# Robust Bayesian Transfer Learning between Autoregressive Inference Tasks

Alec Barber<sup>a</sup> barberalec2@gmail.com

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

Abstract—Bayesian transfer learning typically relies on a complete stochastic dependence specification between source and target learners. We advocate a solution to the Bayesian transfer learning paradigm which adopts Fully Probabilistic Design (FPD) to search for an optimal choice of distribution constrained by probabilistic source knowledge. Using this optimal decisionmaking strategy, an algorithm for accepting source knowledge is identified but is found to be effectively insensitive to source uncertainty. Therefore, we propose an adaptation of the FPD framework which results in a robust transfer learning algorithm.

Experimental evidence gathered via synthetic data shows enhanced performance when employing both optimal algorithms in a low source data predictor variance regime. In a high source data predictor variance setting, only our adapted FPD-optimal algorithm achieves robustness.

*Index Terms*—Autoregressive (AR) model, Bayesian transfer learning, data-predictive transfer, FPD, robust transfer.

### I. INTRODUCTION

Transfer learning (multi-task learning) [1] [2] is the process of adopting knowledge from a set of source learning tasks to increase the learning rate of a given target task [3]. This paper is concerned with transfer learning in a Bayesian context and proposes two algorithms for knowledge transfer between a source and target autoregressive (AR) inference task via the sharing of a one-step-ahead data predictor or a conditional one-step-ahead mean data predictor (probability distributions). The objective of Bayesian Transfer Learning (BTL) is to identify an optimal target distribution M°, describing the targets beliefs in its generative parameters  $\theta$ , conditioned on provided source knowledge F<sub>S</sub> and n locally observed data  $z_n$ . We are further concerned with the robustness property, i.e. positive knowledge transfer when quality source data is available and mitigating negative transfer otherwise.

BTL in literature is typically approached by the design of a prior for the target via knowledge provided by source learners [4], or by the assumption/design of a joint model between the target and source learners [5] which provides the setting for Bayesian conditioning. Prior design for knowledge transfer can prove effective but is restrictive due to its inability to incorporate synchronous, on-line learning. We advocate that any requirement for the design or assumption of a full model

The research has been supported by GAČR grant 18-15970S

Anthony Quinn <sup>a,b</sup> aquinn@tcd.ie

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



Fig. 1. Source and target inference tasks process their local data  $z_{S,n}$  and  $z_n$  to learn AR parameters  $a_S$ ,  $r_S$  and a, r respectively. At each time step the source makes available a distribution  $F_S$  to the target. The target elicits an FPD-optimal,  $F_S$ -conditioned, posterior predictive distribution which improves its parameter inference performance on a and r.

between target and sources is a restrictive form of transfer learning.

This paper avoids the specification of a complete model between target and source by adopting principles of Fully Probabilistic Design (FPD) [6] [7] to elicit a source knowledge conditioned distribution on the target's parameters. FPD is as an axiomatically justified [8] optimal distribution design framework which is rooted in the minimum cross-entropy principle. FPD has previously been applied to transfer learning applications such as Gaussian process regression [9], Kalman filtering [10] and Student-t filtering [11].

This paper extends the application of FPD-based transfer learning to an AR setting. The AR model is favoured for its ability to describe relationships between lagged realisations of temporal data. AR models are frequently adopted in the domains of speech analysis [12] and finance [13].

The remainder of this paper is organised as follows. Section II gives an overview of Bayesian parameter inference for an isolated AR setting. Sections III and IV describe the principles of FPD and applies them to the AR model to achieve knowledge transfer between a source and target inference task. Section V identifies an alternative FPD-optimal robust algorithm for the AR transfer learning paradigm. Sections VI and VII outline the experimental setting for validations of the

algorithms and provides a discussion of the findings. Section VIII concludes the paper.

## II. REVIEW OF BAYESIAN PARAMETER INFERENCE FOR THE AUTOREGRESSION MODEL

The objective when adopting an AR model of order p (AR(p)) to describe n signal data  $z_n$ , is the inference of unknown static generative parameters  $\theta = \{a, r\}$ . Adopting inferences on the parameters, any desired moments or marginals on objects of interest can be derived (e.g. data predictors).

Data is typically made available to the model at discrete steps indexed by i = 1, 2, 3, ... Equations (1) describe a univariate time-invariant AR(p) model where  $e_n$  is the Gaussian innovation at step n,  $z_n \in \mathbb{R}$  is the scalar observation and  $\psi_n \in \mathbb{R}^p$  is the autoregression vector (p previous observations). This paper adopts <sup>T</sup> to denote the transposition operator.

$$\begin{aligned} \mathbf{z}_{n} &= \mu_{\mathbf{z}_{n}} + r\mathbf{e}_{n} = \boldsymbol{\psi}_{n}^{\mathsf{T}} \mathbf{a} + r\mathbf{e}_{n}, \\ & \mathbf{a} \in \mathbb{R}^{\mathsf{p}}, \mathsf{r} \in \mathbb{R}^{+}, \\ & \mathbf{e}_{n} \sim \mathcal{N}_{\mathsf{e}}(0, 1), \\ & \boldsymbol{\psi}_{n} = \begin{bmatrix} \mathsf{z}_{\mathsf{n}-1} \dots \mathsf{z}_{\mathsf{n}-\mathsf{p}} \end{bmatrix}^{\mathsf{T}} \end{aligned}$$
(1)

AR(p) parameter estimation is classically solved via the normal equations which adopt the Wiener criterion [15]. In the Bayesian setting, we adopt a distribution to describe stochastic belief in  $\{a, r\}$  which can be subsequently updated after each datum observation. A key advantage of the Bayesian method is that parameter (or other moments) estimates come equipped with a first principled measure of uncertainty.

We will now identify a solution for the processing of data  $z_n$  in both an online and offline inference setting with respect to maintaining computational tractability (i.e. adopting conjugate priors on the model parameters).

$$\mathsf{F}(\mathsf{z}_{\mathsf{n}} \mid \mathbf{a}, \mathsf{r}, \mathbf{z}_{\mathsf{n}-1}) \equiv \mathcal{N}_{\mathsf{z}_{\mathsf{n}}}(\mu_{\mathsf{z}_{\mathsf{n}}}, \mathsf{r}) = \mathcal{N}_{\mathsf{z}_{\mathsf{n}}}(\boldsymbol{\psi}_{\mathsf{n}}^{\mathsf{T}} \mathbf{a}, \mathsf{r}) \qquad (2)$$

$$\begin{split} F(\boldsymbol{a},r \mid \boldsymbol{z}_n) &= F(\boldsymbol{a},r \mid \boldsymbol{V}_n, \nu_n) = \mathcal{N}\mathrm{i}\mathrm{G}_{\boldsymbol{a},r}(\boldsymbol{V}_n, \nu_n), \\ \boldsymbol{V}_n &= \begin{bmatrix} v_{11,n} & v_{a1,n}^T \\ v_{a1,n} & \boldsymbol{V}_{aa,n} \end{bmatrix} \end{split} \tag{3}$$

The Gaussian AR generative model (2) is a member of the exponential family [14] and therefore is guaranteed to have a conjugate prior distribution and sufficient statistics. The conjugate prior for (2) is identified as the Normal inverse-Gamma (NiG) (4) distribution [14]. Conjugacy is important in signal processing applications as it facilitates a fully tractable, recursive, and online computational inference flow [18].

$$\mathcal{N}iG_{\mathbf{a},\mathbf{r}}(\mathbf{V}_{\mathbf{n}},\nu_{\mathbf{n}}) = \mathbf{K}_{\mathbf{n}}^{-1}\mathbf{r}^{-\frac{\nu_{\mathbf{n}}}{2}} \exp\left[-\frac{1}{2\mathbf{r}}\begin{bmatrix}-1 & \mathbf{a}^{T}\end{bmatrix}\mathbf{V}_{\mathbf{n}}\begin{bmatrix}-1\\\mathbf{a}\end{bmatrix}\right],$$

$$\lambda_{\mathbf{n}} = \mathbf{v}_{11,\mathbf{n}} - \mathbf{v}_{\mathbf{a}1,\mathbf{n}}^{\mathsf{T}}\mathbf{V}_{\mathbf{a}\mathbf{a},\mathbf{n}}^{-1}\mathbf{v}_{\mathbf{a}1,\mathbf{n}}$$

$$\mathbf{K}_{\mathbf{n}} = \Gamma\left(\frac{\nu_{\mathbf{n}}}{2}\right)\lambda_{\mathbf{n}}^{-\frac{\nu_{\mathbf{n}}}{2}}|\mathbf{V}_{\mathbf{a}\mathbf{a},\mathbf{n}}|^{-\frac{1}{2}}(2\pi)^{\frac{1}{2}}.$$
(4)

Equation (4) describes the sufficient statistic parameterisation of  $\mathcal{N}iG$  which is adopted in this paper.  $\mathbf{V}_n \in \mathbb{R}^{(p+1) \times (p+1)}$  is the Extended Information Matrix (EIM),  $\nu_n \in \mathbb{R}^+$  is the degrees of freedom parameter,  $\Gamma(\cdot)$  denotes the gamma function

and  $K_n$  is the normalisation constant. A requirement of the EIM is that it must be symmetric and positive semi-definite. Equation (3) shows the partitioning of the EIM into its subcomponents, isolating the scalar  $v_{11,n}$ , the vector  $\mathbf{v}_{a1,n} \in \mathbb{R}^p$ and the matrix  $\mathbf{V}_{aa,n} \in \mathbb{R}^{p \times p}$ . The degrees of freedom parameter  $\nu_n$  can be interpreted informally as a counter of observed data (note  $\nu_n$  can be any positive, non-zero, real value).  $\lambda_n$ is the total Bayesian squared error. Provided in (5) below are some important moments of  $\mathcal{N}$ iG.

$$\begin{split} \mathsf{E}[\mathbf{a} \mid \mathbf{z}_{\mathsf{n}}] &= \mathsf{E}[\mathbf{a} \mid \mathsf{r}, \mathbf{z}_{\mathsf{n}}] = \mathbf{V}_{\mathsf{a}\mathsf{a},\mathsf{n}}^{-1} \mathbf{v}_{\mathsf{a}\mathsf{1},\mathsf{n}}, \\ \mathsf{E}[\mathsf{r} \mid \mathbf{z}_{\mathsf{n}}] &= \frac{\lambda_{\mathsf{n}}}{\nu_{\mathsf{n}} - \mathsf{p} + 2}, \qquad \nu_{\mathsf{n}} > \mathsf{p} - 2 \end{split}$$
(5)

It can be shown that the following are the fixed and finite dimensional updates of the Bayesian sufficient statistics prescribed by (2) and (3) after observing  $z_n$  [15].

$$\mathbf{V}_{n} = \mathbf{V}_{n-1} + \begin{bmatrix} \mathbf{z}_{n} \\ \boldsymbol{\psi}_{n} \end{bmatrix} \begin{bmatrix} \mathbf{z}_{n} & \boldsymbol{\psi}_{n}^{\mathsf{T}} \end{bmatrix} = \mathbf{V}_{0} + \sum_{i=1}^{n} \begin{bmatrix} \mathbf{z}_{i} \\ \boldsymbol{\psi}_{i} \end{bmatrix} \begin{bmatrix} \mathbf{z}_{i} & \boldsymbol{\psi}_{i}^{\mathsf{T}} \end{bmatrix}$$
$$\nu_{n} = \nu_{n-1} + 1 = \nu_{0} + \mathbf{n}$$
(6)

The EIM  $V_{n-1}$  is updated by the accumulation of the outer product of the vector  $z_n$  extended by  $\psi_n$ . Of particular note is that this isolated update procedure reappears in similar forms in both sections IV and V when conducting transfer learning.  $V_0$  and  $\nu_0$  in (6) represent the prior belief in the parameters before any data is considered. If no preference is provided, a diffuse distribution is typically adopted (i.e.  $V_0 = \epsilon I_{p+1}, \nu_0 = \epsilon$ , where  $\epsilon$  is a small positive real number).

The one-step-ahead data predictor for the NiG model is identified in (7) as the 3 parameter Student-t distribution.

$$\begin{split} \mathsf{F}(\mathsf{z}_{\mathsf{n}+1} \mid \mathbf{z}_{\mathsf{n}}) &= \int_{\mathbf{a},\mathsf{r}} \mathsf{F}(\mathsf{z}_{\mathsf{n}+1} \mid \mathbf{a},\mathsf{r},\boldsymbol{\psi}_{\mathsf{n}+1}) \mathsf{F}(\mathbf{a},\mathsf{r} \mid \mathbf{z}_{\mathsf{n}}) d\mathbf{a} d\mathsf{r} \\ &= \mathsf{St-t}_{\mathsf{z}_{\mathsf{n}+1}} \left( \boldsymbol{\psi}_{\mathsf{n}+1}^{\mathsf{T}} \mathbf{V}_{\mathsf{a}\mathsf{a},\mathsf{n}}^{-1} \mathbf{v}_{\mathsf{a}\mathsf{1},\mathsf{n}}, \lambda_{\mathsf{n}} \frac{1 + \boldsymbol{\psi}_{\mathsf{n}+1}^{\mathsf{T}} \mathbf{V}_{\mathsf{a}\mathsf{a},\mathsf{n}}^{-1} \boldsymbol{\psi}_{\mathsf{n}+1}}{\nu_{\mathsf{n}}}, \nu_{\mathsf{n}} \right) \end{split}$$
(7)

We will now provide a description of the FPD distribution design framework in the context of transfer learning which we will adopt to optimally choose a source knowledge conditioned object.

#### III. FPD-OPTIMAL BAYESIAN TRANSFER LEARNING

This paper is concerned with BTL in an incompletely modeled setting (i.e. where no explicit relationship between source and target is assumed or designed). This incompletely modeled setting is shown in figure 1 where both source and target inference machines conduct independent modelling of local parameters with no explicit association. We desire a distribution on the target's parameters  $M \in M$  that is conditioned on source knowledge  $F_S$ . The choice of M is not uniquely defined and is a decision making task. In this paper we adopt an optimal choice of distribution  $M^{\circ}$  via the Bayesian minimum risk decision making strategy known as FPD [7].

Adoption of FPD design principles require the specifications of two key elements: 1) a user-specified choice of an ideal, zero-loss choice distribution  $M_1$ , 2) a description of how knowledge constrains the set of possible candidate distributions  $M \in \mathbf{M}$ . Note that this paper adopts M and F as variational (i.e. to be designed) and fixed-form distributions respectively.

$$M^{\circ} \equiv \mathop{\arg\min}_{M \in \textbf{M}} \mathcal{D}\left(M || M_{I}\right) \tag{8}$$

FPD dictates that the optimal distribution choice  $M^{\circ}$  should be chosen as the candidate  $M \in \mathbf{M}$  which minimises the Kullback-Leibler divergence to the user-specified ideal  $M_1$ . When no restrictions are placed on  $\mathbf{M}$ , the optimal zeroloss choice in distribution is always  $M_1$ . The Kullback-Leibler divergence is defined in (9) as the expected log-odds between the candidate and ideal distributions [18].

$$\mathcal{D}\left(\mathsf{M}||\mathsf{M}_{\mathsf{I}}\right) = \mathsf{E}_{\mathsf{M}}\left[\log\left(\frac{\mathsf{M}}{\mathsf{M}_{\mathsf{I}}}\right)\right] \tag{9}$$

We will now instantiate equations (8) and (9) for the NiG-AR tasks (2) and (3) where the object being transferred is the source's one-step-ahead data predictor transfer (7).

### IV. FPD-OPTIMAL BTL BETWEEN AR INFERENCE TASKS

Equation (10) outlines the choice in ideal distribution as the joint model of **a**, r and the next unknown datum  $z_{n+1}$ . This joint model expressed in terms of the target's fixed-form generative model (2) and the posterior predictive distribution of the parameters (3). This choice in ideal distribution is motivated by its complete description of the state of the inference task in the absence of external knowledge.

$$\mathsf{M}_{\mathsf{I}}(\mathbf{a},\mathsf{r},\mathsf{z}_{\mathsf{n}+1}\mid\mathbf{z}_{\mathsf{n}}) \equiv \mathsf{F}(\mathsf{z}_{\mathsf{n}+1}\mid\mathbf{a},\mathsf{r},\boldsymbol{\psi}_{\mathsf{n}+1})\mathsf{F}\left(\mathbf{a},\mathsf{r}\mid\mathbf{z}_{\mathsf{n}}\right) \quad (10)$$

$$\mathsf{M}(\mathbf{a},\mathsf{r},\mathsf{z}_{\mathsf{n}+1} \mid \mathsf{F}_{\mathsf{S}}, \mathbf{z}_{\mathsf{n}}) \equiv \mathsf{F}_{\mathsf{S}}(\mathsf{z}_{\mathsf{n}+1} \mid \mathbf{z}_{\mathsf{S},\mathsf{n}})\mathsf{M}\left(\mathbf{a},\mathsf{r} \mid \mathsf{F}_{\mathsf{S}}, \mathbf{z}_{\mathsf{n}}\right) \quad (11)$$

The variational candidate distribution (11) is a conditional on the source's one-step-ahead data predictor  $F_S$  (described in (7)). (11) is presumed unavailable, i.e. unknown functional form. We desire to find an optimal choice of this variational object which has been constrained by source knowledge  $F_S$ .

We assert in (11) that  $F_S$  is descriptive of the target's next datum  $z_{n+1}$ . Moreover, a conditional independence assumption between  $\{a, r\}$  and  $z_{S,n+1}$  given  $F_S$  is assumed. Equation (12) illustrates the conditional independence assumption made in (11).

$$\mathsf{M}(\mathsf{z}_{\mathsf{n}+1} \mid \mathsf{F}_{\mathsf{S}}, \mathbf{a}, \mathsf{r}, \mathbf{z}_{\mathsf{n}}) \equiv \mathsf{M}(\mathsf{z}_{\mathsf{n}+1} \mid \mathsf{F}_{\mathsf{S}}, \boldsymbol{\psi}_{\mathsf{n}+1}) \equiv \\ \mathsf{F}_{\mathsf{S}}(\mathsf{z}_{\mathsf{n}+1} \mid \boldsymbol{\psi}_{\mathsf{n}+1}, \mathbf{z}_{\mathsf{S},\mathsf{n}})$$
(12)

$$\begin{split} \mathsf{M} \in \mathbf{M} &\equiv \{ \text{models (11) with } \mathsf{F}_{\mathsf{S}}(\mathsf{z}_{\mathsf{n}+1}) \\ & \text{fixed and } \mathsf{M}\left(\mathsf{a},\mathsf{r} \mid \mathsf{F}_{\mathsf{S}},\mathsf{z}_{\mathsf{n}}\right) \text{variational} \} \end{split} \tag{13}$$

With regards to brevity,  $\psi_{n+1}$  will be omitted for the rest of this paper. We now wish to solve the FPD optimisation task (8) instantiated with (10), (11) and (13).

**Proposition 1.** The unknown model belongs to the knowledge constrained set,  $M \in M$  (13), and the ideal model  $M_I$  is (10). Then the FPD-optimal model is

$$\mathsf{M}^{\circ}(\mathbf{a},\mathsf{r},\mathsf{z}_{\mathsf{n}+1}|\mathsf{F}_{\mathsf{S}},\mathbf{z}_{\mathsf{n}}) \propto \mathsf{F}_{\mathsf{S}}(\mathsf{z}_{\mathsf{n}+1})\mathsf{M}^{\circ}(\mathbf{a},\mathsf{r}|\mathsf{F}_{\mathsf{S}},\mathbf{z}_{\mathsf{n}}) \tag{14}$$

where

$$M^{\circ}(\mathbf{a}, \mathbf{r} \mid \mathbf{F}_{\mathsf{S}}, \mathbf{z}_{\mathsf{n}}) \propto F(\mathbf{a}, \mathbf{r} \mid \mathbf{z}_{\mathsf{n}}) \exp\left[\int_{\mathbf{z}_{\mathsf{n}+1}} \ln(F(\mathbf{z}_{\mathsf{n}+1} \mid \mathbf{a}, \mathbf{r})) F_{\mathsf{S}}(\mathbf{z}_{\mathsf{n}+1}) d\mathbf{z}_{\mathsf{n}+1}\right]$$
(15)

Proof. See Appendix A.

The result achieved in proposition 1 is the canonical solution as reported in literature. A full description of a recursive and computationally tractable algorithm for the AR(p) model is identified in proposition 2.

**Proposition 2.** At step n, the FPD-optimal source knowledge conditioned model is (15) and the conjugate prior is of the form (4). Then the appropriate recursive hyperparameter update is

$$\mathbf{V}_{n}^{\circ} = \mathbf{V}_{n} + \begin{bmatrix} \mathbf{m}_{S,n+1} \\ \psi_{n+1} \end{bmatrix} \begin{bmatrix} \mathbf{m}_{S,n+1} & \psi_{n+1}^{\mathsf{T}} \end{bmatrix} + \begin{bmatrix} \mathbf{w}_{S,n+1} & 0 & \cdots & 0 \\ 0 & 0 & \vdots \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & 0 \end{bmatrix}$$
(16)  
$$\nu_{n}^{\circ} = \nu_{n} + 1$$
(17)

Where conjugacy between the prior and posterior is conserved and, the first and seconds moments of  $F_S$ ,  $m_{S,n+1}$  and  $w_{S,n+1}$ respectively, are given in (7).

The resulting computational flow adopts equations (6) for the data-step update procedure with the replacement of  $V_{n-1}$ and  $\nu_{n-1}$  with the optimal posterior  $V_{n-1}^{\circ}$  and  $\nu_{n-1}^{\circ}$  at step n-1. The source knowledge transfer step adopts equations (16) and (17).

**Remark 1.** The source distribution acceptance procedure outlined in proposition 2 processes both the first and second moments of  $F_S$ . This is an unexpected result as the second order moment transfer does not happen in similar FPD instantiations in literature [10]. Despite this, the expected values of the parameters **a** and future data realisations are invariant to the source predictor's uncertainty  $w_{S,n+1}$ , although it does have an effect on the target's confidence in these objects.

We will now investigate an alternative algorithm for the transfer of knowledge with motivations to ensure robustness to source knowledge uncertainty.

## V. ROBUST FPD-OPTIMAL BTL

A new specification of ideal distribution  $M_I$  is defined in (18). Similar to (10), it is a joint model on the parameters  $\{a, r\}$  and the target's next datum  $z_{n+1}$ . In contrast to (10), we choose to design our ideal model  $M_I$  as a conditional on the source provided object  $F_S$  where we define  $F_S$  as the source's conditional one-step ahead mean data predictor  $F_S(\mu_{z_{n+1}}|\mathbf{z}_{S,n}, r_{CE,n})$  (20).  $\mu_{z_{n+1}}$  is the first parameter (mean) of

the AR(p) generative model in (1). F<sub>S</sub> should not be confused with the one step ahead data predictor described in (7) and adopted in section IV.

$$\mathsf{M}_{\mathsf{I}}(\mathbf{a},\mathsf{r},\mathsf{z}_{\mathsf{n}+1} \mid \mathsf{F}_{\mathsf{S}},\mathbf{z}_{\mathsf{n}}) \equiv \mathsf{F}_{\mathsf{S}}(\mathsf{z}_{\mathsf{n}+1})\mathsf{F}\left(\mathbf{a},\mathsf{r} \mid \mathbf{z}_{\mathsf{n}}\right) \qquad (18)$$

$$\mathsf{M}(\mathbf{a},\mathsf{r},\mathsf{z}_{n+1}\mid\mathsf{F}_{\mathsf{S}},\mathbf{z}_{n}) \equiv \mathsf{F}(\mathsf{z}_{n+1}\mid\mathbf{a},\mathsf{r})\mathsf{M}\left(\mathbf{a},\mathsf{r}\mid\mathsf{F}_{\mathsf{S}},\mathbf{z}_{n}\right) \quad (19)$$

$$F_{\mathsf{S}}(\mu_{\mathsf{z}_{\mathsf{n}+1}} \mid \mathbf{z}_{\mathsf{S},\mathsf{n}},\mathsf{r}_{\mathsf{CE},\mathsf{n}}) = \mathcal{N}_{\mu_{\mathsf{z}_{\mathsf{n}+1}}}(\psi_{\mathsf{n}+1}^{\mathsf{T}} \mathbf{V}_{\mathsf{a}\mathsf{a},\mathsf{n}}^{-1} \mathbf{v}_{\mathsf{a}\mathsf{1},\mathsf{n}}, \frac{\lambda_{\mathsf{n}}}{\nu_{\mathsf{n}} - \mathsf{p} + 2} \psi_{\mathsf{n}+1}^{\mathsf{T}} \mathbf{V}_{\mathsf{a}\mathsf{a}}^{-1} \psi_{\mathsf{n}+1})$$
(20)

We constrain our beliefs in M by asserting that  $z_{n+1}$  is described by the target's generative model only. When comparing (18) and (19) to (10) and (11), the effective difference is swapping the role of  $F_{S}$  and the generative model respectively.

$$M \in \mathbf{M} \equiv \{ \text{models (19) with } \mathsf{F}(\mathsf{z}_{\mathsf{n}+1} \mid \mathbf{a}, \mathsf{r}) \\ \text{fixed and } \mathsf{M}(\mathbf{a}, \mathsf{r} \mid \mathsf{F}_\mathsf{S}, \mathsf{z}_\mathsf{n}) \text{ variational} \}$$
(21)

The objective is to once again optimally choose M° from the set of distributions (21) via the minimisation of (8) with the user-specified ideal (18).

**Proposition 3.** The unknown model belongs to the knowledge constrained set,  $M \in M$  (21), and the ideal model  $M_1$  is (18). Then the FPD-optimal model is

$$\mathsf{M}^{\circ}(\mathbf{a},\mathsf{r},\mathsf{z}_{\mathsf{n}+1} | \mathsf{F}_{\mathsf{S}},\mathsf{z}_{\mathsf{n}}) \propto \mathsf{F}_{\mathsf{S}}(\mathsf{z}_{\mathsf{n}+1}) \mathsf{M}^{\circ}(\mathbf{a},\mathsf{r} | \mathsf{F}_{\mathsf{S}},\mathsf{z}_{\mathsf{n}}) \tag{22}$$

where

$$\begin{split} \mathsf{M}^{\circ}(\mathbf{a},\mathsf{r} \mid \mathsf{F}_{\mathsf{S}}, \mathbf{z}_{\mathsf{n}}) &\propto \\ \mathsf{F}(\mathbf{a},\mathsf{r} \mid \mathbf{z}_{\mathsf{n}}) \exp\left[-\mathcal{D}\Big(\mathsf{F}(\mathsf{z}_{\mathsf{n}+1} \mid \mathbf{a},\mathsf{r})\Big\|\mathsf{F}_{\mathsf{S}}(\mathsf{z}_{\mathsf{n}+1})\Big)\right] \quad \ \ (23) \\ \textit{pof. See Appendix C.} & \Box \end{split}$$

Proof. See Appendix C.

Proposition 3 outlines a conical solution for transfer learning problem where both the ideal and candidate distributions are conditioned on the source knowledge F<sub>S</sub>. The next task is to identify a computationally tractable inference flow for the updating of belief at each data step. Unfortunately, proposition 3 does not result in a conjugate updating procedure when using the  $\mathcal{N}$ iG distribution. We now propose the use of an extended  $\mathcal{N}$ iG distribution to model belief in the parameters as follows.

$$M(\mathbf{a}, \mathbf{r} \mid \mathbf{V}, \mathbf{B}, \nu, \mathbf{k}) \propto \mathbf{r}^{-\frac{\nu}{2}} \exp \left[ -\frac{1}{2\mathbf{r}} \begin{bmatrix} -1 & \mathbf{a}^{\mathsf{T}} \end{bmatrix} \mathbf{V} \begin{bmatrix} -1 \\ \mathbf{a} \end{bmatrix} -\frac{1}{2} \begin{bmatrix} -1\mathbf{a}^{\mathsf{T}} \end{bmatrix} \mathbf{B} \begin{bmatrix} -1 \\ \mathbf{a} \end{bmatrix} -\mathbf{r}\mathbf{k} \right]$$
(24)

A full investigation of (24) is out of scope of this paper. Regardless, it possess a number of key properties.

- 1) When diffuse choices of hyper-parameters **B** and k are adopted ( $\epsilon I_p$  and  $\epsilon$  respectively, where  $\epsilon$  is a small positive real number), the extended NiG is equivalent to a standard  $\mathcal{N}$ iG with hyper parameters **V**,  $\nu$ .
- 2) When conditioning on a known r, the extended NiGreduces to a multivariate Normal distribution.
- 3) The hyper-parameter update procedure when processing local data  $\mathbf{z}_n$  is unchanged (see (6)).

Considering Property 1 and Property 3, if the extended NiGis adopted in sequential updating in an isolated learner when diffuse priors are adopted, the inference task is identical to traditional inference described in section II.

For the purposes of deriving a robust transfer learning algorithm, we are predominately concerned with producing a marginal distribution on a. Unfortunately, this object is believed to be analytically intractable. This issue is addressed via the adoption of a certainty equivalent choice of r (i.e.  $r_{CE}$ ) which adopts E[r] (5) as a known parameter. Considering Property 2, we can derive a r<sub>CE</sub> conditional on **a** and is provided in (25) where  $\mathbf{\bar{V}}_{n} = \mathbf{V}_{n} + r_{CE}\mathbf{B}_{n}$ .

$$\begin{split} \mathsf{M}(\mathbf{a} \mid \mathsf{r}_{\mathsf{CE}}, \mathbf{\bar{V}}_{n}, \mathbf{B}_{n}, \nu_{n}, \mathsf{k}_{n}) &= \mathcal{N}(\mathbf{\bar{V}}_{\mathsf{aa},n}^{-1} \mathbf{v}_{\mathsf{a1},n}, \mathsf{r}_{\mathsf{CE}} \mathbf{\bar{V}}_{\mathsf{aa},n}^{-1}) \\ \mathsf{E}[\mathsf{z}_{\mathsf{n}+1} \mid \mathsf{r}_{\mathsf{CE}}] &= \psi_{\mathsf{n}+1}^{\mathsf{T}} \mathbf{\bar{V}}_{\mathsf{aa},n}^{-1} \mathbf{v}_{\mathsf{a1},n} \\ \mathsf{VAR}[\mathsf{z}_{\mathsf{n}+1} \mid \mathsf{r}_{\mathsf{CE}}] &= \mathsf{r}_{\mathsf{CE}} \psi_{\mathsf{n}+1}^{\mathsf{T}} \mathbf{\bar{V}}_{\mathsf{aa},n}^{-1} \psi_{\mathsf{n}+1} \end{split} \tag{25}$$

Note that r<sub>CE</sub> must be generated at each step using isolated target knowledge only. Equation (25) is not a complete solution as it does not provide a description on r but is sufficient for making inferences on the parameters **a** and generating the expected values of future data. We now propose a recursive and tractable update procedure for the extended NiG distribution.

Proposition 4. At step n, the optimal source knowledge conditioned model is (23) and the conjugate prior is (24). Then the appropriate recursive hyper-parameter update is

$$\begin{split} \boldsymbol{\nu}_{n}^{\circ} &= \boldsymbol{\nu}_{n} + 2\\ \boldsymbol{k}_{n}^{\circ} &= \boldsymbol{k}_{n} + \frac{1}{\boldsymbol{w}_{\mathsf{S},\mathsf{n}+1}}\\ \boldsymbol{V}_{n}^{\circ} &= \boldsymbol{V}_{n} \end{split} \tag{26} \\ \boldsymbol{B}_{n}^{\circ} &= \boldsymbol{B}_{n} + \frac{1}{\boldsymbol{w}_{\mathsf{S},\mathsf{n}+1}} \begin{bmatrix} \boldsymbol{m}_{\mathsf{S},\mathsf{n}+1} \\ \boldsymbol{\psi}_{\mathsf{n}+1} \end{bmatrix} \begin{bmatrix} \boldsymbol{m}_{\mathsf{S},\mathsf{n}+1} & \boldsymbol{\psi}_{\mathsf{n}+1}^{\mathsf{T}} \end{bmatrix} \end{split}$$

Where  $m_{S,n+1}$  and  $w_{S,n+1}$  are the first and second moments of (20) respectively.

*Proof.* See Appendix D. 
$$\Box$$

Adopting Property 1 of the extended NiG model, the resulting computational flow is defined in (6) for the data-step update with the replacement of  $V_{n-1}$  and  $\nu_{n-1}$  with the optimal posterior  $\mathbf{V}_{n-1}^{\circ}$  and  $\nu_{n-1}^{\circ}$ . Equation (26) describes the source knowledge update. (25) is adopted for deducing moments.

## VI. EXPERIMENTATION

To illustrate the performance of the hitherto described FPDoptimal BTL algorithms, we design an experimental setting with a target AR(8) inference machine and a single source making available a distribution on the targets next datum at every step. We compare the following methods: No Transfer (NT); FPD-Optimal BTL (FPD-BTL); Robust FPD-Optimal BTL (RFPD-BTL). n = 30 local target data are produced using the AR(8) parameters below with a driving Gaussian innovation with variance r = 5.

$$\mathbf{a} = \begin{bmatrix} 0.175, & -0.126, & 0.067, & -0.035, \\ & \dots 0.014, & -0.007, & 0.003, & -0.0001 \end{bmatrix}$$



Fig. 2. Mean Absolute Prediction Error of the target AR(8) inference machine over the last 10 data with varying source prediction accuracy  $\rho_{\text{S}}$ . Methods compared: No Transfer (NT), FPD optimal BTL (FPD-BTL) and Robust FPD optimal BTL (RFPD-BTL). 30 target data considered with variance r = 5.

At each step the target processes its local datum, then accepts a distribution  $F_S = \mathcal{N}(\mu_{z_{n+1}}, \rho_S)$  from the source (simulating the one-step-ahead data predictor and conditional one-step-ahead mean data predictor as appropriate).  $\mu_{z_{n+1}}$  is sampled at each step from  $\mathcal{N}(\psi_{n+1}^T \mathbf{a}, \rho_S)$ .  $\rho_S$  is the operating condition which determines the source predictor variance and a substitute for the second parameter of the transferred distribution (20) and the second moment of (7).  $\rho_S$  dictates the quality of the source knowledge being transferred.

For each value of  $\rho_{S}$ , we plot in figure 1 the Mean Absolute Prediction Error (MAPE) of the target's one-step-ahead data predictor on the last 10 data points, over 5000 Monte-Carlo trials. The MAPE for a single experiment is defined in equation (27) where  $\hat{z}_{i+1}$  is the predicted value of the next datum by the model at step i and  $z_{i+1}$  is the true realisation.

$$\mathsf{MAPE} = \frac{1}{10} \sum_{i=n-10}^{n} |\hat{\mathbf{z}}_{i+1} - \mathbf{z}_{i+1}| \tag{27}$$

#### VII. DISCUSSION

The isolated NT AR task acts a baseline for accessing the performance of the BTL algorithms. Unsurprisingly, its MAPE is invariant with  $\rho_S$ . In the low source data predictor variance regime ( $\rho_S < 20$ ), both FPD algorithms perform well and provide positive knowledge transfer. For ( $\rho_S < r$ ), RFPD-BTL outperforms FPD-BTL due to its ability to accept source knowledge weighted above locally sourced data. For  $\rho_S \approx r$ , FPD-BTL marginally outperforms RFPD-BTL. This is because the two algorithms are effectively equivalent when  $\rho_S = r$  (with respect to expected data prediction) except RBTL-FPD has the added overhead of estimating  $r_{CE}$  (directly proportional to source knowledge acceptance weight). In the stressful regime of  $\rho_S > 20$ , RFPD-BTL achieves good levels of robustness while FPD-BTL does not. The RFPD-BTL algorithm accepts the expected value of the source mean data predictor as a locally observed datum weighted by  $\frac{r_{CE}}{w_S}$ . When  $r_{CE} \ll w_S$ , we assign little weight to the source knowledge. When  $r_{CE} > w_S$ , we accept the source knowledge with a weight greater than 1 (i.e. weighting it higher than a locally observed datum). A key disadvantage to RFPD-BTL is the believed analytical intractability of the extended  $\mathcal{N}iG$  which omits the possibility for a full analytical Bayesian description of the source conditioned posterior.

The BTL technique presented in this paper has potentially widespread application in contexts of distributed and/or multioutput filtering and prediction, assuming that individual observation channels are amenable to all-pole modelling [16]. This includes microphone arrays and other networks of sensors in (low-noise) acoustic and vibrational fields. Recall that our framework involves knowledge transfer between local AR tasks where the interaction model does not need to be specified. It can be replace and robustify complete-modelling approaches to multi-output AR processing, such as the conventional vector-AR [17] process. More work will be required to assess the performance of our BTL framework for transfer from multiple AR sources. It will be interesting to test whether BTL can avoid sensitivity to (spatial) model mismatch.

#### VIII. CONCLUSION

This paper has presented an investigation of Bayesian Transfer Learning (BTL) in an incompletely modelled scenario. BTL is typically approached via the assumption or design of a complete model, allowing Bayesian conditioning between learners. We focus on the incomplete model scenario where this complete model is not designed or assumed. Two algorithms are derived for the AR BTL task via FPD principles, for the transfer of source data predictive knowledge. FPD-BTL adopts a standard FPD setting but is found to be effectively invariant to source uncertainty. We propose the RFPD-BTL instantiation which results in a robust algorithm and is validated in Monte-Carlo simulations. In future work, a formalisation of the RFPD-BTL algorithm will be investigated which does not require the adoption of a certainty equivalent,  $r_{CE}$ , in (25).

## APPENDIX

### A. Proof Proposition 1

with differential entropy of  $\mathsf{F}_{\mathsf{S}}$  and normalising constant  $c_{\mathsf{M}^\circ}$ 

$$\mathcal{H}_{F_{S}} = -\int F_{S}\left(z_{n+1} \mid \mathbf{z}_{S,n}\right) \ln F_{S}\left(z_{n+1} \mid \mathbf{z}_{S,n}\right) dz_{n+1}$$

$$\begin{split} c_{M^{\circ}} &= \int F\left(\boldsymbol{a}, r \mid \boldsymbol{z}_{n}\right) \\ &\times \exp\left\{\int \ln F\left(z_{n+1} \mid \boldsymbol{a}, r\right) F_{S}\left(z_{n+1} \mid \boldsymbol{z}_{S, n}\right) dz_{n+1}\right\} d\boldsymbol{a} dr \end{split}$$

### B. Proof Proposition 2

$$\begin{split} \mathsf{M}^{\circ}(\mathbf{a},\mathsf{rF}_{\mathsf{S}},\!\mathbf{z}_{\mathsf{n}}) &\propto \mathcal{N}\mathrm{iG}_{\mathsf{a},\mathsf{r}}(\mathbf{V}_{\mathsf{n}},\nu_{\mathsf{n}}) \\ & \exp\left[\int_{\mathsf{Z}_{\mathsf{n}+1}} \ln \mathcal{N}_{\mathsf{Z}_{\mathsf{n}+1}}(\boldsymbol{\psi}_{\mathsf{n}+1}^{\mathsf{T}}\mathbf{a},\mathsf{r})\mathsf{St}_{\mathsf{Z}_{\mathsf{n}+1}}(\cdot)\mathsf{d}\mathsf{Z}_{\mathsf{n}+1}\right] \quad (28) \end{split}$$

Considering the Boltzmann modulation term.

$$\begin{split} &\int_{z_{n+1}} \left( \ln \frac{1}{\sqrt{2\pi r}} - \frac{1}{2r} (z_{n+1} - \psi_{n+1}^{\mathsf{T}} \mathbf{a})^2 \right) \mathsf{St}_{z_{n+1}} (\cdot) \mathsf{d} z_{n+1} \\ &= \ln \frac{1}{\sqrt{2\pi r}} - \frac{1}{2r} \Big( \psi_{n+1}^{\mathsf{T}} \mathbf{a} \mathbf{a}^{\mathsf{T}} \psi_{n+1} - 2 \psi_{n+1}^{\mathsf{T}} \mathbf{a} \mathsf{m}_{\mathsf{S}, n+1} \\ &+ \int_{z_{n+1}} z_{n+1}^2 \mathsf{St}_{z_{n+1}} (\cdot) \mathsf{d} z_{n+1} \Big) \\ &= \ln \frac{1}{\sqrt{2\pi r}} - \frac{1}{2r} \Big( (\psi_{n+1}^{\mathsf{T}} \mathbf{a} - \mathsf{m}_{\mathsf{S}, n+1})^2 + \mathsf{w}_{\mathsf{S}, n+1} \Big) \end{split}$$
(29)

Adopting (28), (29) and (4) where  $m_{S,n+1}$  and  $w_{S,n+1}$  are moments of  $St_{z_{n+1}}$  and can be generated via (7).

$$\begin{split} \mathsf{M}^{\circ}(\boldsymbol{a},\boldsymbol{r}|\;\mathsf{F}_{\mathsf{S}},\boldsymbol{z}_{\mathsf{n}}) &\propto \boldsymbol{r}^{-\frac{\nu_{\mathsf{n}}+1}{2}} \exp\left(-\frac{1}{2\mathsf{r}}\begin{bmatrix}-1 & \boldsymbol{a}^{\mathsf{T}}\end{bmatrix}\boldsymbol{\mathsf{V}}_{\mathsf{n}}\begin{bmatrix}-1\\\boldsymbol{a}\end{bmatrix}\right) \\ &-\frac{1}{2\mathsf{r}}\bigg((\boldsymbol{\psi}_{\mathsf{n}+1}^{\mathsf{T}}\boldsymbol{a}-\mathsf{m}_{\mathsf{S},\mathsf{n}})^{2}+\mathsf{w}_{\mathsf{S},\mathsf{n}+1}\bigg)\bigg) \end{split}$$

## C. Proof of Proposition 3

D. Proof of Proposition 4

$$\begin{split} & \mathsf{M}^{\circ}(\boldsymbol{a},\boldsymbol{r}\mid\mathsf{F}_{\mathsf{S}},\boldsymbol{z}_{n}) \propto \\ & \mathsf{F}(\boldsymbol{a},\boldsymbol{r}\mid\boldsymbol{z}_{n}) \exp\left[-\mathcal{D}\big(\mathsf{F}(\boldsymbol{z}_{n+1}\mid\boldsymbol{a},\boldsymbol{r})\|\mathsf{F}_{\mathsf{S}}(\boldsymbol{z}_{n+1})\big)\right] \end{split} \tag{30}$$

Considering the exponential modulation term only.

$$\exp\left[-\frac{1}{2\mathsf{w}_{\mathsf{S},\mathsf{n}+1}}\left(\mathsf{r}+(\psi_{\mathsf{n}+1}^{\mathsf{T}}\mathbf{a}-\mathsf{m}_{\mathsf{S},\mathsf{n}+1})^{2}\right)-\ln(\mathsf{r})\right]$$
$$=\mathsf{r}^{-1}\exp\left[-\frac{1}{2}\left(\mathsf{r}\frac{1}{\mathsf{w}_{\mathsf{S},\mathsf{n}+1}}+\right.\right.$$
$$\left[-1\quad\mathbf{a}^{\mathsf{T}}\right]\frac{1}{\mathsf{w}_{\mathsf{S},\mathsf{n}+1}}\left[\frac{\mathsf{m}_{\mathsf{S},\mathsf{n}+1}}{\psi_{\mathsf{n}+1}}\right]\left[\mathsf{m}_{\mathsf{S},\mathsf{n}+1}\quad\psi_{\mathsf{n}+1}^{\mathsf{T}}\right]\left[-1\atop \mathbf{a}\right]\right)\right]$$

Adopting (24) and (30).

$$\begin{split} \mathsf{M}^{\circ}(\mathbf{a},\mathbf{r} \mid \mathsf{F}_{\mathsf{S}},\mathbf{z}_{\mathsf{n}}) &\propto \\ \mathsf{r}^{-\frac{\nu+2}{2}} \exp\left[-\frac{1}{2\mathsf{r}} \begin{bmatrix} -1 & \mathbf{a}^{\mathsf{T}} \end{bmatrix} \mathbf{V}_{\mathsf{n}} \begin{bmatrix} -1 \\ \mathbf{a} \end{bmatrix} - \frac{1}{2} \begin{bmatrix} -1 & \mathbf{a}^{\mathsf{T}} \end{bmatrix} \mathbf{B}_{\mathsf{n}} \begin{bmatrix} -1 \\ \mathbf{a} \end{bmatrix} \\ &-\frac{1}{2} \begin{bmatrix} -1 & \mathbf{a}^{\mathsf{T}} \end{bmatrix} \frac{1}{\mathsf{w}_{\mathsf{S},\mathsf{n}+1}} \begin{bmatrix} \mathsf{m}_{\mathsf{S},\mathsf{n}+1} \\ \psi_{\mathsf{n}+1} \end{bmatrix} \begin{bmatrix} \mathsf{m}_{\mathsf{S},\mathsf{n}+1} & \psi_{\mathsf{n}+1}^{\mathsf{T}} \end{bmatrix} \begin{bmatrix} -1 \\ \mathbf{a} \end{bmatrix} \\ &-\frac{1}{\mathsf{w}_{\mathsf{S},\mathsf{n}+1}} \mathsf{r} \end{bmatrix} \end{split}$$

#### REFERENCES

- Torrey, L., & Shavlik, J. (2010). Transfer learning. In Handbook of research on machine learning applications and trends: algorithms, methods, and techniques (pp. 242-264). IGI global.
- [2] Pan, S. J., & Yang, Q. (2009). A survey on transfer learning. IEEE Transactions on knowledge and data engineering, 22(10), 1345-1359.
- [3] Weiss, K., Khoshgoftaar, T. M., & Wang, D. (2016). A survey of transfer learning. Journal of Big data, 3(1), 9.
- [4] Marx, Z., Rosenstein, M. T., Kaelbling, L. P., & Dietterich, T. G. (2005). Transfer learning with an ensemble of background tasks. Inductive Transfer, 10.
- [5] Karbalayghareh, A., Qian, X., & Dougherty, E. R. (2018). Optimal Bayesian transfer learning. IEEE Transactions on Signal Processing, 66(14), 3724-3739.
- [6] Kárný, M., & Guy, T. V. (2006). Fully probabilistic control design. Systems & Control Letters, 55(4), 259-265.
- [7] Kárný, M. (1996). Towards fully probabilistic control design. Automatica, 32(12), 1719-1722.
- [8] Kárný, M., & Kroupa, T. (2012). Axiomatisation of fully probabilistic design. Information Sciences, 186(1), 105-113.
- [9] Papež, M., & Quinn, A. (2019, December). Bayesian transfer learning between Gaussian process regression tasks. In 2019 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT) (pp. 1-6). IEEE.
- [10] Foley, C., & Quinn, A. (2017). Fully probabilistic design for knowledge transfer in a pair of Kalman filters. IEEE Signal Processing Letters, 25(4), 487-490.
- [11] Papež, M., & Quinn, A. (2020). Bayesian transfer learning between Student-t filters. Signal Processing, 107624.
- [12] Ephraim, Y., & Roberts, W. J. (2005). Revisiting autoregressive hidden Markov modeling of speech signals. IEEE Signal processing letters, 12(2), 166-169.
- [13] Pacurar, M. (2008). Autoregressive conditional duration models in finance: a survey of the theoretical and empirical literature. Journal of economic surveys, 22(4), 711-751.
- [14] Bernardo, J. M., & Smith, A. F. (2009). Bayesian theory (Vol. 405). John Wiley & Sons.
- r [15] Kay, S. M. (1993). Fundamentals of statistical signal processing. Prentice Hall PTR.
  - [16] Tiao, G. C., & Box, G. E. (1981). Modeling multiple time series with applications. journal of the American Statistical Association, 76(376), 802-816.
  - [17] Corander, J., & Villani, M. (2006). A Bayesian approach to modelling graphical vector autoregressions. Journal of Time Series Analysis, 27(1), 141-156.
  - [18] Šmídl, V., & Quinn, A. (2006). The variational Bayes method in signal processing. Springer Science & Business Media.