

Review: Bayesian Linear Regression

$$D = \{(x_i, y_i)\}_{i=1}^n, \quad X = [x_1, \dots, x_n] \in \mathbb{R}^{n \times d}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

$$P(y|x, D) = \int P(y|x, w) P(w|D) dw \propto \exp\left(-\frac{1}{2}(y - U_{y|x,D})^T \Sigma_{y|x,D}^{-1} (y - U_{y|x,D})\right)$$

$$U_{y|x,D} = \frac{X^T (XX^T + \frac{\sigma^2}{\alpha} I)^{-1} X y}{1 - X^T (XX^T + \frac{\sigma^2}{\alpha} I)^{-1} X}$$

$$\Sigma_{y|x,D} = \frac{\sigma^2}{1 - X^T (XX^T + \frac{\sigma^2}{\alpha} I)^{-1} X}$$

Woodbury matrix Identity

$$\Rightarrow \left(\frac{1}{\alpha} I - \frac{X (XX^T + \frac{\sigma^2}{\alpha} I)^{-1} X}{\sigma^2} \right)^{-1} = \frac{1}{\alpha} I + X^T (XX^T + \alpha I)^{-1} X$$

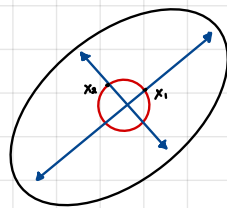
$$\Rightarrow (A + UCV)^T = A^{-1} - A^{-1} U (C^{-1} + V A^{-1} U)^{-1} V A^{-1}$$

$$\text{Var}[y|x, D] = \sigma^2 (1 + X^T \underbrace{(XX^T + \alpha I)^{-1}}_{\in \mathbb{R}^{d \times d}} X) \quad (\alpha = \frac{\sigma^2}{\alpha})$$

$$(XX^T + \alpha I)^{-1}$$

$$\Rightarrow \underbrace{(X I X^T + \alpha I)^{-1}}_{\substack{\tilde{A} \tilde{C} \tilde{V} \\ \tilde{A}}} \leftarrow \text{Using Woodbury matrix Identity}$$

$$= \frac{1}{\alpha} I - \frac{1}{\alpha^2} X (I + \frac{1}{\alpha} X^T X)^{-1} X^T$$



- x_1 방향이 x_2 방향보다 크다.

$$\Rightarrow \text{Var}(x_1) > \text{Var}(x_2)$$

$$\text{Var}[y|x, D] = \sigma^2 (1 + X^T (\frac{1}{\alpha} I - \frac{1}{\alpha^2} X (I + \frac{1}{\alpha} X^T X)^{-1} X^T) X)$$

$$= \sigma^2 + \frac{\sigma^2}{\alpha} X^T (I - X (\alpha I + X^T X)^{-1} X^T) X \quad \leftarrow \text{Small } X$$

$$= \sigma^2 + \frac{\sigma^2}{\alpha} X^T X - \frac{\sigma^2}{\alpha} X^T X (XX^T + \alpha I)^{-1} X^T X$$

Inner product

$\in \mathbb{R}^{n \times n}$, Data의 4 / Dimension에 따라 더 가까운 수는 사용하는 것이 가능하다.

$$XX^T = [x_1, \dots, x_n]^T [x_1, \dots, x_n] = \begin{bmatrix} x_1^T x_1 & x_1^T x_2 & \dots & x_1^T x_n \\ x_2^T x_1 & x_2^T x_2 & \dots & x_2^T x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n^T x_1 & x_n^T x_2 & \dots & x_n^T x_n \end{bmatrix} \leftarrow \text{Training dataset}$$

$$E[y|x, D] = \frac{X^T (XX^T + \frac{\sigma^2}{\alpha} I)^{-1} X y}{1 - X^T (XX^T + \frac{\sigma^2}{\alpha} I)^{-1} X}$$

→ 아래의 관계 때문에 변화 없이 생각하자.
→ 같은 방정식이다.

$$(1 + X^T (XX^T + \alpha I)^{-1} X) X^T (XX^T + \alpha I)^{-1} X y \leftarrow E[y|x, D] \text{ 분배 } (1 + X^T (XX^T + \alpha I)^{-1} X) \text{를 분배하자}$$

분배 미리 재배치 가능

$$\underbrace{\left(\frac{X X^T}{\alpha} + XX^T + \alpha I \right)^{-1}}_{\substack{I = C \\ \tilde{U} \tilde{V} \\ \tilde{A}}} = (XX^T + \alpha I)^{-1} - (XX^T + \alpha I)^{-1} X (1 + X^T (XX^T + \alpha I)^{-1} X)^{-1} X^T (XX^T + \alpha I)^{-1}$$

$$(1 + X^T (XX^T + \alpha I)^{-1} X) X^T (XX^T + \alpha I)^{-1} X y = (1 + \cancel{X^T (XX^T + \alpha I)^{-1} X}) \cdot \cancel{X^T (XX^T + \alpha I)^{-1} X} \cdot \cancel{(XX^T + \alpha I)^{-1} X} \cdot \cancel{(1 + X^T (XX^T + \alpha I)^{-1} X)^{-1} X^T (XX^T + \alpha I)^{-1} X y}$$

$$= X^T (XX^T + \alpha I)^{-1} X y = X^T X (XX^T + \alpha I)^{-1} y$$

$\in \mathbb{R}^{d \times d}$

$\in \mathbb{R}^{n \times n}$

이 항을 빼

$$\begin{aligned} (XX^T + \alpha I) X &= XX^T X + \alpha X = X (X^T X + \alpha I) \\ X (X^T X + \alpha I)^{-1} &= (XX^T + \alpha I)^{-1} X, \quad X = (XX^T + \alpha I)^{-1} X (X^T X + \alpha I) \end{aligned}$$

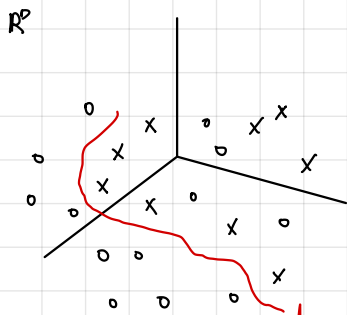
← 항을 빼는 것

$$X^T (XX^T + \alpha I)^{-1} X y = X^T U w \Rightarrow E[y|x, D] = X^T U w$$

$$\text{Var}[y|x, D] = \sigma^2 + \frac{\sigma^2}{\alpha} X^T (XX^T + \alpha I)^{-1} X = \sigma^2 + X^T \Sigma_w X$$

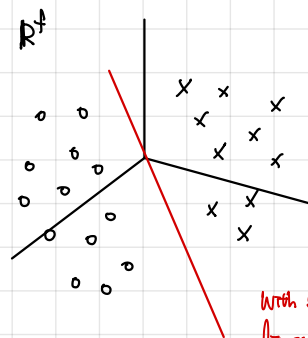
kernel

$$y = wx + b$$



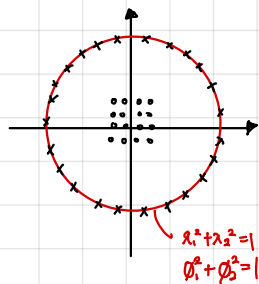
linearly non-separable
(Data space)

mapping data
 $\phi(x) \in \mathbb{R}^d$
 $f \gg D$



With good mapping,
linearly separable
(Feature space)

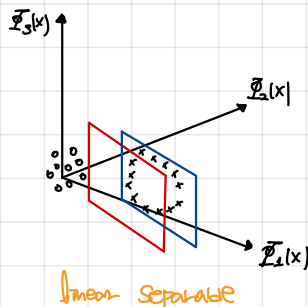
Ex)



$$\phi_1^2 + \phi_2^2 = 1$$

$$\phi(x) \in \mathbb{R}^3$$

$$\phi(x) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}$$



$$\phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

$$X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, Y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

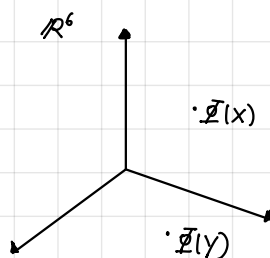
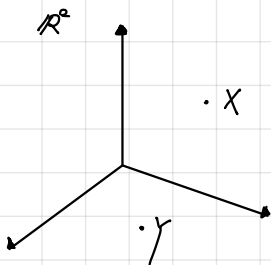
$$\phi(x)^T \phi(y) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}^T \begin{pmatrix} y_1^2 \\ y_2^2 \\ \sqrt{2}y_1y_2 \end{pmatrix} = x_1^2y_1^2 + x_2^2y_2^2 + 2x_1x_2y_1y_2 = (x_1y_1 + x_2y_2)^2 = \underbrace{(x^Ty)^2}_{\in \mathbb{R}} \equiv K(X, Y)$$

Ex) polynomial kernel

$$\text{kernel func: } K(x, y) = (x^Ty + 1)^d \in \mathbb{R}$$

$$\text{let } d=2 \quad (x, y \in \mathbb{R}^2)$$

$$\Rightarrow K(x, y) = (x_1y_1 + x_2y_2 + 1)^2 = x_1^2y_1^2 + x_2^2y_2^2 + 2x_1x_2y_1y_2 + 2x_1y_1 + 2x_2y_2 + 1 = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \frac{1}{\sqrt{2}}x_1x_2 \\ \frac{1}{\sqrt{2}}x_1 \\ \frac{1}{\sqrt{2}}x_2 \end{pmatrix} \begin{pmatrix} y_1^2 \\ y_2^2 \\ \frac{1}{\sqrt{2}}y_1y_2 \\ \frac{1}{\sqrt{2}}y_1 \\ \frac{1}{\sqrt{2}}y_2 \end{pmatrix} = \phi(x)^T \phi(y)$$



Review

Predictive distribution

$$P(y|x, D) = \mathcal{N}(y; \mu_{y|x,D}, \Sigma_{y|x,D})$$

$$\mu_{y|x,D} = x^T (XX^T + \frac{\sigma^2}{N} I)^{-1} X y = x^T \underbrace{\mu_w}_{\text{Posterior mean}}$$

$$\Sigma_{y|x,D} = \sigma^2 (1 + x^T (XX^T + \frac{\sigma^2}{N} I)^{-1} x) = \underbrace{\sigma^2}_{\text{Noise, Randomness (Model \& 학습 error)}} + x^T \underbrace{\Sigma_w}_{\text{Posterior Covariance}} x$$

$$y = w^T x + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

Epistemic Uncertainty

: Uncertainty in the model itself
This can be reduced to 0
as $N \rightarrow \infty$

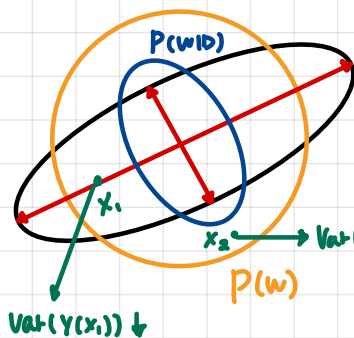
$$\Sigma_w = \frac{\sigma^2}{N} (XX^T + \frac{\sigma^2}{N} I)^{-1}$$

$$= \frac{\sigma^2}{N} \sum_{i=1}^N x_i x_i^T, \text{ if } N \rightarrow \infty \quad \frac{\sigma^2}{N} \sum_{i=1}^N x_i x_i^T \text{ becomes constant, Model } x \Rightarrow \lim_{N \rightarrow \infty} \frac{\sigma^2}{N} (\text{constant}) = 0$$

Aleatoric Uncertainty

: Uncertainty in data that cannot be reduced even if $N \rightarrow \infty$

Data Covariance

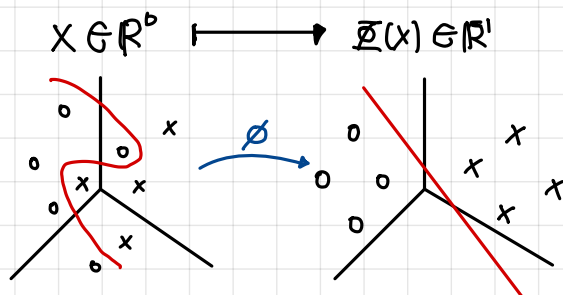


$$\text{Prior } P(w) = \mathcal{N}(0, \frac{\sigma^2}{N} I)$$

If Data Covariance $\uparrow \rightarrow$ Posterior Covariance \downarrow
If Data Covariance $\downarrow \rightarrow$ Posterior Covariance \uparrow

Posterior distribution 모델의 불확실성

—: Predictive Model

Expansion to model nonlinear w.r.t x 

$$y(x) = w^T x \mapsto y(x) = \alpha^T \phi(x)$$

kernel function

$$k(x, y) = \phi(x)^T \phi(y)$$

Example of kernel function

1. Polynomial kernel

$$K(x, y) = (x^T y + 1)^d, \quad d \in \mathbb{N}$$

$$= \phi(x)^T \phi(y)$$

2. Gaussian kernel

$$K(x, y) = \exp\left(-\frac{\|x - y\|^2}{2}\right) \in \mathbb{R}$$

$$= \exp\left(-\frac{x^T x + y^T y - 2x^T y}{2}\right)$$

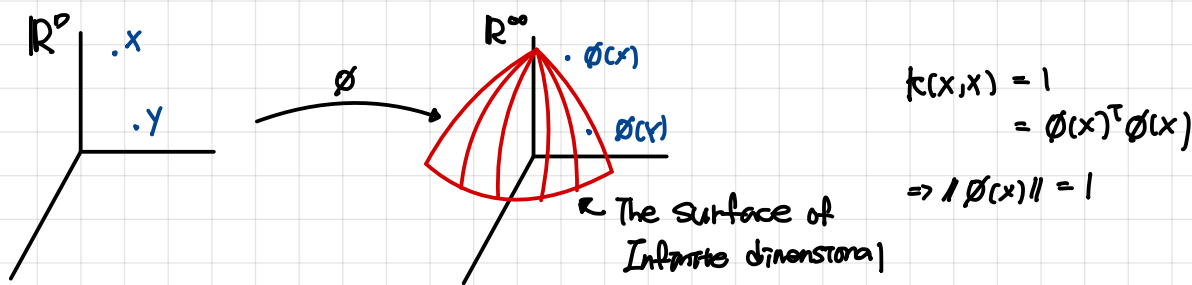
$$= \exp\left(-\frac{x^T x}{2}\right) \cdot \exp\left(-\frac{y^T y}{2}\right) \cdot \exp(x^T y)$$

$$\exp(t) = 1 + t + \frac{t^2}{2!} + \frac{t^3}{3!} + \dots$$

$$x^T y = x_1 y_1 + x_2 y_2 + \dots + x_n y_n$$

$$\exp(x^T y) = \frac{1 + x^T y + \frac{1}{2!} (x^T y)^2 + \frac{1}{3!} (x^T y)^3 + \dots}{\sqrt{1 + x^T x + \frac{1}{2!} (x^T x)^2 + \dots} \sqrt{1 + y^T y + \frac{1}{2!} (y^T y)^2 + \dots}} = \phi(x)^T \phi(y)$$

$$\phi(x) = \frac{1}{\sqrt{1 + x^T x + \frac{1}{2!} (x^T x)^2 + \dots}} \begin{pmatrix} 1 \\ x_1 \\ \frac{x_1^2}{\sqrt{2}} \\ x_2 \\ \frac{x_1 x_2}{\sqrt{2}} \\ \frac{x_1^3}{\sqrt{6}} \\ \vdots \end{pmatrix}$$



Back to Bayesian Linear Regression, Another version for $\mathcal{U}_{Y|X,D}$, $\Sigma_{Y|X,D}$

$$\mathcal{U}_{Y|X,D} = X^T X (X^T X + \frac{\sigma_0^2}{\sigma^2} I)^{-1} y = K^T (K + \frac{\sigma_0^2}{\sigma^2} I)^{-1} y$$

$$\Sigma_{Y|X,D} = \sigma^2 + \sigma_0^2 X^T X - \sigma_0^2 X^T X (X^T X + \frac{\sigma_0^2}{\sigma^2} I)^{-1} X^T X = \sigma^2 + \sigma_0^2 K^T (K + \frac{\sigma_0^2}{\sigma^2} I) K$$

kernelize

$$x \mapsto \phi(x)$$

$$X = [x_1 \dots x_n] \mapsto \Phi = [\phi(x_1) \dots \phi(x_n)]$$

$$x^T x \mapsto \phi(x)^T \phi(x) = k(x, x)$$

$$X^T X \mapsto [\phi(x)^T \phi(x_1) \quad \phi(x)^T \phi(x_2) \quad \dots \quad \phi(x)^T \phi(x_n)]$$

$$= [k(x, x_1) \quad k(x, x_2) \quad \dots \quad k(x, x_n)]$$

$$= K^T G R^{1 \times n} \quad \text{and} \quad K_i = k(x, x_i)$$

Small K

$$X^T X \mapsto \Phi^T \Phi$$

$$= \begin{bmatrix} \phi(x_1)^T \phi(x_1) & \dots & \phi(x_1)^T \phi(x_n) \\ \vdots & & \vdots \\ \phi(x_n)^T \phi(x_1) & \dots & \phi(x_n)^T \phi(x_n) \end{bmatrix} = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix} \equiv K G R^{n \times n}$$

Large K

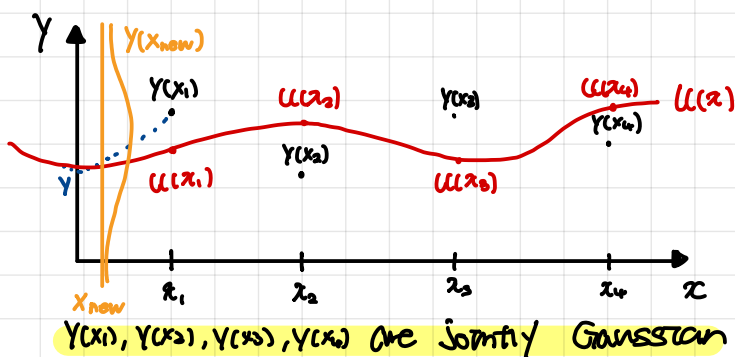
$K_{ii} = k(x_i, x_i)$

In fact, the above process is an example of the Gaussian processes
Gaussian Processes (= Infinite dimensional Gaussian)

Recall:

$$\mathcal{U}_{a|b} = \mathcal{U}_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mathcal{U}_b)$$

$$\Sigma_{a|b} = \Sigma_a - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$



$\mathcal{U}(x)$: mean function

$\Sigma(x_1, x_2)$: Covariance function
 $= \text{COV}(y(x_1), y(x_2))$

GP(\mathcal{U}, Σ): Gaussian process of
 mean function $\mathcal{U}(x)$ and
 Covariance matrix $\Sigma(x_1, x_2)$

Assume that $y(x_1) \dots y(x_n)$ are known for one realization.

We'd like to know the prediction for $y(x)$ for a new x