

1 Gaussian Process

[1] introduces Gaussian Process regression as an alternative view to Bayesian regression.

1.1 Bayesian Regression

In linear regression setup, we assume output $y \in \mathbb{R}$ is a linear function of inputs $x \in \mathbb{R}^d$, corrupted with additive iid normal noise,

$$y = f(x) + \epsilon \quad \text{where} \quad f(x) = w^T x \quad \text{and} \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2) \quad (1)$$

The Bayesian setup considers $w \in \mathbb{R}^d$ as a random variable, endowed with prior $w \sim \mathcal{N}(0, \Sigma_p)$. Using Bayes rule, we can find the posterior of weights given data, which is again a normal random variable $p(w | X, y) = \mathcal{N}(w; A^{-1}b, A^{-1})$ where $A = \frac{1}{\sigma_n^2} X^T X + \Sigma_p^{-1}$ and $b = \frac{1}{\sigma_n^2} X^T y$ and $X \in \mathbb{R}^{n \times d}$, $y \in \mathbb{R}^{n \times 1}$ are design matrices. For test point x_* , the predictive distribution of $f_* = f(x_*)$ is the average likelihood of f_* under model $f(x; w)$ with respect to posterior of w .

$$p(f_* | x_*, X, y) = \int p(f_* | x_*, w) p(w | X, y) dw \quad (2)$$

We can think of the predictive distribution as a linear function $f_* = x_*^T w$ of weights, a normal random variable, and therefore is normal. Therefore, $f_* | x_*, X, y \sim \mathcal{N}(x_*^T A^{-1} b, x_*^T A^{-1} x_*)$. The natural extension to Bayesian linear regression is to kernelize it, for example assume a linear model in some feature space $f(x) = \phi(x)^T w$ for some feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$. We can write predictive distribution as

$$f_* | x_*, X, y \sim \mathcal{N}(k(x_*, X)(k(X, X) + \sigma_n^2 I)^{-1} y, \quad (3)$$

$$k(x_*, x_*) - k(x_*, X)(k(X, X) + \sigma_n^2 I)^{-1} k(X, x_*)) \quad (4)$$

where $k(X, X') = \Phi \Sigma_p \Phi^T \in \mathbb{R}^{n \times n'}$ and $\Phi \in \mathbb{R}^{n \times D}$ are the feature vectors.

1.2 Gaussian Process

Definition 1. (Gaussian Process) A Gaussian process is a stochastic process $\{X_t\}_{t \in T}$ such that, for every finite subset of indices $t_1, \dots, t_k \in T$, $(X_{t_1}, \dots, X_{t_k})$ is multivariate normal.

A Gaussian process $f \sim \mathcal{GP}(m, k)$ over \mathcal{X} is fully specified by the mean function $m : \mathcal{X} \rightarrow \mathbb{R}$ and covariance function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ where

$$m(x) = \mathbb{E}[f(x)] \quad k(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x')))] \quad (5)$$

The covariance function determines function's behavior, for example its stationarity, smoothness, and periodicity etc. For example, the squared exponential covariance function $k_{SE}(x, x') = \exp(-\frac{1}{2\ell^2} \|x - x'\|_2^2)$ enforces the prior knowledge that functions are smooth, i.e. inputs are close in the Euclidean sense will have similar outputs. See Figure (1) for some examples of samples from Gaussian process.

Example 1. (Intuition about covariance matrix for \mathcal{GP}) First consider $(y_1, y_2) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ where $\Sigma_{11} = \Sigma_{22} = 1$ and $\Sigma_{12} = \Sigma_{21} = \rho_{12}$. We know $y_1 | y_2 = a \sim \mathcal{N}(\rho_{12}a, 1 - \rho_{12}^2)$. When $\text{Cov}(y_1, y_2) = \rho_{12} \uparrow 1$, y_1 's samples conditioned on $y_2 = a$ fall close to a with high probability. When $\rho_{12} \downarrow 0$, $y_1 | y_2 = a$ will distribute like a unit normal, regardless of values that y_2 take. Extend this intuition to gaussian process: whenever $k(x_i, x_j)$ is large, y_i, y_j are correlated and so observing $y_i = a$ provides a strong prior on how y_j will behave, or in other words, reduce the uncertainty of values that y_j can take dramatically.

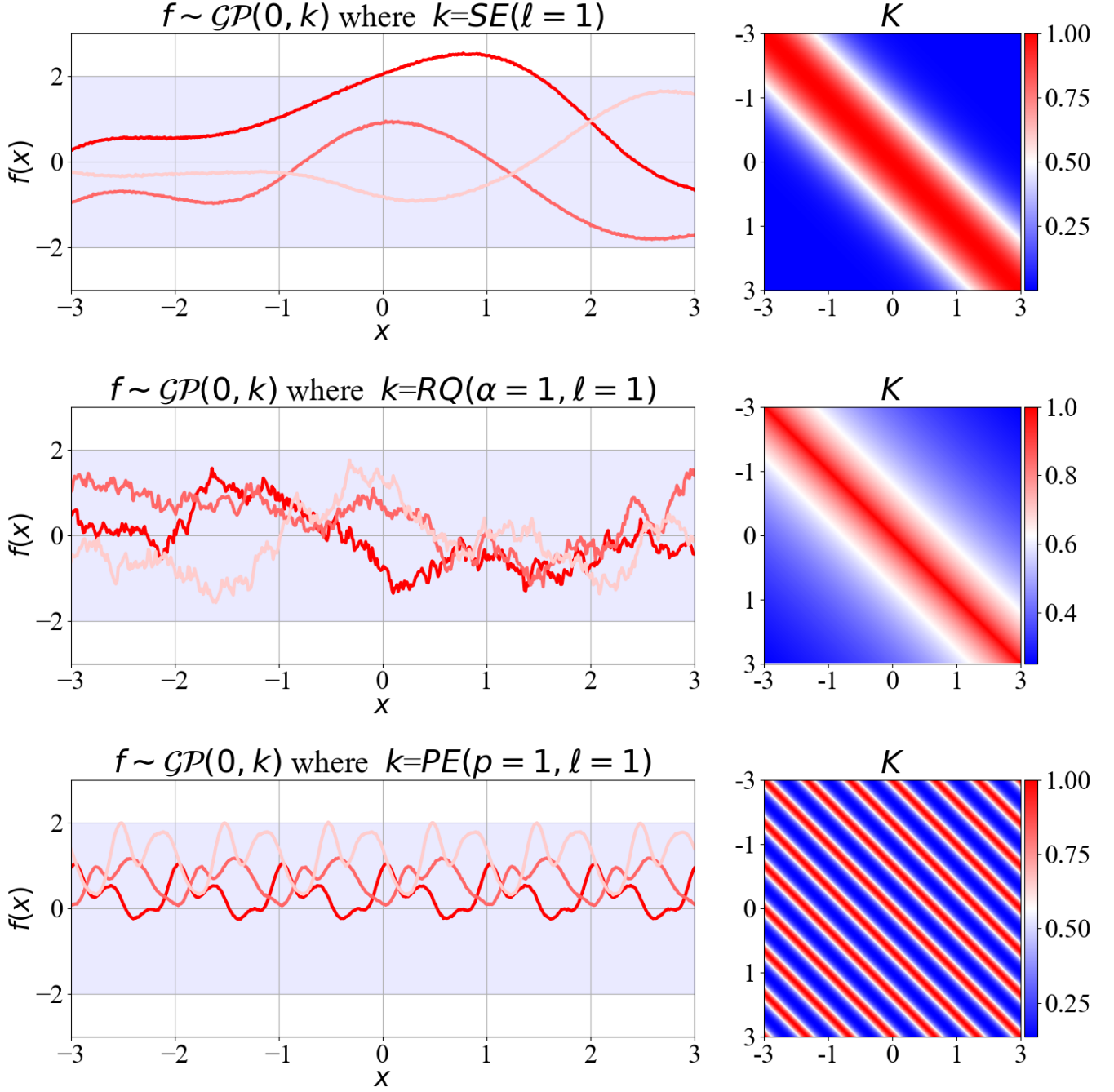


Figure 1: (Left) Samples from Gaussian process prior and (Right) covariance matrix at test locations.

1.3 Gaussian Process Regression

Instead of placing a prior over weights $p(w)$ to quantify randomness in function $f(x) = \phi(x)^T w$, we model function directly as a Gaussian process, $f \sim \mathcal{GP}(0, k)$. There is a one-to-one correspondence between the two views. For example, $f(x) = \phi(x)^T w$ with prior $w \sim \mathcal{N}(0, \Sigma_p)$ used in kernel Bayesian regression has

$$\mathbb{E}[f(x)] = \phi(x)^T \mathbb{E}[w] = 0 \quad \mathbb{E}[f(x)f(x')] = \phi(x)^T \Sigma_p \phi(x') \quad (6)$$

Therefore, $f \sim \mathcal{GP}(0, k)$ where $k(x, x') = \phi(x)^T \Sigma_p \phi(x')$. Note k is in fact a valid kernel. (Since Σ_p is psd, $\Sigma_p = UDU^T$ by SVD. We can write $k(x, x') = \langle \psi(x), \psi(x') \rangle$ where $\psi(x) = \Sigma_p^{1/2} \phi(x)$ and $\Sigma_p^{1/2} = UD^{1/2}U^T$). Conversely, any valid kernel used in kernel Bayesian regression can be used to parameterize the covariance function of the Gaussian process model over f .

Since $y = f(x) + \epsilon$, we have $\mathbf{y} = \mathbf{f} + \sigma_n^2 I \sim \mathcal{N}(0, k(X, X) + \sigma_n^2 I)$ where $\mathbf{y}, \mathbf{f} \in \mathbb{R}^{n \times 1}$. We can write the joint distribution of observed values and function values at some test locations X_* as

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} k(X, X) + \sigma_n^2 I & k(X, X_*) \\ k(X_*, X) & k(X_*, X_*) \end{bmatrix} \right) \quad (7)$$

We can derive the predictive distribution for $\mathbf{f}_* | \mathbf{y}$ by simply apply conditional distribution formula

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(k(X_*, X)(k(X, X) + \sigma_n^2 I)^{-1} \mathbf{y} \quad (8)$$

$$k(X_*, X_*) - k(X_*, X)(k(X, X) + \sigma_n^2 I)^{-1} k(X, X_*)) \quad (9)$$

which has exact form compared to (4). More compactly, for a single test point x_* , $\mu_{\mathbf{f}_*} = k_*^T (K + \sigma_n^2 I)^{-1} \mathbf{y}$ and $\text{Var}(\mathbf{f}_*) = k(x_*, x_*) - k_*^T (K + \sigma_n^2 I)^{-1} k_*$ where $k_* = k(X, x_*) \in \mathbb{R}^{n \times 1}$. See Figure (2) for examples of Gaussian process regression with varying data size and hyperparameters.

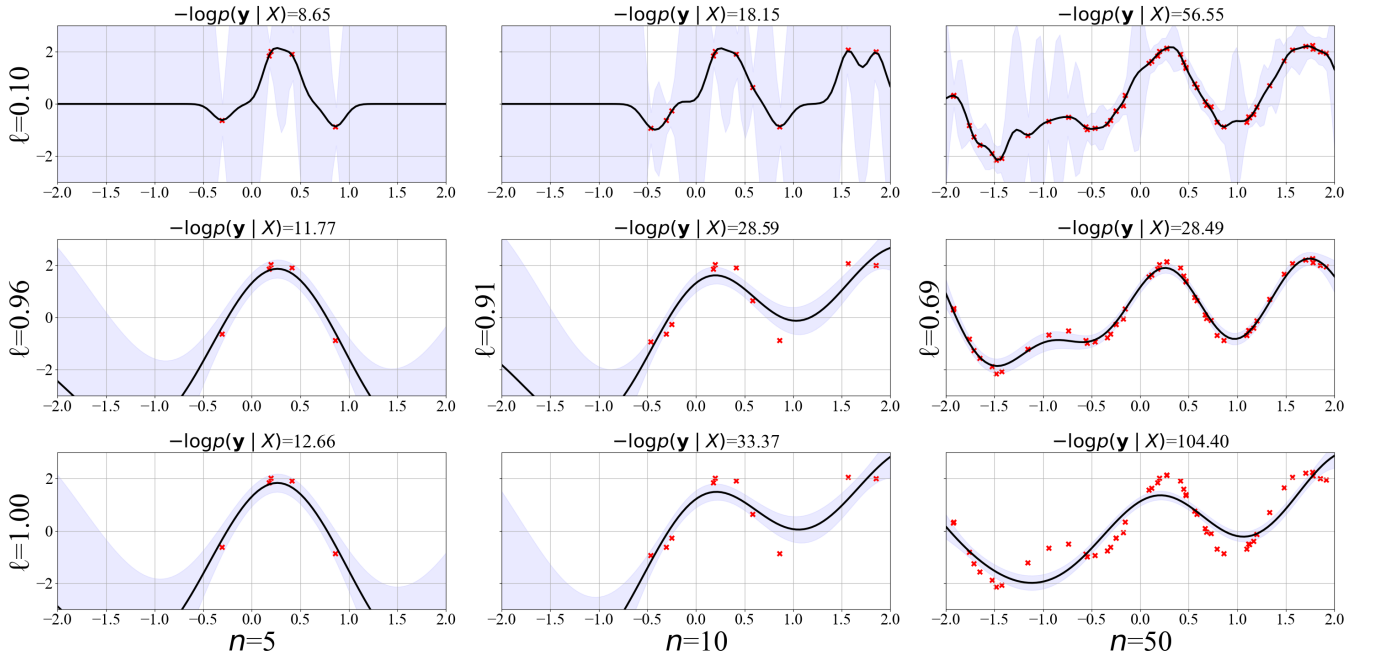


Figure 2: This plot shows mean (black) and 95% confidence interval (light blue) for predictive distribution $p(\mathbf{f}_* | X, \mathbf{y}, X_*)$ fit using Gaussian process regression assuming a SE prior over $f \sim \mathcal{GP}(0, k_{SE})$ of varying lengthscale ℓ and observed sample sizes n given groundtruth noise $\sigma_n = 0.1$. The middle row's hyperparameters is taken to be the empirical Bayes estimate $\ell, \sigma_n = \arg \max p(\mathbf{y} | X, \ell, \sigma_n)$ where $p(\mathbf{y} | X, \ell, \sigma_n)$ given by Equation (11), optimized via gradient descent. Note estimating $\log(\ell), \log(\sigma_n)$ makes the optimization problem much easier.

1.4 Model Selection

Model selection for Gaussian process regression involves picking the form and the hyperparameters of the covariance function. Given data (X, \mathbf{y}) , marginal likelihood $p(\mathbf{y} | X)$ quantifies how likely data is observed under our additive noise model $\mathbf{y} | \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma_n^2 I)$ on average with respect to latent function values (or parameters of our model) $\mathbf{f} | X \sim \mathcal{N}(0, K)$ (due to assumption of Gaussian process prior over \mathbf{f}),

$$p(\mathbf{y} | X) = \int p(\mathbf{y} | \mathbf{f}, X) p(\mathbf{f} | X) d\mathbf{f} \quad (10)$$

This formulation is analogous to that in Bayesian linear regression where the marginal likelihood marginalizes over weights $p(\mathbf{y} | X) = \int p(\mathbf{y} | X, w) p(w) dw$. We can obtain a closed form expression by reading off (7), i.e. $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, K + \sigma_n^2 I)$.

$$\log p(\mathbf{y} | X) = -\frac{1}{2} \mathbf{y}^T (K + \sigma_n^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{n}{2} \log 2\pi \quad (11)$$

In Bayesian model selection, model hyperparameters for Gaussian process regression $\theta = \{\sigma_n, \ell\}$ can be found by maximizing the marginal likelihood or the type II likelihood $\theta^* = \arg \max_{\theta} \log p(\mathbf{y} | X, \theta)$.

2 Approximate Gaussian Process

2.1 Inducing Variables (DIC/FITC)

The unifying view paper [2] gives a pretty good explanation/interpretation of sparse GP methods based on inducing points till circa 2005. The basic idea is to augment existing model $p(\mathbf{y}, \mathbf{f})$ with m additional inducing variables $\mathbf{u} = \{\mathbf{u}_i\}_{i=1}^m$ that are instantiation from the same GP prior, corresponding to a set of inducing locations $X_{\mathbf{u}}$. Inducing variables act as an information channel between training/testing latent function values \mathbf{f}, \mathbf{f}_* , alternatively we can think of \mathbf{u} as a sufficient statistics of potentially redundant latent function values \mathbf{f} . This translates to a conditional independence assumption $\mathbf{f} \perp\!\!\!\perp \mathbf{f}_* | \mathbf{u}$, which implies the following factorization,

$$p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}, \mathbf{f}_*, \mathbf{u}) = p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \mathbf{u}) p(\mathbf{f}_* | \mathbf{u}) p(\mathbf{u}) \quad (12)$$

Without additional assumptions, we can think of joint priors factorizes akin to a latent variable model,

$$\mathbf{u} \sim \mathcal{N}(0, K_{\mathbf{uu}}) \quad (13)$$

$$\mathbf{f} | \mathbf{u} \sim \mathcal{N}(K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mathbf{u}, K_{\mathbf{uu}} - Q_{\mathbf{ff}}) \quad (14)$$

where $Q_{\mathbf{ff}} = K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} K_{\mathbf{uf}}$ is Nyström approximation of $K_{\mathbf{ff}}$. Various approximation schemes corresponds to approximations to the conditional prior, $q(\mathbf{f} | \mathbf{u}), q(\mathbf{f}_* | \mathbf{u})$. For exact GP, inference using the posterior GP $p(\mathbf{f} | \mathbf{y})$ requires inversion of a large kernel matrix; The introduction of inducing variables modifies the form that posterior GP $q(\mathbf{f} | \mathbf{y})$ takes.

To start, the Subset of Regressor (SoR) or Deterministic Inducing Conditional (DIC) approximation assumes that the latent function values as a linear mapping of inducing variables, e.g. $\mathbf{f} = K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mathbf{u}$. Note that the linear projection is one which preserves the mean of conditional prior in Equation (14). This assumption gives to an approximate conditional prior with zero covariance, $q_{\text{SoR}}(\mathbf{f} | \mathbf{u}) = \mathcal{N}(K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mathbf{u}, 0)$. Because of the linear mapping, $\mathbb{E}_q[\mathbf{f}] = 0$ and $\text{Cov}_q(\mathbf{f}) = K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} K_{\mathbf{uu}} K_{\mathbf{uu}}^{-1} K_{\mathbf{uf}} = Q_{\mathbf{ff}}$, similarly for $q(\mathbf{f}_*)$, and therefore,

$$q_{\text{DIC}}(\mathbf{f}, \mathbf{f}_*) = \mathcal{N}\left(0, \begin{bmatrix} Q_{\mathbf{ff}} & Q_{\mathbf{f}^*} \\ Q_{\mathbf{f}^*} & Q_{**} \end{bmatrix}\right) \quad (15)$$

Inference with the DIC approximation is equivalent to exact inference with a modified covariance function $k_{\text{DIC}}(x_i, x_j) = k(x_i, \mathbf{u})K_{\mathbf{uu}}^{-1}k(\mathbf{u}, x_j)$. The predictive distribution is then given by

$$q_{\text{DIC}}(\mathbf{f}_* | \mathbf{y}) \sim \mathcal{N}(Q_{*\mathbf{f}}(Q_{\mathbf{ff}} + \sigma_n^2 I)^{-1}\mathbf{y}, Q_{**} - Q_{*\mathbf{f}}(Q_{\mathbf{ff}} + \sigma_n^2 I)^{-1}Q_{\mathbf{f}*}) \quad (16)$$

Inversion of a $n \times n$ matrix is replaced with inversion of $m \times m$ matrix, making computation faster.

Sparse Gaussian Process using Pseudo-inputs (SGPP) or Fully Independent Training Conditional (FITC) approximation is another popular approximate GP method [3]. Simply, FITC assumes $\mathbf{f}_i \perp \mathbf{f}_j | \mathbf{u}$ for all $i, j \in [n]$, henceforth the conditional priors are fully (conditionally) independent.

$$q_{\text{FITC}}(\mathbf{f} | \mathbf{u}) = \prod_{i=1}^n p(f_i | \mathbf{u}) = \mathcal{N}(K_{\mathbf{fu}}K_{\mathbf{uu}}^{-1}\mathbf{u}, \text{diag}[K_{\mathbf{uu}} - Q_{\mathbf{ff}}]) \quad (17)$$

Compare this to Equation (14), we see that FITC replaces approximate covariances with exact covariances on the diagonal. Equivalently, we can think of \mathbf{f} as a noisy linear map of \mathbf{u} where the noise has a particular form, e.g. $\mathbf{f} = K_{\mathbf{fu}}K_{\mathbf{uu}}^{-1}\mathbf{u} + \gamma$ where $\gamma \sim \mathcal{N}(0, \text{diag}[K_{\mathbf{uu}} - Q_{\mathbf{ff}}])$. We can derive the covariance for the approximate posterior $\text{Cov}_q(\mathbf{f}) = \text{Cov}(K_{\mathbf{fu}}K_{\mathbf{uu}}^{-1}\mathbf{u}) + \text{Cov}(\gamma) = Q_{\mathbf{ff}} + \text{diag}[K_{\mathbf{uu}} - Q_{\mathbf{ff}}]$. Therefore,

$$q_{\text{FITC}}(\mathbf{f}, \mathbf{f}_*) = \mathcal{N}\left(0, \begin{bmatrix} Q_{\mathbf{ff}} - \text{diag}[Q_{\mathbf{ff}} - K_{\mathbf{ff}}] & Q_{\mathbf{f}*} \\ Q_{*\mathbf{f}} & K_{**} \end{bmatrix}\right) \quad (18)$$

Inference with FITC approximation is equivalent to exact inference with a modified covariance function $k_{\text{FITC}}(x_i, x_j) = k_{\text{DIC}}(x_i, x_j) + \delta_{i,j}(k(x_i, x_j) - k_{\text{DIC}}(x_i, x_j))$. The predictive distribution follows,

$$q_{\text{FITC}}(\mathbf{f}_* | \mathbf{y}) \sim \mathcal{N}(Q_{*\mathbf{f}}(Q_{\mathbf{ff}} + \Lambda)^{-1}\mathbf{y}, Q_{**} - Q_{*\mathbf{f}}(Q_{\mathbf{ff}} + \Lambda)^{-1}Q_{\mathbf{f}*}) \quad (19)$$

where $\Lambda = \text{diag}[K_{\mathbf{ff}} - Q_{\mathbf{ff}} + \sigma_n^2 I]$.

One popular methods for selecting inducing inputs \mathbf{X}_u is simply by maximizing the marginal likelihood

$$q(\mathbf{y} | X_u) = \iint p(\mathbf{y} | \mathbf{f})q(\mathbf{f} | \mathbf{u})p(\mathbf{u} | X_u) d\mathbf{u} d\mathbf{f} \quad (20)$$

2.2 Variational Formulation (VFE/SVGP)

Sparse GP with inducing variables is inflicted with the problem that modeling assumptions, e.g. choice of covariance function, is intertwined with approximation assumptions. Specifically, we see that inference using approximation schemes like DIC and FITC can be interpreted as performing exact inference on GP model with a modified covariance function, where inducing locations X_u are part of kernel hyperparameters. This is undesirable.

[4] proposes to use variational inference to learning inducing locations. [5] extends the variational bound for stochastic gradient optimization, [6] extends the variational bound to account for non-conjugate likelihood. The accompanying technical report for the 2009 paper [7], the gentle tutorial [8], another report by Thang Bui [9], Mark van der Wilk's PhD thesis [10] are helpful in understanding the derivations.

[4] proposes Variational Free Energy (VFE) which applies variational inference to learn inducing variables in sparse GP, alleviating the above concern. Similar to inducing variable approximations, the model is augmented with inducing variables and assume the same conditional independence relationship so that the joint density factorizes $p(\mathbf{y}, \mathbf{f}, \mathbf{u}) = p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{u})p(\mathbf{u})$. Here \mathbf{y} is observed while $\{\mathbf{u}, \mathbf{f}\}$ are latent variables. In contrast to inducing variable methods which makes approximating assumptions about the conditional priors, the paper introduce a variational density that preserves the true conditional prior with a freely varying prior over the inducing variables, e.g. $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f} | \mathbf{u})q(\mathbf{u})$ for some $q(\mathbf{u}) \sim \mathcal{N}(\mu_u, \Sigma_u)$

that depends on inducing locations $X_{\mathbf{u}}$ as well as kernel hyperparameters. Together with Equation (14), we can derive an approximate posterior for \mathbf{f} ,

$$q(\mathbf{f}) = \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) d\mathbf{u} = \mathcal{N}(K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mu_{\mathbf{u}}, K_{\mathbf{ff}} - Q_{\mathbf{ff}} + K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \Sigma_{\mathbf{u}} K_{\mathbf{uu}}^{-1} K_{\mathbf{uf}}) \quad (21)$$

The paper then proposes to find variational parameters $\{\mu_n, \Sigma_n, X_{\mathbf{u}}, \theta\}$ by minimizing the distance between the approximate and true posterior, $q^* = \arg \min_{q(\mathbf{u})} \text{KL}(q(\mathbf{f}, \mathbf{u}) | p(\mathbf{f}, \mathbf{u} | \mathbf{y}))$. As usual in variational inference, we derive the ELBO to make objective tractable for optimization,

$$\text{KL}(q(\mathbf{f}, \mathbf{u}) | p(\mathbf{f}, \mathbf{u} | \mathbf{y})) = \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) \log \left(\frac{p(\mathbf{f} | \mathbf{u}) q(\mathbf{u})}{p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \mathbf{u}) p(\mathbf{u})} \right) d\mathbf{f} d\mathbf{u} + \log p(\mathbf{y}) \quad (22)$$

$$\log p(\mathbf{y}) \geq \int q(\mathbf{u}) \int p(\mathbf{f} | \mathbf{u}) \log p(\mathbf{y} | \mathbf{f}) d\mathbf{f} d\mathbf{u} - \text{KL}(q(\mathbf{u}) || p(\mathbf{u})) \quad (23)$$

With $\mathbb{E}[\mathbf{f} | \mathbf{u}] = K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mathbf{u}$ and $\text{Cov}(\mathbf{f} | \mathbf{u}) = K_{\mathbf{ff}} - Q_{\mathbf{ff}}$, we can simplify,

$$\langle \log p(\mathbf{y} | \mathbf{f}) \rangle_{p(\mathbf{f} | \mathbf{u})} = \int p(\mathbf{f} | \mathbf{u}) \left[-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{tr}(\mathbf{y}\mathbf{y}^T - 2\mathbf{y}\mathbf{f}^T + \mathbf{f}\mathbf{f}^T) \right] d\mathbf{f} \quad (24)$$

$$= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{tr}(\mathbf{y}\mathbf{y}^T - 2\mathbf{y}(\mathbb{E}[\mathbf{f} | \mathbf{u}])^T + \text{Cov}(\mathbf{f} | \mathbf{u}) + \mathbb{E}[\mathbf{f} | \mathbf{u}] \mathbb{E}[\mathbf{f} | \mathbf{u}]^T) \quad (25)$$

$$= \log \mathcal{N}(\mathbf{f} | K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mathbf{u}, \sigma^2 I) - \frac{1}{2\sigma^2} \text{tr}(K_{\mathbf{ff}} - Q_{\mathbf{ff}}) \quad (26)$$

Now with $P := K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1}$, $\mathbb{E}_q[\mathbf{u}] = \mu_{\mathbf{u}}$, and $\text{Cov}_q(\mathbf{u}) = \Sigma_{\mathbf{u}} + \mu_{\mathbf{u}} \mu_{\mathbf{u}}^T$, then

$$\langle \log p(\mathbf{y} | \mathbf{f}) \rangle_{q(\mathbf{f}, \mathbf{u})} = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{tr}(\mathbf{y}\mathbf{y}^T - 2\mathbf{y}(P\mu_{\mathbf{u}})^T + P(\Sigma_{\mathbf{u}} + \mu_{\mathbf{u}} \mu_{\mathbf{u}}^T)P^T + K_{\mathbf{uu}} - Q_{\mathbf{uu}}) \quad (27)$$

$$= \log \mathcal{N}(\mathbf{f} | K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mu_{\mathbf{u}}, \sigma^2 I) - \frac{1}{2\sigma^2} \text{tr}(K_{\mathbf{ff}} - Q_{\mathbf{ff}}) - \frac{1}{2\sigma^2} \text{tr}(K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \Sigma_{\mathbf{u}} K_{\mathbf{uu}}^{-1} K_{\mathbf{uf}}) \quad (28)$$

[4] finds the optimal variational distribution $q(\mathbf{u})$ and substitute back into the lower bound, resulting in a bound where the first term is simply the marginal likelihood for sparse GP with DIC/DTC approximation while the second term regularizes approximate kernel be close to the true kernel matrix,

$$\mathcal{L}_{\text{VFE}} := \log \mathcal{N}(\mathbf{f} | K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} \mathbf{u}, Q_{\mathbf{ff}} - \sigma^2 I) - \frac{1}{2\sigma^2} \text{tr}(K_{\mathbf{ff}} - Q_{\mathbf{ff}}) \quad (29)$$

The bound is strictly better with additional inducing variables and becomes tight when $X_{\mathbf{u}} = X$, recovering the full GP marginal likelihood. Maximizing \mathcal{L}_{VFE} with respect to kernel hyperparameters and inducing locations is tractable as the objective requires $\mathcal{O}(nm^2)$ to compute.

[5] proposes Stochastic Variational Gaussian Process (SVGP) that extends the bound (23) to allow for stochastic variational inference. The key observation is that with $p(\mathbf{y} | \mathbf{f}) = \prod_{i=1}^n p(y_i | f_i)$, the lower bound factors with respect to data points. Start with Equation (23),

$$\log p(\mathbf{y}) \geq \langle \log p(\mathbf{y} | \mathbf{f}) \rangle_{q(\mathbf{f})} - \text{KL}(q(\mathbf{u}) || p(\mathbf{u})) = \sum_{i=1}^n \langle \log p(y_i | f_i) \rangle_{q(f_i)} - \text{KL}(q(\mathbf{u}) || p(\mathbf{u})) \quad (30)$$

$$\langle \log p(y_i | f_i) \rangle_{q(f_i)} = \langle -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y_i^2 - 2y_i f_i + f_i^2) \rangle_{q(f_i)} \quad (31)$$

$$= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y_i^2 - 2y_i \mathbb{E}_q[f_i] + \text{Var}_q(f_i) + \mathbb{E}_q[f_i]^2) \quad (32)$$

$$= \log \mathcal{N}(y_i | K_{f_i \mathbf{u}} K_{\mathbf{uu}}^{-1} \mu_{\mathbf{u}}, \sigma^2) - \frac{1}{2\sigma^2} (K_{f_i f_i} - Q_{f_i f_i} - K_{f_i \mathbf{u}} K_{\mathbf{uu}}^{-1} \Sigma_{\mathbf{u}} K_{\mathbf{uu}}^{-1} K_{\mathbf{u} f_i}) \quad (33)$$

where we have used fact that $q(f_i) \sim \mathcal{N}(P_i \mu_{\mathbf{u}}, K_{f_i f_i} - Q_{f_i f_i} - P_i \Sigma_{\mathbf{u}} P_i^T)$ from (21) and write $P_i := K_{f_i \mathbf{u}} K_{\mathbf{u} \mathbf{u}}^{-1}$. This de-coupled term when summed over data points is exactly same as the coupled term in Equation (28). Now we can write down the lower bound which can be optimized with stochastic optimization,

$$\mathcal{L}_{\text{SVGP}} := \sum_{i=1}^n \left(\log \mathcal{N}(y_i | P_i \mu_{\mathbf{u}}, \sigma^2) - \frac{1}{2\sigma^2} (K_{f_i f_i} - Q_{f_i f_i} - P_i \Sigma_{\mathbf{u}} P_i^T) \right) - \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) \quad (34)$$

where computation of KL term has $\mathcal{O}(m^3)$ cost

$$\text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) = \frac{1}{2} \left(\text{tr}(K_{\mathbf{u} \mathbf{u}}^{-1} \Sigma_{\mathbf{u}}) + \mu_{\mathbf{u}}^T K_{\mathbf{u} \mathbf{u}}^{-1} \mu_{\mathbf{u}} - m + \log \frac{|K_{\mathbf{u} \mathbf{u}}|}{|\Sigma_{\mathbf{u}}|} \right) \quad (35)$$

3 Multitask Learning

3.1 MTGP

Multitask Gaussian Process (MTGP) regression [11] is a method to do multitask learning using Gaussian process. Given design $X \in \mathbb{R}^{N \times D}$, $Y \in \mathbb{R}^{N \times M}$, we want to learn a vector valued regressor $f : \mathbb{R}^D \rightarrow \mathbb{R}^M$ that fits data well. We can put a GP prior over $f \sim \mathcal{GP}(0, k)$ where the covariance function models both the relationship between inputs via k_x and similarity between tasks / output coordinate via k_t , specifically define k as tensor product $k_x \otimes k_t$ over $\mathcal{X} \times \mathcal{T}$ where $\mathcal{T} = \{1, 2, \dots, M\}$ is space of tasks,

$$k((x, t), (x', t')) = k_x(x, x') k_t(t, t') \quad (36)$$

Consider a likelihood model with task specific noise variance $y_t \sim \mathcal{N}(f_t(x), \sigma_t^2)$. Denote Y_t as t -th column in Y , then $Y_t \sim \mathcal{N}(0, k_t(t, t)K + \sigma_t^2 I)$ where $[K]_{ij} = k_x(x_i, x_j)$. Furthermore, if we let $\mathbf{y} = \text{vec}(Y) = (Y_1, \dots, Y_M)$, then $\mathbf{y} \sim \mathcal{N}(0, \Sigma)$ where

$$\Sigma = K^t \otimes K + D \otimes I = \begin{bmatrix} K_{11}^t K + \sigma_1^2 I & K_{12}^t K & \dots & K_{1M}^t K \\ K_{21}^t K & K_{22}^t K + \sigma_2^2 I & \dots & K_{2M}^t K \\ \vdots & \vdots & \ddots & \vdots \\ K_{M1}^t K & K_{M2}^t K & \dots & K_{MM}^t K + \sigma_M^2 I \end{bmatrix} \quad (37)$$

Note this is the covariance matrix with no missing outputs, which leads to efficient inference methods. The task similarity matrix K_t ($[K_t]_{ij} = k_t(i, j)$) induce correlation between tasks. Intuitively, if two tasks t, t' are related in the sense an optimal regressor for $y_t, y_{t'}$ vary together in some systematic manner, then observation at some location (x, y_t) constrains what the value $y_{t'}$ can take. The degree to which we can reduce the uncertainty that $y_{t'}$ can takes depends on $k_t(t, t')$. Therefore, a learner which utilizes information obtained from related tasks can be learnt more efficiently than if learnt without knowledge of other tasks, in which case $K_t = I$. [11] proposes to use EM to first impute missing latent variables \mathbf{f} from noisy observations \mathbf{y} , then find hyperparameters $\theta = \{\theta_{k_x}, K^t\}$ by maximizing the full data likelihoods $p(\mathbf{y}, \mathbf{f} | \theta_k)$. Alternatively, we can simply can learn the hyperparameters by applying gradient descent to maximize type-II likelihood $\theta^* = \arg \max_{\theta} \log p(\mathbf{y} | X, \theta)$. See Figure (3) for an example of MTGP in case where $D = 1$ and $M = 2$.

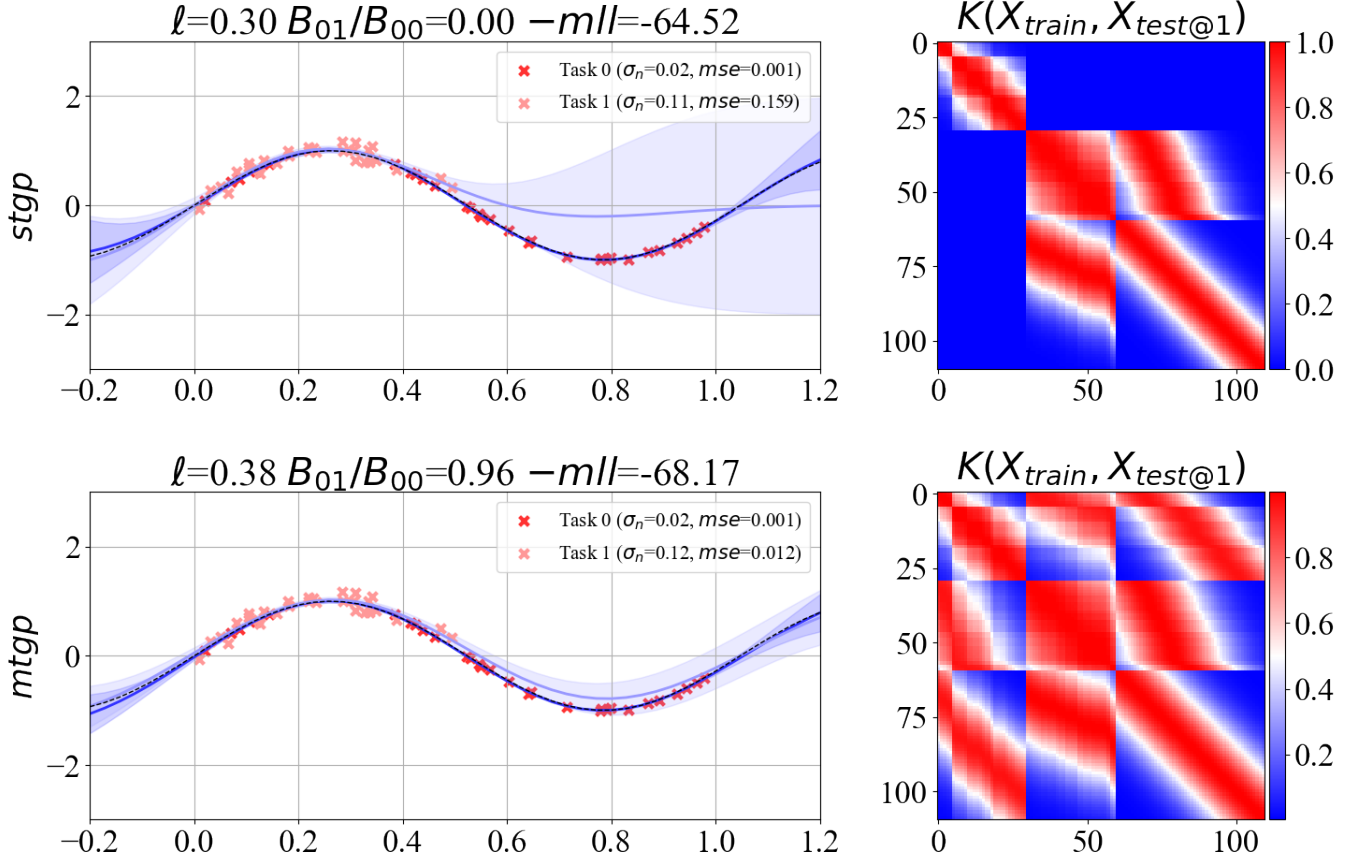


Figure 3: Observations for the two tasks are generated $y_0(x) = \sin(6x) + .03\epsilon$ and $y_1(x) = \sin(6x + \text{shift}) + .1\epsilon$ where $\epsilon \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ and $\text{shift} = .3$. Here for task 1 we only sample points in $[0, .5]$ while for task 0 we sample points in $[0, 1]$. (Top) Two Gaussian process regressor are learnt to fit data independently. (Bottom) A MTGP is fit to data whereby the task similarity matrix $B := K^t$ is learnt by maximizing marginal log likelihood. We see that by taking into account of relationship between tasks, samples from task 0 reduces uncertainty that task 1 has at test locations far from its training support.

3.2 Asymmetric MTGP

When the goal is to improve performance of a target/main task given the other tasks, optimizing for maximum marginal likelihood of all tasks is suboptimal, as shown in Figure (4). There are extension of MTGP to the asymmetric case by assuming that auxiliary task functions are addition of a shared main task function and task specific function [12], however it still uses the full marginal likelihood for all tasks to find task similarity and kernel hyperparameters. Ideally, we want

1. An objective which depends on likelihood of main task only
2. Automatically learn the task similarity such that auxiliary task or side information can improve learning of main task. If a side task is totally unrelated to the main task, it should have minimal influence on learning of the main task.

Given design $X \in \mathbb{R}^{N \times D}$, $Y \in \mathbb{R}^{N \times 1}$, $S \in \mathbb{R}^{N \times M}$, we want to learn a scalar valued regressor $f : \mathbb{R}^D \rightarrow \mathbb{R}$ that fits the data well. We can alternatively try to learn the scalar valued regressor taking into account of relevant side information $f : \mathbb{R}^D \times \mathbb{R}^M \rightarrow \mathbb{R}$. Relevant side information that helps with learning the main

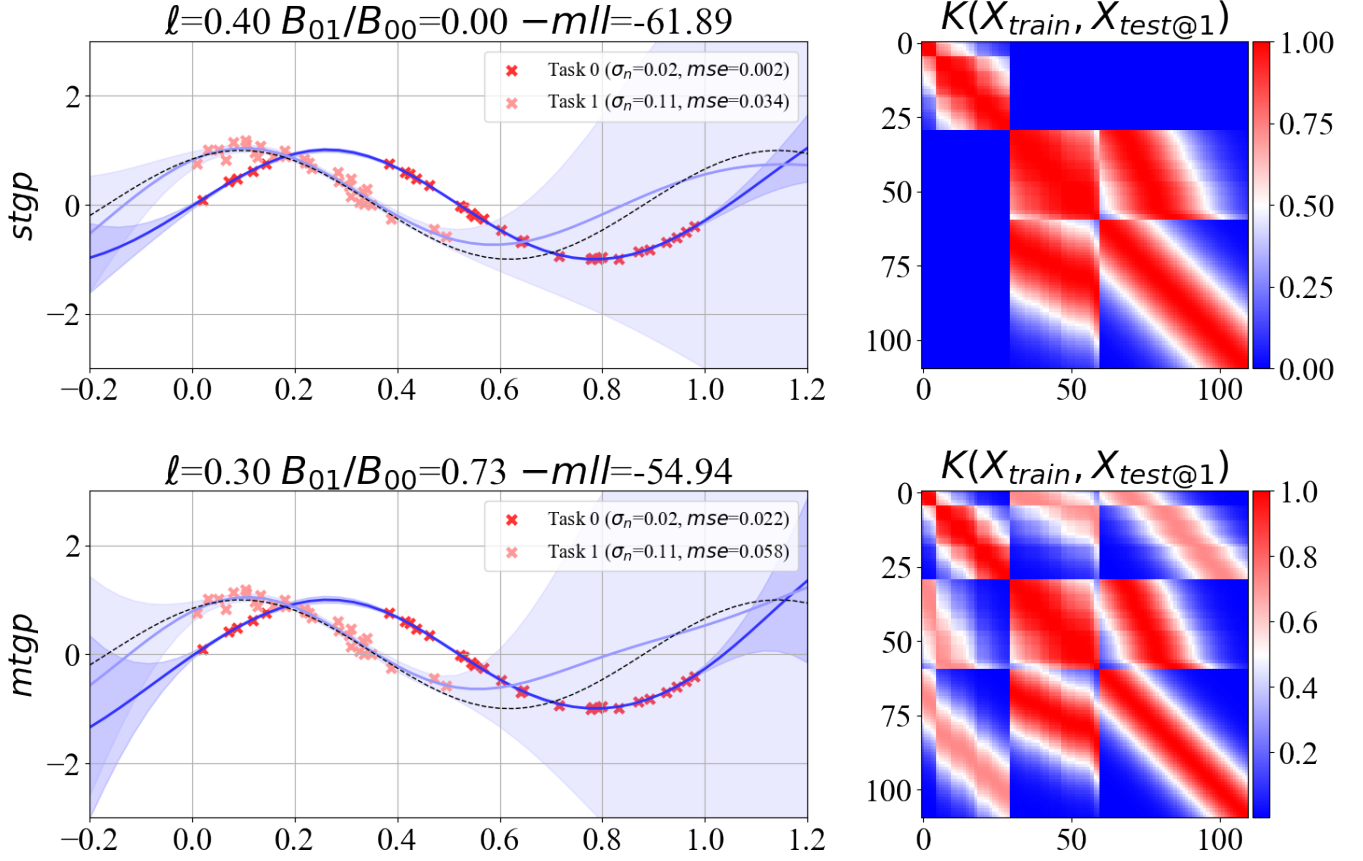


Figure 4: Observations for the two tasks are generated $y_0(x) = \sin(6x) + .03\epsilon$ and $y_1(x) = \sin(6x + \text{shift}) + .1\epsilon$ where $\epsilon \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ and $\text{shift} = 1.5$. Here for task 1 we only sample points in $[0, .5]$ while for task 0 we sample points in $[0, 1]$. There are 4 times training points for task 0 than task 1 (Top) Two Gaussian process regressor are learnt to fit data independently. (Bottom) A MTGP is fit to data whereby the task similarity matrix $B := K^t$ is learnt by maximizing marginal log likelihood. We see in case of data imbalance, which often happens in multitask learning, we arrive at an optimal marginal likelihood for all data points (a larger marginal likelihood), but renders prediction suboptimal for task 1 (a higher MSE)

task should be selected from potentially irrelevant side information. A simple thing one can do is to define a kernel that captures similarity between inputs via k_x and similarity of side information k_s , similar to the setup of MTGP. Let k be a tensor product $k_x \otimes k_s$ over $\mathcal{X} \times \mathcal{S}$,

$$k((x, s), (x', s')) = k_x(x, x')k_s(s, s') \quad (38)$$

Instead of task similarity matrix, correlation between any pair y, y' depends on value of $k_s(s, s')$. In the case where s describes a particular sub-population of our training dataset, it gives rise to a stratification dependent kernel. In health care domain, we might want to provide stratification dependent predictions given our side information of scanner, patient orientation, and existence of comorbidities. In this setting, automatic selection of relevance can be achieved via an Automatic Relevance Determination (ARD) kernel

$$k_s(s, s') = \sigma^2 \exp \left(-\frac{1}{2} \sum_{l=1}^M \frac{1}{\ell_l^2} (s_l - s'_l)^2 \right) \quad (39)$$

If ARD kernel is used, current formulation captures the prior that main task regressor is smooth with respect to some stratification of dataset - at test time the model will give similar prediction to those that appears to be likely in the same training stratification, as shown in Figure (7). Hyperparameters $\{\{\ell_l\}_{l=1}^M, \theta_{k_x}\}$ can be learnt by maximizing marginal likelihood.

One complication is missing side information during training, no information exists for some components of s . A more serious issue is at test time, we would want our parameterized function f to make prediction based on x_* only, absent of any side information. Naively this is not possible, since our regressor is supported over $\mathcal{X} \times \mathcal{S}$. A easy way to fix this is to pre-train a model with the goal of predicting side information from x - we first impute the missing side information, then proceed with both training and testing. A more principled way to do this is to train some GP model and use its predictive distribution over the side information as inputs for training and testing of f . This is related to GP with uncertain inputs and there is closed form solution for inference.

For a toy case where $D = 1, M = 1$, we see that the asymmetric MTGP described just now is able to utilize side information when they help Figure (5) and ignore side information when they are irrelevant Figure (6). A summary of model performance for main task regressor with varying side information relevance is in Table (1). A summary of how much **amtgp** depends on side information is in Table (2).

shift	stgp	mtgp	amtgp
0.000	0.579	0.021	0.181
0.100	0.660	0.030	0.106
0.200	0.735	0.026	0.088
0.300	0.802	0.027	0.162
0.400	0.857	0.041	0.246
0.500	0.898	0.070	0.256
1.000	0.837	0.139	0.296

Table 1: Mean squared error for main task (task=1) regressor at test locations along the grid (`linspace(0,1,100)`) for the three different methods **stgp**, **mtgp**, **amtgp** trained on data of varying task similarity. Larger values of **shift** implies the two task functions are more unrelated. **mtgp** is more flexible than **stgp** and therefore gives a higher marginal likelihood assuming the task similarity matrix can be found properly. **amtgp** is able to utilize side information when they are really relevant **shift** = 0, and does not suffer from performance drop when the two tasks are irrelevant **shift** = 1.0

shift	$\frac{\ell_s}{\ell_x}$
0.0	0.1
0.1	0.7
0.5	8.8
1.0	23.4
1.5	29.4

Table 2: Ratio of lengthscale for k_s over length scale for k_x , e.g. $\frac{\ell_s}{\ell_x}$ for main task (task=1) regressor trained via **amtgp**. A larger ratio implies the model relies more on x and less on side information s

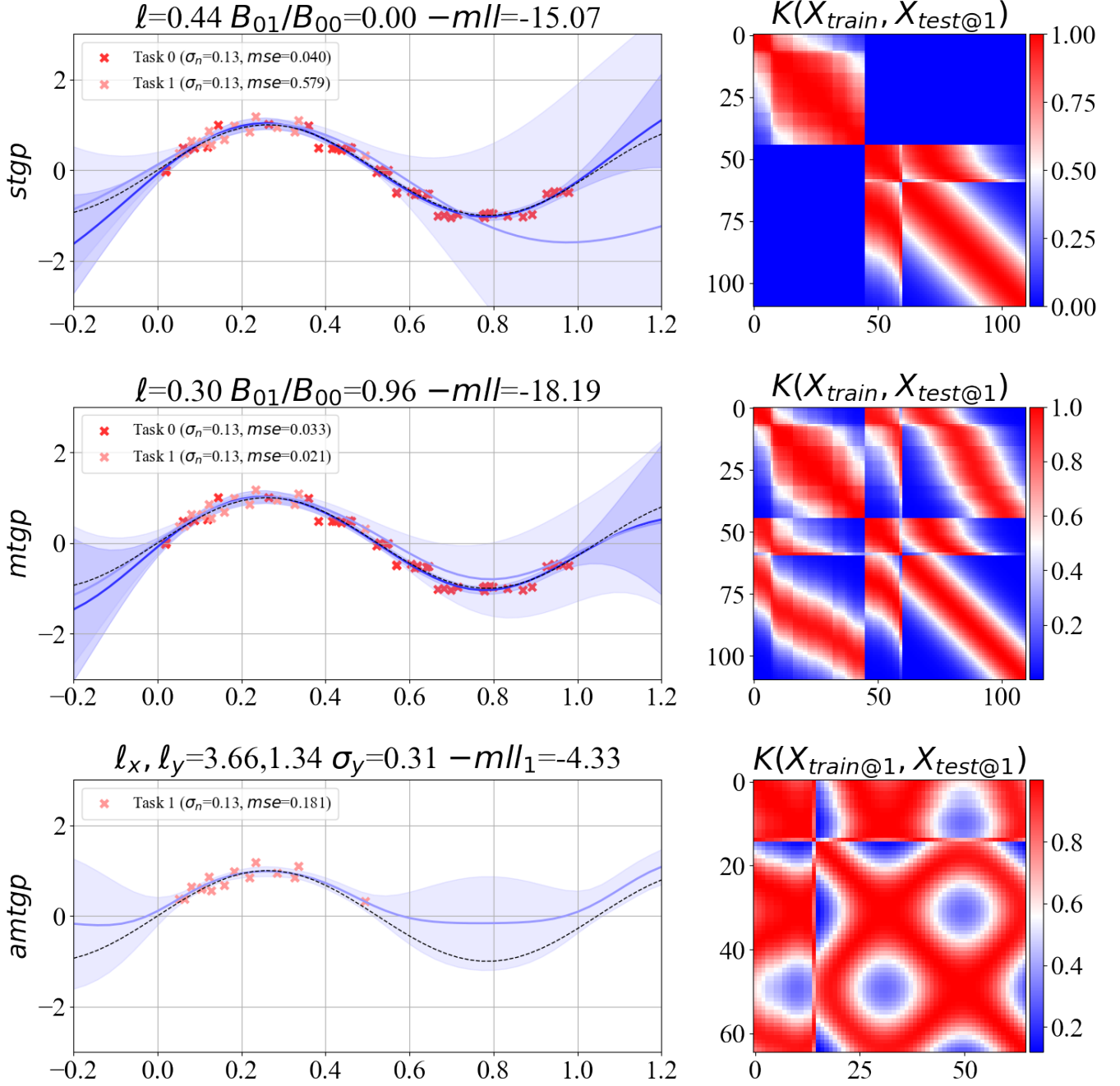


Figure 5: (Top) **stgp** GP regressors are learnt independently (Middle) **mtgp** GP regressors learnt via MTGP (Bottom) **amtgp** GP regressor for task 1 (main task) where the GP prior has kernel defined in Equation (38, 39). Missing side information is imputed from the predictive mean of a pre-trained **stgp** GP regressor for task 0. Bandwidth for the squared exponential kernel k_x and ARD kernel k_s , e.g. $\{\ell_x, \ell_s\}$, is optimized by maximizing the marginal likelihood of task 1 data points. We see, in case where the tasks are related, as measured by a large **shift** = 0. between these two task functions. We see **mtgp** is able to utilize side information really effectively while **amtgp** also learns to rely on side information, as seen from a small ℓ_s , and reduce its uncertainty of prediction far away from its training support.

References

- [1] Carl Edward Rasmussen and Williams Christopher. *Gaussian Process for Machine Learning*. MIT Press, 2006.
- [2] Joaquin Quiñonero-Candela and Carl Edward Rasmussen. “A Unifying View of Sparse Approximate Gaussian Process Regression”. In: *The Journal of Machine Learning Research* 6 (Dec. 1, 2005), pp. 1939–1959. ISSN: 1532-4435.
- [3] Edward Snelson and Zoubin Ghahramani. “Sparse Gaussian Processes using Pseudo-inputs”. In: *Advances in Neural Information Processing Systems* 18 (2005), pp. 1257–1264. URL: <https://proceedings.neurips.cc/paper/2005/hash/4491777b1aa8b5b32c2e8666dbe1a495-Abstract.html> (visited on 02/05/2021).
- [4] Michalis Titsias. “Variational Learning of Inducing Variables in Sparse Gaussian Processes”. In: *Artificial Intelligence and Statistics*. Artificial Intelligence and Statistics. ISSN: 1938-7228. PMLR, Apr. 15, 2009, pp. 567–574. URL: <http://proceedings.mlr.press/v5/titsias09a.html> (visited on 02/21/2021).
- [5] James Hensman, Nicolo Fusi, and Neil D. Lawrence. “Gaussian Processes for Big Data”. In: *arXiv:1309.6835 [cs, stat]* (Sept. 26, 2013). arXiv: [1309.6835](https://arxiv.org/abs/1309.6835). URL: <http://arxiv.org/abs/1309.6835> (visited on 03/07/2021).
- [6] James Hensman, Alex Matthews, and Zoubin Ghahramani. “Scalable Variational Gaussian Process Classification”. In: *arXiv:1411.2005 [stat]* (Nov. 7, 2014). arXiv: [1411.2005](https://arxiv.org/abs/1411.2005). URL: <http://arxiv.org/abs/1411.2005> (visited on 02/23/2021).
- [7] Michalis K Titsias. “Variational Model Selection for Sparse Gaussian Process Regression”. In: (2009), p. 20.
- [8] Yarin Gal and Mark van der Wilk. “Variational Inference in Sparse Gaussian Process Regression and Latent Variable Models - a Gentle Tutorial”. In: *arXiv:1402.1412 [stat]* (Sept. 29, 2014). arXiv: [1402.1412](https://arxiv.org/abs/1402.1412). URL: <http://arxiv.org/abs/1402.1412> (visited on 03/08/2021).
- [9] Thang Bui and Richard Turner. “Sparse Approximations for Non-Conjugate Gaussian Process Regression”. In: (2014), p. 7. URL: https://thangbui.github.io/docs/reports/tr_sparseNonConj.pdf.
- [10] “Sparse Gaussian Process Approximations and Applications”. In: (2018), p. 188. URL: <https://markvdw.github.io/vanderwilk-thesis.pdf>.
- [11] Edwin V Bonilla, Kian Ming A Chai, and Christopher K I Williams. “Multi-task Gaussian Process Prediction”. In: (2008), p. 8.
- [12] Gayle Leen, Jaakko Peltonen, and Samuel Kaski. “Focused multi-task learning in a Gaussian process framework”. In: *Machine Learning* 89.1 (Oct. 1, 2012), pp. 157–182. ISSN: 1573-0565. DOI: [10.1007/s10994-012-5302-y](https://doi.org/10.1007/s10994-012-5302-y). URL: <https://doi.org/10.1007/s10994-012-5302-y> (visited on 02/15/2021).

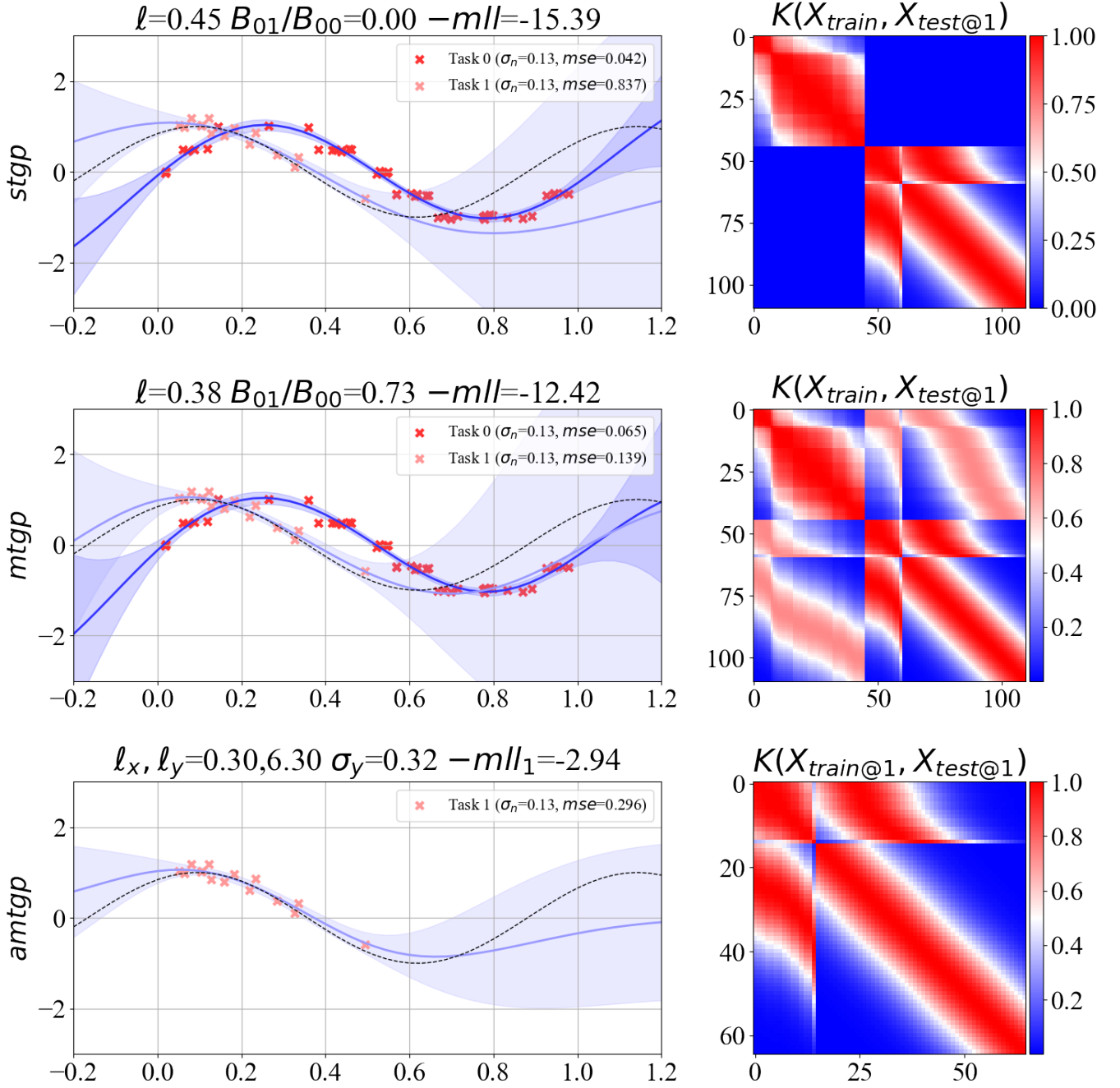


Figure 6: (Top) **stgp** GP regressors are learnt independently (Middle) **mtgp** GP regressors learnt via MTGP (Bottom) **amtgp** GP regressor for task 1 (main task) where the GP prior has kernel defined in Equation (38, 39). Missing side information is imputed from the predictive mean of a pre-trained **stgp** GP regressor for task 0. Bandwidth for the squared exponential kernel k_x and ARD kernel k_s , e.g. $\{\ell_x, \ell_s\}$, is optimized by maximizing the marginal likelihood of task 1 data points. We see, in case where the tasks are unrelated and perhaps even conflicting, as measured by a large $\text{shift} = 1.0$ between these two task functions, we see **mtgp** prediction is more uncertain in regions where task 1 function is supported and task 0 function is not. This is a result of **mtgp** able to use data points from task 0 observations even if the two functions are unrelated. **amtgp** learns to ignore side information as seen from large values of ℓ_s , and reduces to **stgp** regression

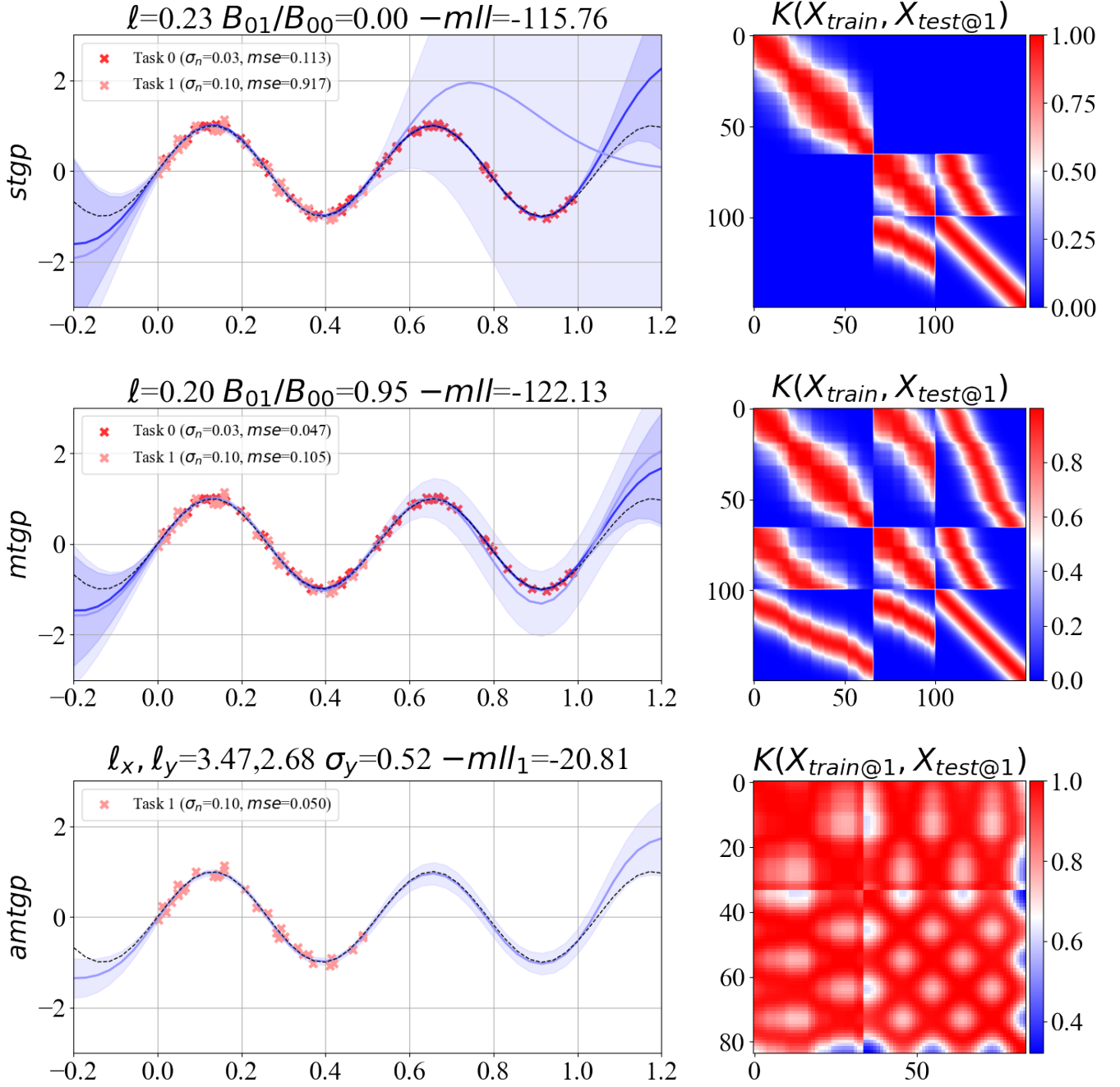


Figure 7: (Top) **stgp** GP regressors are learnt independently (Middle) **mtgp** GP regressors learnt via MTGP (Bottom) **amtgp** GP regressor for task 1 (main task) where the GP prior has kernel defined in Equation (38, 39). We see **amtgp** can extrapolate training support by utilizing prediction of side task.