

# Automatic Construction and Natural-Language Description of Nonparametric Regression Models

**James Robert Lloyd**  
Department of Engineering  
University of Cambridge

**David Duvenaud**  
Department of Engineering  
University of Cambridge

**Roger Grosse**  
Brain and Cognitive Sciences  
Massachusetts Institute of Technology

**Joshua B. Tenenbaum**  
Brain and Cognitive Sciences  
Massachusetts Institute of Technology

**Zoubin Ghahramani**  
Department of Engineering  
University of Cambridge

## Abstract

This paper presents the beginnings of an automatic statistician, focusing on regression problems. Our system explores an open-ended space of possible statistical models to discover a good explanation of the data, and then produces a detailed report with figures and natural-language text.

Our approach treats unknown functions nonparametrically using Gaussian processes, which has two important consequences. First, Gaussian processes model functions in terms of high-level properties (e.g. smoothness, trends, periodicity, changepoints). Taken together with the compositional structure of our language of models, this allows us to automatically describe functions through a decomposition into additive parts. Second, the use of flexible nonparametric models and a rich language for composing them in an open-ended manner also results in state-of-the-art extrapolation performance evaluated over 13 real time series data sets from various domains.

## 1 Introduction

Automating the process of statistical modeling would have a tremendous impact on fields that currently rely on expert statisticians, machine learning researchers, and data scientists. While fitting simple models (such as linear regression) is largely automated by standard software packages, there has been little work on the automatic construction of flexible but interpretable models. What are the ingredients required for an artificial intelligence system to be able to perform statistical modeling automatically? In this paper we conjecture that the following ingredients may be useful for building an AI system for statistics, and we develop a working system which incorporates them:

- **An open-ended language of models** expressive enough to capture many of the modeling assumptions and model composition techniques applied by human statisticians to capture real-world phenomena
- **A search procedure** to efficiently explore the space of models spanned by the language
- **A principled method for evaluating models** in terms of their complexity and their degree of fit to the data
- **A procedure for automatically generating reports** which explain and visualize different factors underlying the data, make the chosen modeling assumptions explicit,

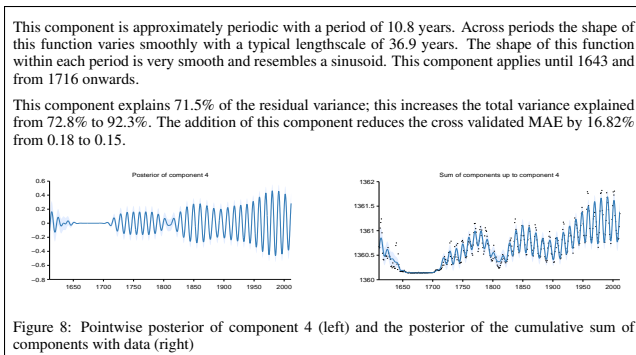


Figure 8: Pointwise posterior of component 4 (left) and the posterior of the cumulative sum of components with data (right)

Figure 1: Extract from an automatically-generated report describing the model components discovered by automatic model search. This part of the report isolates and describes the approximately 11-year sunspot cycle, also noting its disappearance during the 16th century, a time known as the Maunder minimum ([Lean, Beer, and Bradley, 1995](#)).

and quantify how each component improves the predictive power of the model

In this paper, we introduce a system for modeling time-series data containing the above ingredients which we call the Automatic Bayesian Covariance Discovery (ABCD) system. The system defines an open-ended language of Gaussian process models via a compositional grammar. The space is searched greedily, using marginal likelihood and the Bayesian Information Criterion (BIC) to evaluate models. The compositional structure of the language allows us to develop a method for automatically translating components of the model into natural-language descriptions of patterns in the data.

We show examples of automatically generated reports which highlight interpretable features discovered in a variety of data sets (e.g. figure 1). The supplementary material to this paper includes 13 complete reports automatically generated by ABCD<sup>1</sup>.

Good statistical modeling requires not only interpretability but predictive accuracy. We compare ABCD against ex-

<sup>1</sup>A link to these reports will be maintained at <http://mlg.eng.cam.ac.uk/lloyd/>

isting model construction techniques in terms of predictive performance at extrapolation, and we find state-of-the-art performance on 13 time series. In the remainder of this paper we describe the components of ABCD in detail.

## 2 A language of regression models

The general problem of regression consists of learning a function  $f$  mapping from some input space  $\mathcal{X}$  to some output space  $\mathcal{Y}$ . We would like an expressive language which can represent both simple parametric forms of  $f$  such as linear, polynomial, etc. and also complex nonparametric functions specified in terms of properties such as smoothness, periodicity, etc. Fortunately, Gaussian processes (GPs) provide a very general and analytically tractable way of capturing both simple and complex functions.

Gaussian processes are distributions over functions such that any finite subset of function evaluations,  $(f(x_1), f(x_2), \dots, f(x_N))$ , have a joint Gaussian distribution (Rasmussen and Williams, 2006). A GP is completely specified by its mean function,  $\mu(x) = \mathbb{E}(f(x))$  and kernel (or covariance) function  $k(x, x') = \text{Cov}(f(x), f(x'))$ . It is common practice to assume zero mean, since marginalizing over an unknown mean function can be equivalently expressed as a zero-mean GP with a new kernel. The structure of the kernel captures high-level properties of the unknown function,  $f$ , which in turn determines how the model generalizes or extrapolates to new data. We can therefore define a language of regression models by specifying a language of kernels.

The elements of this language are a set of base kernels capturing different function properties, and a set of composition rules which combine kernels to yield other valid kernels. Our base kernels are white noise (WN), constant (C), linear (LIN), squared exponential (SE) and periodic (PER), which on their own encode for uncorrelated noise, constant functions, linear functions, smooth functions and periodic functions respectively<sup>2</sup>. The composition rules are addition and multiplication:

$$k_1 + k_2 = k_1(x, x') + k_2(x, x') \quad (2.1)$$

$$k_1 \times k_2 = k_1(x, x') \times k_2(x, x') \quad (2.2)$$

Combining kernels using these operations can yield kernels encoding for richer structures such as approximate periodicity ( $\text{SE} \times \text{PER}$ ) or smooth functions with linear trends ( $\text{SE} + \text{LIN}$ ).

This kernel composition framework (with different base kernels) was described by Duvenaud et al. (2013). We extend and adapt this framework in several ways. In particular, we have found that incorporating changepoints into the language is essential for realistic models of time series (e.g. figure 1). Changepoints can be defined through addition and multiplication with sigmoidal functions:

$$\text{CP}(k_1, k_2) = k_1 \times \sigma + k_2 \times \bar{\sigma} \quad (2.3)$$

where  $\sigma = \sigma(x)\sigma(x')$  and  $\bar{\sigma} = (1 - \sigma(x))(1 - \sigma(x'))$ . Changewindows  $\text{CW}(\cdot, \cdot)$  can be defined similarly by replacing  $\sigma(x)$  with a product of two sigmoids.

<sup>2</sup>Definitions of kernels are in the supplementary material.

We also expanded and reparametrised the set of base kernels so that they were more amenable to automatic description and to extend the number of common regression models included in the language. Table 1 lists common regression models that can be expressed by our language.

Regression model	Kernel
GP smoothing	SE + WN
Linear regression	C + LIN + WN
Multiple kernel learning	$\sum \text{SE} + \text{WN}$
Trend, cyclical, irregular	$\sum \text{SE} + \sum \text{PER} + \text{WN}$
Fourier decomposition	$\text{C} + \sum \cos + \text{WN}$
Sparse spectrum GPs	$\sum \cos + \text{WN}$
Spectral mixture	$\sum \text{SE} \times \cos + \text{WN}$
Changepoints	e.g. $\text{CP}(\text{SE}, \text{SE}) + \text{WN}$
Heteroscedasticity	e.g. $\text{SE} + \text{LIN} \times \text{WN}$

Table 1: Common regression models expressible in our language.  $\cos$  is a special case of our reparametrised PER.

## 3 Model Search and Evaluation

As in Duvenaud et al. (2013) we explore the space of regression models using a greedy search. We use the same search operators, but also include additional operators to incorporate changepoints; a complete list is contained in the supplementary material.

After each model is proposed its kernel parameters are optimised by conjugate gradient descent. We evaluate each optimized model,  $M$ , using the Bayesian Information Criterion (BIC) (Schwarz, 1978):

$$\text{BIC}(M) = -2 \log p(D | M) + p \log n \quad (3.1)$$

where  $p$  is the number of kernel parameters,  $\log p(D|M)$  is the marginal likelihood of the data,  $D$ , and  $n$  is the number of data points. BIC trades off model fit and complexity and implements what is known as ‘‘Bayesian Occam’s Razor’’ (e.g. Rasmussen and Ghahramani, 2001; MacKay, 2003).

## 4 Automatic description of regression models

**Overview** In this section, we describe how ABCD generates natural-language descriptions of the models found by the search procedure. There are two main features of our language of GP models that allow description to be performed automatically.

First, the sometimes complicated kernel expressions found can be simplified into a sum of products. A sum of kernels corresponds to a sum of functions so each product can be described separately. Second, each kernel in a product modifies the resulting model in a consistent way. Therefore, we can choose one kernel to be described as a noun, with all others described using adjectives.

**Sum of products normal form** We convert each kernel expression into a standard, simplified form. We do this by first distributing all products of sums into a sum of products.

Next, we apply several simplifications to the kernel expression: The product of two SE kernels is another SE with different parameters. Multiplying WN by any stationary kernel (C, WN, SE, or PER) gives another WN kernel. Multiplying any kernel by C only changes the parameters of the original kernel.

After applying these rules, the kernel can as be written as a sum of terms of the form:

$$K \prod_m \text{LIN}^{(m)} \prod_n \sigma^{(n)}, \quad (4.1)$$

where  $K$  is one of WN, C, SE,  $\prod_k \text{PER}^{(k)}$  or  $\text{SE} \prod_k \text{PER}^{(k)}$  and  $\prod_i k^{(i)}$  denotes a product of kernels, each with different parameters.

**Sums of kernels are sums of functions** Formally, if  $f_1(x) \sim \text{GP}(0, k_1)$  and independently  $f_2(x) \sim \text{GP}(0, k_2)$  then  $f_1(x) + f_2(x) \sim \text{GP}(0, k_1 + k_2)$ . This lets us describe each product of kernels separately.

**Each kernel in a product modifies a model in a consistent way** This allows us to describe the contribution of each kernel in a product as an adjective, or more generally as a modifier of a noun. We now describe how each kernel modifies a model and how this can be described in natural language:

- **Multiplication by SE** removes long range correlations from a model since  $\text{SE}(x, x')$  decreases monotonically to 0 as  $|x - x'|$  increases. This can be described as making an existing model’s correlation structure ‘local’ or ‘approximate’.
- **Multiplication by LIN** is equivalent to multiplying the function being modeled by a linear function. If  $f(x) \sim \text{GP}(0, k)$ , then  $xf(x) \sim \text{GP}(0, k \times \text{LIN})$ . This causes the standard deviation of the model to vary linearly without affecting the correlation and can be described as e.g. ‘with linearly increasing standard deviation’.
- **Multiplication by  $\sigma$**  is equivalent to multiplying the function being modeled by a sigmoid which means that the function goes to zero before or after some point. This can be described as e.g. ‘from [time]’ or ‘until [time]’.
- **Multiplication by PER** modifies the correlation structure in the same way as multiplying the function by an independent periodic function. Formally, if  $f_1(x) \sim \text{GP}(0, k_1)$  and  $f_2(x) \sim \text{GP}(0, k_2)$  then

$$\text{Cov}[f_1(x)f_2(x), f_1(x')f_2(x')] = k_1(x, x')k_2(x, x').$$

This can be loosely described as e.g. ‘modulated by a periodic function with a period of [period] [units]’.

**Constructing a complete description of a product of kernels** We choose one kernel to act as a noun which is then described by the functions it encodes for when unmodified e.g. ‘smooth function’ for SE. Modifiers corresponding to the other kernels in the product are then appended to this description, forming a noun phrase of the form:

Determiner + Premodifiers + Noun + Postmodifiers

As an example, a kernel of the form  $\text{SE} \times \text{PER} \times \text{LIN} \times \sigma$  could be described as an

$$\underbrace{\text{SE}}_{\text{approximately}} \times \underbrace{\text{PER}}_{\text{periodic function}} \times \underbrace{\text{LIN}}_{\text{with linearly growing amplitude}} \times \underbrace{\sigma}_{\text{until 1700.}}$$

where PER has been selected as the head noun.

In principle, any assignment of kernels in a product to these different phrasal roles is possible, but in practice we found certain assignments to produce more interpretable phrases than others. The head noun is chosen according to the following ordering:

$$\text{PER} > \text{WN}, \text{SE}, \text{C} > \prod_m \text{LIN}^{(m)} > \prod_n \sigma^{(n)}$$

i.e. PER is always chosen as the head noun when present.

**Ordering additive components** The reports generated by ABCD attempt to present the most interesting or important features of a data set first. As a heuristic, we order components by always adding next the component which most reduces the 10-fold cross-validated mean absolute error.

#### 4.1 Worked example

Suppose we start with a kernel of the form

$$\text{SE} \times (\text{WN} \times \text{LIN} + \text{CP}(\text{C}, \text{PER})).$$

This is converted to a sum of products:

$$\text{SE} \times \text{WN} \times \text{LIN} + \text{SE} \times \text{C} \times \sigma + \text{SE} \times \text{PER} \times \bar{\sigma}.$$

which is simplified to

$$\text{WN} \times \text{LIN} + \text{SE} \times \sigma + \text{SE} \times \text{PER} \times \bar{\sigma}.$$

To describe the first component, the head noun description for WN, ‘uncorrelated noise’, is concatenated with a modifier for LIN, ‘with linearly increasing standard deviation’. The second component is described as ‘A smooth function with a lengthscale of [lengthscale] [units]’, corresponding to the SE, ‘which applies until [changepoint]’, which corresponds to the  $\sigma$ . Finally, the third component is described as ‘An approximately periodic function with a period of [period] [units] which applies from [changepoint]’.

### 5 Example descriptions of time series

We demonstrate the ability of our procedure to discover and describe a variety of patterns on two time series. Full automatically-generated reports for 13 data sets are provided as supplementary material.

#### 5.1 Summarizing 400 Years of Solar Activity

We show excerpts from the report automatically generated on annual solar irradiation data from 1610 to 2011 (figure 2). This time series has two pertinent features: a roughly 11-year cycle of solar activity, and a period lasting from 1645 to 1715 with much smaller variance than the rest of the dataset. This flat region corresponds to the Maunder minimum, a period in which sunspots were extremely rare (Lean, Beer, and

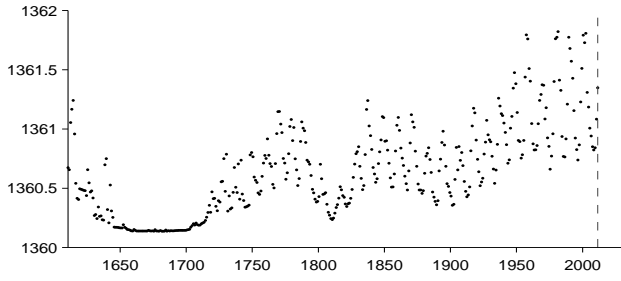


Figure 2: Solar irradiance data.

The structure search algorithm has identified eight additive components in the data. The first 4 additive components explain 92.3% of the variation in the data as shown by the coefficient of determination ( $R^2$ ) values in table 1. The first 6 additive components explain 99.7% of the variation in the data. After the first 5 components the cross validated mean absolute error (MAE) does not decrease by more than 0.1%. This suggests that subsequent terms are modelling very short term trends, uncorrelated noise or are artefacts of the model or search procedure. Short summaries of the additive components are as follows:

- A constant.
- A constant. This function applies from 1643 until 1716.
- A smooth function. This function applies until 1643 and from 1716 onwards.
- An approximately periodic function with a period of 10.8 years. This function applies until 1643 and from 1716 onwards.

Figure 3: Automatically generated descriptions of the components discovered by ABCD on the solar irradiance data set. The dataset has been decomposed into diverse structures with simple descriptions.

Bradley, 1995). ABCD clearly identifies these two features, as discussed below.

Figure 3 shows the natural-language summaries of the top four components chosen by ABCD. From these short summaries, we can see that our system has identified the Maunder minimum (second component) and 11-year solar cycle (fourth component). These components are visualized in figures 4 and 1, respectively. The third component corresponds to long-term trends, as visualized in figure 5.

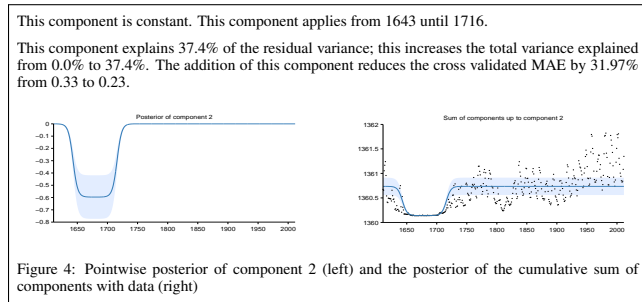


Figure 4: One of the learned components corresponds to the Maunder minimum.

## 5.2 Finding heteroscedasticity in air traffic data

Next, we present the analysis generated by our procedure on international airline passenger data (figure 6). The model constructed by ABCD has four components:  $LIN + SE \times PER \times LIN + SE + WN \times LIN$ , with descriptions given in

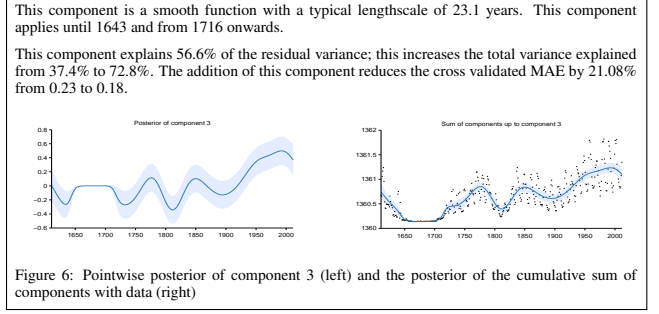


Figure 5: Characterizing the medium-term smoothness of solar activity levels. By allowing other components to explain the periodicity, noise, and the Maunder minimum, ABCD can isolate the part of the signal best explained by a slowly-varying trend.

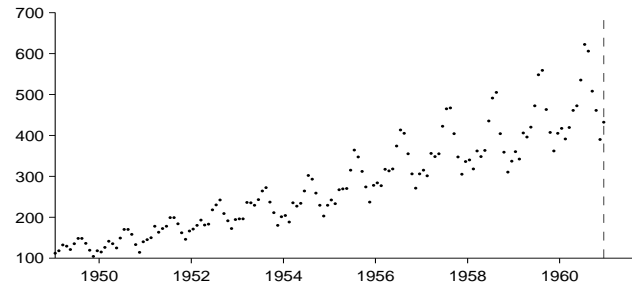


Figure 6: International airline passenger monthly volume (e.g. Box, Jenkins, and Reinsel, 2013).

figure 7.

The second component (figure 8) is accurately described as approximately (SE) periodic (PER) with linearly growing amplitude (LIN). By multiplying a white noise kernel by a linear kernel, the model is able to express heteroscedasticity (figure 9).

## 5.3 Comparison to equation learning

We now compare the descriptions generated by ABCD to parametric functions produced by an equation learning system. We show equations produced by Eureka (Nuttonian, 2011) for the data sets shown above, using the default mean absolute error performance metric.

The learned function for the solar irradiance data is

$$\text{Irradiance}(t) = 1361 + \alpha \sin(\beta + \gamma t) \sin(\delta + \epsilon t^2 - \zeta t)$$

where  $t$  is time and constants are replaced with symbols for brevity. This equation captures the constant offset of the data, and models the long-term trend with a product of sinusoids, but fails to capture the solar cycle or the Maunder minimum.

The learned function for the airline passenger data is

$$\text{Passengers}(t) = \alpha t + \beta \cos(\gamma - \delta t) \log(\epsilon t - \zeta) - \eta$$

which captures the approximately linear trend, and the periodic component with approximately linearly (logistic) in-



The structure search algorithm has identified four additive components in the data. The first 2 additive components explain 98.5% of the variation in the data as shown by the coefficient of determination ( $R^2$ ) values in table 1. The first 3 additive components explain 99.8% of the variation in the data. After the first 3 components the cross validated mean absolute error (MAE) does not decrease by more than 0.1%. This suggests that subsequent terms are modelling very short term trends, uncorrelated noise or are artefacts of the model or search procedure. Short summaries of the additive components are as follows:

- A linearly increasing function.
- An approximately periodic function with a period of 1.0 years and with linearly increasing amplitude.
- A smooth function.
- Uncorrelated noise with linearly increasing standard deviation.

#	$R^2$ (%)	$\Delta R^2$ (%)	Residual $R^2$ (%)	Cross validated MAE	Reduction in MAE (%)
-	-	-	-	280.30	-
1	85.4	85.4	85.4	34.03	87.9
2	98.5	13.2	89.9	12.44	63.4
3	99.8	1.3	85.1	9.10	26.8
4	100.0	0.2	100.0	9.10	0.0

Figure 7: Short descriptions and summary statistics for the four components of the airline model.

## 2.2 Component 2 : An approximately periodic function with a period of 1.0 years and with linearly increasing amplitude

This component is approximately periodic with a period of 1.0 years and varying amplitude. Across periods the shape of this function varies very smoothly. The amplitude of the function increases linearly. The shape of this function within each period has a typical lengthscale of 6.0 weeks.

This component explains 89.9% of the residual variance; this increases the total variance explained from 85.4% to 98.5%. The addition of this component reduces the cross validated MAE by 63.45% from 34.03 to 12.44.

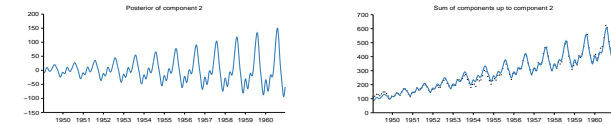


Figure 4: Pointwise posterior of component 2 (left) and the posterior of the cumulative sum of components with data (right)

Figure 8: Capturing non-stationary periodicity in the airline data

creasing amplitude. However, the annual cycle is heavily approximated by a sinusoid and the model does not capture heteroscedasticity.

## 6 Related work

**Building Kernel Functions** Rasmussen and Williams (2006) devote 4 pages to manually constructing a composite kernel to model a time series of carbon dioxide concentrations. In the supplementary material, we include a report automatically generated by ABCD for this dataset; our procedure chose a model similar to the one they constructed by hand. Other examples of papers whose main contribution is to manually construct and fit a composite GP kernel are Klenske (2012) and Lloyd (2013).

Diosan, Rogozan, and Pecuchet (2007); Bing et al. (2010) and Kronberger and Kommenda (2013) search over a similar space of models as ABCD using genetic algorithms but do not interpret the resulting models. Our procedure is based on the model construction method of Duvenaud et al. (2013) which automatically decomposed models but components were interpreted manually and the space of models searched over was smaller than that in this work.

## 2.4 Component 4 : Uncorrelated noise with linearly increasing standard deviation

This component models uncorrelated noise. The standard deviation of the noise increases linearly.

This component explains 100.0% of the residual variance; this increases the total variance explained from 99.8% to 100.0%. The addition of this component reduces the cross validated MAE by 0.00% from 9.10 to 9.10. This component explains residual variance but does not improve MAE which suggests that this component describes very short term patterns, uncorrelated noise or is an artefact of the model or search procedure.

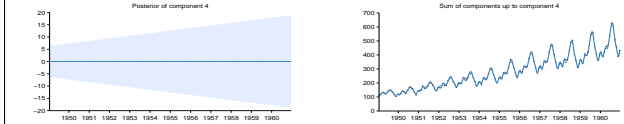


Figure 8: Pointwise posterior of component 4 (left) and the posterior of the cumulative sum of components with data (right)

Figure 9: Modeling heteroscedasticity

**Kernel Learning** Sparse spectrum GPs (Lázaro-Gredilla et al., 2010) approximate the spectral density of a stationary kernel function using delta functions which corresponds to kernels of the form  $\sum \cos$ . Similarly, Wilson and Adams (2013) introduce spectral mixture kernels which approximate the spectral density using a scale-location mixture of Gaussian distributions corresponding to kernels of the form  $\sum SE \times \cos$ . Both demonstrate, using Bochner’s theorem (Bochner, 1959), that these kernels can approximate any stationary covariance function. Our language of kernels includes both of these kernel classes (see table 1).

There is a large body of work attempting to construct rich kernels through a weighted sum of base kernels called multiple kernel learning (MKL) (e.g. Bach, Lanckriet, and Jordan, 2004). These approaches find the optimal solution in polynomial time but only if the component kernels and parameters are pre-specified. We compare to a Bayesian variant of MKL in section 7 which is expressed as a restriction of our language of kernels.

**Equation learning** Todorovski and Dzeroski (1997), Washio et al. (1999) and Schmidt and Lipson (2009) learn parametric forms of functions specifying time series, or relations between quantities. In contrast, ABCD learns a parametric form for the covariance, allowing it to model functions without a simple parametric form.

**Searching over open-ended model spaces** This work was inspired by previous successes at searching over open-ended model spaces: matrix decompositions (Grosse, Salakhutdinov, and Tenenbaum, 2012) and graph structures (Kemp and Tenenbaum, 2008). In both cases, the model spaces were defined compositionally through a handful of components and operators, and models were selected using criteria which trade off model complexity and goodness of fit. Our work differs in that our procedure automatically interprets the chosen model, making the results accessible to non-experts.

**Natural-language output** To the best of our knowledge, our procedure is the first example of automatic description of nonparametric statistical models. However, systems with

natural language output have been built in the areas of video interpretation (Barbu et al., 2012) and automated theorem proving (Ganesalingam and Gowers, 2013).

## 7 Predictive Accuracy

In addition to our demonstration of the interpretability of ABCD, we compared the predictive accuracy of various model-building algorithms at interpolating and extrapolating time-series. ABCD outperforms the other methods on average.

**Data sets** We evaluate the performance of the algorithms listed below on 13 real time-series from various domains from the time series data library (Hyndman, Accessed summer 2013); plots of the data can be found at the beginning of the reports in the supplementary material.

**Algorithms** We compare ABCD to equation learning using Eureqa (Nuttonian, 2011) and six other regression algorithms: linear regression, GP regression with a single SE kernel (squared exponential), a Bayesian variant of multiple kernel learning (MKL) (e.g. Bach, Lanckriet, and Jordan, 2004), change point modeling (e.g. Garnett et al., 2010; Saatçi, Turner, and Rasmussen, 2010; Fox and Dunson, 2013), spectral mixture kernels (Wilson and Adams, 2013) (spectral kernels) and trend-cyclical-irregular models (e.g. Lind et al., 2006).

We use the default mean absolute error criterion when using Eureqa. All other algorithms can be expressed as restrictions of our modeling language (see table 1) so we perform inference using the same search methodology and selection criterion<sup>3</sup> with appropriate restrictions to the language. For MKL, trend-cyclical-irregular and spectral kernels, the greedy search procedure of ABCD corresponds to a forward-selection algorithm. For squared exponential and linear regression the procedure corresponds to marginal likelihood optimisation. More advanced inference methods are typically used for changepoint modeling but we use the same inference method for all algorithms for comparability.

We restricted to regression algorithms for comparability; this excludes models which regress on previous values of times series, such as autoregressive or moving-average models (e.g. Box, Jenkins, and Reinsel, 2013). Constructing a language for this class of time-series model would be an interesting area for future research.

**Interpretability versus accuracy** BIC trades off model fit and complexity by penalizing the number of parameters in a kernel expression. This can result in ABCD favoring kernel expressions with nested products of sums, producing descriptions involving many additive components. While these models have good predictive performance the large number of components can make them less interpretable. We experimented with distributing all products over addition during the search, causing models with many additive components

<sup>3</sup>We experimented with using unpenalised marginal likelihood as the search criterion but observed overfitting, as is to be expected.

to be more heavily penalized by BIC. We call this procedure ABCD-interpretability, in contrast to the unrestricted version of the search, ABCD-accuracy.

**Extrapolation** To test extrapolation we trained all algorithms on the first 90% of the data, predicted the remaining 10% and then computed the root mean squared error (RMSE). The RMSEs are then standardised by dividing by the smallest RMSE for each data set so that the best performance on each data set will have a value of 1.

Figure 10 shows the standardised RMSEs across algorithms. ABCD-accuracy outperforms ABCD-interpretability but both versions have lower quartiles than all other methods.

Overall, the model construction methods with greater capacity perform better: ABCD outperforms trend-cyclical-irregular, which outperforms Bayesian MKL, which outperforms squared exponential. Despite searching over a rich model class, Eureqa performs relatively poorly, since very few datasets are parsimoniously explained by a parametric equation.

Not shown on the plot are large outliers for spectral kernels, Eureqa, squared exponential and linear regression with values of 11, 493, 22 and 29 respectively. All of these outliers occurred on a data set with a large discontinuity (see the call centre data in the supplementary material).

**Interpolation** To test the ability of the methods to interpolate, we randomly divided each data set into equal amounts of training data and testing data. The results are similar to those for extrapolation and are included in the supplementary material.

## 8 Conclusion

Towards the goal of automating statistical modeling we have presented a system which constructs an appropriate model from an open-ended language and automatically generates detailed reports that describe patterns in the data captured by the model. We have demonstrated that our procedure can discover and describe a variety of patterns on several time series. Our procedure’s extrapolation and interpolation performance on time-series are state-of-the-art compared to existing model construction techniques. We believe this procedure has the potential to make powerful statistical model-building techniques accessible to non-experts.

**Source Code** Source code to perform all experiments is available on github<sup>4</sup>.

## References

Bach, F. R.; Lanckriet, G. R.; and Jordan, M. I. 2004. Multiple kernel learning, conic duality, and the SMO algo-

<sup>4</sup><http://www.github.com/jamesrobertlloyd/gpss-research>. All GP parameter optimisation was performed by automated calls to the GPML toolbox available at <http://www.gaussianprocess.org/gpml/code/>.

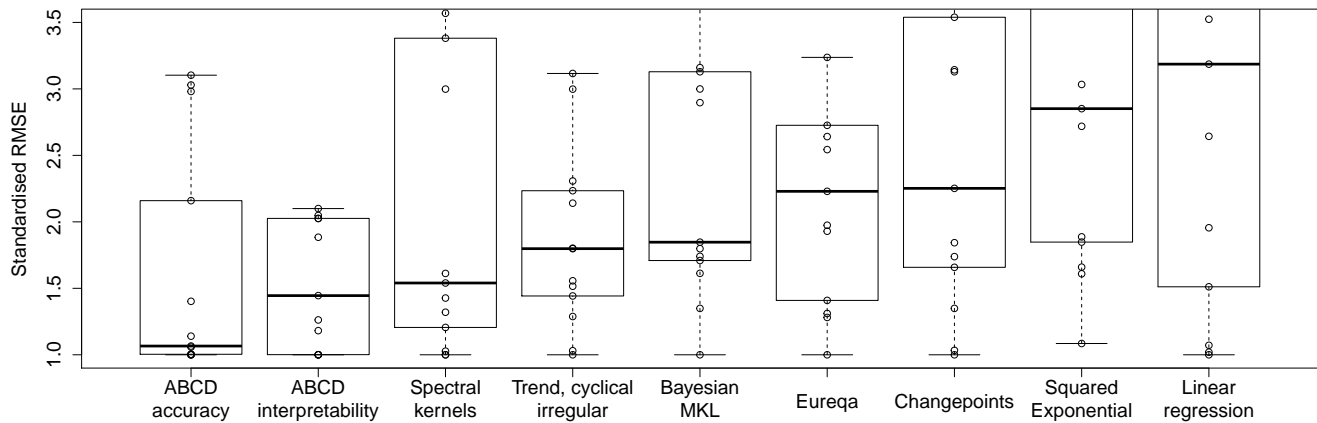


Figure 10: Raw data, and box plot (showing median and quartiles) of standardised extrapolation RMSE (best performance = 1) on 13 time-series. The methods are ordered by median.

- rithm. In *Proceedings of the twenty-first international conference on Machine learning*, 6. ACM.
- Barbu, A.; Bridge, A.; Burchill, Z.; Coroian, D.; Dickinson, S.; Fidler, S.; Michaux, A.; Mussman, S.; Narayanaswamy, S.; Salvi, D.; Schmidt, L.; Shangquan, J.; Siskind, J.; Waggoner, J.; Wang, S.; Wei, J.; Yin, Y.; and Zhang, Z. 2012. Video in sentences out. In *Conference on Uncertainty in Artificial Intelligence*.
- Bing, W.; Wen-qiong, Z.; Ling, C.; and Jia-hong, L. 2010. A GP-based kernel construction and optimization method for RVM. In *International Conference on Computer and Automation Engineering (ICCAE)*, volume 4, 419–423.
- Bochner, S. 1959. *Lectures on Fourier integrals*, volume 42. Princeton University Press.
- Box, G. E.; Jenkins, G. M.; and Reinsel, G. C. 2013. *Time series analysis: forecasting and control*. Wiley. com.
- Diosan, L.; Rogozan, A.; and Pecuchet, J. 2007. Evolving kernel functions for SVMs by genetic programming. In *Machine Learning and Applications, 2007*, 19–24. IEEE.
- Duvenaud, D.; Lloyd, J. R.; Grosse, R.; Tenenbaum, J. B.; and Ghahramani, Z. 2013. Structure discovery in nonparametric regression through compositional kernel search. In *Proceedings of the 30th International Conference on Machine Learning*.
- Fox, E., and Dunson, D. 2013. Multiresolution Gaussian Processes. In *Neural Information Processing Systems 25*. MIT Press.
- Ganesalingam, M., and Gowers, W. T. 2013. A fully automatic problem solver with human-style output. *CoRR* abs/1309.4501.
- Garnett, R.; Osborne, M. A.; Reece, S.; Rogers, A.; and Roberts, S. J. 2010. Sequential bayesian prediction in the presence of changepoints and faults. *The Computer Journal* 53(9):1430–1446.
- Grosse, R.; Salakhutdinov, R.; and Tenenbaum, J. 2012. Exploiting compositionality to explore a large space of model structures. In *Uncertainty in Artificial Intelligence*.
- Hyndman, R. J. Accessed summer 2013. Time series data library.
- Kemp, C., and Tenenbaum, J. 2008. The discovery of structural form. *Proceedings of the National Academy of Sciences* 105(31):10687–10692.
- Klenske, E. 2012. *Nonparametric System Identification and Control for Periodic Error Correction in Telescopes*. Ph.D. Dissertation, University of Stuttgart.
- Kronberger, G., and Kommenda, M. 2013. Evolution of covariance functions for gaussian process regression using genetic programming. *arXiv preprint arXiv:1305.3794*.
- Lázaro-Gredilla, M.; Quiñonero-Candela, J.; Rasmussen, C. E.; and Figueiras-Vidal, A. R. 2010. Sparse spectrum gaussian process regression. *The Journal of Machine Learning Research* 99:1865–1881.
- Lean, J.; Beer, J.; and Bradley, R. 1995. Reconstruction of solar irradiance since 1610: Implications for climate change. *Geophysical Research Letters* 22(23):3195–3198.
- Lind, D. A.; Marchal, W. G.; Wathen, S. A.; and Magazine, B. W. 2006. *Basic statistics for business and economics*. McGraw-Hill/Irwin Boston.
- Lloyd, J. R. 2013. GEFCom2012 hierarchical load forecasting: Gradient boosting machines and gaussian processes. *International Journal of Forecasting*.
- MacKay, D. J. 2003. *Information theory, inference and learning algorithms*. Cambridge university press.
- Nutonian. 2011. Eureka.
- Rasmussen, C., and Ghahramani, Z. 2001. Occam’s razor. In *Advances in Neural Information Processing Systems*.

- Rasmussen, C., and Williams, C. 2006. *Gaussian Processes for Machine Learning*. The MIT Press, Cambridge, MA, USA.
- Saatçi, Y.; Turner, R. D.; and Rasmussen, C. E. 2010. Gaussian process change point models. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, 927–934.
- Schmidt, M., and Lipson, H. 2009. Distilling free-form natural laws from experimental data. *Science* 324(5923):81–85.
- Schwarz, G. 1978. Estimating the dimension of a model. *The Annals of Statistics* 6(2):461–464.
- Todorovski, L., and Dzeroski, S. 1997. Declarative bias in equation discovery. In *International Conference on Machine Learning*, 376–384.
- Washio, T.; Motoda, H.; Niwa, Y.; et al. 1999. Discovering admissible model equations from observed data based on scale-types and identity constraints. In *International Joint Conference On Artificial Intelligence*, volume 16, 772–779.
- Wilson, A. G., and Adams, R. P. 2013. Gaussian process covariance kernels for pattern discovery and extrapolation. In *Proceedings of the 30th International Conference on Machine Learning*.

## Appendices

### A Kernels

#### A.1 Base kernels

For scalar-valued inputs, the white noise (WN), constant (C), linear (LIN), squared exponential (SE), and periodic kernels (PER) are defined as follows:

$$\text{WN}(x, x') = \sigma^2 \delta_{x, x'} \quad (\text{A.1})$$

$$\text{C}(x, x') = \sigma^2 \quad (\text{A.2})$$

$$\text{LIN}(x, x') = \sigma^2 (x - \ell)(x' - \ell) \quad (\text{A.3})$$

$$\text{SE}(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right) \quad (\text{A.4})$$

$$\text{PER}(x, x') = \sigma^2 \frac{\exp\left(\frac{\cos\left(\frac{2\pi(x-x')}{\ell^2}\right)}{\ell^2}\right) - I_0\left(\frac{1}{\ell^2}\right)}{\exp\left(\frac{1}{\ell^2}\right) - I_0\left(\frac{1}{\ell^2}\right)} \quad (\text{A.5})$$

where  $\delta_{x, x'}$  is the Kronecker delta function,  $I_0$  is the modified Bessel function of the first kind of order zero and other symbols are parameters of the kernel functions.

#### A.2 Changepoints and changewindows

The changepoint,  $\text{CP}(\cdot, \cdot)$  operator is defined as follows:

$$\begin{aligned} \text{CP}(k_1, k_2)(x, x') = & \sigma(x)k_1(x, x')\sigma(x') \\ & + (1 - \sigma(x))k_2(x, x')(1 - \sigma(x')) \end{aligned} \quad (\text{A.6})$$

where  $\sigma(x) = 0.5 \times (1 + \tanh(\frac{\ell-x}{s}))$ . This can also be written as

$$\text{CP}(k_1, k_2) = \sigma k_1 + \bar{\sigma} k_2 \quad (\text{A.7})$$

where  $\sigma(x, x') = \sigma(x)\sigma(x')$  and  $\bar{\sigma}(x, x') = (1 - \sigma(x))(1 - \sigma(x'))$ .

Changewindow,  $\text{CW}(\cdot, \cdot)$ , operators are defined similarly by replacing the sigmoid,  $\sigma(x)$ , with a product of two sigmoids.

#### A.3 Properties of the periodic kernel

A simple application of l'Hôpital's rule shows that

$$\text{PER}(x, x') \rightarrow \sigma^2 \cos\left(\frac{2\pi(x-x')}{p}\right) \quad \text{as } \ell \rightarrow \infty. \quad (\text{A.8})$$

This limiting form is written as the cosine kernel (cos).

## B Model construction / search

### B.1 Overview

The model construction phase of ABCD starts with the kernel equal to the noise kernel, WN. New kernel expressions are generated by applying search operators to the current kernel. When new base kernels are proposed by the search operators, their parameters are randomly initialised with several restarts. Parameters are then optimized by conjugate gradients to maximise the likelihood of the data conditioned on the kernel parameters. The kernels are then scored by the Bayesian information criterion and the top scoring kernel is



selected as the new kernel. The search then proceeds by applying the search operators to the new kernel i.e. this is a greedy search algorithm.

In all experiments, 10 random restarts were used for parameter initialisation and the search was run to a depth of 10.

## B.2 Search operators

ABCD is based on a search algorithm which used the following search operators

$$\mathcal{S} \rightarrow \mathcal{S} + \mathcal{B} \quad (\text{B.1})$$

$$\mathcal{S} \rightarrow \mathcal{S} \times \mathcal{B} \quad (\text{B.2})$$

$$\mathcal{B} \rightarrow \mathcal{B}' \quad (\text{B.3})$$

where  $\mathcal{S}$  represents any kernel subexpression and  $\mathcal{B}$  is any base kernel within a kernel expression i.e. the search operators represent addition, multiplication and replacement.

To accommodate changepoint/window operators we introduce the following additional operators

$$\mathcal{S} \rightarrow \text{CP}(\mathcal{S}, \mathcal{S}) \quad (\text{B.4})$$

$$\mathcal{S} \rightarrow \text{CW}(\mathcal{S}, \mathcal{S}) \quad (\text{B.5})$$

$$\mathcal{S} \rightarrow \text{CW}(\mathcal{S}, \mathcal{C}) \quad (\text{B.6})$$

$$\mathcal{S} \rightarrow \text{CW}(\mathcal{C}, \mathcal{S}) \quad (\text{B.7})$$

where  $\mathcal{C}$  is the constant kernel. The last two operators result in a kernel only applying outside or within a certain region.

Based on experience with typical paths followed by the search algorithm we introduced the following operators

$$\mathcal{S} \rightarrow \mathcal{S} \times (\mathcal{B} + \mathcal{C}) \quad (\text{B.8})$$

$$\mathcal{S} \rightarrow \mathcal{B} \quad (\text{B.9})$$

$$\mathcal{S} + \mathcal{S}' \rightarrow \mathcal{S} \quad (\text{B.10})$$

$$\mathcal{S} \times \mathcal{S}' \rightarrow \mathcal{S} \quad (\text{B.11})$$

where  $\mathcal{S}'$  represents any other kernel expression. Their introduction is currently not rigorously justified.

## C Predictive accuracy

**Interpolation** To test the ability of the methods to interpolate, we randomly divided each data set into equal amounts of training data and testing data. We trained each algorithm on the training half of the data, produced predictions for the remaining half and then computed the root mean squared error (RMSE). The values of the RMSEs are then standardised by dividing by the smallest RMSE for each data set i.e. the best performance on each data set will have a value of 1.

Figure 11 shows the standardised RMSEs for the different algorithms. The box plots show that all quartiles of the distribution of standardised RMSEs are lower for both versions of ABCD. The median for ABCD-accuracy is 1; it is the best performing algorithm on 7 datasets. The largest outliers of ABCD and spectral kernels are similar in value.

Changepoints performs slightly worse than MKL despite being strictly more general than Changepoints. The introduction of changepoints allows for more structured models, but it introduces parametric forms into the regression models (i.e. the sigmoids expressing the changepoints). This results in worse interpolations at the locations of the change

points, suggesting that a more robust modeling language would require a more flexible class of changepoint shapes or improved inference (e.g. fully Bayesian inference over the location and shape of the changepoint).

Eureqa is not suited to this task and performs poorly. The models learned by Eureqa tend to capture only broad trends of the data since the fine details are not well explained by parametric forms.

## C.1 Tabeles of standardised RMSEs

See table 2 for raw interpolation results and table 3 for raw extrapolation results. The rows follow the order of the datasets in the rest of the supplementary material. The following abbreviations are used: ABCD-accuracy (ABCD-acc), ABCD-interpretability ((ABCD-int), Spectral kernels (SP), Trend-cyclical-irregular (TCI), Bayesian MKL (MKL), Eureqa (EL), Changepoints (CP), Squared exponential (SE) and Linear regression (Lin).

## D Guide to the automatically generated reports

Additional supplementary material to this paper is 13 reports automatically generated by ABCD. A link to these reports will be maintained at <http://mlg.eng.cam.ac.uk/lloyd/>. We recommend that you read the report for ‘01-airline’ first and review the reports that follow afterwards more briefly. ‘02-solar’ is discussed in the main text. ‘03-mauna’ analyses a dataset mentioned in the related work. ‘04-wheat’ demonstrates changepoints being used to capture heteroscedasticity. ‘05-temperature’ extracts an exactly periodic pattern from noisy data. ‘07-call-centre’ demonstrates a large discontinuity being modeled by a changepoint. ‘10-sulphuric’ combines many changepoints to create a highly structured model of the data. ‘12-births’ discovers multiple periodic components.

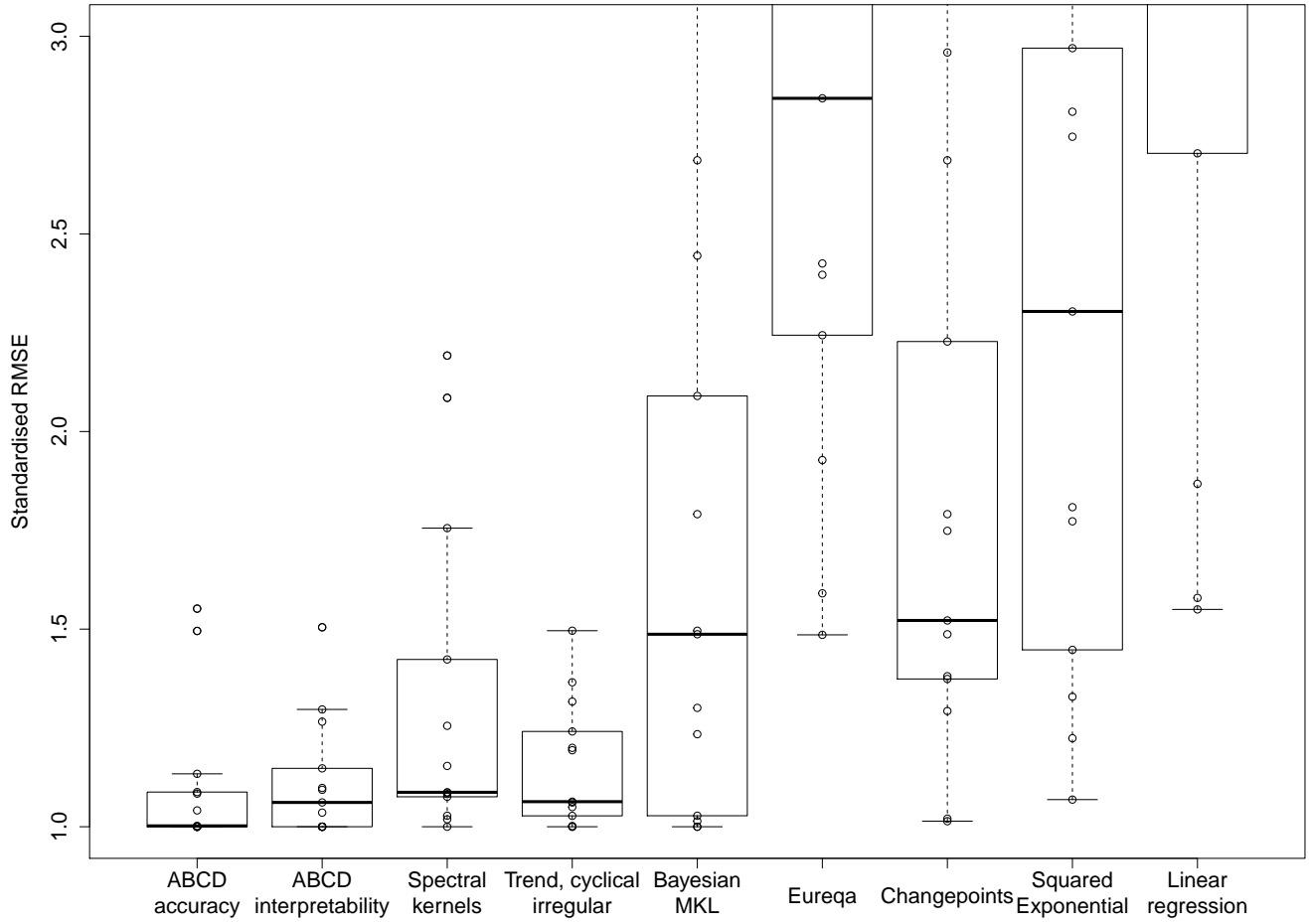


Figure 11: Box plot of standardised RMSE (best performance = 1) on 13 interpolation tasks.

ABCD-acc	ABCD-int	SP	TCI	MKL	EL	CP	SE	Lin
1.04	1.00	2.09	1.32	3.20	5.30	3.25	4.87	5.01
1.00	1.27	1.09	1.50	1.50	3.22	1.75	2.75	3.26
1.00	1.00	1.09	1.00	2.69	26.20	2.69	7.93	10.74
1.09	1.04	1.00	1.00	1.00	1.59	1.37	1.33	1.55
1.00	1.06	1.08	1.06	1.01	1.49	1.01	1.07	1.58
1.50	1.00	2.19	1.37	2.09	7.88	2.23	6.19	7.36
1.55	1.50	1.02	1.00	1.00	2.40	1.52	1.22	6.28
1.00	1.30	1.26	1.24	1.49	2.43	1.49	2.30	3.20
1.00	1.09	1.08	1.06	1.30	2.84	1.29	2.81	3.79
1.08	1.00	1.15	1.19	1.23	42.56	1.38	1.45	2.70
1.13	1.00	1.42	1.05	2.44	3.29	2.96	2.97	3.40
1.00	1.15	1.76	1.20	1.79	1.93	1.79	1.81	1.87
1.00	1.10	1.03	1.03	1.03	2.24	1.02	1.77	9.97

Table 2: Interpolation standardised RMSEs

ABCD-acc	ABCD-int	SP	TCI	MKL	EL	CP	SE	Lin
1.14	2.10	1.00	1.44	4.73	3.24	4.80	32.21	4.94
1.00	1.26	1.21	1.03	1.00	2.64	1.03	1.61	1.07
1.40	1.00	1.32	1.29	1.74	2.54	1.74	1.85	3.19
1.07	1.18	3.00	3.00	3.00	1.31	1.00	3.03	1.02
1.00	1.00	1.03	1.00	1.35	1.28	1.35	2.72	1.51
1.00	2.03	3.38	2.14	4.09	6.26	4.17	4.13	4.93
2.98	1.00	11.04	1.80	1.80	493.30	3.54	22.63	28.76
3.10	1.88	1.00	2.31	3.13	1.41	3.13	8.46	4.31
1.00	2.05	1.61	1.52	2.90	2.73	3.14	2.85	2.64
1.00	1.45	1.43	1.80	1.61	1.97	2.25	1.08	3.52
2.16	2.03	3.57	2.23	1.71	2.23	1.66	1.89	1.00
1.06	1.00	1.54	1.56	1.85	1.93	1.84	1.66	1.96
3.03	4.00	3.63	3.12	3.16	1.00	5.83	5.35	4.25

Table 3: Extrapolation standardised RMSEs