

FETISH: Functional Embedding of Terms in a Spatial Hierarchy

Alex Grabanski* and Lily Wilk†
Case Western Reserve University

October 2, 2020

Abstract

Motivated by the problem of program induction (learning of a computer program from data), we describe an incremental Bayesian method to derive embeddings of terms in a simply-typed combinatory calculus as probability distributions over finite-dimensional vector spaces. Our method, which we call "Functional Embedding of Terms in a Spatial Hierarchy", or "FETISH", assigns a space of embeddings to every type of the underlying language, and assigns an embedding to every term which has been evaluated by the interpreter. As the interpreter is issued commands from a down-stream task, these term embeddings are kept updated in such a way that they reflect the best-available knowledge about term behavior. To do so, we propose novel Bayesian inference routines for random-feature kernel regressors situated in a hierarchy of typed functions. Finally, we discuss ways that supervised, cumulative, multi-task, and reinforcement learning could potentially benefit from the proposed framework.

1 Introduction

The problem of *inductive reasoning*, or reasoning about hypotheses from data, is the central problem of machine learning. While many modern machine learning techniques typically restrict the collection of hypotheses to well-defined and manageable classes, in the spirit of radicalism, we choose instead to jump back to the 1964 theory Ray Solomonoff advanced about induction [16], which was fundamentally a process of induction over *computer programs*. In some sense,

this is the most general setting for the problem of induction in machine learning, since the hypothesis class of computer programs *is* the most general class of hypotheses expressible by machine. Unfortunately, induction over the space of all programs is not computable. In fact, this is incredibly unfortunate, since under the definition of "artificial general intelligence" ("AGI") adopted by AIXI [5], we only need optimal inductive inference to achieve provably-optimal AGI. While approximations to the conceptual process of Bayesian inference over all computer programs have been previously advanced, such as the one in AIXI-tl [5], they have exhibited very poor performance to date.

Motivated by this problem, and by the utility of embeddings for dealing with otherwise-unwieldy objects such as natural-language words in machine learning, we devise a method to embed the terms of a simply-typed combinatory calculus into separate spaces for each type. However, instead of simple vector embeddings, we embed functional terms as *probability distributions* over vectors in the embedding space. In this way, we hope to provide an immensely practical Bayesian inference method over program terms which is not only computable, but also has a very intuitive spatial structure which is readily consumed by down-stream tasks.

2 Prior Research

While what we will describe in this paper does not involve any kind of process for program search, since we leave that open to future developments, the closest area we can draw a reference point to is that of *program synthesis*. A good high-level survey of this area is given in ([4]). In particular, the sub-field of *neural program synthesis* has some relatively recent works [13] [11] which focus on the compositionality of programs and program embeddings using neural-network based encoders as

*e-mail: ajg137@case.edu

†She is a Papillon, and so she does not have an e-mail address. All inquiries should be directed to her owner Alex Grabanski, and if necessary for the query, the two will go on a walk to think on it.

part of their gradient-based optimization pipeline to derive programs from examples. However, unlike neural-network-based frameworks such as AlphaNPI [13], our framework for deriving embeddings is *not* based on interpretation in an imperative language, but instead a typed combinatory calculus, which we would expect to have better compositional properties. In addition, unlike the aforementioned neural program synthesis approaches, we do not require any gradients in our derivation of embeddings, which means that we are free to use an actual interpreter instead of a differentiable surrogate.

Even closer, very recently, ([1]) proposed a neural-network-based framework for gradually growing and abstracting a library of functional programs in a multi-task learning environment. However, unlike the approach here, they do not at any point generate embeddings for program terms, but instead, a discriminative neural-network model on programs which encodes how likely they are to solve the task at hand.

Another research area which is tangentially related is that of *neural architecture search*, a survey of which is provided in ([18]). In neural architecture search, building-blocks of neural networks are composed to attempt to achieve optimal architecture, which is quite similar to the task of composing a program from combinators, with the caveat that neural architecture search is intended to run before a full optimization over all parameters in the network. While this research area is largely divergent from what we consider here, we do note that [8] did consider a Bayesian Optimization framework for neural architecture search which leveraged Gaussian Processes, which we also do here. However, the internal representation in their framework differs radically from the one we are about to present, as their representation is not compositional in nature.

3 Language and Interpreter

In order to make the problem of program inference from data more tractable, we restrict our attention to a very minimal language based on combinatory calculus consisting entirely of constants and combinators under simple typing. While this paradigm may seem restrictive at first, we will describe how these choices are in fact justified, and the resulting language is more powerful than it may seem at first glance.

3.1 Types

The combinatory language we consider has only two broad categories of types: primitive types, and function types.

3.1.1 Primitive Types

A primitive type consists of a *label* [which is an arbitrary identifying token] and an $n \in \mathbb{N}$ which is referred to as the *dimensionality* of the primitive type. We assume that an arbitrary, but finite number of primitive types with distinct labels are defined for the language, which comprise a set we denote as \mathcal{P} . While the labels must be globally unique, any two primitive types with different labels may share the same dimensionality. In the rest of this paper, we will refer to types by their labels.

Intuitively, a primitive type of dimension n may be thought of as the vector space \mathbb{R}^n together with a label describing its preferred interpretation (e.g: spatial vectors, images, voxel bitmaps, word embeddings, etc.) While the restriction of primitive types to finite-dimensional vector spaces rules out a large collection of sequence types (e.g: Strings, arbitrary-length Lists), it notably allows many of the data formats which are found commonly in ML literature.

3.1.2 Function Types

A function type, denoted $(X \rightarrow Y)$ consists of a *domain* type X and a *codomain* type Y . Unlike nearly every other programming language, our language does *not* realize a function type for every possible domain and codomain. Instead, we only recognize function types belonging to a particular finite set \mathcal{F} as being well-formed. The only restriction we place on the finite set \mathcal{F} is that it is *closed*, in the sense that if $(X \rightarrow Y) \in \mathcal{F}$, then $X \in \mathcal{P} \cup \mathcal{F}$ and $Y \in \mathcal{P} \cup \mathcal{F}$.

While the requirement that the set of function types is a finite set may seem restrictive at first, it's important to note that in actual practice, human programmers rarely use functions with high arity nor higher-order functions with overly-complex functional arguments. In practice, this means that in any given language, we could restrict the nesting of the function type constructor to a sufficiently large finite value and lose absolutely nothing of practical import. Since our language has only a finite number of primitive types, such a nesting restriction would be sufficient to guarantee that \mathcal{F} is a finite set.

3.2 Terms

With the types defined above, we may now describe the terms of the language.

3.2.1 Primitive Terms

For every primitive type with label L and dimensionality n , and any arbitrary vector $v \in \mathbb{R}^n$, $L(v) : L$.

3.2.2 Function Terms

For every function type $F = (X \rightarrow Y)$, a term $f : F$ is a procedure which takes any term of type X and yields a term of type Y in finite time. While we do require f to be total, we do not require that f is deterministic, nor do we require that there is a uniform time bound on evaluation of f across all inputs. We will write function application of term $f : F$ to term $x : X$ by $f(x)$.

To define the collection of all function terms, we suppose that a finite set of *primitive function terms* is provided, and take the closure of this set under function application of function terms to terms of the appropriate domain type.

The requirement that every function term is associated with a total procedure rules out the class of general recursive functions from consideration under this framework, but once again, this is not a severe restriction in practice, since we are allowed to define function terms for procedures whose runtime is dependent on characteristics of the input.

3.3 Interpreter

With the definitions of the constructs of the language given, we can now describe the operation of the interpreter and its interface, as exposed to downstream tasks.

In our interpreter, the unique terms of functional types will be assigned their own *address* for long-term storage. Upon initialization, only the primitive function terms will be assigned addresses, but new addresses will be incrementally assigned as new function terms are derived as the return values from evaluated term applications. Vector terms, in contrast, are never assigned addresses, and simply exist inline as arrays of floating-point values. To be able to refer to both situations under the same framework, we say that a *term reference* is either the address of a function term, or a bare array of floating-point values representing a vector term.

We only expose a single mechanism to interact with the interpreter in downstream tasks: Evaluating a function $f : X \rightarrow Y$ corresponding to a given address on an argument $x : X$ for a given term reference, and returning a resulting term reference to $f(x)$ to the downstream task. In addition, as a side-effect of doing this, if $f(x)$ is a function term which does not already have an address, we allocate space for it and assign it to a new address.

Although this mechanism for interacting with the interpreter is extremely minimal, it still allows the evaluation of arbitrarily-nested applications of terms to terms, since the code invoking the interpreter could take an arbitrary expression tree and evaluate it using the strict evaluation order, with intermediate expressions always represented by their term references.

4 Term Embeddings

We are now ready to discuss the embeddings associated with each type of the language. Recall that an embedding is a mapping from objects of some class to a vector space, such that similar objects map to vectors which are close in norm. Given this context, for us to proceed, it's crucial for us to define the notion of "similarity" for terms of primitive and function types.

In the case of primitive terms, since terms of a primitive type L of dimension n are directly identifiable with vectors in \mathbb{R}^n , we may simply define an embedding by $\phi : L \rightarrow \mathbb{R}^n$ by $\phi(L(v)) = v$ for each such primitive type.

However, a mechanism for deriving embeddings of function terms is far less obvious. Function spaces on real vector spaces are inherently infinite-dimensional, so we need to either apply a technique which reduces the dimensionality to something that is dependent on the number of data points [as in kernel regression], or we'll need to choose a restricted class of functions which may be parametrized by finite-dimensional vectors. Generally, when faced with such a problem, ML practitioners will advance either kernel-based techniques or deep neural networks to achieve a solution.

Unfortunately, neural networks are ill-suited to this domain due to the heavy computational expense of training via gradient descent. Our aim is to provide an incremental algorithm for deriving term embeddings which exploits the totality of knowledge about all terms after each evaluation step, and in order to make this work with neural networks as models of

functions, we would need to fully train networks at every evaluation step of the interpreter, the cost of which is prohibitive.

On the other hand, kernel regression has the clear advantage in this domain of having closed-form minimizers of the squared error when modeling a particular collection of data-points. However, due to a technical difficulty, we cannot apply kernel regression to this case, either: since we require that the embedding for terms of a given function type targets a space of *fixed* dimension, we cannot apply the kernel trick here. However, we can *approximate* kernel regression with a fixed-dimensional function space by using a finite, pre-determined collection of random features. In this setting, we only need to perform a linear regression on feature-mapped input-output pairs, which fits our requirement for incrementality quite well. Our reference implementation supports sketched linear features, randomized quadratic features [10], and Random Fourier Features [12].

To be more concrete, the above proposal would represent the embedding of a function of type $F = (X \rightarrow Y)$ by the vectorization of the $t \times r$ coefficient matrix A in the best-fit multivariate regression model:

$$y = A * \psi_F(x) + \epsilon$$

where $x \in \mathbb{R}^s$ is the embedding of the input in X , $y \in \mathbb{R}^t$ is the embedding of the output in Y , $\epsilon \in \mathbb{R}^t$ is a noise term, and ψ_F is a feature mapping from the input embedding space \mathbb{R}^s to the feature space \mathbb{R}^r ¹.

However, simply representing the embedding of a function by the maximum-likelihood ordinary regression model is insufficient for our purposes upon further reflection. In particular, the proposed model above takes on a recursive character in the case where either X or Y are themselves function types. Since estimation of function embeddings is data-dependent, the accuracy of a fit at the level of F would implicitly depend on the accuracy of fits at the levels of X or Y . This motivates a consideration of the error in the best coefficients associated with the regression in addition to the best regression coefficients themselves.

¹The choice of ψ_F is fixed across all functions in F , but we do have the freedom to adapt this choice according to the identities of the types X and Y . In particular, this leaves open the possibility of choosing features based on the label of X , if X is a primitive type. The authors can imagine great utility in tailoring these choices by e.g. specifying wavelet-based features when X represents an image, or using features originating from a pre-trained neural network, but such design decisions are beyond the scope of this paper.

4.1 Schmears and Extended Embeddings

To approach the problem of representing uncertainty in our knowledge of function embeddings, we first introduce some definitions which will greatly simplify this task.

First, a *schmear* is our term for the representation of [limited information about] a probability distribution by its first two moments [mean and covariance]. In other words, an n -dimensional schmear is a pair (μ, Σ) where $\mu \in \mathbb{R}^n$ and Σ is a n -by- n positive semi-definite matrix. We denote the set of all n -dimensional schmears by \mathcal{S}^n .

With the definition of "schmear" in hand, we can now extend our previous notion of an "embedding" to accurately model uncertainty. An *extended embedding* for the type X into \mathbb{R}^n is a map:

$$\phi : X \rightarrow \mathcal{S}^n$$

which associates each term reference of that type to n -dimensional schmears. Moreover, while not a requirement of our definition, we would like our extended embeddings to send operationally-similar terms to schmears that are "close" to one another.

Once again, deriving an extended embedding in the case where X is the primitive type L is straightforward. Vectors are vectors, and have no uncertainty about them, so in the spirit of the definition, we may simply take:

$$\phi(L(v)) = (v, \mathbf{0})$$

4.2 Extended Embeddings for Function Types

In line with our motivating sections above, we will use regression models to represent embeddings of function terms. However, in order to get the extended embeddings we desire, it is necessary for us to consider full-blown *Bayesian* multivariate regression, not just formulas for the maximum-likelihood regression coefficients.

Recall that for the multivariate linear regression model:

$$y = A * z + \epsilon$$

for y a t -dimensional vector and z a r -dimensional vector, where we assume that $\epsilon \sim_{iid} N(\mathbf{0}_{t \times t}, \Sigma_\epsilon)$ for noise covariance matrix Σ_ϵ , we can put a conjugate *matrix-normal inverse-Wishart* prior over A and Σ_ϵ [15], meaning that we suppose:

$$\Sigma_\epsilon \sim \mathcal{IW}_t(V, v)$$

$$B|\Sigma_\epsilon \sim \mathcal{MN}_{t,r}(B, \Lambda^+, \Sigma_\epsilon)$$

Where $\mathcal{IW}_t(V, v)$ denotes the *inverse-Wishart* distribution with $t \times t$ scale matrix V and v degrees of freedom, and $\mathcal{MN}_{t,r}(B, \Lambda^+, \Sigma_\epsilon)$ denotes the *matrix normal* distribution with $t \times r$ mean B , $r \times r$ input covariance Λ^+ (for precision Λ)² and $t \times t$ output covariance Σ_ϵ . We refer the reader to [2] or [14] for the probability density functions and some additional properties of these distributions. Here, we simply emphasize the fact that if

$$X \sim \mathcal{MN}_{t,r}(M, U, V)$$

then

$$\text{vec}(X) \sim \mathcal{N}_{t*r}(\text{vec}(M), V \otimes U)$$

where $\text{vec}(-)$ denotes vectorization of a matrix, $\mathcal{N}_n(-, -)$ denotes the multivariate normal distribution parameterized by its n -dimensional mean and $n \times n$ covariance, respectively, and \otimes denotes the Kronecker product.

In what follows, we will use the shorthand $\mathcal{MN}\mathcal{IW}_{t,r}(B, \Lambda, V, v)$ for the matrix-normal inverse-Wishart prior as described above. Since the mean of $\mathcal{IW}_t(V, v)$ is $\frac{V}{v-t-1}$ and the mean of the matrix-normal inverse-Wishart prior is independent of the noise covariance, by an application of the law of total covariance, we can note that if

$$X \sim \mathcal{MN}\mathcal{IW}_{t,r}(B, \Lambda, V, v)$$

then

$$E[\text{vec}(X)] = B, \quad \text{Cov}[\text{vec}(X)] = \frac{V}{v-t-1} \otimes \Lambda^+$$

As a consequence, we can now define our extended embeddings for function terms. All we need to do is let $z = \psi_F(x)$ in the regression model above, and assume that our current knowledge of the regression parameters is matrix-normal inverse-Wishart distributed as above, which results in the extended embedding:

$$\phi(f) = (B_f, \frac{V_f}{v_f - t - 1} \otimes \Lambda_f^+)$$

where $f : F = (X \rightarrow Y)$, and $(B_f, \Lambda_f, V_f, v_f)$ is the current parameter tuple for the $t \times r$ matrix-normal inverse-Wishart model placed on f .

²Here and elsewhere in this paper, we opt to use the Moore-Penrose pseudoinverse instead of the matrix inverse, because we frequently translate between covariance and precision matrices which are sometimes of deficient rank since they often originate from sums of outer products in our inference routines.

5 Bayesian Updates

With our assumption of matrix-normal inverse-Wishart priors on our models of functions, our task is now to keep these models updated with the current best-available information as we evaluate new term applications in our interpreter.

5.1 Data Updates

As with any typical supervised regression model, we need to utilize information from the input/output pairings which are provided to us in order to determine the regression coefficients. However, unlike a typical regression model, the input/output pairings used in our case to update the function type $F = (X \rightarrow Y)$ are not in general input/output pairings of data points, but rather, of input/output pairings of *schmears*. As a result, we cannot directly employ the typical expressions for incremental updates to a multivariate linear regression.

However, we can *reduce* this problem of ours to a form which is expressible as weighted input/output data point updates. To accomplish this, we first greatly simplify matters by ignoring the input covariance entirely, so we are left with just a pairing of an observed input mean and the output schmar. Our reason for doing this is that in general, accounting for uncertainty in inputs in a Bayesian regression framework typically adds significant complexity [2]. On the other hand, uncertainty in outputs in a Bayesian regression framework is relatively simple to deal with.

After ignoring the spread of the input distribution, we may then proceed using a trick from the field of control theory to also eliminate the need to consider the output distribution in favor of instead considering a collection of points with the same mean and covariance. For a given schmar $(\mu, \Sigma) \in \mathcal{S}^n$, we can derive the *symmetric set of $2n+1$ sigma points* [7] as:

$$\{\mu\} \cup \{\mu \pm [\sqrt{n\Sigma}]_i \mid i \in [0..n]\}$$

Where $[\sqrt{n\Sigma}]_i$ denotes the i th column of the unique symmetric matrix square root of $n\Sigma$. Straightforward calculation may be used to verify that this set of points has the same mean and covariance as the original schmar. Consequently, we can represent our data-point/schmar update as a collection of $2n+1$ data-point/data-point updates. In particular, if (μ_X, Σ_X) is our input schmar, and $s_Y = (\mu_Y, \Sigma_Y)$ is our output schmar, we perform $2n+1$ data-point updates, each with weight $\frac{1}{2n+1}$ using the input-output pairings

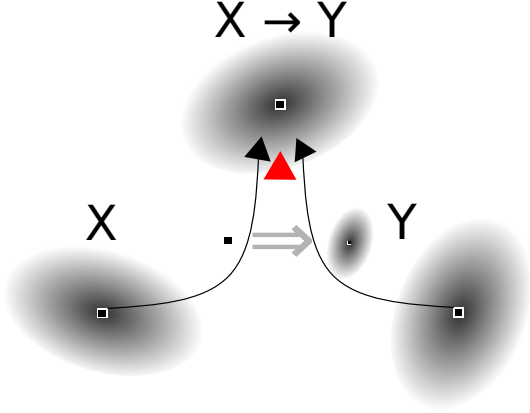


Figure 1: Schematic diagram for data updates, which utilize the mean of the input schmear and the entirety of the output schmear. In this and in the other diagrams of this paper, red triangles will be used to indicate which distribution the update(s) are being performed on, and black arrows will be used to indicate how information is propagated for the update(s).

$\{(\mu_X, \sigma_{Yi}) \mid i \in [0..2n+1]\}$, where σ_{Yi} is the i th sigma-point derived from s_Y .

Now that we have reduced the problem of keeping our models updated with respect to input/output schmears to the problem of keeping our models updated with respect to input/output data-points, we may simply use the usual expressions for the Bayesian update of the matrix-normal inverse-Wishart distribution $\mathcal{MN}\mathcal{IW}(B, \Lambda, V, v)$ with respect to observation of a w -weighted input-output vector pairing (z, y) .

$$\begin{aligned}
 B &\Rightarrow (B\Lambda + w * yz^T)\Lambda_*^+ &= B_* \\
 \Lambda &\Rightarrow \Lambda + w * zz^T &= \Lambda_* \\
 v &\Rightarrow v + w \\
 V &\Rightarrow V + (B - B_*)\Lambda(B - B_*)^T \\
 &\quad + w * \left(\frac{yz^T}{z^T z} - B_*\right)zz^T\left(\frac{yz^T}{z^T z} - B_*\right)^T
 \end{aligned}$$

Note that using the above expressions, we may also undo an update by simply replacing $w \Rightarrow -w$. In practice, we will not use the above formulas directly, but will instead opt to also track the value of Λ^+ across updates utilizing the Sherman-Morrison rank-1 inverse update formula for greater computa-

tional efficiency [9].

5.2 Prior Updates

While we are now relatively well-equipped, since the previously-described process is able to account for observed input/output pairings resulting from evaluation of a function application, we still have not exploited all information available to us. To see this, we will consider a worked example involving partial application of a binary function.

Suppose that we have a function term $f : X \rightarrow (Y \rightarrow Z)$ and two terms $x_1, x_2 : X$, and we have already directed the interpreter to evaluate $f(x_1) = f_1$ and $f(x_2) = f_2$. Furthermore, suppose that we have a large collection of terms $y_1, y_2, \dots, y_n : Y$ for which we have evaluated $f_1(y_i)$ for each i , but we have not yet evaluated f_2 on anything. In this situation, if we only use the above update process, the embedding of f_1 will be close to the maximum-likelihood regression model for the pairings $\{(y_i, f_1(y_i)) \mid i \in [1, \dots, n]\}$. On the other hand, the embedding for f_2 will be entirely determined by the prior we adopt over regression coefficients. From these observations, we can immediately notice a conflict: If we take the limit as the embeddings of x_1 and x_2 approach each other, there's a sharp discontinuity in the behavior we purport to observe for f on embeddings. Fundamentally, this discontinuity originates from the fact that the estimates of the embeddings for f_1 and f_2 are not linked, despite the fact that both stem from partial evaluation of f .

To resolve this problem at a conceptual level, we can note that our observations of f_1 and f_2 will also result in a data-based update to the embedding of f . In the situation above, if we could only propagate the information we have gained about f to f_2 (ignoring the information stemming from the embedding of f_2 itself), we could encode the intuition that f_2 should be "close" to f_1 .

We follow exactly this intuition in proposing another update routine which we call a *prior update*, since it in some sense updates the prior knowledge we apply to partially-evaluated functions. Moving beyond the particular example above, consider the general scenario where $f : X \rightarrow Y$ where $Y = (Z \rightarrow W)$ and $x : X, g = f(x) : (Z \rightarrow W)$. We want to use the embeddings of f and x to impute an update to the embedding of g . Conceptually, we break down this routine into a series of steps: First, we impute a schmear in the embedding space of $(Z \rightarrow W)$ which represents our best current estimate of the schmear for the application of f applied to x ,

where we modify our information about f to ignore any data update applied to it stemming from x . Then, once we have obtained our best-estimate schmear in the embedding space of $(Z \rightarrow W)$, we translate this schmear into a Bayesian update to perform on the matrix-normal inverse-Wishart model for g .

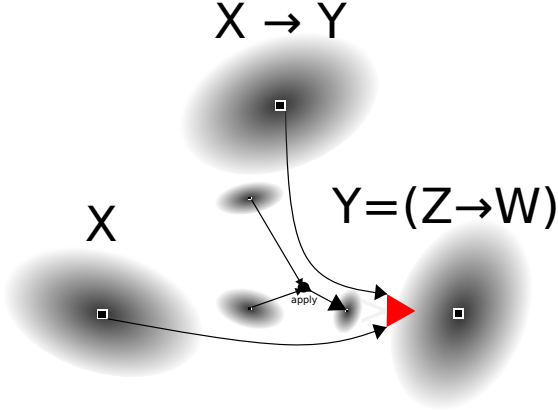


Figure 2: Schematic diagram for prior updates. The embedded schmear for the function and the embedded schmear for an argument are applied to one another to impute an output schmear, which is then translated into an update for the model of the output term.

5.2.1 Imputing the Output Schmear

For the moment, consider the straightforward linear model $\tilde{A}\tilde{z} + \epsilon = \tilde{y}$ where \tilde{A} is $t \times r$ and further suppose that we employ a $\mathcal{MN}\mathcal{IW}_{t,r}(\tilde{B}, \tilde{\Lambda}, \tilde{V}, \tilde{v})$ prior over the usual parameters and suppose that $E[\tilde{z}] = \mu$ and $Cov[\tilde{z}] = \Sigma$. Then,

$$E[\tilde{y}] = \tilde{B}\mu$$

$$Cov[\tilde{y}] = [\langle \Sigma, \tilde{\Lambda}^+ \rangle_F + \mu^T \tilde{\Lambda}^+ \mu] \frac{\tilde{V}}{\tilde{v} - t - 1} + \tilde{B}\Sigma\tilde{B}^T$$

where $\langle -, - \rangle_F$ denotes the Frobenius inner product of matrices. An elementary derivation of this fact is provided in Appendix A.

While the above formulas may be used to get an output schmear from the matrix-normal inverse-Wishart distribution and the input distribution for a *linear* model, we are interested in the case of a potentially *nonlinear* model where $\tilde{z} = \phi(\tilde{x})$ for some feature mapping $\phi : \mathbb{R}^s \rightarrow \mathbb{R}^r$. Once again, we borrow a trick from the field of control theory, since from a schmear over \tilde{x} , we may use an *unscented transform* [7] to obtain an estimate of the schmear over \tilde{z} by using a collection of sigma-points over the input distribution, transforming them through ϕ , and then computing the empirical mean and covariance of the transformed points. [TODO: diagram of the whole process here]

By composing an unscented transform on the argument-space schmear to obtain a feature-space schmear with the above description of the output schmear obtained by applying a given linear model to a feature-space schmear, we are now able to straightforwardly impute output schmeared from arbitrary term applications.

5.2.2 Updating the Output Model

For the second stage of a prior update, we need to utilize the obtained output-space schmear over the embedding space $(Z \rightarrow W)$ to update our information about the model $w = A * \phi(z) + \epsilon$, which we assume has a $\mathcal{MN}\mathcal{IW}_{t,r}(B, \Lambda, V, v)$ prior over the relevant model parameters.

First, it is important to note that there is an inherent ambiguity in the interpretation of the covariance of a schmear in the embedding space for $(Z \rightarrow W)$ with respect to our usual model. Fundamentally, this is due to the fact that the covariance in the linear operator obtained from a matrix-normal inverse-Wishart distribution is only expressible as a *Kronecker product* of relevant input/output covariances, which is in general invariant with respect to alterations of the relative scaling of the input covariance and output covariance factors. [3]. Luckily, in our case, we may bypass this difficulty by assuming that the schmear for the prior update carries *no information* about the residual noise ϵ . This is a reasonable thing for us to do, since we only can get a sense of the "typical" size of the residual error from data updates, not prior updates.

After adopting the assumption that we gain no information about the noise covariance, we can note that the $\mathcal{MN}\mathcal{IW}_{t,r}$ likelihood may be separated back into its constituent matrix-normal and inverse-Wishart components, of which we only want to update the matrix-normal component, since the inverse-Wishart component entirely pertains to the

likelihood for the noise covariance. Recall that in our models, the matrix normal component of the likelihood is given by $\mathcal{MN}_{t,r}(B, \Lambda^+, \Sigma_\epsilon)$, where Σ_ϵ is the noise covariance. Ignoring scaling factors, the log-likelihood of this distribution is given by:

$$\mathcal{L}_{\mathcal{MN}}(X) = -\frac{1}{2} \langle \Sigma_\epsilon^+, (X - B) \Lambda (X - B)^T \rangle_F + C$$

for some normalizing constant $C \in \mathbb{R}$. Now, suppose that we also have a schmear $(\mu, \Sigma) \in \mathcal{S}^{t \times r}$ over the space that $\text{vec}(X)$ belongs to. We would like to obtain a reasonable log-likelihood function from the schmear which is amenable to manipulation when combined with the above log-likelihood for our pre-existing matrix-normal model. To do so, since we are assuming that our update will yield no additional information about the error covariance, a reasonable functional form for such a log-likelihood is given by another matrix normal log-likelihood:

$$\mathcal{L}_{\text{schmear}}(X) = -\frac{1}{2} \langle \Sigma_\epsilon^+, (X - M) \Lambda_\sigma (X - M)^T \rangle_F + C$$

for some $r \times r$ matrix Λ_σ which we will derive from Σ , where $\text{vec}(M) = \mu$ and C is once again a normalizing constant.

To obtain a suitable Λ_σ , we simply find where $\frac{V}{v-t-1} \otimes \Lambda_\sigma^+$ is closest to Σ in the Frobenius norm. This yields the simple expression:

$$\Lambda_\sigma = \frac{V \Sigma_T^+}{v - t - 1}$$

where Σ_T is the result of re-interpreting Σ as 4-tensor laid out in matrix form as $(t \times r) \times (t \times r)$ and transposing dimensions to re-order as $(t \times t) \times (r \times r)$, which is once again laid out in matrix form.

Once we have done this, we can convert the likelihood of the update to a form which is not conditional on the error covariance by packaging it as a matrix-normal inverse-Wishart distribution whose inverse-Wishart component is completely non-informative. The result of doing this is the distribution $\mathcal{MNIW}_{t,r}(B, \Lambda, V, v)$ where:

$$\begin{aligned} B &= M \\ \Lambda &= \Lambda_\sigma \\ v &= t \\ V &= \mathbf{0} \end{aligned}$$

Using the formulas expressed in [CITE own refl], we note that if we have two matrix-normal inverse-Wishart priors $\mathcal{MNIW}_{t,r}(B_1, \Lambda_1, V_1, v_1)$ and $\mathcal{MNIW}_{t,r}(B_2, \Lambda_2, V_2, v_2)$, we may combine the information from these priors into a new matrix-normal inverse-Wishart prior $\mathcal{MNIW}_{t,r}(B, \Lambda, V, v)$ using the \mathcal{MNIW} -sum as follows:

$$\begin{aligned} B &= (B_1 \Lambda_1 + B_2 \Lambda_2) \Lambda^+ \\ \Lambda &= \Lambda_1 + \Lambda_2 \\ v &= v_1 + v_2 + k \\ V &= V_1 + V_2 + (B_1 - B) \Lambda_1 (B_1 - B)^T + (B_2 - B) \Lambda_2 (B_2 - B)^T \end{aligned}$$

Consequently, we may use the above formulas to incorporate the information from the imputed \mathcal{MNIW} distribution into the model for the output term. If we want to later remove the information added to the model in this way (e.g: because we have more accurate information to apply instead), we can simply take the \mathcal{MNIW} distribution used in the update and add its *additive inverse* to the model parameters instead. In line with [CITE own reference], the additive inverse of a matrix-normal inverse-Wishart prior $\mathcal{MNIW}_{t,r}(B_+, \Lambda_+, V_+, v_+)$ is given by $\mathcal{MNIW}_{t,r}(B_-, \Lambda_-, V_-, v_-)$ where:

$$\begin{aligned} B_- &= B_+ \\ \Lambda_- &= -\Lambda_+ \\ v_- &= 2t - v_+ \\ V_- &= -V_+ \end{aligned}$$

5.3 Update Ordering

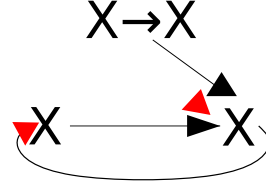
With our previously-described mechanisms for updating the models for function terms with respect to data updates and prior updates, we are now faced with the task of scheduling these updates during the regular operation of the interpreter. Since, for a function type $X \rightarrow Y$ where $Y = Z \rightarrow W$, data updates will depend on paired terms in X and Y , and prior updates will be issued back to the terms in Y which were outputs of the function. Consequently, there is an inherent cyclicity in the update process. Ideally, we would deal with this by deriving explicit expressions for the fixed point of the update process. However, given the intricacies of the processes we have described, it is highly unlikely that there is

such a closed-form expression for the fixed point. Instead, we will simply settle for an update process which will gradually approach a fixed point in the large-sample limit, as the number of interpreter-evaluated terms goes to infinity.

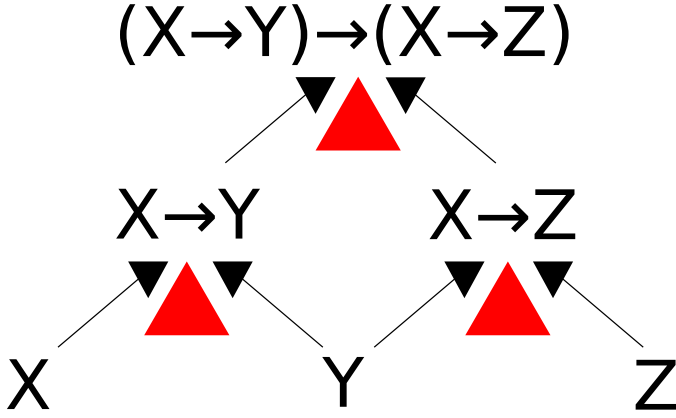
To do so, we first note that after each term application is evaluated by the interpreter, we would ideally like to only update those terms which need updating as a consequence of the just-performed evaluation. To this end, we define two kinds of passes over the structure defined by all term applications which have been evaluated in the past.

First, we define a *data update pass* as a pass where we take the newly-updated term applications from a previous pass, and update the models for the functions. Then, we repeatedly take the terms that were modified, and find all term applications involving them as either the argument or the result of an evaluation until there are no more terms to propagate data updates for. The following is a schematic for this pass:

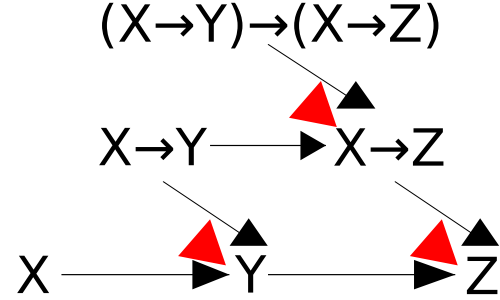
if we propagated updates whenever the argument embedding changed, we could easily run into an infinite loop, as in the following figure.



The following figure is a schematic of the prior update pass.



For the prior updates, we also define a *prior update pass* as a pass where we take all function terms modified in a previous pass and compute prior updates on their output spaces, if their output type is a function type. We then iterate this until we run out of terms to propagate prior updates for. It's important to note that we only perform prior updates when the *function* embedding changes. This is necessary because



In our reference implementation, after each batch of interpreter evaluations, we perform one data update pass followed by one prior update pass.

6 Additional Considerations

With the previous sections, we have elaborated upon the main content of the proposed framework for incrementally deriving

term embeddings in FETISH, and so a reader who is only interested in the intuition behind FETISH may skip straight to the "Conclusions and Further Research" section. However, in the exposition above, we have left a few small loose ends lying around which bear tying down.

6.1 Dimensionality Reduction

First, the previously-described scheme is actually intractable as written for higher-order function types, since e.g: if we assign a fixed number of feature dimensions for each input dimension, if embeddings for X have dimensionality on the order of t , then embeddings for $X \rightarrow X$ have dimensionality on the order of t^2 , embeddings for $(X \rightarrow X) \rightarrow (X \rightarrow X)$ have dimensionality on the order of t^4 , embeddings for $((X \rightarrow X) \rightarrow (X \rightarrow X)) \rightarrow ((X \rightarrow X) \rightarrow (X \rightarrow X))$ have dimensionality on the order of t^8 , and so on, in a pattern of exponentially-bigger spaces with each increase in function order.

Not only is this incredibly parameter and time-inefficient, but it is also in many ways counter to intuition about the kind of inductive biases we would expect for sufficiently high-order function terms. In particular, an argument from the vantage point of Wadler's "Theorems for Free" ([17]) is particularly appealing to the authors: Despite the fact that unlike in Wadler's situation, we do not have universally quantified function types, it still seems likely that the collection of "useful" or "natural" higher-order functions belongs to a much smaller sub-space than the exponentially-sized spaces which appear as a consequence of the previously-described framework.

Since we conjecture that the fundamental behavior of higher-order types is captured in a low-dimensional subspace of the "true" embedding space for such functions, we also conjecture that random projections [6] will do an excellent job of maintaining the salient features of such spaces. In particular, if we pick the target space of each random projection of the embeddings to be logarithmically-sized in the unprojected size of the embedding, we may completely side-step the exponential growth observed previously.

In particular, this modification (which is in our reference implementation) associates every function type $F = (X \rightarrow Y)$ where the embedding space of X is s -dimensional and the embedding space of Y is t -dimensional with a random projection operator $\pi : \mathbb{R}^{r \times t} \rightarrow \mathbb{R}^m$, where r is the number of features on the input space and m is the dimension of the projected

embedding space for F . To curb exponential growth in dimension, we propose setting $m \in O(\log(r) + \log(t))$. Then, whenever a schmear (μ, Σ) in the embedding space $\mathbb{R}^{r \times t}$ of F is mentioned in the above exposition, we simply replace it with the projected schmear:

$$\pi((\mu, \Sigma)) = (\pi\mu, \pi\Sigma\pi^T)$$

. If we ever need to re-expand the projected schmear into a full $r \times t$ -dimensional schmear, we can use the (pseudo-)inverse operation:

$$\pi^+((\mu_\pi, \Sigma_\pi)) = (\pi^+\mu_\pi, \pi^+\Sigma(\pi^+)^T)$$

In particular, a consequence of this replacement process is that when X and Y are function types, models in F are no longer mappings from embeddings of X to embeddings of Y , but are instead mappings from *projections of* embeddings in X to *projections of* mappings in Y . We refer the curious reader to our reference implementation for the nuts-and-bolts of this tweak.

6.2 Choice of Prior Distributions

While we have everywhere adopted a Bayesian framework for statistical inference, up to this point, we have made no mention of the priors which we propose to adopt for the involved \mathcal{MNTW} models of function embeddings. To resolve this quandary, we point out the following desirable properties of a choice of priors in our setting:

1. Model embedding schmears are well-defined for every model with at least one data-point.
2. Priors are isotropic over their embedding space and have zero mean to encourage small regression coefficients.
3. The prior distribution for embeddings in Y is close to the distribution over applications $f(x)$ where f is drawn from the prior over models for $X \rightarrow Y$ and x is drawn from the prior over embeddings for X .

To satisfy the first condition, since we are interested in priors expressible as $\mathcal{MNTW}_{t,r}(B, \Lambda, V, v)$, we simply note that schmears are well-defined for the distribution so long as $v > t + 1$, so for the prior, we may simply set $v = t + 1$ to satisfy the condition while minimizing the informativeness of our prior.

For the second condition, we may simply set $B = \mathbf{0}_{t \times r}$, $\Lambda = \sigma_\Lambda^{-1} I_r$, and $V = \sigma_V I_t$ where we are free in our choices of the scalars σ_Λ^{-1} and σ_V . For simplicity, our reference implementation always sets $\sigma_\Lambda = \sigma_V = \sigma$, which entails that we would roughly expect the covariance of the model schmeur to be near $\sigma^2(I_t \otimes I_r)$ after a single observation.

The third condition takes a bit more work to satisfy, but it ultimately reduces to some elementary algebraic manipulation. In particular, we assume that every embedding space’s prior covariances [after a single observation] will be of the above form, meaning that their overall variance is proportional to σ^2 for a constant σ which encodes the “expected” variance of embeddings in every space. However, since the models we consider for the function space $F = (X \rightarrow Y)$ involve an in general nonlinear feature mapping $\psi_F : \mathbb{R}^s \rightarrow \mathbb{R}^r$, satisfying the third condition above for our model $A * \psi_F(x) + \epsilon = y$ is a bit trickier, since we require the total variances of $x \in \mathbb{R}^s$, $y \in \mathbb{R}^t$, and the $t \times r$ matrix A to all be of the same scale. To do so, we propose applying a post-scaling to ψ_F dependent on its functional form to achieve the desired property. Our choices of scaling factors by feature type are listed in the table below, with derivations given in Appendix B:

Features	Scaling Factor
Linear	$\frac{1}{\sigma} \sqrt{\frac{2}{s}}$
Quadratic	$\frac{\sqrt{2}}{s\sigma^2}$
Fourier	$\sqrt{\frac{2}{r}}$

7 Conclusions and Directions for Further Research

We have presented an incremental procedure to keep distributional embeddings of terms in a simply-typed combinatory calculus updated in a Bayesian fashion. The authors hope that this procedure yields the inspiration necessary for an entirely new class of machine-learning systems to take root.

But before discussing such potential applications, the authors think it prudent to suggest extensions to the basic framework presented in this paper. In particular, the framework we have discussed is far from unique – different choices could have been made at many steps along the way. For

instance, instead of the matrix-normal inverse-Wishart conjugate Bayesian models we have adopted for the functional embeddings, a more expressive class of distributions may have been chosen instead. For another example, instead of using the unscented transform to propagate uncertainty through nonlinear functions, we could have instead used a local Taylor series approximation, drawn empirical samples from the distribution instead of using sigma points, or even regularized the estimates to a degree to reduce the risk of overconfidence in our estimates. However, we largely consider such things to be incremental improvements.

More intriguing to the authors is the potential for extending our framework to languages with type systems with more expansive capabilities. A very simple, yet impactful extension that the authors have in mind is to extend the framework to allow for types of finite, unbounded-length sequences, such as lists and strings. However, the authors also wonder if the framework may be extended to universal types, which could provide better generalization behavior over types, or if the framework could be extended even further in some fashion to allow some analogue of dependent typing.

On applications, the authors believe that the most immediate task at hand is to leverage the proposed framework to tackle supervised learning problems. In particular, the authors believe that some variant of Monte Carlo Tree Search which takes into account the probability distributions defined by the embeddings could be a fruitful avenue in this pursuit. Here, the benefits of the proposed framework really show, since if supervised learning is tackled under this framework, the nature of this framework virtually guarantees that multi-task learning will follow soon after.

Taking a longer view, the authors are particularly interested in applications of the proposed framework to general reinforcement learning problems. On the purely-theoretical side, the authors ask whether there is a version of the universal artificial general intelligence procedure AIXI to the setting we have described, of a combinatorial base language for which the distributional information maintained over programs is fundamentally spatial and kept in the form of embeddings through probability distributions. On the practical side, the authors are also curious as to whether there are any immediate and pragmatic fusions of the proposed framework with typical algorithms from the reinforcement learning setting. We leave this rich area for exploration open to other daring researchers who wish to investigate one of the most important topics of our time leveraging the framework which we have described.

Our reference implementation in Rust is available under the MIT License at <https://github.com/bubble-07/FETISH-RS>.

References

- [1] Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sable-Meyer, Luc Cary, Lucas Morales, Luke Hewitt, Armando Solar-Lezama, and Joshua B Tenenbaum. Dream-coder: Growing generalizable, interpretable knowledge with wake-sleep bayesian program learning. *arXiv preprint arXiv:2006.08381*, 2020.
- [2] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. CRC press, 2013.
- [3] Hunter Glanz and Luis Carvalho. An expectation-maximization algorithm for the matrix normal distribution. *arXiv preprint arXiv:1309.6609*, 2013.
- [4] Sumit Gulwani, Oleksandr Polozov, Rishabh Singh, et al. Program synthesis. *Foundations and Trends® in Programming Languages*, 4(1-2):1–119, 2017.
- [5] Marcus Hutter. A theory of universal artificial intelligence based on algorithmic complexity. *arXiv preprint cs/0004001*, 2000.
- [6] William B Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics*, 26(189-206):1, 1984.
- [7] Simon J Julier and Jeffrey K Uhlmann. New extension of the kalman filter to nonlinear systems. In *Signal processing, sensor fusion, and target recognition VI*, volume 3068, pages 182–193. International Society for Optics and Photonics, 1997.
- [8] Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric P Xing. Neural architecture search with bayesian optimisation and optimal transport. In *Advances in neural information processing systems*, pages 2016–2025, 2018.
- [9] Kaare Brandt Petersen and Michael Syskind Pedersen. The matrix cookbook, nov 2012. URL <http://www2.imm.dtu.dk/pubdb/p.php, 3274:14>, 2012.
- [10] Ninh Pham and Rasmus Pagh. Fast and scalable polynomial kernels via explicit feature maps. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 239–247, 2013.
- [11] Thomas Pierrot, Guillaume Ligner, Scott E Reed, Olivier Sigaud, Nicolas Perrin, Alexandre Laterre, David Kas, Karim Beguir, and Nando de Freitas. Learning compositional neural programs with recursive tree search and planning. In *Advances in Neural Information Processing Systems*, pages 14673–14683, 2019.
- [12] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.
- [13] Scott Reed and Nando De Freitas. Neural programmer-interpreters. *arXiv preprint arXiv:1511.06279*, 2015.
- [14] Peter E Rossi, Greg M Allenby, and Rob McCulloch. *Bayesian statistics and marketing*. John Wiley & Sons, 2012.
- [15] Stanley Sawyer. Wishart distributions and inverse-wishart sampling. URL: www.math.wustl.edu/sawyer/hmhandouts/Whishart.pdf, 2007.
- [16] Ray J Solomonoff. A formal theory of inductive inference. part i. *Information and control*, 7(1):1–22, 1964.
- [17] Philip Wadler. Theorems for free! In *Proceedings of the fourth international conference on Functional programming languages and computer architecture*, pages 347–359, 1989.
- [18] Martin Wistuba, Amrith Rawat, and Tejaswini Pedapati. A survey on neural architecture search. *arXiv preprint arXiv:1905.01392*, 2019.

Appendix A Derivation for Output Schmeat Imputation

Consider the linear model $Az + \epsilon = y$ where A is $t \times r$ and we employ a $MNTW_{t,r}(B, \Lambda, V, v)$. Suppose further that we have an input schmeat given by $E[z] = \mu$ and $Cov[z] = \Sigma$.

First, note that using an observation we previously made about the \mathcal{MNTW} distribution, we can notice that:

$$E[vec(A)] = vec(B)$$

$$Cov[vec(A)] = \frac{V}{v-t-1} \otimes \Lambda^+$$

Through elementary manipulation [see, e.g: here], we can see that:

$$Cov[Az]_{kl} = \sum_i \sum_j Cov[A_{ki}, A_{jl}] (\Sigma_{ij} + \mu_i \mu_j) + B_{ki} B_{lj} \Sigma_{ij}$$

Since $Cov[A_{ki}, A_{jl}] = \frac{1}{v-t-1} V_{kl} \Lambda_{ij}^+$. Re-expressing the outer double-summation in matrix operations, we arrive at the claimed result

$$Cov[Az] = [\langle \Sigma, \Lambda^+ \rangle_F + \mu^T \Lambda^+ \mu] \frac{V}{v-t-1} + B \Sigma B^T$$

Appendix B Derivations of Scaling Factors for Common Feature Maps

Before deriving the formulas for scaling factors, we first note the elementary result that if M is a random $t \times r$ matrix whose entries are i.i.d. $N(0, \sigma_M^2)$, and x is a random r -dimensional vector whose entries are i.i.d. $N(0, \sigma_x^2)$, then

$$E[Mx] = 0, \quad Cov[Mx] = \frac{r}{2} \sigma_M^2 \sigma_x^2 I_t$$

B.1 Linear Features

For the case of linear features, we can directly apply the above formula, where we want:

$$Cov[\alpha Mx] = \alpha \frac{r}{2} \sigma^4 I_t = \sigma^2 I_t$$

Hence, since for this case $r = s$ we obtain the formula in the table,

$$\alpha = \frac{1}{\sigma} \sqrt{\frac{2}{s}}$$

B.2 Quadratic Features

For the case of the quadratic features $x \mapsto xx^T$, we may use the above formula to note that:

$$Cov[vec(xx^T)] = \frac{1}{2} \sigma_x^4 I_{s*s}$$

Upon another application of the above formula, we see that the condition we want is:

$$Cov[\alpha Mvec(xx^T)] = \alpha \frac{t^2}{2} \sigma^6 = \sigma^2$$

Which implies that we need to set:

$$\alpha = \frac{\sqrt{2}}{s\sigma^2}$$

B.3 Fourier Features

For the case of random Fourier features, we suppose that we're using the fourier features which are composed of pairings of $\sin(w\theta)$ with $\cos(w\theta)$. Using the fact that the squares of these functions add to one, we obtain the simple formula in the table for the scaling factor:

$$\alpha = \sqrt{\frac{2}{r}}$$