

# Chapter 1

## Dropout in Gaussian processes

### 1.1 Dropout in Gaussian processes

Dropout is a method for regularizing neural networks (Hinton et al., 2012; Srivastava, 2013). Training with dropout entails randomly and independently “dropping” (setting to zero) some proportion  $p$  of features or inputs, in order to improve the robustness of the resulting network by reducing co-dependence between neurons. To maintain similar overall activation levels, weights are multiplied by  $1/p$  at test time. Alternatively, feature activations are multiplied by  $1/p$  during training. At test time, the set of models defined by all possible ways of dropping-out neurons is averaged over, usually in an approximate way.

Baldi and Sadowski (2013) and Wang and Manning (2013) analyzed dropout in terms of the effective prior induced by this procedure in several models, such as linear and logistic regression. In this section, we examine the priors on functions that result from performing dropout in the one-hidden-layer neural network implicitly defined by a GP (equation (??)).

#### 1.1.1 Dropout on feature activations

First, we examine the prior that results from randomly dropping features from  $\mathbf{h}(\mathbf{x})$  with probability  $p$ . If these features have a weight distribution with finite moments  $\mathbb{E}[\alpha_i] = \mu$ ,  $\mathbb{V}[\alpha_i] = \sigma^2$ , then the distribution of weights after each one has been dropped out with probability  $p$  is:

$$r_i \stackrel{\text{iid}}{\sim} \text{Ber}(p) \quad \mathbb{E}[r_i \alpha_i] = p\mu, \mathbb{V}[r_i \alpha_i] = p^2 \sigma^2. \quad (1.1)$$

However, after we increase the remaining activations to maintain the same expected activation by multiplying them by  $1/p$ , the resulting moments are again:

$$\mathbb{E} \left[ \frac{1}{p} r_i \alpha_i \right] = \frac{p}{p} \mu = \mu, \quad \mathbb{V} \left[ \frac{1}{p} r_i \alpha_i \right] = \frac{p^2}{p^2} \sigma^2 = \sigma^2. \quad (1.2)$$

Thus, dropping out features of an infinitely-wide MLP does not change the model at all, since no individual feature can have more than infinitesimal contribution to the network output.

### 1.1.2 Dropping out inputs

In a GP with kernel  $k(\mathbf{x}, \mathbf{x}') = \prod_{d=1}^D k_d(\mathbf{x}_d, \mathbf{x}'_d)$ , exact averaging over all possible ways of dropping out inputs with probability  $1/2$  results in a mixture of GPs, each depending on only a subset of the inputs:

$$p(f(\mathbf{x})) = \frac{1}{2^D} \sum_{\mathbf{r} \in \{0,1\}^D} \text{GP} \left( 0, \prod_{d=1}^D k_d(\mathbf{x}_d, \mathbf{x}'_d)^{r_d} \right) \quad (1.3)$$

We present two results ways to gain intuition about this model.

First, if the kernel on each dimension has the form  $k_d(\mathbf{x}_d, \mathbf{x}'_d) = g\left(\frac{\mathbf{x}_d - \mathbf{x}'_d}{w_d}\right)$ , as does the SE kernel (??), then any input dimension can be dropped out by setting its lengthscale  $w_d$  to  $\infty$ . Thus, dropout on the inputs of a GP corresponds to a spike-and-slab prior on lengthscales, with each dimension independently having  $w_d = \infty$  with probability  $1/2$ .

Another way to understand the resulting prior is to note that the dropout mixture (1.3) has the same covariance as

$$f(\mathbf{x}) \sim \text{GP} \left( 0, \frac{1}{2^{-2D}} \sum_{\mathbf{r} \in \{0,1\}^D} \prod_{d=1}^D k_d(\mathbf{x}_d, \mathbf{x}'_d)^{r_d} \right) \quad (1.4)$$

A GP whose covariance is a sum of kernels corresponds to a sum of functions, each distributed according to a GP. Therefore, (1.4) describes a sum of  $2^D$  functions, each depending on a different subset of the inputs.

In this chapter, we introduce a Gaussian process model of functions which are *additive*. An additive function is one which decomposes into a sum of low-dimensional functions, each depending on only a subset of the input variables. Additive GPs generalize both Generalized Additive Models, and the standard GP models which use squared-exponential kernels. Hyperparameter learning in this model can be seen as Bayesian

Hierarchical Kernel Learning (HKL). We introduce an expressive but tractable parameterization of the kernel function, which allows efficient evaluation of all input interaction terms, whose number is exponential in the input dimension. The additional structure discoverable by this model results in state-of-the-art predictive power in regression tasks, as well as increased interpretability.

## 1.2 Introduction

Most statistical regression models in use today are of the form:  $g(y) = f(x_1) + f(x_2) + \dots + f(x_D)$ . Popular examples include logistic regression, linear regression, and Generalized Linear Models (Nelder and Wedderburn, 1972). This family of functions, known as Generalized Additive Models (GAM) (Hastie and Tibshirani, 1990), are typically easy to fit and interpret. Some extensions of this family, such as smoothing-splines ANOVA (Wahba, 1990), add terms depending on more than one variable. However, such models generally become intractable and difficult to fit as the number of terms increases.

At the other end of the spectrum are kernel-based models, which typically allow the response to depend on all input variables simultaneously. These have the form:  $y = f(x_1, x_2, \dots, x_D)$ . A popular example would be a Gaussian process model using a squared-exponential (or Gaussian) kernel. We denote this model as SE-GP. This model is much more flexible than the GAM, but its flexibility makes it difficult to generalize to new combinations of input variables.

In this paper, we introduce a Gaussian process model that generalizes both GAMs and the SE-GP. This is achieved through a kernel which allow additive interactions of all orders, ranging from first order interactions (as in a GAM) all the way to  $D$ th-order interactions (as in a SE-GP). Although this kernel amounts to a sum over an exponential number of terms, we show how to compute this kernel efficiently, and introduce a parameterization which limits the number of hyperparameters to  $O(D)$ . A Gaussian process with this kernel function (an additive GP) constitutes a powerful model that allows one to automatically determine which orders of interaction are important. We show that this model can significantly improve modeling efficacy, and has major advantages for model interpretability.

We can see that a GP with a first-order additive kernel is an example of a GAM: Each function drawn from this model is a sum of orthogonal one-dimensional functions. Compared to functions drawn from the higher-order GP, draws from the first-order GP have more long-range structure.

We can expect many natural functions to depend only on sums of low-order interactions. For example, the price of a house or car will presumably be well approximated by a sum of prices of individual features, such as a sun-roof. Other parts of the price may depend jointly on a small set of features, such as the size and building materials of a house. Capturing these regularities will mean that a model can confidently extrapolate to unseen combinations of features.

### 1.3 Additive Kernels

We now give a precise definition of additive kernels. We first assign each dimension  $i \in \{1 \dots D\}$  a one-dimensional *base kernel*  $k_i(x_i, x'_i)$ . We then define the first order, second order and  $n$ th order additive kernel as:

$$k_{add_1}(\mathbf{x}, \mathbf{x}') = \sigma_1^2 \sum_{i=1}^D k_i(x_i, x'_i) \quad (1.5)$$

$$k_{add_2}(\mathbf{x}, \mathbf{x}') = \sigma_2^2 \sum_{i=1}^D \sum_{j=i+1}^D k_i(x_i, x'_i) k_j(x_j, x'_j) \quad (1.6)$$

$$k_{add_n}(\mathbf{x}, \mathbf{x}') = \sigma_n^2 \sum_{1 \leq i_1 < i_2 < \dots < i_n \leq D} \left[ \prod_{d=1}^n k_{i_d}(x_{i_d}, x'_{i_d}) \right] \quad (1.7)$$

$$k_{add_D}(\mathbf{x}, \mathbf{x}') = \sigma_n^2 \sum_{1 \leq i_1 < i_2 < \dots < i_D \leq D} \left[ \prod_{d=1}^D k_{i_d}(x_{i_d}, x'_{i_d}) \right] = \sigma_D^2 \prod_{d=1}^D k_d(x_d, x'_d) \quad (1.8)$$

where  $D$  is the dimension of our input space, and  $\sigma_n^2$  is the variance assigned to all  $n$ th order interactions. The  $n$ th covariance function is a sum of  $\binom{D}{n}$  terms. In particular, the  $D$ th order additive covariance function has  $\binom{D}{D} = 1$  term, a product of each dimension's covariance function:

$$k_{add_D}(\mathbf{x}, \mathbf{x}') = \sigma_D^2 \prod_{d=1}^D k_d(x_d, x'_d) \quad (1.9)$$

In the case where each base kernel is a one-dimensional squared-exponential kernel, the  $D$ th-order term corresponds to the multivariate squared-exponential kernel:

$$k_{add_D}(\mathbf{x}, \mathbf{x}') = \sigma_D^2 \prod_{d=1}^D k_d(x_d, x'_d) = \sigma_D^2 \prod_{d=1}^D \exp\left(-\frac{(x_d - x'_d)^2}{2l_d^2}\right) = \sigma_D^2 \exp\left(-\sum_{d=1}^D \frac{(x_d - x'_d)^2}{2l_d^2}\right) \quad (1.10)$$

also commonly known as the Gaussian kernel. The full additive kernel is a sum of the additive kernels of all orders.

### 1.3.1 Parameterization

The only design choice necessary in specifying an additive kernel is the selection of a one-dimensional base kernel for each input dimension. Any parameters (such as length-scales) of the base kernels can be learned as usual by maximizing the marginal likelihood of the training data.

In addition to the hyperparameters of each dimension-wise kernel, additive kernels are equipped with a set of  $D$  hyperparameters  $\sigma_1^2 \dots \sigma_D^2$  controlling how much variance we assign to each order of interaction. These “order variance” hyperparameters have a useful interpretation: the  $d$ th order variance hyperparameter controls how much of the target function’s variance comes from interactions of the  $d$ th order. Table 1.1 shows examples of the variance contributed by the different orders of interaction, learned on real datasets.

Table 1.1 Relative variance contribution of each order in the additive model, on different datasets. Here, the maximum order of interaction is set to 10, or smaller if the input dimension less than 10. Values are normalized to sum to 100.

Dataset	Order of interaction									
	1st	2nd	3rd	4th	5th	6th	7th	8th	9th	10th
pima	0.1	0.1	0.1	0.3	1.5	<b>96.4</b>	1.4	0.0		
liver	0.0	0.2	<b>99.7</b>	0.1	0.0	0.0				
heart	<b>77.6</b>	0.0	0.0	0.0	0.1	0.1	0.1	0.1	0.1	22.0
concrete	<b>70.6</b>	13.3	13.8	2.3	0.0	0.0	0.0	0.0		
pumadyn-8nh	0.0	0.1	0.1	0.1	0.1	0.1	0.1	<b>99.5</b>		
servo	<b>58.7</b>	27.4	0.0	13.9						
housing	0.1	0.6	<b>80.6</b>	1.4	1.8	0.8	0.7	0.8	0.6	12.7

On different datasets, the dominant order of interaction estimated by the additive model varies widely. An additive GP with all of its variance coming from the 1st order is equivalent to a GAM; an additive GP with all its variance coming from the  $D$ th order is equivalent to a SE-GP.

Because the hyperparameters can specify which degrees of interaction are important, the additive GP is an extremely general model. If the function we are modeling is

decomposable into a sum of low-dimensional functions, our model can discover this fact and exploit it (see figure ??). If this is not the case, the hyperparameters can specify a suitably flexible model.

### 1.3.2 Interpretability

As noted by Plate (1999), one of the chief advantages of additive models such as GAM is their interpretability. Plate also notes that by allowing high-order interactions as well as low-order interactions, one can trade off interpretability with predictive accuracy. In the case where the hyperparameters indicate that most of the variance in a function can be explained by low-order interactions, it is useful and easy to plot the corresponding low-order functions, as in Figure 1.1.

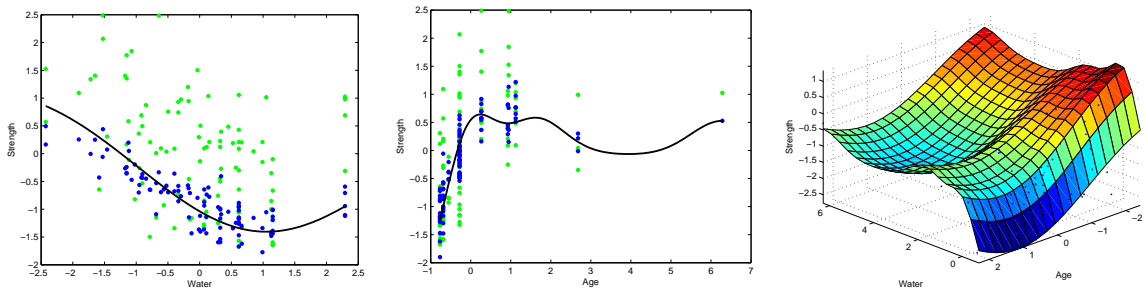


Fig. 1.1 Low-order functions on the concrete dataset. Left, centre: by considering only first-order terms of the additive kernel, we recover a form of Generalized Additive Model, and can plot the corresponding 1-dimensional functions. Green points indicate the original data, blue points are data after the mean contribution from the other dimensions' first-order terms has been subtracted. The black line is the posterior mean of a GP with only one term in its kernel. Right: the posterior mean of a GP with only one second-order term in its kernel.

### 1.3.3 Efficient Evaluation of Additive Kernels

An additive kernel over  $D$  inputs with interactions up to order  $n$  has  $O(2^n)$  terms. Naïvely summing over these terms quickly becomes intractable. In this section, we show how one can evaluate the sum over all terms in  $O(D^2)$ .

To efficiently compute the additive kernel, we exploit the fact that the  $n$ th order additive kernel corresponds to the  $n$ th *elementary symmetric polynomial* (Macdonald, 1998) of the base kernels, which we denote  $e_n$ . For example: if  $\mathbf{x}$  has 4 input dimensions

( $D = 4$ ), and if we let  $z_i = k_i(x_i, x'_i)$ , then

$$k_{add_1}(\mathbf{x}, \mathbf{x}') = e_1(z_1, z_2, z_3, z_4) = z_1 + z_2 + z_3 + z_4 \quad (1.11)$$

$$k_{add_2}(\mathbf{x}, \mathbf{x}') = e_2(z_1, z_2, z_3, z_4) = z_1 z_2 + z_1 z_3 + z_1 z_4 + z_2 z_3 + z_2 z_4 + z_3 z_4 \quad (1.12)$$

$$k_{add_3}(\mathbf{x}, \mathbf{x}') = e_3(z_1, z_2, z_3, z_4) = z_1 z_2 z_3 + z_1 z_2 z_4 + z_1 z_3 z_4 + z_2 z_3 z_4 \quad (1.13)$$

$$k_{add_4}(\mathbf{x}, \mathbf{x}') = e_4(z_1, z_2, z_3, z_4) = z_1 z_2 z_3 z_4 \quad (1.14)$$

The Newton-Girard formulae give an efficient recursive form for computing these polynomials. If we define  $s_k$  to be the  $k$ th power sum:  $s_k(z_1, z_2, \dots, z_D) = \sum_{i=1}^D z_i^k$ , then

$$k_{add_n}(\mathbf{x}, \mathbf{x}') = e_n(z_1, \dots, z_D) = \frac{1}{n} \sum_{k=1}^n (-1)^{(k-1)} e_{n-k}(z_1, \dots, z_D) s_k(z_1, \dots, z_D) \quad (1.15)$$

Where  $e_0 \triangleq 1$ . The Newton-Girard formulae have time complexity  $O(D^2)$ , while computing a sum over an exponential number of terms.

Conveniently, we can use the same trick to efficiently compute all of the necessary derivatives of the additive kernel with respect to the base kernels. We merely need to remove the kernel of interest from each term of the polynomials:

$$\frac{\partial k_{add_n}}{\partial z_j} = e_{n-1}(z_1, \dots, z_{j-1}, z_{j+1}, \dots, z_D) \quad (1.16)$$

This trick allows us to optimize the base kernel hyperparameters with respect to the marginal likelihood.

### 1.3.4 Computation

The computational cost of evaluating the Gram matrix of a product kernel (such as the SE kernel) is  $O(N^2 D)$ , while the cost of evaluating the Gram matrix of the additive kernel is  $O(N^2 D R)$ , where  $R$  is the maximum degree of interaction allowed (up to  $D$ ). In higher dimensions, this can be a significant cost, even relative to the fixed  $O(N^3)$  cost of inverting the Gram matrix. However, as our experiments show, typically only the first few orders of interaction are important for modeling a given function; hence if one is computationally limited, one can simply limit the maximum degree of interaction without losing much accuracy.

Additive Gaussian processes are particularly appealing in practice because their use requires only the specification of the base kernel. All other aspects of GP inference

remain the same. All of the experiments in this chapter were performed using the standard GPML toolbox<sup>1</sup>; code to perform all experiments is available at the author’s website.<sup>2</sup>

## 1.4 Related Work

Plate (1999) constructs a form of additive GP, but using only the first-order and  $D$ th order terms. This model is motivated by the desire to trade off the interpretability of first-order models, with the flexibility of full-order models. Our experiments show that often, the intermediate degrees of interaction contribute most of the variance.

A related functional ANOVA GP model Kaufman and Sain (2010) decomposes the *mean* function into a weighted sum of GPs. However, the effect of a particular degree of interaction cannot be quantified by that approach. Also, computationally, the Gibbs sampling approach used in Kaufman and Sain (2010) is disadvantageous.

Christoudias et al. (2009) previously showed how mixtures of kernels can be learnt by gradient descent in the Gaussian process framework. They call this *Bayesian localized multiple kernel learning*. However, their approach learns a mixture over a small, fixed set of kernels, while our method learns a mixture over all possible products of those kernels.

### 1.4.1 Hierarchical Kernel Learning

We note that a similar breakthrough has recently been made by Bach (2009) called Hierarchical Kernel Learning (HKL). HKL explores a similar class of models, and sidesteps the possibly exponential number of interaction terms by cleverly selecting only a tractable subset. However, this method suffers considerably from the fact that cross-validation must be used to set hyperparameters. In addition, the machinery necessary to train these models is immense. Finally, on real datasets, HKL is outperformed by the standard SE-GP (Bach, 2009).

Bach (2009) uses a regularized optimization framework to learn a weighted sum over an exponential number of kernels which can be computed in polynomial time. The subsets of kernels considered by this method are restricted to be a *hull* of kernels.<sup>3</sup> Given

<sup>1</sup>Available at <http://www.gaussianprocess.org/gpml/code/>

<sup>2</sup>Example code available at: <http://mlg.eng.cam.ac.uk/duvenaud/>

<sup>3</sup>In this setting, a hull is defined as a subset of all terms such that if term  $\prod_{j \in J} k_j(\mathbf{x}, \mathbf{x}')$  is included in the subset, then so are all terms  $\prod_{j \in J/i} k_j(\mathbf{x}, \mathbf{x}')$ , for all  $i \in J$ . For details, see (Bach, 2009).



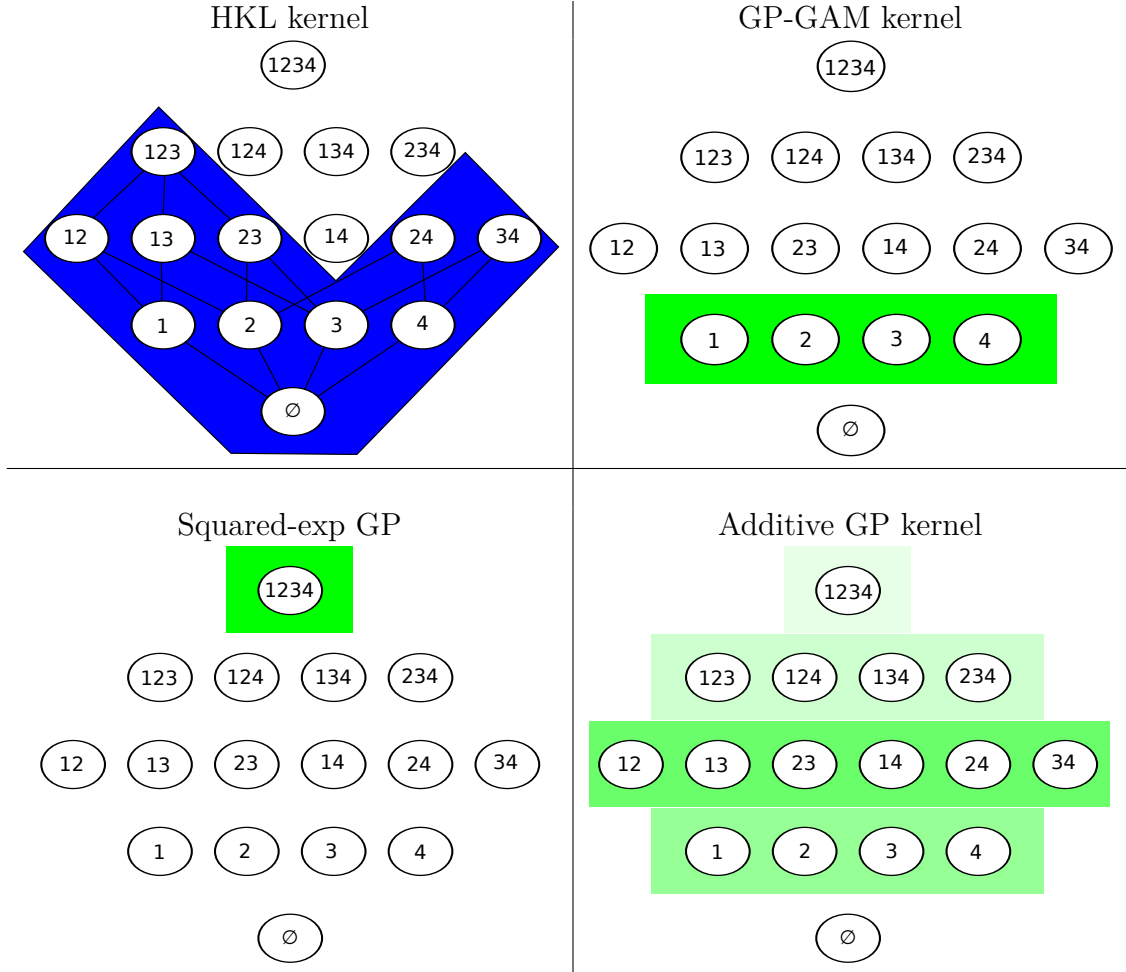


Fig. 1.2 A comparison of different additive model classes. Nodes represent different interaction terms, ranging from first-order to fourth-order interactions. Far left: HKL can select a hull of interaction terms, but must use a pre-determined weighting over those terms. Far right: the additive GP model can weight each order of interaction separately. Neither HKL nor the additive model dominate one another in terms of flexibility, however the GP-GAM and the SE-GP are special cases of additive GPs.

each dimension’s kernel, and a pre-defined weighting over all terms, HKL performs model selection by searching over hulls of interaction terms. In [Bach \(2009\)](#), Bach also fixes the relative weighting between orders of interaction with a single term  $\alpha$ , computing the sum over all orders by:

$$k_a(\mathbf{x}, \mathbf{x}') = v_D^2 \prod_{d=1}^D (1 + \alpha k_d(x_d, x'_d)) \quad (1.17)$$

which has computational complexity  $O(D)$ . However, this formulation forces the weight of all  $n$ th order terms to be weighted by  $\alpha^n$ .

Figure 1.2 contrasts the HKL hull-selection method with the Additive GP hyperparameter-learning method. Neither method dominates the other in flexibility. The main difficulty with the approach of [Bach \(2009\)](#) is that hyperparameters are hard to set other than by cross-validation. In contrast, our method optimizes the hyperparameters of each dimension’s base kernel, as well as the relative weighting of each order of interaction.

## 1.4.2 ANOVA Procedures

Vapnik [Vapnik \(1998\)](#) introduces the support vector ANOVA decomposition, which has the same form as our additive kernel. However, they recommend approximating the sum over all  $D$  orders with only one term “of appropriate order”, presumably because of the difficulty of setting the hyperparameters of an SVM. [Stitson et al. \(1999\)](#) performed experiments which favourably compared the support vector ANOVA decomposition to polynomial and spline kernels. They too allowed only one order to be active, and set hyperparameters by cross-validation.

A closely related procedure from [Wahba \(1990\)](#) is smoothing-splines ANOVA (SS-ANOVA). An SS-ANOVA model is estimated as a weighted sum of splines along each dimension, plus a sum of splines over all pairs of dimensions, all triplets, etc, with each individual interaction term having a separate weighting parameter. Because the number of terms to consider grows exponentially in the order, in practice, only terms of first and second order are usually considered. Learning in SS-ANOVA is usually done via penalized-maximum likelihood with a fixed sparsity hyperparameter.

In contrast to these procedures, our method can easily include all  $D$  orders of interaction, each weighted by a separate hyperparameter. As well, we can learn kernel hyperparameters individually per input dimension, allowing automatic relevance determination to operate.

### 1.4.3 Non-local Interactions

By far the most popular kernels for regression and classification tasks on continuous data are the squared exponential (Gaussian) kernel, and the Matérn kernels. These kernels depend only on the scaled Euclidean distance between two points, both having the form:

$$k(\mathbf{x}, \mathbf{x}') = f\left(\sum_{d=1}^D (x_d - x'_d)^2 / l_d^2\right) \quad (1.18)$$

Bengio et al. (2006) argue that models based on squared-exponential kernels are particularly susceptible to the *curse of dimensionality*. They emphasize that the locality of the kernels means that these models cannot capture non-local structure. They argue that many functions that we care about have such structure. Methods based solely on local kernels will require training examples at all combinations of relevant inputs.

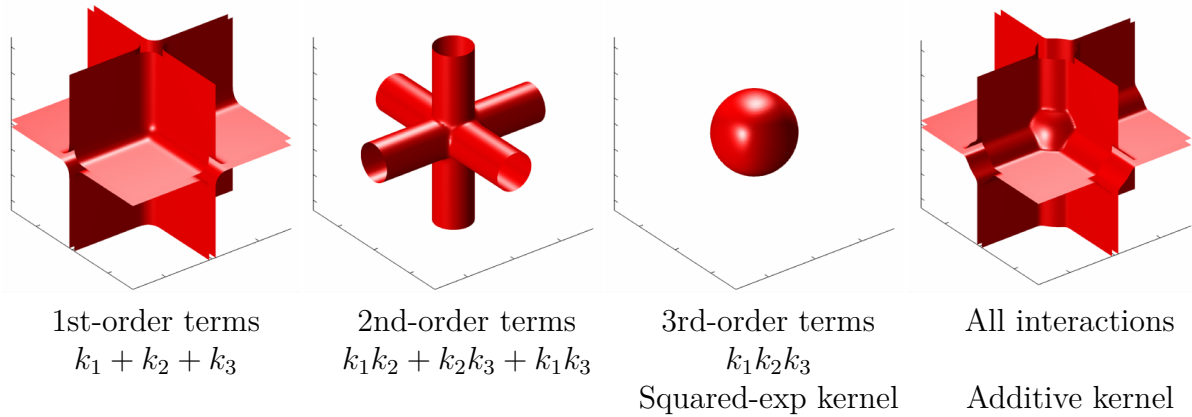


Fig. 1.3 Isocontours of additive kernels in 3 dimensions. The third-order kernel only considers nearby points relevant, while the lower-order kernels allow the output to depend on distant points, as long as they share one or more input value.

Additive kernels have a much more complex structure, and allow extrapolation based on distant parts of the input space, without spreading the mass of the kernel over the whole space. For example, additive kernels of the second order allow strong non-local interactions between any points which are similar in any two input dimensions. Figure 1.3 provides a geometric comparison between squared-exponential kernels and additive kernels in 3 dimensions.

## 1.5 Experiments

### Overfitting

A multivariate squared exponential kernel has  $D + 1$  hyperparameters, while an additive function with one hyperparameter per dimension (e.g., a scale parameter) on each dimension’s covariance function will have  $D \times 2$  hyperparameters. Because each additional hyperparameter increases the tendency to overfit, we recommend allowing only one hyperparameter per input dimension.

Results were averaged over 10 cross-validation folds. To limit the possibility of overfitting, we parametrized each of our one-dimensional kernel functions with only the lengthscale.

On a diverse collection of datasets, we compared five different models. In the results tables below, GP Additive refers to a GP using the additive kernel with squared-exp base kernels. For speed, we limited the maximum order of interaction in the additive kernels to 10. GP-GAM denotes an additive GP model with only first-order interactions. GP Squared-Exp is a GP model with a squared-exponential ARD kernel. HKL<sup>4</sup> was run using the all-subsets kernel, which corresponds to the same set of kernels as considered by the additive GP with a squared-exp base kernel.

For all GP models, we fit hyperparameters by the standard method of maximizing training-set marginal likelihood, using L-BFGS [Nocedal \(1980\)](#) for 500 iterations, allowing five random restarts. In addition to learning kernel hyperparameters, we fit a constant mean function to the data. In the classification experiments, GP inference was done using Expectation Propagation [Minka \(2001\)](#).

### 1.5.1 Bach Synthetic Dataset

In addition to standard UCI repository datasets, we generated a synthetic dataset following the same recipe as [Bach \(2009\)](#). This dataset was designed by [Bach \(2009\)](#) to demonstrate the advantages of HKL over GP-ARD. From a covariance matrix drawn from a Wishart distribution with 1024 degrees of freedom, we select 8 variables. We then construct the non-linear function  $f(X) = \sum_{i=1}^4 \sum_{j=1+1}^4 X_i X_j + \epsilon$ , which sums all 2-way products of the first 4 variables, and adds Gaussian noise  $\epsilon$ . This dataset is one which can be predicted well by a kernel which is a sum of two-way interactions over the first 4 variables, ignoring the extra 4 noisy copies.

<sup>4</sup>Code for HKL available at <http://www.di.ens.fr/~fbach/hkl/>

If the dataset is large enough, HKL can construct a hull around only those subsets of cross terms that are optimal for predicting the output. GP-ARD, in contrast, can only learn to ignore the noisy copy variables, but cannot learn to ignore the higher-term interactions between the predictive variables. However, a GP with an additive kernel can learn both to ignore irrelevant variables, and to ignore certain orders of interaction. In this example, the additive GP is able to recover the relevant structure.

## 1.5.2 Results

Tables 1.2, 1.3, 1.4 and 1.5 show mean performance across 10 train-test splits. Because HKL does not specify a noise model, it could not be included in the likelihood comparisons.

Table 1.2 Regression Mean Squared Error

Method	bach	concrete	pumadyn-8nh	servo	housing
Linear Regression	1.031	0.404	0.641	0.523	0.289
GP GAM	1.259	0.149	0.598	0.281	0.161
HKL	<b>0.199</b>	0.147	0.346	0.199	0.151
GP Squared-exp	<b>0.045</b>	0.157	<b>0.317</b>	<b>0.126</b>	<b>0.092</b>
GP Additive	<b>0.045</b>	<b>0.089</b>	<b>0.316</b>	<b>0.110</b>	<b>0.102</b>

Table 1.3 Regression Negative Log Likelihood

Method	bach	concrete	pumadyn-8nh	servo	housing
Linear Regression	2.430	1.403	1.881	1.678	1.052
GP GAM	1.708	0.467	1.195	0.800	0.457
GP Squared-exp	<b>-0.131</b>	0.398	<b>0.843</b>	0.429	<b>0.207</b>
GP Additive	<b>-0.131</b>	<b>0.114</b>	<b>0.841</b>	<b>0.309</b>	<b>0.194</b>

The model with best performance on each dataset is in bold, along with all other models that were not significantly different under a paired  $t$ -test. The additive model never performs significantly worse than any other model, and sometimes performs significantly better than all other models. The difference between all methods is larger in the case of regression experiments. The performance of HKL is consistent with the results in Bach (2009), performing competitively but slightly worse than SE-GP.

Table 1.4 Classification percent error

Method	breast	pima	sonar	ionosphere	liver	heart
Logistic Regression	7.611	24.392	26.786	16.810	45.060	<b>16.082</b>
GP GAM	<b>5.189</b>	<b>22.419</b>	<b>15.786</b>	<b>8.524</b>	<b>29.842</b>	<b>16.839</b>
HKL	<b>5.377</b>	24.261	<b>21.000</b>	9.119	<b>27.270</b>	<b>18.975</b>
GP Squared-exp	<b>4.734</b>	<b>23.722</b>	<b>16.357</b>	<b>6.833</b>	<b>31.237</b>	<b>20.642</b>
GP Additive	<b>5.566</b>	<b>23.076</b>	<b>15.714</b>	<b>7.976</b>	<b>30.060</b>	<b>18.496</b>

Table 1.5 Classification negative log-likelihood

Method	breast	pima	sonar	ionosphere	liver	heart
Logistic Regression	0.247	0.560	4.609	0.878	0.864	0.575
GP GAM	<b>0.163</b>	<b>0.461</b>	<b>0.377</b>	<b>0.312</b>	<b>0.569</b>	<b>0.393</b>
GP Squared-exp	<b>0.146</b>	0.478	<b>0.425</b>	<b>0.236</b>	<b>0.601</b>	0.480
GP Additive	<b>0.150</b>	<b>0.466</b>	<b>0.409</b>	<b>0.295</b>	<b>0.588</b>	<b>0.415</b>

The additive GP performed best on datasets well-explained by low orders of interaction, and approximately as well as the SE-GP model on datasets which were well explained by high orders of interaction (see table 1.1). Because the additive GP is a superset of both the GP-GAM model and the SE-GP model, instances where the additive GP performs slightly worse are presumably due to over-fitting, or due to the hyperparameter optimization becoming stuck in a local maximum. Additive GP performance can be expected to benefit from integrating out the kernel hyperparameters.

### 1.5.3 Future Work

Since the non-local structure capturable by additive kernels is necessarily axis-aligned, we can naturally consider that combining the hyperparameter optimization with an initial rotation of the input space might allow us to recover non-axis aligned additivity in functions.

Note that we are free to choose a different covariance function along each dimension.

## 1.6 Conclusion

We present additive Gaussian processes: a simple family of models which generalizes two widely-used classes of models. Additive GPs also introduce a tractable new type of structure into the GP framework. Our experiments indicate that such additive structure is present in real datasets, allowing our model to perform better than standard GP models. In the case where no such structure exists, our model can recover arbitrarily flexible models, as well.

In addition to improving modeling efficacy, the additive GP also improves model interpretability: the order variance hyperparameters indicate which sorts of structure are present in our model.

Compared to HKL, which is the only other tractable procedure able to capture the same types of structure, our method benefits from being able to learn individual kernel hyperparameters, as well as the weightings of different orders of interaction. Our experiments show that additive GPs are a state-of-the-art regression model.

# References

- Francis Bach. High-dimensional non-linear variable selection through hierarchical kernel learning. *CoRR*, abs/0909.0844, 2009. (pages 8, 10, 12, and 13)
- Pierre Baldi and Peter J Sadowski. Understanding dropout. In *Advances in Neural Information Processing Systems*, pages 2814–2822, 2013. (page 1)
- Yoshua Bengio, Olivier Delalleau, and Nicolas Le Roux. The curse of highly variable functions for local kernel machines. *Advances in neural information processing systems*, 18:107–114, 2006. ISSN 1049-5258. (page 11)
- M. Christoudias, R. Urtasun, and T. Darrell. Bayesian localized multiple kernel learning. *Technical report, EECS Department, University of California, Berkeley*, 2009. (page 8)
- T.J. Hastie and R.J. Tibshirani. *Generalized additive models*. Chapman & Hall/CRC, 1990. (page 3)
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012. (page 1)
- C.G. Kaufman and S.R. Sain. Bayesian functional anova modeling using Gaussian process prior distributions. *Bayesian Analysis*, 5(1):123–150, 2010. (page 8)
- I.G. Macdonald. *Symmetric functions and Hall polynomials*. Oxford University Press, USA, 1998. ISBN 0198504500. (page 6)
- T.P. Minka. Expectation propagation for approximate bayesian inference. In *Uncertainty in Artificial Intelligence*, volume 17, pages 362–369. Citeseer, 2001. (page 12)



- J.A. Nelder and R.W.M. Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society. Series A (General)*, 135(3):370–384, 1972. (page 3)
- J. Nocedal. Updating quasi-newton matrices with limited storage. *Mathematics of computation*, 35(151):773–782, 1980. (page 12)
- T.A. Plate. Accuracy versus interpretability in flexible modeling: Implementing a trade-off using Gaussian process models. *Behaviormetrika*, 26:29–50, 1999. ISSN 0385-7417. (pages 6 and 8)
- Nitish Srivastava. *Improving neural networks with dropout*. PhD thesis, University of Toronto, 2013. (page 1)
- M. Stitson, A. Gammerman, V. Vapnik, V. Vovk, C. Watkins, and J. Weston. Support vector regression with ANOVA decomposition kernels. *Advances in kernel methods: Support vector learning*, pages 285–292, 1999. (page 10)
- V.N. Vapnik. *Statistical learning theory*, volume 2. Wiley New York, 1998. (page 10)
- G. Wahba. *Spline models for observational data*. Society for Industrial Mathematics, 1990. ISBN 0898712440. (pages 3 and 10)
- Sida Wang and Christopher Manning. Fast dropout training. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pages 118–126, 2013. (page 1)