

GP-select: Accelerating EM using adaptive subspace preselection

Jacquelyn A. Shelton¹, Jan Gasthaus^{2, 3}, Zhenwen Dai^{3, 4}, Jörg Lücke⁵, and Arthur Gretton²

¹ Technical University Berlin, Marchstrasse 23, 10587 Berlin, Germany.

² University College London, Gatsby Unit, 25 Howland Street, London W1T 4JG, UK.

³ Amazon Development Center, Karl-Liebknecht-Str. 5, 10178 Berlin, Germany.*

⁴ University of Sheffield, Western Bank, Sheffield, South Yorkshire S10 2TN, UK.

⁵ University of Oldenburg, Ammerländer Heerstr. 114, 26129 Oldenburg, Germany.

Keywords: Approximate inference, generative graphical models, latent variable models, Expectation Maximization, EM acceleration, variable preselection

Abstract

We propose a nonparametric procedure to achieve fast inference in generative graphical models when the number of latent states is very large. The approach is based on iterative

*This work was done prior to joining Amazon.

latent variable preselection, where we alternate between learning a 'selection function' to reveal the relevant latent variables, and using this to obtain a compact approximation of the posterior distribution for EM; this can make inference possible where the number of possible latent states is e.g. exponential in the number of latent variables, whereas an exact approach would be computationally infeasible. We learn the selection function entirely from the observed data and current EM state via Gaussian process regression. This is by contrast with earlier approaches, where selection functions were manually-designed for each problem setting. We show that our approach performs as well as these bespoke selection functions on a wide variety of inference problems: in particular, for the challenging case of a hierarchical model for object localization with occlusion, we achieve results that match a customized state-of-the-art selection method, at a far lower computational cost.

1 Introduction

Inference in probabilistic graphical models can be challenging in situations where there are a large number of hidden variables, each of which may take on one of several state values. The Expectation Maximization (EM) algorithm is widely applied to learn model parameters when hidden variables are present, however inference can quickly become intractable as the dimensionality of hidden states increases. Consider, for instance, the floor of a nursery populated with different toys, and images of this floor large enough to contain a number of toys. A nursery easily contains a hundred different toys and any subset of these hundred toys may appear in any image. For one hundred toys there

is therefore a combinatorics of 2^{100} different combinations of toys that can make up an image. An inference task may now be to infer, for any given image, the toys it contains. If we approached this task using a probabilistic graphical model, we would define a basic such model using a set of one hundred hidden variables (one for each toy). Given a specific image, inference would then take the form of computing the posterior probability for any combination of toys, and from this, e.g., the probability of each toy to be in the image can be computed. If done exactly, this process needs to evaluate all the 2^{100} different toy combinations which easily exceeds currently available computational resources.

While there are also many tasks for which graphical models with few latent variables are sufficient, the requirement for many hidden variables (as in the toy example) is typical for visual, auditory and many other types of data with very rich structure. Graphical models for such data are often a central building block for tasks such as denoising (Elad and Aharon, 2006; Titsias and Lázaro-Gredilla, 2011), inpainting (Mairal et al., 2009b,a; Titsias and Lázaro-Gredilla, 2011), classification (Raina et al., 2007), or collaborative filtering (Titsias and Lázaro-Gredilla, 2011). Typically, the performance in these tasks improves with the number of latent variables that can be used (and which is usually limited by computational demands).

Expectation truncation (ET) (Lücke and Eggert, 2010) is an approximate EM algorithm for accelerating inference and learning in graphical models with many latent variables. Its basic idea is to restrict the inference performed during the E-step to an “interesting” subset of states of the latent variables, chosen per data point according to a *selection function*. This subspace reduction can lead to a significant decrease in com-

putational demand with very little loss of accuracy (compared with the full model). To provide an intuition: For the toy example, we could for instance first analyze the colors contained in a given image. If the image did not contain the color "red", we could already assume red toys or partly red toys to be absent. Only in a second step would we then consider the combinatorics of the remaining toys. More features and more refined features would allow for a reduction to still smaller sets of toys until the combinatorics of these selected toys becomes computationally tractable. The selection function of expectation truncation mathematically models the process of selecting the relevant hidden variables (the relevant toys); while truncated posterior distributions then models their remaining combinatorics (see further below).

In previous work, functions to select states of high posterior mass were derived individually for each graphical model of interest, e.g., by taking upper bounds or noiseless limits (Lücke and Eggert, 2010; Shelton et al., 2012; Bornschein et al., 2013; Henniges et al., 2014; Sheikh et al., 2014). The crucial underlying assumption remains that when EM has converged, the posterior mass is concentrated in small volumes of the latent state space (see, e.g., Lücke and Eggert, 2010; Sheikh et al., 2014, for discussions). We can expect the approximation to be accurate only if restricting the combinatorics (e.g., combinations of a restricted number of toys) does not miss large parts of posterior mass. This property is observed to hold, however, for many types of data in the auditory, visual or general pattern recognition domains.

The definition of appropriate selection functions for basic graphical models (such as the nursery floor example) is already non-trivial. For models incorporating more detailed data properties, the definition of selections functions becomes still more demand-

ing. For visual data, e.g., models that also capture mutual object occlusions (Henniges et al., 2014) and/or the object position (Dai and Lücke, 2014), the design of suitable selection functions is extremely challenging: it requires both expert knowledge on the problem domain and considerable computational resources to implement (indeed, the design of such functions for particular problems has been a major contribution in previous work on the topic).

In the present work, we propose a generalization of the ET approach, where we completely avoid the challenge of problem-specific selection function design. Instead, we learn selection functions adaptively and non-parametrically from the data, while learning the model parameters simultaneously using EM. We emphasize that the selection function is used only to "guide" the underlying base inference algorithm to regions of high posterior probability, but is not itself used as an approximation to the posterior distribution. As such, the learned function does not have to be a completely accurate indication of latent variable predictivity, as long as the *relative importance* of the latent states likely to contribute posterior probability mass is preserved. We use Gaussian process regression (Rasmussen and Williams, 2005) to learn the selection function – by regressing the expected values of the latent variables onto the observed data – though other regression techniques could also be applied. The main advantage of GPs is that they do not need to be re-trained when only the output changes, as long as the inputs remain the same. This makes adaptive learning of a changing target function (given fixed inputs) computationally trivial. We term this part of our approach *GP-select*. Our nonparametric generalization of ET may be applied as a black-box meta algorithm for accelerating inference in generative graphical models, with no expert knowledge re-

quired.

Our approach is the first to make ET a general purpose algorithm for discrete latent variables, whereas previously, ET had to be modified by hand for each latent variable model addressed. For instance, in Section 5.3 we will show that preselection is crucial for efficient inference in complex models. Although ET has already been successful in some models, this work shows that more complex models will crucially depend on an improved selection step and focuses on automating this step.

For empirical evaluation, we have applied GP-select in a number of experimental settings. First, we considered the case of sparse coding models (binary sparse coding, spike-and-slab, nonlinear spike-and-slab), where the relationship between the observed and latent variables is known to be complex and nonlinear.[†] We show that GP-select can produce results with equal performance to the respective manually-derived selection functions. Interestingly, we find it can be essential to use nonlinear GP regression in the spike-and-slab case, and that simple linear regression is not sufficiently flexible in modeling the posterior shape. Second, we illustrate GP-select on a simple Gaussian mixture model, where we can provide intuition and explicitly visualize the form of the learned regression function. We find that even for a simple model, it can be essential to learn a nonlinear mapping. Finally, we present results for a recent hierarchical model for translation invariant occlusive components analysis (Dai and Lücke, 2014). The performance of our inference algorithm matches that of the complex hand-engineered selection function of the previous work, while being straightforward to implement and having a far lower computational cost.

[†]Note that even when linear relations exist between the latents and outputs, a nonlinear regression may still be necessary in finding relevant variables, as a result of explaining away.

2 Related work

The general idea of aiding inference in graphical models by learning a function that maps from the observed data to a property of the latent variables is quite old. Early work includes the Helmholtz machine (Dayan et al., 1995) and its bottom-up connections trained using the wake-sleep algorithm (Hinton et al., 1995). More recently, the idea has surfaced in the context of learning variational distributions with neural networks (Kingma and Welling, 2014). A two-stage inference procedure has been discussed in the context of computer vision (Yuille and Kersten, 2006) and neural inference (Körner et al., 1999). Recently, researchers (Mnih and Gregor, 2014) have generalized this idea to learning in arbitrary graphical models by training an “inference network” that efficiently implements sampling from the posterior distribution.

GPs have recently been widely used to ”learn” the results of complicated models in order to accelerate inference and parameter selection. GP approximations have been used in lieu of solving complex partial differential equations (Sacks et al., 1989; Currin et al., 1991), to learn data-driven kernel functions for recommendation systems (Schwaighofer et al., 2004), and recently for quantum chemistry (Rupp et al., 2012). Other work has used GPs to simplify computations in approximate Bayesian computation (ABC) methods: namely to model the likelihood function for inference (Wilkinson, 2014), to aid in making Metropolis-Hastings (MH) decisions (Meeds and Welling, 2014), and to model the discrepancies between simulated/observed data in parameter space simplification (Gutmann and Corander, 2015). Recently, instead of the typical choice of GPs for large scale Bayesian optimization, neural networks have been used to learn an adaptive set of basis functions for Bayesian linear regression (Snoek et al.,

2015).

Our work follows the same high level philosophy in that we use GPs to approximate complex/intractable probabilistic models. None of the cited prior work address our problem setting, namely the selection of relevant latent variables by learning a non-parametric relevance function, for use in expectation truncation (ET).

3 Variable selection for accelerated inference

Notation. We denote the observed data by the $D \times N$ matrix $\mathbf{Y} = (\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(N)})$, where each vector $\mathbf{y}^{(n)} = (y_1^{(n)}, \dots, y_D^{(n)})^T$ is the n th observation in a D -dimensional space. Similarly we define corresponding binary latent variables by the matrix $\mathbf{S} = (\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)}) \in \{0, 1\}^{H \times N}$ where each $\mathbf{s}^{(n)} = (s_1^{(n)}, \dots, s_H^{(n)})^T \in \{0, 1\}^H$ is the n th vector in the H -dimensional latent space, and for each individual hidden variable $h = 1, \dots, H$, the vector $\mathbf{s}_h = (s_h^{(1)}, \dots, s_h^{(N)}) \in \{0, 1\}^N$. The number of dimensions in the reduced latent space is denoted by H' , where $H' \ll H$. Note that although we restrict ourselves to binary latent variables here, the procedure could in principle be generalized to variables with higher cardinality (e.g. see (Exarchakis et al., 2012)). We denote the prior distribution over the latent variables as $p(\mathbf{s}|\theta)$ and the likelihood of the data as $p(\mathbf{y}|\mathbf{s}, \theta)$. Using these expressions, the posterior distribution over latent variables is

$$p(\mathbf{s}^{(n)}|\mathbf{y}^{(n)}, \Theta) = \frac{p(\mathbf{s}^{(n)}|\Theta) p(\mathbf{y}^{(n)}|\mathbf{s}^{(n)}, \Theta)}{\sum_{\mathbf{s}'} p(\mathbf{s}'|\Theta) p(\mathbf{y}^{(n)}|\mathbf{s}', \Theta)}. \quad (1)$$

3.1 Selection via Expectation Truncation in EM

Expectation Maximization (EM) is an iterative algorithm to optimize the model parameters of a given graphical model (see e.g. (Dempster et al., 1977; Neal and Hinton, 1998)). EM iteratively optimizes a lower bound on the data likelihood by inferring the posterior distribution over hidden variables given the current parameters (the E-step), and then adjusting the parameters to maximize the likelihood of the data averaged over this posterior (the M-step). When the number of latent states to consider is large (e.g. exponential in the number of latent variables), the computation of the posterior distribution in the E-step becomes intractable and approximations are required.

Expectation truncation (ET) is a meta algorithm, which improves convergence of the expectation maximization (EM) algorithm (Lücke and Eggert, 2010). The main idea underlying ET is that the posterior probability mass is concentrated in a small subspace of the full latent space. This is the case, for instance, if for a given data point $\mathbf{y}^{(n)}$ only a subset of the H latent variables $s_1^{(n)}, s_2^{(n)}, \dots, s_H^{(n)}$ are relevant. Even when the probability mass is supported everywhere, it may still be largely concentrated on a small number of the latents.

A *selection function* can be used to identify a subset $I \subseteq \{1, 2, \dots, H\}$ of size H' ($H' \ll H$) of salient variables, which in turn is used to define a subset $\mathcal{K}_n = \{\mathbf{s} \mid \mathbf{s} \in 2^H \wedge \forall h \notin I : s_h = 0\} \subseteq \{0, 1\}^H$ of the possible state configurations of the latent variables $\mathbf{s}^{(n)}$ for each data point. This subset contains only state configurations where the values of the variables not identified to be relevant by the selection function are fixed to 0. The posterior distribution in Equation (1) can then be approximated by a truncated posterior distribution, computed on the reduced support,

$$\begin{aligned}
& p(\mathbf{s}^{(n)} | \mathbf{y}^{(n)}, \Theta) \\
& \approx q_n(\mathbf{s}^{(n)}; \Theta) = \frac{p(\mathbf{s}^{(n)}, \mathbf{y}^{(n)} | \Theta) \mathbb{I}(\mathbf{s}^{(n)} \in \mathcal{K}_n)}{\sum_{\mathbf{s}' \in \mathcal{K}_n} p(\mathbf{s}', \mathbf{y}^{(n)} | \Theta)},
\end{aligned} \tag{2}$$

where $\mathbb{I}(\mathbf{s} \in \mathcal{K}_n) = 1$ if $\mathbf{s} \in \mathcal{K}_n$ is true, and 0 otherwise. In other words, Equation (2) is proportional to Equation (1) if $\mathbf{s}^{(n)} \in \mathcal{K}_n$ (and zero otherwise), so that the approximate posterior $q_n(\mathbf{s}^{(n)}; \Theta)$ assigns zero mass to states $\mathbf{s}^{(n)} \notin \mathcal{K}_n$. The set \mathcal{K}_n contains only states for which $s_h = 0$ for all h that are not selected, i.e. all states where $s_h = 1$ for some non-selected h are assigned zero probability. This means that there are fewer terms in the denominator of Equation (2) compared with Equation (1), thus reducing the computational complexity. Equation (2) still remains proportional to Equation (1) for the remaining states $\mathbf{s} \in \mathcal{K}_n$, however. As there are only $2^{H'}$ terms in the sum over \mathcal{K}_n , computing this posterior approximation is much more efficient than computing the exact normalizing constant for the full posterior (containing 2^H terms). The number of latent dimensions to select, H' , is chosen based on the compute resources available: i.e. as large as resources allow in order to be closer to true EM, although empirically it has been shown that much smaller values suffice (see e.g. Sheikh et al., 2014, App. B on complexity-accuracy trade-offs).

3.2 ET with affinity

One way of constructing a selection function is by first ranking the latent variables according to an *affinity function* $f_h(\mathbf{y}^{(n)}) : \mathbb{R}^D \mapsto \mathbb{R}$ which directly reflects the relevance of the latent variable $s_h^{(n)}$. A natural choice for such a function is the one that

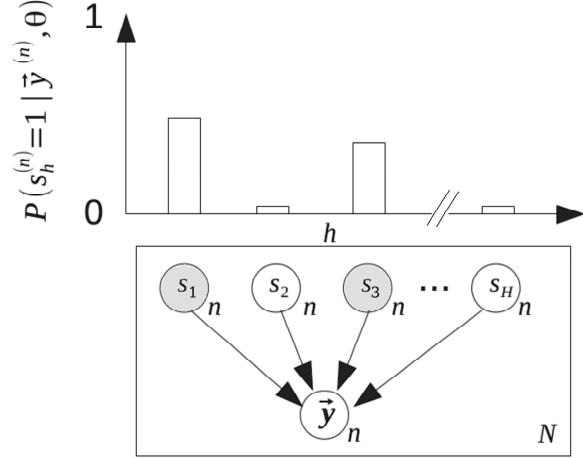


Figure 1: Illustration of the affinity function for selection. The affinity approximates the marginal posterior probability of each $h = 1, \dots, H$ latent variable (top), which corresponds to the most relevant variables for a given data point $\mathbf{y}^{(n)}$ (bottom). Here, the variables s_1 and s_3 yield high affinity and would thus be considered relevant for $\mathbf{y}^{(n)}$.

approximates the marginal posterior probability of each variable, e.g. we try to learn f as follows:

$$f_h(\mathbf{y}^{(n)}) = \hat{p}_h^{(n)} \approx p_h^{(n)} \equiv p(s_h^{(n)} = 1 | \mathbf{y}^{(n)}, \Theta), \quad (3)$$

meaning that the relevant variables will have greater marginal posterior probability $p_h^{(n)}$.

See Figure 1 for a simplified illustration. When the latent variables $s_{h=1}^{(n)}, \dots, s_H^{(n)}$ in the marginal posterior probability $\hat{\mathbf{p}}^{(n)} = \hat{p}_{h=1}^{(n)}, \dots, \hat{p}_H^{(n)}$ are conditionally independent given a data point $\mathbf{y}^{(n)}$, this affinity function correctly isolates the most relevant variables in the posterior. To see this, consider the full joint $p(s_1, \dots, s_h | \mathbf{y}, \Theta)$ in the case when a subset of latents has values clamped to zero, i.e., $s_h = 0$ for all $h \notin I$ (compare Equation (2)). We can then ask what the overall joint posterior mass is in this case. If we suppose the latents to be conditionally independent, this total mass is given by a

product of marginals as follows:

$$\sum_{\mathbf{s} \text{ with } s_h=0 \text{ for all } h \notin I} p(s_1, \dots, s_H | \mathbf{y}, \Theta) = 1 - \prod_{h \notin I} p(s_h = 1 | \mathbf{y}, \Theta). \quad (4)$$

We want this mass to be as large as possible as its complement is the posterior mass that we discard with our approximation. If the affinity function correctly estimates the marginals $p(s_h = 1 | \mathbf{y}, \Theta)$, then discarding those $(H - H')$ marginal with lowest values is equivalent to discarding the space with the least posterior mass (compared to discarding w.r.t. all alternative choices with the same number of latents). Even when this strong assumption does not hold in practice (which is often the case), however, the affinity can still correctly highlight relevant variables, and has been empirically shown to be quite effective when dependencies exist (see e.g. the source separation tasks in (Sheikh et al., 2014)).

Next, using all $\hat{p}_{h=1}^{(n)}, \dots, \hat{p}_H^{(n)}$ from the affinity function $\mathbf{f}(\mathbf{y}^{(n)}) = (f_1(\mathbf{y}^{(n)}), \dots, f_H(\mathbf{y}^{(n)}))$, we define $\gamma(\hat{\mathbf{p}}^{(n)})$ to simultaneously sort the indices of the latent variables in descending order [of probability $\hat{p}^{(n)}$] and *reduce* the sorted set to the H' highest (most relevant) variables' indices. $\gamma(\hat{\mathbf{p}}^{(n)})$ thus returns the H' selected variable indices I chosen by the affinity to be relevant to the n th data point. To ensure that there is a non-zero probability of selecting each variable per EM iteration, 10% of the H' indices are uniformly chosen from H at random. This prevents the possible propagation of errors from $q_{(n)}$ continuously assigning small probabilities to a variable s_h in early EM iterations. The rationale for this is that the optimization of $q_{(n)}$ in early iterations of EM starts from randomly initialized s_h . If the affinity function itself is based on the posterior approximation (as

it will be in the algorithm described in Section 4), it has a tendency to not select indices that were previously not selected. Thus in order for selection-based EM to not “get stuck”, it is important to select a few extra hidden indices randomly, to give the algorithm an opportunity to evaluate possibly unused variables which might be relevant for $\mathbf{y}^{(n)}$.

Finally, using the indices I from γ , we define $\mathcal{I}(I)$ to return an H' -dimensional subset of selected relevant latent states \mathcal{K}_n for each data point $\mathbf{y}^{(n)}$. All ‘non-relevant’ variable states s_h for all variables $h \notin I$ are effectively set to 0 in Equation (2) by not being present in the state set \mathcal{K}_n . For example, let’s say that there are five s_h , where $h \in \{1, \dots, 5\}$. We consider the case where only s_1 and s_2 are selected. The \mathcal{I} function will then return zeros for s_3 , s_4 , and s_5 , but will return both allowed possibilities 0 or 1 for s_1 and s_2 . Thus a valid setting for the entire vector \mathbf{s} can be $\mathbf{s} = [01000]$, but not $\mathbf{s} = [01100]$.

Using \mathbf{f} , \mathcal{I} , and γ , we can define a *selection function* $\mathcal{S} : \mathbb{R}^D \mapsto 2^{\{1, \dots, H\}}$ to select subsets \mathcal{K}_n per data point $\mathbf{y}^{(n)}$. Again, the goal is for the states \mathcal{K}_n to contain most of the probability mass $p(\mathbf{s} | \mathbf{y})$ and to be significantly smaller than the entire latent space. The *affinity based selection function* to obtain the set of states \mathcal{K}_n can be expressed as

$$\mathcal{S}(\mathbf{y}^{(n)}) = \mathcal{I}[\gamma[\mathbf{f}(\mathbf{y}^{(n)})]] = \mathcal{K}_n. \quad (5)$$

To summarize, the main task is to formulate a general data-driven function to identify relevant latent variables and to select the corresponding set of states \mathcal{K}_n . This is performed using GP regression in order to compute the truncated posterior Equation (2)

on the reduced support \mathcal{K}_n . With the combined effort of the above utility functions, we have concisely defined the function $\mathcal{S}(\mathbf{y}^{(n)})$ in Equation (5) to perform this selection.

3.3 Inference in EM with selection

In each iteration of EM, the following occurs: prior to the E-step, the selection function $\mathcal{S}(\mathbf{y}^{(n)})$ in (5) is computed to select the most relevant states \mathcal{K}_n , which are then used to compute the truncated posterior distribution $q_n(\mathbf{s})$ in (2). The truncated posterior can be computed using any standard inference method, such as exact inference or e.g. Gibbs sampling from $q(\mathbf{s})$ if inference is still intractable or further computational acceleration is desired. The result of the E-step is then used to update the model parameters with maximum likelihood in the M-step.

4 GP-Select

In previous work, the selection function $\mathcal{S}(\mathbf{y}^{(n)})$ was a deterministic function derived individually for each model (see e.g. Shelton et al., 2011, 2012; Dai and Lücke, 2012a,b; Bornschein et al., 2013; Sheikh et al., 2014; Shelton et al., 2015), specific examples of which will be shown in Section 5.1. We now generalize the selection approach: instead of predefining the form of \mathcal{S} for variable selection, we want to learn it in a black-box and model-free way based on the data. We learn \mathcal{S} using Gaussian process (GP) regression (e.g. Rasmussen and Williams, 2005), which is a flexible nonparametric model and scales cubically[‡] with the number of data points N but linearly with the number of latent

[‡]If the scaling with N is still too expensive, an incomplete Cholesky approximation is used, with cost linear in N and quadratic in the rank Q of the approximation (see Section 5.3 for details).

variables H . We define the affinity function f_h as being drawn from a Gaussian process model: $f_h(\mathbf{y}^{(n)}) \sim \text{GP}(0, k(\cdot, \cdot))$, where $k(\cdot, \cdot)$ is the covariance kernel, which can be flexibly parameterized to represent the relationship between variables. Again, we use f_h to approximate the marginal posterior probability p_h that $s_h^{(n)} = 1$. A nice property of Gaussian processes is that the kernel matrix K need only be computed once (until the kernel function hyperparameters are updated) to approximate $p_h^{(n)}$ for the entire $H \times N$ set of latent variables \mathbf{S} .

Thus, prior to each E-step in each EM iteration, within each calculation of the selection function, we calculate the affinity using a GP to regress the expected values of the latent variables $\langle \mathbf{S} \rangle$ from the observed data \mathbf{Y} . Specifically, we train on p_h from the previous EM iteration (where p_h is equal to $\langle s_h \rangle$), for training data of $\mathcal{D} = \{(\mathbf{y}^{(n)}, \langle \mathbf{s}^{(n)} \rangle_{q_n(\mathbf{s}^{(n)})}) | n = 1, \dots, N\}$, where we recall that $q_n(\mathbf{s}^{(n)})$ is the approximate posterior distribution for $\mathbf{s}^{(n)}$ in Equation (2). Note that we do not use a sigmoid link, hence this is clearly not a correct estimate of a probability (it can be negative, or greater than one). From the selection perspective, however, it is not necessary to avoid these pathologies, as we only want an ordering of the variables. A correct GP classification approach with a properly defined likelihood will no longer have a marginal Gaussian distribution, and we would no longer be able to trivially express the posterior means of different functions with the same inputs, without considerable extra computation.

In the first EM iteration, the expectations $\langle \mathbf{s}^{(n)} \rangle_q$ are initialized randomly; in each subsequent EM iteration, the expectations w.r.t. the \mathcal{K}_n -truncated posterior $q(\mathbf{s})$ are used. The EM algorithm is run for T iterations and the hyperparameters of the kernel

are optimized by maximum likelihood every T^* EM iterations.

For each data point n and latent variable s_h we compute the predicted mean of the GP by leaving this data point out of the training set and considering all others, which is called leave-one-out (LOO) prediction. It can be shown that this can be implemented efficiently (see Section 5.4.2 in Rasmussen and Williams, 2005), and we use this result to update the predicted affinity as follows:

$$\hat{p}_h^{(n)} \leftarrow \langle s_h^{(n)} \rangle_{q_n} - \frac{[K^{-1} \langle \mathbf{s}_h \rangle_{q_n}]_{nn}}{[K^{-1}]_{nn}}. \quad (6)$$

Equation (6) can be efficiently implemented for all latent variables $h = 1, \dots, H$ and all data points $n = 1, \dots, N$ using matrix operations, thereby requiring only one kernel matrix inversion for the entire dataset.

Substituting Equation (6) for \mathbf{f} in the affinity based selection function (5) ,

$$\mathcal{S}(\mathbf{y}^{(n)}) = \mathcal{I} \left[\gamma \left[\langle s_h^{(n)} \rangle_{q_n} - \frac{[K^{-1} \langle \mathbf{s}_h \rangle_{q_n}]_{nn}}{[K^{-1}]_{nn}} \right] \right] = \mathcal{I} [\gamma [\mathbf{f}(\mathbf{y}^{(n)})]] = \mathcal{K}_n$$

we call the entire process *GP-select*. An outline is shown in Algorithm 1.

Algorithm 1 GP-Select to accelerate inference in Expectation Maximization

```

for EM iterations  $t = 1, \dots, T$  do
    initialize all latent variables expectations  $\langle s_h^{(n)} \rangle_{q_{n,t}}$ 
    for data point  $n = 1, \dots, N$  do
        compute affinity of all latent variables  $\hat{\mathbf{p}}_t^{(n)}$ : Eq. (6)
        compute subset of relevant states  $\mathcal{S}$ : Eq. (5)
        compute truncated posterior  $q_{n,t}(\mathbf{s}^{(n)})$  in E-step: Eq. (2)
        update model parameters in M-step, e.g. as in Sec. 5.1
        store  $\langle s_h^{(n)} \rangle_{q_{n,t}}$  for  $\mathbf{p}^{(n)}$  in EM iteration  $t + 1$ 
    end for
    optimize kernel hyperparams every  $T^*$  EM iterations
end for
```

5 Experiments

We apply our GP-select inference approach to five different probabilistic generative models. First, we considered three sparse coding models (binary sparse coding, spike-and-slab, nonlinear spike-and-slab), where the relationship between the observed and latent variables is known to be complex and nonlinear. Second, we apply GP-select to a simple Gaussian mixture model, to both provide functional intuition of approach and to explicitly visualize the form of the learned regression function. Finally, we apply our approach to a recent hierarchical model for translation invariant occlusive components analysis (Dai and Lücke, 2012a; Dai et al., 2013; Dai and Lücke, 2014).

5.1 Sparse coding models

Many types of natural data are composed of potentially many component types, but any data point often only contains a very small number of this potentially large set of components. For the introductory example of toys on the nursery floor, for instance, there are many different toys that can *potentially* be in a given image but there is typically only a relatively small number of toys *actually* appearing in any one image. Another example is a sound played by a piano at a given time t . While the sound can contain waveforms generated by pressing any of the 88 piano keys, there are only relatively few keys (typically much smaller than ten) that actually generated the sound. Sparse Coding algorithms model such data properties by providing a large number of hidden variables (potential data components) but assigning non-zero (or significantly different from zero) values only to a small subset of components (those actually appearing). Sparse coding algorithms are typically used for tasks such as denoising (Elad

and Aharon, 2006; Mairal et al., 2009b), inpainting (Mairal et al., 2009b,a; Titsias and Lázaro-Gredilla, 2011), classification (LeCun, NEC; Titsias and Lázaro-Gredilla, 2011; Raina et al., 2007, e.g. MNIST dataset <http://yann.lecun.com/exdb/mnist/>, the flowers dataset <http://www.robots.ox.ac.uk/~vgg/data/flowers/>), transfer learning (Raina et al., 2007), collaborative filtering (Titsias and Lázaro-Gredilla, 2011) and are important models for neuro-sensory processing (Olshausen and Field, 1997; Zylberberg et al., 2011; Bornschein et al., 2013; Sheikh et al., 2014, and many more). A variety of sparse coding models have been successfully scaled to high-dimensional latent spaces with the use of selection (Henniges et al., 2010; Bornschein et al., 2013; Sheikh et al., 2014) or selection combined with Gibbs sampling (Shelton et al., 2011, 2012, 2015) inference approaches. Latent variables were selected in these earlier works using selection functions that were individually defined for each model. In order to demonstrate our method of autonomously learned selection functions, we consider three of these sparse generative models, and perform inference in EM with our GP-select approach instead of a hand-crafted selection function. The models are relevant for different tasks such as classification (e.g., binary sparse coding), source separations and denoising (linear spike-and-slab sparse coding) or sparse encoding and extraction of interpretable image components (nonlinear sparse coding). Note that when it is obvious from context, we drop the notation referring to each data point n in order to make the equations more concise.

The models and their parameters are:

A. Binary sparse coding:

latents: $\mathbf{s} \sim \text{Bern}(\mathbf{s}|\pi) = \prod_{h=1}^H \pi^{s_h} (1 - \pi)^{1-s_h}$

observations: $\mathbf{y} \sim \mathcal{N}(\mathbf{y}; W\mathbf{s}, \sigma^2 I)$

parameters: $W = \left(\sum_{n=1}^N \mathbf{y}^{(n)} \langle \mathbf{s} \rangle_{q_n}^T \right) \left(\sum_{n=1}^N \langle \mathbf{s} \mathbf{s}^T \rangle_{q_n} \right)^{-1}$

$$\sigma^2 = \frac{1}{ND} \sum_n \left\langle \|\mathbf{y}^{(n)} - W\mathbf{s}\|^2 \right\rangle_{q_n}$$

$$\pi = \frac{1}{N} \sum_n |\langle \mathbf{s} \rangle_{q_n}|, \text{ where } |\mathbf{x}| = \frac{1}{H} \sum_h x_h$$

where $W \in \mathbb{R}^{D \times H}$ denotes the components / dictionary elements and π parameterizes the sparsity (see e.g. (Henniges et al., 2010)).

B. Spike-and-slab sparse coding:

latents: $\mathbf{s} = \mathbf{b} \odot \mathbf{z}$ where $\mathbf{b} \sim \text{Bern}(\mathbf{b}|\pi)$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mu, \Sigma_h)$

observations: $\mathbf{y} \sim \mathcal{N}(\mathbf{y}; W\mathbf{s}, \sigma^2 I)$

parameters: $W = \frac{\sum_{n=1}^N \mathbf{y}^{(n)} \langle \mathbf{s} \odot \mathbf{z} \rangle_n^T}{\sum_{n=1}^N \langle (\mathbf{s} \odot \mathbf{z})(\mathbf{s} \odot \mathbf{z})^T \rangle_n}$

$$\pi = \frac{1}{N} \sum_{n=1}^N \langle \mathbf{s} \rangle_n$$

$$\sigma^2 = \sum_{n=1}^N \left[\langle (\mathbf{s} \odot \mathbf{z})(\mathbf{s} \odot \mathbf{z})^T \rangle_n - \langle \mathbf{s} \mathbf{s}^T \rangle_n \odot \mu \mu^T \right] \odot \left(\sum_{n=1}^N \left[\langle \mathbf{s} \mathbf{s}^T \rangle_n \right] \right)^{-1}$$

$$\mu_{pr} = \frac{\sum_{n=1}^N \langle \mathbf{s} \odot \mathbf{z} \rangle_n}{\sum_{n=1}^N \langle \mathbf{s} \rangle_n}$$

$$\sigma_{pr}^2 = \frac{1}{N} \sum_{n=1}^N \left[\mathbf{y}^{(n)} (\mathbf{y}^{(n)})^T - W \left[\langle (\mathbf{s} \odot \mathbf{z}) \rangle_n \langle (\mathbf{s} \odot \mathbf{z}) \rangle_n^T \right] W^T \right]$$

where the point-wise multiplication of the two latent vectors, i.e., $(\mathbf{s} \odot \mathbf{z})_h = s_h z_h$ generates a ‘spike-and-slab’ distributed variable $(\mathbf{s} \odot \mathbf{z})$, that has either continuous values or exact zero entries (e.g. (Titsias and Lázaro-Gredilla, 2011; Goodfellow

et al., 2013; Sheikh et al., 2014)).

C. Nonlinear Spike-and-slab sparse coding:

$$\text{latents: } \mathbf{s} = \mathbf{b} \odot \mathbf{z} \quad \text{where} \quad \mathbf{b} \sim \text{Bern}(\mathbf{b}|\pi)$$

$$\text{and } \mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mu_{\text{pr}}, \sigma^2)$$

$$\text{observations: } \mathbf{y} \sim \mathcal{N}(\mathbf{y}; \max_h \{s_h \mathbf{W}_h\}, \sigma^2 I)$$

$$\text{parameters: } \hat{W}_{hd} = \frac{\langle s_h y_d \rangle^*}{\langle s_d^2 \rangle^*} \quad \hat{\pi} = \langle \mathbb{I}(s) \rangle$$

$$\hat{\sigma}^2 = \left\langle W_{dh} s_h - y_d^{(n)} \right\rangle^*$$

$$\hat{\mu}_{\text{pr}} = \langle s_h \rangle^* \quad \hat{\sigma}_{\text{pr}}^2 = \langle (s_h - \hat{\mu}_{\text{pr}})^2 \rangle^*$$

Where expectations $\langle . \rangle^*$ mean:

$$\langle f(s) \rangle^* = \sum_n \frac{\int_s p(\mathbf{s}|\mathbf{y}^{(n)}, \Theta) f(\mathbf{s}) \mathbb{I}(\mathbf{h} \text{ is max}) ds}{\int_s p(\mathbf{s}|\mathbf{y}^{(n)}, \Theta) \mathbb{I}(\mathbf{h} \text{ is max}) ds}$$

where \mathbb{I} is the indicator function denoting the domain to integrate over, namely where h is the maximum. Using $\langle f(s) \rangle^*$ allows for the condensed expression of the update equations shown above. The mean of the Gaussian for each $\mathbf{y}^{(n)}$ is centered at $\max_h \{s_h \mathbf{W}_h\}$, where \max_h is a nonlinearity that considers all H latent components and takes the h yielding the maximum value for $s_h \mathbf{W}_h$ (Lücke and Sahani, 2008; Shelton et al., 2012; Bornschein et al., 2013; Shelton et al., 2015), instead of centering the data at the linear combination of $\sum_h s_h \mathbf{W}_h = W \mathbf{s}$.

In the above models, inference with the truncated posterior of Equation (2) using hand-crafted selection functions $\mathcal{S}_h(\mathbf{y}^{(n)})$ to obtain the subset of states \mathcal{K}_n [of selected relevant variables $\mathbf{s}(\mathbf{y}^{(n)})$], shown in Equation (5), has yielded results as good

or more robust performance than exact inference (converging less frequently to local optima than exact inference; see earlier references for details). For models **A** and **C**, the hand-constructed function approximating $f(\mathbf{y}^{(n)})$, for substitution in Equation (5), was the cosine similarity between the weights \mathbf{W}_h (e.g. dictionary elements, components, etc.) associated with each latent variable s_h and each data point $\mathbf{y}^{(n)}$: $f(\mathbf{y}^{(n)}) = (\mathbf{W}_h^T / \|\mathbf{W}_h\|) \mathbf{y}^{(n)}$. For model **B**, the constructed affinity function was the data likelihood given a singleton state: $f(\mathbf{y}^{(n)}) = p(\mathbf{y}^{(n)} | \mathbf{s} = \mathbf{s}_h, \Theta)$, where \mathbf{s}_h represents a singleton state in which only the entry h is non-zero. The goal of these experiments is to demonstrate the performance of GP-select and the effects/benefits of using different selection functions. To do this, we consider artificial data generated according to each sparse coding model, and thus with known ground-truth parameters. As discussed above, we could also apply the sparse coding models using GP-select to other application domains listed, but that is not the focus of these experiments. We generate $N = 2,000$ data points consisting of $D = 5 \times 5 = 25$ observed dimensions and $H = 10$ latent components according to each of the models **A-C**: N images of randomly selected overlapping 'bars' with varying intensities for models **B** and **C**, and additive Gaussian noise parameterized by ground-truth $\sigma^2 = 2$ and we choose $H' = 5$, (e.g. following the spike-and-slab prior). On average, each data point contains 2 bars, i.e. ground-truth is $\pi H = 2$, and we choose $H' = 5$. With this choice, we can select sufficiently many latents for virtually all data points.

For each of the models considered, we run 10 repetitions of each of the following set of experiments: (1) selection using the respective hand-crafted selection function, (2) GP-select using a linear covariance kernel, (3) GP-select using an RBF covariance

kernel, and (4) GP-select using a kernel composed by adding the following kernels: RBF, linear, bias and white noise kernels, which we will term the *composition kernel*. As hyperparameters of kernels are learned, the composition kernel (4) can adapt itself to the data and only use the kernel components required. See Rasmussen and Williams (2005, Chapter 4, Section 4.2.4) for a discussion on kernel adaptation. Kernel parameters were model-selected via maximum marginal likelihood every 10 EM iterations. For models **A** and **B**, inference was performed exactly using the truncated posterior (2), but as exact inference is analytically intractable in model **C**, inference was performed by drawing Gibbs samples from the truncated space (Shelton et al., 2011, 2012, 2015).

We run all models until convergence.

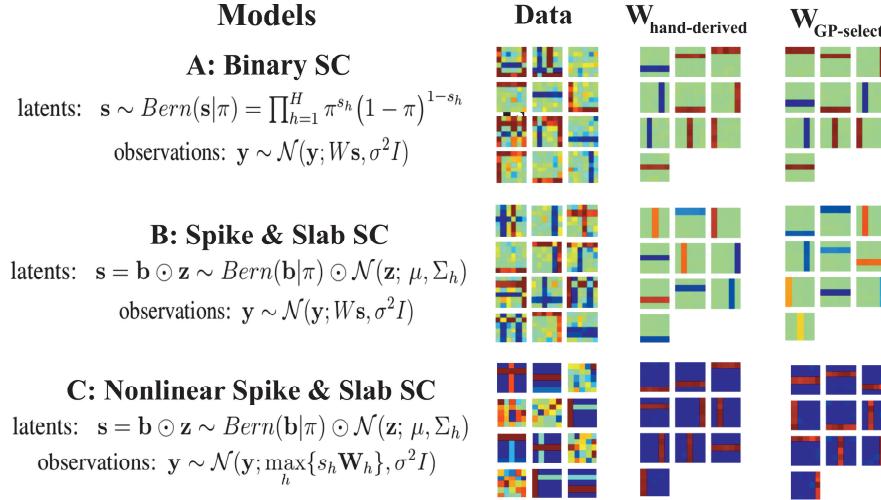


Figure 2: Sparse coding models results comparing GP-select with a successful hand-derived selection function. Results are shown on artificial ground-truth data with $H = 10$ latent variables and $H' = 5$ preselected variables for: **A** Binary sparse coding, **B** Spike-and-slab sparse coding, and **C** Nonlinear spike-and-slab sparse coding. First column: Example data points $\mathbf{y}^{(n)}$ generated by each of the models. Middle column: Converged dictionary elements \mathbf{W} learned by the hand-crafted selection functions. Third column: Converged dictionary elements \mathbf{W} learned by GP-select with $H' = 5$ using the kernel with best performance (matching that of inference with hand-crafted selection function). In all cases, the model using the GP-select function converged to the ground-truth solution, just as the hand-crafted selection functions did.

Results are shown in Figure 2. In all experiments, the GP-select approach was able to infer ground-truth parameters as well as the hand-crafted function. For models where the cosine similarity was used (in **A** and **C**), GP regression with a linear kernel quickly learned the ground-truth parameters, and hence fewer iterations of EM were necessary. In other words, even without providing GP-select explicit weights W as required for the hand-crafted function, its affinity function using GP regression (6) learned a similar enough function to quickly yield identical results. Furthermore, in the model with a less straight-forward hand-crafted function (in the spike-and-slab model of **B**), only GP regression with an RBF kernel was able to recover ground-truth parameters. In this case (model **B**), GP-select using an RBF kernel recovered the ground-truth 'bars' in 7 out of 10 repetitions, whereas the hand-crafted function recovered the bases in 8 instances. For the remaining models, GP-select converged to the ground-truth parameters with the same average frequency as the hand-crafted functions.

Finally, we have observed empirically that the composition kernel is flexible enough to subsume all other kernels: the variance of the irrelevant kernels dropped to zero in simulations. This suggests the composition kernel is a good choice for general use.

5.2 Gaussian mixture model

Next, we apply GP-select to a simple example, a Gaussian mixture model, where the flexibility of the approach can be easily and intuitively visualized. Furthermore, the GMMs flexibility allows us to explicitly visualize the effect of different selection functions. A demonstration and code for the GMM application is provided in (Dai, 2016).

The model of the data likelihood is

$$p(\mathbf{y}^{(n)} | \mu_{\mathbf{c}}, \sigma_{\mathbf{c}}, \pi) = \sum_{c=1}^C \mathcal{N}(\mathbf{y}^{(n)}; \mu_{\mathbf{c}}, \sigma_{\mathbf{c}}) \pi_c, \quad (7)$$

where C is the number of mixture components; the task is to assign each data point to its latent cluster.

The training data used for GP regression was $\mathcal{D} = \{(\mathbf{y}^{(n)}, \langle s_h^{(n)} \rangle_{q_n}) | n = 1, \dots, N\}$, where the targets were the expected cluster responsibilities (posterior probability distribution for each cluster) for all data points, $\langle s_h \rangle_q$, and we use one-hot encoding for cluster identity. With this, we apply our GP-select approach to this model, computing the selection function according to Equation (5) with affinity f defined by GP regression (6) and following the approximate EM approach as in the previous experiments. In these experiments we consider two scenarios for EM learning of the data likelihood in Equation (7): GP-select with an RBF covariance kernel and a linear covariance kernel. We do not include the composition kernel suggested (based on experiments) in Section 4.1, as the goal of the current experiments is to show the effects of using the 'wrong' kernel. These effects would further support the use of the flexible composition kernel in general, as it can subsume both kernels considered in the current experiments (RBF and linear).

To easily visualize the output, we generate 2-dimensional observed data $(\mathbf{y}^{(n)} \in \mathbb{R}^{D=2})$ from $C = 3$ clusters – first with randomly assigned cluster means, and second such that the means of the clusters lie roughly on a line. In the GP-select experiments, we select $C' = 2$ clusters from the full set, and run 40 EM iterations for both kernel choices (linear and RBF). Note that for mixture models, the notation of C' selected clusters of the C set is analogous to the H' selected latent variables from the H full set,

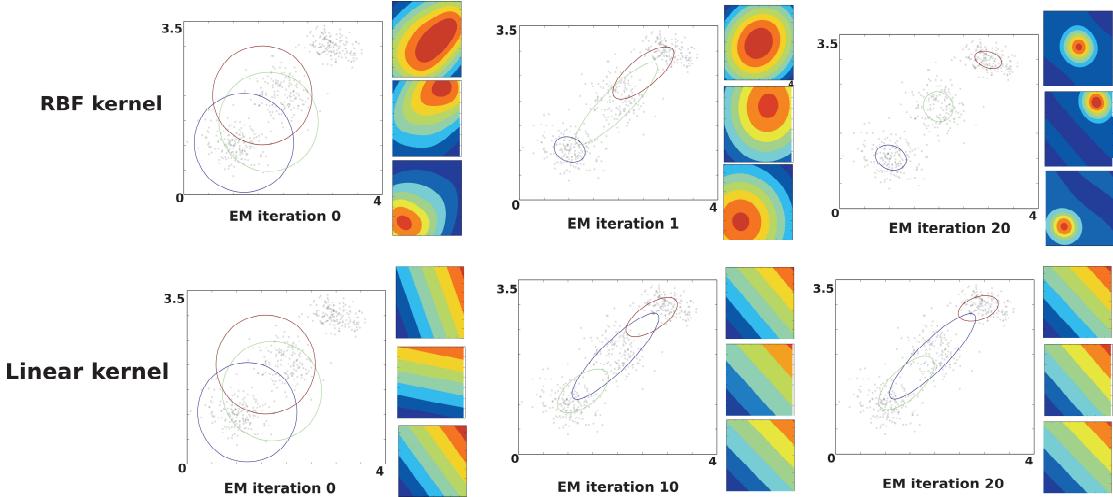


Figure 3: Gaussian mixture model results using GP-select (selection of $C' = 2$ in a $C = 3$ class scenario) for inference. Progress of the inference is shown using (row one) an RBF covariance kernel in the regression, and (row two) a linear covariance kernel. For each iteration shown, we see (1) the observed data and their inferred cluster assignments and (2) the C corresponding GP regression functions learned/used for GP-select in that iteration. Different iterations are pictured due to different convergence rates. As shown, inference with GP-select using a linear kernel is unable to assign the data points to the appropriate clusters, whereas GP-select with an RBF kernel succeeds.

as described in the non-mixture model setting, and the GP-select algorithm proceeds unchanged. We randomly initialize the variance of the clusters σ_c and initialize the cluster means μ_c at randomly selected data points. Results are shown in Figure 3.

on these data, the linear GP regression prediction cannot correctly assign the data to their clusters (as seen in Figure 3B), but the nonlinear approach successfully and easily finds the ground-truth clusters (Figure 3A). Furthermore, even when both approaches were initialized in the optimal solution, the cluster assignments from GP-select with a linear kernel quickly wandered away from the optimal solution and were identical to random initialization, converging to the same result shown in iteration 20 of Figure 3B). The RBF kernel cluster assignments remained at the optimal solution even with number of selected clusters set to $C' = 1$.

These experiments demonstrate that the selection function needs to be flexible even for very simple models, and that nonlinear selection functions are an essential tool even in such apparently straightforward cases.

5.3 Translation Invariant Occlusive models

Now that we have verified that GP-select can be applied to various generative graphical models and converge to ground-truth parameters, we consider a more challenging model that addresses a problem in computer vision: *translations of objects in a scene*.

Model. Translation invariant models address the problem that, e.g., visual objects can appear in any location of an image. Probabilistic models for translation invariance are particularly appealing as they allow to separately infer object positions and object type, making them very interpretable and powerful tools for image processing.

Translation invariant models are particularly difficult to optimize, however, because they must consider a massive latent variable space: evaluating multiple objects and locations in a scene leads a *latent space complexity of the number of locations exponentiated by the number of objects*. Inference in such a massive latent space heavily relies on the idea of variable selection to reduce the number of candidate objects and locations. In particular, hand-engineered selection functions that consider *translational invariance* have been successfully applied to this type of model (Dai and Lücke, 2012b, 2014; Dai et al., 2013). The selection function used so far for reducing latent space complexity in this model has been constructed as follows. First, the candidate locations of all the objects in the model are predicted. Then a subset of candidate objects that might appear in the image are selected according to those predicted locations. Next,

the subset of states \mathcal{K}_n is constructed according to the combinations of the possible locations and numbers of candidate objects. The posterior distribution is then computed following Equation (2).

This selection system is very costly: the selection function has parameters which need to be hand-tuned, e.g., the number of representative features, and it needs to scan through the entire image, considering all possible locations, which becomes computationally demanding for large-scale experiments. To maximally exploit the capabilities of the GP-selection function, we directly use the GP regression model to *predict the possible locations* of a component *without introducing any knowledge of translation invariance* into the selection function. In this work, a GP regression model is fitted from the input image to marginal posterior probabilities of individual components appearing at all possible locations. Therefore, the input to the GP-selection function is the image to be inferred and the output is a score for each possible location of each component in the model. For example, when learning 10 components in a $D = 30 \times 30$ pixel image patch, the output dimensionality of GP-select is 9000. This task is computationally feasible, since GP models scale linearly with output dimensionality. The inference of components' locations with GP-select is significantly faster than the selection function in the original work, as it avoids explicitly scanning through the image.

Although there are additional computations necessary for an automatic selection function like GP-select, for instance due to the adjustment of its parameters, there are many options to reduce computational costs. First, we may approximate the full $N \times N$ Gram matrix by an incomplete Cholesky approximation (Fine and Scheinberg, 2001) resulting in a cost of $O(N \times Q)$, where $Q \ll N$ is the rank of the Cholesky approxi-



Figure 4: COIL Dataset (Nene et al., 1996): A handful of data points used in experiments with the Translation Invariant Occlusive (InvECA) model, showing the occluding objects to be learned.

mation. Second, we may reduce the update frequency of the kernel hyperparameters to be computed only every T^* EM iterations, where a $T^* > 1$ represents a corresponding computation reduction. The combination of the Cholesky approximation plus infrequent updates will have the following benefits: a factor of five speedup for infrequent updates, and a factor of $(N - Q)^2$ speedup from incomplete Cholesky, where Q is the rank of the Cholesky approximation and N is the number of original data points.

COIL Dataset. We apply our GP-selection function to the *Invariant Exclusive Component Analysis (InvECA) model* (Dai and Lücke, 2012b; Dai et al., 2013). For our experiments, we consider an image dataset used in previous work: data were generated using objects from the COIL-100 image dataset (Nene et al., 1996), taking 16 different objects, downscaled to $D = 10 \times 10$ pixels and segmented out from the black background. A given image was generated by randomly selecting a subset of the 16 objects, where each object has a probability of 0.2 of appearing. The appearing objects were placed at random positions on a 30×30 black image. When the objects overlap, they occlude each other with a different random depth order for each image. In total, $N = 2000$ images were generated in the dataset (examples shown in Figure 4). The task of the InvECA model is to discover the visual components (i.e. the images of 16

objects) from the image set without any label information. We compare the visual components learned by using *four different selection functions in the InvECA model*: the hand-crafted selection function used in the original work by Dai and Lücke (2012b), GP-select updated every iteration, GP-select updated every $T^* = 5$ iterations, and GP-select with incomplete Cholesky decomposition updated every iteration, or $T^* = 1$ (in this manner we isolate the improvements due to Cholesky from those due to infrequent updates). In these experiments, the parameters of GP-select are optimized at the end of each T^* EM iteration(s), using a maximum of 20 gradient updates. The number of objects to be learned is $H = 20$ and the algorithm pre-selects $H' = 5$ objects for each data point. The kernel used was the composition kernel, as suggested in Section 4.1, although after fitting the hyperparameters only the RBF kernel remained with large variance (i.e. a linear kernel alone would not have produced good variable selection, thus the flexible composition kernel was further shown to be a good choice).

Results. All four versions of the InvECA model using each of the selection functions considered successfully recover each individual objects in our modified COIL image set. The learned object representations with GP-select are shown in Figure 5. Four additional components developed into representations, however these all had very low mask values, allowing them to easily be distinguished from other true components.

Next, we compare the accuracy of the four selection functions. For this, we collected the object locations (pixels) indicated by each selection function after all EM iterations, applied the selection functions (for the GP selection functions, this was using the final function learned after all EM iterations) to the entire image dataset again, then compared these results with the ground-truth location of all of the objects in the dataset.

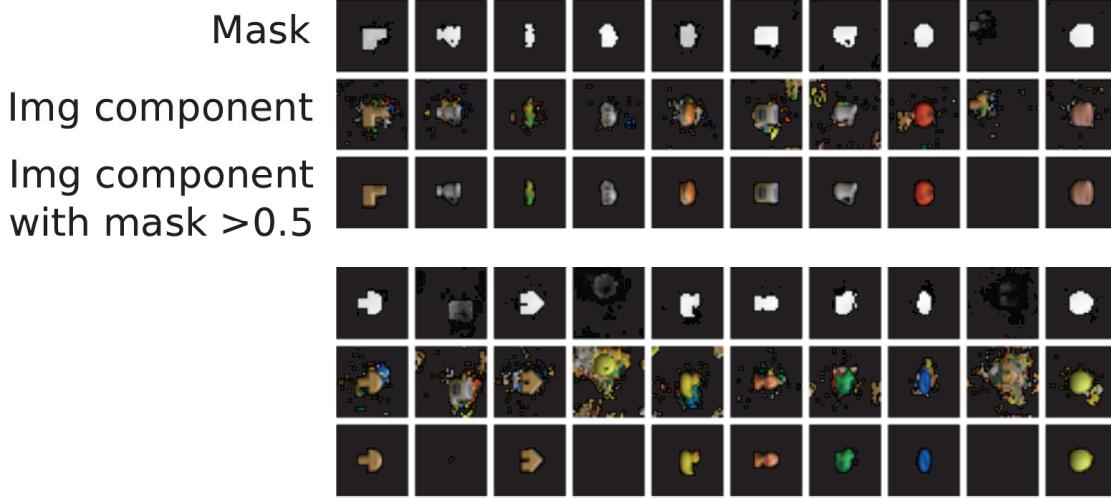


Figure 5: Image components and their masks learned by GP-select with the Translation Invariant model. GP-select learned all objects in the dataset. The first row shows the mask of each component, the second row shows the learned image components, and the third row shows only the area of the learned components that had a mask > 0.5 . For the second three-row block of images, the same titles of the first three-row block hold.

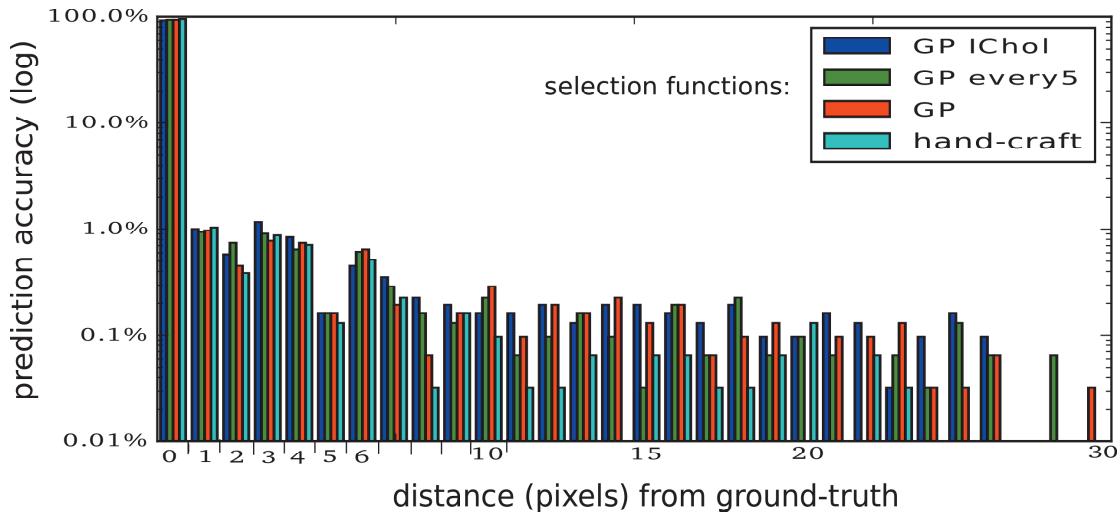


Figure 6: Prediction accuracy of the four selection functions in the InvECA model. Functions depicted in the figures: GP-select with no modifications (GP, red), the incomplete Cholesky decomposition (GP IChol, blue), with updated kernel hyperparameters every 5 EM iterations (GP every5, green), and with hand-crafted selection (hand-craft, cyan). Shown: the log-scale histogram of the prediction accuracy for the four selection functions, measured by the distance each function’s predicted object location was to the ground-truth object location. All bars of the selection functions show very similar accuracy for the various distances.

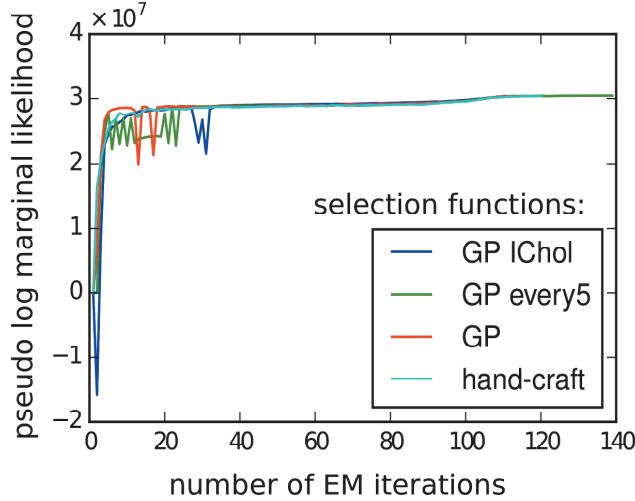


Figure 7: Baseline comparison of the four selection functions in the InvECA model. Functions depicted in the figures are identical to those in Figure 6. Shown: the convergence of the pseudo log marginal likelihood [of the model parameters learned at each EM iteration] for the four selection functions over all EM iterations. After about 40 EM iterations, all selection function versions of the algorithm converge to the same likelihood solution. Simultaneously, the GP-select approaches exhibit no loss of accuracy compared to the hand-crafted function, and 'GP IChol' represents a factor of 100 speedup vs. 'GP', and 'GP every5' represents a factor of 5 speedup.

The accuracy of the predicted locations was then computed by comparing the distance of all ground-truth object location to the location of the top candidate locations from each selection function. See Figure 6 for a histogram of these distances and the corresponding accuracy for all selection functions. Note that the percentages in the histogram are plotted in log scale. Also, as a baseline verification, we computed and compared the pseudo log likelihood (Dai et al., 2013) of the original selection function to the three GP-select based ones. The pseudo log likelihood for all selection functions is shown in Figure 7. Figures 6-7 show that all four selection functions can very accurately predict the locations of all the objects in the dataset – the GP-select selection functions yields no loss in inference performance in comparison to the original hand-engineered selection function. Even those using speed-considerate approximations (incomplete Cholesky

decomposition of the kernel matrix (GP IChol) and updating kernel hyperparameters only every 5 EM iterations (GP every5)) have indistinguishable prediction accuracy on the task.

An analysis of the benefits indicate that, as GP-select avoids explicitly scanning through the image, the time to infer the location of an object is significantly reduced compared to the hand-crafted function. GP-select requires 22.1 seconds on a single CPU core to infer the locations of objects across the whole image set, while the hand-crafted function requires 1830.9 seconds. In the original work, the selection function was implemented with GPU acceleration and parallelization. Although we must compute the kernel hyperparameters for GP-select, it is important to note that the hyperparameters need not be fit perfectly each iteration – for the purposes of our approach, a decent approximation suffices for excellent variable selection. In this experiment, updating the parameters of GP-select with 10 gradient steps took about 390 seconds for the full-rank kernel matrix. When we compute the incomplete Cholesky decomposition while inverting the covariance matrix, compute time was reduced to 194 seconds (corresponding to the $(N - Q)^2$ speedup, where Q is the rank of the Cholesky approximation), with minimal loss in accuracy. Furthermore, when updating the GP-select hyperparameters only every 5 iterations, average compute time was reduced by another one fifth, again without loss in accuracy.

6 Discussion

We have proposed a means of achieving fast EM inference in Bayesian generative models, by learning an approximate selection function to determine relevant latent variables for each observed variable. The process of learning the relevance functions is interleaved with the EM steps, and these functions are used in obtaining an approximate posterior distribution in the subsequent EM iteration. The functions themselves are learned via Gaussian process regression, and do not require domain-specific engineering, unlike previous selection functions. In experiments on mixtures and sparse coding models with interpretable output, the learned selection functions behaved in accordance with our expectations for the posterior distribution over the latents.

The significant benefit we show empirically is that by learning the selection function in a general and flexible nonparametric way, we can avoid using potentially expensive hand-engineered selection functions. Cost reduction is both in terms of required expertise in the problem domain, and computation time in identifying the relevant latent variables. Inference using our approach required 22.1 seconds on a single CPU core, versus 1830.9 seconds with the original hand-crafted function for the complex hierarchical model of (Dai et al., 2013).

A major area where further performance gains might be expected is in improving computational performance, since we expect the greatest advantages of GP-select to occur for complex models at large scale. For instance, kernel ridge regression may be parallelized (Zhang et al., 2014), or the problem may be solved in the primal via random Fourier features (Le et al., 2013). Furthermore, there are many recent developments regarding the scaling up of GP inference to large-scale problems, e.g., sparse

GP approximation (Lawrence et al., 2002), stochastic variational inference (Hensman et al., 2013, 2012), using parallelization techniques and GPU acceleration (Dai et al., 2014), or in combination with stochastic gradient descent (Bottou and Bousquet, 2008). For instance, for very large datasets where the main model is typically trained with mini-batch learning, stochastic variational inference can be used for GP fitting as in (Hensman et al., 2013) and the kernel parameters can be efficiently updated each (or only every T * few) iteration with respect to a mini-batch.

Acknowledgments

We acknowledge funding by the German Research Foundation (DFG) under grants LU 1196/4-2 (JS), LU 1196/5-1 (JL), by the Cluster of Excellence EXC 1077/1 "Hearing4all" (JL), and by the RADIANT and WYSIWYD (EU FP7-ICT) projects (ZD).

References

- Bornschein, J., Henniges, M., and Lücke, J. (2013). Are V1 simple cells optimized for visual occlusions? A comparative study. *PLoS Computational Biology*, 9(6):e1003062.
- Bottou, L. and Bousquet, O. (2008). The tradeoffs of large scale learning. In Platt, J. C., Koller, D., Singer, Y., and Roweis, S. T., editors, *Advances in Neural Information Processing Systems 20 (NIPS)*, Vancouver, British Columbia, Canada, pages 161–168. MIT Press.
- Currin, C., Mitchell, T., Morris, M., and Ylvisaker, D. (1991). Bayesian prediction

of deterministic functions, with applications to the design and analysis of computer experiments. *J. American Statistical Association*, 86(416):953–963.

Dai, Z. (2016). GP-select Demo on Gaussian mixture models.

https://github.com/fatflake/GP-select-Code/GMM_demo.ipynb.

Dai, Z., Damianou, A., Hensman, J., and Lawrence, N. (2014). Gaussian process models with parallelization and gpu acceleration. *NIPS 2014, Workshop on Modern non-parametrics: automating the learning pipeline*. In *NIPS Workshop on Modern non-parametrics: automating the learning pipeline*.

Dai, Z., Exarchakis, G., and Lücke, J. (2013). What are the invariant occlusive components of image patches? a probabilistic generative approach. In Burges, C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K., editors, *Advances in Neural Information Processing Systems*, pages 243–251.

Dai, Z. and Lücke, J. (2012a). Autonomous cleaning of corrupted scanned documents – a generative modeling approach. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3338–3345.

Dai, Z. and Lücke, J. (2012b). Unsupervised learning of translation invariant occlusive components. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2400–2407.

Dai, Z. and Lücke, J. (2014). Autonomous document cleaning – a generative approach to reconstruct strongly corrupted scanned texts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(10):1950–1962.

- Dayan, P., Hinton, G. E., Neal, R. M., and Zemel, R. S. (1995). The helmholtz machine. *Neural Computation*, 7:889–904.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society B*, 39:1–38.
- Elad, M. and Aharon, M. (2006). Image denoising via sparse and redundant representations over learned dictionaries. *Trans. Img. Proc.*, 15(12):3736–3745.
- Exarchakis, G., Henniges, M., Eggert, J., and Lücke, J. (2012). Ternary sparse coding. In *LVA/ICA*, Lecture Notes in Computer Science, pages 204–212. Springer.
- Fine, S. and Scheinberg, K. (2001). Efficient SVM training using low-rank kernel representations. *Journal of Machine Learning Research*, 2:243–264.
- Goodfellow, I. J., Courville, A., and Bengio, Y. (2013). Scaling up spike-and-slab models for unsupervised feature learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1902–1914.
- Gutmann, M. U. and Corander, J. (2015). Bayesian Optimization for Likelihood-Free Inference of Simulator-Based Statistical Models. Technical report, University of Helsinki. <http://arxiv.org/abs/1501.03291>.
- Henniges, M., Puertas, G., Bornschein, J., Eggert, J., and Lücke, J. (2010). Binary Sparse Coding. In *Proceedings LVA/ICA*, LNCS 6365, pages 450–57. Springer.
- Henniges, M., Turner, R. E., Sahani, M., Eggert, J., and Lücke, J. (2014). Efficient occlusive components analysis. *JMLR*, 15:2689–2722.

Hensman, J., Fusi, N., and Lawrence, N. (2013). Gaussian processes for big data. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence, UAI 2013, Bellevue, WA, USA, August 11-15, 2013*.

Hensman, J., Rattray, M., and Lawrence, N. D. (2012). Fast Variational Inference in the Conjugate Exponential Family. In Bartlett, P. L., Pereira, F. C. N., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 25 (NIPS), Lake Tahoe, Nevada, United States*, pages 2897–2905.

Hinton, G. E., Dayan, P., Frey, B. J., and Neal, R. M. (1995). The ‘wake-sleep’ algorithm for unsupervised neural networks. *Science*, 268:1158 – 1161.

Kingma, D. P. and Welling, M. (2014). Efficient gradient-based inference through transformations between bayes nets and neural nets. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pages 1782–1790.

Körner, E., Gewaltig, M. O., Körner, U., Richter, A., and Rodemann, T. (1999). A model of computation in neocortical architecture. *Neural Networks*, 12:989 – 1005.

Lawrence, N., Seeger, M., and Herbrich, R. (2002). Fast sparse gaussian process methods: The informative vector machine. In *Advances in Neural Information Processing Systems 15 (NIPS), Vancouver, British Columbia, Canada*, pages 609–616. MIT Press.

Le, Q., Sarlos, T., and Smola, A. J. (2013). Fastfood — computing hilbert space ex-

pansions in loglinear time. In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, pages 244–252.

LeCun, Y. (NEC). MNIST database of handwritten digits.

<http://yann.lecun.com/exdb/mnist/>.

Lücke, J. and Eggert, J. (2010). Expectation truncation and the benefits of preselection in training generative models. *Journal of Machine Learning Research*, 11:2855–900.

Lücke, J. and Sahani, M. (2008). Maximal causes for non-linear component extraction. *Journal of Machine Learning Research*, 9:1227–67.

Mairal, J., Bach, F., Ponce, J., and Sapiro, G. (2009a). Online dictionary learning for sparse coding. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009*, volume 25, pages 689–696.

Mairal, J., Bach, F., Ponce, J., Sapiro, G., and Zisserman, A. (2009b). Non-local sparse models for image restoration. In *IEEE 12th International Conference on Computer Vision, ICCV 2009, Kyoto, Japan, September 27 - October 4, 2009*, volume 25, pages 2272–2279.

Meeds, E. and Welling, M. (2014). GPS-ABC: gaussian process surrogate approximate bayesian computation. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence (UAI), Quebec City, Quebec, Canada, July 23-27, 2014*, pages 593–602.

Mnih, A. and Gregor, K. (2014). Neural variational inference and learning in belief

networks. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pages 1791–1799.

Neal, R. and Hinton, G. (1998). A view of the EM algorithm that justifies incremental, sparse, and other variants. In Jordan, M. I., editor, *Learning in Graphical Models*. Kluwer.

Nene, S. A., Nayar, S. K., and Murase, H. (1996). Columbia object image library (coil-100). Technical report, CUCS-006-96.

Olshausen, B. and Field, D. (1997). Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, 37(23):3311–3325.

Raina, R., Battle, A., Lee, H., Packer, B., and Ng, A. Y. (2007). Self-taught learning: Transfer learning from unlabeled data. In *Proceedings of the 24th International Conference on Machine Learning, ICML '07*, pages 759–766, New York, NY, USA. ACM.

Rasmussen, C. E. and Williams, C. K. I. (2005). *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press.

Rupp, M., Tkatchenko, A., Müller, K.-R., and von Lilienfeld, O. A. (2012). Fast and accurate modeling of molecular atomization energies with machine learning. *Phys. Rev. Lett.*, 108:058301.

Sacks, J., Welch, W. J., Mitchell, T. J., and Wynn, H. P. (1989). Design and analysis of computer experiments. *Statist. Sci.*, 4(4):433–435.

Schwaighofer, A., Tresp, V., and Yu, K. (2004). Learning gaussian process kernels via hierarchical bayes. In *Advances in Neural Information Processing Systems 17 (NIPS), Vancouver, British Columbia, Canada*, pages 1209–1216. MIT Press.

Sheikh, A.-S., Shelton, J., and Lücke, J. (2014). A truncated EM approach for spike-and-slab sparse coding. *Journal of Machine Learning Research*, 15:2653–2687.

Shelton, J., Bornschein, J., Sheikh, A.-S., Berkes, P., and Lücke, J. (2011). Select and sample - A model of efficient neural inference and learning. In *Advances in Neural Information Processing Systems 24 (NIPS), Granada, Spain.*, pages 2618–2626.

Shelton, J., Sterne, P., Bornschein, J., Sheikh, A.-S., and Lücke, J. (2012). Why MCA? nonlinear sparse coding with spike-and-slab prior for neurally plausible image encoding. In *Advances in Neural Information Processing Systems 25 (NIPS), Lake Tahoe, Nevada, United States*, pages 2285–2293.

Shelton, J. A., Sheikh, A.-S., Bornschein, J., Sterne, P., and Lücke, J. (2015). Nonlinear spike-and-slab sparse coding for interpretable image encoding. *PLoS ONE*, 10(5):1–25.

Snoek, J., Rippel, O., Swersky, K., Kiros, R., Satish, N., Narayanan Sundaram, M., Ali, M., Patwary, P., and Adams, R. (2015). Scalable bayesian optimization using deep neural networks. Technical report, Harvard University.
<http://arxiv.org/abs/1502.05700>.

Titsias, M. and Lázaro-Gredilla, M. (2011). Spike and slab variational inference for

multi-task and multiple kernel learning. In *Advances in Neural Information Processing Systems 24 (NIPS), Granada, Spain*, pages 2510–2518.

Wilkinson, R. D. (2014). Accelerating ABC methods using gaussian processes. Technical report, University of Sheffield. <http://arxiv.org/abs/1401.1436>.

Yuille, A. and Kersten, D. (2006). Vision as Bayesian inference: analysis by synthesis? *Trends in Cognitive Sciences*, 10(7):301–308.

Zhang, Y., Duchi, J. C., and Wainwright, M. J. (2014). Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. Technical report, University of California, Berkeley. <http://arxiv.org/abs/1305.5029>.

Zylberberg, J., Murphy, J., and Deweese, M. (2011). A Sparse Coding Model with Synaptically Local Plasticity and Spiking Neurons Can Account for the Diverse Shapes of V1 Simple Cell Receptive Fields. *PLoS Computational Biology*, 7(10):e1002250.