

PyTorch Hyperparameter Tuning — A Tutorial for spotPython

Version 0.2.0

Thomas Bartz-Beielstein
bartzbeielstein@gmail.com

<https://orcid.org/0000-0002-5938-5158>

May, 27th 2023

1 Hyperparameter Tuning

Hyperparameter tuning is an important, but often difficult and computationally intensive task. Changing the architecture of a neural network or the learning rate of an optimizer can have a significant impact on the performance.

The goal of hyperparameter tuning is to optimize the hyperparameters in a way that improves the performance of the machine learning or deep learning model. The simplest, but also most computationally expensive, approach uses manual search (or trial-and-error (Meignan et al. 2015)). Commonly encountered is simple random search, i.e., random and repeated selection of hyperparameters for evaluation, and lattice search (“grid search”). In addition, methods that perform directed search and other model-free algorithms, i.e., algorithms that do not explicitly rely on a model, e.g., evolution strategies (Bartz-Beielstein et al. 2014) or pattern search (Lewis, Torczon, and Trosset 2000) play an important role. Also, “hyperband”, i.e., a multi-armed bandit strategy that dynamically allocates resources to a set of random configurations and uses successive bisections to stop configurations with poor performance (Li et al. 2016), is very common in hyperparameter tuning. The most sophisticated and efficient approaches are the Bayesian optimization and surrogate model based optimization methods, which are based on the optimization of cost functions determined by simulations or experiments.

We consider below a surrogate model based optimization-based hyperparameter tuning approach based on the Python version of the SPOT (“Sequential Parameter Optimization Toolbox”) (Bartz-Beielstein, Lasarczyk, and Preuss 2005), which is suitable for situations where only limited resources are available. This may be due to limited availability and cost of hardware, or due to the fact that confidential data may only be processed locally, e.g., due to legal

requirements. Furthermore, in our approach, the understanding of algorithms is seen as a key tool for enabling transparency and explainability. This can be enabled, for example, by quantifying the contribution of machine learning and deep learning components (nodes, layers, split decisions, activation functions, etc.). Understanding the importance of hyperparameters and the interactions between multiple hyperparameters plays a major role in the interpretability and explainability of machine learning models. SPOT provides statistical tools for understanding hyperparameters and their interactions. Last but not least, it should be noted that the SPOT software code is available in the open source `spotPython` package on github¹, allowing replicability of the results. This tutorial describes the Python variant of SPOT, which is called `spotPython`. The R implementation is described in Bartz et al. (2022). SPOT is an established open source software that has been maintained for more than 15 years (Bartz-Beielstein, Lasarczyk, and Preuss 2005) (Bartz et al. 2022).

This tutorial is structured as follows. The concept of the hyperparameter tuning software `spotPython` is described in Section 2. Section 3 describes the integration of `spotPython` into the PyTorch training workflow and presents the results. Finally, Section 4 presents a summary and an outlook.

Note

The corresponding `.ipynb` notebook (Bartz-Beielstein 2023) is updated regularly and reflects updates and changes in the `spotPython` package. It can be downloaded from https://github.com/sequential-parameter-optimization/spotPython/blob/main/notebooks/14_spot_ray_hpt_torch_cifar10.ipynb.

2 The Hyperparameter Tuning Software SPOT

Surrogate model based optimization methods are common approaches in simulation and optimization. SPOT was developed because there is a great need for sound statistical analysis of simulation and optimization algorithms. SPOT includes methods for tuning based on classical regression and analysis of variance techniques. It presents tree-based models such as classification and regression trees and random forests as well as Bayesian optimization (Gaussian process models, also known as Kriging). Combinations of different meta-modeling approaches are possible. SPOT comes with a sophisticated surrogate model based optimization method, that can handle discrete and continuous inputs. Furthermore, any model implemented in `scikit-learn` can be used out-of-the-box as a surrogate in `spotPython`.

SPOT implements key techniques such as exploratory fitness landscape analysis and sensitivity analysis. It can be used to understand the performance of various algorithms, while simultaneously giving insights into their algorithmic behavior. In addition, SPOT can be used as an

¹<https://github.com/sequential-parameter-optimization>

optimizer and for automatic and interactive tuning. Details on SPOT and its use in practice are given by Bartz et al. (2022).

A typical hyperparameter tuning process with `spotPython` consists of the following steps:

1. Loading the data (training and test datasets), see Section 3.3.
2. Specification of the preprocessing model, see Section 3.4. This model is called `prep_model` (“preparation” or pre-processing). The information required for the hyperparameter tuning is stored in the dictionary `fun_control`. Thus, the information needed for the execution of the hyperparameter tuning is available in a readable form.
3. Selection of the machine learning or deep learning model to be tuned, see Section 3.5. This is called the `core_model`. Once the `core_model` is defined, then the associated hyperparameters are stored in the `fun_control` dictionary. First, the hyperparameters of the `core_model` are initialized with the default values of the `core_model`. As default values we use the default values contained in the `spotPython` package for the algorithms of the `torch` package.
4. Modification of the default values for the hyperparameters used in `core_model`, see Section 3.7.1. This step is optional.
 1. numeric parameters are modified by changing the bounds.
 2. categorical parameters are modified by changing the categories (“levels”).
5. Selection of target function (loss function) for the optimizer, see Section 3.8.
6. Calling SPOT with the corresponding parameters, see Section 3.9. The results are stored in a dictionary and are available for further analysis.
7. Presentation, visualization and interpretation of the results, see Section 3.10.

3 Hyperparameter Tuning for PyTorch With `spotPython`

In this tutorial, we will show how `spotPython` can be integrated into the PyTorch training workflow. It is based on the tutorial “Hyperparameter Tuning with Ray Tune” from the PyTorch documentation (PyTorch 2023a), which is an extension of the tutorial “Training a Classifier” (PyTorch 2023b) for training a CIFAR10 image classifier.

This document refers to the following software versions:

- python: 3.10.10
- torch: 2.0.1
- torchvision: 0.15.0
- spotPython: 0.2.0

`spotPython` can be installed via `pip`².

²Alternatively, the source code can be downloaded from gitHub: <https://github.com/sequential-parameter-optimization/spotPython>.

```
!pip install spotPython
```

Results that refer to the Ray Tune package are taken from https://PyTorch.org/tutorials/beginner/hyperparameter_tuning_tutorial.html³.

3.1 Setup

Before we consider the detailed experimental setup, we select the parameters that affect run time, initial design size and the device that is used.

```
MAX_TIME = 60
INIT_SIZE = 20
DEVICE = "cpu" # "cuda:0"
```

3.2 Initialization of the fun_control Dictionary

spotPython uses a Python dictionary for storing the information required for the hyperparameter tuning process. This dictionary is called `fun_control` and is initialized with the function `fun_control_init`. The function `fun_control_init` returns a skeleton dictionary. The dictionary is filled with the required information for the hyperparameter tuning process. It stores the hyperparameter tuning settings, e.g., the deep learning network architecture that should be tuned, the classification (or regression) problem, and the data that is used for the tuning. The dictionary is used as an input for the SPOT function.

```
fun_control = fun_control_init()
```

3.3 Data Loading

The data loading process is implemented in the same manner as described in the Section “Data loaders” in PyTorch (2023a). The data loaders are wrapped into the function `load_data`. A global data directory is used, which allows sharing the data directory between different trials.

```
def load_data(data_dir="./data"):
    transform = transforms.Compose([
        transforms.ToTensor(),
        transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))
    ])
```

³We were not able to install Ray Tune on our system. Therefore, we used the results from the PyTorch tutorial.

```

trainset = torchvision.datasets.CIFAR10(
    root=data_dir, train=True, download=True, transform=transform)

testset = torchvision.datasets.CIFAR10(
    root=data_dir, train=False, download=True, transform=transform)

return trainset, testset

```

The test and train data are added to the dictionary `fun_control`.

```

train, test = load_data()
n_samples = len(train)
# add the dataset to the fun_control
fun_control.update({
    "train": train,
    "test": test,
    "n_samples": n_samples})

```

3.4 Specification of the Preprocessing Model

After the training and test data are specified and added to the `fun_control` dictionary, `spotPython` allows the specification of a data preprocessing pipeline, e.g., for the scaling of the data or for the one-hot encoding of categorical variables. The preprocessing model is called `prep_model` (“preparation” or pre-processing) and includes steps that are not subject to the hyperparameter tuning process. The preprocessing model is specified in the `fun_control` dictionary. The preprocessing model can be implemented as a `sklearn` pipeline. The following code shows a typical preprocessing pipeline:

```

categorical_columns = ["cities", "colors"]
one_hot_encoder = OneHotEncoder(handle_unknown="ignore",
                                sparse_output=False)

prep_model = ColumnTransformer(
    transformers=[
        ("categorical", one_hot_encoder, categorical_columns),
    ],
    remainder=StandardScaler(),
)

```

Because the Ray Tune (`ray[tune]`) hyperparameter tuning as described in PyTorch (2023a) does not use a preprocessing model, the preprocessing model is set to `None` here.

```

prep_model = None
fun_control.update({"prep_model": prep_model})

```

3.5 Select algorithm and core_model_hyper_dict

The same neural network model as implemented in the section “Configurable neural network” of the PyTorch tutorial (PyTorch 2023a) is used here. We will show the implementation from PyTorch (2023a) in Section 3.5.1 first, before the extended implementation with `spotPython` is shown in Section 3.5.2.

3.5.1 Implementing a Configurable Neural Network With Ray Tune

We used the same hyperparameters that are implemented as configurable in the PyTorch tutorial. We specify the layer sizes, namely 11 and 12, of the fully connected layers:

```

class Net(nn.Module):
    def __init__(self, l1=120, l2=84):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(3, 6, 5)
        self.pool = nn.MaxPool2d(2, 2)
        self.conv2 = nn.Conv2d(6, 16, 5)
        self.fc1 = nn.Linear(16 * 5 * 5, l1)
        self.fc2 = nn.Linear(l1, l2)
        self.fc3 = nn.Linear(l2, 10)

    def forward(self, x):
        x = self.pool(F.relu(self.conv1(x)))
        x = self.pool(F.relu(self.conv2(x)))
        x = x.view(-1, 16 * 5 * 5)
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x

```

The learning rate, i.e., `lr`, of the optimizer is made configurable, too:

```

optimizer = optim.SGD(net.parameters(), lr=config["lr"], momentum=0.9)

```

3.5.2 Implementing a Configurable Neural Network With spotPython

spotPython implements a class which is similar to the class described in the PyTorch tutorial. The class is called `Net_CIFAR10` and is implemented in the file `netcifar10.py`.

```
import spotPython.torch.netcore as netcore
class Net_CIFAR10(netcore.Net_Core):
    def __init__(self, l1, l2, lr_mult, batch_size, epochs, k_folds, patience):
        super(Net_CIFAR10, self).__init__(
            lr_mult=lr_mult, batch_size=batch_size, epochs=epochs, k_folds=k_folds,
            patience=patience
        )
        self.conv1 = nn.Conv2d(3, 6, 5)
        self.pool = nn.MaxPool2d(2, 2)
        self.conv2 = nn.Conv2d(6, 16, 5)
        self.fc1 = nn.Linear(16 * 5 * 5, 11)
        self.fc2 = nn.Linear(11, 12)
        self.fc3 = nn.Linear(12, 10)

    def forward(self, x):
        x = self.pool(F.relu(self.conv1(x)))
        x = self.pool(F.relu(self.conv2(x)))
        x = x.view(-1, 16 * 5 * 5)
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x
```

`Net_CIFAR10` inherits from the class `Net_Core` which is implemented in the file `netcore.py`. It implements the additional attributes that are common to all neural network models. The attributes are the learning rate multiplier `lr_mult`, the batch size `batch_size`, the number of epochs `epochs`, the number of folds `k_folds` for the cross validation, and the patience `patience` for the early stopping. The class `Net_Core` is shown below.

```
from torch import nn

class Net_Core(nn.Module):
    def __init__(self, lr_mult, batch_size, epochs, k_folds, patience):
        super(Net_Core, self).__init__()
        self.lr_mult = lr_mult
        self.batch_size = batch_size
```

```
self.epochs = epochs
self.k_folds = k_folds
self.patience = patience
```

3.5.3 Comparison of the Approach Described in the PyTorch Tutorial With spotPython

Comparing the class `Net` from the PyTorch tutorial and the class `Net_CIFAR10` from `spotPython`, we see that the class `Net_CIFAR10` has additional attributes and does not inherit from `nn` directly. It adds an additional class, `Net_core`, that takes care of additional attributes that are common to all neural network models, e.g., the learning rate multiplier `lr_mult` or the batch size `batch_size`.

`spotPython`'s `core_model` implements an instance of the `Net_CIFAR10` class. In addition to the basic neural network model, the `core_model` can use these additional attributes. `spotPython` provides methods for handling these additional attributes to guarantee 100% compatibility with the PyTorch classes. The method `add_core_model_to_fun_control` adds the hyper-parameters and additional attributes to the `fun_control` dictionary. The method is shown below.

```
core_model = Net_CIFAR10
fun_control = add_core_model_to_fun_control(core_model=core_model,
                                           fun_control=fun_control,
                                           hyper_dict=TorchHyperDict,
                                           filename=None)
```

Note

In addition to the class `Net` from the PyTorch tutorial, the class `Net_CIFAR10` has additional attributes, e.g.:

- learning rate (`lr`),
- batch size (`batch_size`),
- epochs (`epochs`),
- k_folds (`k_folds`), and
- early stopping criterion “patience” (`patience`)

Further attributes can be easily added to the class, e.g., `optimizer`.

3.6 The Search Space

In Section 3.6.1, we first describe how to configure the search space with `ray[tune]` (as shown in PyTorch (2023a)) and then how to configure the search space with `spotPython` in Section 3.6.2.

3.6.1 Configuring the Search Space With Ray Tune

Ray Tune’s search space can be configured as follows (PyTorch 2023a):

```
config = {
    "l1": tune.sample_from(lambda _: 2**np.random.randint(2, 9)),
    "l2": tune.sample_from(lambda _: 2**np.random.randint(2, 9)),
    "lr": tune.loguniform(1e-4, 1e-1),
    "batch_size": tune.choice([2, 4, 8, 16])
}
```

The `tune.sample_from()` function enables the user to define sample methods to obtain hyperparameters. In this example, the `l1` and `l2` parameters should be powers of 2 between 4 and 256, so either 4, 8, 16, 32, 64, 128, or 256. The `lr` (learning rate) should be uniformly sampled between 0.0001 and 0.1. Lastly, the batch size is a choice between 2, 4, 8, and 16.

At each trial, `ray[tune]` will randomly sample a combination of parameters from these search spaces. It will then train a number of models in parallel and find the best performing one among these. `ray[tune]` uses the `ASHAScheduler` which will terminate bad performing trials early.

3.6.2 Configuring the Search Space With spotPython

3.6.2.1 The `hyper_dict` Hyperparameters for the Selected Algorithm

`spotPython` uses simple JSON files for the specification of the hyperparameters. Users can specify their individual JSON files, or they can use the JSON files provided by `spotPython`. The JSON file for the `core_model` is called `torch_hyper_dict.json`.

In contrast to `ray[tune]`, `spotPython` can handle numerical, boolean, and categorical hyperparameters. They can be specified in the JSON file in a similar way as the numerical hyperparameters as shown below:

```

"factor_hyperparameter": {
  "levels": ["A", "B", "C"],
  "type": "factor",
  "default": "B",
  "transform": "None",
  "core_model_parameter_type": "str",
  "lower": 0,
  "upper": 2},

```

Each entry in the JSON file represents one hyperparameter with the following structure: `type`, `default`, `transform`, `lower`, and `upper`. The corresponding entries for the `Net_CIFAR10` class are shown below.

```

{"Net_CIFAR10":
  {
    "l1": {
      "type": "int",
      "default": 5,
      "transform": "transform_power_2_int",
      "lower": 2,
      "upper": 9},
    "l2": {
      "type": "int",
      "default": 5,
      "transform": "transform_power_2_int",
      "lower": 2,
      "upper": 9},
    "lr_mult": {
      "type": "float",
      "default": 1.0,
      "transform": "None",
      "lower": 0.1,
      "upper": 10},
    "batch_size": {
      "type": "int",
      "default": 4,
      "transform": "transform_power_2_int",
      "lower": 1,
      "upper": 4},
    "epochs": {
      "type": "int",
      "default": 3,

```

```

        "transform": "transform_power_2_int",
        "lower": 1,
        "upper": 4},
    "k_folds": {
        "type": "int",
        "default": 2,
        "transform": "None",
        "lower": 2,
        "upper": 3},
    "patience": {
        "type": "int",
        "default": 5,
        "transform": "None",
        "lower": 2,
        "upper": 10},
    "optimizer": {
        "levels": ["Adadelata",
                    "Adagrad",
                    "Adam",
                    "AdamW",
                    "SparseAdam",
                    "Adamax",
                    "ASGD",
                    "LBFGS",
                    "NAdam",
                    "RAdam",
                    "RMSprop",
                    "Rprop",
                    "SGD"],
        "type": "factor",
        "default": "SGD",
        "transform": "None",
        "class_name": "torch.optim",
        "core_model_parameter_type": "str",
        "lower": 0,
        "upper": 12},
    "sgd_momentum": {
        "type": "float",
        "default": 0.0,
        "transform": "None",
        "lower": 0.0,
        "upper": 1.0}

```

```
}  
}
```

3.7 Modifying the Hyperparameters

Ray tune (PyTorch 2023a) does not provide a way to change the specified hyperparameters without re-compilation. However, `spotPython` provides functions for modifying the hyperparameters, their bounds and factors as well as for activating and de-activating hyperparameters without re-compilation of the Python source code. These functions are described in the following.

3.7.1 Modify `hyper_dict` Hyperparameters for the Selected Algorithm aka `core_model`

After specifying the model, the corresponding hyperparameters, their types and bounds are loaded from the JSON file `torch_hyper_dict.json`. After loading, the user can modify the hyperparameters, e.g., the bounds. `spotPython` provides a simple rule for de-activating hyperparameters: If the lower and the upper bound are set to identical values, the hyperparameter is de-activated. This is useful for the hyperparameter tuning, because it allows to specify a hyperparameter in the JSON file, but to de-activate it in the `fun_control` dictionary. This is done in the next step.

3.7.2 Modify Hyperparameters of Type numeric and integer (boolean)

Since the hyperparameter `k_folds` is not used in the PyTorch tutorial, it is de-activated here by setting the lower and upper bound to the same value. Note, `k_folds` is of type “integer”.

```
fun_control = modify_hyper_parameter_bounds(fun_control, "batch_size", bounds=[1, 5])  
fun_control = modify_hyper_parameter_bounds(fun_control, "k_folds", bounds=[0, 0])  
fun_control = modify_hyper_parameter_bounds(fun_control, "patience", bounds=[3, 3])  
fun_control["core_model_hyper_dict"]
```

3.7.3 Modify Hyperparameter of Type factor

In a similar manner as for the numerical hyperparameters, the categorical hyperparameters can be modified. New configurations can be chosen by adding or deleting levels. For example, the hyperparameter `optimizer` can be re-configured as follows:

In the following setting, two optimizers ("SGD" and "Adam") will be compared during the `spotPython` hyperparameter tuning. The hyperparameter `optimizer` is active.

```
fun_control = modify_hyper_parameter_levels(fun_control, "optimizer", ["SGD", "Adam"])
```

The hyperparameter `optimizer` can be de-activated by choosing only one value (level), here: "SGD".

```
fun_control = modify_hyper_parameter_levels(fun_control, "optimizer", ["SGD"])
```

As discussed in Section 3.7.4, there are some issues with the LBFGS optimizer. Therefore, the usage of the LBFGS optimizer is not deactivated in `spotPython` by default. However, the LBFGS optimizer can be activated by adding it to the list of optimizers. `Rprop` was removed, because it does perform very poorly (as some pre-tests have shown). However, it can also be activated by adding it to the list of optimizers. Since `SparseAdam` does not support dense gradients, `Adam` was used instead. Therefore, there are 10 default optimizers:

```
fun_control = modify_hyper_parameter_levels(fun_control, "optimizer", ["Adadelata",
    "Adagrad", "Adam", "AdamW", "Adamax", "ASGD", "NAdam", "RAdam",
    "RMSprop", "SGD"])
```

3.7.4 Optimizers

Table 1 shows some of the optimizers available in PyTorch:

Table 1: Optimizers available in PyTorch (selection). “mom” denotes **momentum**, “weight” **weight_decay**, “damp” **dampening**, “nest” **nesterov**, “lr_sc” **learning rate for scaling delta**, “mom_dec” for **momentum_decay**, and “step_s” for **step_sizes**. The default values are shown in the table.

Optimizer	lr	mom	weight	damp	nest	rho	lr_sc	lr_decay	lambd	alpha	mom_dec	decay	step_s
Adadelata	-	0.	-	-	0.9	1.0	-	-	-	-	-	-	-
Adagrad	1e-2	-	0.	-	-	-	-	0.	-	-	-	-	-
Adam	1e-3	-	0.	-	-	-	-	-	(0.9,0.999)	-	-	-	-
AdamW	1e-3	-	1e-2	-	-	-	-	-	(0.9,0.999)	-	-	-	-
SparseAdam	1e-3	-	-	-	-	-	-	-	(0.9,0.999)	-	-	-	-
Adamax	2e-3	-	0.	-	-	-	-	-	(0.9, 0.999)	-	-	-	-

Optimizer	lr	momentum	weight_decay	dampening	nest	rho	lr_scheduler	lr_decay	betas	lambda	alpha	momentum_decay	step_size
ASGD	1e-2	0.9	0.	-	False	-	-	-	-	1e-4	0.75	-	-
LBFGS	1.	-	-	-	-	-	-	-	-	-	-	-	-
NAdam	2e-3	-	0.	-	-	-	-	-	(0.9,0.999)	-	0	-	-
RAdam	1e-3	-	0.	-	-	-	-	-	(0.9,0.999)	-	-	-	-
RMSprop	1e-2	0.	0.	-	-	-	-	-	(0.9,0.999)	-	-	-	-
Rprop	1e-2	-	-	-	-	-	-	-	-	(0.5,1.2)	1e-6, 50)	-	-
SGD	required	0.	0.	0.	False	-	-	-	-	-	-	-	-

`spotPython` implements an `optimization` handler that maps the optimizer names to the corresponding `PyTorch` optimizers.

i A note on LBFGS

We recommend deactivating `PyTorch`'s LBFGS optimizer, because it does not perform very well. The `PyTorch` documentation, see <https://pytorch.org/docs/stable/generated/torch.optim.LBFGS.html#torch.optim.LBFGS>, states:

This is a very memory intensive optimizer (it requires additional `param_bytes * (history_size + 1)` bytes). If it doesn't fit in memory try reducing the history size, or use a different algorithm.

Furthermore, the LBFGS optimizer is not compatible with the `PyTorch` tutorial. The reason is that the LBFGS optimizer requires the `closure` function, which is not implemented in the `PyTorch` tutorial. Therefore, the LBFGS optimizer is recommended here.

Since there are 10 optimizers in the portfolio, it is not recommended tuning the hyperparameters that effect one single optimizer only.

i A note on the learning rate

`spotPython` provides a multiplier for the default learning rates, because optimizers use different learning rates. Using a multiplier for the learning rates might enable a simultaneous tuning of the learning rates for all optimizers. However, this is not recommended, because the learning rates are not comparable across optimizers. Therefore, we recom-

mend fixing the learning rate for all optimizers if multiple optimizers are used. This can be done by setting the lower and upper bounds of the learning rate multiplier to the same value as shown below.

Thus, the learning rate, which affects the SGD optimizer, will be set to a fixed value. We choose the default value of $1e-3$ for the learning rate, because it is used in other PyTorch examples (it is also the default value used by `spotPython` as defined in the `optimizer_handler()` method). We recommend tuning the learning rate later, when a reduced set of optimizers is fixed. Here, we will demonstrate how to select in a screening phase the optimizers that should be used for the hyperparameter tuning.

For the same reason, we will fix the `sgd_momentum` to 0.9.

```
fun_control = modify_hyper_parameter_bounds(fun_control, "lr_mult", bounds=[1.0, 1.0])
fun_control = modify_hyper_parameter_bounds(fun_control, "sgd_momentum", bounds=[0.9, 0.9])
```

3.8 Evaluation

The evaluation procedure requires the specification of two elements:

1. the way how the data is split into a train and a test set and
2. the loss function (and a metric).

3.8.1 Hold-out Data Split and Cross-Validation

As a default, `spotPython` provides a standard hold-out data split and cross validation.

3.8.1.1 Hold-out Data Split

If a hold-out data split is used, the data will be partitioned into a training, a validation, and a test data set. The split depends on the setting of the `eval` parameter. If `eval` is set to `train_hold_out`, one data set, usually the original training data set, is split into a new training and a validation data set. The training data set is used for training the model. The validation data set is used for the evaluation of the hyperparameter configuration and early stopping to prevent overfitting. In this case, the original test data set is not used. The following splits are performed in the hold-out setting: $\{\text{train}_0, \text{test}\} \rightarrow \{\text{train}_1, \text{validation}_1, \text{test}\}$, where $\text{train}_1 \cup \text{validation}_1 = \text{train}_0$.

Note

`spotPython` returns the hyperparameters of the machine learning and deep learning models, e.g., number of layers, learning rate, or optimizer, but not the model weights. Therefore, after the SPOT run is finished, the corresponding model has to be trained again with the best hyperparameter configuration. The training is performed on the training data set. The test data set is used for the final evaluation of the model.

Summarizing, the following splits are performed in the hold-out setting:

1. Run `spotPython` with `eval` set to `train_hold_out` to determine the best hyperparameter configuration.
2. Train the model with the best hyperparameter configuration on the training data set:
 - `train_tuned(model_spot, train, "model_spot.pt")`.
3. Test the model on the test data:
 - `test_tuned(model_spot, test, "model_spot.pt")`

These steps will be exemplified in the following sections.

In addition to this **hold-out** setting, `spotPython` provides another hold-out setting, where an explicit test data is specified by the user that will be used as the validation set. To choose this option, the `eval` parameter is set to `test_hold_out`. In this case, the training data set is used for the model training. Then, the explicitly defined test data set is used for the evaluation of the hyperparameter configuration (the validation).

3.8.1.2 Cross-Validation

The cross validation setting is used by setting the `eval` parameter to `train_cv` or `test_cv`. In both cases, the data set is split into k folds. The model is trained on $k - 1$ folds and evaluated on the remaining fold. This is repeated k times, so that each fold is used exactly once for evaluation. The final evaluation is performed on the test data set. The cross validation setting is useful for small data sets, because it allows to use all data for training and evaluation. However, it is computationally expensive, because the model has to be trained k times.

Note

Combinations of the above settings are possible, e.g., cross validation can be used for training and hold-out for evaluation or *vice versa*. Also, cross validation can be used for training and testing. Because cross validation is not used in the PyTorch tutorial (PyTorch 2023a), it is not considered further here.

3.8.1.3 Overview of the Evaluation Settings

3.8.1.3.1 Settings for the Hyperparameter Tuning

Table 2 provides an overview of the training evaluations.

Table 2: Overview of the evaluation settings.

eval	train	test	function	comment
"train_hold_out"	✓		train_hold_out_split() validate_fold_train_hold_out() for early stopping	It's the train_hold_out() set internally
"test_hold_out"	✓	✓	train_hold_out_split() validate_fold_test_hold_out() for early stopping	the test validate_fold_test_hold_out() validate_fold_or_hold_out()
"train_cv"	✓		evaluate_cv(model, train)	CV, using the train data set
"test_cv"		✓	evaluate_cv(model, test)	CV, using the test data set . Identical to "train_cv", uses only test data.

- "train_cv" and "test_cv" use sklearn.model_selection.KFold() internally.

3.8.1.4 Detailed Description of the "train_hold_out" Setting

The "train_hold_out" setting is used by default. It uses the loss function specified in fun_control and the metric specified in fun_control.

1. First, the method HyperTorch().fun_torch is called.
2. fun_torch() calls spotPython.torch.traintest.evaluate_hold_out() as follows:

```
df_eval, _ = evaluate_hold_out(  
    model,  
    train_dataset=fun_control["train"],  
    shuffle=self.fun_control["shuffle"],
```

```

    loss_function=self.fun_control["loss_function"],
    metric=self.fun_control["metric_torch"],
    device=self.fun_control["device"],
    show_batch_interval=self.fun_control["show_batch_interval"],
    path=self.fun_control["path"],
)

```

Note: Only the data set `fun_control["train"]` is used for training and validation. It is used as follows:

```

trainloader, valloader = create_train_val_data_loaders(
    dataset=train_dataset, batch_size=batch_size_instance, shuffle=shuffle
)

```

`create_train_val_data_loaders()` splits the `train_dataset` into `trainloader` and `valloader` using `torch.utils.data.random_split()` as follows:

```

def create_train_val_data_loaders(dataset, batch_size, shuffle, num_workers=0):
    test_abs = int(len(dataset) * 0.6)
    train_subset, val_subset = random_split(dataset, [test_abs, len(dataset) - test_abs])
    trainloader = torch.utils.data.DataLoader(
        train_subset, batch_size=int(batch_size), shuffle=shuffle, num_workers=num_workers
    )
    valloader = torch.utils.data.DataLoader(
        val_subset, batch_size=int(batch_size), shuffle=shuffle, num_workers=num_workers
    )
    return trainloader, valloader

```

The optimizer is set up as follows:

```

lr_mult_instance = net.lr_mult
optimizer = optimizer_handler(optimizer_name=optimizer_instance, params=net.parameters(), lr_mult=lr_mult_instance)

```

3. `evaluate_hold_out()` sets the `net` attributes such as `epochs`, `batch_size`, `optimizer`, and `patience`. For each epoch, the methods `train_hold_out()` and `validate_fold_or_hold_out()` are called, the former for training and the latter for validation and early stopping. The validation loss from the last epoch (not the best validation loss) is returned from `evaluate_hold_out`.
4. The method `train_hold_out()` is implemented as follows:

```

def train_hold_out(net, trainloader, batch_size, loss_function, optimizer, device, show_batch_interval):
    running_loss = 0.0
    epoch_steps = 0
    for i, data in enumerate(trainloader, 0):
        inputs, labels = data
        inputs, labels = inputs.to(device), labels.to(device)
        optimizer.zero_grad()
        outputs = net(inputs)
        loss = loss_function(outputs, labels)
        loss.backward()
        torch.nn.utils.clip_grad_norm_(net.parameters(), max_norm=1.0)
        optimizer.step()
        running_loss += loss.item()
        epoch_steps += 1
        if i % show_batch_interval == (show_batch_interval - 1): # print every show_batch_interval
            print(
                "Batch: %5d. Batch Size: %d. Training Loss (running): %.3f"
                % (i + 1, int(batch_size), running_loss / epoch_steps)
            )
            running_loss = 0.0
    return loss.item()

```

5. The method `validate_fold_or_hold_out()` is implemented as follows:

```

def validate_fold_or_hold_out(net, valloader, loss_function, metric, device):
    val_loss = 0.0
    val_steps = 0
    metric.reset()
    for i, data in enumerate(valloader, 0):
        with torch.no_grad():
            inputs, labels = data
            inputs, labels = inputs.to(device), labels.to(device)
            outputs = net(inputs)
            _, predicted = torch.max(outputs.data, 1)
            metric_value = metric(predicted, labels).to(device)
            loss = loss_function(outputs, labels)
            val_loss += loss.cpu().numpy()
            val_steps += 1
    loss = val_loss / val_steps
    metric_value = metric.compute()
    return metric_value, loss

```

3.8.1.5 Detailed Description of the "test_hold_out" Setting

It uses the loss function specified in `fun_control` and the metric specified in `fun_control`.

1. First, the method `HyperTorch().fun_torch` is called.
2. `fun_torch()` calls `spotPython.torch.traintest.evaluate_hold_out()` similar to the "train_hold_out" setting with one exception: It passes an additional `test` data set to `evaluate_hold_out()` as follows:

```
test_dataset=fun_control["test"]
```

`evaluate_hold_out()` calls `create_train_test_data_loaders` instead of `create_train_val_data_loaders` as follows: The two data sets are used in `create_train_test_data_loaders` as follows:

```
def create_train_test_data_loaders(dataset, batch_size, shuffle, test_dataset, num_workers=0):
    trainloader = torch.utils.data.DataLoader(
        dataset, batch_size=int(batch_size), shuffle=shuffle, num_workers=num_workers
    )
    testloader = torch.utils.data.DataLoader(
        test_dataset, batch_size=int(batch_size), shuffle=shuffle, num_workers=num_workers
    )
    return trainloader, testloader
```

3. The following steps are identical to the "train_hold_out" setting. Only a different data loader is used for testing.

3.8.1.6 Detailed Description of the "train_cv" Setting

It uses the loss function specified in `fun_control` and the metric specified in `fun_control`.

1. First, the method `HyperTorch().fun_torch` is called.
2. `fun_torch()` calls `spotPython.torch.traintest.evaluate_cv()` as follows:

```
df_eval, _ = evaluate_cv(
    model,
    dataset=fun_control["train"],
    shuffle=self.fun_control["shuffle"],
    device=self.fun_control["device"],
    show_batch_interval=self.fun_control["show_batch_interval"],
)
```

Note: Only the data set `fun_control["train"]` is used for CV. 3. In `evaluate_cv()`, the following steps are performed: The optimizer is set up as follows:

```
lr_instance = net.lr
optimizer = optimizer_handler(optimizer_name=optimizer_instance, params=net.parameters(), lr=
```

`evaluate_cv()` sets the `net` attributes such as `epochs`, `batch_size`, `optimizer`, and `patience`. CV is implemented as follows:

```
kfold = KFold(n_splits=k_folds_instance, shuffle=shuffle)
for fold, (train_ids, val_ids) in enumerate(kfold.split(dataset)):
    train_subsampler = torch.utils.data.SubsetRandomSampler(train_ids)
    val_subsampler = torch.utils.data.SubsetRandomSampler(val_ids)
    trainloader = torch.utils.data.DataLoader(
        dataset, batch_size=batch_size_instance, sampler=train_subsampler, num_workers=num_w
    )
    valloader = torch.utils.data.DataLoader(
        dataset, batch_size=batch_size_instance, sampler=val_subsampler, num_workers=num_wor
    )
    reset_weights(net)
    # Train fold for several epochs:
    train_fold(
        net,
        trainloader,
        epochs_instance,
        loss_function,
        optimizer,
        device,
        show_batch_interval=show_batch_interval,
    )
    # Validate fold: use only loss for tuning
    metric_values[fold], loss_values[fold] = validate_fold_or_hold_out(net, valloader, loss_
df_eval = sum(loss_values.values()) / len(loss_values.values())
```

4. The method `train_fold()` is implemented as follows:

```
def train_fold(net, trainloader, epochs, loss_function, optimizer, device, show_batch_interv
    for epoch in range(epochs):
        print(f"Epoch: {epoch + 1}")
        running_loss = 0.0
        epoch_steps = 0
        for i, data in enumerate(trainloader, 0):
            inputs, labels = data
            inputs, labels = inputs.to(device), labels.to(device)
            optimizer.zero_grad()
```

```

        outputs = net(inputs)
        loss = loss_function(outputs, labels)
        loss.backward()
        torch.nn.utils.clip_grad_norm_(net.parameters(), max_norm=1.0)
        optimizer.step()
        # the following is for printing the statistic only
        running_loss += loss.item()
        epoch_steps += 1
        if i % show_batch_interval == (show_batch_interval - 1): # print every show_batch_interval
            print("Batch: %5d. Training Loss (running): %.3f" % (i + 1, running_loss / epoch_steps))
            running_loss = 0.0

```

5. The method `validate_fold_or_hold_out()` is implemented as shown above. In contrast to the hold-out setting, it is called for each of the k folds. The results are stored in a dictionaries `metric_values` and `loss_values` as follows:

```

# Validate fold: use only loss for tuning
metric_values[fold], loss_values[fold] = validate_fold_or_hold_out(net, valloader, loss_function)
df_eval = sum(loss_values.values()) / len(loss_values.values())

```

The results are averaged over the k folds and returned as `df_eval`.

3.8.1.7 Detailed Description of the "test_cv" Setting

It uses the loss function specified in `fun_control` and the metric specified in `fun_control`.

1. First, the method `HyperTorch().fun_torch` is called.
2. `fun_torch()` calls `spotPython.torch.traintest.evaluate_cv()` as follows:

```

df_eval, _ = evaluate_cv(
    model,
    dataset=fun_control["test"],
    shuffle=self.fun_control["shuffle"],
    device=self.fun_control["device"],
    show_batch_interval=self.fun_control["show_batch_interval"],
)

```

Note: The data set `fun_control["test"]` is used for CV. The rest is the same as for the "train_cv" setting.

3.8.1.8 Settings for the Final Evaluation of the Tuned Architecture

3.8.1.8.1 Training of the Tuned Architecture

`train_tuned(model, train)`: train the model with the best hyperparameter configuration (or simply the default) on the training data set. It splits the `traindata` into new `train` and `validation` sets using `create_train_val_data_loaders()`, which calls `torch.utils.data.random_split()` internally. Currently, 60% of the data is used for training and 40% for validation. The `train` data is used for training the model with `train_hold_out()`. The `validation` data is used for early stopping using `validate_fold_or_hold_out()` on the `validation` data set.

`train_tuned()` is just a wrapper to `evaluate_hold_out` using the `train` data set. It is implemented as follows:

```
def train_tuned(net, train_dataset, shuffle, loss_function, metric, device=None, show_batch_interval=10, path=None):
    evaluate_hold_out(
        net=net,
        train_dataset=train_dataset,
        shuffle=shuffle,
        test_dataset=None,
        loss_function=loss_function,
        metric=metric,
        device=device,
        show_batch_interval=show_batch_interval,
        path=path,
    )
```

Note: During training, `shuffle` is set to `True`, whereas during testing, `shuffle` is set to `False`.

3.8.1.8.2 Testing of the Tuned Architecture

`test_tuned(model, test)`: test the model on the test data set. No data splitting is performed. The (trained) model is evaluated using the `validate_fold_or_hold_out()` function.

Note: During training, `shuffle` is set to `True`, whereas during testing, `shuffle` is set to `False`.

```
def test_tuned(net, shuffle, test_dataset=None, loss_function=None, metric=None, device=None, show_batch_interval=10, path=None):
    batch_size_instance = net.batch_size
    removed_attributes, net = get_removed_attributes_and_base_net(net)
    if path is not None:
        net.load_state_dict(torch.load(path))
    net.eval()
```

```

try:
    device = getDevice(device=device)
    if torch.cuda.is_available():
        device = "cuda:0"
        if torch.cuda.device_count() > 1:
            print("We will use", torch.cuda.device_count(), "GPUs!")
            net = nn.DataParallel(net)
    net.to(device)
    valloader = torch.utils.data.DataLoader(
        test_dataset, batch_size=int(batch_size_instance), shuffle=shuffle, num_workers=dev
    )
    metric_value, loss = validate_fold_or_hold_out(
        net, valloader=valloader, loss_function=loss_function, metric=metric, device=device
    )
    df_eval = loss
    df_metric = metric_value
    df_preds = np.nan
except Exception as err:
    print(f"Error in Net_Core. Call to test_tuned() failed. {err=}, {type(err)=}")
    df_eval = np.nan
    df_metric = np.nan
    df_preds = np.nan
add_attributes(net, removed_attributes)
print(f"Final evaluation: Validation loss: {df_eval}")
print(f"Final evaluation: Validation metric: {df_metric}")
print("-----")
return df_eval, df_preds, df_metric

```

3.8.2 Loss Functions and Metrics

The key "loss_function" specifies the loss function which is used during the optimization. There are several different loss functions under PyTorch's `nn` package. For example, a simple loss is `MSELoss`, which computes the mean-squared error between the output and the target. In this tutorial we will use `CrossEntropyLoss`, because it is also used in the PyTorch tutorial.

```

from torch.nn import CrossEntropyLoss
loss_function = CrossEntropyLoss()
fun_control.update({"loss_function": loss_function})

```

In addition to the loss functions, `spotPython` provides access to a large number of metrics. The key "metric_sklearn" is used for metrics that follow the `scikit-learn` conventions.

The key "river_metric" is used for the river based evaluation (Montiel et al. 2021) via `eval_oml_iter_progressive`, and the key "metric_torch" is used for the metrics from TorchMetrics. TorchMetrics is a collection of more than 90 PyTorch metrics⁴. Because the PyTorch tutorial uses the accuracy as metric, we use the same metric here. Currently, accuracy is computed in the tutorial's example code. We will use TorchMetrics instead, because it offers:

- A standardized interface to increase reproducibility
- Reduces Boilerplate
- Distributed-training compatible
- Rigorously tested
- Automatic accumulation over batches
- Automatic synchronization between multiple devices

Therefore, we set

```
metric_torch = torchmetrics.Accuracy(task="multiclass", num_classes=10)
```

```
loss_function = CrossEntropyLoss()
weights = 1.0
metric_torch = torchmetrics.Accuracy(task="multiclass", num_classes=10)
shuffle = True
eval = "train_hold_out"
device = DEVICE
show_batch_interval = 100_000
path="torch_model.pt"

fun_control.update({
    "data_dir": None,
    "checkpoint_dir": None,
    "horizon": None,
    "oml_grace_period": None,
    "weights": weights,
    "step": None,
    "log_level": 50,
    "weight_coeff": None,
    "metric_torch": metric_torch,
    "metric_river": None,
    "metric_sklearn": None,
    "loss_function": loss_function,
    "shuffle": shuffle,
```

⁴<https://torchmetrics.readthedocs.io/en/latest/>.

```

    "eval": eval,
    "device": device,
    "show_batch_interval": show_batch_interval,
    "path": path,
    })

```

3.9 Calling the SPOT Function

Now, the dictionary `fun_control` contains all information needed for the hyperparameter tuning. Before the hyperparameter tuning is started, it is recommended to take a look at the experimental design. The method `gen_design_table` generates a design table as follows:

```
print(gen_design_table(fun_control))
```

This allows to check if all information is available and if the information is correct. Table 3 shows the experimental design for the hyperparameter tuning. Hyperparameter transformations are shown in the column “transform”, e.g., the `l1` default is 5, which results in the value $2^5 = 32$ for the network, because the transformation `transform_power_2_int` was selected in the JSON file. The default value of the `batch_size` is set to 4, which results in a batch size of $2^4 = 16$.

Table 3: Experimental design for the hyperparameter tuning. The table shows the hyperparameters, their types, default values, lower and upper bounds, and the transformation function. The transformation function is used to transform the hyperparameter values from the unit hypercube to the original domain. The transformation function is applied to the hyperparameter values before the evaluation of the objective function.

name	type	default	lower	upper	transform
<code>l1</code>	int	5	2	9	<code>transform_power_2_int</code>
<code>l2</code>	int	5	2	9	<code>transform_power_2_int</code>
<code>lr</code>	float	0.001	0.001	0.001	None
<code>batch_size</code>	int	4	1	5	<code>transform_power_2_int</code>
<code>epochs</code>	int	3	3	4	<code>transform_power_2_int</code>
<code>k_folds</code>	int	2	0	0	None
<code>patience</code>	int	5	3	3	None
<code>optimizer</code>	factor	SGD	0	9	None

The objective function `fun_torch` is selected next. It implements an interface from PyTorch’s training, validation, and testing methods to `spotPython`.

```
fun = HyperTorch().fun_torch
```

The `spotPython` hyperparameter tuning is started by calling the `Spot` function. Here, we will run the tuner for approximately 30 minutes (`max_time`). Note: the initial design is always evaluated in the `spotPython` run. As a consequence, the run may take longer than specified by `max_time`, because the evaluation time of initial design (here: `init_size`, 10 points) is performed independently of `max_time`.

```
spot_tuner = spot.Spot(fun=fun,
                        lower = lower,
                        upper = upper,
                        fun_evals = inf,
                        fun_repeats = 1,
                        max_time = MAX_TIME,
                        noise = False,
                        tolerance_x = np.sqrt(np.spacing(1)),
                        var_type = var_type,
                        var_name = var_name,
                        infill_criterion = "y",
                        n_points = 1,
                        seed=123,
                        log_level = 50,
                        show_models= False,
                        show_progress= True,
                        fun_control = fun_control,
                        design_control={"init_size": INIT_SIZE,
                                       "repeats": 1},
                        surrogate_control={"noise": True,
                                          "cod_type": "norm",
                                          "min_theta": -4,
                                          "max_theta": 3,
                                          "n_theta": len(var_name),
                                          "model_optimizer": differential_evolution,
                                          "model_fun_evals": 10_000,
                                          "log_level": 50
                                          })

spot_tuner.run(X_start=X_start)
```

During the run, the following output is shown:

```
config: {'l1': 4, 'l2': 64, 'lr_mult': 1.0, 'batch_size': 16, 'epochs': 16, 'k_folds': 0, 'p
Epoch: 1
Loss on hold-out set: 1.602842689704895
Accuracy on hold-out set: 0.4006
Metric value on hold-out data: 0.40059998631477356
Epoch: 2
Loss on hold-out set: 1.4648857820034027
Accuracy on hold-out set: 0.47685
Metric value on hold-out data: 0.4768500030040741
Epoch: 3
Loss on hold-out set: 1.403354868555069
Accuracy on hold-out set: 0.482
Metric value on hold-out data: 0.4819999933242798
Epoch: 4
Loss on hold-out set: 1.384560032081604
Accuracy on hold-out set: 0.49065
Metric value on hold-out data: 0.4906499981880188
Epoch: 5
Loss on hold-out set: 1.4326466094970702
Accuracy on hold-out set: 0.4809
Metric value on hold-out data: 0.48089998960494995
Epoch: 6
Loss on hold-out set: 1.3759961807250976
Accuracy on hold-out set: 0.4995
Metric value on hold-out data: 0.49950000643730164
Epoch: 7
Loss on hold-out set: 1.3684927892208099
Accuracy on hold-out set: 0.50695
Metric value on hold-out data: 0.5069500207901001
Epoch: 8
Loss on hold-out set: 1.3642385012149811
Accuracy on hold-out set: 0.506
Metric value on hold-out data: 0.5059999823570251
Epoch: 9
Loss on hold-out set: 1.3157437609672546
Accuracy on hold-out set: 0.5304
Metric value on hold-out data: 0.5303999781608582
Epoch: 10
Loss on hold-out set: 1.3481314319610596
Accuracy on hold-out set: 0.5268
Metric value on hold-out data: 0.5267999768257141
Epoch: 11
```

```

Loss on hold-out set: 1.3608774542331696
Accuracy on hold-out set: 0.51525
Metric value on hold-out data: 0.515250027179718
Epoch: 12
Loss on hold-out set: 1.359324642753601
Accuracy on hold-out set: 0.52355
Metric value on hold-out data: 0.5235499739646912
Early stopping at epoch 11
Returned to Spot: Validation loss: 1.359324642753601
-----

```

3.10 Results

After the hyperparameter tuning run is finished, the progress of the hyperparameter tuning can be visualized. The following code generates the progress plot from Figure 2.

```
spot_tuner.plot_progress(log_y=False, filename="./figures" + experiment_name+"_progress.png")
```

Figure 1: ?(caption)

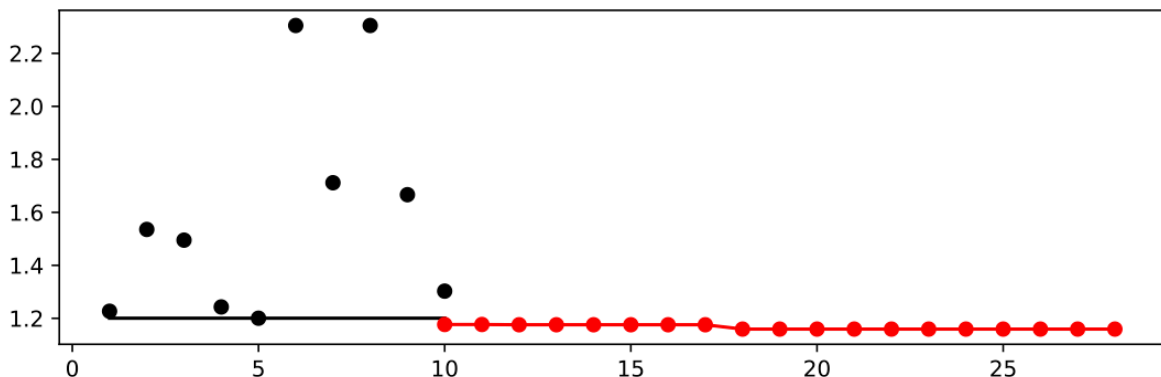


Figure 2: Progress plot. Black dots denote results from the initial design. Red dots illustrate the improvement found by the surrogate model based optimization (surrogate model based optimization).

Figure 2 shows a typical behaviour that can be observed in many hyperparameter studies (Bartz et al. 2022): the largest improvement is obtained during the evaluation of the initial design. The surrogate model based optimization refines the results. Figure 2 also illustrates one major difference between `ray[tune]` as used in PyTorch

(2023a) and `spotPython`: the `ray[tune]` uses a random search and will generate results similar to the *black* dots, whereas `spotPython` uses a surrogate model based optimization and presents results represented by *red* dots in Figure 2. The surrogate model based optimization is considered to be more efficient than a random search, because the surrogate model guides the search towards promising regions in the hyperparameter space.

In addition to the improved (“optimized”) hyperparameter values, `spotPython` allows a statistical analysis, e.g., a sensitivity analysis, of the results. We can print the results of the hyperparameter tuning, see Table 4.

```
print(gen_design_table(fun_control=fun_control, spot=spot_tuner))
```

Table 4: Results of the hyperparameter tuning. The table shows the hyperparameters, their types, default values, lower and upper bounds, and the transformation function. The column “tuned” shows the tuned values. The column “importance” shows the importance of the hyperparameters. The column “stars” shows the importance of the hyperparameters in stars. The importance is computed by the SPOT software.

name	type	default	lower	upper	tuned	transform	importance	stars
l1	int	5	2.0	9.0	7.0	pow_2_int	100.00	***
l2	int	5	2.0	9.0	3.0	pow_2_int	96.29	***
lr_mult	float	1.0	0.1	10.0	0.1	None	0.00	
batchsize	int	4	1.0	5.0	4.0	pow_2_int	0.00	
epochs	int	3	3.0	4.0	4.0	pow_2_int	4.18	*
k_folds	int	2	0.0	0.0	0.0	None	0.00	
patience	int	5	3.0	3.0	3.0	None	0.00	
optimizer	factor	SGD	0.0	9.0	3.0	None	0.16	.

To visualize the most important hyperparameters, `spotPython` provides the function `plot_importance`. The following code generates the importance plot from Figure 3.

```
spot_tuner.plot_importance(threshold=0.025, filename="./figures" + experiment_name+"_importance")
```

3.11 Get SPOT Results

The architecture of the `spotPython` model can be obtained by the following code:

```
X = spot_tuner.to_all_dim(spot_tuner.min_X.reshape(1,-1))
model_spot = get_one_core_model_from_X(X, fun_control)
model_spot
```

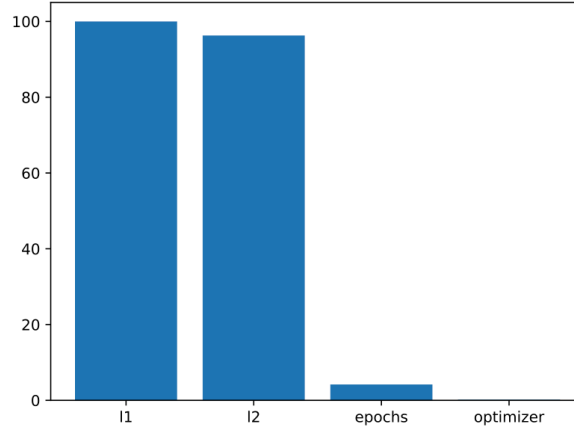


Figure 3: Variable importance

First, the numerical representation of the hyperparameters are obtained, i.e., the numpy array `X` is generated. This array is then used to generate the model `model_spot` by the function `get_one_core_model_from_X`. The model `model_spot` has the following architecture:

```
Net_CIFAR10(
  (conv1): Conv2d(3, 6, kernel_size=(5, 5), stride=(1, 1))
  (pool): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
  (conv2): Conv2d(6, 16, kernel_size=(5, 5), stride=(1, 1))
  (fc1): Linear(in_features=400, out_features=64, bias=True)
  (fc2): Linear(in_features=64, out_features=32, bias=True)
  (fc3): Linear(in_features=32, out_features=10, bias=True)
)
```

3.12 Get Default Hyperparameters

In a similar manner as in Section 3.11, the default hyperparameters can be obtained.

```
# fun_control was modified, we generate a new one with the original default hyperparameters
fc = copy.deepcopy(fun_control)
fc.update({"core_model_hyper_dict": hyper_dict[fun_control["core_model"].__name__]})
model_default = get_one_core_model_from_X(X_start, fun_control=fc)
```

The corresponding default model has the following architecture:

```

Net_CIFAR10(
    (conv1): Conv2d(3, 6, kernel_size=(5, 5), stride=(1, 1))
    (pool): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
    (conv2): Conv2d(6, 16, kernel_size=(5, 5), stride=(1, 1))
    (fc1): Linear(in_features=400, out_features=32, bias=True)
    (fc2): Linear(in_features=32, out_features=32, bias=True)
    (fc3): Linear(in_features=32, out_features=10, bias=True)
)

```

3.13 Evaluation of the Tuned Architecture

The method `train_tuned` takes a model architecture without trained weights and trains this model with the train data. The train data is split into train and validation data. The validation data is used for early stopping. The trained model weights are saved as a dictionary.

This evaluation is similar to the final evaluation in PyTorch (2023a).

```

train_tuned(net=model_default, train_dataset=train, shuffle=True,
            loss_function=fun_control["loss_function"],
            metric=fun_control["metric_torch"],
            device = DEVICE, show_batch_interval=1_000,)
test_tuned(net=model_default, test_dataset=test,
            loss_function=fun_control["loss_function"],
            metric=fun_control["metric_torch"],
            shuffle=False,
            device = DEVICE)

```

The following code trains the model `model_spot`. If `path` is set to a filename, e.g., `path = "model_spot_trained.pt"`, the weights of the trained model will be saved to this file.

```

train_tuned(net=model_spot, train_dataset=train,
            loss_function=fun_control["loss_function"],
            metric=fun_control["metric_torch"],
            shuffle=True,
            device = DEVICE,
            path=None)

```

```

Loss on hold-out set: 1.2267619131326675
Accuracy on hold-out set: 0.58955
Early stopping at epoch 13

```


If `path` is set to a filename, e.g., `path = "model_spot_trained.pt"`, the weights of the trained model will be loaded from this file.

```
test_tuned(net=model_spot, test_dataset=test,
           shuffle=False,
           loss_function=fun_control["loss_function"],
           metric=fun_control["metric_torch"],
           device = DEVICE)
```

Loss on hold-out set: 1.242568492603302

Accuracy on hold-out set: 0.5957

3.14 Comparison with Default Hyperparameters and Ray Tune

Table 5 shows the loss and accuracy of the default model, the model with the hyperparameters from SPOT, and the model with the hyperparameters from `ray[tune]`.

Table 5: Comparison of the loss and accuracy of the default model, the model with the hyperparameters from SPOT, and the model with the hyperparameters from `ray[tune]`. `ray[tune]` only shows the validation loss, because training loss is not reported by `ray[tune]`.

Model	Validation Loss	Validation Accuracy	Loss	Accuracy
Default	2.1221	0.2452	2.1182	0.2425
<code>spotPython</code>	1.2268	0.5896	1.2426	0.5957
<code>ray[tune]</code>	1.1815	0.5836	-	0.5806

3.15 Detailed Hyperparameter Plots

The contour plots in this section visualize the interactions of the three most important hyperparameters, `l1`, `l2`, and `epochs`, and `optimizer` of the surrogate model used to optimize the hyperparameters. Since some of these hyperparameters take factorial or integer values, sometimes step-like fitness landscapes (or response surfaces) are generated. SPOT draws the interactions of the main hyperparameters by default. It is also possible to visualize all interactions. For this, again refer to the notebook (Bartz-Beielstein 2023).

Figure 5 to Figure 10 show the contour plots of the loss as a function of the hyperparameters. These plots are very helpful for benchmark studies and for understanding neural networks. `spotPython` provides additional tools for a visual inspection of the results and give valuable

```

threshold = 0.025
impo = spot_tuner.print_importance(threshold=threshold, print_screen=True)
var_plots = [i for i, x in enumerate(impo) if x[1] > threshold]
min_z = min(spot_tuner.y)
max_z = max(spot_tuner.y)
n = spot_tuner.k
for i in var_plots:
    for j in var_plots:
        if j > i:
            filename = "./figures" + experiment_name+"_contour_"+str(i)+"_"+str(j)+".png"
            spot_tuner.plot_contour(i=i, j=j, min_z=min_z, max_z = max_z, filename=filename)

```

Figure 4: ?(caption)

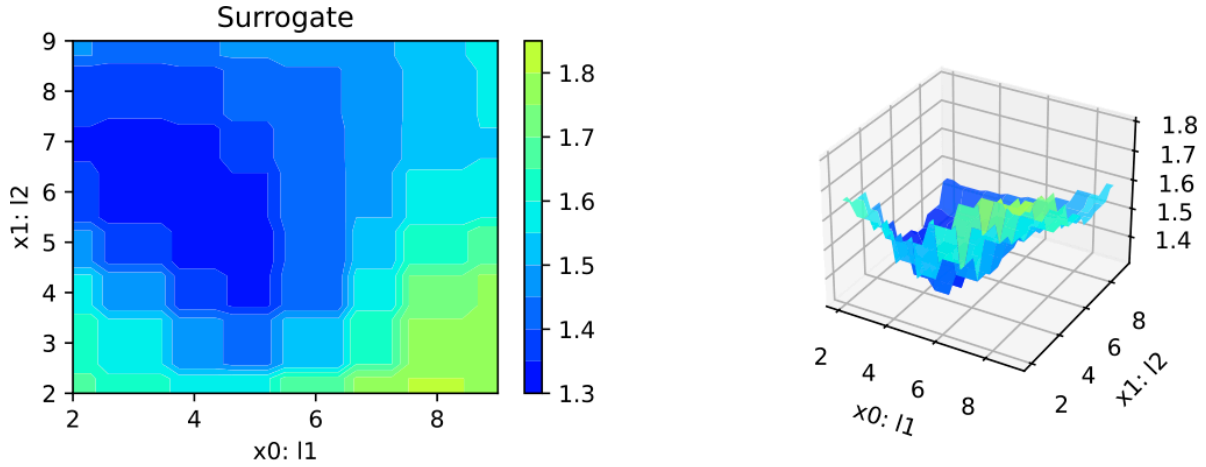


Figure 5: Contour plot of the loss as a function of l1 and l2, i.e., the number of neurons in the layers.

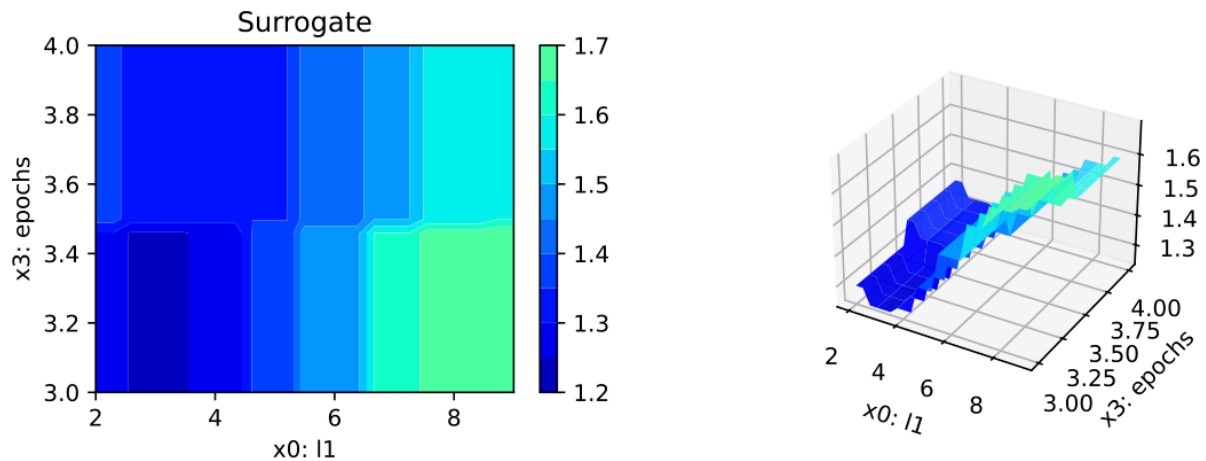


Figure 6: Contour plot of the loss as a function of the number of epochs and the neurons in layer 11.

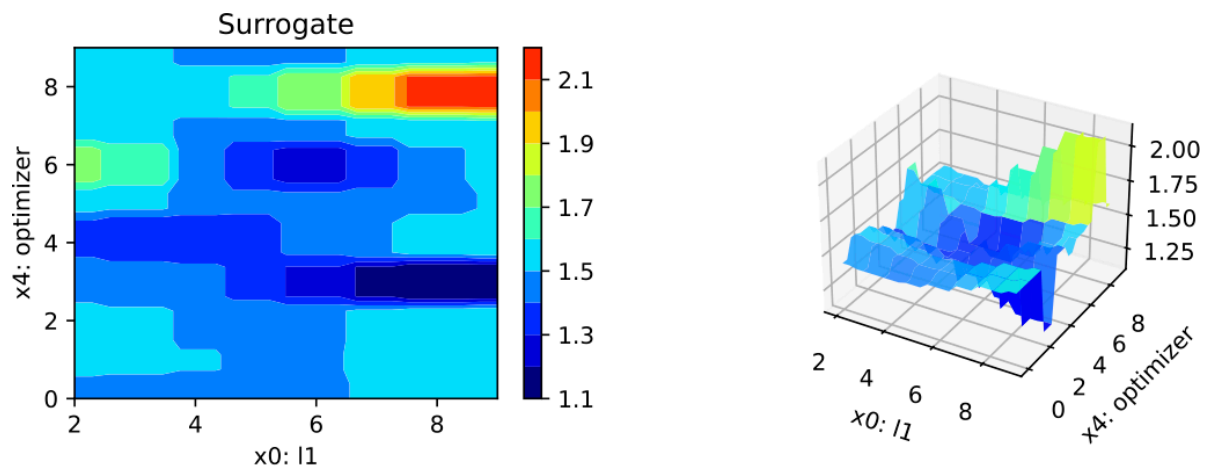


Figure 7: Contour plot of the loss as a function of the optimizer and the neurons in layer 11.

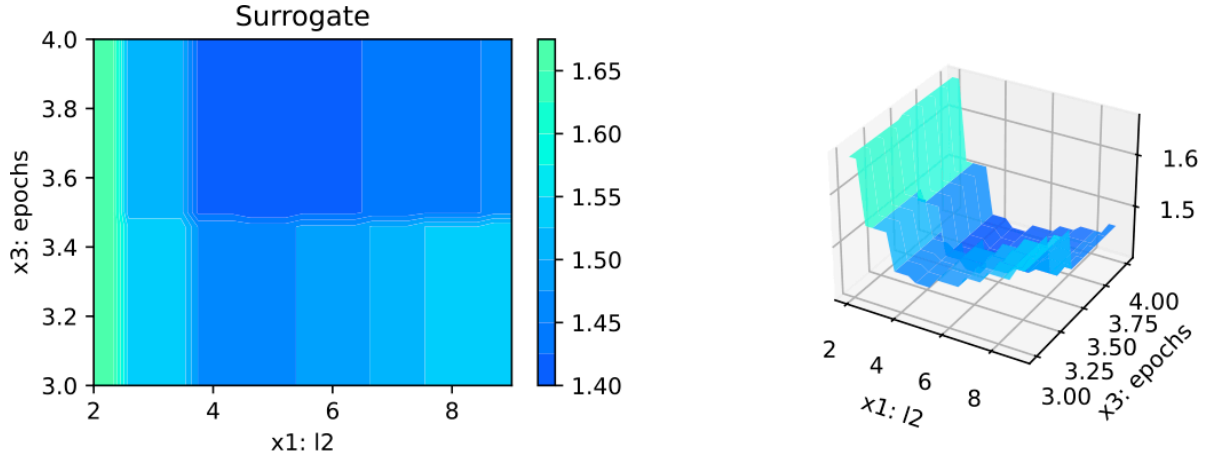


Figure 8: Contour plot of the loss as a function of the number of epochs and the neurons in layer 12.

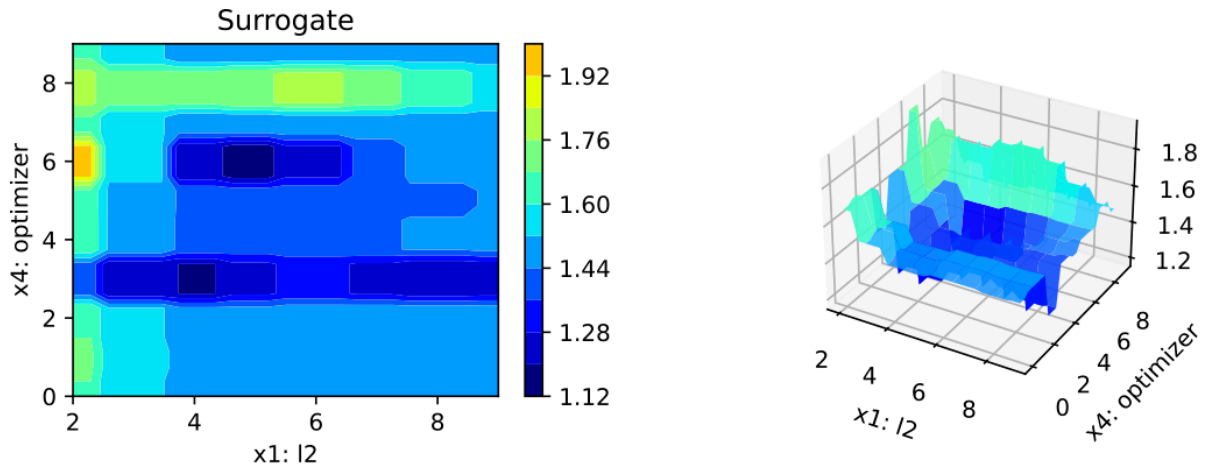


Figure 9: Contour plot of the loss as a function of the optimizer and the neurons in layer 12.

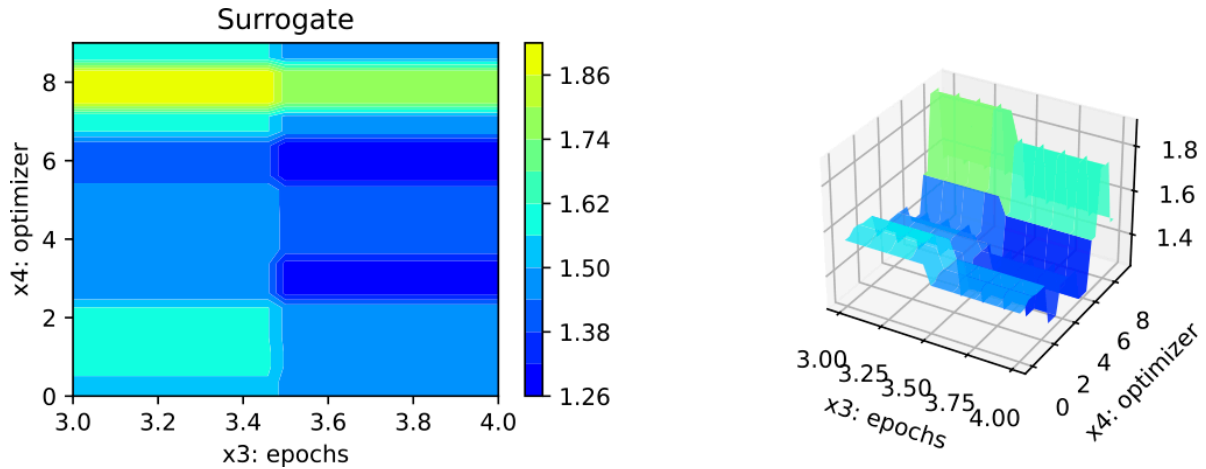


Figure 10: Contour plot of the loss as a function of the optimizer and the number of epochs.

insights into the hyperparameter tuning process. This is especially useful for model explainability, transparency, and trustworthiness. In addition to the contour plots, Figure 12 shows the parallel plot of the hyperparameters.

```
spot_tuner.parallel_plot()
```

Figure 11: ?(caption)

4 Summary and Outlook

This tutorial presents the hyperparameter tuning open source software `spotPython` for PyTorch. To show its basic features, a comparison with the “official” PyTorch hyperparameter tuning tutorial (PyTorch 2023a) is presented. Some of the advantages of `spotPython` are:

- Numerical and categorical hyperparameters.
- Powerful surrogate models.
- Flexible approach and easy to use.
- Simple JSON files for the specification of the hyperparameters.
- Extension of default and user specified network classes.
- Noise handling techniques.

Currently, only rudimentary parallel and distributed neural network training is possible, but these capabilities will be extended in the future. The next version of `spotPython` will also include a more detailed documentation and more examples.

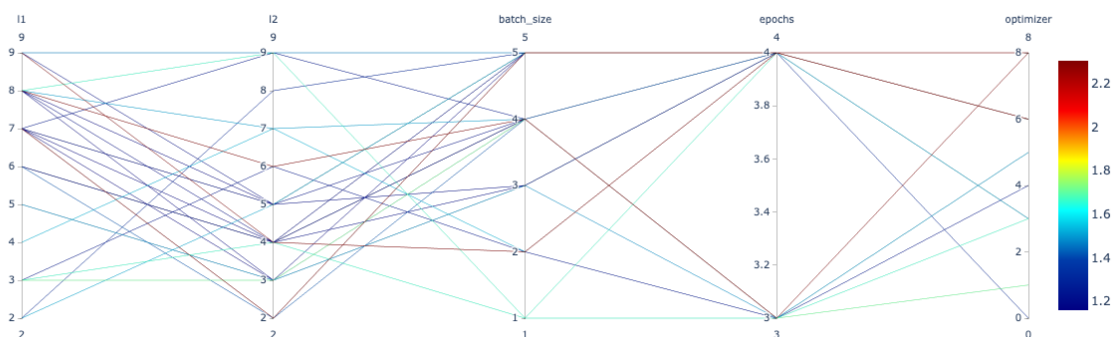


Figure 12: Parallel plot

! Important

Important: This tutorial does not present a complete benchmarking study (Bartz-Beielstein et al. 2020). The results are only preliminary and highly dependent on the local configuration (hard- and software). Our goal is to provide a first impression of the performance of the hyperparameter tuning package `spotPython`. To demonstrate its capabilities, a quick comparison with `ray[tune]` was performed. `ray[tune]` was chosen, because it is presented as “an industry standard tool for distributed hyperparameter tuning.” The results should be interpreted with care.

5 Appendix

Sample Output From Ray Tune’s Run

The output from `ray[tune]` could look like this (PyTorch 2023b):

Number of trials: 10 (10 TERMINATED)

l1	l2	lr	batch_size	loss	accuracy	training_iteration
64	4	0.00011629	2	1.87273	0.244	2
32	64	0.000339763	8	1.23603	0.567	8
8	16	0.00276249	16	1.1815	0.5836	10
4	64	0.000648721	4	1.31131	0.5224	8
32	16	0.000340753	8	1.26454	0.5444	8

	8		4		0.000699775		8		1.99594		0.1983		2	
	256		8		0.0839654		16		2.3119		0.0993		1	
	16		128		0.0758154		16		2.33575		0.1327		1	
	16		8		0.0763312		16		2.31129		0.1042		4	
	128		16		0.000124903		4		2.26917		0.1945		1	

Best trial config: {'l1': 8, 'l2': 16, 'lr': 0.00276249, 'batch_size': 16, 'data_dir': '...'}
Best trial final validation loss: 1.181501
Best trial final validation accuracy: 0.5836
Best trial test set accuracy: 0.5806

References

- Bartz, Eva, Thomas Bartz-Beielstein, Martin Zaefferer, and Olaf Mersmann, eds. 2022. *Hyperparameter Tuning for Machine and Deep Learning with R - A Practical Guide*. Springer.
- Bartz-Beielstein, Thomas. 2023. “PyTorch Hyperparameter Tuning with SPOT: Comparison with Ray Tuner and Default Hyperparameters on CIFAR10.” https://github.com/sequential-parameter-optimization/spotPython/blob/main/notebooks/14_spot_ray_hpt_torch_cifar10.ipynb.
- Bartz-Beielstein, Thomas, Jürgen Branke, Jörn Mehnen, and Olaf Mersmann. 2014. “Evolutionary Algorithms.” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 4 (3): 178–95.
- Bartz-Beielstein, Thomas, Carola Doerr, Jakob Bossek, Sowmya Chandrasekaran, Tome Eftimov, Andreas Fischbach, Pascal Kerschke, et al. 2020. “Benchmarking in Optimization: Best Practice and Open Issues.” arXiv. <https://arxiv.org/abs/2007.03488>.
- Bartz-Beielstein, Thomas, Christian Lasarczyk, and Mike Preuss. 2005. “Sequential Parameter Optimization.” In *Proceedings 2005 Congress on Evolutionary Computation (CEC’05), Edinburgh, Scotland*, edited by B McKay et al., 773–80. Piscataway NJ: IEEE Press.
- Lewis, R M, V Torczon, and M W Trosset. 2000. “Direct search methods: Then and now.” *Journal of Computational and Applied Mathematics* 124 (1–2): 191–207.
- Li, Lisha, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. 2016. “Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization.” *arXiv e-Prints*, March, arXiv:1603.06560.
- Meignan, David, Sigrid Knust, Jean-Marc Frayet, Gilles Pesant, and Nicolas Gaud. 2015. “A Review and Taxonomy of Interactive Optimization Methods in Operations Research.” *ACM Transactions on Interactive Intelligent Systems*, September.
- Montiel, Jacob, Max Halford, Saulo Martiello Mastelini, Geoffrey Bolmier, Raphael Sourty, Robin Vaysse, Adil Zouitine, et al. 2021. “River: Machine Learning for Streaming Data in Python.”
- PyTorch. 2023a. “Hyperparameter Tuning with Ray Tune.” https://pytorch.org/tutorials/beginner/hyperparameter_tuning_tutorial.html.
- . 2023b. “Training a Classifier.” https://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html.