

# The Automated Statistician for Gaussian Process Classification

Nikola Mrkšić  
Trinity College



**UNIVERSITY OF  
CAMBRIDGE**

*A dissertation submitted to the University of Cambridge  
in partial fulfilment of the requirements for the Part III of  
the Computer Science Tripos*

May 7, 2014

University of Cambridge  
Computer Laboratory  
William Gates Building  
15 JJ Thomson Avenue  
Cambridge CB3 0FD  
UNITED KINGDOM



# Declaration

I Nikola Mrkšić of Trinity College, being a candidate for the Part III in Computer Science, hereby declare that this report and the work described in it are my own work, unaided except as may be specified below, and that the report does not contain material that has already been used to any substantial extent for a comparable purpose.

Total word count: 0

**Signed:**

**Date:**

This dissertation is copyright ©2014 Nikola Mrkšić.

All trademarks used in this dissertation are hereby acknowledged.



# Abstract



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Building an Automated Statistician . . . . .	1
1.2	Contribution . . . . .	1
<b>2</b>	<b>Gaussian Processes</b>	<b>2</b>
2.1	The Role of Kernels for Gaussian Processes . . . . .	2
2.2	Gaussian Process Classification . . . . .	2
2.2.1	The Laplace Approximation to Marginal Likelihood . .	3
2.2.2	Expectation Propagation . . . . .	3
2.2.3	Variational Bayes . . . . .	3
<b>3</b>	<b>Related Work</b>	<b>4</b>
3.1	Kernel learning . . . . .	4
3.2	Kernel structure discovery for regression . . . . .	4
3.3	Additive Gaussian Processes . . . . .	5
<b>4</b>	<b>Kernel Structure Discovery for GP Classification</b>	<b>6</b>
4.1	Defining the Kernel Grammar . . . . .	6
4.1.1	The Search Operators . . . . .	6
4.2	Model selection . . . . .	6
4.3	Optimising the Hyperparameters . . . . .	6
4.3.1	Overfitting . . . . .	7
4.4	Guiding the Structure Search . . . . .	7
4.4.1	Bayesian and Akaike Information Criteria . . . . .	7
4.4.2	The Number of Effective Hyperparameters . . . . .	7
4.4.3	Cross-validated training accuracy . . . . .	8
4.5	Adapting the likelihood function . . . . .	8
4.5.1	Dealing with Outliers . . . . .	8
4.6	Providing Interpretability . . . . .	9
4.6.1	Pima . . . . .	9
4.6.2	Breast . . . . .	11
4.6.3	Heart and Liver . . . . .	13

4.7	Bayesian Model Averaging . . . . .	13
4.7.1	BMA for Predictive Performance . . . . .	13
4.7.2	BMA/Cross-validated model for Model Selection . . . .	15
4.7.3	Liver . . . . .	15
<b>5</b>	<b>Evaluation</b>	<b>17</b>
5.1	Experiments with Synthetic Data . . . . .	17
5.1.1	Adding Salt and Pepper Noise . . . . .	19
5.2	Experiments on Real World Data Sets . . . . .	23
<b>6</b>	<b>Summary and Conclusions</b>	<b>24</b>
6.1	Further Work . . . . .	24



# Chapter 1

## Introduction

### 1.1 Building an Automated Statistician

### 1.2 Contribution

Built a system for automatic kernel discovery for Gaussian Process classification...

# Chapter 2

## Gaussian Processes

1. Define GPs, explain how it is fully specified by the mean and covariance function... Discuss covariance functions. Talk about function as infinite vector to generalize from multinomial Gaussian to this.
2. GP as a prior on function, show some function draws from the prior...
3. Marginal likelihood as criterion for model selection and fitting the hyper-parameters. write out nlml formula and discuss that the first term penalises poor fit, second model complexity, third is normalization. Say this is type II ML.

### 2.1 The Role of Kernels for Gaussian Processes

Describe different types of kernels.

### 2.2 Gaussian Process Classification

Explain why it is harder, non-Gaussian likelihood...

**2.2.1 The Laplace Approximation to Marginal Likelihood**

**2.2.2 Expectation Propagation**

**2.2.3 Variational Bayes**

# Chapter 3

## Related Work

### 3.1 Kernel learning

General review of work done, with different models.

### 3.2 Kernel structure discovery for regression

Unsupervised structure discovery - mention that paper, maybe initially, as a similar type of work... 1 paragraph?

Talk about the fact that these methods fix the structure of the kernel beforehand. Do they? Or they lack interpretability.

As a critique, discuss different information criteria, i.e. the fact that RQs can be penalised more than an SE... preferring products to additive components, say that this is something we evaluate thoroughly and try to improve on.

Write an intro to doing this for classification, why it is harder (less information than regression, harder interpretability, non-Gaussian likelihoods and the need to resort to nlml approximations).

### 3.3 Additive Gaussian Processes

Find other relevant work - Dave mentioned some of these. Read through all the papers' relevant work sections.

# Chapter 4

## Kernel Structure Discovery for GP Classification

### 4.1 Defining the Kernel Grammar

Draw the basic kernels, describe additive and product kernels.

Present challenges for classification, difference from regression, lack of clear component interpretability, as opposed to i.e. time-series data.

#### 4.1.1 The Search Operators

Adding or multiplying with a base kernel.

### 4.2 Model selection

### 4.3 Optimising the Hyperparameters

Subsampling the training data to optimize the process.

Discuss parallelisation. Maybe add a table of running times and numbers of restarts.

### 4.3.1 Overfitting

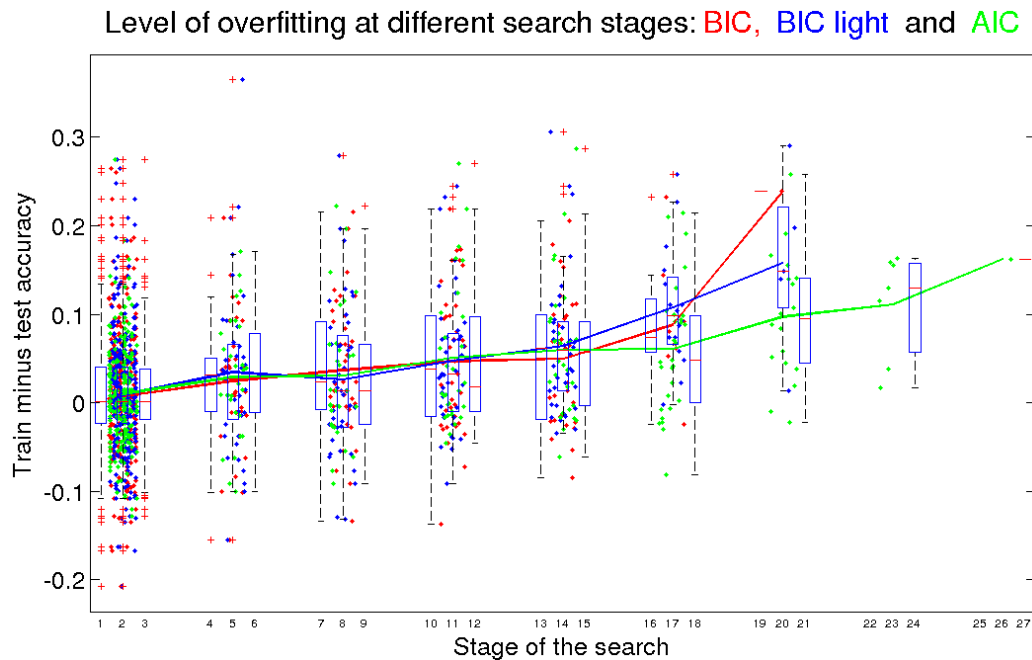
Dealing with very small and very large hyperparameters.

## 4.4 Guiding the Structure Search

### 4.4.1 Bayesian and Akaike Information Criteria

### 4.4.2 The Number of Effective Hyperparameters

Outline BIC light, present it as middle ground between full BIC and AIC. Insert the number of search steps figure that shows overfitting.



### 4.4.3 Cross-validated training accuracy

## 4.5 Adapting the likelihood function

### 4.5.1 Dealing with Outliers

a problem with classification and BIC is that estimating the number of effective hyperparameters is hard: not all of them affect the classification boundary equally. The only thing important for classification is the zero-crossings of the function (where it is greater, and where it is smaller than 0, for the purposes of prediction). Hence, smaller length-scales are the most important factor (as they determine the predictive mean and where the crossings exactly are). Larger lengthscales equate to multiplication with a constant function, so less beneficial. Signal variances are useful for computing the predictive variances, but we only care about the predictive means for determining the target label. Not counting signal frequencies might add a quick boost to quality of BIC as a model discriminator (i.e. not number of product terms, but 0 - wont decrease by 2, but might help our choice-making).

The reason that GPs out of the box are easily outperformed by classification specific schemes such as SVMs is that soft-margin SVMs are less certain of their predictions, and hence are less certain about points far away from the high density regions of the data. SVMs are never more certain than the estimated level of pepper noise (outliers in the sense of noise in the target labels of the data) whereas GPs would be very certain about points close to one of the clusters...

As a potential solution, we might want to use a @likMix which adds the pepper noise into the likelihood function, determining the level of the outliers and thus having the likelihood (i.e. sigmoid) which doesnt range from 0 to 1, but from  $\alpha$  to  $(1 - \alpha)$ , where  $\alpha$  is the level of pepper noise, which is, again, to be learnt from the data? This can be implemented in GPML using @likMix, and is something we can first evaluate w.r.t. synthetic data we previously

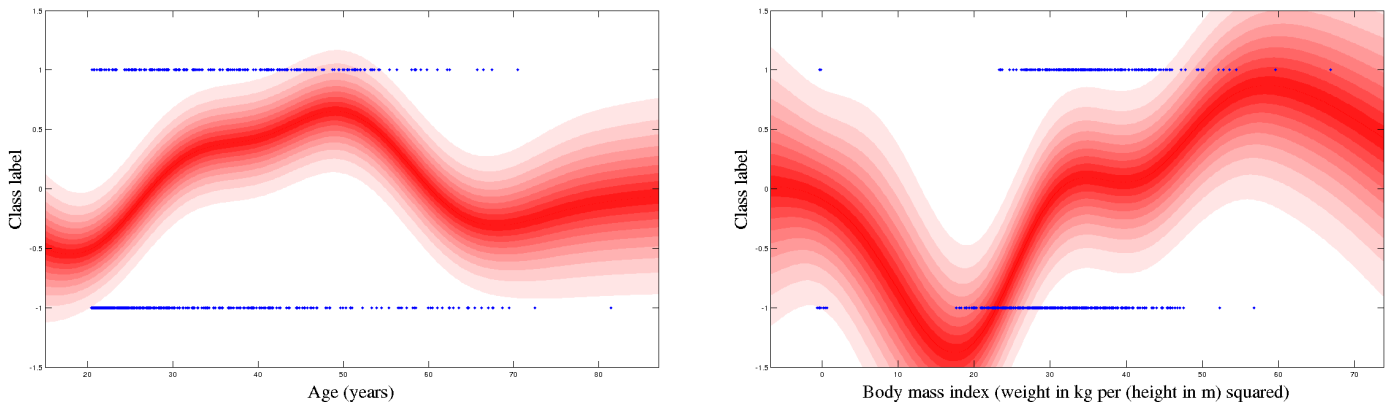


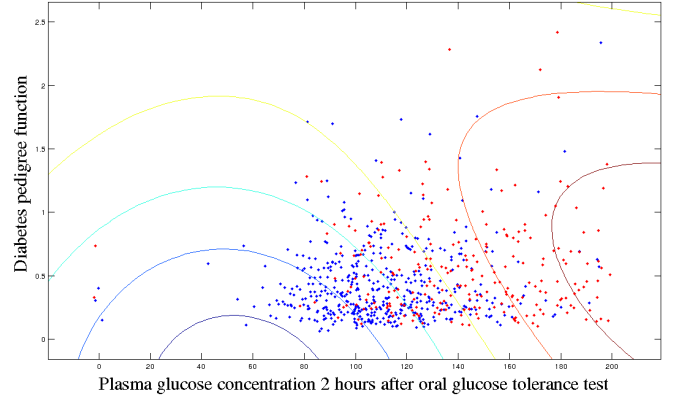
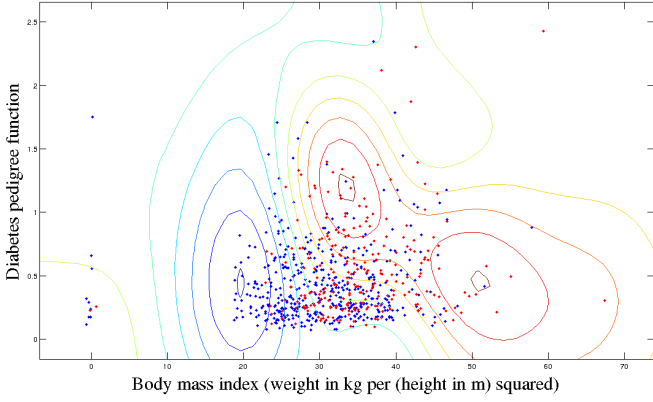
generated. The only difference is that now we add additional salt and pepper noise (random outliers) into the data and we compare performance (on this dataset) of likMix vs likErf on its own.

## 4.6 Providing Interpretability

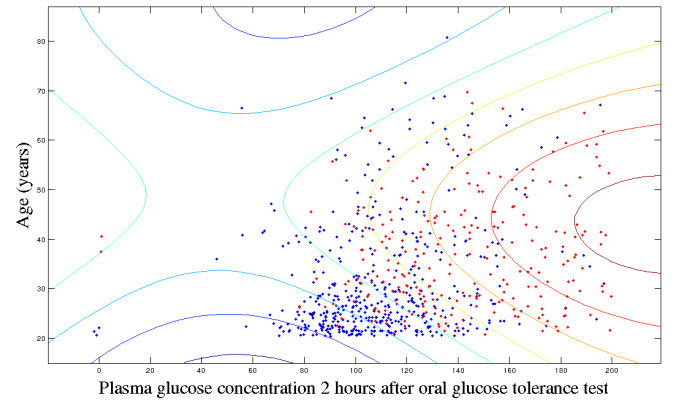
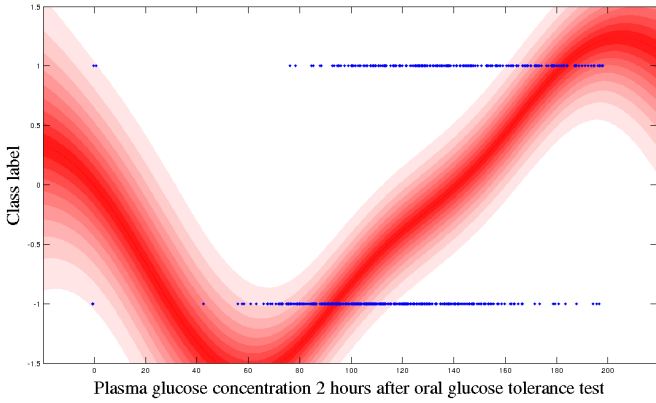
### 4.6.1 Pima

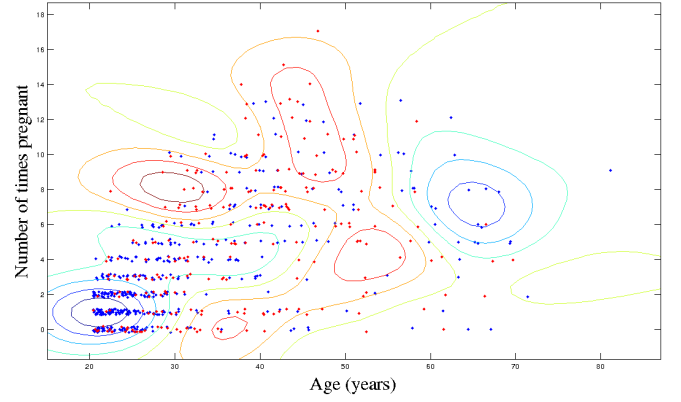
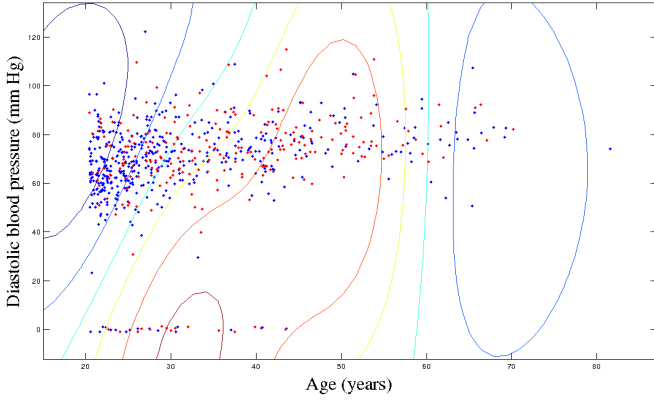
The Pima Indian Diabetes dataset contains medical measurements of 768 individuals, together with an indicator of whether they later developed diabetes or not. There are 8 attributes in total, 4 of which feature in the kernel structures returned by the kernel discovery procedure. These are dimensions 2, 6, 7 and 8: plasma glucose concentration, body mass index (BMI), diabetes pedigree function and the age of the patient. The remaining four attributes seem to carry less information regarding the likelihood of developing diabetes, at least in a dataset of this size and using this search procedure. These four features are the number of pregnancies, diastolic blood pressure, triceps skin fold thickness and 2-hour serum insulin measures. The structure returned most often (4 times using 10-fold cross validation for computing classification accuracy) was  $[2 \times 7] + [6] + [8]$ , closely followed by  $[2] + [8] + [6 \times 7]$ , which occurred for 3 folds. The four posterior plots indicating the class boundaries determined with respect to these dimensions are shown below.





In addition to these additive components, the remaining three folds' structures contained occurrences of  $[2]$ ,  $[2 \times 3 \times 7]$  and  $[2 \times 3 \times 7]$ . Showing the three-way dependencies in 2D is not straightforward. Below, we show the posterior plots for  $[2]$ ,  $[2 \times 8]$  and  $[3 \times 8]$ , showing that interactions of these dimensions produce points well separated by the classifiers produced (with the possible exception of the interaction between age and diastolic blood pressure). The two way interaction between the number of pregnancies and age is also shown: this dependency was found during various stages of the search but discarded, even though the classifier seems to observe the general pattern present in the data set.



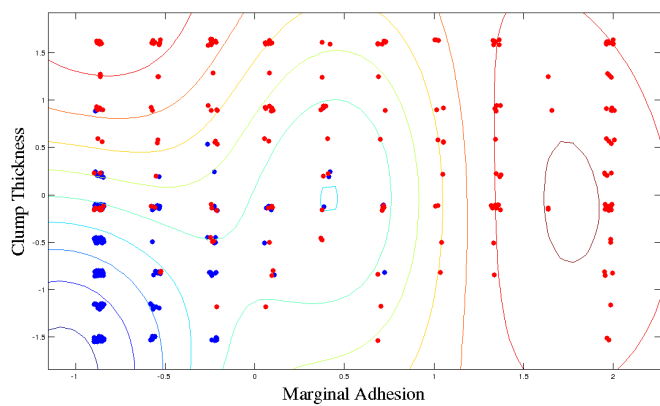
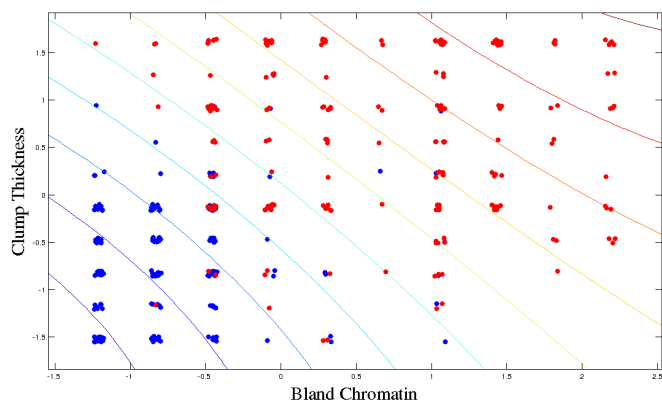
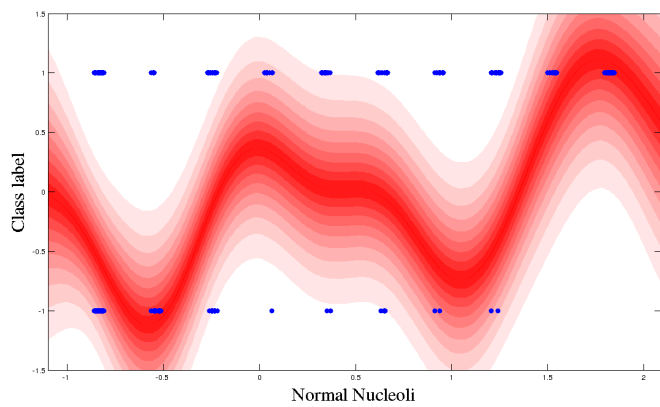
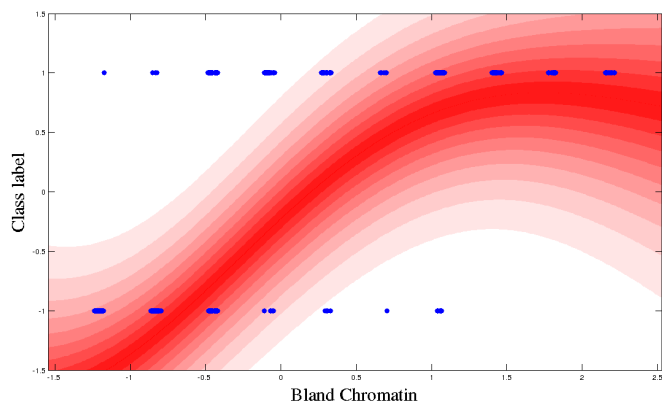
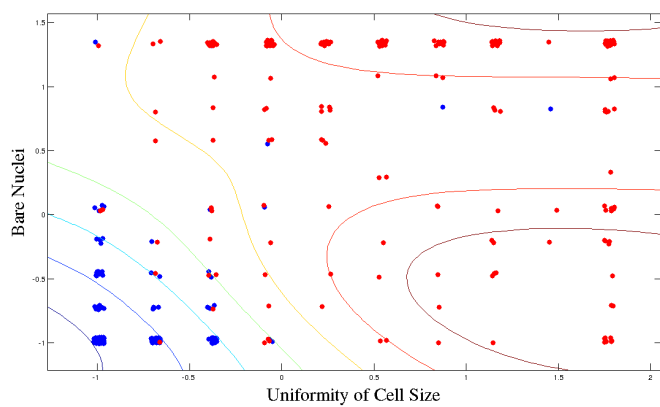
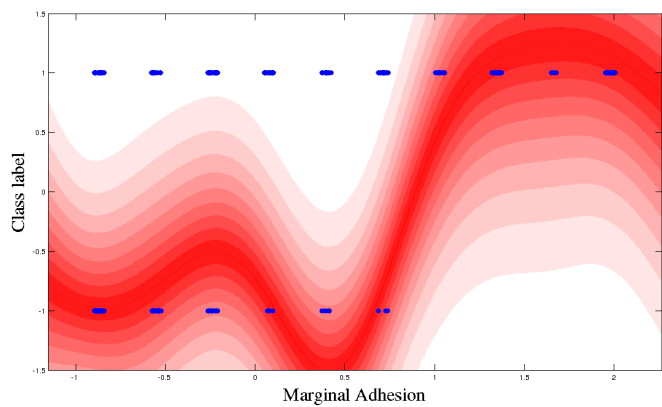


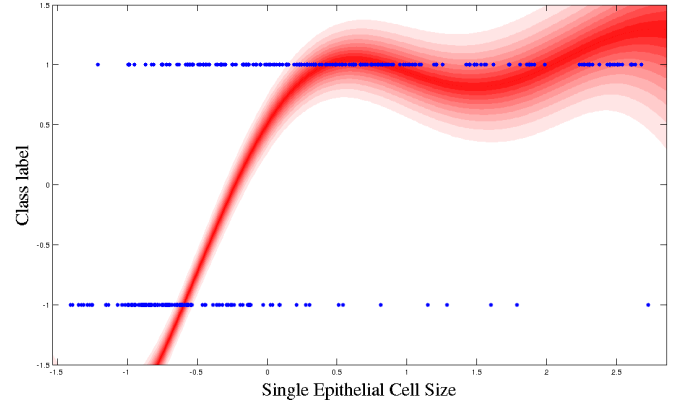
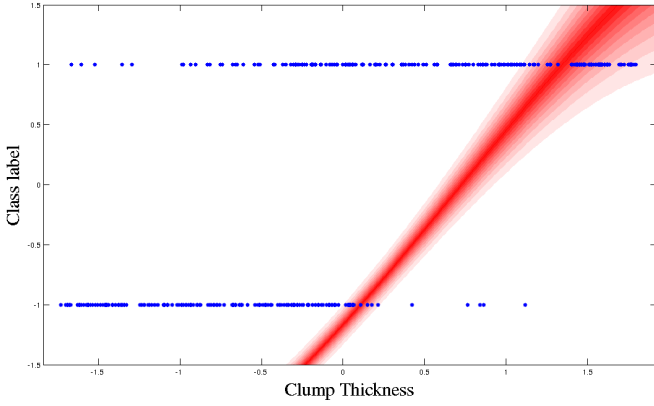
### 4.6.2 Breast

Wisconsin Diagnostic Breast Cancer Dataset is a 9-dimensional dataset which contains measured attributes of 449 individuals, together with an indicator of whether or not they later developed breast cancer. Across the 10 folds generated, there were nine different kernel structures that the procedure returned, with seven out of the nine attributes utilised in these structures. The dimensions most frequently found were 1, 2, 4, 6 and 8. Most kernels returned represent some combination of these five base kernels. Dimension 5 features once, and dimension 7 three times. Interestingly, whenever  $SE_7$  was introduced, the accuracy increased far above the levels attained in the other kernels (ranging from 97% to 100%), but it was part of the final kernel in only three out of the ten folds evaluated.

Below, we show those components of the 10 final structures which represent 1 or 2-way interactions.  $[1 \times 2 \times 6]$  features in three folds, as do some other high dimensional interactions, but we are still unable to visualise those.

The first six plots show individual components of the final kernel structures obtained using the structure search (not necessarily across the same fold). The final row shows dimensions 1 and 5 using a single SE kernel in the respective dimension.





### 4.6.3 Heart and Liver

These two data sets require dealing with representation of discretised data. Plots don't make much sense (yet). Adding jitter might alleviate these issues.

## 4.7 Bayesian Model Averaging

### 4.7.1 BMA for Predictive Performance

Averaging assigned labels or probabilities converges to the same thing when many models are used.

Discuss assigning different alphas, make a plot showing how the accuracy changes.

Currently averaging class labels across different models, and  $\alpha$  shows the BIC values scaling factor used to assign the weight factors to different classes.

From the set of models in the layers above, in and below the current model, we use their BIC values  $b_1, b_2, \dots$  to assign each model a weight:

$$w_i = \frac{e^{-\alpha b_i}}{\sum_j e^{-\alpha b_j}}$$

Table 4.1: Bayesian Model Averaging, classification errors across folds.

Liver			Pima		
Best model	$\alpha = 0.5$	$\alpha = 1$	Best model	$\alpha = 0.1$	$\alpha = 1$
23.53%	<b>20.59%</b>	<b>20.59%</b>	21.05%	22.37%	21.05%
20.00%	22.86%	20.00%	23.38%	23.38%	24.68%
37.14%	37.14%	37.14%	15.58%	19.48%	16.88%
40.00%	<b>31.43%</b>	<b>34.29%</b>	18.18%	19.48%	18.18%
28.57%	31.43%	28.57%	19.48%	20.78%	22.08%
26.47%	<b>23.53%</b>	26.47%	29.87%	<b>28.57%</b>	<b>28.57%</b>
23.53%	<b>20.59%</b>	26.47%	29.87%	29.87%	36.36%
26.47%	<b>23.53%</b>	26.47%	19.48%	<b>16.88%</b>	<b>14.29%</b>
25.71%	28.57%	28.57%	25.00%	<b>19.74%</b>	<b>23.68%</b>
23.53%	<b>17.65%</b>	<b>17.65%</b>	20.78%	<b>19.48%</b>	<b>19.48%</b>
<b>27.50%</b>	<b>25.73%</b>	<b>26.62%</b>	<b>22.27%</b>	<b>22.00%</b>	22.53%

Heart			Breast		
Best model	$\alpha = 0.5$	$\alpha = 1$	Best model	$\alpha = 0.5$	$\alpha = 1$
17.24%	20.69%	17.24%	2.27%	2.27%	2.27%
13.33%	53.33%	53.33%	4.55%	4.55%	4.55%
20.69%	20.69%	24.14%	11.11%	<b>8.89%</b>	<b>8.89%</b>
6.90%	<b>3.45%</b>	<b>3.45%</b>	15.56%	<b>11.11%</b>	<b>11.11%</b>
16.67%	<b>10.00%</b>	<b>10.00%</b>	2.27%	2.27%	2.27%
26.67%	26.67%	30.00%	10.87%	<b>6.52%</b>	<b>6.52%</b>
30.00%	53.33%	53.33%	0.00%	0.00%	0.00%
20.00%	23.33%	23.33%	2.22%	2.22%	2.22%
16.67%	20.00%	20.00%	4.35%	4.35%	4.35%
10.00%	16.67%	16.67%	11.11%	<b>8.89%</b>	11.11%
<b>17.82%</b>	<b>24.82%</b>	<b>25.15%</b>	<b>6.43%</b>	<b>5.11%</b>	<b>5.33%</b>

### 4.7.2 BMA/Cross-validated model for Model Selection

Choosing the most frequently occurring model structure, then optimising hyperparameters on each training fold and evaluating on the respective test fold gets the following figures:

### 4.7.3 Liver

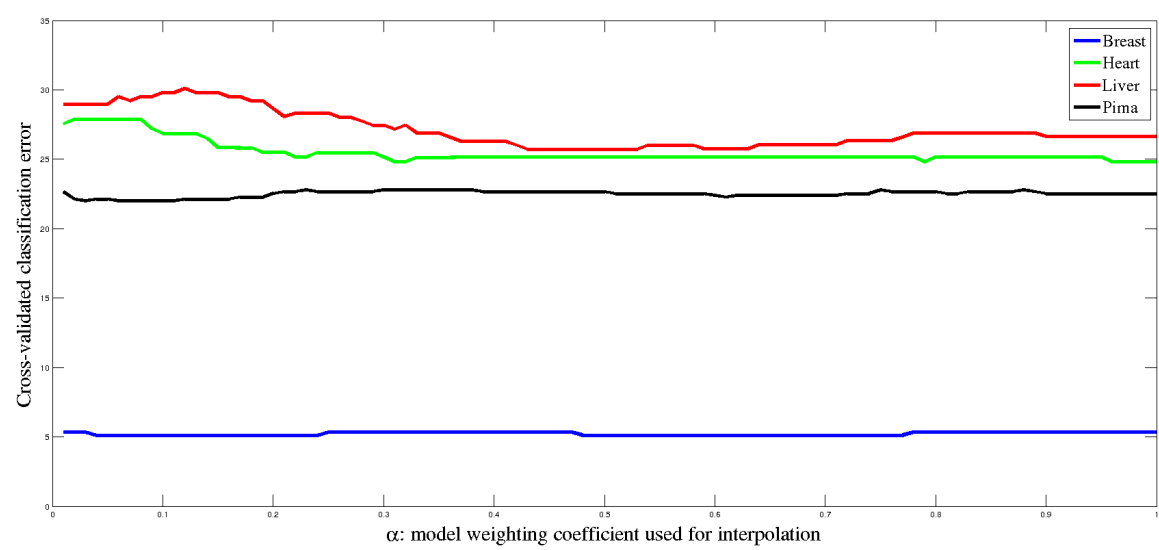
Kernel chosen for this data set is:  $SE_1 + SE_3 \times SE_4 \times SE_5 + SE_6$ , and the testing errors on the ten folds are:

Liver	
Per-fold models	'Best' model
23.53%	23.53%
20.00%	20.00%
37.14%	37.14%
40.00%	<b>34.29%</b>
28.57%	34.29%
26.47%	<b>23.53%</b>
23.53%	26.47%
26.47%	26.47%
25.71%	25.71%
23.53%	<b>14.71%</b>
27.50%	<b>26.61%</b>

For breast, picking  $[1 \times 2 \times 6 + 7]$  leads to the following performance, yielding an improvement of 0.9% over the standard BIC light version.

The average error is 26.61%, which improves on the previous value of 27.5%, but falls short of using weighted model averaging.

Breast
'Best' model
2.27%
4.55%
11.11%
11.11%
2.27%
10.87%
0.00%
2.22%





# Chapter 5

## Evaluation

In this section, we consider the results achieved by the structure discovery procedure presented in the previous section. We first provide a proof of concept for the procedure by showing that it is able to extract correct kernel structure from synthetic data generated using squared exponential kernels. We then proceed to demonstrate that the performance of the procedure on real world data sets is on par with other state of the art methods such as additive Gaussian Processes and Random Forests. Finally, we present the kernel decompositions on the real world data sets which can be used to uncover and visualise the underlying data patterns which dictate class membership of the data points.

### 5.1 Experiments with Synthetic Data

To prove that the greedy structure search procedure is able to extract structure from data, the algorithm was first applied to data drawn from a single GP prior. If the BIC guiding criterion indeed picked the correct model based on marginal likelihood, the procedure should be able to recover the original kernel used to generate the data. The amount of data available to the procedure, as well as the signal-to-noise ratio in the data were varied across

experiments. As we decrease the noise levels and add more data points, the structure search should get closer to the underlying truth, that is the original kernel used to generate the data.

True Kernel	N	Kernel recovered (SNR = 1)	Kernel recovered (SNR = 100)
$[1]$ <i>3 dimensions</i>	100	$[1]$	$[1]$
	300	$[1]$	$[1]$
	500	$[1]$	$[1]$
$[2] + [2] + [2]$ <i>3 dimensions</i>	100	$[2]$	$[2]$
	300	$[2]$	$[2] + [2]$
	500	$[2]$	$[2] + [2]$
$[2 \times 3]$ <i>3 dimensions</i>	100	$[2 \times 3]$	$[2 \times 3]$
	300	$[2 \times 3]$	$[2 \times 3]$
	500	$[2 \times 3]$	$[2 \times 3]$
$[1] + [2 \times 3] + [4]$ <i>4 dimensions</i>	100	$[1] + [4]$	$[1] + [2] + [4]$
	300	$[1 \times 4] + [2 \times 3]$	$[1] + [2 \times 3] + [4]$
	500	$[1 \times 4] + [2 \times 3]$	$[1] + [2 \times 3] + [4]$
$[1] + [2 \times 3] + [4]$ <i>10 dimensions</i>	100	$[1] + [4]$	$[1] + [2] + [4]$
	300	$[1 \times 4] + [2 \times 3]$	$[1] + [2 \times 3] + [4]$
	500	$[1 \times 4] + [2 \times 3]$	$[1] + [2 \times 3] + [4]$
$[1] + [2 \times 3] + [4] + [5 \times 6]$ <i>10 dimensions</i>	100	$[4] + [5]$	$[1 \times 9] + [4]$
	300	$[1] + [2] + [4] + [5 \times 6 \times 7]$	$[1] + [1 \times 5 \times 6] + [2 \times 3] + [4]$
	500	$[1] + [1 \times 4] + [2 \times 3] + [5 \times 6]$	$[1] + [1 \times 4 \times 5 \times 6] + [2 \times 3] + [4]$
$[3 \times 5 \times 7]$ <i>10 dimensions</i>	100	$[4 \times 7]$	$[3 \times 5 \times 7]$
	300	$[2 \times 3 \times 5 \times 7]$	$[3 \times 5 \times 7]$
	500	$[2 \times 3 \times 5 \times 7]$	$[3 \times 5 \times 7]$
$[1] + [3 \times 5 \times 7] + [10]$ <i>10 dimensions</i>	100	$[1 \times 3] + [10]$	$[1 \times 10]$
	300	$[1] + [10]$	$[1] + [1 \times 10] + [3 \times 5 \times 7]$
	500	$[1] + [10]$	$[1] + [3 \times 5 \times 7] + [10]$
$[3 \times 5 \times 7 \times 9]$ <i>10 dimensions</i>	100	$[3 \times 7]$	$[1] + [7 \times 9]$
	300	$[7 \times 9]$	$[3 \times 5 \times 7 \times 9]$
	500	$[3 \times 5 \times 7 \times 9]$	$[3 \times 5 \times 7 \times 9]$
$[1] + [3 \times 5 \times 7 \times 9] + [10]$ <i>10 dimensions</i>	100	$[9 + 10]$	$[10]$
	300	$[1] + [10]$	$[1 \times 5] + [7] + [10]$
	500	$[1] + [3 \times 5 \times 9] + [10]$	$[1] + [3 \times 5 \times 7 \times 9] + [10]$

Kernel optimal rates give a bound on the performance of the structure search: on average, the search is (by design) unable to construct kernels that are better than the one used to generate the data. The Bayes optimal rates give a theoretical limit on what any model could achieve, as they compare the actual function (y values before adding noise) to the data after the noise is

added. Kernel rate should always be below the Bayes rate...

		SNR = 1			SNR = 100		
True Kernel	N	Test Accuracy	Kernel	Function	Test Accuracy	Kernel	Function
<b>[1]</b> <i>3 dimensions</i>	100	<b>81.0%</b>	80.0%		<b>98.0%</b>	97.0%	
	300	<b>82.7%</b>	82.0%	83.6%	<b>99.3%</b>	98.7%	99.0%
	500	83.6%	83.8%		<b>98.4%</b>	98.4%	
<b>[2] + [2] + [2]</b> <i>3 dimensions</i>	100	74.0%	75.0%		<b>98.0%</b>	98.0%	
	300	<b>77.3%</b>	77.0%	78.2%	<b>98.7%</b>	98.7%	98.8 %
	500	<b>78.2%</b>	78.2%		<b>98.6%</b>	98.6%	
<b>[2 × 3]</b> <i>3 dimensions</i>	100	67.0%	68.0%		<b>89.0%</b>	89.0%	
	300	69.7%	70.3%	74.8%	<b>94.7%</b>	91.3%	98.0 %
	500	71.4%	72.6%		<b>95.8%</b>	93.4%	
<b>[1] + [2 × 3] + [4]</b> <i>3 dimensions</i>	100	71.0%	80.0%		82.0%	87.0%	
	300	73.0%	76.3%	76.4%	<b>94.7%</b>	91.7%	97.4%
	500	73.4%	74.8%		<b>94.2%</b>	91.8%	
<b>[1] + [2 × 3] + [4]</b> <i>10 dimensions</i>	100	71.0%	80.0%		82.0%	87.0%	
	300	73.0%	76.3%	76.4%	<b>94.7%</b>	91.7%	97.4%
	500	73.4%	74.8%		<b>94.2%</b>	91.8%	
<b>[1] + [2 × 3] + [4] + [5 × 6]</b> <i>10 dimensions</i>	100	<b>61.0%</b>	60.0%		70.0%	84.0%	
	300	67.7%	70.7%	76.6%	91.0%	92.3%	97.2%
	500	<b>73.4%</b>	73.0%		<b>94.0%</b>	92.0%	
<b>[3 × 5 × 7]</b> <i>10 dimensions</i>	100	53.0%	64.0%		<b>87.0%</b>	83.0%	
	300	69.3%	71.0%	79.0%	<b>92.7%</b>	89.0%	98.4%
	500	73.2%	73.6%		<b>93.6%</b>	91.6%	
<b>[1] + [3 × 5 × 7] + [10]</b> <i>10 dimensions</i>	100	<b>71.0%</b>	70.0%		87.0%	92.0%	
	300	<b>73.0%</b>	71.3%	75.0%	<b>94.3%</b>	94.3%	97.6%
	500	70.0%	72.2%		<b>93.8%</b>	92.8%	
<b>[3 × 5 × 7 × 9]</b> <i>10 dimensions</i>	100	57.0%	59.0%		65.0%	77.0%	
	300	59.3%	70.3%	77.0%	<b>83.7%</b>	82.7%	97.0 %
	500	<b>71.0%</b>	69.0%		<b>85.6%</b>	82.6%	
<b>[1] + [3 × 5 × 7 × 9] + [10]</b> <i>10 dimensions</i>	100	68.0%	74.0%		77.0%	81.0%	
	300	70.7%	73.3%	78.2%	81.7%	89.0%	96.8%
	500	72.0%	72.4%		<b>90.6%</b>	89.8%	

### 5.1.1 Adding Salt and Pepper Noise

We validated our method’s ability to recover known structure on a set of synthetic datasets. For several composite kernel expressions, we constructed synthetic data by first sampling 100, 300 and 500 points uniformly at random,

then sampling function values at those points from a GP prior. We then added i.i.d. Gaussian noise to the functions, at various signal-to-noise ratios (SNR), as well as different amounts of salt and pepper noise (random outliers in the data set).

Table 5.1 lists the true kernels we used to generate the data. Subscripts indicate which dimension each kernel was applied to. Subsequent columns show the dimensionality  $D$  of the input space, and the kernels chosen by our search for different SNRs and different amounts of added salt and pepper noise. We also show the kernel optimal rates (the accuracy the kernel used to generate the data achieves on the noisy test set) and the function optimal rates (the rate a classifier which knew the *exact* function used to generate the data achieves on the noisy test data set).

Table 5.1: True kernel:  $\text{SE}_1 + \text{SE}_2 + \text{SE}_3$ ,  $D = 3$ .

Data size	SNR	sp_noise	Kernel chosen	Test accuracy	Kernel rate	Bayes rate
100	100	0%	$\text{SE}_1 + \text{SE}_1 \times \text{SE}_3 + \text{SE}_2$	87.0%	91.0%	97.4%
300	100	0%	$\text{SE}_1 + \text{SE}_2 + \text{SE}_3$	94.0%	95.7%	97.4%
500	100	0%	$\text{SE}_1 + \text{SE}_2 + \text{SE}_3$	95.8%	95.4%	97.4%
100	100	5%	$\text{SE}_1 + \text{SE}_2 + \text{SE}_3$	77.0%	80.0%	91.6%
300	100	5%	$\text{SE}_1 \times \text{SE}_3 + \text{SE}_2$	87.0%	85.7%	91.6%
500	100	5%	$\text{SE}_1 \times \text{SE}_2 \times \text{SE}_3$	89.8%	89.8%	91.6%
100	100	20%	$\text{SE}_1 \times \text{SE}_3$	69.0%	69.0%	82.0%
300	100	20%	$\text{SE}_1 \times \text{SE}_3 + \text{SE}_2$	75.3%	73.0%	82.0%
500	100	20%	$\text{SE}_1 \times \text{SE}_3 + \text{SE}_2$	77.6%	74.0%	82.0%
100	1	0%	$\text{SE}_1 + \text{SE}_3$	64.0%	72.0%	77.4%
300	1	0%	$\text{SE}_1 + \text{SE}_3$	74.3%	75.0%	77.4%
500	1	0%	$\text{SE}_1 + \text{SE}_3$	75.6%	76.6%	77.4%
100	1	5%	$\text{SE}_1 + \text{SE}_3$	63.0%	63.0%	74.4%
300	1	5%	$\text{SE}_1 \times \text{SE}_3$	70.7%	68.3%	74.4%
500	1	5%	$\text{SE}_1 \times \text{SE}_3$	72.6%	72.6%	74.4%
100	1	20%	$\text{SE}_1 \times \text{SE}_3$	53.0%	60.0%	68.8%
300	1	20%	$\text{SE}_1 \times \text{SE}_3$	65.3%	65.3%	68.8%
500	1	20%	$\text{SE}_1 \times \text{SE}_3$	66.2%	67.8%	68.8%

Table 5.2: True kernel:  $\text{SE}_1 + \text{SE}_2 \times \text{SE}_3 + \text{SE}_4$ ,  $D = 4$ .

Data size	SNR	sp_noise	Kernel chosen	Test accuracy	Kernel rate	Bayes rate
100	100	0%	$\text{SE}_1 + \text{SE}_2 \times \text{SE}_3 + \text{SE}_4$	87.0%	92.0%	97.4%
300	100	0%	$\text{SE}_1 + \text{SE}_2 \times \text{SE}_3 + \text{SE}_4$	94.0%	94.7%	97.4%
500	100	0%	$\text{SE}_1 + \text{SE}_2 \times \text{SE}_3 + \text{SE}_4$	95.6%	96.2%	97.4%
100	100	5%	$\text{SE}_1 + \text{SE}_2$	81.0%	76.0%	92.0%
300	100	5%	$\text{SE}_1 + \text{SE}_2 + \text{SE}_3 \times \text{SE}_4$	85.7%	84.0%	92.0%
500	100	5%	$\text{SE}_1 \times \text{SE}_4 + \text{SE}_2 \times \text{SE}_3 + \text{SE}_3$	87.6%	88.6%	92.0%
100	100	20%	$\text{SE}_2 \times \text{SE}_4$	67.0%	67.0%	82.0%
300	100	20%	$\text{SE}_2 \times \text{SE}_3 + \text{SE}_4$	76.0%	73.7%	82.0%
500	100	20%	$\text{SE}_2 + \text{SE}_3 \times \text{SE}_4$	77.0%	79.8%	82.0%
100	1	0%	$\text{SE}_2$	68.0%	67.0%	76.0%
300	1	0%	$\text{SE}_1 + \text{SE}_2 \times \text{SE}_3$	72.3%	70.3%	76.0%
500	1	0%	$\text{SE}_1 + \text{SE}_2 \times \text{SE}_3$	72.2%	73.2%	76.0%
100	1	5%	$\text{SE}_2$	67.0%	58.0%	72.2%
300	1	5%	$\text{SE}_1 \times \text{SE}_2$	71.0%	64.3%	72.2%
500	1	5%	$\text{SE}_1 \times \text{SE}_2 \times \text{SE}_3$	70.6%	68.0%	72.2%
100	1	20%	$\text{SE}_2$	59.0%	61.0%	69.0%
300	1	20%	$\text{SE}_2 \times \text{SE}_3 \times \text{SE}_4$	65.3%	62.3%	69.0%
500	1	20%	$\text{SE}_2 \times \text{SE}_3 \times \text{SE}_4$	64.8%	64.8%	69.0%

Table 5.3: True kernel:  $\text{SE}_1 + \text{SE}_3 \times \text{SE}_7 + \text{SE}_{10}$ ,  $D = 10$ .

Data size	SNR	sp_noise	Kernel chosen	Test accuracy	Kernel rate	Bayes rate
100	100	0%	$\text{SE}_1 \times \text{SE}_9 + \text{SE}_{10}$	61.0%	88.0%	96.0%
300	100	0%	$\text{SE}_1 + \text{SE}_1 \times \text{SE}_{10} + \text{SE}_3 \times \text{SE}_7$	92.0%	92.7%	96.0%
500	100	0%	$\text{SE}_1 + \text{SE}_1 \times \text{SE}_3 \times \text{SE}_7 \times \text{SE}_{10} + \text{SE}_{10}$	94.2%	94.6%	96.0%
100	100	5%	$\text{SE}_1 \times \text{SE}_9 + \text{SE}_{10}$	53.0%	71.0%	91.8%
300	100	5%	$\text{SE}_1 + \text{SE}_3 \times \text{SE}_7 + \text{SE}_6 \times \text{SE}_{10}$	82.0%	81.3%	91.8%
500	100	5%	$\text{SE}_1 \times \text{SE}_3 \times \text{SE}_7 \times \text{SE}_{10} + \text{SE}_{10}$	85.0%	86.2%	91.8%
100	100	20%	$\text{SE}_1$	49.0%	64.0%	79.8%
300	100	20%	$\text{SE}_1 + \text{SE}_{10}$	60.0%	70.0%	79.8%
500	100	20%	$\text{SE}_1 \times \text{SE}_3 \times \text{SE}_7 \times \text{SE}_{10}$	74.2%	75.2%	79.8%
100	1	0%	$\text{SE}_{10}$	59.0%	70.0%	74.4%
300	1	0%	$\text{SE}_1 \times \text{SE}_3 \times \text{SE}_7 \times \text{SE}_{10} + \text{SE}_{10}$	71.3%	72.7%	74.4%
500	1	0%	$\text{SE}_1 \times \text{SE}_{10} + \text{SE}_3 \times \text{SE}_7 + \text{SE}_9$	72.0%	71.4%	74.4%
100	1	5%	$\text{SE}_{10}$	55.0%	66.0%	71.4%
300	1	5%	$\text{SE}_1 \times \text{SE}_{10}$	58.7%	68.7%	71.4%
500	1	5%	$\text{SE}_1 + \text{SE}_{10}$	60.6%	69.4%	71.4%
100	1	20%	$\text{SE}_3$	55.0%	56.0%	65.4%
300	1	20%	$\text{SE}_{10}$	58.0%	61.7%	65.4%
500	1	20%	$\text{SE}_1 \times \text{SE}_{10}$	58.2%	62.0%	65.4%

Table 5.4: True kernel:  $\text{SE}_1 + \text{SE}_3 \times \text{SE}_5 \times \text{SE}_7 + \text{SE}_9$ ,  $D = 10$ .

Data size	SNR	sp_noise	Kernel chosen	Test accuracy	Kernel rate	Bayes rate
100	100	0%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	85.0%	86.0%	97.0%
300	100	0%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7 + \text{SE}_9$	93.7%	93.0%	97.0%
500	100	0%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	91.4%	92.2%	97.0%
100	100	5%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	78.0%	76.0%	91.6%
300	100	5%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	84.0%	83.7%	91.6%
500	100	5%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	86.2%	83.6%	91.6%
100	100	20%	$\text{SE}_8$	49.0%	59.0%	82.0%
300	100	20%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	68.3%	66.0%	82.0%
500	100	20%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	72.2%	66.0%	82.0%
100	1	0%	$\text{SE}_1 \times \text{SE}_3 \times \text{SE}_4 \times \text{SE}_5 + \text{SE}_7$	59.0%	66.0%	74.2%
300	1	0%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7 + \text{SE}_9$	71.7%	72.7%	74.2%
500	1	0%	$\text{SE}_1 + \text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	73.0%	70.6%	74.2%
100	1	5%	$\text{SE}_1 \times \text{SE}_3 \times \text{SE}_4 \times \text{SE}_5 + \text{SE}_7$	55.0%	62.0%	70.8%
300	1	5%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	64.3%	68.7%	70.8%
500	1	5%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_7$	70.4%	67.4%	70.8%
100	1	20%	$\text{SE}_3 \times \text{SE}_5 \times \text{SE}_9$	52.0%	64.0%	66.4%
300	1	20%	$\text{SE}_3 \times \text{SE}_7 \times \text{SE}_8$	55.7%	61.7%	66.4%
500	1	20%	$\text{SE}_3 \times \text{SE}_7$	56.4%	62.6%	66.4%

## 5.2 Experiments on Real World Data Sets

In this section, we compare the performance of models constructed using our algorithm with related methods and show that the performance of our structurally simpler models is on par with more complicated models such as additive GPs [1] and Hierarchical Kernel Learning. We also compare the performance of structure search using different information criteria (BIC, AIC, BIClight), as well as the search guided by cross-validated test accuracy. We also show the performance of the kernel using a likelihood mixture to account for outliers.

The table below contains the mean classification error across 10 train-test splits between different methods. The best performing model is shown in bold, together with all other models that were not significantly different from it, according to the paired t-test for statistical significance. In addition to the structure search, we show the performance of the random forest method, which constructs 1000 decision trees using the training data and then uses the mode of the classifications produced by these trees to label the test set data. This method was intended to be a *ceiling* performance for our methods, as its focus is just predictive performance: it does not contribute to interpretability or our understanding of the data set considered.

Table 5.5: Classification Percent Error

Method	breast	pima	liver	heart
Logistic Regression	7.611	24.392	45.060	<b>16.082</b>
GP GAM	<b>5.189</b>	<b>22.419</b>	<b>29.842</b>	<b>16.839</b>
HKL	<b>5.377</b>	24.261	<b>27.270</b>	<b>18.975</b>
GP Squared-exp	<b>4.734</b>	<b>23.722</b>	<b>31.237</b>	<b>20.642</b>
GP Additive	<b>5.566</b>	<b>23.076</b>	<b>30.060</b>	<b>18.496</b>
GPSS (AIC)	6.430	<b>22.529</b>	28.924	19.860
GPSS (BIC)	5.980	23.440	37.010	<b>18.150</b>
GPSS (BIC light)	6.430	<b>22.270</b>	<b>27.500</b>	<b>17.820</b>
GPSS (likMix)	<b>11.240</b>	<b>23.180</b>	<b>28.370</b>	<b>16.460</b>
GPSS (crossValGuide)	<b>5.090</b>	23.700	-	<b>17.160</b>
Random Forest	<b>4.220</b>	<b>23.440</b>	<b>24.030</b>	<b>17.130</b>

## Chapter 6

# Summary and Conclusions

### 6.1 Further Work



# Bibliography

- [1] D. Duvenaud, H. Nickisch, and C.E. Rasmussen. Additive Gaussian processes. In *Advances in Neural Information Processing Systems*, 2011.