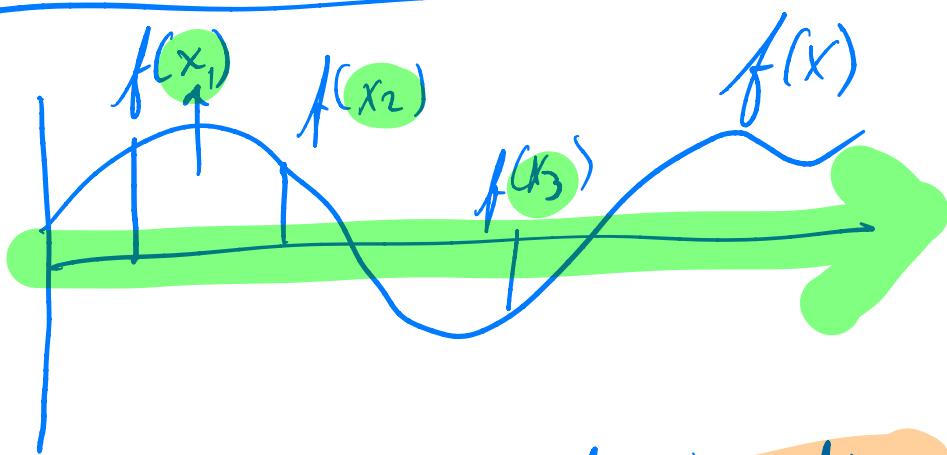
one measurement

$\xrightarrow{\text{normal dist}} f(x) \sim N(m(x), k(x, x))$

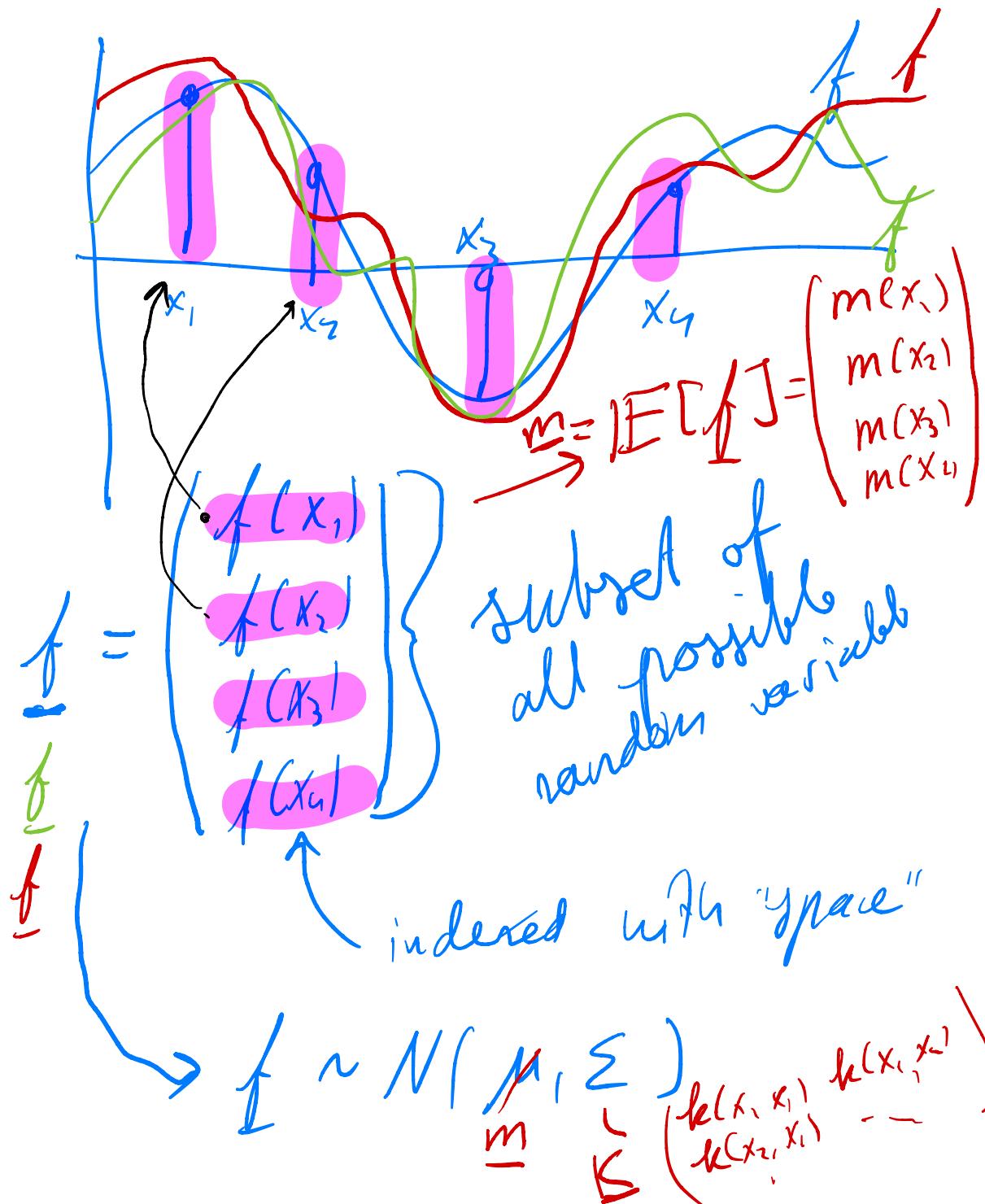
temperature measured at fixed
location x . One measurement is
random. But actually the entire
func is random

$$f(\cdot) \sim GP(m(\cdot), k(\cdot, \cdot))$$

$E = \left(\begin{array}{c} f_1 \\ f_2 \\ f_3 \\ \vdots \\ f_N \end{array} \right) \in \mathbb{R}^N$ finite dim vector space
 vector elements are indexed with integer index



functions are infinite dim. vectors. Its values are continuously indexed with indices over the real line \mathbb{R}



Gaussian Processes

- Think of functions as infinite dimensional vectors

Sampling function $f(\cdot)$

In contrast:

Sampling random weights $\underline{w} \sim p(\underline{w}) \rightarrow f(x) = \phi(x)^T \underline{w}$

This is what we did
before in the
parametric scenario.

- Sample **finite-dimensional** vectors $\mathbf{f} \in \mathbb{R}^N$ from **Gaussians** $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with mean vector $\boldsymbol{\mu} \in \mathbb{R}^N$, and covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{N \times N}$

$$\text{Recall property of Gaussians: If } \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \end{pmatrix} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \text{ then finite subset is also Gaussian. E.g. } \begin{pmatrix} f_2 \\ f_4 \\ f_5 \end{pmatrix} \sim N(\boldsymbol{\mu}_{2,4,5}, \boldsymbol{\Sigma}_{2,4,5})$$

- Sample infinite-dimensional vectors (functions) $f: \mathbb{R} \rightarrow \mathbb{R}$ from **Gaussian processes** $GP(m(\cdot), k(\cdot, \cdot))$, defined by mean function $m(\cdot)$ and covariance kernel $k(\cdot, \cdot)$

$$\text{Gaussian process: If } f(\cdot) \sim GP(m(\cdot), k(\cdot, \cdot)), \text{ then finite subset is also Gaussian. E.g. } \begin{pmatrix} f(x_1) \\ f(x_4) \\ f(x_5) \end{pmatrix} \sim N(\boldsymbol{\mu}_{2,4,5}, \boldsymbol{\Sigma}_{2,4,5})$$

Functional Viewpoint, why is this a GP?

- Take any finite set $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ with corresponding random variables $\{f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)\}$ then

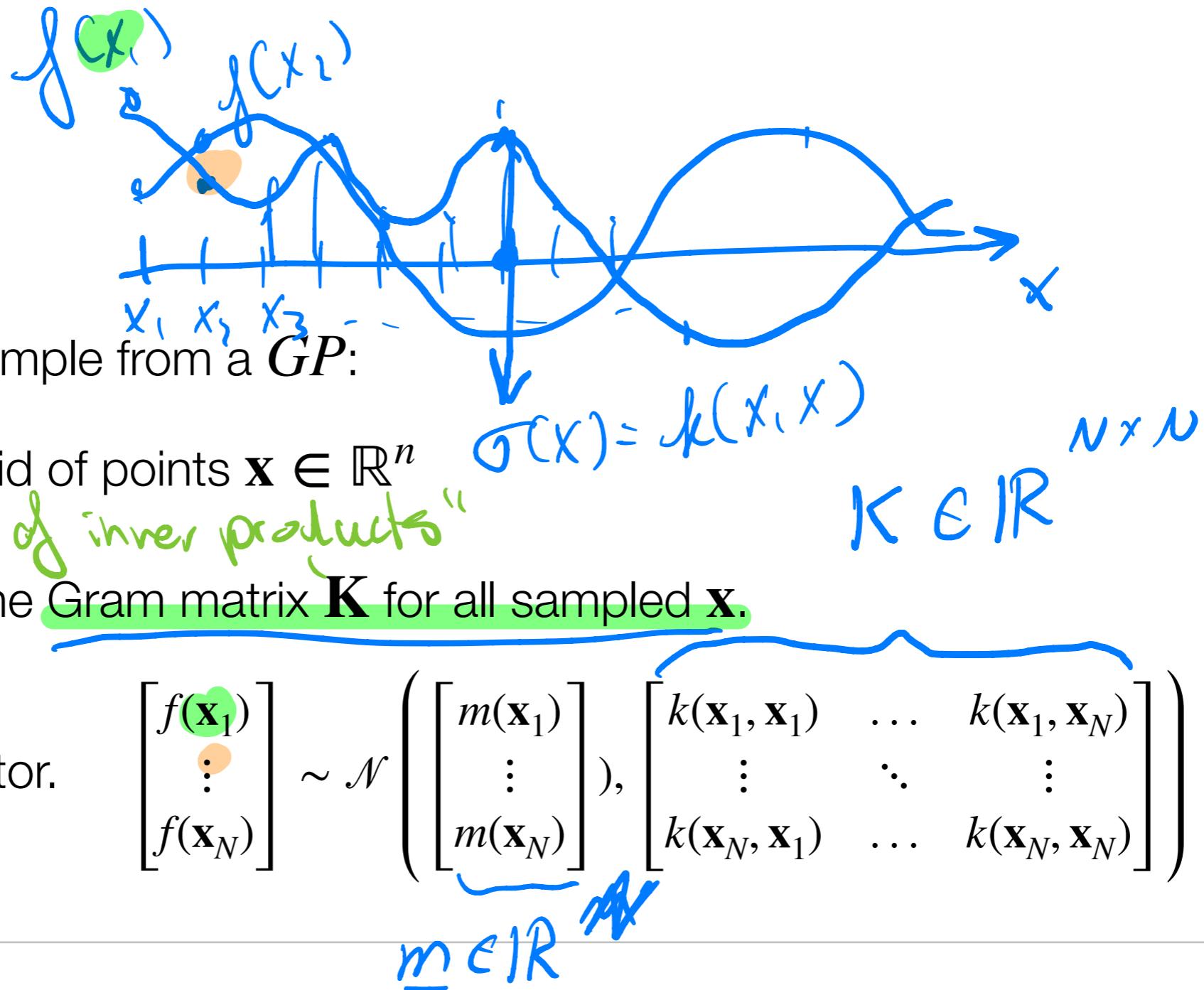
$$p\left(\begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_N) \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} m(\mathbf{x}_1) \\ \vdots \\ m(\mathbf{x}_N) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}\right)$$

- Consistency requirement:** any subset of $\{f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)\}$ should also be Gaussian distributed.
- This is true because of Gaussian marginalization property:

$$p\left(\begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix}, \begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix}\right) \rightarrow p(\mathbf{f}_1) = \mathcal{N}(\mathbf{m}_1, K_{11})$$

Functions as vectors

- Think of a function $f(\cdot)$ drawn from a GP as an extremely high-dimensional vector drawn from an extremely high-dimensional multivariate Gaussian distribution



Example: Bayesian Linear Regression

- Bayesian linear models:

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

- Prior on \mathbf{w} :

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_p)$$

- Then $f(\mathbf{x})$ is a Gaussian process

$$m(x) \simeq \mathbb{E}[f(\mathbf{x})] = \boldsymbol{\phi}(\mathbf{x})^T \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

$$\begin{aligned} k(x, x') &= \text{cov}(f(\mathbf{x}), f(\mathbf{x}')) = \mathbb{E}[f(\mathbf{x})f(\mathbf{x}')^T] = \boldsymbol{\phi}(\mathbf{x})^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \boldsymbol{\phi}(\mathbf{x}') \\ &= \boldsymbol{\phi}(\mathbf{x})^T \Sigma_p \boldsymbol{\phi}(\mathbf{x}') \end{aligned}$$

- Thus $f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)$ for any N are jointly Gaussian!

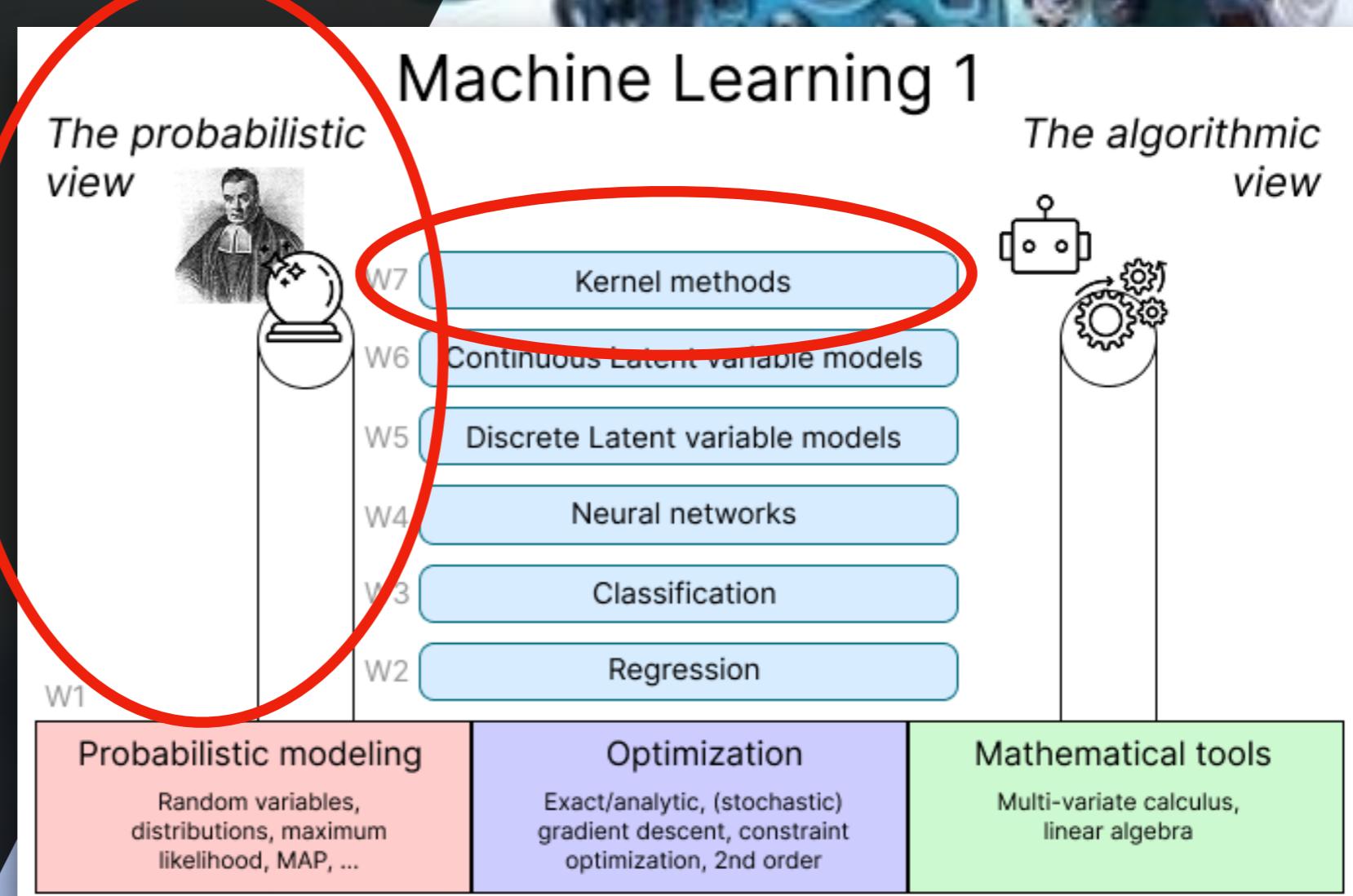
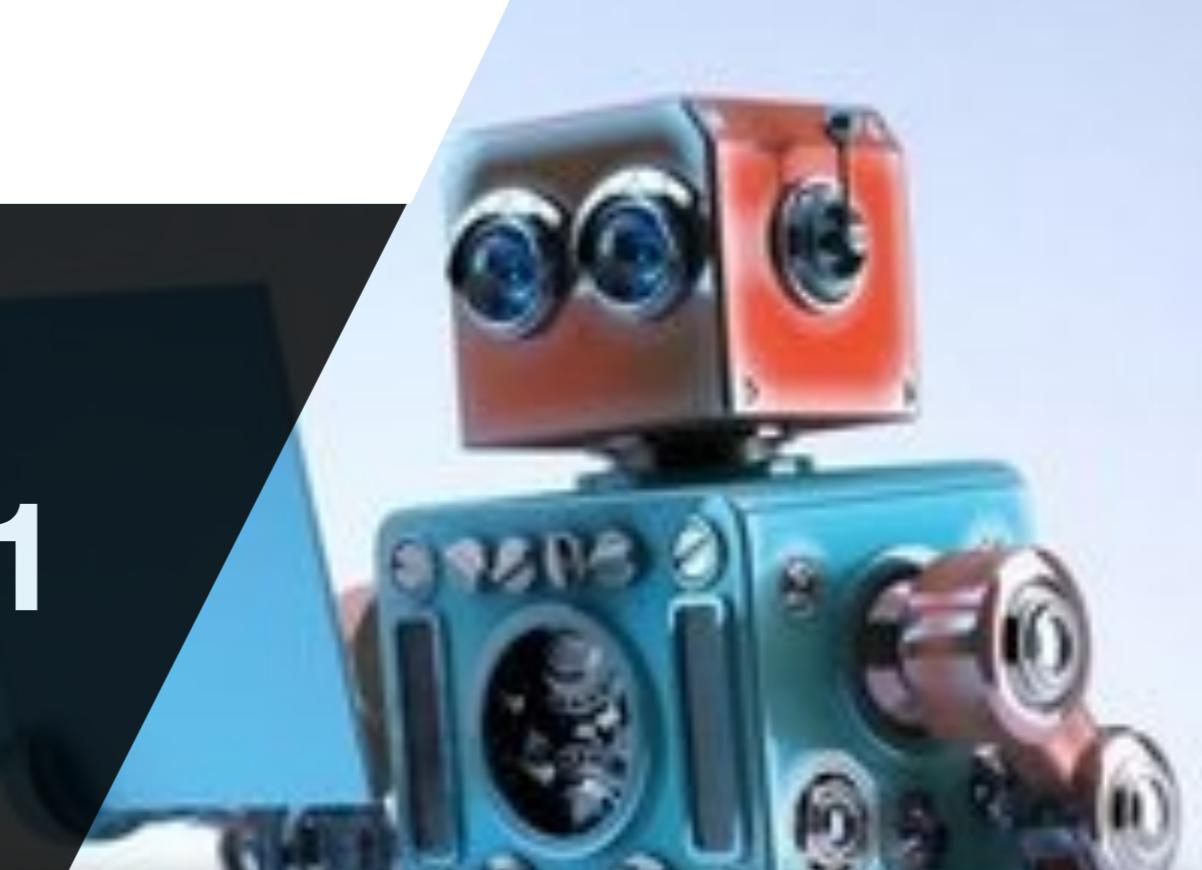
Machine Learning 1

Lecture 12.4 - Kernel Methods

Gaussian Processes - With Exponential
Kernels

Erik Bekkers

(Bishop 6.4.2)



Drawing functions from GP's

- Specifying a kernel determines the characteristics over functions drawn from the *GP*
- For simplicity, let's take

$$\mathbf{f} = \begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_N) \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \right)$$

- We consider this kernel

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

exp

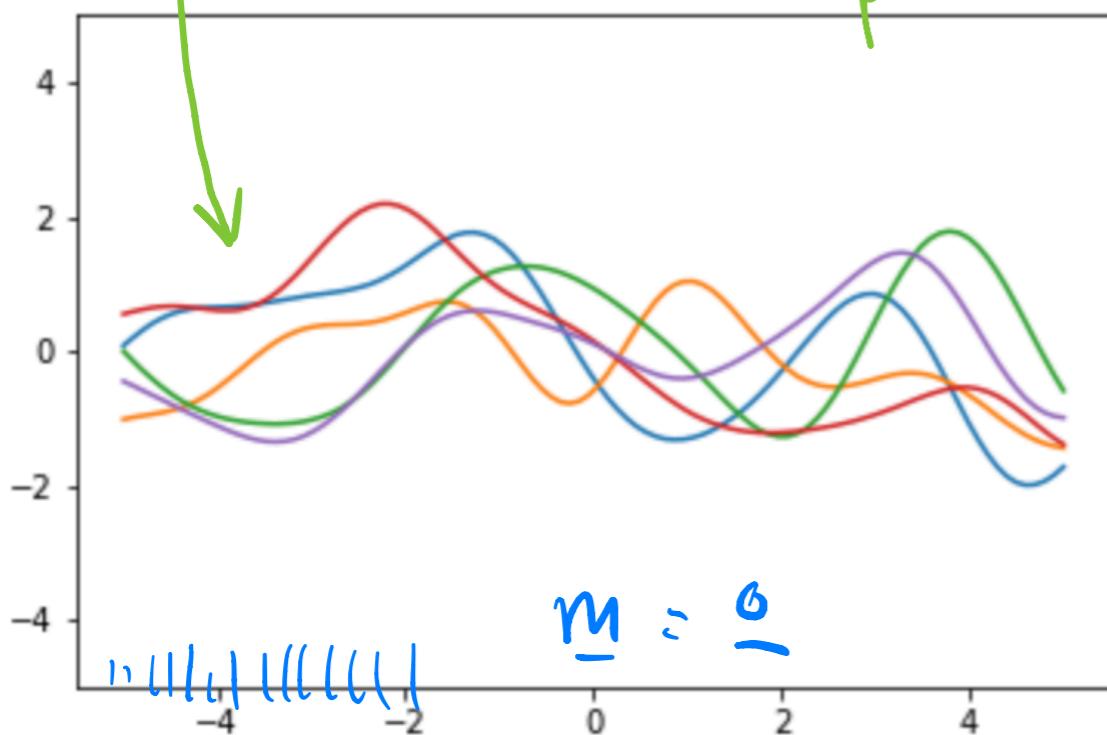
radial basis
func.

linear kernel

Drawing functions from GP's

- Sample fine grid of points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in [-5, 5]$
- Compute $\underline{\mathbf{K}}$ Gram matrix $\in \mathbb{R}^{N \times N}$
- Reparametrization trick for sampling:
 - Compute $\mathbf{K} = \mathbf{L}\mathbf{L}^T$
 - Sample random vector of size N: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_N)$
 - Sample $\mathbf{f} = \begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_N) \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$ by computing $\mathbf{f} = \mathbf{L}\mathbf{z}$

draw 5 samples and plot them



Using kernel:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

$$\theta_0 = 1$$

$$\theta_1 = 1$$

$$\theta_2 = 0$$

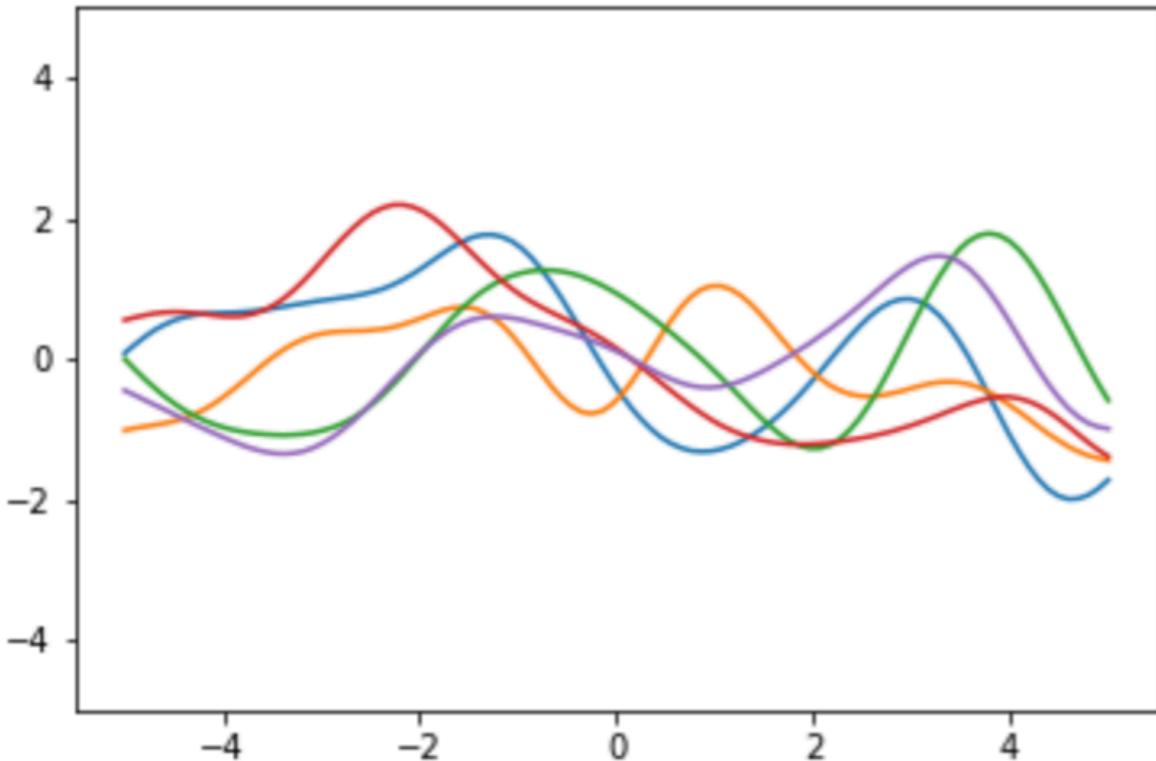
$$\theta_3 = 0$$

Varying the pre-factor θ_0

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

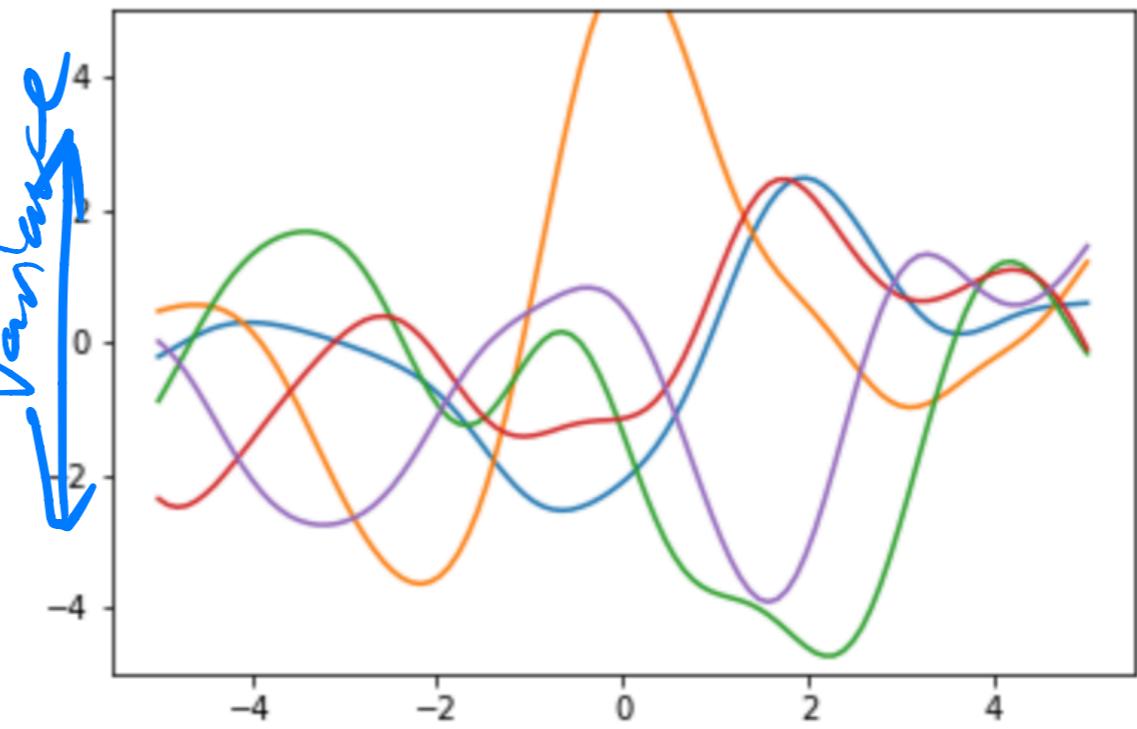
Variance regardless of $\underline{\mathbf{x}}_n, \underline{\mathbf{x}}_m$

$$\theta_0 = 1 \quad \theta_1 = 1 \quad \theta_2 = 0 \quad \theta_3 = 0$$



increases
variance

$$\theta_0 = 4 \quad \theta_1 = 1 \quad \theta_2 = 0 \quad \theta_3 = 0$$

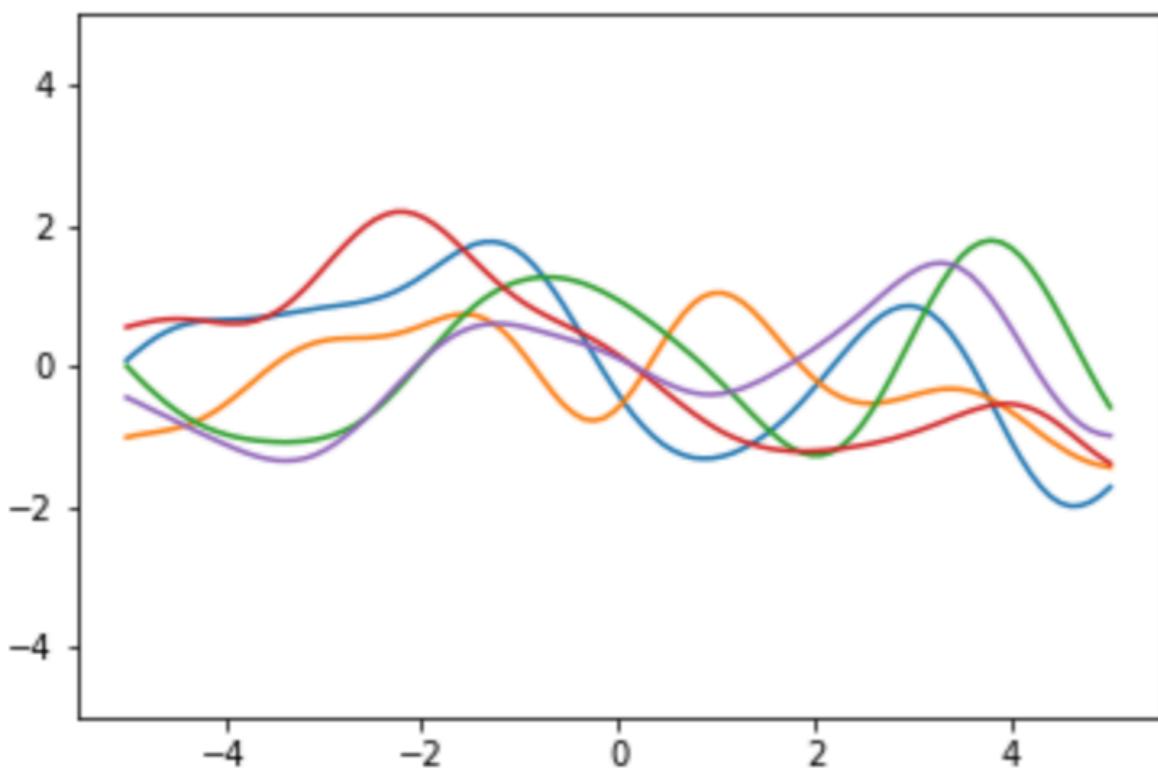


Varying the length scale θ_1

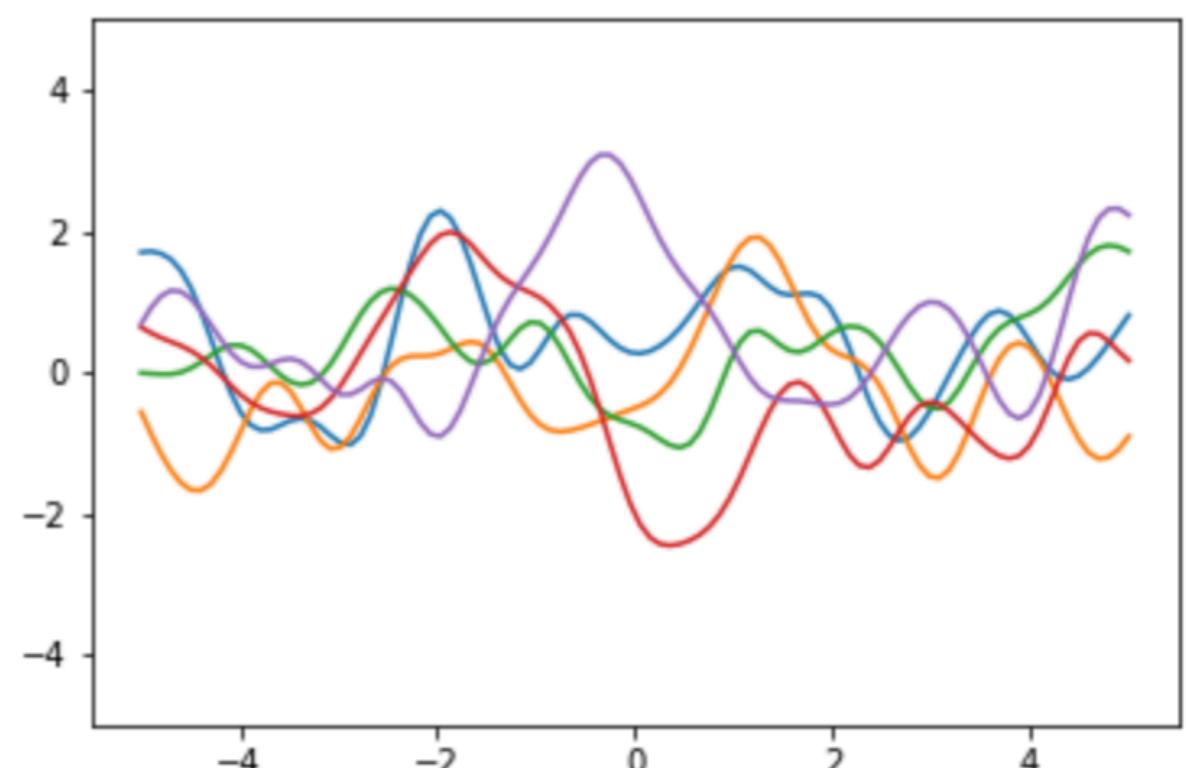
$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

parametrizes how quickly "similarity" decays with distance

$$\theta_0 = 1 \quad \theta_1 = 1 \quad \underline{\theta_2 = 0} \quad \underline{\theta_3 = 0}$$

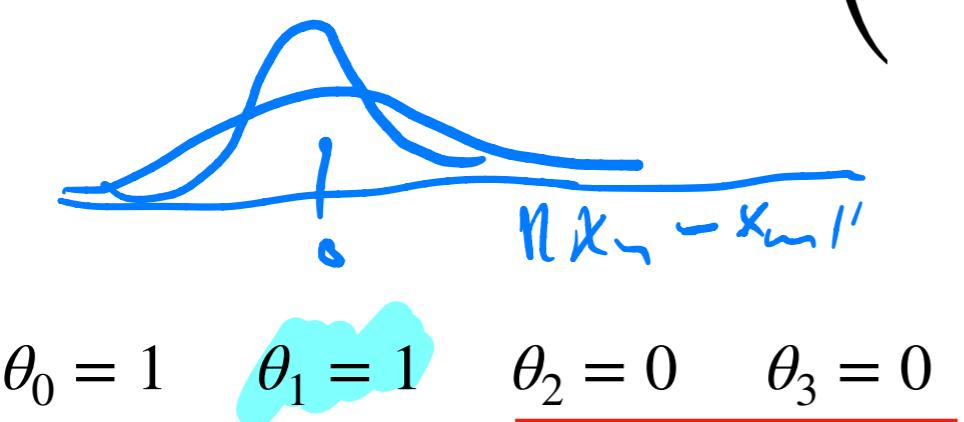


$$\theta_0 = 1 \quad \theta_1 = 0.5 \quad \underline{\theta_2 = 0} \quad \underline{\theta_3 = 0}$$



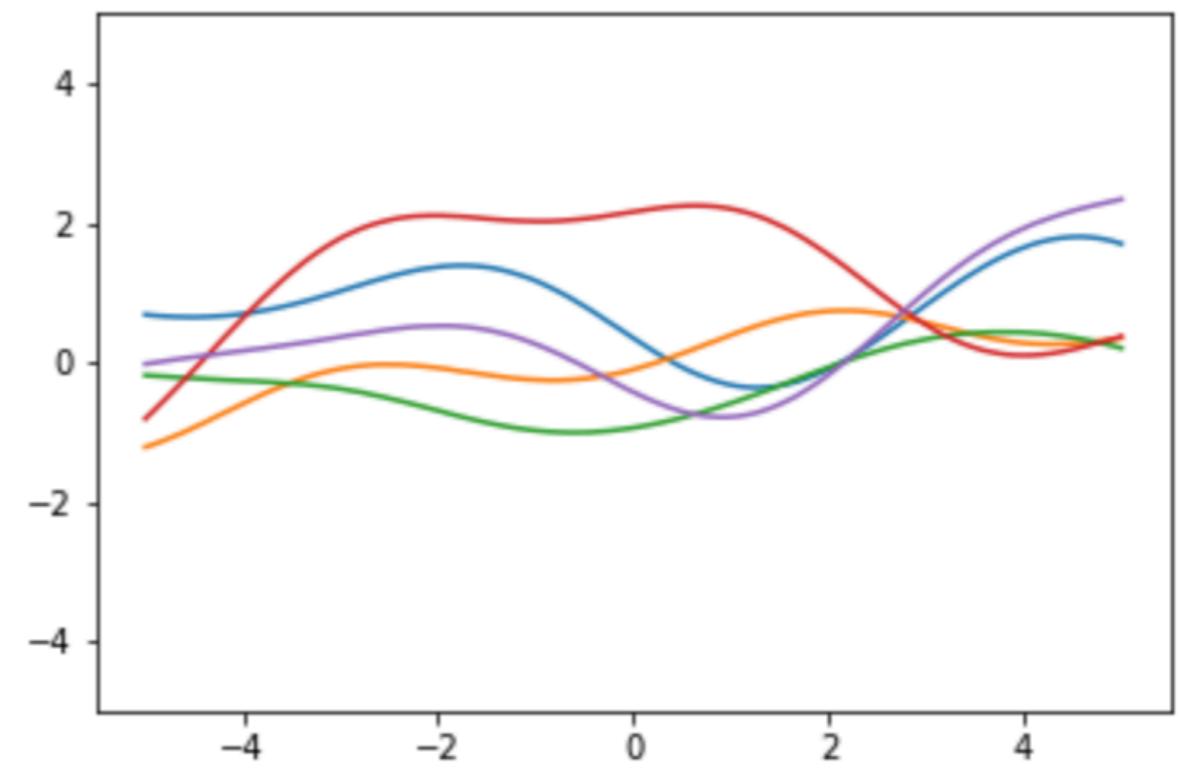
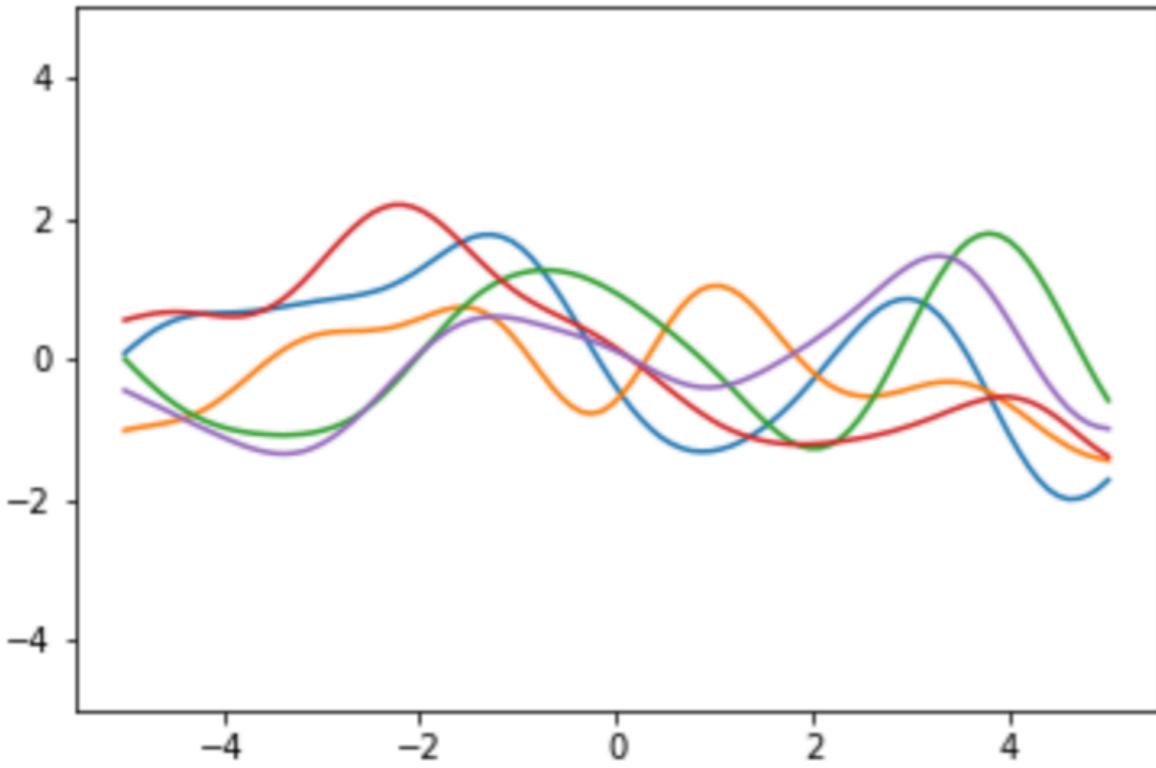
Varying the length scale θ_1

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$



large θ_2 slow decay
and thus even distant points
vary together

$\theta_0 = 1 \quad \theta_1 = 2 \quad \theta_2 = 0 \quad \theta_3 = 0$

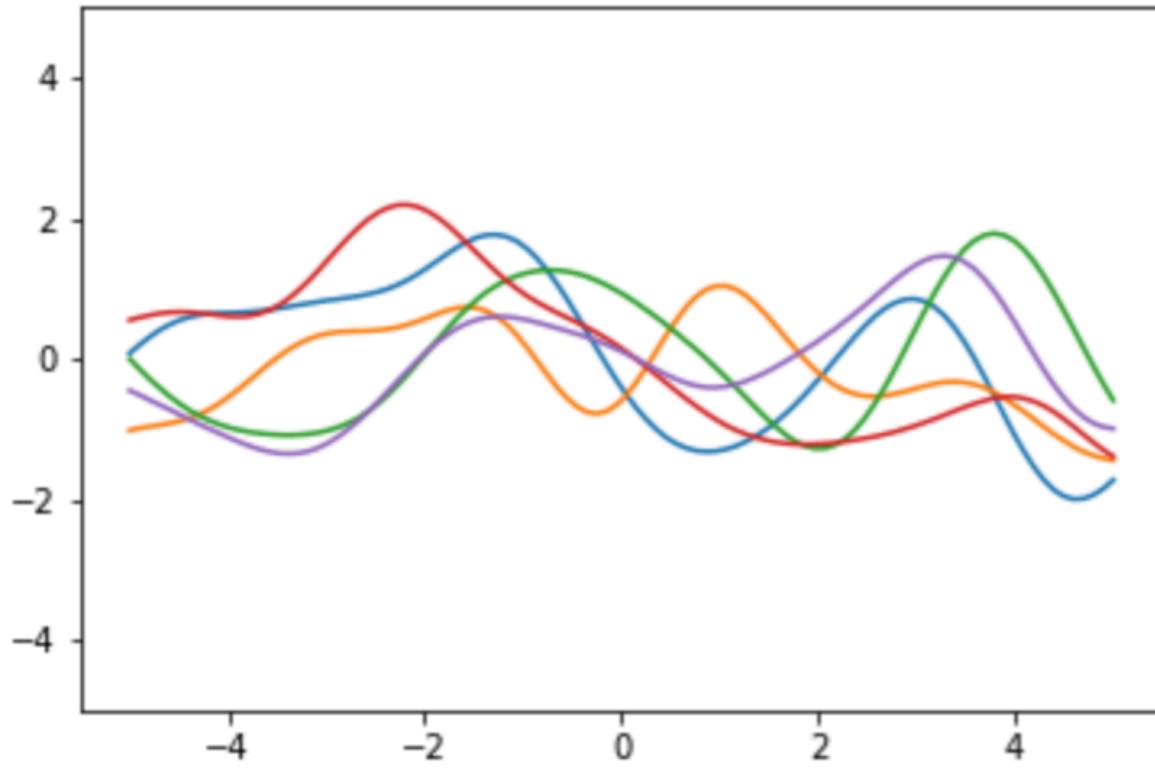


Varying the offset θ_2

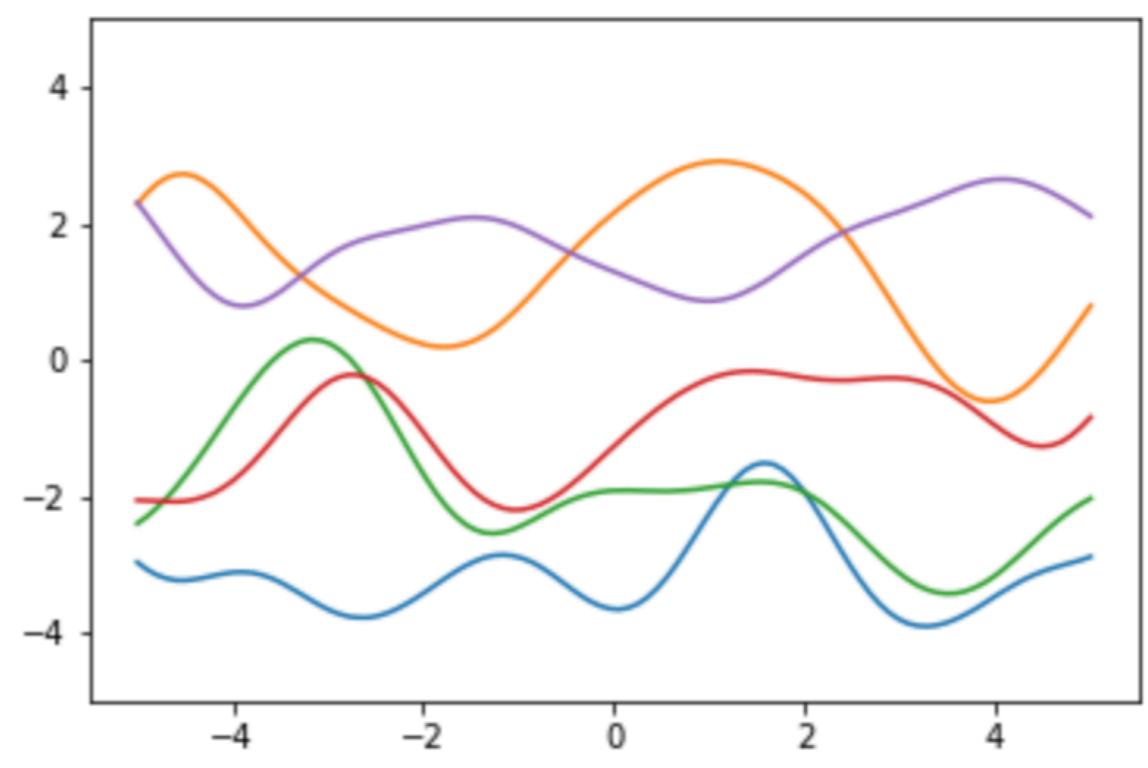
Pos. After independent
Correlation
heat std

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

$$\theta_0 = 1 \quad \theta_1 = 1 \quad \theta_2 = 0 \quad \theta_3 = 0$$



$$\theta_0 = 1 \quad \theta_1 = 1 \quad \theta_2 = 5 \quad \theta_3 = 0$$

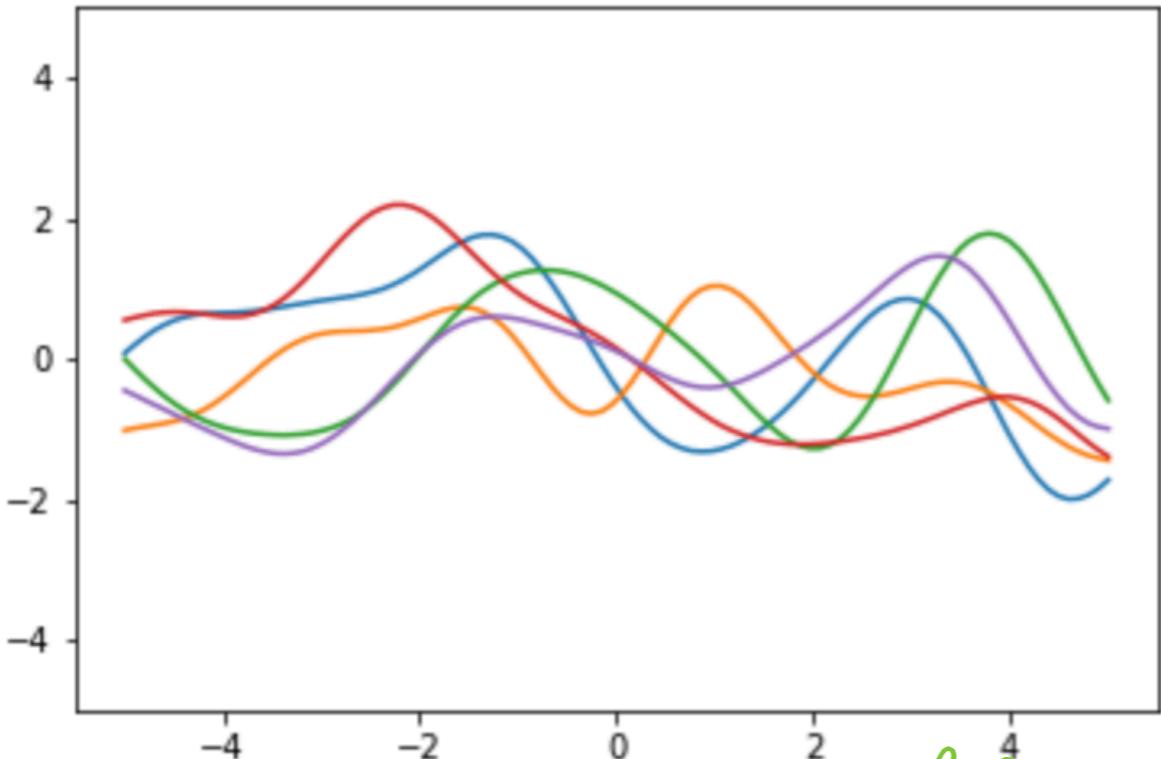


Varying the offset θ_2

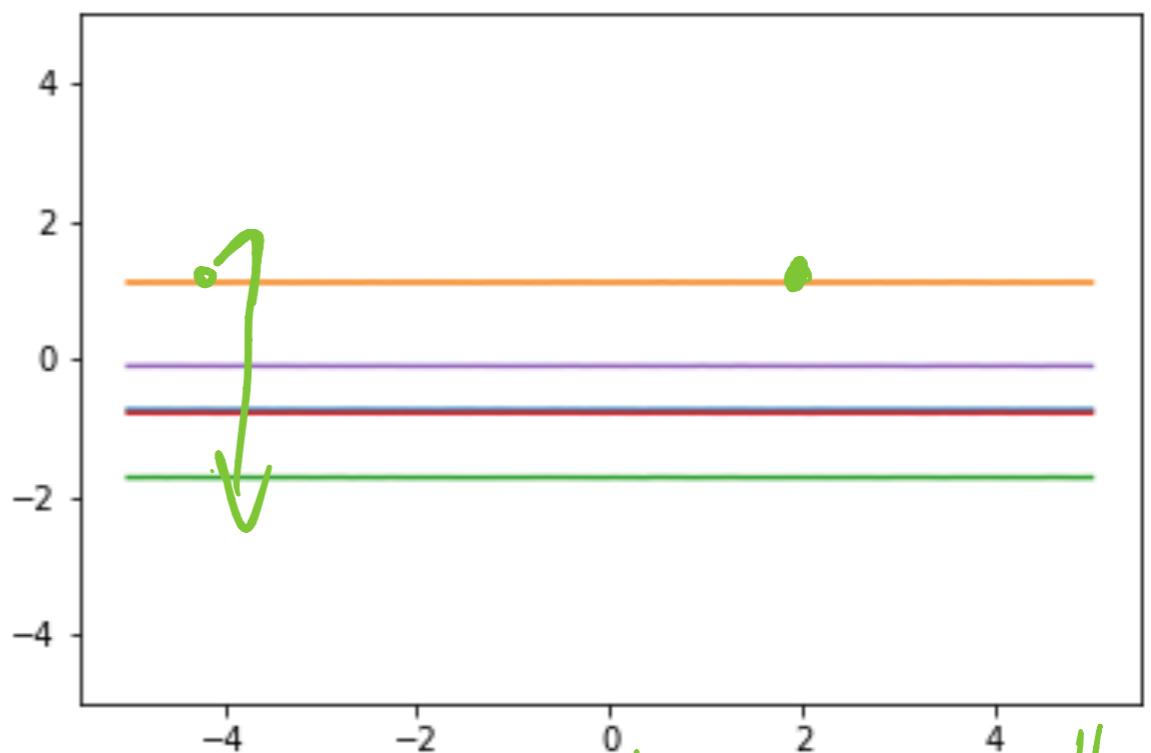
$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

perfect correlation
as
↓

$$\theta_0 = 1 \quad \theta_1 = 1 \quad \underline{\theta_2 = 0} \quad \theta_3 = 0$$



$$\theta_0 = 0 \quad \theta_1 = 0 \quad \theta_2 = 5 \quad \theta_3 = 0$$



each $f(\mathbf{x}_m), f(\mathbf{x}_n)$ pair varies in exactly
the same way

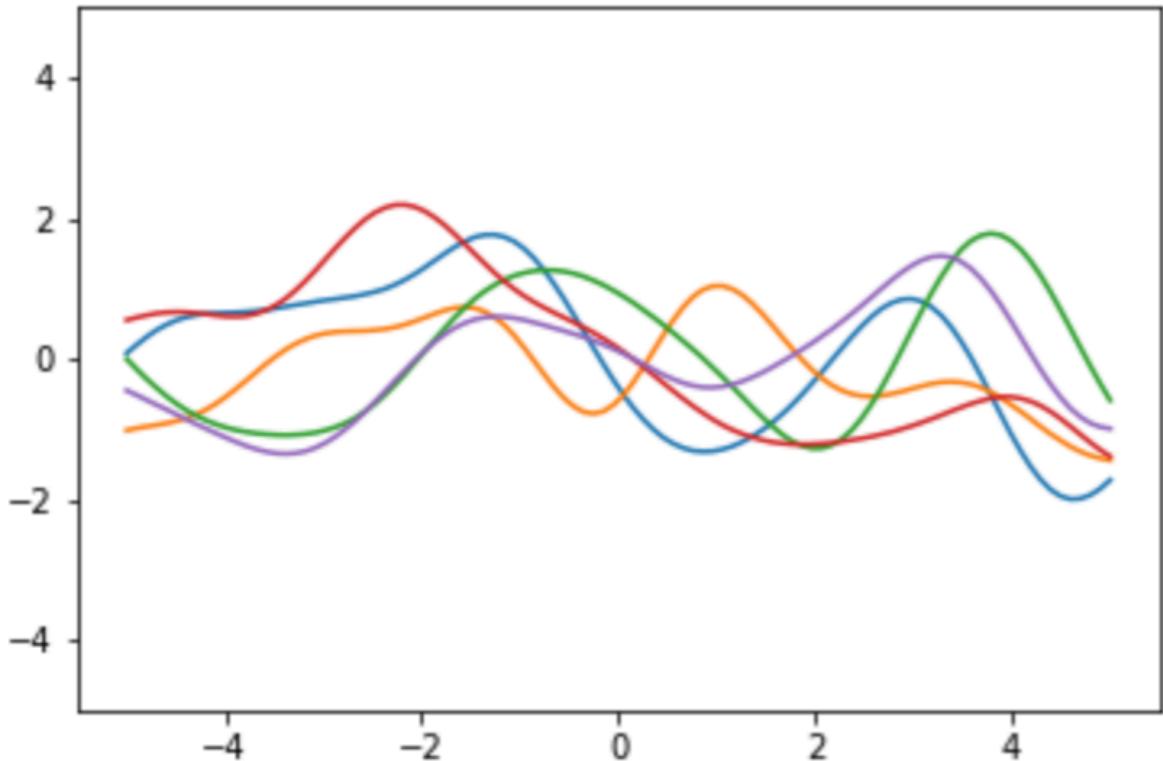
Varying the linear term θ_3

$$\delta^2 \ell(x) = \mathcal{K}(x, x)$$

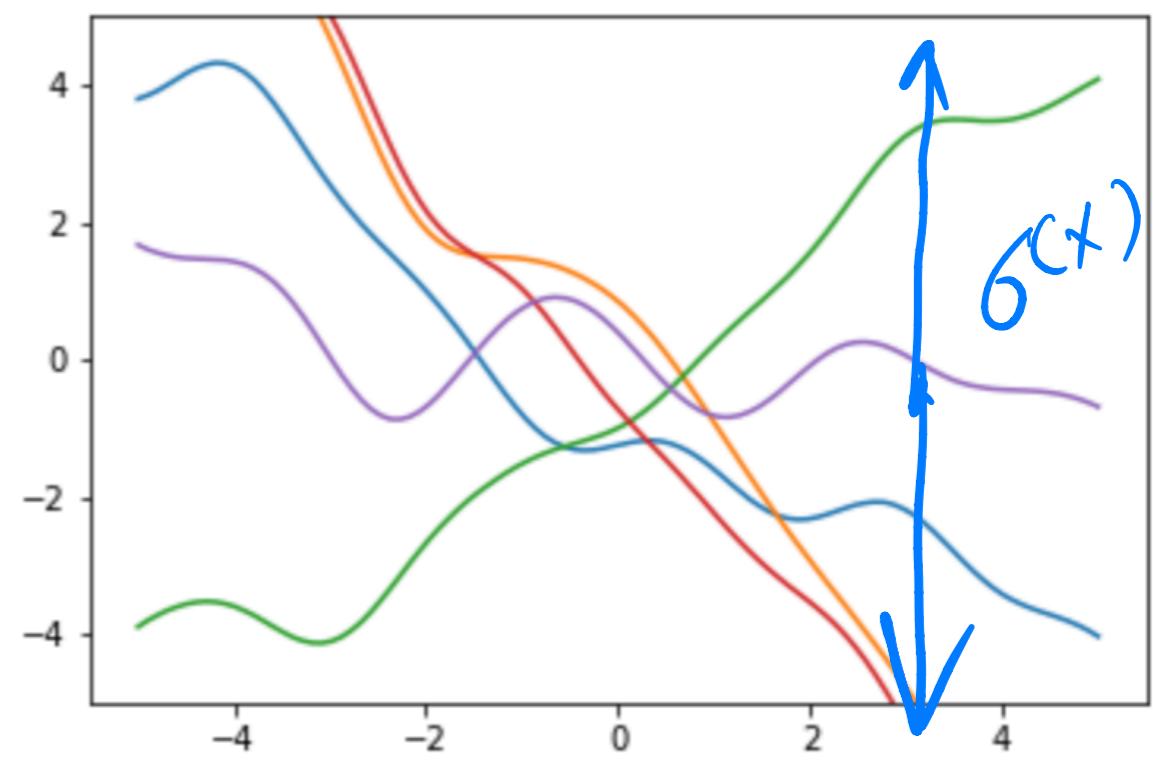
$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \underline{\theta_2} + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

Sampled functions $f(\cdot)$ tend to look linear.

$$\theta_0 = 1 \quad \theta_1 = 1 \quad \underline{\theta_2 = 0} \quad \theta_3 = 0$$



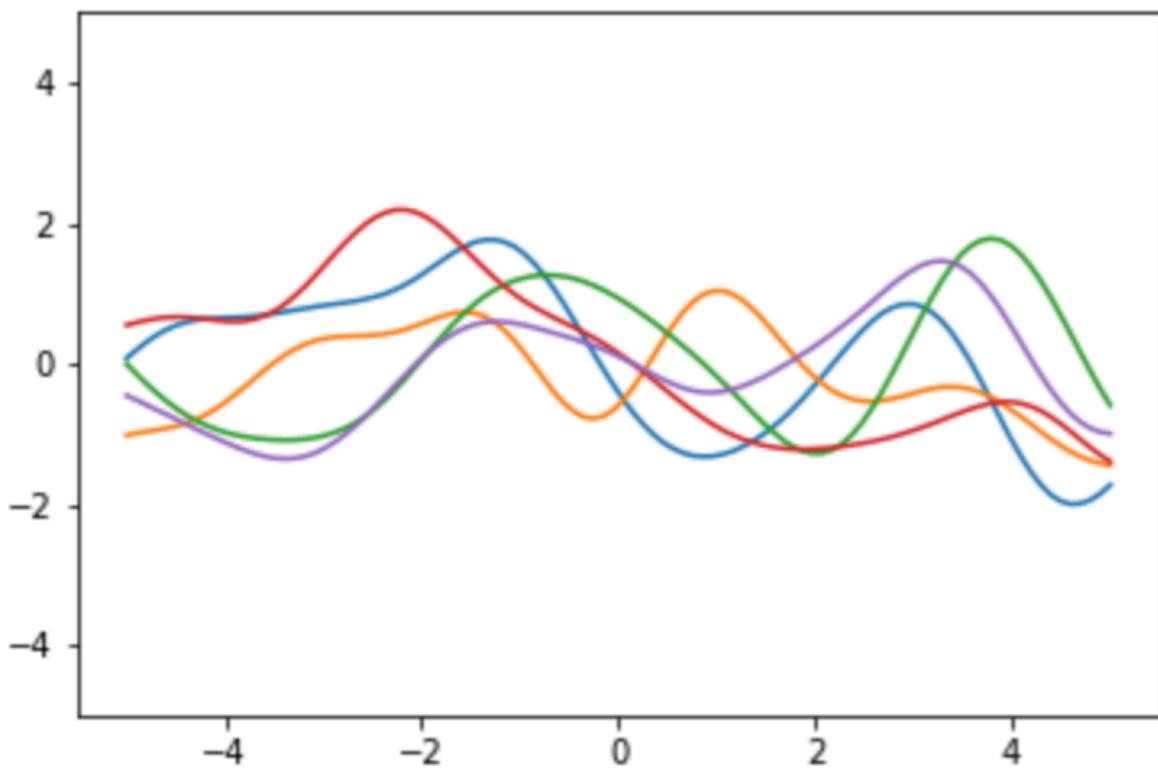
$$\theta_0 = 1 \quad \theta_1 = 1 \quad \underline{\theta_2 = 0} \quad \theta_3 = 1$$



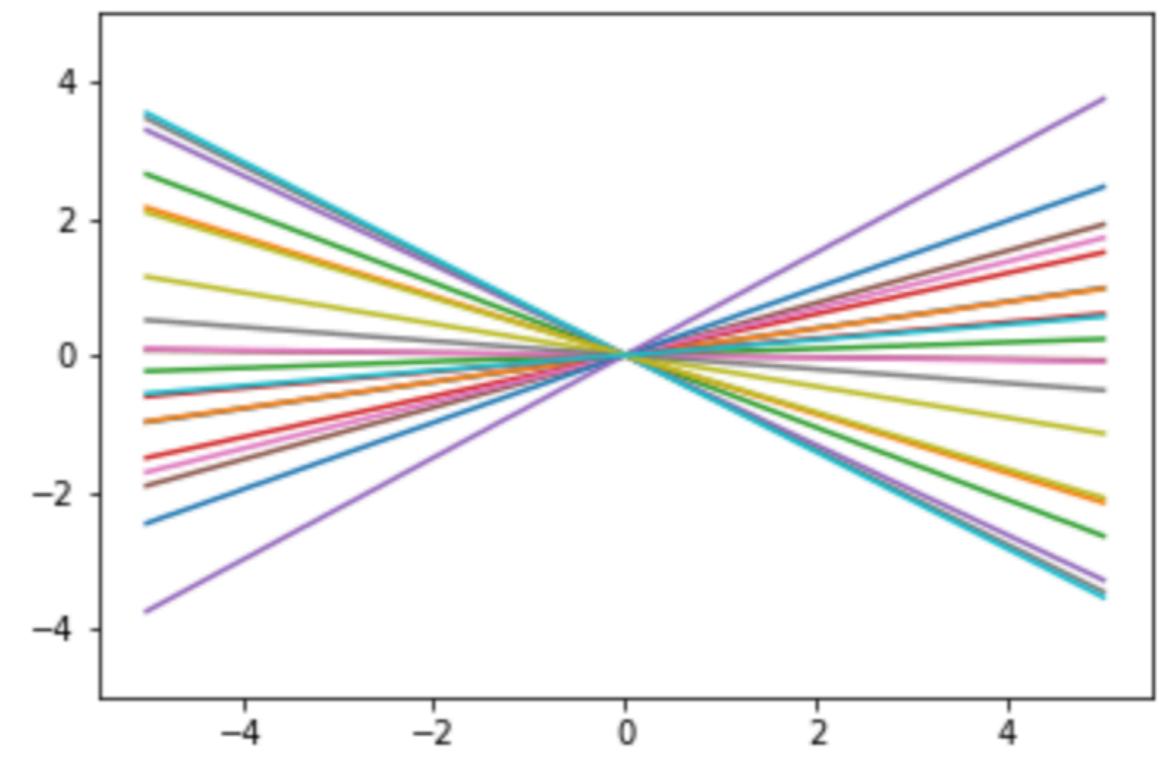
Varying the linear term θ_3

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{1}{2\theta_1^2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

$$\theta_0 = 1 \quad \theta_1 = 1 \quad \theta_2 = 0 \quad \theta_3 = 0$$



$$\theta_0 = 0 \quad \theta_1 = 0 \quad \theta_2 = 0 \quad \theta_3 = 0.2$$

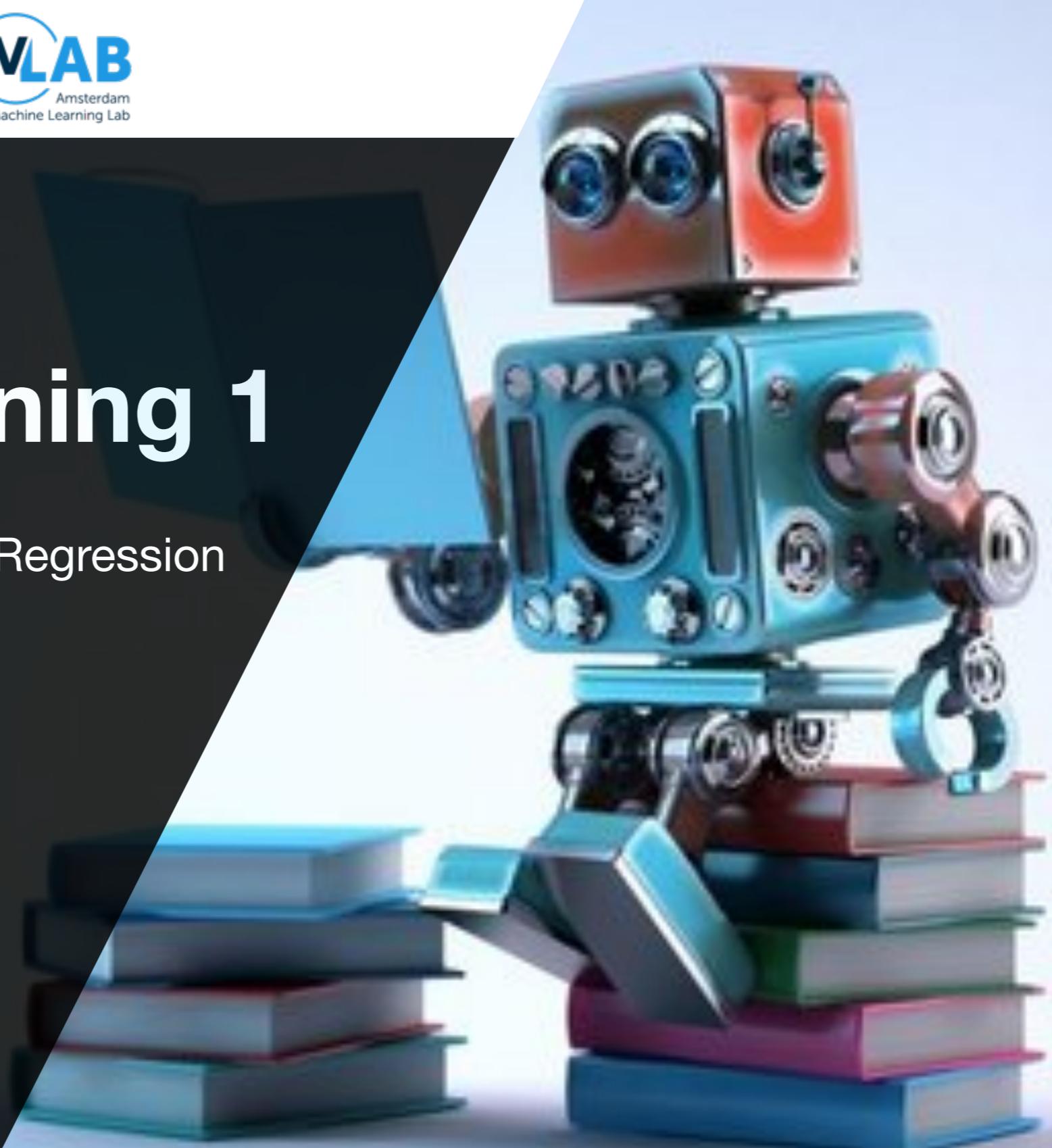


Machine Learning 1

Lecture 12.5 - Kernel Methods
Gaussian Processes - Bayesian Regression

Erik Bekkers

(Bishop 6.4.2, 6.4.3)



Regression with GP's

- We have observed $\{(\mathbf{x}_i, f_i)\}_{i=1}^N$ where we assume
no explicit parametrization, assume drawn from a GP

$$f_i = f(\mathbf{x}_i) = y(\mathbf{x}_i) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$

- Assume we have a *GP* for $y(\mathbf{x})$, so any vector of observations is a Gaussian random variable

Like a prior distribution

$$\mathbf{y} = \begin{bmatrix} y(\mathbf{x}_1) \\ \vdots \\ y(\mathbf{x}_N) \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, K \in \mathbb{R}^{N \times N} \right)$$

Choice *choice determin.*

class of func.
E.g. Smoother
functions are obtained with large θ_i

- Then $f(\cdot)$ is also a *GP* since $\mathbf{f} = \mathbf{y} + \boldsymbol{\varepsilon}$; the sum of two independent random variables is also Gaussian distributed

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, K(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I})$$

Predictions with GP's

- The joint distribution of test points \mathbf{f}' (at \mathbf{X}') and \mathbf{f} (train points), according to our *GP*, is given by

Conditioning property

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}' \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I} & \mathbf{K}(\mathbf{X}, \mathbf{X}') \\ \mathbf{K}(\mathbf{X}', \mathbf{X}) & \mathbf{K}(\mathbf{X}', \mathbf{X}') + \beta^{-1} \mathbf{I} \end{bmatrix} \right)$$

Then

posterior

$$p(\mathbf{f}' | \mathbf{X}', \mathbf{X}, \mathbf{f}) = \mathcal{N}(\mu', \Sigma')$$

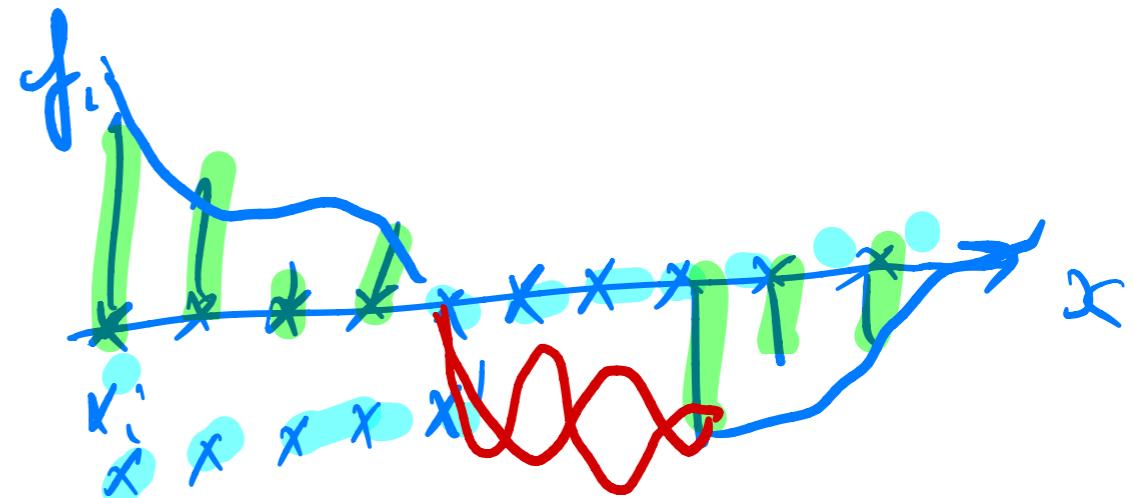
with (property of Gaussian conditionals,
see beginning of lecture)

$$\mu' = \mathbf{K}(\mathbf{X}', \mathbf{X}) (\mathbf{K}(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I})^{-1} \mathbf{f}$$

$$\Sigma' = \mathbf{K}(\mathbf{X}', \mathbf{X}') + \beta^{-1} \mathbf{I} - \mathbf{K}(\mathbf{X}', \mathbf{X}) (\mathbf{K}(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I})^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}')$$

Inverse required!

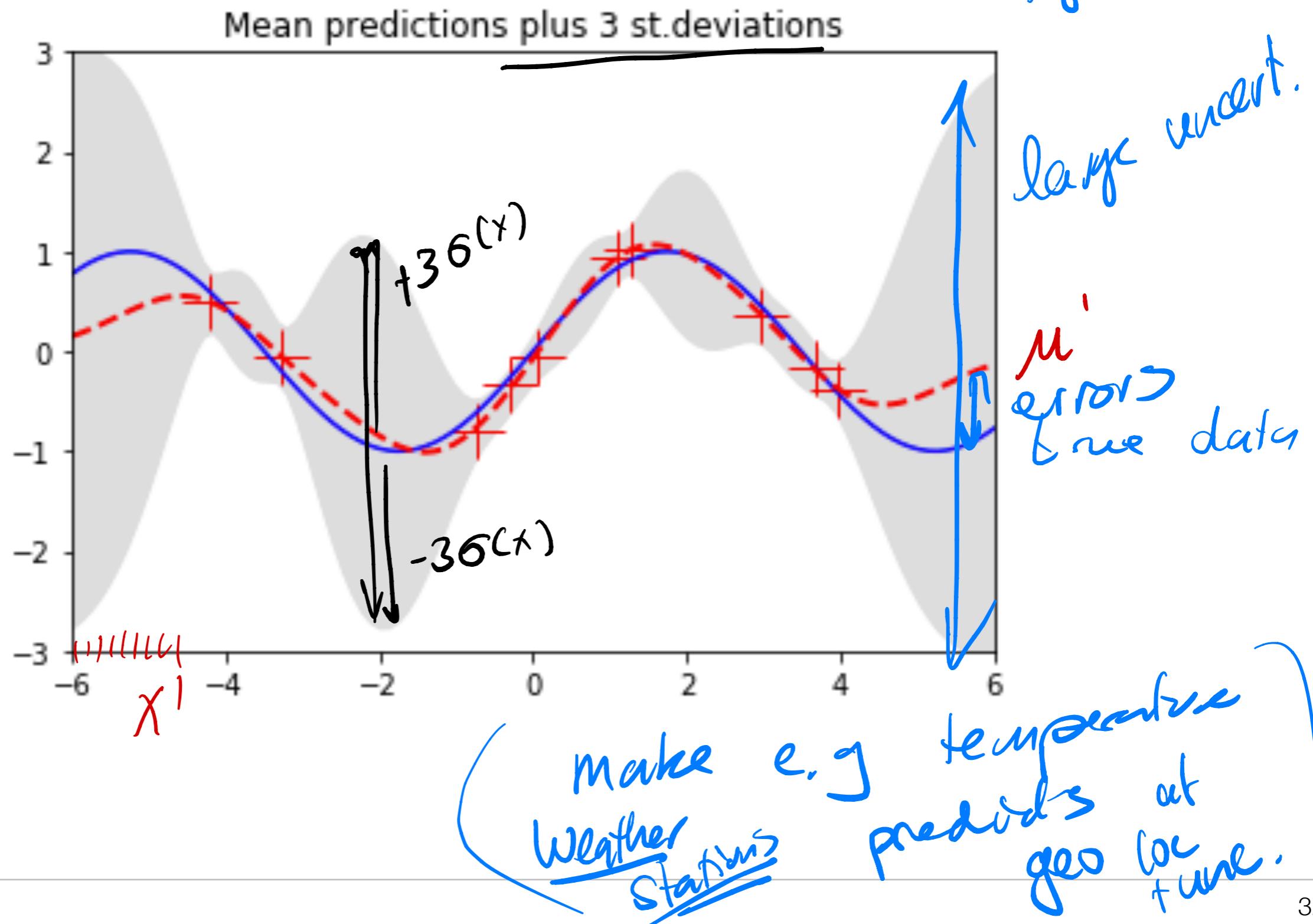
\mathbf{X}'



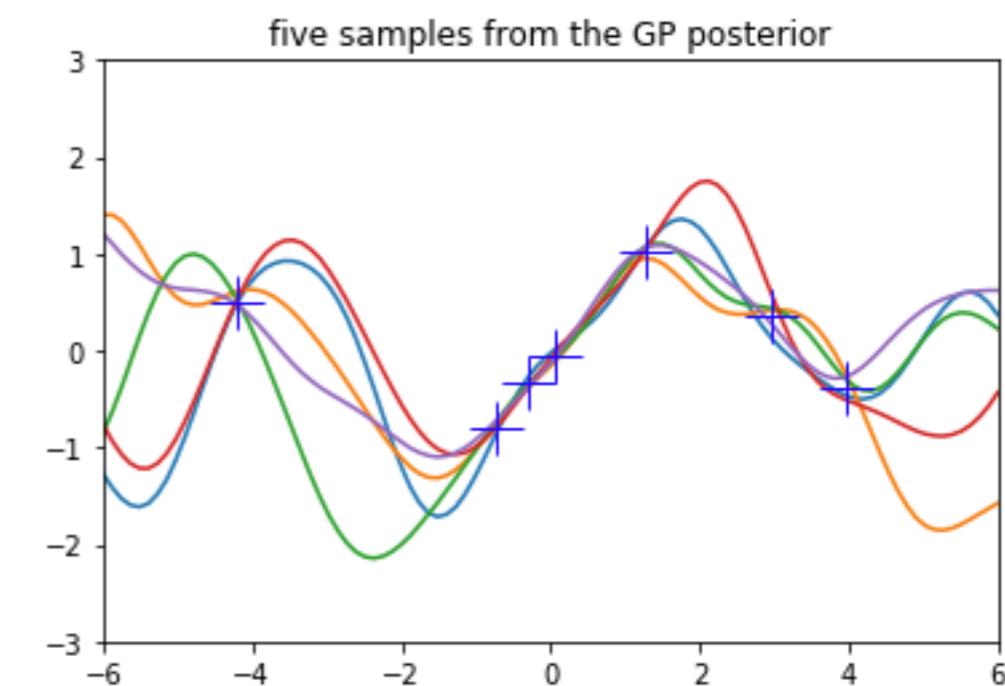
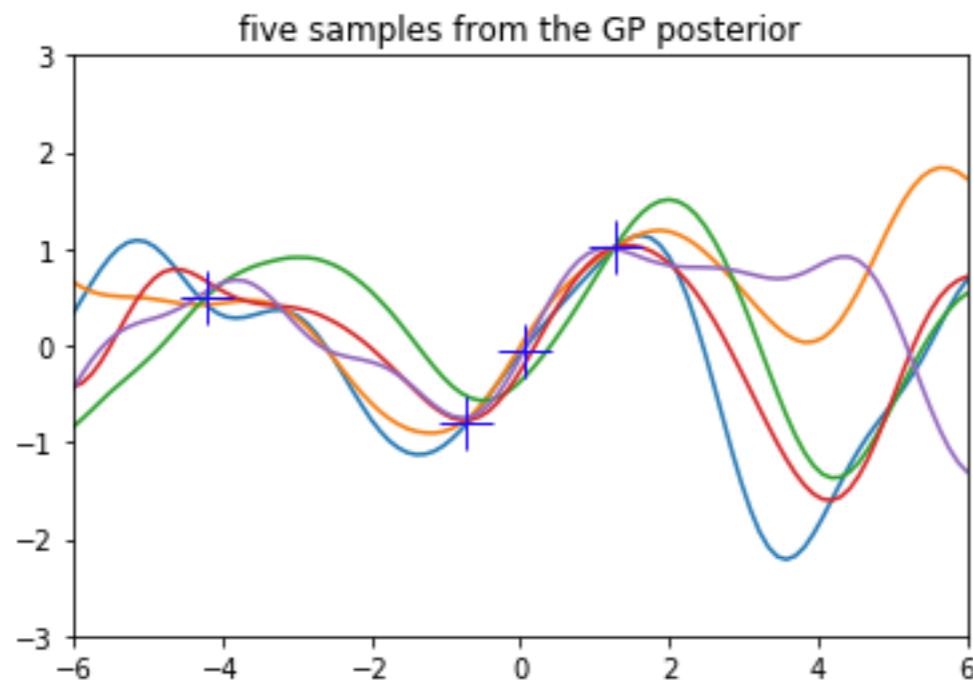
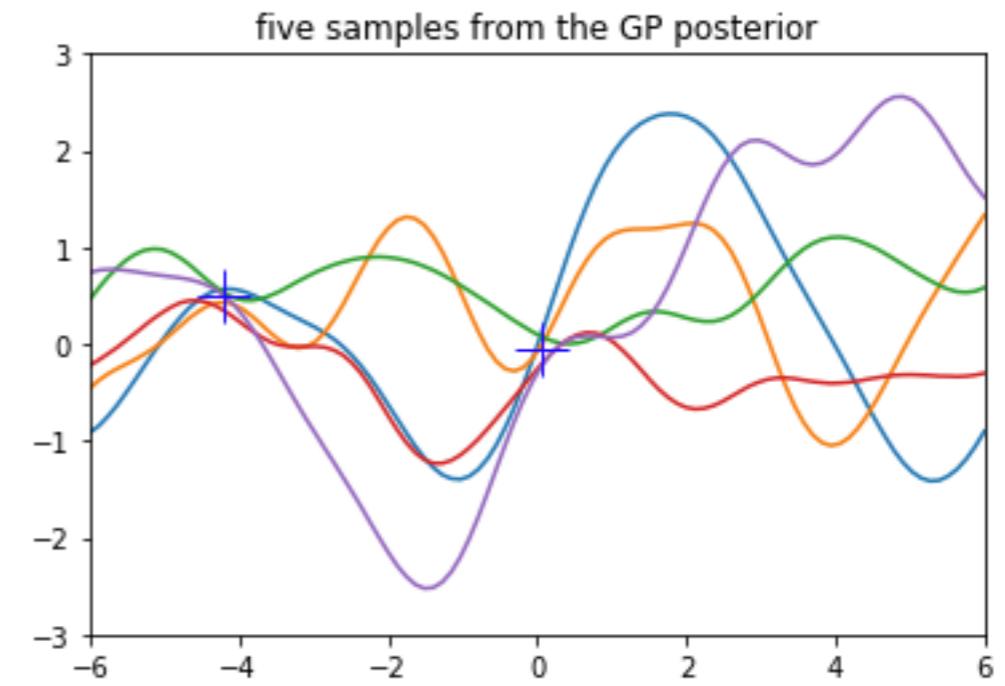
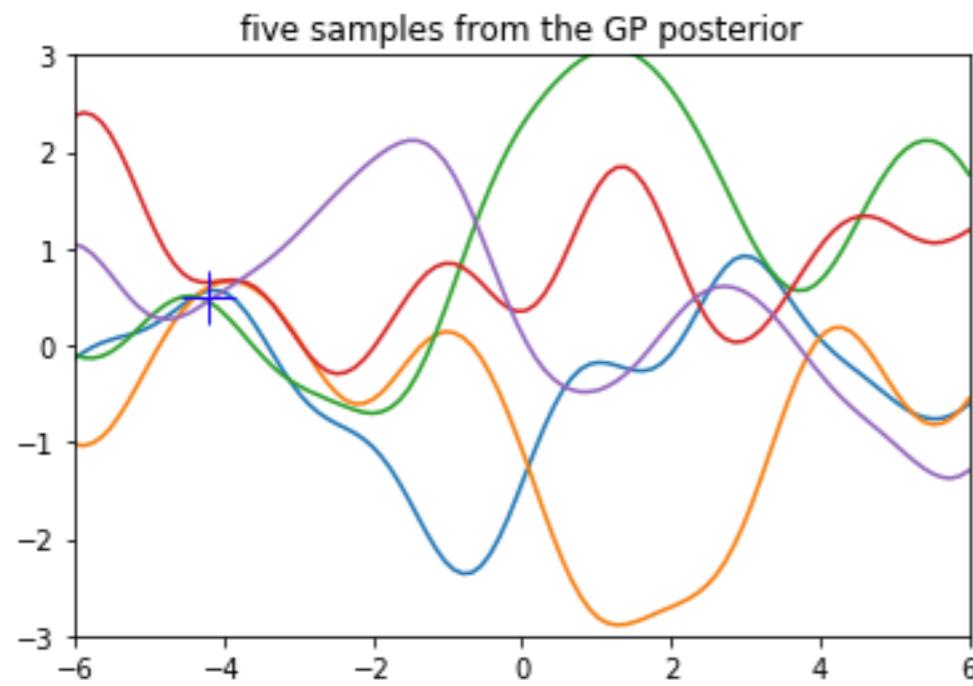
Predictions with GP's

active learning

- identify uncertain regions
- collect data in that region



Drawing functions from GP posterior



How to choose kernel parameters?

- The kernel parameters $\theta_0, \theta_1, \theta_2, \theta_3$ are hyperparameters
- Simplest approach: take training observations, for which we know

*known
data,
see slide
3.1.*

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, C_{\theta}(\mathbf{X}, \mathbf{X})) = \frac{1}{(2\pi)^{N/2} |C_{\theta}|^{1/2}} \exp\left(-\frac{1}{2}\mathbf{f}^T \mathbf{C}_{\theta}^{-1} \mathbf{f}\right)$$

with $C_{\theta}(\mathbf{X}, \mathbf{X}) = K(\mathbf{X}, \mathbf{X}) + \beta^{-1}I$

- Make a maximum likelihood estimate

$$\max_{\theta} \ln p(\mathbf{f} | \mathbf{X}, \theta) = \max_{\theta} -\frac{1}{2} \ln |\mathbf{C}_{\theta}| - \frac{1}{2} \mathbf{f}^T \mathbf{C}_{\theta}^{-1} \mathbf{f} - \frac{N}{2} \ln 2\pi$$

- Solve numerically for θ

E.g. with SGD