

# Comprehensive analysis of gradient-based hyperparameter optimization algorithms

O. Y. Bakhteev · V. V. Strijov

Received: date / Accepted: date

**Abstract** The paper investigates hyperparameter optimization problem. Hyperparameters are the parameters of model parameter distribution. The adequate choice of hyperparameter values prevents model overfit and allows it to obtain higher predictive performance. Neural network models with large amount of hyperparameters are analyzed. The hyperparameter optimization for models is computationally expensive. The paper proposes modifications of various gradient-based methods to simultaneously optimize many hyperparameters. The paper compares the experiment results with the random search. The main impact of the paper is hyperparameter optimization algorithms analysis for the models with high amount of parameters. To select precise and stable models the authors suggest to use two model selection criteria: cross-validation and evidence lower bound. The experiments show that the models optimized using the evidence lower bound give higher error rate than the models obtained using cross-validation. These models also show greater stability when data is noisy. The evidence lower bound usage is preferable when the model tends to overfit or when the cross-validation is computationally expensive. The algorithms are evaluated on regression and classification datasets.

**Keywords** gradient descent · hyperparameter optimization · model selection · neural networks · classification · regression

The research was made possible by Government of the Russian Federation (Agreement 05.Y09.21.0018) and FASIE project 44116. This paper contains results of the project Statistical methods of machine learning, which is carried out within the framework of the Program "Center of Big Data Storage and Analysis" of the National Technology Initiative Competence Center. It is supported by the Ministry of Science and Higher Education of the Russian Federation according to the agreement between the M.V. Lomonosov Moscow State University and the Foundation of project support of the National Technology Initiative from 11.12.2018, No 13/1251/2018.

---

O. Y. Bakhteev  
Moscow Institute of Physics and Technology  
Tel.: +74991354163  
E-mail: bakhteev@phystech.edu

V. V. Strijov  
Moscow Institute of Physics and Technology  
FRCCSC of the Russian Academy of Sciences

## 1 Introduction

The paper analyzes hyperparameter optimization problem. *Model* is a function superposition, which solves a classification or regression problem. *Model hyperparameters* are the parameters of model parameter distribution.

The paper focuses on the neural network models. The parameter optimization problem is computationally expensive because of training loss non-convexity. The amount of model parameters can reach millions [27] and such models optimization requires multiple days [31]. The hyperparameter values can influence significantly the quality of models [15, 24]. These two facts make the hyperparameter optimization problem very important. The adequate choice of hyperparameter values also helps to prune redundant parameters to make the model more compact and stable [11, 16]. Comprehensive analysis of the hyperparameters and hyperparameter optimization algorithms leads to challenges addressed by operations research.

In this paper gradient-based algorithms are analyzed. They optimize a large amount of hyperparameters within a short time in contrast to gradient-free algorithms [23]. The complexity of hyperparameter optimization is comparable to the model parameter optimization complexity, therefore, the parameters and hyperparameters are optimized in a single procedure. As a baseline, a random search algorithm is used. The pros and cons of all the algorithms are illustrated in Table 1. Opposing to papers [10, 26, 21] this research does not focus on the hold-out cross-validation and analyzes algorithms in general. The analyzed algorithms are implemented as a toolbox [2]. The authors propose a number of analyzed algorithms modifications: we extend all the algorithms to work with validation loss that can contain hyperparameters, for DrMad algorithm [10] we add an additional parameter that regularizes a length of analyzed parameter update trajectory, for HOAG algorithm [26] we use stochastic gradient descent for large linear system solving because of high dimension of the linear system. The main impact of the paper is the analysis of algorithm behavior with various validation loss functions and a quality and stability analysis of the resulting models. The experiments are conducted using cross-validation and evidence lower bound as a model selection criterion. Opposing to the paper [26], where the optimization algorithms comparison can also be found, this paper presents the results on the large datasets, such as WISDM [19] and MNIST [20], where the number of hyperparameters is significant. The experimental results show that the gradient-based algorithms are significantly more effective than random search-based algorithms when the number of hyperparameters is large.

*Related works.* The papers [4, 6] proposed hyperparameters optimization strategies based on the random search. The current gold-standard methods [29, 14] proposed hyperparameter optimizations based on the probabilistic models construction for the optimal hyperparameter value prediction. These methods are ineffective whenever the number of hyperparameters is large [23].

The gradient-based hyperparameter optimization methods were suggested in [23, 8, 10, 26, 21]. The papers [10, 23] proposed the gradient-based methods that use reverse differentiation method. This method is similar to backpropagation. The main idea of this approach is to restore the full history of parameter updates in order to optimize hyperparameters using this history. In general this procedure is computationally expensive. The authors proposed to use only the last update of

Algorithm	Type	Pros	Cons
Random search	Stochastic	Easy to implement	Ineffective in high dimensions (curse of dimension)
Greedy [21]	Gradient-based	Can be used in inner model parameter optimization	Greedy and non-optimal: the algorithm is correct whenever the loss function Hessian is identical.
HOAG [26]	Gradient-based	Fast convergence	The algorithm requires additional parameter configuration: it is required to solve the linear system, which depends on the Hessian of a loss function. The error in linear system solution influences the final model quality.
DrMAD [10]	Gradient-based	Considers the trajectory of parameter updates and the optimization algorithm.	Can be instable because of gradient explosion or vanishing. The algorithm is correct when the trajectory of parameter update is linear, which is mostly not true for the complex non-convex models.

Table 1: The analyzed algorithms properties.

the parameters at each optimization iteration in [21]. This method can be considered as a greedy hyperparameter optimization algorithm. The authors optimized the model parameters using stochastic gradient descent with momentum in [23]. This allowed to use less memory for storing the parameter update history. The authors of the paper [10] proposed to linearize the history trajectory to effectively restore the update history. The authors of the paper [26] used an approximation of the exact hyperparameter gradient, which can be calculated when the problem satisfies some regular conditions.

For the hyperparameter optimization various model selection criteria can be used [22, 7]. One of them is the marginal likelihood or *evidence* [22, 7, 18, 30]. The authors of the papers [18, 30] considered the problem of model selection and hyperparameter optimization for the linear regression problem. One of the marginal likelihood approximation methods is an *evidence lower bound*, which can be obtained using variational inference [7]. A stochastic version of the variational inference was proposed in [13]. The authors of the paper [28] analyzed the relation between the gradient-based evidence lower bound estimations and MCMC methods. An alternative model selection criterion is a minimum description length [12] that characterizes statistical complexity of the model. An evidence lower bound estimation method for the neural networks was proposed in [11]. The authors analyzed the relation between the evidence lower bound and minimum description length.

The other model selection criterion is the cross-validation [1, 18]. One of this criterion problem is its high computational cost [32, 17]. The authors of the pa-

per [3] analyzed the variance of cross-validation model performance estimation. In [18] the authors compared hyperparameter values obtained during hyperparameter optimization using various model selection criteria.

## 2 Problem statement

There is given a dataset

$$\mathfrak{D} = \{(\mathbf{x}_i, y_i)\}, i = 1, \dots, m. \quad (1)$$

It contains the object matrix  $[\mathbf{x}_1, \dots, \mathbf{x}_m]^\top = \mathbf{X}$  and the label vector  $[y_1, \dots, y_m]^\top = \mathbf{y}$ . The label  $y_i$  is in finite set,  $y_i \in \mathbb{Y} = \{1, \dots, Z\}$ , in case of  $Z$ -class classification problem. For the regression problem the label  $y_i$  is in the real-value subset  $y_i \in \mathbb{Y} \subset \mathbb{R}$ .

A model  $\mathbf{f}(\mathbf{w}, \mathbf{x})$  predicts the variable  $y_i$ :

$$f : \mathbb{R}^u \times \mathbb{R}^n \rightarrow \mathbb{Y}, \quad \mathbf{w} \in \mathbb{R}^u.$$

It is differentiable with respect to parameters  $\mathbf{w}$ .

Introduce the probabilistic interpretation of the model  $f(\mathbf{w}, \mathbf{x})$  for the classification and regression problems.

*Regression problem.* Suppose the dependent variable  $y$  is normally distributed:

$$\mathbf{y} = \mathcal{N}(\mathbf{f}, \mathbf{I}), \quad \text{where } \mathbf{f} = \mathbf{f}(\mathbf{w}, \mathbf{X}) = [f(\mathbf{w}, \mathbf{x}_1), \dots, f(\mathbf{w}, \mathbf{x}_m)]^\top. \quad (2)$$

Define the likelihood function  $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$ :

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\frac{m}{2} \log(2\pi) - \frac{1}{2} (\mathbf{y} - \mathbf{f}(\mathbf{w}, \mathbf{X}))^\top \mathbf{I} (\mathbf{y} - \mathbf{f}(\mathbf{w}, \mathbf{X})).$$

*Classification problem.* For the 2-class classification problem suppose the dependent variable  $y$  is distributed binomially:

$$\mathbf{y} = \mathcal{B}(1 - \mathbf{f}, \mathbf{f}), \quad (3)$$

where  $\mathbf{f}$  is the probability that objects from the matrix  $\mathbf{X}$  belong to the first class. For the multiclass classification problem the dependent variable  $y$  is distributed multinomially and  $r$ -th component  $f_r$  is the probability that objects from the matrix  $\mathbf{X}$  belong to the class  $r$ . Define the likelihood function

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \sum_{\mathbf{x}, \mathbf{y} \in \mathbf{X}, \mathbf{y}} \sum_{r=1}^Z [y = r] \log f_r(\mathbf{w}, \mathbf{x}).$$

For both classification (3) and regression (2) problems define the prior distribution  $p(\mathbf{w}|\mathbf{A})$ :

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}), \quad (4)$$

where  $\mathbf{A}^{-1} = \text{diag}[\alpha_1, \dots, \alpha_u]^{-1}$  is the covariance diagonal matrix. The hypotheses (2), (3) and (4) do not contradict each other since the normal distribution is unbounded [9].

In order to describe the general scheme of model optimization procedure denote a vector of model parameters by  $\mathbf{w}$ . Denote by  $\boldsymbol{\theta} \in \mathbb{R}^s$  a vector of parameters of model parameters distribution. In the simple case the  $\boldsymbol{\theta}$  is a concatenation of vectors of parameters for multiple model instances. The details of  $\boldsymbol{\theta}$  construction are given in the model selection methods description.

Introduce an example of hyperparameter optimization using the coherent Bayesian inference. Optimize model parameters  $\boldsymbol{\theta} = \mathbf{w}$  according to the model parameter distribution  $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \mathbf{A})$ :

$$\hat{\boldsymbol{\theta}} = \arg \max(-L(\boldsymbol{\theta}, \mathbf{A})) = p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \mathbf{A}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{A})}{p(\mathbf{y}|\mathbf{X}, \mathbf{A})}. \quad (5)$$

Calculate the hyperparameter  $\mathbf{A}$  posterior distribution:

$$p(\mathbf{A}|\mathbf{X}, \mathbf{y}) \propto p(\mathbf{y}|\mathbf{X}, \mathbf{A})p(\mathbf{A}),$$

where  $\propto$  means proportionality.

Suppose  $p(\mathbf{A})$  is uniform on a large interval. The hyperparameter optimization problem is:

$$Q(\boldsymbol{\theta}, \mathbf{A}) = p(\mathbf{y}|\mathbf{X}, \mathbf{A}) = \int_{\mathbf{w} \in \mathbb{R}^u} p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{A}) \rightarrow \max_{\text{diag}(\mathbf{A})=[\alpha_1, \dots, \alpha_u] \in \mathbb{R}^n}. \quad (6)$$

Formulate the hyperparameter optimization problem. Given functions  $L$  and  $Q$  that characterize the model selection method. The loss function  $L$  optimizes parameters  $\boldsymbol{\theta}$ . The validation function  $Q$  evaluates the model predictive performance. The problem is to optimize the parameters  $\boldsymbol{\theta}$  and the hyperparameters  $\mathbf{A}$ :

$$\hat{\mathbf{A}} = \arg \max_{[\alpha_1, \dots, \alpha_u] \in \mathbb{R}^n} Q(\hat{\boldsymbol{\theta}}(\mathbf{A}), \mathbf{A}), \quad (7)$$

$$\hat{\boldsymbol{\theta}}(\mathbf{A}) = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^s} L(\boldsymbol{\theta}, \mathbf{A}). \quad (8)$$

Introduce the functions  $L, Q$  and  $\boldsymbol{\theta}$  for various model selection methods.

*Basic method.* Optimize the model parameters  $\boldsymbol{\theta}$  and hyperparameters  $\mathbf{A}$  using the whole dataset  $\mathcal{D}$  and the same function for both optimization and model evaluation with one model instance,  $L = -Q$ :

$$Q(\boldsymbol{\theta}, \mathbf{A}) = -L(\boldsymbol{\theta}, \mathbf{A}) = \log p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{A}) = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}|\mathbf{A}). \quad (9)$$

Describe two model selection methods that prevent model overfitting.

*Cross-validation [1].* Split the dataset  $\mathfrak{D}$  into  $k$  equal parts:  $\mathfrak{D} = \mathfrak{D}_1 \sqcup \dots \sqcup \mathfrak{D}_k$ .

Run  $k$  model optimizations for each subset. The vector  $\theta$  is a concatenation of model parameters  $\mathbf{w}_k$  for each optimization instance:

$$\theta = [\mathbf{w}_1, \dots, \mathbf{w}_k].$$

Construct the loss function  $L$  as an average negative log posterior probability over the remaining part of each split  $\mathfrak{D}_1, \dots, \mathfrak{D}_k$ :

$$L(\theta, \mathbf{A}) = -\frac{1}{k} \sum_{c=1}^k \left( \frac{k}{k-1} \log p(\mathbf{y} \setminus \mathbf{y}_c | \mathbf{X} \setminus \mathbf{X}_c, \mathbf{w}_c) + \log p(\mathbf{w}_c | \mathbf{A}) \right). \quad (10)$$

Construct the validation function  $Q$  as an average log likelihood over  $k$  splits:

$$Q(\theta, \mathbf{A}) = \frac{1}{k} \sum_{c=1}^k k \log p(\mathbf{y}_c | \mathbf{X}_c, \mathbf{w}_c). \quad (11)$$

*Evidence Lower Bound [11].* As in the basic method description let  $L = -Q$ . It is an evidence lower bound:

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{X}, \mathbf{A}) &\geq -D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{A})) + \int_{\mathbf{w}} q(\mathbf{w}) \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \mathbf{A}) d\mathbf{w} \approx \quad (12) \\ &\approx \sum_{i=1}^m \log p(y_i | \mathbf{x}_i, \mathbf{w}_i) - D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{A})) = -L(\theta, \mathbf{A}) = Q(\theta, \mathbf{A}), \end{aligned}$$

where  $q$  is an auxiliary distribution called variational distribution. Let  $q$  be a normal distribution:

$$q \sim \mathcal{N}(\boldsymbol{\mu}_q, \mathbf{A}_q^{-1}), \quad (13)$$

where  $\mathbf{A}_q = \text{diag}[\alpha_1^q, \dots, \alpha_u^q]^{-1}$  is a diagonal covariance matrix and  $\boldsymbol{\mu}_q$  is a mean vector. The divergence  $D_{\text{KL}}$  between 2 Gaussian variables is

$$D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{f})) = \frac{1}{2} (\text{Tr}[\mathbf{A} \mathbf{A}_q^{-1}] + (\boldsymbol{\mu} - \boldsymbol{\mu}_q)^\top \mathbf{A} (\boldsymbol{\mu} - \boldsymbol{\mu}_q) - u + \ln |\mathbf{A}^{-1}| - \ln |\mathbf{A}_q^{-1}|).$$

The vector of optimized parameters  $\theta$  is a vector of variational distribution  $q$  parameters:

$$\theta = [\text{diag}(\mathbf{A}_q), \mu_1, \dots, \mu_u].$$

### 3 Gradient-based hyperparameter optimization methods

In this section we analyze hyperparameter optimization methods based on gradient descent. Suppose the parameters  $\theta$  are also optimized using gradient-based methods.

**Definition 1** A stochastic gradient descent operator  $T$  estimates the parameters  $\theta'$  using their previous values  $\theta$  with random subset of dataset:

$$\theta' = T(\theta, \mathbf{A}, \mathfrak{D}) = \theta - \gamma \nabla L(\theta, \mathbf{A})_{\mathfrak{D}=\hat{\mathfrak{D}}},$$

where  $\hat{\mathfrak{D}}$  is a random subset of  $\mathfrak{D}$ .

Algorithm	Type	Optimization iteration complexity	Correctness suppositions
Random search	stochastic	$O(\eta s  \hat{\mathcal{D}} )$	-
Greedy [21]	gradient-based	$O(\eta u s  \hat{\mathcal{D}} )$	$\mathbf{H}(\boldsymbol{\theta}) = \mathbf{I}$
HOAG [26]	gradient-based	$O(\eta u s  \hat{\mathcal{D}}  + o)$ , where $o$ is a complexity of linear equation solution	first derivatives of $Q$ and second derivatives of $L$ are Lipschitz functions; $\det \mathbf{H} \neq 0$ ;
DrMAD [10]	gradient-based	$O(\eta u s  \hat{\mathcal{D}} )$	Parameter trajectory $\boldsymbol{\theta} = \boldsymbol{\theta}_0, \dots, \boldsymbol{\theta}_\eta$ is linear

Table 2: The analyzed algorithms complexity and correctness. The model is a multilayer perceptron optimized using backpropagation with the basic model selection method (9).

Optimize the parameters  $\boldsymbol{\theta}$  with  $\eta$  steps of stochastic gradient descent:

$$\hat{\boldsymbol{\theta}} = T \circ T \circ \dots \circ T(\boldsymbol{\theta}_0, \mathbf{A}) = T^\eta(\boldsymbol{\theta}_0, \mathbf{A}), \quad T(\boldsymbol{\theta}, \mathbf{A}) = \boldsymbol{\theta} - \gamma \nabla L(\boldsymbol{\theta}, \mathbf{A}), \quad (14)$$

where  $\gamma$  is the learning rate,  $\boldsymbol{\theta}_0$  is the initial values of the vector  $\boldsymbol{\theta}$ .

Redefine the optimization problem according to the definition of operator  $T$ :

$$\hat{\mathbf{A}} = \arg \max_{\text{diag}(\mathbf{A}) \in \mathbb{R}^n} Q(T^\eta(\boldsymbol{\theta}_0, \mathbf{A})). \quad (15)$$

The general scheme of the hyperparameter optimization is the following:

1. in range from 1 to  $l$ , where  $l$  is the number of the hyperparameter optimization iterations:
2. initialize parameters  $\boldsymbol{\theta}$ ;
3. solve the optimization problem (15) and obtain the new hyperparameter values  $\hat{\mathbf{A}}$ ;
4. set  $\mathbf{A} = \hat{\mathbf{A}}$ .

One of the methods the problem (15) solution is to use a gradient descent not only for the parameters  $\boldsymbol{\theta}$ , but also for the hyperparameters  $\mathbf{A}$ . The exact gradient calculation  $\nabla_{\mathbf{A}} Q(T^\eta(\boldsymbol{\theta}_0, \mathbf{A}))$  is intractable because of the inner optimization operator  $T$ . Therefore the number of approximated solutions were presented. Their main properties are listed in Table 2. An example of the parameter and hyperparameter update with different hyperparameter optimization algorithms is illustrated in Fig. 1.

*Greedy algorithm.* Update the hyperparameter  $\mathbf{A}$  using gradient descent, which depends only on the last update of parameters  $\boldsymbol{\theta}$ . Optimize the hyperparameters and parameters in a single optimization procedure. Update the hyperparameter  $\mathbf{A}$  at every iteration:

$$\hat{\mathbf{A}} = \mathbf{A} + \gamma_{\mathbf{A}} \nabla_{\mathbf{A}} Q(T(\boldsymbol{\theta}, \mathbf{A}), \mathbf{A}) = \mathbf{A} + \gamma_{\mathbf{A}} \nabla_{\mathbf{A}} Q(\boldsymbol{\theta} - \gamma \nabla L(\boldsymbol{\theta}, \mathbf{A}), \mathbf{A}),$$

where  $\gamma_{\mathbf{A}}$  is the learning rate for the hyperparameter optimization.

*HOAG.* Obtain approximate hyperparameter gradient values  $\nabla_{\mathbf{A}}Q(T^\eta(\boldsymbol{\theta}_0, \mathbf{A}))$  using the following formula approximation of the of the gradient:

$$\nabla_{\mathbf{A}}Q(T^\eta(\boldsymbol{\theta}_0, \mathbf{A})) = \nabla_{\mathbf{A}}Q(\boldsymbol{\theta}, \mathbf{A}) - (\nabla_{\mathbf{A}}\nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta}, \mathbf{A}))^\top \mathbf{H}(\boldsymbol{\theta})^{-1} \nabla_{\boldsymbol{\theta}}Q(\boldsymbol{\theta}, \mathbf{A}),$$

where  $\mathbf{H}$  is the Hessian of the function  $L$  with respect to the parameters  $\boldsymbol{\theta}$ .

The algorithm of hyperparameter gradient  $\nabla_{\mathbf{A}}Q$  approximation:

1. Optimize the parameters  $\boldsymbol{\theta} = T^\eta(\boldsymbol{\theta}_0, \mathbf{A})$ .
2. Solve linear system for a vector  $\boldsymbol{\lambda}$ :  $\mathbf{H}(\boldsymbol{\theta})\boldsymbol{\lambda} = \nabla_{\boldsymbol{\theta}}Q(\boldsymbol{\theta}, \mathbf{A})$ .
3. Compute the approximate hyperparameter gradient:  $\hat{\nabla}_{\mathbf{A}}Q = \nabla_{\mathbf{A}}Q(\boldsymbol{\theta}, \mathbf{A}) - \nabla_{\mathbf{A}}\nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta}, \mathbf{A})^\top \boldsymbol{\lambda}$ .

The final update rule is

$$\hat{\mathbf{A}} = \mathbf{A} + \gamma_{\mathbf{A}} \hat{\nabla}_{\mathbf{A}}Q. \quad (16)$$

The computational cost of the Hessian  $\mathbf{H}(\boldsymbol{\theta})$  evaluation from step 2 is high. In this paper we use stochastic gradient descent for solving this linear system.

*DrMad.* For the hyperparameter gradient evaluation restore the full history of  $\eta$  parameters updates starting from the initial value  $\boldsymbol{\theta}_0$ . For the simplification of this procedure suppose that the trajectory of parameters  $\boldsymbol{\theta}$  update is linear:

$$\boldsymbol{\theta}^\tau = \boldsymbol{\theta}_0 + \frac{\tau}{\eta} T(\boldsymbol{\theta}). \quad (17)$$

The algorithm of approximate hyperparameter gradient calculation:

1. Optimize parameters  $\boldsymbol{\theta} = T^\eta(\boldsymbol{\theta}_0, \mathbf{A})$ .
2. Let  $\hat{\nabla}_{\mathbf{A}} = \nabla_{\mathbf{A}}Q(\boldsymbol{\theta}, \mathbf{A})$ .
3. Let  $\hat{\nabla}_{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}}Q(\boldsymbol{\theta}, \mathbf{A})$ .
4. In range from  $\tau = \eta$  to 1:
5. Compute  $\boldsymbol{\theta}^\tau$  (17).
6.  $d\mathbf{v} = \gamma \hat{\nabla}_{\boldsymbol{\theta}}$ .
7.  $\hat{\nabla}_{\mathbf{A}} = \hat{\nabla}_{\mathbf{A}} + d\mathbf{v} \nabla_{\mathbf{A}} \nabla_{\boldsymbol{\theta}} L$ .
8.  $\hat{\nabla}_{\boldsymbol{\theta}} = \hat{\nabla}_{\boldsymbol{\theta}} - d\mathbf{v} \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} L$ .

The final update rule is analogous to (16). The algorithm is instable when the learning rate  $\gamma$  is large [10]. For the stability of the algorithm we discard first 5% of values  $\boldsymbol{\theta}$  and calculate only each  $\tau_{\text{step}}$  step of parameter updates history:

$$\boldsymbol{\theta}^\tau = \boldsymbol{\theta}_{\tau_0} + \frac{\tau}{\eta} T(\boldsymbol{\theta}), \quad \tau \in \{\tau_0, \tau_0 + \tau_{\text{step}}, \dots, \eta - \tau_{\text{step}}, \eta\}, \quad (18)$$

where  $\tau_0 = \lceil 0.05 \cdot \eta \rceil$ .



## 4 Experiments

For the evaluation of the analyzed algorithms we conduct a series of computational experiments. A synthetic dataset [2], WISDM [19] and MNIST [20] datasets are used. All datasets are split into Train and Test subsets. The algorithms are optimized on the Train subsets and evaluated on the Test subsets. For each dataset and each hyperparameter optimization algorithm we run the optimizations 5 times. The results are averaged. We use the following evaluation criteria.

1. Quality: the best value of  $Q$ :  $\hat{Q} = \max_{j \in \{1, \dots, l\}} Q^j$ .
2. Convergence: the number of iterations to have a validation value greater than 99% of the best value  $\hat{Q}$ :

$$\arg \min_j : \frac{Q^j - Q^0}{\hat{Q} - Q^0} \geq 0.99,$$

where  $Q^0$  is the value of  $Q$  before the hyperparameter optimization.

3. Error function  $E$ :

$$E = \text{RMSE}(\mathbf{X}_{\text{test}}, \mathbf{y}_{\text{test}}) = \left( \frac{1}{|\mathbf{X}_{\text{test}}|} \sum_{\mathbf{x} \in \mathbf{X}_{\text{test}}, y \in \mathbf{y}_{\text{test}}} (f(\mathbf{x}, \mathbf{w}) - y) \right)^{\frac{1}{2}}$$

for the regression problem and

$$E = \text{Acc}(\mathbf{X}_{\text{test}}, \mathbf{y}_{\text{test}}) = 1 - \frac{1}{|\mathbf{X}_{\text{test}}|} \sum_{\mathbf{x} \in \mathbf{X}_{\text{test}}, y \in \mathbf{y}_{\text{test}}} [f(\mathbf{x}, \mathbf{w}) \neq y]$$

for the classification problem.

4. The error function  $E_\sigma$  for the model with noisy dataset:

$$\text{RMSE}_\sigma = \text{RMSE}(\mathbf{X}_{\text{test}} + \boldsymbol{\varepsilon}, \mathbf{y}_{\text{test}}),$$

$$\text{Acc}_\sigma = \text{Acc}(\mathbf{X}_{\text{test}} + \boldsymbol{\varepsilon}, \mathbf{y}_{\text{test}}), \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma \mathbf{I}).$$

The random search is used as a baseline. We set the number of the random search iterations to the number  $l$  of gradient-based hyperparameter optimization algorithms iterations. It is set to  $l = 50$  for the synthetic dataset and WISDM dataset,  $l = 25$  for the MNIST dataset. We use two types of loss function  $L$  and validation function  $Q$ :

1. evidence lower bound for both  $L$  and  $Q$  (12);
2. cross-validation functions (10),(11).

For the neural network models we set all the hyperparameters equal for each layer when use cross-validation (10). We use the full diagonal parameterization of the hyperparameters for all the linear models and for the neural network models when we optimize evidence lower bound (12).

The hyperparameters are initialized with uniform distribution  $\mathcal{U}(a, b)$ , where  $a = -2, b = 10$  for the synthetic dataset and  $a = -4, b = 10$  for the WISDM and MNIST datasets.

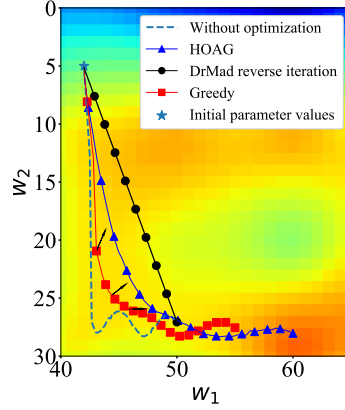


Fig. 1: An example of parameter update trajectories. The color displays the value of the validation function  $Q$ . Greedy algorithm optimizes hyperparameter during the parameter optimization, therefore it has the light blue trajectory between the optimized green trajectory of HOAG algorithm and the dark blue trajectory of parameters without hyperparameter optimization. DrMAD uses a linearized dashed parameter trajectory during hyperparameter optimization procedure.

We calibrate the hyperparameter learning rate  $\gamma_{\mathbf{A}}$  for each algorithm using grid search:  $\{r \cdot 10^s, s \leq 1, r \in \{1, 25, 50, 75\}\}$ . The largest value of  $\gamma_{\mathbf{A}}$  is chosen if the final hyperparameter value  $\mathbf{A}$  satisfied the condition:

$$a_{\min} \leq \min(\mathbf{A}), \quad \max(\mathbf{A}) \leq b_{\max},$$

where  $a_{\min} = -2.5, b_{\max} = 10.5$  for the synthetic dataset and  $a_{\min} = -5, b_{\max} = 11$  for the WISDM and MNIST datasets. We calibrate  $\gamma_{\mathbf{A}}$  with a small amount of iterations:  $l = 50$  for the synthetic dataset,  $l = 10$  for the WISDM and  $l = 5$  for the MNIST dataset. Whenever each algorithm shows instability in further experiments (gradient explosion or overflow) we lower the value of  $\gamma_{\mathbf{A}}$ . The parameter  $\tau_{\text{step}}$  is set to 1 for synthetic dataset and WISDM and  $\tau_{\text{step}} = 10$  for the MNIST dataset. The overall results are presented in Table 3. The implementation of the proposed algorithms is available at [2].

*Synthetic dataset.* The synthetic dataset is generated according to the rule:

$$\mathbf{y} = \mathbf{x} + \varepsilon, \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where  $m = 40, \quad n = 1$ .

The regression model  $\mathbf{f}$  uses the following features:  $\{\mathbf{x}^0, \dots, \mathbf{x}^9, \sin(\mathbf{x}), \cos(\mathbf{x})\}$ . The unregularized regression model with such ratio of object number and feature number tends to overfit significantly. The goal of this experiment is to analyze if the hyperparameter optimization could prevent such models from overfitting.

Fig. 2 plots the resulting polynomials. Since the model is overparametrized, the optimization without any overfitting control tends to overfit drastically. The

real dependency between  $\mathbf{y}$  and  $\mathbf{X}$  is linear. As we can see, the evidence lower bound criterion (12) gives non-overfitted polynomials, which parameters are close to the linear models.

*WISDM.* The WISDM dataset contains a set of records from accelerometer. Each record has three coordinates that correspond to the accelerometer axes. As a set of features we use first 199 records from each 200 sequential records. We use  $l_2$  norm of the 200th records as a label and required to predict.

A neural network with 10 neurons on the hidden layer is used:

$$\mathbf{f} = \mathbf{W}_2 \cdot \text{RELU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2,$$

where  $\mathbf{W}_1, \mathbf{b}_1$  — hidden layer parameters,  $\mathbf{W}_2, \mathbf{b}_2$  — output layer parameters,

$$\text{RELU}(\mathbf{x}) = \max(\mathbf{0}, \mathbf{x}).$$

Fig. 3 shows the plots with RMSE and best validation loss value  $\hat{Q}$  for the WISDM dataset. DrMad and HOAG algorithms shows significantly lower results than greedy algorithm. The random search shows good results in case of cross-validation, when the number of hyperparameters is small. When the number of hyperparameters is large the greedy and HOAG algorithm shows good results. HOAG requires more iteration number  $l$  than the greedy algorithm for the convergence.

*MNIST.* The MNIST dataset contains a set of images of handwritten digits. We use a neural network with 300 neurons on the hidden layer and softmax output.

Fig. 4 plots the error  $E$  and quality  $\hat{Q}$ . For the evidence lower bound criterion, the models with highest  $\hat{Q}$  value have also highest error  $E$ . Nevertheless, these models show the best results for the stability, when the noise in dataset is high. Fig. 5 plots the error  $E_\sigma$ . Models with highest evidence lower bound have the lowest error when the Test dataset is noisy. We interpret this as an ability to select the model with highest generalization ability.

The experiments show that the gradient-based methods give better results than the random search when the number of hyperparameters is large. The best results were obtained by the greedy algorithm. We find that the hyperparameter optimization is instable when using DrMad algorithm, therefore its learning rate  $\gamma_{\mathbf{A}}$  is set to low values.

## 5 Discussion

The experiments show that each algorithm perform effectively and therefore the appropriate hyperparameter optimization method should rely on the amount of hyperparameters and the specific of the problem.

When dealing with small amount of hyperparameters random search shows the best results since search procedure can be employed more effectively than gradient-based optimization in low-dimensional hyperparameters space. For the high-dimensional hyperparameter space both HOAG and greedy algorithms show good performance. HOAG algorithm is more preferable if the model optimization

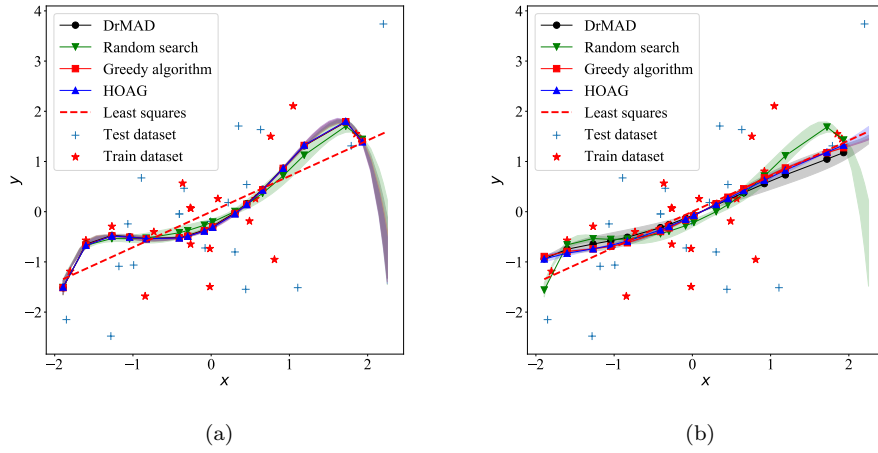


Fig. 2: Resulting models for the synthetic dataset: a — cross-validation, b — evidence lower bound

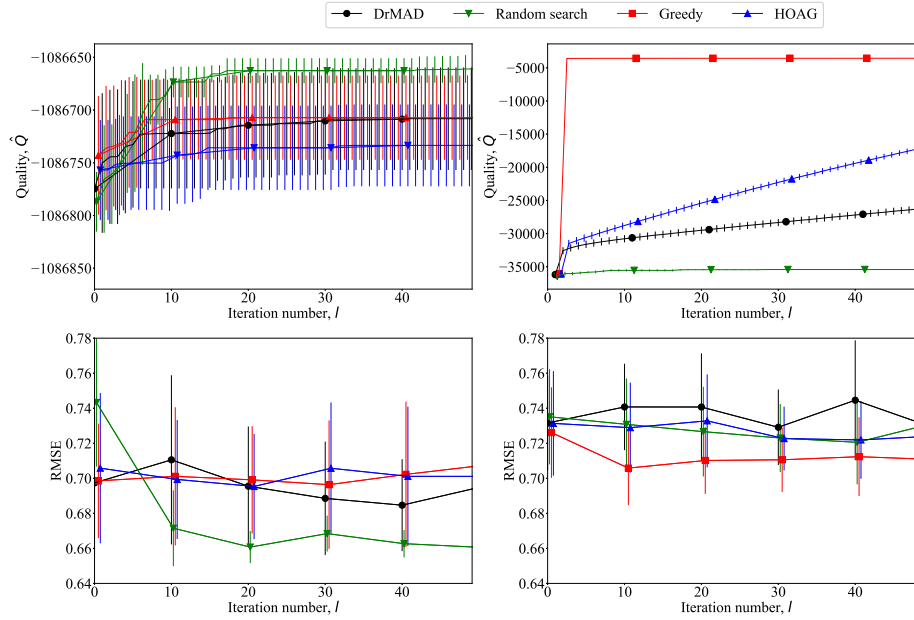


Fig. 3: WISDM, best validation value  $\hat{Q}$  and RMSE for cross-validation (left) and evidence lower bound (right)

Algorithm	$L, Q$	$Q(\theta, \mathbf{A})$	Convergence	E	$E_{0.25}$	$E_{0.5}$
<i>Synthetic</i>						
Random search	(10)	<b>-171.6</b>	<b>26.2</b> $\pm$ <b>20.0</b>	<b>1.367</b>	<b>1.410</b>	<b>1.555</b>
Greedy	(10)	-172.5	30.0 $\pm$ 24.5	1.421	1.439	1.536
DrMAD	(10)	-174.1	40.2 $\pm$ 16.1	1.403	1.424	<b>1.512</b>
HOAG	(10)	-174.7	29.4 $\pm$ 24.0	<b>1.432</b>	1.463	1.553
Random Search	(12)	-63.5	32.4 $\pm$ 18.7	1.368	1.426	1.546
Greedy	(12)	-25.5	<b>1.2</b> $\pm$ <b>0.4</b>	1.161	1.174	1.193
DrMAD	(12)	<b>-25.1</b>	10.6 $\pm$ 0.8	1.157	1.163	1.184
HOAG	(12)	-25.8	10.8 $\pm$ 1.5	<b>1.141</b>	<b>1.149</b>	<b>1.177</b>
<i>WISDM</i>						
Random search	(10)	<b>-1086661.1</b>	22.0 $\pm$ 19.3	<b>0.660</b>	<b>0.670</b>	<b>0.690</b>
Greedy	(10)	-1086707.1	<b>15.4</b> $\pm$ <b>17.2</b>	0.707	0.723	0.769
DrMAD	(10)	-1086708.2	29.2 $\pm$ 8.0	0.694	0.708	0.742
HOAG	(10)	-1086733.5	28.2 $\pm$ 7.13	0.701	0.724	0.753
Random search	(12)	-35420.4	14.4 $\pm$ 7.8	0.732	0.755	0.785
Greedy	(12)	<b>-3552.9</b>	<b>1.0</b> $\pm$ <b>0.0</b>	<b>0.702</b>	<b>0.730</b>	<b>0.767</b>
DrMAD	(12)	-26091.4	50.0 $\pm$ 0.0	0.729	0.753	0.816
HOAG	(12)	-16566.6	49.0 $\pm$ 0.0	0.733	0.755	0.801
<i>MNIST</i>						
Random search	(10)	-3236.4	7.8 $\pm$ 1.9	0.981	<b>0.966</b>	<b>0.866</b>
Greedy	(10)	<b>-3416.7</b>	10.8 $\pm$ 10.4	0.979	0.962	0.860
DrMAD	(10)	-3469.0	17.0 $\pm$ 5.6	<b>0.982</b>	0.962	0.831
HOAG	(10)	-3748.6	<b>8.6</b> $\pm$ <b>7.3</b>	0.980	0.961	0.853
Random search	(12)	-1304556.4	14.2 $\pm$ 5.7	<b>0.982</b>	0.943	0.814
Greedy	(12)	<b>-11136.2</b>	<b>1.0</b> $\pm$ <b>0.0</b>	0.977	<b>0.952</b>	<b>0.884</b>
DrMAD	(12)	-1305432.9	24.6 $\pm$ 0.5	<b>0.982</b>	0.941	0.813
HOAG	(12)	-280061.6	24.0 $\pm$ 0.0	0.981	0.943	0.819

Table 3: Experiment results

problem is expensive. On the other hand we can schedule greedy hyperparameter optimization to make it less expensive as in [21].

DrMad algorithm showed rather poor results on the MNIST and WISDM datasets. Perhaps, it is so because of high learning rate  $\gamma$  used in experiments. The large value of the learning rate can make DrMad algorithm instable. Two improvements can be proposed. We can use more stable optimization like Adam or AdaGrad for both parameter and hyperparameter optimization. The second improvement was proposed in [10]: we can use more complicated parameter trajectory approximation to make it more similar to the original parameter trajectory. Opposing to HOAG and greedy algorithms, DrMad optimization has prerequisites for not only hyperparameters optimization but also the metaparameters, i.e., the optimization procedure parameters. The opportunity of such optimization using reversed differentiation was shown in [10].

The other interesting aspect of our experiments is the relation between the model error (RMSE or Accuracy) and the value of validation loss  $Q$ . The models obtained by the evidence lower bound showed higher error rate than the models

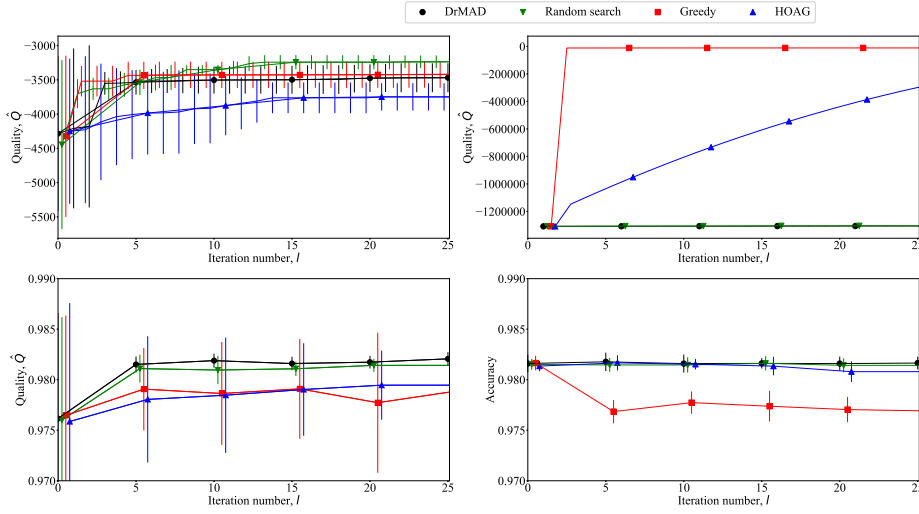


Fig. 4: MNIST, best validation value  $\hat{Q}$  and Accuracy for cross-validation (left) and evidence lower bound (right)

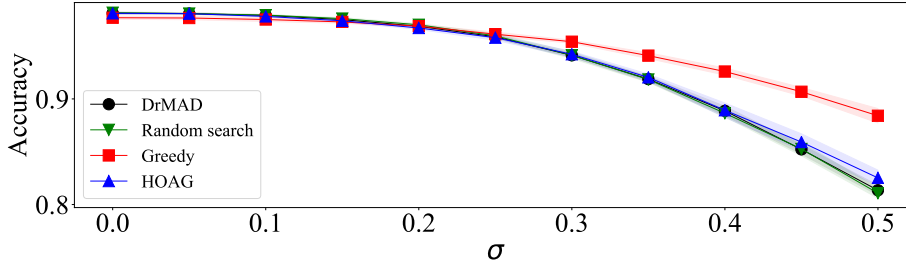


Fig. 5: MNIST, Model accuracy with noise in the Test dataset. The hyperparameters were optimized with evidence lower bound criterion.

obtained using cross-validation on the MNIST and WISDM datasets. These models also show greater stability when noise was added to the Test datasets. The evidence lower bound showed significantly better results on the synthetic dataset, when the amount of objects in the Train dataset is small. Therefore we conclude that the evidence lower bound usage is preferable when the model tends to overfit or when the cross-validation usage is too computationally expensive. In [11] it was noted that the evidence lower bound optimization required more iterations for the convergence. In our experiments we used the same number of iterations both for the cross-validation and evidence lower bound. The more accurate iteration number calibration can improve the final quality of these models.

## 6 Conclusion and future work

The paper analyzes the gradient-based hyperparameter optimization algorithms. We adapt the analyzed algorithms for general validation functions and evaluated their performance on the MNIST and WISDM datasets. Two model selection criteria are compared: the cross-validation and evidence lower bound.

The experiments show that the gradient-based algorithms are effective when the number of hyperparameters is large. The results show that models obtained using evidence lower bound have higher error rates than models obtained using cross-validation, but they are also more stable when the test dataset contains a lot of noise.

The authors implement these algorithms as a toolbox available at [2]. The toolbox is developed in Python using Theano [5] and Numpy [25] libraries.

In our future work we are planning to develop the analyzed algorithm and to extend gradient-based algorithms to optimize not only hyperparameters, but also the model optimization parameters. The other object of our future research is the difference between the cross-validation and evidence lower bound and the theoretical aspects of their properties for the models with large amount of parameters.

## References

1. Arlot, S., Celisse, A.: A survey of cross-validation procedures for model selection. *Statistics Surveys* **4**, 40–79 (2010). DOI 10.1214/09-SS054. URL <http://dx.doi.org/10.1214/09-SS054>
2. Bakhteev, O.: pyfos: a library for hyperparameter optimization. URL <https://git.io/fjZks>
3. Bengio, Y., Grandvalet, Y.: No unbiased estimator of the variance of k-fold cross-validation. *Journal of machine learning research* **5**, 1089–1105 (2004)
4. Bergstra, J., Bengio, Y.: Random search for hyper-parameter optimization. *Journal of Machine Learning Research* **13**(Feb), 281–305 (2012)
5. Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., Turian, J., Warde-Farley, D., Bengio, Y.: Theano: a cpu and gpu math expression compiler. In: *Proceedings of the Python for scientific computing conference (SciPy)*, vol. 4 (2010)
6. Bergstra, J.S., Bardenet, R., Bengio, Y., Kégl, B.: Algorithms for hyper-parameter optimization. In: *Advances in Neural Information Processing Systems*, pp. 2546–2554 (2011)
7. Bishop, C.M.: *Pattern Recognition and Machine Learning* (Information Science and Statistics). Springer-Verlag New York, Inc., Secaucus, NJ, USA (2006)
8. Domke, J.: Generic methods for optimization-based modeling. In: *AISTATS*, vol. 22, pp. 318–326 (2012)
9. Farcomeni, A.: Bayesian constrained variable selection. *Statistica Sinica* pp. 1043–1062 (2010)
10. Fu, J., Luo, H., Feng, J., Low, K.H., Chua, T.S.: Drmad: distilling reverse-mode automatic differentiation for optimizing hyperparameters of deep neural networks. In: *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, pp. 1469–1475 (2016)
11. Graves, A.: Practical variational inference for neural networks. In: *Advances in Neural Information Processing Systems* 24, pp. 2348–2356 (2011)
12. Grünwald, P.: A tutorial introduction to the minimum description length principle. In: *Advances in Minimum Description Length: Theory and Applications*. MIT Press (2005)
13. Hoffman, M.D., Blei, D.M., Wang, C., Paisley, J.: Stochastic variational inference. *Journal of machine learning research* **14**(1), 1303–1347 (2013)
14. Hutter, F., Hoos, H.H., Leyton-Brown, K.: Sequential model-based optimization for general algorithm configuration. In: *International Conference on Learning and Intelligent Optimization*, pp. 507–523 (2011)
15. Hwang, S., Jeong, M.K.: Robust relevance vector machine for classification with variational inference. *Annals of Operations Research* **263**(1-2), 21–43 (2018)

16. Karaletsos, T., Rätsch, G.: Automatic relevance determination for deep generative models. *stat* **1050**, 26 (2015)
17. Krstajic, D., Buturovic, L.J., Leahy, D.E., Thomas, S.: Cross-validation pitfalls when selecting and assessing regression and classification models. *Journal of Cheminformatics* **6**(1), 1–15 (2014)
18. Kuznetsov, M., Tokmakova, A., Strijov, V.: Analytic and stochastic methods of structure parameter estimation. *Informatica* **27**(3), 607–624 (2016). DOI <http://www.mii.lt/informatica/pdf/INFO1109.pdf>
19. Kwapisz, J.R., Weiss, G.M., Moore, S.A.: Activity recognition using cell phone accelerometers. *ACM SigKDD Explorations Newsletter* **12**(2), 74–82 (2011)
20. LeCun, Y.: The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/> (1998)
21. Luketina, J., Berglund, M., Greff, K., Raiko, T.: Scalable gradient-based tuning of continuous regularization hyperparameters. In: *International Conference on Machine Learning*, pp. 2952–2960 (2016)
22. MacKay, D.J.C.: *Information Theory, Inference & Learning Algorithms*. Cambridge University Press, New York, NY, USA (2002)
23. Maclaurin, D., Duvenaud, D., Adams, R.: Gradient-based hyperparameter optimization through reversible learning. In: *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pp. 2113–2122 (2015)
24. Nystrup, P., Boyd, S., Lindström, E., Madsen, H.: Multi-period portfolio selection with drawdown control. *Annals of Operations Research* pp. 1–27 (2017)
25. Oliphant, T.E.: *A guide to NumPy*, vol. 1. Trelgol Publishing USA (2006)
26. Pedregosa, F.: Hyperparameter optimization with approximate gradient. In: *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, vol. 48, pp. 737–746 (2016)
27. Salakhutdinov, R., Hinton, G.: Learning a nonlinear embedding by preserving class neighbourhood structure. In: *Journal of Machine Learning Research - Proceedings Track*, vol. 2, pp. 412–419 (2007)
28. Salimans, T., Kingma, D.P., Welling, M.: Markov chain monte carlo and variational inference: Bridging the gap. In: *ICML*, vol. 37, pp. 1218–1226 (2015)
29. Snoek, J., Rippel, O., Swersky, K., Kiros, R., Satish, N., Sundaram, N., Patwary, M.M.A., Prabhat, M., Adams, R.P.: Scalable bayesian optimization using deep neural networks. In: *ICML*, pp. 2171–2180 (2015)
30. Strijov, V.V.: *Model genetation and selection for regression and classification problems (dsc thesis)* (2014)
31. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8–13 2014, Montreal, Quebec, Canada*, pp. 3104–3112 (2014)
32. Vishnu, A., Narasimhan, J., Holder, L., Kerbyson, D.J., Hoisie, A.: Fast and accurate support vector machines on large scale systems. In: *2015 IEEE International Conference on Cluster Computing, CLUSTER 2015, Chicago, IL, USA, September 8–11, 2015*, pp. 110–119 (2015)