

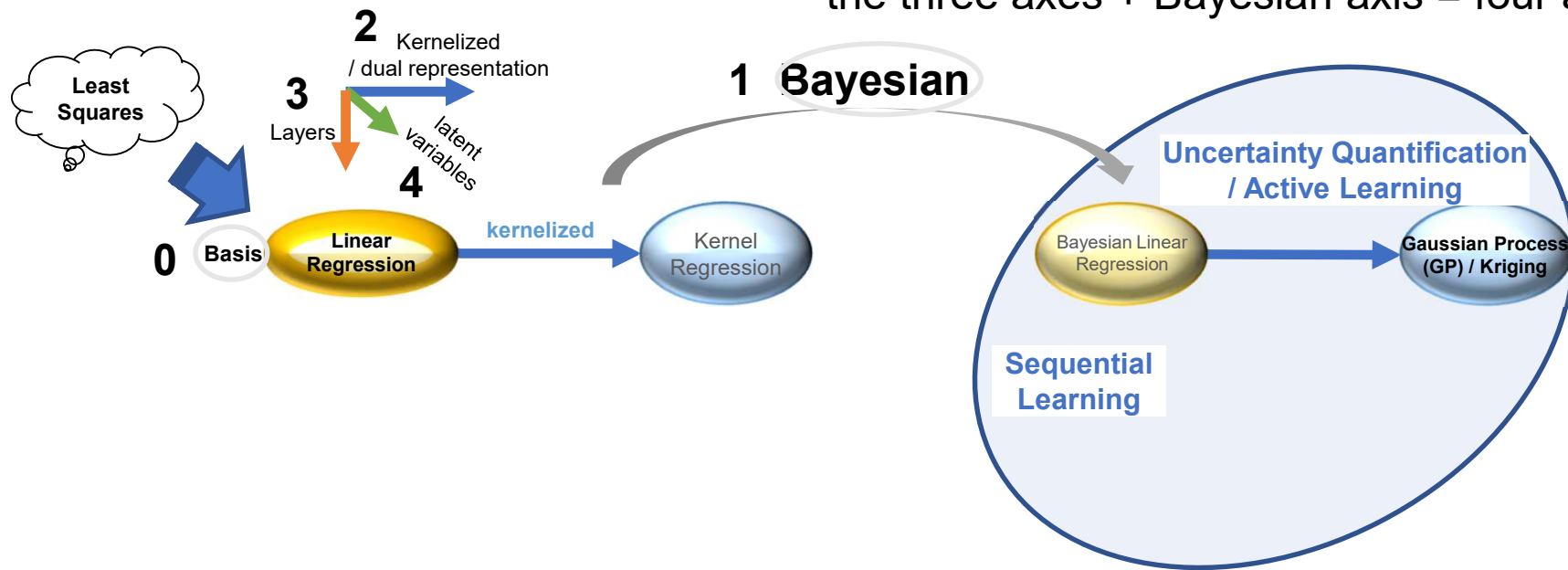


Dr. Daigo Maruyama
Prof. Dr. Ali Elham

Prof. Dr. Ali Elham

Key Components

the three axes + Bayesian axis = four axes



Introduction of Books (for Gaussian Processes)

Gaussian Process

(as perspectives from [machine learning](#)) – free pdf available officially

- C. E. Rasmussen and K. I. Williams, “Gaussian Processes for Machine Learning”, the MIT Press, 2006.

Kriging

(as perspectives from [geostatistics](#)) – no free pdf available

- Alexander Forrester, Andras Sobester, Andy Keane, “Engineering Design via Surrogate Modelling: A Practical Guide”, Wiley, 2008.
- Hans Wackernagel, “Multivariate Geostatistics”, Springer, 1995, 1998.

Please compare the description in these books from these two fields

Lecture content

- Review of Bayesian approach
- Introduction of kernel
- Dual representation (Introduction to Gaussian processes)

The lecture of this time partially follows the Chapter 6 and Section 2.3.3 and 3.3 of the book:

Christopher M. Bishop "Pattern Recognition And Machine Learning" Springer-Verlag (2006)
The name of this book is shown as "PRML" when it is referred in the slides.

The lecture slides contains many original contents in the context apart from the above sections in the book.

Lecture content

1. Bayesian approach (review)



Bayesian Approach

- Principally no overfitting
 - not because of prior information
 - because of considering all the possibilities of the parameter \mathbf{w}

➡ \mathbf{w} is not in the prediction anymore in the Bayesian approach

- **Deterministic (point estimate)**

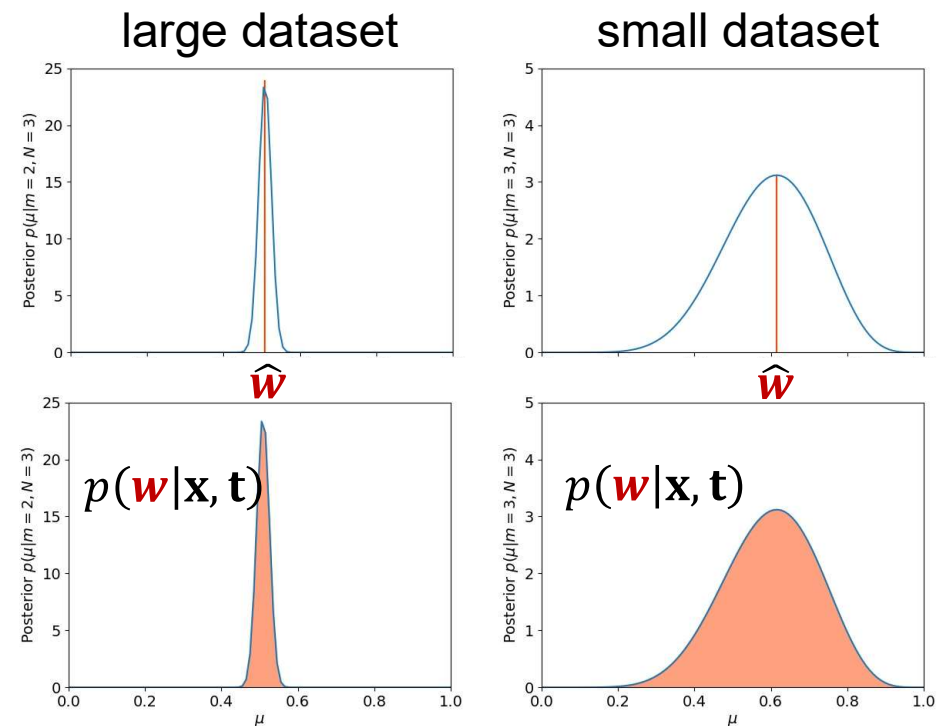
$$p(t|x, \hat{\mathbf{w}})$$

We believe that there is a true \mathbf{w}
(but learned by finite data in hand).

- **Stochastic (Bayesian)**

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t})d\mathbf{w}$$

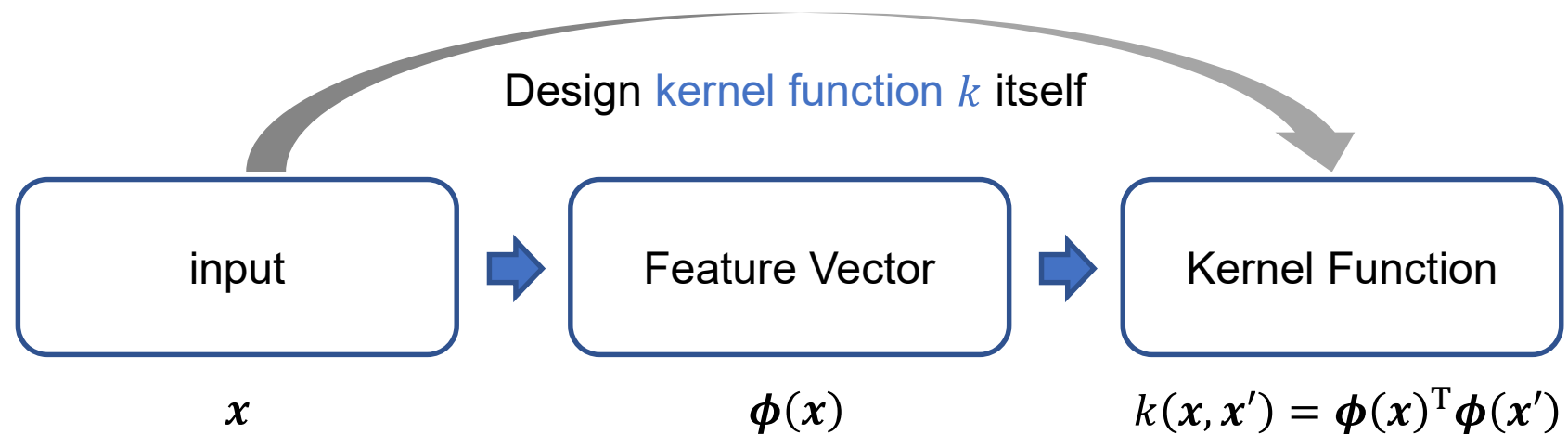
\mathbf{w} is integrated out by
reflecting all the info of \mathbf{w} .



Lecture content

2. Kernel approach

Feature Vector to Kernel Function



- Nonlinear function
- String
- Graph
- etc.

dot product
= similarity between x and x'

Linear regression

Input Vector

$$y(x, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

Linear regression (**simple**)



point estimate (MLE)

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{ND} \end{pmatrix}$$

Please see Lecture 4

Prediction by “a line” (or hyperplane in general)

Feature Vector

$$y(x, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x) \quad \text{Linear regression (general)}$$



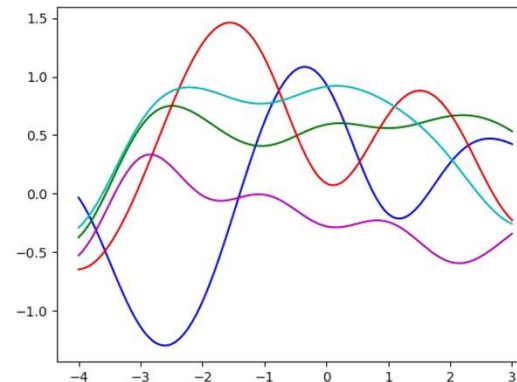
point estimate (MLE)

$$\hat{\mathbf{w}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t}$$

$$\boldsymbol{\Phi} = \begin{pmatrix} \boldsymbol{\phi}(x_1) \\ \vdots \\ \boldsymbol{\phi}(x_N) \end{pmatrix} = \begin{pmatrix} \phi_1(x_1) & \cdots & \phi_M(x_1) \\ \vdots & \ddots & \vdots \\ \phi_1(x_N) & \cdots & \phi_M(x_N) \end{pmatrix}$$

Please see Lecture 4

Prediction by “a curve” (curve fitting problem)



Kernel Function

$$y(x, \mathbf{w}) = \mathbf{w}^T \mathbf{k}(x) = \sum_{n=1}^N w_n k(x, x_n)$$



point estimate (MLE)

$$\begin{aligned} \hat{\mathbf{w}} &= (\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T \mathbf{t} \\ &= \mathbf{K}^{-1} \mathbf{t} \end{aligned}$$

Kernel regression

RBF (radial basis function) interpolation

The dimensionality (degree of freedom) of \mathbf{w} is always adjusted to the sample size N .



Guaranteed to pass through all the sample points.

Linear regression

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

$$\sigma = 1$$

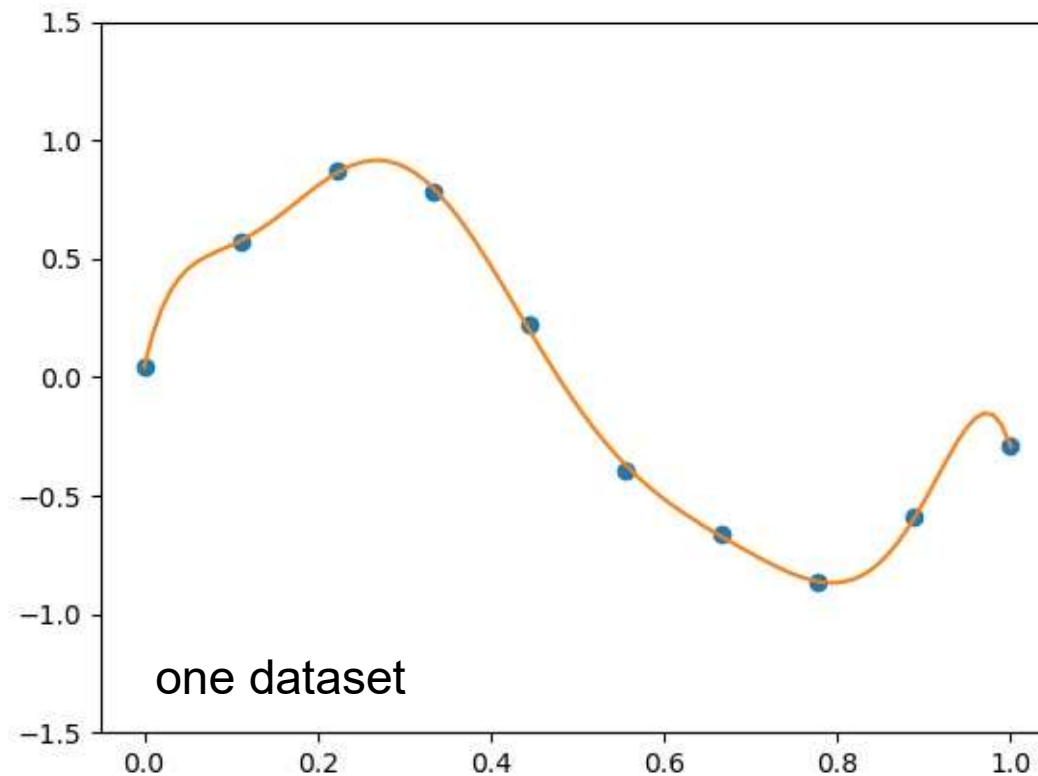
$$\mathbf{K} = \begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_N) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \cdots & k(x_N, x_N) \end{pmatrix}$$

Kernel Function

Kernel regression: examples

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

$$\sigma = 1$$

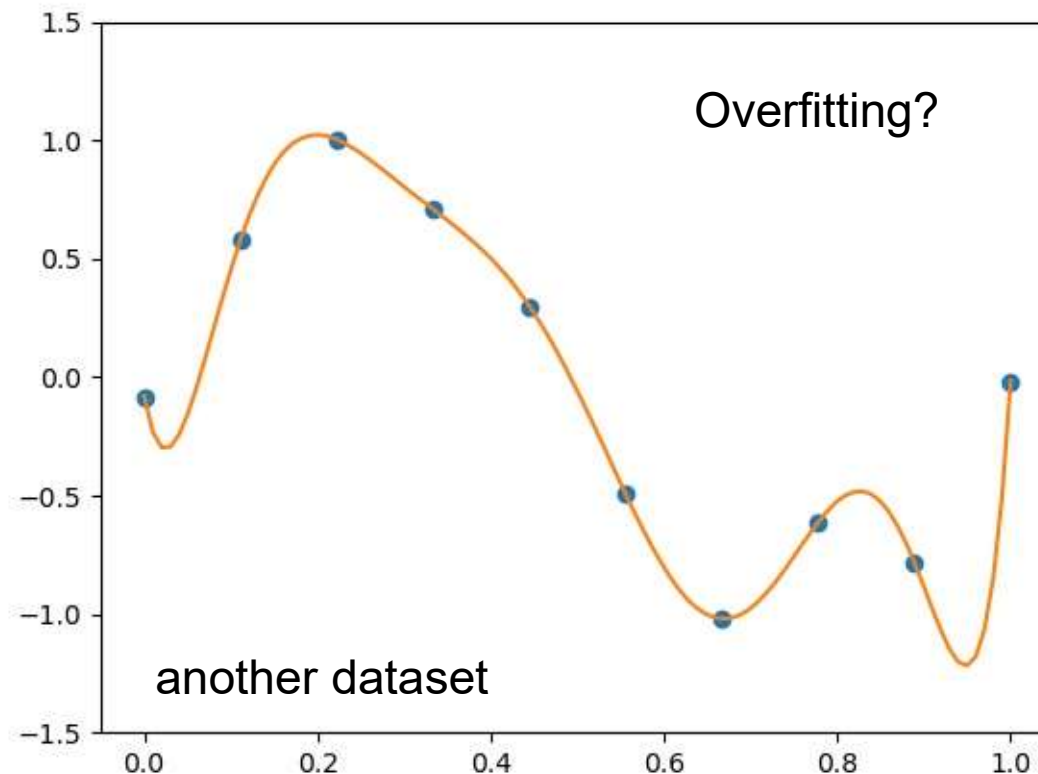


Kernel Function

Kernel regression: examples

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

$$\sigma = 1$$



$\sigma = 1$ and it cannot be tuned by MLE.



dual representation can solve this.

Design of Kernel Function

Necessary and sufficient conditions for a function $k(\mathbf{x}, \mathbf{x}')$

Gram matrix \mathbf{K} to be positive (semi)definite

the same property as the covariance matrix Σ of multiple Gaussian distributions $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$

$$\mathbf{Y} = (\mathbf{y}(\mathbf{x}_1), \dots, \mathbf{y}(\mathbf{x}_N))^T$$
$$\mathbf{Y}_C = \mathbf{Y} - E[\mathbf{Y}] \quad \text{centered}$$

$$\Sigma = \frac{1}{N} \boxed{\mathbf{Y}_C \mathbf{Y}_C^T} \quad \text{dot product}$$

$$= \begin{pmatrix} \text{cov}[\mathbf{y}(\mathbf{x}_1), \mathbf{y}(\mathbf{x}_1)] & \cdots & \text{cov}[\mathbf{y}(\mathbf{x}_1), \mathbf{y}(\mathbf{x}_N)] \\ \vdots & \ddots & \vdots \\ \text{cov}[\mathbf{y}(\mathbf{x}_N), \mathbf{y}(\mathbf{x}_1)] & \cdots & \text{cov}[\mathbf{y}(\mathbf{x}_N), \mathbf{y}(\mathbf{x}_N)] \end{pmatrix} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Definition of covariance

$$\text{cov}[\mathbf{y}, \mathbf{y}'] = E[(\mathbf{y} - E[\mathbf{y}])(\mathbf{y}' - E[\mathbf{y}'])^T]$$

$$\Sigma[\mathbf{Y}, \mathbf{Y}] = E[(\mathbf{Y} - E[\mathbf{Y}])(\mathbf{Y} - E[\mathbf{Y}])^T]$$

Arbitrary covariance matrix Σ can be a Gram matrix \mathbf{K} .

Design of Kernel Function

Kernel trick

Gram matrix \mathbf{K} to be **positive (semi)definite**

➡ As far as this is satisfied, any kernel functions $k(\mathbf{x}, \mathbf{x}')$ can be designed.

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Design of Kernel Function

An example of kernel trick

$$k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right) \quad \text{Gaussian kernel}$$

We did not have to explicitly specify what the feature vector $\boldsymbol{\phi}(\mathbf{x})$ is.

➡ $\boldsymbol{\phi}(\mathbf{x})$ in this case is originally a nonlinear function of infinity dimension.

An exercise to get familiar with kernel functions see PRML, 294-295

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

➡ What is $\boldsymbol{\phi}(\mathbf{x})$?

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

$$= \dots$$

$$= (x_1^2, \sqrt{2}x_1x_2, x_2^2)(z_1^2, \sqrt{2}z_1z_2, z_2^2)^T$$

Design of Kernel Function

Once valid kernels are created:

Sum of kernel = kernel

Product of kernel = kernel

...

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$$

...

Practical design of kernel functions:

Existing kernel functions which are already guaranteed that the covariance matrix composed by them are positive (semi)definite.



Read books and papers, then free to compose new kernels based on them.

This can be naturally understood by considering some covariance matrices:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix} + \hat{\sigma}^2 \begin{pmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{pmatrix}$$

$$k(\mathbf{x}, \mathbf{x}') = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}' \\ 0 & \text{else} \end{cases}$$

3. Kernels from dual representation

- A bridge between the Bayesian linear regression and Gaussian processes
- Introduction to Gaussian processes in the next lectures

Bayesian Linear Regression (REVIEW)

Let's try to apply this concept to **the curve fitting problem**.

Probabilistic model

$$p(t|x, \mathbf{w}, \sigma) = \mathcal{N}(t|y(x, \mathbf{w}), \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{\{t - y(x, \mathbf{w})\}^2}{2\sigma^2} \right\}$$

Likelihood function

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma) = \prod_{i=1}^N \mathcal{N}(t_i|y(x_i, \mathbf{w}), \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{\{t_i - y(x_i, \mathbf{w})\}^2}{2\sigma^2} \right]$$

➡ **Posterior distribution** $p(\mathbf{w}, \sigma|\mathbf{x}, \mathbf{t}) = \textit{complicated}$

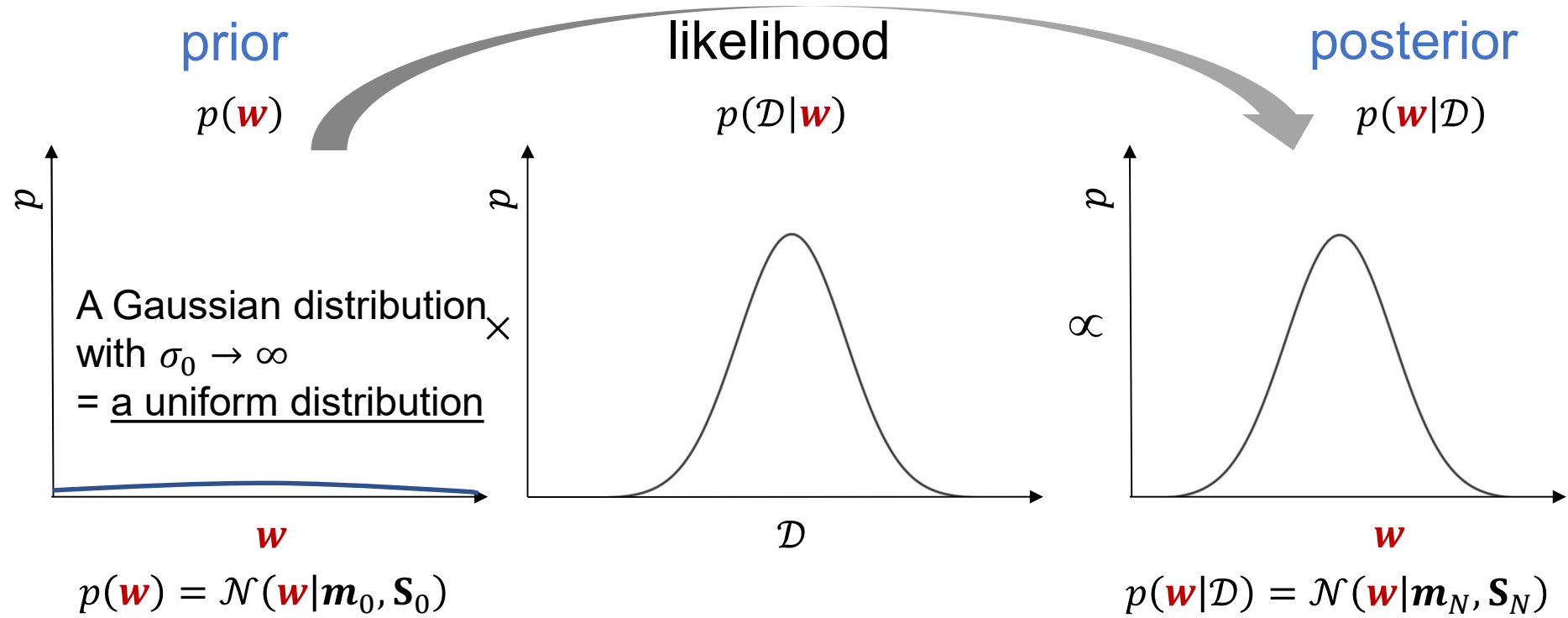
This is a Gaussian distribution wrt t_i ,
but NOT a Gaussian distribution wrt \mathbf{w}, σ .

Predictive distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \iint p(t|x, \mathbf{w}, \sigma) p(\mathbf{w}, \sigma|\mathbf{x}, \mathbf{t}) d\mathbf{w} d\sigma = \textit{complicated}$$

probabilistic model \times **posterior**

Conjugate Prior + Linear Regression (REVIEW)



**Conjugate prior
+ Linear regression**

Analytical solutions of

- Posterior distribution
- Predictive distribution

Bayesian Linear Regression

Let's try to apply this concept to **the curve fitting problem**.

Probabilistic model

$$p(t|x, \mathbf{w}) = \mathcal{N}(t|y(x, \mathbf{w}), \hat{\sigma}^2) = \textit{Gaussian}$$

Likelihood function

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \textit{Gaussian}$$

If **Prior distribution** $p(\mathbf{w}) = \textit{Gaussian}$

➡ **Posterior distribution** $p(\mathbf{w}|\mathbf{x}, \mathbf{t}) = \textit{Gaussian}$

Predictive distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t})d\mathbf{w} = \textit{Gaussian}$$

\mathbf{w} is already integrated out.

Directly design the
predictive distribution

please read 2.3.3 in PRML to follow the formulations.
But the concept described here is important.

Bayesian Linear Regression

This has to be a linear regression model.

1. We have (defined):

- A probabilistic model (of t)

$$p(t|\mathbf{w}) = \mathcal{N}(t|\mathbf{w}^T\boldsymbol{\phi}, \hat{\Sigma})$$

- A prior distribution (of \mathbf{w})

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

a familiar probabilistic model (of t)

$$p(t|\mathbf{w}) = \mathcal{N}(t|\mathbf{w}^T\boldsymbol{\phi}(x), \hat{\sigma}^2)$$

$$p(\mathbf{t}|\mathbf{w}) = \mathcal{N}\left(\begin{matrix} t_1 \\ t_2 \\ \vdots \end{matrix} \middle| \begin{matrix} \mathbf{w}^T\boldsymbol{\phi}(x_1) \\ \mathbf{w}^T\boldsymbol{\phi}(x_2) \\ \vdots \end{matrix}, \begin{pmatrix} \hat{\sigma}^2 & 0 & \dots \\ 0 & \hat{\sigma}^2 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}\right)$$

$$\hat{\Sigma} = \hat{\sigma}^2 \mathbf{I}$$

2. Then we obtain:

- A posterior distribution (of \mathbf{w})

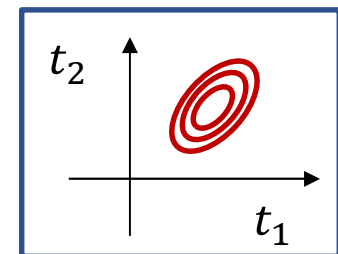
$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

- A predictive distribution (of t)

$$p(t) = \mathcal{N}(t|\mathbf{m}_N^T\boldsymbol{\phi}, \hat{\Sigma} + \boldsymbol{\phi}\mathbf{S}_N\boldsymbol{\phi}^T)$$



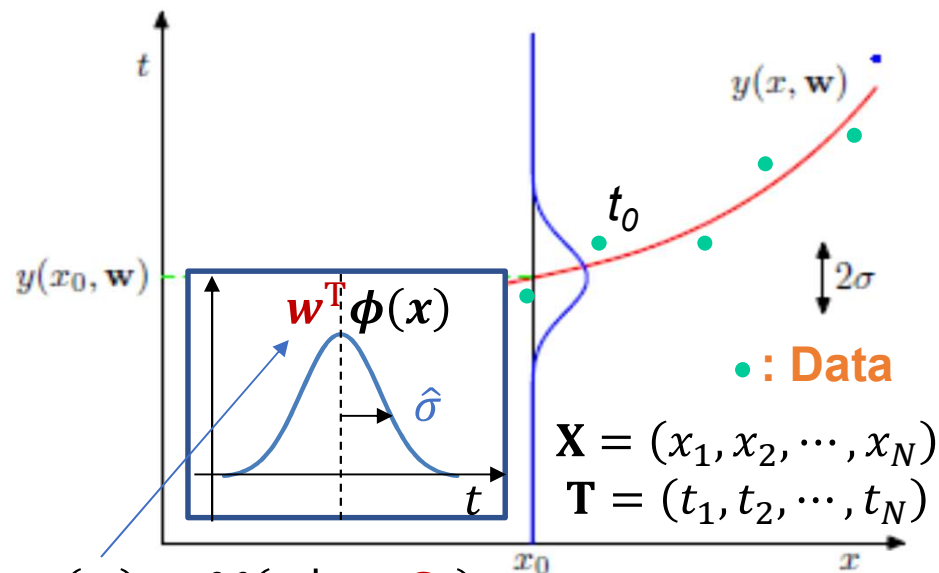
A multivariate Gaussian distribution
in general



Dual Representation

Probabilistic model:
An isotropic Gaussian distribution

$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$



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

Prior of \mathbf{w}

Start (traditional)

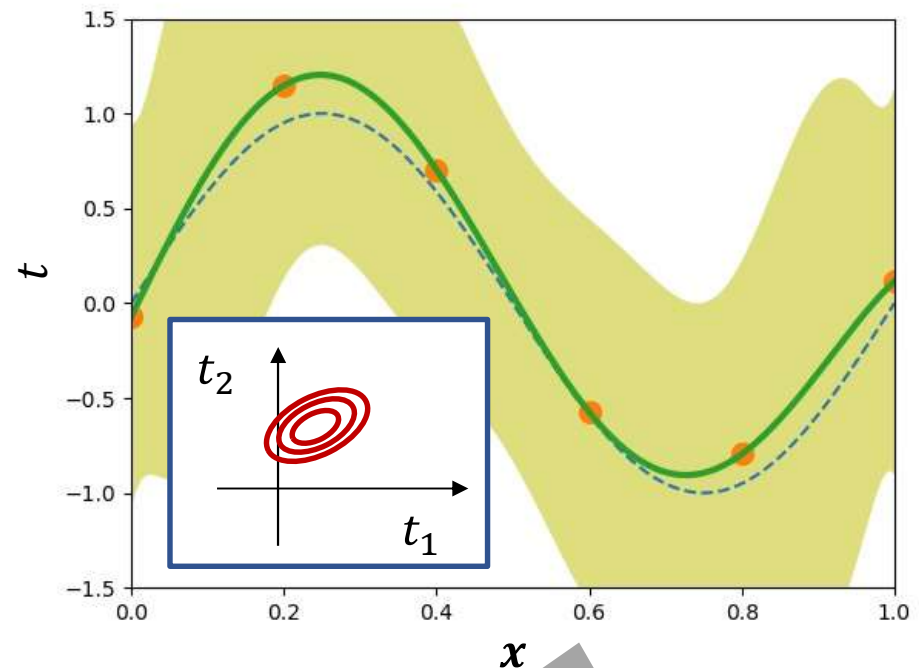
Arbitrary data: \mathbf{X}, \mathbf{T}

Goal (traditional)

Predictive distribution:

A multivariate Gaussian distribution

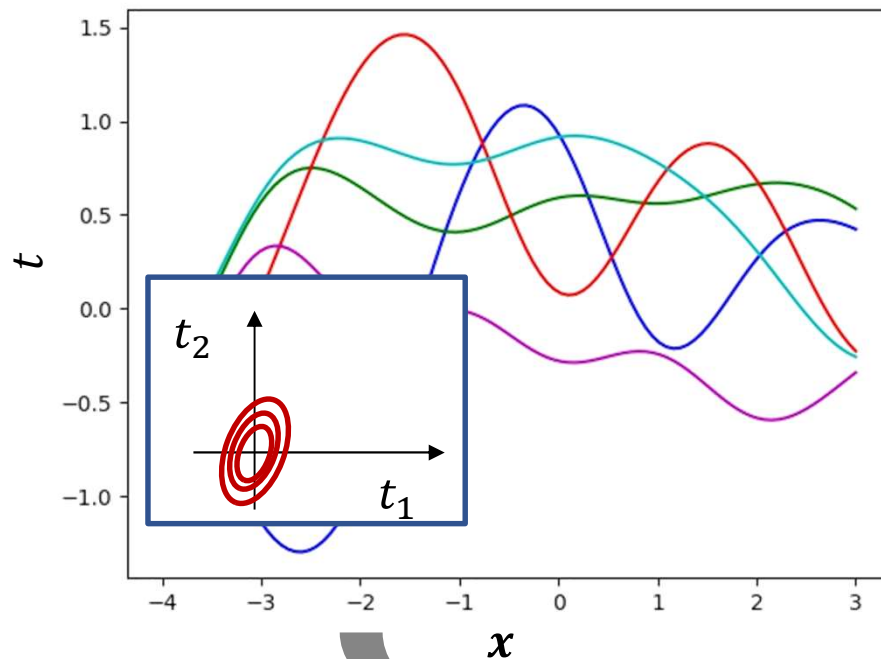
$$p(\mathbf{t} | \mathbf{x}, \mathbf{X}, \mathbf{T}) = \mathcal{N}(\mathbf{t} | \mathbf{m}_N^T \boldsymbol{\Phi}, \hat{\sigma}^2 \mathbf{I} + \boldsymbol{\Phi} \mathbf{S}_N \boldsymbol{\Phi}^T)$$



Dual Representation

Predictive distribution:
A multivariate Gaussian distribution

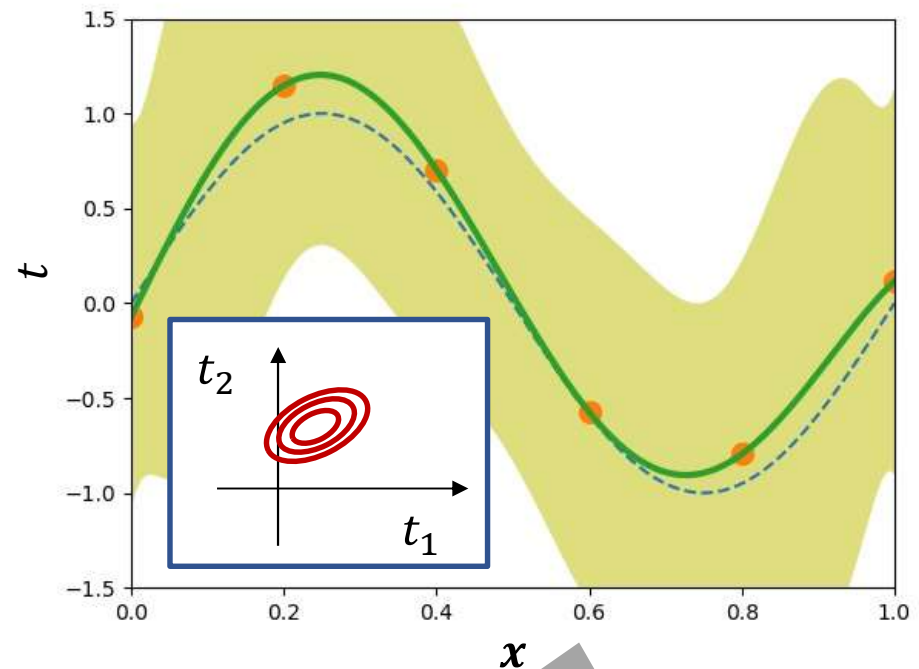
$$p(\mathbf{t}|\mathbf{x}) = \mathcal{N}(\mathbf{t} | \mathbf{m}_0^T \Phi, \hat{\sigma}^2 \mathbf{I} + \Phi \mathbf{S}_0 \Phi^T)$$



Start (Dual)

Predictive distribution:
A multivariate Gaussian distribution

$$p(\mathbf{t}|\mathbf{x}, \mathbf{X}, \mathbf{T}) = \mathcal{N}(\mathbf{t} | \mathbf{m}_N^T \Phi, \hat{\sigma}^2 \mathbf{I} + \Phi \mathbf{S}_N \Phi^T)$$



Goal (Dual)

Arbitrary data: \mathbf{X}, \mathbf{T}



Dual Representation

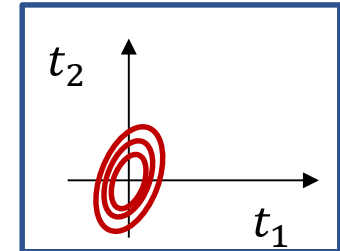
Direct definition of a predictive distribution:
A multivariate Gaussian distribution
(represented by kernels)

1. We have (defined):

- A predictive distribution (of \mathbf{t})

$$p(\mathbf{t}) = \mathcal{N}(\mathbf{t} | \mathbf{m}^T \Phi, \hat{\Sigma} + \Phi \mathbf{S} \Phi^T)$$

Since we know the form.



- A prior distribution (of \mathbf{w})

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

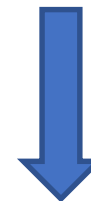
$$p(\mathbf{t} | \mathbf{x}) = \mathcal{N}(\mathbf{t} | \mathbf{m}_0^T \Phi, \hat{\Sigma} + \Phi \mathbf{S}_0 \Phi^T)$$

$\begin{matrix} \text{= 0} & \text{\(\equiv K\)} \end{matrix}$

Prior

2. Then we obtain:

- A posterior distribution (of \mathbf{w})



Data \mathbf{X}, \mathbf{T}
is given.

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{X}, \mathbf{T}) = \mathcal{N}(\mathbf{t} | \mathbf{m}_N^T \Phi, \hat{\Sigma} + \Phi \mathbf{S}_N \Phi^T)$$

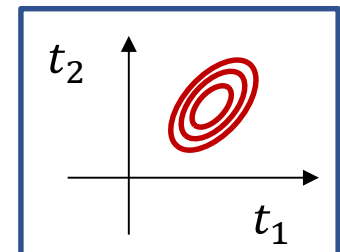
Posterior

\mathbf{x}, \mathbf{t} : prediction

\mathbf{X}, \mathbf{T} : data

$$\mathbf{m}_N = \mathbf{m}_N(\mathbf{X}, \mathbf{T})$$

$$\mathbf{S}_N = \mathbf{S}_N(\mathbf{X})$$



Dual Representation

Direct definition of a predictive distribution:
A multivariate Gaussian distribution
(represented by kernels)

1. We have (defined):

- A predictive distribution (of \mathbf{t})

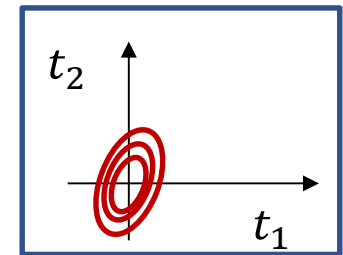
$$p(\mathbf{t}) = \mathcal{N}(\mathbf{t} | \mathbf{m}^T \Phi, \hat{\Sigma} + \Phi \mathbf{S} \Phi^T) \quad \text{Since we know the form.}$$

- A prior distribution (of \mathbf{w})

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

$$p(\mathbf{t} | \mathbf{x}) = \mathcal{N}(\mathbf{t} | \mathbf{0}, \mathbf{K}_0)$$

Prior



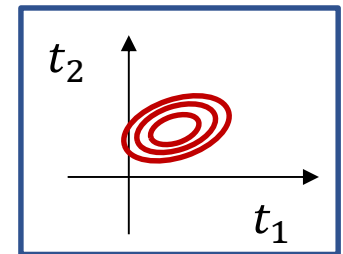
2. Then we obtain:

- A posterior distribution (of \mathbf{w})

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

$$p(\mathbf{t} | \mathbf{x}) = \mathcal{N}(\mathbf{t} | \mathbf{0}, \mathbf{K}_N)$$

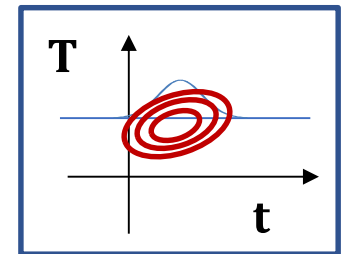
← Data \mathbf{X}, \mathbf{T}
is given.



Please remind you of **the conditional Gaussian distributions** in general introduced in Lecture 4.

$\mathbf{K}_N \neq \mathbf{K}_N(\mathbf{X}, \mathbf{T})$

Posterior



$$p(\mathbf{t} | \mathbf{x}, \mathbf{X}, \mathbf{T}) = \mathcal{N}(\mathbf{t} | \mathbf{k}_N(\mathbf{x})^T \mathbf{K}_N^{-1} \mathbf{t}, \mathbf{K}_N(\mathbf{x}, \mathbf{x}) - \mathbf{k}_N(\mathbf{x})^T \mathbf{K}_N^{-1} \mathbf{k}_N(\mathbf{x}))$$