

# Bayesian transfer learning between Gaussian process regression tasks

Milan Papež<sup>a</sup> and Anthony Quinn<sup>a,b</sup>

<sup>a</sup> Institute of Information Theory and Automation, Czech Academy of Sciences, Czech Republic

<sup>b</sup> Department of Electronic and Electrical Engineering, Trinity College Dublin, the University of Dublin, Ireland

**Abstract**—Bayesian knowledge transfer in supervised learning scenarios often relies on a complete specification and optimization of the stochastic dependence between source and target tasks. This is a critical requirement of completely modelled settings, which can often be difficult to justify. We propose a strategy to overcome this. The methodology relies on fully probabilistic design to develop a target algorithm which accepts source knowledge in the form of a probability distribution. We present this incompletely modelled setting in the supervised learning context where the source and target tasks are to perform Gaussian process regression. Experimental evaluation demonstrates that the transfer of the source distribution substantially improves prediction performance of the target learner when recovering a distorted nonparametric function realization from noisy data.

**Index Terms**—Bayesian transfer learning, supervised learning, fully probabilistic design, incomplete modelling, Gaussian process regression.

## I. INTRODUCTION

Transfer learning [1]—also known as multi-task learning [2]—has become a fundamental part of machine learning [3] and artificial intelligence [4]. The root principle behind transfer learning is to use knowledge provided by one (or more) source task(s) in order to improve performance and generalization capabilities of one (or more) related target task(s). In this paper, we restrict to Bayesian transfer learning [5] and propose an algorithm for transferring knowledge between supervised learning tasks. The central aim of supervised learning is to learn behaviour of a latent function from labelled pairs of input-output data, and then make predictions about its future values based on unlabelled input data [6], [7]. We specify our framework in the case where source and target tasks rely on a Gaussian process (GP) [8] to model their latent functions.

Bayesian transfer learning with GP priors has been successfully applied in diverse areas, including Bayesian optimization [9], medical time-series analysis [10], terrain modelling [11], and robot pose estimation [12]. The central question of transfer learning is how to specify the dependence structure between the tasks. GP-based approaches define a GP prior for each of the tasks and then model their correlations based on a linear dependence structure with fixed coefficients [13]–[15] or varying, input-dependent, coefficients [16], [17], which can be generalized under the framework of convolving processes [18]. To capture a nonlinear dependence structure, and foster a richer representation of relationships between the tasks, an additional GP layer can be applied to model the correlations [19].

The aforementioned techniques are all instances of *complete modelling*, where the inference of the target unknown quantities is improved by conditioning on knowledge given by crude

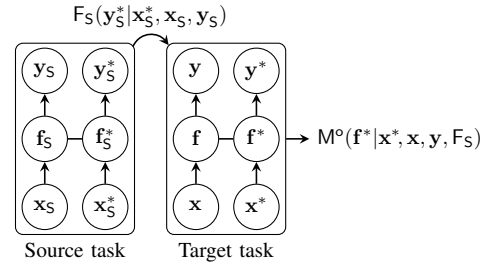


Fig. 1: Supervised source and target tasks process their source ( $x_S, y_S, x_S^*$ ) and target ( $x, y, x^*$ ) known values to learn source ( $f_S, f_S^*, y_S^*$ ) and target ( $f, f^*, y^*$ ) unknown values. Here,  $x, y$ , and  $f$  are labelled inputs, outputs, and function evaluations, respectively, and the superscript \* denotes their unlabelled variants. The source task transfers the posterior predictive distribution of unlabelled outputs,  $F_S$ . The target task improves its performance by using an FPD-optimal,  $F_S$ -conditioned, posterior predictive distribution of unlabelled function values,  $M^o$ .

source data. This requires a dependence structure between the source and target tasks to be specified. In contrast, this paper presents an approach where the target unknown quantities are additionally conditioned on knowledge in the form of a source probability distribution. Therefore, the challenge is to find a probabilistic model over a source probability distribution, which is assumed to be *unavailable*. This is an instance of *incomplete modelling*. The fully probabilistic design (FPD) [20], [21] is an optimal model-completion strategy rooted in the minimum cross-entropy principle [22]. Its primary purpose is to condition unknown quantities on knowledge—a source probability distribution in this paper—for which there is no probabilistic model [23]. The conditioning on a source probability distribution via the FPD approach is the key feature that releases the requirement of any dependence structure between source and target tasks to be specified.

The FPD-optimal model completion was originally used to design the target prior of unknown parameters conditioned on the source output predictor [24] and was later extended to various signal processing problems [25]–[28]. The present paper casts this methodology into the supervised learning context and designs the target posterior predictor of unlabelled function evaluations conditioned on the source posterior predictor of unlabelled outputs. We evaluate the prediction performance and the ability to reject imprecise source knowledge of different methods. The experimental results demonstrate that the proposed approach provides competitive performance and lower computational requirements compared to an alternative strategy relying on explicit dependence assumptions.

## II. FPD KNOWLEDGE TRANSFER BETWEEN SUPERVISED SOURCE AND TARGET TASKS

Let us consider a supervised learning problem specified by a conditional likelihood function in the following form:

$$F(\mathbf{y}|\mathbf{f}, \mathbf{x}), \quad (1)$$

where  $\mathbf{y} \equiv (y_i \in \mathbb{R})_{i=1}^n$  are output data,  $\mathbf{x} \equiv (x_i \in \mathbb{R}^{n_x})_{i=1}^n$  are input data, and  $\mathbf{f} \equiv (f(x_i))_{i=1}^n$  are evaluations of a latent function,  $f: \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ , at the inputs,  $\mathbf{x}$ . One of the primary objectives of supervised learning is to predict the unlabelled function values,  $\mathbf{f}^*$ , and outputs,  $\mathbf{y}^*$ , based on the labelled input-output data,  $(\mathbf{x}, \mathbf{y})$ , and the unlabelled inputs,  $\mathbf{x}^*$ . That is, we seek the joint posterior predictor,

$$F(\mathbf{y}^*, \mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}). \quad (2)$$

In this paper, we consider source and target tasks as shown by Fig. 1, where  $S$  denotes the source variables. Note that, to simplify notation, we do not use any decorator to distinguish the target variables. We assume that the target task does not have any direct knowledge about the quantities processed by the source task. The target task only receives the posterior predictor from the source task,  $F_S(\mathbf{y}_S^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S)$ . The central inference aim is to extend the joint posterior model of the target task, (2), to accommodate additional knowledge in the form of the source posterior predictor,

$$M(\mathbf{y}^*, \mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S). \quad (3)$$

Nevertheless, a probability distribution over this distributional knowledge,  $F_S$ , is assumed to be unavailable. Therefore, the functional form of (3) is unknown, and our goal is to find a mechanism to condition (2) on  $F_S$ . In this paper,  $F$  and  $M$  denote specified (fixed-form) and unspecified (variational-form) distributions, respectively.

To transfer the source posterior predictor of unlabelled outputs,  $F_S$ , we constrain the functional form of the unknown joint model (3) as follows:

$$M(\mathbf{y}^*, \mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) \equiv F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S)M(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S), \quad (4)$$

where we apply the conditional independence, and we restrict the unknown  $F_S$ -conditioned target posterior predictor of unlabelled outputs,  $\mathbf{y}^*$ , to be the source posterior predictor of unlabelled outputs computed at  $\mathbf{y}_S^* \equiv \mathbf{y}^*$ , that is,

$$M(\mathbf{y}^*|\mathbf{f}^*, \mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) \equiv F_S(\mathbf{y}_S^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S)|_{\mathbf{y}_S^* \equiv \mathbf{y}^*}. \quad (5)$$

This specification leaves  $M(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S)$  as the only variational quantity in (4), and enables us to delineate the knowledge-constrained set of possible models, as follows:

$$M \in \mathbf{M} \equiv \{\text{models (4) with } F_S(\mathbf{f}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) \text{ fixed and } M(\mathbf{f}|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) \text{ variational}\}. \quad (6)$$

The joint posterior model (2) reflects behaviour of a supervised learning algorithm in the absence of any source knowledge.

Therefore, we use it as the ideal (prescriptive) design

$$\begin{aligned} M_I(\mathbf{y}^*, \mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) &\equiv F(\mathbf{y}^*, \mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) \\ &= F(\mathbf{y}^*|\mathbf{f}^*, \mathbf{x}^*)F(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}). \end{aligned} \quad (7)$$

Here, we rely on the fact that  $(\mathbf{x}, \mathbf{y})$  does not provide more information about  $\mathbf{y}^*$  when  $(\mathbf{f}^*, \mathbf{x}^*)$  is known.

FPD is an optimal tool to condition a probability distribution on another probability distribution. FPD makes this possible by seeking an optimal model,  $M^\circ$ , that belongs to the knowledge-constrained set,  $M \in \mathbf{M}$  (6), and incorporates preferences about  $M$  defined by an ideal model  $M_I$  (7). The distribution that is closest to  $M_I$  in the minimum Kullback-Leibler divergence (KLD) [29] sense,

$$M^\circ(\mathbf{y}^*, \mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) \equiv \underset{M \in \mathbf{M}}{\operatorname{argmin}} \mathcal{D}(M||M_I), \quad (8)$$

is the FPD-optimal design,  $M^\circ \in \mathbf{M}$ , which we seek. Here, the KLD from  $M$  to  $M_I$  is

$$\mathcal{D}(M||M_I) = E_M \left[ \log \left( \frac{M}{M_I} \right) \right],$$

and the expected value under  $M$  is  $E_M$ .

**Proposition 1.** *The unknown model is a member of the knowledge constrained set,  $M \in \mathbf{M}$  (6), and the ideal model  $M_I$  is (7). Then, the solution of (8)—the FPD-optimal model—is*

$$\begin{aligned} M^\circ(\mathbf{y}^*, \mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) &= \\ &F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S)M^\circ(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S), \end{aligned} \quad (9)$$

where

$$\begin{aligned} M^\circ(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) &\propto F(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) \\ &\times \exp \left\{ \int \log F(\mathbf{y}^*|\mathbf{f}^*, \mathbf{x}^*)F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) d\mathbf{y}^* \right\}. \end{aligned} \quad (10)$$

*Proof.* See Appendix A.  $\square$

Proposition 1 defines the approach to incorporate the source posterior predictor,  $F_S$ , into the target learning procedure. Note that (10) uses  $F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S)$  to perform the update from the posterior  $F(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y})$  to the FPD-optimal,  $F_S$ -conditioned, posterior  $M^\circ(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S)$ . A specific computational flow for producing  $F(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y})$  is known as the *data learning* step. From this perspective, (10) can be seen as an additional step after the data learning step, which we refer to as the *transfer learning* step.

## III. FPD KNOWLEDGE TRANSFER BETWEEN GAUSSIAN PROCESS TASKS

This section is concerned with the transfer learning scenario where the source and target task is to perform the GP regression. We specify the conditional likelihood function (1) as

$$F(\mathbf{y}|\mathbf{f}, \mathbf{x}) \equiv \mathcal{N}(\mathbf{y}; \mathbf{f}, \sigma^2 I_n), \quad (11)$$

where  $\sigma^2$  is the variance of the noise corrupting the (independent and identically distributed) outputs, and  $I_n$  is the  $n$ -dimensional identity matrix. To find the joint posterior model

---

**Algorithm 1:** The FPD-optimal GP regression transfer

---

**Input:**  $m(x), k(x, x'), \sigma^2, \mathbf{x}, \mathbf{x}^*, \mathbf{y}, m_{S,n}(\mathbf{x}_S^*)$ 

1 Data learning step:

2 Use  $m, k, \sigma^2, \mathbf{x}, \mathbf{x}^*$ , and  $\mathbf{y}$  in (15) to compute  $m_n(\mathbf{x}^*)$  and  $k_n(\mathbf{x}^*, \mathbf{x}^*)$ .

3 Transfer learning step:

4 Use  $\mathbf{x}^*, m_n(\mathbf{x}^*), k_n(\mathbf{x}^*, \mathbf{x}^*)$ , and  $m_{S,n}(\mathbf{x}_S^*)$  in (18) to compute  $m_n^\circ(\mathbf{x}^*)$  and  $k_n^\circ(\mathbf{x}^*, \mathbf{x}^*)$ .**Output:**  $m_n^\circ(\mathbf{x}^*), k_n^\circ(\mathbf{x}^*, \mathbf{x}^*)$ 

---

(2), we need to specify a prior distribution over the latent function values,  $F(\mathbf{f}|\mathbf{x})$ . We choose the Gaussian process prior for this purpose.

A stochastic process over a random function,  $f : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ , defines a probability distribution,  $f \sim F$ , such that any finite collection of evaluation points,  $\mathbf{x}$ , induces a joint probability distribution over the function values,  $\mathbf{f}$ . In the case of Gaussian process, we define  $f(x) \sim \mathcal{GP}(m(x), k(x, x'))$ , where  $m(x)$  and  $k(x, x')$  are the mean and covariance functions, respectively, and the induced joint probability distribution is the multivariate Gaussian distribution,

$$F(\mathbf{f}|\mathbf{x}) \equiv \mathcal{N}(\mathbf{f}; m(\mathbf{x}), k(\mathbf{x}, \mathbf{x})). \quad (12)$$

Here,  $m(\mathbf{x})$  is an  $n$ -dimensional vector, and  $k(\mathbf{x}, \mathbf{x})$  is an  $n \times n$ -dimensional matrix. Specific cases of  $m$  and  $k$  allow us to express prior beliefs in the behaviour of  $f$ , including smoothness and periodicity assumptions.

The data learning step of the source and target algorithms is based on the joint posterior model (2). We will require its marginals, which we now recall.

**Lemma 1.** *The probability distribution of output data is (11), and the prior distribution of latent function evaluations is (12). Then, the joint posterior model (2) admits the following marginal distributions:*

$$F(\mathbf{y}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{y}^*; m_n(\mathbf{x}^*), k_n(\mathbf{x}^*, \mathbf{x}^*) + \sigma^2 I_{n_*}), \quad (13)$$

$$F(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{f}^*; m_n(\mathbf{x}^*), k_n(\mathbf{x}^*, \mathbf{x}^*)), \quad (14)$$

where

$$m_n(x) = m(x) + k(x, \mathbf{x}) \times (k(\mathbf{x}, \mathbf{x}) + \sigma^2 I_n)^{-1}(\mathbf{y} - m(\mathbf{x})), \quad (15a)$$

$$k_n(x, x') = k(x, x') - k(x, \mathbf{x}) \times (k(\mathbf{x}, \mathbf{x}) + \sigma^2 I_n)^{-1}k(\mathbf{x}, x'), \quad (15b)$$

with  $k(x, \mathbf{x})$  denoting an  $n$ -dimensional vector.

*Proof.* The proof follows from standard calculus with Gaussian distributions, see, e.g. [6].  $\square$

The target algorithm then additionally proceeds through the transfer learning step given by Proposition 1. We instantiate this in the present context in Proposition 2.

**Proposition 2.** *The target posterior predictor of unlabelled function values is (14), and the source posterior predictor of*

*unlabelled outputs has the same functional form as (13), i.e.,*

$$F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) \equiv \mathcal{N}(\mathbf{y}^*; m_{S,n}(\mathbf{x}_S^*), k_{S,n}(\mathbf{x}_S^*, \mathbf{x}_S^*) + \sigma_S^2 I_{n_*}). \quad (16)$$

Then, the FPD-optimal posterior predictor (10) is

$$M^\circ(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) = \mathcal{N}(\mathbf{f}^*; m_n^\circ(\mathbf{x}^*), k_n^\circ(\mathbf{x}^*, \mathbf{x}^*)), \quad (17)$$

where

$$m_n^\circ(x) = m_n(x) + k_n(x, \mathbf{x}^*) \times (k_n(\mathbf{x}^*, \mathbf{x}^*) + \sigma^2 I_{n_*})^{-1}(m_{S,n}(\mathbf{x}_S^*) - m_n(\mathbf{x}^*)), \quad (18a)$$

$$k_n^\circ(x, x') = k_n(x, x') - k_n(x, \mathbf{x}^*) \times (k_n(\mathbf{x}^*, \mathbf{x}^*) + \sigma^2 I_{n_*})^{-1}k_n(\mathbf{x}^*, x'). \quad (18b)$$

*Proof.* See Appendix B.  $\square$

Proposition 2 implies a computational procedure for the FPD-optimal *target* task as summarized in Algorithm 1. The *source* procedure is based only on the data step, i.e., line 2 of Algorithm 1, with the associated quantities and the GP shaping functions decorated by  $S$ .

## IV. EXPERIMENTS

To illustrate the key features of the proposed FPD-optimal algorithm, we focus on a simple example where the source and target tasks are to regress a nonlinear function observed in Gaussian noise. We compare the following methods: Source algorithm with *No Transfer (SNT)*; target algorithm with *No Transfer (NT)*; target algorithm based on the *Linear Model of Coregionalization (LMC)* method [14], where the correlation coefficient,  $\rho$ , between the source and target tasks is specified, see [30] for details; and target algorithm using the *FPD-optimal transfer (FPD)* given in Algorithm 1. We consider that the source and target output data are generated according to

$$F(\mathbf{y}_S|\mathbf{f}_S, \mathbf{x}_S) \equiv \prod_{i=1}^{n_S} \mathcal{N}(y_{S,i}; f_S(x_{S,i}), \sigma_S^2), \\ F(\mathbf{y}|\mathbf{f}, \mathbf{x}) \equiv \prod_{i=1}^n \mathcal{N}(y_i; f(x_i), \sigma^2),$$

where  $f_S(x) = \frac{x}{2} \sin(x)$ ,  $\sigma_S^2 = 0.01$ ,  $f(x) = \frac{x^{0.95}}{1.9} \sin(x)$ , and  $\sigma^2 = 1$ . The source and target input data are simulated by

$$\mathbf{x} \sim \mathcal{U}_n(-10, 10), \\ \mathbf{x}_S \sim \mathcal{U}_{n_S}(-10, 10).$$

Here,  $\mathcal{U}_n(a, b)$  is the uniform distribution on the  $n$ -fold open interval  $(a, b)$ . This scenario reflects a situation where  $f$  is a moderately distorted version of  $f_S$ , which is observed in a higher noise. In such a case, we would like to know if the source procedure with high-quality data can assist the target procedure with low-quality data. The prior mean functions,  $m, m_S$ , are given by the zero vector, and the covariance functions,  $k, k_S$ , are both specified as the squared exponential kernel [8],

$$k(x, x') = \sigma_f^2 \exp \left\{ -\frac{1}{2} \frac{\|x - x'\|^2}{l^2} \right\},$$

where  $\sigma_f^2$  is the signal variance,  $l$  is the length-scale, and  $\|\cdot\|$  denotes the euclidean norm. We set  $\sigma_f^2 = \sigma_f^2 = 2$  and  $l_S = l = 0.5$  for all algorithms. We are interested in the prediction

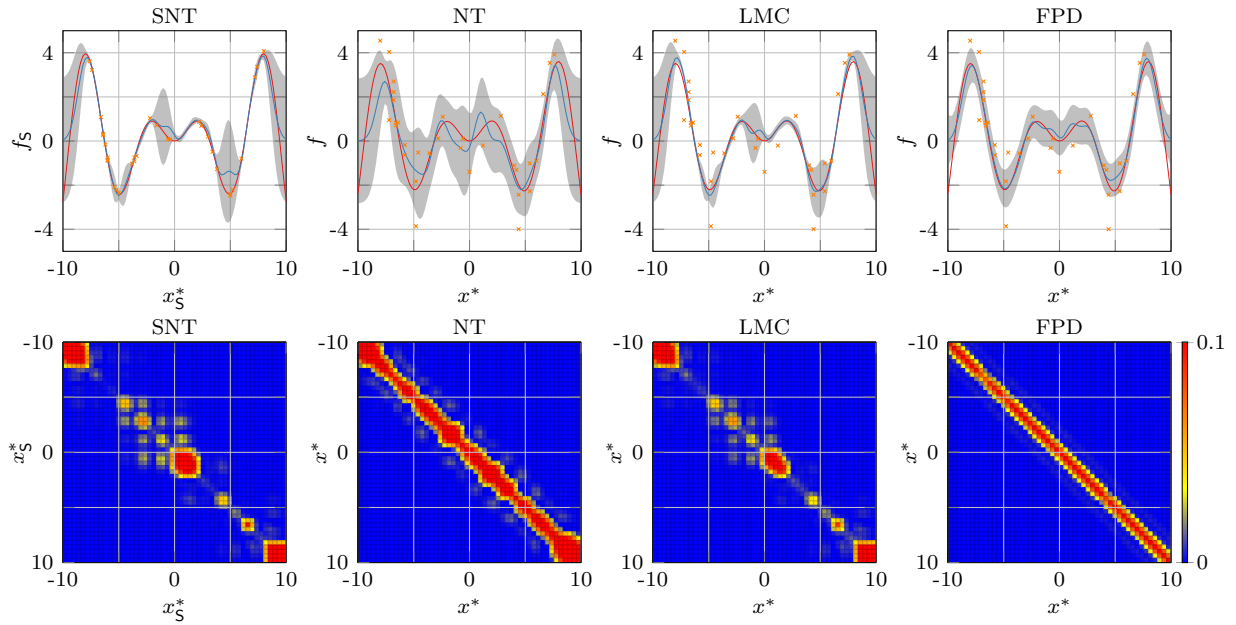


Fig. 2: Top: the prediction performance at the unlabelled input data  $\mathbf{x}^*$ . Here, (—) is the true function,  $f(x^*)$ , ( $\circ$ ) are the labelled (noisy) outputs,  $y$ , (—) is the posterior predictive mean function,  $m_n(x^*)$ , and (shaded) is the posterior predictive  $2\sigma$ -region,  $\pm 2\sqrt{k_n(x^*, x^*)}$ . Bottom: the covariance function evaluated at the unlabelled inputs for the SNT algorithm,  $k_{S,n}(x_S^*, x_S^*)$ , the NT and LMC algorithms,  $k_n(x^*, x^*)$ , and the FPD algorithm,  $k_n^o(x^*, x^*)$ .

performance of the various methods for unlabelled input data  $\mathbf{x}^*$  generated on the closed interval  $[-10, 10]$  with the step-size 0.2. We evaluate the error norm (EN) between the true function values at  $\mathbf{x}^*$ ,  $\mathbf{f}^* \equiv (f(x_i^*))_{i=1}^{n_S}$ , and their posterior predictive mean estimate,  $m_n(\mathbf{x}^*)$ ,

$$\text{EN} \equiv \|\mathbf{f}^* - m_n(\mathbf{x}^*)\|.$$

In the first experiment, whose results are depicted in Fig. 2, we illustrate the prediction performance of the compared methods for  $n \equiv n_S \equiv 32$ . The top row of Fig. 2 reveals that—despite processing the imprecise target output data—the FPD algorithm takes advantage of the source posterior predictive distribution,  $F_S$ , and—compared to the NT algorithm—it significantly recovers the shape of the latent function and reduces the  $2\sigma$ -region. The LMC algorithm—with its correlation coefficient set to  $\rho = 1$  (extreme value)—performs moderately better than the FPD algorithm. The bottom row of Fig. 2 presents the covariance matrices associated with the results above, providing more details of the uncertainty representation.

Transfer learning is useful mainly when a source algorithm processes either more precise data or simply more data of the same precision, compared to the target algorithm. We demonstrate this with two experiments depicted in the top row of Fig. 3, where the EN of various methods changes as a function of the variance of the source output data,  $\sigma_S^2$  (top-left), and the number of source data,  $n_S$  (top-right). The NT algorithm does not depend on these parameters and thus defines a reference EN level for comparing the remaining methods. If any algorithm produces an EN below or above this reference level, we say that it provides *positive* or *negative*

transfer, respectively. When any algorithm saturates at the EN level of the NT algorithm, we say that it achieves *robust* transfer. The FPD algorithm switches its behaviour near the intersection points  $\sigma_S^2 = \sigma^2$  and  $n_S = n$ , yielding positive transfer for (approximately)  $\sigma_S^2 < \sigma^2$  and  $n_S > n$ , and negative transfer for  $\sigma_S^2 > \sigma^2$  and  $n_S < n$ . We explain this behaviour in Section V. The LMC algorithm gives positive transfer for  $\sigma_S^2 < 10^2$  and all  $n_S$ , achieving robust transfer for  $\sigma_S^2 > 10^2$ .

The next experiment (bottom-left of Fig. 3) evaluates the computational time of the compared methods by changing the number of target data according to  $n = 2^i$  for  $i = 1, \dots, 9$ . We see that the proposed FPD algorithm is computationally more efficient and offers an improved prediction performance compared to the LMC algorithm (again, with  $\rho = 1$ ) when  $i = 1, \dots, 7$ . For  $i = 8$  and  $i = 9$ , the FPD and LMC algorithms perform similarly.

The final experiment (bottom-right of Fig. 3) demonstrates the impact of the correlation coefficient,  $\rho$ , on the EN of the LMC algorithm for source knowledge of varying quality. We observe that values of  $\rho$  other than  $\rho \rightarrow 1$  cause the EN of the LMC algorithm to be worse compared to the FPD algorithm. More importantly, we see that the LMC algorithm yields negative transfer for  $\rho < 0$  even for high-quality source knowledge. Note that if  $\rho = 0$  (uncorrelated source and target tasks), the LMC algorithm recovers the EN of the NT algorithm. The proposed FPD algorithm does not depend on any structural ( $\rho$ -like) assumptions, and yet it is able to deliver very competitive performance compared to the LMC algorithm when the source knowledge is of high quality.

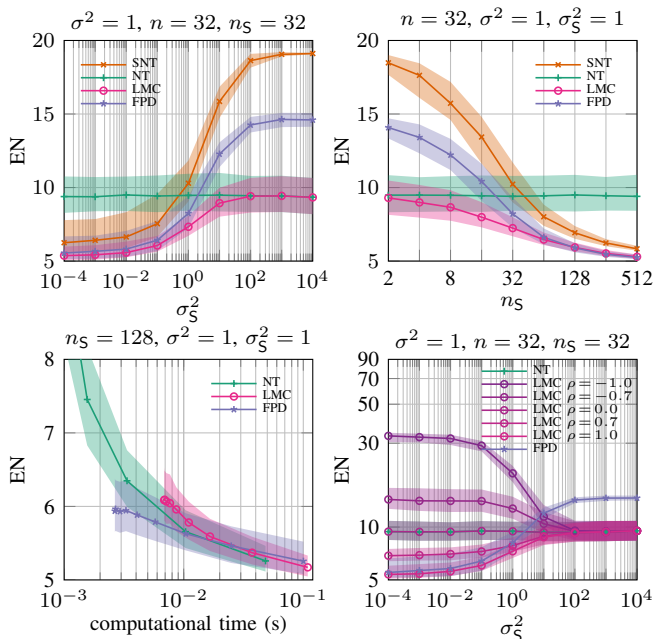


Fig. 3: Top-left: the error norm (EN) versus the source output data variance,  $\sigma_S^2$ . Top-right: the EN versus the number of source data,  $n_S$ . Bottom-left: the EN versus the computational time in seconds (influenced by the number of target data). Bottom-right: the EN versus the source output data variance,  $\sigma_S^2$ , for various correlation coefficients of the LMC algorithm,  $\rho$ . The results are averaged over 1000 independent simulation epochs, with the solid line and shaded area being the median and interquartile range, respectively.

## V. DISCUSSION

In the first part of this section, we would like to discuss the following question: if the true function lies in the  $2\sigma$ -region for the SNT and NT algorithms—as shown on top of Fig. 2—why is this behaviour not fully preserved by the LMC and FPD algorithms? In the case of the LMC algorithm, surprisingly, this happens close to the points where both the SNT and NT algorithms process output data, i.e., close to  $x^* = -8$  and  $x^* = 8$ . Although, such behaviour occurred only occasionally during our experiments—and is still consistent due to the fact that the true function is certainly in the  $3\sigma$ -region—this situation demonstrates that the LMC algorithm can deliver overconfident predictions compared to the SNT and NT algorithms. We observed that this is often the case when  $\rho \rightarrow 1$ , which is, however, the value where the LMC algorithm provides the best EN, see the bottom-right of Fig. 3. Importantly, this issue of overconfident predictions diminishes for lower  $\rho$ , but the FPD algorithm then outperforms the LMC algorithm (Fig. 3). In the case of the FPD algorithm, the true function lies outside the  $2\sigma$ -interval where neither the SNT nor NT algorithm processes any output data, i.e., close to  $x^* = -10$  and  $x^* = 10$ . This is more sensible behaviour than in the LMC case. Moreover, a closer look at bottom of Fig. 2 reveals that the uncertainty representation of the FPD algorithm is rather uniform compared to the LMC algorithm which adjusts its uncertainty according to the number of data processed. This behaviour of the FPD algorithm is a result of the loss of the covariance function

of the source posterior predictor,  $k_{S,n}$ , during the transfer learning step. The source covariance function is lost when computing the exponential structure of (10), see also (19) in Appendix B. The consequence of this can best be seen in (18), where only the source mean function,  $m_{S,n}$ , is present but the source covariance function,  $k_{S,n}$ , is missing. The same reasoning explains the fact that the FPD algorithm suffers from negative transfer for (approximately)  $\sigma_S^2 > \sigma^2$  and  $n_S < n$ , as illustrated in the top row of Fig. 3. The kernel hyperparameters were set to fixed values for all the investigated methods in order to ensure a fair comparison. However, we expect that optimizing these hyperparameters can address the problems associated with uncertainty representation of both the LMC and FPD algorithms. This will form the agenda of future work.

Next, we comment on the choice of the correlation coefficient,  $\rho$ , of the LMC method, as illustrated in bottom-right of Fig. 2. We see that this choice is critical, leading not only to a slightly worse EN compared to the FPD algorithm, but, crucially, to negative transfer even for high-quality source knowledge. In practice,  $\rho$  can be treated as an additional hyperparameter which needs to be optimized [14]. There is no such hyperparameter dependence in the proposed FPD algorithm. Notwithstanding this, the FPD algorithm is competitive with the LMC algorithm. Recall that the LMC algorithm requires complete modelling of the source and target dependence, and directly processes the raw source data. In contrast, our FPD algorithm (resulting from the Bayesian transfer in Fig. 1) only transfers the sufficient statistics of  $F_S$ . Furthermore, the essence of this approach is to avoid explicit joint modelling assumptions between the source and target tasks, which, anyway, are very hard to propose in most cases. In addition,  $\rho$  will increase the computational cost of marginalizing the hyperparameters in the LMC algorithm via Markov chain Monte Carlo methods. This added complexity is avoided in the FPD algorithm.

## VI. CONCLUSION

This paper has presented FPD-optimal Bayesian transfer learning in the context of Gaussian process regression. The proposed FPD algorithm is a consequence of optimal model completion in the context where the stochastic dependence between the source and target tasks is not specified. As illustrated by the experimental evidence, the FPD algorithm can offer competitive (and sometimes even better) prediction performance, lower computational requirements, and fewer hyperparameters to tune compared to the LMC approach. The latter relies on explicit—and often brittle—dependence assumptions. However, the FPD algorithm suffers from negative transfer when the source knowledge is imprecise compared to the target knowledge. We encountered similar behaviour when developing transfer learning algorithms in Bayesian filtering applications [25], [26]. There, we were successful in resolving this issue by introducing an auxiliary variable augmentation [27]. Our preliminary investigations suggest that the same solution can also be successful in the present context. We will provide more details in a future paper.

A. Proof of Proposition 1

After substituting (4) and (7) into (8), we obtain:

$$\begin{aligned} \mathcal{D}(M||M_I) &= \int F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) M(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) \\ &\times \log \left( \frac{F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) M(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S)}{F(\mathbf{y}^*|\mathbf{f}^*, \mathbf{x}^*) F(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y})} \right) d\mathbf{y}^* d\mathbf{f}^* \\ &= \int M(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S) \\ &\times \log \left( \frac{M(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S)}{M^o(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}, F_S)} \right) d\mathbf{f}^* - \mathcal{H}_{F_S} - \log c_{M^o}, \end{aligned}$$

where  $\mathcal{H}_{F_S}$  is the differential entropy of  $F_S$ ,

$$\mathcal{H}_{F_S} = - \int F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) \log F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) d\mathbf{y}^*,$$

and  $c_{M^o}$  is the normalizing constant,

$$\begin{aligned} c_{M^o} &= \int F(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) \\ &\times \exp \left\{ \int \log F(\mathbf{y}^*|\mathbf{f}^*, \mathbf{x}^*) F_S(\mathbf{y}^*|\mathbf{x}_S^*, \mathbf{x}_S, \mathbf{y}_S) d\mathbf{y}^* \right\} d\mathbf{f}^*. \quad \square \end{aligned}$$

B. Proof of Proposition 2

Substituting (11) and (16) into the exponential term of (10) leads to

$$\exp \{E_{F_S}[\log F(\mathbf{y}^*|\mathbf{f}^*, \mathbf{x}^*)]\} \propto \mathcal{N}(m_{S,n}; \mathbf{f}^*, \sigma^2 I_{n_*}). \quad (19)$$

Then, the product of (14) and (19) yields

$$\mathcal{N} \left( \begin{bmatrix} m_{S,n} \\ \mathbf{f}^* \end{bmatrix}; \begin{bmatrix} m_n(\mathbf{x}^*) \\ m_n(\mathbf{x}^*) \end{bmatrix}, \begin{bmatrix} k_n(\mathbf{x}^*, \mathbf{x}^*) + \sigma^2 I_{n_*} & k_n(\mathbf{x}^*, \mathbf{x}^*) \\ k_n(\mathbf{x}^*, \mathbf{x}^*) & k_n(\mathbf{x}^*, \mathbf{x}^*) \end{bmatrix} \right),$$

which admits the conditional distribution (17) with the moments (18).  $\square$

REFERENCES

[1] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.

[2] Y. Zhang and Q. Yang, "A survey on multi-task learning," *arXiv preprint arXiv:1707.08114*, 2018.

[3] K. Weiss, T. M. Khoshgoftaar, and D. D. Wang, "A survey of transfer learning," *Journal of Big Data*, vol. 3, no. 1, pp. 9, 2016.

[4] J. Lu, V. Behbood, P. Hao, H. Zuo, S. Xue, and G. Zhang, "Transfer learning using computational intelligence: A survey," *Knowledge-Based Systems*, vol. 80, pp. 14–23, 2015.

[5] L. Torrey and J. Shavlik, "Transfer learning," in *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques*, pp. 242–264. IGI Global, 2010.

[6] K. P. Murphy, *Machine learning: A probabilistic perspective*, MIT Press, 2012.

[7] C. M. Bishop, *Pattern recognition and machine learning*, Springer, 2006.

[8] C. Rasmussen and C. Williams, *Gaussian processes for machine learning*, MIT Press, 2006.

[9] K. Swersky, J. Snoek, and R. P. Adams, "Multi-task Bayesian optimization," in *Advances in neural information processing systems*, 2013, pp. 2004–2012.

[10] L.-F. Cheng, G. Darnell, B. Dumitrescu, C. Chivers, M. E. Draugelis, K. Li, and B. E. Engelhardt, "Sparse multi-output Gaussian processes for medical time series prediction," *arXiv preprint arXiv:1703.09112*, 2017.

[11] S. Vasudevan, F. Ramos, E. Nettleton, and H. Durrant-Whyte, "Non-stationary dependent Gaussian processes for data fusion in large-scale terrain modeling," in *2011 IEEE International Conference on Robotics and Automation*. IEEE, 2011, pp. 1875–1882.

[12] K. M. A. Chai, Christopher K. I. Williams, S. Klanke, and S. Vijayakumar, "Multi-task Gaussian process learning of robot inverse dynamics," in *Advances in neural information processing systems*, 2008, pp. 265–272.

[13] Y. W. Teh, M. Seeger, and M. I. Jordan, "Semiparametric latent factor models," in *AISTATS 2005-Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics*, 2005.

[14] E. V. Bonilla, K. M. Chai, and C. Williams, "Multi-task Gaussian process prediction," in *Advances in neural information processing systems*, 2008, pp. 153–160.

[15] T. V. Nguyen and E. V. Bonilla, "Collaborative multi-output Gaussian processes," in *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2014, pp. 643–652.

[16] A. G. Wilson, D. A. Knowles, and Z. Ghahramani, "Gaussian process regression networks," in *Proceedings of the 29th International Conference on Machine Learning*. Omnipress, 2012, pp. 1139–1146.

[17] T. Nguyen and E. V. Bonilla, "Efficient variational inference for Gaussian process regression networks," in *Artificial Intelligence and Statistics*, 2013, pp. 472–480.

[18] M. A. Álvarez and N. D. Lawrence, "Computationally efficient convolved multiple output Gaussian processes," *Journal of Machine Learning Research*, vol. 12, no. May, pp. 1459–1500, 2011.

[19] A. Boustati and R. S. Savage, "Multi-task learning in deep Gaussian processes with multi-kernel layers," 2019.

[20] M. Kárný, "Towards fully probabilistic control design," *Automatica*, vol. 32, no. 12, pp. 1719–1722, 1996.

[21] M. Kárný and T. Kroupa, "Axiomatisation of fully probabilistic design," *Information Sciences*, vol. 186, no. 1, pp. 105–113, 2012.

[22] J. Shore and R. Johnson, "Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross-entropy," *IEEE Transactions on Information Theory*, vol. 26, no. 1, pp. 26–37, 1980.

[23] A. Quinn, M. Kárný, and T. V. Guy, "Fully probabilistic design of hierarchical Bayesian models," *Information Sciences*, vol. 369, pp. 532–547, 2016.

[24] A. Quinn, M. Kárný, and T. V. Guy, "Optimal design of priors constrained by external predictors," *International Journal of Approximate Reasoning*, vol. 84, pp. 150–158, 2017.

[25] C. Foley and A. Quinn, "Fully probabilistic design for knowledge transfer in a pair of Kalman filters," *IEEE Signal Processing Letters*, vol. 25, no. 4, pp. 487–490, 2018.

[26] M. Papež and A. Quinn, "Dynamic Bayesian knowledge transfer between a pair of Kalman filters," in *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2018.

[27] M. Papež and A. Quinn, "Robust Bayesian transfer learning between Kalman filters," in *2019 IEEE 29th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2019.

[28] L. Jirsa, L. Pavelkova, and A. Quinn, "Knowledge transfer in a pair of uniformly modelled Bayesian filters," in *2019 16th International Conference on Informatics in Control, Automation and Robotics (ICINCO)*. IEEE, 2019.

[29] S. Kullback and R. A. Leibler, "On information and sufficiency," *The Annals of Mathematical Statistics*, vol. 22, no. 1, pp. 79–86, 1951.

[30] K. M. Chai, "Generalization errors and learning curves for regression with multi-task Gaussian processes," in *Advances in neural information processing systems*, 2009, pp. 279–287.