
Raiders of the Lost Architecture: Kernels for Bayesian Optimization in Conditional Parameter Spaces

Abstract

In practical Bayesian optimization, we must often search over structures with differing numbers of parameters. For instance, we may wish to search over neural network architectures with an unknown number of layers. To relate performance data gathered for different architectures, we define a new kernel for conditional parameter spaces that explicitly includes information about which parameters are relevant in a given structure. We show that this kernel improves model quality and Bayesian optimization results over several simpler baseline kernels.

1 Introduction

Bayesian optimization (BO) is an efficient approach for solving blackbox optimization problems of the form $\arg \min_{x \in \mathcal{X}} f(x)$ (see (Brochu et al., 2010) for a detailed overview), where f is expensive to evaluate. It employs a prior distribution $p(f)$ over functions that is updated as new information on f becomes available. The most common choice of prior distribution are Gaussian processes (GPs (Rasmussen & Williams, 2006)), as they are powerful and flexible models for which the marginal and conditional distributions can be computed efficiently.¹ However, some problem domains remain challenging to model well with GPs, and the efficiency and effectiveness of Bayesian optimization suffers as a result. In this paper, we tackle the common problem of input dimensions that are only relevant if other inputs take certain values (Hutter, 2009; Bergstra et al., 2011). This is a general problem in algorithm configuration (Hutter, 2009) that occurs in many ma-

¹There are prominent exceptions to this rule, though. In particular, tree-based models, such as random forests, can be a better choice if there are many data points (and GPs thus become computationally inefficient), if the input dimensionality is high, if the noise is not normally distributed, or if there are non-stationarities (Taddy et al., 2011; Hutter et al., 2011; Bergstra et al., 2011).

chine learning contexts, such as, e.g., in deep neural networks (Hinton et al., 2006); flexible computer vision architectures (Bergstra et al., 2013); and the combined selection and hyperparameter optimization of machine learning algorithms (Thornton et al., 2013). We detail the case of deep neural networks below.

Bayesian optimization has recently been applied successfully to deep neural networks (Snoek et al., 2012; Bergstra et al., 2011) to optimize high level model parameters and optimization parameters, which we will refer to collectively as *hyperparameters*. Deep neural networks represent the state-of-the-art on multiple machine learning benchmarks such as object recognition (Krizhevsky et al., 2012), speech recognition (Hinton et al., 2012a), natural language processing (Mikolov et al., 2010) and more. They are multi-layered models by definition, and each layer is typically parameterized by a unique set of hyperparameters, such as regularization parameters and the layer capacity or number of hidden units. Thus adding additional layers introduces additional hyperparameters to be optimized. The result is a complex hierarchical conditional parameter space, which is difficult to search over. Historically, practitioners have simply built a separate model for each type of architecture or used non-GP models (Bergstra et al., 2011), or assumed a fixed architecture (Snoek et al., 2012). If there is any relation between networks with different architectures, separately modelling each is wasteful.

GPs with standard kernels fail to model the performance of architectures with such conditional hyperparameters. To remedy this, the contribution of this paper is the introduction of a kernel that allows observed information to be shared across architectures when this is appropriate. We demonstrate the effectiveness of this kernel on a GP regression task and a Bayesian optimization task using a feed-forward classification neural network.

2 A Kernel for Conditional Parameter Spaces

GPs employ a positive-definite kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ to model the covariance between function values. Typical GP models cannot, however, model the covariance between function values whose inputs have different (pos-

sibly overlapping) sets of relevant variables.

In this section, we construct a kernel between points in a space that may have dimensions which are irrelevant under known conditions (further details are available in (Hutter & Osborne, 2013)). As an explicit example, we consider a deep neural network: if we set the network depth to 2 we know that the 3rd layer’s hyperparameters do not have any effect (as there is no 3rd layer).

Formally, we aim to do inference about some function f with domain \mathcal{X} . $\mathcal{X} = \prod_{i=1}^D \mathcal{X}_i$ is a D -dimensional input space, where each individual dimension is bounded real, that is, $\mathcal{X}_i = [l_i, u_i] \subset \mathbb{R}$ (with lower and upper bounds l_i and u_i , respectively). We define functions $\delta_i: \mathcal{X} \rightarrow \{\text{true}, \text{false}\}$, for $i \in \{1, \dots, D\}$. $\delta_i(\mathbf{x})$ stipulates the relevance of the i th feature x_i to $f(\mathbf{x})$.

2.1 The problem

As an example, imagine trying to model the performance of a neural network having either one or two hidden layers, with respect to the regularization parameters for each layer, x_1 and x_2 . If y represents the performance of a one layer-net with regularization parameters x_1 and x_2 , then the value x_2 doesn’t matter, since there is no second layer to the network. Below, we’ll write an input triple as $(x_1, \delta_2(\mathbf{x}), x_2)$ and assume that $\delta_1(\mathbf{x}) = \text{true}$; that is, the regularization parameter for the first layer is always relevant.

In this setting, we want a kernel k to be dependent on which parameters are relevant, and the values of relevant parameters for both points. For example, consider first-layer parameters x_1 and x'_1 :

- If we are comparing two points for which the same parameters are relevant, the value of any unused parameters shouldn’t matter,

$$\begin{aligned} & k((x_1, \text{false}, x_2), (x'_1, \text{false}, x'_2)) \\ &= k((x_1, \text{false}, x''_2), (x'_1, \text{false}, x'''_2)), \forall x_2, x'_2, x''_2, x'''_2; \end{aligned} \quad (1)$$

- The covariance between a point using both parameters and a point using only one should again only depend on their shared parameters,

$$\begin{aligned} & k((x_1, \text{false}, x_2), (x'_1, \text{true}, x'_2)) \\ &= k((x_1, \text{false}, x''_2), (x'_1, \text{true}, x'''_2)), \forall x_2, x'_2, x''_2, x'''_2. \end{aligned} \quad (2)$$

Put another way, in the absence of any other information, this specification encodes our prior ignorance about the irrelevant (missing) parameters while still allowing us to model correlations between relevant parameters.

Figure 1: A demonstration of the embedding giving rise to the pseudo-metric. All points for which $\delta_2(x) = \text{false}$ are mapped onto a line varying only along x_1 . Points for which $\delta_2(x) = \text{true}$ are mapped to the surface of a semicylinder, depending on both x_1 and x_2 . This embedding gives a constant distance between pairs of points which have differing values of δ but the same values of x_1 .

165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

Figure 2: The box embedding (placeholder)

220		275
$g(0, 0, \text{true})$		276
221		277
222		278
223		279
$g(0, 1, \text{true})$		280
224		281
225		282
226		283
227		284
228		285
229		286
230		287
231		288
232		289
233		290
234		291
235		292
236		293
237		294
238		295
239		296
240		297
241		298
242		299
243		300
244		301
245		302
246		303
247		304
248		305
249		306
250		307
251		308
252		309
253		310
254		311
255		312
256		313
257		314
258		315
259		316
260		317
261		318
262		319
263		320
264		321
265		322
266		323
267		324
268		325
269		326
270		327
271		328
272		329
273		
274		

GP, as in (Bergstra et al., 2011). In the case where all data comes from the same architecture, the hierarchical GP model makes slightly different assumptions than a standard GP, because it places the data on a semi-circle. To check whether this alternate assumption affects performance, we also include a Separate-Hierarchical GP model. We also compare against a simpler embedding method, or “Poor Man’s Embedding”. In this procedure, we combine data from all models into a single model, replacing any unused parameters with a heuristic constant, set to -1 in these experiments.

The results are averaged over 10-fold train/test splits. Kernel parameters were inferred using slice sampling (Murray & Adams, 2010). As the errors lie between 0 and 1 with many distributed toward the lower end, it can be beneficial to take the log of the outputs before modelling them with a GP. We experiment with both the original and transformed outputs.

Data. 951 datapoints were generated by performing Bayesian optimization on MNIST. Results are averaged over 10-fold train/test splits.

Table 1: Normalized Mean Squared Error on MNIST Bayesian optimization data

Method	Original data	Log outputs
Separate Linear	0.812 ± 0.045	0.737 ± 0.049
Separate GP	0.546 ± 0.038	0.446 ± 0.041
Separate Arc GP	0.535 ± 0.030	0.440 ± 0.031
Linear	0.876 ± 0.043	0.834 ± 0.047
GP	0.481 ± 0.031	0.401 ± 0.028
Arc GP	0.421 ± 0.033	0.335 ± 0.028

Results. Table 1 shows that a GP using the arc kernel performs favourably to a GP that ignores the relevance information of each point. The “separate” categories apply a different model to each layer and therefore do not take advantage of dependencies between layers. Interestingly, the separate Arc GP, which is effectively just a standard GP with additional embedding, performs comparably to a standard GP, suggesting that the embedding doesn’t limit the expressiveness of the model.

4.2 Bayesian Optimization Experiments

In this experiment, we test the ability of Bayesian optimization to tune the hyperparameters of each layer of a deep neural network. We allow the neural networks for these problems to use up to 5 hidden layers (or no hidden layer). We optimize over learning rates, L2 weight constraints, dropout rates (Hinton et al., 2012b), and the number of hidden units per layer leading to a total of up to 23 hyperparameters and 6 architectures. On MNIST, most effort is

spent improving the error by a fraction of a percent, therefore we optimize this dataset using the log-classification error. For CIFAR-10, we use classification error as the objective. We use the Deepnet² package, and each function evaluation took approximately 1000 to 2000 seconds to run on NVIDIA GTX Titan GPUs. Note that when a network of depth n is tested, all hyperparameters from layers $n + 1$ onward are deemed irrelevant.

Experimental Setup. For Bayesian optimization, we follow the methodology of (Snoek et al., 2012), using slice sampling and the expected improvement heuristic. In this methodology, the acquisition function is optimized by first selecting from a pre-determined grid of points lying in $[0, 1]^{23}$, distributed according to a Sobol sequence. Our baseline is a standard Gaussian process over this space that is agnostic to whether particular dimensions are irrelevant for a given point.

Results. Figure 4 shows that on these datasets, using the arc kernel consistently reaches good solutions faster than the naive baseline, or it finds a better solution. In the case of MNIST, the best discovered model achieved 1.19% test error using 50000 training examples. By comparison, (Wan et al., 2013) achieved 1.28% test error using a similar model and 60000 training examples. Similarly, our best model for CIFAR-10 achieved 21.1% test error using 45000 training examples and 400 features. For comparison, a support vector machine using 1600 features with the same feature pipeline and 50000 training examples achieves 22.1% error. Figure 4c shows the proportion of function evaluations spent on each architecture size for the CIFAR-10 experiments. Interestingly, the baseline tends to favour smaller models while a GP using the arc kernel distributes it’s efforts amongst deeper architectures that tend to yield better results.

5 Conclusion

We introduced the arc kernel for conditional parameter spaces that facilitates modelling the performance of deep neural network architectures by enabling the sharing of information across architectures where useful. Empirical results show that this kernel improves GP model quality and GP-based Bayesian optimization results over several simpler baseline kernels. Allowing information to be shared across architectures improves the efficiency of Bayesian optimization and removes the need to manually search for good architectures. The resulting models perform favourably compared to established benchmarks by domain experts.

²<https://github.com/nitishsrivastava/deepnet>

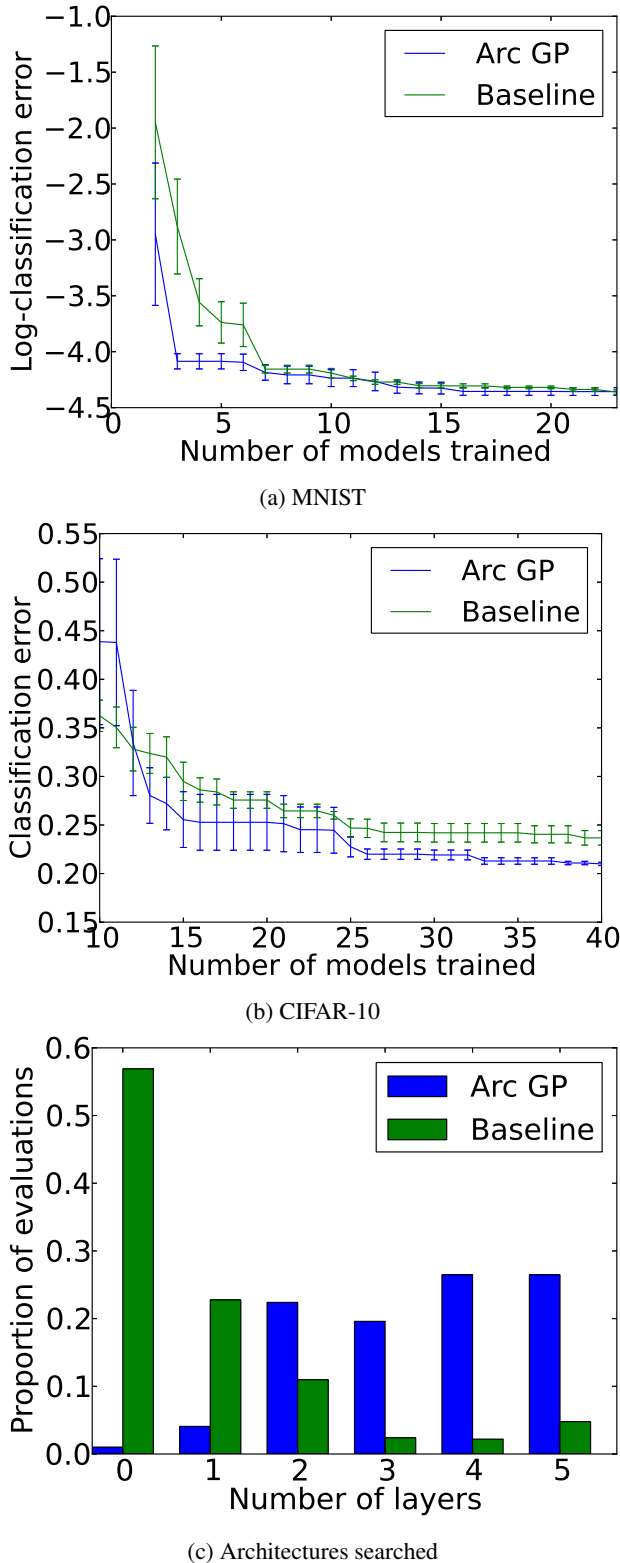


Figure 4: Bayesian optimization results using the arc kernel.

6 Acknowledgements

References

- Bergstra, James, Bardenet, Rémi, Bengio, Yoshua, Kégl, Balázs, et al. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, 2011.
- Bergstra, James, Yamins, Daniel, and Cox, David. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *International Conference on Machine Learning*, pp. 115–123, 2013.
- Brochu, Eric, Brochu, Tyson, and de Freitas, Nando. A Bayesian interactive optimization approach to procedural animation design. In *ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, 2010.
- Coates, Adam, Lee, Honglak, and Ng, Andrew Y. An analysis of single-layer networks in unsupervised feature learning. *Artificial Intelligence and Statistics*, 2011.
- Hinton, Geoffrey E., Osindero, Simon, and Teh, Yee-Whye. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, July 2006. ISSN 0899-7667.
- Hinton, Geoffrey E., Deng, Li, Yu, Dong, Dahl, George E., rahman Mohamed, Abdel, Jaitly, Navdeep, Senior, Andrew, Vanhoucke, Vincent, Nguyen, Patrick, Sainath, Tara N., and Kingsbury, Brian. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Process. Mag.*, 29(6):82–97, 2012a.
- Hinton, Geoffrey E., Srivastava, Nitish, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012b.
- Hutter, Frank. *Automated Configuration of Algorithms for Solving Hard Computational Problems*. PhD thesis, University of British Columbia, Department of Computer Science, Vancouver, Canada, October 2009.
- Hutter, Frank and Osborne, Michael A. A kernel for hierarchical parameter spaces, 2013. arXiv:1310.5738.
- Hutter, Frank, Hoos, Holger H., and Leyton-Brown, Kevin. Sequential model-based optimization for general algorithm configuration. In *Proc. of LION-5*, pp. 507–523, 2011.
- Krizhevsky, Alex. Learning multiple layers of features from tiny images. *Technical report, Department of Computer Science, University of Toronto*, 2009.

550	Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoff. Im-	605
551	agenet classification with deep convolutional neural net-	606
552	works. In <i>Advances in Neural Information Processing</i>	607
553	<i>Systems</i> . 2012.	608
554		609
555	Lecun, Yann, Bottou, Lon, Bengio, Yoshua, and Haffner,	610
556	Patrick. Gradient-based learning applied to document	611
557	recognition. In <i>Proc. of the IEEE</i> , pp. 2278–2324, 1998.	612
558		613
559	Mikolov, Tomas, Karafiát, Martin, Burget, Lukas, Cer-	614
560	nocký, Jan, and Khudanpur, Sanjeev. Recurrent neu-	615
561	ral network based language model. In <i>Interspeech</i> , pp.	616
562	1045–1048, 2010.	617
563		618
564	Murray, Iain and Adams, Ryan P. Slice sampling covari-	619
565	ance hyperparameters of latent Gaussian models. In <i>Ad-</i>	620
566	<i>vances in Neural Information Processing Systems</i> , 2010.	621
567		622
568	Rasmussen, Carl E. and Williams, Christopher K.I. Gaus-	623
569	sian Processes for Machine Learning. <i>The MIT Press,</i>	624
570	<i>Cambridge, MA, USA</i> , 2006.	625
571		626
572	Snoek, Jasper, Larochelle, Hugo, and Adams,	627
573	Ryan Prescott. Practical Bayesian optimization of	628
574	machine learning algorithms. In <i>Advances in Neural</i>	629
575	<i>Information Processing Systems</i> , 2012.	630
576		631
577	Taddy, Matthew A., Gramacy, Robert B., and Polson,	632
578	Nicholas G. Dynamic trees for learning and design.	633
579	<i>Journal of the American Statistical Association</i> , 106	634
580	(493):109–123, 2011.	635
581		636
582	Thornton, Chris, Hutter, Frank, Hoos, Holger H., and	637
583	Leyton-Brown, Kevin. Auto-WEKA: Combined selec-	638
584	tion and hyperparameter optimization of classification	639
585	algorithms. In <i>Proc. of KDD'13</i> , pp. 847–855, 2013.	640
586		641
587	Wan, Li, Zeiler, Matthew, Zhang, Sixin, Cun, Yann L, and	642
588	Fergus, Rob. Regularization of neural networks using	643
589	dropconnect. In <i>International Conference on Machine</i>	644
590	<i>Learning</i> , 2013.	645
591		646
592		647
593		648
594		649
595		650
596		651
597		652
598		653
599		654
600		655
601		656
602		657
603		658
604		659