

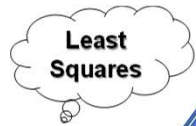
Scientific Machine Learning

Lecture 4: Linear Regression

Dr. Daigo Maruyama

Prof. Dr. Ali Elham

Current Position



Basis

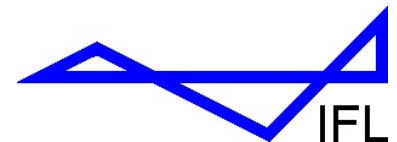


Today at Lecture 4



Technische
Universität
Braunschweig

Dr. Daigo Maruyama | Scientific Machine Learning: Lecture 4 | Slide 2



Lecture content

- Gaussian Distribution
 - Sampling Methods (Design of Experiments)
 - Linear Regression
- } continued from
Lecture 3

The lecture of this time basically follows the 2nd and 3rd chapters 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.

Gaussian Distribution (Normal Distribution)



Carl Friedrich Gauss
(1777-1855)

Born in Braunschweig

Collegium Carolinum at TUBS

Some important topics
related to Gaussian distributions

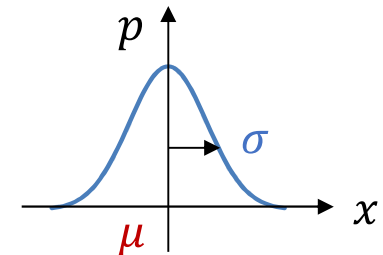
- Least square method
- Central limit theorem
- Gaussian Process

Gaussian Distribution (Normal Distribution)

Used in Regression

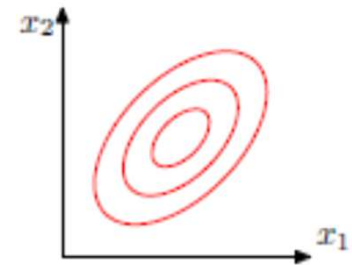
Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}$$



Multivariate Gaussian distribution

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{(\sqrt{2\pi})^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right\}$$



Some useful properties

- conditional distribution
 - marginal distribution
- ➡ also Gaussian distributions



Gaussian Distribution (Normal Distribution)

Multivariate Gaussian distribution

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(\sqrt{2\pi})^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

$$\mathbf{x} = (x_1, \dots, x_D)^T$$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_D)^T$$

D parameters

\mathbf{x} : random variable (vector)

$\boldsymbol{\mu}$: mean (vector)

$\boldsymbol{\Sigma}$: covariance (matrix)

$$\boldsymbol{\Sigma} = \begin{pmatrix} \text{var}[1] & \cdots & \text{cov}[1, D] \\ \vdots & \ddots & \vdots \\ \text{cov}[D, 1] & \cdots & \text{var}[D] \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \cdots & \sigma_{1,D} \\ \vdots & \ddots & \vdots \\ \sigma_{D,1} & \cdots & \sigma_D^2 \end{pmatrix}$$

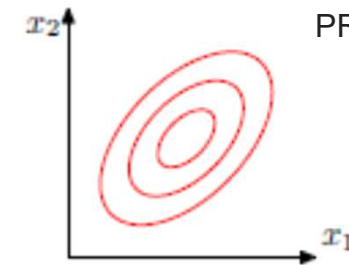
$\frac{D(D+1)}{2}$ parameters

➡ MLE?

$$\text{cov}[i, j] = \text{cov}[j, i]$$

Important properties of the covariance matrix $\boldsymbol{\Sigma}$:

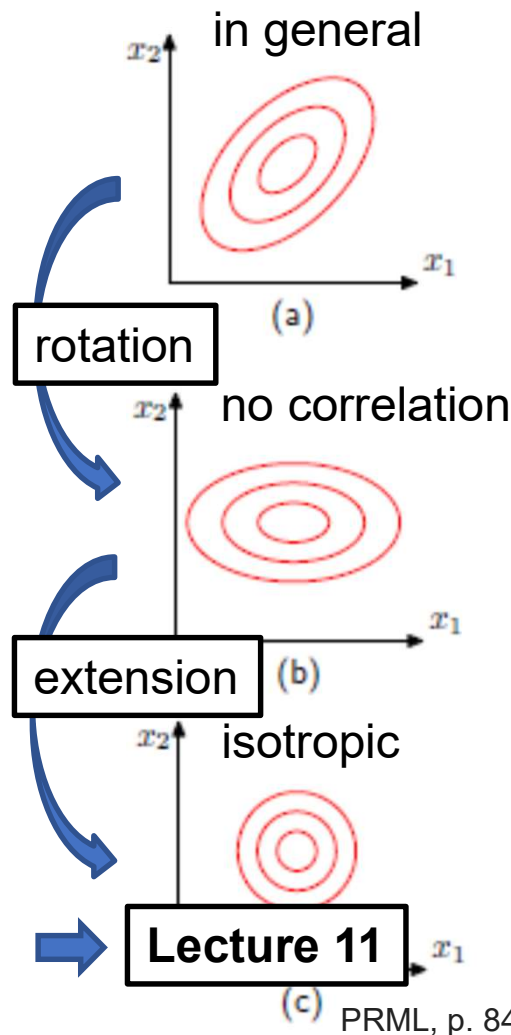
- Symmetric
- Positive (semi)definite



PRML, p. 84

Contour map of a multivariate Gaussian distribution when $D = 2$

Gaussian Distribution (Normal Distribution)



$$\Sigma_{(a)} = \begin{pmatrix} \sigma_1^2 & \cdots & \sigma_{1,D} \\ \vdots & \ddots & \vdots \\ \sigma_{D,1} & \cdots & \sigma_D^2 \end{pmatrix}$$

$$\mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$\frac{D(D+1)}{2}$ parameters

$$\Sigma_{(b)} = \begin{pmatrix} \sigma_1^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_D^2 \end{pmatrix}$$

$$\mathcal{N}(y|\mu, \sigma_i^2) \quad i = 1, \dots, M$$

for each axis

D parameters

$$\Sigma_{(c)} = \sigma^2 \mathbf{I} = \begin{pmatrix} \sigma^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma^2 \end{pmatrix}$$

$$\mathcal{N}(y|\mu, \sigma^2)$$

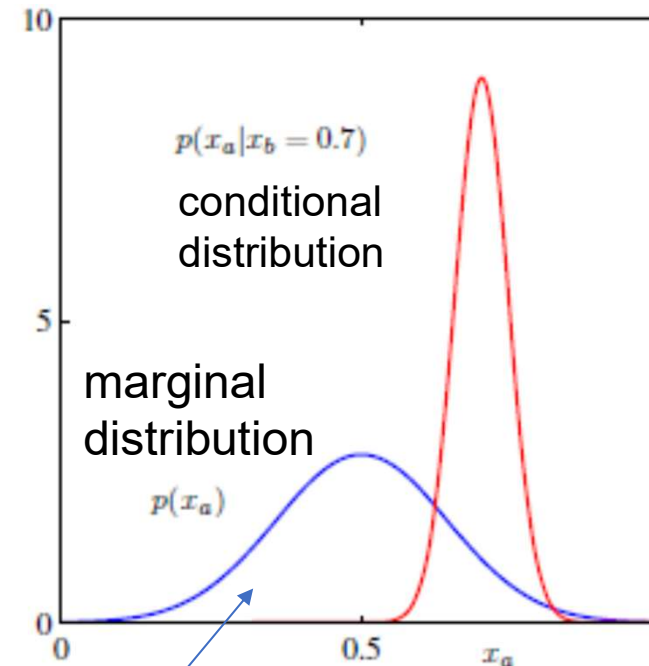
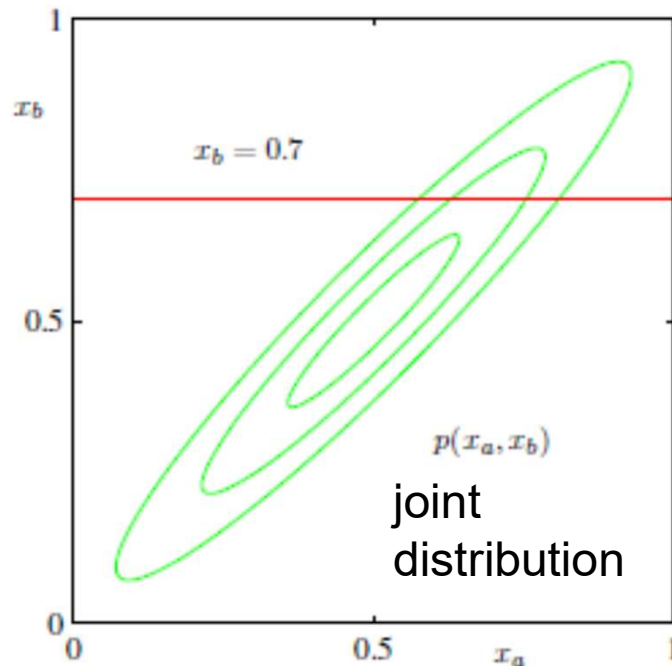
for each axis

1 parameter

PRML, p. 84

Gaussian Distribution (Normal Distribution)

Please imagine a mountain from the top view and a side view
(but the view has to be scaled to satisfy the area=1)



PRML, p. 90

Important property: Both of them become Gaussian distributions.

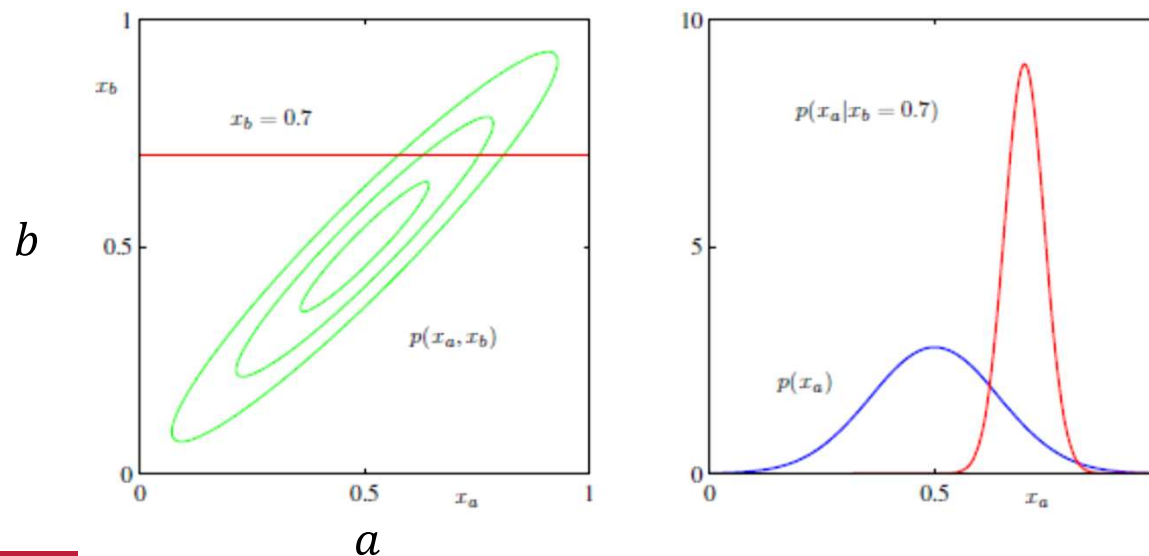
Gaussian Distribution (Normal Distribution)

We have now a multivariate Gaussian distribution:

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

When we partition \mathbf{x} into two disjoint subsets \mathbf{x}_a and \mathbf{x}_b :

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$



Gaussian Distribution (Normal Distribution)

We have now a multivariate Gaussian distribution:

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

When we partition \mathbf{x} into two disjoint subsets \mathbf{x}_a and \mathbf{x}_b :

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

What if :

- **Marginal distribution** $p(\mathbf{x}_a)$
- **Conditional distribution** $p(\mathbf{x}_a|\mathbf{x}_b)$
- **(Joint distribution** $p(\mathbf{x}) = p(\mathbf{x}_a, \mathbf{x}_b) = p(\mathbf{x}_a|\mathbf{x}_b)p(\mathbf{x}_b)$)

In general, there is no guarantee that $p(\mathbf{x}_a)$ and $p(\mathbf{x}_a|\mathbf{x}_b)$ can be represented by the same distributions $p(\mathbf{x})$.

Gaussian Distribution (Normal Distribution)

We have now a multivariate Gaussian distribution:

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

When we partition \mathbf{x} into two disjoint subsets \mathbf{x}_a and \mathbf{x}_b :

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

Marginal distribution $p(\mathbf{x}_a)$

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

simple and intuitive

pretending that we did not see any information from the subset \mathbf{x}_b .



The basic concept of marginal distributions

Gaussian Distribution (Normal Distribution)

We have now a multivariate Gaussian distribution:

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

When we partition \mathbf{x} into two disjoint subsets \mathbf{x}_a and \mathbf{x}_b :

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

Conditional distribution $p(\mathbf{x}_a|\mathbf{x}_b)$

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\mathbf{x}_b, \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$

where,

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$$

No need to remember but the concept is important
in Gaussian Processes (Lectures 6-8)

Other Probability Distributions

- **Laplace distribution**

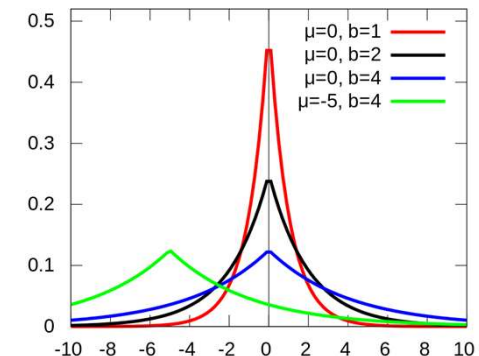
$$p(x|\mu, b) = \frac{1}{2b} \exp\left\{-\frac{|x - \mu|}{b}\right\}$$

- **Cauchy distribution** (a special case of **Student's t-distribution**)

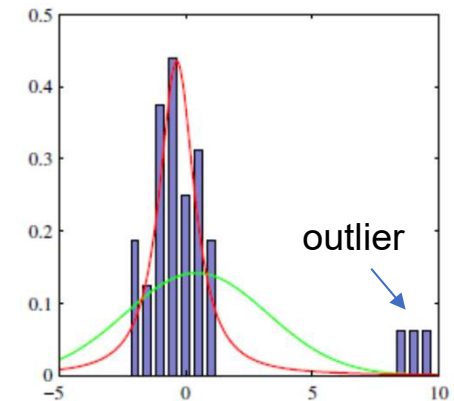
$$p(x|x_0, \gamma) = \frac{1}{\pi} \frac{\gamma}{(x - x_0)^2 + \gamma^2} \quad \text{used to treat outlier}$$

Most of the introduced probability distributions are categorized in Exponential Family.

$$p(x|\eta) = h(x)g(\eta)\exp\{\eta^T u(x)\}$$



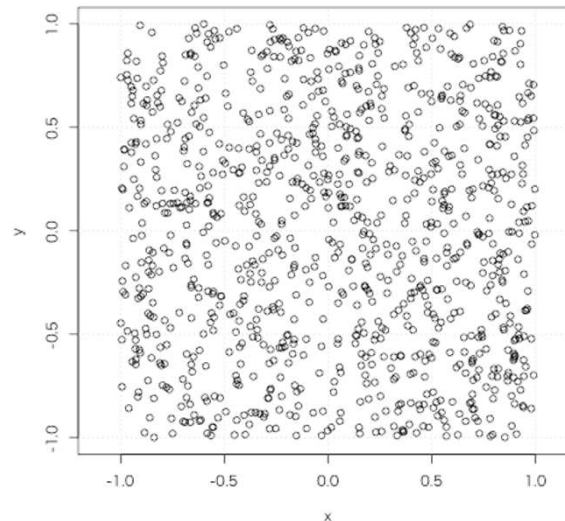
wikipedia



PRML, p. 104

red: Gaussian
Green: student's-t

3. Sampling Methods (Design of Experiments)



Curse of Dimensionality

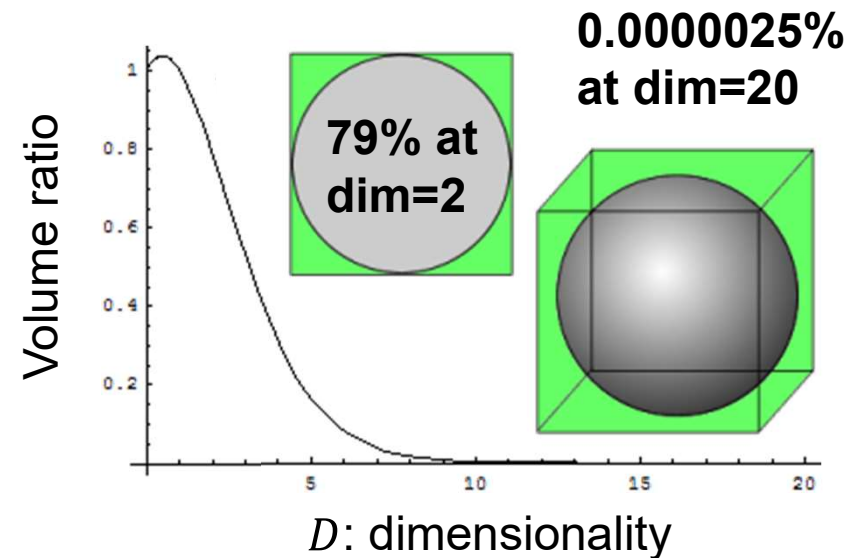
In high-dimensional space

The volume ratio between the cube and the sphere is counterintuitive.

almost skin (tiny volume)

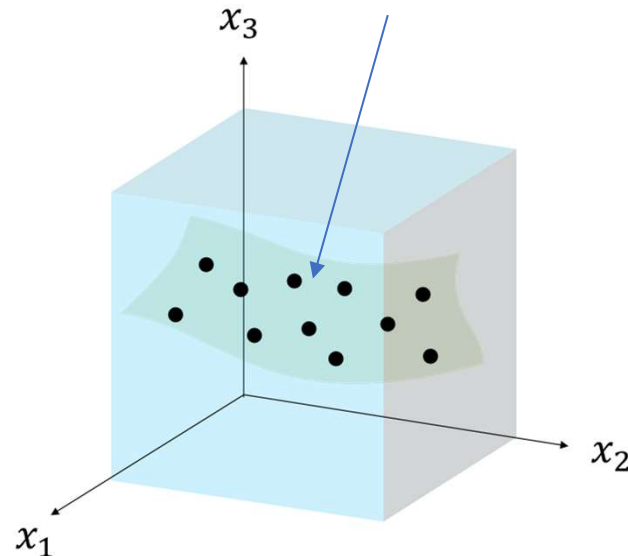
$$V_D = \frac{\pi^{D/2}}{(D/2)!} = \frac{\pi^{D/2}}{\Gamma(D/2 + 1)}$$

$$\frac{\text{volume of hypersphere}}{\text{volume of hypercube}} = \frac{V_D}{2^D} = \frac{\pi^{D/2}}{2^D (D/2)!} \longrightarrow 0$$



Curse of Dimensionality

a low-dimensional space - manifold



only two input actually:

- Rotation
- Translation

10,000 dims to 2 dims

➡ **Lecture 11**

Dimensionality Reduction (data is lying on a low-dimensional space - manifold)

- Big data: plentiful data in hand
- **Traditional engineering design: data is produced ➡ Design or Experiments**

The parameters (design variables) are selected from engineering viewpoints.

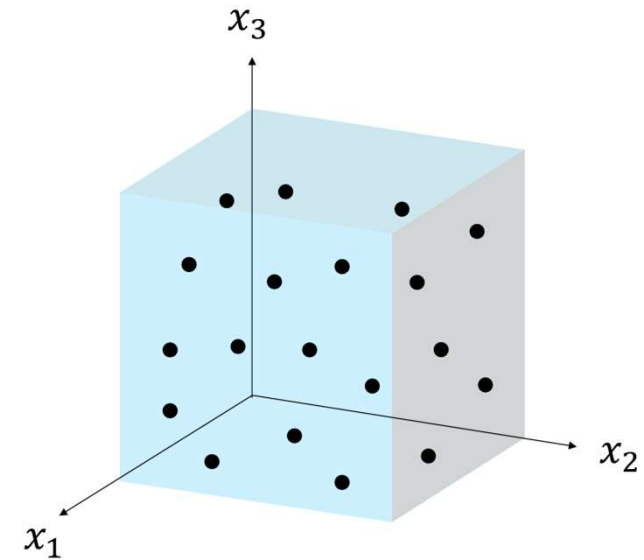
Design of Experiments (DoE)

Uniform distributions (**Design of Experiments**)

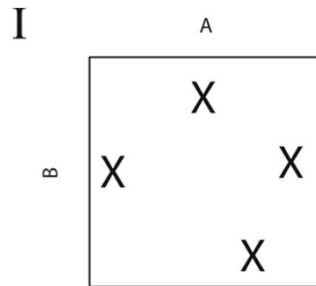
- Monte Carlo (MC) sampling method
- Latin Hypercube sampling (LHS) method
- Quasi Monte Carlo (QMC) sampling method
 - Halton sequence
 - Sobol sequence

Arbitrary distributions (**Lecture 12**)

- Markov-Chain Monte Carlo (MCMC)



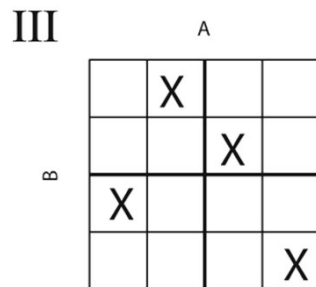
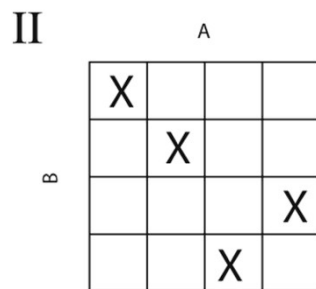
Design of Experiments (DoE)



Latin hypercube sampling

Properties in practical use:

- The partition has to be defined first (the sample size defined).
- New sample points cannot be added.



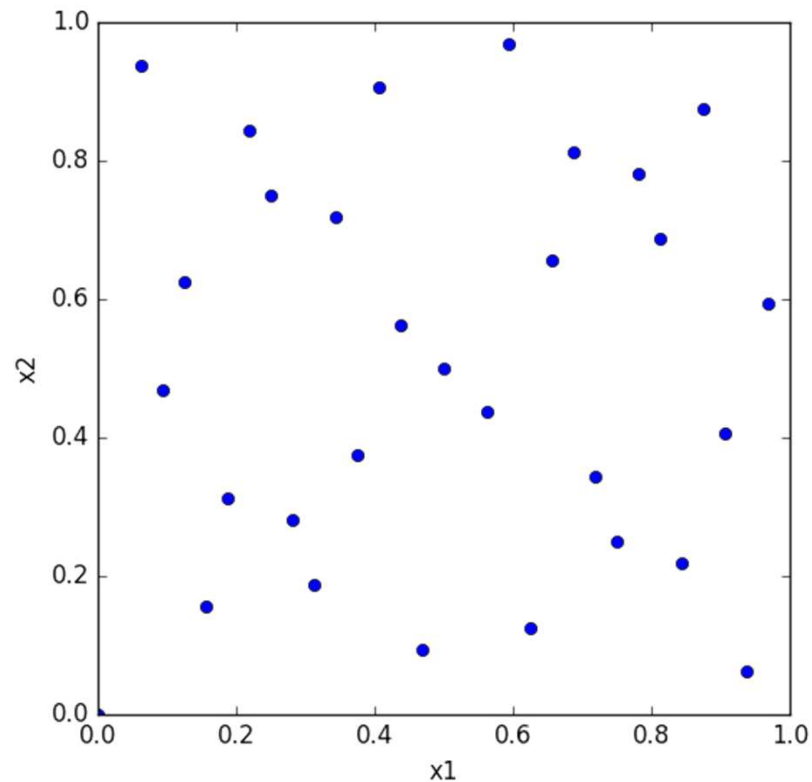
wikipedia:
Latin hypercube sampling

Design of Experiments (DoE)

```
[[ 0. 0. ]
 [ 0.5 0.5 ]
 [ 0.75 0.25 ]
 [ 0.25 0.75 ]
 [ 0.375 0.375 ]
 [ 0.875 0.875 ]
 [ 0.625 0.125 ]
 [ 0.125 0.625 ]
```

```
[ 0.1875 0.3125 ]
 [ 0.6875 0.8125 ]
 [ 0.9375 0.0625 ]
 [ 0.4375 0.5625 ]
 [ 0.3125 0.1875 ]
 [ 0.8125 0.6875 ]
 [ 0.5625 0.4375 ]
 [ 0.0625 0.9375 ]
 [ 0.09375 0.46875 ]
 [ 0.59375 0.96875 ]
 [ 0.84375 0.21875 ]
 [ 0.34375 0.71875 ]
 [ 0.46875 0.09375 ]
 [ 0.96875 0.59375 ]
 [ 0.71875 0.34375 ]
 [ 0.21875 0.84375 ]
 [ 0.15625 0.15625 ]
 [ 0.65625 0.65625 ]
 [ 0.90625 0.40625 ]
 [ 0.40625 0.90625 ]
 [ 0.28125 0.28125 ]
 [ 0.78125 0.78125 ]
 [ 0.53125 0.03125 ]
 [ 0.03125 0.53125 ]
```

```
....
 [ 0.05273438 0.23242188]
 [ 0.55273438 0.73242188]
 [ 0.80273438 0.48242188]
 [ 0.30273438 0.98242188]
 [ 0.42773438 0.35742188]
 [ 0.92773438 0.85742188]
 [ 0.67773438 0.10742188]
 [ 0.17773438 0.60742188]
 [ 0.24023438 0.41992188]
 [ 0.74023438 0.91992188]
 [ 0.99023438 0.16992188]
 [ 0.49023438 0.66992188]
```



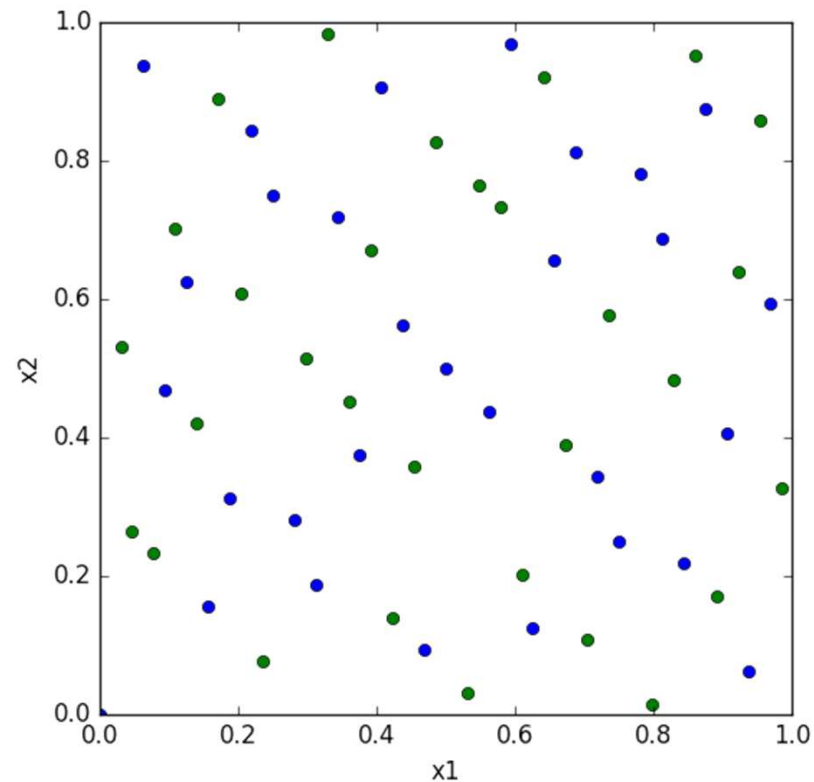
Sobol sequence (Quasi Monte Carlo sampling)

- Low-discrepancy sequence (uniformity)
- Reproducibility
- The uniformity is kept in high dimensions (especially Sobol).

Design of Experiments (DoE)

```
[[ 0.5 0.5 ]
 [ 0.75 0.25 ]
 [ 0.25 0.75 ]
 [ 0.375 0.375 ]
 [ 0.875 0.875 ]
 [ 0.625 0.125 ]
 [ 0.125 0.625 ]
 [ 0.1875 0.3125 ]
 [ 0.6875 0.8125 ]
 [ 0.9375 0.0625 ]
 [ 0.4375 0.5625 ]
 [ 0.3125 0.1875 ]
 [ 0.8125 0.6875 ]
 [ 0.5625 0.4375 ]
 [ 0.0625 0.9375 ]
 [ 0.09375 0.46875 ]
 [ 0.59375 0.96875 ]
 [ 0.84375 0.21875 ]
 [ 0.34375 0.71875 ]
 [ 0.46875 0.09375 ]
 [ 0.96875 0.59375 ]
 [ 0.71875 0.34375 ]
 [ 0.21875 0.84375 ]
 [ 0.15625 0.15625 ]
 [ 0.65625 0.65625 ]
 [ 0.90625 0.40625 ]
 [ 0.40625 0.90625 ]
 [ 0.28125 0.28125 ]
 [ 0.78125 0.78125 ]
 [ 0.53125 0.03125 ]
 [ 0.03125 0.53125 ]
....
```

```
[ 0.05273438 0.23242188]
 [ 0.55273438 0.73242188]
 [ 0.80273438 0.48242188]
 [ 0.30273438 0.98242188]
 [ 0.42773438 0.35742188]
 [ 0.92773438 0.85742188]
 [ 0.67773438 0.10742188]
 [ 0.17773438 0.60742188]
 [ 0.24023438 0.41992188]
 [ 0.74023438 0.91992188]
 [ 0.99023438 0.16992188]
 [ 0.49023438 0.66992188]]
```

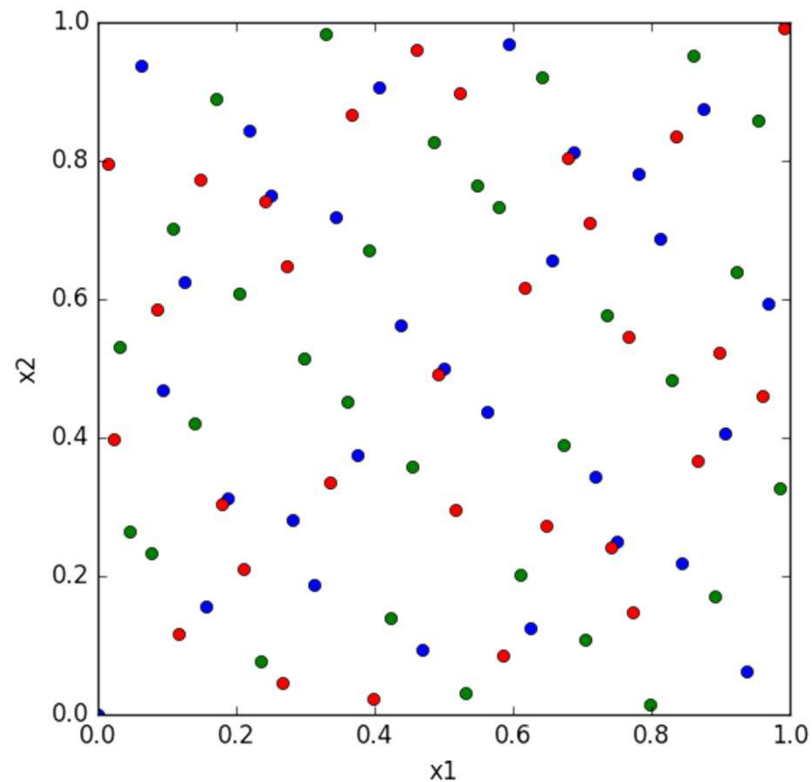


New sample points can be added keeping the uniformity.

Design of Experiments (DoE)

```
[[ 0.5 0.5 ]
 [ 0.75 0.25 ]
 [ 0.25 0.75 ]
 [ 0.375 0.375 ]
 [ 0.875 0.875 ]
 [ 0.625 0.125 ]
 [ 0.125 0.625 ]
 [ 0.1875 0.3125 ]
 [ 0.6875 0.8125 ]
 [ 0.9375 0.0625 ]
 [ 0.4375 0.5625 ]
 [ 0.3125 0.1875 ]
 [ 0.8125 0.6875 ]
 [ 0.5625 0.4375 ]
 [ 0.0625 0.9375 ]
 [ 0.09375 0.46875 ]
 [ 0.59375 0.96875 ]
 [ 0.84375 0.21875 ]
 [ 0.34375 0.71875 ]
 [ 0.46875 0.09375 ]
 [ 0.96875 0.59375 ]
 [ 0.71875 0.34375 ]
 [ 0.21875 0.84375 ]
 [ 0.15625 0.15625 ]
 [ 0.65625 0.65625 ]
 [ 0.90625 0.40625 ]
 [ 0.40625 0.90625 ]
 [ 0.28125 0.28125 ]
 [ 0.78125 0.78125 ]
 [ 0.53125 0.03125 ]
 [ 0.03125 0.53125 ]
....
```

```
[ 0.05273438 0.23242188]
 [ 0.55273438 0.73242188]
 [ 0.80273438 0.48242188]
 [ 0.30273438 0.98242188]
 [ 0.42773438 0.35742188]
 [ 0.92773438 0.85742188]
 [ 0.67773438 0.10742188]
 [ 0.17773438 0.60742188]
 [ 0.24023438 0.41992188]
 [ 0.74023438 0.91992188]
 [ 0.99023438 0.16992188]
 [ 0.49023438 0.66992188]]
```

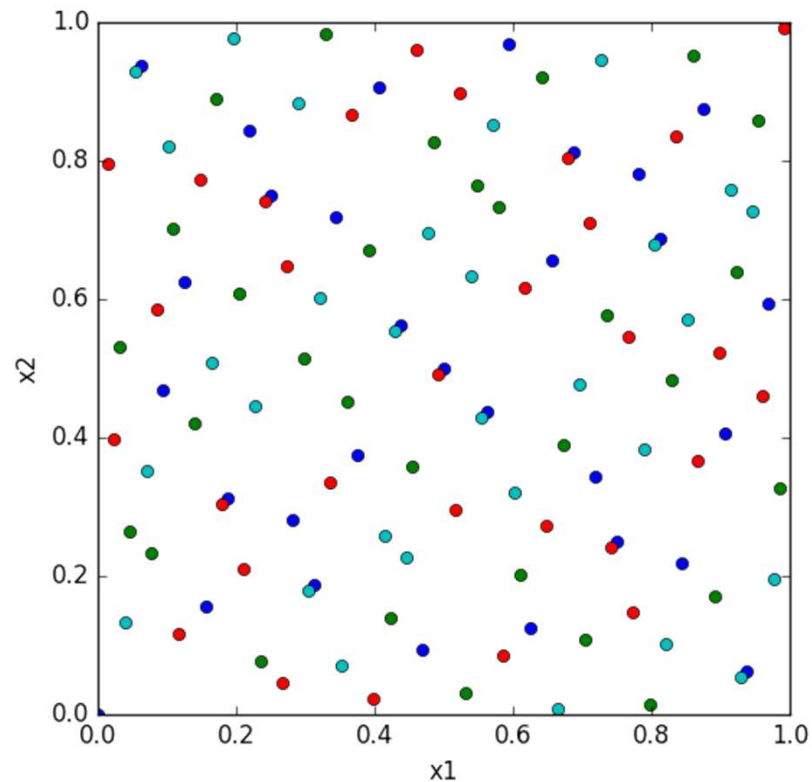


New sample points can be added keeping the uniformity.

Design of Experiments (DoE)

```
[[ 0.5 0.5 ]
[ 0.75 0.25 ]
[ 0.25 0.75 ]
[ 0.375 0.375 ]
[ 0.875 0.875 ]
[ 0.625 0.125 ]
[ 0.125 0.625 ]
[ 0.1875 0.3125 ]
[ 0.6875 0.8125 ]
[ 0.9375 0.0625 ]
[ 0.4375 0.5625 ]
[ 0.3125 0.1875 ]
[ 0.8125 0.6875 ]
[ 0.5625 0.4375 ]
[ 0.0625 0.9375 ]
[ 0.09375 0.46875 ]
[ 0.59375 0.96875 ]
[ 0.84375 0.21875 ]
[ 0.34375 0.71875 ]
[ 0.46875 0.09375 ]
[ 0.96875 0.59375 ]
[ 0.71875 0.34375 ]
[ 0.21875 0.84375 ]
[ 0.15625 0.15625 ]
[ 0.65625 0.65625 ]
[ 0.90625 0.40625 ]
[ 0.40625 0.90625 ]
[ 0.28125 0.28125 ]
[ 0.78125 0.78125 ]
[ 0.53125 0.03125 ]
[ 0.03125 0.53125 ]
....
```

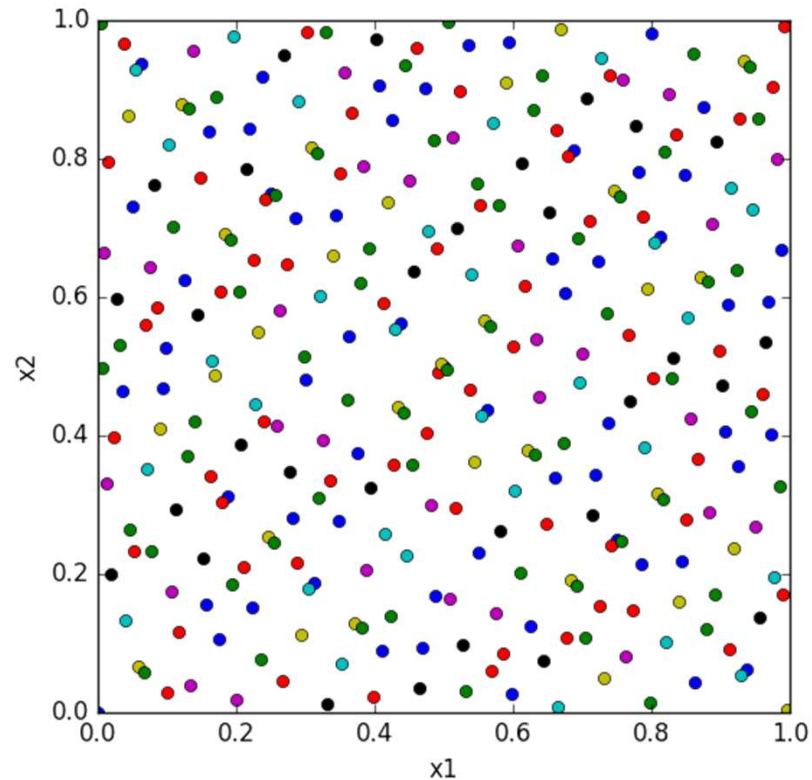
```
[ 0.05273438 0.23242188]
[ 0.55273438 0.73242188]
[ 0.80273438 0.48242188]
[ 0.30273438 0.98242188]
[ 0.42773438 0.35742188]
[ 0.92773438 0.85742188]
[ 0.67773438 0.10742188]
[ 0.17773438 0.60742188]
[ 0.24023438 0.41992188]
[ 0.74023438 0.91992188]
[ 0.99023438 0.16992188]
[ 0.49023438 0.66992188]]
```



New sample points can be added keeping the uniformity.

Design of Experiments (DoE)

```
[[ 0.5 0.5 ]
[ 0.75 0.25 ]
[ 0.25 0.75 ]
[ 0.375 0.375 ]
[ 0.875 0.875 ]
[ 0.625 0.125 ]
[ 0.125 0.625 ]
[ 0.1875 0.3125 ]
[ 0.6875 0.8125 ]
[ 0.9375 0.0625 ]
[ 0.4375 0.5625 ]
[ 0.3125 0.1875 ]
[ 0.8125 0.6875 ]
[ 0.5625 0.4375 ]
[ 0.0625 0.9375 ]
[ 0.09375 0.46875 ]
[ 0.59375 0.96875 ]
[ 0.84375 0.21875 ]
[ 0.34375 0.71875 ]
[ 0.46875 0.09375 ]
[ 0.96875 0.59375 ]
[ 0.71875 0.34375 ]
[ 0.21875 0.84375 ]
[ 0.15625 0.15625 ]
[ 0.65625 0.65625 ]
[ 0.90625 0.40625 ]
[ 0.40625 0.90625 ]
[ 0.28125 0.28125 ]
[ 0.78125 0.78125 ]
[ 0.53125 0.03125 ]
[ 0.03125 0.53125 ]
....
[ 0.05273438 0.23242188]
[ 0.55273438 0.73242188]
[ 0.80273438 0.48242188]
[ 0.30273438 0.98242188]
[ 0.42773438 0.35742188]
[ 0.92773438 0.85742188]
[ 0.67773438 0.10742188]
[ 0.17773438 0.60742188]
[ 0.24023438 0.41992188]
[ 0.74023438 0.91992188]
[ 0.99023438 0.16992188]
[ 0.49023438 0.66992188]]
```



New sample points can be added keeping the uniformity.

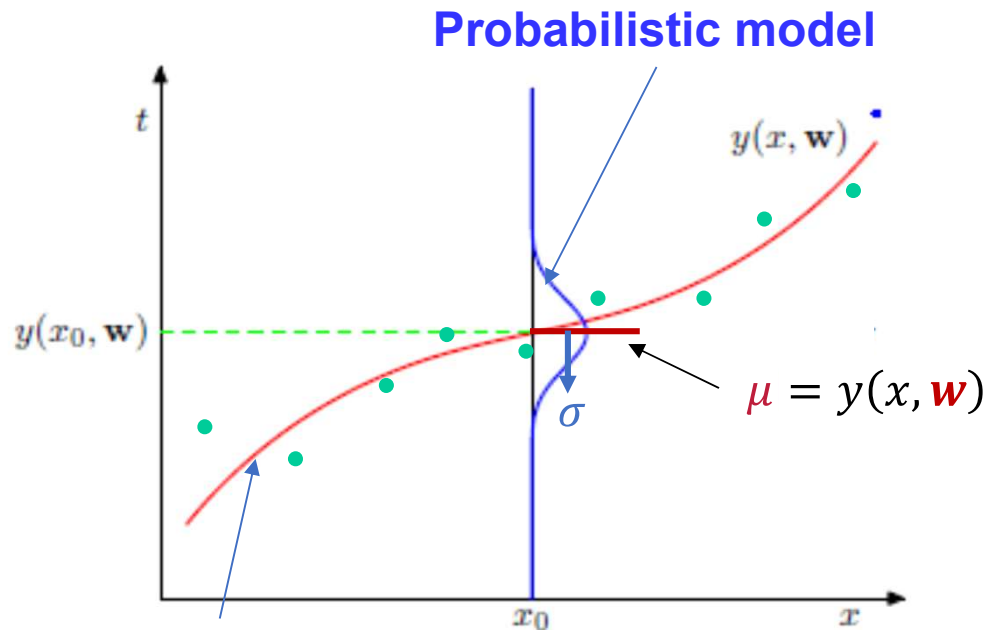
Lecture content

3. Linear Regression



Regression Model

In the curve fitting problem, we defined a probabilistic model.



Regression model

based on PRML, p. 29

x : deterministic variable
 t : random variable

Probabilistic model $p(t|x)$

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

e.g. $p(t|x)$ is a Gaussian distribution.

Regression model $E[t|x]$

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

e.g. $y(x, \mathbf{w})$ is a polynomial function.

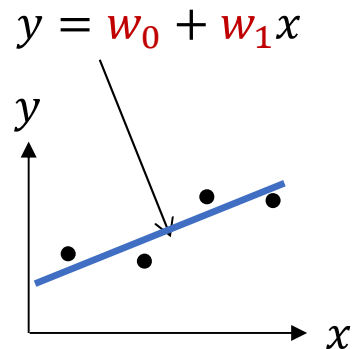
The regression model is a part of the probabilistic model*.

Regression model \in **Probabilistic model**

*There are a few exceptions (e.g. support vector machine).

Regression Model

The simplest linear regression



$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w})$$

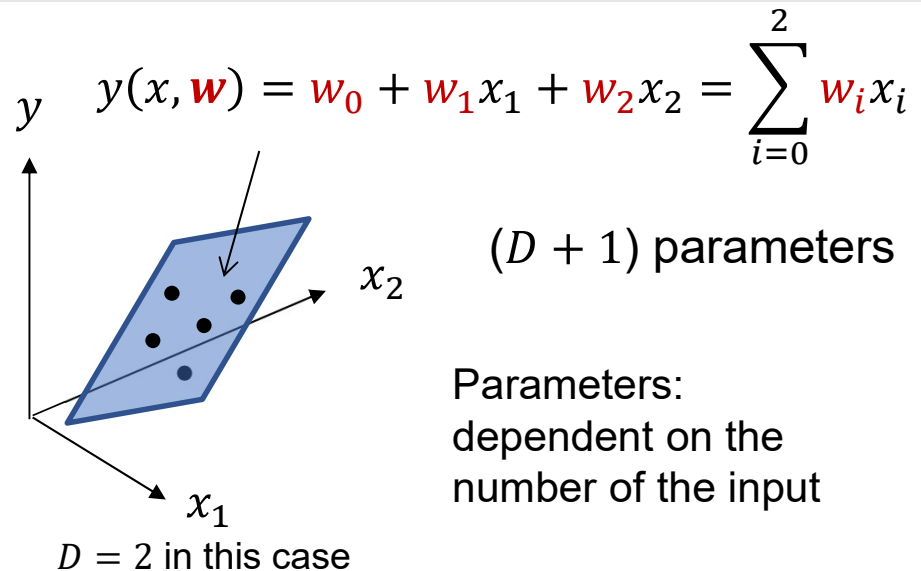
$$E(\mathbf{w}) = \sum_{i=1}^N \{y_i - y(x_i, \mathbf{w})\}^2$$

$$\begin{cases} \frac{\partial E(\mathbf{w})}{\partial w_0} = 0 \\ \frac{\partial E(\mathbf{w})}{\partial w_1} = 0 \end{cases}$$



$\hat{\mathbf{w}}$ is analytically solved.

Please confirm this by yourself.



The components of x

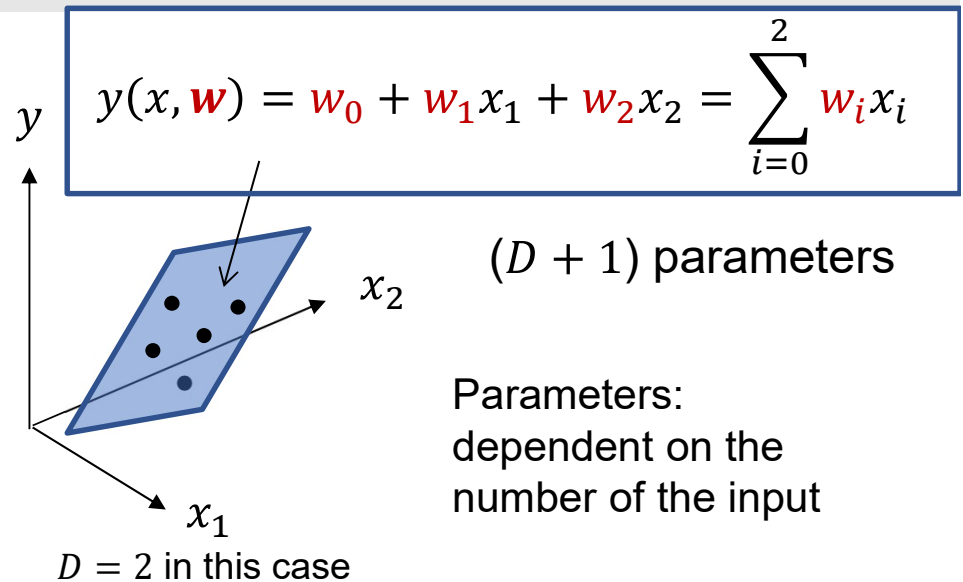
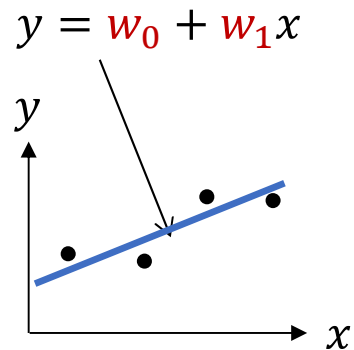
$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1D} \\ 1 & \vdots & \ddots & \vdots \\ 1 & x_{N1} & \cdots & x_{ND} \end{pmatrix} \quad \left. \begin{matrix} x_{i1}, \dots, x_{iD} \\ \text{sample} \\ x_1, \dots, x_N \end{matrix} \right\}$$

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

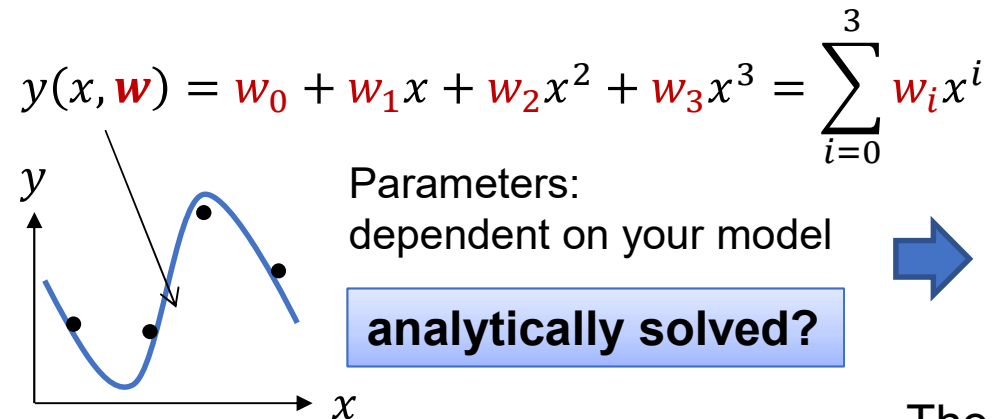
analytically solved

Regression Model

The simplest linear regression



Nonlinear model (**still linear regression**)



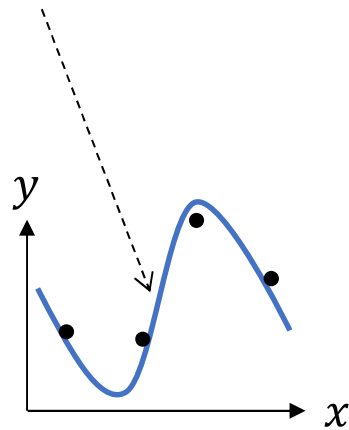
- ↓
- D-dimensional input
 - M-degrees

The parameters \mathbf{w} increase exponentially.

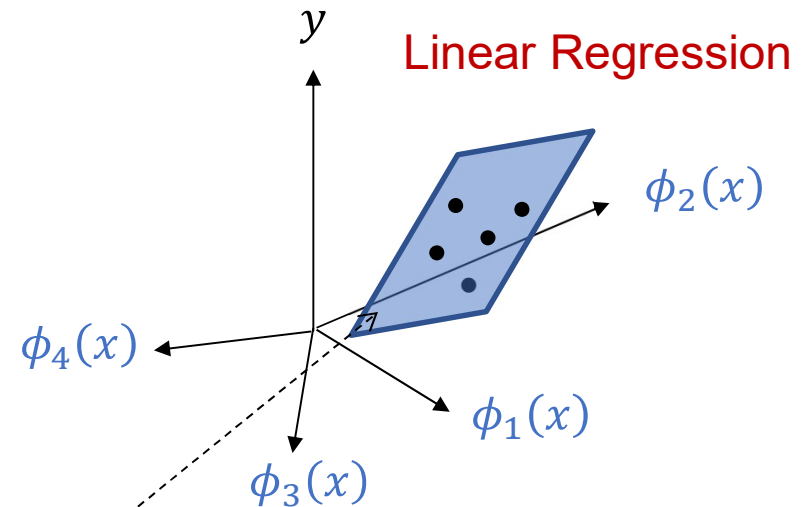
Linear Regression

Shifting to the simplest linear regression by mapping $\phi: x \rightarrow s$ $s = \phi(x)$

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 = \sum_{i=0}^3 w_i x^i$$



from 1D to 4D
in this example



In this example:

$$\begin{aligned}\phi(x) &= (\phi_0(x), \phi_1(x), \phi_2(x), \phi_3(x))^T \\ &= (x^0, x^1, x^2, x^3)^T\end{aligned}$$

$$y(\mathbf{s}, \mathbf{a}) = a_0s_0 + a_1s_1 + a_2s_2 + a_3s_3 = \sum_{i=0}^3 a_i s_i$$

Linear Regression

Examples of $\phi(x)$:

$$\phi_i(x) = \exp\left\{-\frac{(x - \mu_i)^2}{2\sigma^2}\right\}$$

Gaussian function

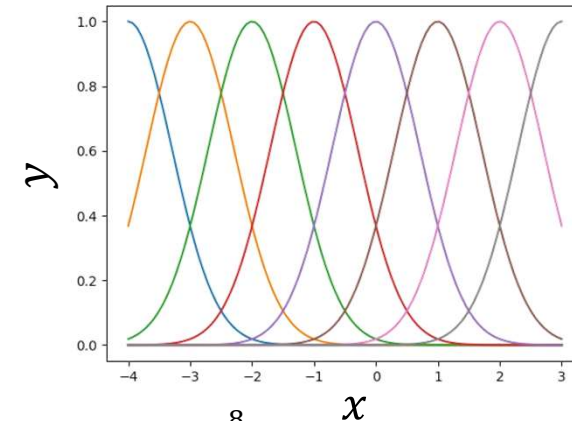
$\mu = (\mu_1, \mu_2, \dots, \mu_i, \dots, \mu_8)$ and σ are user-defined.

We can use any nonlinear functions for $\phi(x)$.

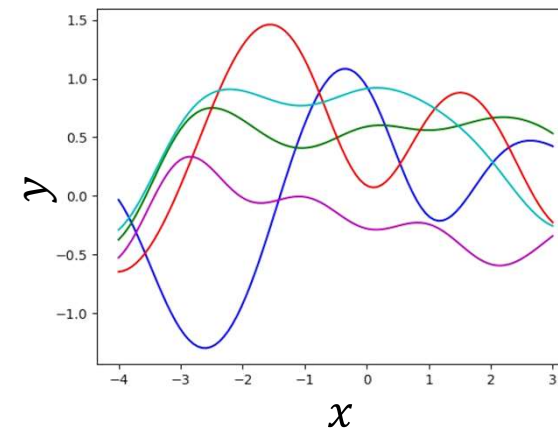
5 examples from randomly generated w

w is composed of 8 elements.

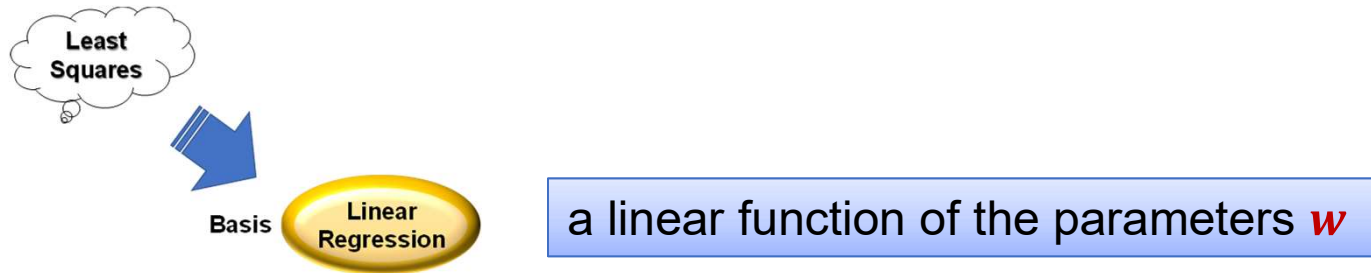
8 Gaussian functions $\phi(x)$



$$\sum_{i=1}^8 w_i \phi_i(x) = w^T \phi(x)$$



Linear Regression

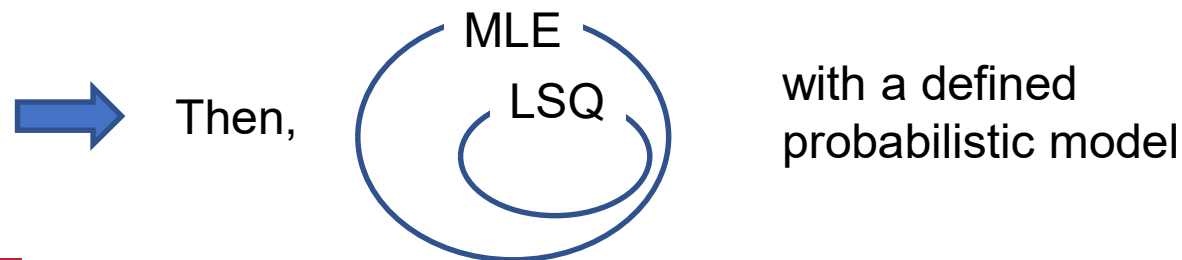


Linear regression (in general)

$$y(x, \mathbf{w}) = w_0 + w_1 \phi_1(x) + w_2 \phi_2(x) + \cdots + w_M \phi_M(x) = \sum_{i=0}^M w_i \phi_i(x) = [\mathbf{w}^T \boldsymbol{\phi}(x)]$$

M can be freely defined.

- $(M + 1)$ parameters $\mathbf{w} = (w_0, w_1, w_2, \dots, w_M)^T$
- no matter how many dimensionality the original input x has



Linear Regression

Likelihood function

$$L(\mathbf{w}, \sigma) \equiv -\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma^2) = \frac{1}{2\sigma^2} \sum_{i=1}^N \underbrace{\{y_i - y(x_i, \mathbf{w})\}^2}_{\text{the least square term}} + \frac{N}{2} \ln(2\pi\sigma^2)$$

negative log of the likelihood function

In Linear Regression: $y(x, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x)$

Maximum Likelihood Estimation (MLE)

$$\hat{\mathbf{w}}, \hat{\sigma} = \underset{\mathbf{w}, \sigma}{\operatorname{argmin}} L(\mathbf{w}, \sigma)$$

Need optimization algorithm?

Linear Regression

M : number of the parameter \mathbf{w}
 N : sample size

$$\frac{\partial L(\mathbf{w}, \sigma)}{\partial \mathbf{w}} = 0 \text{ or } \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = 0$$



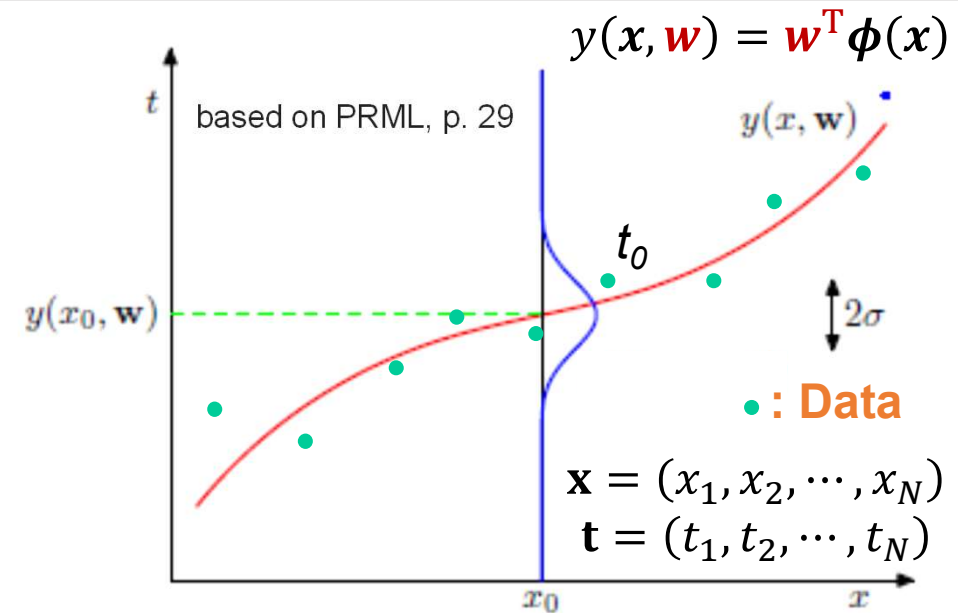
analytically solved

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

$M \times M$ square matrix

$$\Phi^{*-1} \equiv (\Phi^T \Phi)^{-1} \Phi^T$$

$$\Phi^* \hat{\mathbf{w}} = \mathbf{t} \quad \text{a linear system}$$



The components of $\phi(x)$

$$\phi_1(x), \dots, \phi_M(x)$$

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

sample
 x_1, \dots, x_N

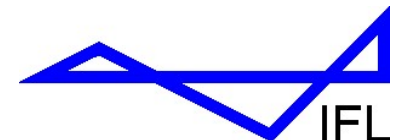
$N \times M$ matrix



Technische
 Universität
 Braunschweig

Dr. Daigo Maruyama | Scientific Machine Learning: Lecture 4 | Slide 32

Pseudo inverse matrix



Linear Regression

M : number of the parameter \mathbf{w}
 N : sample size

Important properties of Linear Regression

- The regression model is a linear function of the parameters $\mathbf{w} = (w_0, w_1, w_2, \dots, w_M)^T$
- The function $\boldsymbol{\phi}(\mathbf{x}) = (\phi_0(\mathbf{x}), \phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_M(\mathbf{x}))^T$ can be nonlinear of the input \mathbf{x} .



- The number of the parameters \mathbf{w} is M .
- Therefore, $\hat{\mathbf{w}}$ can be obtained analytically.



(Under the assumption that the probabilistic model is an isotropic Gaussian distribution)



Linear Regression

Regularization

Penalty on the parameter \mathbf{w} to avoid overfitting

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w})$$

$$\text{s.t. } \|\mathbf{w}\|^2 \leq \eta$$



$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} E_{\text{reg}}(\mathbf{w}) = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t}$$

$$\text{where, } E_{\text{reg}}(\mathbf{w}) = E(\mathbf{w}) + \lambda \|\mathbf{w}\|^2$$

regularization term

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



$$E(\mathbf{w}) = \sum_{i=1}^N \{y_i - y(x_i, \mathbf{w})\}^2$$

Review:

The error function $E(\mathbf{w})$

= the simplest expression of the negative log likelihood.

still analytically solved

Linear Regression

Other regularization techniques

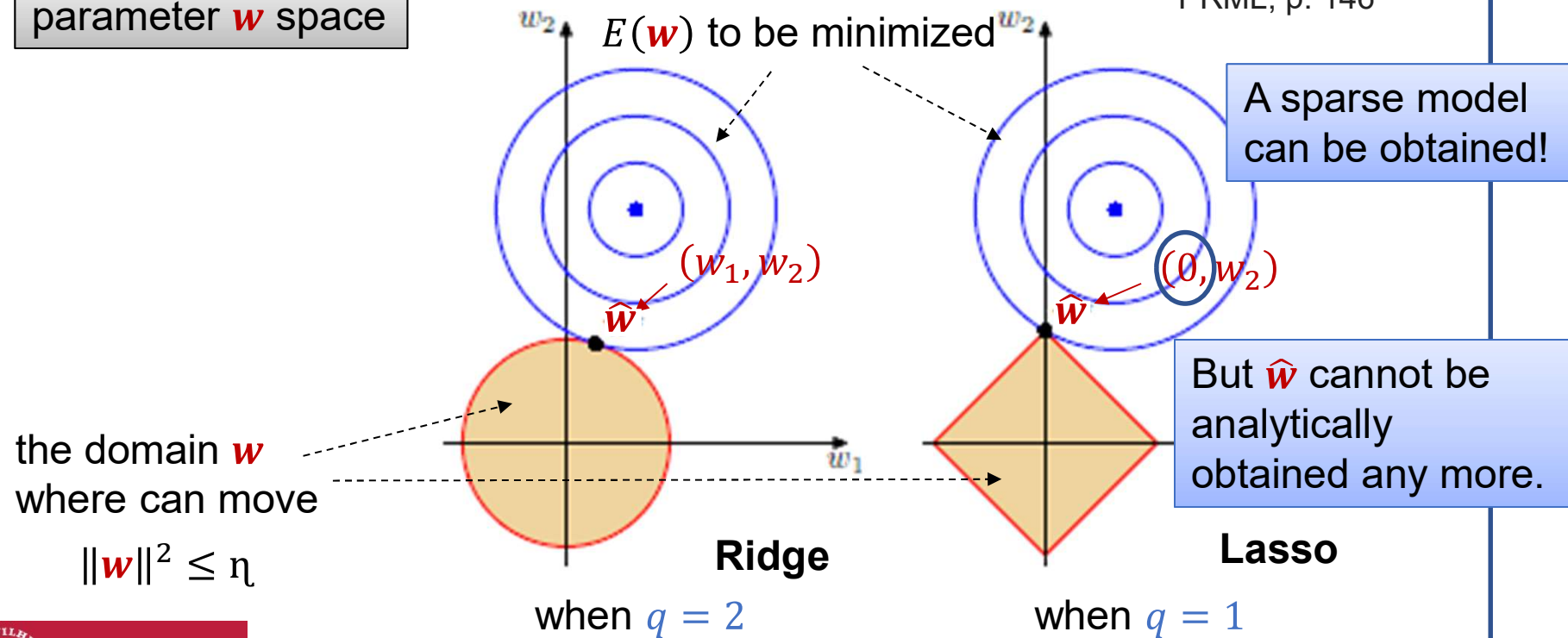
regularization term

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} E_{reg}(\mathbf{w})$$

where, $E_{reg}(\mathbf{w}) = E(\mathbf{w}) + \lambda \|\mathbf{w}\|^q$

parameter \mathbf{w} space

PRML, p. 146



Summary until Linear Regression (1/2)

A larger framework than the least square method was introduced based on perspectives of the probability theory
(Some examples will be introduced in other lectures).

The concept (a procedure):

1. **Define a probabilistic model (using probability distributions)**
2. **Maximize the likelihood function determined by dataset**
3. **(if prior distributions are set, maximize the posterior distribution)**

Compute the log of the likelihood/posterior to avoid numerical errors

error function the least square is one of the error functions.

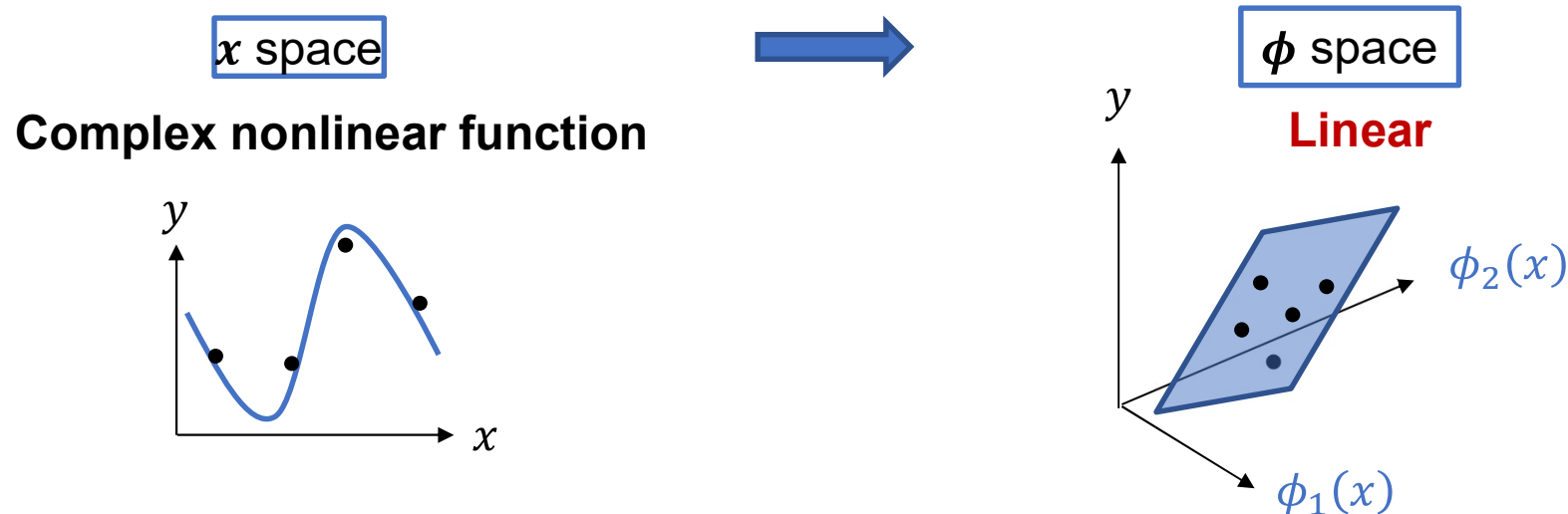
- The concept is further extended to be generalized until the lectures of “Bayes”.

Summary until Linear Regression (2/2)

A regression model can be defined in the process of the definition of the probabilistic model.

- Linear regression model is one of the regression models.
 - The model is a linear function of the parameters w

➡ The 2nd process: MLE becomes easy (analytically obtained).



This concept will be important.

Current Position to Next



➔ The theory is extended by introducing Bayesian perspective by from the next lecture using two slots of lectures

The required tools will be all equipped.