

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

## Gaussian Process

### 1. Regression

$\Rightarrow$  making inferences about the relationship between inputs and targets  
(conditional distribution of the targets given the inputs.)

#### ① Weight space view:

- standard linear model

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

$\Rightarrow$  likelihood ( $p(y | x, \omega)$ )

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

$\Rightarrow$  prior ( $p(\omega)$ )

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

$\Rightarrow$  posterior ( $p(\omega | y, X)$ )

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

$\downarrow$  MAP estimate  $\approx$  penalized MLE (ridge regression)

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

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

- kernel trick

$\Rightarrow$  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  $f(x) = \phi(x)^T w$ )

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

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

where  $\underline{k} + \sigma_n^2 I = \phi(x)^T \Sigma_p \phi(x) + \sigma_n^2 I \quad (n \times n) \rightarrow$  computationally relieved.

$$\underline{k(x_*, x)} = \phi(x_*)^T \Sigma_p \phi(x) = \left(\sum_i^k \phi(x_i)\right) \cdot \left(\sum_j^k \phi(x_j)\right) = \psi(x_*) \cdot \psi(x)$$

$\downarrow$  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  $m(x) = \mathbb{E}[f(x)]$ ,  $k(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x'))^T]$

$\Rightarrow$  marginalization property (consistency)

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

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

$\Rightarrow$  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

$\Rightarrow$  noise free observations

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

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

$\Rightarrow$  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

- linear combination of observation  $y$ ;  $\rightarrow$  weight  $\approx$  (equivalent kernel Nadaraya-Watson estimator)
- linear combination of kernel  $k(x_*, x_i)$   $\rightarrow$  representer theorem

\* Interpretation on covariance prediction

- do not depend on observation  $y$ :

$\Rightarrow$  marginal likelihood ( $p(y | X)$ )

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

### • Squared Exponential kernel

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

$\rightarrow$  signal variance ( $\sigma_f^2$ )

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

$\rightarrow$  noise variance ( $\sigma_n^2$ )  $\downarrow$  distance among all pairs

### • Decision theory for regression

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

$$\rightarrow \tilde{R}_L(y_{\text{guess}} | X_*) = \int L(y_*, y_{\text{guess}}) \cdot p(y_* | X_*, y_{\text{guess}}) dy_*$$

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

$$(L(y_*, y_{\text{guess}}) = |y_{\text{guess}} - y_*| \implies y_{\text{optimal}} | X_* = \text{median}(p(y_* | X_*, y_{\text{guess}})))$$

$$(L(y_*, y_{\text{guess}}) = (y_{\text{guess}} - y_*)^2 \implies y_{\text{optimal}} | X_* = \text{mean}(p(y_* | X_*, y_{\text{guess}})))$$

$$(L(y_*, y_{\text{guess}}), p(y_* | X_*, y_{\text{guess}}): \text{symmetric} \implies y_{\text{optimal}} | X_* = \text{mean}(p(y_* | X_*, y_{\text{guess}})))$$

### • Non-zero mean Gaussian Process

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

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

$$\rightarrow p(g_* | X_*, y, X) = N(\underline{m(X_*)} + k(X_*, X) (k(X, X) + \sigma_n^2 I)^{-1} (y - \underline{m(X)}), \text{Cov}(f_*))$$

$\Rightarrow$  with a few fixed basis functions ( $h(\cdot)$ )

: global linear model with residuals modelled by Gaussian Process

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

$$\rightarrow g(x) \sim GP(h(x)^T b, k(x, x') + h(x)^T B h(x'))$$

$$\rightarrow p(g_* | X_*, y, X) = N(\underline{m(X_*)} + R^T \bar{\beta}, \text{Cov}(f_*) + R^T (B^{-1} + h(X)(k(X, X) + \sigma_n^2 I)^{-1} h(X)^T)^{-1} R)$$

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

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

$\Rightarrow$  marginal likelihood ( $p(y | X)$ )

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

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

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

## Sparse Gaussian Process

### 1. A unifying view of sparse approximate Gaussian Process Regression

: understand the theoretical foundations of the various approximations  
 ⇒ every algorithm approximates the joint prior  $p(f, f)$  by conditional independence

$$(p(f, f) \approx q(f, f) = \int q(f|u) q(f|u) p(u) du)$$

① GP

$$\rightarrow \begin{bmatrix} f \\ f_* \end{bmatrix} \sim N \left[ 0, \begin{bmatrix} k_{ff} & k_{f*} \\ k_{*f} & k_{**} \end{bmatrix} \right]$$

② Subset of Regressor (SoR)

$$\begin{cases} f_* = k_{*u} w_u, \quad p(w_u) = N(0, k_{uu}^{-1}) \\ u = k_{uu}^{-1} w_u, \quad p(u) = N(0, k_{uu}) \end{cases}$$

$$\Rightarrow q(f|u) = N(k_{fu} k_{uu}^{-1} u, 0), \quad q(f_*|u) = N(k_{*u} k_{uu}^{-1} u, 0)$$

$$\Rightarrow q(f, f_*) = N \left( 0, \begin{bmatrix} Q_{ff} & Q_{f*} \\ Q_{*f} & Q_{**} \end{bmatrix} \right)$$

$$(\because \text{Cov}(f, f_*) = \mathbb{E}[(k_{fu} k_{uu}^{-1} u - 0)(k_{*u} k_{uu}^{-1} u - 0)^T] = \mathbb{E}[k_{fu} k_{uu}^{-1} u \cdot u^T k_{*u} k_{uu}^{-1}] = Q_{f*}]$$

③ Deterministic Training Conditional (DTC)  $\Leftrightarrow$  PLV, PP

$$\Rightarrow p(y|f) \approx q(y|f) = N(k_{fu} k_{uu}^{-1} u, \sigma^2 I)$$

$$\Rightarrow q(f|u) = N(k_{fu} k_{uu}^{-1} u, 0), \quad q(f_*|u) = p(f_*|u)$$

$$\Rightarrow q(f, f_*) = N \left( 0, \begin{bmatrix} Q_{ff} & Q_{f*} \\ Q_{*f} & K_{**} \end{bmatrix} \right) \rightarrow \text{violating consistency}$$

$$(\because \text{Cov}(f, f_*) = \mathbb{E}[(f_* - 0)(f_* - 0)^T] = K_{**})$$

④ Fully Independent Training Conditional (FITC)  $\Leftrightarrow$  SPGP

$$\Rightarrow p(y|f) \approx q(y|f) = N(k_{fu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_{ff}] + \sigma^2 I)$$

$$\Rightarrow q(f|u) = N(k_{fu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_{ff}]), \quad q(f_*|u) = p(f_*|u)$$

$$\Rightarrow q(f, f_*) = N \left( 0, \begin{bmatrix} Q_{ff} + \text{diag}[k_{ff} - Q_{ff}] & Q_{f*} \\ Q_{*f} & K_{**} \end{bmatrix} \right) \rightarrow \text{violating consistency}$$

$$(\because \text{Cov}(f, f_*) = \mathbb{E}[(k_{fu} k_{uu}^{-1} u)(k_{*u} k_{uu}^{-1} u)^T] + \text{diag}[k_{ff} - Q_{ff}])$$

⑤ Fully Independent Conditional (FIC)

$$\Rightarrow p(y|f) \approx q(y|f) = N(k_{fu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_{ff}] + \sigma^2 I)$$

$$\Rightarrow \begin{cases} q(f|u) = N(k_{fu} k_{uu}^{-1} u, \text{diag}[k_{ff} - Q_{ff}]) \\ q(f_*|u) = N(k_{*u} k_{uu}^{-1} u, \text{diag}[k_{**} - Q_{**}]) \end{cases}$$

$$\Rightarrow q(f, f_*) = N \left( 0, \begin{bmatrix} Q_{ff} + \text{diag}[k_{ff} - Q_{ff}] & Q_{f*} \\ Q_{*f} & Q_{**} + \text{diag}[k_{**} - Q_{**}] \end{bmatrix} \right)$$

$$(\because \text{Cov}(f, f_*) = \mathbb{E}[(k_{fu} k_{uu}^{-1} u)(k_{*u} k_{uu}^{-1} u)^T] + \text{diag}[k_{**} - Q_{**}])$$

⑥ Partially Independent Training Conditional (PITC)

$$\Rightarrow p(y|f) \approx q(y|f) = N(k_{fu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_{ff}] + \sigma^2 I)$$

$$\Rightarrow q(f|u) = N(k_{fu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_{ff}]), \quad q(f_*|u) = p(f_*|u)$$

$$\Rightarrow q(f, f_*) = N \left( 0, \begin{bmatrix} Q_{ff} + \text{blockdiag}[k_{ff} - Q_{ff}] & Q_{f*} \\ Q_{*f} & K_{**} \end{bmatrix} \right) \rightarrow \text{violating consistency}$$

$$(\because \text{Cov}(f, f_*) = \mathbb{E}[(k_{fu} k_{uu}^{-1} u)(k_{*u} k_{uu}^{-1} u)^T] + \text{blockdiag}[k_{**} - Q_{**}])$$

⑦ Partially Independent Conditional (PIC)

$$\Rightarrow p(y|f) \approx q(y|f) = N(k_{fu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_{ff}] + \sigma^2 I)$$

$$\Rightarrow \begin{cases} q(f|u) = N(k_{fu} k_{uu}^{-1} u, \text{blockdiag}[k_{ff} - Q_{ff}]) \\ q(f_*|u) = N(k_{*u} k_{uu}^{-1} u, \text{blockdiag}[k_{**} - Q_{**}]) \end{cases}$$

$$\Rightarrow q(f, f_*) = N \left( 0, \begin{bmatrix} Q_{ff} + \text{blockdiag}[k_{ff} - Q_{ff}] & Q_{f*} \\ Q_{*f} & Q_{**} + \text{blockdiag}[k_{**} - Q_{**}] \end{bmatrix} \right)$$

$$(\because \text{Cov}(f, f_*) = \mathbb{E}[(k_{fu} k_{uu}^{-1} u)(k_{*u} k_{uu}^{-1} u)^T] + \text{blockdiag}[k_{**} - Q_{**}])$$

### 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

⇒ predictive distribution ( $p(f_*|y)$ )

$$\rightarrow p(f_*|y) = \int p(f_*|f) p(f|y) df$$

$$= \int p(f_*|f, u) p(f|y, u) p(u) df du \text{ where } p(u) = N(0, k_{uu})$$

$$= \int p(f_*|u) p(f|u) p(u) df du$$

$$= \int p(f_*|u) p(u) du = \int p(f_*, u|y) du$$

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

⇒ variational distribution ( $q(f_*)$ )

$$\rightarrow q(f_*) = \int p(f_*|u) p(f|u) \phi(u) df du \text{ where } \phi(u) = N(\mu, A)$$

$$= \int p(f_*|u) \phi(u) du = \int q(f_*, u) du$$

$$= N(K_{*u}^T K_u^{-1} \mu, K_{**} - K_{*u}^T K_u^{-1} K_{*u} + K_u^T K_u^{-1} A K_u^{-1} K_{*u})$$

⇒ Variational learning of  $z$  and kernel hyperparameters

$$\rightarrow \arg \min_{z, \mu, A} \text{KL}(q(f_*) \| p(f_*|y)) \Leftrightarrow \arg \min_{z, \mu, A} \text{KL}(q(f_*, u) \| p(f_*, u|y))$$

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

$$\geq \int q(f_*, u) \log \frac{p(f_*, u, y)}{q(f_*, u)} df_* du$$

$$= \int p(f_*) \phi(u) \log \frac{p(y|f_*) p(u)}{\phi(u)} df_* du$$

$$= \int \phi(u) \left\{ \int p(f_*|u) \log p(y|f_*) df_* + \log \frac{p(u)}{\phi(u)} \right\} du$$

$$\begin{aligned} & \int p(f_*|u) \left[ -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \text{Tr}[y y^T - 2y f_u^T + f_u f_u^T] \right] df_* \\ &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \text{Tr}[y y^T - 2y f_u^T + f_u f_u^T + \alpha x^T + K_f - Q_f] \\ &= \log N(\alpha, \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}[K_f - Q_f] \end{aligned}$$

where  $\alpha = K_u^T K_u^{-1} u$ ,  $Q_f = K_u^T K_u^{-1} K_u$

$$= \int \phi(u) \log \frac{N(\alpha, \sigma^2 I) p(u)}{\phi(u)} du - \frac{1}{2\sigma^2} \text{Tr}[K_f - Q_f]$$

$$\geq \log \int N(\alpha, \sigma^2 I) p(u) du - \frac{1}{2\sigma^2} \text{Tr}[K_f - Q_f] \quad (\text{Jensen's inequality})$$

$$= \log N(0, K_u^T K_u^{-1} K_u + \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}[K_f - Q_f]$$

$$M_{\text{opt}} = \sigma^2 K_u (K_u + \sigma^2 K_u^T K_u^{-1})^{-1} K_u^T y$$

$$A_{\text{opt}} = K_u (K_u + \sigma^2 K_u^T K_u^{-1})^{-1} K_u$$

### 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

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

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

⇒ Orthogonal decomposition

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

$$\rightarrow f = f_{||} + f_\perp \text{ where } f_{||} \in V, f_\perp \perp V$$

$$\left( \begin{array}{l} f_{||} = k_{uf}^T k_u^{-1} u \sim N(0, k_{uf}^T k_u^{-1} k_{uf}) \\ f_\perp = f - f_{||} \sim N(0, k_f - k_{uf}^T k_u^{-1} k_{uf}) \end{array} \right)$$

$$\Rightarrow \text{joint distribution } (p(y, u, f_\perp, v))$$

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

$$\Rightarrow \text{variational distribution } (q(u, f_\perp, v))$$

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

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

$$\left( \begin{array}{l} m_{f_\perp} = C_{vf}^T C_v^{-1} m_v \\ S_{f_\perp} = C_{ff} - C_{vf}^T C_v^{-1} (C_v - S_v) C_v^{-1} C_{vf} \end{array} \right) \text{ where } C_{ab} = k_{ab} - k_{ua}^T k_u^{-1} k_{ub}$$

⇒ variational learning of  $z, \sigma, m_v, S_v$  and kernel hyperparameters

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

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

↓

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

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

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

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

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

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

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

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

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