

* Multivariate Gaussian dist of $y = f(x) \rightarrow$ "prior": mean=0, covariance= $k(x, x')$

Gaussian Process

1. Regression

⇒ Making inferences about the relationship between inputs and targets
(conditional distribution of the targets given the inputs.)

① Weights space view.

• standard linear model

$$f(x) = x^T w, \quad y = f(x) + \epsilon \quad \text{where } \epsilon \sim N(0, \sigma_n^2)$$

⇒ likelihood ($p(y|X, \omega)$)

$$\begin{aligned} \rightarrow p(y|X, \omega) &= \prod_{i=1}^n p(y_i|x_i, \omega) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(y_i - x_i^T \omega)^2}{2\sigma_n^2}\right) \\ &= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2} (y - X^T \omega)^T\right) = N(X^T \omega, \sigma_n^2 I) \end{aligned}$$

⇒ prior ($p(\omega)$)

$$\rightarrow p(\omega) = N(0, \Sigma_p)$$

⇒ posterior ($p(\omega|y, X)$)

$$\begin{aligned} \rightarrow p(\omega|y, X) &\propto p(y|X, \omega) \cdot p(\omega) = N(X^T \omega, \sigma_n^2 I) \cdot N(0, \Sigma_p) \\ &\propto \exp\left(-\frac{1}{2\sigma_n^2} (y - X^T \omega)^T (y - X^T \omega)\right) \cdot \exp\left(-\frac{1}{2} \omega^T \Sigma_p^{-1} \omega\right) \\ &\propto \exp\left(-\frac{1}{2} (\omega - \bar{\omega})^T \left(\frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1}\right) (\omega - \bar{\omega})\right) \end{aligned}$$

$$\text{where } \bar{\omega} = \frac{1}{\sigma_n^2} \left(\frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1} \right)^{-1} X^T y$$

→ $p(\omega|y, X) = N\left(\frac{1}{\sigma_n^2} A^{-1} X^T y, A^{-1}\right)$ where $A = \frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1}$ (DxD)

MAP estimate ≈ penalized MLE (ridge regression)

⇒ predictive posterior ($p(f_*|x_*, y, X)$)

$$\begin{aligned} \rightarrow p(f_*|x_*, y, X) &= \int p(f_*|x_*, \omega) \cdot p(\omega|y, X) d\omega \\ &= N\left(\frac{1}{\sigma_n^2} Z_*^T A^{-1} X^T y, Z_*^T A^{-1} Z_*\right) \end{aligned}$$

• kernel trick

⇒ project the inputs into some high dimensional space using a set of basis functions
($\phi: \mathbb{R}^D \rightarrow \mathbb{R}^N$ where $D < N$ and $f(x) = \phi(x)^T w$)

⇒ predictive posterior ($p(f_*|x_*, y, X)$)

$$\begin{aligned} \rightarrow p(f_*|x_*, y, X) &= N\left(\frac{1}{\sigma_n^2} \phi(x_*)^T A^{-1} \phi(X) y, \phi(x_*)^T A^{-1} \phi(X)\right) \\ \text{where } A &= \frac{1}{\sigma_n^2} \phi(X) \phi(X)^T + \Sigma_p^{-1} \quad (\text{NXN}) \rightarrow \text{computationally expensive} \\ \rightarrow p(f_*|x_*, y, X) &= N\left(\phi(x_*)^T \sum_p \phi(x) \left(K + \sigma_n^2 I\right)^{-1} y, \right. \\ &\quad \left. \phi(x_*)^T \sum_p \phi(x) - \phi(x_*)^T \sum_p \phi(x) \left(K + \sigma_n^2 I\right)^{-1} \phi(x) \sum_p \phi(x)\right) \\ &= N\left(K(x_*, X) (K(X, X) + \sigma_n^2 I)^{-1} y, \right. \\ &\quad \left. K(x_*, x_*) - K(x_*, X) (K(X, X) + \sigma_n^2 I)^{-1} K(X, x_*)\right) \end{aligned}$$

where $K + \sigma_n^2 I = \phi(X)^T \sum_p \phi(X) + \sigma_n^2 I \quad (n \times n) \rightarrow$ computationally relieved

$$k(x_*, x_n) = \phi(x_*)^T \sum_p \phi(x_n) = \left(\sum_p \phi(x_*)\right) \cdot \left(\sum_p \phi(x_n)\right) = \psi(x_*) \cdot \psi(x_n)$$

covariance function, kernel

② Function-space view

• Gaussian Process

• collection of random variables, any finite number of which have a joint Gaussian distribution
 $\left(f(x) \sim GP(m(x), k(x, x')) \right)$
where $m(x) = \mathbb{E}[f(x)]$, $k(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x'))^T]$

⇒ marginalization property (consistency)

$$\rightarrow (y_1, y_2) \sim N(\mu, \Sigma) \Rightarrow y_1 \sim N(\mu_1, \Sigma_{11}) \text{ and } y_2 \sim N(\mu_2, \Sigma_{22})$$

(examination of larger set of variables does not change the distribution of smaller set)

⇒ generating samples from $GP(m, k)$

1) Cholesky decomposition on covariance matrix : $K = LL^T$

(L is a lower triangle matrix)

2) Sample $u \sim N(0, I)$

3) Compute $x = m + L \cdot u$

$$(\mathbb{E}[x] = 0, \text{Cov}(x, x') = \mathbb{E}[L u (L u)^T] = L \cdot \mathbb{E}[u u^T] L^T = K)$$

• zero mean Gaussian Process

⇒ noise free observations

$$\rightarrow \begin{bmatrix} f \\ f_* \end{bmatrix} \sim GP \left[0, \begin{bmatrix} k(x, x) & k(x, x_*) \\ k(x_*, x) & k(x_*, x_*) \end{bmatrix} \right]$$

$$\rightarrow p(f_*|x_*, f, X) = N(k(x_*, X) k(X, X)^{-1} f, k(x_*, x_*) - k(x_*, X) k(X, X)^{-1} k(X, x_*))$$

⇒ noisy observations

$$\rightarrow \begin{bmatrix} y \\ f_* \end{bmatrix} \sim GP \left[0, \begin{bmatrix} k(x, x) + \sigma_n^2 I & k(x, x_*) \\ k(x_*, x) & k(x_*, x_*) \end{bmatrix} \right]$$

$$\rightarrow p(f_*|x_*, y, X) = N(k(x_*, X) (k(X, X) + \sigma_n^2 I)^{-1} y, k(x_*, x_*) - k(x_*, X) (k(X, X) + \sigma_n^2 I)^{-1} k(X, x_*))$$

* Interpretation on mean prediction

1) linear combination of observation y : → weight ≈ (equivalent kernel)
Matern-3/2 estimator

2) linear combination of kernel $k(x_n, x_i)$ → representer theorem

* Interpretation on covariance prediction

1) do not depend on observation y

⇒ marginal likelihood ($p(y|x)$)

$$\begin{aligned} \rightarrow p(y|X) &= \int p(y|f, X) p(f|X) df \\ &= \int N(f, \sigma_n^2 I) \cdot N(0, k(X, X)) df \\ &= N(0, k(X, X) + \sigma_n^2 I) \end{aligned}$$

• Squared Exponential kernel

$$\rightarrow k(x_p, x_q) = \sigma_p^2 \exp\left(-\frac{1}{2\ell^2} (x_p - x_q)^2\right) + \sigma_n^2 \delta_{pq}$$

→ signal variance (σ_p^2)

→ length scale (ℓ) → Median Heuristics : $\ell^* = \frac{\text{Median}}{2 \log(n+1)}$

→ noise variance (σ_n^2)

• Decision theory for regression

: minimize the expected loss w.r.t. the model on what the truth might be

$$\rightarrow \tilde{R}_L(y_{\text{guess}}|x) = \int L(y_{\text{true}}, y_{\text{guess}}) p(y_{\text{true}}|x, y, X) dy$$

$$\rightarrow y_{\text{optimal}}|x = \arg \min_{y_{\text{guess}}} \tilde{R}_L(y_{\text{guess}}|x)$$

$$(L(y_{\text{true}}, y_{\text{guess}}) = |y_{\text{true}} - y_{\text{guess}}| \rightarrow y_{\text{optimal}}|x = \text{median}(p(y_{\text{true}}|x, y, X)))$$

$$(L(y_{\text{true}}, y_{\text{guess}}) = (y_{\text{true}} - y_{\text{guess}})^2 \rightarrow y_{\text{optimal}}|x = \text{mean}(p(y_{\text{true}}|x, y, X)))$$

$$(L(y_{\text{true}}, y_{\text{guess}}), p(y_{\text{true}}|x, y, X); \text{symmetric} \rightarrow y_{\text{optimal}}|x = \text{mean}(p(y_{\text{true}}|x, y, X)))$$

• Non-zero mean Gaussian Process

⇒ with a fixed mean function ($m(\cdot)$)

$$\rightarrow g(x) \sim GP(m(x), k(x, x))$$

$$\rightarrow p(g_*|x_*, y, X) = N(m(x_*) + k(x_*, X) k(X, X)^{-1} (y - m(X)), \text{Cov}(f_*))$$

⇒ with a few fixed basis functions ($b(\cdot)$)

: global linear model with residuals modelled by Gaussian Process

$$\rightarrow g(x) = f(x) + h(x)^T \beta \quad \text{where } f(x) \sim GP(0, k(x, x')), \beta \sim N(b, B)$$

$$\rightarrow g(x) \sim GP(h(x)^T b, k(x, x) + h(x)^T B h(x))$$

$$\rightarrow p(g_*|x_*, y, X) = N\left(\mathbb{E}[f_*] + R^T \bar{\beta}, \text{Cov}(f_*) + R^T (B^T + h(x)(k(X, X) + \sigma_n^2 I)^{-1} h(X)^T)^T R\right)$$

$$\text{where } \bar{\beta} = (B^{-1} + h(x)(k(X, X) + \sigma_n^2 I)^{-1} h(X)^T)^{-1} (h(X)(k(X, X) + \sigma_n^2 I)^{-1} h(X)^T)^T R$$

$$R = h(x_*) - h(x) (k(x, x) + \sigma_n^2 I)^{-1} k(x, x_*)$$

⇒ marginal likelihood ($p(y|x)$)

$$\rightarrow p(y|x) = \int p(y|g, X) p(g|x) dg$$

$$= \int N(g, \sigma_n^2 I) \cdot N(h(x)^T b, k(x, X) + h(x)^T B h(x)) dg$$

$$= N(h(x)^T b, k(x, x) + h(x)^T B h(x) + \sigma_n^2 I)$$

Sparse Gaussian Process

1. A unifying view of sparse approximate Gaussian Process Regression
: understand the theoretical foundations of the various approximations
 \Rightarrow every algorithm approximates the joint prior $p(f, f)$ by conditional independence
- $$p(f_x, f) \approx q(f_x | f) = \int q(f_x | u) q(f | u) p(u) du$$
- ① GP
- $$\rightarrow \begin{bmatrix} f \\ f_x \end{bmatrix} \sim N \left[0, \begin{bmatrix} k_{xx} & k_{xu} \\ k_{ux} & k_x \end{bmatrix} \right]$$
- ② Subset of Regressor (SoR)
- $$\begin{cases} f_x = k_{uu} w_u, \quad p(w_u) = N(0, k_{uu}^{-1}) \\ u = k_{uu} w_u, \quad p(u) = N(0, k_{uu}) \end{cases}$$
- $$\Rightarrow q(f_x | u) = N(k_{uu} k_{uu}^{-1} u, 0), \quad q(f_x | u) = N(k_{uu} k_{uu}^{-1} u, 0)$$
- $$\Rightarrow q(f_x | f) = N \left(0, \begin{bmatrix} Q_{ff} & Q_{fx} \\ Q_{xf} & K_{xx} \end{bmatrix} \right)$$
- $$(\because \text{Cov}(f, f_x) = \mathbb{E}[(k_{uu} k_{uu}^{-1} u - 0)(k_{uu} k_{uu}^{-1} u - 0)^T] = \mathbb{E}[k_{uu} k_{uu}^{-1} u \cdot u^T k_{uu} k_{uu}^{-1}] = Q_{fx})$$
- ③ Deterministic Training Conditional (DTC) \Leftrightarrow PLV, PP
- $$\Rightarrow p(y|f) \approx q(y|f) = N(k_{uu} k_{uu}^{-1} u, \sigma^2 I)$$
- $$\Rightarrow q(f|u) = N(k_{uu} k_{uu}^{-1} u, 0), \quad q(f_x|u) = p(f_x|u)$$
- $$\Rightarrow q(f_x | f) = N \left(0, \begin{bmatrix} Q_{ff} & Q_{fx} \\ Q_{xf} & K_{xx} \end{bmatrix} \right) \rightarrow \text{violating consistency}$$
- $$(\because \text{Cov}(f, f_x) = \mathbb{E}[(f - 0)(f_x - 0)^T] = k_{xx})$$
- ④ Fully Independent Training Conditional (FITC) \Leftrightarrow SPGP
- $$\Rightarrow p(y|f) \approx q(y|f) = N(k_{uu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_f] + \sigma^2 I)$$
- $$\Rightarrow q(f|u) = N(k_{uu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_f]), \quad q(f_x|u) = p(f_x|u)$$
- $$\Rightarrow q(f_x | f) = N \left(0, \begin{bmatrix} Q_{ff} + \text{diag}[k_{ff} - Q_f] & Q_{fx} \\ Q_{xf} & K_{xx} \end{bmatrix} \right) \rightarrow \text{violating consistency}$$
- $$(\because \text{Cov}(f, f_x) = \mathbb{E}[(k_{uu} k_{uu}^{-1} u)(k_{uu} k_{uu}^{-1} u)^T] + \text{diag}[k_{ff} - Q_f])$$
- ⑤ Fully Independent Conditional (FIC)
- $$\Rightarrow p(y|f) \approx q(y|f) = N(k_{uu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_f] + \sigma^2 I)$$
- $$\Rightarrow \begin{cases} q(f|u) = N(k_{uu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_f]) \\ q(f_x|u) = N(k_{uu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_f]) \end{cases}$$
- $$\Rightarrow q(f_x | f) = N \left(0, \begin{bmatrix} Q_{ff} + \text{diag}[k_{ff} - Q_f] & Q_{fx} \\ Q_{xf} & Q_{xx} + \text{diag}[k_{ff} - Q_f] \end{bmatrix} \right)$$
- $$(\because \text{Cov}(f, f_x) = \mathbb{E}[(k_{uu} k_{uu}^{-1} u)(k_{uu} k_{uu}^{-1} u)^T] + \text{diag}[k_{ff} - Q_f])$$
- ⑥ Partially Independent Training Conditional (PITC)
- $$\Rightarrow p(y|f) \approx q(y|f) = N(k_{uu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_f] + \sigma^2 I)$$
- $$\Rightarrow q(f|u) = N(k_{uu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_f]), \quad q(f_x|u) = p(f_x|u)$$
- $$\Rightarrow q(f_x | f) = N \left(0, \begin{bmatrix} Q_{ff} + \text{blockdiag}[k_{ff} - Q_f] & Q_{fx} \\ Q_{xf} & K_{xx} \end{bmatrix} \right) \rightarrow \text{violating consistency}$$
- $$(\because \text{Cov}(f, f_x) = \mathbb{E}[(k_{uu} k_{uu}^{-1} u)(k_{uu} k_{uu}^{-1} u)^T] + \text{blockdiag}[k_{ff} - Q_f])$$
- ⑦ Partially Independent Conditional (PIC)
- $$\Rightarrow p(y|f) \approx q(y|f) = N(k_{uu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_f] + \sigma^2 I)$$
- $$\Rightarrow \begin{cases} q(f|u) = N(k_{uu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_f]) \\ q(f_x|u) = N(k_{uu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_f]) \end{cases}$$
- $$\Rightarrow q(f_x | f) = N \left(0, \begin{bmatrix} Q_{ff} + \text{blockdiag}[k_{ff} - Q_f] & Q_{fx} \\ Q_{xf} & Q_{xx} + \text{blockdiag}[k_{ff} - Q_f] \end{bmatrix} \right)$$
- $$(\because \text{Cov}(f, f_x) = \mathbb{E}[(k_{uu} k_{uu}^{-1} u)(k_{uu} k_{uu}^{-1} u)^T] + \text{blockdiag}[k_{ff} - Q_f])$$

2. Variational Learning of Inducing Variables

: inducing variables and hyperparameters are jointly selected by maximizing the lower bound of the marginal likelihood to minimize the KL divergence between a variational GP and the true posterior GP

\Rightarrow predictive distribution $(p(f_x | y))$

$$\begin{aligned} \rightarrow p(f_x | y) &= \int p(f_x | f) p(f | y) df \\ &= \int p(f_x | f, u) p(f | y, u) p(u) df du \quad \text{where } p(u) = N(0, k_u) \\ &= \int p(f_x | u) p(f | u) p(u) df du \\ &= \int p(f_x | u) p(u) du = \int p(f_x, u | y) du \end{aligned}$$

(u is assumed to be a sufficient statistic for f)

\Rightarrow variational distribution $(q(f_x))$

$$\begin{aligned} \rightarrow q(f_x) &= \int p(f_x | u) p(f | u) \phi(u) df du \quad \text{where } \phi(u) = N(\mu, A) \\ &= \int p(f_x | u) \phi(u) du = \int q(f_x, u) du \\ &= N(k_u^\top K_u^{-1} \mu, k_u^\top K_u^{-1} k_u + k_u^\top A K_u^{-1} k_u) \end{aligned}$$

\Rightarrow Variational learning of ϕ and kernel hyperparameters

$$\begin{aligned} \rightarrow \arg \min_{\phi, \mu, A} \text{KL}(q(f_x) || p(f_x | y)) &\Leftrightarrow \arg \min_{\phi, \mu, A} \text{KL}(q(f_x, u) || p(f_x, u | y)) \\ \rightarrow \log p(y | x) &= \int q(f_x, u) \log \frac{p(f_x, u | y)}{q(f_x, u)} df_x du \\ &= \int p(f_x | u) \phi(u) \log \frac{p(y | f_x) p(u)}{\phi(u)} df_x du \\ &= \int \phi(u) \left\{ \int p(f_x | u) \log p(y | f_x) df_x + \log \frac{p(u)}{\phi(u)} \right\} du \\ &\quad \downarrow \\ &\int p(f_x | u) \left[-\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \text{Tr}[y_y^\top T - 2y_x^\top f_x + f_x^\top f_x] \right] df_x \\ &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \text{Tr}[y_y^\top T - 2y_x^\top f_x + \alpha Q_f + k_f - Q_f] \\ &= \log N(\alpha, \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}[k_f - Q_f] \\ &\quad \text{where } \alpha = k_u^\top K_u^{-1} u, \quad Q_f = K_u^\top K_u^{-1} k_f \\ &= \int \phi(u) \log \frac{N(\alpha, \sigma^2 I) p(u)}{\phi(u)} du - \frac{1}{2\sigma^2} \text{Tr}[k_f - Q_f] \\ &\geq \log \int N(\alpha, \sigma^2 I) p(u) du - \frac{1}{2\sigma^2} \text{Tr}[k_f - Q_f] \quad (\because \text{Jensen's inequality}) \\ &= \log N(0, K_u^\top K_u^{-1} k_f + \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}[k_f - Q_f] \\ &\quad \text{Mope} = \sigma^2 k_u (k_u + \sigma^2 K_u^\top K_u^{-1})^\top K_u^\top y \\ &\quad \text{Aope} = k_u (k_u + \sigma^2 K_u^\top K_u^{-1})^\top k_u \end{aligned}$$

3. Sparse Orthogonal Variational Inference for Gaussian Process

: Introduce another set of inducing variables for the orthogonal complement, which is beneficial for the computational cost

⇒ Reinterpreting SVGP

$$\rightarrow p(f|u) = N(k_{uf}^T k_u^{-1} u, k_f - k_{uf}^T k_u^{-1} k_{uf})$$

$$\rightarrow f = k_{uf}^T k_u^{-1} u + f_1 \text{ where } f_1 \sim N(0, k_f - k_{uf}^T k_u^{-1} k_{uf})$$

⇒ Orthogonal decomposition

$$\rightarrow V = \left\{ \sum_{j=1}^m \alpha_j k(z_j, \cdot) \right\} : \text{linear span of kernel basis function}$$

$$\rightarrow f = f_H + f_L \text{ where } f_H \in V, f_L \perp V$$

$$\begin{cases} f_H = k_{uf}^T k_u^{-1} u \sim N(0, k_f - k_{uf}^T k_u^{-1} k_{uf}) \\ f_L = f - f_H \sim N(0, k_f - k_{uf}^T k_u^{-1} k_{uf}) \end{cases}$$

⇒ joint distribution ($p(y, u, f_L, v)$)

$$\rightarrow p(y, u, f_L, v) = p(y|k_{uf}^T k_u^{-1} u + f_L) p(u) p(f_L|v) p(v) \text{ where } p(v) = N(0, C_v)$$

⇒ variational distribution ($q(u, f_L, v)$)

$$\rightarrow q(u, f_L, v) = q(u) p(f_L|v) q(v) \text{ where } q(v) = N(m_v, S_v)$$

$$\rightarrow q(f_L) = \int p(f_L|v) q(v) dv = N(m_{f_L}, S_{f_L})$$

$$\begin{cases} m_{f_L} = C_{vf}^T C_v^{-1} m_v \\ S_{f_L} = C_{ff} - C_{vf}^T C_v^{-1} (C_v - S_v) C_v^{-1} C_{vf} \end{cases}$$

⇒ variational learning of z, σ, m_v, S_v and kernel hyperparameters

$$\rightarrow \log p(y|x)$$

$$\geq \mathbb{E}_{q(u)q(f_L)} [\log p(y|k_{uf}^T k_u^{-1} u + f_L)] - KL[q(u)\|p(u)] - KL[q(v)\|p(v)]$$

↓

$$= \mathbb{E}_{q(u)q(f_L)} [\log N(k_{uf}^T k_u^{-1} u + f_L, \sigma^2 I)]$$

$$= \mathbb{E}_{q(u)q(f_L)} \left[-\frac{N}{2} \log 2\pi - \frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y - k_{uf}^T k_u^{-1} u - f_L)^T (y - k_{uf}^T k_u^{-1} u - f_L) \right]$$

$$= \mathbb{E}_{q(u)} [\log N(k_{uf}^T k_u^{-1} u + m_{f_L}, \sigma^2 I)] - \frac{1}{2\sigma^2} \text{Tr}(S_{f_L})$$

$$\geq \int \log N(k_{uf}^T k_u^{-1} u + m_{f_L}, \sigma^2 I) p(u) du - \frac{1}{2\sigma^2} \text{Tr}(S_{f_L}) - KL[q(v)\|p(v)]$$

$$= \log N(m_{f_L}, k_{uf}^T k_u^{-1} k_{uf} + \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}(S_{f_L}) - KL[N(m_v, S_v)\|N(0, C_v)]$$

$$\begin{cases} m_v^{\text{opt}} = C_v [C_v + C_{vf} (k_{uf}^T k_u^{-1} k_{uf} + \sigma^2 I)^{-1} C_{vf}^T]^{-1} C_{vf} (k_{uf}^T k_u^{-1} k_{uf} + \sigma^2 I)^{-1} y \\ S_v^{\text{opt}} = C_v [C_v + \sigma^{-2} C_{vf} C_{vf}^T]^{-1} C_v \end{cases}$$

★ If $m_v = 0, S_v = C_v$, then it reduces to SVGP

★ If $S_v = C_v$, then it reduces to follow

$$\log N(m_{f_L}, k_{uf}^T k_u^{-1} k_{uf} + \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}(C_f) - \frac{1}{2} m_v^T C_v^{-1} m_v$$

Gaussian Process Extrapolation

1. Gaussian Process Kernels for pattern discovery and extrapolation
 - : derive kernel by modeling a spectral density with a gaussian mixture

