

---

# **baconian Documentation**

***Release 0.2.6***

**Linsen Dong**

**May 22, 2020**



<b>1</b>	<b>Installation Guide</b>	<b>1</b>
<b>2</b>	<b>Step by step guide to run a RL experiment</b>	<b>3</b>
<b>3</b>	<b>Best practice and core concepts of Baconian</b>	<b>9</b>
<b>4</b>	<b>Examples</b>	<b>15</b>
<b>5</b>	<b>Logging and Visualization</b>	<b>37</b>
<b>6</b>	<b>How to implement a new algorithm</b>	<b>43</b>
<b>7</b>	<b>How to implement a new environment</b>	<b>51</b>
<b>8</b>	<b>How to implement a new dynamics model</b>	<b>55</b>
<b>9</b>	<b>System Overview of Baconian</b>	<b>59</b>
<b>10</b>	<b>Baconian API Reference</b>	<b>63</b>
<b>11</b>	<b>Contributions</b>	<b>79</b>
<b>12</b>	<b>Indices and tables</b>	<b>81</b>
	<b>Python Module Index</b>	<b>83</b>
	<b>Index</b>	<b>85</b>



# CHAPTER 1

---

## Installation Guide

---

Baconian is easy to install. We offer the pip requirement file to install the required packages. Make sure your machine has python 3.5, 3.6, or 3.7 with Ubuntu 16.04, 18.04 (recommend python 3.5 and ubuntu 16.04).

1. We recommend you to use anaconda to manage your package and environment, since installing the required packages may overwrite some of already installed packages with a different version.

```
source activate your_env
```

2. Install by pip or clone the source code and install the packages with requirement file:

```
// pip install:
pip install baconian
// source code install
git clone git@github.com:cap-ntu/baconian-project.git baconian
cd ./baconian
pip install pip -U
pip install tensorflow==1.15.2 // or pip install tensorflow-gpu==1.15.2
pip install -e .
```

Then you are free to go. You can either use the Baconian as a third party package and import it into your own project, or directly modify it.

After you finish the above installation, you are able to run the following environments of Gym:

- algorithmic
- toy\_text
- classic\_control

3. Support for Gym full environments

If you want to use the full environments of gym, please refer to [gym](#) to obtain the license and library.

Mostly you will need to install requirements for mujoco environments and box-2d environments.

4. Support for DeepMind Control Suit

For DeepMind control suit, you should install it by:

```
pip install git+git://github.com/deepmind/dm_control.git
```

DeepMind control suit also relies on the mujoco engine which is the same as the mujoco-py environments in gym. They provide similar tasks with slightly differences.

And for the different default mujoco key and mujoco binaries for mujoco-py and DeepMind control suit, please follow the setting of mujoco-py and we will take care of the setting for DeepMind Control Suit at runtime.

5. We have implemented many examples, you may try them first at [Examples](#)
6. If you prefer to conduct a new experiments by yourself, you can follow the tutorial here [How to implement](#)

---

## Step by step guide to run a RL experiment

---

This is a step by step guide of how to compose an model-based RL experiments in Baconian, we will take the example of Dyna algorithm, which is a very typical model-based RL architecture proposed by Sutton in 1990.

In this method, we need the environment and the task we aim to solve, a model-free method, an approximator as dynamics for real world environment. We also need to specific the hyper-parameters in Dyna architecture, e.g., how many samples we should get from real environment to update the dynamics model, and how many samples from real environment and dynamics model to serve the training of the model-free algorithms.

For more information on the modules used in the following codes, please see *Best Practices* and *API references*

---

**Note:** Complete codes can be found at [here](#)

---

### 2.1 Create the tasks and DDPG algorithms

Create environment Pendulum-v0, and the DDPG algorithms to solve the task

```
1 def task_fn():
2     # create the gym environment by make function
3     env = make('Pendulum-v0')
4     # give your experiment a name which is used to generate the log path etc.
5     name = 'demo_exp'
6     # construct the environment specification
7     env_spec = EnvSpec(obs_space=env.observation_space,
8                        action_space=env.action_space)
9     # construct the neural network to approximate q function of DDPG
10    mlp_q = MLPQValueFunction(env_spec=env_spec,
11                             name_scope=name + '_mlp_q',
12                             name=name + '_mlp_q',
13                             mlp_config=[
14                                 {
15                                     "ACT": "RELU",
```

(continues on next page)

(continued from previous page)

```

16         "B_INIT_VALUE": 0.0,
17         "NAME": "1",
18         "N_UNITS": 16,
19         "TYPE": "DENSE",
20         "W_NORMAL_STDDEV": 0.03
21     },
22     {
23         "ACT": "LINEAR",
24         "B_INIT_VALUE": 0.0,
25         "NAME": "OUPTUT",
26         "N_UNITS": 1,
27         "TYPE": "DENSE",
28         "W_NORMAL_STDDEV": 0.03
29     }
30 ]])
31 # construct the neural network to approximate policy for DDPG
32 policy = DeterministicMLPPolicy(env_spec=env_spec,
33                                name_scope=name + '_mlp_policy',
34                                name=name + '_mlp_policy',
35                                mlp_config=[
36                                    {
37                                        "ACT": "RELU",
38                                        "B_INIT_VALUE": 0.0,
39                                        "NAME": "1",
40                                        "N_UNITS": 16,
41                                        "TYPE": "DENSE",
42                                        "W_NORMAL_STDDEV": 0.03
43                                    },
44                                    {
45                                        "ACT": "LINEAR",
46                                        "B_INIT_VALUE": 0.0,
47                                        "NAME": "OUPTUT",
48                                        "N_UNITS": env_spec.flat_action_dim,
49                                        "TYPE": "DENSE",
50                                        "W_NORMAL_STDDEV": 0.03
51                                    }
52                                ],
53                                reuse=False)
54 # construct the DDPG algorithms
55 ddpq = DDPG(
56     env_spec=env_spec,
57     config_or_config_dict={
58         "REPLAY_BUFFER_SIZE": 10000,
59         "GAMMA": 0.999,
60         "CRITIC_LEARNING_RATE": 0.001,
61         "ACTOR_LEARNING_RATE": 0.001,
62         "DECAY": 0.5,
63         "BATCH_SIZE": 50,
64         "TRAIN_ITERATION": 1,
65         "critic_clip_norm": 0.1,
66         "actor_clip_norm": 0.1,
67     },
68     value_func=mlp_q,
69     policy=policy,
70     name=name + '_ddpg',
71     replay_buffer=None
72 )

```



## 2.2 Create a global dynamics model with MLP network

```

1  # construct a neural network based global dynamics model to approximate the state_
  ↪ transition of environment
2  mlp_dyna = ContinuousMLPGlobalDynamicsModel(
3      env_spec=env_spec,
4      name_scope=name + '_mlp_dyna',
5      name=name + '_mlp_dyna',
6      learning_rate=0.01,
7      state_input_scaler=RunningStandardScaler(dims=env_spec.flat_obs_dim),
8      action_input_scaler=RunningStandardScaler(dims=env_spec.flat_action_dim),
9      output_delta_state_scaler=RunningStandardScaler(dims=env_spec.flat_obs_dim),
10     mlp_config=[
11         {
12             "ACT": "RELU",
13             "B_INIT_VALUE": 0.0,
14             "NAME": "l1",
15             "L1_NORM": 0.0,
16             "L2_NORM": 0.0,
17             "N_UNITS": 16,
18             "TYPE": "DENSE",
19             "W_NORMAL_STDDEV": 0.03
20         },
21         {
22             "ACT": "LINEAR",
23             "B_INIT_VALUE": 0.0,
24             "NAME": "OUPTUT",
25             "L1_NORM": 0.0,
26             "L2_NORM": 0.0,
27             "N_UNITS": env_spec.flat_obs_dim,
28             "TYPE": "DENSE",
29             "W_NORMAL_STDDEV": 0.03
30         }
31     ])

```

## 2.3 Create the Dyna architecture as algorithms

Create the Dyna algorithms by passing the ddpq and dynamics model, and wrap by the agent

```

1  # finally, construct the Dyna algorithms with a model free algorithm DDGP, and a_
  ↪ NN model.
2  algo = Dyna(env_spec=env_spec,
3      name=name + '_dyna_algo',
4      model_free_algo=ddpg,
5      dynamics_model=mlp_dyna,
6      config_or_config_dict=dict(
7          dynamics_model_train_iter=10,
8          model_free_algo_train_iter=10
9      ))
10 # To make the NN based dynamics model a proper environment so be a sampling_
  ↪ source for DDGP, reward function and
11 # terminal function need to be set.
12
13 # For examples only, we use random reward function and terminal function with_
  ↪ fixed episode length.

```

(continues on next page)

(continued from previous page)

```

14     algo.set_terminal_reward_function_for_dynamics_env(
15         terminal_func=FixedEpisodeLengthTerminalFunc(max_step_length=env.unwrapped._
↪max_episode_steps,
16                                                         step_count_fn=algo.dynamics_env.
↪total_step_count_fn),
17         reward_func=RandomRewardFunc())
18     # construct agent with additional exploration strategy if needed.
19     agent = Agent(env=env, env_spec=env_spec,
20                  algo=algo,
21                  algo_saving_scheduler=PeriodicalEventSchedule(
22                      t_fn=lambda: get_global_status_collect() ('TOTAL_AGENT_TRAIN_
↪SAMPLE_COUNT'),
23                      trigger_every_step=20,
24                      after_t=10),
25                      name=name + '_agent',
26                      exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
↪space,
27                                                         init_random_prob=0.5))

```

## 2.4 Configure the workflow and experiments object

Create the dyna-like workflow, and the experiments object, and run the experiment within your *task\_fn*

```

1     # construct the training flow, called Dyna flow. It defines how the training_
↪proceed, and the terminal condition
2     flow = create_dyna_flow(
3         train_algo_func=(agent.train, (), dict(state='state_agent_training')),
4         train_algo_from_synthesized_data_func=(agent.train, (), dict(state='state_
↪agent_training')),
5         train_dynamics_func=(agent.train, (), dict(state='state_dynamics_training')),
6         test_algo_func=(agent.test, (), dict(sample_count=1)),
7         test_dynamics_func=(agent.algo.test_dynamics, (), dict(sample_count=10,
↪env=env)),
8         sample_from_real_env_func=(agent.sample, (), dict(sample_count=10,
9                                                         env=agent.env,
10                                                         store_flag=True)),
11         sample_from_dynamics_env_func=(agent.sample, (), dict(sample_count=10,
12                                                         env=agent.algo.dynamics_
↪env,
13                                                         store_flag=True)),
14         train_algo_every_real_sample_count_by_data_from_real_env=40,
15         train_algo_every_real_sample_count_by_data_from_dynamics_env=40,
16         test_algo_every_real_sample_count=40,
17         test_dynamics_every_real_sample_count=40,
18         train_dynamics_every_real_sample_count=20,
19         start_train_algo_after_sample_count=1,
20         start_train_dynamics_after_sample_count=1,
21         start_test_algo_after_sample_count=1,
22         start_test_dynamics_after_sample_count=1,
23         warm_up_dynamics_samples=1
24     )
25     # construct the experiment
26     experiment = Experiment(
27         tuner=None,

```

(continues on next page)

(continued from previous page)

```

28     env=env,
29     agent=agent,
30     flow=flow,
31     name=name + '_exp'
32 )
33 # run!
34 experiment.run()
35
36
37 from baconian.core.experiment_runner import *
```

**Note:** Don't confuse the workflow with the Dyna algorithm itself. The flow only specifies how the algorithms interact with environments, and how to update and evaluate the ddpg (model-free method) and dynamics model.

## 2.5 Set the global configuration and launch the experiment

Set some global config if needed and wrap all the above task into a function *task\_fn*, and pass into experiment runner.

```

1  # set some global configuration here
2
3  # set DEFAULT_EXPERIMENT_END_POINT to indicate when to stop the experiment.
4  # one usually used is the TOTAL_AGENT_TRAIN_SAMPLE_COUNT, i.e., how many samples/
5  ↪timesteps are used for training
6  GlobalConfig().set('DEFAULT_EXPERIMENT_END_POINT', dict(TOTAL_AGENT_TRAIN_SAMPLE_
7  ↪COUNT=200))
8
9  # set the logging path to write log and save model checkpoints.
10 GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')
11
12 # feed the task into a exp runner.
13 single_exp_runner(task_fn, del_if_log_path_existed=True)
```

## 2.6 Results logging/ Visualization

Please refer to *Logging and Visualization*.



---

## Best practice and core concepts of Baconian

---

Here we introduce the core ideas and features about Baconian, to make sure you utilize the code correctly. As for the detailed usage of different algorithms, dynamics, please refer to the [API](#) page

### 3.1 Put your task into a function

Here we introduce the basic usage of Baconian, and introduce how it can help you to set up the model-based RL experiments.

First of all, whenever you want to run some algorithms, or any codes within the baconian, simply you need to wrap your code into a function, and pass this function to `single_exp_runner` or `duplicate_exp_runner`. In this method, Baconian will do some internal initialization of logging, experiment set-up etc.

`single_exp_runner` will run your function for once. As for `duplicate_exp_runner`, it is designed for running multiple experiments in a row, because in RL experiments, we usually run the experiment with a certain set of parameters but with different seeds to get a more stable results. So use `duplicate_exp_runner` can easily help you to achieve this, and the log file will be stored into sub-directory under your home log directory respectively.

Specifically, you can do it by:

```
from baconian.core.experiment_runner import single_exp_runner, duplicate_exp_runner
# Define your function first.
def your_function():
    a = 1
    b = 2
    print(a + b)
# Then pass the function object to single_exp_runner, then it will set up everything,
↳ and run your code.
single_exp_runner(your_function)
# Or call duplicate_exp_runner to run multiple experiments in a row. 10 is the number,
↳ of experiments:
duplicate_exp_runner(10, your_function)
```

## 3.2 Global Configuration

The global configuration offer the setting including default log path, log level, and some other system related default configuration. We implement the global configuration module with singleton method, and you can utilize it by following examples:

```
from baconian.config.global_config import GlobalConfig
from baconian.core.experiment_runner import single_exp_runner, duplicate_exp_runner
def your_function():
    a = 1;
    b = 2;
    print(a + b)
# Use GlobalConfig() to access the instance of GlobalConfig
# anywhere you want, and set the log path by yourself
# First argument is key you want to set, e.g., DEFAULT_LOG_PATH
# Second argument is the value.
GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')
single_exp_runner(task_fn, del_if_log_path_existed=True)
```

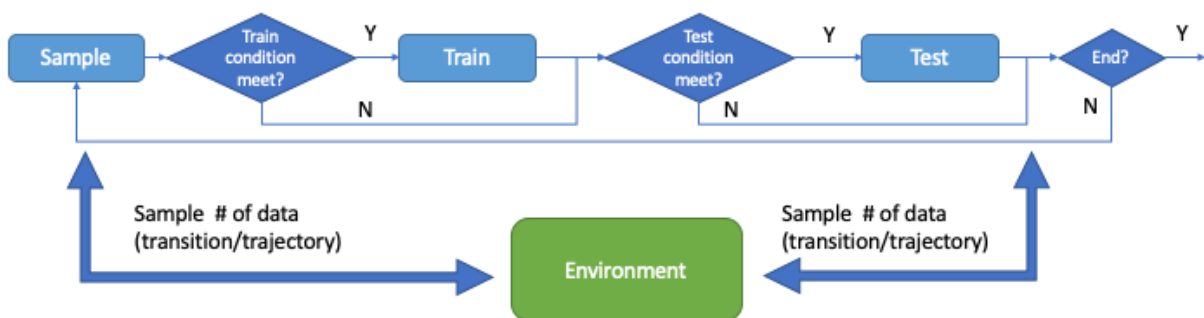
During the time task is running, the global configuration will be frozen, if you try to change it, an error will be raised.

## 3.3 Workflow for RL Experiments

In Baconian, the control flow of the experiments is delegated to an independent module `baconian.core.flow.train_test_flow:Flow` which is an abstract class. The reason to do so is to improve the flexibility and extensibility of framework. Two typical flow are implemented. One is `baconian.core.flow.train_test_flow:TrainTestFlow`, which corresponds to the pipeline of most model-free algorithms, which is sampling-training-testing pipeline. The other one is `baconian.core.flow.dyna_flow.py:DynaFlow`, which is the flow in Dyna algorithm [Sutton, 1992].

A typical train-test flow follows the diagram below:

Train Test Flow Flowchart



- Train/Test condition:  $(t > x)$  and  $(t - \text{last train/test time step} > n)$
- End condition:  $t > N$
- $t$  is the current time step, called a timer.
  - usually is the number of data the agent gets from environment.
  - can be changed to other timer if it's monotonic, e.g., total times the agent execute training.
- $x$  represents start the train/test at least after  $x$  time steps
- $n$  represents trigger the train/test every  $n$  time steps.
- $N$  is the max value for the time steps.

---

**Note:** The reason why we add this module with these parameters settings is to enable the user to fully customize the experiment execution process.

---

For setting the end point N in the flow, see the section you can achieve by setting the global configuration in following code example. For using other status as end point, see section [Built-in Global Status/Counter](#) for more built-in status.

A flow module is required to pass in one step function as the timer to control the process. More detailed explanations of the time step function are given below.

For `baconian.core.flow.dyna_flow.py:DynaFlow`, comparing to `baconian.core.flow.train_test_flow:TrainTestFlow`, there are three more processes are added: `train_agent_from_model`, `train_dynamics_model`, and `test_dynamics_model`. They follows the similar ideas on how to control these processes, user can refer to its usage and example codes for more details.

## 3.4 Time Step Function

RL experiment often relies on a timer/counter to indicate the progress of the experiment. It can be used to schedule the parameters from modules like action noise, exploration strategy. It can also be used to timestamp a recorded log like the evaluated performance of agent so you can know the changes of the performance along with the different timestep.

Naturally, the number of samples generated from environment for training purpose is used as the time step function. This value can be retrieved from following code. Or you can built the time step function using any status value from any objects (i.e., agent, env, algorithm) as long as it is a monotonic counter.

For using built-in agent status counter as the time step function, you can see more in the following section.

## 3.5 Built-in Global Status/Counter

We include some recorded values as a global shared status which can be accessed during the experiments.

The all built-in recorded values are listed below:

- `TOTAL_AGENT_TRAIN_SAMPLE_COUNT`: the timesteps/samples used by agent for training
- `TOTAL_AGENT_TRAIN_SAMPLE_FUNC_COUNT`: the times of sampling function called by agent during training
- `TOTAL_AGENT_TEST_SAMPLE_COUNT`: the timesteps/samples used by agent for testing
- `TOTAL_AGENT_UPDATE_COUNT`: the times of training function called by agent
- `TOTAL_ENV_STEP_TRAIN_SAMPLE_COUNT`: the timesteps used by environment during training, it differs a little from `TOTAL_AGENT_TRAIN_SAMPLE_COUNT`
- `TOTAL_ENV_STEP_TEST_SAMPLE_COUNT`: the timesteps used by environment during testing, it differs a little from `TOTAL_AGENT_TEST_SAMPLE_COUNT`

User can access these values from anywhere they want or register new global status into by following code snippet after you start to execute the experiment:

```
from baconian.core.status import get_global_status_collect
# access directly
print(get_global_status_collect() ('TOTAL_AGENT_TRAIN_SAMPLE_COUNT'))
# wrap in to a function
```

(continues on next page)

(continued from previous page)

```
train_sample_count_func=lambda: get_global_status_collect() ('TOTAL_AGENT_TRAIN_SAMPLE_
↪COUNT')

# register new status
get_global_status_collect().register_info_key_status(
    # object that hold the source_
    ↪value
    obj=object,
    # which key the object used
    info_key='predict_counter',
    # under which status
    under_status='TRAIN',
    # name used to store in global_
    ↪status
    return_name='TOTAL_AGENT_TRAIN_
    ↪SAMPLE_COUNT'
)
```

Example Usages:

1. Use it to decide when to end the experiment using DEFAULT\_EXPERIMENT\_END\_POINT in Global Config User can set by:

```
# only set one value with its limit as the end point.
# once the recorded value TOTAL_AGENT_TRAIN_SAMPLE_COUNT exceed 200, the experiment_
↪will end.
GlobalConfig().set('DEFAULT_EXPERIMENT_END_POINT', dict(TOTAL_AGENT_TRAIN_SAMPLE_
↪COUNT=200))
```

2. Access the value as the scheduler parameters' clock:

See Page *Scheduler Parameters*.

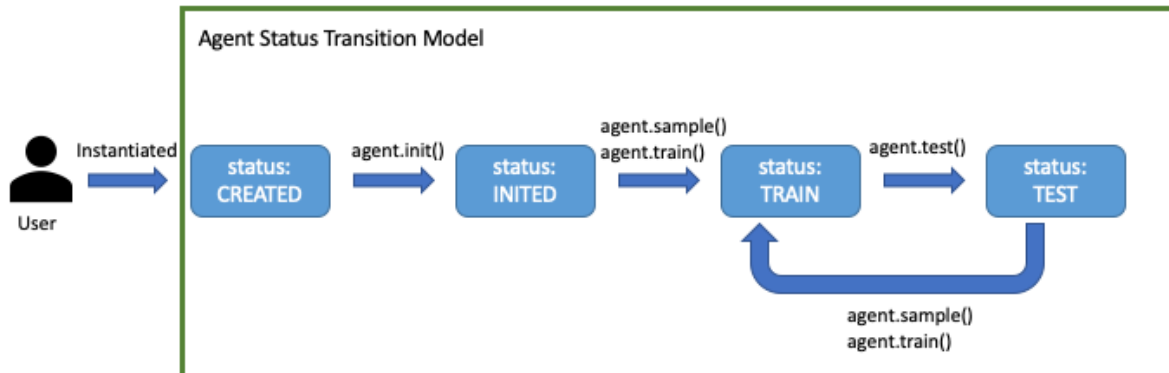
## 3.6 Stateful Behaviour

Status control is a must for DRL experiments. For instance, off-policy DRL methods need to switch between behavior policy and target policy during sampling and testing or decay the exploration action noise w.r.t the training progress status.

The following figure shows the status transition of agent and how it control the behaviour of agent.



## Status Cycle of Agent



- `agent.init()`: initialize neural network parameters, warm up the algorithm (e.g., sample a batch of data to get the mean and variance of states)
- `agent.sample()`, `agent.train()`: sample the data under the TRAIN status, and append the data into the memory. Train agent's algorithm using data from memory. Under this status, agent will enable the exploration strategy, action noise, etc.
- `agent.test()`: test the agent, the exploration strategy/action noise, etc. will be disabled.

Every class that inherited from `baconian.core.core:Basic` will have two class attributes: `STATUS_LIST` which contains all status of this class or module.

You can call `set_status` method to change the status of one instance. You can call `get_status` method to get the current status of an instance, which is a dict type. The return value not only contains the status (i.e., TRAIN, TEST) but also other extra information that is specially added in the code. Such as, you can register a counter of a function by doing so:

```

1      else:
2          if self.noise_adder and not self.is_testing:
3              res = self.env_spec.action_space.clip(self.noise_adder(self.algo.
↪predict(**kwargs)))
4          else:
5              res = self.algo.predict(**kwargs)
6              self.recorder.append_to_obj_log(obj=self, attr_name='action', status_
↪info=self.get_status(), value=res)
7          return res
8
9      @register_counter_info_to_status_decorator(increment=1, info_key='sample_counter',
↪under_status=('TRAIN', 'TEST'),
10                                              ignore_wrong_status=True)
11      def sample(self, env, sample_count: int, in_which_status: str = 'TRAIN', store_
↪flag=False,
12              sample_type: str = 'transition') -> (
13          TransitionData, TrajectoryData):
14          """
15          sample a certain number of data from environment
16
17          :param env: environment to sample
18          :param sample_count: int, sample count
19          :param in_which_status: string, environment status
20          :param store_flag: to store environment samples or not, default False
  
```

The counter for calling function `predict` is added as one attribute of status, which will be returned with the key

`predict_counter`

For detailed usage of these methods, you can find it in API reference.

We offer some examples on how to better use Baconian in a more practical and efficient way.

## 4.1 DQN with Acrobot-v1

```
from baconian.algo.dqn import DQN
from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
from baconian.algo.value_func.mlp_q_value import MLPQValueFunction
from baconian.core.agent import Agent
from baconian.algo.misc import EpsilonGreedy
from baconian.core.experiment import Experiment
from baconian.core.flow.train_test_flow import create_train_test_flow
from baconian.config.global_config import GlobalConfig
from baconian.common.schedules import LinearScheduler
from baconian.core.status import get_global_status_collect

def task_fn():
    env = make('Acrobot-v1')
    name = 'demo_exp'
    env_spec = EnvSpec(obs_space=env.observation_space,
                       action_space=env.action_space)
    mlp_q = MLPQValueFunction(env_spec=env_spec,
                              name_scope=name + '_mlp_q',
                              name=name + '_mlp_q',
                              mlp_config=[
                                  {
                                      "ACT": "TANH",
                                      "B_INIT_VALUE": 0.0,
                                      "NAME": "1",
                                      "N_UNITS": 64,
                                      "TYPE": "DENSE",
```

(continues on next page)

(continued from previous page)

```

        "W_NORMAL_STDDEV": 0.03
    },
    {
        "ACT": "TANH",
        "B_INIT_VALUE": 0.0,
        "NAME": "2",
        "N_UNITS": 64,
        "TYPE": "DENSE",
        "W_NORMAL_STDDEV": 0.03
    },
    {
        "ACT": "RELU",
        "B_INIT_VALUE": 0.0,
        "NAME": "3",
        "N_UNITS": 256,
        "TYPE": "DENSE",
        "W_NORMAL_STDDEV": 0.03
    },
    {
        "ACT": "LINEAR",
        "B_INIT_VALUE": 0.0,
        "NAME": "OUPTUT",
        "N_UNITS": 1,
        "TYPE": "DENSE",
        "W_NORMAL_STDDEV": 0.03
    }
]
dqn = DQN(env_spec=env_spec,
          config_or_config_dict=dict(REPLAY_BUFFER_SIZE=50000,
                                     GAMMA=0.99,
                                     BATCH_SIZE=32,
                                     LEARNING_RATE=0.001,
                                     TRAIN_ITERATION=1,
                                     DECAY=0),
          name=name + '_dqn',
          value_func=mlp_q)
agent = Agent(env=env, env_spec=env_spec,
              algo=dqn,
              name=name + '_agent',
              exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
→space,
                                              prob_scheduler=LinearScheduler(
                                              t_fn=lambda: get_global_
→status_collect() (
                                              'TOTAL_AGENT_TRAIN_
→SAMPLE_COUNT'),
                                              schedule_timesteps=int(0.1 *
→100000),
                                              initial_p=1.0,
                                              final_p=0.02),
                                              init_random_prob=0.1),
              noise_adder=None)

flow = create_train_test_flow(
    test_every_sample_count=100,
    train_every_sample_count=1,
    start_test_after_sample_count=0,

```

(continues on next page)

(continued from previous page)

```

start_train_after_sample_count=1000,
sample_func_and_args=(agent.sample, (), dict(sample_count=1,
                                              env=agent.env,
                                              store_flag=True)),

train_func_and_args=(agent.train, (), dict()),
test_func_and_args=(agent.test, (), dict(sample_count=3)),
)
experiment = Experiment(
    tuner=None,
    env=env,
    agent=agent,
    flow=flow,
    name=name
)
experiment.run()

from baconian.core.experiment_runner import *

GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')
single_exp_runner(task_fn, del_if_log_path_existed=True)

```

## 4.2 DDPG with Pendulum-v0

```

"""
A simple example to show how to build up an experiment with ddpq training and testing_
on Pendulum-v0
"""

from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
from baconian.algo.value_func.mlp_q_value import MLPQValueFunction
from baconian.algo.ddpg import DDPG
from baconian.algo.policy import DeterministicMLPPolicy
from baconian.core.agent import Agent
from baconian.algo.misc import EpsilonGreedy
from baconian.core.experiment import Experiment
from baconian.core.flow.train_test_flow import create_train_test_flow
from baconian.config.global_config import GlobalConfig
from baconian.core.status import get_global_status_collect
from baconian.common.schedules import PeriodicalEventSchedule
import baconian.common.log_data_loader as loader
from pathlib import Path

def task_fn():
    env = make('Pendulum-v0')
    name = 'demo_exp'
    env_spec = EnvSpec(obs_space=env.observation_space,
                      action_space=env.action_space)

    mlp_q = MLPQValueFunction(env_spec=env_spec,
                             name_scope=name + '_mlp_q',
                             name=name + '_mlp_q',
                             mlp_config=[

```

(continues on next page)

(continued from previous page)

```

        {
            "ACT": "RELU",
            "B_INIT_VALUE": 0.0,
            "NAME": "1",
            "N_UNITS": 16,
            "TYPE": "DENSE",
            "W_NORMAL_STDDEV": 0.03
        },
        {
            "ACT": "LINEAR",
            "B_INIT_VALUE": 0.0,
            "NAME": "OUPTUT",
            "N_UNITS": 1,
            "TYPE": "DENSE",
            "W_NORMAL_STDDEV": 0.03
        }
    ])
    policy = DeterministicMLPPolicy(env_spec=env_spec,
                                   name_scope=name + '_mlp_policy',
                                   name=name + '_mlp_policy',
                                   mlp_config=[
                                       {
                                           "ACT": "RELU",
                                           "B_INIT_VALUE": 0.0,
                                           "NAME": "1",
                                           "N_UNITS": 16,
                                           "TYPE": "DENSE",
                                           "W_NORMAL_STDDEV": 0.03
                                       },
                                       {
                                           "ACT": "LINEAR",
                                           "B_INIT_VALUE": 0.0,
                                           "NAME": "OUPTUT",
                                           "N_UNITS": env_spec.flat_action_dim,
                                           "TYPE": "DENSE",
                                           "W_NORMAL_STDDEV": 0.03
                                       }
                                   ],
                                   reuse=False)

    ddpq = DDPG(
        env_spec=env_spec,
        config_or_config_dict={
            "REPLAY_BUFFER_SIZE": 10000,
            "GAMMA": 0.999,
            "CRITIC_LEARNING_RATE": 0.001,
            "ACTOR_LEARNING_RATE": 0.001,
            "DECAY": 0.5,
            "BATCH_SIZE": 50,
            "TRAIN_ITERATION": 1,
            "critic_clip_norm": 0.1,
            "actor_clip_norm": 0.1,
        },
        value_func=mlp_q,
        policy=policy,
        name=name + '_ddpg',
        replay_buffer=None
    )

```

(continues on next page)

(continued from previous page)

```

    )
    agent = Agent(env=env, env_spec=env_spec,
                  algo=ddpg,
                  algo_saving_scheduler=PeriodicalEventSchedule(
                      t_fn=lambda: get_global_status_collect() ('TOTAL_AGENT_TRAIN_
↪SAMPLE_COUNT'),
                      trigger_every_step=20,
                      after_t=10),
                  name=name + '_agent',
                  exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
↪space,
                                                    init_random_prob=0.5))

    flow = create_train_test_flow(
        test_every_sample_count=10,
        train_every_sample_count=10,
        start_test_after_sample_count=5,
        start_train_after_sample_count=5,
        train_func_and_args=(agent.train, (), dict()),
        test_func_and_args=(agent.test, (), dict(sample_count=1)),
        sample_func_and_args=(agent.sample, (), dict(sample_count=100,
                                                    env=agent.env,
                                                    store_flag=True))
    )

    experiment = Experiment(
        tuner=None,
        env=env,
        agent=agent,
        flow=flow,
        name=name
    )
    experiment.run()

from baconian.core.experiment_runner import *

GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')
single_exp_runner(task_fn, del_if_log_path_existed=True)

```

## 4.3 PPO with Pendulum-v0

```

# Date: 3/30/19
# Author: Luke
# Project: baconian-internal

"""
A simple example to show how to build up an experiment with ppo training and testing_
↪on Pendulum-v0
"""

from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
from baconian.algo.value_func import MLPVValueFunc
from baconian.algo.ppo import PPO

```

(continues on next page)

(continued from previous page)

```

from baconian.algo.policy.normal_distribution_mlp import NormalDistributionMLPPolicy
from baconian.core.agent import Agent
from baconian.algo.misc import EpsilonGreedy
from baconian.core.experiment import Experiment
from baconian.core.flow.train_test_flow import create_train_test_flow
from baconian.config.global_config import GlobalConfig
from baconian.core.status import get_global_status_collect
from baconian.common.schedules import PeriodicalEventSchedule

def task_fn():
    env = make('Pendulum-v0')
    name = 'demo_exp_'
    env_spec = EnvSpec(obs_space=env.observation_space,
                       action_space=env.action_space)

    mlp_v = MLPVValueFunc(env_spec=env_spec,
                          name_scope=name + 'mlp_v',
                          name=name + 'mlp_v',
                          mlp_config=[
                              {
                                  "ACT": "RELU",
                                  "B_INIT_VALUE": 0.0,
                                  "NAME": "1",
                                  "N_UNITS": 16,
                                  "L1_NORM": 0.01,
                                  "L2_NORM": 0.01,
                                  "TYPE": "DENSE",
                                  "W_NORMAL_STDDEV": 0.03
                              },
                              {
                                  "ACT": "LINEAR",
                                  "B_INIT_VALUE": 0.0,
                                  "NAME": "OUPTUT",
                                  "N_UNITS": 1,
                                  "TYPE": "DENSE",
                                  "W_NORMAL_STDDEV": 0.03
                              }
                          ])

    policy = NormalDistributionMLPPolicy(env_spec=env_spec,
                                         name_scope=name + 'mlp_policy',
                                         name=name + 'mlp_policy',
                                         mlp_config=[
                                             {
                                                 "ACT": "RELU",
                                                 "B_INIT_VALUE": 0.0,
                                                 "NAME": "1",
                                                 "N_UNITS": 16,
                                                 "L1_NORM": 0.01,
                                                 "L2_NORM": 0.01,
                                                 "TYPE": "DENSE",
                                                 "W_NORMAL_STDDEV": 0.03
                                             },
                                             {
                                                 "ACT": "LINEAR",
                                                 "B_INIT_VALUE": 0.0,

```

(continues on next page)



(continued from previous page)

```

        "NAME": "OUPTUT",
        "N_UNITS": env_spec.flat_action_dim,
        "TYPE": "DENSE",
        "W_NORMAL_STDDEV": 0.03
    }
],
reuse=False)

ppo = PPO(
    env_spec=env_spec,
    config_or_config_dict={
        "gamma": 0.995,
        "lam": 0.98,
        "policy_train_iter": 10,
        "value_func_train_iter": 10,
        "clipping_range": None,
        "beta": 1.0,
        "eta": 50,
        "value_func_memory_size": 10,
        "log_var_init": -1.0,
        "kl_target": 0.003,
        "policy_lr": 0.01,
        "value_func_lr": 0.01,
        "value_func_train_batch_size": 10,
        "lr_multiplier": 1.0
    },
    value_func=mlp_v,
    stochastic_policy=policy,
    name=name + 'ppo'
)
agent = Agent(env=env, env_spec=env_spec,
              algo=ppo,
              algo_scheduler=PeriodicalEventSchedule(
                  t_fn=lambda: get_global_status_collect() ('TOTAL_AGENT_TRAIN_
↪SAMPLE_COUNT'),
                  trigger_every_step=20,
                  after_t=10),
              name=name + 'agent',
              exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
↪space,
                                              init_random_prob=0.5))

flow = create_train_test_flow(
    test_every_sample_count=10,
    train_every_sample_count=10,
    start_test_after_sample_count=5,
    start_train_after_sample_count=5,
    train_func_and_args=(agent.train, (), dict()),
    test_func_and_args=(agent.test, (), dict(sample_count=10)),
    sample_func_and_args=(agent.sample, (), dict(sample_count=100,
                                                  env=agent.env,
                                                  sample_type='trajectory',
                                                  store_flag=True))
)

experiment = Experiment(
    tuner=None,
    env=env,

```

(continues on next page)

(continued from previous page)

```

        agent=agent,
        flow=flow,
        name=name
    )
    experiment.run()

from baconian.core.experiment_runner import single_exp_runner

GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')
single_exp_runner(task_fn, del_if_log_path_existed=True)

```

## 4.4 MPC with Pendulum-v0

```

"""
A simple example to show how to build up an experiment with ddpq training and testing.
↳ on MountainCarContinuous-v0
"""
from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
from baconian.core.agent import Agent
from baconian.algo.misc import EpsilonGreedy
from baconian.core.experiment import Experiment
from baconian.core.flow.train_test_flow import create_train_test_flow
from baconian.algo.mpc import ModelPredictiveControl
from baconian.algo.dynamics.terminal_func.terminal_func import RandomTerminalFunc
from baconian.algo.dynamics.reward_func.reward_func import RandomRewardFunc
from baconian.algo.policy import UniformRandomPolicy
from baconian.algo.dynamics.mlp_dynamics_model import ContinuousMLPGlobalDynamicsModel
from baconian.config.global_config import GlobalConfig

def task_fn():
    env = make('Pendulum-v0')
    name = 'demo_exp'
    env_spec = EnvSpec(obs_space=env.observation_space,
                       action_space=env.action_space)

    mlp_dyna = ContinuousMLPGlobalDynamicsModel(
        env_spec=env_spec,
        name_scope=name + '_mlp_dyna',
        name=name + '_mlp_dyna',
        learning_rate=0.01,
        mlp_config=[
            {
                "ACT": "RELU",
                "B_INIT_VALUE": 0.0,
                "NAME": "1",
                "L1_NORM": 0.0,
                "L2_NORM": 0.0,
                "N_UNITS": 16,
                "TYPE": "DENSE",
                "W_NORMAL_STDDEV": 0.03
            }
        ],

```

(continues on next page)

(continued from previous page)

```

        {
            "ACT": "LINEAR",
            "B_INIT_VALUE": 0.0,
            "NAME": "OUPTUT",
            "L1_NORM": 0.0,
            "L2_NORM": 0.0,
            "N_UNITS": env_spec.flat_obs_dim,
            "TYPE": "DENSE",
            "W_NORMAL_STDDEV": 0.03
        }
    ])
    algo = ModelPredictiveControl(
        dynamics_model=mlp_dyna,
        env_spec=env_spec,
        config_or_config_dict=dict(
            SAMPLED_HORIZON=2,
            SAMPLED_PATH_NUM=5,
            dynamics_model_train_iter=10
        ),
        name=name + '_mpc',

        policy=UniformRandomPolicy(env_spec=env_spec, name='uni_policy')
    )
    algo.set_terminal_reward_function_for_dynamics_env(reward_
↳ func=RandomRewardFunc(name='reward_func'),
                                                    terminal_
↳ func=RandomTerminalFunc(name='random_terminal'), )
    agent = Agent(env=env, env_spec=env_spec,
                  algo=algo,
                  name=name + '_agent',
                  exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
↳ space,
                                                    init_random_prob=0.5))

    flow = create_train_test_flow(
        test_every_sample_count=10,
        train_every_sample_count=10,
        start_test_after_sample_count=5,
        start_train_after_sample_count=5,
        train_func_and_args=(agent.train, (), dict()),
        test_func_and_args=(agent.test, (), dict(sample_count=10)),
        sample_func_and_args=(agent.sample, (), dict(sample_count=100,
                                                    env=agent.env,
                                                    store_flag=True))
    )
    experiment = Experiment(
        tuner=None,
        env=env,
        agent=agent,
        flow=flow,
        name=name
    )
    experiment.run()

from baconian.core.experiment_runner import single_exp_runner

GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')

```

(continues on next page)

(continued from previous page)

```
single_exp_runner(task_fn, del_if_log_path_existed=True)
```

## 4.5 Dyna with Pendulum-v0

```
"""
A simple example to show how to build up an experiment with Dyna training and testing.
↳ on Pendulum-v0
"""
from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
from baconian.algo.value_func.mlp_q_value import MLPQValueFunction
from baconian.algo.ddpg import DDPG
from baconian.algo.policy import DeterministicMLPPolicy
from baconian.core.agent import Agent
from baconian.algo.misc import EpsilonGreedy
from baconian.core.experiment import Experiment
from baconian.config.global_config import GlobalConfig
from baconian.core.status import get_global_status_collect
from baconian.common.schedules import PeriodicalEventSchedule
from baconian.algo.dynamics.mlp_dynamics_model import ContinuousMLPGlobalDynamicsModel
from baconian.algo.dyna import Dyna
from baconian.algo.dynamics.reward_func.reward_func import RandomRewardFunc
from baconian.algo.dynamics.terminal_func.terminal_func import _
↳ FixedEpisodeLengthTerminalFunc
from baconian.core.flow.dyna_flow import create_dyna_flow
from baconian.common.data_pre_processing import RunningStandardScaler

def task_fn():
    # create the gym environment by make function
    env = make('Pendulum-v0')
    # give your experiment a name which is used to generate the log path etc.
    name = 'demo_exp'
    # construct the environment specification
    env_spec = EnvSpec(obs_space=env.observation_space,
                       action_space=env.action_space)
    # construct the neural network to approximate q function of DDPG
    mlp_q = MLPQValueFunction(env_spec=env_spec,
                              name_scope=name + '_mlp_q',
                              name=name + '_mlp_q',
                              mlp_config=[
                                  {
                                      "ACT": "RELU",
                                      "B_INIT_VALUE": 0.0,
                                      "NAME": "1",
                                      "N_UNITS": 16,
                                      "TYPE": "DENSE",
                                      "W_NORMAL_STDDEV": 0.03
                                  },
                                  {
                                      "ACT": "LINEAR",
                                      "B_INIT_VALUE": 0.0,
                                      "NAME": "OUPTUT",
                                      "N_UNITS": 1,
```

(continues on next page)

(continued from previous page)

```

        "TYPE": "DENSE",
        "W_NORMAL_STDDEV": 0.03
    }
    ])

# construct the neural network to approximate policy for DDPG
policy = DeterministicMLPPolicy(env_spec=env_spec,
                                name_scope=name + '_mlp_policy',
                                name=name + '_mlp_policy',
                                mlp_config=[
                                    {
                                        "ACT": "RELU",
                                        "B_INIT_VALUE": 0.0,
                                        "NAME": "1",
                                        "N_UNITS": 16,
                                        "TYPE": "DENSE",
                                        "W_NORMAL_STDDEV": 0.03
                                    },
                                    {
                                        "ACT": "LINEAR",
                                        "B_INIT_VALUE": 0.0,
                                        "NAME": "OUPTUT",
                                        "N_UNITS": env_spec.flat_action_dim,
                                        "TYPE": "DENSE",
                                        "W_NORMAL_STDDEV": 0.03
                                    }
                                ],
                                reuse=False)

# construct the DDPG algorithms
ddpg = DDPG(
    env_spec=env_spec,
    config_or_config_dict={
        "REPLAY_BUFFER_SIZE": 10000,
        "GAMMA": 0.999,
        "CRITIC_LEARNING_RATE": 0.001,
        "ACTOR_LEARNING_RATE": 0.001,
        "DECAY": 0.5,
        "BATCH_SIZE": 50,
        "TRAIN_ITERATION": 1,
        "critic_clip_norm": 0.1,
        "actor_clip_norm": 0.1,
    },
    value_func=mlp_q,
    policy=policy,
    name=name + '_ddpg',
    replay_buffer=None
)

# construct a neural network based global dynamics model to approximate the state_
→transition of environment
mlp_dyna = ContinuousMLPGlobalDynamicsModel(
    env_spec=env_spec,
    name_scope=name + '_mlp_dyna',
    name=name + '_mlp_dyna',
    learning_rate=0.01,
    state_input_scaler=RunningStandardScaler(dims=env_spec.flat_obs_dim),
    action_input_scaler=RunningStandardScaler(dims=env_spec.flat_action_dim),
    output_delta_state_scaler=RunningStandardScaler(dims=env_spec.flat_obs_dim),
    mlp_config=[

```

(continues on next page)

(continued from previous page)

```

        {
            "ACT": "RELU",
            "B_INIT_VALUE": 0.0,
            "NAME": "l",
            "L1_NORM": 0.0,
            "L2_NORM": 0.0,
            "N_UNITS": 16,
            "TYPE": "DENSE",
            "W_NORMAL_STDDEV": 0.03
        },
        {
            "ACT": "LINEAR",
            "B_INIT_VALUE": 0.0,
            "NAME": "OUPTUT",
            "L1_NORM": 0.0,
            "L2_NORM": 0.0,
            "N_UNITS": env_spec.flat_obs_dim,
            "TYPE": "DENSE",
            "W_NORMAL_STDDEV": 0.03
        }
    ])
    # finally, construct the Dyna algorithms with a model free algorithm DDGP, and a
    ↪ NN model.
    algo = Dyna(env_spec=env_spec,
                name=name + '_dyna_algo',
                model_free_algo=ddpg,
                dynamics_model=mlp_dyna,
                config_or_config_dict=dict(
                    dynamics_model_train_iter=10,
                    model_free_algo_train_iter=10
                ))
    # To make the NN based dynamics model a proper environment so be a sampling
    ↪ source for DDPG, reward function and
    # terminal function need to be set.

    # For examples only, we use random reward function and terminal function with
    ↪ fixed episode length.
    algo.set_terminal_reward_function_for_dynamics_env(
        terminal_func=FixedEpisodeLengthTerminalFunc(max_step_length=env.unwrapped._
    ↪ max_episode_steps,
                                                    step_count_fn=algo.dynamics_env.
    ↪ total_step_count_fn),
        reward_func=RandomRewardFunc())
    # construct agent with additional exploration strategy if needed.
    agent = Agent(env=env, env_spec=env_spec,
                  algo=algo,
                  algo_saving_scheduler=PeriodicalEventSchedule(
                      t_fn=lambda: get_global_status_collect() ('TOTAL_AGENT_TRAIN_
    ↪ SAMPLE_COUNT'),
                      trigger_every_step=20,
                      after_t=10),
                  name=name + '_agent',
                  exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
    ↪ space,
                                                    init_random_prob=0.5))

    # construct the training flow, called Dyna flow. It defines how the training
    ↪ proceed, and the terminal condition

```

(continues on next page)

(continued from previous page)

```

flow = create_dyna_flow(
    train_algo_func=(agent.train, (), dict(state='state_agent_training')),
    train_algo_from_synthesized_data_func=(agent.train, (), dict(state='state_
↪agent_training')),
    train_dynamics_func=(agent.train, (), dict(state='state_dynamics_training')),
    test_algo_func=(agent.test, (), dict(sample_count=1)),
    test_dynamics_func=(agent.algo.test_dynamics, (), dict(sample_count=10,
↪env=env)),
    sample_from_real_env_func=(agent.sample, (), dict(sample_count=10,
                                                    env=agent.env,
                                                    store_flag=True)),
    sample_from_dynamics_env_func=(agent.sample, (), dict(sample_count=10,
                                                    env=agent.algo.dynamics_
↪env,
                                                    store_flag=True)),
    train_algo_every_real_sample_count_by_data_from_real_env=40,
    train_algo_every_real_sample_count_by_data_from_dynamics_env=40,
    test_algo_every_real_sample_count=40,
    test_dynamics_every_real_sample_count=40,
    train_dynamics_every_real_sample_count=20,
    start_train_algo_after_sample_count=1,
    start_train_dynamics_after_sample_count=1,
    start_test_algo_after_sample_count=1,
    start_test_dynamics_after_sample_count=1,
    warm_up_dynamics_samples=1
)
# construct the experiment
experiment = Experiment(
    tuner=None,
    env=env,
    agent=agent,
    flow=flow,
    name=name + '_exp'
)
# run!
experiment.run()

from baconian.core.experiment_runner import *

# set some global configuration here

# set DEFAULT_EXPERIMENT_END_POINT to indicate when to stop the experiment.
# one usually used is the TOTAL_AGENT_TRAIN_SAMPLE_COUNT, i.e., how many samples/
↪timesteps are used for training
GlobalConfig().set('DEFAULT_EXPERIMENT_END_POINT', dict(TOTAL_AGENT_TRAIN_SAMPLE_
↪COUNT=200))

# set the logging path to write log and save model checkpoints.
GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')

# feed the task into a exp runner.
single_exp_runner(task_fn, del_if_log_path_existed=True)

```

## 4.6 Gaussian Process Dynamics

```

"""
This gives a simple example on how to use Gaussian Process (GP) to approximate the
↳ Gym environment Pendulum-v0
We use gpflow package to build the Gaussian Process.
"""

from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
import numpy as np
from baconian.common.sampler.sample_data import TransitionData
from baconian.algo.policy import UniformRandomPolicy
from baconian.algo.dynamics.gaussian_process_dynamics_model import
↳ GaussianProcessDynamicsModel
from baconian.algo.dynamics.dynamics_model import DynamicsEnvWrapper
from baconian.algo.dynamics.terminal_func.terminal_func import RandomTerminalFunc
from baconian.algo.dynamics.reward_func.reward_func import RandomRewardFunc

env = make('Pendulum-v0')
name = 'demo_exp'
env_spec = EnvSpec(obs_space=env.observation_space,
                   action_space=env.action_space)
data = TransitionData(env_spec=env_spec)
policy = UniformRandomPolicy(env_spec=env_spec)
# Do some initial sampling here to train GP model
st = env.reset()
for i in range(100):
    ac = policy.forward(st)
    new_st, re, _, _ = env.step(ac)
    data.append(state=st, new_state=new_st, action=ac, reward=re, done=False)
    st = new_st

gp = GaussianProcessDynamicsModel(env_spec=env_spec, batch_data=data)
gp.init()
gp.train()

dyna_env = DynamicsEnvWrapper(dynamics=gp)
# Since we only care about the prediction here, so we pass the terminal function and
↳ reward function setting with
# random one
dyna_env.set_terminal_reward_func(terminal_func=RandomTerminalFunc(),
                                reward_func=RandomRewardFunc())

st = env.reset()
real_state_list = []
dynamics_state_list = []
test_sample_count = 100
for i in range(test_sample_count):
    ac = env_spec.action_space.sample()
    gp.reset_state(state=st)
    new_state_dynamics, _, _, _ = dyna_env.step(action=ac, allow_clip=True)
    new_state_real, _, done, _ = env.step(action=ac)
    real_state_list.append(new_state_real)
    dynamics_state_list.append(new_state_dynamics)
    st = new_state_real
    if done is True:
        env.reset()

```

(continues on next page)



(continued from previous page)

```
l1_loss = np.linalg.norm(np.array(real_state_list) - np.array(dynamics_state_list),
↳ord=1)
l2_loss = np.linalg.norm(np.array(real_state_list) - np.array(dynamics_state_list),
↳ord=2)
print('l1 loss is {}, l2 loss is {}'.format(l1_loss, l2_loss))
```

## 4.7 Early Stopping Flow (DDPG with Pendulum-v0)

```
"""
This script show the example for adding an early stopping feature so when the agent_
↳can't increase its received average
reward for evaluation, the experiment will end early.

To do so in a extensible and modular way. We can implement a new flow called_
↳EarlyStoppingFlow that implement a special
ending condition detections by accessing the agent's evaluation reward (with built-in_
↳modular to access). Such mechanism
can be re-used by all algorithms, which avoid the redundant coding for users.
"""
from baconian.config.dict_config import DictConfig
from baconian.core.flow.train_test_flow import TrainTestFlow

from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
from baconian.algo.value_func.mlp_q_value import MLPQValueFunction
from baconian.algo.ddpg import DDPG
from baconian.algo.policy import DeterministicMLPPolicy
from baconian.core.agent import Agent
from baconian.algo.misc import EpsilonGreedy
from baconian.core.experiment import Experiment
from baconian.core.status import get_global_status_collect
from baconian.common.schedules import PeriodicalEventSchedule

class EarlyStoppingFlow(TrainTestFlow):
    required_key_dict = {
        **TrainTestFlow.required_key_dict,
        'USE_LAST_K_EVALUATION_REWARD': 10
    }

    def __init__(self, train_sample_count_func, config_or_config_dict: (DictConfig,
↳dict), func_dict: dict, agent):
        super().__init__(train_sample_count_func, config_or_config_dict, func_dict)
        self.agent = agent

    def _is_ended(self):
        test_reward = sorted(self.agent.recorder.get_log(attr_name='sum_reward',
↳filter_by_status=dict(status='TEST')),
        key=lambda x: x['sample_counter'])
        if len(test_reward) >= self.parameters('USE_LAST_K_EVALUATION_REWARD') * 2:
            last_reward = test_reward[-self.parameters('USE_LAST_K_EVALUATION_REWARD'
↳')]
            pre_reward = test_reward[-self.parameters('USE_LAST_K_EVALUATION_REWARD')
↳* 2: -self.parameters(
```

(continues on next page)

(continued from previous page)

```

        'USE_LAST_K_EVALUATION_REWARD'])
    last_reward = np.mean([r['value'] for r in last_reward])
    pre_reward = np.mean([r['value'] for r in pre_reward])
    if last_reward < pre_reward:
        ConsoleLogger().print('info', 'training ended because last {} step_
↪reward: {} < previous {} step reward {}'.format(self.parameters('USE_LAST_K_
↪EVALUATION_REWARD'), last_reward, self.parameters('USE_LAST_K_EVALUATION_REWARD'),
↪pre_reward))
        return True
    return super()._is_ended()

def create_early_stopping_flow(test_every_sample_count, train_every_sample_count,
↪start_train_after_sample_count,
                                start_test_after_sample_count, train_func_and_args,
↪test_func_and_args,
                                sample_func_and_args,
                                agent,
                                use_last_k_evaluation_reward,
                                train_samples_counter_func=None):
    config_dict = dict(
        TEST_EVERY_SAMPLE_COUNT=test_every_sample_count,
        TRAIN_EVERY_SAMPLE_COUNT=train_every_sample_count,
        START_TRAIN_AFTER_SAMPLE_COUNT=start_train_after_sample_count,
        START_TEST_AFTER_SAMPLE_COUNT=start_test_after_sample_count,
        USE_LAST_K_EVALUATION_REWARD=use_last_k_evaluation_reward
    )

    def return_func_dict(s_dict):
        return dict(func=s_dict[0],
                    args=s_dict[1],
                    kwargs=s_dict[2])

    func_dict = dict(
        train=return_func_dict(train_func_and_args),
        test=return_func_dict(test_func_and_args),
        sample=return_func_dict(sample_func_and_args),
    )
    if train_samples_counter_func is None:
        def default_train_samples_counter_func():
            return get_global_status_collect()('TOTAL_AGENT_TRAIN_SAMPLE_COUNT')

        train_samples_counter_func = default_train_samples_counter_func

    return EarlyStoppingFlow(config_or_config_dict=config_dict,
                             train_sample_count_func=train_samples_counter_func,
                             agent=agent,
                             func_dict=func_dict)

def task_fn():
    env = make('Pendulum-v0')
    name = 'demo_exp'
    env_spec = EnvSpec(obs_space=env.observation_space,
                       action_space=env.action_space)

    mlp_q = MLPQValueFunction(env_spec=env_spec,

```

(continues on next page)

(continued from previous page)

```

        name_scope=name + '_mlp_q',
        name=name + '_mlp_q',
        mlp_config=[
            {
                "ACT": "RELU",
                "B_INIT_VALUE": 0.0,
                "NAME": "l",
                "N_UNITS": 16,
                "TYPE": "DENSE",
                "W_NORMAL_STDDEV": 0.03
            },
            {
                "ACT": "LINEAR",
                "B_INIT_VALUE": 0.0,
                "NAME": "OUPTUT",
                "N_UNITS": 1,
                "TYPE": "DENSE",
                "W_NORMAL_STDDEV": 0.03
            }
        ])
    policy = DeterministicMLPPolicy(env_spec=env_spec,
                                   name_scope=name + '_mlp_policy',
                                   name=name + '_mlp_policy',
                                   mlp_config=[
                                       {
                                           "ACT": "RELU",
                                           "B_INIT_VALUE": 0.0,
                                           "NAME": "l",
                                           "N_UNITS": 16,
                                           "TYPE": "DENSE",
                                           "W_NORMAL_STDDEV": 0.03
                                       },
                                       {
                                           "ACT": "LINEAR",
                                           "B_INIT_VALUE": 0.0,
                                           "NAME": "OUPTUT",
                                           "N_UNITS": env_spec.flat_action_dim,
                                           "TYPE": "DENSE",
                                           "W_NORMAL_STDDEV": 0.03
                                       }
                                   ],
                                   reuse=False)

    ddpq = DDPG(
        env_spec=env_spec,
        config_or_config_dict={
            "REPLAY_BUFFER_SIZE": 10000,
            "GAMMA": 0.999,
            "CRITIC_LEARNING_RATE": 0.001,
            "ACTOR_LEARNING_RATE": 0.001,
            "DECAY": 0.5,
            "BATCH_SIZE": 50,
            "TRAIN_ITERATION": 1,
            "critic_clip_norm": 0.1,
            "actor_clip_norm": 0.1,
        },
        value_func=mlp_q,

```

(continues on next page)

(continued from previous page)

```

        policy=policy,
        name=name + '_ddpg',
        replay_buffer=None
    )
    agent = Agent(env=env, env_spec=env_spec,
                  algo=ddpg,
                  algo_saving_scheduler=PeriodicalEventSchedule(
                      t_fn=lambda: get_global_status_collect() ('TOTAL_AGENT_TRAIN_
↪SAMPLE_COUNT'),
                      trigger_every_step=20,
                      after_t=10),
                  name=name + '_agent',
                  exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
↪space,
                                                    init_random_prob=0.5))

    flow = create_early_stopping_flow(
        agent=agent,
        use_last_k_evaluation_reward=5,
        test_every_sample_count=10,
        train_every_sample_count=10,
        start_test_after_sample_count=5,
        start_train_after_sample_count=5,
        train_func_and_args=(agent.train, (), dict()),
        test_func_and_args=(agent.test, (), dict(sample_count=1)),
        sample_func_and_args=(agent.sample, (), dict(sample_count=100,
                                                    env=agent.env))
    )

    experiment = Experiment(
        tuner=None,
        env=env,
        agent=agent,
        flow=flow,
        name=name
    )
    experiment.run()

from baconian.core.experiment_runner import *

GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')
GlobalConfig().set('DEFAULT_EXPERIMENT_END_POINT',
                   dict(TOTAL_AGENT_TRAIN_SAMPLE_COUNT=2000,
                         TOTAL_AGENT_TEST_SAMPLE_COUNT=None,
                         TOTAL_AGENT_UPDATE_COUNT=None))
single_exp_runner(task_fn, del_if_log_path_existed=True)

```

## 4.8 Environment Wrapper (MountainCarContinuous-v0, RewardWrapper)

## 4.9 Use scheduler module in experiments

```

"""
In this example, we demonstrate how to utilize the scheduler module to dynamically
↪ setting the
learning rate of your algorithm, or epsilon-greedy probability
"""

from baconian.algo.dqn import DQN
from baconian.core.core import EnvSpec
from baconian.envs.gym_env import make
from baconian.algo.value_func.mlp_q_value import MLPQValueFunction
from baconian.core.agent import Agent
from baconian.algo.misc import EpsilonGreedy
from baconian.core.experiment import Experiment
from baconian.core.flow.train_test_flow import create_train_test_flow
from baconian.config.global_config import GlobalConfig
from baconian.common.schedules import LinearScheduler, PiecewiseScheduler,
↪ PeriodicalEventSchedule
from baconian.core.status import get_global_status_collect

def task_fn():
    env = make('Acrobot-v1')
    name = 'example_scheduler_'
    env_spec = EnvSpec(obs_space=env.observation_space,
                       action_space=env.action_space)

    mlp_q = MLPQValueFunction(env_spec=env_spec,
                             name_scope=name + '_mlp_q',
                             name=name + '_mlp_q',
                             mlp_config=[
                                 {
                                     "ACT": "RELU",
                                     "B_INIT_VALUE": 0.0,
                                     "NAME": "1",
                                     "N_UNITS": 16,
                                     "TYPE": "DENSE",
                                     "W_NORMAL_STDDEV": 0.03
                                 },
                                 {
                                     "ACT": "LINEAR",
                                     "B_INIT_VALUE": 0.0,
                                     "NAME": "OUPTUT",
                                     "N_UNITS": 1,
                                     "TYPE": "DENSE",
                                     "W_NORMAL_STDDEV": 0.03
                                 }
                             ])

    dqn = DQN(env_spec=env_spec,
              config_or_config_dict=dict(REPLAY_BUFFER_SIZE=1000,
                                         GAMMA=0.99,

```

(continues on next page)

(continued from previous page)

```

                                BATCH_SIZE=10,
                                LEARNING_RATE=0.001,
                                TRAIN_ITERATION=1,
                                DECAY=0.5),

    name=name + '_dqn',
    value_func=mlp_q)
agent = Agent(env=env, env_spec=env_spec,
              algo=dqn,
              name=name + '_agent',
              algo_saving_scheduler=PeriodicalEventSchedule(
                  t_fn=lambda: get_global_status_collect() ('TOTAL_AGENT_TRAIN_
↪SAMPLE_COUNT'),
                  trigger_every_step=20,
                  after_t=10),
              exploration_strategy=EpsilonGreedy(action_space=env_spec.action_
↪space,
                                              prob_
↪scheduler=PiecewiseScheduler(
                                              t_fn=lambda: get_global_
↪status_collect() (
                                              'TOTAL_AGENT_TRAIN_
↪SAMPLE_COUNT'),
                                              endpoints=((10, 0.3), (100,
↪0.1), (200, 0.0)),
                                              outside_value=0.0
                                              ),
                                              init_random_prob=0.5))

    flow = create_train_test_flow(
        test_every_sample_count=10,
        train_every_sample_count=10,
        start_test_after_sample_count=5,
        start_train_after_sample_count=5,
        train_func_and_args=(agent.train, (), dict()),
        test_func_and_args=(agent.test, (), dict(sample_count=10)),
        sample_func_and_args=(agent.sample, (), dict(sample_count=100,
                                                    env=agent.env,
                                                    store_flag=True))

    )
    experiment = Experiment(
        tuner=None,
        env=env,
        agent=agent,
        flow=flow,
        name=name + 'experiment_debug'
    )

    dqn.parameters.set_scheduler(param_key='LEARNING_RATE',
                                scheduler=LinearScheduler(
                                    t_fn=experiment.TOTAL_AGENT_TRAIN_SAMPLE_COUNT,
                                    schedule_timesteps=GlobalConfig().DEFAULT_
↪EXPERIMENT_END_POINT[
                                    'TOTAL_AGENT_TRAIN_SAMPLE_COUNT'],
                                    final_p=0.0001,
                                    initial_p=0.01))

    experiment.run()

```

(continues on next page)

(continued from previous page)

```
from baconian.core.experiment_runner import single_exp_runner

GlobalConfig().set('DEFAULT_LOG_PATH', './log_path')
single_exp_runner(task_fn, del_if_log_path_existed=True)
```

## 4.10 Visualization of experiments log

Below is the code of Baconian visualisation module. For its usage, please refer to :doc:`Logging and Visualization <how\_to\_log>`.





---

## Logging and Visualization

---

In this part, we will introduce how to record and save everything you want during an experiments. The contents are organized as follows:

- How to understand the log file after the experiment finished.
- How to visualize with log file after experiment finished.
- How the logging module of Baconian works.

### 5.1 Log Directory Explanation

After one experiment finished, lots of built-in logging information will be recorded and saved to file so you can visualize and analysis the results. A typical logging file directory looks like this:

```
.
├── console.log           # log output to console during experiments
├── final_status.json     # final status of all modules when the experiment
└─ finished
├── global_config.json   # config you use for this experiment include default
└─ config and customized config.
├── record               # record folder contains all raw log of experiments.
├── └── agent            # agent level log (training samples' reward, test samples
└─ ' reward etc.)
    ├── └── CREATED.json  # agent's log under CREATED status
    ├── └── INITED.json   # agent's log after CREATED and into INITED status
    ├── └── TRAIN.json    # agent's log under training status (e.g., reward
└─ received when sample for training, action during training sampling)
    ├── └── TEST.json     # agent's log under test status (e.g., reward received
└─ when sample for evaluation/test)
    ├── └── algo          # algorithm level log (loss, KL divergence, etc.)
    ├── └── CREATED.json  # algorithm's log under CREATED status
    ├── └── INITED.json   # algorithm's log after CREATED and into INITED status
    ├── └── TRAIN.json    # algorithm's log under training status (e.g., training
└─ loss)
```

(continues on next page)

(continued from previous page)

```
└─┬─ TEST.json          # algorithm's log under test status
└─ model_checkpoints     # store the checkpoints of tensorflow models if you have_
→ add the saving scheduler in your experiments.
```

## 5.2 Visualize and analyze the log

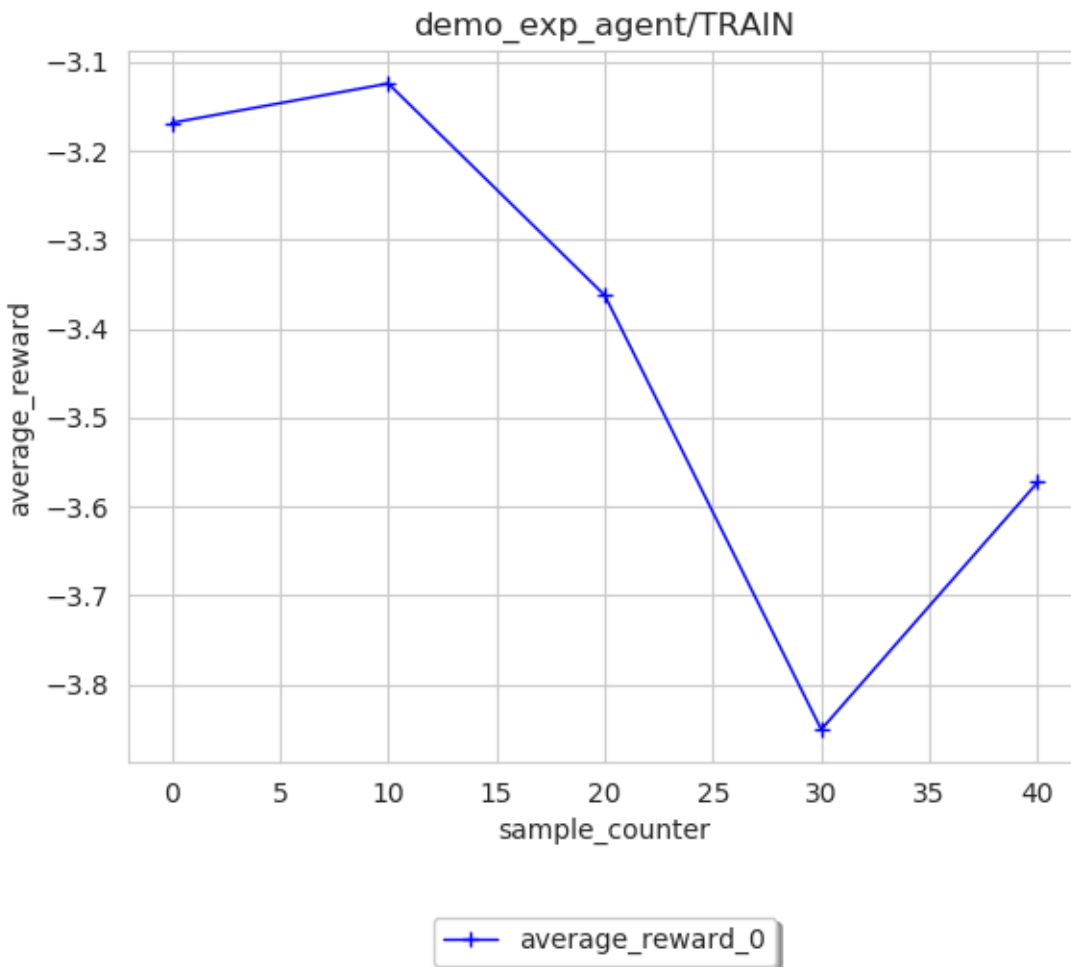
Three plot modes: `line`, `scatter` and `histogram` are available within `plot_res` function in `log_data_loader`. `plotter` will use `Matplotlib` package to plot a graph with `data_new` `DataFrame` as input. Users may utilise `visualisation.py` for fast plotting, with existing single or multiple experiment results.

Please refer to the following examples for single or multiple experiment visualisation:

- Single Experiment Visualisation

We use the log generated by running the *Dyna examples* as example.

Following code snippet is to draw a `line plot` of `average_reward` that agent received during training, using data sampled from environment as index and averaged every 10 rewards for smooth line.

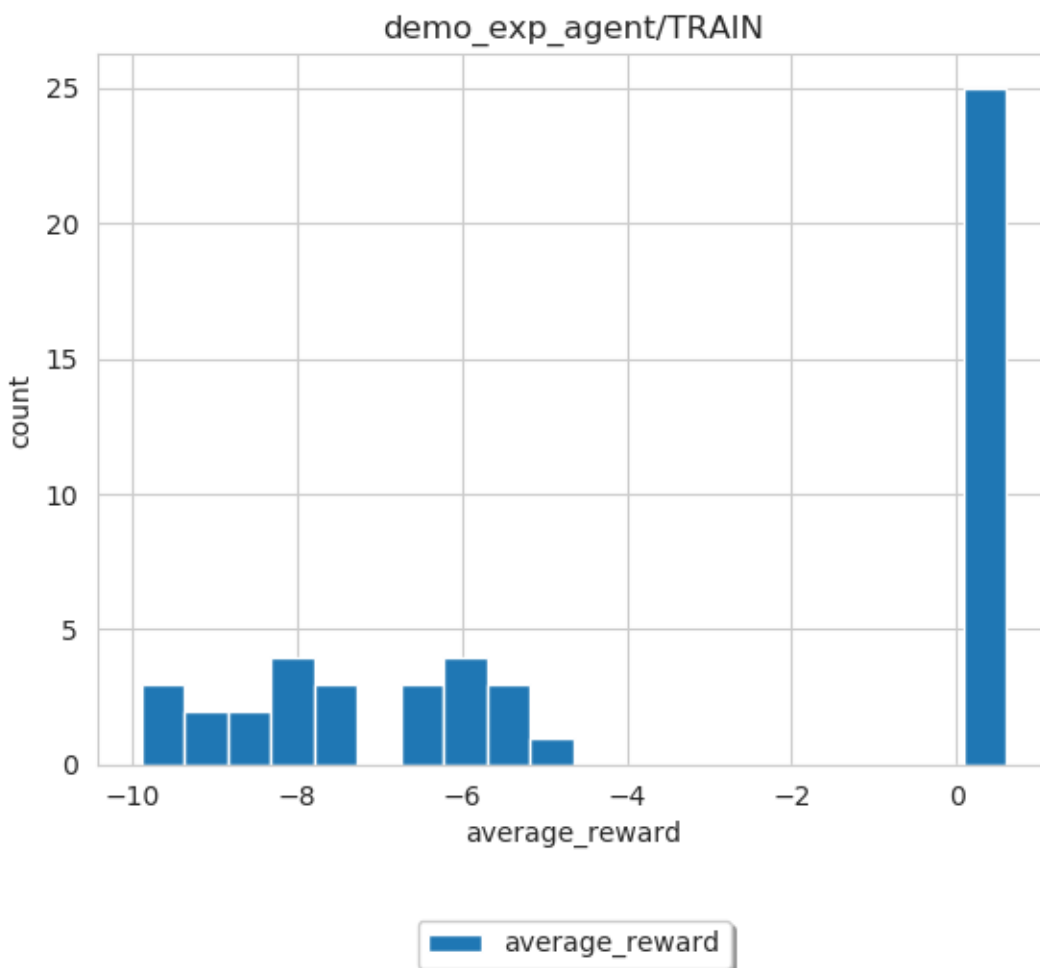


**Note:** `sub_log_dir_name` should include the COMPLETE directory in between the `log_path` directory and

`json.log`. For example, if you have a log folder structured as `/path/to/log/record/demo_exp_agent/TEST/log.json`, then the `sub_log_dir_name` should be `demo_exp_agent/TEST/` and `exp_root_dir` should be `/path/to/log/`.

Please note that histogram plot mode is a bit different from the other two modes, in terms of data manipulation. Although `index` is unnecessary under histogram mode, but currently user still should pass in one for internal data processing.

```
image = loader.SingleExpLogDataLoader(path)
image.plot_res(sub_log_dir_name='demo_exp_agent/TRAIN',
               key="average_reward",
               index='sample_counter',
               mode='histogram',
               save_format='pdf',
               save_flag=True,
               file_name='average_reward_histogram'
               )
```



- Multiple Experiment Visualisation

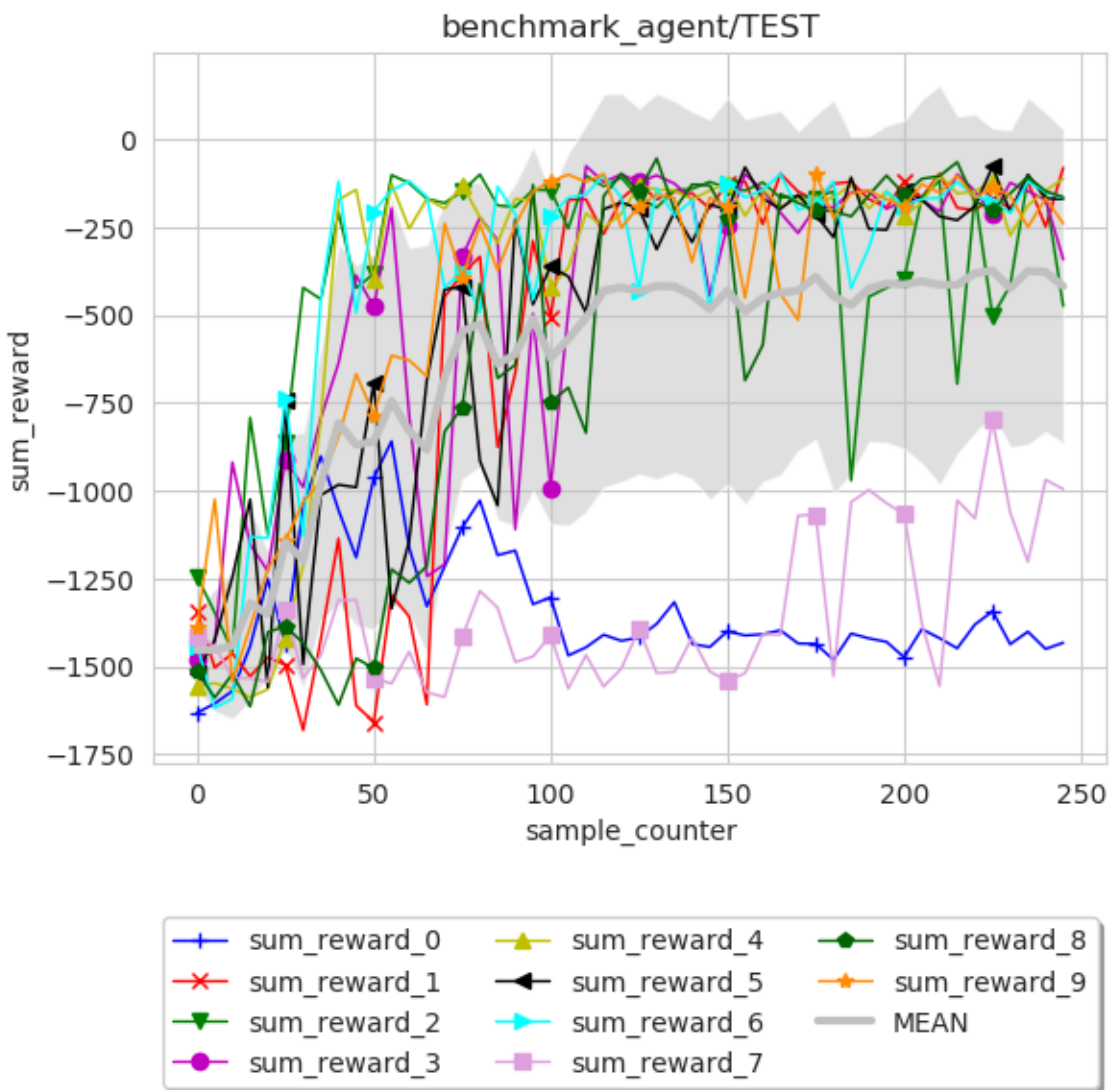
Visualize the results from multiple runs can give a more reliable analysis of the RL methods, by plotting the mean and variance over these results. Such can be done by `MultipleExpLogDataLoader`

We use the DDPG benchmark experiments as example, use can found the script under the source code

baconian-project/baconian/benchmark/run\_benchmark.py

Following code snippet is to draw a line plot of sum\_reward in benchmark\_agent/TEST as a result of 10 times of DDPG benchmark experiments.

```
path = '/path/to/log' # under the path, there should be 10 sub folders, each contains_
↳ 1 experiment results.
image = loader.MultipleExpLogDataLoader(path)
image.plot_res(sub_log_dir_name='benchmark_agent/TEST',
               key="sum_reward",
               index='sample_counter',
               mode='line',
               save_flag=True,
               average_over=10,
               save_format='pdf'
               )
```



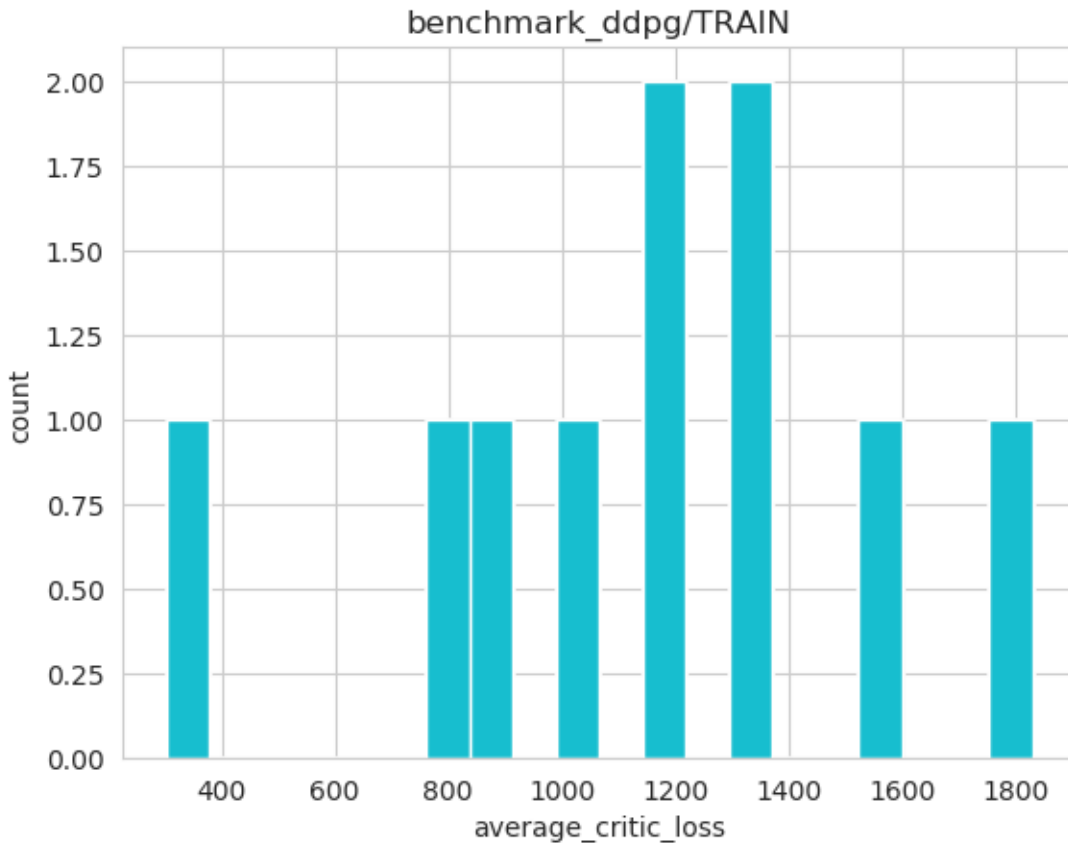
We can see from the results that DDPG is not quite stable as 2 out of 10 runs failed to converge.

When plotting multiple experiment results in histogram mode, figure will reflect the histogram/data distribution using all experiments' data.

```

path = '/path/to/log'
image = loader.MultipleExpLogDataLoader(path)
image.plot_res(sub_log_dir_name='benchmark_ddpg/TRAIN',
               key="average_critic_loss",
               index='train',
               mode='histogram',
               file_name='average_critic_loss_benchmark',
               save_format='pdf',
               save_flag=True,
               )

```



We can use the action distribution to analyze and diagnose algorithms.

### 5.3 How the logging module of Baconian works

There are two important modules of Baconian: `Logger` and `Recorder`, `Recorder` is coupled with every module or class you want to record something during training or testing, for such as DQN, Agent or Environment. It will record the information like loss, gradient or reward in a way that you specified. While `Logger` will take charge of these recorded information, group them in a certain way and output them into file, console etc.



---

## How to implement a new algorithm

---

In this section, we will walk through the implementation of Deep Deterministic Policy Gradient (DDPG) algorithm, available at `baconian/algo/ddpg.py`. It utilizes many functionalities provided by the framework, which we will describe below.

- The `ModelFreeAlgo` and `OffPolicyAlgo` Classes

For the algorithms in Baconian project, we have written many abstract classes to indicate the characteristics of the algorithm, in `baconian/algo/rl_algo.py`. The DDPG class inherits from `ModelFreeAlgo` and `OffPolicyAlgo` classes' `ModelFreeAlgo`, `OffPolicyAlgo` and other classes in `baconian/algo/rl_algo.py` inherit `Algo` class to categorize DRL algorithms.

```
from baconian.algo.algo import Algo
from baconian.algo.dynamics.dynamics_model import DynamicsModel
from baconian.core.core import EnvSpec
from baconian.common.logging import record_return_decorator
import numpy as np

class ModelFreeAlgo(Algo):
    def __init__(self, env_spec: EnvSpec, name: str = 'model_free_algo', warm_up_
↳trajectories_number=0):
        super(ModelFreeAlgo, self).__init__(env_spec, name, warm_up_trajectories_
↳number)

class OnPolicyAlgo(Algo):
    pass

class OffPolicyAlgo(Algo):
    pass

class ValueBasedAlgo(Algo):
    pass
```

(continues on next page)

(continued from previous page)

```

class PolicyBasedAlgo(Algo):
    pass

class ModelBasedAlgo(Algo):
    def __init__(self, env_spec, dynamics_model: DynamicsModel, name: str = 'model_
↳based_algo'):
        super(ModelBasedAlgo, self).__init__(env_spec, name)
        self._dynamics_model = dynamics_model
        self.dynamics_env = self._dynamics_model.return_as_env()

    def train_dynamics(self, *args, **kwargs):
        pass

    @record_return_decorator(which_recorder='self')
    def test_dynamics(self, env, sample_count, *args, **kwargs):
        self.set_status('TEST')
        env.set_status('TEST')
        st = env.reset()
        real_state_list = []
        dyanmics_state_list = []
        for i in range(sample_count):
            ac = self.env_spec.action_space.sample()
            self._dynamics_model.reset_state(state=st)
            new_state_dynamics, _, _ = self.dynamics_env.step(action=ac, )
            new_state_real, _, done, _ = env.step(action=ac)
            real_state_list.append(new_state_real)
            dyanmics_state_list.append(new_state_dynamics)
            st = new_state_real
            if done is True:
                env.reset()
        l1_loss = np.linalg.norm(np.array(real_state_list) - np.array(dyanmics_state_
↳list), ord=1)
        l2_loss = np.linalg.norm(np.array(real_state_list) - np.array(dyanmics_state_
↳list), ord=2)
        return dict(dynamics_test_l1_error=l1_loss, dynamics_test_l2_error=l2_loss)

    def set_terminal_reward_function_for_dynamics_env(self, terminal_func, reward_
↳func):
        self.dynamics_env.set_terminal_reward_func(terminal_func, reward_func)

```

Each new algorithm should implement the methods and attributes defined in Algo class (baconian/algo/algorithm.py).

```

from baconian.core.core import Basic, EnvSpec, Env
from baconian.core.status import StatusWithSubInfo
import abc
from typeguard import typechecked
from baconian.common.logging import Recorder
from baconian.core.parameters import Parameters
from baconian.common.sampler.sample_data import TrajectoryData

class Algo(Basic):

```

(continues on next page)



(continued from previous page)

```

"""
Abstract class for algorithms
"""
STATUS_LIST = ['CREATED', 'INITED', 'TRAIN', 'TEST']
INIT_STATUS = 'CREATED'

@typechecked
def __init__(self, env_spec: EnvSpec, name: str = 'algo', warm_up_trajectories_
↪number=0):
    """
    Constructor

    :param env_spec: environment specifications
    :type env_spec: EnvSpec
    :param name: name of the algorithm
    :type name: str
    :param warm_up_trajectories_number: how many trajectories used to warm up the_
↪training
    :type warm_up_trajectories_number: int
    """

    super().__init__(status=StatusWithSubInfo(obj=self), name=name)
    self.env_spec = env_spec
    self.parameters = Parameters(dict())
    self.recorder = Recorder(default_obj=self)
    self.warm_up_trajectories_number = warm_up_trajectories_number

def init(self):
    """
    Initialization method, such as network random initialization in Tensorflow

    :return:
    """
    self._status.set_status('INITED')

def warm_up(self, trajectory_data: TrajectoryData):
    """
    Use some data to warm up the algorithm, e.g., compute the mean/std-dev of the_
↪state to perform normalization.
    Data used in warm up process will not be added into the memory
    :param trajectory_data: TrajectoryData object
    :type trajectory_data: TrajectoryData

    :return: None
    """
    pass

def train(self, *arg, **kwargs) -> dict:
    """
    Training API, specific arguments should be defined by each algorithms itself.

    :return: training results, e.g., loss
    :rtype: dict
    """

    self._status.set_status('TRAIN')
    return dict()

```

(continues on next page)

(continued from previous page)

```

def test(self, *arg, **kwargs) -> dict:
    """
    Testing API, most of the evaluation can be done by agent instead of
    ↪ algorithms, so this API can be skipped

    :return: test results, e.g., rewards
    :rtype: dict
    """

    self._status.set_status('TEST')
    return dict()

@abc.abstractmethod
def predict(self, *arg, **kwargs):
    """
    Predict function, given the obs as input, return the action, obs will be read
    ↪ as the first argument passed into
    this API, like algo.predict(obs=x, ...)

    :return: predicted action
    :rtype: np.ndarray
    """
    raise NotImplementedError

@abc.abstractmethod
def append_to_memory(self, *args, **kwargs):
    """
    For off-policy algorithm, use this API to append the data into replay buffer.
    ↪ samples will be read as the first
    argument passed into this API, like algo.append_to_memory(samples=x, ...)

    """
    raise NotImplementedError

@property
def is_training(self):
    """
    A boolean indicate the if the algorithm is in training status

    :return: True if in training
    :rtype: bool
    """
    return self.get_status()['status'] == 'TRAIN'

@property
def is_testing(self):
    """
    A boolean indicate the if the algorithm is in training status

    :return: True if in testing
    :rtype: bool
    """
    return self.get_status()['status'] == 'TEST'

```

- The MultiPlaceholderInput Class

The algorithms in Baconian project are mostly implemented with TensorFlow, similar in the process of saving and

loading the parameters. Hence, parameters are stored in the format of TensorFlow variables by PlaceholderInput and MultiPlaceholderInput classes.

```
class DDPG(ModelFreeAlgo, OffPolicyAlgo, MultiPlaceholderInput):

    # ...

    @record_return_decorator(which_recorder='self')
    def save(self, global_step, save_path=None, name=None, **kwargs):
        save_path = save_path if save_path else GlobalConfig().DEFAULT_MODEL_
        ↪CHECKPOINT_PATH
        name = name if name else self.name
        MultiPlaceholderInput.save(self, save_path=save_path, global_step=global_step,
        ↪name=name, **kwargs)
        return dict(check_point_save_path=save_path, check_point_save_global_
        ↪step=global_step,
                    check_point_save_name=name)

    @record_return_decorator(which_recorder='self')
    def load(self, path_to_model, model_name, global_step=None, **kwargs):
        MultiPlaceholderInput.load(self, path_to_model, model_name, global_step,
        ↪**kwargs)
        return dict(check_point_load_path=path_to_model, check_point_load_global_
        ↪step=global_step,
                    check_point_load_name=model_name)
```

- Constructor

```
class DDPG(ModelFreeAlgo, OffPolicyAlgo, MultiPlaceholderInput):
    required_key_dict = DictConfig.load_json(file_path=GlobalConfig().DEFAULT_DDPG_
    ↪REQUIRED_KEY_LIST)

    @typechecked()
    def __init__(self,
                  env_spec: EnvSpec,
                  config_or_config_dict: (DictConfig, dict),
                  value_func: MLPQValueFunction,
                  policy: DeterministicMLPPolicy,
                  schedule_param_list=None,
                  name='ddpg',
                  replay_buffer=None):

        """
        :param env_spec: environment specifications, like action apace or observation_
        ↪space
        :param config_or_config_dict: configuraion dictionary, like learning rate or_
        ↪decay, if any
        :param value_func: value function
        :param policy: agent policy
        :param schedule_param_list:
        :param name: name of algorithm class instance
        :param replay_buffer: replay buffer, if any
        """
        ModelFreeAlgo.__init__(self, env_spec=env_spec, name=name)
        config = construct_dict_config(config_or_config_dict, self)

        self.config = config
        self.actor = policy
```

(continues on next page)

(continued from previous page)

```

        self.target_actor = self.actor.make_copy(name_scope='{}_target_actor'.
↪format(self.name),
                                                    name='{}_target_actor'.format(self.
↪name),
                                                    reuse=False)

        self.critic = value_func
        self.target_critic = self.critic.make_copy(name_scope='{}_target_critic'.
↪format(self.name),
                                                    name='{}_target_critic'.
↪format(self.name),
                                                    reuse=False)

        self.state_input = self.actor.state_input

        if replay_buffer:
            assert issubclass(replay_buffer, BaseReplayBuffer)
            self.replay_buffer = replay_buffer
        else:
            self.replay_buffer = UniformRandomReplayBuffer(limit=self.config('REPLAY_
↪BUFFER_SIZE'),
                                                    action_shape=self.env_spec.
↪action_shape,
                                                    observation_shape=self.env_
↪spec.obs_shape)

        self.parameters = ParametersWithTensorflowVariable(tf_var_list=[],
                                                            rest_parameters=dict(),
                                                            to_scheduler_param_
↪tuple=schedule_param_list,
                                                            name='ddpg_param',
                                                            source_config=config,
                                                            require_snapshot=False)

        """
        self.parameters contains all the parameters (variables) of the algorithm
        """
        self._critic_with_actor_output = self.critic.make_copy(reuse=True,
                                                                name='actor_input_{}'.
↪format(self.critic.name),
                                                                state_input=self.state_
↪input,
                                                                action_input=self.
↪actor.action_tensor)
        self._target_critic_with_target_actor_output = self.target_critic.make_
↪copy(reuse=True,
                                                                name='target_critic_with_target_actor_output_{}'.format(
↪self.critic.name),
                                                                action_input=self.target_actor.action_tensor)

        with tf.variable_scope(name):
            self.reward_input = tf.placeholder(shape=[None, 1], dtype=tf.float32)
            self.next_state_input = tf.placeholder(shape=[None, self.env_spec.flat_
↪obs_dim], dtype=tf.float32)
            self.done_input = tf.placeholder(shape=[None, 1], dtype=tf.bool)
            self.target_q_input = tf.placeholder(shape=[None, 1], dtype=tf.float32)

```

(continues on next page)

(continued from previous page)

```

        done = tf.cast(self.done_input, dtype=tf.float32)
        self.predict_q_value = (1. - done) * self.config('GAMMA') * self.target_q_
↪ input + self.reward_input
        with tf.variable_scope('train'):
            self.critic_loss, self.critic_update_op, self.target_critic_update_op,
↪ self.critic_optimizer, \
                self.critic_grads = self._setup_critic_loss()
            self.actor_loss, self.actor_update_op, self.target_actor_update_op,
↪ self.action_optimizer, \
                self.actor_grads = self._set_up_actor_loss()

        var_list = get_tf_collection_var_list(
            '{}/train'.format(name)) + self.critic_optimizer.variables() + self.
↪ action_optimizer.variables()
        self.parameters.set_tf_var_list(tf_var_list=sorted(list(set(var_list))),
↪ key=lambda x: x.name))
        MultiPlaceholderInput.__init__(self,
                                       sub_placeholder_input_list=[dict(obj=self.
↪ target_actor,
                                       attr_name=
↪ 'target_actor',
                                       ),
                                       dict(obj=self.
↪ actor,
                                       attr_name=
↪ 'actor'),
                                       dict(obj=self.
↪ critic,
                                       attr_name=
↪ 'critic'),
                                       dict(obj=self.
↪ target_critic,
                                       attr_name=
↪ 'target_critic')
                                       ],
                                       parameters=self.parameters)

```



---

## How to implement a new environment

---

Similar to algorithms, environments in Baconian project also should implement the methods and attributes defined in Env class `baconian/core/core.py`, inheriting gym Env class.

```
class Env(gym.Env, Basic):
    """
    Abstract class for environment
    """
    key_list = ()
    STATUS_LIST = ('JUST_RESET', 'JUST_INITED', 'TRAIN', 'TEST', 'NOT_INIT')
    INIT_STATUS = 'NOT_INIT'

    @typechecked
    def __init__(self, name: str = 'env'):
        super(Env, self).__init__(status=StatusWithSubInfo(obj=self), name=name)
        self.action_space = None
        self.observation_space = None
        self.step_count = None
        self.recorder = Recorder()
        self._last_reset_point = 0
        self.total_step_count_fn = lambda: self._status.group_specific_info_key(info_
↪key='step', group_way='sum')

        @register_counter_info_to_status_decorator(increment=1, info_key='step', under_
↪status=('TRAIN', 'TEST'),
                                                    ignore_wrong_status=True)

        def step(self, action):
            pass

        @register_counter_info_to_status_decorator(increment=1, info_key='reset', under_
↪status='JUST_RESET')
        def reset(self):
            self._status.set_status('JUST_RESET')
            self._last_reset_point = self.total_step_count_fn()
```

(continues on next page)

(continued from previous page)

```

@register_counter_info_to_status_decorator(increment=1, info_key='init', under_
↪ status='JUST_INITED')
def init(self):
    self._status.set_status('JUST_INITED')

def get_state(self):
    raise NotImplementedError

def seed(self, seed=None):
    return self.unwrapped.seed(seed=seed)

```

We use STATUS to record and control the status of an environment, register\_counter\_info\_to\_status\_decorator is a decorator that counts the times of initialization and reset of an environment.

```

def register_counter_info_to_status_decorator(increment, info_key, under_status: (str,
↪ tuple) = None,
                                           ignore_wrong_status=False):

    def wrap(fn):
        if under_status:
            assert isinstance(under_status, (str, tuple))
            if isinstance(under_status, str):
                final_st = tuple([under_status])
            else:
                final_st = under_status

        else:
            final_st = (None,)

        @wraps(fn)
        def wrap_with_self(self, *args, **kwargs):
            # todo record() called in fn will lost the just appended info_key at the_
↪ very first
            obj = self
            if not hasattr(obj, '_status') or not isinstance(getattr(obj, '_status'),
↪ StatusWithInfo):
                raise ValueError(
                    ' the object {} does not not have attribute StatusWithInfo_
↪ instance or hold wrong type of Status'.format(
                        obj))

            assert isinstance(getattr(obj, '_status'), StatusWithInfo)
            obj_status = getattr(obj, '_status')
            for st in final_st:
                obj_status.append_new_info(info_key=info_key, init_value=0, under_
↪ status=st)
            res = fn(self, *args, **kwargs)
            for st in final_st:
                if st and st != obj.get_status()['status'] and not ignore_wrong_
↪ status:
                    raise ValueError('register counter info under status: {} but got_
↪ status {}'.format(st,
                        obj.get_status()['status']))

```

(continues on next page)



(continued from previous page)

```

        obj_status.update_info(info_key=info_key, increment=increment,
                               under_status=obj.get_status()['status'])

        return res

    return wrap_with_self

return wrap

```

The class `EnvSpec` stores and regulates the environment specifications, e.g. data type of observation space and action space in an environment.

```

class EnvSpec(object):
    @init_func_arg_record_decorator()
    @typechecked
    def __init__(self, obs_space: Space, action_space: Space):
        self._obs_space = obs_space
        self._action_space = action_space
        self.obs_shape = tuple(np.array(self.obs_space.sample()).shape)
        if len(self.obs_shape) == 0:
            self.obs_shape = (1,)
        self.action_shape = tuple(np.array(self.action_space.sample()).shape)
        if len(self.action_shape) == 0:
            self.action_shape = ()

    @property
    def obs_space(self):
        return self._obs_space

    @property
    def action_space(self):
        return self._action_space

    @property
    def flat_obs_dim(self) -> int:
        return int(flat_dim(self.obs_space))

    @property
    def flat_action_dim(self) -> int:
        return int(flat_dim(self.action_space))

    @staticmethod
    def flat(space: Space, obs_or_action: (np.ndarray, list)):
        return flatten(space, obs_or_action)

    def flat_action(self, action: (np.ndarray, list)):
        return flatten(self.action_space, action)

    def flat_obs(self, obs: (np.ndarray, list)):
        return flatten(self.obs_space, obs)

```



---

## How to implement a new dynamics model

---

New dynamics model in Baconian project are supposed to implement the methods and attributes defined in DynamicsModel class (baconian/algo/dynamics/dynamics\_model.py).

```

1  STATUS_LIST = ('CREATED', 'INITED')
2  INIT_STATUS = 'CREATED'
3
4  def __init__(self, env_spec: EnvSpec, parameters: Parameters = None, init_
↪state=None, name='dynamics_model',
5      state_input_scaler: DataScaler = None,
6      action_input_scaler: DataScaler = None,
7      state_output_scaler: DataScaler = None):
8
9      """
10
11      :param env_spec: environment specifications, such as observation space and
↪action space
12      :type env_spec: EnvSpec
13      :param parameters: parameters
14      :type parameters: Parameters
15      :param init_state: initial state of dynamics model
16      :type init_state: str
17      :param name: name of instance, 'dynamics_model' by default
18      :type name: str
19      :param state_input_scaler: data preprocessing scaler of state input
20      :type state_input_scaler: DataScaler
21      :param action_input_scaler: data preprocessing scaler of action input
22      :type action_input_scaler: DataScaler
23      :param state_output_scaler: data preprocessing scaler of state output
24      :type state_output_scaler: DataScaler
25      """
26      super().__init__(name=name)
27      self.env_spec = env_spec
28      self.state = init_state
29      self.parameters = parameters

```

(continues on next page)

(continued from previous page)

```

30     self.state_input = None
31     self.action_input = None
32     self.new_state_output = None
33     self.recorder = Recorder(flush_by_split_status=False, default_obj=self)
34     self._status = StatusWithSingleInfo(obj=self)
35     self.state_input_scaler = state_input_scaler if state_input_scaler else
↳ IdenticalDataScaler(
36         dims=env_spec.flat_obs_dim)
37     self.action_input_scaler = action_input_scaler if action_input_scaler else
↳ IdenticalDataScaler(
38         dims=env_spec.flat_action_dim)
39     self.state_output_scaler = state_output_scaler if state_output_scaler else
↳ IdenticalDataScaler(
40         dims=env_spec.flat_obs_dim)
41
42     def init(self, *args, **kwargs):
43         self.set_status('INITED')
44         self.state = self.env_spec.obs_space.sample()
45
46     @register_counter_info_to_status_decorator(increment=1, info_key='step_counter')
47     def step(self, action: np.ndarray, state=None, allow_clip=False, **kwargs_for_
↳ transit):
48
49         """
50         State transition function (only support one sample transition instead of
↳ batch data)
51
52         :param action: action to be taken
53         :type action: np.ndarray
54         :param state: current state, if None, will use stored state (saved from last
↳ transition)
55         :type state: np.ndarray
56         :param allow_clip: allow clip of observation space, default False
57         :type allow_clip: bool
58         :param kwargs_for_transit: extra kwargs for calling the _state_transit, this
↳ is typically related to the
59                                 specific mode you used
60         :type kwargs_for_transit:
61         :return: new state after step
62         :rtype: np.ndarray
63         """
64         state = np.array(state).reshape(self.env_spec.obs_shape) if state is not None
↳ else self.state
65         action = action.reshape(self.env_spec.action_shape)
66         if allow_clip is True:
67             if state is not None:
68                 state = self.env_spec.obs_space.clip(state)
69                 action = self.env_spec.action_space.clip(action)
70             if self.env_spec.action_space.contains(action) is False:
71                 raise StateOrActionOutOfBoundError(
72                     'action {} out of bound of {}'.format(action, self.env_spec.action_
↳ space.bound()))
73             if self.env_spec.obs_space.contains(state) is False:
74                 raise StateOrActionOutOfBoundError(
75                     'state {} out of bound of {}'.format(state, self.env_spec.obs_space.
↳ bound()))
76             new_state = self._state_transit(state=state, action=self.env_spec.flat_
↳ action(action),

```

(continues on next page)

(continued from previous page)

```

77                                     **kwargs_for_transit)
78         if allow_clip is True:
79             new_state = self.env_spec.obs_space.clip(new_state)
80         if self.env_spec.obs_space.contains(new_state) is False:
81             raise StateOrActionOutOfBoundError(
82                 'new state {} out of bound of {}'.format(new_state, self.env_spec.obs_
↪space.bound()))
83         self.state = new_state
84         return new_state
85
86     @abc.abstractmethod
87     def _state_transit(self, state, action, **kwargs) -> np.ndarray:
88         """
89
90         :param state: original state
91         :type state: np.ndarray
92         :param action: action taken by agent
93         :type action: np.ndarray
94         :param kwargs:
95         :type kwargs:
96         :return: new state after transition
97         :rtype: np.ndarray
98         """
99         raise NotImplementedError
100
101     def copy_from(self, obj) -> bool:
102         """
103
104         :param obj: object to copy from
105         :type obj:
106         :return: True if successful else raise an error
107         :rtype: bool
108         """
109         if not isinstance(obj, type(self)):
110             raise TypeError('Wrong type of obj %s to be copied, which should be %s' %_
↪(type(obj), type(self)))
111         return True
112
113     def make_copy(self):
114         """ Make a copy of parameters and environment specifications."""
115         raise NotImplementedError
116
117     def reset_state(self, state=None):
118         """
119
120         :param state: original state
121         :type state: np.ndarray
122         :return: a random sample space in observation space
123         :rtype: np.ndarray
124         """
125         if state is not None:
126             assert self.env_spec.obs_space.contains(state)
127             self.state = state
128         else:
129             self.state = self.env_spec.obs_space.sample()
130
131     def return_as_env(self) -> Env:

```

(continues on next page)

(continued from previous page)

```
132     """
133
134     :return: an environment with this dynamics model
135     :rtype: DynamicsEnvWrapper
136     """
137     return DynamicsEnvWrapper(dynamics=self,
138                               name=self._name + '_env')
139
140
```

Similar to algorithms, dynamics models are categorized in `baconian/algo/dynamics/dynamics_model.py`, such as `GlobalDynamicsModel` and `DifferentiableDynamics`.

## System Overview of Baconian

The system overview of Baconian is shown in below figure. We design Baconian with the objective to minimize users' coding effort on developing and testing MBRL algorithms. With Baconian, the user can easily setup a MBRL experiment by configuring the target algorithms and modules without the need for understanding the inner mechanisms.

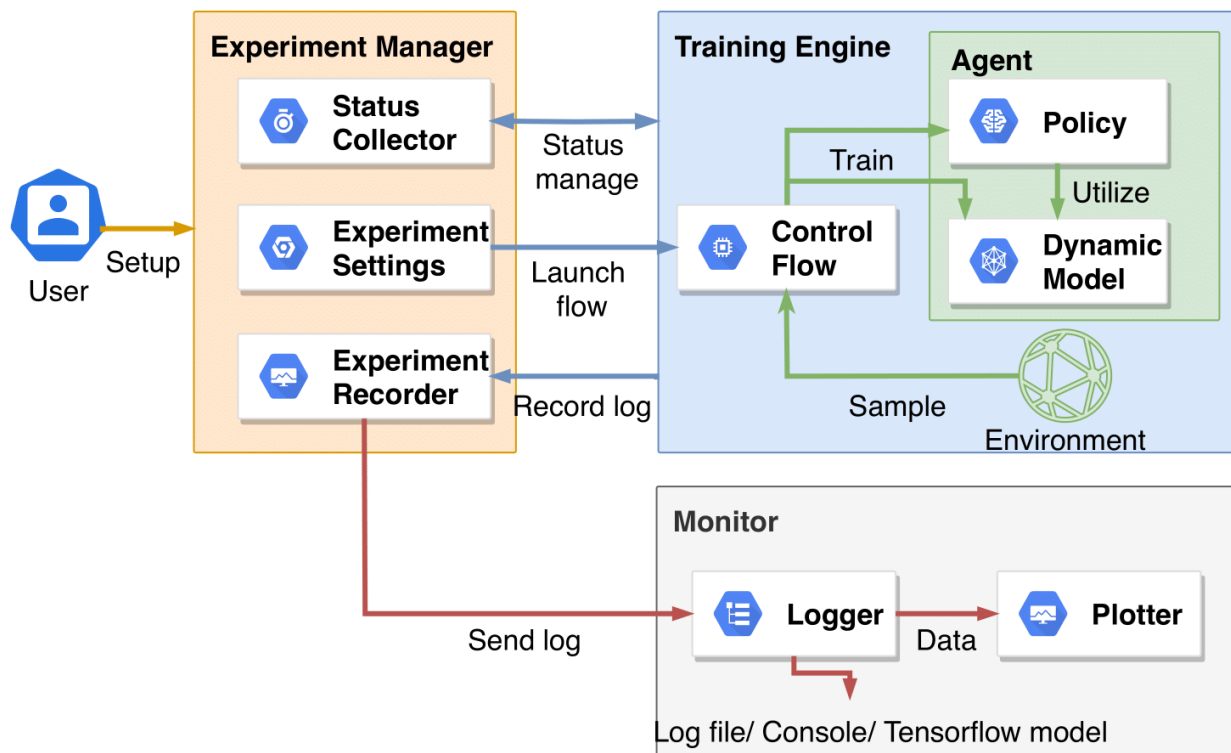


Fig. 1: System Overview of Baconian.

Following figure shows the experiment creation flow of Baconian.

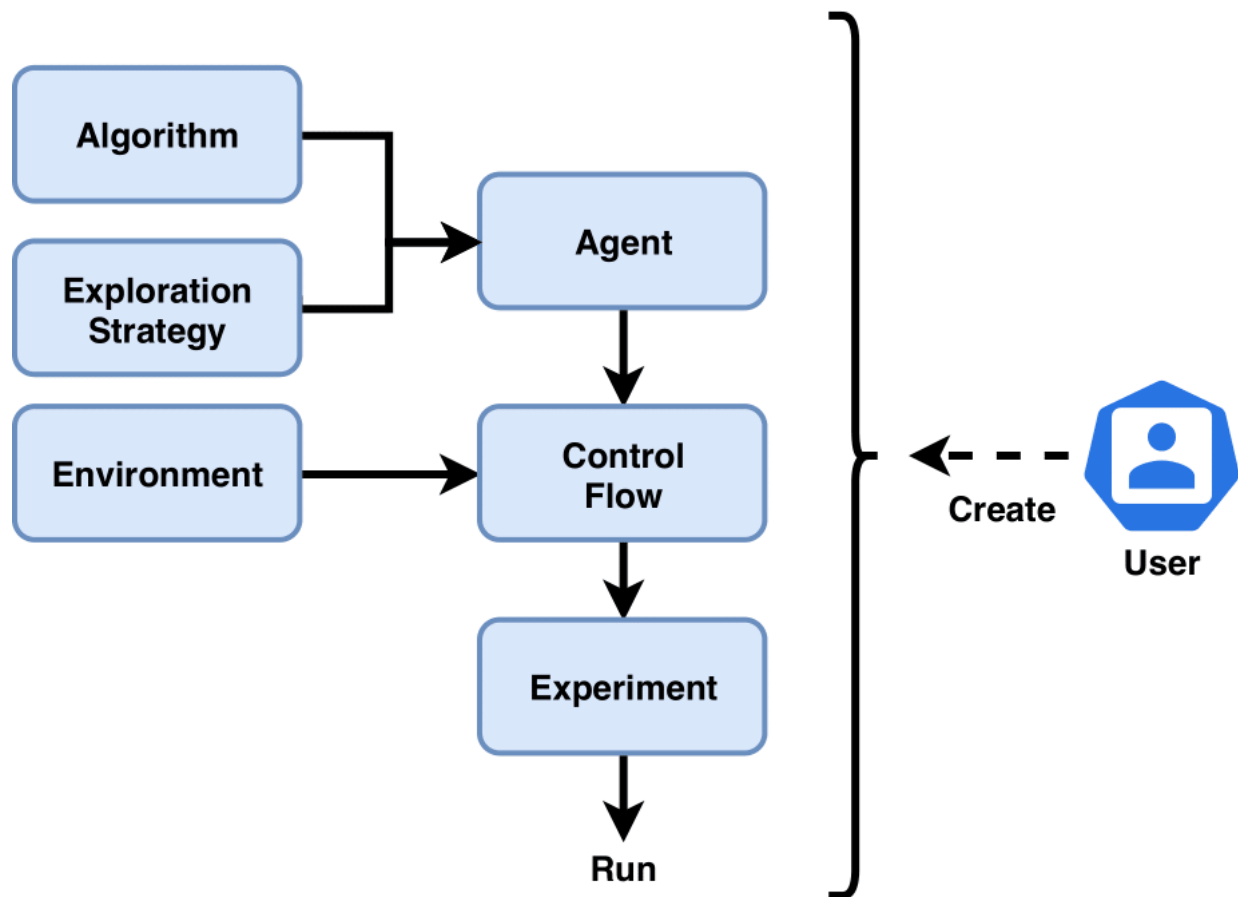


Fig. 2: Baconian experiment creation process.



First, the user should create a environment and a RL algorithm module with necessary hyper-parameters configured, e.g., neural network size, learning rate. Algorithm module is usually composed by a policy module and a dynamics model module depending on different algorithms. Then, the user needs to create an agent module by passing the algorithm module, and the exploration strategy module if needed, into it. Second, the user should create a control flow module that defines how the experiments should be proceeded and the stopping conditions. This includes defining how much samples should be collected for training at each step, and what condition to indicate the completion of an experiment, etc. Some typical control flows have already been implemented in Baconian to meet the users' needs.

Finally, an experiment module should be created by passing the agent, environment, control flow modules into it, and then launched. After that, the Baconian will handle the experiment running, monitoring, and results saving/logging, etc.

For more information on system design of Baconian, you can read our demo paper at: <https://arxiv.org/abs/1904.10762>.



## 10.1 Core Module

### 10.1.1 `baconian.core.core.Basic`

**class** `baconian.core.core.Basic` (*name: str, status=None*)

Basic class within the whole framework

**INIT\_STATUS** = `None`

**STATUS\_LIST** = (`'TRAIN'`, `'TEST'`)

**\_\_init\_\_** (*name: str, status=None*)

Init a new Basic instance.

#### Parameters

- **name** (*str*) – name of the object, can be determined to generate log path, handle tensor-flow name scope etc.
- **status** (*Status*) – A status instance `Status` to indicate the status of the object

**allow\_duplicate\_name** = `False`

**get\_status** () → dict

Return the object's status, a dictionary.

**init** (*\*args, \*\*kwargs*)

Initialize the object

**load** (*\*args, \*\*kwargs*)

Load the parameters from training checkpoints.

**name**

The name(id) of object, a string.

**required\_key\_dict** = ()

**save** (\*args, \*\*kwargs)  
Save the parameters in training checkpoints.

**set\_status** (val)  
Set the object's status.

**status\_list**  
Status list of the object, ('TRAIN', 'TEST').

### 10.1.2 baconian.core.core.Env

**class** `baconian.core.core.Env` (*name: str = 'env', copy\_from\_env=None*)  
Abstract class for environment

**INIT\_STATUS** = 'CREATED'

**STATUS\_LIST** = ('JUST\_RESET', 'INITED', 'TRAIN', 'TEST', 'CREATED')

**\_\_init\_\_** (*name: str = 'env', copy\_from\_env=None*)

**get\_state** ()  
Get the status of the environment.

**init** ()  
Set the status to 'INITED'.

**key\_list** = ()

**reset** ()  
Set the status to 'JUST\_RESET', and update new reset point

**seed** (*seed=None*)

**Parameters** **seed** (*int*) – seed to generate random number

**Returns** seed of the unwrapped environment

**Return type** `int`

**step** (*action*)

**Parameters** **action** (*method*) – agent's action, the environment will react responding to action

### 10.1.3 baconian.core.core.EnvSpec

**class** `baconian.core.core.EnvSpec` (*obs\_space: baconian.common.spaces.base.Space, action\_space: baconian.common.spaces.base.Space*)

**\_\_init\_\_** (*obs\_space: baconian.common.spaces.base.Space, action\_space: baconian.common.spaces.base.Space*)

**action\_shape**

**action\_space**

**Returns** Action space of environment

**Return type** `Space`

```

static flat (space:      baconian.common.spaces.base.Space,      obs_or_action:      (<class
                        'numpy.ndarray'>, <class 'list'>))
    flat the input obs or action :param space: space of environment :type space: Space :param obs_or_action:
    action or observation space :type obs_or_action: (np.ndarray, list) :return: flatten action or observation
    space :rtype: Space

flat_action (action: (<class 'numpy.ndarray'>, <class 'list'>))
    Parameters action ((np.ndarray, list)) – action taken by agent
    Returns flatten action parameter
    Return type np.ndarray

flat_action_dim
    Returns the dimension(length) of flatten action space
    Return type int

flat_obs (obs: (<class 'numpy.ndarray'>, <class 'list'>))
    Parameters obs ((np.ndarray, list)) – observation of the agent
    Returns flatten observation parameter
    Return type np.ndarray

flat_obs_dim
    Returns the dimension(length) of flatten observation space
    Return type int

obs_shape
obs_space
    Returns Observation space of environment
    Return type Space

```

## 10.2 Agent module

### 10.2.1 baconian.core.agent

```

class baconian.core.agent.Agent (name,      env:      (<class      'baconian.core.core.Env'>,
<class
                        'baconian.envs.env_wrapper.Wrapper'>),
    algo:      baconian.algo.algo.Algo,      env_spec:      baco-
    nian.core.core.EnvSpec,      sampler:      baco-
    nian.common.sampler.sampler.Sampler = None, noise_adder:
    baconian.common.noise.AgentActionNoiseWrapper      =
    None,      reset_noise_every_terminal_state=False,      re-
    set_state_every_sample=False,      exploration_strategy:
    baconian.algo.misc.epsilon_greedy.ExplorationStrategy
    =      None,      algo_saving_scheduler:      baco-
    nian.common.schedules.EventScheduler = None)

INIT_STATUS = 'CREATED'

STATUS_LIST = ('CREATED', 'INITED', 'TRAIN', 'TEST')

```

```
__init__(name, env: (<class 'baconian.core.core.Env'>, <class 'baconian.envs.env_wrapper.Wrapper'>), algo: baconian.algo.algo.Algo, env_spec: baconian.core.core.EnvSpec, sampler: baconian.common.sampler.sampler.Sampler = None, noise_adder: baconian.common.noise.AgentActionNoiseWrapper = None, reset_noise_every_terminal_state=False, reset_state_every_sample=False, exploration_strategy: baconian.algo.misc.epsilon_greedy.ExplorationStrategy = None, algo_saving_scheduler: baconian.common.schedules.EventScheduler = None)
```

#### Parameters

- **name** (*str*) – the name of the agent instance
- **env** (*Env*) – environment that interacts with agent
- **algo** (*Algo*) – algorithm of the agent
- **env\_spec** (*EnvSpec*) – environment specifications: action apace and environment space
- **sampler** (*Sampler*) – sampler
- **reset\_noise\_every\_terminal\_state** (*bool*) – reset the noise every sampled trajectory
- **reset\_state\_every\_sample** (*bool*) – reset the state everytime perofrm the sample/rollout
- **noise\_adder** (*AgentActionNoiseWrapper*) – add action noise for exploration in action space
- **exploration\_strategy** (*ExplorationStrategy*) – exploration strategy in action space
- **algo\_saving\_scheduler** (*EventSchedule*) – control the schedule the varying parameters in training process

**init()**

Initialize the algorithm, and set status to 'INITED'.

**is\_testing**

Check whether the agent is testing. Return a boolean value.

**Returns** true if the agent is testing

**Return type** bool

**is\_training**

Check whether the agent is training. Return a boolean value.

**Returns** true if the agent is training

**Return type** bool

**predict** (*\*\*kwargs*)

predict the action given the state

**Parameters** **kwargs** – rest parameters, include key: obs

**Returns** predicted action

**Return type** numpy ndarray

**required\_key\_dict** = {}

**reset\_on\_terminal\_state()**

**sample** (*env*, *sample\_count*: int, *in\_which\_status*: str = 'TRAIN', *store\_flag*=False, *sample\_type*: str = 'transition') -> (<class 'baconian.common.sampler.sample\_data.TransitionData'>, <class 'baconian.common.sampler.sample\_data.TrajectoryData'>)  
sample a certain number of data from environment

#### Parameters

- **env** – environment to sample
- **sample\_count** – int, sample count
- **in\_which\_status** – string, environment status
- **store\_flag** – to store environment samples or not, default False
- **sample\_type** – the type of sample, 'transition' by default

**Returns** sample data from environment

**Return type** some subclass of SampleData: TrajectoryData or TransitionData

**store\_samples** (*samples*: *baconian.common.sampler.sample\_data.SampleData*)  
store the samples into memory/replay buffer if the algorithm that agent hold need to do so, like DQN, DDPG

**Parameters** **samples** (*SampleData*) – sample data of the experiment

**test** (*sample\_count*) → *baconian.common.sampler.sample\_data.SampleData*  
test the agent

**Parameters** **sample\_count** (*int*) – how many trajectories used to evaluate the agent's performance

**Returns** SampleData object.

**train** (*\*args*, *\*\*kwargs*)  
train the agent

**Returns** True for successfully train the agent, false if memory buffer did not have enough data.

**Return type** bool

## 10.3 Experiment modules

### 10.3.1 baconian.core.experiment

For experiments, its functionality should include: 1. experiment and config set up 2. logging control 3. hyper-param tuning etc. 4. visualization 5. any related experiment utility ...

```
class baconian.core.experiment.Experiment (name: str, agent: baconian.core.agent.Agent,
                                           env: baconian.core.core.Env, flow: baconian.core.flow.train_test_flow.Flow,
                                           tuner: baconian.core.tuner.Tuner = None, register_default_global_status=True)
```

```
INIT_STATUS = 'CREATED'
```

```
STATUS_LIST = ('CREATED', 'INITED', 'RUNNING', 'FINISHED', 'CORRUPTED')
```

```
TOTAL_AGENT_TEST_SAMPLE_COUNT ()
```

```
TOTAL_AGENT_TRAIN_SAMPLE_COUNT ()
```

```
TOTAL_AGENT_UPDATE_COUNT()
```

```
TOTAL_ENV_STEP_TEST_SAMPLE_COUNT()
```

```
TOTAL_ENV_STEP_TRAIN_SAMPLE_COUNT()
```

```
__init__(name: str, agent: baconian.core.agent.Agent, env: baconian.core.core.Env, flow: baconian.core.flow.train_test_flow.Flow, tuner: baconian.core.tuner.Tuner = None, register_default_global_status=True)
```

#### Parameters

- **name** (*str*) – name of experiment
- **agent** (*Agent*) – agent of experiment
- **env** (*Env*) – environment of experiment
- **flow** (*Flow*) – control flow to experiment
- **tuner** (*Tuner*) – hyper-parameter tuning method, currently in development
- **register\_default\_global\_status** (*bool*) – register info key and status into global status collection

```
init()
```

Create a new TensorFlow session, and set status to ‘INITED’.

```
required_key_dict = {}
```

```
run()
```

Run the experiment, and set status to ‘RUNNING’.

### 10.3.2 baconian.core.experiment\_runner

```
baconian.core.experiment_runner.single_exp_runner(task_fn,
                                                    auto_choose_gpu_flag=False,
                                                    gpu_id: int = 0, seed=None,
                                                    del_if_log_path_existed=False,
                                                    **task_fn_kwargs)
```

#### Parameters

- **task\_fn** (*method*) – task function defined bu users
- **auto\_choose\_gpu\_flag** (*bool*) – auto choose gpu, default False
- **gpu\_id** (*int*) – gpu id, default 0
- **seed** (*int*) – seed generated by system time

:param del\_if\_log\_path\_existed: delete obsolete log file path if existed, by default False :type del\_if\_log\_path\_existed: bool :param task\_fn\_kwargs: :type task\_fn\_kwargs: :return: :rtype:

```
baconian.core.experiment_runner.duplicate_exp_runner(num, task_fn,
                                                       auto_choose_gpu_flag=False,
                                                       gpu_id: int = 0,
                                                       seeds: list = None,
                                                       del_if_log_path_existed=False,
                                                       **task_fn_kwargs)
```

#### Parameters

- **num** (*int*) – the number of multiple experiments



- **task\_fn** (*method*) – task function, defined by users
- **auto\_choose\_gpu\_flag** (*bool*) – auto choose gpu, default False
- **gpu\_id** (*int*) – gpu id, default 0
- **seeds** (*list*) – seeds generated by system time
- **del\_if\_log\_path\_existed** (*bool*) – delete the existing log path, default False
- **task\_fn\_kwargs** –

**Returns**

**Return type**

## 10.4 Environment Module

### 10.4.1 baconian.envs.gym\_env

Gym environment wrapping module

`baconian.envs.gym_env.GymEnv.name`

The name(id) of object, a string.

`baconian.envs.gym_env.GymEnv.status_list`

Status list of the object, ('TRAIN', 'TEST').

`baconian.envs.gym_env.GymEnv.unwrapped`

**Returns** original unwrapped gym environment

**Return type** gym env

`baconian.envs.gym_env.GymEnv.unwrapped_gym`

**Returns** gym environment, depend on attribute 'unwrapped'

**Return type** gym env

## 10.5 Dynamics Module

### 10.5.1 `baconian.algo.dynamics.dynamics_model.DynamicsModel`

```
class baconian.algo.dynamics.dynamics_model.DynamicsModel (env_spec: baconian.core.core.EnvSpec,  
                                                         parameters: baconian.core.parameters.Parameters  
                                                         = None,  
                                                         init_state=None,  
                                                         name='dynamics_model',  
                                                         state_input_scaler:  
                                                         baconian.common.data_pre_processing.DataScaler  
                                                         = None, action_input_scaler:  
                                                         baconian.common.data_pre_processing.DataScaler  
                                                         = None,  
                                                         state_output_scaler:  
                                                         baconian.common.data_pre_processing.DataScaler  
                                                         = None)
```

```
INIT_STATUS = 'CREATED'
```

```
STATUS_LIST = ('CREATED', 'INITED')
```

```
__init__ (env_spec: baconian.core.core.EnvSpec, parameters: baconian.core.parameters.Parameters  
          = None, init_state=None, name='dynamics_model', state_input_scaler: baconian.common.data_pre_processing.DataScaler = None,  
          action_input_scaler: baconian.common.data_pre_processing.DataScaler = None, state_output_scaler: baconian.common.data_pre_processing.DataScaler = None)
```

#### Parameters

- **env\_spec** (*EnvSpec*) – environment specifications, such as observation space and action space
- **parameters** (*Parameters*) – parameters
- **init\_state** (*str*) – initial state of dynamics model
- **name** (*str*) – name of instance, 'dynamics\_model' by default
- **state\_input\_scaler** (*DataScaler*) – data preprocessing scaler of state input
- **action\_input\_scaler** (*DataScaler*) – data preprocessing scaler of action input
- **state\_output\_scaler** (*DataScaler*) – data preprocessing scaler of state output

```
copy_from (obj) → bool
```

**Parameters** *obj* – object to copy from

**Returns** True if successful else raise an error

**Return type** **bool**

```
init (*args, **kwargs)
```

**make\_copy()**

Make a copy of parameters and environment specifications.

**reset\_state** (*state=None*)

**Parameters** *state* (*np.ndarray*) – original state

**Returns** a random sample space in observation space

**Return type** *np.ndarray*

**return\_as\_env** () → *baconian.core.core.Env*

**Returns** an environment with this dynamics model

**Return type** *DynamicsEnvWrapper*

**step** (*action: numpy.ndarray, state=None, allow\_clip=False, \*\*kwargs\_for\_transit*)

State transition function (only support one sample transition instead of batch data)

**Parameters**

- **action** (*np.ndarray*) – action to be taken
- **state** (*np.ndarray*) – current state, if *None*, will use stored state (saved from last transition)
- **allow\_clip** (*bool*) – allow clip of observation space, default *False*
- **kwargs\_for\_transit** – extra kwargs for calling the *\_state\_transit*, this is typically related to the specific mode you used

**Returns** new state after step

**Return type** *np.ndarray*

## 10.5.2 Abstraction class for different types of dynamics model:

```
class baconian.algo.dynamics.dynamics_model.LocalDyanmicsModel (env_spec: baco-
    nian.core.core.EnvSpec,
    paramete-
    rs: baco-
    nian.core.parameters.Parameters
    = None,
    init_state=None,
    name='dynamics_model',
    state_input_scaler:
    baco-
    nian.common.data_pre_processing.Dataa
    = None, ac-
    tion_input_scaler:
    baco-
    nian.common.data_pre_processing.Dataa
    = None,
    state_output_scaler:
    baco-
    nian.common.data_pre_processing.Dataa
    = None)
```

```
class baconian.algo.dynamics.dynamics_model.GlobalDynamicsModel (env_spec:
                                                                    baco-
                                                                    nian.core.core.EnvSpec,
                                                                    param-
                                                                    eters:    baco-
                                                                    nian.core.parameters.Parameters
                                                                    =        None,
                                                                    init_state=None,
                                                                    name='dynamics_model',
                                                                    state_input_scaler:
                                                                    baco-
                                                                    nian.common.data_pre_processing.Dat
                                                                    = None, ac-
                                                                    tion_input_scaler:
                                                                    baco-
                                                                    nian.common.data_pre_processing.Dat
                                                                    =        None,
                                                                    state_output_scaler:
                                                                    baco-
                                                                    nian.common.data_pre_processing.Dat
                                                                    = None)

class baconian.algo.dynamics.dynamics_model.TrainableDyanmicsModel

    train (*args, **kwargs)

class baconian.algo.dynamics.dynamics_model.DifferentiableDynamics (input_node_dict:
                                                                    dict, out-
                                                                    put_node_dict:
                                                                    dict)

    __init__ (input_node_dict: dict, output_node_dict: dict)

    grad_on_input (key_or_node: (<class 'str'>, <sphinx.ext.autodoc.importer._MockObject object at
                                0x7f604e7de2e8>), order=1, batch_flag=False)

    split_and_hessian (out_node, innode)
```

## 10.6 Algorithm Module

### 10.6.1 baconian.algo.algo.Algo

```
class baconian.algo.algo.Algo (env_spec:  baconian.core.core.EnvSpec, name:  str = 'algo',
                                warm_up_trajectories_number=0)

    Abstract class for algorithms

    INIT_STATUS = 'CREATED'

    STATUS_LIST = ['CREATED', 'INITED', 'TRAIN', 'TEST']

    __init__ (env_spec:      baconian.core.core.EnvSpec,      name:      str      =      'algo',
              warm_up_trajectories_number=0)
        Constructor

        Parameters

        • env_spec (EnvSpec) – environment specifications
```

- **name** (*str*) – name of the algorithm
- **warm\_up\_trajectories\_number** (*int*) – how many trajectories used to warm up the training

**append\_to\_memory** (*\*args, \*\*kwargs*)

For off-policy algorithm, use this API to append the data into replay buffer. samples will be read as the first argument passed into this API, like algo.append\_to\_memory(samples=x, ...)

**init** ()

Initialization method, such as network random initialization in Tensorflow

**Returns**

**is\_testing**

A boolean indicate the if the algorithm is in training status

**Returns** True if in testing

**Return type** bool

**is\_training**

A boolean indicate the if the algorithm is in training status

**Returns** True if in training

**Return type** bool

**predict** (*\*arg, \*\*kwargs*)

Predict function, given the obs as input, return the action, obs will be read as the first argument passed into this API, like algo.predict(obs=x, ...)

**Returns** predicted action

**Return type** np.ndarray

**test** (*\*arg, \*\*kwargs*) → dict

Testing API, most of the evaluation can be done by agent instead of algorithms, so this API can be skipped

**Returns** test results, e.g., rewards

**Return type** dict

**train** (*\*arg, \*\*kwargs*) → dict

Training API, specific arguments should be defined by each algorithms itself.

**Returns** training results, e.g., loss

**Return type** dict

**warm\_up** (*trajectory\_data: baconian.common.sampler.sample\_data.TrajectoryData*)

Use some data to warm up the algorithm, e.g., compute the mean/std-dev of the state to perform normalization. Data used in warm up process will not be added into the memory :param trajectory\_data: TrajectoryData object :type trajectory\_data: TrajectoryData

**Returns** None

## 10.7 Flow Module

### 10.7.1 `baconian.core.flow.train_test_flow.Flow`

**class** `baconian.core.flow.train_test_flow.Flow` (*func\_dict*)

Interface of experiment flow module, it defines the workflow of the reinforcement learning experiments.

**\_\_init\_\_** (*func\_dict*)

Constructor for Flow.

**Parameters** **func\_dict** (*dict*) – the function and its arguments that will be called in the Flow

**\_call\_func** (*key*, *\*\*extra\_kwargs*)

Call a function that is pre-defined in self.func\_dict

**Parameters**

- **key** (*str*) – name of the function, e.g., train, test, sample.
- **extra\_kwargs** – some extra kwargs you may want to be passed in the function calling

**Returns** actual return value of the called function if self.func\_dict has such function otherwise None.

**Return type**

**\_launch** () → bool

Abstract method to be implemented by subclass for a certain workflow.

**Returns** True if the flow correctly executed and finished

**Return type** bool

**launch** () → bool

Launch the flow until it finished or catch a system-allowed errors (e.g., out of GPU memory, to ensure the log will be saved safely).

**Returns** True if the flow correctly executed and finished

**Return type** bool

**required\_func** = ()

**required\_key\_dict** = {}

**class** `baconian.core.flow.train_test_flow.TrainTestFlow` (*train\_sample\_count\_func*,  
*config\_or\_config\_dict*:  
(*<class 'baconian.config.dict\_config.DictConfig'>*,  
*<class 'dict'>*), *func\_dict*:  
*dict*)

A typical sampling-training and testing workflow, that used by most of model-free/model-based reinforcement learning method. Typically, it repeat the sampling(saving to memory if off policy)->training(from memory if off-policy, from samples if on-policy)->test

**\_\_init\_\_** (*train\_sample\_count\_func*, *config\_or\_config\_dict*: (*<class 'baconian.config.dict\_config.DictConfig'>*, *<class 'dict'>*), *func\_dict*: *dict*)

Constructor of TrainTestFlow

**Parameters**

- **train\_sample\_count\_func** (*method*) – a function indicates how much training samples the agent has collected currently.
- **config\_or\_config\_dict** (*Config or dict*) – a `Config` or a `dict` should have the keys: (`TEST EVERY SAMPLE COUNT`, `TRAIN EVERY SAMPLE COUNT`, `START TRAIN AFTER SAMPLE COUNT`, `START TEST AFTER SAMPLE COUNT`)
- **func\_dict** (*dict*) – function dict, holds the keys: ‘sample’, ‘train’, ‘test’. each item in the dict as also should be a dict, holds the keys ‘func’, ‘args’, ‘kwargs’

**\_\_is\_ended()**

**Returns** True if an experiment is ended

**Return type** bool

**\_\_launch()** → bool

Launch the flow until it finished or catch a system-allowed errors (e.g., out of GPU memory, to ensure the log will be saved safely).

**Returns** True if the flow correctly executed and finished

**Return type** bool

**required\_func** = ('train', 'test', 'sample')

**required\_key\_dict** = {'START\_TEST\_AFTER\_SAMPLE\_COUNT': 1, 'START\_TRAIN\_AFTER\_SAMPLE\_COUNT': 1}

## 10.7.2 baconian.core.flow.dyna\_flow.DynaFlow

```
class baconian.core.flow.dyna_flow.DynaFlow(train_sample_count_func, config_or_config_dict: (<class 'baconian.config.dict_config.DictConfig'>, <class 'dict'>), func_dict: dict)
```

A typical flow for utilizing the model-based algo, it is not restricted to Dyna algorithms, but can be utilized by others.

```
__init__(train_sample_count_func, config_or_config_dict: (<class 'baconian.config.dict_config.DictConfig'>, <class 'dict'>), func_dict: dict)
```

**Parameters**

- **train\_sample\_count\_func** (*method*) – a function indicates how much training samples the agent has collected currently.
- **config\_or\_config\_dict** (*Config or dict*) – a `Config` or a `dict` should have the keys: (`TEST EVERY SAMPLE COUNT`, `TRAIN EVERY SAMPLE COUNT`, `START TRAIN AFTER SAMPLE COUNT`, `START TEST AFTER SAMPLE COUNT`)
- **func\_dict** (*dict*) – function dict, holds the keys: ‘sample’, ‘train’, ‘test’. each item in the dict as also should be a dict, holds the keys ‘func’, ‘args’, ‘kwargs’

**\_\_is\_ended()**

**Returns** True if an experiment is ended

**Return type** bool

**\_\_launch()** → bool

Launch the flow until it finished or catch a system-allowed errors (e.g., out of GPU memory, to ensure the log will be saved safely).

**Returns** True if the flow correctly executed and finished

**Return type** bool

```
required_func = ('train_algo', 'train_algo_from_synthesized_data', 'train_dynamics', 'train_
```

```
required_key_dict = {'START_TEST_ALGO_AFTER_SAMPLE_COUNT': 1, 'START_TEST_DYNAMICS_AFT
```

## 10.8 Common module

### 10.8.1 baconian.common

```
class baconian.common.sampler.sampler.Sampler
```

Sampler module that handle the sampling procedure for training/testing of the agent.

```
static sample (env:      baconian.core.core.Env,      agent,      sample_count:      int,      sam-
ple_type='transition',      reset_at_start=None)      ->      (<class      'baco-
nian.common.sampler.sample_data.TransitionData'>,      <class      'baco-
nian.common.sampler.sample_data.TrajectoryData'>)
```

a static method of sample function

## Parameters

- **env** – environment object to sample from.
- **agent** – agent object to offer the sampling policy
- **in\_which\_status** – a string, “TEST” or “TRAIN” indicate this sample is used for training or testing (evaluation)
- **sample\_count** – number of samples. If the sample\_type == “transition”, then this value means the number of transitions, usually for off-policy method like DQN, DDPG. If the sample\_type == “trajectory”, then this value means the numbers of trajectories.
- **sample\_type** – a string, “transition” or “trajectory”.
- **reset\_at\_start** – A bool, if True, will reset the environment at the beginning, if False, continue sampling based on previous state (this is useful for certain tasks that you need to preserve previous state to reach the terminal goal state). If None, for sample\_type == “transition”, it will set to False, for sample\_type == “trajectory”, it will set to True.

**Returns** SampleData object.

## 10.9 Visualisation Module

### 10.9.1 baconian.common.plotter

```
class baconian.common.plotter.Plotter
```

```
color_list = ['b', 'r', 'g', 'm', 'y', 'k', 'cyan', 'plum', 'darkgreen', 'darkorange',
```

```
markers = ('+', 'x', 'v', 'o', '^', '<', '>', 's', 'p', '*', 'h', 'H', 'D', 'd', 'P',
```

```
static plot_any_key_in_log(data, key, index, exp_num=1, sub_log_dir_name=None,
                           scatter_flag=False, histogram_flag=False, save_flag=False,
                           save_path=None, save_format='png', file_name=None, separate_exp_flag=True, mean_stddev_flag=False, path_list=None)
```



### Parameters

- **data** – a pandas DataFrame containing (index and) key columns
- **key** – in y-axis, the variable to plot, assigned by user
- **index** – in x-axis, the argument assigned by user
- **sub\_log\_dir\_name** – the sub-directory which the log file to plot is saved in
- **exp\_num** – [optional] the number of experiments to visualize
- **scatter\_flag** – [optional] draw scatter plot if true
- **histogram\_flag** – [optional] draw histogram if true
- **save\_flag** – [optional] save the figure to a file if true
- **save\_path** – [optional] save path of figure, assigned by user, the directory of log\_path by default
- **save\_format** – [optional] format of figure to save, png by default
- **file\_name** – [optional] the file name of the file to save, key\_VERUS\_index by default
- **separate\_exp\_flag** – [optional] plot the results of each experiment separately if true
- **mean\_stddev\_flag** – [optional] plot the mean value of multiple experiment results and standard deviation
- **path\_list** – [optional] the list of save paths assigned by users, figure file will be saved to each path

### Returns

```
static plot_any_scatter_in_log (res_dict, res_name, file_name, key, index, op, scatter_flag=False, save_flag=False, save_path=None, fig_id=4, label="", restrict_dict=None)  
plot_fig (fig_num, col_id, x, y, title, x_label, y_label, label=' ', marker='*')
```



# CHAPTER 11

---

## Contributions

---

We hope everyone that is working on Model-Based RL research can contribute to this project to better boost the development of Model-Based RL research.

Please submit an issue in the GitHub repository if you run into any problem or bug. Also feel free to contact me: [linsen001@e.ntu.edu.sg](mailto:linsen001@e.ntu.edu.sg).



## CHAPTER 12

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



### b

`baconian.core.experiment`, [67](#)

`baconian.core.experiment_runner`, [68](#)

`baconian.envs.gym_env.GymEnv`, [69](#)





## Symbols

- `__init__()` (*baconian.algo.algo.Algo* method), 72
  - `__init__()` (*baconian.algo.dynamics.dynamics\_model.DifferentiableDynamics* method), 72
  - `__init__()` (*baconian.algo.dynamics.dynamics\_model.DynamicsModel* method), 70
  - `__init__()` (*baconian.core.agent.Agent* method), 65
  - `__init__()` (*baconian.core.core.Basic* method), 63
  - `__init__()` (*baconian.core.core.Env* method), 64
  - `__init__()` (*baconian.core.core.EnvSpec* method), 64
  - `__init__()` (*baconian.core.experiment.Experiment* method), 68
  - `__init__()` (*baconian.core.flow.dyna\_flow.DynaFlow* method), 75
  - `__init__()` (*baconian.core.flow.train\_test\_flow.Flow* method), 74
  - `__init__()` (*baconian.core.flow.train\_test\_flow.TrainTestFlow* method), 74
  - `_call_func()` (*baconian.core.flow.train\_test\_flow.Flow* method), 74
  - `_is_ended()` (*baconian.core.flow.dyna\_flow.DynaFlow* method), 75
  - `_is_ended()` (*baconian.core.flow.train\_test\_flow.TrainTestFlow* method), 75
  - `_launch()` (*baconian.core.flow.dyna\_flow.DynaFlow* method), 75
  - `_launch()` (*baconian.core.flow.train\_test\_flow.Flow* method), 74
  - `_launch()` (*baconian.core.flow.train\_test\_flow.TrainTestFlow* method), 75
- ## A
- `action_shape` (*baconian.core.core.EnvSpec* attribute), 64
  - `action_space` (*baconian.core.core.EnvSpec* attribute), 64
- ## B
- Agent* (class in *baconian.core.agent*), 65
  - Algo* (class in *baconian.algo.algo*), 72
  - `allow_duplicate_name` (*baconian.core.core.Basic* attribute), 63
  - `append_to_memory()` (*baconian.algo.algo.Algo* method), 73
- ## C
- `color_list` (*baconian.common.plotter.Plotter* attribute), 76
  - `copy_from()` (*baconian.algo.dynamics.dynamics\_model.DynamicsModel* method), 70
- ## D
- DifferentiableDynamics* (class in *baconian.algo.dynamics.dynamics\_model*), 72
  - `duplicate_exp_runner()` (in module *baconian.core.experiment\_runner*), 68
  - DynaFlow* (class in *baconian.core.flow.dyna\_flow*), 75
  - DynamicsModel* (class in *baconian.algo.dynamics.dynamics\_model*), 70
- ## E
- Env* (class in *baconian.core.core*), 64
  - EnvSpec* (class in *baconian.core.core*), 64
  - Experiment* (class in *baconian.core.experiment*), 67
- ## F
- `flat()` (*baconian.core.core.EnvSpec* static method), 64
  - `flat_action()` (*baconian.core.core.EnvSpec* method), 65

flat\_action\_dim (*baconian.core.core.EnvSpec attribute*), 65  
 flat\_obs() (*baconian.core.core.EnvSpec method*), 65  
 flat\_obs\_dim (*baconian.core.core.EnvSpec attribute*), 65  
 Flow (*class in baconian.core.flow.train\_test\_flow*), 74

## G

get\_state() (*baconian.core.core.Env method*), 64  
 get\_status() (*baconian.core.core.Basic method*), 63  
 GlobalDynamicsModel (*class in baconian.algo.dynamics.dynamics\_model*), 71  
 grad\_on\_input() (*baconian.algo.dynamics.dynamics\_model.DifferentiableDynamics method*), 72

## I

init() (*baconian.algo.algo.Algo method*), 73  
 init() (*baconian.algo.dynamics.dynamics\_model.DynamicsModel method*), 70  
 init() (*baconian.core.agent.Agent method*), 66  
 init() (*baconian.core.core.Basic method*), 63  
 init() (*baconian.core.core.Env method*), 64  
 init() (*baconian.core.experiment.Experiment method*), 68  
 INIT\_STATUS (*baconian.algo.algo.Algo attribute*), 72  
 INIT\_STATUS (*baconian.algo.dynamics.dynamics\_model.DynamicsModel attribute*), 70  
 INIT\_STATUS (*baconian.core.agent.Agent attribute*), 65  
 INIT\_STATUS (*baconian.core.core.Basic attribute*), 63  
 INIT\_STATUS (*baconian.core.core.Env attribute*), 64  
 INIT\_STATUS (*baconian.core.experiment.Experiment attribute*), 67  
 is\_testing (*baconian.algo.algo.Algo attribute*), 73  
 is\_testing (*baconian.core.agent.Agent attribute*), 66  
 is\_training (*baconian.algo.algo.Algo attribute*), 73  
 is\_training (*baconian.core.agent.Agent attribute*), 66

## K

key\_list (*baconian.core.core.Env attribute*), 64

## L

launch() (*baconian.core.flow.train\_test\_flow.Flow method*), 74  
 load() (*baconian.core.core.Basic method*), 63  
 LocalDyanmicsModel (*class in baconian.algo.dynamics.dynamics\_model*), 71

## M

make\_copy() (*baconian.algo.dynamics.dynamics\_model.DynamicsModel*

*method*), 70

markers (*baconian.common.plotter.Plotter attribute*), 76

## N

name (*baconian.core.core.Basic attribute*), 63  
 name (*in module baconian.envs.gym\_env.GymEnv*), 69

## O

obs\_shape (*baconian.core.core.EnvSpec attribute*), 65  
 obs\_space (*baconian.core.core.EnvSpec attribute*), 65

## P

plot\_any\_key\_in\_log() (*baconian.common.plotter.Plotter static method*), 76  
 plot\_any\_scatter\_in\_log() (*baconian.common.plotter.Plotter static method*), 76  
 plot\_fig() (*baconian.common.plotter.Plotter method*), 77  
 Plotter (*class in baconian.common.plotter*), 76  
 predict() (*baconian.algo.algo.Algo method*), 73  
 predict() (*baconian.core.agent.Agent method*), 66

## R

required\_func (*baconian.core.flow.dyna\_flow.DynaFlow attribute*), 76  
 required\_func (*baconian.core.flow.train\_test\_flow.Flow attribute*), 74  
 required\_func (*baconian.core.flow.train\_test\_flow.TrainTestFlow attribute*), 75  
 required\_key\_dict (*baconian.core.agent.Agent attribute*), 66  
 required\_key\_dict (*baconian.core.core.Basic attribute*), 63  
 required\_key\_dict (*baconian.core.experiment.Experiment attribute*), 68  
 required\_key\_dict (*baconian.core.flow.dyna\_flow.DynaFlow attribute*), 76  
 required\_key\_dict (*baconian.core.flow.train\_test\_flow.Flow attribute*), 74  
 required\_key\_dict (*baconian.core.flow.train\_test\_flow.TrainTestFlow attribute*), 75  
 reset() (*baconian.core.core.Env method*), 64  
 reset\_on\_terminal\_state() (*baconian.core.agent.Agent method*), 66

`reset_state()` (*baconian.algo.dynamics.dynamics\_model.DynamicsModel method*), 71

`TOTAL_ENV_STEP_TEST_SAMPLE_COUNT()` (*baconian.core.experiment.Experiment method*), 68

`return_as_env()` (*baconian.algo.dynamics.dynamics\_model.DynamicsModel method*), 71

`TOTAL_ENV_STEP_TRAIN_SAMPLE_COUNT()` (*baconian.core.experiment.Experiment method*), 68

`run()` (*baconian.core.experiment.Experiment method*), 68

`train()` (*baconian.algo.algo.Algo method*), 73

`train()` (*baconian.algo.dynamics.dynamics\_model.TrainableDynamicsModel method*), 72

`train()` (*baconian.core.agent.Agent method*), 67

`TrainableDynamicsModel` (class in *baconian.algo.dynamics.dynamics\_model*), 72

`TrainTestFlow` (class in *baconian.core.flow.train\_test\_flow*), 74

## S

`sample()` (*baconian.common.sampler.sampler.Sampler static method*), 76

`sample()` (*baconian.core.agent.Agent method*), 66

`Sampler` (class in *baconian.common.sampler.sampler*), 76

`save()` (*baconian.core.core.Basic method*), 63

`seed()` (*baconian.core.core.Env method*), 64

`set_status()` (*baconian.core.core.Basic method*), 64

`single_exp_runner()` (in module *baconian.core.experiment\_runner*), 68

`split_and_hessian()` (*baconian.algo.dynamics.dynamics\_model.DifferentiableDynamics method*), 72

`STATUS_LIST` (*baconian.algo.algo.Algo attribute*), 72

`STATUS_LIST` (*baconian.algo.dynamics.dynamics\_model.DynamicsModel attribute*), 70

`STATUS_LIST` (*baconian.core.agent.Agent attribute*), 65

`STATUS_LIST` (*baconian.core.core.Basic attribute*), 63

`status_list` (*baconian.core.core.Basic attribute*), 64

`STATUS_LIST` (*baconian.core.core.Env attribute*), 64

`STATUS_LIST` (*baconian.core.experiment.Experiment attribute*), 67

`status_list` (in module *baconian.envs.gym\_env.GymEnv*), 69

`step()` (*baconian.algo.dynamics.dynamics\_model.DynamicsModel method*), 71

`step()` (*baconian.core.core.Env method*), 64

`store_samples()` (*baconian.core.agent.Agent method*), 67

## U

`unwrapped` (in module *baconian.envs.gym\_env.GymEnv*), 69

`unwrapped_gym` (in module *baconian.envs.gym\_env.GymEnv*), 69

## W

`warm_up()` (*baconian.algo.algo.Algo method*), 73

## T

`test()` (*baconian.algo.algo.Algo method*), 73

`test()` (*baconian.core.agent.Agent method*), 67

`TOTAL_AGENT_TEST_SAMPLE_COUNT()` (*baconian.core.experiment.Experiment method*), 67

`TOTAL_AGENT_TRAIN_SAMPLE_COUNT()` (*baconian.core.experiment.Experiment method*), 67

`TOTAL_AGENT_UPDATE_COUNT()` (*baconian.core.experiment.Experiment method*), 67