
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 $\arg \min_{x \in X} f(x)$ (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 machine learning contexts, such as, e.g., in deep neural net-

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).

works (Hinton et al., 2006); flexible computer vision architectures (Bergstra et al.); 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 (possibly 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,

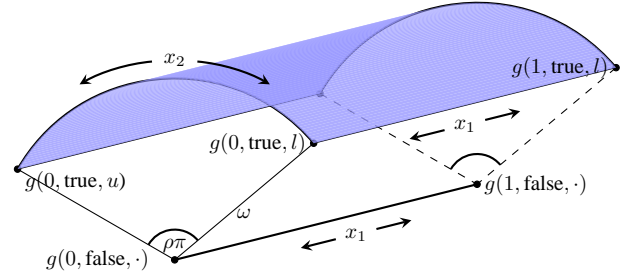
$$\begin{aligned} & k((x_1, \text{false}, x_2), (x'_1, \text{true}, x'_2)) \\ &= k((x_1, \text{false}, x'_2), (x'_1, \text{true}, x'_2)), \forall x_2, x'_2, x''_2, x'''_2. \end{aligned} \quad (2)$$

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.

2.2 Cylindrical Embedding

We can build a kernel with these properties for each possibly irrelevant input dimension i by embedding our points

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 .



into a Euclidean space. Specifically, the embedding we use is

$$g_i(\mathbf{x}) = \begin{cases} [0, 0]^\top & \text{if } \delta_i(\mathbf{x}) = \text{false} \\ \omega_i [\sin \pi \rho_i \frac{x_i}{u_i - l_i}, \cos \pi \rho_i \frac{x_i}{u_i - l_i}]^\top & \text{otherwise.} \end{cases} \quad (3)$$

Where $\omega_i \in \mathbb{R}^+$ and $\rho_i \in [0, 1]$.

Figure 2.2 shows a visualization of the embedding of points $(x_1, \delta_2(\mathbf{x}), x_2)$ into \mathbb{R}^3 . In this space, we have the Euclidean distance,

$$\begin{aligned} d_i(\mathbf{x}, \mathbf{x}') &= \|g_i(\mathbf{x}) - g_i(\mathbf{x}')\|_2 \\ &= \begin{cases} 0 & \text{if } \delta_i(\mathbf{x}) = \delta_i(\mathbf{x}') = \text{false} \\ \omega_i & \text{if } \delta_i(\mathbf{x}) \neq \delta_i(\mathbf{x}') \\ \omega_i \sqrt{2} \sqrt{1 - \cos(\pi \rho_i \frac{x_i - x'_i}{u_i - l_i})} & \text{if } \delta_i(\mathbf{x}) = \delta_i(\mathbf{x}') = \text{true.} \end{cases} \end{aligned} \quad (4)$$

We can use this to define a covariance over our original space. In particular, we consider the class of covariances that are functions only of the Euclidean distance Δ between points. There are many examples of such covariances. Popular examples are the exponentiated quadratic, for which $\kappa(\Delta) = \sigma^2 \exp(-\frac{1}{2}\Delta^2)$, or the rational quadratic, for which $\kappa(\Delta) = \sigma^2 (1 + \frac{1}{2\alpha}\Delta^2)^{-\alpha}$. We can simply take (4) in the place of Δ , returning a valid covariance that satisfies all desiderata above.

Explicitly, note that as desired, if i is irrelevant for both \mathbf{x} and \mathbf{x}' , d_i specifies that $g(\mathbf{x})$ and $g(\mathbf{x}')$ should not differ owing to differences between x_i and x'_i . Secondly, if i is relevant for both \mathbf{x} and \mathbf{x}' , the difference between $f(\mathbf{x})$ and $f(\mathbf{x}')$ due to x_i and x'_i increases monotonically with increasing $|x_i - x'_i|$. The parameter ρ_i controls whether differing in the relevance of i contributes more or less to the distance than differing in the value of x_i , should i be

(a) MNIST

(b) CIFAR-10

(c) Architectures searched

Figure 3: Bayesian optimization results using the arc kernel.

Results. Figure 3 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 3c 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.

4 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.

5 Acknowledgements

References

- Bergstra, James, Yamins, Daniel, and Cox, David. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures.
- Bergstra, James, Bardenet, Rémi, Bengio, Yoshua, Kégl, Balázs, et al. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, 2011.
- 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

440	nets. <i>Neural Computation</i> , 18(7):1527–1554, July 2006.	495
441	ISSN 0899-7667.	496
442		497
443	Hinton, Geoffrey E., Deng, Li, Yu, Dong, Dahl, George E.,	498
444	rahman Mohamed, Abdel, Jaitly, Navdeep, Senior, An-	499
445	drew, Vanhoucke, Vincent, Nguyen, Patrick, Sainath,	500
446	Tara N., and Kingsbury, Brian. Deep neural networks	501
447	for acoustic modeling in speech recognition: The shared	502
448	views of four research groups. <i>IEEE Signal Process.</i>	503
449	<i>Mag.</i> , 29(6):82–97, 2012a.	504
450	Hinton, Geoffrey E., Srivastava, Nitish, Krizhevsky, Alex,	505
451	Sutskever, Ilya, and Salakhutdinov, Ruslan. Improving	506
452	neural networks by preventing co-adaptation of feature	507
453	detectors. <i>arXiv preprint arXiv:1207.0580</i> , 2012b.	508
454		509
455	Hutter, Frank. <i>Automated Configuration of Algorithms for</i>	510
456	<i>Solving Hard Computational Problems</i> . PhD thesis, Uni-	511
457	versity of British Columbia, Department of Computer	512
458	Science, Vancouver, Canada, October 2009.	513
459		514
460	Hutter, Frank and Osborne, Michael A. A kernel for hier-	515
461	archical parameter spaces, 2013. <i>arXiv:1310.5738</i> .	516
462	Hutter, Frank, Hoos, Holger H., and Leyton-Brown, Kevin.	517
463	Sequential model-based optimization for general algo-	518
464	rithm configuration. In <i>Proc. of LION-5</i> , pp. 507–523,	519
465	2011.	520
466		521
467	Krizhevsky, Alex. Learning multiple layers of features	522
468	from tiny images. <i>Technical report, Department of Com-</i>	523
469	<i>puter Science, University of Toronto</i> , 2009.	524
470		525
471	Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoff. Im-	526
472	agenet classification with deep convolutional neural net-	527
473	works. In <i>Advances in Neural Information Processing</i>	528
474	<i>Systems</i> . 2012.	529
475		530
476	Lecun, Yann, Bottou, Lon, Bengio, Yoshua, and Haffner,	531
477	Patrick. Gradient-based learning applied to document	532
478	recognition. In <i>Proc. of the IEEE</i> , pp. 2278–2324, 1998.	533
479		534
480	Mikolov, Tomas, Karafiát, Martin, Burget, Lukas, Cer-	535
481	nocký, Jan, and Khudanpur, Sanjeev. Recurrent neu-	536
482	ral network based language model. In <i>Interspeech</i> , pp.	537
483	1045–1048, 2010.	538
484		539
485	Murray, Iain and Adams, Ryan P. Slice sampling covari-	540
486	ance hyperparameters of latent Gaussian models. In <i>Ad-</i>	541
487	<i>vances in Neural Information Processing Systems</i> , 2010.	542
488		543
489	Rasmussen, Carl E. and Williams, Christopher K.I. Gaus-	544
490	sian Processes for Machine Learning. <i>The MIT Press,</i>	545
491	<i>Cambridge, MA, USA</i> , 2006.	546
492		547
493	Snoek, Jasper, Larochelle, Hugo, and Adams,	548
494	Ryan Prescott. Practical Bayesian optimization of	549
	machine learning algorithms. In <i>Advances in Neural</i>	
	<i>Information Processing Systems</i> , 2012.	