

Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms

Chris Thornton Frank Hutter Holger H. Hoos Kevin Leyton-Brown

Department of Computer Science, University of British Columbia
201-2366 Main Mall, Vancouver BC, V6T 1Z4, Canada
{cwtornt, hutter, hoos, kevinlb}@cs.ubc.ca

ABSTRACT

Many different machine learning algorithms exist; taking into account each algorithm’s hyperparameters, there is a staggeringly large number of possible alternatives overall. We consider the problem of simultaneously selecting a learning algorithm and setting its hyperparameters, going beyond previous work that addresses these issues in isolation. We show that this problem can be addressed by a fully automated approach, leveraging recent innovations in Bayesian optimization. Specifically, we consider a wide range of feature selection techniques (combining 3 search and 8 evaluator methods) and all classification approaches implemented in WEKA, spanning 2 ensemble methods, 10 meta-methods, 27 base classifiers, and hyperparameter settings for each classifier. On each of 21 popular datasets from the UCI repository, the KDD Cup 09, variants of the MNIST dataset and CIFAR-10, we show classification performance often much better than using standard selection/hyperparameter optimization methods. We hope that our approach will help non-expert users to more effectively identify machine learning algorithms and hyperparameter settings appropriate to their applications, and hence to achieve improved performance.

Keywords

model selection, hyperparameter optimization, WEKA

1. INTRODUCTION

Increasingly, users of machine learning tools are non-experts who require off-the-shelf solutions. The machine learning community has much aided such users by making available a wide variety of sophisticated learning algorithms and feature selection methods through open source packages, such as WEKA [13] and PyBrain [22]. Each of these packages asks a user to make two kinds of choices: selecting a learning algorithm and customizing it by setting its hyperparameters (which also control any feature selection being performed). It can be challenging to make the right choice when faced with these degrees of freedom, leaving many users to select algorithms based on reputation or intuitive appeal, and/or to leave hyperparameters set to default values. Of course, this approach can yield performance far worse than that of the best method and hyperparameter settings.

This suggests a natural challenge for machine learning: given a dataset, to automatically and simultaneously choose a learning algorithm and set its hyperparameters to optimize empirical performance. We dub this the *combined algorithm selection and hyperparameter optimization problem* (short:

CASH); we formally define it in Section 3. Despite the practical importance of this problem, we are surprised to find no evidence that it has previously been considered in the literature. A likely explanation is that the combined space of learning algorithms and their hyperparameters is very challenging to search: the response function is noisy and the space is high dimensional, involves both categorical and continuous choices, and contains hierarchical dependencies (*e.g.*, the hyperparameters of a learning algorithm are only meaningful if that algorithm is chosen; the algorithm choices in an ensemble method are only meaningful if that ensemble method is chosen; etc). In contrast, we do note that there has been considerable past work separately addressing model selection [*e.g.*, 1, 5, 6, 7, 9, 20, 21, 29] and hyperparameter optimization [*e.g.*, 2, 3, 4, 12, 25, 27].

In what follows, we demonstrate that CASH can be viewed as a single hierarchical hyperparameter optimization problem, in which even the choice of algorithm itself is considered a hyperparameter. We also show that — based on this problem formulation — recent Bayesian optimization methods can obtain high quality results in reasonable time and with minimal human effort. After discussing some preliminaries (Section 2), we define the CASH problem and discuss methods for tackling it (Section 3). We then define a concrete CASH problem encompassing the full range of classifiers and feature selectors in the open source package WEKA (Section 4), and show that a search in the combined space of algorithms and hyperparameters yields better-performing models than standard algorithm selection/hyperparameter optimization methods (Section 5). More specifically, we show that the recent Bayesian optimization procedures TPE [3] and SMAC [14] find combinations of algorithms and hyperparameters that often outperform existing baseline methods, especially on large datasets.

2. PRELIMINARIES

This work focuses on classification problems: learning a function $f : \mathcal{X} \mapsto \mathcal{Y}$ with finite \mathcal{Y} . A *learning algorithm* A maps a set $\{d_1, \dots, d_n\}$ of training data points $d_i = (\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ to such a function, which is often expressed via a vector of *model parameters*. Most learning algorithms A further expose *hyperparameters* $\lambda \in \Lambda$, which change the way the learning algorithm A_λ itself works. For example, hyperparameters are used to describe a description-length penalty, the number of neurons in a hidden layer, the number of data points that a leaf in a decision tree must contain to be eligible for splitting, etc. These hyperparameters are typically optimized in an “outer loop” that evaluates the

performance of each hyperparameter configuration using cross-validation.

2.1 Model Selection

Given a set of learning algorithms \mathcal{A} and a limited amount of training data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, the goal of model selection is to determine the algorithm $A^* \in \mathcal{A}$ with optimal generalization performance. Generalization performance is estimated by splitting \mathcal{D} into disjoint training and validation sets $\mathcal{D}_{\text{train}}^{(i)}$ and $\mathcal{D}_{\text{valid}}^{(i)}$, learning functions f_i by applying A^* to $\mathcal{D}_{\text{train}}^{(i)}$, and evaluating the predictive performance of these functions on $\mathcal{D}_{\text{valid}}^{(i)}$. This allows for the model selection problem to be written as:

$$A^* \in \operatorname{argmin}_{A \in \mathcal{A}} \frac{1}{k} \sum_{i=1}^k \mathcal{L}(A, \mathcal{D}_{\text{train}}^{(i)}, \mathcal{D}_{\text{valid}}^{(i)}),$$

where $\mathcal{L}(A, \mathcal{D}_{\text{train}}^{(i)}, \mathcal{D}_{\text{valid}}^{(i)})$ is the loss (here: misclassification rate) achieved by A when trained on $\mathcal{D}_{\text{train}}^{(i)}$ and evaluated on $\mathcal{D}_{\text{valid}}^{(i)}$. We use k -fold cross-validation [18], which splits the training data into k equal-sized partitions $\mathcal{D}_{\text{valid}}^{(1)}, \dots, \mathcal{D}_{\text{valid}}^{(k)}$, and sets $\mathcal{D}_{\text{train}}^{(i)} = \mathcal{D} \setminus \mathcal{D}_{\text{valid}}^{(i)}$ for $i = 1, \dots, k$.¹

2.2 Hyperparameter Optimization

The problem of optimizing the hyperparameters $\lambda \in \Lambda$ of a given learning algorithm A is conceptually similar to that of model selection. Some key differences are that hyperparameters are often continuous, that hyperparameter spaces are often high dimensional, and that we can exploit correlation structure between different hyperparameter settings $\lambda_1, \lambda_2 \in \Lambda$. Given n hyperparameters $\lambda_1, \dots, \lambda_n$ with domains $\Lambda_1, \dots, \Lambda_n$, the hyperparameter space Λ is a subset of the crossproduct of these domains: $\Lambda \subset \Lambda_1 \times \dots \times \Lambda_n$. This subset is often strict, such as when certain settings of one hyperparameter render other hyperparameters inactive. For example, the parameters determining the specifics of the third layer of a deep belief network are not relevant if the network depth is set to one or two. Likewise, the parameters of a support vector machine’s polynomial kernel are not relevant if we use a different kernel instead.

More formally, following [15], we say that a hyperparameter λ_i is *conditional* on another hyperparameter λ_j , if λ_i is only active if hyperparameter λ_j takes values from a given set $V_i(j) \subseteq \Lambda_j$; in this case we call λ_j a *parent* of λ_i . Conditional hyperparameters can in turn be parents of other conditional hyperparameters, giving rise to a tree-structured space [3] or, in some cases, a directed acyclic graph (DAG) [15]. Given such a structured space Λ , the (hierarchical) hyperparameter optimization problem can be written as:

$$\lambda^* \in \operatorname{argmin}_{\lambda \in \Lambda} \frac{1}{k} \sum_{i=1}^k \mathcal{L}(A_\lambda, \mathcal{D}_{\text{train}}^{(i)}, \mathcal{D}_{\text{valid}}^{(i)}).$$

3. COMBINED ALGORITHM SELECTION AND HYPERPARAMETER OPTIMIZATION (CASH)

Given a set of algorithms $\mathcal{A} = \{A^{(1)}, \dots, A^{(k)}\}$ with associated hyperparameter spaces $\Lambda^{(1)}, \dots, \Lambda^{(k)}$, we define the

¹K-fold cross-validation is not the only available method for estimating generalization performance. We also experimented with the technique of repeated random subsampling validation [18], with similar results.

Algorithm 1 SMBO

- 1: initialise model \mathcal{M}_L ; $\mathcal{H} \leftarrow \emptyset$
 - 2: **while** time budget for optimization has not been exhausted **do**
 - 3: $\lambda \leftarrow$ candidate configuration from \mathcal{M}_L
 - 4: Compute $c = \mathcal{L}(A_\lambda, \mathcal{D}_{\text{train}}^{(i)}, \mathcal{D}_{\text{valid}}^{(i)})$
 - 5: $\mathcal{H} \leftarrow \mathcal{H} \cup \{(\lambda, c)\}$
 - 6: Update \mathcal{M}_L given \mathcal{H}
 - 7: **end while**
 - 8: **return** λ from \mathcal{H} with minimal c
-

combined algorithm selection and hyperparameter optimization problem (CASH) as computing

$$A^* \lambda^* \in \operatorname{argmin}_{A^{(j)} \in \mathcal{A}, \lambda \in \Lambda^{(j)}} \frac{1}{k} \sum_{i=1}^k \mathcal{L}(A_\lambda^{(j)}, \mathcal{D}_{\text{train}}^{(i)}, \mathcal{D}_{\text{valid}}^{(i)}). \quad (1)$$

We note that this problem can be reformulated as a single combined hierarchical hyperparameter optimization problem with parameter space $\Lambda = \Lambda^{(1)} \cup \dots \cup \Lambda^{(k)} \cup \{\lambda_r\}$, where λ_r is a new root-level hyperparameter that selects between algorithms $A^{(1)}, \dots, A^{(k)}$. The root-level parameters of each subspace $\Lambda^{(i)}$ are made conditional on λ_r being instantiated to A_i .

In principle, Problem 1 can be tackled in various ways. A promising approach is Bayesian Optimization [8], and in particular Sequential Model-Based Optimization [SMBO; 14], a versatile stochastic optimization framework that can work explicitly with both categorical and continuous hyperparameters, and that can exploit hierarchical structure stemming from conditional parameters. SMBO (outlined in Algorithm 1) first builds a model \mathcal{M}_L that captures the dependence of loss function \mathcal{L} on hyperparameter settings λ (line 1 in Algorithm 1). It then iterates the following steps: use \mathcal{M}_L to determine a promising candidate configuration of hyperparameters λ to evaluate next (line 3); evaluate the loss c of λ (line 4); and update the model \mathcal{M}_L with the new data point (λ, c) thus obtained (lines 5–6).

In order to select its next hyperparameter configuration λ using model \mathcal{M}_L , SMBO uses a so-called *acquisition function* $a_{\mathcal{M}_L} : \Lambda \rightarrow \mathbb{R}$, which uses the predictive distribution of model \mathcal{M}_L at arbitrary hyperparameter configurations $\lambda \in \Lambda$ to quantify (in closed form) how useful knowledge about λ would be. SMBO then simply maximizes this function over Λ to select the most useful configuration λ to evaluate next. Several prominent acquisition functions exist [17, 23, 26] that all aim to automatically trade off exploitation (locally optimizing hyperparameters in regions known to perform well) versus exploration (trying hyperparameters in a relatively unexplored region of the space) in order to avoid premature convergence. In this work, we use one of the most prominent acquisition functions, the *positive expected improvement (EI)* attainable over an existing given error rate c_{\min} [23]. Let $c(\lambda)$ denote the error rate of hyperparameter configuration λ . Then, the positive improvement function over c_{\min} is defined as:

$$I_{c_{\min}}(\lambda) := \max\{c_{\min} - c(\lambda), 0\}.$$

Of course, we do not know $c(\lambda)$. We can, however, compute its expectation with respect to the current model \mathcal{M}_L :

$$\mathbb{E}_{\mathcal{M}_L}[I_{c_{\min}}(\lambda)] := \int_{-\infty}^{c_{\min}} \max\{c_{\min} - c, 0\} \cdot p_{\mathcal{M}_L}(c | \lambda) dc. \quad (2)$$

One main difference between existing SMBO algorithms lies in the model class they employ. We now review the two SMBO algorithms whose models can handle hierarchical hyperparameters and that are thus suitable for the task of combined algorithm selection and hyperparameter optimization.

3.1 Sequential Model-based Algorithm Configuration (SMAC)

Sequential model-based algorithm configuration [SMAC; 14] supports a variety of models $p(c | \lambda)$ to capture the dependence of the loss function c on hyper-parameters λ , including approximate Gaussian processes and random forests. In this paper we use random forest models, since they tend to perform well with discrete and high-dimensional input data. SMAC handles conditional parameters by instantiating inactive conditional parameters in λ to default values for model training and prediction. This allows the individual decision trees to include splits of the kind “is hyperparameter λ_i active?”, allowing them to focus on active hyperparameters. While random forests are not usually treated as probabilistic models, SMAC obtains a predictive mean μ_λ and variance σ_λ^2 of $p(c | \lambda)$ as frequentist estimates over the predictions of its individual trees for λ ; it then models $p_{\mathcal{M}_L}(c | \lambda)$ as a Gaussian $\mathcal{N}(\mu_\lambda, \sigma_\lambda^2)$.

SMAC uses the expected improvement criterion defined in Equation 2, instantiating c_{min} to the error rate of the best hyperparameter configuration measured so far. Under SMAC’s predictive distribution $p_{\mathcal{M}_L}(c | \lambda) = \mathcal{N}(\mu_\lambda, \sigma_\lambda^2)$, this expectation can be computed by the closed-form expression

$$\mathbb{E}_{\mathcal{M}_L}[I_{c_{min}}(\lambda)] = \sigma_\lambda \cdot [u \cdot \Phi(u) + \varphi(u)],$$

where $u = \frac{c_{min} - \mu_\lambda}{\sigma_\lambda}$, and φ and Φ denote the probability density function and cumulative distribution function of a standard normal distribution, respectively [17].

SMAC is designed for robust optimization under noisy function evaluations, and as such implements special mechanisms to keep track of its best known configuration and assure high confidence in its estimate of that configuration’s performance. This robustness against noisy function evaluations can be exploited in combined algorithm selection and hyperparameter optimization, since the function to be optimized in Equation (1) is a mean over a set of loss terms (each corresponding to one pair of $\mathcal{D}_{train}^{(i)}$ and $\mathcal{D}_{valid}^{(i)}$ constructed from the training set). A key idea in SMAC is to make progressively better estimates of this mean by evaluating these terms one at a time, thus trading off accuracy against computational cost. In order for a new configuration to become a new incumbent, it must outperform the previous incumbent in every comparison made: considering only one fold, two folds, and so on up to the total number of folds previously used to evaluate the incumbent. (Furthermore, every time the incumbent survives such a comparison, it is evaluated on a new fold, up to the total number available, meaning that the number of folds used to evaluate the incumbent grows over time.) This means that a poorly performing configuration can be discarded after considering as little as a single fold.

Finally, SMAC also implements a diversification mechanism to achieve robust performance even when its model is misled: every second configuration is selected at random. Because of the evaluation procedure just described, the over-

head required by this safeguard is less than it might appear.

3.2 Tree-structured Parzen Estimator (TPE)

While SMAC models $p(c | \lambda)$ explicitly, the Tree-structure Parzen Estimator [TPE; 3] uses separate models for $p(c)$ and $p(\lambda | c)$. Specifically, it models $p(\lambda | c)$ as one of two density estimates, conditional on whether c is greater or less than a given threshold value c^* :

$$p(\lambda | c) = \begin{cases} \ell(\lambda), & \text{if } c < c^*. \\ g(\lambda), & \text{if } c \geq c^*. \end{cases} \quad (3)$$

Here, c^* is chosen as the γ -quantile of the losses TPE obtained so far (where γ is an algorithm parameter with a default value of $\gamma = 0.15$), $\ell(\cdot)$ is a density estimate learned from all previous hyperparameters λ with corresponding loss smaller than c^* , and $g(\cdot)$ is a density estimate learned from all previous hyperparameters λ with corresponding loss greater than or equal to c^* . Intuitively, this creates a probabilistic density estimator $\ell(\cdot)$ for hyperparameters that appear to do ‘well’, and a different density estimator $g(\cdot)$ for hyperparameters that appear ‘poor’ with respect to the threshold. Bergstra et al. [3] showed that the expected improvement $\mathbb{E}_{\mathcal{M}_L}[I_{c_{min}}(\lambda)]$ from Equation 2 is proportional to a quantity that can be computed in closed-form from γ , $g(\lambda)$, and $\ell(\lambda)$:

$$\mathbb{E}[I_{c_{min}}(\lambda)] \propto \left(\gamma + \frac{g(\lambda)}{\ell(\lambda)} \cdot (1 - \gamma) \right)^{-1}.$$

TPE maximizes this expression by generating many candidate hyperparameter configurations at random and picking λ with the smallest value of $g(\lambda)/\ell(\lambda)$.

The density estimators $\ell(\cdot)$ and $g(\cdot)$ have a hierarchical structure with discrete, continuous, and conditional variables reflecting the hyperparameters and their dependence relationships. For each node in this tree structure, a 1-D Parzen estimator is created to model the density of the node’s corresponding hyperparameter. For a given hyperparameter configuration λ that is being added to either ℓ or g , only the 1-D estimators corresponding to active hyperparameters in λ are updated. For continuous hyperparameters, these 1-D estimators are constructed by placing density in the form of a Gaussian at each hyperparameter value λ_i , with standard deviation set to the larger of each point’s left and right neighbour. Discrete hyperparameters are estimated with probabilities proportional to the number of times that a particular choice occurred in the set of observations. To evaluate a candidate hyperparameter λ ’s probability estimate, TPE starts at the root of the tree and descends into the leaves by following paths that only use active hyperparameters. At each node in this traversal, the probability of the corresponding hyperparameter is computed according to its 1-D estimator, and the individual probabilities are combined on a pass back up to the root of the tree. Note that this means that TPE assumes independence for hyperparameters that do not appear together along any path from the tree’s root to one of its leaves.

4. AUTO-WEKA

To demonstrate the feasibility of an automatic approach to solving the CASH problem, we built a tool, *Auto-WEKA*, that solves this problem for all classification algorithms and feature selectors/evaluators implemented in the WEKA package [13]. Note that while we have focused on classification

Table 1: Classifiers in Auto-WEKA: Classifiers marked with * are meta-methods, which in addition to their own parameters take one ‘base’ classifier and its parameters. Classifiers marked with + are ensemble methods that take as input up to 5 ‘base’ classifiers and their parameters. *Categorical* and *Numeric* refer to the number of hyperparameters of each kind for each classifier.

Classifier	Categorical Numeric	
	Categorical	Numeric
BAYES NET	2	0
NAIVE BAYES	2	0
NAIVE BAYES MULTINOMIAL	0	0
GAUSSIAN PROCESS	3	6
LINEAR REGRESSION	2	1
LOGISTIC REGRESSION	0	1
SINGLE-LAYER PERCEPTRON	5	2
STOCHASTIC GRADIENT DESCENT	3	2
SVM	4	6
SIMPLE LINEAR REGRESSION	0	0
SIMPLE LOGISTIC REGRESSION	2	1
VOTED PERCEPTRON	1	2
KNN	4	1
K-STAR	2	1
DECISION TABLE	4	0
RIPPER	3	1
M5 RULES	3	1
1-R	0	1
PART	2	2
0-R	0	0
DECISION STUMP	0	0
C4.5 DECISION TREE	6	2
LOGISTIC MODEL TREE	5	2
M5 TREE	3	1
RANDOM FOREST	2	3
RANDOM TREE	4	4
REP TREE	2	3
LOCALLY WEIGHTED LEARNING*	3	0
ADABOOST M1*	2	2
ADDITIVE REGRESSION*	1	2
ATTRIBUTE SELECTED*	2	0
BAGGING*	1	2
CLASSIFICATION VIA REGRESSION*	0	0
LOGITBOOST*	4	4
MULTICLASS CLASSIFIER*	3	0
RANDOM COMMITTEE*	0	1
RANDOM SUBSPACE*	0	2
VOTING ⁺	1	0
STACKING ⁺	0	0

algorithms in WEKA, there is no obstacle to extending our approach to other settings.

Table 1 provides a list of all 39 WEKA classification algorithms. Of these models, 27 are considered ‘base’ classifiers (which can be used independently), 10 of the remaining classifiers are meta methods (which take a single base classifier and its parameters as an input), and the final 2 ensemble classifiers can take any number of base classifiers as input. We allowed the meta-methods to use any base classifier with any hyperparameter settings, and allowed ensemble methods to use up to five base classifiers, again with any hyperparameter settings. Not all classifiers are applicable on all datasets (e.g., due to a classifier’s inability to handle missing data). For a given dataset, our Auto-WEKA implementation automatically only considers the subset of applicable classifiers.

Table 2 provides a list of WEKA’s 3 feature search methods, as well its 8 feature evaluators, and their respective number of subparameters (up to 5 for search; up to 4 for evaluators). To perform feature selection, a search method is combined

Table 2: Feature Search/Evaluator methods in Auto-WEKA: Methods marked with * are search methods, which require one feature evaluator that is used to determine the importance of a feature. *Categorical* and *Numeric* refer to the number of hyperparameters of each method.

Feature Method	Categorical Numeric	
	Categorical	Numeric
BEST FIRST*	1	1
GREEDY STEPWISE*	3	2
RANKER*	0	1
CFS SUBSET EVAL	2	0
PEARSON CORRELATION EVAL	0	0
GAIN RATIO EVAL	0	0
INFO GAIN EVAL	2	0
1-R EVAL	1	2
PRINCIPAL COMPONENTS EVAL	2	2
RELIEF EVAL	1	2
SYMMETRICAL UNCERTAINTY EVAL	1	0

with a feature evaluator, and the subparameters of both of them need to be instantiated. Feature selection is run as a pre-processing phase before building any classifier.

The algorithms in Table 1 and 2 have a wide variety of hyperparameters, which take on values from continuous intervals, from ranges of integers, and from other discrete sets. We associated either a uniform or log uniform prior with each numerical parameter, depending on its semantics. For example, we set a log uniform prior for the ridge regression penalty, and a uniform prior for the maximum depth for a tree in a random forest. Auto-WEKA works with continuous hyperparameter values directly; nevertheless, to give a sense of the size of the hypothesis space we studied, we note that discretizing hyperparameter domains to a maximum of 10 values each would give rise to over 10^{47} hyperparameter settings. We emphasize that this space is *much* larger than a simple union of the base learners’ hypothesis spaces (whose size is roughly 10^8), since the ensemble methods allow up to 5 *independent* base learners, giving rise to a space with roughly $(10^8)^5 = 10^{40}$ elements. The feature selection part gives rise to another independent decision between roughly 10^6 choices, and several parameters on the meta and ensemble level contribute another order of magnitude to the total size of AutoWEKA’s hypothesis space.

Auto-WEKA can be understood as a single learning algorithm with a highly conditional parameter space, as depicted in Figure 1. Auto-WEKA has two top-level Boolean parameters, the first of which is `is_base` that selects among single base classifiers and ensemble or meta-classifiers. If `is_base` is `true`, then the parameter `base` determines which of the 27 base classifiers are to be used. If `is_base` is `false`, then `class` indicates either an ensemble or a meta-classifier. If `class` is a meta-classifier, then the parameter `meta_base` is chosen to be one of the 27 base classifiers. In the event that `class` is an ensemble classifier, an additional parameter `num_classes` is an integer chosen from $\{1, \dots, 5\}$. `base_i` variables are then selected according to the value of `num_classes`, which again select which of the 27 base classifiers to use. For each `*base*` parameter, conditional hyperparameters for every model are attached.

The second top level Boolean parameter `feat_sel` indicates if one of the feature selection methods is going to be applied. If `feat_sel` is `false`, then Auto-WEKA passes the unmodified dataset to the classifier. If it is `true`, then

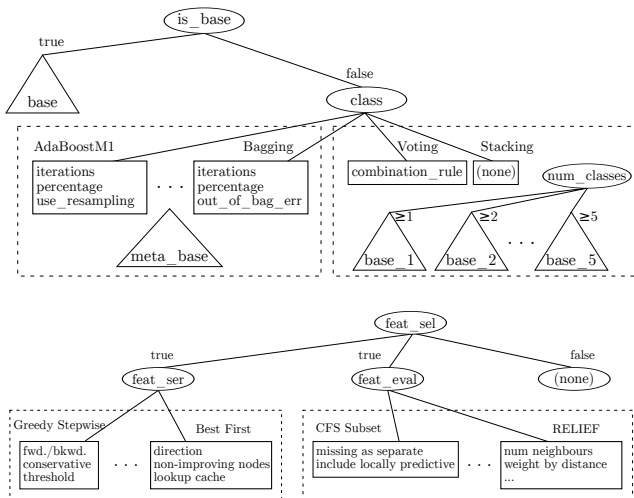


Figure 1: Auto-WEKA’s parameter space. Top: first top level Boolean, concerning Auto-WEKA’s classification methods. The triangular items represent a parameter that selects one of the 27 base classifiers, and adds conditional classifier hyperparameters accordingly. Bottom: second top level Boolean, concerning Auto-WEKA’s feature selection methods.

`feat_search` selects the choice of feature search method, and `feat_eval` selects the choice of feature evaluator. This results in a very wide tree that captures all the hierarchical nature of the model hyperparameters, and allows the creation of a single hyperparameter optimization problem with four hierarchical layers of a total of 786 parameters.

Auto-WEKA is agnostic to the choice of optimizer, so we implemented variants leveraging SMAC and TPE, respectively – the two Bayesian optimization algorithms that can handle hierarchical parameter spaces (see Sections 3.1 and 3.2.² We defined two Auto-WEKA variants, based on SMAC and TPE, respectively. We made both of these Auto-WEKA versions available to the public at www.cs.ubc.ca/labs/beta/Projects/autoweka/ and are committed to provide support for their widespread use in practice.

Both TPE and SMAC have their own parameters that influence their performance (such as TPE’s choice of the γ -quantile indicating ‘good’ or ‘bad’ performance, or the parameters of SMAC’s random forest model). In Auto-WEKA, we used the defaults for these meta-hyperparameters, as set by the authors (we assume further small improvements may be possible by meta-hyperparameter optimization, but a separate process with a meta-training/validation set split would be required to guard against over-fitting, and we did not attempt this).

Finally, we note that both TPE and SMAC are randomized algorithms and thus expected to produce different results based on the random seed provided. As demonstrated in [16], this allows for trivial, yet effective parallelization of the optimization process: simply perform k independent runs of the optimization method in parallel and select the result of the run with the lowest cross-validation error to return.³ We

²We thank the authors of SMAC and TPE for giving us access to their respective implementations.

³Other, more sophisticated methods for the parallelization of Bayesian optimization exist [16, 3, 10, 25], but to date, there

Table 3: Datasets Used; *Num. Discr.* and *Num. Cont.* refer to the number of discrete and continuous attributes of elements in the dataset, respectively.

Name	Num Discr.	Num Cont.	Num Classes	Num Training	Num Test
DEXTER	20 000	0	2	420	180
GERMANCREDIT	13	7	2	700	300
DOROTHEA	100 000	0	2	805	345
YEAST	0	8	10	1 038	446
AMAZON	10 000	0	49	1 050	450
SECOM	0	591	2	1 096	471
SEMEION	256	0	10	1 115	478
CAR	6	0	4	1 209	519
MADELON	500	0	2	1 820	780
KR-VS-KP	37	0	2	2 237	959
ABALONE	1	7	28	2 923	1 254
WINE QUALITY	0	11	11	3 425	1 469
WAVEFORM	0	40	3	3 500	1 500
GISETTE	5 000	0	2	4 900	2 100
CONVEX	0	784	2	8 000	50 000
CIFAR-10-SMALL	3 072	0	10	10 000	10 000
MNIST BASIC	0	784	10	12 000	50 000
ROT. MNIST + BI	0	784	10	12 000	50 000
SHUTTLE	9	0	7	43 500	14 500
KDD09-APPENTENCY	190	40	2	35 000	15 000
CIFAR-10	3 072	0	10	50 000	10 000

verified experimentally that the more parallel runs, the faster this process identifies high-quality configurations; nevertheless, we restricted the version of Auto-WEKA studied here to run only 4 parallel processes, in order to study a setting typical for commonly used workstations.

5. EVALUATING AUTO-WEKA

We now describe an experimental study of the performance that can be achieved by Auto-WEKA on various datasets. After specifying our experiment environment, we demonstrate the importance of addressing the algorithm selection and the CASH problems, and establish baselines for them (Section 5.2). We then demonstrate Auto-WEKA’s ability to search its enormous hyperparameter space effectively to find algorithms and hyperparameters with low cross-validation error (Section 5.3). Then, we analyse its test performance and address concerns regarding overfitting (Section 5.4). Finally, we provide a synopsis of the classifiers and feature search/evaluators Auto-WEKA chose in our experiments (Section 5.5).

5.1 Experimental setup

We evaluated Auto-WEKA on 21 prominent benchmark datasets (see Table 3): 15 sets from the UCI repository [11]; the ‘convex’, ‘MNIST basic’ and ‘rotated MNIST with background images’ tasks used in [4]; the appentency task from the KDD Cup ’09; and two versions of the CIFAR-10 image classification task [19] (CIFAR-10-Small is a subset of CIFAR-10, where only the first 10 000 training data points are used rather than the full 50 000.) For datasets with a predefined training/test split, we used that split. Otherwise, we randomly split the dataset into 70% training and 30% test data. The test data was never seen by any optimization method; it was only used once in an offline analysis stage to evaluate the models found by the various optimization

is no empirical evidence that these methods outperform the simple approach we use here when the cost of evaluating hyperparameter configurations varies across the space.

methods. We denote datasets with at least 10 000 training data points as ‘large’ and all others as ‘small’.

All of our experiments were run on Linux machines with Intel Xeon X5650 six-core processors, running at 2.66GHz. All datasets had a RAM limit of 3GB for classification; if training a classifier ever exceeded this memory limit, the classifier job was terminated, returning a misclassification rate of 100%. An additional 1GB of RAM was allocated for the SMBO method. While these limits are somewhat arbitrary, we believe them to be reasonably close to the resource limitations faced by any user of machine learning algorithms. We also limited the training time for each evaluation of a learning algorithm on each fold, to ensure that the optimization method had a chance to explore the search space. Once this training budget for a fold is consumed, Auto-WEKA sends an interrupt to the learning algorithm to terminate as soon as possible, and the (partially) trained model is then evaluated on the validation set to determine the error estimate of the fold. This timeout was set to 150 minutes for classification and 15 minutes for feature search and evaluation in our experiments.⁴ For each dataset, we ran Auto-WEKA with each hyperparameter optimization algorithm with a total time budget of 30 hours. For each method, we performed 25 runs of this process with different random seeds and then – in order to simulate parallelization on a typical workstation – used bootstrap sampling to repeatedly select 4 random runs and report the performance of the one with best cross-validation performance.

In early experiments, we observed a few cases in which Auto-WEKA’s SMBO method picked hyperparameters that had excellent training performance, but turned out to generalize poorly. To enable Auto-WEKA to detect cases of overfitting, we partitioned its training set into two subsets: 70% for use inside the SMBO method, and 30% of validation data that is only looked at after the SMBO method has finished.

5.2 Algorithm Selection and CASH: Baseline Methods

Auto-WEKA aims to aid non-expert users of machine learning techniques. The simplest approach for selecting a classifier that is widely adopted amongst non-experts is to use a classifier merely based on its popularity or intuitive appeal, without any empirical consideration of alternatives. To quantify the difference such a choice can make, we consider the 39 WEKA classifiers (each with default hyperparameter settings) for each dataset, train each on the training set, and measure its accuracy on the test set. For each dataset, the second and third columns in Table 4 present the best and worst “oracle performance” of these classifiers on the test set. We observe that the gap between the best and worst classifier was huge, *e.g.* misclassification rates of 4.93% *vs* 99.24% on the Dorothea dataset. Even when the set of classifiers was restricted to a few popular ones (we considered neural networks, random forests, SVMs, AdaBoost, C4.5 decision trees, logistic regression, and KNN), this gap still exceeded 20% on 14 out of the 21 datasets. Furthermore, there was

⁴In preliminary experiments, only few models exceeded this timeout for the datasets studied here. [24] presents a promising approach for using runtime predictions in the expected improvement calculation to automatically drive the search away from excessively expensive models. We plan to incorporate this approach into future versions of Auto-WEKA.

no single method that achieved good performance across all datasets: every method was at least 22% worse than the best for at least one data set. We conclude that some form of algorithm selection is essential for achieving good performance.

A straight-forward algorithm selection method is to perform exhaustive 10-fold cross-validation on the training set and to return the classifier with the smallest average misclassification error across folds. We will refer to this method applied to the set of 39 WEKA classifiers as *Ex-Def*; it is the best choice that can be made among the 39 WEKA classifiers (with their default hyperparameters) based on an exhaustive cross-validation and will serve as a baseline to compare an optimal solution to the algorithm selection problem in WEKA to our solution for the CASH problem in WEKA.

More experienced users of machine learning algorithms would not only select between a fixed set of default algorithms, but would also consider different hyperparameter settings — for example by performing a grid search over the hyperparameter space of a single classifier (as, *e.g.*, implemented in WEKA⁵). Since different learning algorithms perform well for different problems, more experienced users optimally would also want to consider different hyperparameter settings for more than one learning algorithm. Therefore, a stronger baseline we will use is an approach that — in addition to the 39 WEKA default classifiers — considers various hyperparameter settings for all of WEKA’s 27 base classifiers. More precisely, this baseline considers a grid of hyperparameter settings for each of these 27 base classifier, and performs a random search [4] in the union of these grids (plus the 39 WEKA default classifiers). We refer to this baseline as *Random Grid* and note that — as an optimization approach in the joint space of algorithms and hyperparameter settings — it is a simple CASH algorithm. We executed this Random Grid search for all our datasets in parallel, using 400 CPU hours on average per dataset (at least 120 hours for each). Table 4 (columns 4 and 5) shows the best and worst “oracle performance” on the test set across these Random Grid classifiers. Comparing these performances to the default performance, we note that in most cases even WEKA’s best default algorithm could be improved by selecting better hyperparameter settings, sometimes rather substantially so: *e.g.*, in the rotated MNIST with background images task, random grid search offered a 6% improvement over Ex-Def. We conclude that choosing hyperparameter settings appropriately can lead to substantial differences in performance and that because of this fact even a relatively simple (albeit computationally expensive) approach for CASH can outperform algorithm selection by itself.

5.3 Results for Cross-Validation Performance

With 786 hierarchical hyperparameters, Auto-WEKA’s combined algorithm / hyperparameter space is very complex. We now study how effectively SMAC and TPE could search this space to optimize 10-fold cross-validation performance, and compare their performance to that of the Ex-Def and Random Grid methods defined in the previous section.

The middle portion of Table 4 reports the results. First, we note that random grid search over the hyperparameters of all base-classifiers yielded better results than Ex-Def in 14/21 cases (and tied in the remaining seven), which underlines

⁵This is implemented in WEKA’s `CVParameterSelection` class, see weka.wikispaces.com/Optimizing+parameters

Table 4: Performance on both 10-fold cross-validation and test data. Random Grid Search was run once for an average of 400 hours per dataset; for SMAC and TPE, we performed 25 runs of 30 hours each. We report results as median percent error rate across 100 000 bootstrap samples simulating 4 parallel runs. Ex-Def is deterministic. Test error rates are determined by training the selected model/hyperparameters on the entire 70% training data and computing the accuracy on the previously unused 30% test data. Boldface indicates the lowest error within a block of comparable methods. SC denotes correlation coefficients (see Section 5.4).

Dataset	Oracle Perf. (%)				10-Fold C.V. Performance (%)				Test Performance (%)				SC	
	EX-DEF		RAND. GRID		EX-DEF	RAND. GRID	AUTO-WEKA		EX-DEF	RAND. GRID	AUTO-WEKA		TPE	SMAC
	BEST	WORST	BEST	WORST			TPE	SMAC			TPE	SMAC		
DEXTER	7.78	52.78	5.00	58.33	10.20	7.48	9.90	5.48	8.89	5.00	9.44	7.22	0.82	0.25
GERMANCREDIT	26.00	38.00	25.00	63.67	22.45	22.45	21.43	19.59	27.33	27.33	27.67	28.33	0.31	0.20
DOROTHEA	4.93	99.24	4.93	99.24	6.03	6.03	6.93	5.52	6.96	6.96	6.96	6.38	0.95	0.40
YEAST	40.00	68.99	37.53	68.99	39.43	38.87	35.03	36.27	40.45	40.90	41.12	40.45	0.36	0.49
AMAZON	28.44	99.33	28.44	99.33	43.94	43.94	48.43	48.30	28.44	28.44	37.56	37.56	0.92	0.97
SECOM	7.87	14.26	7.66	40.64	6.25	6.12	6.25	5.34	8.09	8.30	7.87	7.87	-0.10	-0.56
SEMEION	8.18	92.45	6.08	92.45	6.52	6.52	6.91	4.86	8.18	8.18	8.18	5.03	0.84	0.73
CAR	0.77	29.15	0.19	31.66	2.71	1.54	0.94	0.71	0.77	0.19	0.00	0.58	0.12	0.75
MADELON	17.05	50.26	17.05	51.03	25.98	24.26	24.26	20.87	21.38	20.77	20.77	21.15	0.44	0.43
KR-VS-KP	0.31	48.96	0.21	51.04	0.89	0.70	0.45	0.32	0.31	0.52	0.52	0.31	0.22	0.32
ABALONE	73.18	84.04	72.55	89.23	73.33	72.45	72.20	71.76	73.18	72.79	72.71	73.02	0.15	0.10
WINE QUALITY	36.35	60.99	36.08	81.62	38.94	37.28	35.94	34.74	37.51	36.08	33.56	33.70	0.73	0.85
WAVEFORM	14.27	68.80	14.20	68.80	12.73	12.73	12.57	11.71	14.40	14.40	14.20	14.40	0.36	0.26
GISETTE	2.52	50.91	2.38	50.91	3.62	3.27	3.70	2.42	2.81	2.38	2.57	2.24	0.69	0.79
CONVEX	25.96	50.00	25.96	50.57	28.68	28.50	29.04	24.70	25.96	26.76	25.45	22.05	0.98	0.84
CIFAR-10-SMALL	65.91	90.00	64.54	90.00	66.59	65.11	57.97	57.76	65.91	64.54	56.65	55.93	0.93	0.80
MNIST BASIC	5.19	88.75	3.79	88.75	5.12	4.00	13.64	3.64	5.19	3.79	18.03	3.56	1.00	0.87
ROT. MNIST + BI	63.14	88.88	57.28	90.96	66.15	59.75	73.04	59.61	63.14	58.16	69.86	55.84	0.50	0.95
SHUTTLE	0.0138	20.8414	0.0069	20.8414	0.0328	0.0263	0.0230	0.0230	0.0138	0.0276	0.0069	0.0069	0.60	0.73
KDD09-APPENTENCY	1.74	6.97	1.64	54.08	1.88	1.88	1.88	1.75	1.75	1.77	1.74	1.74	0.89	1.00
CIFAR-10	64.27	90.00	64.27	90.00	65.54	65.54	66.68	63.21	64.27	64.27	64.80	62.39	0.33	0.69

the importance of not only choosing the right algorithm but of also setting its hyperparameters well. However, we note that this performance of random grid search is based on a very large time budget of an average of 400 CPU hours per dataset (650 CPU hours on average for each of the large datasets), making it a somewhat unrealistic alternative in practice. In contrast, Auto-WEKA was only run for 4×30 CPU hours per dataset, but still yielded substantially better performance than random grid search, outperforming it in 20/21 cases (and performing worse in one⁶). Comparing the two Auto-WEKA variants, SMAC outperformed TPE in 19/21 cases, with one tie. We note that sometimes Auto-WEKA’s performance improvements over the other methods were substantial, with relative reductions of the cross-validation error rate exceeding 15% in 12/21 cases.

We conclude that by searching Auto-WEKA’s combined algorithm/hyperparameter space, we can effectively find models with much better cross-validation performance than by grid search over WEKA’s base classifiers.

5.4 Results for Test Performance

The results just shown demonstrate that Auto-WEKA is effective at optimizing its given objective function; however, this is not sufficient to allow us to conclude that it fits models that generalize well. As the hypothesis space of a machine learning algorithm grows, so does its potential for overfitting. The use of cross-validation substantially increases Auto-WEKA’s robustness against overfitting, but

⁶For the Amazon data set, WEKA’s default implementation of support vector machines yielded a very strong error rate of 44%, which was below that of SMAC’s median performance. One of SMAC’s 25 runs actually reached an error rate of 36%, indicating that Auto-WEKA would be competitive given more time.

since its hypothesis space is much larger than that of standard classification algorithms, it is important to carefully study whether (and to what extent) overfitting poses a problem.

To evaluate generalization, we determined a combination of algorithm and hyperparameter settings A_λ by running Auto-WEKA as before (cross-validating on the training set), trained A_λ on the entire training set, and then evaluated the resulting model on the test set. The right portion of Table 4 reports the test performance obtained with all methods. Broadly speaking, similar trends held as for cross-validation performance: random grid search performed better than Ex-Def, and Auto-WEKA in turn outperformed random grid search. However, the performance differences were less pronounced: random grid search only yielded better results than Ex-Def in 9/21 cases, with 6/21 ties and 6/21 cases in which Ex-Def performed better. Auto-WEKA continued to outperform random grid search and Ex-Def in 15/21 cases, with 3 ties and 3 losses. Notably, Auto-WEKA performed best on all of the 11 largest datasets; we attribute this to the fact that the risk of overfitting decreases with dataset size. Sometimes, Auto-WEKA’s performance improvements over the other methods were substantial, with relative reductions of the test error rate exceeding 15% in 5/21 cases. Comparing the different Auto-WEKA variants, SMAC outperformed TPE on 12 datasets and tied on 3, with TPE performing better on 6. Even when compared to the unrealistic “oracle best” Random Grid classifier (which has access to the test set!), Auto-WEKA found algorithms/hyperparameters with a smaller error on 9/21 datasets, tied on 2, and was outperformed on the remaining 10. For the 10 largest datasets, it performed better in 8 cases, tied in 1, and only lost in 1.

As mentioned in our experimental setup, Auto-WEKA only used 70% of its training set during the optimization of cross-validation performance, reserving the remaining 30%

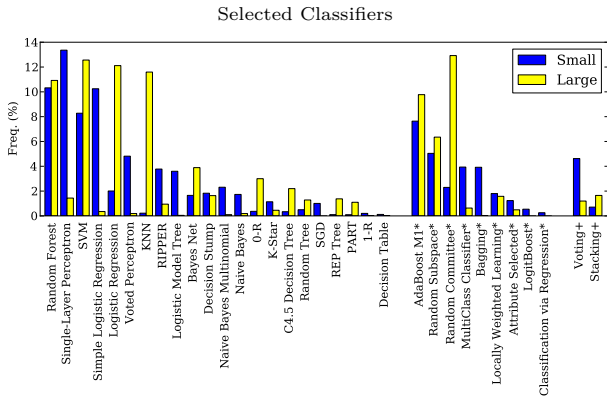


Figure 2: Distribution of chosen classifiers across the small and large datasets, aggregated across TPE, and SMAC, ranked on their frequency of being selected. Meta methods are marked by a * suffix, ensemble methods by a + suffix.

for assessing the risk of overfitting. At any point in time Auto-WEKA’s SMBO method keeps track of its *incumbent* (the hyperparameter configuration with the lowest cross-validation error rate seen so far). After its SMBO method has finished, Auto-WEKA extracts a trajectory of these incumbents from it and computes their generalization performance on the withheld 30% validation data. It then computes the Spearman rank coefficient between the sequence of training performances (evaluated by the SMBO method through cross-validation) and this generalization performance. The rightmost columns in Table 4 (labelled SC) show the average correlation coefficient for each run of Auto-WEKA. We note a general trend: as the absolute gap between cross-validation and test performance grows, this correlation coefficient decreases. The GermanCredit dataset is a good example where Auto-WEKA can signal that it only has low confidence in how well its chosen hyperparameters will generalize. We do note, however, that this weak signal has to be used with caution: there is no guarantee that large correlation coefficients yield a small gap and vice versa.

5.5 Classifiers Selected by Auto-WEKA

Figure 2 shows the distribution of classifiers chosen by our two Auto-WEKA variants (aggregated across runs and datasets - both TPE and SMAC produced similar results when considered individually). We note that no single classifier clearly dominated the others: the most frequently used classifiers (random forests, the single layer perceptron, and SVMs) were only selected in roughly 12% of all cases each, and most classifiers were selected in at least a few percent of the cases. Furthermore, the selected methods differed considerably between the large and small datasets, demonstrating the need for dataset-specific methods; for example, the large datasets benefitted more from meta methods than the small ones. A more detailed investigation of the top two meta-methods in Figure 3 (left) shows which base methods were chosen. Note that AdaBoostM1 frequently used the single layer perceptron on the small datasets, but never for the large ones, while the REP tree was highly popular for the large datasets. In the random subspace, the two most prominent methods were naive Bayes and the decision table. It is interesting to note that these two methods, as well as the REP tree frequently selected by AdaBoost, were not often

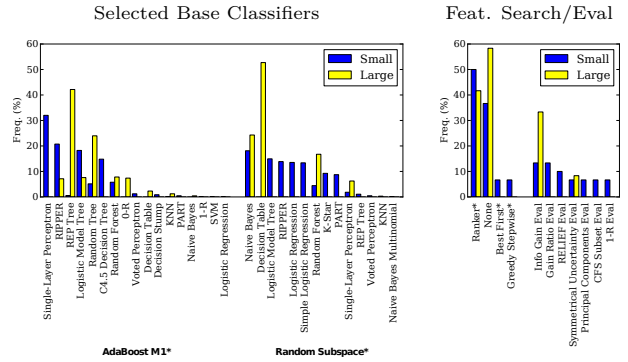


Figure 3: Left: distribution of chosen base classifiers for the two most frequently selected meta methods: AdaBoostM1 and random subspace. Right: distribution of chosen feature search and evaluator methods. Both plots are aggregated across TPE and SMAC, ranked on their frequency of being selected; None indicates that no feature selection was performed.

selected as a base classifier on their own. This underlines the importance of searching Auto-WEKA’s entire parameter space instead of, *e.g.*, restricting one’s attention to a small number of favourite base classifiers.

Figure 3 (right) provides a breakdown of the feature search and evaluation methods Auto-WEKA selected. Overall, it used these feature selection methods more often on the smaller datasets than on the larger ones, and if it did use a feature selection method it clearly favored the ranker method. All feature evaluators were used with roughly the same frequency for small datasets; in contrast, if Auto-WEKA performed feature selection for a large dataset it clearly favored the information gain evaluator. We note that Auto-WEKA’s data-dependent choices (based on its internal cross-validation evaluation) allow it to use feature selection as a regularization method for small data sets, while at the same time using all features to construct more complex hypotheses for large datasets.

6. CONCLUSION AND FUTURE WORK

In this work, we have shown that the daunting problem of combined algorithm selection and hyperparameter optimization (short: CASH) can be solved by a practical, fully automated tool. This is made possible by the use of recent Bayesian optimization techniques that iteratively build models of the algorithm/hyperparameter landscape and leverage these models to identify new points in the space that deserve investigation.

We built a tool, Auto-WEKA, that utilizes the full range of classification algorithms in WEKA and makes it easy for non-experts to build high-quality classifiers for given application scenarios. An extensive empirical comparison on 21 prominent datasets showed that Auto-WEKA often outperformed standard algorithm selection/hyperparameter optimization methods, especially on large datasets. We empirically compared two different optimizers for searching Auto-WEKA’s 786-dimensional parameter space and in the end recommend an Auto-WEKA variant based on the Bayesian optimization method SMAC [14]. We have written a freely-downloadable software package to make Auto-WEKA easy for end-users to access; it is available

at www.cs.ubc.ca/labs/beta/Projects/autoweika/.

We see several promising avenues for future work. First, Auto-WEKA still shows larger improvements in cross-validation performance than on test data, suggesting the investigation of more sophisticated methods for detecting and avoiding overfitting than our simple correlation-based approach. Second, we see potential value in extending our current approach to allow parameter sharing between classifiers used within ensemble methods, likely increasing their chance of being selected by Auto-WEKA. Finally, we could use our approach as an inner loop for training ensembles of machine learning algorithms by iteratively adding algorithms with maximal marginal contribution (this idea is conceptually related to the Hydra approach for constructing algorithm selectors [28]).

References

- [1] M. Adankon and M. Cheriet. Model selection for the LS-SVM. application to handwriting recognition. *Pattern Recognition*, 42(12):3264–3270, 2009.
- [2] Y. Bengio. Gradient-based optimization of hyperparameters. *Neural Computation*, 12(8):1889–1900, 2000.
- [3] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for Hyper-Parameter Optimization. In *Proc. of NIPS-11*, 2011.
- [4] J. Bergstra and Y. Bengio. Random search for hyperparameter optimization. *JMLR*, 13:281–305, 2012.
- [5] A. Biem. A model selection criterion for classification: Application to HMM topology optimization. In *Proc. of ICDAR-03*, pages 104–108. IEEE, 2003.
- [6] H. Bozdogan. Model selection and Akaike’s information criterion (AIC): The general theory and its analytical extensions. *Psychometrika*, 52(3):345–370, 1987.
- [7] P. Brazdil, C. Soares, and J. Da Costa. Ranking learning algorithms: Using IBL and meta-learning on accuracy and time results. *Machine Learning*, 50(3):251–277, 2003.
- [8] E. Brochu, V. M. Cora, and N. de Freitas. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. Technical Report UBC TR-2009-23 and arXiv:1012.2599v1, Department of Computer Science, University of British Columbia, 2009.
- [9] O. Chapelle, V. Vapnik, and Y. Bengio. Model selection for small sample regression. *Machine Learning*, 2001.
- [10] T. Desautels, A. Krause, and J. Burdick. Parallelizing exploration-exploitation tradeoffs with gaussian process bandit optimization. In *Proc. of ICML-12*, 2012.
- [11] A. Frank and A. Asuncion. UCI machine learning repository, 2010. URL: <http://archive.ics.uci.edu/ml>. University of California, Irvine, School of Information and Computer Sciences.
- [12] X. Guo, J. Yang, C. Wu, C. Wang, and Y. Liang. A novel LS-SVMs hyper-parameter selection based on particle swarm optimization. *Neurocomputing*, 71(16):3211–3215, 2008.
- [13] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. Witten. The WEKA data mining software: an update. *ACM SIGKDD Explorations Newsletter*, 11(1):10–18, 2009.
- [14] F. Hutter, H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration. *Proc. of LION-5*, pages 507–523, 2011.
- [15] F. Hutter, H. Hoos, K. Leyton-Brown, and T. Stütze. ParamILS: an automatic algorithm configuration framework. *JAIR*, 36(1):267–306, 2009.
- [16] F. Hutter, H. H. Hoos, and K. Leyton-Brown. Parallel algorithm configuration. In *Proc. of LION-6*, 2012.
- [17] D. R. Jones, M. Schonlau, and W. J. Welch. Efficient global optimization of expensive black box functions. *Journal of Global Optimization*, 13:455–492, 1998.
- [18] R. Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proc. of IJCAI-95*, pages 1137–1145, 1995.
- [19] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Master’s thesis, Department of Computer Science, University of Toronto*, 2009.
- [20] O. Maron and A. Moore. Hoeffding races: Accelerating model selection search for classification and function approximation. In *Proc. of NIPS-94*, pages 59–66, 1994.
- [21] A. McQuarrie and C. Tsai. *Regression and time series model selection*. World Scientific, 1998.
- [22] T. Schaul, J. Bayer, D. Wierstra, Y. Sun, M. Felder, F. Sehnke, T. Rückstieß, and J. Schmidhuber. PyBrain. *JMLR*, 2010.
- [23] M. Schonlau, W. J. Welch, and D. R. Jones. Global versus local search in constrained optimization of computer models. In N. Flournoy, W. Rosenberger, and W. Wong, editors, *New Developments and Applications in Experimental Design*, volume 34, pages 11–25. Institute of Mathematical Statistics, Hayward, California, 1998.
- [24] J. Snoek, H. Larochelle, and R. Adams. Opportunity cost in Bayesian optimization. In *NIPS Workshop on Bayesian Optimization, Sequential Experimental Design, and Bandits*, 2011. Published online.
- [25] J. Snoek, H. Larochelle, and R. P. Adams. Practical bayesian optimization of machine learning algorithms. In *Proc. of NIPS-12*, 2012.
- [26] N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proc. of ICML-10*, pages 1015–1022, 2010.
- [27] V. Strijov and G. Weber. Nonlinear regression model generation using hyperparameter optimization. *Computers & Mathematics with Applications*, 60(4):981–988, 2010.
- [28] L. Xu, H. H. Hoos, and K. Leyton-Brown. Hydra: Automatically configuring algorithms for portfolio-based selection. In *Proc. of AAAI-10*, pages 210–216, 2010.
- [29] P. Zhao and B. Yu. On model selection consistency of lasso. *JMLR*, 7:2541–2563, Dec. 2006.