# Hyperparameter optimization in deep multi-target prediction

Dimitrios Iliadis[1][0000−0002−3676−5940], Marcel Wever[2][0000−0001−9782−6818], Bernard De Baets[1][0000−0002−3876−620X], and Willem Waegeman[1][0000−−0002−5950−3003]

[1] KERMIT, Department of Data Analysis and Mathematical Modelling, Ghent University, Coupure links 653, B-9000 Ghent, Belgium
{dimitrios.iliadis,bernard.debaets,willem.waegeman}@ugent.be
[2] Department of Computer Science, Ludwig-Maximilians-University Munich, Akademiestr. 7, 80799 Munich, Germany
marcel.wever@ifi.lmu.de

**Abstract.** As a result of the ever increasing complexity of configuring and fine-tuning machine learning models, the field of automated machine learning (AutoML) has emerged over the past decade. However, software implementations like Auto-WEKA and Auto-sklearn typically focus on classical machine learning (ML) tasks such as classification and regression. Our work can be seen as the first attempt at offering a single AutoML framework for most problem settings that fall under the umbrella of multi-target prediction, which includes popular ML settings such as multi-label classification, multivariate regression, multi-task learning, dyadic prediction, matrix completion, and zero-shot learning. Automated problem selection and model configuration are achieved by extending DeepMTP, a general deep learning framework for MTP problem settings, with popular hyperparameter optimization (HPO) methods. Our extensive benchmarking across different datasets and MTP problem settings identifies cases where specific HPO methods outperform others.

**Keywords:** Multi-target prediction · automated machine learning · hyperparameter optimization · multi-label classification · multivariate regression · matrix completion · multi-task learning · dyadic prediction

## 1 Introduction

The past decade of AI research has been dominated by significant advances in deep learning. From convolutional neural networks (CNNs) [20,32] to generative adversarial networks (GANs) [10] and transformers [36], deep learning architectures have enabled major breakthroughs in several application areas, such as computer vision, speech recognition, and protein folding [17]. However, the configuration of these increasingly complex architectures requires multiple design decisions that are not standardized. Selecting the appropriate architecture and hyperparameters for the optimal neural network is usually reserved for highly

arXiv:2211.04362v1 [cs.LG] 8 Nov 2022

experienced users who navigate the configuration space through trial and error. This process becomes dull in a per-dataset case and essentially infeasible when one wants to offer a model in a software tool that thousands of inexperienced users will potentially use.

The increased demand for machine learning applications and the limited availability of expertise has led to the emergence of the field of automated machine learning (AutoML) [15], which is concerned with automating the process of engineering machine learning applications. In particular, this field aims to develop methods that help to move away from the tedious task of manually configuring machine learning algorithms to a data-driven approach that is able to efficiently navigate through the space of potential solution candidates while maintaining near-optimal performance. An essential sub-task that needs to be tackled to achieve this performance deals with the optimization of hyperparameters. Since we focus on hyperparameter optimization (HPO) in this work, we refer to several AutoML-related surveys [3, 6, 13] for a more thorough introduction and overview. Especially in the case of neural networks, hyperparameter optimization plays a significant role as the hyperparameters greatly affect the computational complexity and the generalization performance.

A subarea of AutoML research, also known as neural architecture search (NAS) [6], is solely concerned with hyperparameter optimization of a particular class of models, namely neural networks. While NAS has demonstrated promising performance [41] and efficiency [28] improvements, most existing work has focused on the challenging, yet narrow task of image classification, leaving other types of equally interesting learning tasks largely unexplored. A similar trend can be seen at the software level, as most published tools are designed for single-target classification and regression, and only a limited number of them focus on other types of learning tasks. One of those less explored areas involves the simultaneous prediction of multiple targets. Despite the broad applicability potential of the area of multi-target prediction (MTP), only a few tools have been proposed for specific subareas of multi-label classification. A more detailed review of such tools will be given in Section 4.

The possibility of utilizing an automated tool for the majority of sub-areas that fall under the umbrella of MTP is an exciting idea. Typical examples of MTP settings are multi-label classification, multivariate regression, multi-task learning, dyadic prediction, zero-shot learning, network inference, and matrix completion. A first attempt in this direction was introduced by the DeepMTP framework [16], which will be reviewed in Section 2. This framework makes it possible for non-expert users to automatically select the most appropriate MTP problem setting by answering a handful of questions. After selecting the most appropriate MTP problem setting, the DeepMTP framework utilizes a flexible two-branch architecture that can be adjusted for specific MTP settings. The experiments presented in [16] showcased DeepMTP as a competitive approach, compared to other baseline methods across multiple MTP problem settings and datasets. In terms of HPO, a standard grid search was used for all the comparisons. Even though this is acceptable in a research environment, it is certainly not

practical for a user-centered software package. This is the primary motivation behind the work presented in this paper. We intend to increase the practical usability of the DeepMTP framework with an extension that utilizes efficient HPO methods. These HPO methods will be described in Section 3, and the benchmarking results will be presented in Section 5.

## 2   A short review of DeepMTP

### 2.1   Automated selection of the most suitable MTP problem setting

As mentioned in the introduction, multi-target prediction comprises various sub-areas of machine learning, such as multi-label classification, multivariate regression, multi-task learning, dyadic prediction, zero-shot learning, network inference, and matrix completion. All these problem settings display specific characteristics that have resulted in the development of specific machine learning methods. At the same time, they also share a significant commonality, i.e., the prediction of multiple target variables. An example of a multi-label classification problem, the detection of dog breeds from images of mixed-breed dogs, can be seen in Fig. 1. The main difference in this task is that every dog can be associated with multiple breeds simultaneously (multi-label) instead of just one of them (single-label). Since a detailed formal definition of every MTP problem setting is out of scope for this work, we only present the general definition of MTP. For more details about the other MTP problem settings, we refer the reader to the survey by Waegeman et al. [38]. A software package for DeepMTP is already available online[3] and a manuscript with a detailed explanation of the capabilities and general functionality is under review.

**Definition 1.** *A **multi-target prediction problem** is characterized by an instance set $\mathcal{X}$, a target set $\mathcal{T}$ and a score set $\mathcal{Y}$ with the following properties:*

*(P1) A training dataset $\mathcal{D}$ contains triplets $(\mathbf{x}_i, \mathbf{t}_j, y_{ij})$, where $\mathbf{x}_i \in \mathcal{X}$ represents an instance, $\mathbf{t}_j \in \mathcal{T}$ represents a target, and $y_{ij} \in \mathcal{Y}$ is the score that quantifies the relationship between an instance and a target, with $i \in \{1, \ldots, n\}$ and $j \in \{1, \ldots, m\}$. The scores can be arranged in an $n \times m$ matrix $\mathbf{Y}$ that is usually incomplete.*
*(P2) The score set $\mathcal{Y}$ consists of nominal, ordinal or real values.*
*(P3) During testing, the objective is to predict the score for any unobserved instance-target couple $(\mathbf{x}, \mathbf{t}) \in \mathcal{X} \times \mathcal{T}$.*

This definition is very general, as it intends to cover all the MTP problem settings that were considered by Waegeman et al. [38]. It describes three basic properties, so every MTP problem setting can be defined by adding more custom properties to the primary three. This is clear in the formal definition of multi-label classification, which requires four additional properties. Property **P4** defines the generalization objective since one expects to make only predictions

---

[3] `https://github.com/diliadis/DeepMTP`

| | beagle | cocker | fox terrier | Jack russell | poodle | shepherd |
|---|---|---|---|---|---|---|
|  | 0 | 1 | 0 | 0 | 0 | 0 |
|  | 1 | 0 | 0 | 0 | 0 | 1 |
|  | 0 | 1 | 0 | 0 | 1 | 0 |
|  | 1 | 0 | 1 | 1 | 0 | 0 |

**Fig. 1.** Example of a multi-label classification task in which the goal is to identify the breeds of mixed-dogs. The instances correspond to images of dogs and the labels/targets to the breed names.

for new targets during testing. Property **P5** specifies the absence of side information (commonly known as input features) for the targets. Finally, properties **P6** and **P7** inform about the state of the score matrix, mainly that the interaction values for all possible (instance, target) pairs $(\mathbf{x}_i, \mathbf{t}_j)$ in our training set are known and of binary type. These properties form the basis for the questionnaire used by the DeepMTP framework.

**Q1**: Is it expected to encounter novel instances during testing? **(yes/no)**
**Q2**: Is it expected to encounter novel targets during testing? **(yes/no)**
**Q3**: Is there side information available for the instances? **(yes/no)**
**Q4**: Is there side information available for the targets? **(yes/no)**
**Q5**: Is the score matrix fully observed? **(yes/no)**
**Q6**: What is the type of the target variable? **(binary/nominal/ordinal/real-valued)**

Question **Q2** can be mapped to property **P4** as it relates to generalizing to new targets. Question **Q4** and property **P5** refer to the existence of side information for the targets. Questions **Q1** and **Q2** were generated from similar properties **P4** and **P5** that refer to the instances. Finally, questions **Q5** and **Q6** are derived from properties **P6** and **P7**, respectively (state of score matrix and target variable type).

**Table 1.** A snapshot of specific answers to the DeepMTP questionnaire and the corresponding MTP problem setting. "-" denotes a wildcard.

| Q1 | Q2 | Q3 | Q4 | Q5 | Q6 | MTP method |
|----|----|----|----|----|----|-----------|
| yes | no | yes | no | yes | binary | Multi-label classification |
| yes | no | yes | no | yes | real-valued | Multivariate regression |
| yes | no | yes | no | no | - | Multi-task learning |
| yes | no | yes | yes (hierarchy) | yes | binary | Hierarchical multi-label classification |
| yes | no | yes | yes | no | - | Dyadic prediction |
| yes | yes | yes | yes | no | - | Zero-shot learning |
| no | no | no | no | no | - | Matrix completion |
| no | no | yes | yes | no | - | Hybrid matrix completion |
| yes | yes | yes | yes | no | - | Cold-start collaborative filtering |
| yes | no | yes | no | yes | nominal/categorical | Multi-dimensional classification |

Table 1 shows how specific combinations of answers lead to individual MTP problem settings. The mapping is not based on a data-driven method, but the in-house expertise from the research team that developed DeepMTP [16, 38]. Furthermore, the selection of a specific MTP problem setting does not affect the configuration of the neural network architecture used by the DeepMTP framework, but intends to guide the user to the most appropriate literature. Configuration-related decisions are made based on the respective questions that comprise the questionnaire. The practical aspect of the framework moves away from the individual MTP settings to a more general view that is based on three principles. These include the generalization objectives for instances and targets, the existence of side information for instances and targets, and finally, the target variable type.

The questionnaire can be answered manually by (inexperienced) users or even automatically if the dataset is provided. There are at most three possible datasets that can be required by any of the MTP problem settings. Two of them contain the side information for the instances and targets, and the third one contains the score matrix, a prerequisite for every MTP problem setting. If these datasets are supplied to the framework, the task of answering the questionnaire becomes trivial with a simple computer program. If a user uploads the side information files for the instances and targets, questions **Q3** and **Q4** are answered trivially. Question **Q5** can be answered by detecting missing values in the supplied dataset. To automatically determine the answers for questions **Q1** and **Q2**, one can compare the relations between instances (targets) in the training and test files, thus arriving at one of four generalization (validation) settings. Setting A involves the prediction of missing values inside the interaction matrix. In Setting B, the goal is to make predictions for new instances, while Setting C involves the prediction for new targets. Setting D requires the prediction for novel pairs of instances and targets. In conclusion, we show that in the first stage of the DeepMTP framework, the selection of the most appropriate MTP problem setting can be automated by combining the original questionnaire with basic programming.
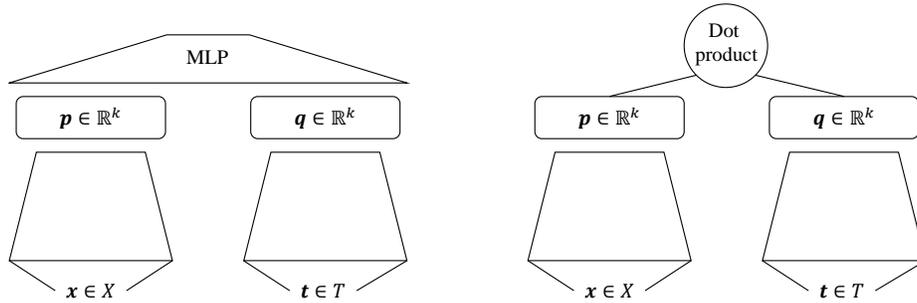
**Fig. 2.** Simplified view of the two two-branch neural network architectures.

### 2.2    The neural network architecture behind DeepMTP

The DeepMTP framework utilizes a two-branch architecture that has gained popularity in the field of collaborative filtering [12]. However, the same architecture can be easily modified to achieve competitive performance across multiple MTP problem settings. The architecture (see Fig. 2 left) features two branches that take as input any available side information for the instances and targets and then output two embedding vectors $\mathbf{p_x}$ and $\mathbf{q_t}$, respectively. The embeddings are then concatenated and the resulting vector is used as input to a series of fully connected layers that terminate at a single output node, as follows:

$$
\begin{aligned}
\mathbf{z}_1 = \phi_1(\mathbf{p_x}, \mathbf{q_t}) &= \begin{bmatrix} \mathbf{p_x} \\ \mathbf{q_t} \end{bmatrix}, \\
\phi_2(\mathbf{z}_1) &= \alpha_2(\mathbf{W}_2^T \mathbf{z}_1 + \mathbf{b}_2), \\
&\vdots \\
\phi_L(\mathbf{z}_{L-1}) &= \alpha_L(\mathbf{W}_L^T \mathbf{z}_{L-1} + \mathbf{b}_L), \\
\hat{y}_{\mathbf{xt}} &= \sigma(\mathbf{h}^T \phi_L(\mathbf{z}_{L-1})),
\end{aligned}
\tag{1}
$$

where $\mathbf{W}$, $\mathbf{b}$ and $\alpha$ represent the weight matrix, bias vector, and activation function of the final multi-layer perceptron (MLP) layer, respectively. Alternatively, a seemingly more straightforward yet also less expressive approach skips the final series of fully-connected layers and instead computes the dot product (see Fig. 2 right) of the two embedding vectors in the following way:

$$
\hat{y}_{\mathbf{xt}} = \sigma(\mathbf{p_x} \cdot \mathbf{q_t}). \tag{2}
$$

Even though the architecture that uses the MLP was initially proposed as a more powerful approximator [12], subsequent work [4,29] has argued that the reality of training the more complex model results in practical disadvantages. Our preliminary experiments seem to agree with that observation, so all the experiments presented below will use the dot product version. Any modifications to this architecture are guided by the answers to the aforementioned questionnaire in the following way:

- The combination of answers for **Q1** and **Q2** determines the validation setting the user expects. This can be explicitly requested by the user, especially if more than one validation setting is available given the datasets, or inferred by the relation of the instance and target IDs in the train and test datasets.
- The answers to questions **Q3** and **Q4** play multiple roles. A negative answer for any of the two questions will lead to the use of one-hot encoded vectors as the input to the corresponding branch. Furthermore, these answers determine the feasibility of the generalization of the user requests (through **Q1** and **Q2**).
- Question **Q5** is used to distinguish between MTP problem settings (e.g., multi-label classification and multi-task learning).
- The answer to question **Q6** dictates the type of loss function that is used during training (binary cross-entropy loss for classification tasks and squared error loss for regression tasks).

To conclude, the neural network used in the DeepMTP framework is capable of adapting to the various characteristics that MTP problem settings exhibit, from the existence of input features to the type of task that they represent. This functionality, combined with the purpose-made questionnaire and a basic user interface, can make multi-target prediction a more accessible area of research.

## 3 Hyperparameter optimization methods

In this section, we give an overview of the hyperparameter optimization (HPO) methods that we consider for the benchmarking experiments in Section 5. These techniques are usually grouped into two main categories: black box and multi-fidelity techniques.

**Grid Search & Random Search:** Conceptually, grid search is generally considered the simplest HPO method, as the optimal configuration is identified by brute-forcing the evaluation of all possible configurations of a user-defined hyperparameter space [25]. The main weakness of this approach is its high computational cost, as the curse of dimensionality means that the number of configuration evaluations grows exponentially with the size of the configuration space. Furthermore, because continuous hyperparameters need to be discretized, an increase in their resolution can quickly lead to an explosion in the resulting number of configuration evaluations.

Instead of exhaustively searching over the hyperparameter space, a random search samples configurations at random until a user-defined budget is exhausted [2]. This approach is usually able to outperform grid search in cases where some hyperparameters are more important than others, as it may find configurations that are left out when discretizing. In the literature, random search is a standard baseline when comparing HPO methods, which is also why we choose to include a random search in the experimental section of this paper.

**Sequential Model-based Algorithm Configuration:** The sequential model-based algorithm configuration (SMAC) approach [14] is a Bayesian optimization method that uses random forests as a surrogate model. The main advantage over other options such as Gaussian processes is that random forests can naturally support categorical hyperparameters, handle larger search spaces, and scale better as the number of training samples increases.

After an initialization phase of randomly sampled observations, SMAC fits a random forest to the collected observations, serving as the surrogate model in the framework of Bayesian optimization. Subsequently, SMAC alternates between determining the next hyperparameter configuration to evaluate and updating the surrogate model with the newly evaluated observation. In the former step, the surrogate model is employed to estimate the usefulness of evaluating a hyperparameter configuration $\theta$. Assessing the usefulness of an evaluation requires tackling the so-called exploration-exploitation dilemma, which is traditionally implemented via a so-called acquisition function, e.g., the expected improvement $\mathbb{EI}(\theta)$, comparing the potential improvement of some configuration over the best configuration observed so far:

$$\mathbb{EI}(\theta) = \sigma_\theta[u \cdot \Phi(u) + \phi(u)], \quad u = \frac{o_{min} - \mu_\theta}{\sigma_\theta} \tag{3}$$

where $\Phi$ denotes the cumulative distribution function of the standard normal distribution, $\phi$ is the corresponding probability density function, and $o_{min}$ is the loss of the best performing configuration. This approach has proven to be quite successful and is at the heart of several AutoML software packages (e.g., Auto-Weka [34], auto-sklearn [8]). For our experiments, we used the SMAC3 implementation [22] that is available online[4].

**Multi-fidelity optimization - Hyperband:**  One of the core steps in any standard HPO method is the performance evaluation of a given configuration. This can be manageable for simple models that are relatively cheap to train and test, but can become a significant bottleneck for more complex models that need hours or even days to train. This is particularly evident in deep learning, as big neural networks with millions of parameters trained on increasingly larger datasets can deem traditional black-box HPO methods impractical.

Addressing this issue, multi-fidelity HPO methods have been devised to discard unpromising hyperparameter configurations already at an early stage. To this end, the evaluation procedure is adapted to support cheaper evaluations of hyperparameter configurations, such as evaluating on sub-samples (feature-wise or instance-wise) of the provided data set or executing the training procedure only for a certain number of epochs in the case of iterative learners. The more promising candidates are subsequently evaluated on increasing budgets until a maximum assignable budget is reached.

A popular representative of such methods is Hyperband [21]. Hyperband builds upon Successive Halving (SH) [18], where a set of $n$ candidates is first

---

[4] `https://github.com/automl/SMAC3`

evaluated on a small budget. Based on these *low-fidelity* performance estimates, the $\frac{n}{\eta}$ ($\eta \geq 2$) best candidates are preserved, while the remaining configurations are already discarded. Iteratively increasing the evaluation budget and re-evaluating the remaining candidates with the increased budget while discarding the inferior candidates results in fewer resources wasted on inferior candidates. In return, one focuses more on the promising candidates.

Despite the efficiency of the successive halving strategy, it is well known that it suffers from the exploration-exploitation trade-off. In simple terms, a static budget $\mathcal{B}$ means that the user has to manually decide whether to explore a number of configurations $n$ or give each configuration a sufficient budget to develop. An incorrect decision can lead to an inadequate exploration of the search space (small $n$) or the early rejection of promising configurations (large $n$). Hyperband overcomes the exploration-exploitation trade-off by repeating the successive halving strategy with different initializations of SH, varying the budget and the number of initial candidate configurations.

**Bayesian Optimization with Hyperband:** Despite the advantages that Hyperband displays compared to baseline methods like grid search and random search, it is still restricted by the random sampling of configurations at the beginning of each iteration of the successive halving routine. Learning from past sampled configurations has the potential to provide improvements in the final performance compared to standard, model-free Hyperband. The Bayesian optimization with Hyperband method [7] tries to improve over this by replacing random sampling with a model-based approach that utilizes Bayesian optimization.

More specifically, at the beginning of each bracket, Hyperband determines the number of configurations, and the Bayesian optimization component decides which configurations to consider. While initially the latter will suggest configurations randomly, once sufficiently many observations have been made, a surrogate model is fitted to those observations, and it is used to suggest new hyperparameter configurations that maximize the expected improvement acquisition function. The experiments presented by Falkner et al. [7] support that this model-based approach outperforms other model-free baselines such as random search and grid search.

## 4   Related Work

Recent advances in the theoretical and methodological aspects of AutoML have been closely followed by the publication of accompanying software. These open-source frameworks work as a test bench of AutoML theory in the real world. A software package called Auto-WEKA [34] was one of the first attempts to offer an implementation to tackle the combined algorithm selection and hyperparameter optimization (CASH) problem. The package optimizes over the classification models and feature selectors offered by the original WEKA package using the SMAC optimizer detailed above. Hyperopt-Sklearn [19] is another project designed for the CASH problem. The optimization component uses the

Hyperopt library, another implementation of Bayesian optimization, and the baseline models are provided by the popular scikit-learn library. Hyperopt is designed to optimize over the search domain that is generated by the combinations of scikit-learn's preprocessing, classification and regression modules. Similar to Auto-WEKA, the search domain can be comprised of random variables that are sampled and then mapped by an objective function to a scalar score. The score can then be minimized by any of the supported optimizers (Random Search, TPE, Gaussian Process Trees). Auto-sklearn [9] is another AutoML system that uses scikit-learn's implemented preprocessing, classification, and regression modules to define the configuration space. In this implementation, SMAC is used for optimizing over a hypothesis space. The model-based optimization approach was designed to use performance data from similar datasets and to construct ensembles of baseline models evaluated during the optimization state.

All of the tools mentioned above are designed for single-target classification and regression problems. Generalizing to multiple targets, several frameworks have been published and gained attention in their respective fields. MULAN is a popular Java library [35] built on top of the well-known WEKA platform and provides a wide range of multi-label classification and multivariate regression algorithms. Another open-source library called scikit-multilearn [33] was published more recently, providing fewer methods compared to MULAN, but it is written using the more popular Python language. In the area of hierarchical multi-label classification, Clus [37] is a decision tree-based open-source framework that provides a Java interface. Despite the popularity of these tools in specific sub-areas of MTP, none of them offer any automation options in the algorithm selection or hyperparameter optimization steps. In the area of multi-label classification, the work of de Sa et al. [30, 31] uses genetic algorithms in the first attempt at automating the task. Furthermore, Wever et al. proposed an extension of ML-Plan [39], an approach that combines hierarchical task network planning with a best-first search, to configure multi-label classifiers and managed to outperform other baselines in multi-label classification benchmarks, including the ones proposed by de Sa et al. [30, 31]. Finally, CascadeML [26] proposed a cascade neural network that utilizes label associations and requires minimal hyperparameter tuning as another viable benchmark for multi-label classification datasets.

In recent years, machine learning frameworks like Pytorch [27] and Tensorflow [1] have gained considerable popularity in the area of deep learning. As a result, AutoML libraries suited specifically for these deep learning frameworks are now being released. Auto-Net is one of the first attempts at automatically tuning neural network architectures. Its first major release, Auto-Net 1.0 [23], uses the SMAC optimizer and Lasagne [5] as the deep learning framework, a seemingly powerful combination, as it was one of the first to outperform expert users on competition datasets [11]. Auto-Net 2.0, first described in a book chapter [24], was able to bring performance improvements over the first release by replacing SMAC with BOHB. Those ideas were further improved in [40] with the introduction of Auto-Pytorch, a framework that combines ensembling with

**Table 2.** Basic information about the datasets using in the experiments across five MTP problem settings.

| | | # instances | # targets | # instance features | # target features |
|---|---|---|---|---|---|
| **Multi-label classification** | Bibtex | 7395 | 159 | 1836 | 159 |
| | Corel5k | 5000 | 374 | 499 | 374 |
| **Multivariate regression** | Rf2 | 9125 | 8 | 576 | 8 |
| | scm1d | 9803 | 16 | 280 | 16 |
| **Multi-task learning** | dog | 800 | 52 | 3*224*224 | 52 |
| | bird | 2000 | 65 | 3*224*224 | 65 |
| **Matrix completion** | Movielens100k | 943 | 1682 | 943 | 1682 |
| | Movielens1M | 6040 | 3706 | 6040 | 3706 |
| **Dyadic prediction** | ERN | 1164 | 154 | 445 | 445 |
| | SRN | 1821 | 9884 | 113 | 1685 |

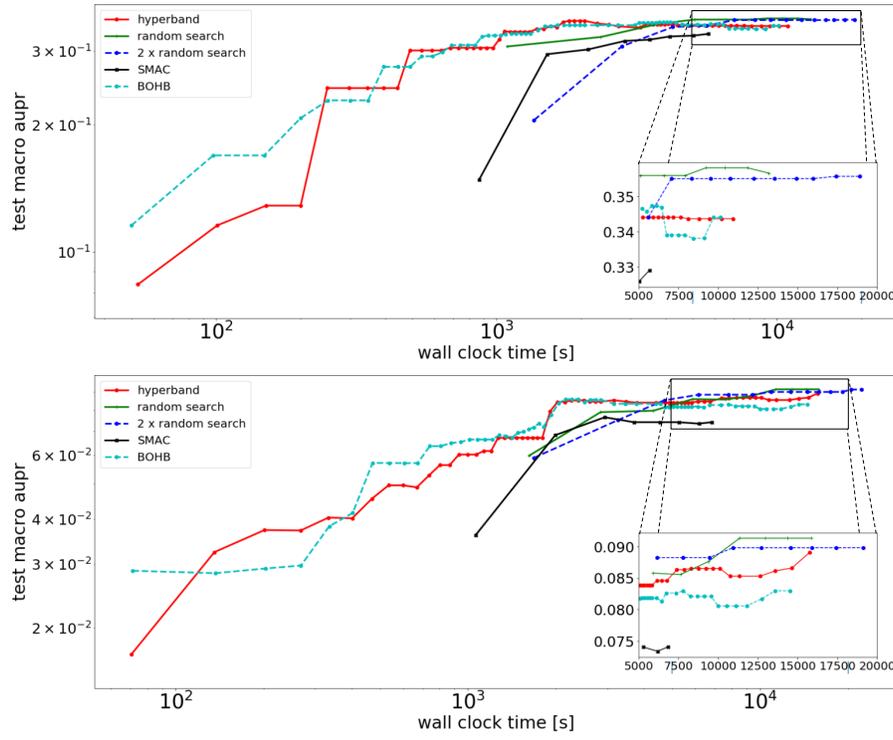multi-fidelity optimization and meta-learning and, as a result, brings significant efficiencies.

To conclude, the aforementioned methods and software packages show that multi-target prediction is a largely unexplored area, which has seen some recent interest in one of its most popular problem settings. In this work, we attempt to set the beginning stages of a similar software package specifically adapted for all the problem settings that fall under the umbrella of MTP. This is achieved by using the automatically answered questionnaire to select the most appropriate MTP setting, and then deploy one of the HPO methods that we benchmark in the next section, on a flexible two-branch neural network architecture.

# 5    Evaluation of HPO methods for DeepMTP

This section's goal is to compare Random Search, Random Search with a doubled budget, Hyperband, BOHB, and SMAC3 across different MTP settings, task types (classification, regression), dataset sizes, and types.

## 5.1    Experiment setup

Basic information about all the datasets used for benchmarking is available in Table 2. Every dataset is split into training and test sets (80-20%). We also randomly sample 20% from the training dataset to form an internal validation set that we use for early stopping and to determine the best configurations while optimizing the network. Also, every HPO run is repeated 5 times, and we report the average performance. The maximum budget allowed for the HPO methods like Hyperband, BOHB, and SMAC is different for every dataset. This parameter is determined by calculating the average best epoch from 20 randomly selected configurations tested before the HPO methods are benchmarked. For Hyperband and BOHB, every other parameter is set to the default values. The configuration space is similar for most of the experiments, with some additional restrictions introduced in cases where a branch encodes one-hot encoded vectors (only one layer allowed). A primary goal of this work is to identify potential cases where

**Fig. 3.** Average test macro-AUPR of every HPO method across time for Bibtex (top) and Corel5k (bottom)

one of the HPO methods could provide a clear advantage. That information can then be used to determine the HPO method suggested by the DeepMTP framework, further minimizing the number of inputs a regular user has to provide. Space constraints do not allow us to present detailed plots for every MTP problem setting and dataset, so in this section, we showcase a limited, yet representative number of examples. Detailed information about the hyperparameter spaces used for every dataset and all additional results and extensive visualizations are publicly available in a web-based application[5].

### 5.2 Comparing performance across time

For the multi-label classification problem setting, we selected two of the largest datasets available in the Mulan repository [35]. The maximum budget allocated for both datasets was set to 27 epochs. In Fig. 3, we observe that random search is quite competitive with Hyperband and BOHB, while SMAC shows the worst

---

[5] https://share.streamlit.io/anonymousmlresearcher/deepmtp_hpo_results/
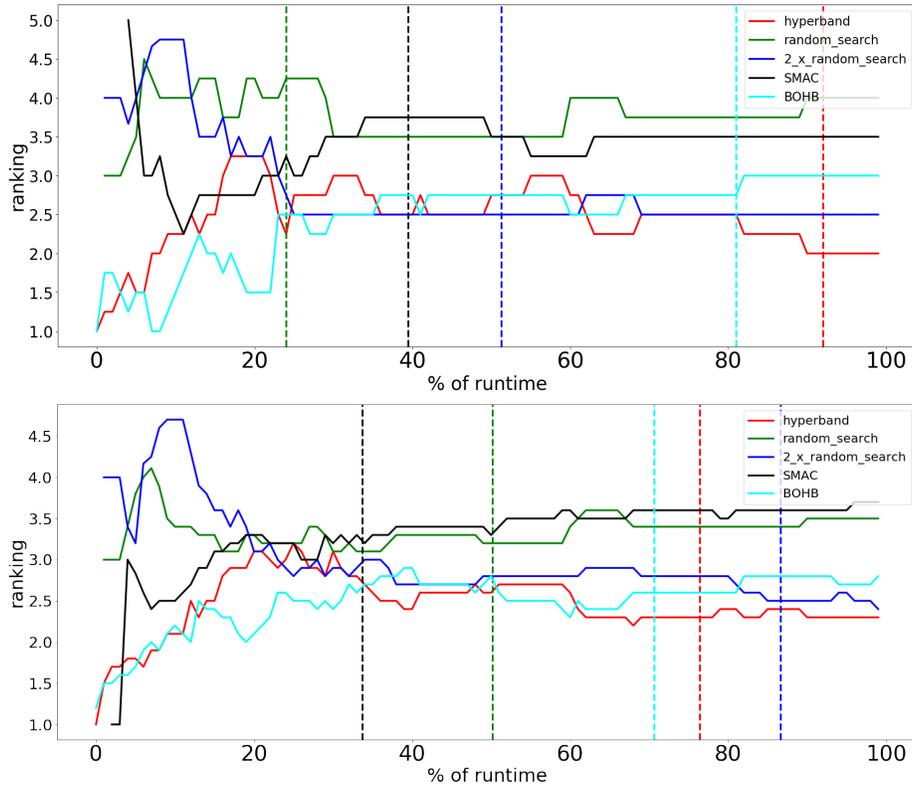main/streamlit_interface.py

performance out of the five candidates. Hyperband provides an early speed-up on the bibtex dataset (2.2 times compared to random search and 5.7 times compared to random search with double the budget) as well as corel5k dataset (1.8 times compared to random search and 2.3 times compared to random search with double the budget). Another interesting fact is that BOHB fails to outperform the standard Hyperband approach, something that might be attributed to the relatively small budget.

In order to generalize across MTP problem settings and provide a more informative comparison between the HPO methods considered, we generated a global ranking plot, shown in the bottom panel of Fig. 4. Based on this ranking, we observe that Hyperband and BOHB outperform the other three approaches for the first 30-40% of the runtime. The same behavior is observed in the average ranking at the classification (multi-label classification, multi-task learning, and dyadic prediction) and regression level (multivariate regression and matrix completion), shown in the top panel of Fig. 4, with Hyperband and BOHB outranking the others. Diving deeper into the rankings per MTP problem setting, Hyperband ranks in the top 2 across all five settings in the early stages of runtime. After 30-40% of the runtime has passed, random search shows a similar performance as Hyperband and BOHB, sometimes outperforming them. This effect can be seen in Fig. 4, as random search achieves a similar ranking as Hyperband at the very end of the total runtime. Because the ranking plots do not inform us about the actual performance difference between the best HPO methods, we decided to generate plots that quantify them. The results show that in the classification and regression settings, the actual performance difference of the top-2 ranked HPO methods is relatively small. Space limitation does not allow us to present these plots here, so we make them available in the aforementioned web-based application[5].

Compared to Hyperband, the similar performance of the two more advanced Bayesian approaches, BOHB and SMAC, likely results from the relatively small budget we assign. Specifically for SMAC, the available budget ranges from eight to 10 configurations across all experiments, which is insufficient for the underlying optimizer. Even though SMAC and BOHB can show potential improvements by increasing the maximum budget, we argue that the standing comparison is fair and more representative of the time constraints that users impose in a real-world environment. The importance of getting relatively good results fairly quickly is the main advantage of Hyperband in the results we have obtained. If the user favors performance over runtime, a random search with a doubled budget can provide a similar or marginally better performance in some MTP problem settings.

## 6   Conclusions

The goal of this paper was to present a fully-automated deep learning pipeline for multi-target prediction by extending the DeepMTP framework with hyperparameter optimization methods. We benchmarked the most popular HPO meth-

**Fig. 4.** Average ranking of every HPO method across the percentage of the respective runtimes for the MTP setting. In the y axis the lower values indicate higher ranking (lower is better). The vertical dotted lines represent the average end-points of every HPO approach. In terms of the performance metrics used to calculate the metrics, we decided to use the most frequently used metrics and averaging schemes for every MTP problem setting (macro-AUPR for the multi-label classification and multi-task learning datasets, micro-AUPR for the dyadic prediction datasets, macro-RRMSE for the multivariate regression datasets, and micro-RMSE for the matrix completion datasets). The top panel visualizes the average ranking over the regression datasets and the bottom panel the global ranking over all datasets.

ods on ten different benchmark datasets from five different MTP settings. An important finding from our experiments is that BOHB and SMAC, two seemingly more advanced methods, do not result in any performance improvements compared to the standard Hyperband. Based on the results, we also conclude that Hyperband is a viable option for the majority of the MTP problem settings that our framework considers, as it provides significant speed-ups compared to the other baselines. Despite the competitiveness of random search, Hyperband is able to return results early in the optimization process, thus providing the

option to stop the optimization if the performance is deemed adequate by the user.

## Declarations

**Funding** This research received funding from the Flemish Government under the "Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen" programme.

**Conflicts of interest** The authors declare that they have no conflct of interest.

**Ethics approval** Not applicable.

**Consent to participate** Not applicable.

**Consent for publication** Not applicable.

**Availability of data and material** The data used for the experiments are available online, see Section 6 for more details.

**Code availability** The code used to run the experiments can be found on github[6].

**Authors' contributions** The first author implemented the python package. All four authors contributed equally to the manuscript.

## References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: TensorFlow: A system for Large-Scale machine learning. In: 12th USENIX symposium on OSDI 16. pp. 265–283 (2016)
2. Bergstra, J., Bengio, Y.: Random search for hyper-parameter optimization. J Mach Learn Res **13**(2) (2012)
3. Bischl, B., Binder, M., Lang, M., Pielok, T., Richter, J., Coors, S., Thomas, J., Ullmann, T., Becker, M., Boulesteix, A., Deng, D., Lindauer, M.: Hyperparameter optimization: Foundations, algorithms, best practices and open challenges. CoRR **abs/2107.05847** (2021), `https://arxiv.org/abs/2107.05847`
4. Dacrema, M.F., Boglio, S., Cremonesi, P., Jannach, D.: A troubling analysis of reproducibility and progress in recommender systems research. ACM TOIS **39**(2), 1–49 (2021)
5. Dieleman, S., Schlüter, J., Raffel, C., Olson, E., Sønderby, S.K., Nouri, D., Maturana, D., Thoma, M., Battenberg, E., Kelly, J., Fauw, J.D., Heilman, M., diogo149, McFee, B., Weideman, H., takacsg84, peterderivaz, Jon, instagibbs, Rasul, D.K., CongLiu, Britefury, Degrave, J.: Lasagne: First release. (Aug 2015). https://doi.org/10.5281/zenodo.27878, `https://doi.org/10.5281/zenodo.27878`
6. Elsken, T., Metzen, J.H., Hutter, F.: Neural architecture search: A survey. J Mach Learn Res **20**(1), 1997–2017 (2019)
7. Falkner, S., Klein, A., Hutter, F.: Bohb: Robust and efficient hyperparameter optimization at scale. In: ICML. pp. 1437–1446. PMLR (2018)
8. Feurer, M., Klein, A., Eggensperger, K., Springenberg, J., Blum, M., Hutter, F.: Efficient and robust automated machine learning. In: NeurIPS. pp. 2962–2970 (2015)

---

[6] `https://github.com/diliadis/DeepMTP`

9. Feurer, M., Klein, A., Eggensperger, K., Springenberg, J.T., Blum, M., Hutter, F.: Auto-sklearn: Efficient and Robust Automated Machine Learning, pp. 113–134. Springer International Publishing, Cham (2019)

10. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. NeurIPS **27** (2014)

11. Guyon, I., Bennett, K., Cawley, G., Escalante, H.J., Escalera, S., Ho, T.K., Macià, N., Ray, B., Saeed, M., Statnikov, A., et al.: Design of the 2015 chalearn automl challenge. In: 2015 International Joint Conference on Neural Networks (IJCNN). pp. 1–8. IEEE (2015)

12. He, X., Liao, L., Zhang, H., Nie, L., Hu, X., Chua, T.S.: Neural collaborative filtering. In: Proceedings of the 26th WWW. pp. 173–182 (2017)

13. He, X., Zhao, K., Chu, X.: Automl: A survey of the state-of-the-art. Knowl. Based. Syst. **212**, 106622 (2021)

14. Hutter, F., Hoos, H.H., Leyton-Brown, K.: Sequential model-based optimization for general algorithm configuration. In: LION. pp. 507–523. Springer (2011)

15. Hutter, F., Kotthoff, L., Vanschoren, J. (eds.): Automated Machine Learning - Methods, Systems, Challenges. The Springer Series on Challenges in Machine Learning, Springer (2019). https://doi.org/10.1007/978-3-030-05318-5, `https://doi.org/10.1007/978-3-030-05318-5`

16. Iliadis, D., De Baets, B., Waegeman, W.: Multi-target prediction for dummies using two-branch neural networks. Mach. Learn pp. 1–34 (2022)

17. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., et al.: Highly accurate protein structure prediction with alphafold. Nature **596**(7873), 583–589 (2021)

18. Karnin, Z.S., Koren, T., Somekh, O.: Almost optimal exploration in multi-armed bandits. In: ICML 2013, Atlanta, GA, USA, 16-21 June 2013. JMLR, vol. 28, pp. 1238–1246. JMLR.org (2013), `http://proceedings.mlr.press/v28/karnin13.html`

19. Komer, B., Bergstra, J., Eliasmith, C.: Hyperopt-sklearn: automatic hyperparameter configuration for scikit-learn. In: ICML workshop on AutoML. vol. 9, p. 50. Citeseer (2014)

20. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. NeurIPS **25** (2012)

21. Li, L., Jamieson, K.G., DeSalvo, G., Rostamizadeh, A., Talwalkar, A.: Hyperband: A novel bandit-based approach to hyperparameter optimization. J. Mach. Learn. Res. **18**, 185:1–185:52 (2017), `http://jmlr.org/papers/v18/16-558.html`

22. Lindauer, M., Eggensperger, K., Feurer, M., Biedenkapp, A., Deng, D., Benjamins, C., Ruhkopf, T., Sass, R., Hutter, F.: Smac3: A versatile bayesian optimization package for hyperparameter optimization (2021)

23. Mendoza, H., Klein, A., Feurer, M., Springenberg, J.T., Hutter, F.: Towards automatically-tuned neural networks. In: Workshop on Automatic Machine Learning. pp. 58–65. PMLR (2016)

24. Mendoza, H., Klein, A., Feurer, M., Springenberg, J.T., Urban, M., Burkart, M., Dippel, M., Lindauer, M., Hutter, F.: Towards automatically-tuned deep neural networks. In: Automated machine learning, pp. 135–149. Springer, Cham (2019)

25. Montgomery, D.C.: Design and analysis of experiments. John wiley & sons (2017)

26. Pakrashi, A., Mac Namee, B.: Cascademl: An automatic neural network architecture evolution and training algorithm for multi-label classification (best technical paper). In: SGAI. pp. 3–17. Springer (2019)

27. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. NeurIPS **32** (2019)
28. Pham, H., Guan, M., Zoph, B., Le, Q., Dean, J.: Efficient neural architecture search via parameters sharing. In: ICML. pp. 4095–4104. PMLR (2018)
29. Rendle, S., Krichene, W., Zhang, L., Anderson, J.: Neural collaborative filtering vs. matrix factorization revisited. In: Fourteenth ACM Conference on RecSys. pp. 240–248 (2020)
30. de Sá, A.G., Freitas, A.A., Pappa, G.L.: Automated selection and configuration of multi-label classification algorithms with grammar-based genetic programming. In: PPSN. pp. 308–320. Springer (2018)
31. de Sá, A.G., Pappa, G.L., Freitas, A.A.: Towards a method for automatically selecting and configuring multi-label classification algorithms. In: GECCO. pp. 1125–1132 (2017)
32. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
33. Szymański, P., Kajdanowicz, T.: A scikit-based Python environment for performing multi-label classification. ArXiv e-prints (Feb 2017)
34. Thornton, C., Hutter, F., Hoos, H.H., Leyton-Brown, K.: Auto-weka: Combined selection and hyperparameter optimization of classification algorithms. In: Proceedings of the 19th ACM SIGKDD. pp. 847–855 (2013)
35. Tsoumakas, G., Spyromitros-Xioufis, E., Vilcek, J., Vlahavas, I.: Mulan: A java library for multi-label learning. J Mach. Learn. Res. **12**, 2411–2414 (2011)
36. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. NeurIPS **30** (2017)
37. Vens, C., Struyf, J., Schietgat, L., Džeroski, S., Blockeel, H.: Decision trees for hierarchical multi-label classification. Machine learning **73**(2), 185–214 (2008)
38. Waegeman, W., Dembczyński, K., Hüllermeier, E.: Multi-target prediction: a unifying view on problems and methods. Data Min Knowl Discov **33**(2), 293–324 (2019)
39. Wever, M.D., Mohr, F., Tornede, A., Hüllermeier, E.: Automating multi-label classification extending ml-plan (2019)
40. Zimmer, L., Lindauer, M., Hutter, F.: Auto-pytorch: multi-fidelity metalearning for efficient and robust autodl. IEEE PAMI **43**(9), 3079–3090 (2021)
41. Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. arXiv preprint arXiv:1611.01578 (2016)