

# Gaussian Processes and Bayesian Linear Regression

---

Yung-Kyun Noh

---

---



# GP as an Infinite Dimensional Gaussians

$$\text{Ex)} \quad x_{t+1} = \alpha x_t + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma^2)$$

$$E[\varepsilon_t] = 0$$

$$E[\varepsilon_t^2] = \sigma^2$$

$$E[\varepsilon_t \varepsilon_{t'}] = 0 \quad t \neq t'$$

$$x_t = \alpha x_{t-1} + \varepsilon_{t-1}$$

$$= \alpha(\alpha x_{t-2} + \varepsilon_{t-2}) + \varepsilon_{t-1}$$

$$= \alpha^2 x_{t-2} + \alpha \varepsilon_{t-2} + \varepsilon_{t-1}$$

$$= \alpha^3 x_{t-3} + \alpha^2 \varepsilon_{t-3} + \alpha \varepsilon_{t-2} + \varepsilon_{t-1}$$

= ...

$$= \sum_{i=0}^{\infty} \alpha^i \varepsilon_{t-1-i}$$

$$E[x_t] = \sum_{i=0}^{\infty} \alpha^i E[\varepsilon_{t-1-i}] = 0$$

$$E[x_t x_{t+dt}] = E\left(\sum_{i=0}^{\infty} \alpha^i \varepsilon_{t-1-i}\right) \left(\sum_{j=0}^{\infty} \alpha^j \varepsilon_{t+dt-1-j}\right)$$

$$t-1-i = t+dt-j \Rightarrow j = i+dt$$

$$\begin{pmatrix} \dots & \alpha^2 \varepsilon_{t-3} & \alpha \varepsilon_{t-2} & \varepsilon_{t-1} \\ \dots & \alpha^{t+2} \varepsilon_{t-3} & \alpha^{t+1} \varepsilon_{t-2} & \alpha^t \varepsilon_{t-1} & \dots & \alpha \varepsilon_{t+dt-2} & \varepsilon_{t+dt-1} \end{pmatrix}$$

$$= \sum_{i=0}^{\infty} \alpha^i \alpha^{i+dt} E[\varepsilon_{t-1-i}^2] = \sum_{i=0}^{\infty} \alpha^{2i+dt} \sigma^2 = \frac{\alpha^{dt} \sigma^2}{1-\alpha^2}$$

1)

$$m(t) = 0, \quad k(t, t') = \frac{a^{|t-t'|} b^2}{1-a^2}$$

$$k(t_1, t_2) = \frac{a^{|t_1-t_2|} b^2}{1-a^2} = \begin{pmatrix} \vdots & \\ 0 & \\ \vdots & \\ 0 & \\ \vdots & \\ t_1 \leftarrow 0 & \\ t_1-1 \leftarrow b & \\ t_1-2 \leftarrow ab & \\ t_1-3 \leftarrow a^2b & \\ \vdots & \end{pmatrix}^T \begin{pmatrix} \vdots & \\ 0 & \\ \vdots & \\ 0 & \\ \vdots & \\ t_2 \leftarrow 0 & \\ t_2-1 \leftarrow b & \\ t_2-2 \leftarrow ab & \\ t_2-3 \leftarrow a^2b & \\ \vdots & \end{pmatrix}$$

$$= \phi(t_1)^T \phi(t_2)$$

$$\phi(t) = \begin{pmatrix} \vdots & \vdots \\ 0 & \\ 0 \leftarrow t & \\ b \leftarrow t-1 & \\ ab \leftarrow t-2 & \\ a^2b & \\ \vdots & \end{pmatrix}$$

Covariance

= Inner product of  
two vectors in  
 $\phi$ -space.

Note:

Inner product matrix is a p.d. matrix

$$K = \begin{pmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \cdots & \phi(x_1)^T \phi(x_N) \\ \vdots & & & \\ \phi(x_N)^T \phi(x_1) & \cdots & \phi(x_N)^T \phi(x_N) \end{pmatrix}$$

For a nonzero  $\vec{c} \in \mathbb{R}^N$ ,

$$\begin{aligned}\vec{c}^T K \vec{c} &= \sum_{i,j=1}^N c_i K_{ij} c_j = \sum_{i,j=1}^N c_i \phi(x_i)^T \phi(x_j) c_j \\ &= \left( \sum_{i=1}^N c_i \phi(x_i) \right)^T \left( \sum_{j=1}^N c_j \phi(x_j) \right) > 0\end{aligned}$$

$\therefore$  Inner product matrix  $K$  is p.d.

$$E[y|x, \mathcal{D}] = x^\top (XX^\top + \frac{\sigma^2}{\sigma_0^2} I)^{-1} X y \quad \dots \textcircled{1}$$

$$(XX^\top + \frac{\sigma^2}{\sigma_0^2} I) X = XX^\top X + \frac{\sigma^2}{\sigma_0^2} X$$

$$\underbrace{A''}_{A^{-1}} = X \underbrace{(X^\top X + \frac{\sigma^2}{\sigma_0^2} I)}_{= B}$$

$$AX = XB$$

$$A^{-1} A X B^{-1} = A^{-1} X B B^{-1}$$

$$XB^{-1} = A^{-1} X$$

$$\Rightarrow X(X^\top X + \frac{\sigma^2}{\sigma_0^2} I)^{-1} = (XX^\top + \frac{\sigma^2}{\sigma_0^2} I)^{-1} X$$

$$\therefore \textcircled{1} = x^\top X (X^\top X + \frac{\sigma^2}{\sigma_0^2} I)^{-1} y$$

$$= k^\top (K + \frac{\sigma^2}{\sigma_0^2} I)^{-1} y \quad \underbrace{\qquad \qquad \qquad}_{\leftarrow \text{Same as GP mean}}$$

If  $K_{ij} = x_i^\top x_j$ ,  $k_i = x_i^\top x_i$ ,

When we consider GP with covariance

$$k(t_1, t_2) = \frac{\alpha^{|t_1 - t_2|} \sigma^2}{1 - \alpha^2}, \quad \text{GP mean after}$$

observation is the predictive mean of the linear Bayesian model in the  $\phi(t)$  space, with  $y = w^\top \phi(x)$  and Gaussian prior on  $w$ .

Note: Inner product of two  $\phi(x_1), \phi(x_2)$  is the covariance between  $y(x_1)$  and  $y(x_2)$ ,

$$k(x_1, x_2) = E[y(x_1)y(x_2)] - E[y(x_1)]E[y(x_2)]$$

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

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

$$= \exp\left(-\frac{1}{2\sigma^2} \{x^2 + z^2 - 2x^\top z\}\right)$$

$$= \frac{\exp\left(\frac{x^\top z}{\sigma^2}\right)}{\exp\left(\frac{x^2}{\sigma^2}\right)^{\frac{1}{2}} \exp\left(\frac{z^2}{\sigma^2}\right)^{\frac{1}{2}}}$$

$$= \frac{1 + \frac{x^T z}{6^2} + \frac{1}{2!} \left( \frac{x^T z}{6^2} \right)^2 + \dots}{\sqrt{1 + \frac{x^2}{6^2} + \frac{1}{2!} \left( \frac{x^2}{6^2} \right)^2 + \dots} \sqrt{1 + \frac{z^2}{6^2} + \frac{1}{2!} \left( \frac{z^2}{6^2} \right)^2 + \dots}}$$

$$= \frac{\psi(x)^T \psi(z)}{\|\psi(x)\| \|\psi(z)\|}$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \\ x_0 \end{pmatrix} \quad \psi(x) = \begin{pmatrix} 1 \\ x_1/6 \\ x_2/6 \\ \vdots \\ x_n/6 \\ x_0^2/\sqrt{2}6^2 \\ x_1^2/\sqrt{2}6^2 \\ \vdots \\ x_n^2/\sqrt{2}6^2 \\ x_1 x_2/6^2 \\ x_1 x_3/6^2 \\ \vdots \\ x_{n-1} x_n/6^2 \\ x_1^3/\sqrt{3!} 6^3 \\ \vdots \\ x_n^3/\sqrt{3!} 6^3 \\ \vdots \end{pmatrix}$$

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

$$\phi(x) = \frac{\psi(x)}{\|\psi(x)\|}$$

$$\begin{aligned} (x^T z)^2 &= (x_1 z_1 + \dots + x_0 z_0)^2 \\ &= x_1^2 z_1^2 + \dots + x_0^2 z_0^2 \\ &\quad + 2x_1 x_2 z_1 z_2 + \\ &\quad \dots + 2x_{n-1} x_n z_{n-1} z_n \end{aligned}$$

$$\begin{aligned} (x^T z)^3 &= (x_1 z_1 + \dots + x_0 z_0)^3 \\ &= x_1^3 z_1^3 + \dots + x_0^3 z_0^3 \\ &\quad + 3x_1^2 x_2 z_1^2 z_2 + \\ &\quad \dots + 6x_1 x_2 x_3 z_1 z_2 z_3 \\ &\quad + \dots \\ &\quad + 6x_{n-2} x_{n-1} x_n z_{n-2} z_{n-1} z_n \end{aligned}$$

$\sigma^2 \uparrow$  The element values decay fast

$\Rightarrow \phi(x)$  maps data to effectively low-dimensional space.

$\Rightarrow$  Flexibility  $\downarrow$ .

$\sigma^2 \downarrow$  The element values decay slow.

$\Rightarrow \phi(x)$  maps data to effectively high-dimensional space

$\Rightarrow$  Flexibility  $\uparrow$

Easy to overfit,

Data away to the training data becomes orthogonal to the training data.

$$k(x_i, x_j) = \phi^\top(x_i) \phi(x_j) = 0.$$

Linear Bayesian predictive density:

$$E[y|x, \mathcal{D}] = k^T \left( K + \frac{\sigma^2}{\sigma_0^2} I \right)^{-1} y$$

$$\text{Var}[y|x, \mathcal{D}] = \sigma^2 + \sigma_0^2 k(x, x) - \sigma_0^2 k^T \left( K + \frac{\sigma^2}{\sigma_0^2} I \right)^{-1} k$$

GPs

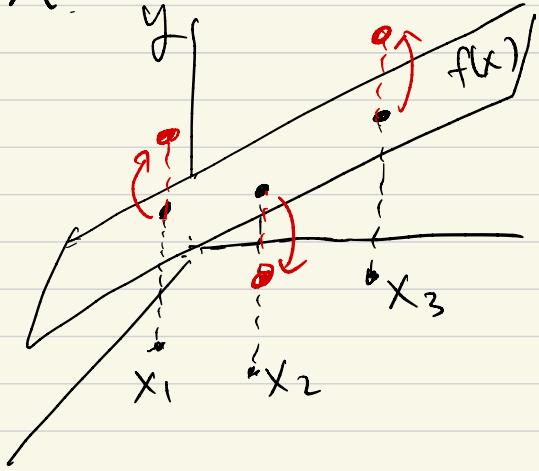
$$\mu(x|\mathcal{D}) = k^T K^{-1} y \quad \hookrightarrow \mu \text{ at } x$$

$$k(x, x|\mathcal{D}) = k(x, x) - k^T K^{-1} k$$

$$k_{GP}(x, z) = \sigma_0^2 k_{BL}(x, z) + \sigma^2 \delta_{x,z}$$

$$\delta_{x,z} = \begin{cases} 1 & \text{if } x = z \\ 0 & \text{otherwise} \end{cases}$$

# Linear Regression.

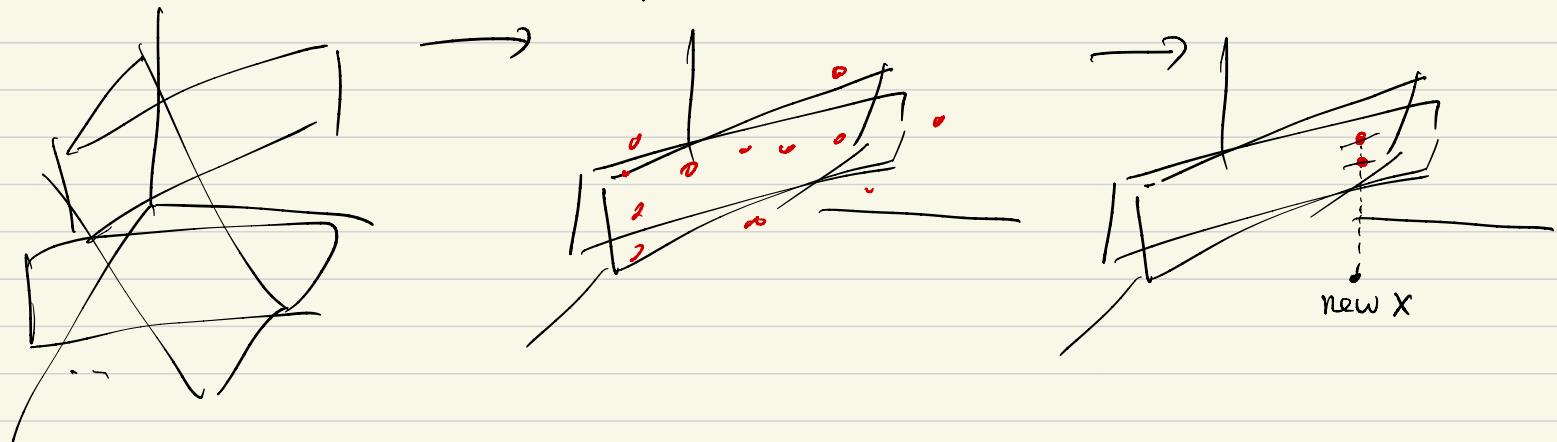


$$y = w^T x + \epsilon$$

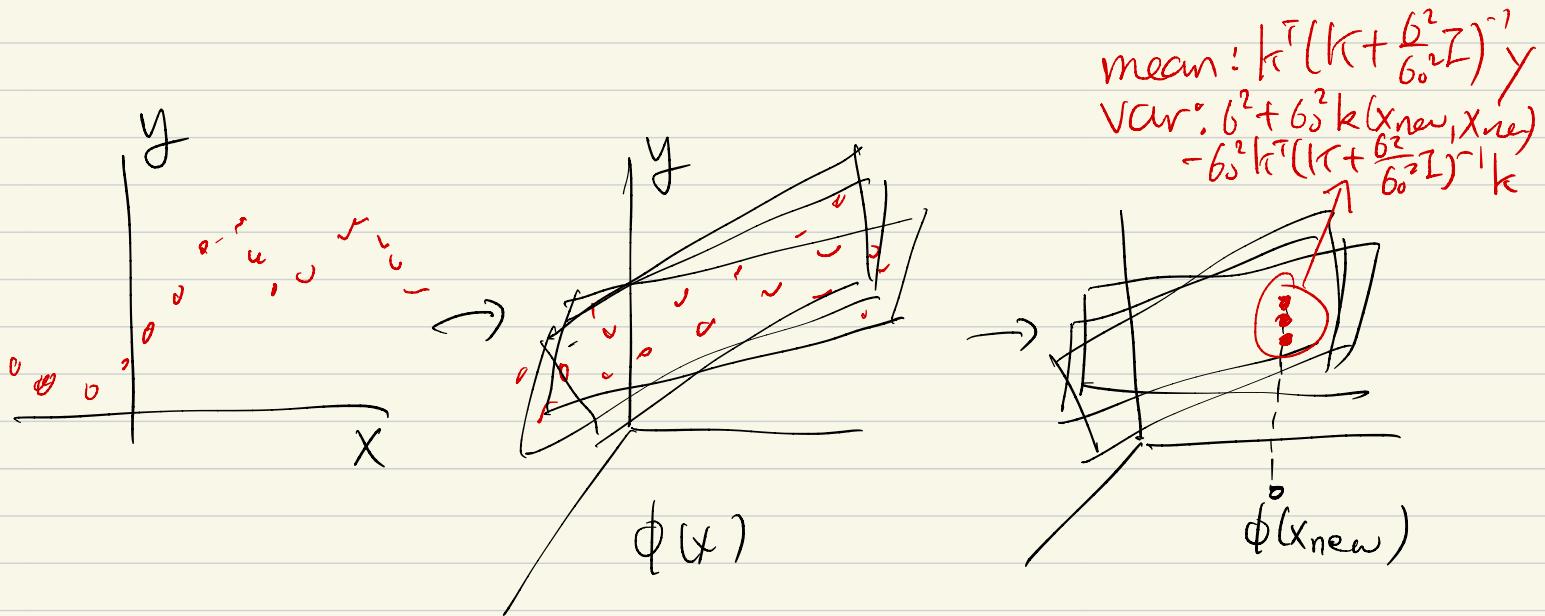
$$= f(x) + \epsilon$$

# Bayesian Linear Regression.(BLR)

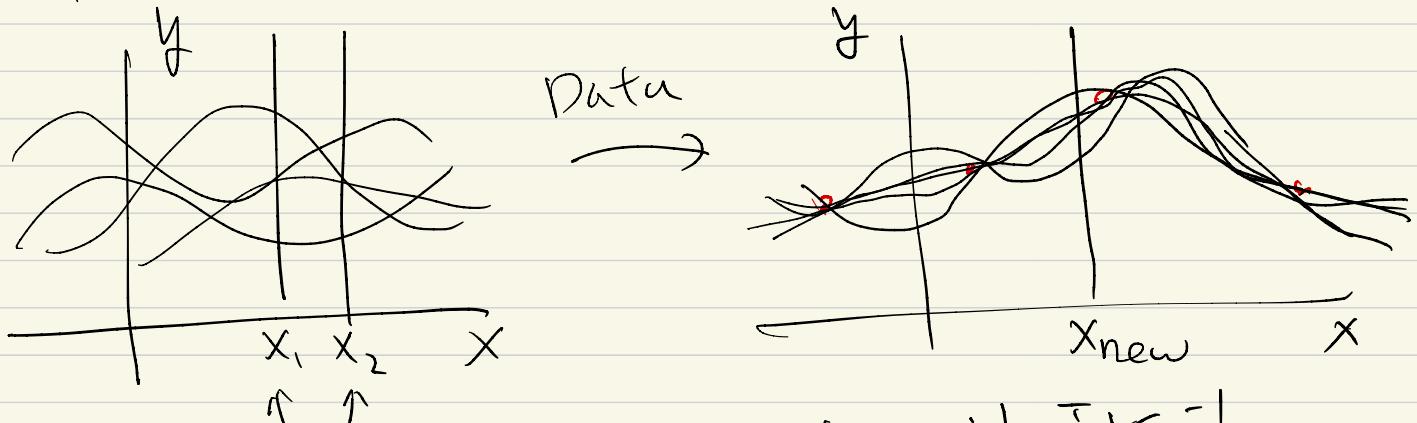
Prior for  $w$       Data      Posterior for  $w$       Prediction,



# BLR with kernels.



GPs,



$$\text{Cov}(y(x_i), y(x_j)) = k_{GP}(x_i, x_j)$$

$$k_{GP}(x_i, x_j)$$

$$= \sigma_0^2 k_{BLR}(x_i, x_j) + \sigma^2 S_{x_i, x_j}$$

$$\text{mean: } k_{GP}^T K_{GP}^{-1} y$$

$$\text{Var: } k_{GP}(x_{\text{new}}, x_{\text{new}}) - k_{GP}^T K_{GP}^{-1} k_{GP}$$

