

# Supplementary material for ‘Regularized Multi-output Gaussian Convolution Process with Domain Adaptation’

Xinming Wang, Chao Wang, Xuan Song, Levi Kirby, Jianguo Wu



## APPENDIX A

### DERIVATION OF COVARIANCE FUNCTION IN CONVOLUTION PROCESS

For the convolution process:

$$f_i(\mathbf{x}) = g_i(\mathbf{x}) * Z(\mathbf{x}) = \int_{-\infty}^{\infty} g_i(\mathbf{x} - \mathbf{u})Z(\mathbf{u})d\mathbf{u},$$

If  $Z(\mathbf{x})$  is a commonly used white Gaussian noise process, i.e.,  $\text{cov}(Z(\mathbf{x}), Z(\mathbf{x}')) = \delta(\mathbf{x} - \mathbf{x}')$  and  $\mathbb{E}(Z(\mathbf{x})) = 0$ , then the cross covariance is derived as:

$$\begin{aligned} \text{cov}_{ij}^f(\mathbf{x}, \mathbf{x}') &= \text{cov}\{g_i(\mathbf{x}) * Z(\mathbf{x}), g_j(\mathbf{x}') * Z(\mathbf{x}')\} \\ &= \mathbb{E}\left\{\int_{-\infty}^{\infty} g_i(\mathbf{x} - \mathbf{u})Z(\mathbf{u})d\mathbf{u} \int_{-\infty}^{\infty} g_j(\mathbf{x}' - \mathbf{u}')Z(\mathbf{u}')d\mathbf{u}'\right\} \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g_i(\mathbf{u})g_j(\mathbf{u}')\mathbb{E}\{Z(\mathbf{x} - \mathbf{u})Z(\mathbf{x}' - \mathbf{u}')\}d\mathbf{u}d\mathbf{u}' \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g_i(\mathbf{u})g_j(\mathbf{u}')\delta(\mathbf{x} - \mathbf{u} - \mathbf{x}' + \mathbf{u}')d\mathbf{u}d\mathbf{u}' \\ &= \int_{-\infty}^{\infty} g_i(\mathbf{u})g_j(\mathbf{u} - \mathbf{v})d\mathbf{u}, \end{aligned} \quad (1)$$

where  $\mathbf{v} = \mathbf{x} - \mathbf{x}'$  and the last equality is based on the property of Dirac function that  $\int g(\mathbf{u}')\delta(\mathbf{u}' - \mathbf{x})d\mathbf{u}' = g(\mathbf{x})$ .

For our MGCP structure:

$$\begin{aligned} y_i(\mathbf{x}) &= f_i(\mathbf{x}) + \epsilon_i(\mathbf{x}) = g_{ii}(\mathbf{x}) * Z_i(\mathbf{x}) + \epsilon_i(\mathbf{x}), i \in \mathcal{I}^S \\ y_t(\mathbf{x}) &= f_t(\mathbf{x}) + \epsilon_t(\mathbf{x}) = \sum_{j \in \mathcal{I}} g_{jt}(\mathbf{x}) * Z_j(\mathbf{x}) + \epsilon_t(\mathbf{x}), \end{aligned}$$

the source-target covariance function can be calculated as:

$$\begin{aligned} \text{cov}_{it}^f(\mathbf{x}, \mathbf{x}') &= \text{cov}(f_i(\mathbf{x}), f_t(\mathbf{x}')) \\ &= \text{cov}\left\{g_{ii}(\mathbf{x}) * Z_i(\mathbf{x}), \sum_{j \in \mathcal{I}} g_{jt}(\mathbf{x}') * Z_j(\mathbf{x}')\right\} \\ &= \sum_{j \in \mathcal{I}} \text{cov}\{g_{ii}(\mathbf{x}) * Z_i(\mathbf{x}), g_{jt}(\mathbf{x}') * Z_j(\mathbf{x}')\} \\ &= \int_{-\infty}^{\infty} g_{ii}(\mathbf{u})g_{it}(\mathbf{u} - \mathbf{v})d\mathbf{u}, \quad i \in \mathcal{I}^S \end{aligned} \quad (2)$$

where  $\mathbf{v} = \mathbf{x} - \mathbf{x}'$ . In the same way, we can derive the auto-covariance as

$$\begin{aligned} \text{cov}_{ii}^f(\mathbf{x}, \mathbf{x}') &= \int_{-\infty}^{\infty} g_{ii}(\mathbf{u})g_{ii}(\mathbf{u} - \mathbf{v})d\mathbf{u}, i \in \mathcal{I}^S \\ \text{cov}_{tt}^f(\mathbf{x}, \mathbf{x}') &= \sum_{j \in \mathcal{I}} \int_{-\infty}^{\infty} g_{jj}(\mathbf{u})g_{jt}(\mathbf{u} - \mathbf{v})d\mathbf{u}. \end{aligned}$$

## APPENDIX B

### PROOF OF THEOREM 1

Suppose that  $g_{it}(\mathbf{x}) = 0, \forall i \in \mathcal{U} \subseteq \mathcal{I}^S$  for all  $\mathbf{x} \in \mathcal{X}$ . For notational convenience, suppose  $\mathcal{U} = \{1, 2, \dots, h|h \leq q\}$ , then the predictive distribution of the model at any new input  $\mathbf{x}_*$  is unrelated with  $\{f_1, f_2, \dots, f_h\}$  and is reduced to:

$$\begin{aligned} p(y_t(\mathbf{x}_*)|\mathbf{y}) &= \mathcal{N}(\mathbf{k}_+^T \mathbf{C}_+^{-1} \mathbf{y}_+, \\ &\quad \text{cov}_{tt}^f(\mathbf{x}_*, \mathbf{x}_*) + \sigma_t^2 - \mathbf{k}_+^T \mathbf{C}_+^{-1} \mathbf{k}_+), \end{aligned}$$

where  $\mathbf{k}_+ = (\mathbf{K}_{h+1,*}^T, \dots, \mathbf{K}_{q,*}^T, \mathbf{K}_{t,*}^T)^T$ ,  $\mathbf{y}_+ = (\mathbf{y}_{h+1}^T, \dots, \mathbf{y}_q^T, \mathbf{y}_t^T)^T$ , and

$$\mathbf{C}_+ = \begin{pmatrix} \mathbf{C}_{h+1,h+1} & \cdots & \mathbf{0} & \mathbf{C}_{h+1,t} \\ \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \cdots & \mathbf{C}_{q,q} & \mathbf{C}_{q,t} \\ \mathbf{C}_{h+1,t}^T & \cdots & \mathbf{C}_{q,t}^T & \mathbf{C}_{t,t} \end{pmatrix}.$$

**Proof.** Recall that

$$\begin{aligned} \text{cov}_{jt}^y(\mathbf{x}, \mathbf{x}') &= \text{cov}_{jt}^f(\mathbf{x}, \mathbf{x}') \\ &= \int_{-\infty}^{\infty} g_{jj}(\mathbf{u})g_{jt}(\mathbf{u} - \mathbf{v})d\mathbf{u}, \\ \text{cov}_{tt}^y(\mathbf{x}, \mathbf{x}') &= \text{cov}_{tt}^f(\mathbf{x}, \mathbf{x}') + \sigma_t^2 \delta(\mathbf{x} - \mathbf{x}') \\ &= \sum_{h \in \mathcal{I}} \int_{-\infty}^{\infty} g_{hh}(\mathbf{u})g_{ht}(\mathbf{u} - \mathbf{v})d\mathbf{u} + \sigma_t^2 \delta(\mathbf{x} - \mathbf{x}'), \end{aligned}$$

for all  $j \in \{1, 2, \dots, q\}$ , so  $g_{it}(\mathbf{x}) = 0, i \in \{1, 2, \dots, h|h \leq q\}$  implies that  $\text{cov}_{it}^y(\mathbf{x}, \mathbf{x}') = 0$  for all  $i \in \{1, 2, \dots, h\}$  and

$$\text{cov}_{tt}^y(\mathbf{x}, \mathbf{x}') = \sum_{i=h+1}^t \int_{-\infty}^{\infty} g_{ii}(\mathbf{u})g_{it}(\mathbf{u} - \mathbf{v})d\mathbf{u} + \sigma_t^2 \delta(\mathbf{x} - \mathbf{x}').$$

Therefore, we have that  $\mathbf{C}_{i,t} = 0, i \in \{1, 2, \dots, h\}$  and partition covariance matrix  $\mathbf{C} = \begin{pmatrix} \mathbf{C}_- & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_+ \end{pmatrix}$ , where  $\mathbf{C}_- =$

$$\begin{pmatrix} \mathbf{C}_{1,1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{2,2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{C}_{h,h} \end{pmatrix}.$$

The predictive distribution at point  $\mathbf{x}_*$  is

$$y_t(\mathbf{x}_*) \sim \mathcal{N}(\mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{y}, \text{cov}_{tt}^f(\mathbf{x}_*, \mathbf{x}_*) + \sigma_t^2 - \mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{K}_*).$$

Also, based on that  $\text{cov}_{it}^y(\mathbf{x}, \mathbf{x}') = 0$  for all  $i \in \{1, 2, \dots, h\}$ , we have that  $\mathbf{K}_* = (\mathbf{0}, \mathbf{k}_+^T)^T$ . Let  $\mathbf{y}_- = (\mathbf{y}_1^T, \dots, \mathbf{y}_h^T)^T$ , then  $\mathbf{y} = (\mathbf{y}_-^T, \mathbf{y}_+^T)^T$ . Therefore,

$$\begin{aligned} \mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{y} &= (\mathbf{0}, \mathbf{k}_+^T) \begin{pmatrix} \mathbf{C}_- & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_+ \end{pmatrix}^{-1} (\mathbf{y}_-^T, \mathbf{y}_+^T)^T \\ &= (\mathbf{0}, \mathbf{k}_+^T) \begin{pmatrix} \mathbf{C}_-^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_+^{-1} \end{pmatrix} (\mathbf{y}_-^T, \mathbf{y}_+^T)^T \\ &= \mathbf{k}_+^T \mathbf{C}_+^{-1} \mathbf{y}_+, \\ \mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{K}_* &= (\mathbf{0}, \mathbf{k}_+^T) \begin{pmatrix} \mathbf{C}_- & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_+ \end{pmatrix}^{-1} (\mathbf{0}, \mathbf{k}_+^T)^T \\ &= \mathbf{k}_+^T \mathbf{C}_+^{-1} \mathbf{k}_+. \end{aligned}$$

Note that the auto-covariance matrix of target output  $f_t$ ,  $\mathbf{C}_{tt}$ , is also unrelated with observed data  $\{\mathbf{X}_i | i = 1, 2, \dots, h\}$  which from source output  $\{f_i | i = 1, 2, \dots, h\}$ . As a result, the predictive distribution is totally independent on these outputs. Proof completes.

## APPENDIX C REGULARITY CONDITIONS

In this part, we state the regularity conditions for the consistency theorem of the MLE  $\hat{\boldsymbol{\theta}}_{\#}$ , which are formulated in [34].

Denote  $\mathbf{y}$  with total  $N$  observations as  $\mathbf{y}^N$ , and let

$$p_k(\boldsymbol{\theta}) = \frac{p(\mathbf{y}^k | \boldsymbol{\theta})}{p(\mathbf{y}^{k-1} | \boldsymbol{\theta})}$$

for each  $k$ . Assume  $p_k(\boldsymbol{\theta})$  is twice differentiable with respect to  $\boldsymbol{\theta}$  in a neighborhood of  $\boldsymbol{\theta}^*$ . Also assume that the support of  $p(\mathbf{y}^N | \boldsymbol{\theta})$  is independent of  $\boldsymbol{\theta}$  in the neighborhood. Define  $\phi_k(\boldsymbol{\theta}) = \log p_k(\boldsymbol{\theta})$ , and its first derivative  $\phi_k'(\boldsymbol{\theta})$ , second derivative  $\phi_k''(\boldsymbol{\theta})$ .

For simplicity and without loss of generality, we only consider the conditions for one-dimensional case. Define  $\phi_k^{*'} = \phi_k'(\boldsymbol{\theta}^*)$  and  $\phi_k^{*''} = \phi_k''(\boldsymbol{\theta}^*)$ . Let  $\mathcal{F}_N$  be the  $\sigma$ -field generated by  $y_j, 1 \leq j \leq N$ , and  $\mathcal{F}_0$  be the trivial  $\sigma$ -field. Define the random variable  $i_k^* = \text{var}(\phi_k^{*'} | \mathcal{F}_{k-1}) = \mathbb{E}[(\phi_k^{*'})^2 | \mathcal{F}_{k-1}]$  and  $I_N^* = \sum_{k=1}^N i_k^*$ . Define  $S_N = \sum_{k=1}^N \phi_k^{*'}$  and  $S_N^* = \sum_{k=1}^N \phi_k^{*''} + I_N^*$ . If the following conditions hold:

- (c1)  $\phi_k(\boldsymbol{\theta})$  is thrice differentiable in the neighborhood of  $\boldsymbol{\theta}^*$ . Let  $\phi_k^{*'''} = \phi_k'''(\boldsymbol{\theta}^*)$  be the third derivative,
- (c2) Twice differentiation of  $\int p(\mathbf{y}^N | \boldsymbol{\theta}) d\mu^N(\mathbf{y}^N)$  with respect to  $\boldsymbol{\theta}$  exists in the neighborhood of  $\boldsymbol{\theta}^*$ ,
- (c3)  $\mathbb{E}|\phi_k^{*''}| < \infty$  and  $\mathbb{E}|\phi_k^{*'''} + (\phi_k^{*'})^2| < \infty$ .
- (c4) There exists a sequence of constants  $K(N) \rightarrow \infty$  as  $N \rightarrow \infty$  such that:

- (i)  $K(N)^{-1} S_N \xrightarrow{P} 0$ ,
- (ii)  $K(N)^{-1} S_N^* \xrightarrow{P} 0$ ,
- (iii) there exists  $a(\boldsymbol{\theta}^*) > 0$  such that  $\forall \epsilon > 0$ ,  $P[K(N)^{-1} I_N^* \geq 2a(\boldsymbol{\theta}^*)] \geq 1 - \epsilon$  for all  $N \geq N(\epsilon)$ ,
- (iv)  $K(N)^{-1} \sum_{k=1}^N \mathbb{E}|\phi_k^{*'''}| < M < \infty$  for all  $N$ ,

then the MLE  $\hat{\boldsymbol{\theta}}_{\#}$  is consistent for  $\boldsymbol{\theta}^*$ . There exists a sequence  $r_N$  such that  $r_N \rightarrow \infty$  as  $N \rightarrow \infty$ , i.e.,

$$\|\hat{\boldsymbol{\theta}}_{\#} - \boldsymbol{\theta}^*\| = O_P(r_N^{-1}).$$

## APPENDIX D PROOF OF THEOREM 2

Suppose that the MLE for  $L(\boldsymbol{\theta} | \mathbf{y})$ ,  $\hat{\boldsymbol{\theta}}_{\#}$ , is  $r_N$  consistent, i.e., satisfying Eq. (19). If  $\max\{|\mathbb{P}_{\gamma}''(\boldsymbol{\theta}_{i0}^*)| : \boldsymbol{\theta}_{i0}^* \neq \mathbf{0}\} \rightarrow 0$ , then there exists a local maximizer  $\hat{\boldsymbol{\theta}}$  of  $L_{\mathbb{P}}(\boldsymbol{\theta} | \mathbf{y})$  s.t.  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\| = O_P(r_N^{-1} + r_0)$ , where  $r_0 = \max\{|\mathbb{P}_{\gamma}'(\boldsymbol{\theta}_{i0}^*)| : \boldsymbol{\theta}_{i0}^* \neq \mathbf{0}\}$ .

**Proof.** Recall the assumptions in Section 3.2. For the unpenalized log-likelihood  $L(\boldsymbol{\theta})$ , the MLE  $\hat{\boldsymbol{\theta}}_{\#}$  is  $r_N$  consistent where  $r_N$  is a sequence such that  $r_N \rightarrow \infty$  as  $N \rightarrow \infty$ . And we have that  $L'(\boldsymbol{\theta}^*) = O_P(r_N)$  and  $\mathbf{I}_N(\boldsymbol{\theta}^*) = O_P(r_N^2)$ , which are the standard argument based on the consistency of estimator. Based on that, we aim to study the asymptotic properties of the penalized likelihood  $L_{\mathbb{P}}(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - r_N^2 \mathbb{P}_{\gamma}(\boldsymbol{\theta}_0)$ . Here we multiply the penalty function by  $r_N^2$  to avoid that penalty term degenerates as  $N \rightarrow \infty$ . The following proof is similar to that of Fan and Li [33] but based on dependent observations.

To prove theorem 2, we need to show that for any given  $\epsilon > 0$ , there exists a large constant  $U$  such that:

$$P \left\{ \sup_{\|\mathbf{u}\|=U} L_{\mathbb{P}}(\boldsymbol{\theta}^* + r_N^+ \mathbf{u}) < L_{\mathbb{P}}(\boldsymbol{\theta}^*) \right\} \geq 1 - \epsilon, \quad (3)$$

where  $r_N^+ = r_N^{-1} + r_0$ . This implies that with probability at least  $1 - \epsilon$  there exists a local maximum in the ball  $\{\boldsymbol{\theta}^* + r_N^+ \mathbf{u} : \|\mathbf{u}\| \leq U\}$ . So the local maximizer  $\hat{\boldsymbol{\theta}}$  satisfies that  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\| = O_P(r_N^+)$ .

By  $\mathbb{P}_{\gamma}(0) = 0$ , we have

$$\begin{aligned} &L_{\mathbb{P}}(\boldsymbol{\theta}^* + r_N^+ \mathbf{u}) - L_{\mathbb{P}}(\boldsymbol{\theta}^*) \\ &\leq L(\boldsymbol{\theta}^* + r_N^+ \mathbf{u}) - L(\boldsymbol{\theta}^*) \\ &\quad - r_N^2 \sum_{i=h+1}^q [\mathbb{P}_{\gamma}(|\boldsymbol{\theta}_{i0}^* + r_N^+ u_{i0}|) - \mathbb{P}_{\gamma}(|\boldsymbol{\theta}_{i0}^*|)], \end{aligned}$$

where  $h$  and  $q$  are the number of zero components and all components in  $\boldsymbol{\theta}_{i0}^*$ , and  $u_{i0}$  is the element corresponding to  $\boldsymbol{\theta}_{i0}$  in  $\mathbf{u}$ . Let  $\mathbf{I}_N(\boldsymbol{\theta}^*)$  be the finite and positive definite information matrix at  $\boldsymbol{\theta}^*$  with  $N$  observations. Applying a Taylor expansion on the likelihood function, we have that

$$\begin{aligned} &L_{\mathbb{P}}(\boldsymbol{\theta}^* + r_N^+ \mathbf{u}) - L_{\mathbb{P}}(\boldsymbol{\theta}^*) \\ &\leq r_N^+ L'(\boldsymbol{\theta}^*)^T \mathbf{u} - \frac{1}{2} (r_N^+)^2 \mathbf{u}^T \mathbf{I}_N(\boldsymbol{\theta}^*) \mathbf{u} [1 + o_P(1)] \\ &\quad - r_N^2 \sum_{i=h+1}^q \left\{ r_N^+ \mathbb{P}_{\gamma}'(|\boldsymbol{\theta}_{i0}^*|) \text{sign}(\boldsymbol{\theta}_{i0}^*) u_{i0} \right. \\ &\quad \left. + \frac{1}{2} (r_N^+)^2 \mathbb{P}_{\gamma}''(|\boldsymbol{\theta}_{i0}^*|) u_{i0}^2 [1 + o_P(1)] \right\}, \quad (4) \end{aligned}$$

Note that  $\|L'(\boldsymbol{\theta}^*)\| = O_P(r_N)$  and  $\mathbf{I}_N(\boldsymbol{\theta}^*) = O_P(r_N^2)$ . so the first term on the right-hand side of Eq. (4) is on the order  $O_P(r_N^+ r_N)$ , while the second term is  $O_P((r_N^+ r_N)^2)$ . By choosing a sufficient large  $U$ , the first term can be dominated by the second term uniformly in  $\|\mathbf{u}\| = U$ . Besides, the absolute value of the third term is bounded by

$$\sqrt{q - h} r_N^2 r_N^+ r_0 \|\mathbf{u}\| + (r_N r_N^+)^2 \max\{|\mathbb{P}_{\gamma}''(\boldsymbol{\theta}_{i0}^*)| : \boldsymbol{\theta}_{i0}^* \neq \mathbf{0}\} \|\mathbf{u}\|^2,$$

which is also dominated by second term as it is on the order of  $o_P((r_N r_N^+)^2)$ . Thus, Eq. (3) holds and the proof completes.

## APPENDIX E PROOF OF THEOREM 3

Let  $\theta_{10}^*$  and  $\theta_{20}^*$  contain the zero and non-zero components in  $\theta_0^*$  respectively. Assume the conditions in Theorem 2 also hold, and  $\hat{\theta}$  is  $r_N$  consistent by choosing proper  $\gamma$  in  $\mathbb{P}_\gamma(\theta_0)$ . If  $\liminf_{N \rightarrow \infty} \liminf_{\theta \rightarrow 0^+} \gamma^{-1} \mathbb{P}'_\gamma(\theta) > 0$  and  $(r_N \gamma)^{-1} \rightarrow 0$ , then

$$\lim_{N \rightarrow \infty} P(\hat{\theta}_{10} = \mathbf{0}) = 1.$$

**Proof.** To prove this theorem, we only need to prove that for a small  $\epsilon_N = U r_N$ , where  $U$  is a given constant and  $i = 1, \dots, s$ ,

$$\frac{\partial L_{\mathbb{P}}(\theta)}{\partial \theta_{i0}} \theta_{i0} < 0, 0 < |\theta_{i0}| < \epsilon_N. \quad (5)$$

By Taylor's expansion,

$$\begin{aligned} \frac{\partial L_{\mathbb{P}}(\theta)}{\partial \theta_{i0}} &= \frac{\partial L(\theta)}{\partial \theta_{i0}} - r_N^2 \mathbb{P}'_\gamma(|\theta_{i0}|) \text{sign}(\theta_{i0}) \\ &= \frac{\partial L(\theta^*)}{\partial \theta_{i0}} + \left[ \partial \left( \frac{\partial L(\theta^*)}{\partial \theta_{i0}} \right) / \partial \theta \right]^T (\theta - \theta^*) [1 + o_P(1)] \\ &\quad - r_N^2 \mathbb{P}'_\gamma(|\theta_{i0}|) \text{sign}(\theta_{i0}). \end{aligned}$$

As  $\frac{\partial L(\theta)}{\partial \theta_{i0}} = O_P(r_N)$ ,  $\partial \left( \frac{\partial L(\theta^*)}{\partial \theta_{i0}} \right) / \partial \theta_j = O_P(r_N^2)$  by the standard argument for  $r_N$  consistent estimator, thus

$$\begin{aligned} \frac{\partial L_{\mathbb{P}}(\theta)}{\partial \theta_{i0}} &= O_P(r_N) - r_N^2 \mathbb{P}'_\gamma(|\theta_{i0}|) \text{sign}(\theta_{i0}) \\ &= r_N^2 \gamma \left( O_P\left(\frac{1}{r_N \gamma}\right) - \gamma^{-1} \mathbb{P}'_\gamma(|\theta_{i0}|) \text{sign}(\theta_{i0}) \right). \end{aligned}$$

Because that  $\liminf_{N \rightarrow \infty} \liminf_{\theta \rightarrow 0^+} \gamma^{-1} \mathbb{P}'_\gamma(\theta) > 0$  and  $(r_N \gamma)^{-1} \rightarrow 0$ ,  $\frac{\partial L_{\mathbb{P}}(\theta)}{\partial \theta_{i0}}$  will be positive while  $\theta_{i0}$  is negative and vice versa. As a result, Eq. (5) follows. Proof completes.

## APPENDIX F INTERPRETATION OF THE BENCHMARK: MGCP-RF

The illustration of MGCP-RF is shown in Fig. 1.

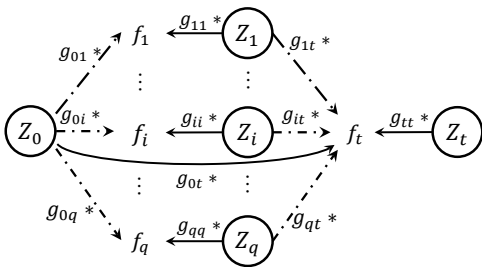


Fig. 1. The structure of MGCP-RF

In this structure, target  $f_t$  is generated by three kinds of latent process:  $Z_0(x)$ ,  $\{Z_i(x)\}_{i=1}^q$  and  $Z_t(x)$ . As  $Z_0(x)$  is the common process shared by sources, the covariance matrix blocks between source  $f_i$  and the other outputs are zero only when the scale parameters in  $g_{0i}(x)$  and  $g_{it}(x)$

are zero simultaneously. Thus, the marginalized covariance matrix  $C_+$  in Theorem 1 will be:

$$C_+ = \begin{pmatrix} C_{h+1,h+1} & \cdots & C_{h+1,q} & C_{h+1,t} \\ \vdots & \ddots & \vdots & \vdots \\ C_{h+1,q}^T & \cdots & C_{q,q} & C_{q,t} \\ C_{h+1,t}^T & \cdots & C_{q,t}^T & C_{t,t} \end{pmatrix}.$$

The difference to MGCP-R is that covariance among the remaining sources  $\{f_i\}_{i=h+1}^q$  can be modeled. This structure is indeed more comprehensive but with the cost of a half more parameters than MGCP-R. The cost will increase if we use more latent process to model the correlation among sources.

To realize the effect of shrinking  $g_{0i}(x)$  and  $g_{it}(x)$  at the same time, group-L1 penalty is used and the penalized log-likelihood function is:

$$\max_{\theta} L_{\mathbb{P}}(\theta|\mathbf{y}) = L(\theta|\mathbf{y}) - \gamma \sum_{i=1}^q \sqrt{\alpha_{0i}^2 + \alpha_{it}^2},$$

## APPENDIX G INFLUENCE OF TUNING-PARAMETER

To test the influence of the tuning-parameter  $\gamma$  in our model, we conduct the following experiment. Based on the same dataset in the 1D example of simulation case I, we construct MGCP-R model only with sources  $f_1$  and  $f_2$ , and let  $\gamma$  vary from 0 to 10 at a step of 1. Note that MGCP-T is equal to the model with  $\gamma = 0$ . The boxplot of MAE with respect to different values of  $\gamma$  is shown in Fig. 2. The estimated value of  $\alpha_{1t}$ ,  $\alpha_{2t}$  in one repetition is presented in Fig. 3. It can be seen that as  $\gamma$  increases, source  $f_2$  will be excluded from the prediction of target, leading to an increased prediction error. In practice, cross-validation can be used to select an optimal tuning-parameter.

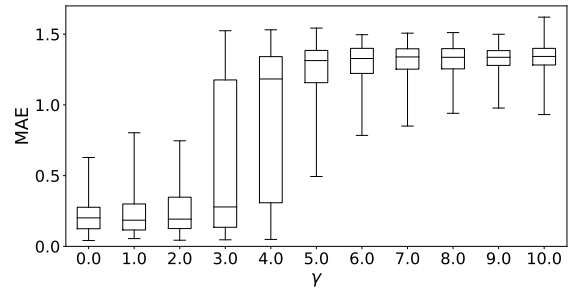


Fig. 2. Prediction error with different  $\gamma$  in 100 repetition.

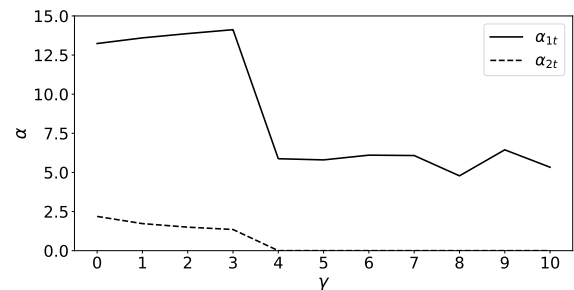


Fig. 3. Estimated values of  $\alpha_{1t}$ ,  $\alpha_{2t}$  in one repetition.