

Automatic machine learning model selection

David Laredo¹

¹School of Mechanical Engineering, University of California, Merced

Abstract

Neural networks and deep learning are changing the way that artificial intelligence is being done. Efficiently choosing a suitable model (including hyperparameters) for a specific problem is a time-consuming task. Choosing among the many different combinations of neural networks available gives rise to a staggering number of possible alternatives overall. Here we address this problem by proposing a fully automated framework for efficiently selecting a neural network model given a specific problem (whether it is classification or regression). Our proposal focuses on a distributed decision-making algorithm for keeping the most promising models among a pool of possible models. We hope that this approach will help non-expert users to more effectively identify neural network based models and hyperparameter settings appropriate to their applications, and hence to achieve improved performance.

Index terms— artificial neural networks, model selection, hyperparameter tuning, distributed computing, evolutionary algorithms

1. INTRODUCTION

Machine learning studies automatic algorithms that improve themselves through experience. Given the large amounts of data currently available in many fields such as engineering, biomedical, finance, etc, and the increasingly computing power available machine learning is now practiced by people with very diverse backgrounds. Increasingly, users of machine learning tools are non-experts who require off-the-shelf solutions. The machine learning community has aided these users by making available a variety of easy to use learning algorithms and feature selection methods as WEKA [1] and PyBrain [2]. Nevertheless, the user still needs to make some choices which not may be obvious or intuitive (selecting a learning algorithm, hyperparameters, features, etc).

Recently, neural networks have gained a lot of attention due to the newer models (CNN, RNN, Deep Learning, etc.) and their flexibility and generality for solving a large number of problems: regression, classification, natural language processing, recommendation systems, just to mention a few. Furthermore, there are a lot software libraries which makes their implementations easy to use (tensorflow, keras, kaffe, etc.). Nevertheless, the task of picking the right neural network model (hyperparameters included) can be even more complicated than that of other algorithms.

Given the popularity of neural networks, specially among non computer scientist we will focus our efforts in this study to them and leave other algorithms for future work.

Usually, the process of selecting a suitable machine learning model for a particular problem is done in an iterative manner. First, an input dataset must be transformed from a domain specific format to features which are predictive of the field of interest. Once features have engineered users must pick a learning setting appropriate to their problem, e.g. regression, classification or recommendation. Next users must pick an appropriate model, such as support vector machines (SVM), logistic regression, any flavor of neural networks (NN). Each model family has a number of hyperparameters, such as regularization degree, learning rate, number of neurons, and each pf these must be tuned to achieve optimal results. Finally, users must pick a software package that can train their model, configure one or more machines to execute the training and evaluate the model's quality. It can be challenging to make the right choice when faced with so many degrees of freedom, leaving many users to select a model based on intuition or randomness and/or leave hyperparameters set to default. Certainly this approach will usually yield suboptimal results.

This suggests a natural challenge for machine learning: given a dataset, to automatically and simultaneously chose a learning algorithm and set its hyperparameters to optimize performance. As mentioned in [1] the combined space of learning algorithm and hyperparemeters 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. hyperparameters of the algorithm are only meaningful if that algorithm is chosen). Thus, identifying a high quality model is typically costly (in the sense that entails a lot of computational effort) and time consuming.

Distributed and cloud computing provide a compelling way to accelerate this process, but also present additional challenges. Though parallel storage and processing techniques enable users to train models on massive datasets and accelerate the search process by training multiple models at once, the distributed setting forces several more decisions upon users: what parallel execution strategy to use, how big a cluster to provision, how to efficiently distribute computation across it, and what machine learning framework to use. These decisions are onerous, particularly for users who are experts in their own field but inexperienced in machine learning and distributed systems.

To address this challenges we propose NeuroTuner a flexible and scalable system to automate the process of selecting artificial neural network models.

2. MODEL SELECTION

In this section we introduce and formally describe the model selection problem, for this section we borrow the definitions given in [3]. This work focuses on supervised learning: 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 $d_i = (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{X} \times \mathcal{Y}$ to such a function. Most learning algorithms A further expose hyperparameters $\lambda \in \Lambda$, which change the way the learning algorithm A_λ works. One example of hyperparameters is the number of neurons in a hidden layer of an ANN, another common example is the learning rate α of a neural network. 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 $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_n, \mathbf{y}_n)\}$, 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}_{train}^i and $\mathcal{D}_{validation}^i$, learning functions f_i by applying A^* to \mathcal{D}_{train}^i , and evaluating the predictive performance of these functions $\mathcal{D}_{validation}^i$. This allows for the model selection to be written as:

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

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

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: $\Lambda \subset \Lambda_1 \cdots \Lambda_n$. 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. More formally, following [4], 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 [5] or, in some cases, a directed acyclic graph (DAG) [4]. Given such a structured space Λ , the (hierarchical) hyperparameter optimization problem can be written as:

$$\lambda^* \in \operatorname{argmin}_{\lambda \in \Lambda} \frac{1}{k} \sum_{i=1}^k \mathcal{L} \left(A_\lambda, \mathcal{D}_{train}^i, \mathcal{D}_{validation}^i \right), \quad (2)$$

In this study we consider the more general combined algorithm selection and hyperparameter optimization (CASH). That is we intend to optimize both problems at the same time (maybe through a scalarization of both objectives, maybe through multi-objective optimization).

3. LITERATURE REVIEW

Automatic model selection has been of research interest since the uprising of deep learning. This is no surprise since selecting an effective combination of algorithm and hyperparameter values is

currently a challenging task requiring both deep machine learning knowledge and repeated trials. This is not only beyond the capability of layman users with limited computing expertise, but also often a non-trivial task even for machine learning experts [6].

To make machine learning accessible to non-expert users, researchers have proposed various automatic selection methods for machine learning algorithms and/or hyperparameter values for a given supervised machine learning problem. These methods' goal is to find, within a pre-specified resource limit (usually specified in terms of time, number of algorithms and/or combinations of hyperparameter values), an effective algorithm and/or combination of hyperparameter values that maximize the accuracy measure on the given machine learning problem and data set. Using an automatic selection method, the machine learning practitioner can skip the manual and iterative process of selecting an efficient combination of hyperparameter values and neural network model, which is high labor intensive and requires a high skill set in machine learning.

In the recent years a number of tools have been made available for users to automate the model selection and/or hyperparameter tuning, in the following we present a brief survey of the most popular methods.

3.1. AutoWEKA

Auto-WEKA [7] is a system designed to help machine learning users by automatically searching through the joint space of WEKA's learning algorithms and their respective hyperparameter settings to maximize performance using a state-of-the-art Bayesian optimization method. AutoWEKA addresses the CASH problem by treating all of WEKA as a single, highly parametric machine learning framework, and using Bayesian optimization to find a strong instantiation for a given dataset. AutoWEKA also natively supports parallel runs (on a single machine) to find good configurations faster and save the N best configurations of each run instead of just the single best. AutoWEKA is tightly integrated with WEKA and does provide support for Multilayer Perceptrons (MLP).

3.2. Auto-sklearn

Auto-sklearn [8] is Auto-WEKA's sister package, it uses the same Bayesian optimizer but comprises a smaller space of models and hyperparameters, however it includes additional meta-learning techniques.

3.3. TuPAQ

TuPAQ [6] is a system designed to efficiently and scalably automate the process of training predictive models. One of its main features is a planning algorithm which decides on an efficient parallel execution strategy during model training while identifying new hyperparameter configurations and proactively eliminating models which are unlikely to provide good results. TuPAQ is aimed at large scale machine learning, it builds on top of the well known Apache Spark. TuPAQ only focuses on classification problems and considers only three model families (Support Vector Machines, Logistic Regression and nonlinear SVMs), each with several hyperparameters. TuPAQ

performs batching to train multiple models simultaneously and deploys bandit resource allocation to allocate more resources to the most promising models. TuPAQ does not provide support for neural networks.

4. OUR PROPOSAL

While there is a number of methods for automatic model selection and hyperparameter tuning, the most popular ones still have room for improvement. In the case of AutoWEKA and Auto-sklearn they do not provide good support for large machine learning problems, nor provide support for distributed computing. TuPAQ on the other hand, provides wide support for distributed computing, maximizing the use of computational resources through the use of sophisticated optimizations, nevertheless its restricted to only classification problems and does not provide support for neural networks.

We propose to implement a system for automatically selecting the most fitting neural network architecture (only fully connected networks in the first stage) for a given problem, whether it is classification or regression. Furthermore, we plan that the system should be scalable and should be able to be used in distributed computing environments, allowing it to be usable for large datasets and complex models. To achieve the latter we propose to build our system using Ray [9] which is a distributed system designed with large scale distributed machine learning in mind.

Our proposal includes the comparisson of an evolutionary algorithm against sequential model based optimization for the selection of an optimal model, along with sampling and pruning techniques to only keep the most promising models and discard the rest of them. Furthermore, we propose to train multiple models simultaneously using Ray and perform bandit allocation resource to perform efficient resource allocation towards the most promising models. Our results will be first tested using the CMAPSS dataset [10] for regression and the MNIST dataset [11] for classification. A comparisson against AutoWEKA and Auto-sklearn (on the neural networks) will also be provided.

Specific milestones are presented next:

1. (1 week) Install and run Ray, do some basic computations using the framework and the available resources (CPU, GPU, multiple machines, etc.)
2. (1 week) Define suitable metrics for each of the problems at hand e.g., a composed metric (precision and recall) for classification.
3. (1 week) Select an appropriate evolutionary algorithm (genetic algorithm, differential evolution) as a solver for the optimization problem of choosing the most fit neural network.
4. (1 week) Define an appropriate encoding for the neural network model to be used by the evolutionary algorithm (array of bits, array of integers, array of real and integer values, etc.)
5. (3 weeks) Implement the proposed framework in python and test it in a single computer
6. (2 weeks) Implement the distributed version of the framework (using ray) and test it.
7. (2 week) Evaluate the performance of the framework, synthesize the results and compare against existing frameworks.

REFERENCES

- [1] 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.
- [2] T. Schaul, J. Bayer, D. Wierstra, Y. Sun, M. Felder, F. Sehnke, T. Rucksties, and J. Schmidhuber. Pybrain. *JMLR*, 11:743–746, 2010.
- [3] Thornton C., Hutter F., Hoos H., and Leyton-Brown K. Auto-weka: Combined selection and hyperparameter optimization of classification algorithms. In *KDD*, 2013.
- [4] Hutter F., Hoos H., Leyton-Brown K, and Stutzle T. Paramils: and automatic algorithm configuration framework. *JAIR*, 36(1):267–306, 2009.
- [5] Bergstra J., Bardenet R., Bengio Y., and Kegl B. Algorithms for hyper-parameter optimization. In *NIPS*, 2011.
- [6] Sparks ER., Talwalkar A., Smith V., Kottalam J., Pan X, and Gonzales JE. Automated model search for large scale machine learning. In *SoCC*, pages 368–380, 2015.
- [7] Thornton C., Hutter F., Hoos H., Leyton-Brown K, and Kotthoff L. Auto-weka 2.0: Automatic model selection and hyperparameter optimization in weka. *JMLR*, 2016.
- [8] Feurer M., Klein A., Eggenberger K., Springenberg J., Blum M., and Hutter F. Efficient and robust automated machine learning. In *NIPS*, volume 17, pages 1–5, 2015.
- [9] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, William Paul, Michael I. Jordan, and Ion Stoica. Ray: A distributed framework for emerging AI applications. *CoRR*, abs/1712.05889, 2017.
- [10] A. Saxena and K. Goebel. Phm08 challenge data set. [Online] Available at: <https://ti.arc.nasa.gov/tech/dash/groups/pcoe/prognostic-data-repository/>.
- [11] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. -, 2010.