

Scientific Machine Learning

Lecture 6: Bayesian Statistics (2/2) – Bayesian Linear Regression

Dr. Daigo Maruyama

Prof. Dr. Ali Elham

Where are we going now?

We are going to learn:

- **Gaussian Processes** →
- **Neural Networks** →

If one sentence is used to explain them:

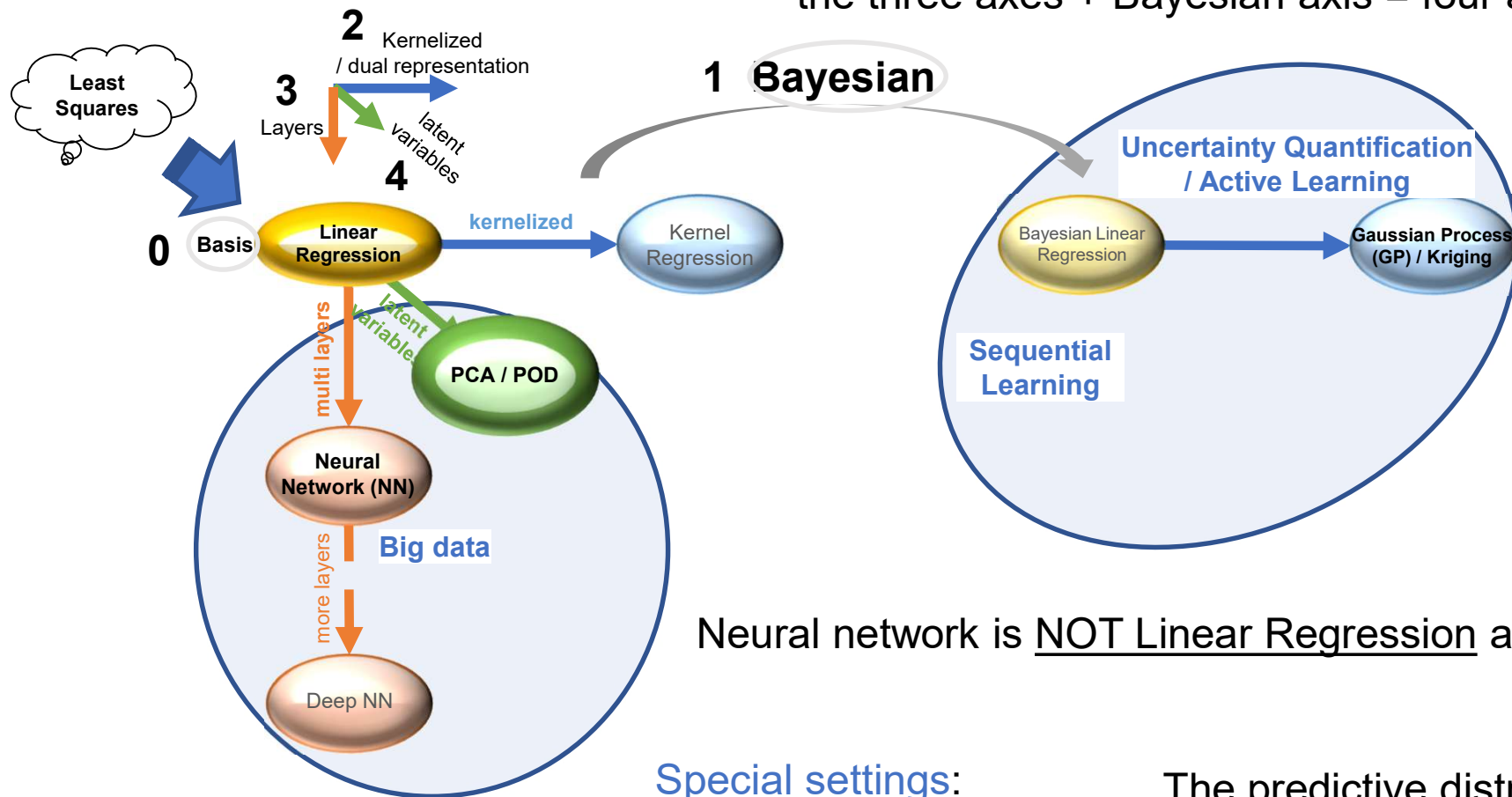
The probabilistic model is
a multivariate Gaussian distribution.

Nonlinear regression

by learning tools now.

Key Components

the three axes + Bayesian axis = four axes



Neural network is NOT Linear Regression anymore.

Special settings:

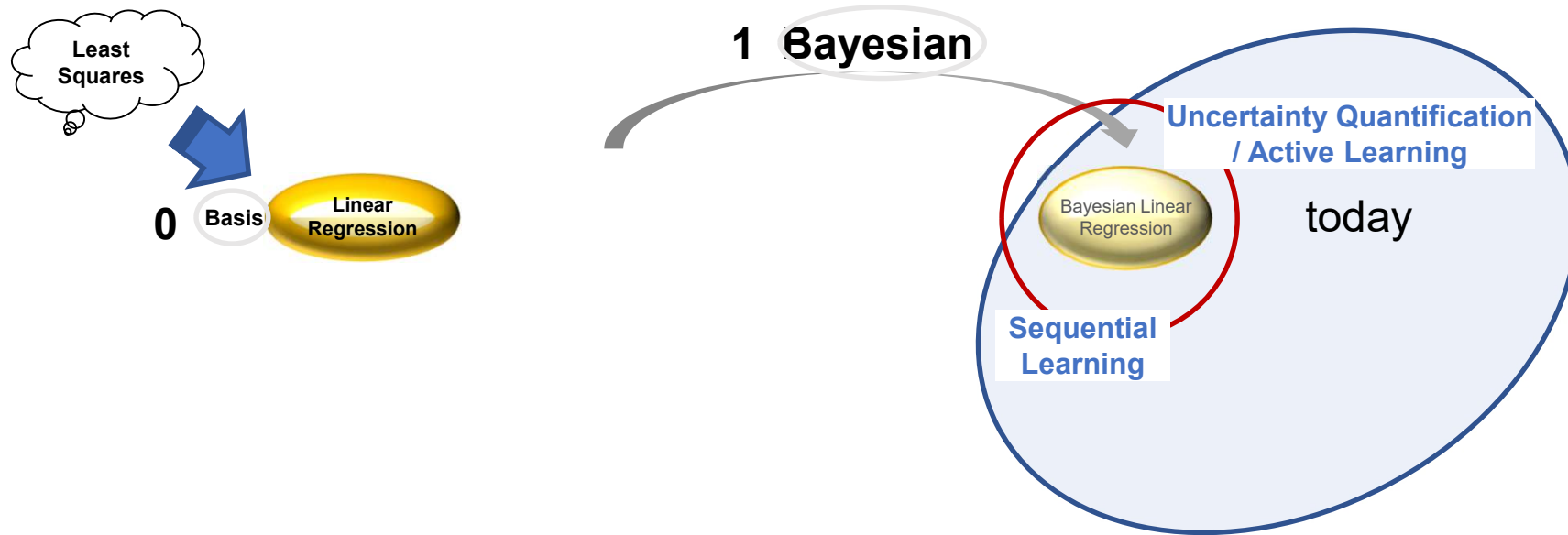
- Conjugate prior
- Linear regression



The predictive distribution (the goal) can be analytically obtained.

This can be also a preparation of Gaussian Process.

Current Position



Special settings:

- Conjugate prior
- Linear regression



The predictive distribution (the goal) can be analytically obtained.

This can be also a preparation of Gaussian Process.



Technische
Universität
Braunschweig

Dr. Daigo Maruyama |



Lecture content

- Bayesian approach (review)
- Bayesian linear regression
- Bayesian sequential learning
- Uncertainty due to data (active learning)

The lecture of this time basically follows the Sections 3.3 and 2.3.6 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.

Dr. Daigo Maruyama | Scientific Machine Learning: Lecture 6 | Slide 5

Lecture content

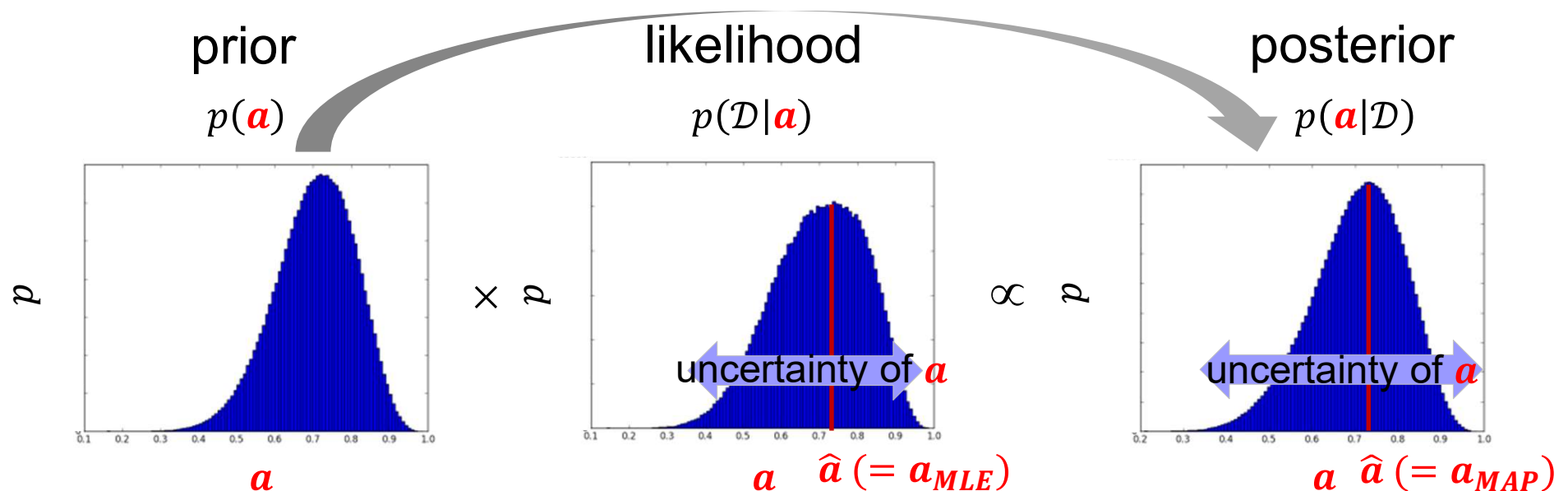
- Bayesian approach (review)



Bayesian Approach - Review

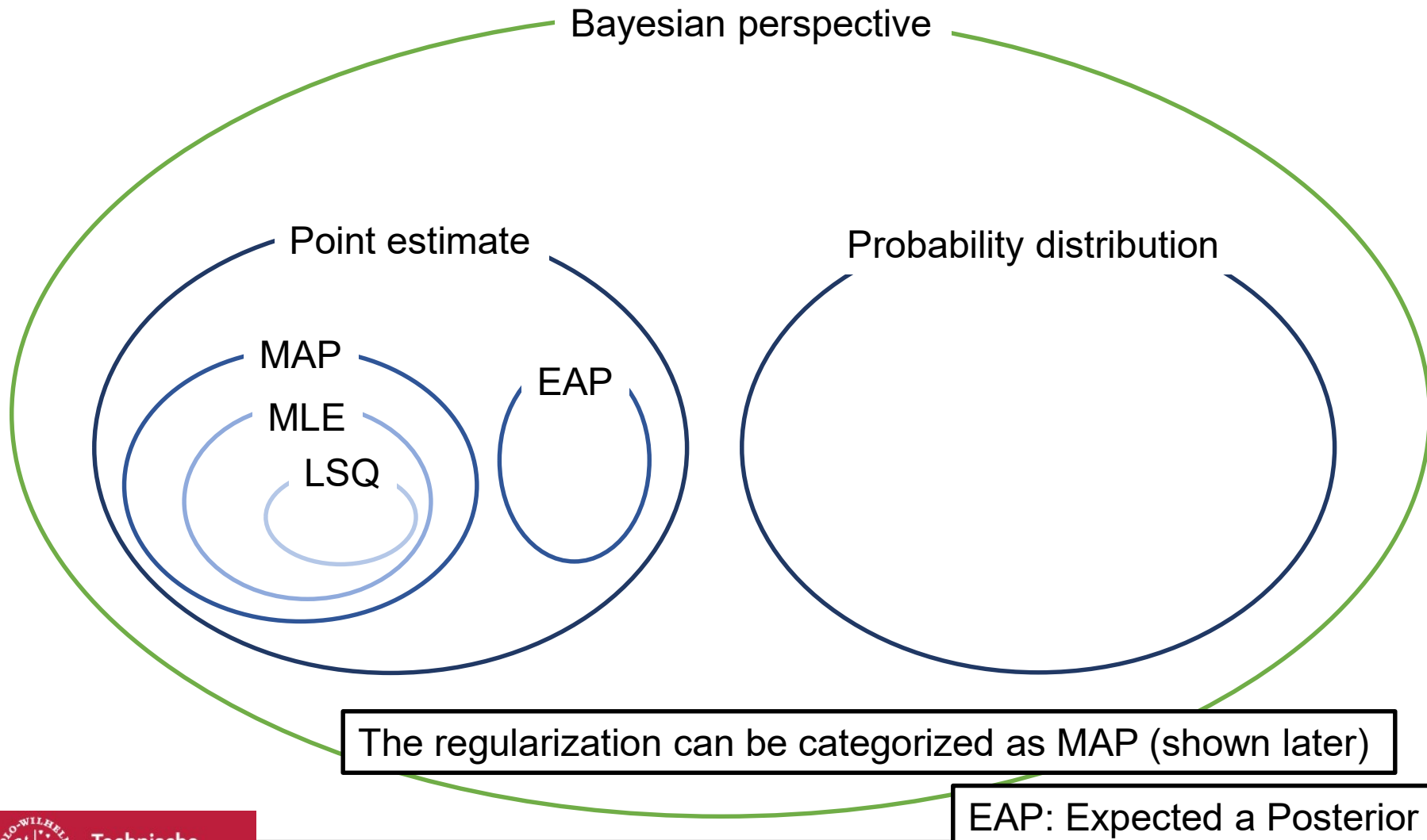
Generalized Process

1. Define a probabilistic model
2. Then, **compute the posterior**
 - (Define a prior distribution)
 - **Point estimate** (Deterministic)
 - **Probability distribution** (Stochastic) ← computations hard



$$\text{prior} \times \text{likelihood} \propto \text{posterior}$$

Bayesian Approach – Generalized Perspective



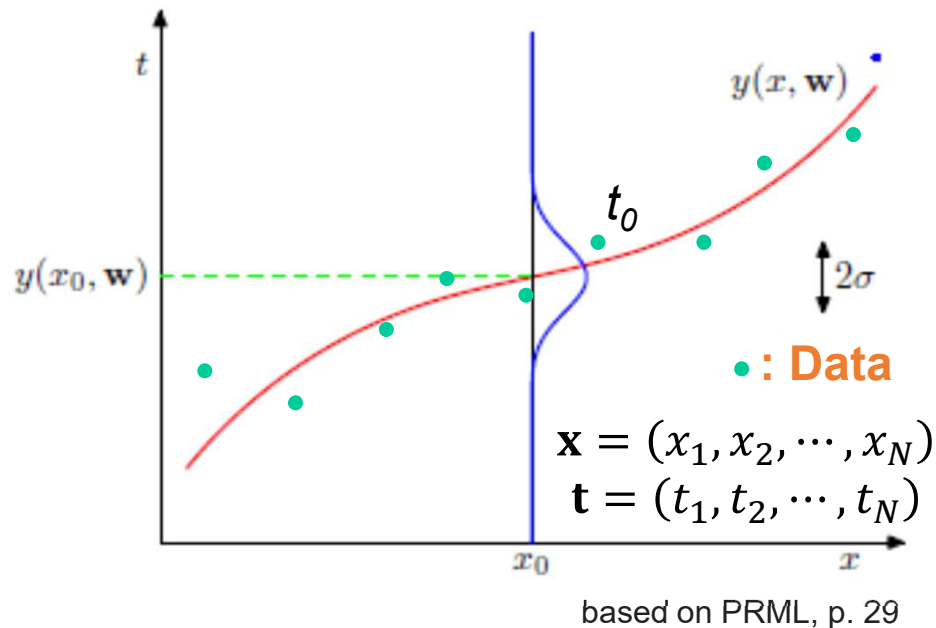
Lecture content

- Bayesian linear regression



Curve Fitting Revisiting

The least square method and the regularization method are summarized in perspectives based in the **probability theory**.



Define **a Probabilistic model**

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

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

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

x : deterministic variable

t : random variable

\mathbf{w}, σ : random variable

e.g. $y(x, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x)$ is a polynomial function.

Curve Fitting Revisiting (MLE)

Likelihood

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma) = \prod_{i=1}^N \mathcal{N}(t_i | y(x_i, \mathbf{w}), \sigma^2)$$

The likelihood function is uniquely determined by the data \mathbf{x}, \mathbf{t} .



MLE

$$\hat{\mathbf{w}}, \hat{\sigma} = \underset{\mathbf{w}, \sigma}{\operatorname{argmax}} p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma)$$

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^N \{t_i - y(x_i, \mathbf{w})\}^2 = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N \{t_i - y(x_i, \hat{\mathbf{w}})\}^2$$

$$\text{If } y(x, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x)$$

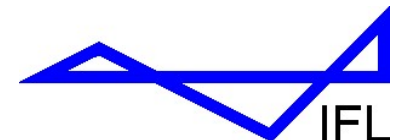
(If **Linear regression**)

x : deterministic variable

t : random variable

\mathbf{w}, σ : **deterministic variable**

$$(-\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma)) = \frac{1}{2\sigma^2} \sum_{i=1}^N \{t_i - y(x_i, \mathbf{w})\}^2 + \frac{N}{2} \ln(2\pi\sigma^2)$$



Curve Fitting Revisiting (Bayesian Perspective)

Prior

$$p(\mathbf{w}, \sigma)$$

You can set it.

Likelihood

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma) = \prod_{i=1}^N \mathcal{N}(t_i | y(x_i, \mathbf{w}), \sigma^2)$$

The likelihood function is uniquely determined by the data \mathbf{x}, \mathbf{t} .

Posterior

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

by the Bayes' theorem

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

If the prior $p(\mathbf{w}, \sigma)$ is a uniform distribution, $p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \sigma) \propto p(\mathbf{w}, \sigma | \mathbf{x}, \mathbf{t})$

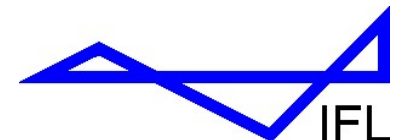
x : deterministic variable \leftarrow Please do not consider x

t : random variable

\mathbf{w}, σ : random variable

$$\mathbf{a} \equiv \mathbf{w}, \sigma$$

$$p(\mathbf{a}) p(\mathbf{t} | \mathbf{a}) \propto p(\mathbf{a} | \mathbf{t})$$



Bayesian Approach

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**

Bayesian Approach

There are so many difficulties in the Bayesian approach in general.

Difficulty 1: Computing **the posterior distribution** (the point estimate is fine)

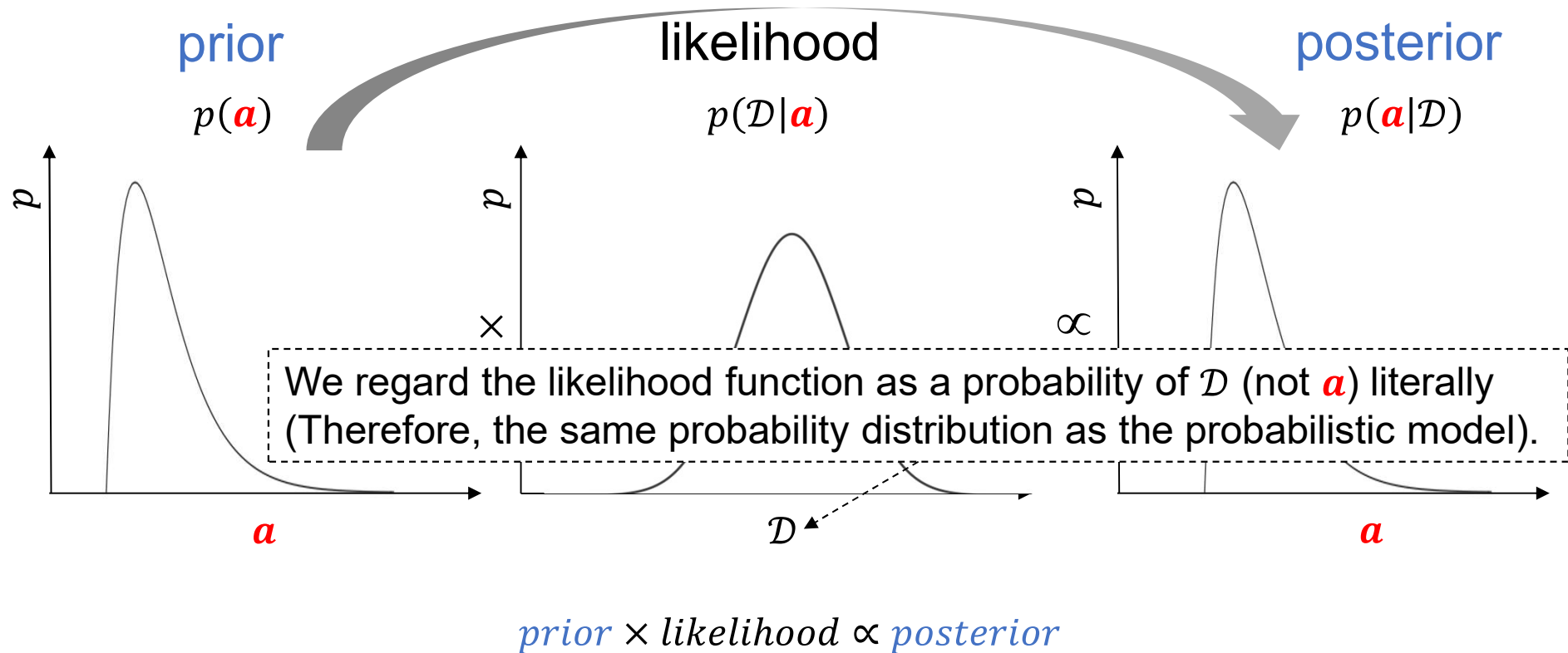
Difficulty 2: Computing **the predictive distribution** (the goal)

Computations are hard.

It is rather **rare** if we can obtain analytical solutions of them.

Conjugate Prior

Special cases



The posterior becomes the same type of probability distribution of the prior.

Conjugate Prior

Likelihood Function	Unknown Parameters	Conjugate Prior	Predictive Distribution
Bernoulli	μ	Beta	Bernoulli
Binomial	μ	Beta	Beta-Binomial
Categorical	$\boldsymbol{\mu}$	Dirichlet	Categorical
Multinomial	$\boldsymbol{\mu}$	Dirichlet	Dirichlet-Multinomial
Poisson	λ	Gamma	Negative Binomial
Gaussian	μ (σ^2 known)	Gaussian	Gaussian
Gaussian	σ^2 (μ known)	Gamma	Student-t
Gaussian	μ, σ^2	Gauss-Gamma	Student-t
Multivariate Gaussian	$\boldsymbol{\mu}$ ($\boldsymbol{\Sigma}$ known)	Multivariate Gaussian	Multivariate Gaussian
Multivariate Gaussian	$\boldsymbol{\Sigma}$ ($\boldsymbol{\mu}$ known)	Wishart	Multivariate Student-t
Multivariate Gaussian	$\boldsymbol{\mu}, \boldsymbol{\Sigma}$	Gaussian-Wishart	Multivariate Student-t

Conjugate Prior

Likelihood Function	Unknown Parameters	Conjugate Prior	Predictive Distribution
Bernoulli	<p>When we can have cases like “σ^2 known”?</p> <p>Please think about such cases by yourself.</p> <p>The answers were already given in the previous slides.</p>		Bernoulli
Binomial			Binomial
Categorical			Categorical
Multinomial			Multinomial
Poisson			Negative Binomial
Gaussian		Gaussian	Gaussian
Gaussian		Gamma	Student-t
Gaussian	μ, σ^2	Gauss-Gamma	Student-t
Multivariate Gaussian	μ (Σ known)	Multivariate Gaussian	Multivariate Gaussian
Multivariate Gaussian	Σ (μ known)	Wishart	Multivariate Student-t
Multivariate Gaussian	μ, Σ	Gaussian-Wishart	Multivariate Student-t

No need to remember!

Bayesian Linear Regression

PRML, p.152-158

There is an important thing which is not emphasized in the book:

- σ is known and constant.

P. 152: For the moment, we shall treat the noise precision parameter β as a known constant.

- ✓ We have been doing MLE for μ, σ both of which are unknown.

➡ $\hat{\mu}, \hat{\sigma}$ obtained

- ✓ We will extend this to the Bayesian approach but only for μ as unknown, σ is known.

A possible idea: e.g. $\hat{\sigma}$ is used.

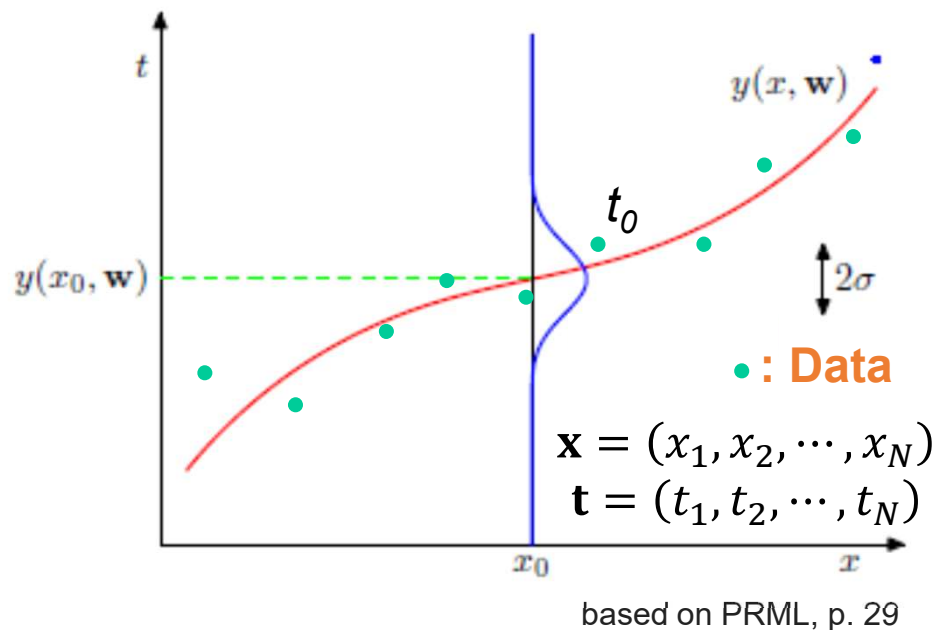
➡ Then, only μ is considered as probability as $p(\mu)$.

Conjugate Prior

Likelihood Function	Unknown Parameters	Conjugate Prior	Predictive Distribution
Bernoulli	μ	Beta	Bernoulli
Binomial	μ	Beta	Beta-Binomial
Categorical	$\boldsymbol{\mu}$	Dirichlet	Categorical
Multinomial	$\boldsymbol{\mu}$	Dirichlet	Dirichlet-Multinomial
Poisson	λ	Gamma	Negative Binomial
Gaussian	μ (σ^2 known)	Gaussian	Gaussian
Gaussian	σ^2 (μ known)	Gamma	Student-t
Gaussian	μ, σ^2	Gauss-Gamma	Student-t
Multivariate Gaussian	$\boldsymbol{\mu}$ ($\boldsymbol{\Sigma}$ known)	Multivariate Gaussian	Multivariate Gaussian
Multivariate Gaussian	$\boldsymbol{\Sigma}$ ($\boldsymbol{\mu}$ known)	Wishart	Multivariate Student-t
Multivariate Gaussian	$\boldsymbol{\mu}, \boldsymbol{\Sigma}$	Gaussian-Wishart	Multivariate Student-t

Bayesian Linear Regression

The least square method and the regularization method are summarized in perspectives based in the **probability theory**.



$x, \hat{\sigma}$: deterministic variable

t : random variable

\mathbf{w} : random variable

Define **a Probabilistic model**

$$p(t|x, \mu) = \mathcal{N}(t|\mu(x), \hat{\sigma}^2)$$

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



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

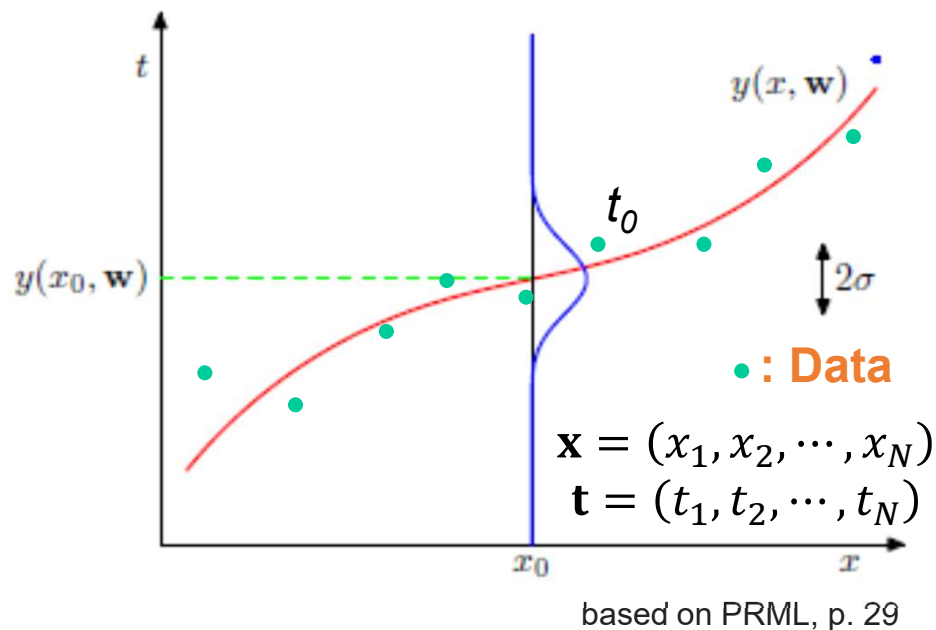
Likelihood can be then determined.

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(t_i|y(x_i, \mathbf{w}), \hat{\sigma}^2)$$

Gaussian distributions wrt t_i

Bayesian Linear Regression

The least square method and the regularization method are summarized in perspectives based in the **probability theory**.

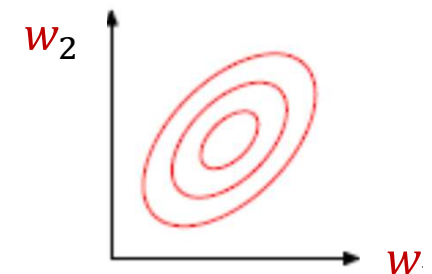


$x, \hat{\sigma}$: deterministic variable
 t : random variable
 \mathbf{w} : random variable

Define **a prior**

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

your setting



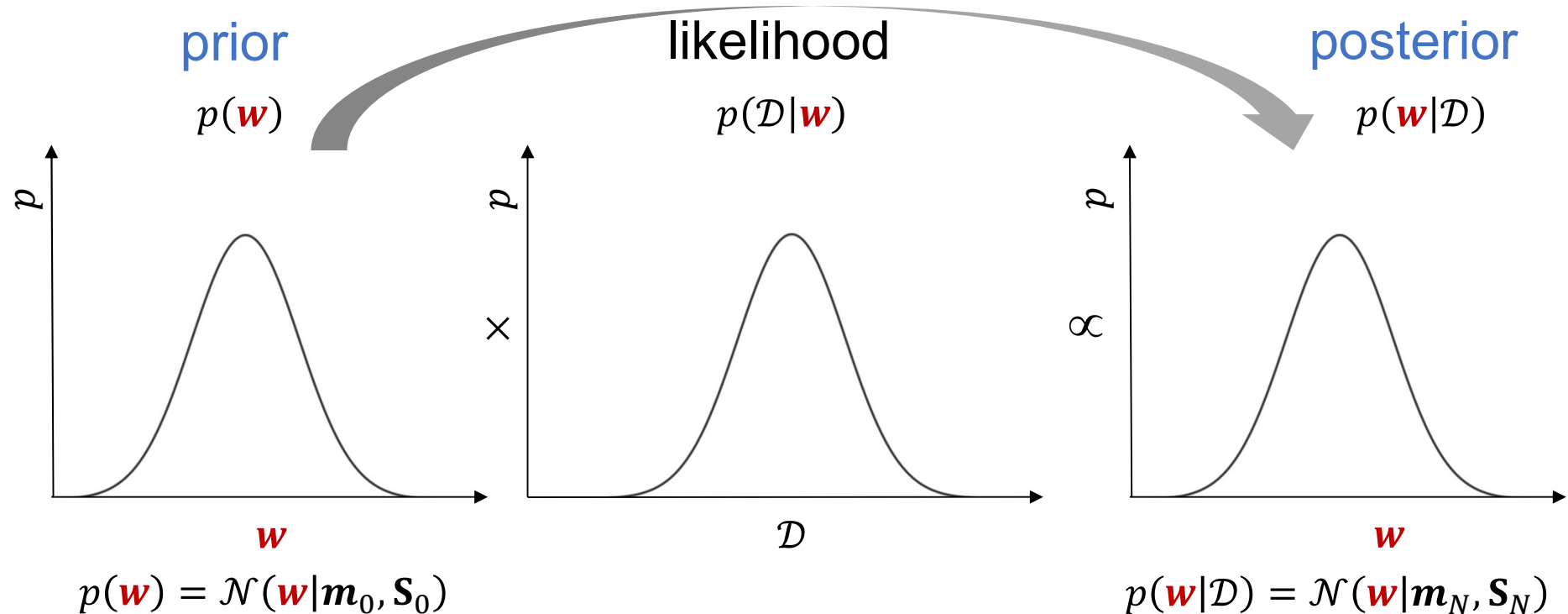
A Gaussian distribution in general



a posterior using the likelihood

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

Bayesian Linear Regression

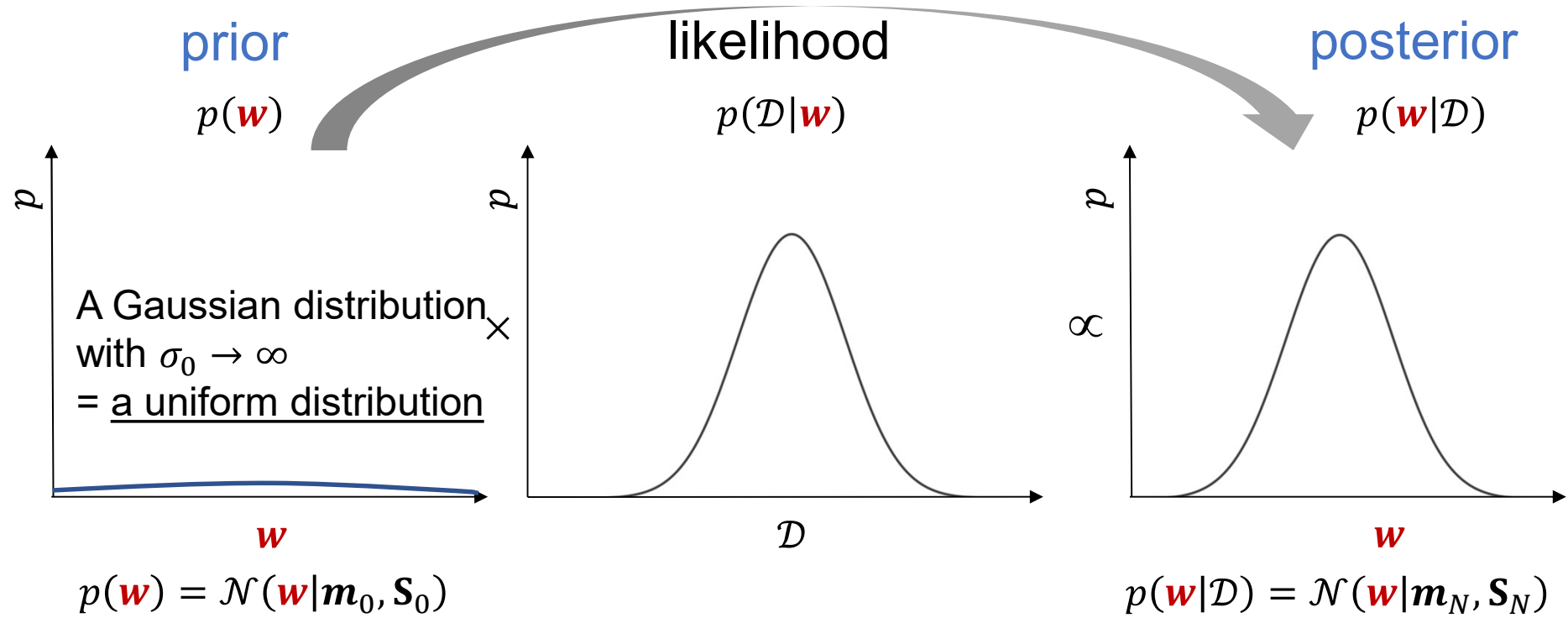


**Conjugate prior
+ Linear regression**

Analytical solutions of

- Posterior distribution
- Predictive distribution

Bayesian Linear Regression



**Conjugate prior
+ Linear regression**

Analytical solutions of

- Posterior distribution
- Predictive distribution

Bayesian Linear Regression

Details of the analytical solutions:
See PRML, p.152

Prior (your setting)

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \sigma_0^2 \mathbf{I})$$

Special settings:

- Conjugate prior
 - and all Gaussian distributions
- Linear regression

Posterior

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

$\mathbf{m}_N, \mathbf{S}_N$: analytically obtained

By Bayes' theorem

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

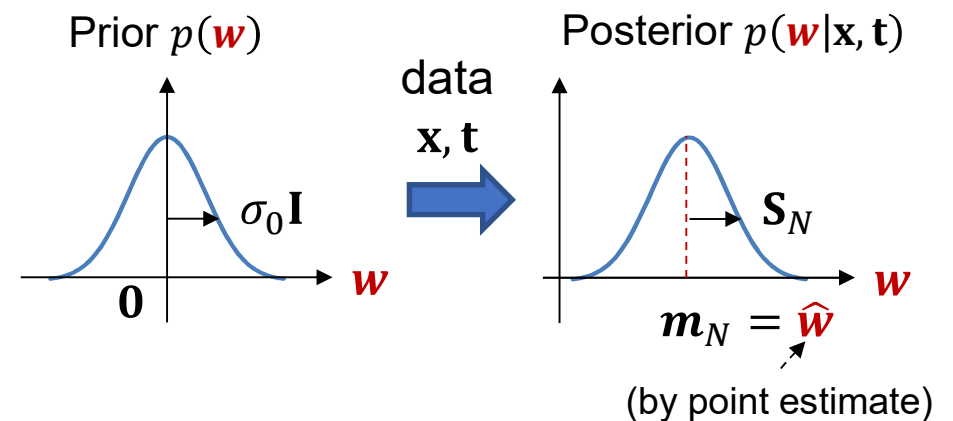
$$\text{where, } p(\mathbf{t} | \mathbf{x}, \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(t_i | y(x_i, \mathbf{w}), \hat{\sigma}^2)$$

By the way,

How does the prior work in this setting?



Relationship with
regularization



Curve Fitting Revisiting (Review from Lecture 3)

e.g. $p(\mathbf{w})$: a Gaussian distribution around $\mathbf{0}$

$$p(\mathbf{w}|\sigma_w) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma_w^2 \mathbf{I})$$

Bayes' theorem $p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$

➡ $-\ln p(\mathbf{w}|\mathcal{D}) \propto -\ln p(\mathcal{D}|\mathbf{w}) - \ln p(\mathbf{w})$

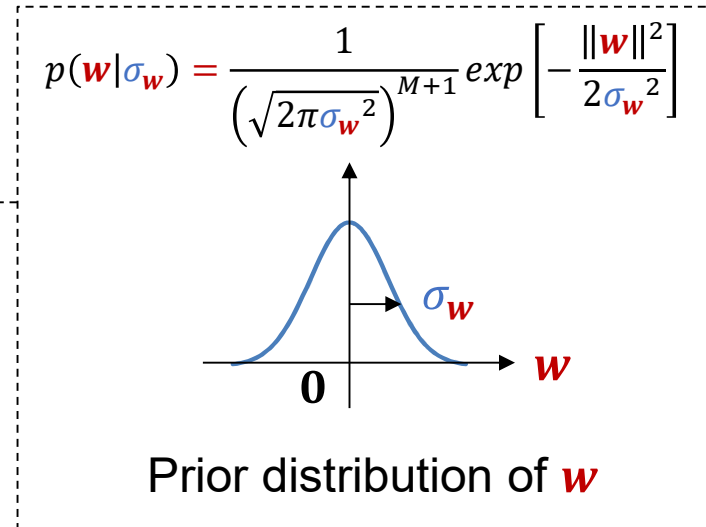
$$= \frac{1}{2\hat{\sigma}^2} \sum_{i=1}^N \{t_i - y(x_i, \mathbf{w})\}^2 + \frac{N}{2} \ln(2\pi\hat{\sigma}^2) + \frac{1}{2\sigma_w^2} \|\mathbf{w}\|^2$$

$$-\ln p(\mathbf{w}|\mathcal{D}) \propto \sum_{i=1}^N \{t_i - y(x_i, \mathbf{w})\}^2 + \left(\frac{\hat{\sigma}}{\sigma_w}\right)^2 \|\mathbf{w}\|^2 = \underbrace{E(\mathbf{w}) + \lambda \|\mathbf{w}\|^2}_{\text{regularization}}$$

where,

$$\lambda \equiv \left(\frac{\hat{\sigma}}{\sigma_w}\right)^2$$

The same function as that for the regularization



Linear Regression (Review from Lecture 4)

If the regression $y(x_i, \mathbf{w})$ is a linear regression form $\mathbf{w}^T \boldsymbol{\phi}(x)$:

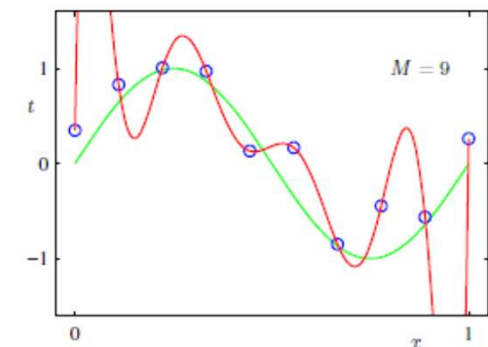
$$\text{i.e. } y(x_i, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x)$$

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

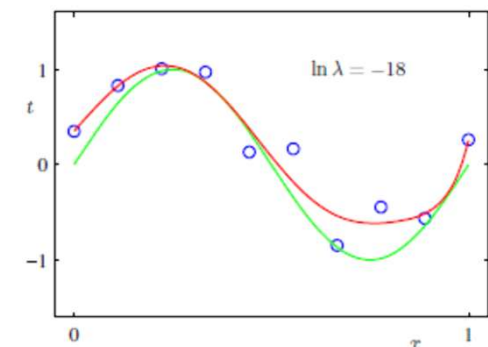
The **regularization** was (first introduced in Lecture 2):
The least square method with penalty term of \mathbf{w}

The **regularization** is now:
The point estimate with a prior distribution of \mathbf{w}

= Maximizing a posterior (MAP)



↓ regularization



The prediction process is not decided only by the given data, but also with prior information.

Linear Regression (Review from Lecture 4)

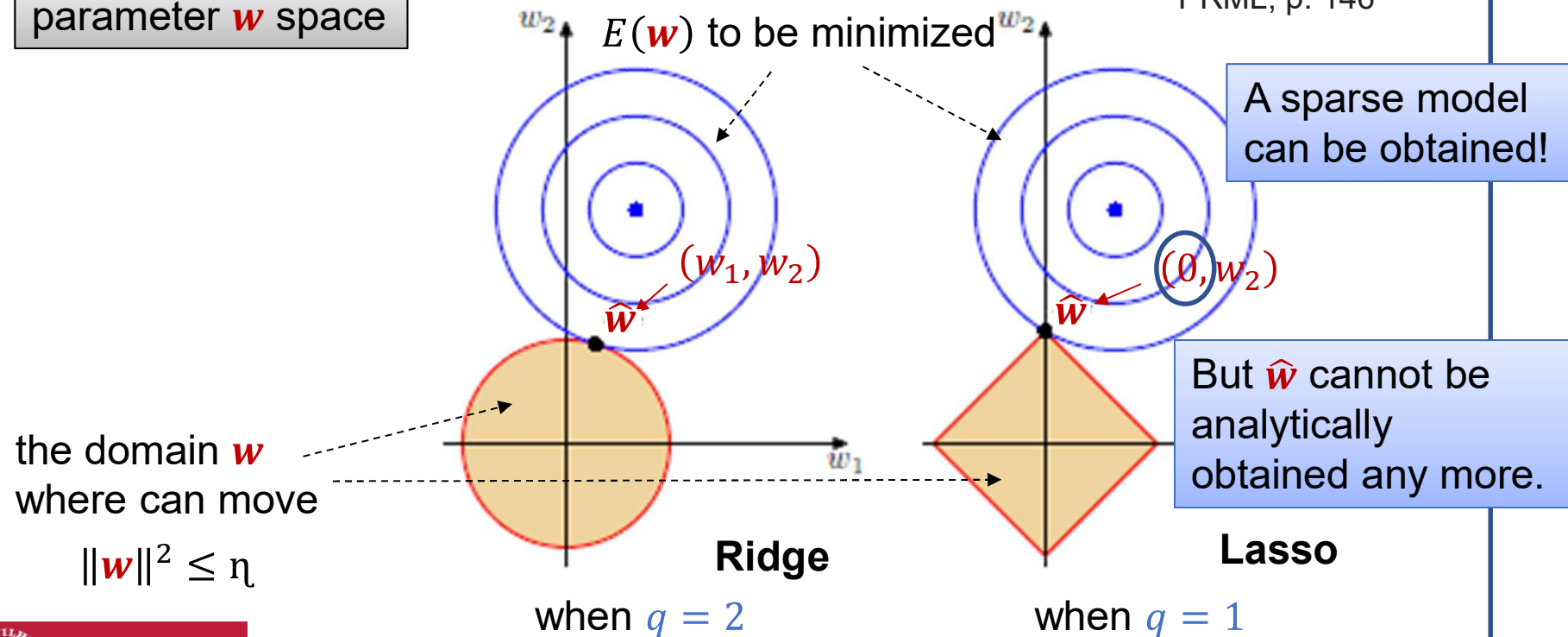
Other regularization techniques

regularization term

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} E_{reg}(\mathbf{w}) \quad \text{where,} \quad E_{reg}(\mathbf{w}) = E(\mathbf{w}) + \lambda \|\mathbf{w}\|^q$$

parameter \mathbf{w} space

PRML, p. 146



Regularization - LASSO Regression

Prior (your setting)

$$p(\mathbf{w}|\sigma_w) = \mathcal{N}(\mathbf{w}|0, \sigma_w^2 \mathbf{I})$$

when the prior is a Gaussian distribution



Ridge

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

Another Prior (your setting)

$$p(\mathbf{w}|b_w) = \text{Laplace}(\mathbf{w}|0, b_w)$$

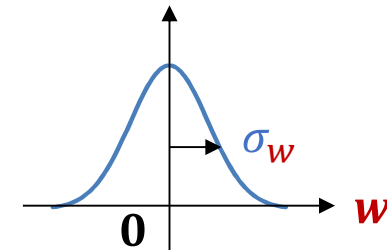
when is the prior is a Laplace distribution



Lasso

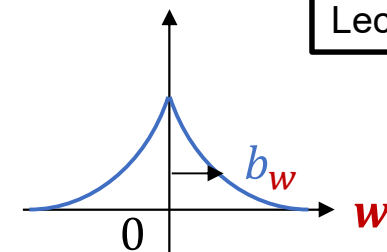
$\hat{\mathbf{w}}$: no analytical solution

$$p(\mathbf{w}|\sigma_w) = \frac{1}{(\sqrt{2\pi}\sigma_w)^{M+1}} \exp\left[-\frac{\|\mathbf{w}\|^2}{2\sigma_w^2}\right]$$



Prior distribution of \mathbf{w}

$$p(\mathbf{w}|b_w) = \frac{1}{(2b_w)^{M+1}} \exp\left\{-\frac{\|\mathbf{w}\|}{b_w}\right\}$$



Prior distribution of \mathbf{w}

Lecture 4

Bayesian Linear Regression (Note)

Details of the analytical solutions:
See PRML, p.152

Prior (your setting)

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \sigma_0^2 \mathbf{I})$$

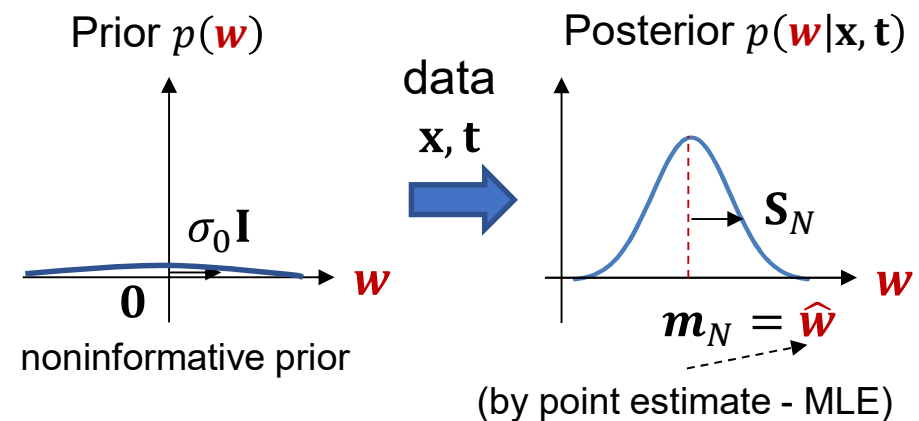
Special settings:

- Conjugate prior
 - and all Gaussian distributions
- Linear regression

Posterior

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

$\mathbf{m}_N, \mathbf{S}_N$: analytically obtained



$$\mathbf{m}_N = \hat{\sigma}^{-2} \mathbf{S}_N \Phi^T \mathbf{t} \xrightarrow{\sigma_0 \rightarrow \infty} \hat{\sigma}^{-2} (\hat{\sigma}^{-2} \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = \hat{\mathbf{w}} \text{ by MLE}$$

$$\mathbf{S}_N^{-1} = \sigma_0^{-2} \mathbf{I} + \hat{\sigma}^{-2} \Phi^T \Phi \xrightarrow{\sigma_0 \rightarrow \infty} \hat{\sigma}^{-2} \Phi^T \Phi$$

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

likelihood \propto posterior

Lecture content

- Bayesian sequential learning



Bayesian Sequential Learning

Bayes' theorem

$$p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) \quad \text{wrt } \mathbf{w} \qquad p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

A dataset \mathcal{D} divided into N datasets as $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_N$.

$$p(\mathbf{w}|\mathcal{D}_1) \propto p(\mathcal{D}_1|\mathbf{w})p(\mathbf{w})$$

next prior

$$p(\mathbf{w}|\mathcal{D}_1, \mathcal{D}_2) \propto p(\mathcal{D}_2|\mathbf{w})p(\mathbf{w}|\mathcal{D}_1)$$

next prior

...

$$p(\mathbf{w}|\mathcal{D}) \propto \prod_{i=1}^N p(\mathcal{D}_i|\mathbf{w}) p(\mathbf{w})$$



$$p(\mathbf{w}|\mathcal{D}_1, \mathcal{D}_2) \propto p(\mathcal{D}_2|\mathbf{w})p(\mathcal{D}_1|\mathbf{w})p(\mathbf{w})$$

$$p(\mathbf{w}|\mathcal{D}_1, \mathcal{D}_2) = \frac{p(\mathcal{D}_2|\mathbf{w})p(\mathbf{w}|\mathcal{D}_1)}{p(\mathcal{D}_2)}$$

Sequential view
Bayesian sequential learning

Bayesian Sequential Learning

$$p(\mathbf{w}|\mathcal{D}) \propto \prod_{i=1}^N p(\mathcal{D}_i|\mathbf{w}) p(\mathbf{w})$$

The likelihood of data in batch itself

Sequential view
Bayesian sequential learning

This can be applied to any problems in which the observed data are assumed to be **i.i.d.**

Review of Lecture 2:

Independent and identically distributed (i.i.d.)

The likelihood function

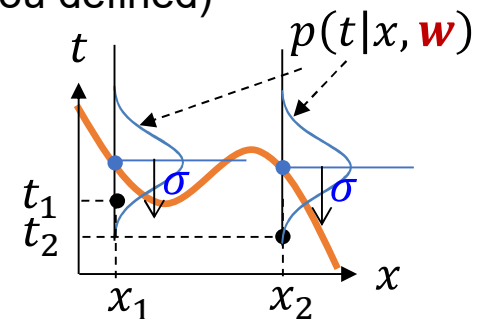
$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(t_i|y(x_i, \mathbf{w}), \hat{\sigma}^2)$$

Each sample point is i.i.d.

independently generated from a Gaussian distribution (the probabilistic model that you defined)

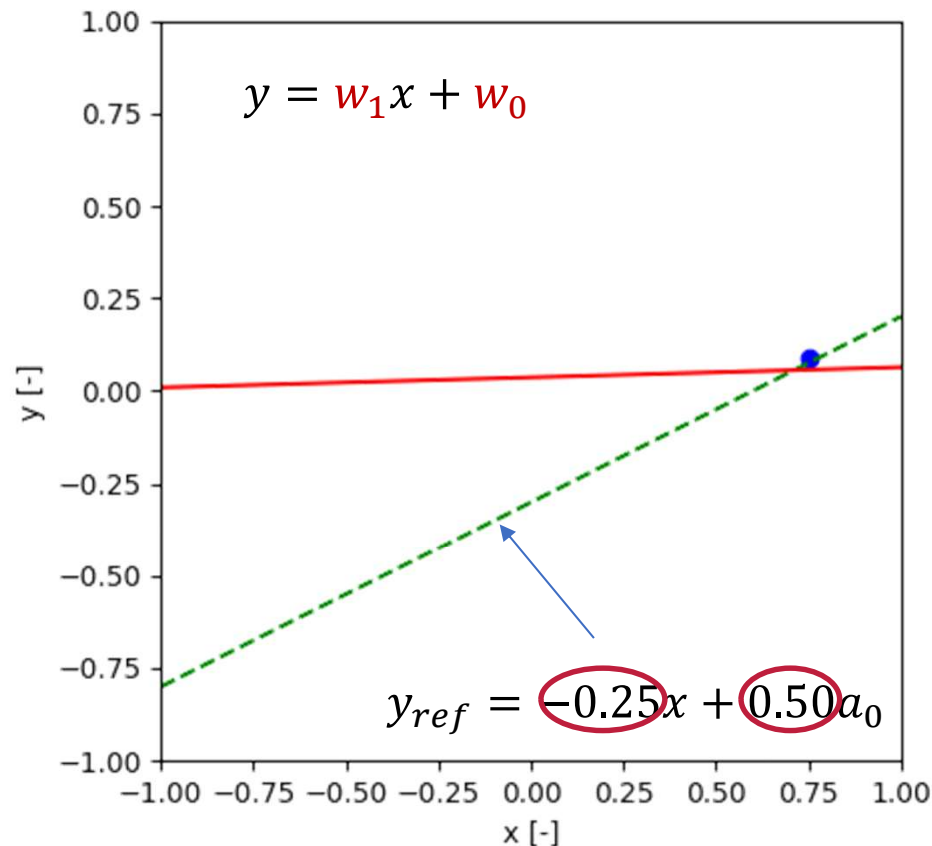
independent

$$p(x_1, x_2) = p(x_1)p(x_2) = \prod_{i=1}^2 p(x_i)$$

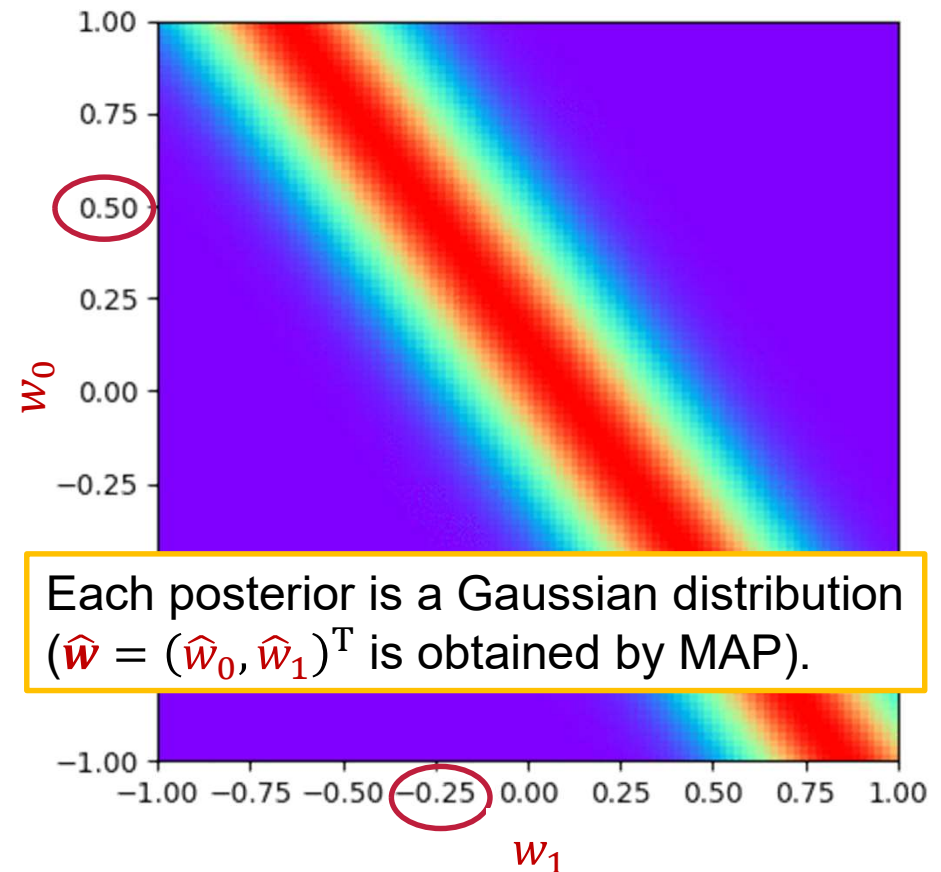


Bayesian Sequential Learning

input - output space



parameter space $p(\mathbf{w}|\mathcal{D})$



Each posterior is a Gaussian distribution ($\hat{\mathbf{w}} = (\hat{w}_0, \hat{w}_1)^T$ is obtained by MAP).

The parameters w_0 and w_1 are being clarified to be around (-0.25, 0.50).

Lecture content

- Uncertainty due to data (active learning)



Bayesian Linear Regression

Details of the analytical solutions:
See PRML, p.152

Prior (your setting)

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \sigma_0^2 \mathbf{I})$$

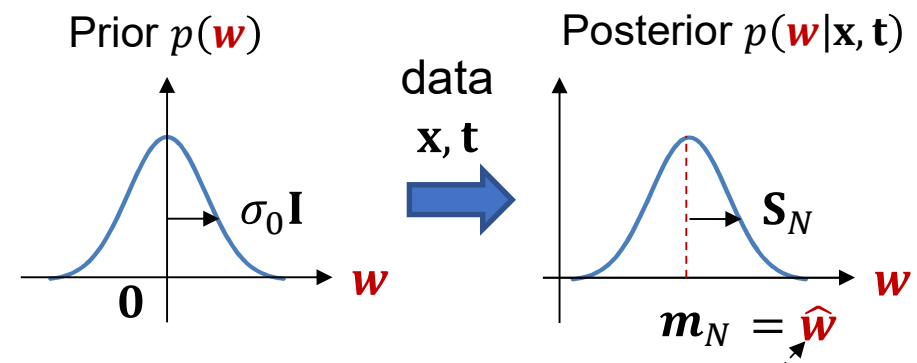
Special settings:

- Conjugate prior
 - and all Gaussian distributions
- Linear regression

Posterior

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

$\mathbf{m}_N, \mathbf{S}_N$: analytically obtained



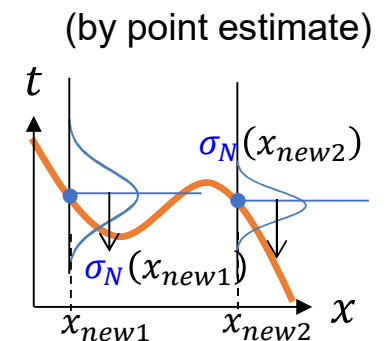
Predictive distribution (the goal)

$$\Rightarrow p(t | x, \mathbf{x}, \mathbf{t}) = \int p(t | x, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w} = \mathcal{N}(t | \boxed{\mathbf{m}_N^T \phi(x)}, \underbrace{\sigma_N^2(x)}_{\text{the linear regression itself}})$$

The predictive distribution result contains the result of “point estimate”.

$$\hat{\mathbf{w}}^T \phi(x)$$

What is σ_N ?



Bayesian Linear Regression

Note:

The reason why the variance can be decomposed clearly like this is again due to the properties of Gaussian distributions in general.

$$\begin{aligned}\sigma_N^2(\mathbf{x}) &= \hat{\sigma}^2 + \cancel{\sigma_N'^2(\mathbf{x})} \xrightarrow{N \rightarrow \infty} 0 \\ &= \text{noise} + \text{uncertainty associated with } \mathbf{w}\end{aligned}$$

What does this mean?

where, $\sigma_N'^2(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$

The meanings of these two components are totally different from each other.

Bayesian Linear Regression

(MLE) Least square method

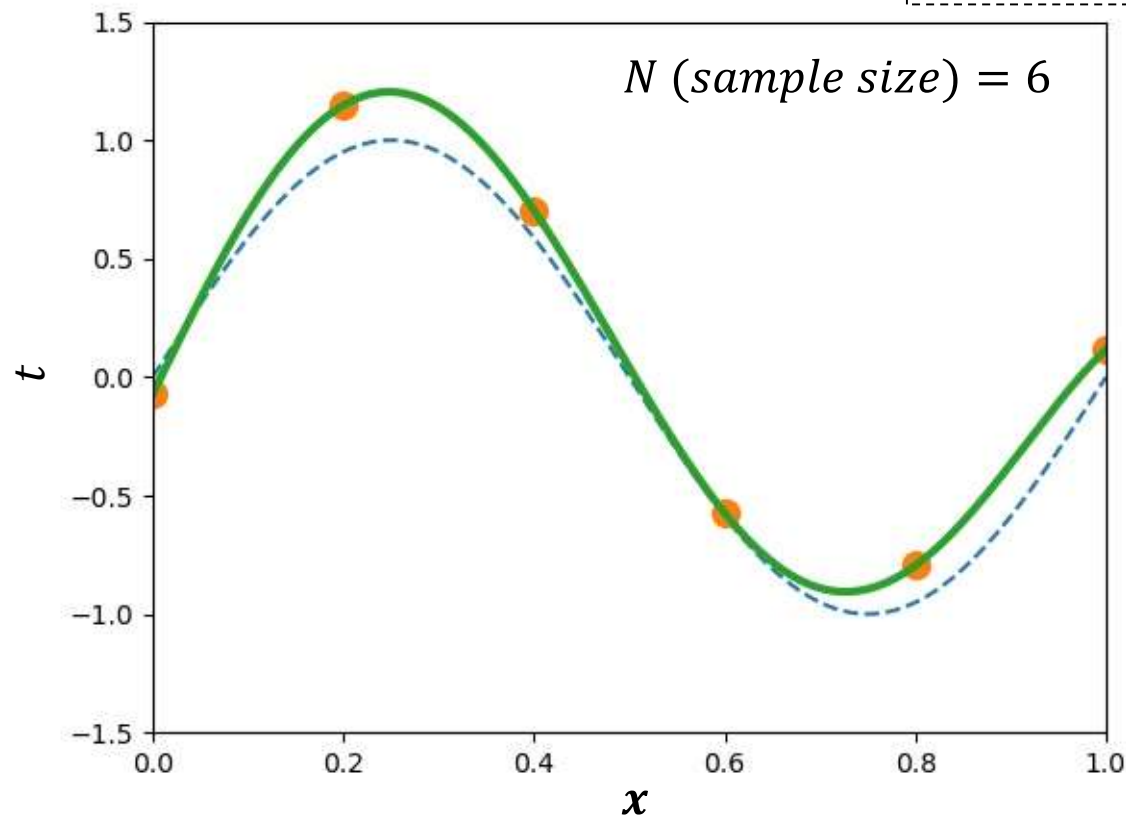
Examples

The probabilistic model

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

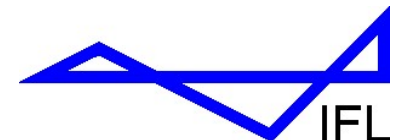
$\boldsymbol{\phi}(x)$ = polynomials ($M = 9$)

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



Technische
Universität
Braunschweig

Dr. Daigo Maruyama | Scientific Machine Learning: Lecture 6 | Slide 37



Bayesian Linear Regression

The probabilistic model

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

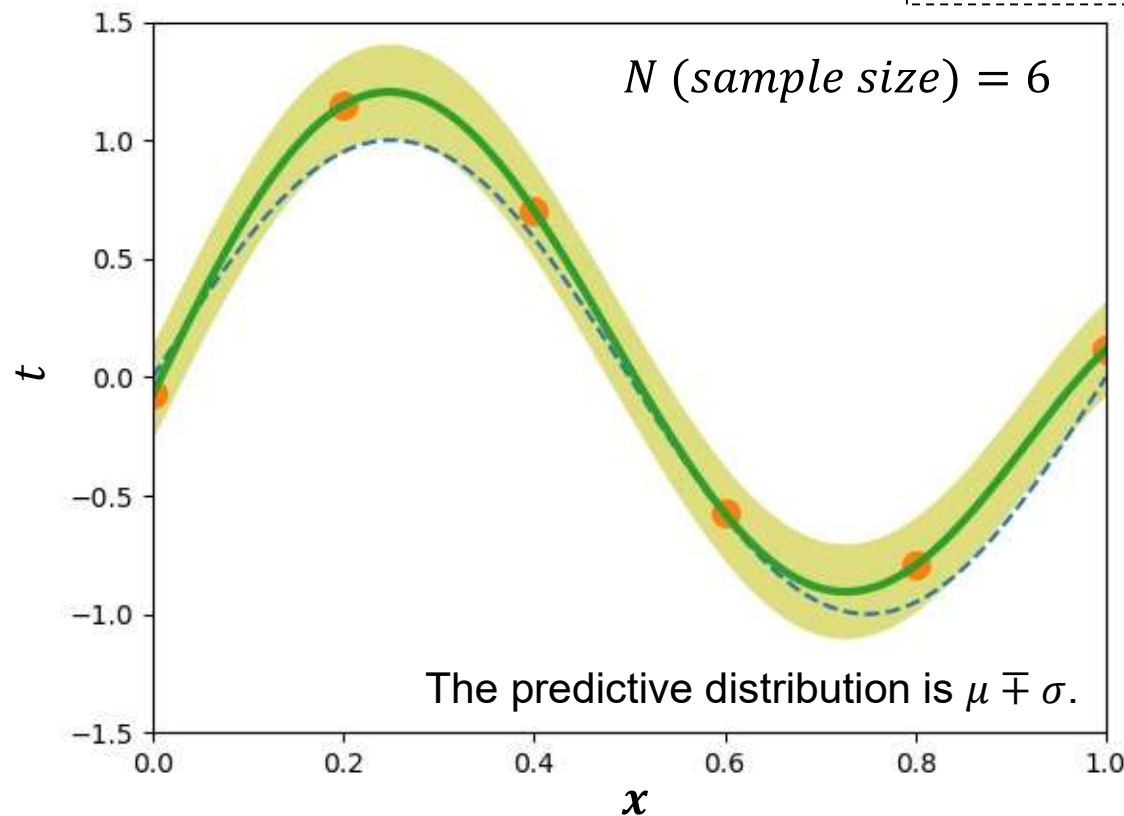
$\boldsymbol{\phi}(x)$ = polynomials ($M = 9$)

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

(MLE) Least square method

Examples

$\hat{\sigma}$



$\hat{\sigma}$: constant noise

Bayesian Linear Regression

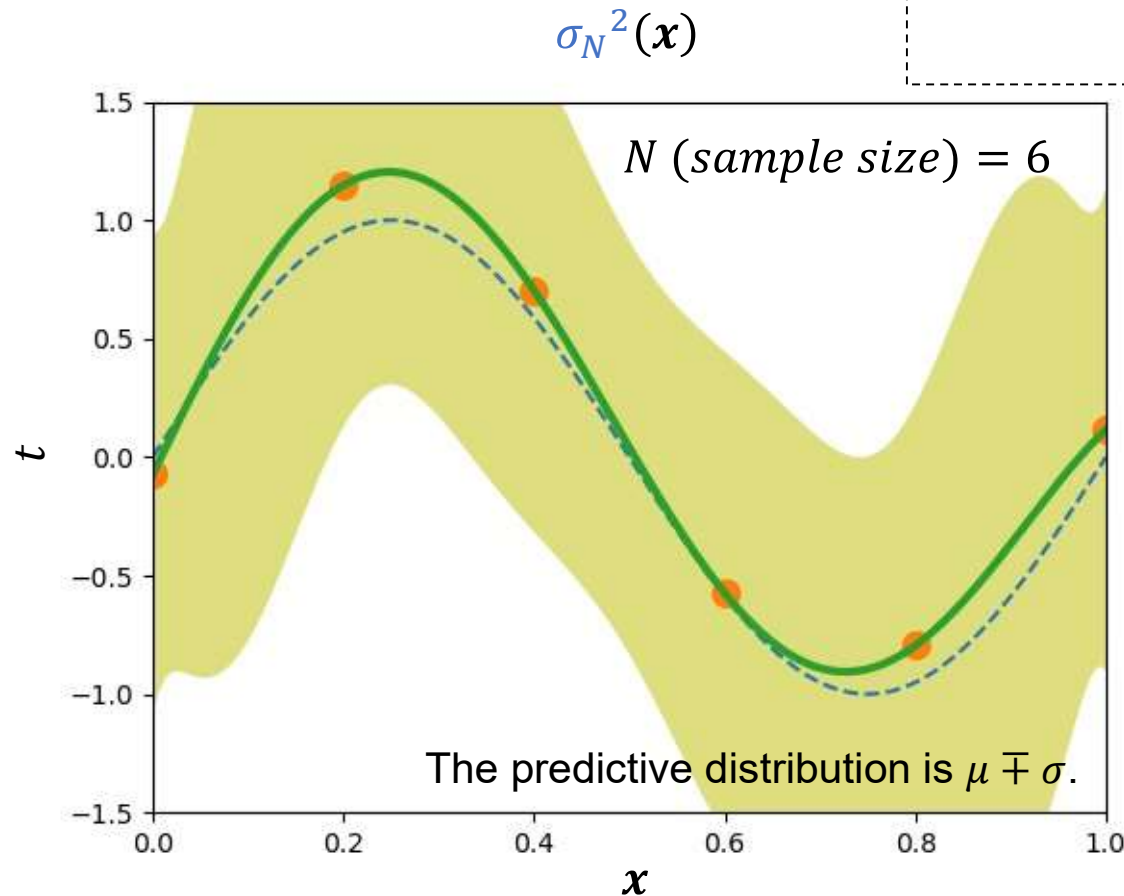
The probabilistic model

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

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

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

Examples



$\sigma_N^2(x)$: constant noise + uncertainty associated with \mathbf{w}

Bayesian Linear Regression

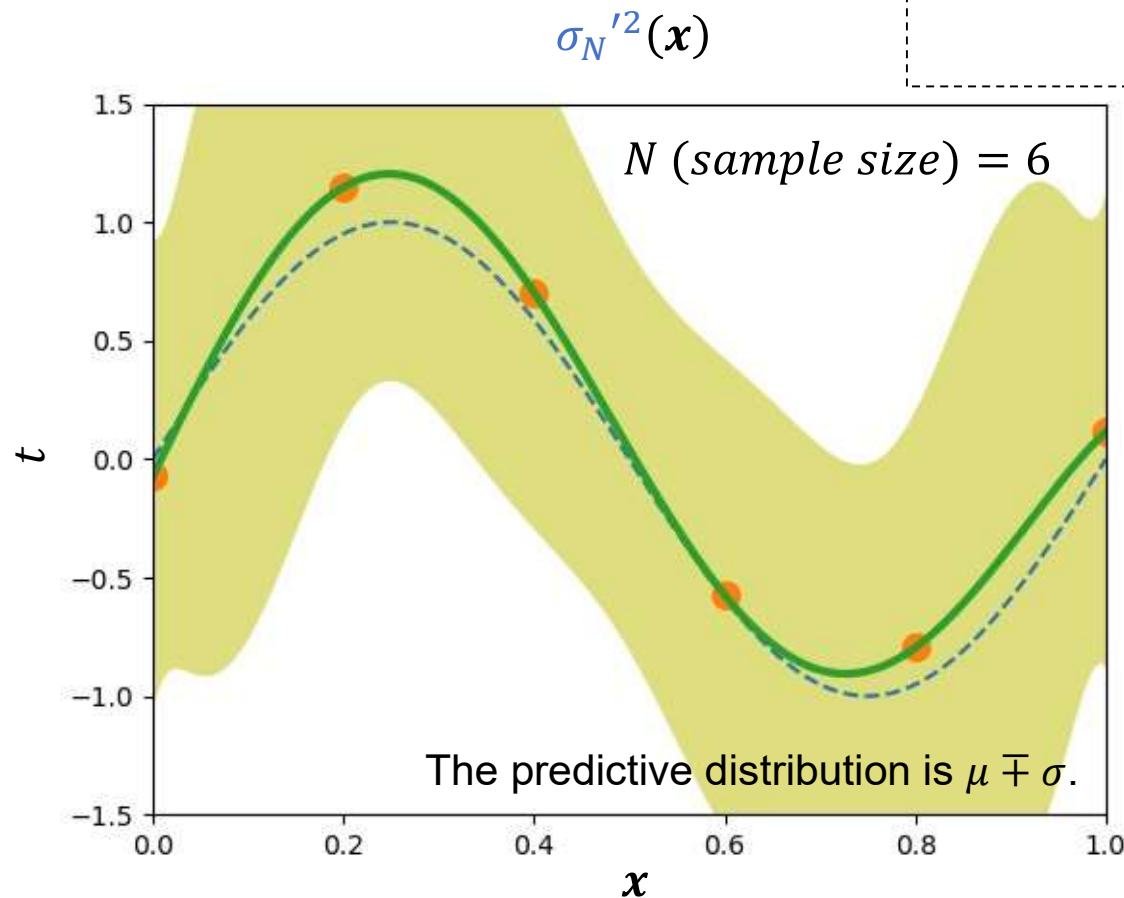
The probabilistic model

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

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

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

Examples



$\sigma_N'^2(x)$: uncertainty associated with \mathbf{w}



Technische
Universität
Braunschweig

Bayesian Linear Regression

The probabilistic model

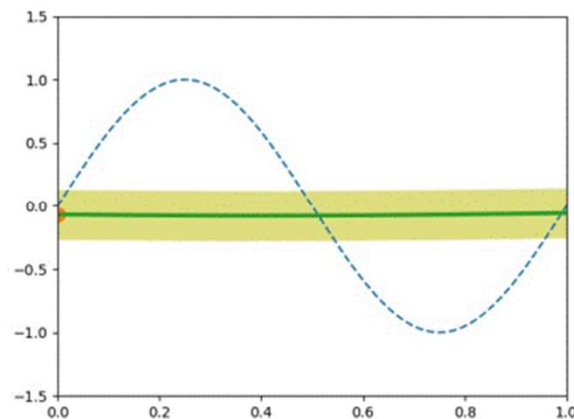
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

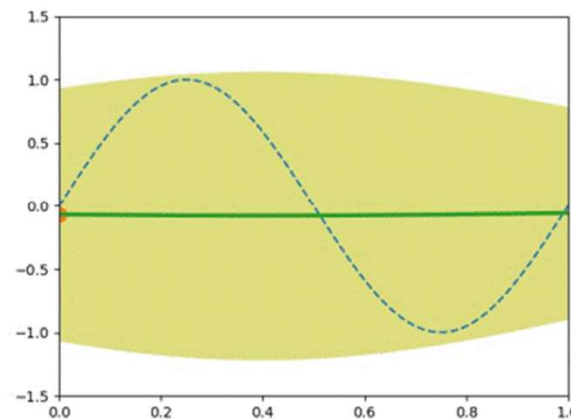
Examples

N (sample size) increasing $N = 1$



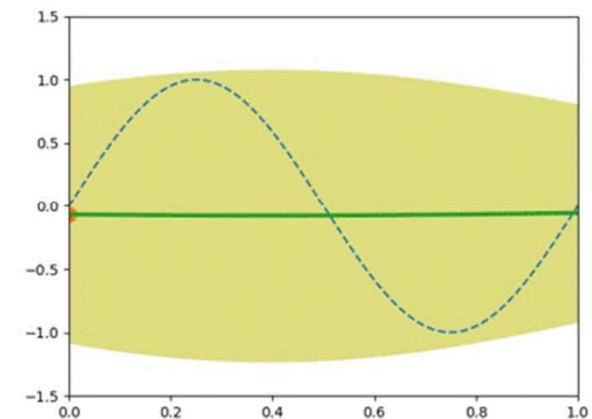
$\hat{\sigma}^2$

+



$\sigma_N'^2(x)$

=



$\sigma_N^2(x)$

$\sigma_N'^2(x)$: uncertainty associated with \mathbf{w}

Note: The predictive distribution is $\mu \mp \sigma$.

Bayesian Linear Regression

The probabilistic model

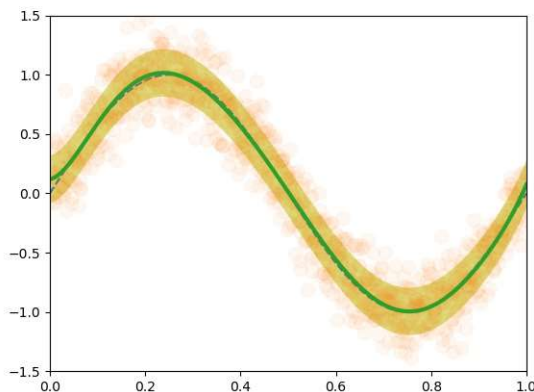
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

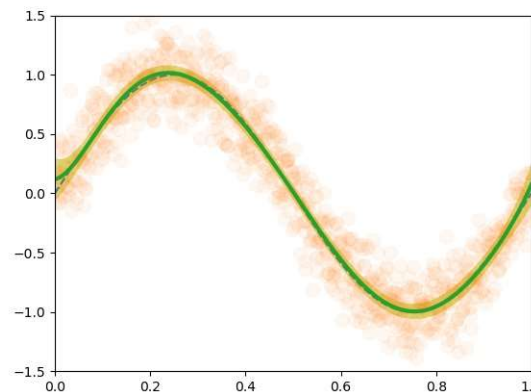
Examples

$N = 1000$



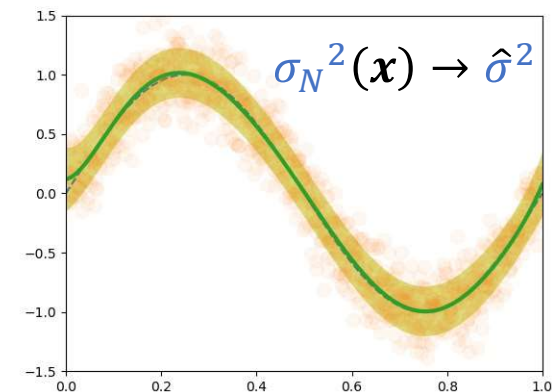
$\hat{\sigma}^2$

+



$\sigma_N'^2(x)$ ⁰

=



$\sigma_N^2(x) \rightarrow \hat{\sigma}^2$

$\sigma_N^2(x)$

$\sigma_N'^2(x)$: uncertainty associated with \mathbf{w}

Note: The predictive distribution is $\mu \mp \sigma$.

Bayesian Linear Regression

Question: How can we use this information?

Active learning

if evaluating the output t (for a given input x) is **expensive**

Request new sample points on:

- where σ_N is large.
- where a criterion defined by σ_N and the optimum location among the current sample points is large.
- etc. (σ_N + other information)

Annotation in classification

Bayesian optimization

With a small amount of learning datasets, the output t is **efficiently** predicted.



Gaussian Processes (Lectures 7-9) are more often used (since weak assumption on the probabilistic model can be used).

Bayesian Linear Regression

The probabilistic model

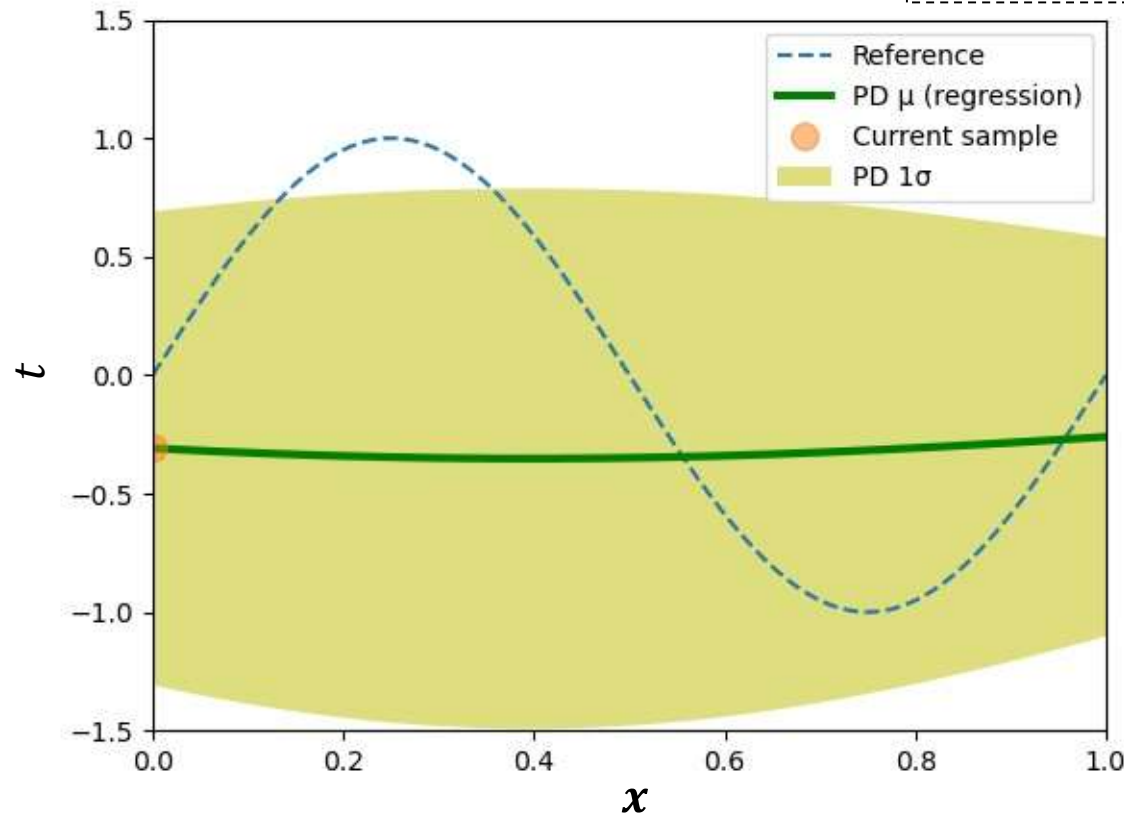
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



starting from N (sample size) = 1



Technische
Universität
Braunschweig

Bayesian Linear Regression

The probabilistic model

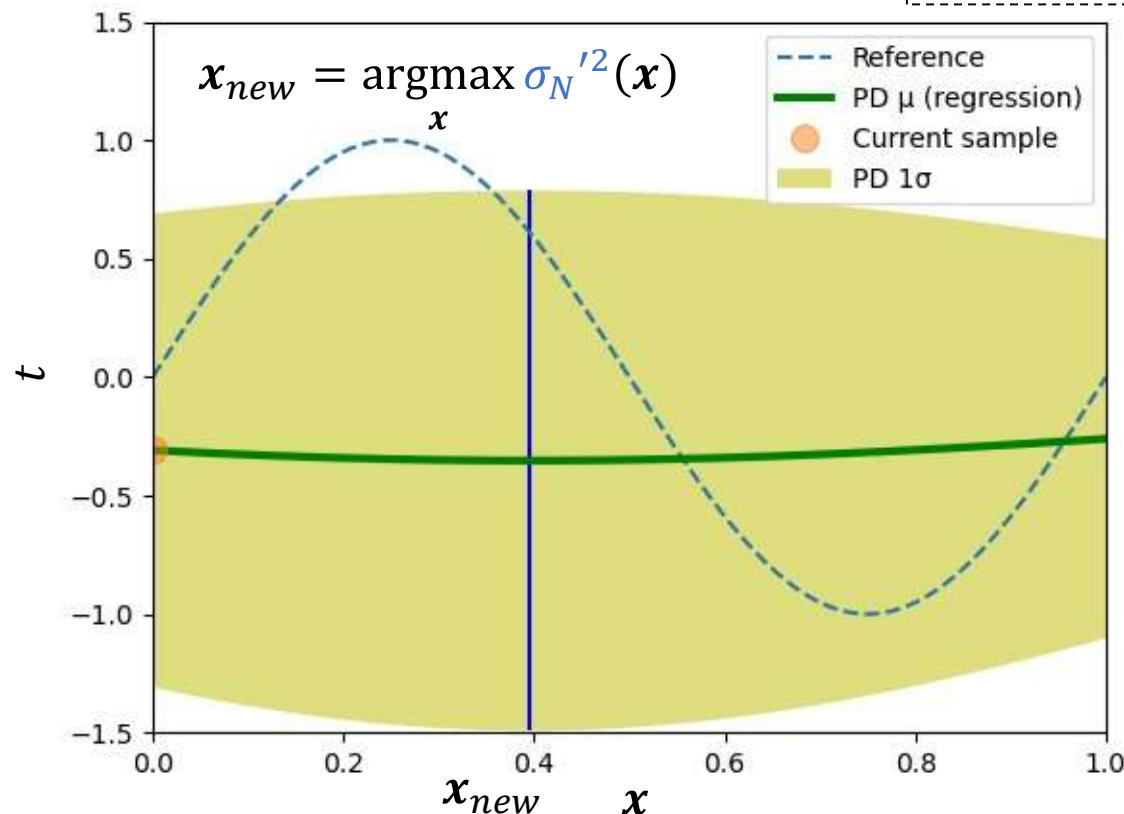
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Ask to provide the output t for the location x_{new}



Technische
Universität
Braunschweig

Bayesian Linear Regression

The probabilistic model

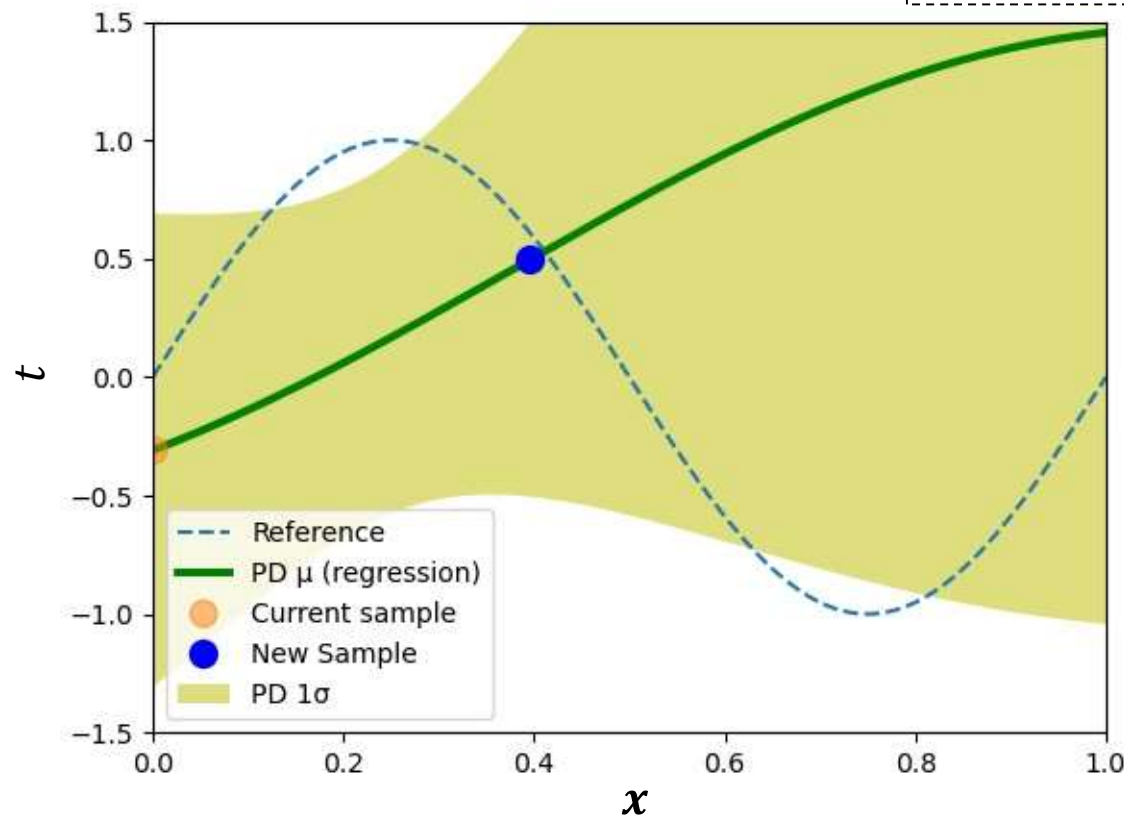
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



A new sample point has been added **actively/automatically**.

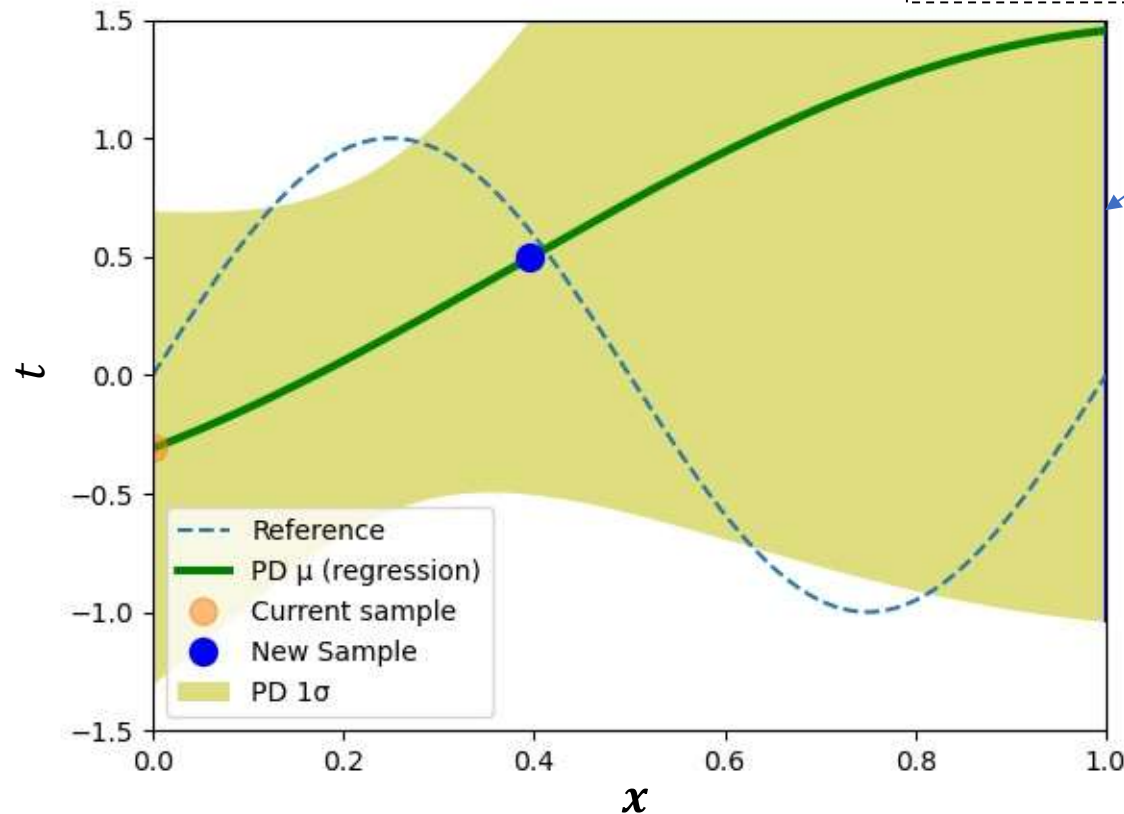
Bayesian Linear Regression

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)



Continue the iterative process...

Bayesian Linear Regression

The probabilistic model

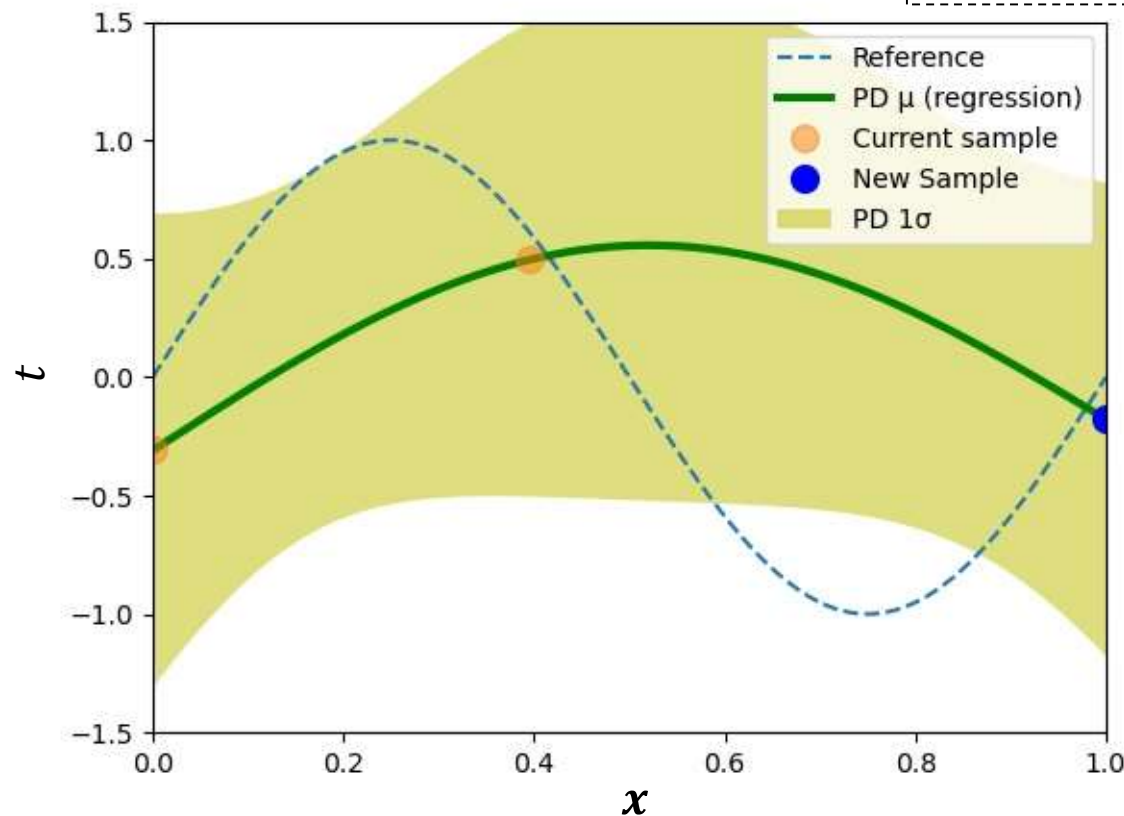
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

Examples

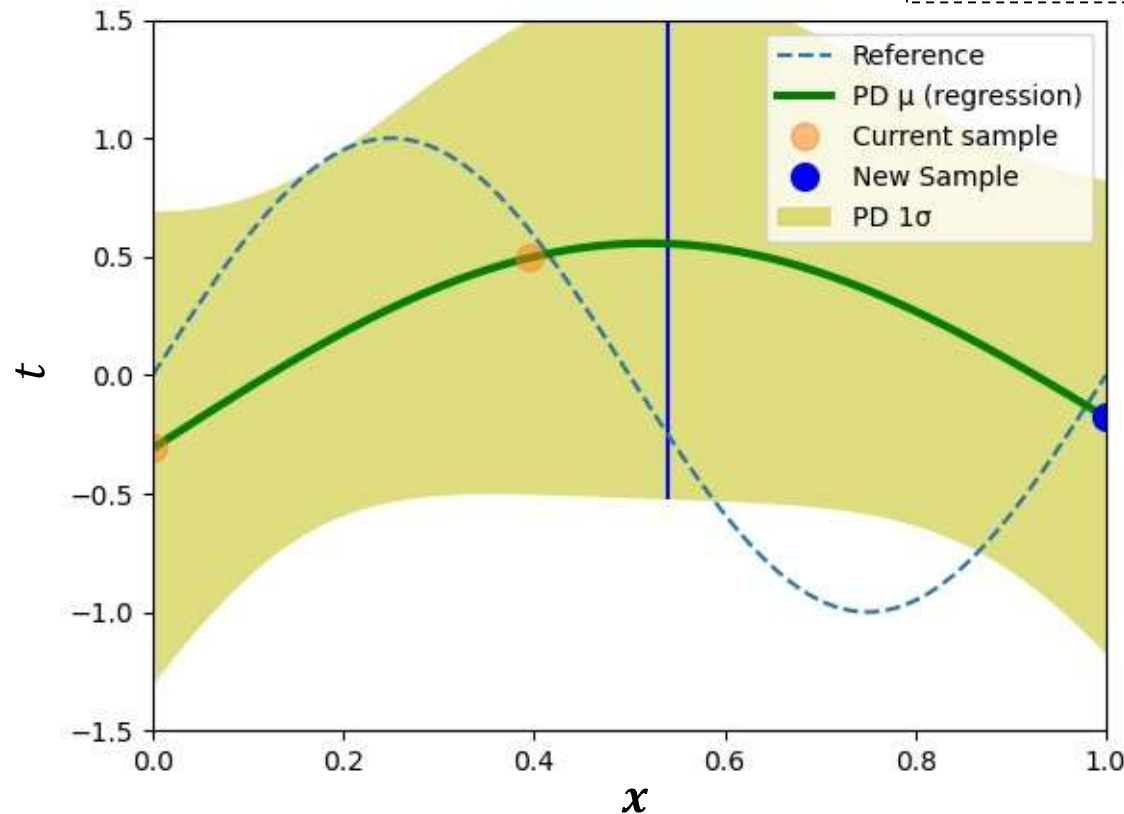
Adding a new point at the location x where $\sigma_N'^2(x)$ (or $\sigma_N^2(x)$) is max

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)



Bayesian Linear Regression

The probabilistic model

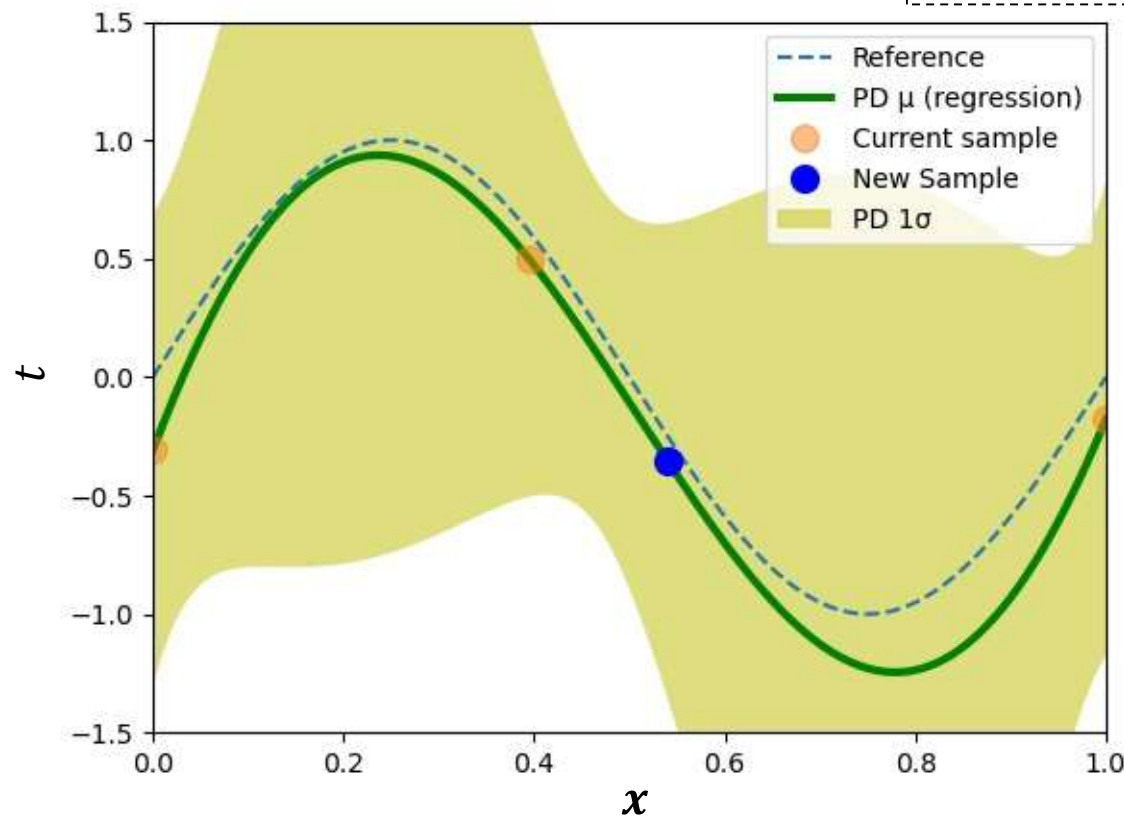
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

Examples

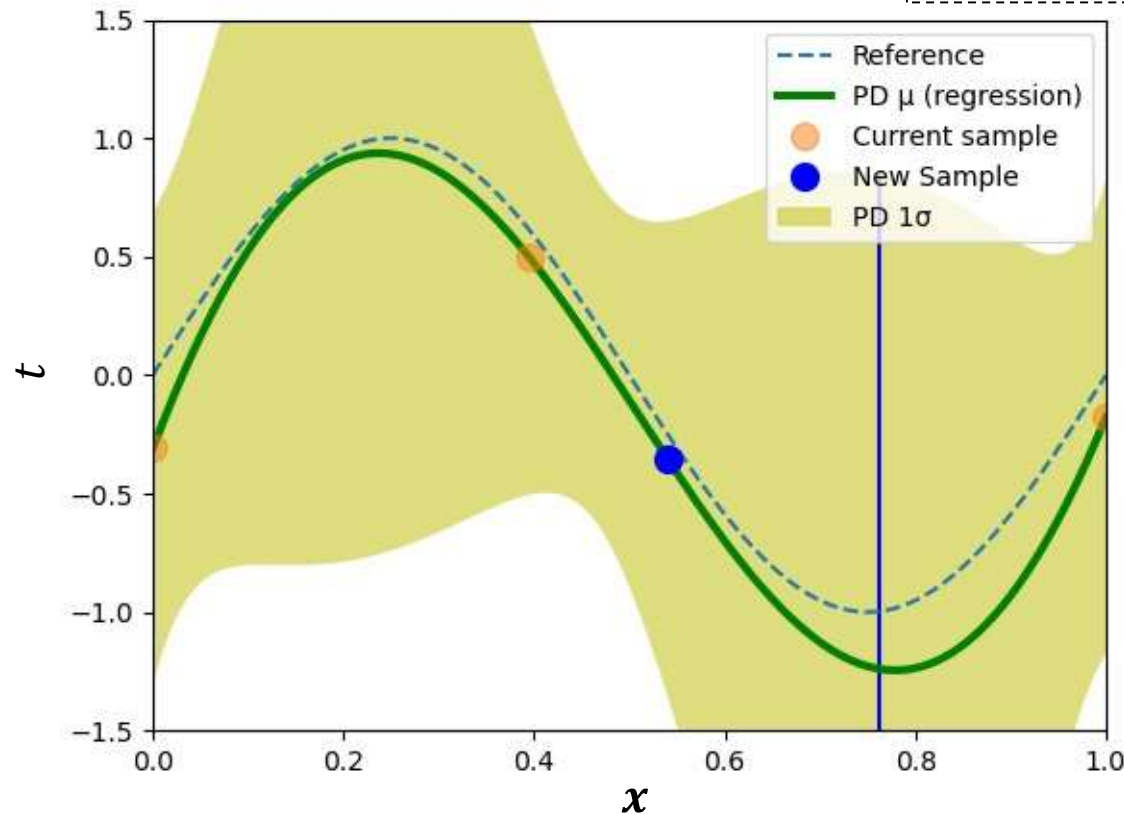
Adding a new point at the location x where $\sigma_N'^2(x)$ (or $\sigma_N^2(x)$) is max

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)



Bayesian Linear Regression

The probabilistic model

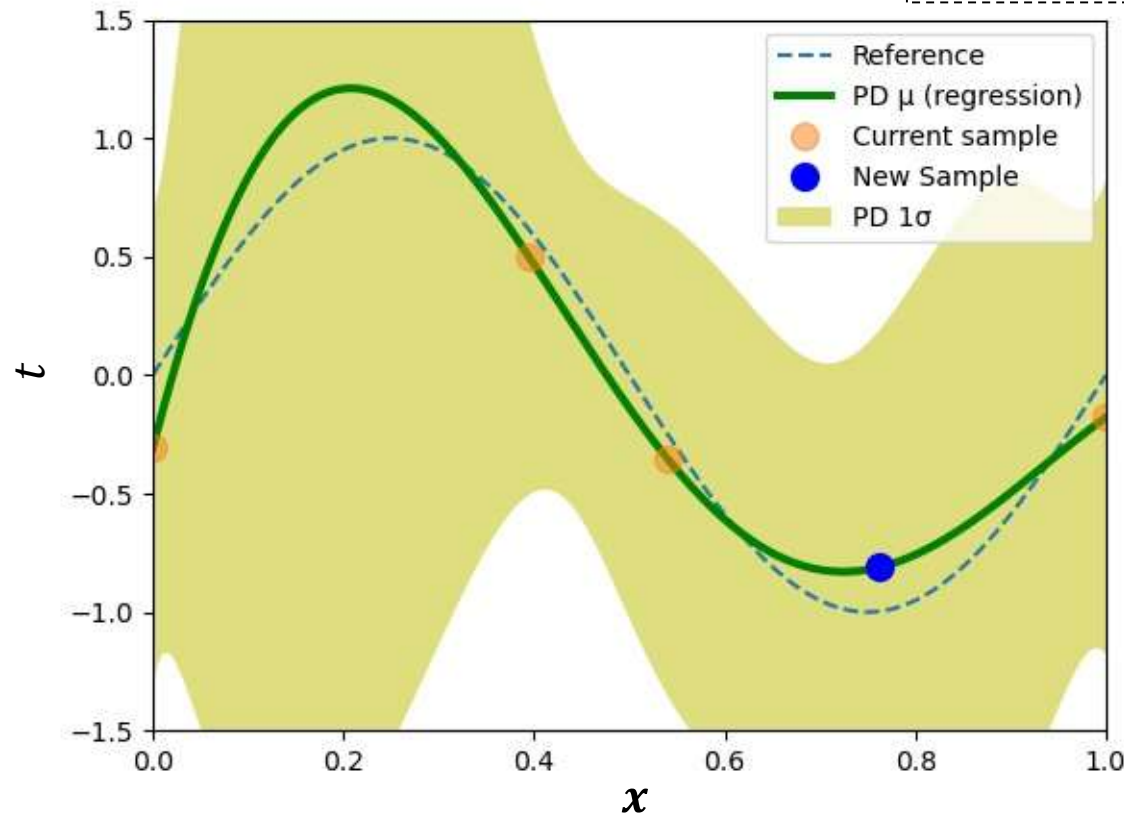
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

The probabilistic model

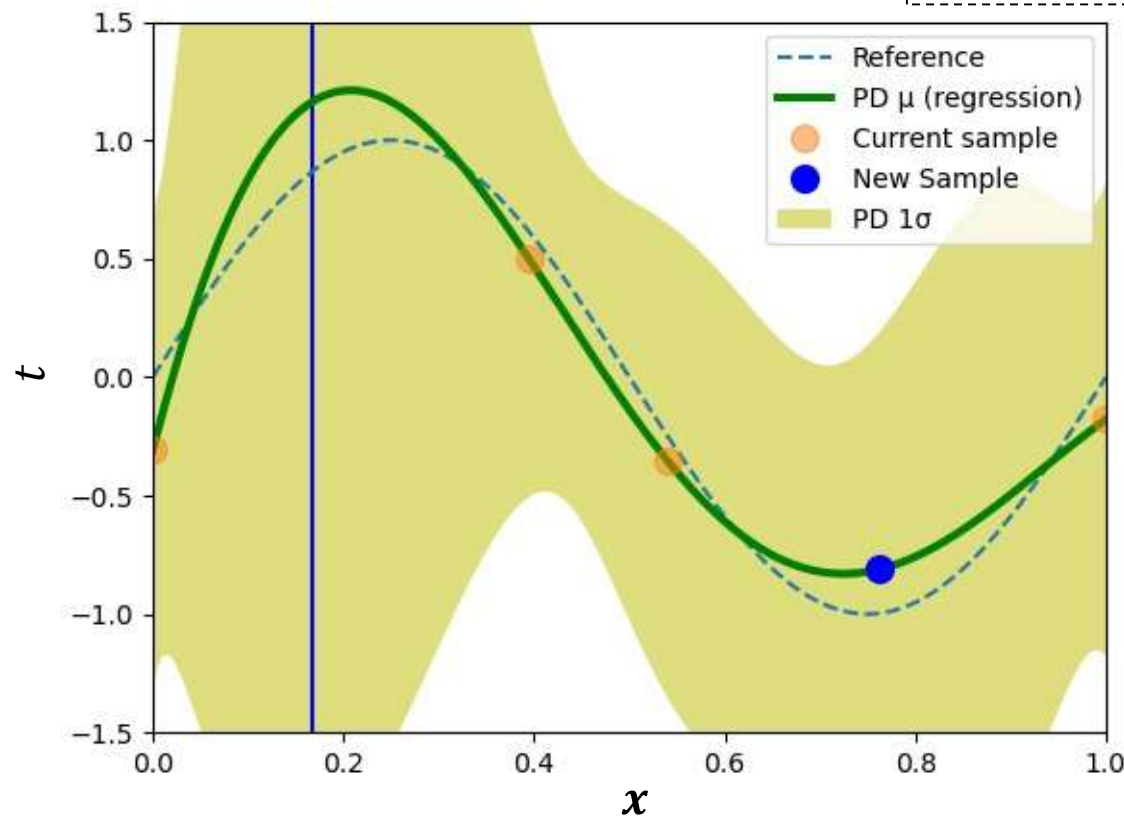
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

The probabilistic model

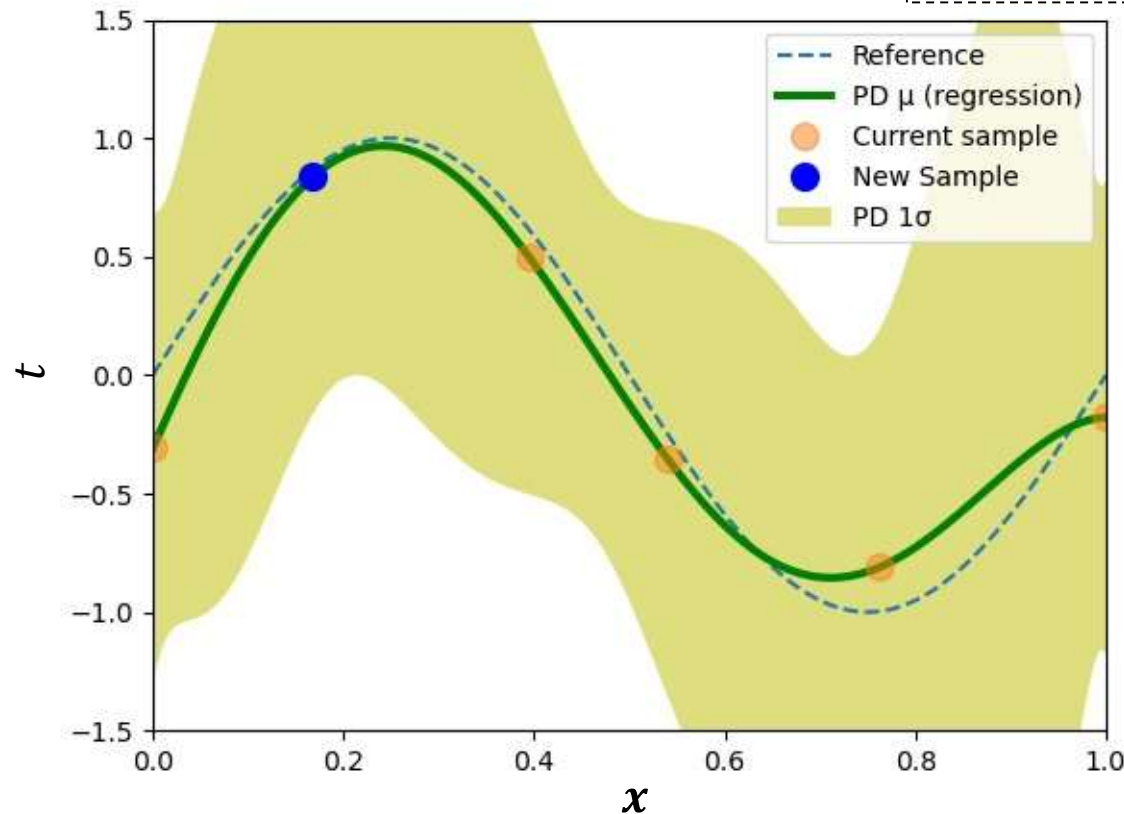
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

Examples

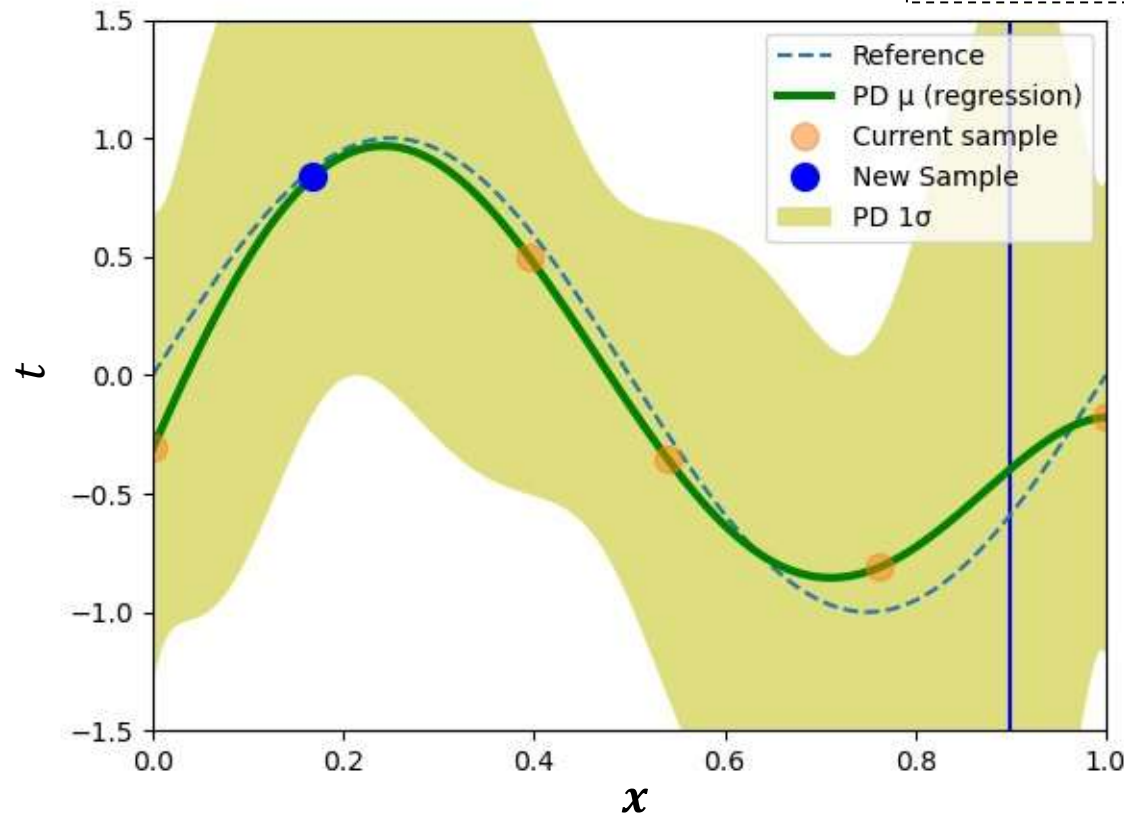
Adding a new point at the location x where $\sigma_N'^2(x)$ (or $\sigma_N^2(x)$) is max

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)



Technische
Universität
Braunschweig

Bayesian Linear Regression

The probabilistic model

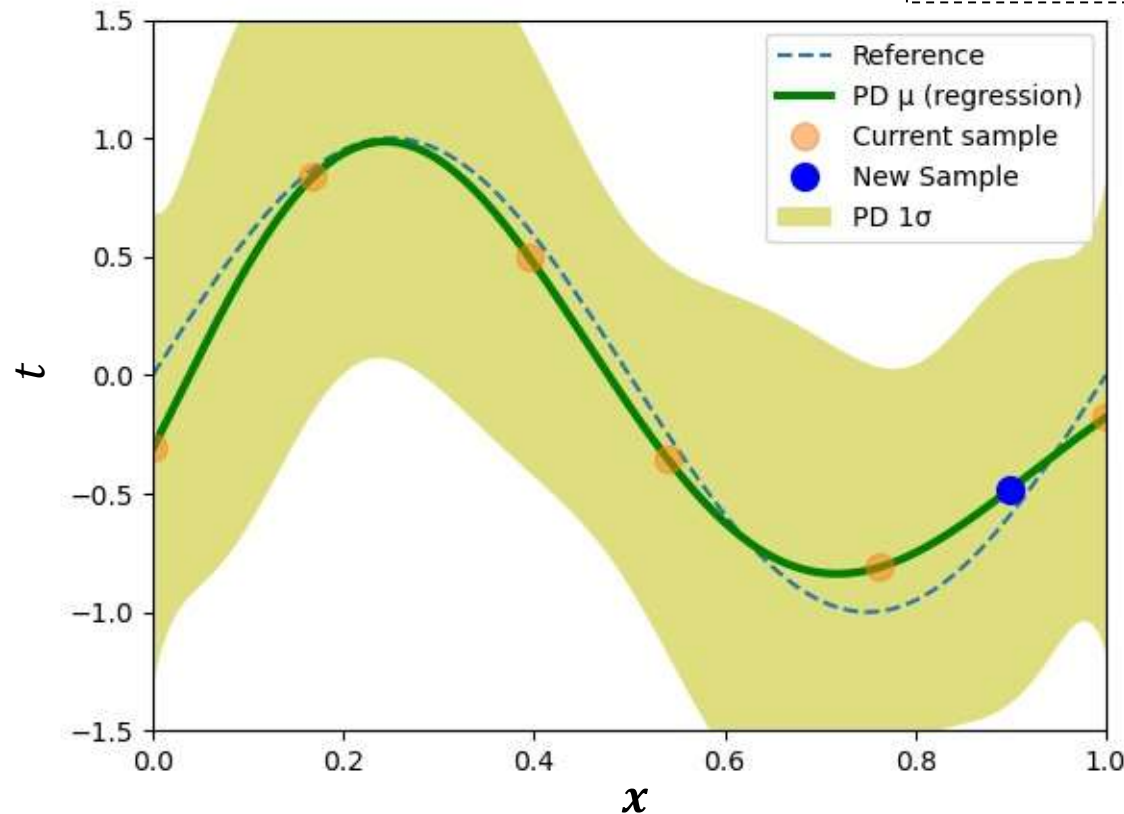
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

The probabilistic model

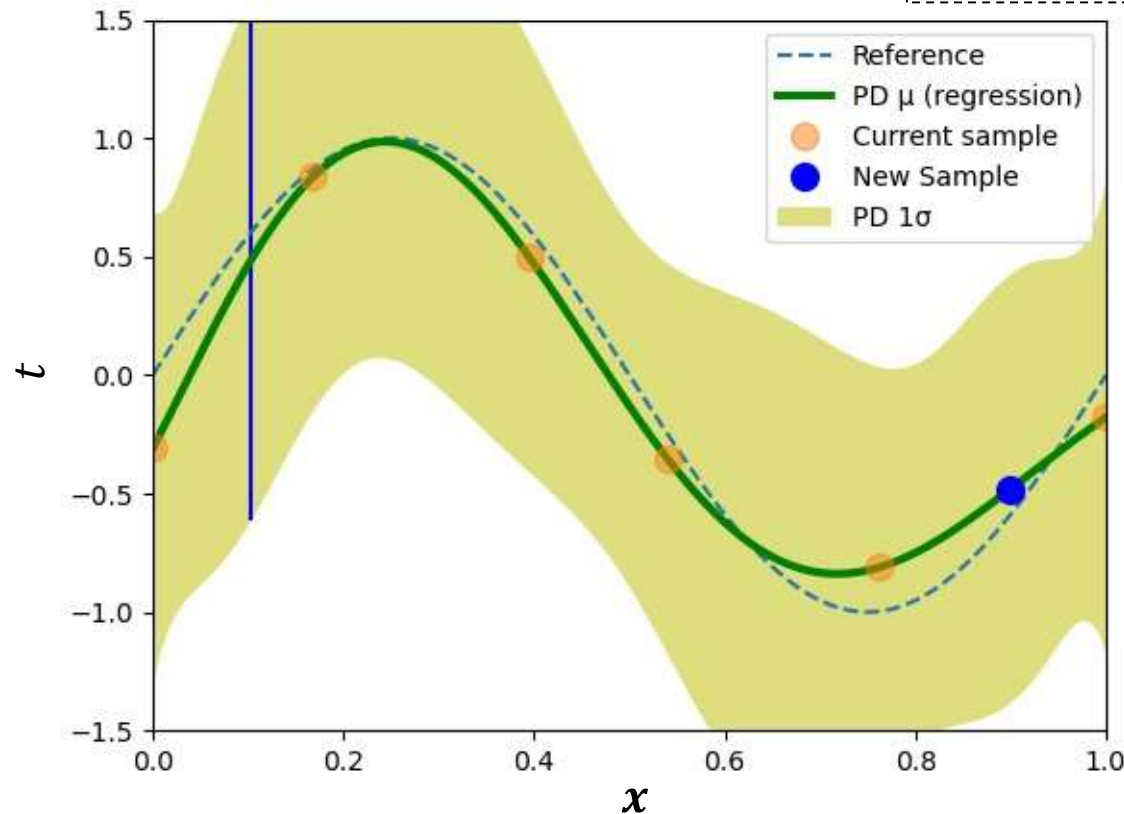
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

The probabilistic model

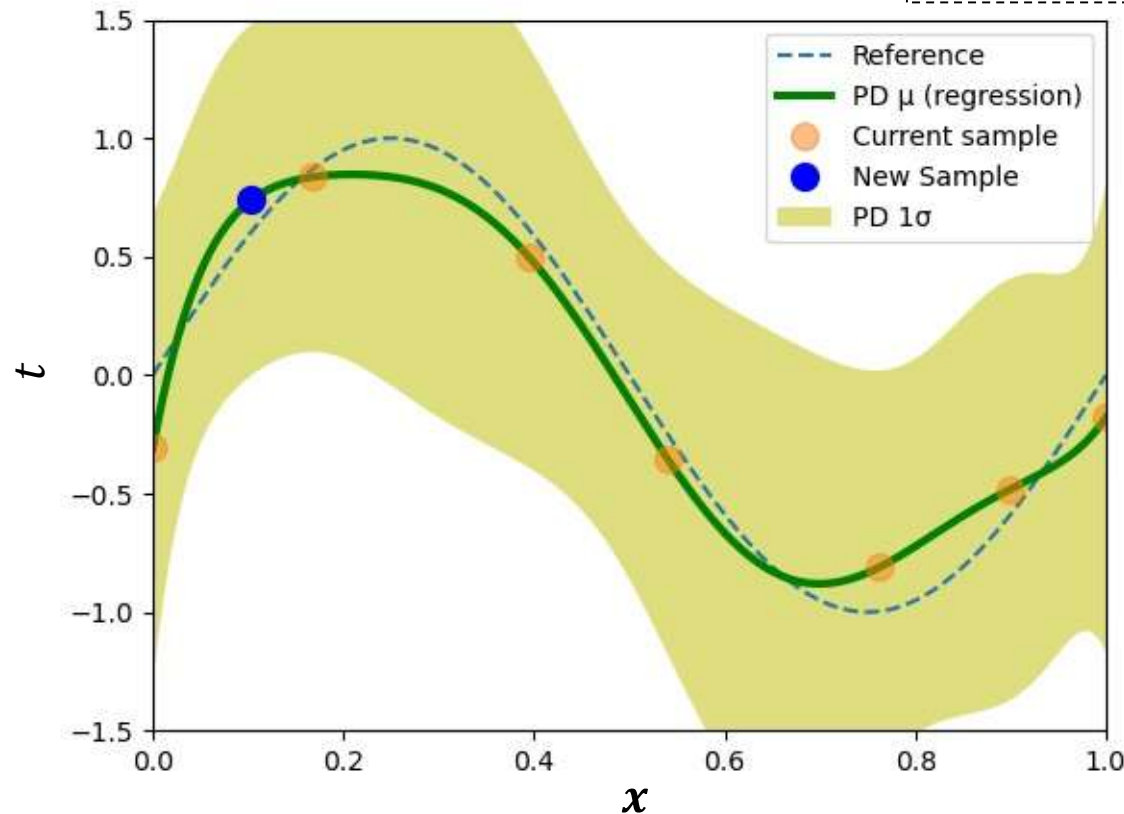
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

The probabilistic model

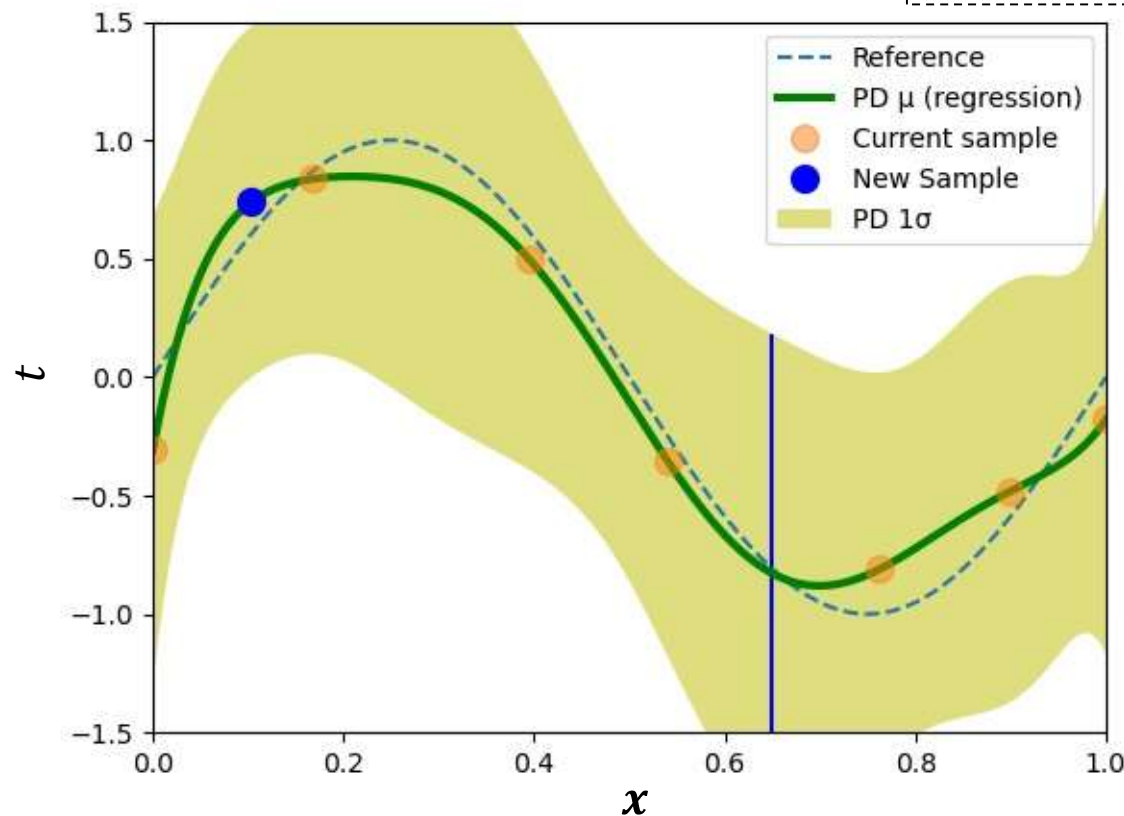
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Bayesian Linear Regression

The probabilistic model

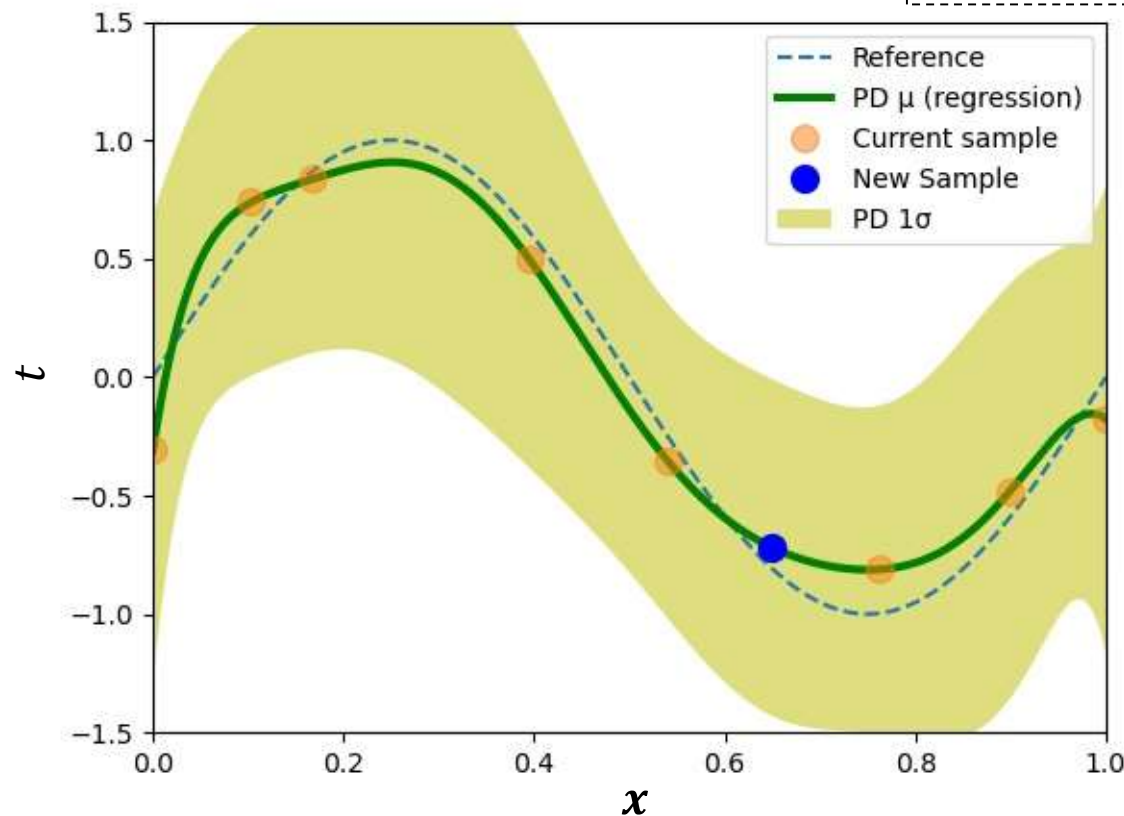
$$p(t|x, \mathbf{w}) = \mathcal{N}(t | \mathbf{w}^T \boldsymbol{\phi}(x), \hat{\sigma}^2)$$

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

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

Examples

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



Prior in General

Prior:

- Anything is fine.
 - Evidence in past
 - Regularization
 - Imaginary sample data
 - Your belief
 - etc.



Summary

Generalized Process

1. Define a probabilistic model
2. Then, **compute the posterior**
 - (Define a prior distribution)
 - **Point estimate** (Deterministic)
 - **Probability distribution** (Stochastic) ← computations hard

By using special cases,

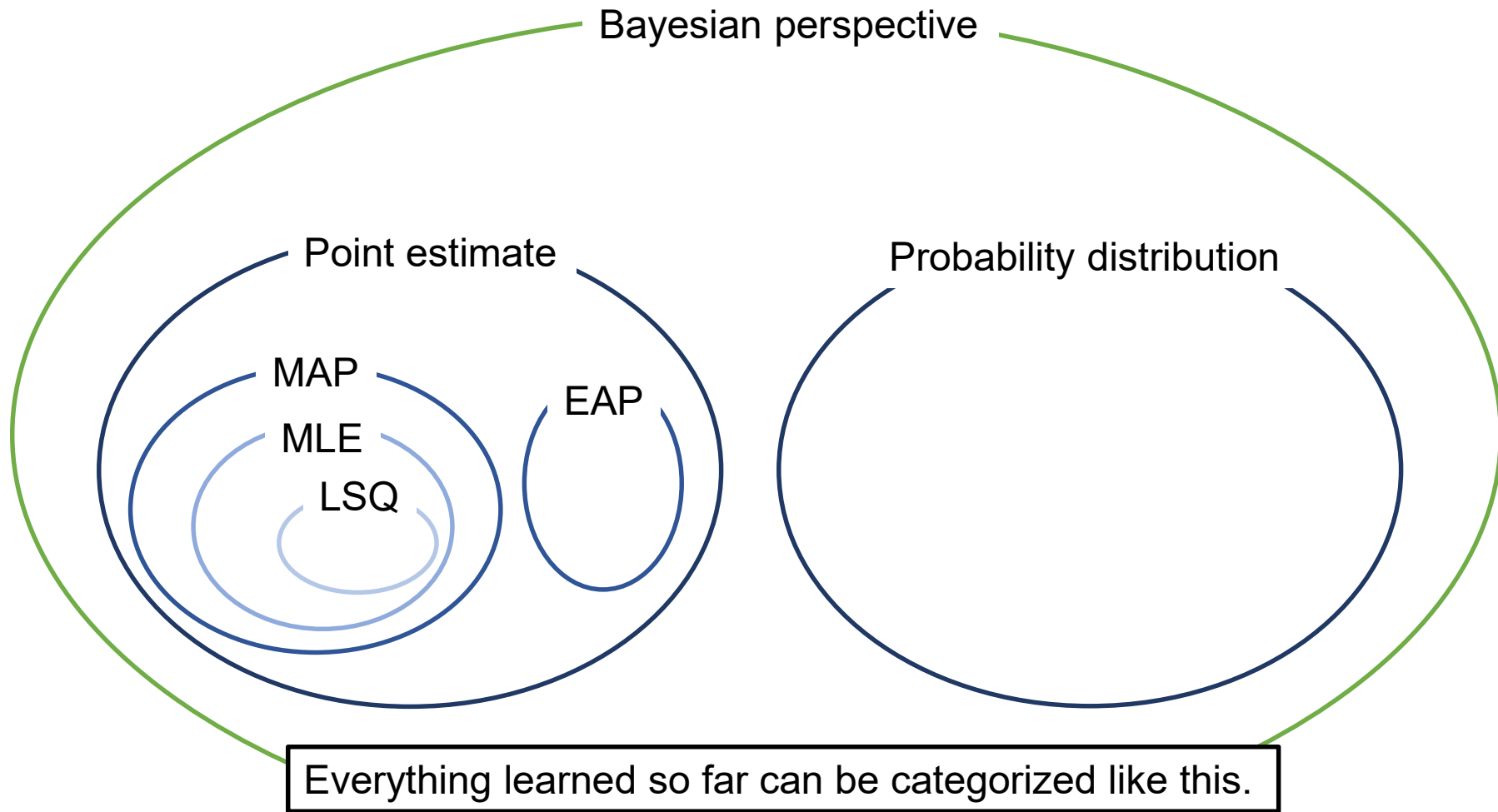
We could see what we can do by the Bayesian perspectives

➡ We could know: Generalized perspectives + Various possibilities

There are numerical techniques – e.g. MCMC (Lecture 12) to compute the difficult parts:

- Posterior distribution
- Predictive distribution (as the goal)

Bayesian Approach – Generalized Perspective



- Regularization
- Bayesian sequential learning
- ...