

✓ Introduction

- The following notebook focuses on applying Optuna in Reinforcement Learning.
- Source / Reference of this notebook: https://colab.research.google.com/github/araffin/tools-for-robotic-rl-icra2022/blob/main/notebooks/optuna_lab.ipynb#scrollTo=4UU17YpjymPr
- To study the application of Optuna in Deep Learning, refer to: <https://www.geeksforgeeks.org/hyperparameter-tuning-with-optuna-in-pytorch/>

Quick guides on the step needed.

1. Define a config.
2. Define a search space or a function that returns the parameters that define the models.
3. Define an objective score function that will return the objective score function for the sampled set of hyperparameter.
4. Create a optimization loop
 - For each trial:
 - a) sample a hyperparameter
 - b) Use the hyperparameter to train. At intervals, evaluate the performance of the models and decided if to prune this trial.

✓ Step 0: Library Installation and Import

```
1 # Install optuna library
2 !pip install optuna
```

 [Show hidden output](#)

```
1 # Install Stable Baseline 3
2 !pip install stable-baselines3==2.0.0a5
```

 [Show hidden output](#)

```
1 # Install gymnasium
2
3 !pip install swig
4 !pip install gymnasium[box2d]
```

 [Show hidden output](#)

```
1 import optuna
2 from optuna.pruners import MedianPruner
3 from optuna.samplers import TPESampler
4 from optuna.visualization import plot_optimization_history, plot_param_importances
5
6
7 from stable_baselines3.common.callbacks import EvalCallback
8 from stable_baselines3.common.env_util import make_vec_env
9 from stable_baselines3.a2c import A2C
10 from stable_baselines3 import PPO
11
12 import gymnasium as gym
13
14 import torch as th
15 import torch.nn as nn
```

✓ DEMO: Tune a A2C agent that plays CartPole-v1

✓ Step 1: Create config / key parameters

- Terminology:
 1. **TRIALS** - Each **TRIAL** is initiated with different sampled set of hyperparameter. If needed, multiple **JOBS** can be initiated in parallel for each trial. Each **TRIAL** will involve training the agent / model for **N_TIMESTEPS**.
 2. **EVAL_EPISODES** - During the training that last for **N_TIMESTEPS**, at an interval of **EVAL_FREQ**, evaluation will be performed. For each evaluation, **N_EVAL_EPISODES** of evaluation episodes will be sampled and reviewed. This may help the scheduler to decide whether to

prune early.

```

1 # Config
2 # Hyperparameter Optimization Loop
3 N_TRIALS = 100 # Maximum number of trials during Hyperparameter Optimization Loop
4 N_JOBS = 1 # Number of parallel jobs to run during each trials in Hyperparameter Optimization Loop
5 N_STARTUP_TRIALS = 5 # Number of trials to perform random sampling (without relying on sampler) during the Hyperparameter Optimization Loop
6 TIMEOUT = int(60*15) # Maximum number of times (in seconds) the entire loop is allowed up to.
7
8 # Evaluation Parameter for each set of hyperparameter
9 N_TIMESTEPS = int(2e4) # Training budget - Number of time steps in one FULL TRIAL for each set of hyperparameter.
10 N_EVALUATIONS = 2 # Number of intermediate evaluations performed in one FULL TRIAL for each set of hyperparameter.
11 EVAL_FREQ = int(N_TIMESTEPS / N_EVALUATIONS) # Step interval for each intermediate evaluations during one FULL TRIAL.
12 N_EVAL_EPISODES = 10 # Number of episodes to be sampled for each evaluation.
13
14
15 # Environment Parameter
16 N_EVAL_ENVS = 5 # Number of environments used in parallel during evaluation.
17
18 ENV_ID = "CartPole-v1" # ID of the environments, to be initiated with gym.make()
19
20 # Algorithm Parameter
21 ALGO_NAME = "A2C"

```

✓ Step 2: Define the search space

```

1 def sample_a2c_params(trial):
2
3     """
4     Sample a set of hyperparameter to be trial'd.
5
6     Args:
7         trial (optuna.Trial) : An Optuna trial object.
8
9     Returns:
10         params (dict): The hyperparameters to be trial'd. Its key matches the keywords used in defining models.
11     """
12
13     # Refer to this link: https://stable-baselines3.readthedocs.io/en/master/modules/a2c.html#example
14     # To study the hyperparameter to be updated.
15
16     #####
17     # Discount factor #
18     #####
19     # suggest.float -> sample from a continuous space (float)
20     # "gamma" - name (to be showcased in the final plot)
21     # log - means sample from log space
22     gamma = 1 - trial.suggest_float("one_minus_gamma", 0.0001, 0.1, log = True)
23
24     # Create another attribute to store the actual gamma value
25     trial.set_user_attr("gamma", gamma)
26
27     #####
28     # Maximum value for gradient clipping #
29     #####
30     max_grad_norm = trial.suggest_float("max_grad_norm", 0.3, 0.5, log=True)
31
32     #####
33     # Number of steps to run for each environment per update #
34     #####
35     n_steps = 2 ** trial.suggest_int("exponent_n_steps", 3, 10)
36
37     # Create another attribute to store the actual n_steps
38     trial.set_user_attr("n_steps", n_steps)
39
40     #####
41     # Learning_rate #
42     #####
43     learning_rate = trial.suggest_float("learning_rate", 1e-5, 1, log=True)
44
45     #####
46     # Network architecture #
47     #####
48     # https://stable-baselines3.readthedocs.io/en/master/modules/a2c.html#a2c-policies
49
50     net_arch = trial.suggest_categorical("net_arch", ["tiny", "small"])
51
52     # The net_arch expects a list, that is why it is wrapped in a list
53     net_arch = [{"pi": [64], "vf": [64]} if net_arch == 'tiny' \

```

```

54         else {"pi" : [64, 64], "vf" : [64, 64]]]
55
56 #####
57 # Activation Function #
58 #####
59 activation_fn = trial.suggest_categorical("activation_fn", ["tanh", "relu"])
60
61 activation_fn = {"tanh": nn.Tanh, "relu": nn.ReLU}[activation_fn]
62
63
64
65 # Note: The key used in this dictionary match the key used in defining the models
66 # Therefore, the naming convention must be followed.
67 params = {"n_steps": n_steps,
68           "gamma": gamma,
69           "learning_rate": learning_rate,
70           "max_grad_norm": max_grad_norm,
71           "policy_kwargs": {"net_arch": net_arch,
72                             "activation_fn": activation_fn}}
73
74 return params
75

```

▼ Step 3: Define objective

- A custom callback function is defined to report the results of periodic evaluations.

```

1 class TrialEvalCallback(EvalCallback):
2
3     """
4     Callback used for evaluating and reporting a trial.
5
6     Args:
7         eval_env (gym.env): An evaluation environment.
8         trial (Optuna.trial): An Optuna trial object.
9         n_eval_episodes (int): Number of evaluation episodes for each evaluation.
10        eval_freq (int): Step interval for an intermediate evaluation in each trial.
11        deterministic (boolean): Whether the evaluation should use stochastic or deterministic policy.
12        verbose (int):
13
14    Returns:
15        out (boolean):
16    """
17
18    def __init__(self, eval_env, trial, n_eval_episodes, eval_freq, deterministic, verbose = 0):
19
20        super().__init__(eval_env = eval_env, n_eval_episodes = n_eval_episodes,
21                        eval_freq = eval_freq, deterministic = deterministic,
22                        verbose = verbose)
23
24        self.trial = trial
25        self.eval_idx = 0
26        self.is_pruned = False
27
28    def _on_step(self):
29        if self.eval_freq > 0 and self.n_calls % self.eval_freq == 0:
30            super()._on_step()
31            self.eval_idx += 1
32
33            # Send report to optuna
34            self.trial.report(self.last_mean_reward, self.eval_idx)
35
36            # Prune trial if needed
37            if self.trial.should_prune():
38                self.is_pruned = True
39                return False
40            return True

```

- The true objective function.

```

1 def objective(trial):
2
3     """
4     A function that returns the objective score that decides the quality of a set of hyperparameter.
5
6     Args:
7         trial (optuna.Trial): An Optuna trial object.
8

```

```

9 Returns:
10     objective_score (float): The score that represents the quality of this set of hyperparameter.
11     """
12
13 # Create the default keyword arguments (those that wasnt defined in the hyperparameter sampling function)
14 kwargs = {"policy": "MlpPolicy",
15           "env": ENV_ID}
16
17 # Update with the inclusion of the sampled hyperparameter
18 kwargs.update(sample_a2c_params(trial))
19
20 # Create a model using the sampled hyperparameter
21 model = A2C(**kwargs)
22
23 # Create the environments
24 eval_envs = make_vec_env(env_id = ENV_ID,
25                           n_envs = N_EVAL_ENVS)
26
27 # Create the call back for reporting evaluation results
28 eval_callback = TrialEvalCallback(eval_env = eval_envs,
29                                  trial = trial,
30                                  n_eval_episodes = N_EVAL_EPISODES,
31                                  eval_freq = EVAL_FREQ,
32                                  deterministic = True)
33
34 nan_encountered = False
35 try:
36     model.learn(N_TIMESTEPS, callback = eval_callback)
37 except AssertionError as e:
38     # Sometimes, randomly sampled error can lead to NaN
39     print(e)
40     nan_encountered = True
41 finally:
42     # At the end of training or if error is encountered
43     # Free Memory
44     model.env.close()
45     eval_envs.close()
46
47 # Inform the optimizer that a non-valid hyperparameter is sampled
48 if nan_encountered:
49     return float('nan')
50
51 if eval_callback.is_pruned:
52     raise optuna.exceptions.TrialPruned()
53
54 return eval_callback.last_mean_reward
55
56
57

```

✓ Step 4: Define Hyperparameter Optimization Loop

```

1 # Set PyTorch num threads to 1 for faster training
2 # Parallel environment will demand heavy use of CPU.
3 # Therefore, this line limits to the usage of cpu for PyTorch to be only 1 line.
4 th.set_num_threads(1)
5
6 # Select a sampler
7 # https://optuna.readthedocs.io/en/stable/reference/generated/optuna.samplers.TPESampler.html
8 # n_startup_trials -> Number of trials at the beginning that sample a set of hyperparameter randomly instead of using the algorithm
9 # This allows the creation of initial database.
10 sampler = TPESampler(n_startup_trials = N_STARTUP_TRIALS)
11
12 # Select a scheduler / pruner
13 # https://optuna.readthedocs.io/en/stable/reference/generated/optuna.pruners.MedianPruner.html
14 # n_startup_trials -> Pruning is disabled at the beginning for this many trials for initial database creation.
15 # n_warmup_steps -> Number of steps in each trial that disable the pruning.
16 pruner = MedianPruner(n_startup_trials = N_STARTUP_TRIALS,
17                       n_warmup_steps = N_TIMESTEPS // 3)
18
19 # Create a study for Hyperparameter Optimization
20 # https://optuna.readthedocs.io/en/stable/reference/generated/optuna.create_study.html
21 study = optuna.create_study(sampler = sampler,
22                             pruner = pruner,
23                             direction = "maximize")
24
25 try:
26     # https://optuna.readthedocs.io/en/stable/reference/generated/optuna.study.Study.html#optuna.study.Study.optimize
27     study.optimize(objective,
28                   n_trials = N_TRIALS,

```

```
29         timeout = TIMEOUT,
30         n_jobs = N_JOBS)
31 except KeyboardInterrupt:
32     pass
33
34 # Print the meta info for the hyperparameter optimization process
35 print(f"Number of finished trials: {len(study.trials)}")
36 trial = study.best_trial
37 print(f"Best trial: {trial.value}")
38
39 print("Params: ")
40 for key, value in trial.params.items():
41     print(f"  {key}: {value}")
42
43 print("User Attributes: ")
44 for key, value in trial.user_attrs.items():
45     print(f"  {key}: {value}")
46
47
48 # Write report
49 study.trials_dataframe().to_csv(f"study_result_{ALGO_NAME}_{ENV_ID}.csv")
50
51 # Show plot
52 fig1 = plot_optimization_history(study)
53 fig2 = plot_param_importances(study)
54
55 fig1.show()
56 fig2.show()
57
```

[I 2025-06-12 09:35:47,739] A new study created in memory with name: no-name-813ee568-ff50-443a-8067-43acc1197b30
/usr/local/lib/python3.11/dist-packages/stable_baselines3/common/policies.py:460: UserWarning:

As shared layers in the mlp_extractor are removed since SB3 v1.8.0, you should now pass directly a dictionary and not a list (net_

```
[I 2025-06-12 09:36:19,525] Trial 0 finished with value: 9.2 and parameters: {'one_minus_gamma': 0.012924707618275728, 'max_grad_r
[I 2025-06-12 09:36:49,605] Trial 1 finished with value: 9.2 and parameters: {'one_minus_gamma': 0.038723746654246494, 'max_grad_r
[I 2025-06-12 09:37:18,442] Trial 2 finished with value: 9.3 and parameters: {'one_minus_gamma': 0.002021039321497278, 'max_grad_r
[I 2025-06-12 09:37:45,863] Trial 3 finished with value: 91.7 and parameters: {'one_minus_gamma': 0.0008731735219335531, 'max_grad_r
[I 2025-06-12 09:38:14,860] Trial 4 finished with value: 116.7 and parameters: {'one_minus_gamma': 0.0012304139554903581, 'max_grad_r
[I 2025-06-12 09:38:43,782] Trial 5 finished with value: 122.5 and parameters: {'one_minus_gamma': 0.0001086993609073862, 'max_grad_r
[I 2025-06-12 09:39:24,908] Trial 6 finished with value: 84.9 and parameters: {'one_minus_gamma': 0.000103139578231042, 'max_grad_r
[I 2025-06-12 09:39:54,671] Trial 7 finished with value: 481.1 and parameters: {'one_minus_gamma': 0.00010095783296398445, 'max_grad_r
[I 2025-06-12 09:40:23,040] Trial 8 finished with value: 463.6 and parameters: {'one_minus_gamma': 0.00035606177672072123, 'max_grad_r
[I 2025-06-12 09:40:55,476] Trial 9 finished with value: 500.0 and parameters: {'one_minus_gamma': 0.008262752774391171, 'max_grad_r
[I 2025-06-12 09:41:33,076] Trial 10 finished with value: 9.5 and parameters: {'one_minus_gamma': 0.007654797065370172, 'max_grad_r
[I 2025-06-12 09:42:04,835] Trial 11 finished with value: 138.9 and parameters: {'one_minus_gamma': 0.08561534553697434, 'max_grad_r
[I 2025-06-12 09:42:35,400] Trial 12 finished with value: 392.9 and parameters: {'one_minus_gamma': 0.007621459534012113, 'max_grad_r
[I 2025-06-12 09:43:05,790] Trial 13 finished with value: 500.0 and parameters: {'one_minus_gamma': 0.0004338676441916004, 'max_grad_r
[I 2025-06-12 09:43:39,324] Trial 14 finished with value: 44.6 and parameters: {'one_minus_gamma': 0.00042913443902935054, 'max_grad_r
[I 2025-06-12 09:44:08,172] Trial 15 finished with value: 9.5 and parameters: {'one_minus_gamma': 0.003222336477326619, 'max_grad_r
[I 2025-06-12 09:44:41,767] Trial 16 finished with value: 149.9 and parameters: {'one_minus_gamma': 0.004097294434259803, 'max_grad_r
[I 2025-06-12 09:45:10,447] Trial 17 finished with value: 118.7 and parameters: {'one_minus_gamma': 0.023388281427777047, 'max_grad_r
[I 2025-06-12 09:45:42,328] Trial 18 finished with value: 266.5 and parameters: {'one_minus_gamma': 0.00039643173176913117, 'max_grad_r
[I 2025-06-12 09:46:14,635] Trial 19 finished with value: 208.1 and parameters: {'one_minus_gamma': 0.0011274364044843253, 'max_grad_r
[I 2025-06-12 09:46:49,696] Trial 20 finished with value: 110.1 and parameters: {'one_minus_gamma': 0.0002099399167953067, 'max_grad_r
[I 2025-06-12 09:47:21,125] Trial 21 finished with value: 294.8 and parameters: {'one_minus_gamma': 0.0001840475379181496, 'max_grad_r
[I 2025-06-12 09:47:55,061] Trial 22 finished with value: 428.9 and parameters: {'one_minus_gamma': 0.0006772044746225003, 'max_grad_r
[I 2025-06-12 09:48:27,116] Trial 23 finished with value: 172.5 and parameters: {'one_minus_gamma': 0.0002309687834281705, 'max_grad_r
[I 2025-06-12 09:49:01,142] Trial 24 finished with value: 406.1 and parameters: {'one_minus_gamma': 0.005463870698583285, 'max_grad_r
[I 2025-06-12 09:49:32,196] Trial 25 finished with value: 500.0 and parameters: {'one_minus_gamma': 0.0015837290597985752, 'max_grad_r
[I 2025-06-12 09:50:01,660] Trial 26 finished with value: 9.1 and parameters: {'one_minus_gamma': 0.002124909567287219, 'max_grad_r
[I 2025-06-12 09:50:40,218] Trial 27 finished with value: 65.6 and parameters: {'one_minus_gamma': 0.001853162648692653, 'max_grad_r
[I 2025-06-12 09:51:11,973] Trial 28 finished with value: 9.2 and parameters: {'one_minus_gamma': 0.016521266890261135, 'max_grad_r
```

Number of finished trials: 29

Best trial: 500.0

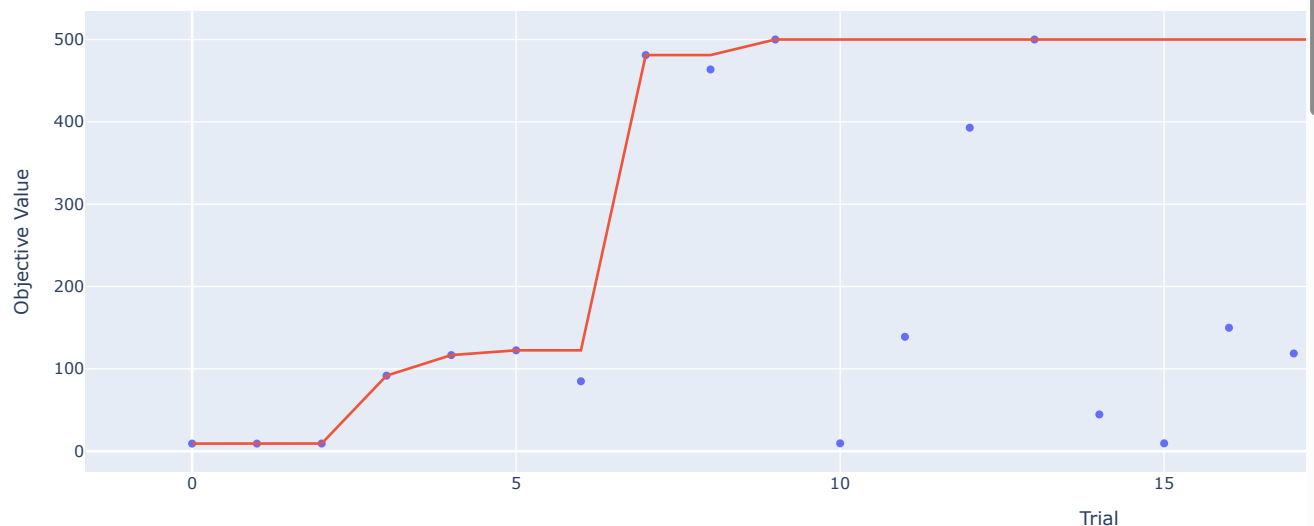
Params:

```
one_minus_gamma: 0.008262752774391171
max_grad_norm: 0.4988993250029
exponent_n_steps: 5
learning_rate: 0.0009036800602866176
net_arch: small
activation_fn: tanh
```

User Attributes:

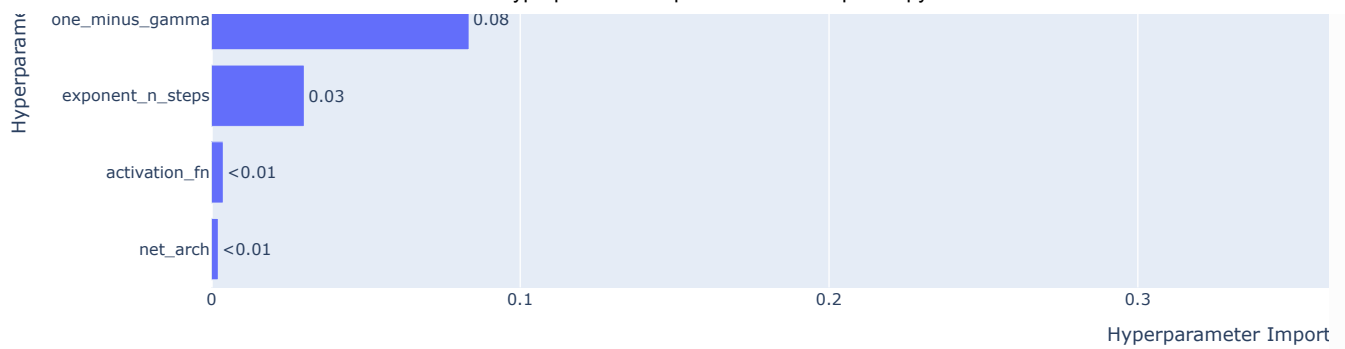
```
gamma: 0.9917372472256089
n_steps: 32
```

Optimization History Plot



Hyperparameter Importances





✓ Practice: Tune a PPO agent that plays LunarLander-v2

✓ Step 1: Create a config

```

1
2
3 #####
4 # Step 1: Create a config #
5 #####
6 # Hyperparameter Optimization Loop
7 N_TRIALS = 100 # Maximum number of trials during Hyperparameter Optimization Loop
8 N_JOBS = 1 # Number of parallel jobs to run during each trials in Hyperparameter Optimization Loop
9 N_STARTUP_TRIALS = 5 # Number of trials to perform random sampling (without relying on sampler) during the Hyperparameter Optimization Loop
10 TIMEOUT = int(60*15) # Maximum number of times (in seconds) the entire loop is allowed up to.
11
12 # Evaluation Parameter for each set of hyperparameter
13 N_TIMESTEPS = int(1e4) # Training budget - Number of time steps in one FULL TRIAL for each set of hyperparameter.
14 N_EVALUATIONS = 2 # Number of intermediate evaluations performed in one FULL TRIAL for each set of hyperparameter.
15 EVAL_FREQ = int(N_TIMESTEPS / N_EVALUATIONS) # Step interval for each intermediate evaluations during one FULL TRIAL.
16 N_EVAL_EPISODES = 10 # Number of episodes to be sampled for each evaluation.
17
18
19 # Environment Parameter
20 N_EVAL_ENVS = 16 # Number of environments used in parallel during evaluation.
21
22 ENV_ID = "LunarLander-v2" # ID of the environments, to be initiated with gym.make()
23
24 # Algorithm Parameter
25 ALGO_NAME = "PPO"

```

✓ Step 2: Define the search space

```

1
2 #####
3 # Step 2: Define a function that samples the hyperparameter #
4 #####
5
6 def sample_ppo_params(trial):
7
8     """
9     Sample a set of hyperparameter to be trial'd.
10
11     Args:
12         trial (optuna.Trial) : An Optuna trial object.
13
14     Returns:
15         params (dict): The hyperparameters to be trial'd. Its key matches the keywords used in defining models.
16     """
17
18     # To study the hyperparameter to be updated.
19
20     #####
21     # Discount factor #
22     #####
23     # suggest.float -> sample from a continuous space (float)
24     # "gamma" - name (to be showcased in the final plot)
25     # log - means sample from log space
26     gamma = 1 - trial.suggest_float("one_minus_gamma", 0.0001, 0.1, log = True)

```

```

27
28 # Create another attribute to store the actual gamma value
29 trial.set_user_attr("gamma", gamma)
30
31 #####
32 # Maximum value for gradient clipping #
33 #####
34 max_grad_norm = trial.suggest_float("max_grad_norm", 0.3, 0.5, log=True)
35
36 #####
37 # Number of steps to run for each environment per update #
38 #####
39 n_steps = 2 ** trial.suggest_int("exponent_n_steps", 3, 11)
40
41 # Create another attribute to store the actual n_steps
42 trial.set_user_attr("n_steps", n_steps)
43
44 #####
45 # Learning_rate #
46 #####
47 learning_rate = trial.suggest_float("learning_rate", 1e-5, 1, log=True)
48
49 #####
50 # Network architecture #
51 #####
52 net_arch = trial.suggest_categorical("net_arch", ["tiny", "small"])
53
54 # The net_arch expects a list, that is why it is wrapped in a list
55 net_arch = [{"pi": [64], "vf": [64]} if net_arch == 'tiny' \
56             else {"pi" : [64, 64], "vf" : [64, 64]}]
57
58 #####
59 # Activation Function #
60 #####
61 activation_fn = trial.suggest_categorical("activation_fn", ["tanh", "relu"])
62
63 activation_fn = {"tanh": nn.Tanh, "relu": nn.ReLU}[activation_fn]
64
65
66
67 # Note: The key used in this dictionary match the key used in defining the models
68 # Therefore, the naming convention must be followed.
69 params = {"n_steps": n_steps,
70          "gamma": gamma,
71          "learning_rate": learning_rate,
72          "max_grad_norm": max_grad_norm,
73          "policy_kwargs": {"net_arch": net_arch,
74                           "activation_fn": activation_fn}}
75
76 return params

```

▼ Step 3: Define objective functions

```

1
2
3 #####
4 # Step 3: Objective function #
5 #####
6 # 3.1: Callback
7 class TrialEvalCallback(EvalCallback):
8
9     """
10     Callback used for evaluating and reporting a trial.
11
12     Args:
13         eval_env (gym.env): An evaluation environment.
14         trial (Optuna.trial): An Optuna trial object.
15         n_eval_episodes (int): Number of evaluation episodes for each evaluation.
16         eval_freq (int): Step interval for an intermediate evaluation in each trial.
17         deterministic (boolean): Whether the evaluation should use stochastic or deterministic policy.
18         verbose (int):
19
20     Returns:
21         out (boolean):
22     """
23
24     def __init__(self, eval_env, trial, n_eval_episodes, eval_freq, deterministic, verbose = 0):
25
26         super().__init__(eval_env = eval_env, n_eval_episodes = n_eval_episodes,
27                          eval_freq = eval_freq, deterministic = deterministic,

```



```

28         verbose = verbose)
29     self.trial = trial
30     self.eval_idx = 0
31     self.is_pruned = False
32
33     def _on_step(self):
34         if self.eval_freq > 0 and self.n_calls % self.eval_freq == 0:
35             super()._on_step()
36             self.eval_idx += 1
37
38             # Send report to optuna
39             self.trial.report(self.last_mean_reward, self.eval_idx)
40
41             # Prune trial if needed
42             if self.trial.should_prune():
43                 self.is_pruned = True
44                 return False
45             return True
46
47         # 3.2: The definition of the objective score.
48 def objective(trial):
49     """
50     A function that returns the objective score that decides the quality of a set of hyperparameter.
51
52     Args:
53         trial (optuna.Trial): An Optuna trial object.
54
55     Returns:
56         objective_score (float): The score that represents the quality of this set of hyperparameter.
57     """
58
59     # Creat the default keyword arguments (those that wasnt defined in the hyperparameter sampling function)
60     kwargs = {"policy": "MlpPolicy",
61              "env": ENV_ID}
62
63     # Update with the inclusion of the sampled hyperparameter
64     kwargs.update(sample_ppo_params(trial))
65
66     # Create a model using the sampled hyperparameter
67     model = PPO(**kwargs)
68
69     # Create the environments
70     eval_envs = make_vec_env(env_id = ENV_ID,
71                             n_envs = N_EVAL_ENVS)
72
73     # Create the call back for reporting evaluation results
74     eval_callback = TrialEvalCallback(eval_env = eval_envs,
75                                     trial = trial,
76                                     n_eval_episodes = N_EVAL_EPISODES,
77                                     eval_freq = EVAL_FREQ,
78                                     deterministic = True)
79
80     nan_encountered = False
81     try:
82         model.learn(N_TIMESTEPS, callback = eval_callback)
83     except AssertionError as e:
84         # Sometimes, randomly sampled error can lead to NaN
85         print(e)
86         nan_encountered = True
87     finally:
88         # At the end of training or if error is encountered
89         # Free Memory
90         model.env.close()
91         eval_envs.close()
92
93     # Inform the optimizer that a non-valid hyperparameter is sampled
94     if nan_encountered:
95         return float('nan')
96
97     if eval_callback.is_pruned:
98         raise optuna.exceptions.TrialPruned()
99
100     return eval_callback.last_mean_reward
101

```

▼ Step 4: Hyperparameter Optimization Loop

```

1
2     #####
3     # Step 4: Optimization Loop #

```

```

4 #####
5
6 # Set PyTorch num threads to 1 for faster training
7 # Parallel environment will demand heavy use of CPU.
8 # Therefore, this line limits to the usage of cpu for PyTorch to be only 1 line.
9 th.set_num_threads(1)
10
11 # Select a sampler
12 # https://optuna.readthedocs.io/en/stable/reference/samplers/generated/optuna.samplers.TPESampler.html
13 # n_startup_trials -> Number of trials at the beginning that sample a set of hyperparameter randomly instead of using the algorithm
14 # This allows the creation of initial database.
15 sampler = TPESampler(n_startup_trials = N_STARTUP_TRIALS)
16
17 # Select a scheduler / pruner
18 # https://optuna.readthedocs.io/en/stable/reference/generated/optuna.pruners.MedianPruner.html
19 # n_startup_trials -> Pruning is disabled at the beginning for this many trials for initial database creation.
20 # n_warmup_steps -> Number of steps in each trial that disable the pruning.
21 pruner = MedianPruner(n_startup_trials = N_STARTUP_TRIALS,
22                       n_warmup_steps = N_TIMESTEPS // 3)
23
24 # Create a study for Hyperparameter Optimization
25 # https://optuna.readthedocs.io/en/stable/reference/generated/optuna.create\_study.html
26 study = optuna.create_study(sampler = sampler,
27                             pruner = pruner,
28                             direction = "maximize")
29
30 try:
31     # https://optuna.readthedocs.io/en/stable/reference/generated/optuna.study.Study.html#optuna.study.Study.optimize
32     study.optimize(objective,
33                   n_trials = N_TRIALS,
34                   timeout = TIMEOUT,
35                   n_jobs = N_JOBS)
36 except KeyboardInterrupt:
37     pass
38
39 # Print the meta info for the hyperparameter optimization process
40 print(f"Number of finished trials: {len(study.trials)}")
41 trial = study.best_trial
42 print(f"Best trial: {trial.value}")
43
44 print("Params: ")
45 for key, value in trial.params.items():
46     print(f" {key}: {value}")
47
48 print("User Attributes: ")
49 for key, value in trial.user_attrs.items():
50     print(f" {key}: {value}")
51
52
53 # Write report
54 study.trials_dataframe().to_csv(f"study_result_{ALGO_NAME}_{ENV_ID}.csv")
55
56 # Show plot
57 fig1 = plot_optimization_history(study)
58 fig2 = plot_param_importances(study)
59
60 fig1.show()
61 fig2.show()
62
63
64

```

 Show hidden output

✓ Practice: Tune a DQN agent that plays Atari Games - Space Invader

✓ Install Libraries

```
1 !pip install git+https://github.com/DLR-RM/rl-baselines3-zoo
```

 Show hidden output

```
1 !pip install gymnasium[atari]
2 !pip install gymnasium[accept-rom_license]
```

 Show hidden output

✓ Create a config file

- Save the config file as `dqn.yml`
- Note: The config file is only mostly used as a placeholder. During Hyperparameter optimization, a set of hyperparameter will be randomly sampled, thus replacing these values.
- The type and range of hyperparameter values to be sampled may be referred to: https://github.com/DLR-RM/rl-baselines3-zoo/blob/master/rl_zoo3/hyperparams_opt.py#L222

```

1 # The default config file
2 SpaceInvadersNoFrameskip-v4:
3   env_wrapper:
4     - stable_baselines3.common.atari_wrappers.AtariWrapper
5   frame_stack: 4 #Every 4 frame as 1 input to allow the model to learn the trajectories of the object.
6   policy: 'CnnPolicy'
7   n_timesteps: !!float 1e2 # 1e6 (Recommended, but shortened in this notebook as its only for demo)
8   buffer_size: 100000
9   learning_rate: !!float 1e-4
10  batch_size: 32
11  learning_starts: 100000
12  target_update_interval: 1000
13  train_freq: 4
14  gradient_steps: 1
15  exploration_fraction: 0.1
16  exploration_final_eps: 0.01
17  # If True, you need to deactivate handle_timeout_termination
18  # in the replay_buffer_kwargs
19  optimize_memory_usage: False

```

✓ Access the API for hyperparameter optimization (Built-in)

- Refer to the raw code to find out the usage: https://github.com/DLR-RM/rl-baselines3-zoo/blob/master/rl_zoo3/train.py
- The full command is as following: `!python -m rl_zoo3.train --algo dqn --env SpaceInvadersNoFrameskip-v4 -f logs/ -c dqn.yml -optimize --optimization-log-path logs/optimization --eval-episodes 10 --n-eval-envs 1 --max-total-trials 100 --n-jobs 1 --sampler "tpe" --pruner "median" --n-startup-trials 10 --n-evaluations 2`

```
1 !python -m rl_zoo3.train --algo dqn --env SpaceInvadersNoFrameskip-v4 -f logs/ -c dqn.yml -optimize --optimization-log-path logs/opt
```

 [Show hidden output](#)

✓ Post Notes: How to apply it onto a PyTorch Model?

- A good reference link: <https://www.geeksforgeeks.org/hyperparameter-tuning-with-optuna-in-pytorch/>
- The main difference will be on the definition of the objective function and how to manually report back the intermediate and/or final evaluation to optune.
- The following codes are generated by chatGPT on how to accomplish both. (Note: It has not yet to be tested)

```

1 def objective(trial):
2     # Sample hyperparameters
3     lr = trial.suggest_float("lr", 1e-5, 1e-1, log=True)
4     n_units = trial.suggest_int("n_units", 16, 128)
5
6     # Model

```