
Bayesian Deconditional Kernel Mean Embeddings

Kelvin Hsu^{1 2} Fabio Ramos^{1 3}

Abstract

Conditional kernel mean embeddings form an attractive nonparametric framework for representing conditional means of functions, describing the observation processes for many complex models. However, the recovery of the original underlying function of interest whose conditional mean was observed is a challenging inference task. We formalize deconditional kernel mean embeddings as a solution to this inverse problem, and show that it can be naturally viewed and used as a nonparametric Bayes' rule. Critically, we introduce the notion of task transformed Gaussian processes and establish deconditional kernel means embeddings as their posterior predictive mean. This connection provides Bayesian interpretations and uncertainty estimates for deconditional kernel means, explains their regularization hyperparameters, and provides a marginal likelihood for kernel hyperparameter learning. They further enable practical applications such as learning sparse representations for big data and likelihood-free inference.

1. Introduction

Observations of complex phenomena often lead to likelihoods that are described by a conditional mean. A widely applicable setting where this occurs is collecting observations under uncertain inputs, where the task is to learn a function $f : \mathcal{X} \rightarrow \mathbb{R}$ to model a real-valued response z as a function of inputs $x \in \mathcal{X}$ without being able to query or measure x directly to observe this phenomenon. Instead, another measured input $y \in \mathcal{Y}$ relates to x through $p(x|y)$. Consequently, given y , the response Z has mean $g(y) := \mathbb{E}[f(X)|Y = y]$, where g is called the conditional mean of f . Furthermore, $p(x|y)$ is often only available as sample pairs $\{x_i, y_i\}_{i=1}^n$, from simulations, algorithms, or separate experiments, making recovery of latent functions f from conditional means g a challenging inference task.

¹University of Sydney ²CSIRO, Sydney ³NVIDIA, Seattle.
Correspondence to: Kelvin Hsu <Kelvin.Hsu@sydney.edu.au>.

Our first contribution begins with formulating **deconditional mean embeddings (DMEs)** as solutions to this inference problem by building upon the framework of **conditional mean embeddings (CMEs)** (Song et al., 2013). We show that the **DME** can be established as a nonparametric Bayes' rule in the **reproducing kernel Hilbert space (RKHS)** and used for likelihood-free Bayesian inference. In contrast to **kernel Bayes' rule (KBR)** (Fukumizu et al., 2013) which uses third order tensors that can result in vanishing priors, **DMEs** use second order tensors and avoids this problem.

Together with **CMEs** and **KBR**, **DMEs** form a critical part of the **kernel mean embedding (KME)** (Muandet et al., 2017) framework, where probabilistic rules can be represented nonparametrically as operators that are linear in the **RKHS**. This greatly simplifies probabilistic inference without requiring parametrized distributions and compromising flexibility.

Despite this connection, there are elements unique to the **KME** framework that cannot be interpreted or solved via the parallel between probability rules and **RKHS** mean operations. Similar to empirical forms for **KBR** and **CMEs**, empirical **DMEs** are obtained by replacing expectations in its constituent operators with their empirical means, and introduce regularization for operator inverses to relax **RKHS** assumptions, instead of as the optimal solution to a particular loss. Setting regularization hyperparameters is difficult in practice without an appropriate loss for the inference task. Furthermore, similar to **KBR**, the nonparametric Bayes' rule provided by **DMEs** is a statement between observed (or simulated) variables and not on latent functions or quantities. Consequently, uncertainty estimation in inference of latent functions f still require a separate Bayesian formulation.

Our second contribution establishes a Bayesian view of **DMEs** as posterior predictive means of the **task transformed Gaussian process (TTGP)**, a novel nonparametric Bayesian model that recover latent relationships between variables without observing them jointly. **TTGPs** are so named because we show that they are a type of transformed Gaussian process (Murray-Smith & Pearlmutter, 2005) where the transformations and noise covariances are learned, by transforming one **Gaussian process (GP)** task to another, rather than designed from expert knowledge. We use this connection to derive posterior and predictive uncertainty estimates for **DMEs** and explain their regularization hyperparameters

as a function of noise variance. Finally, we derive marginal likelihoods and their scalable computational forms to learn **DME** hyperparameters, which can also be applied to learn inducing points for sparse representations as a special case. All proofs are in the supplementary material.

2. Kernel Mean Embeddings

We begin with an overview of the **KME** framework from which **DMEs** are built upon. **KMEs** are an arsenal of techniques concerned with representations and transformations of function expectations under highly flexible distributions. They consider functions that lie within **RKHSs** \mathcal{H}_k and \mathcal{H}_ℓ , formed by positive definite kernels $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. The **RKHSs** \mathcal{H}_k and \mathcal{H}_ℓ are the closure span of the features $\phi(x) = k(x, \cdot)$ and $\psi(y) = \ell(y, \cdot)$ across $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ respectively, endowed with the inner products $\langle \cdot, \cdot \rangle_k \equiv \langle \cdot, \cdot \rangle_{\mathcal{H}_k}$ and $\langle \cdot, \cdot \rangle_\ell \equiv \langle \cdot, \cdot \rangle_{\mathcal{H}_\ell}$.

The key object is the mean embedding of a distribution $\mu_X := \mathbb{E}[k(X, \cdot)] \in \mathcal{H}_k$. They encode function expectations in the sense that $\mathbb{E}[f(X)] = \langle \mu_X, f \rangle_k$, due to the reproducing property that $\langle k(x, \cdot), f \rangle_k = f(x)$ for all $f \in \mathcal{H}_k$.

Higher ordered mean embeddings are vital components of the framework. Specifically, second order mean embeddings such as $C_{YY} := \mathbb{E}[\ell(Y, \cdot) \otimes \ell(Y, \cdot)] \in \mathcal{H}_\ell \otimes \mathcal{H}_\ell$ and $C_{XY} := \mathbb{E}[k(X, \cdot) \otimes \ell(Y, \cdot)] \in \mathcal{H}_k \otimes \mathcal{H}_\ell$ can be identified as cross-covariance operators $C_{YY} : \mathcal{H}_\ell \rightarrow \mathcal{H}_\ell$ and $C_{XY} : \mathcal{H}_\ell \rightarrow \mathcal{H}_k$ that serve as building blocks of **CMEs** and **DMEs**.

In practical scenarios where only *iid* samples $\{x_i, y_i\}_{i=1}^n$ that are realizations of $(X_i, Y_i) \sim \mathbb{P}_{XY}$ for $i \in \{1, \dots, n\}$ are available, the **KME** framework becomes attractive for nonparametric inference because core objects only require expectations under distributions. Consequently, they can be estimated via empirical means as $\hat{\mu}_X := \frac{1}{n} \sum_{i=1}^n k(x_i, \cdot)$, $\hat{C}_{YY} := \frac{1}{n} \sum_{i=1}^n \ell(y_i, \cdot) \otimes \ell(y_i, \cdot)$, and $\hat{C}_{XY} := \frac{1}{n} \sum_{i=1}^n k(x_i, \cdot) \otimes \ell(y_i, \cdot)$ (Muandet et al., 2017).

For feature matrices, we stack features by columns $\Phi := [\phi(x_1) \ \dots \ \phi(x_n)]$ and $\Psi := [\psi(y_1) \ \dots \ \psi(y_n)]$. We write gram matrices as $K := \Phi^T \Phi$ and $L := \Psi^T \Psi$, where the (i, j) -th element of $A^T B$ is the inner product of the i -th column of A with the j -th column of B . That is, $K_{ij} = \phi(x_i)^T \phi(x_j)$ and $L_{ij} = \psi(y_i)^T \psi(y_j)$. When columns are elements of **RKHSs** such as when $\phi(x) = k(x, \cdot)$ in Φ and $\psi(y) = \ell(y, \cdot)$ in Ψ , the notation $(\cdot)^T(\cdot)$ is a shorthand for the corresponding **RKHS** inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ when it is clear from context what \mathcal{H} is. For example, $f^T h$ is shorthand for $\langle f, h \rangle_k$ if $f, h \in \mathcal{H}_k$. Another common usage is $\Phi^T f = \{\phi(x_i)^T f\}_{i=1}^n = \{k(x_i, \cdot)^T f\}_{i=1}^n = \{\langle k(x_i, \cdot), f \rangle_k\}_{i=1}^n = \{f(x_i)\}_{i=1}^n =: \mathbf{f}$. For summing outer products, we write $\hat{C}_{YY} = \frac{1}{n} \Psi \Psi^T$ and $\hat{C}_{XY} = \frac{1}{n} \Phi \Psi^T$. Note that we use non-bold letters for single points x and y , even though they are often multivariate in practice.

3. Conditional Kernel Mean Embeddings

We now present **CMEs** in a fashion that focuses on their operator properties. By reviewing **CMEs** this way, parallels and contrast with **DMEs** in the subsequent section 4 become more apparent. Importantly, instead of defining **CMEs** via an explicit form, we begin by forming problem statements.

Definition 3.1 (Conditional Mean Problem Statement). Given a function $f : \mathcal{X} \rightarrow \mathbb{R}$, infer the function $g : \mathcal{Y} \rightarrow \mathbb{R}$ such that $g(y) = \mathbb{E}[f(X)|Y = y] \equiv \mathbb{E}_{X|Y}[f](y)$. We call g the *conditional mean* of f with respect to $\mathbb{P}_{X|Y}$ and write the shorthand $g = \mathbb{E}_{X|Y}[f] = \mathbb{E}[f(X)|Y = \cdot]$.

This naturally leads to the notion of operators that map functions f to their conditional means $g = \mathbb{E}[f(X)|Y = \cdot]$.

Definition 3.2 (Conditional Mean Operators). The *conditional mean operator* (**CMO**) $C_{X|Y} : \mathcal{H}_\ell \rightarrow \mathcal{H}_k$ corresponding to $\mathbb{P}_{X|Y}$ is the operator that satisfies

$$(C_{X|Y})^T f = \mathbb{E}[f(X)|Y = \cdot], \quad \forall f \in \mathcal{H}_k, \quad (3.1)$$

where $(C_{X|Y})^T : \mathcal{H}_k \rightarrow \mathcal{H}_\ell$ denotes the adjoint of $C_{X|Y}$.

Depending on the nature of ℓ , unique solutions exist.

Theorem 3.1 ((Fukumizu et al., 2004)). Assume that $\ell(y, \cdot) \in \text{image}(C_{YY})$ for all $y \in \mathcal{Y}$. The *conditional mean operator* (**CMO**) $C_{X|Y}$ is unique and given by

$$C_{X|Y} = C_{XY} C_{YY}^{-1}. \quad (3.2)$$

The assumption that $\ell(y, \cdot) \in \text{image}(C_{YY})$ for all $y \in \mathcal{Y}$ is commonly relaxed by introducing a regularization hyperparameter $\lambda > 0$ to the inverse, so that the **CMO** is replaced with $C_{XY}(C_{YY} + \lambda I)^{-1}$ (Song et al., 2013).

Contrary to definition 3.2, it is more common in the literature to define the **CMO** as the operator $C_{X|Y}$ that satisfies

$$C_{X|Y} \ell(y, \cdot) = \mathbb{E}[k(X, \cdot)|Y = y], \quad \forall y \in \mathcal{Y}, \quad (3.3)$$

while (3.1) is taken as an immediate property of **CMOs** (Fukumizu et al., 2004). However, due to lemma 3.2, we instead take definition 3.2 as the definition of **CMOs**, emphasizing **CMOs** as solutions to the conditional mean problem, and treat (3.3) as an immediate property.

Lemma 3.2. Statements (3.1) and (3.3) are equivalent.

The **CME** of $\mathbb{P}_{X|Y=y}$ is $\mu_{X|Y=y} := C_{X|Y} \ell(y, \cdot)$, equivalent to querying the **CMO** at a particular input y . Consequently, $\langle \mu_{X|Y=y}, f \rangle_k = \langle C_{X|Y} \ell(y, \cdot), f \rangle_k = \langle \ell(y, \cdot), (C_{X|Y})^T f \rangle_\ell = \langle \ell(y, \cdot), \mathbb{E}_{X|Y}[f] \rangle_\ell = \mathbb{E}_{X|Y}[f](y)$.

Motivated by theorem 3.1, empirical **CMOs** and **CMEs** are defined by estimating their constituents by empirical means.

Definition 3.3 (Empirical Conditional Mean Operator). The empirical **CMO** is $\hat{C}_{X|Y} := \hat{C}_{XY}(\hat{C}_{YY} + \lambda I)^{-1}$, $\lambda > 0$.

Theorem 3.3 ((Song et al., 2009)). *The nonparametric form for $\hat{C}_{X|Y}$ is*

$$\hat{C}_{X|Y} = \Phi(L + n\lambda I)^{-1}\Psi^T. \quad (3.4)$$

The empirical **CME** is then $\hat{\mu}_{X|Y=y} := \hat{C}_{X|Y}\ell(y, \cdot)$. Consequently, with $\ell(y) := \{\ell(y_i, y)\}_{i=1}^n$, an estimate for $\mathbb{E}_{X|Y}[f](y)$ is $\langle f, \hat{\mu}_{X|Y=y} \rangle_k = \langle f, \hat{C}_{X|Y}\ell(y, \cdot) \rangle_k = f^T \Phi(L + n\lambda I)^{-1}\Psi^T \ell(y, \cdot) = \mathbf{f}^T(L + n\lambda I)^{-1}\ell(y)$.

Critically, while empirical **CMOs** (3.4) are estimated from joint samples from the joint distribution \mathbb{P}_{XY} , they only encode the conditional distribution $\mathbb{P}_{X|Y}$. This means that the empirical **CMOs** will encode the same conditional distribution even if the joint distribution \mathbb{P}_{XY} changes but the conditional distribution $\mathbb{P}_{X|Y}$ stays the same. That is, the empirical **CMO** built from joint samples of $p(x, y) = p(x|y)p(y)$ and the empirical **CMO** built from joint samples of $q(x, y) := p(x|y)q(y)$ will encode the same conditional distribution $p(x|y)$ and converge to the same **CMO**.

4. Deconditional Kernel Mean Embeddings

We now present a novel class of **KMEs** referred to as **deconditional mean embeddings** (**DMEs**). They are natural counterparts to **CMEs**. The presentation of definitions and theorems in this section is mainly parallel to section 3. We define the *deconditional mean* problem as the task of recovering latent functions from their conditional means.

Definition 4.1 (Deconditional Mean Problem Statement). Given a function $g : \mathcal{Y} \rightarrow \mathbb{R}$, infer a function $f : \mathcal{X} \rightarrow \mathbb{R}$ such that $g(y) = \mathbb{E}[f(X)|Y = y]$. We call f a *deconditional mean* of g with respect to $\mathbb{P}_{X|Y}$ and write the shorthand $f = \mathbb{E}_{X|Y}^\dagger[g]$.

The deconditional mean of a function g infers the function f whose conditional mean would be g with respect to $\mathbb{P}_{X|Y}$. The corresponding operator that encodes this transformation is the **deconditional mean operator** (**DMO**).

Definition 4.2 (Deconditional Mean Operators). The **deconditional mean operator** (**DMO**) $C'_{X|Y} : \mathcal{H}_k \rightarrow \mathcal{H}_\ell$ corresponding to $\mathbb{P}_{X|Y}$ is the operator that satisfies

$$(C'_{X|Y})^T \mathbb{E}[f(X)|Y = \cdot] = f, \quad \forall f \in \mathcal{H}_k. \quad (4.1)$$

Depending on the nature of ℓ and k , unique solutions exist.

Theorem 4.1. *Assume that $\ell(y, \cdot) \in \text{image}(C_{YY})$ for all $y \in \mathcal{Y}$ and $k(x, \cdot) \in \text{image}(C_{X|Y}C_{YY}(C_{X|Y})^T)$ for all $x \in \mathcal{X}$. The **deconditional mean operator** (**DMO**) $C'_{X|Y}$ is unique and given by*

$$C'_{X|Y} = (C_{X|Y}C_{YY})^T(C_{X|Y}C_{YY}(C_{X|Y})^T)^{-1}. \quad (4.2)$$

Similar to the case with **CMOs** (Song et al., 2013), the assumption that $k(x, \cdot) \in \text{image}(C_{X|Y}C_{YY}(C_{X|Y})^T)$ for

all $x \in \mathcal{X}$ can be relaxed by introducing a regularization hyperparameter $\epsilon > 0$ to the inverse, so that the **DMO** is replaced with $(C_{X|Y}C_{YY})^T(C_{X|Y}C_{YY}(C_{X|Y})^T + \epsilon I)^{-1}$.

Since **DMOs** invert the results of **CMOs**, they can also be understood as pseudo-inverses of **CMOs**.

Theorem 4.2. *If the assumptions in theorem 4.1 hold and further $((C_{X|Y})^T C_{X|Y})^{-1}$ exists such that the pseudo-inverse $C_{X|Y}^\dagger := ((C_{X|Y})^T C_{X|Y})^{-1}(C_{X|Y})^T$ is defined, then **DMOs** are pseudo-inverses of **CMOs** $C'_{X|Y} = C_{X|Y}^\dagger$.*

The **DME** of $\mathbb{P}_{X=x|Y}$ is $\mu'_{X=x|Y} := C'_{X|Y}k(x, \cdot) \in \mathcal{H}_\ell$, equivalent to querying the **DMO** at a particular input x . Consequently, $\langle \mu'_{X=x|Y}, g \rangle_\ell = \langle C'_{X|Y}k(x, \cdot), g \rangle_\ell = \langle k(x, \cdot), (C'_{X|Y})^T g \rangle_k = \langle k(x, \cdot), f \rangle_k = f(x)$.

The form in (4.2) makes it evident that a **DMO** can be fully specified once $C_{X|Y}$ and C_{YY} , encoding the measures $\mathbb{P}_{X|Y}$ and \mathbb{P}_Y respectively, are known. If densities exist, we write them as $p_{X|Y} \equiv p_{X|Y}(\cdot|y)$ and $p_Y \equiv p_Y(\cdot)$, and drop the subscripts in density evaluations as $p(x|y)$ and $p(y)$ whenever the context is clear. Note that $\mathbb{P}_{X=x|Y}$ corresponds to $p_{X|Y}(x|\cdot)$ which is evaluated at x and now a function of y . This is in contrast with $\mathbb{P}_{X|Y=y}$ corresponding to $p_{X|Y}(\cdot|y)$ evaluated at y and now a function of x .

Consider the case where X and Y play the roles of observed and unobserved (latent) variables respectively. The **DMO** considers the conditional $p_{X|Y}$ and the marginal p_Y encoded as $C_{X|Y}$ and C_{YY} (theorem 4.1), and inverts the **CMO** $C_{X|Y}$ (theorem 4.2) with the help of the encoded marginal C_{YY} . This is analogous to the Bayes' rule, where the posterior $p_{Y|X}(\cdot|x) = \frac{p_{X|Y}(x|\cdot)p_Y(\cdot)}{\int_{\mathcal{Y}} p_{X|Y}(x|y)p_Y(y)dy}$ is fully specified by the likelihood $p_{X|Y}$ and prior p_Y . We can then interpret **DMEs** as querying the rule at the observed quantity x while leaving the rule as a function of y for inference. Consequently, we also refer to $C_{X|Y}$ and C_{YY} as the likelihood operator and the prior operator respectively.

The difference between the **DMO** (4.2) and **CMO** (3.2) equations is akin to writing $p_{Y|X}(\cdot|x)$ using Bayes' rule against using the conditional density rule. Compare the **DMO** decomposition (4.2) with the **CMO** decomposition $C_{Y|X} = C_{YX}C_{XX}^{-1} = (C_{XY})^T C_{XX}^{-1}$ in the other direction by reversing the roles of X and Y in (3.2), which would correspond to the posterior $\mathbb{P}_{Y|X}$. The **CMO** is composed of a *joint* operator C_{XY} and an *evidence* operator C_{XX} corresponding to the joint \mathbb{P}_{XY} and evidence \mathbb{P}_X distributions. Similarly, the **DMO** is also composed of a *joint* operator $C_{XY} = C_{X|Y}C_{YY} : \mathcal{H}_\ell \rightarrow \mathcal{H}_k$ and an *evidence* operator $C'_{XX} := C_{X|Y}C_{YY}(C_{X|Y})^T : \mathcal{H}_k \rightarrow \mathcal{H}_k$, but both specified from the likelihood and prior operators.

Motivated by this, we propose to estimate the likelihood and prior operators using separate and independently drawn

samples. The likelihood operator $C_{X|Y}$ is estimated as $\hat{C}_{X|Y}$ (definition 3.3) using *iid* samples $\{x_i, y_i\}_{i=1}^n$, also denoted as $\mathbf{x} := \{x_i\}_{i=1}^n$ and $\mathbf{y} := \{y_i\}_{i=1}^n$. Note that as the likelihood operator is a **CMO**, these joint samples can be from any joint distribution $Q_{XY} \neq P_{XY}$ as long as its conditional distribution is also $P_{X|Y}$. The prior operator C_{YY} is estimated as $\tilde{C}_{YY} := \frac{1}{m} \sum_{j=1}^m \ell(\tilde{y}_j, \cdot) \otimes \ell(\tilde{y}_j, \cdot)$ using another set of *iid* samples $\tilde{\mathbf{y}} := \{\tilde{y}_j\}_{j=1}^m$ from P_Y .

Definition 4.3 (Empirical Deconditional Mean Operator). Let $\epsilon > 0$ be a regularization hyperparameter and define $\hat{C}_{X|Y}$ and \tilde{C}_{YY} as above. The empirical **DMO** is

$$\bar{C}'_{X|Y} := (\hat{C}_{X|Y} \tilde{C}_{YY})^T (\hat{C}_{X|Y} \tilde{C}_{YY} (\hat{C}_{X|Y})^T + \epsilon I)^{-1}. \quad (4.3)$$

The accents denote the set of samples used for estimation. When both sets are used such as in the estimation of the **DMO** $\bar{C}'_{X|Y}$, we denote it with a bar such as $\bar{C}'_{X|Y}$.

Theorem 4.3. The nonparametric form for $\bar{C}'_{X|Y}$ is

$$\bar{C}'_{X|Y} = \tilde{\Psi} [A^T K A + m\epsilon I]^{-1} A^T \Phi^T, \quad (4.4)$$

where $A := (L + n\lambda I)^{-1} \tilde{L}$, $\tilde{L} := \Psi^T \tilde{\Psi}$, and $\tilde{\Psi} := [\psi(\tilde{y}_1) \cdots \psi(\tilde{y}_m)]$.

The empirical **DME** is then $\bar{\mu}'_{X=x|Y} := \bar{C}'_{X|Y} k(x, \cdot)$. Consequently, with $\mathbf{k}(x) := \{k(x_i, x)\}_{i=1}^n$ and $\tilde{\mathbf{g}} := \{g(\tilde{y}_j)\}_{j=1}^m$, an estimate for $\mathbb{E}_{X|Y}^\dagger[g](x)$ is $\langle g, \bar{\mu}'_{X=x|Y} \rangle_\ell = \tilde{\mathbf{g}}^T [A^T K A + m\epsilon I]^{-1} A^T \mathbf{k}(x)$. This motivates the following definitions, where the notation $\tilde{\mathbf{g}}$ is replaced with $\tilde{\mathbf{z}}$, to be interpreted as target observations of g at $\tilde{\mathbf{y}}$.

Definition 4.4 (Nonparametric **DME** Estimator). The nonparametric **DME** estimator, also called the kernel **DME** estimator or the **DME** estimator in function space view, is $\bar{f}(x) = \bar{\alpha}^T \mathbf{k}(x) = \sum_{i=1}^n \bar{\alpha}_i k(x_i, x)$, where $\bar{\alpha} := A [A^T K A + m\epsilon I]^{-1} \tilde{\mathbf{z}}$ and $A := (L + n\lambda I)^{-1} \tilde{L}$. Equivalently, $\bar{f}(x) = \tilde{\mathbf{z}}^T [A^T K A + m\epsilon I]^{-1} A^T \mathbf{k}(x)$. An alternative form is $\bar{f}(x) = \tilde{\mathbf{z}}^T A^T [K A A^T + m\epsilon I]^{-1} \mathbf{k}(x)$.

When features $\phi(x) \in \mathbb{R}^p$ and $\psi(y) \in \mathbb{R}^q$ are finite dimensional, we define the parametric **DME** estimator as follows by rewriting definition 4.4 using the Woodbury identity.

Definition 4.5 (Parametric **DME** Estimator). The parametric **DME** estimator, also called the feature **DME** estimator or the **DME** estimator in weight space view, is $\bar{f}(x) = \bar{\mathbf{w}}^T \phi(x)$, where $\bar{\mathbf{w}} = [\Phi A A^T \Phi^T + m\epsilon I]^{-1} \Phi A \tilde{\mathbf{z}}$ and $A := \Psi^T (\Psi \Psi^T + n\lambda I)^{-1} \tilde{\Psi}$. Equivalently, $\bar{f}(x) = \tilde{\mathbf{z}}^T A^T \Phi^T [\Phi A A^T \Phi^T + m\epsilon I]^{-1} \phi(x)$.

In definition 4.4 (resp. 4.5), computational complexity is dominated by inversions for $L + n\lambda I$ and $A^T K A + m\epsilon I$ (resp. $\Psi \Psi^T + n\lambda I$ and $\Phi A A^T \Phi^T + m\epsilon I$) at $O(n^3)$ and $O(m^3)$ (resp. $O(q^3)$ and $O(p^3)$). For the alternative form in definition 4.4, both inversions are $O(n^3)$, allowing for larger m at $O(m)$ without compromising tractability.

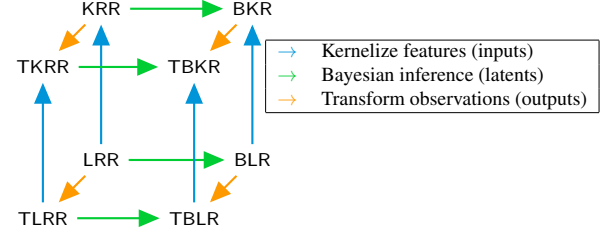


Figure 1. Three dimensions of model extensions (T: Transformed, B: Bayesian, K/L: Kernel/Linear, (R)R: (Ridge) Regression). Kernel extensions (blue) specify feature spaces implicitly through a kernel. Bayesian extensions (green) introduce notions of uncertainty on latent quantities (weights or functions). Finally, transformed extensions (orange) capture indirect function observations.

5. Task Transformed Gaussian Processes

DMEs are constructed as solutions to the task of inferring deconditional means, which are often real-valued functions. Regression problems also address inference of real-valued functions from data. This raises curiosity towards whether **DMEs** can be formulated as solutions to a regression-like problem, and what insights this connection would provide.

In this section, we formulate the task transformed regression problem to provide regression views of **DMEs**. To do this, we first briefly review transformed regression in section 5.1 before we present our contributions in section 5.2.

5.1. Transformed Regression

Standard regression models often assume a Gaussian full data likelihood $p(\mathbf{z}|\mathbf{f}) = \mathcal{N}(\mathbf{z}; \mathbf{f}, \sigma^2 I)$ with targets $\mathbf{z} := \{z_i\}_{i=1}^n \in \mathbb{R}^n$. In the generalized setting when observations of \mathbf{f} at \mathbf{x} are not available but observations of linear combinations thereof are, we can use $p(\tilde{\mathbf{z}}|\mathbf{f}) = \mathcal{N}(\tilde{\mathbf{z}}; M^T \mathbf{f}, \Sigma)$ for some transformation $M \in \mathbb{R}^{n \times m}$ and noise covariance Σ , where $\tilde{\mathbf{z}} := \{\tilde{z}_j\}_{j=1}^m \in \mathbb{R}^m$ are the available observations.

We refer to this setting as *transformed regression*. They can be seen as another dimension of modeling with **linear ridge regression (LRR)** as the base model (fig. 1). **Kernel Ridge Regression (KRR)** is obtained from **LRR** via the kernel trick and Woodbury identity, and they are **maximum a posteriori (MAP)** solutions or predictive means of **Gaussian process regression (GPR)** and **Bayesian linear regression (BLR)** respectively (Rasmussen & Williams, 2006). Consequently, we also refer to **GPR** as **Bayesian kernel regression (BKR)**. Analogous relationships hold between transformed models.

5.2. Task Transformed Regression

We define *task transformed regression (TTR)* as the problem of learning to predict a target variable Z from features X when no direct sample pairs of X and Z are available but instead indirect samples $\{x_i, y_i\}_{i=1}^n$ and $\{\tilde{y}_j, \tilde{z}_j\}_{j=1}^m$ with a

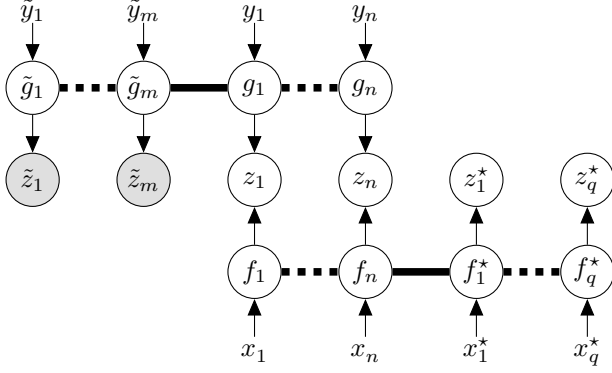


Figure 2. Graphical model (chain graph) for task transformed Gaussian process regression (TTGPR). Circles are random variables. Shaded circles are observed random variables. Undirected edges indicate the GP field, where all the random variables on the field are fully connected to each other (Rasmussen & Williams, 2006). The goal is to infer f^* to predict z^* at x^* , using only a task or original dataset $\{\tilde{y}_j, \tilde{z}_j\}_{j=1}^m$ and a transformation dataset $\{x_i, y_i\}_{i=1}^n$. To connect the two GPs, we posit that the unobserved targets z at x and at y would have been the same if they were observed. Note that like regular GPs, to TTGPs the inputs x and y are not modeled as random variables but treated as index variables instead.

mediating variable Y are available. The name illustrates the idea of transforming the task of regressing Z on Y to learn $g : \mathcal{Y} \rightarrow \mathbb{R}$, using the task or original dataset $\{\tilde{y}_j, \tilde{z}_j\}_{j=1}^m$, to the task of regressing Z on X to learn $f : \mathcal{X} \rightarrow \mathbb{R}$, by mediating the task dataset through the transformation dataset $\{x_i, y_i\}_{i=1}^n$. As the mediating variable Y links the two sets together, y and \tilde{y} control the task transformation.

DME as solution to chained loss We formulate losses for TTR, and establish DMEs as solutions. We begin with the parametric case with $f(x) = \mathbf{w}^T \phi(x)$ and $g(y) = \mathbf{v}^T \psi(y)$.

Theorem 5.1 (Task Transformed LRR (TTLRR)). *The weights of the parametric DME estimator $\hat{f}(x) = \hat{\mathbf{w}}^T \phi(x)$ (definition 4.5) solve chained regularized least square losses,*

$$\begin{aligned} \hat{\mathbf{v}}[\mathbf{w}] &:= \underset{\mathbf{v} \in \mathbb{R}^q}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \phi(x_i) - \mathbf{v}^T \psi(y_i))^2 + \lambda \|\mathbf{v}\|^2, \\ \hat{\mathbf{w}} &:= \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{m} \sum_{j=1}^m (\tilde{z}_j - \hat{\mathbf{v}}[\mathbf{w}]^T \psi(\tilde{y}_j))^2 + \epsilon \|\mathbf{w}\|^2. \end{aligned} \quad (5.1)$$

The notation $\hat{\mathbf{v}}[\mathbf{w}]$ explicitly denotes that $\hat{\mathbf{v}}$ depends on \mathbf{w} . Conceptually, in function space view the first optimization finds g so that f at x best matches with g at y , leading to a solution $\hat{g}[f]$ that is dependent on f . The second finds f so that $\hat{g}[f]$ at \tilde{y} best matches targets \tilde{z} . Using the kernel trick $k(x, x') = \phi(x)^T \phi(x')$, we obtain the nonparametric case.

Lemma 5.2 (Task Transformed KRR (TTKRR)). *The weights of the nonparametric DME estimator $\hat{f}(x) = \hat{\alpha}^T \mathbf{k}(x)$ (definition 4.4) satisfies $\hat{\mathbf{w}} = \Phi \hat{\alpha}$ (the kernel trick).*

DME as posterior predictive mean of TTGP We extend TTLRR and TTKRR to the Bayesian case. This connection reveals that TTR models are transformed regression models with transformations and noise covariances that are learned.

In the parametric case, we have task transformed BLR (TTBLR). We place separate independent Gaussian priors $p(\mathbf{v}) = \mathcal{N}(\mathbf{v}; \mathbf{0}, \beta^2 I)$ and $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \gamma^2 I)$ for g and f respectively. As z is not observed directly from f but only for g , we include noise only for observing g to arrive at likelihoods $p(z|\mathbf{v}) = \mathcal{N}(z; \mathbf{v}^T \psi(y), \sigma^2)$ and $p(z|\mathbf{w}) = \mathcal{N}(z; \mathbf{w}^T \phi(x), 0)$ for g and f respectively.

In the nonparametric case, we have task transformed BKR (TTBKR). We place GP priors $g \sim \mathcal{GP}(0, \ell)$ and $f \sim \mathcal{GP}(0, k)$ on the functions directly. Consequently, TTBKR is also referred to as task transformed Gaussian process regression (TTGPR). Similar to TTBLR, the likelihoods are $p(z|g) = \mathcal{N}(z; g(y), \sigma^2)$ and $p(z|f) = \mathcal{N}(z; f(x), 0)$.

The graphical model for TTGPR is shown in fig. 2. The two GPs for g and f are linked by constraining their targets to be the same at y and x respectively. The GP for g is used to infer the predictive distribution $p(\tilde{z}|\mathbf{z})$, which in turn specifies the overall likelihood $p(\tilde{\mathbf{z}}|\mathbf{f})$ used to infer f . Detailed derivations are provided in the proof of theorem 5.3.

Theorem 5.3 (Task Transformed BLR (TTBLR) and Task Transformed BKR (TTBKR)). (1) *The TTBLR is a transformed BLR (TBLR) with $M = \Psi^T(\Psi\Psi^T + \frac{\sigma^2}{\beta^2}I)^{-1}\tilde{\Psi}$ and $\Sigma = \sigma^2\tilde{\Psi}^T(\Psi\Psi^T + \frac{\sigma^2}{\beta^2}I)^{-1}\tilde{\Psi} + \sigma^2I$ as the transformation and noise covariance.* (2) *The TTBKR is a transformed BKR (TBKR) with transformation $M = (L + \sigma^2I)^{-1}\tilde{L}$ and noise covariance $\Sigma = \tilde{L} + \sigma^2I - \tilde{L}^T(L + \sigma^2I)^{-1}\tilde{L}$.* (3) *The TTBLR and TTBKR marginal likelihoods are $p(\tilde{\mathbf{z}}) = \mathcal{N}(\tilde{\mathbf{z}}; \mathbf{0}, [\Sigma^{-1} - \Sigma^{-1}A^T\Phi^TC\Phi A\Sigma^{-1}]^{-1})$ where $C = [\Phi A\Sigma^{-1}A^T\Phi^T + \frac{1}{\gamma^2}I]^{-1}$ and $p(\tilde{\mathbf{z}}) = \mathcal{N}(\tilde{\mathbf{z}}; \mathbf{0}, A^TKA + \Sigma)$ respectively.* (4) *For both models, when the posterior for g is approximated via MAP, the covariance becomes $\Sigma = \sigma^2I$. In this case, the parametric (resp. nonparametric) DME estimator (definitions 4.4 and 4.5) is the predictive mean of a TTBLR (resp. TTBKR) with $\lambda = \frac{\sigma^2}{n\beta^2}$ and $\epsilon = \frac{\sigma^2}{m\gamma^2}$ (resp. $\lambda = \frac{\sigma^2}{n}$ and $\epsilon = \frac{\sigma^2}{m}$). An alternative TTBKR marginal likelihood is $p(\tilde{\mathbf{z}}) = \mathcal{N}(\tilde{\mathbf{z}}; \mathbf{0}, \sigma^2[I - A^T(KAA^T + \sigma^2I)^{-1}KA]^{-1})$.* (5) *When both posteriors for g and f are approximated via MAP, TTBLR and TTBKR becomes TTLRR and TTKRR respectively with λ and ϵ from (4).*

Importantly, our end goal is to infer f . While this involves inferring g , g is not of direct interest. A simpler alternative is to only perform Bayesian inference on f and approximate g with its MAP solution, simplifying the noise covariance via (4) of theorem 5.3. This establishes a Bayesian interpretation for DMEs as MAP estimates of TTGPs. Critically, by maximizing the TTGP marginal likelihood, we can learn DME hyperparameters of kernels k and ℓ , and also λ and

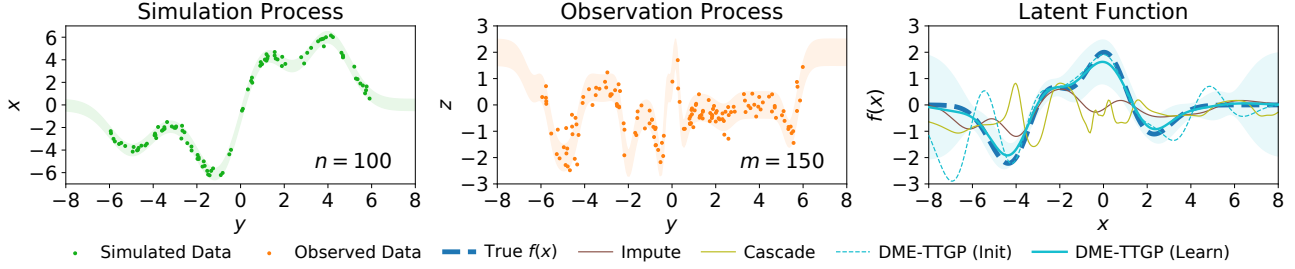


Figure 3. Illustration of latent function recovery with a **task transformed Gaussian process (TTGP)** on non-trivial simulation and observation processes $p(x|y)$ and $p(z|y)$. (Left) The simulation process $p(x|y)$. (Center) The observation process $p(z|y)$. (Right) The true latent function f , the naive solutions using cascaded regressors and imputed data, and the mean and uncertainty bounds of the **TTGP**, also the Bayesian **DME**, with initial and learned hyperparameters. All bounds are 2 standard deviations from the mean.

ε. Furthermore, the computational complexity for alternative marginal likelihood is dominated by inversions that are $O(n^3)$ only, again allowing for larger m at $O(m)$.

In summary, we first establish the **DME** as a nonparametric solution to the **TTR** problem under chained regularized least squares losses that make learning f dependent on learning g . While λ and ϵ are previously seen as numerical adjustments to relax **RKHS** assumptions and stabilize matrix inversions in the **KME** framework, they can now be seen as controlling the amount of function regularization under this loss. Secondly, we present **TTGPs** as nonparametric Bayesian solutions to this regression problem and show that **DMEs** are their posterior predictive means. Again, inference of f is dependent on the inference of g , allowing **GP** uncertainties to propagate through. This connection provides Bayesian interpretations of **DMEs** and enable uncertainty estimation in inferring deconditional means. Critically, we use this to derive marginal likelihoods for hyperparameter learning.

6. Nonparametric Bayes' Rule

While **DMOs** were constructed as solutions to the deconditional mean problem, they also resemble Bayes' rule when we focus solely on considering the encoded relationship between X and Y . This was motivated by theorem 4.1, which revealed that the **DMO** can be fully specified by the **CMO** $C_{X|Y}$ and the second order mean embedding C_{YY} that encoded the likelihood $\mathbb{P}_{Y|X}$ and prior \mathbb{P}_Y respectively. To establish this view, we investigate the conditions for which the **DMO** $C'_{X|Y}$ coincide with the **CMO** $C_{Y|X}$ that encodes the posterior $\mathbb{P}_{Y|X}$, leading to a nonparametric Bayes' rule.

While first class citizens of probability rules are density evaluations, first class citizens of the **KME** framework are expectations. Consequently, instead of relating density evaluations, rules under the **KME** framework relate mean embeddings of distributions at various orders. Importantly, while a distribution $Y \sim \mathbb{P}_Y$ has one simple density evaluation $p_Y(y)$, it can have different **RKHS** representations at different orders such as μ_Y and C_{YY} or higher.

A nonparametric Bayes' rule is a rule which translates Bayes' rule into the **RKHS**, where distributions are represented as **RKHS** operators, alleviating limitations from parametric assumptions such as Gaussian posteriors. It computes a posterior operator $C_{Y|X}$ when given only likelihood operators (e.g. $C_{X|Y}$) and prior operators (e.g. C_{YY}). The **DMO** is appealing as all operators involved are of second order and the same second order likelihood and prior operators are used for both the joint and evidence operator.

However, because C'_{XX} is not necessarily the same as C_{XX} , the **DMO** $C'_{X|Y} = (C_{XY})^T (C'_{XX})^{-1}$ is not necessarily the posterior operator $C_{Y|X} = (C_{XY})^T (C_{XX})^{-1}$. Nevertheless, under certain conditions they coincide with each other.

Theorem 6.1. *If $C_{X|Y} C_{Y|X} C_{XX} = C_{XX}$, then $C'_{XX} = C_{XX}$ and $C'_{X|Y} = C_{Y|X}$.*

A special instance where the assumptions are met is when $X = r(Y)$ where r is not necessarily invertible. Importantly, for empirical **DMOs**, having $x_i = x_j$ for any $y_i = y_j$ suffices, which can be achieved if all y_i are unique. Furthermore, empirical **DMOs** $\bar{C}'_{X|Y}$ can be seen as generalizations of empirical **CMOs** $\hat{C}_{Y|X}$ in the other direction.

Theorem 6.2. *If $m = n$ and $\tilde{y}_i = y_i$ for all $i \in \{1, \dots, n\}$, then the empirical **DMO** corresponding to $\mathbb{P}_{X|Y}$ becomes the empirical **CMO** corresponding to $\mathbb{P}_{Y|X}$ for $\lambda \rightarrow 0^+$,*

$$\lim_{\lambda \rightarrow 0^+} \bar{C}'_{X|Y} = \Psi [K + n\epsilon I]^{-1} \Phi^T = \hat{C}_{Y|X}. \quad (6.1)$$

Intuitively, suppose $\{x_i, y_i\}_{i=1}^n$ are from $p(x, y) := p(x|y)p(y)$ and $\{\tilde{y}_j\}_{j=1}^m = \{y_i\}_{i=1}^n$ is from $p(y)$, then the **DMO** of $p(x|y)$ is equivalent to the **CMO** of $p(y|x) = p(x|y)p(y) / \int_Y p(x|y)p(y)dy$ in the other direction, as per theorem 6.2. In general, however, if $\{x_i, y_i\}_{i=1}^n$ are from $q(x, y) := p(x|y)q(y)$ and $\{\tilde{y}_j\}_{j=1}^m$ is from $p(y)$, then using the joint samples from $q(x, y)$ only to build the **CMO** will yield the **CMO** of $q(y|x) = p(x|y)q(y) / \int_Y p(x|y)q(y)dy$, while using both the joint samples from $q(x, y)$ and marginal samples from $p(y)$ to build the **DMO** will yield the **CMO** corresponding to $p(y|x) = p(x|y)p(y) / \int_Y p(x|y)p(y)dy$. Appendix E details further parallels with probabilistic rules.

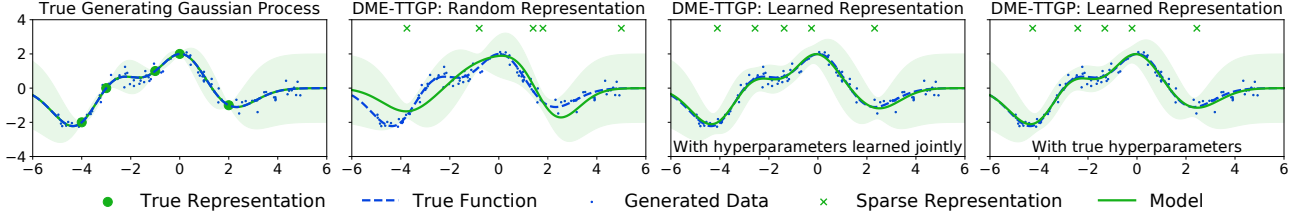


Figure 4. Sparse representation learning on a dataset of 100 points generated by using the toy process from Rasmussen & Williams (2006) as ground truth. (Left) The true function, represented exactly by a GP mean using 5 points. (Left Center) DME using 5 random points. (Right Center) DME with the 5 points and all hyperparameters learned jointly via its marginal likelihood. (Right) DME with the 5 points learned via its marginal likelihood under true hyperparameters. The vertical position of the sparse representation has no meaning.

7. Related Work

DMOs have strong connections to KBR. Both provide a nonparametric Bayes’ rule under the KME framework. In contrast to DMOs where both likelihood and prior operators are of second order, both KBR(a) and KBR(b) (Song et al., 2013; Fukumizu et al., 2013) use a third order likelihood operator $C_{XX|Y}$ and a first order prior embedding μ_Y for the evidence operator $C_{XX} = C_{XX|Y}\mu_Y$. KBR(b) further uses a different third order likelihood operator $C_{XY|Y}$ for the joint operator $C_{XY} = C_{XY|Y}\mu_Y$. Consequently, KBR becomes sensitive to inverse regularizations and effects of prior samples \tilde{y} can vanish. For instance, when $\epsilon \rightarrow 0^+$, KBR(b) degenerate to $\tilde{C}_{Y|X} = \Psi K^{-1} \Phi^T$, which is a CMO that no longer depend on \tilde{y} . Instead, DMOs degenerate to $\tilde{C}'_{X|Y} = \tilde{\Psi} A^T [A A^T]^{-1} K^{-1} \Phi^T$, retaining their original structure. Detailed comparisons are provided in appendix F.

Viewing KMEs as regressors provides valuable insights and interpretations to the framework. CMOs can be established as regressors where the vector-valued targets are also kernel induced features (Grünwälder et al., 2012). In contrast, we establish DMOs as solutions to task transformed regressors which recover latent functions that, together with a likelihood, governs interactions between three variables.

The TTR problem describes the setting of learning from conditional distributions in the extreme case where only one sample of x_i is available for each y_i to describe $p(x|y)$. Dual KMEs (Dai et al., 2017) formulate this setting as a saddle point problem, and employ stochastic approximations to efficiently optimize over the function space. However, without connections to Bayesian models such as TTGPs that admit a marginal likelihood, hyperparameter selection often require inefficient grid search.

Hyperparameter learning of marginal embeddings have been investigated by placing GP priors on the embedding itself to yield a marginal likelihood objective (Flaxman et al., 2016). However, it is unclear how this can be extended to CMEs. Our marginal likelihoods (theorem 5.3 and theorem G.1) provide such objective for DMEs and, due to theorem 6.2, it can also be applied to CMEs as a special case.

8. Applications and Experiments

While DMEs are developed to complement the theoretical framework of KMEs, in this section we describe and demonstrate some of their practical applications with experiments.

8.1. Hyperparameter Learning for TTR

We first illustrate in fig. 3 the TTR problem, the primary application of TTGPs and DMEs. While X and Y are multivariate in general, we use 1D examples to enable visualizations. Although this is a 1D problem, the simulation process $p(x|y)$ and observation process $p(z|y)$ are governed by non-trivial relationships where successful recovery of f requires dealing with difficult multi-modalities in $p(y|x)$. To generate the data, we choose non-trivial functions r and f and generate $X_i = r(Y_i) + \eta_i$ and $\tilde{Z}_j = f(r(\tilde{Y}_j) + \tilde{\eta}_j) + \tilde{\xi}_j$, where $Y_i, \tilde{Y}_j \sim U(-6, 6)$ and $\eta_i, \tilde{\eta}_j, \tilde{\xi}_j \sim \mathcal{N}(0, 0.25^2)$ for all $i \in \{1, \dots, n\}$ and $j \in \{1, \dots, m\}$. In this way, $p(x|y) = \mathcal{N}(x; r(y), 0.25^2)$, $p(z|x) = \mathcal{N}(z; f(x), 0.25^2)$, and $\mathbb{E}[Z|Y = y] = \mathbb{E}[f(r(y) + \tilde{\eta}) + \tilde{\xi}] = \mathbb{E}[f(X)|Y = y]$.

By optimizing the marginal likelihood in theorem 5.3 (3), we see that the DME is able to adapt from its initial hyperparameters to learn the latent function accurately.

We compare this to two naive solutions that one may propose when faced with a TTR problem. The *cascade* method trains separate regressors from X to Y , with the transformation set, and from Y to Z , with the task set. They use the former to predict y^* from x^* and the latter to predict z^* from y^* . The *impute method* trains a regressor from Y to Z with the task set and predicts \mathbf{z}_{fake} at locations \mathbf{y} , and trains a regressor on the dataset $(\mathbf{x}, \mathbf{z}_{\text{fake}})$ to predict z^* from a new x^* . We use GPR means (KRR) for all such regressors. Both methods suffer because uncertainty propagation is lost by training regressors separately. The cascade method suffers further because $p(y|x)$ is usually highly multi-modal such as in this example, so unimodal regressors like GPR from X to Y are unsuitable. This also highlights that while DMEs provides unimodal Gaussian uncertainty on function evaluations and thus Z , they capture multi-modality as a nonparametric Bayes’ rule between X and Y .

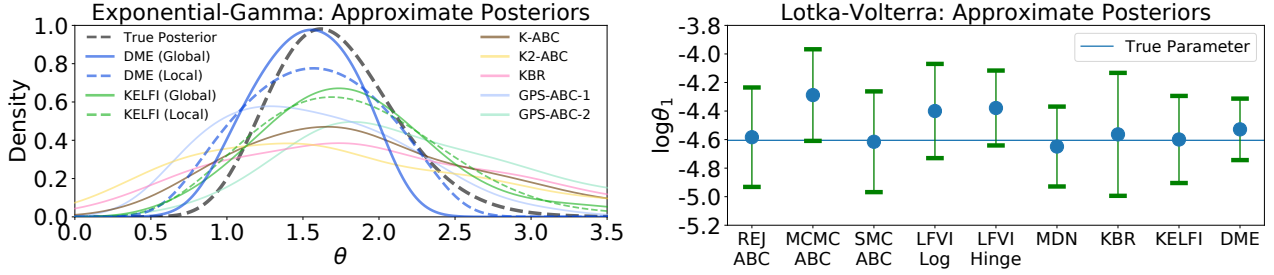


Figure 5. Application to LFI. (Left) Approximate posteriors under kernel based LFI methods for the toy exponential-gamma problem using 100 simulations. ‘Global’ and ‘Local’ refer to the optimality of model hyperparameters with respect to their respective approximate marginal likelihoods. (Right) Approximate posteriors of the first (log) parameter. Error bars represent the middle 95% credible interval.

8.2. Sparse Representation Learning with TTGP

A special case of the TTR problem is to learn sparse representations for big data with trainable inducing points. Continuing with the notations used so far, we are given a large original dataset $(\tilde{\mathbf{y}}, \tilde{\mathbf{z}})$ of size m , with inputs Y and target Z . We let the transformation dataset be a set of n inducing points $\mathbf{x} = \mathbf{y} = \mathbf{u}$ for Y where $n \ll m$. That is, we degenerate to $X = Y$ and $\mathcal{X} = \mathcal{Y}$. We maximize the alternative marginal likelihood in theorem 5.3 (4) with respect to the inducing points and learn the TTGP hyperparameters jointly. For predictive mean we use the alternative computational form in definition 4.4. Similar form exists for the covariance. These alternative forms are suitable for this application because n is small for its $O(n^3)$ inversions and dependence on m is only $O(m)$. We illustrate this process in fig. 4.

8.3. Likelihood-Free Inference with DME

As a nonparametric Bayes’ rule, DMEs can be used for likelihood-free inference (LFI) (Marin et al., 2012) where likelihood evaluations are intractable but sampling from a simulator $\mathbf{x} \sim p(\mathbf{x}|\theta)$ is possible. The simulator takes parameters θ and stochastically generates simulated data that are often summarized into statistics \mathbf{x} . Observed data are also summarized into statistics \mathbf{y} , and discrepancies with \mathbf{x} are often measured by an ϵ -kernel $\kappa_\epsilon(\mathbf{y}, \mathbf{x}) = p_\epsilon(\mathbf{y}|\mathbf{x})$ such that $p_\epsilon(\mathbf{y}|\theta) = \int p_\epsilon(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\theta)d\theta$. This ϵ is not to be confused with the regularization used for C'_{XX} , which we denote as δ for this section only. After selecting a prior $p(\theta)$, the goal is to approximate the posterior $p_\epsilon(\theta|\mathbf{y})$.

Translating notations into the LFI setting, we have $x_i \rightarrow \mathbf{x}_i$, $y_i \rightarrow \theta_i$, $\tilde{y}_j \rightarrow \tilde{\theta}_j$, and $x \rightarrow \mathbf{y}$. We first simulate $\mathbf{x}_i \sim p(\mathbf{x}|\theta_i)$ on parameters $\{\theta_i\}_{i=1}^n \sim \pi(\theta)$ not necessarily from the prior to get $\{\theta_i, \mathbf{x}_i\}_{i=1}^n$ for the likelihood, and sample $\{\tilde{\theta}_j\}_{j=1}^m \sim p(\theta)$ for the prior. We then build the DME $\bar{\mu}_{\Theta|\mathbf{X}=\mathbf{y}}$ and sample it with kernel herding (Chen et al., 2010) for posterior super-samples. This is described in algorithm 8.1. We also provide the approximate marginal likelihood objective \bar{q} to maximize for hyperparameter learning of the DME. Derivations are detailed in appendix G.

Algorithm 8.1 Deconditional Mean Embeddings for LFI

- 1: **Input:** Data \mathbf{y} , simulations $\{\theta_i, \mathbf{x}_i\}_{i=1}^n \sim p(\mathbf{x}|\theta)\pi(\theta)$, prior samples $\{\tilde{\theta}_j\}_{j=1}^m \sim p(\theta)$, query points $\{\theta_r^*\}_{r=1}^R$, kernels k, κ_ϵ, ℓ , and ℓ' , regularization λ and δ
- 2: $L \leftarrow \{\ell(\theta_i, \theta_j)\}_{i,j=1}^{n,n}$, $\tilde{L} \leftarrow \{\ell(\tilde{\theta}_i, \tilde{\theta}_j)\}_{i,j=1}^{m,m}$
- 3: $A \leftarrow (L + n\lambda I)^{-1}\tilde{L}$, $\tilde{L}^* \leftarrow \{\ell(\tilde{\theta}_j, \theta_r^*)\}_{j,r=1}^{m,R}$
- 4: $K \leftarrow \{k(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^{n,n}$, $\mathbf{k}(\mathbf{y}) \leftarrow \{k(\mathbf{x}_i, \mathbf{y})\}_{i=1}^n$
- 5: **DME:** $\mu \leftarrow (\tilde{L}^*)^T A^T [K A A^T + m\delta I]^{-1} \mathbf{k}(\mathbf{y}) \in \mathbb{R}^R$
- 6: **for** $s \in \{1, \dots, S\}$ with $\mathbf{a} \leftarrow \mathbf{0} \in \mathbb{R}^R$ **initialized do**
- 7: $\hat{\theta}_s \leftarrow \theta_{r^*}$ where $r^* \leftarrow \operatorname{argmax}_r \mu_r - (a_r/s)$
- 8: $\mathbf{a} \leftarrow \mathbf{a} + \{\ell'(\theta_{r^*}, \hat{\theta}_s)\}_{r=1}^R$
- 9: **end for**
- 10: **Output:** Posterior super-samples $\{\hat{\theta}_s\}_{s=1}^S$
- 11: **Learning:** $\bar{q} \leftarrow \operatorname{mean}(A^T \kappa_\epsilon)$, $\kappa_\epsilon \leftarrow \{\kappa_\epsilon(\mathbf{y}, \mathbf{x}_i)\}_{i=1}^n$

Figure 5 demonstrates algorithm 8.1 on two standard benchmarks. For the toy exponential-gamma problem we compare directly with other kernel approaches. As simulations are usually very expensive, we show the case with very limited simulations ($n = 100$), leading to most methods producing posteriors wider than the ground truth. Nevertheless, by optimizing \bar{q} in line 11, DMEs can adapt their kernel length scales accordingly. For Lotka-Volterra, the ABC methods used more than 100000 simulations, while MDN used 10000 simulations. To achieve competitive accuracy, kernel approaches such as DMEs, KELFI (Hsu & Ramos, 2019), and KBR used 2000, 2500, and 2500 simulations. Acronyms and experimental details are described in appendix G.

9. Conclusion

The connections of DMEs with CMEs and GPs produce useful insights towards the KME framework, and are important steps towards establishing Bayesian views of KMEs. DMEs provide novel solutions to a class of nonparametric Bayesian regression problems and enable applications such as sparse representation learning and LFI. For future work, relaxing assumptions required for DMOs as a nonparametric Bayes’ rule can have fruitful theoretical and practical implications.

References

- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- Chen, Y., Welling, M., and Smola, A. Super-samples from kernel herding. In *The Twenty-Sixth Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI-10)*, pp. 109–116. AUAI Press, 2010.
- Dai, B., He, N., Pan, Y., Boots, B., and Song, L. Learning from Conditional Distributions via Dual Embeddings. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 1458–1467. PMLR, 20–22 Apr 2017.
- Flaxman, S., Sejdinovic, D., Cunningham, J. P., and Filippi, S. Bayesian learning of kernel embeddings. In *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*, pp. 182–191. AUAI Press, 2016.
- Fukumizu, K., Bach, F. R., and Jordan, M. I. Dimensionality reduction for supervised learning with reproducing kernel Hilbert spaces. *Journal of Machine Learning Research*, 5 (Jan):73–99, 2004.
- Fukumizu, K., Song, L., and Gretton, A. Kernel Bayes’ rule: Bayesian inference with positive definite kernels. *Journal of Machine Learning Research*, 14(1):3753–3783, 2013.
- Grünwälder, S., Lever, G., Baldassarre, L., Patterson, S., Gretton, A., and Pontil, M. Conditional mean embeddings as regressors. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012*, volume 2, pp. 1823–1830, 2012.
- Hastings, W. K. Monte carlo sampling methods using markov chains and their applications. 1970.
- Higham, N. J. *Accuracy and stability of numerical algorithms*. SIAM, 2002.
- Hsu, K. and Ramos, F. Bayesian Learning of Conditional Kernel Mean Embeddings for Automatic Likelihood-Free Inference. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 2631–2640. PMLR, 16–18 Apr 2019.
- Marin, J.-M., Pudlo, P., Robert, C. P., and Ryder, R. J. Approximate Bayesian computational methods. *Statistics and Computing*, 22(6):1167–1180, 2012.
- Meeds, E. and Welling, M. GPS-ABC: Gaussian process surrogate approximate Bayesian computation. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, pp. 593–602. AUAI Press, 2014.
- Muandet, K., Fukumizu, K., Sriperumbudur, B., Schölkopf, B., et al. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017.
- Murray-Smith, R. and Pearlmuter, B. A. Transformations of gaussian process priors. In *Deterministic and Statistical Methods in Machine Learning*, pp. 110–123. Springer, 2005.
- Nakagome, S., Fukumizu, K., and Mano, S. Kernel approximate Bayesian computation in population genetic inferences. *Statistical applications in genetics and molecular biology*, 12(6):667–678, 2013.
- Papamakarios, G. and Murray, I. Fast ε -free inference of simulation models with bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, pp. 1028–1036, 2016.
- Park, M., Jitkrittum, W., and Sejdinovic, D. K2-ABC: Approximate Bayesian Computation with Kernel Embeddings. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pp. 398–407. PMLR, 09–11 May 2016.
- Rasmussen, C. E. and Williams, C. K. I. *Gaussian processes for machine learning*. The MIT Press, 2006.
- Song, L., Huang, J., Smola, A., and Fukumizu, K. Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 961–968. ACM, 2009.
- Song, L., Fukumizu, K., and Gretton, A. Kernel embeddings of conditional distributions: A unified kernel framework for nonparametric inference in graphical models. *IEEE Signal Processing Magazine*, 30(4):98–111, 2013.
- Tran, D., Ranganath, R., and Blei, D. Hierarchical implicit models and likelihood-free variational inference. In *Advances in Neural Information Processing Systems*, pp. 5523–5533, 2017.