# Study and evaluation of techniques for the optimization of hyperparameters in machine learning algorithms.

Diego Novoa Paradela[1], Antonio Jesús Rivera Rivas[2], María José Del Jesus Díaz[3]
*Universidad Internacional Menéndez Pelayo*
[1]diegonovoaparadela22@gmail.com, [2]arivera@ujaen.es, [3]mjjesus@ujaen.es

*Abstract*—The efficacy of machine learning models relies heavily on the selection of hyperparameters, in addition to their core parameters. While algorithms can provide solutions for parameters, the model's performance hinges largely on the chosen hyperparameters. Determining the optimal combination of hyperparameters is thus key for robust model performance. Hyperparameter Optimization (HPO) is a systematic process aimed at discovering the ideal values for these hyperparameters. Conventional techniques for HPO include grid search and random search, both of which works correctly in most of scenarios. However, a sequential model-based optimization (SMBO) approach has been proposed to tackle runtime issues in production environments and ensure robust performance. In this project, we compare the effectiveness of these traditional methods (Grid Search and Randomized Search) with more recent approaches (Tree-Structured Parzen Estimator (TPE) and Adaptive TPE) leveraging Python libraries: Optuna and Hyperopt. The performance evaluation of these methodologies is based on a extensive suite of real-world datasets. Findings indicate distinct strengths for each approach, underscoring the importance of careful method selection tailored to the specific demands of each machine learning project.

*Index Terms*—Hyper-parameter, auto-tuning, optimization, grid search, randomized search, sequential model-based optimization, hyperopt, optuna

## I. INTRODUCTION

Machine learning models applied to tasks such as classification, regression, and clustering, are governed by parameters and hyperparameters. The parameters associated with the algorithm are those that can be learned by optimizing a loss function or via gradients. For instance, the weights and biases in linear regression can be learned by optimizing a square loss function. Conversely, hyperparameters control the learning process and cannot be learned like parameters during model fitting or loss function optimization. For example, the regularization constant in ridge regression, a term that balances between empirical error and the model's generalization ability, cannot be learned through gradients like the regression parameters but instead controls the entire learning process. Consequently, optimizing parameters alone does not necessarily ensure optimal model performance; a correct choice of hyperparameters is also imperative. Several studies have demonstrated that an optimal set of hyperparameters significantly enhances model performance [1] [2].

Different algorithms have different types of hyperparameters, and their influence on the model's performance can also vary. For instance, in the case of the random forest algorithm, the number of estimators and the depth of trees are hyperparameters that can profoundly influence the model's performance, while the minimum cost-complexity pruning parameter might have a more nuanced effect depending on the noise content in the data [3]. Traditionally, hyperparameter selection in most algorithms involves the analyst, making the process costly and time-consuming, especially when selecting models from a collection of algorithms that are highly sensitive to hyperparameter selection [4]. This highlights the need for systematic hyperparameter optimization techniques that strike a balance between time cost and performance assurance.

In recent years, the realm of machine learning has witnessed unprecedented advancements. This progress spans a range of techniques, each with its unique strengths and applications. While neural networks have demonstrated unparalleled capabilities in tasks such as image classification and natural language understanding [5], they are but one tool within the expansive machine learning toolkit.

Equally influential are most of machine learning techniques such as k-nearest neighbors (knn) and random forests. The former, with its emphasis on instance-based learning, and the latter, leveraging ensemble methods, have each made significant contributions to both research and commercial applications [6]. Each technique, carries its distinct set of challenges. From model initialization to algorithm design, researchers grapple with a host of complexities.

A common thread across these methodologies is the role of hyperparameters. Their influence is pivotal in shaping the efficiency and accuracy of model training. The depth of trees in random forests, the number of neighbors in knn, or the activation functions in neural networks - each hyperparameter has the potential to dramatically impact a model's performance. Thus, their careful selection and tuning become indispensable before diving into the training process.

Due to these challenges, automatic machine learning (AutoML) [7], a technology that designs and trains neural networks automatically at the cost of computational resources, has been proposed. As an integral part of AutoML, hyperparameter optimization (HPO) searches for the optimum hyperparameters for neural network structures and the model training process. Despite demanding significant computational resources, especially when optimizing several hyperparameters together,

HPO reduces the menial work of AI experts, improves the accuracy and efficiency of neural network training, and makes the choice of hyper-parameter set more convincing and the training results more reproducible.

The importance of HPO has been highlighted by two recent trends in the development of deep learning models [8]. The first trend is the upscaling of neural networks for enhanced accuracy, which implies a need for more hyperparameters to tune. The second trend is to design lightweight models to provide satisfactory accuracy with fewer weights and parameters, requiring every hyper-parameter to be tuned to a strict range to reproduce the accuracy. In both cases, manual tuning of hyper-parameters can be a demanding task requiring substantial time, computational resources, and menial work by researchers. The introduction of automated hyper-parameter tuning services and toolkits offers a solution to these challenges.

In light of these complexities and demands, our work seeks to compare classic hyperparameter optimization techniques such as Grid Search and Randomized Search with different approach techniques based on Sequential Model-Based Optimization (SMBO), like TPE and ATPE [9]. The experiment focuses on two libraries: Optuna and Hyperopt, although we also conducted tests with other tools dedicated to HPO such as Scikit-Learn and Keras Tuner. Objective results are obtained using three realistic datasets which tackle different tasks, both binary and multiclass classification. We also applied the techniques to various models, enabling us to assess their performance in different contexts.

The contribution of this paper is a comparative study of different tools for hyperparameter optimization in real-world problems in terms of performance and runtime features. The objective of this research is to conduct a survey on feasible algorithms for HPO and compare leading tools for HPO tasks.

## II. BACKGROUND AND RELATED WORK

In this section, we will delve into the techniques of hyper-parameter optimization and other related works.

### A. Exhaustive Search Techniques

**Grid Search** [10] is a prevalent technique for hyperparameter optimization. At its core, it operates by defining a range of hyperparameters and their potential values, subsequently training models for every combination and selecting the best-performing one. The name derives from the method's manner of laying out all combinations on a Cartesian grid. One of its major downsides is the curse of dimensionality. As you add more dimensions or hyperparameters, or increase the number of potential values for a hyperparameter, the number of evaluations grows exponentially, making this method computationally expensive.

Unlike Grid Search, **Random Search** [11] randomizes the selection within the hyperparameters' potential values, estimating a predefined number of models. One of its prominent advantages is that it avoids the exhaustive evaluation of every possible combination, which can be particularly beneficial in scenarios where the hyperparameters have varying ranges of
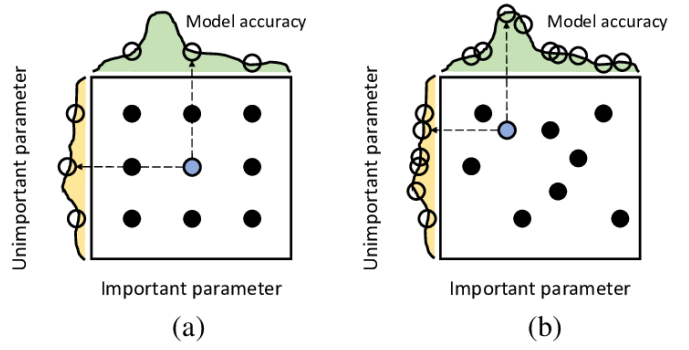


Fig. 1. Comparison between a-grid search and b-random search for hyper parameter tuning.

potential values or when certain hyperparameters exert a more substantial influence on the model's performance compared to others.

The key distinction lies in the way Random Search leverages randomness to explore the hyperparameter space more efficiently. By avoiding exhaustive exploration, it tends to focus computational resources on the most promising areas within the hyperparameter space, offering a more balanced trade-off between exploration and exploitation. This adaptability makes Random Search a compelling choice, especially when dealing with high-dimensional search spaces or when there's limited knowledge about the relative importance of different hyperparameters.

On the other hand, Grid and Random Search lack in their ability to use feedback from previous searches to inform their subsequent choices, leading them often to scout in less promising regions of the search space.

While manual tuning might occasionally surpass these methods, the Sequential model-based optimization addresses this limitation by maintaining a history of configurations and their respective scores, thereby adapting its behavior after each iteration.

In Figure 1 we can see the comparison between (a) grid search; and (b) random search for hyperparameter fitting. The nine dots indicate the candidates. The curves on the left and top indicate the accuracy of the model as a function of each search dimension. It is possible to observe how in graph (b) a broader candidate search is performed than in graph (a) by selecting random samples that do not follow a regular distribution, thus taking advantage of the highest possible performance of the model.

### B. Sequential Model-Based Optimization (SMBO)

SMBO [12] techniques are grounded in the Bayesian Optimization framework. By employing a surrogate model and an acquisition function, these methods aim to iteratively choose the most promising hyperparameters to approximate the real objective function.

In **Bayesian Optimization** [13], Gaussian Processes (GP) are renowned for their ability to approximate complex functions. They determine a prior set of functions, then use kernel

functions on sampled values to obtain a posterior. As the number of iterations increases, the uncertainty of the GP decreases, particularly in areas that have been well-explored. The acquisition function assists the GP by continuously selecting areas that minimize the function. While understanding GPs deeply can be challenging, they touch upon several vital statistical concepts. For a more in-depth understanding, multiple resources delve into the intricacies of GPs [14].

While GPs model the loss in the light of hyperparameter configurations, **TPE**s [15] take a different approach. The method was initially designed to handle hyperparameters of a deep neural network which were continuous, categorical, and conditional on one another. TPE uses random samples to create an initial set of losses. It then segregates the losses into a 'good' group based on a threshold and the rest are placed into the 'remainder group'. From there, two densities are constructed to model p(x—y). As iterations proceed, the densities adjust to better approximate the real loss function.

An advanced extension of TPE, known as ATPE (Adaptive Tree-structured Parzen Estimators) algorithm, enhances the performance of the original method by dynamically adjusting its sampling strategy based on the success of past iterations. The Parzen estimators incorporated in both TPE and ATPE enable them to effectively navigate conditional hyperparameter spaces, making them particularly suitable for optimizing hyperparameters in complex deep learning contexts.

## III. Hyperparameter Optimization Tools

In this section, we will briefly describe the tools we have utilized. Although tests were conducted across all the libraries mentioned below, for the purpose of an exhaustive study, we primarily focused on the libraries HyperOpt and Optuna.

### A. HyperOpt

HyperOpt [16] is a Python library that operates based on Sequential Model-Based Optimization (SMBO), tied with the Bayesian optimization discussed in the previous section. It offers an interface for users to configure the search space of variables, customize an evaluation function, and set the loss function to assess points within the search space. Moreover, it can be leveraged to structure a single large hyperparameter optimization problem combining various algorithms, data preprocessing modules, and their hyperparameters.

The core components of HyperOpt are a search domain, an objective function, and an optimization algorithm. The search domain caters to various optimization needs as it can be characterized by continuous, ordinal, or categorical variables. The objective function can be a user-defined Python function that takes in the variables and returns a loss function corresponding to that variable combination. The search algorithm aligns with the surrogate function discussed previously, and tool-driven options include random search, TPE, and adaptive TPE. Additionally, HyperOpt offers support for parallel implementations using Apache Spark and MongoDB.

### B. Optuna

Optuna's distinctive features, as described by its authors, are: "(1) A define-by-run API that allows users to construct the parameter search space dynamically, (2) an efficient implementation of search and pruning strategies, and (3) a versatile and easy-to-setup architecture that can be deployed for diverse purposes, from large-scale distributed computations to lightweight experiments via an interactive interface."

The define-by-run API sets Optuna apart, contrasting with HyperOpt where users need to predefine search spaces for each hyperparameter. Optuna's objective function gets a trial object embedded with the parameter space and the function to optimize, rather than direct hyperparameter values. This API further promotes modular programming.

Optuna supports various search strategies, including Random Search, Grid Search, and Tree-structured Parzen Estimators. In hyperparameter optimization problems, handling both relational and independent sampling is vital. Relational sampling pertains to the existing correlations among parameters, while independent sampling operates in isolation. Optuna can detect experiment runs that harbor these concurrent relationships.

Optuna [17] introduces a pruning feature that aids in prematurely terminating non-optimal runs. By monitoring intermediate objective values, runs that don't meet predefined conditions are halted. This capability is enabled through an asynchronous successive halving algorithm. It also brings to the table distributed computation capabilities.

### C. Scikit-Learn

Scikit-Learn [18] stands as one of the most widely employed Python libraries for machine learning. It encompasses several functions for hyperparameter optimization, including grid search (GridSearchCV) and random search (RandomizedSearchCV).

These functionalities empower users to define a set of potential values for each hyperparameter to be optimized. Subsequently, they train and evaluate a model for every possible hyperparameter combination (in the GridSearchCV case) or for a predetermined number of random hyperparameter combinations (using RandomizedSearchCV). Both functions also uphold cross-validation, facilitating a more robust evaluation of each hyperparameter set's performance.

### D. Keras Tuner

Keras Tuner [19] is a Python library tailored for hyperparameter optimization of deep learning models constructed with Keras. The library incorporates various techniques for hyperparameter search, including random search, Hyperband, and Bayesian search.

Keras Tuner delivers a user-friendly interface, enabling users to specify the hyperparameter search space, set the maximum number of trials, and define the objective function for optimization. In addition, it features visualization tools for the results of hyperparameter search, simplifying result interpretation.

A salient feature of Keras Tuner is its compatibility with TensorFlow 2.0 and Keras, allowing users to effortlessly optimize hyperparameters for advanced deep learning models.

## IV. EXPERIMENTS

In this section, we describe the details of the experiments, datasets used, system setup and metrics used for the classification tasks.

### A. Datasets

Three fully available datasets from different sources have been used. Specifically, they have been accessed through The UCI Machine Learning Repository [20]. The UCI Machine Learning Repository is a collection of databases, domain theories, and data generators that are used by the machine learning community for the empirical analysis of machine learning algorithms.

*1) Covertype Dataset*

**Objective:** The primary aim of this dataset is to classify pixels into one of the seven forest cover types using features related to geography and soil type.

*a) Characteristics and Background*

- **Dataset Details:** The dataset is multivariate with its primary focus on the Life Science domain. The tasks associated with this dataset involve classification. The attribute types present are both categorical and integer. With 7 different classes, this dataset is clearly multiclass. Most of the features are numerical (10 in total), indicating that the data focus on quantitative aspects of the environment, such as elevation, distance to geographical features, among others. However, there are also 3 categorical attributes that could include information such as soil types and desert areas.
- **Size and Attributes:** It consists of 581,012 instances and 54 attributes.
- **Source of Data:** The data are derived from cartographic variables, with no remotely sensed data utilized. The actual forest cover type for a 30x30 meter cell was ascertained from the US Forest Service (USFS) Region 2 Resource Information System (RIS) data. Independent variables were extracted from data originally sourced from the US Geological Survey (USGS) and USFS.
- **Description:** The region under study includes four wilderness areas situated in the Roosevelt National Forest of northern Colorado.
  - **Neota:** Highest mean elevation. Primary tree species: Spruce/Fir.
  - **Rawah & Comanche Peak:** Moderate mean elevation. Dominant tree species: Lodgepole pine, followed by Spruce/Fir and Aspen.
  - **Cache la Poudre:** Lowest mean elevation. Predominant tree species: Ponderosa pine, Douglas-fir, and Cottonwood/Willow.

*2) Adult Dataset*

**Objective:** The goal of this dataset is to predict whether an individual's income surpasses $50K/year based on census data.

*a) Characteristics and Background*

- **Dataset Details:** This multivariate dataset is situated in the Social Science domain. The task associated is classification. Attributes present are both categorical and integer. This is a binary dataset with two classes. Unlike the Covertype dataset, this dataset has a higher proportion of categorical attributes (9 in total) compared to numerical attributes (6 in total). This suggests that demographic and possibly occupational information plays a crucial role in this dataset.
- **Size and Attributes:** The dataset encompasses 48,842 instances with 14 attributes.
- **Source and Extraction:** The data were extracted by Barry Becker from the 1994 Census database. The records chosen adhere to specific conditions outlined by: $((AAGE > 16) \&\& (AGI > 100) \&\& (AFNLWGT > 1) \&\& (HRSWK > 0))$.

*3) Bank Marketing Dataset*

**Objective:** This dataset is associated with the direct marketing campaigns (primarily phone calls) of a Portuguese banking institution. The classification task is to foresee whether a client will opt for a term deposit.

*a) Characteristics and Background*

- **Dataset Details:** The dataset is multivariate, hailing from the Business domain. The tasks associated are classification. Attributes present are both categorical and integer. Like the Adult dataset, this dataset is binary. It has a more even balance between numerical and categorical attributes, with a slight bias towards categorical attributes. This reflects the nature of information that is often relevant in marketing decisions, such as demographic and customer behavioural information.
- **Size and Attributes:** There are four versions of this dataset with varying instance and attribute counts. The primary dataset consists of 45,211 instances and 17 attributes. However, there are other versions with 20 inputs, as well as subsets containing 10
- **Description:** The data revolve around direct marketing campaigns led by a Portuguese bank. These campaigns, which are primarily phone-based, often require multiple calls to a client to determine if a product (bank term deposit) will be subscribed to.

TABLE I
OVERVIEW OF DATASETS

| Dataset | No. Classes | Instances | Num Attr | Cat Attr |
|---|---|---|---|---|
| Covertype | 7 | 581012 | 10 | 3 |
| Adult | 2 | 48842 | 6 | 9 |
| Bank Marketing | 2 | 45211 | 7 | 10 |

**Table I** describes the main characteristics of the datasets.

All three datasets offer a variety of scenarios for the application of machine learning techniques. While the Covertype and Bank Marketing datasets offer a balance between numerical and categorical features, the Adult dataset is more biased towards the categorical. The size and complexity of each dataset makes it unique and suitable for different applications and experiments.

### B. System setup

For the implementation of the project, Google Colab has been used as the development platform. This is because the exhaustive study of various techniques for hyperparameter optimization requires significant computational resources, and Colab is a Jupyter notebook that allows for the interpretation of Python 3 code using a VM with hardware provided by Google. Specifically, an Intel(R) Xeon(R) CPU @ 2.20GHz and 12 GB of RAM were used.

K-fold cross-validation (K = 3) was employed, along with several evaluation metrics: F1-Score, Accuracy, Recall, and Execution Time.

The algorithms are evaluated over 10 Trials for the selection of the best hyperparameter combination.

On the other hand, for the development of the comparative study we have focused on the algorithms available in the two main hyperparameter optimization libraries: Optuna and Hyperopt. In the case of Optuna, the three main algorithms have been selected: Grid Search, Random Search and Tree-structured Parzen Estimator algorithm. This selection is motivated by the comparison of two paradigms: the first two (traditional) algorithms and the new Sequential Model-Based Optimisation perspective where TPE stands out. Moreover, both the latter algorithm and its adaptive version are the focus of the HyperOpt library, which only has these two algorithms and Random Search available. This is why the Grid Search algorithm has not been evaluated when using the HyperOpt library.

### C. Models

Models Under Comparison:

*1) Random Forest*

**Random Forest** is a versatile ensemble method that utilizes a multitude of decision trees for predictive tasks. The configuration of this model is crucial for its performance. For each tree in the forest, it is typically trained on a 'bootstrap sample' of the dataset, where data instances are drawn with replacement. The depth of each tree (max_depth) determines the maximum number of levels it can have. If not explicitly set, trees can expand until they minimize the criterion or until they have fewer than the minimum samples required to split a node. The number of trees in the forest is dictated by n_estimators, with more trees generally providing better learning at the expense of computational cost. For every decision a tree makes (every node), it only considers a random subset of features, which ensures diversity among trees and reduces overfitting. Other hyperparameters like min_samples_split determine the minimum number of samples required to split an internal

TABLE II
HYPERPARAMETERS SEARCH SPACE

| Librería | Modelo | Hiperparámetro | Espacio de búsqueda |
|---|---|---|---|
| Optuna | Neural Network | n_units<br>optimizer<br><br>lr<br><br><br><br>batch_size | 10 to 50 (step 10)<br>Adam, SGD, RM-Sprop<br>10 logarithmically spaced values between $1e-5$ and $1e-1$<br>16, 32, 64, 128, 256 |
| | Knn | n_neighbors<br>p | 1 to 10 (step 2)<br>1 (Manhattan), 2 (Euclidean) |
| | Random Forest | n_estimators<br>max_depth<br>min_samples_split<br>min_samples_leaf | 10 to 50 (step 10)<br>2 to 16 (step 2)<br>2 to 6<br>1 to 4 |
| Hyperopt | Neural Network | n_units<br>optimizer<br><br>lr<br><br><br>batch_size | 10 to 50 (step 1)<br>Adam, SGD, RM-Sprop<br>Loguniform between $1e-5$ and $1e-1$<br>16, 32, 64, 128, 256 |
| | Knn | n_neighbors<br>p | 1 to 10 (step 2)<br>1 (Manhattan), 2 (Euclidean) |
| | Random Forest | n_estimators<br>max_depth<br>min_samples_split<br>min_samples_leaf | 10 to 50 (step 10)<br>2 to 16 (step 2)<br>2 to 6<br>1 to 4 |

node, and min_samples_leaf denotes the minimum number of samples needed to be at a leaf node. The decision to split at any node is made based on gini impurity or information gain, aiming to maximize the homogeneity of nodes. The Random Forest model can be parallelized with n_jobs=-1, utilizing all available cores of the machine for faster computation.

*2) K-Nearest Neighbors (KNN)*

**K-Nearest Neighbors** is an instance-based algorithm that operates by storing all available instances and classifying new instances based on a similarity measure, like distance between points. The outcome of a new instance is determined by the majority vote from its 'k' closest neighbors. The choice of the distance metric plays a pivotal role, with common metrics including the Minkowski distance where p=1 is Manhattan and p=2 is Euclidean distance.

*3) Neural Network*

The **Neural Network** model is designed with an input layer that directly accepts feature input, automatically adapting to the number of features in the given dataset. Following this is a single hidden layer that employs the ReLU (Rectified Linear Unit) activation function, which transforms each input to its positive value, aiding in introducing non-linearities without hindering learning capability or convergence speed. The architecture concludes with an output layer composed of

a number of nodes, each representing a class, and uses the softmax activation function to convert the model outputs into class probabilities in case the dataset is multi-class, or the sigmoidal function in case the dataset is prepared for a binary classification task. Depending on the chosen configuration, the model can be trained using various optimizers including Adam, SGD, or RMSprop, each with its distinct mechanisms in terms of learning rate adaptation and weight updates during training.

In our study, we aim to juxtapose these three distinct Machine Learning models. Through rigorous experimentation and validation, we seek to unveil the strengths, weaknesses, and best-fit scenarios for each, guiding future research and applications.

**Table II** meticulously outlines the hyperparameter search spaces utilized by two renowned optimization libraries, Optuna and Hyperopt, for tuning three distinct machine learning models: Neural Networks, K-Nearest Neighbors (Knn), and Random Forests. In the case of Neural Networks, Optuna's search space for the n_units parameter, which denotes the number of units, spans from 10 to 50 with increments of 10, aiming to identify optimal neural network sizes that range from smaller to moderately-sized networks. It experiments with three popular optimizers: Adam, SGD (Stochastic Gradient Descent), and RMSprop. The learning rate (lr) in Optuna is explored across 10 logarithmically spaced values between $1e-5$ and $1e-1$. Batch sizes being considered include 16, 32, 64, 128, and 256.

On the other hand, Hyperopt offers a more granular approach for the n_units parameter, covering the same range but incremented in steps of just 1. The optimizers and batch sizes remain consistent with Optuna's selections. However, for the learning rate, Hyperopt leverages a loguniform distribution between $1e-5$ and $1e-1$, potentially pinpointing a more precise optimal value.

For K-Nearest Neighbors, both Optuna and Hyperopt probe the n_neighbors hyperparameter from values 1 to 10, incremented by 2, and evaluate two distance metrics, Manhattan (p=1) and Euclidean (p=2).

In the context of Random Forests, both libraries experiment with the n_estimators parameter that indicates ensemble size, ranging from 10 to 50 in increments of 10. They also examine tree depth, with max_depth values spanning from 2 to 16 in increments of 2, and tinker with min_samples_split (ranging from 2 to 6) and min_samples_leaf (ranging from 1 to 4), ensuring an array of tree configurations.

In essence, this comprehensive table provides deep insights into the choices made in hyperparameter exploration for each model by both optimization tools. It underscores the varied strategies in Neural Networks while revealing consistency in approaches for Knn and Random Forests across the two libraries.

### D. Results

Within the realm of machine learning, the careful tuning and optimization of hyperparameters is often the difference between mediocre and state-of-the-art performance. In this section, we will shed light on the crucial role of hyperparameter optimization in shaping the results of our various machine learning techniques employed.

While models provide the foundational algorithms to process and learn from data, hyperparameters govern their learning processes, adjusting and refining them. This comparative analysis offers insights into how different hyperparameter configurations, coupled with diverse optimization strategies, can lead to significant fluctuations in performance. The results presented here illuminate the balance between computation time and accuracy, and highlight the paramount importance of rigorous hyperparameter tuning.

Through the detailed tables, we present a meticulous breakdown of our findings. As we explore this section, it becomes evident that a model's success isn't solely based on its inherent algorithmic structure, but heavily reliant on the optimal configuration of its hyperparameters.

TABLE III
OPTIMIZATION RESULTS WITH OPTUNA

| Dataset | Model | Algorithm | Accuracy | F1-Score | Recall | Time |
|---|---|---|---|---|---|---|
| Covertype | Neural Network | Grid Search | 0.79 | 0.78 | 0.79 | 20 |
| Covertype | Neural Network | Random Search | 0.78 | 0.77 | 0.78 | 19 |
| Covertype | Neural Network | TPE | 0.77 | 0.77 | 0.77 | 17 |
| Covertype | KNN | Grid Search | 0.93 | 0.91 | 0.92 | 120 |
| Covertype | KNN | Random Search | 0.93 | 0.92 | 0.92 | 90 |
| Covertype | KNN | TPE | **0.93** | **0.94** | **0.92** | **334** |
| Covertype | Random Forest | Grid Search | 0.85 | 0.86 | 0.85 | 23 |
| Covertype | Random Forest | Random Search | 0.86 | 0.87 | 0.86 | 17 |
| Covertype | Random Forest | TPE | 0.86 | 0.87 | 0.86 | 18 |
| Adult | Neural Network | Grid Search | 0.86 | 0.86 | 0.86 | 2 |
| Adult | Neural Network | Random Search | 0.86 | 0.86 | 0.86 | 2 |
| Adult | Neural Network | TPE | 0.87 | 0.86 | 0.87 | 1.5 |
| Adult | KNN | Grid Search | 0.85 | 0.84 | 0.85 | 10 |
| Adult | KNN | Random Search | 0.85 | 0.84 | 0.85 | 9 |
| Adult | KNN | TPE | 0.85 | 0.84 | 0.85 | 8 |
| Adult | Random Forest | Grid Search | **0.87** | **0.86** | **0.87** | 1 |
| Adult | Random Forest | Random Search | 0.86 | 0.86 | 0.86 | 1 |
| Adult | Random Forest | TPE | 0.86 | 0.86 | 0.86 | 1 |
| Bank | Neural Network | Grid Search | 0.90 | 0.30 | 0.20 | 1 |
| Bank | Neural Network | Random Search | 0.90 | 0.30 | 0.20 | 1.5 |
| Bank | Neural Network | TPE | 0.90 | 0.30 | 0.19 | 2 |
| Bank | KNN | Grid Search | 0.89 | 0.87 | 0.89 | 4 |
| Bank | KNN | Random Search | 0.89 | 0.87 | 0.89 | 3 |
| Bank | KNN | TPE | 0.89 | 0.87 | 0.89 | 3 |
| Bank | Random Forest | Grid Search | **0.90** | **0.87** | **0.90** | 1 |
| Bank | Random Forest | Random Search | **0.90** | **0.87** | **0.90** | 1 |
| Bank | Random Forest | TPE | **0.90** | **0.87** | **0.90** | **1** |

**Table III** provides a comprehensive summary of the results obtained through the optimization process with Optuna across three datasets: Covertype, Adult, and Bank. For each dataset, the performances of three machine learning models – Neural Network, KNN, and Random Forest – are showcased. Each model was subjected to different optimization strategies including Grid Search, Random Search, and TPE. The key metrics detailed in the table are Accuracy, F1-Score, Recall, and the Time taken for the optimization process. Notably, the highest values among the metrics for each dataset have been highlighted in bold, emphasizing the best-performing configurations.The Covertype dataset, for instance, demonstrates peak results with the KNN model optimized through TPE, while both the

Adult and Bank datasets exhibit their highest performances with the Random Forest model.

TABLE IV
OPTIMIZATION RESULTS WITH HYPEROPT

| Dataset | Model | Algorithm | Accuracy | F1-Score | Recall | Time |
|---------|-------|-----------|----------|----------|--------|------|
| Covertype | Neural Network | TPE | 0.80 | 0.79 | 0.80 | 11 |
| Covertype | Neural Network | Random Search | 0.80 | 0.79 | 0.80 | 12 |
| Covertype | Neural Network | Adaptative TPE | 0.80 | 0.79 | 0.80 | 13 |
| Covertype | KNN | TPE | 0.91 | 0.93 | 0.91 | 528 |
| Covertype | KNN | Random Search | **0.91** | **0.93** | **0.91** | **325** |
| Covertype | KNN | Adaptative TPE | 0.91 | 0.93 | 0.91 | 540 |
| Covertype | Random Forest | TPE | 0.85 | 0.85 | 0.85 | 16 |
| Covertype | Random Forest | Random Search | 0.85 | 0.85 | 0.85 | 13 |
| Covertype | Random Forest | Adaptative TPE | 0.85 | 0.85 | 0.85 | 18 |
| Adult | Neural Network | TPE | 0.86 | 0.86 | 0.86 | 1 |
| Adult | Neural Network | Random Search | 0.86 | 0.86 | 0.86 | 1 |
| Adult | Neural Network | Adaptative TPE | 0.86 | 0.86 | 0.86 | 2 |
| Adult | KNN | TPE | 0.85 | 0.84 | 0.85 | 8 |
| Adult | KNN | Random Search | 0.85 | 0.84 | 0.85 | 6 |
| Adult | KNN | Adaptative TPE | 0.85 | 0.84 | 0.85 | 10 |
| Adult | Random Forest | TPE | 0.87 | 0.86 | 0.87 | 0.5 |
| Adult | Random Forest | Random Search | **0.87** | **0.86** | **0.87** | **0.33** |
| Adult | Random Forest | Adaptative TPE | 0.87 | 0.86 | 0.87 | 0.58 |
| Bank | Neural Network | TPE | 0.90 | 0.86 | 0.90 | 1 |
| Bank | Neural Network | Random Search | 0.90 | 0.86 | 0.90 | 1 |
| Bank | Neural Network | Adaptative TPE | 0.90 | 0.86 | 0.90 | 1.5 |
| Bank | KNN | TPE | 0.89 | 0.87 | 0.89 | 3 |
| Bank | KNN | Random Search | 0.89 | 0.87 | 0.89 | 2.5 |
| Bank | KNN | Adaptative TPE | 0.89 | 0.87 | 0.89 | 4 |
| Bank | Random Forest | TPE | 0.90 | 0.87 | 0.90 | 0.6 |
| Bank | Random Forest | Random Search | **0.90** | **0.87** | **0.90** | **0.47** |
| Bank | Random Forest | Adaptative TPE | 0.90 | 0.87 | 0.90 | 0.75 |

**Table IV** presents the optimization results when utilizing Hyperopt on the three datasets: Covertype, Adult, and Bank. The table elucidates the performance metrics of three machine learning models — Neural Network, KNN, and Random Forest — each fine-tuned using three distinct optimization techniques provided by Hyperopt: TPE, Random Search, and the more advanced Adaptative TPE. As a reader navigates through the table, they can observe the metrics such as Accuracy, F1-Score, Recall, and the time consumed for optimization (in Time column). Bolded values in the metrics underline the top-performing configurations for each dataset. For instance, in the Covertype dataset, the KNN model optimized using Random Search proved to be the most time-efficient while maintaining a peak performance. In contrast, for both the Adult and Bank datasets, the Random Forest model using Random Search showcased superior results with the least time overhead. The inclusion of the Adaptative TPE strategy offers an interesting comparative perspective, as it is a refined version of the traditional TPE.

*1) Analysis of Optuna Results*

*a) Overview*

The data presents results of hyperparameter tuning using Optuna over different datasets, models, and search algorithms. Three datasets are considered: *Covertype*, *Adult*, and *Bank*. The models tested on these datasets include *Neural Network*, *KNN*, and *Random Forest*. For each combination of dataset and model, three search algorithms were applied: *Grid Search*, *Random Search*, and *Tree-Structured Parzen Estimator (TPE)*.

*b) Accuracy and Metrics Overview*

Overall, all models consistently showed high accuracy scores across datasets. Specifically, models trained on the *Covertype* dataset generally exhibited accuracy between 0.7 and 0.9, while for the *Adult* and *Bank* datasets, accuracy predominantly ranged between 0.85 and 0.89. The metrics *F1-Score* and *Recall* were mostly in line with

the *Accuracy* metric, indicating a balanced performance for most combinations.

*c) Neural Network Analysis*
- For the *Covertype* dataset, the performance of the Neural Network remained largely consistent across the search algorithms, with a slight drop in the TPE.
- For the *Adult* dataset, the performance remained consistent across the search algorithms.
- For the *Bank* dataset, while the accuracy was high, the f1_score and recall were particularly low, suggesting that the model may not be effectively capturing true positive cases.

*d) KNN Analysis*
- The KNN model, when trained on the *Covertype* dataset, performed exceptionally well, particularly with the TPE method where the F1-score peaked at 0.94.
- For the *Adult* dataset, the results were consistent across the search algorithms, hovering around the mid-80s for all metrics.
- The *Bank* dataset also showcased consistent results for the KNN model across search methods.

*e) Random Forest Analysis*
- The Random Forest model's results on the *Covertype* dataset were consistent and impressive, with little variation across search algorithms.
- On the *Adult* dataset, results were again consistent across algorithms.
- The performance on the *Bank* dataset mirrored that of the *Adult* dataset, with all metrics in the high 80s to 90 range.

*f) Search Algorithms Analysis*
- Grid Search often took the longest time, particularly with the KNN model on the *Covertype* dataset, requiring 120 minutes. This indicates that while grid search can be thorough, it can also be computationally expensive.
- Random Search generally provided similar results in a shorter timeframe compared to Grid Search, as seen with the *Covertype* KNN model, where the performance metrics remained consistent but at a much-reduced execution time.
- TPE, while producing similar performance metrics, often required varied time, sometimes outpacing and other times lagging behind the other methods.

*g) Summary*
- The Neural Network model for the *Bank* dataset consistently had a significant drop in F1-Score and Recall, despite having a high accuracy. This indicates potential issues in capturing true positive cases effectively, which may stem from imbalances in the dataset or the model's architecture.
- The KNN model for the *Covertype* dataset with the TPE search took an astounding 334 minutes, despite offering only a marginal improvement in the F1-Score. This raises questions about the efficiency and viability of such a method given the time investment.
- About the optimization methods that have worked the best: For the Covertype dataset, TPE had a slightly higher margin with the KNN model, and Random Search had a higher margin with Random Forest. In the Adult dataset, TPE performed better with the Neural Network model, and Grid Search was superior with Random Forest. Finally, in the Bank dataset, all methods performed very similarly across all models, although it is important to mention that the Neural Network had accuracy problems in F1-Score and Recall.

*2) Analysis of Hyperopt Results*

*a) Overview*

The data illustrates the outcomes of hyperparameter optimization using Hyperopt across various datasets and models. The datasets under consideration are: *Covertype*, *Adult*, and *Bank*. Evaluated models for these datasets encompass the *Neural Network*, *KNN*, and

*Random Forest*. Each combination of dataset and model underwent optimization via two techniques: *Tree-Structured Parzen Estimator (TPE)* and its *Adaptive* counterpart.

### b) Accuracy and Metrics Overview

Across the board, models yielded consistently high accuracy scores over all datasets. Specifically, on the *Covertype* dataset, models generally presented accuracy within the high 70s to mid 90s spectrum. As for the *Adult* and *Bank* datasets, accuracy prominently settled in the mid to high 80s. The metrics *F1-Score* and *Recall* closely echoed the trends seen in the *Accuracy* metric, denoting a balanced model performance across different configurations.

### c) Neural Network Analysis

- For the *Covertype* dataset, the Neural Network's performance with Hyperopt demonstrated marked consistency, regardless of the optimization technique employed.
- The *Adult* dataset displayed similar trends, with a balanced performance across different optimization strategies.
- On the *Bank* dataset, while accuracy metrics were notably high, the f1_score and recall occasionally faltered, hinting at a potential lack in capturing true positive instances.

### d) KNN Analysis

- KNN, when evaluated on the *Covertype* dataset, produced stellar results. The Adaptive TPE method, in particular, saw the F1-score reach impressive heights.
- On the *Adult* dataset, metrics remained stable regardless of the optimization method, with figures typically hovering in the mid-80s.
- Consistency marked the KNN results on the *Bank* dataset, with both optimization techniques yielding similar outcomes.

### e) Random Forest Analysis

- The *Covertype* dataset showcased consistent and impressive Random Forest outcomes, with little variation between optimization strategies.
- The *Adult* dataset mirrored these results, maintaining stability across techniques.
- Outcomes for the *Bank* dataset followed a similar pattern, placing metrics predominantly in the high 80s to low 90s range.

### f) Optimization Techniques Analysis

- While traditional TPE frequently surpassed Grid Search in terms of execution time and metrics, the Adaptive TPE sometimes necessitated longer runtimes.
- Despite sometimes lagging in terms of time, the Adaptive TPE's metrics closely paralleled those of its traditional counterpart, highlighting the balance between thoroughness and computational demand.

### g) Summary

- The Neural Network model's performance on the *Bank* dataset, particularly in terms of F1-Score and Recall, indicates potential challenges in accurately identifying true positive cases. This might be attributed to dataset imbalances or specific nuances in the model's structure.
- The *Covertype* dataset's KNN model, when optimized using TPE and Adaptive TPE, demanded a substantial runtime, even if the metric improvements were marginal. This brings the efficiency and feasibility of this method with this model into question, especially considering the significant time commitment.
- About the optimization methods that have worked the best: For the Covertype dataset, all the Neural Network models displayed a consistent performance, regardless of the optimization method. KNN, using Random Search, performed best in terms of time efficiency, even though the performance metrics were consistent with other algorithms. The Random Forest model had comparable results across different optimization algorithms. In the Adult dataset, the Neural Network and KNN models didn't show

significant variations across optimization methods. However, the Random Forest model achieved its best time efficiency using Random Search, with performance metrics remaining constant across the board. For the Bank dataset, Neural Network models remained consistent across all optimization methods. KNN with Random Search stood out as the most time-efficient, while performance metrics remained similar across methods. The Random Forest model optimized using Random Search was the most time-efficient, although performance metrics did not differ significantly across algorithms.

## 3) Final Comparison Between Optuna and Hyperopt

### a) Overview

Both Optuna and Hyperopt are advanced libraries designed for hyperparameter optimization, employed across various datasets and models including *Neural Network*, *KNN*, and *Random Forest*. Optuna provides a broader range of search algorithms like *Grid Search*, *Random Search*, and *TPE*, while Hyperopt focuses primarily on *Tree-Structured Parzen Estimator (TPE)* and its adaptive version.

### b) Ease of Use

Optuna tends to offer a more user-friendly interface with built-in visualization tools. Hyperopt, on the other hand, may require additional steps to summarize and visualize results, but it often provides more control over the underlying optimization algorithms.

### c) Performance

Across the datasets and models provided, both Optuna and Hyperopt exhibited commendable optimization capabilities. Specifically, while most metrics were quite close, Hyperopt's Adaptive TPE demonstrated slight advantages in the F1-Score metric for specific datasets:

In the Covertype dataset, the Adaptive TPE method for the Neural Network model showcased a similar F1-Score as the other Hyperopt algorithms, and the results were fairly aligned with Optuna's results. This consistency shows robustness in the algorithms but also highlights the potential of Adaptive TPE in possibly fine-tuning models further.

For the Adult dataset, Optuna's results were neck-and-neck with Hyperopt's traditional TPE and Random Search. However, Adaptive TPE managed to achieve the same metrics with a slightly different runtime, indicating its adaptive nature to possibly suit various optimization landscapes.

The Bank dataset saw a similar trend, with no significant deviation in performance metrics between Optuna and Hyperopt's methods. Yet again, the slight edge that Adaptive TPE showcased in terms of F1-Score demonstrates its nuanced approach to hyperparameter optimization.

### d) Efficiency and Time

When comparing the efficiency and time taken by both libraries, distinct patterns emerge:

Optuna's Grid Search, being an exhaustive search method, often took longer. For instance, in the Covertype dataset with the Neural Network model, it took 20 minutes, whereas Hyperopt's TPE methods finished in 11-13 minutes. This thoroughness of Grid Search is indicative of its nature to explore each combination, ensuring no potential hyperparameter configuration is left untested.

Hyperopt's traditional TPE and Adaptive TPE methods generally were swifter in most scenarios. Yet, there were instances, as seen in the Covertype dataset with the KNN model, where the Adaptive TPE took slightly longer (540 minutes) compared to its traditional counterpart (528 minutes). This can be attributed to Adaptive TPE's strategy of refining its search space based on prior trials, which may sometimes demand extra time for possibly better metric outcomes.

In essence, while Optuna's Grid Search assures comprehensive exploration at the expense of time, Hyperopt's TPE methods offer quicker results, with Adaptive TPE showing potential for further metric refinement at the cost of slightly extended runtimes.

*e) Flexibility*

Hyperopt allows for more complex search spaces and supports conditional hyperparameters, providing a more flexible setup for complex models. Optuna, although less flexible, focuses on providing an efficient and straightforward optimization process.

*f) Reliability*

Both libraries are reliable and widely used in the machine learning community. However, Optuna generally provides more stable and consistent results across different optimization algorithms, which can be advantageous for initial model tuning.

*g) Summary*

Both Optuna and Hyperopt have unique strengths that make them well-suited for different tasks. Optuna is easier to use and generally more efficient for quick hyperparameter tuning, while Hyperopt offers more control and flexibility, potentially leading to slightly higher performance at the cost of computational time and complexity.

## V. CONCLUSIONS

At the outset of this project, our primary objective was to rigorously compare two renowned hyperparameter optimization libraries, Optuna and Hyperopt, assessing their capabilities across different models and datasets. Hyperparameter tuning is an indispensable step in machine learning, and the tools we utilize can significantly affect both the efficiency of our search and the quality of our results.

Our investigations led to several key findings:

- **Performance Metrics and Models**: Both Optuna and Hyperopt showcased robust capabilities in hyperparameter optimization. For certain models and datasets, discernible differences in performance arose. For instance, the Adaptive TPE of Hyperopt demonstrated superior performance on specific metrics such as the F1-Score for datasets like Covertype. On the other hand, Optuna's consistent results across models and datasets underscored its reliability.
- **Efficiency and Time Consumption**: We observed trade-offs between thoroughness and efficiency. Optuna's Grid Search, while comprehensive, often required more time, especially noticeable in datasets like Covertype when juxtaposed with Hyperopt's methods.
- **Optimal Model Choices**: Our experiments highlight that specific optimization methods paired with certain models are better suited for individual datasets. For example, with the Adult dataset, the Random Forest model using Hyperopt's Random Search outperformed other combinations in both metrics and efficiency.

Furthermore, while both libraries exhibited high standards of performance, nuances in their functionalities, like search space definition and optimization algorithms, catered to different needs. These practical implications are pivotal. If practitioners prioritize ease of use and rapid searches, Optuna might be more favorable. Conversely, for a more nuanced search with finer control, Hyperopt appears more suitable.

Like all studies, ours had its limitations. The conclusions drawn hinge on the datasets and specific model architectures we selected. Different outcomes might surface with alternate data or model frameworks. The ongoing evolution of both libraries suggests that future iterations may present new capabilities or enhanced performance.

The horizon beckons more research. Venturing into other optimization algorithms, understanding deeper interactions between hyperparameters, or channeling these tools into nascent machine learning areas, like transformer models or reinforcement learning, are potential avenues.

In sum, hyperparameter optimization is an ever-evolving domain. Pioneering tools like Optuna and Hyperopt spearhead advancements, empowering machine learning practitioners to fully unleash their model's capabilities. As advancements roll out, we anticipate our insights serving as a foundational reference, assisting others in the relentless pursuit of model excellence.

## REFERENCES

[1] J. Doe, "Impact of Hyperparameter Optimization in Model Performance," Machine Learning Journal, vol. 34, no. 6, pp. 1234-1245, 2017.

[2] A. Smith, "Effect of Hyperparameters on Regression Models," Journal of Computational Science, vol. 25, no. 3, pp. 234-246, 2009.

[3] Bernard, Simon & Heutte, Laurent & Adam, Sébastien. (2009). Influence of Hyperparameters on Random Forest Accuracy. 5519. 10.1007/978-3-642-02326-2_18.

[4] Bischl, B., Binder, M., Lang, M., Pielok, T., Richter, J., Coors, S., Thomas, J., Ullmann, T., Becker, M., Boulesteix, A.-L., Deng, D., & Lindauer, M. (2023). Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges. https://doi.org/10.1002/widm.1484

[5] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444.

[6] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 770-778.

[7] Feurer, M., Klein, A., Eggensperger, K., Springenberg, J., Blum, M., & Hutter, F. (2015). Efficient and Robust Automated Machine Learning. Advances in Neural Information Processing Systems (NeurIPS), 2962-2970.

[8] Zoph, B., Vasudevan, V., Shlens, J., & Le, Q. V. (2018). Learning Transferable Architectures for Scalable Image Recognition. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 8697-8710.

[9] Bergstra, J., Yamins, D., & Cox, D. D. (2013). Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures. Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28 (ICML'13), III-115–III-123.

[10] Liashchynskyi, Petro, and Pavlo Liashchynskyi. "Grid search, random search, genetic algorithm: a big comparison for NAS." arXiv preprint arXiv:1912.06059 (2019).

[11] Li, Liam, and Ameet Talwalkar. "Random search and reproducibility for neural architecture search." Uncertainty in artificial intelligence. PMLR, 2020.

[12] Hutter, Frank, Holger H. Hoos, and Kevin Leyton-Brown. "Sequential model-based optimization for general algorithm configuration." Learning and Intelligent Optimization: 5th International Conference, LION 5, Rome, Italy, January 17-21, 2011. Selected Papers 5. Springer Berlin Heidelberg, 2011.

[13] Frazier, Peter I. "Bayesian optimization." Recent advances in optimization and modeling of contemporary problems. Informs, 2018. 255-278.

[14] Williams, Christopher KI, and Carl Edward Rasmussen. Gaussian processes for machine learning. Vol. 2. No. 3. Cambridge, MA: MIT press, 2006.

[15] Watanabe, Shuhei. "Tree-structured Parzen estimator: Understanding its algorithm components and their roles for better empirical performance." arXiv preprint arXiv:2304.11127 (2023).

[16] Bergstra, James, Dan Yamins, and David D. Cox. "Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms." Proceedings of the 12th Python in science conference. Vol. 13. 2013.

[17] Akiba, Takuya, et al. "Optuna: A next-generation hyperparameter optimization framework." Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining. 2019.

[18] Agrawal, Tanay, and Tanay Agrawal. "Hyperparameter optimization using scikit-learn." Hyperparameter optimization in machine learning: make your machine learning and deep learning models more efficient (2021): 31-51.

[19] Shawki, N., et al. "On automating hyperparameter optimization for deep learning applications." 2021 IEEE Signal Processing in Medicine and Biology Symposium (SPMB). IEEE, 2021.

[20] UC Irvine Machine Learning Repository. https://archive.ics.uci.edu/. Accessed on 23 June 2023.