

Scientific Machine Learning

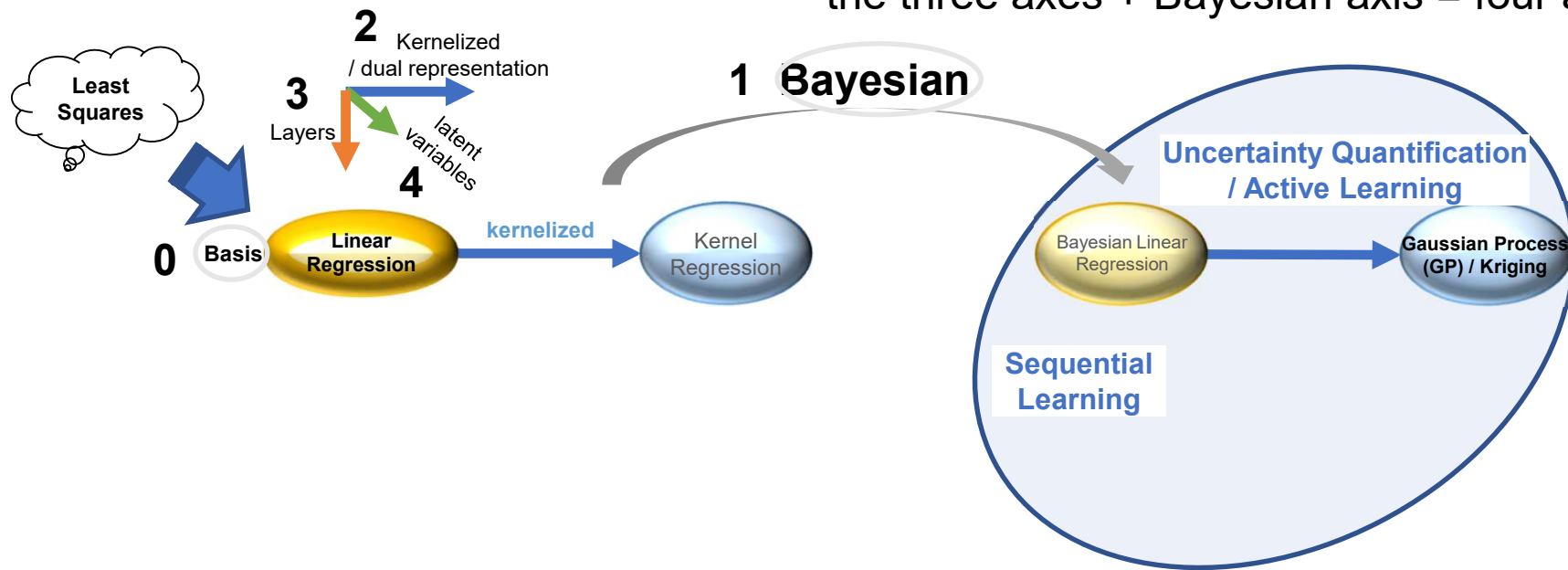
Lecture 9: Gaussian Process (2/2)

Dr. Daigo Maruyama

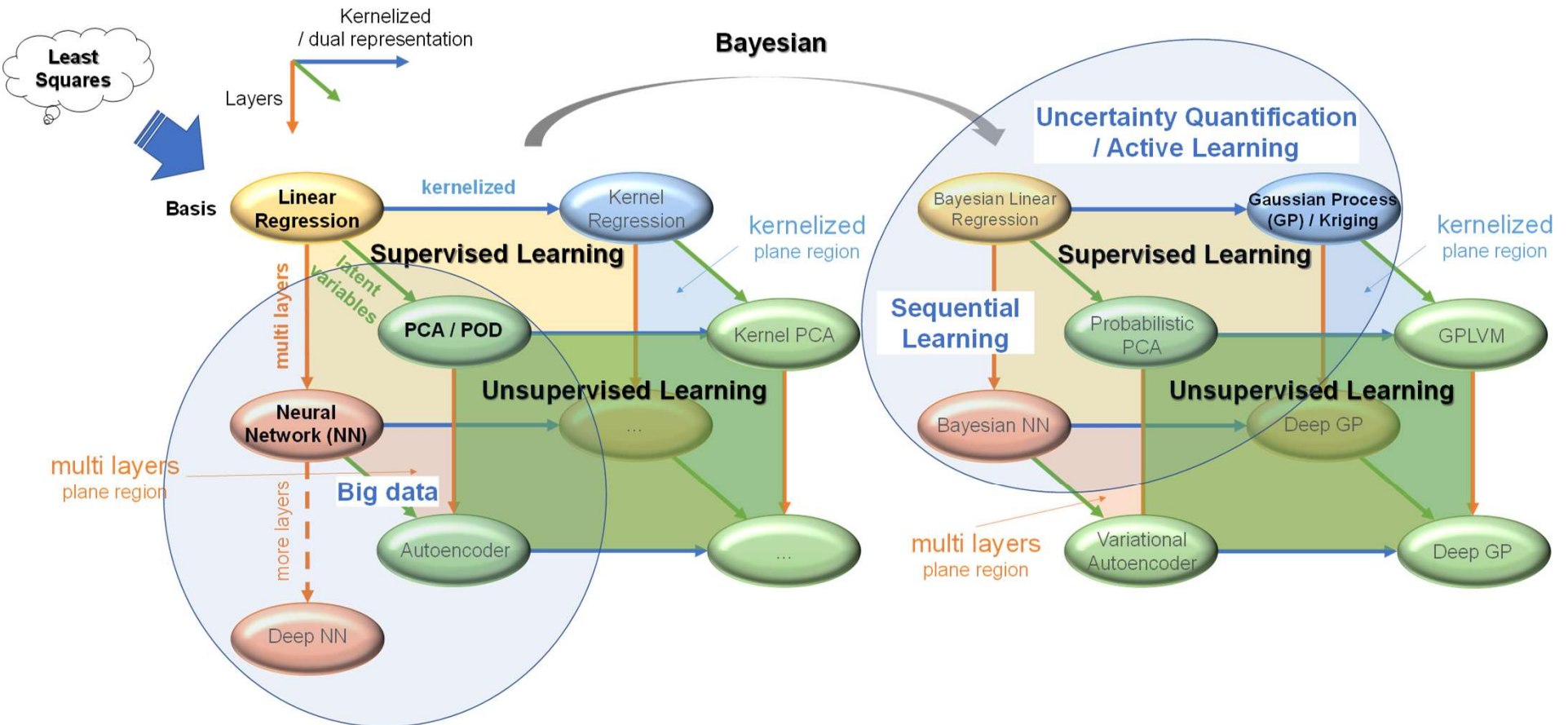
Prof. Dr. Ali Elham

Key Components

the three axes + Bayesian axis = four axes



Structures of the Techniques/Methods



Lecture content

- Further topics of Gaussian processes and introduction of Kriging
- Applications of Gaussian processes
- Advanced Methods/Topics

The lecture of this time partially follows the Section 6.4 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.

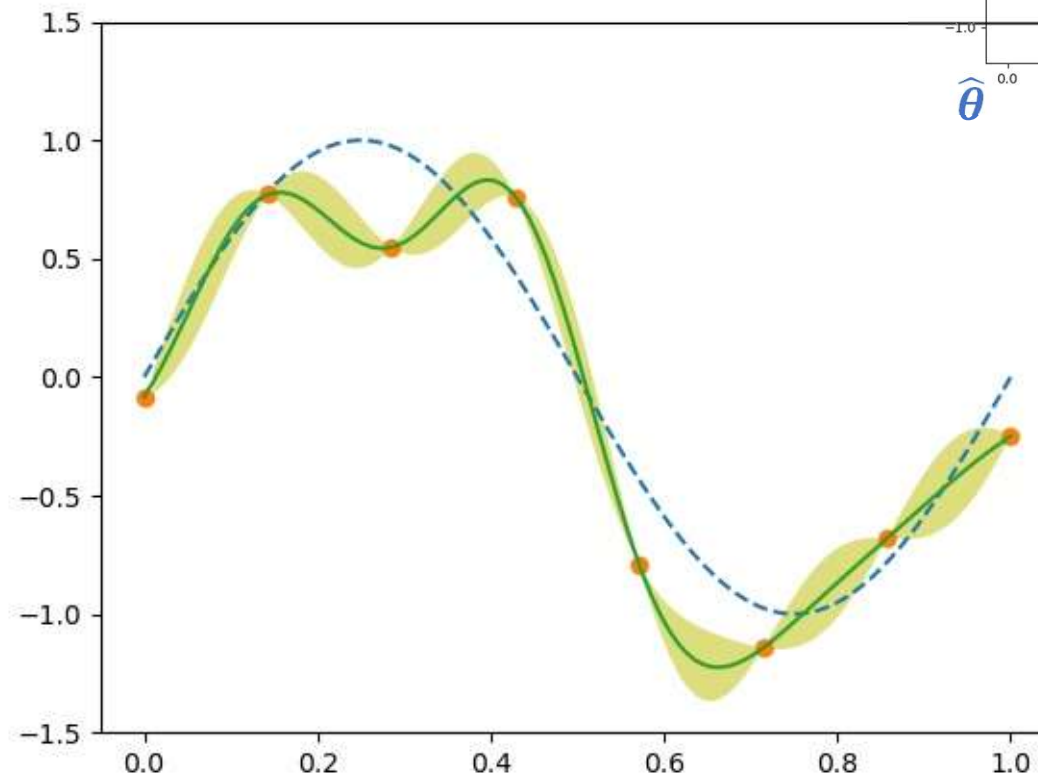
Lecture content

- Further topics of Gaussian processes and introduction of Kriging



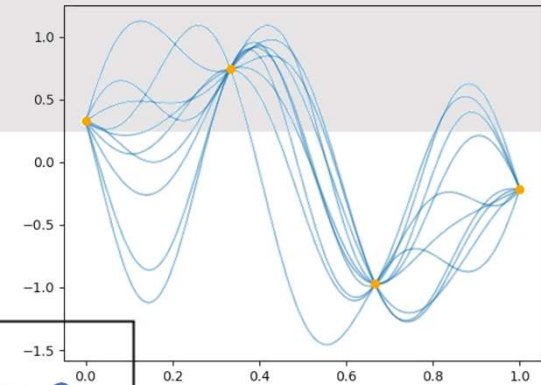
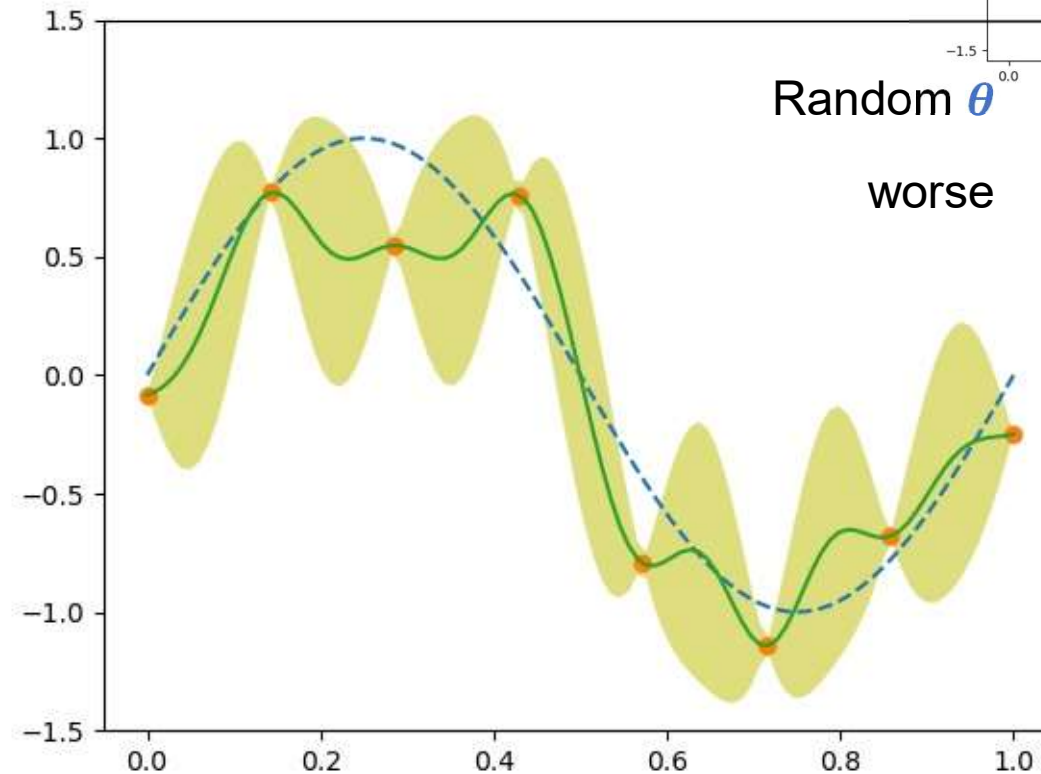
Gaussian Processes

Example ($\hat{\theta}$ by MLE)



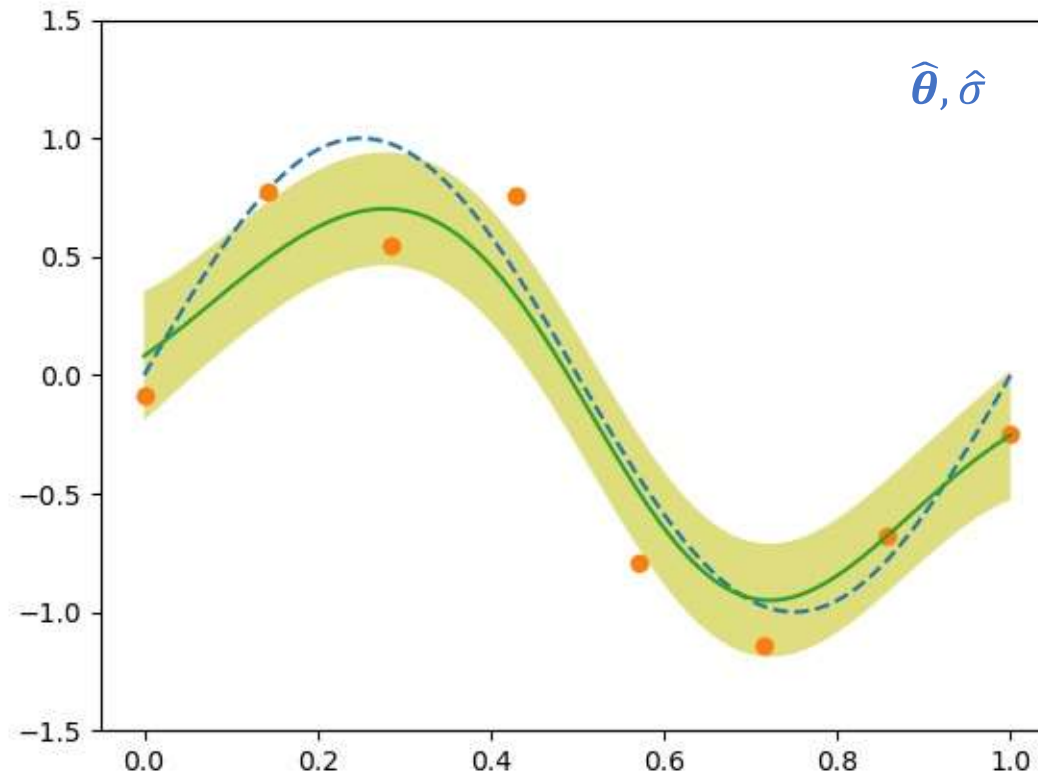
Gaussian Processes

Examples (random θ without learning from the data)



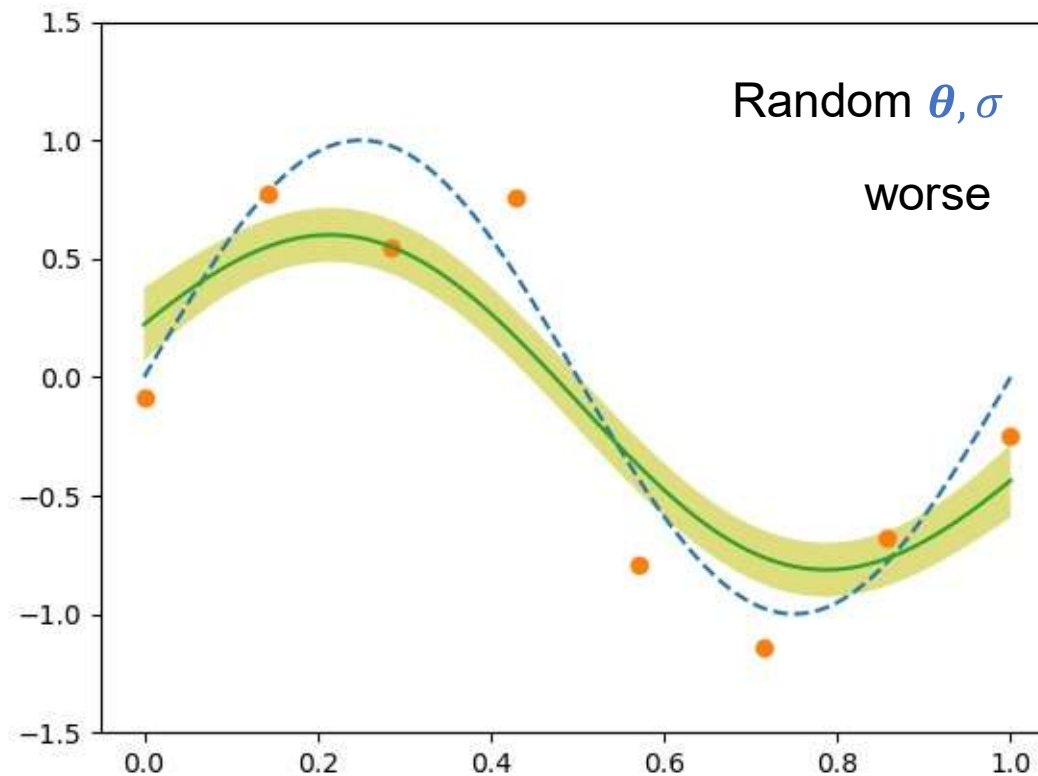
Gaussian Processes

Examples ($\hat{\theta}, \hat{\sigma}$ by MLE)



Gaussian Processes

Examples (random θ, σ without learning from the data)

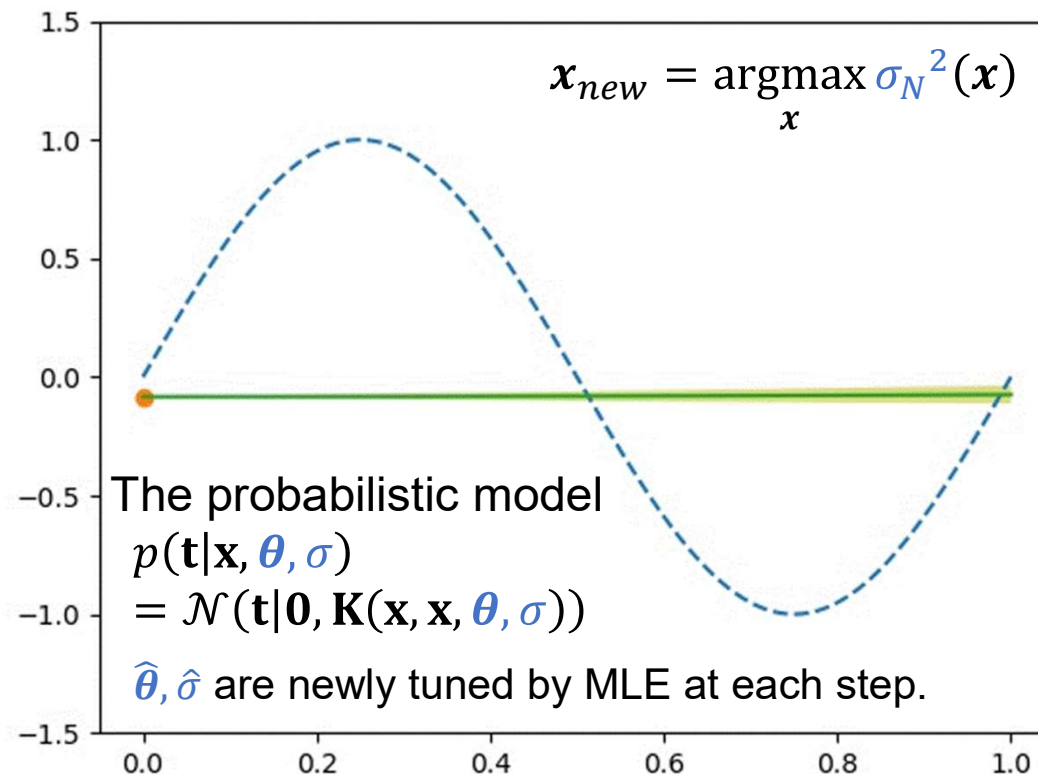


Gaussian Process + adaptive sampling

Examples

Adding a new point at the location x where $\sigma_N^2(x)$ is max

Exactly the same concept as slides 44-60 in Lecture 6

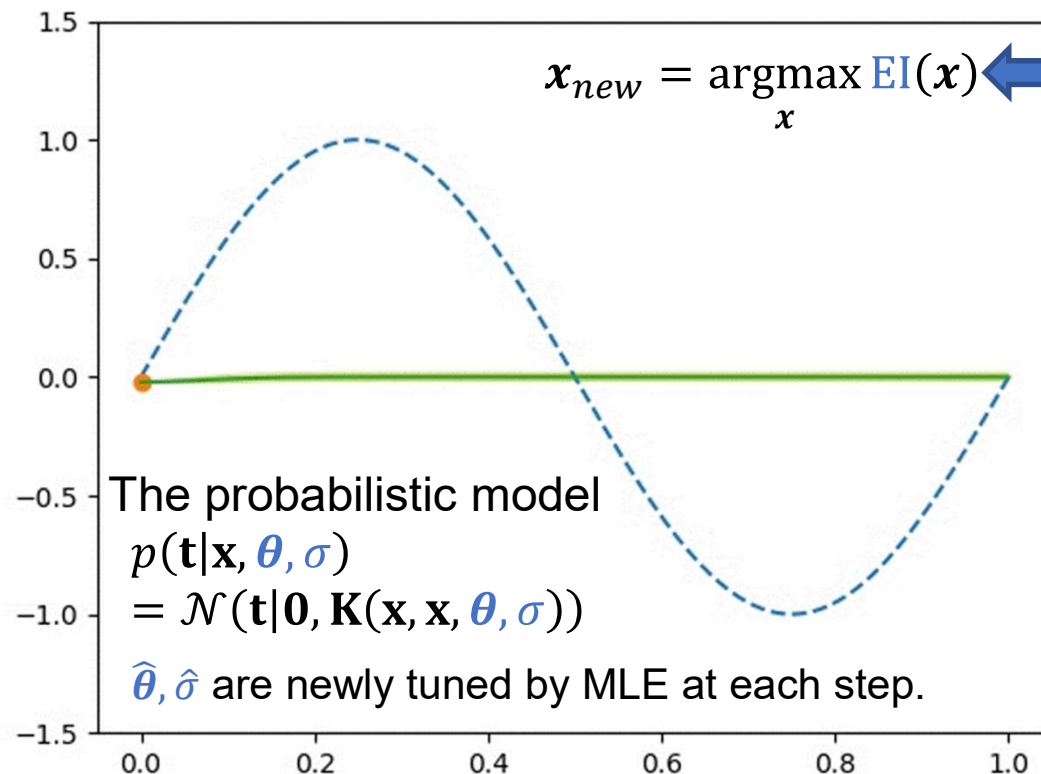


$$\sigma_N^2(x) = k(x, x, \hat{\theta}, \hat{\sigma}) - k(x, \hat{\theta}, \hat{\sigma})^T \mathbf{K}(\hat{\theta}, \hat{\sigma})^{-1} k(x, \hat{\theta}, \hat{\sigma})$$

Gaussian Process + adaptive sampling (for Optimization)

Examples

Adding a new point at the location x **where** a function $EI(x) = f(\mu(x), \sigma_N^2(x), \text{current minimum sample point})$ is **max**



a suitable criterion
to find optimum

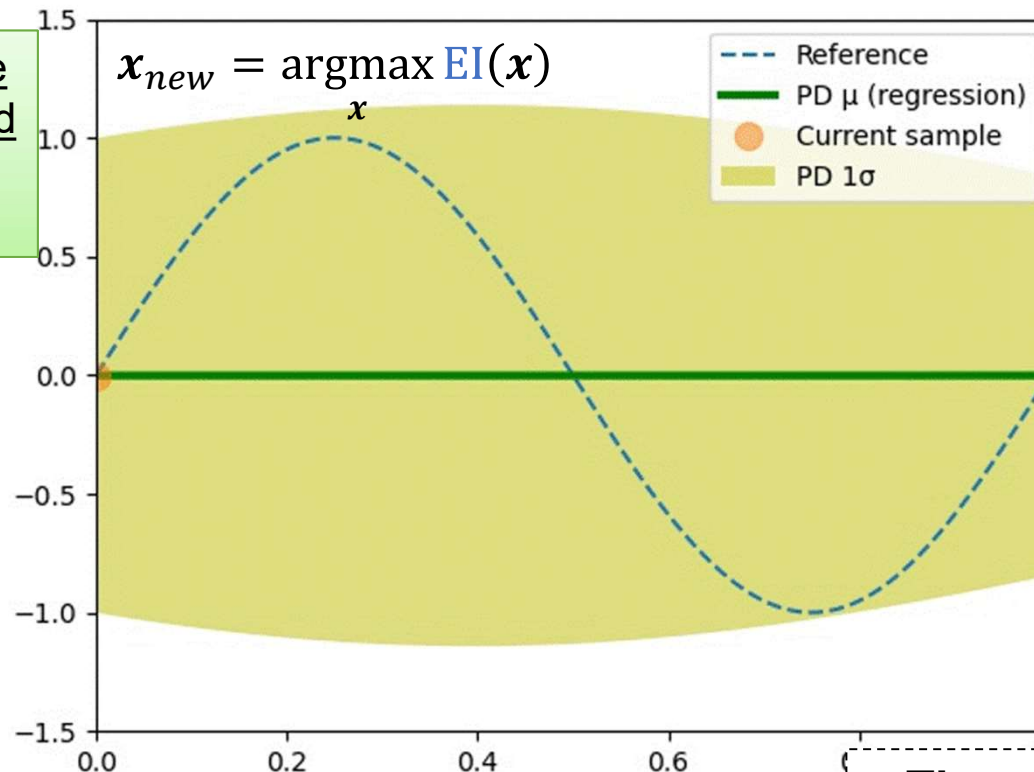
$$EI(x) = f(\mu(x), \sigma_N^2(x), \text{current minimum sample point}) \quad (\mu(x) = \mathbf{k}(x)^T \mathbf{K}(\hat{\boldsymbol{\theta}}, \hat{\sigma})^{-1} \mathbf{T})$$

Bayesian Linear Regression + adaptive sampling (for Opt.)

Examples

Adding a new point at the location x **where** a function $EI(x) = f(\mu(x), \sigma_N^2(x), \text{current minimum sample point})$ is **max**

Of course, the same concept can be used in the Bayesian linear regression.



$$EI(x) = f(\mu(x), \sigma_N^2(x), \text{current minimum sample point})$$

The probabilistic model

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

$\boldsymbol{\phi}(x)$ = polynomials ($M = 9$)

$$\hat{\sigma} = 1e - 3 \text{ (fixed)}$$



Bayesian Linear Regression

Examples

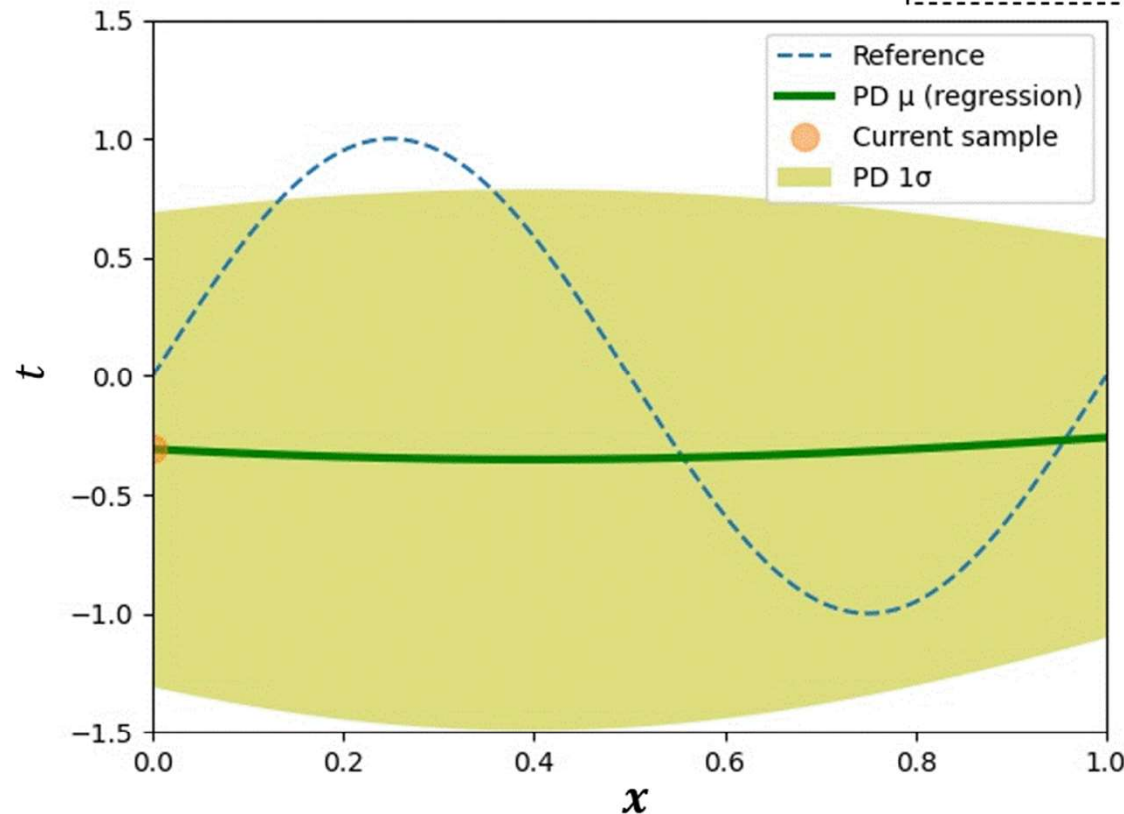
Adding a new point at the location x where $\sigma_N'^2(x)$ (or $\sigma_N^2(x)$) is max

The probabilistic model

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

$\boldsymbol{\phi}(x)$ = polynomials ($M = 9$)

$\hat{\sigma} = 0.2$ (fixed)



starting from N (sample size) = 1

Adaptive Sampling

called acquisition functions in general

1. Variance (Entropy)

$$\mathbf{x}_{new} = \operatorname{argmax}_x \sigma_N^2(\mathbf{x})$$

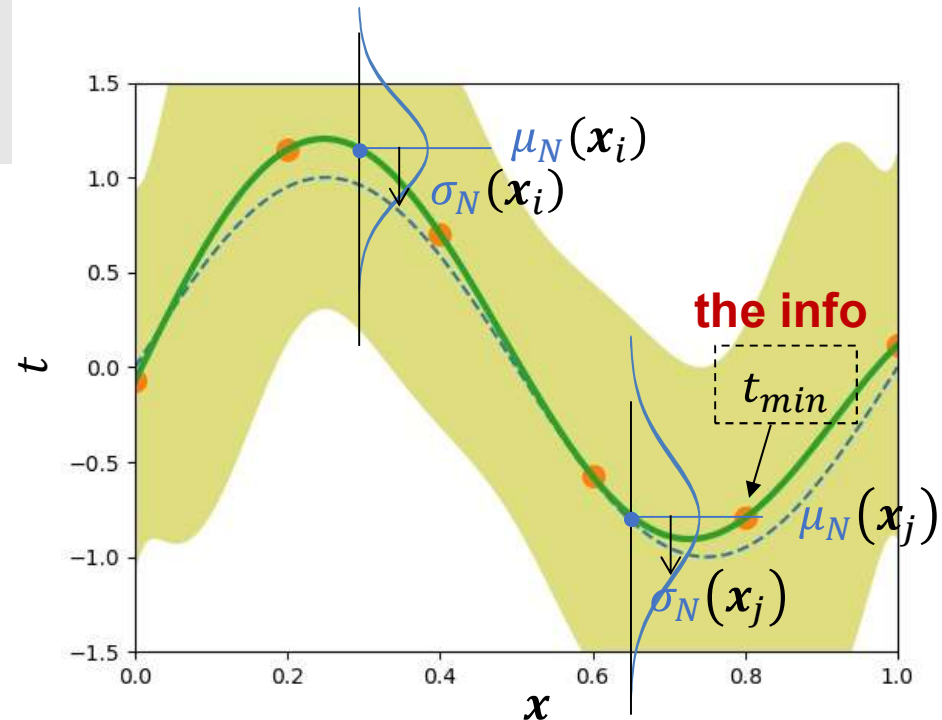
2. Expected Improvement (EI)

$$\mathbf{x}_{new} = \operatorname{argmax}_x \text{EI}(\mathbf{x})$$

$$\text{EI}(\mathbf{x}) = (t_{min} - \mu_N(\mathbf{x}))\text{cdf}(p_{std}(\mathbf{x})) + \sigma_N(\mathbf{x})\text{pdf}(p_{std}(\mathbf{x}))$$
$$p_{std}(\mathbf{x}) = \frac{t_{min} - \mu_N(\mathbf{x})}{\sigma_N(\mathbf{x})}$$

standardized

If you are interested, you can follow the meaning of the equation and can understand it as a probability where min location is expected to be updated.



The predictive distribution (a single input \mathbf{x})

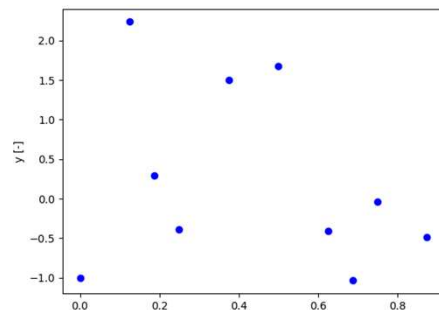
the info

$$p(t|\mathbf{x}) = \mathcal{N}(t|\mu_N(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

History of Gaussian Process and Kriging

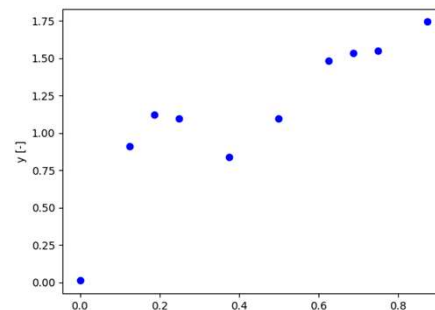
Gaussian Process

“Gaussian process” have been modeled and developed to describe random motions in physics.



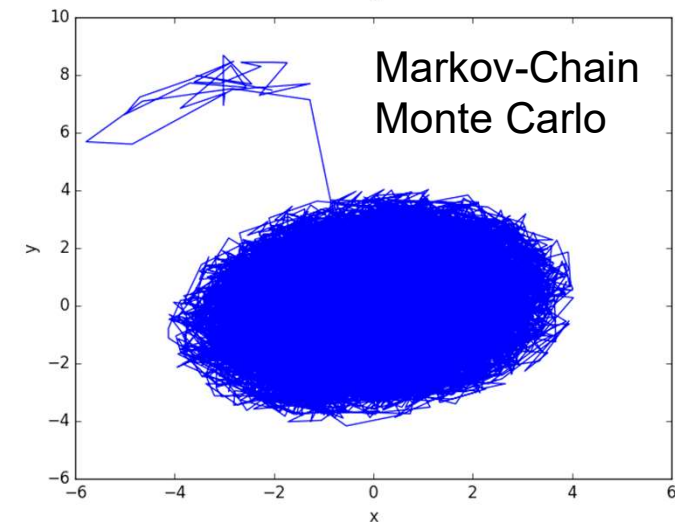
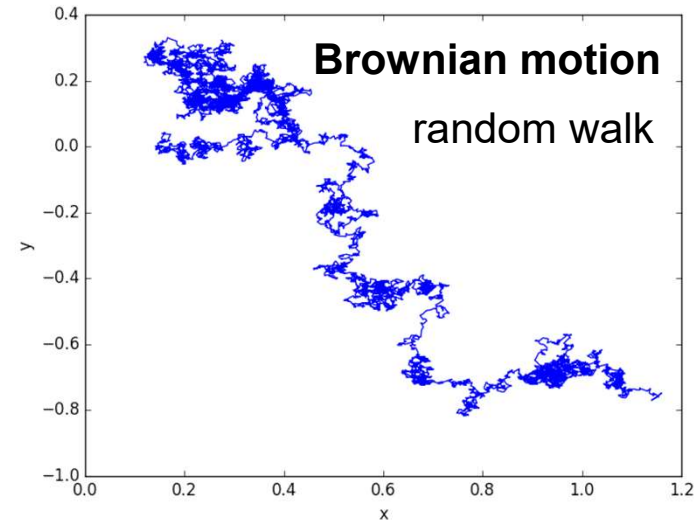
t [time]

$$p(y) = \mathcal{N}(y|0, \sigma^2)$$



t [time]

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K})$$



History of Gaussian Process and Kriging

Kriging

- Spatial statistics
- Geostatistics

to interpolate the value of a random field at an unobserved location from observations of its value at nearby locations

$$\tilde{t}(x) = \mathbf{w}(x)^T \mathbf{T}$$

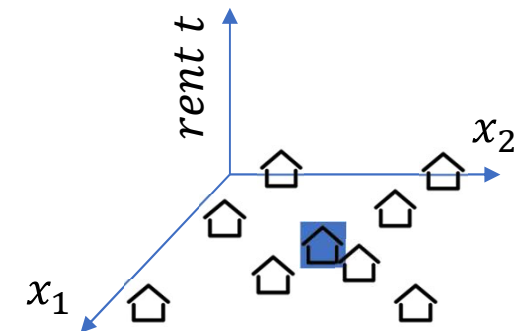
Example:

We want to predict the rent of an apartment.



using:

- The known rents of apartments around it
- Spatial data (in this case the location \mathbf{X})



Spatial autocorrelation: “*closer to each other, closer the properties are.*”

$$\hat{\mathbf{w}}(x) = \min_{\mathbf{w}(x)} MSE \quad MSE = E \left[(t(x) - \tilde{t}(x))^2 \right] \quad \Rightarrow \quad \tilde{t}(x) = \hat{\mathbf{w}}(x)^T \mathbf{T}$$

Prediction in Dual Representation

$$p(\mathbf{t}|\mathbf{x}, \mathcal{D}) = \mathcal{N}(\mathbf{t} | \underbrace{\mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{T}}_{\chi(\mathbf{x})^T}, \underbrace{\mathbf{K}(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})}_{\chi(\mathbf{x})^T})$$

$$\mathcal{D} = (\mathbf{X}, \mathbf{T})$$

Gaussian Process $\tilde{\mathbf{t}}(\mathbf{x}) = \underbrace{\mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{T}}_{\chi(\mathbf{x})^T}$

Kriging

$$\tilde{\mathbf{t}}(\mathbf{x}) = \chi(\mathbf{x})^T \mathbf{T}$$

(Bayesian) linear regression

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

$$\tilde{\mathbf{t}}(\mathbf{x}) = \boldsymbol{\Phi}(\mathbf{x}) \mathbf{w}$$

Prediction using
the training data $\mathcal{D} = (\mathbf{X}, \mathbf{T})$

Prediction using
the learned parameters \mathbf{w}

$$p(\mathbf{t}|\mathbf{x}, \mathcal{D}) = \mathcal{N}(\mathbf{t} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \chi(\mathbf{x})^T \mathbf{T}$$

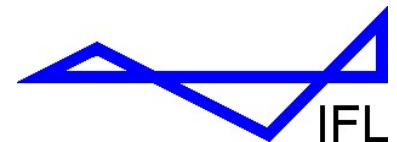
$$\boldsymbol{\Sigma} = \mathbf{K}(\mathbf{x}, \mathbf{x}) - \chi(\mathbf{x})^T \mathbf{k}(\mathbf{x})$$

Note: The covariance can be slightly different in Kriging.

The training data $\mathcal{D} = (\mathbf{X}, \mathbf{T})$ can be discarded in the prediction process.

Lecture content

- Applications of Gaussian processes



Optimization in Engineering Problems

Optimization formulation:

$$\min_x y(x)$$

➡ \hat{x} and $\hat{y} \equiv y(\hat{x})$

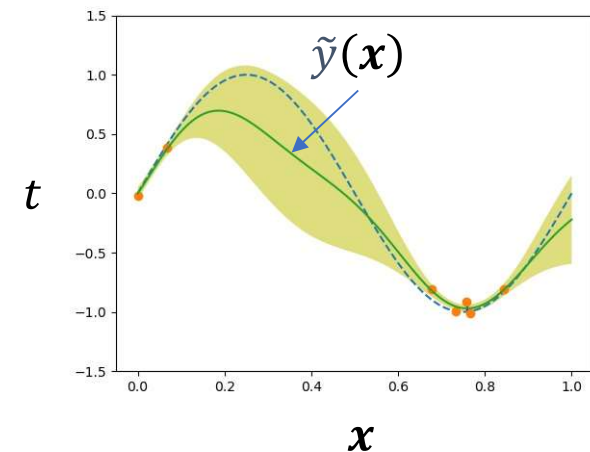
$$\min_x \tilde{y}(x) \quad \tilde{y} \text{ is approximation.}$$

x : input

t : output •

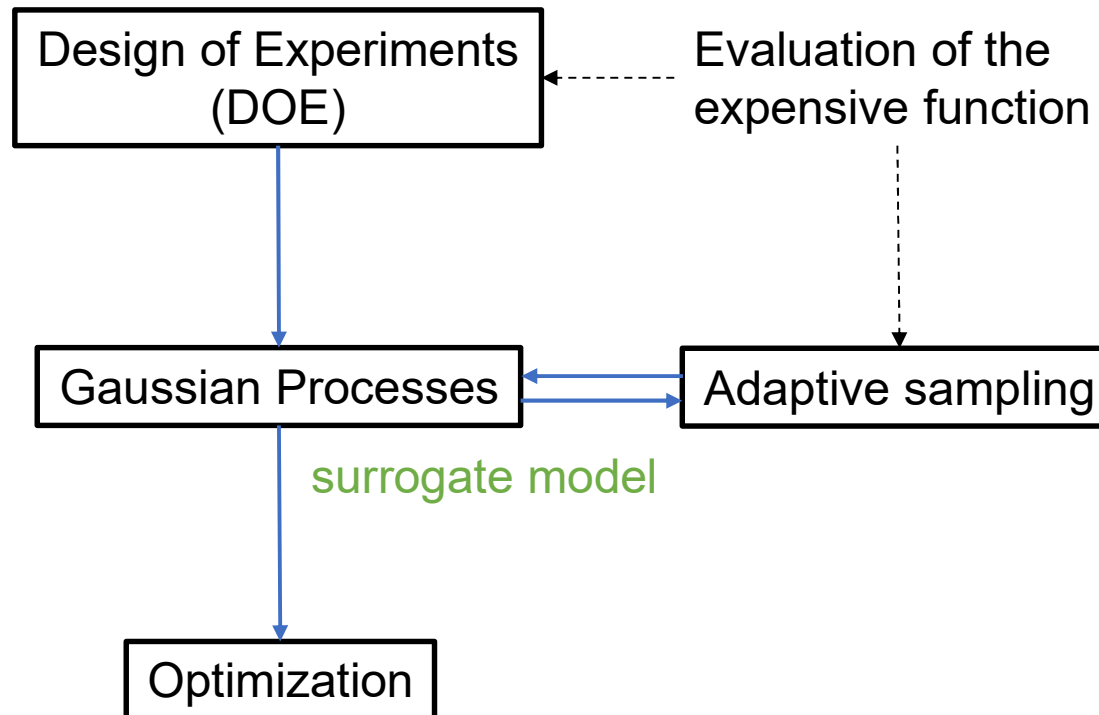
Conditions:

$y(x)$ for arbitrary x is expensive to obtain
(e.g. experimental data).

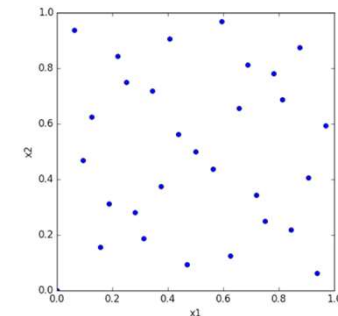


Find \hat{x} accurately but using as small numbers of sample size N as possible

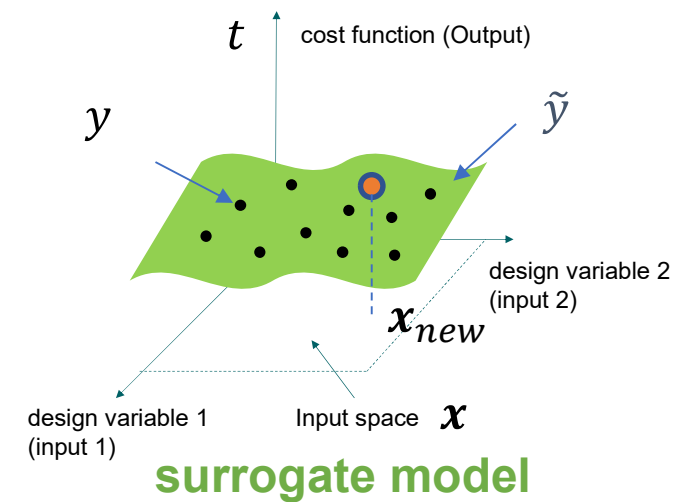
Surrogate-Based Optimization



DoE (see Lecture 4)



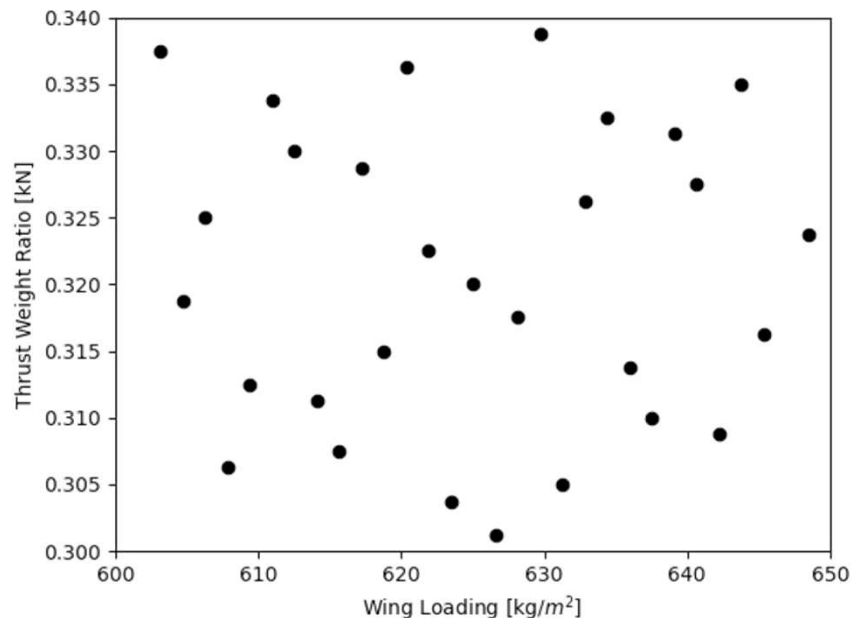
Sampling uniformly in the input space x



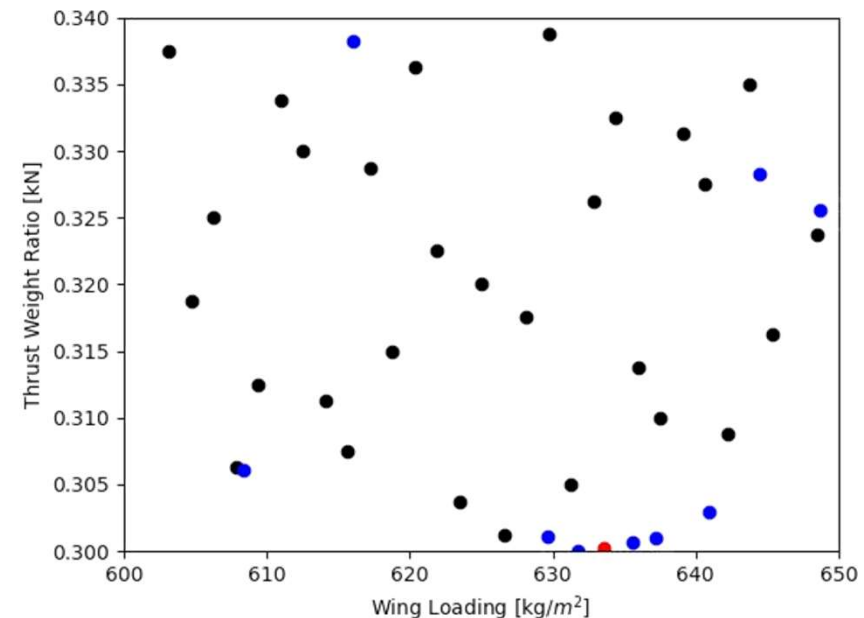
Surrogate-Based Optimization

DOE: Quasi Monte Carlo – Sobol sequence

initial DoE sample



initial DoE sample
with adaptive sample points



Sample points in **input design variables** x space

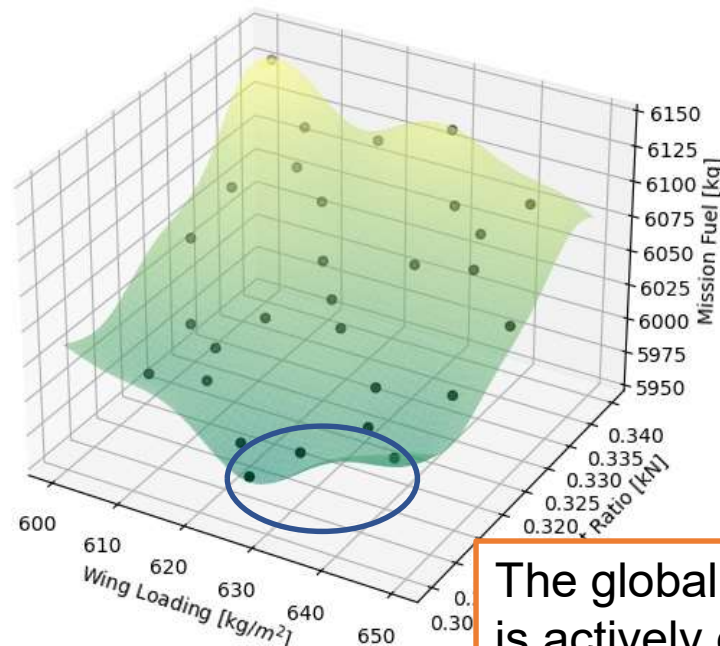
Wing loading $\in [600, 650]$

Thrust weight ratio $\in [0.3, 0.34]$

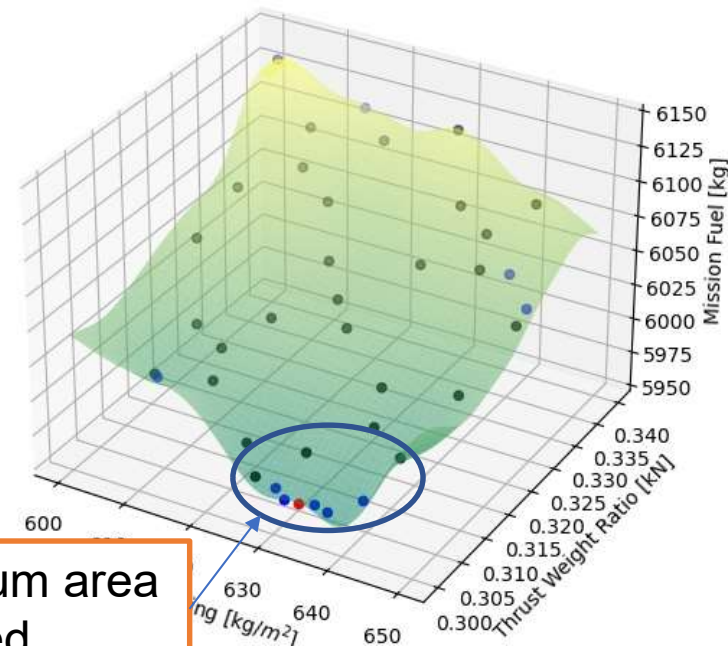
- initial DoE sample
- adaptive sample points
- the optimum point

Surrogate-Based Optimization

initial DoE sample



initial DoE sample
with **adaptive sample points**



The global minimum area
is actively exploited.

The corresponding **cost function values** $y(x)$ on the sample points x ,
and **surrogate models** $\tilde{y}(x)$ (by Gaussian Processes)

Technical Issues

A multivariate Gaussian distribution

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

↑
covariance

$$\begin{aligned} \mathbf{\Sigma}(\mathbf{T}, \mathbf{T}) &= \begin{pmatrix} \text{cov}[t(\mathbf{x}_1), t(\mathbf{x}_1)] & \cdots & \text{cov}[t(\mathbf{x}_1), t(\mathbf{x}_N)] \\ \vdots & \ddots & \vdots \\ \text{cov}[t(\mathbf{x}_N), t(\mathbf{x}_1)] & \cdots & \text{cov}[t(\mathbf{x}_N), t(\mathbf{x}_N)] \end{pmatrix} \\ &= \sigma_T^2 \begin{pmatrix} \text{cor}[t(\mathbf{x}_1), t(\mathbf{x}_1)] & \cdots & \text{cor}[t(\mathbf{x}_1), t(\mathbf{x}_N)] \\ \vdots & \ddots & \vdots \\ \text{cor}[t(\mathbf{x}_N), t(\mathbf{x}_1)] & \cdots & \text{cor}[t(\mathbf{x}_N), t(\mathbf{x}_N)] \end{pmatrix} \\ &= \sigma_T^2 \underbrace{\begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}}_{\text{represented by } \mathbf{x}} = \sigma_T^2 \mathbf{K} \end{aligned}$$

$\text{cov}[T, T'] = \frac{\text{cor}[T, T']}{\sigma_T \sigma_{T'}}$
 $0 \leq \text{cor}[T, T'] \leq 1$

Gram matrix \mathbf{K} : correlation matrix

represented by \mathbf{x} →

\mathbf{x} should be scaled (standardized).

Technical Issues

A multivariate Gaussian distribution

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

↑
covariance

$$\begin{aligned} \Sigma(\mathbf{T}, \mathbf{T}) &= \begin{pmatrix} \text{cov}[t(\mathbf{x}_1), t(\mathbf{x}_1)] & \cdots & \text{cov}[t(\mathbf{x}_1), t(\mathbf{x}_N)] \\ \vdots & \ddots & \vdots \\ \text{cov}[t(\mathbf{x}_N), t(\mathbf{x}_1)] & \cdots & \text{cov}[t(\mathbf{x}_N), t(\mathbf{x}_N)] \end{pmatrix} \\ &= \sigma_T^2 \begin{pmatrix} \text{cor}[t(\mathbf{x}_1), t(\mathbf{x}_1)] & \cdots & \text{cor}[t(\mathbf{x}_1), t(\mathbf{x}_N)] \\ \vdots & \ddots & \vdots \\ \text{cor}[t(\mathbf{x}_N), t(\mathbf{x}_1)] & \cdots & \text{cor}[t(\mathbf{x}_N), t(\mathbf{x}_N)] \end{pmatrix} \\ &= \sigma_T^2 \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1, \boldsymbol{\theta}) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N, \boldsymbol{\theta}) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1, \boldsymbol{\theta}) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N, \boldsymbol{\theta}) \end{pmatrix} = \sigma_T^2 \mathbf{K}(\boldsymbol{\theta}) \end{aligned}$$

$\text{cov}[T, T'] = \frac{\text{cor}[T, T']}{\sigma_T \sigma_{T'}}$
 $0 \leq \text{cor}[T, T'] \leq 1$

Gram matrix $\mathbf{K}(\boldsymbol{\theta})$: correlation matrix

represented by \mathbf{x}



\mathbf{x} should be scaled (standardized).

Technical Issues

A multivariate Gaussian distribution

$$p(\mathbf{t}|\mathbf{x}) = \mathcal{N}(\mathbf{t}|\mathbf{0}, \sigma^2 \mathbf{K}(\boldsymbol{\theta}))$$

Then, the **likelihood function** when data $\mathcal{D} = (\mathbf{X}, \mathbf{T})$ is given,

$$p(\mathbf{T}|\mathbf{X}) = \mathcal{N}(\mathbf{t}|\mathbf{0}, \sigma^2 \mathbf{K}(\boldsymbol{\theta})) = \frac{1}{(2\pi\sigma)^{N/2} |\mathbf{K}(\boldsymbol{\theta})|^{1/2}} \exp\left(-\frac{\mathbf{T}^T \mathbf{K}(\boldsymbol{\theta})^{-1} \mathbf{T}}{2\sigma^2}\right)$$

negative log as usual \rightarrow **Error function** E

$$E(\sigma, \boldsymbol{\theta}) = -\ln p(\mathbf{T}|\mathbf{X}) = \frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \sigma^2 + \frac{1}{2} \ln |\mathbf{K}(\boldsymbol{\theta})| + \frac{\mathbf{T}^T \mathbf{K}(\boldsymbol{\theta})^{-1} \mathbf{T}}{2\sigma^2}$$

$$\frac{\partial E}{\partial \sigma} = 0, \frac{\partial E}{\partial \boldsymbol{\theta}} = 0$$

Please derive $E(\boldsymbol{\theta})$ by yourself.

cautions
in computation!

$$\hat{\boldsymbol{\theta}} = \min_{\boldsymbol{\theta}} E(\boldsymbol{\theta})$$

$$E(\boldsymbol{\theta}) = \frac{N}{2} \ln \hat{\sigma}(\boldsymbol{\theta})^2 + \frac{1}{2} \ln |\mathbf{K}(\boldsymbol{\theta})| + C$$

where,

$$\hat{\sigma}(\boldsymbol{\theta})^2 = \frac{\mathbf{T}^T \mathbf{K}(\boldsymbol{\theta})^{-1} \mathbf{T}}{N}$$



Useful Kernel Functions

- Gaussian kernel (RBF kernel)

$$k(\mathbf{x}, \mathbf{x}', \theta) = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|^2) \quad \text{or} \quad k(\mathbf{x}, \mathbf{x}', \boldsymbol{\theta}) = \exp\left(-\sum_{i=1}^D \theta_i \|x^{(i)} - x'^{(i)}\|^2\right)$$

- Linear kernel

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- Exponential kernel

$$k(\mathbf{x}, \mathbf{x}', \theta) = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|) \quad \longrightarrow \quad k(\mathbf{x}, \mathbf{x}', \theta, p) = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|^p)$$

- Matérn kernel

$$k_\nu(\mathbf{x}, \mathbf{x}', \theta) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\pi}r}{\theta}\right)^\nu K_\nu\left(\frac{\sqrt{2\pi}r}{\theta}\right)$$

$$r = \|\mathbf{x} - \mathbf{x}'\|$$

when $\nu = \frac{1}{2}$, Exponential kernel
when $\nu = \frac{3}{2}$, Matérn3
when $\nu = \frac{5}{2}$, Matérn5
when $\nu = \infty$, Gaussian kernel

Lecture content

- Advanced Methods/Topics

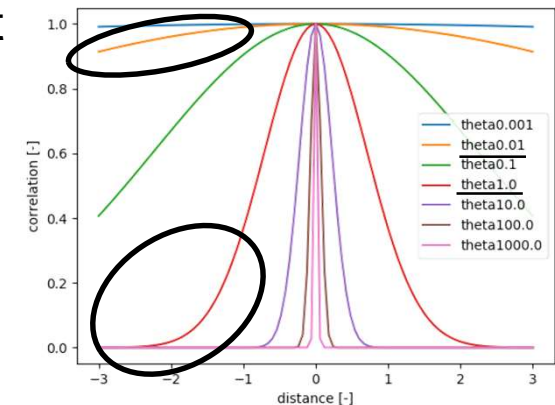


Sensitivity Analysis (Automatic Relevance Determination - ARD)

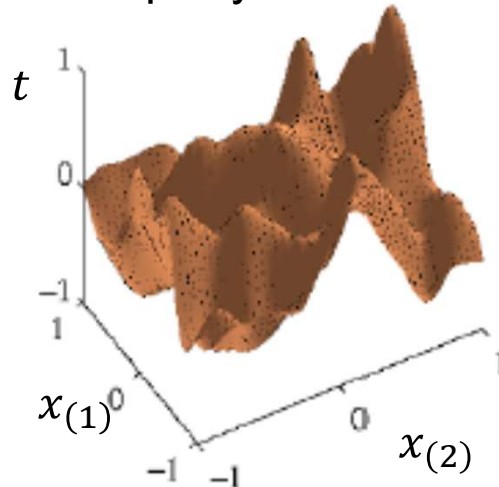
A kernel with hyperparameters at each component of the input

$$k(\mathbf{x}, \mathbf{x}', \boldsymbol{\theta}) = \exp\left(-\sum_{i=1}^2 \theta_i \|x^{(i)} - x'^{(i)}\|^2\right) \quad \text{e.g. 2D case}$$

$\mathbf{x} = (x_{(1)}, x_{(2)})$

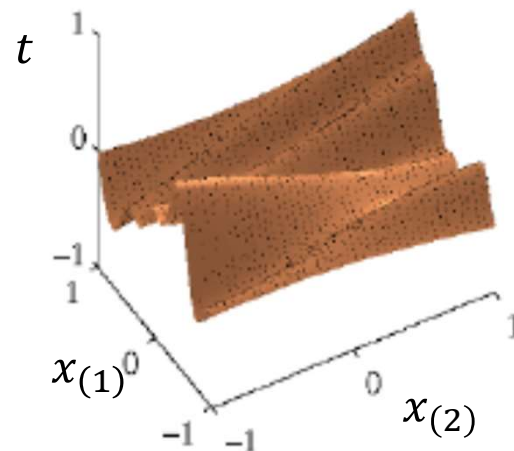


Equally sensitive



$$\boldsymbol{\theta} = (\theta_1, \theta_2) \\ = (1, 1)$$

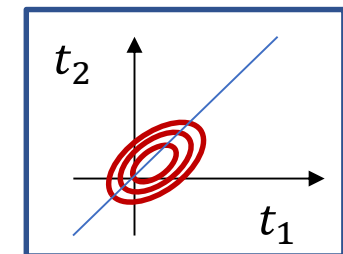
Not that sensitive in $x_{(2)}$



$$\boldsymbol{\theta} = (\theta_1, \theta_2) \\ = (1, 0.01)$$

PRML, p.312

t is correlated strongly in $x_{(2)}$
(even the distance $\|x_i - x_i'\|^2$
is large).



t_1, t_2, t_3, \dots are very similar in $x_{(2)}$.

After MLE with obtaining $\hat{\boldsymbol{\theta}}$, the info of
the sensitivity analysis can be obtained.

Other variations (from Kriging)

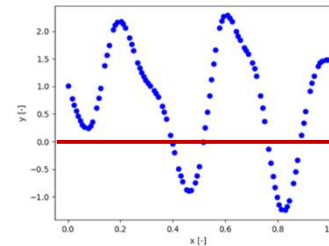
Note:

A little different explanation from that from Kriging in geostatistics

A multivariate Gaussian distribution

$$p(\mathbf{t}|\mathbf{x}) = \mathcal{N}(\mathbf{t}|\boldsymbol{\mu}, \sigma^2 \mathbf{K}(\boldsymbol{\theta}))$$

We can model also the trend.



the trend: $\boldsymbol{\mu}$

When $\boldsymbol{\mu}$ is:

1. $\boldsymbol{\mu} = \mathbf{0}$ (or $\bar{\mathbf{T}}$)
Simple Kriging

2. $\boldsymbol{\mu} = \mathbf{1}\mu$: a constant value
Ordinary Kriging

$$\hat{\mu}(\boldsymbol{\theta}) = (\mathbf{1}^T \mathbf{K}(\boldsymbol{\theta})^{-1} \mathbf{1})^{-1} \mathbf{1}^T \mathbf{K}(\boldsymbol{\theta})^{-1} \mathbf{T}$$

3. $\boldsymbol{\mu} = \boldsymbol{\Phi} \mathbf{w}$: linear regression model
Universal Kriging

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

$\hat{\boldsymbol{\mu}}$ (by MLE) in general has the closed-form (analytically expressed by $\boldsymbol{\theta}$) as **generalized least-squares (GLS)**.

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

$$\boldsymbol{\Sigma} = \sigma^2 \mathbf{K}(\boldsymbol{\theta})$$

when,

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

(see Lecture 4)

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

Least-squares
in linear regression

ordinary least-squares (OLS).

Other variations (from Kriging)

A multivariate Gaussian distribution

$$p(\mathbf{t}|\mathbf{x}) = \mathcal{N}(\mathbf{t}|\boldsymbol{\mu}, \sigma^2 \mathbf{K}(\boldsymbol{\theta}))$$

Basically the same as Slide 25

Then, the **likelihood function** when data $\mathcal{D} = (\mathbf{X}, \mathbf{T})$ is given,

$$p(\mathbf{T}|\mathbf{X}) = \mathcal{N}(\mathbf{t}|\boldsymbol{\mu}, \sigma^2 \mathbf{K}(\boldsymbol{\theta})) = \frac{1}{(2\pi\sigma)^{N/2} |\mathbf{K}(\boldsymbol{\theta})|^{1/2}} \exp\left(-\frac{(\mathbf{T} - \boldsymbol{\mu})^T \mathbf{K}(\boldsymbol{\theta})^{-1} (\mathbf{T} - \boldsymbol{\mu})}{2\sigma^2}\right)$$

negative log as usual \rightarrow **Error function** E

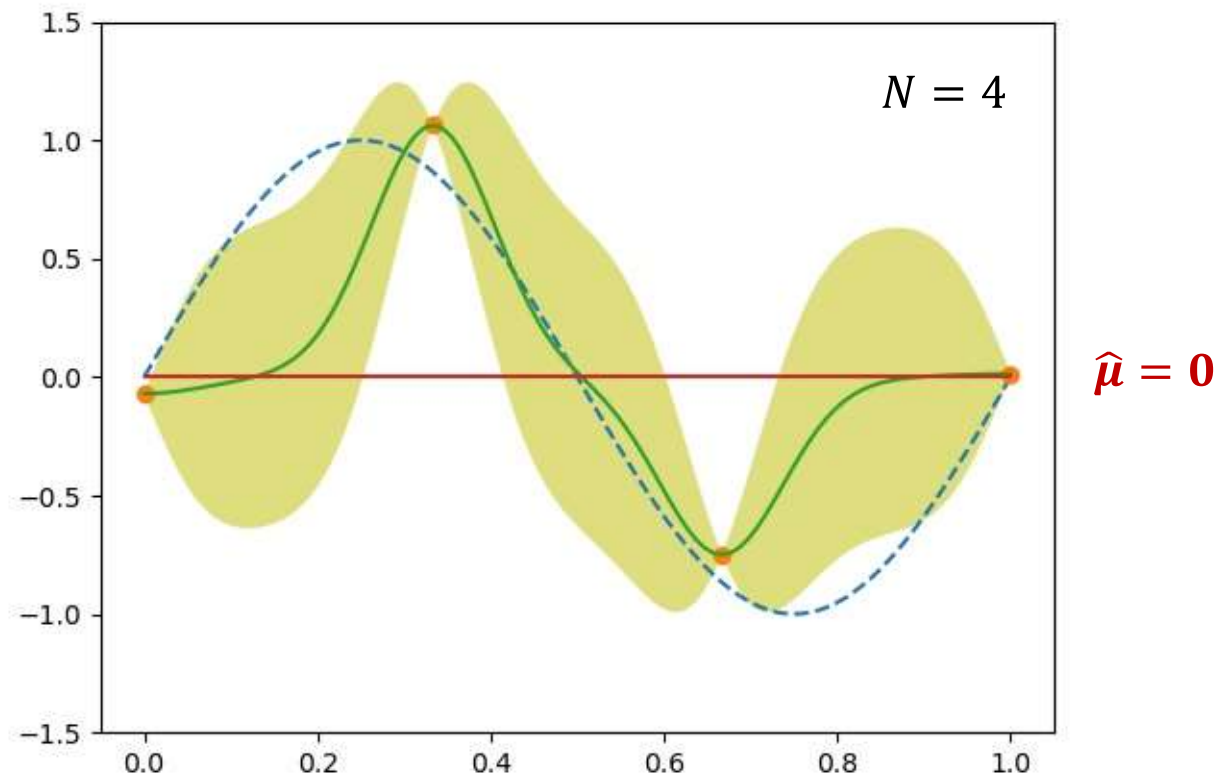
$$E(\boldsymbol{\mu}, \sigma, \boldsymbol{\theta}) = -\ln p(\mathbf{T}|\mathbf{X}) = \frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \sigma^2 + \frac{1}{2} \ln |\mathbf{K}(\boldsymbol{\theta})| + \frac{(\mathbf{T} - \boldsymbol{\mu})^T \mathbf{K}(\boldsymbol{\theta})^{-1} (\mathbf{T} - \boldsymbol{\mu})}{2\sigma^2}$$

$$\frac{\partial E}{\partial \boldsymbol{\mu}} = 0, \frac{\partial E}{\partial \sigma} = 0, \frac{\partial E}{\partial \boldsymbol{\theta}} = 0$$

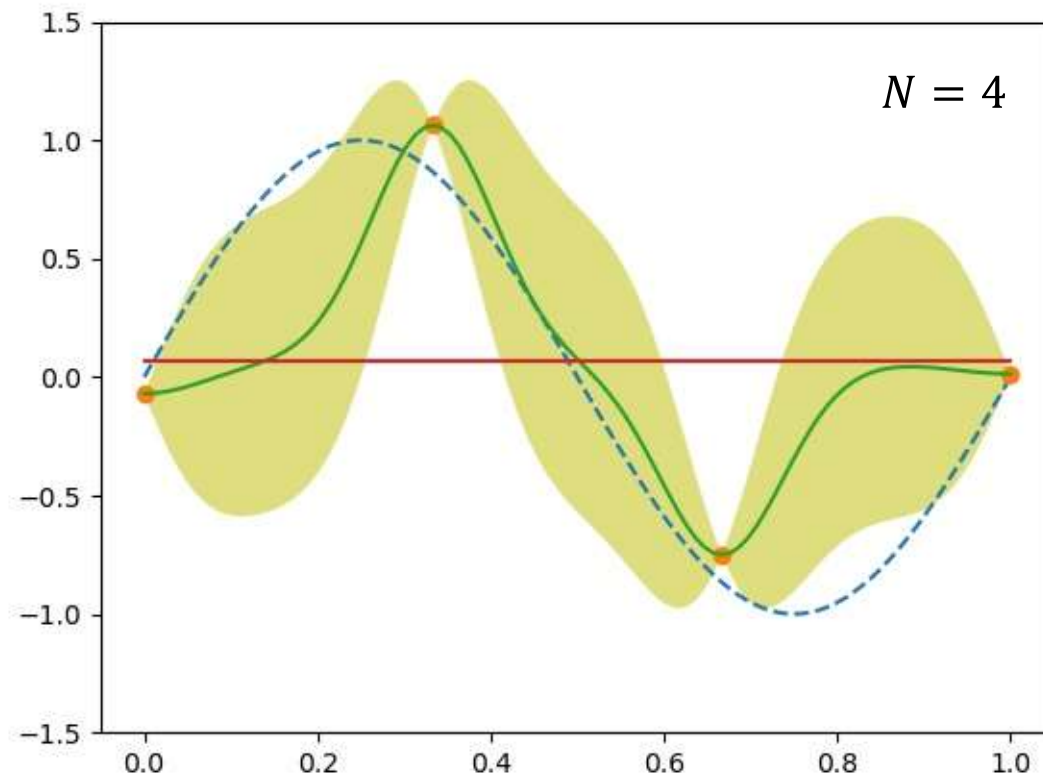
$$\hat{\boldsymbol{\theta}} = \min_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) \quad E(\boldsymbol{\theta}) = \frac{N}{2} \ln \hat{\sigma}(\boldsymbol{\theta})^2 + \frac{1}{2} \ln |\mathbf{K}(\boldsymbol{\theta})| + C$$
$$\hat{\boldsymbol{\mu}} = \text{GLS} \quad \hat{\sigma}(\boldsymbol{\theta})^2 = \frac{(\mathbf{T} - \hat{\boldsymbol{\mu}})^T \mathbf{K}(\boldsymbol{\theta})^{-1} (\mathbf{T} - \hat{\boldsymbol{\mu}})}{N}$$



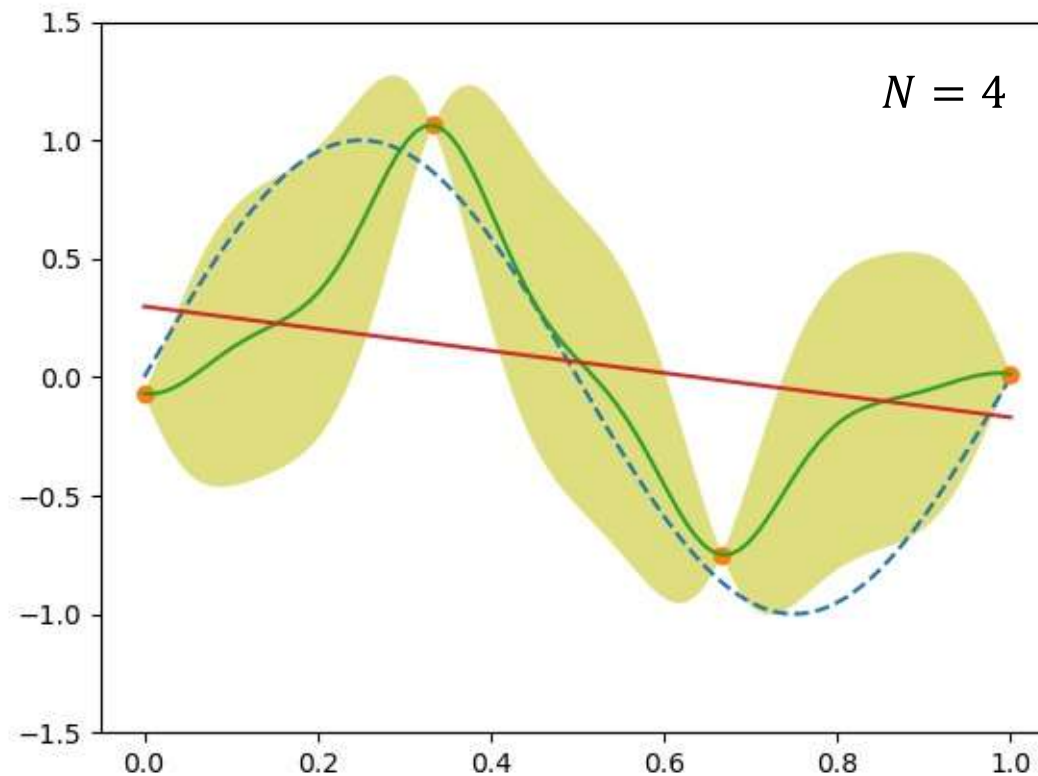
Simple Kriging (Normal Gaussian Process)



Ordinary Kriging



Universal Kriging



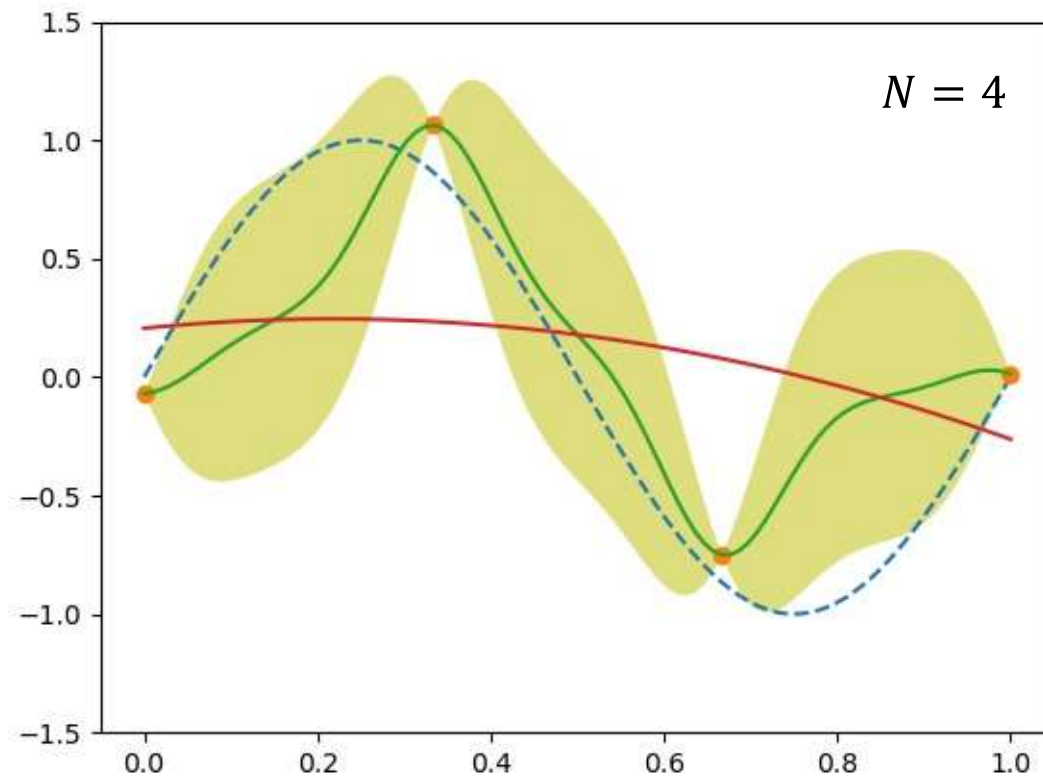
$$\hat{\mu} = \Phi \hat{w}$$

$$\phi(x) = (x^0, x^1)^T$$

Trend:
linear function



Universal Kriging



$$\hat{\mu} = \Phi \hat{w}$$

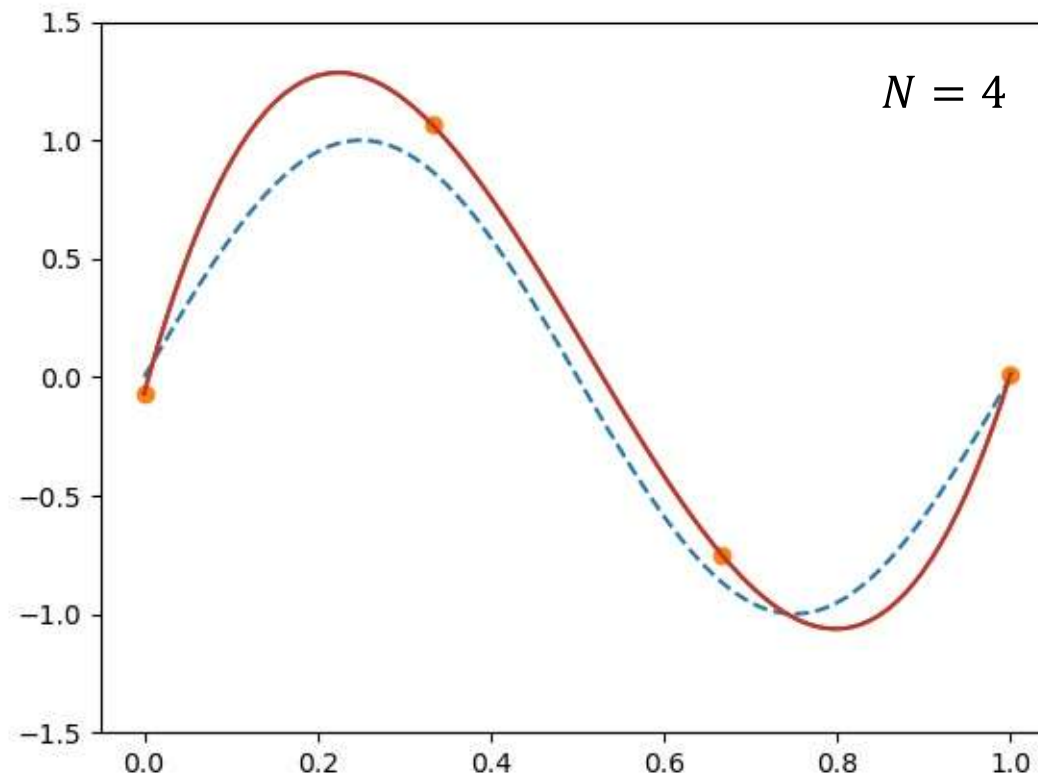
$$\phi(x) = (x^0, x^1, x^2)^T$$

Trend:
quadratic function



Universal Kriging

If you know the characteristic of the trend (it is a cubic function in this case), the universal Kriging works better (not always).

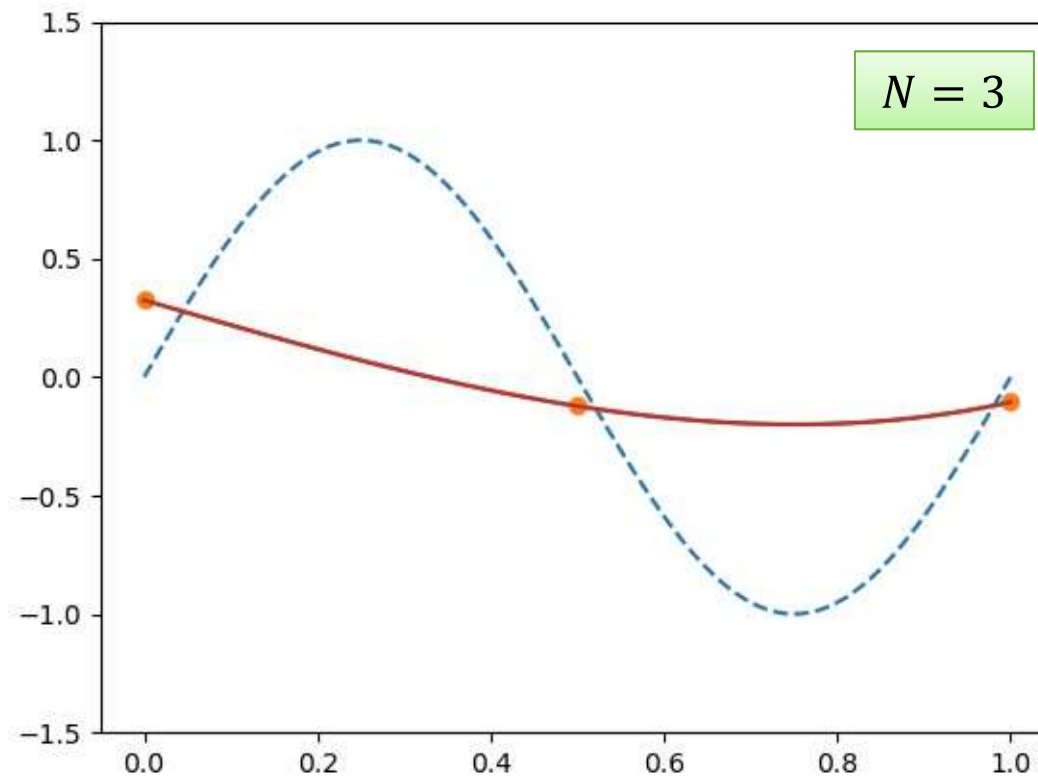


$$\hat{\mu} = \Phi \hat{w}$$

$$\phi(x) = (x^0, x^1, x^2, x^3)^T$$

Trend:
cubic function

Universal Kriging (3/3)



$$\hat{\mu} = \Phi \hat{w}$$

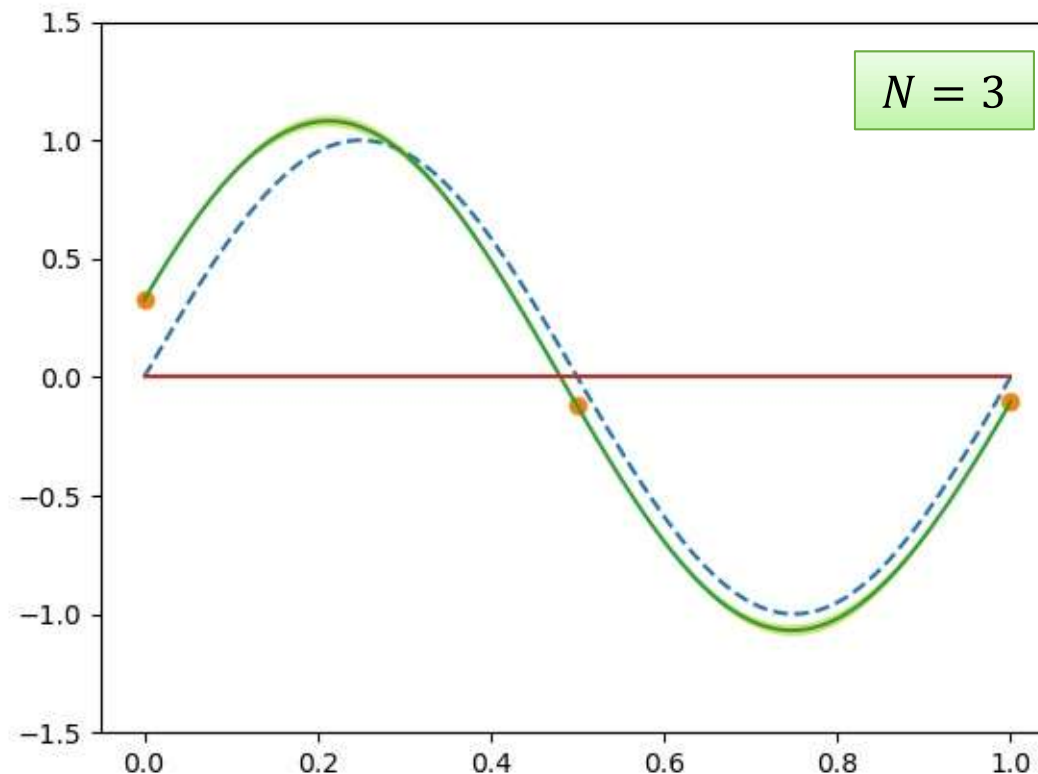
$$\phi(x) = (x^0, x^1, x^2, x^3)^T$$

Trend:
cubic function



Gradient Enhanced

Instead of adding new sample points, using **the gradient information** of the existing sample points



$$\hat{\mu} = 0$$

(based on
Simple Kriging
in this case)

Gradient Enhanced

when the kernel function is differentiable w.r.t. x

e.g. $k(x, x', \theta) = \exp\left(-\sum_{i=1}^D \theta_i \|x^{(i)} - x'^{(i)}\|^2\right)$

Instead of obtaining
new sample points

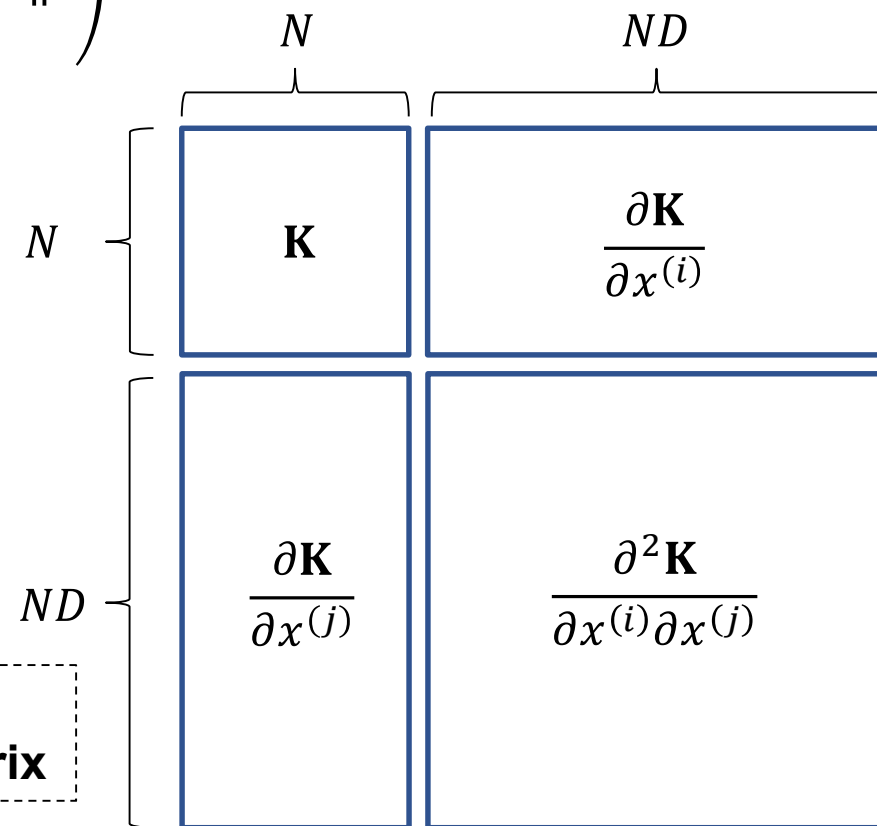
$$\dot{\mathbf{K}} = \begin{pmatrix} \mathbf{K} & \frac{\partial \mathbf{K}}{\partial x^{(i)}} \\ \frac{\partial \mathbf{K}}{\partial x^{(j)}} & \frac{\partial^2 \mathbf{K}}{\partial x^{(i)} \partial x^{(j)}} \end{pmatrix} \quad i, j = 1, \dots, D$$

the size of the matrix

\mathbf{K} : $N \times N$

$\dot{\mathbf{K}}$: $N(1 + D) \times N(1 + D)$

Computational time (during MLE):
cubic proportional to the size of the matrix



Gradient Enhanced

Practical applications:

- When the computations of the gradients are relatively cheap:

- **Automatic differentiation** (the chain rule)

- Forward mode

- Reverse mode

- e.g. single output t , multiple input x

Topics of “calculus”

Note:

AD by the reverse mode is also used in neural networks for different purposes (Lecture 10).

Examples:

- adjoint solvers in computational fluid dynamics (CFD)

If there are tools to efficiently compute the gradients, the gradient-enhanced model is very useful.

Pros and Cons (Gaussian Process / Kriging in General)

Pros

- Weak assumption on using Bayesian approach
 - The model uncertainty can be analytically obtained.



Due to this property, the Gaussian processes are **powerful when data is expensive to obtain.**

In contrast to Big Data

Cons

- Computational cost in large datasets - $O(N^3)$
 - (topics of “linear algebra” and “computer science”)



within 1 sec. to days!

Especially in the cases of high-dimensional input x , the sample size N tends to large.

Technical Issues on Implementation (1/2)

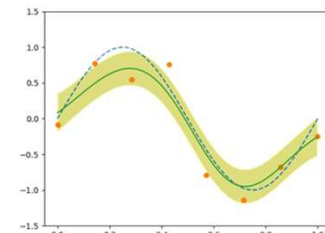
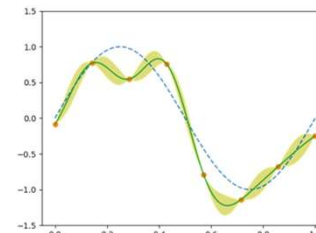
Computation of the likelihood function (error function)

$$E(\boldsymbol{\theta}) = \frac{N}{2} \ln \hat{\sigma}(\boldsymbol{\theta})^2 + \frac{1}{2} \ln |\mathbf{K}(\boldsymbol{\theta})| + C \quad \hat{\sigma}(\boldsymbol{\theta})^2 = \frac{\mathbf{T}^T \mathbf{K}(\boldsymbol{\theta})^{-1} \mathbf{T}}{N}$$

- Effectively use “LU decomposition” and it results
- Use double precision
- If numerically the above process is unstable, a fixed value can be added on the diagonal of $\mathbf{K}(\boldsymbol{\theta})$.

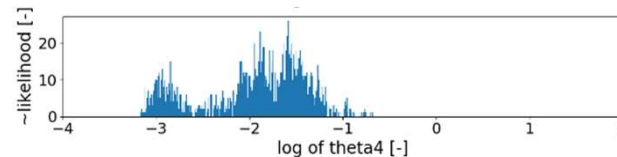
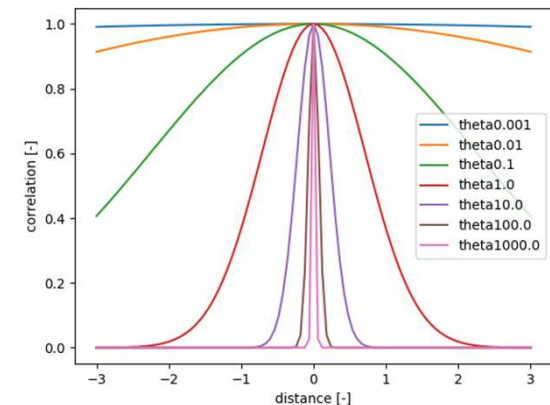
➡ Theoretically this is equivalent to avoid overfitting for noisy data.

See Lecture 8, slides 30-34



Technical Issues on Implementation (2/2)

- The hyperparameter θ should be searched in the **log-scale**.
- The likelihood function (error function) $E(\theta)$ is multimodal.
 - $\frac{\partial E(\theta)}{\partial \theta}$ can be analytically obtained but **global optimization methods** is useful.



MCMC was used to make this fig.

- The input x should be standardized ($0 \leq x \leq 1$).

A common kernel function can be used and ARD can be used.

Technical Issues on Implementation - Advanced

Advanced topics (**research level**):

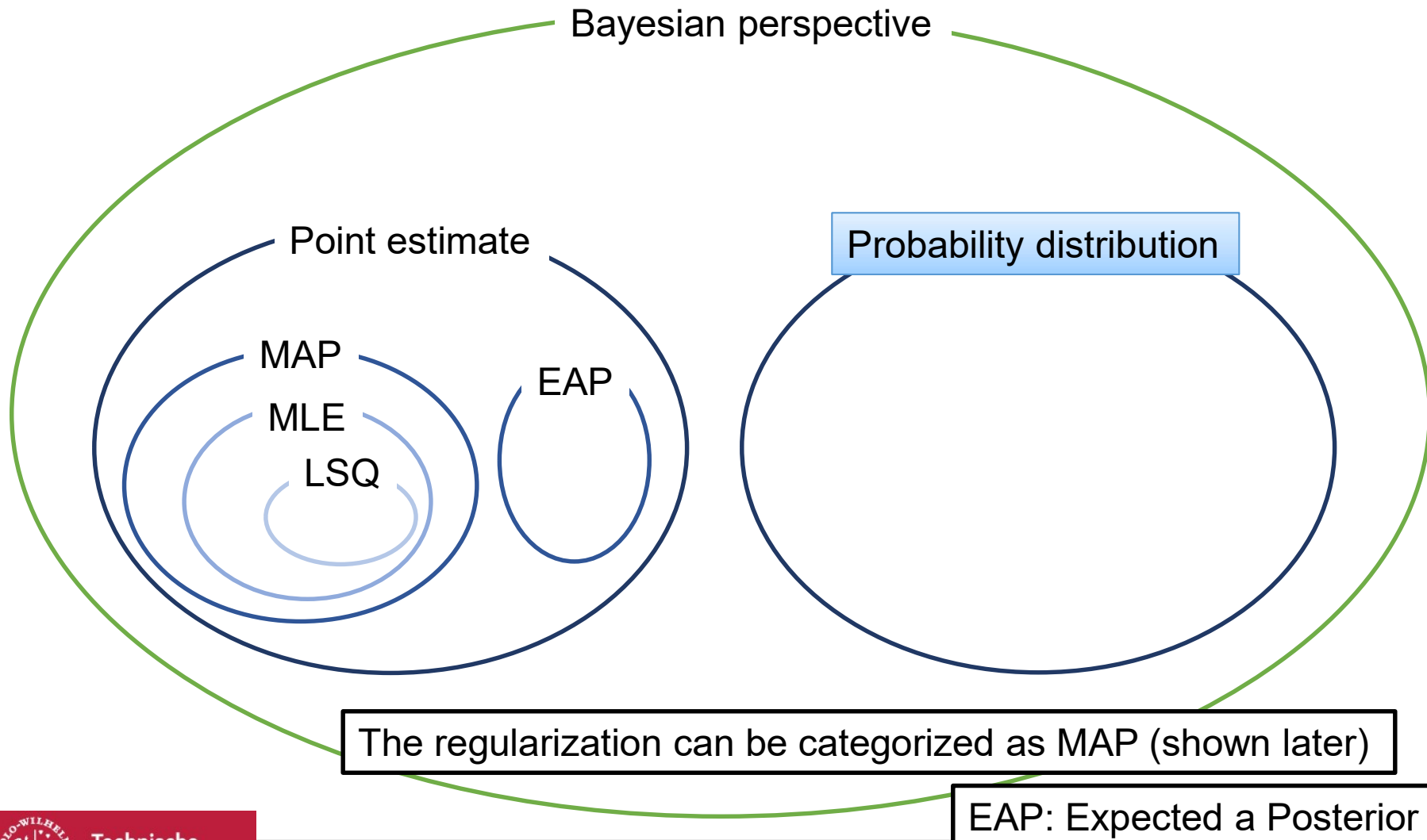
- “Cholesky decomposition” (instead of LU decomposition) for further speed-up
- There are various techniques to reduce the computational costs in large datasets.
 - Basically approximation

Other Advanced Methods

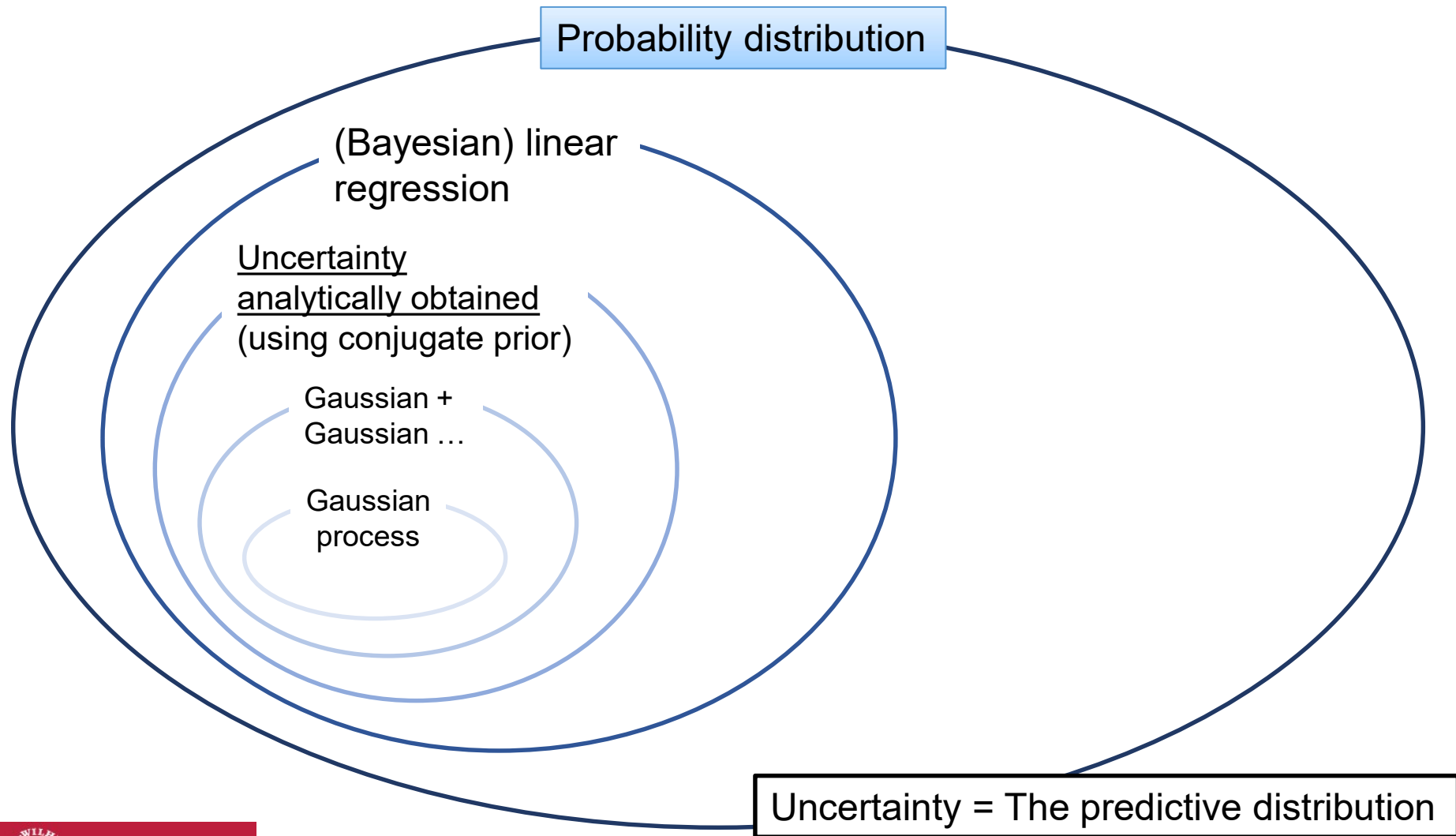
Other advanced topics that were not shown here:

- Cokriging (many variations in different communities)
 - Variable-fidelity methods
- Non-stationary kernel
- Regression Kriging (essentially the same as Universal Kriging)
- Deep Gaussian process

Bayesian Approach – Generalized Perspective (Lecture 6)



Uncertainty (to obtain the predictive distribution)



Uncertainty (to obtain the predictive distribution)

