```python
"""
Created on Mon Mar 01 14:47:27 2024
@author: Bing
The core code of the LBDSAC algorithm
"""
import os
os.environ["CUDA_VISIBLE_DEVICES"] = '0'
import tensorflow as tf
import shutil
import numpy as np
import pandas as pd
import time
import datetime
import seaborn as sns
from logx import EpochLogger
from mpi_tf import sync_all_params, MpiAdamOptimizer
from mpi_tools import mpi_fork, mpi_sum, proc_id, mpi_statistics_scalar, num_procs
from scipy.stats import norm
import globalvar as gl
import json
import argparse
from run_utils import setup_logger_kwargs
import save_data
import matplotlib.pyplot as plt
from env2 import unPower
from model_modules import placeholder,placeholders,cnn_mlp,get_vars,count_vars,gaussian_likelihood,\
    get_target_update,mlp_gaussian_policy,apply_squashing_func,mlp_actor,mlp_critic,Fig
import ReplayBuffer

n_T = gl.get_value('n_T')
n_wind = gl.get_value('n_wind')
n_load = gl.get_value('n_load')
n_gen = gl.get_value('n_gen')
n_wind_data = gl.get_value('n_wind_data')
obs_dim = gl.get_value('state_dim')
act_dim = gl.get_value('action_dim')
P_load = gl.get_value('P_load')
P_wind = gl.get_value('P_wind')
load_capacity = gl.get_value('load_capacity')
wind_capacity = gl.get_value('wind_capacity')
EPS = 1e-8

snow = datetime.datetime.now()
seed_now = str(snow.month)+str(snow.day)+str(snow.hour)
```

```python
45.    seed_now = int(seed_now)
46.
47.    parser = argparse.ArgumentParser()
48.    parser.add_argument('--Train', type=bool, default=True, help='True or False')
49.    parser.add_argument('--Continue', type=bool, default=False, help='True or False')
50.    parser.add_argument('--Portion', type=bool, default=False, help='True or False')
51.    parser.add_argument('--SaveFig', type=bool, default=False, help='True or False')
52.    parser.add_argument('--load_buffer', type=bool, default=False, help='True or False')
53.    parser.add_argument('--Model_version', type=int, default=8000, help='8000')
54.    parser.add_argument('--epochs', type=int, default=8000, help='800')
55.    parser.add_argument('--steps_per_epoch', type=int, default=20, help='20')
56.    parser.add_argument('--update_freq', type=int, default=720, help='720')
57.    parser.add_argument('--save_freq', type=int, default=1000,  help='500')
58.    parser.add_argument('--cost_lim', type=float, default=1e-3, help='1e-3')
59.    parser.add_argument('--lr', type=float, default=1e-3,help='1e-3')
60.    parser.add_argument('--hidden_sizes_actor', type=list, default=[128, 256, 256, 64], help='actor')
61.    parser.add_argument('--hidden_sizes_critic', type=list, default=[128, 256, 64], help='critic')
62.    parser.add_argument('--hidden_sizes_var', type=list, default=[128, 256, 64], help='var')
63.    parser.add_argument('--cnn_actor', type=list, default=[], help='CNNs of policy network actor')
64.    parser.add_argument('--cnn_critic', type=list, default=[], help='CNNs of  network critic')
65.    parser.add_argument('--cnn_var', type=list, default=[], help='CNNs of network var')
66.    parser.add_argument('--hid', type=int, default=512)
67.    parser.add_argument('--l', type=int, default=2)
68.    parser.add_argument('--gamma', type=float, default=0.99)
69.    parser.add_argument('--seed', '-s', type=int, default=seed_now)
70.    parser.add_argument('--exp_name', type=str, default='sac')
71.    parser.add_argument('--cpu', type=int, default=1)
72.    parser.add_argument('--render', default=False, action='store_true')
73.    parser.add_argument('--local_start_steps', default=3600, type=int, help='3600')
74.    parser.add_argument('--local_update_after', default=3600, type=int, help='3600')
75.    parser.add_argument('--batch_size', default=256, type=int)  # 256
76.    parser.add_argument('--fixed_entropy_bonus', default=None, type=float)
77.    parser.add_argument('--entropy_constraint', type=float, default=-1)
78.    parser.add_argument('--fixed_cost_penalty', default=None, type=float)
79.    parser.add_argument('--cost_constraint', type=float, default=None)
80.    parser.add_argument('--lr_s', type=int, default=50)
81.    parser.add_argument('--damp_s', type=int, default=10)
82.    parser.add_argument('--logger_kwargs_str', type=json.loads, default='{"output_dir": "./logger"}')
83.    args = parser.parse_args()
84.
85.    """
86.    Soft Actor-Critic with Logic-Based Benders Decomposition Algorithm
87.    """
88.
```

```python
89.   def LBDSAC(env_fn, actor_fn=mlp_actor, critic_fn=mlp_critic,
90.          ac_kwargs_actor=dict(),ac_kwargs_critic=dict(),ac_kwargs_var=dict(),
91.          seed=0, steps_per_epoch=1200, epochs=101, replay_size=500 * 2 * 720, gamma=0.99, cl=0.5,
92.          polyak=0.995, lr=1e-4, batch_size=1024, local_start_steps=600,
93.          max_ep_len=n_T, logger_kwargs=dict(), save_freq=50, local_update_after=int(1e3),
94.          update_freq=120, render=False,
95.          fixed_entropy_bonus=None, entropy_constraint=-1.0,
96.          fixed_cost_penalty=None, cost_constraint=None, cost_lim=None,
97.          reward_scale=1, lr_scale=1, damp_scale=0,Train=True,Continue_Training=False,
98.          Model_Portion=False,load_buffer=True
99.          ):
100.
101.   use_costs = fixed_cost_penalty or cost_constraint or cost_lim
102.   logger = EpochLogger(**logger_kwargs)
103.   logger.save_config(locals())
104.
105.   # Env instantiation
106.   env, test_env = env_fn(), env_fn()
107.   #Setting seeds
108.   seed += 200 * proc_id()
109.   tf.set_random_seed(seed)
110.   np.random.seed(seed)
111.
112.   # Inputs to computation graph
113.   x_ph, a_ph, x2_ph, r_ph, d_ph, c1_ph, c2_ph = placeholders(obs_dim, act_dim, obs_dim, None, None, None, None
      )
114.
115.   # Main outputs from computation graph
116.   with tf.variable_scope('main'):
117.       mu, pi, logp_pi = actor_fn(x_ph, a_ph, **ac_kwargs_actor)
118.       qr1, qr1_pi = critic_fn(x_ph, a_ph, pi, name='qr1', **ac_kwargs_critic)
119.       qr2, qr2_pi = critic_fn(x_ph, a_ph, pi, name='qr2', **ac_kwargs_critic)
120.       qc1, qc1_pi = critic_fn(x_ph, a_ph, pi, name='qc1', **ac_kwargs_critic)
121.       qc2, qc2_pi = critic_fn(x_ph, a_ph, pi, name='qc2', **ac_kwargs_critic)
122.
123.   with tf.variable_scope('main', reuse=True):
124.       # Additional policy output from a different observation placeholder
125.       _, pi2, logp_pi2 = actor_fn(x2_ph, a_ph, **ac_kwargs_actor)
126.
127.   # Target value network
128.   with tf.variable_scope('target'):
129.       _, qr1_pi_targ = critic_fn(x2_ph, a_ph, pi2, name='qr1', **ac_kwargs_critic)
130.       _, qr2_pi_targ = critic_fn(x2_ph, a_ph, pi2, name='qr2', **ac_kwargs_critic)
131.       _, qc1_pi_targ = critic_fn(x2_ph, a_ph, pi2, name='qc1', **ac_kwargs_critic)
```

```python
132.        _, qc2_pi_targ = critic_fn(x2_ph, a_ph, pi2, name='qc2', **ac_kwargs_critic)
133.
134.    # Entropy bonus
135.    if fixed_entropy_bonus is None:
136.        with tf.variable_scope('entreg'):
137.            soft_alpha = tf.get_variable('soft_alpha',
138.                            initializer=0.0,
139.                            trainable=True,
140.                            dtype=tf.float32)
141.        alpha = tf.nn.softplus(soft_alpha)
142.    else:
143.        alpha = tf.constant(fixed_entropy_bonus)
144.    log_alpha = tf.log(tf.clip_by_value(alpha, 1e-8, 1e8))
145.
146.    # Cost penalty
147.    if use_costs:
148.        if fixed_cost_penalty is None:
149.            with tf.variable_scope('costpen1'):
150.                soft_beta1 = tf.get_variable('soft_beta1',
151.                                initializer=0.0,
152.                                trainable=True,
153.                                dtype=tf.float32)
154.            beta1 = tf.nn.softplus(soft_beta1)
155.            log_beta1 = tf.log(tf.clip_by_value(beta1, 1e-8, 1e8))
156.        else:
157.            beta1 = tf.constant(fixed_cost_penalty)
158.            log_beta1 = tf.log(tf.clip_by_value(beta1, 1e-8, 1e8))
159.    else:
160.        beta1 = 0.0  # costs do not contribute to policy optimization
161.        print('Not using costs')
162.
163.    if use_costs:
164.        if fixed_cost_penalty is None:
165.            with tf.variable_scope('costpen2'):
166.                soft_beta2 = tf.get_variable('soft_beta2',
167.                                initializer=0.0,
168.                                trainable=True,
169.                                dtype=tf.float32)
170.            beta2 = tf.nn.softplus(soft_beta2)
171.            log_beta2 = tf.log(tf.clip_by_value(beta2, 1e-8, 1e8))
172.        else:
173.            beta2 = tf.constant(fixed_cost_penalty)
174.            log_beta2 = tf.log(tf.clip_by_value(beta2, 1e-8, 1e8))
175.    else:
```

```python
176.         beta2 = 0.0  # costs do not contribute to policy optimization
177.         print('Not using costs')
178.
179.     # Experience buffer
180.     replay_buffer = ReplayBuffer(obs_dim=obs_dim, act_dim=act_dim, size=replay_size, load_buffer=load_buffer)
181.
182.     # Count variables
183.     if proc_id() == 0:
184.         var_counts = tuple(count_vars(scope) for scope in
185.                     ['main/pi', 'main/qr1', 'main/qr2', 'main/qc1', 'main/qc2', 'main'])
186.         print((
187.                 '\nNumber of parameters: \t pi: %d, \t qr1: %d, \t qr2: %d, \t qc1: %d, \t qc2: %d, \t total: %d\n') % var_co
        unts)
188.     # Min Double-Q:
189.     min_q_pi = tf.minimum(qr1_pi, qr2_pi)
190.     min_q_pi_targ = tf.minimum(qr1_pi_targ, qr2_pi_targ)
191.
192.     # Targets for Q and V regression
193.     q_backup = tf.stop_gradient(r_ph + gamma * (1 - d_ph) * (min_q_pi_targ - alpha * logp_pi2))
194.     qc1_backup = tf.stop_gradient(c1_ph + gamma * (1 - d_ph) * qc1_pi_targ)
195.     qc2_backup = tf.stop_gradient(c2_ph + gamma * (1 - d_ph) * qc2_pi_targ)
196.
197.     cost_constraint = cost_lim * (1 - gamma ** max_ep_len) / (1 - gamma) / max_ep_len
198.     damp1 = damp_scale * tf.reduce_mean(cost_constraint - qc1)
199.     damp2 = damp_scale * tf.reduce_mean(cost_constraint - qc2)
200.
201.     # LBDSAC losses
202.     pi_loss = tf.reduce_mean(alpha * logp_pi - min_q_pi + (beta1 - damp1) * qc1_pi + (beta2 - damp2) * qc2_pi)
203.     qr1_loss = 0.5 * tf.reduce_mean((q_backup - qr1) ** 2)
204.     qr2_loss = 0.5 * tf.reduce_mean((q_backup - qr2) ** 2)
205.     qc1_loss = 0.5 * tf.reduce_mean((qc1_backup - qc1) ** 2)
206.     qc2_loss = 0.5 * tf.reduce_mean((qc2_backup - qc2) ** 2)
207.     q_loss = qr1_loss + qr2_loss + qc1_loss + qc2_loss
208.
209.     # Loss for alpha
210.     entropy_constraint *= act_dim
211.     pi_entropy = -tf.reduce_mean(logp_pi)
212.     alpha_loss = - alpha * (entropy_constraint - pi_entropy)
213.     print('using entropy constraint', entropy_constraint)
214.
215.     # Loss for beta
216.     if use_costs:
217.         if cost_constraint is None:
218.             cost_constraint = cost_lim * (1 - gamma ** max_ep_len) / (1 - gamma) / max_ep_len
```

```python
219.        print('using cost constraint', cost_constraint)
220.        beta1_loss = beta1 * (cost_constraint - qc1)
221.        beta2_loss = beta2 * (cost_constraint - qc2)
222.
223.    # Policy train op
224.    train_pi_op = MpiAdamOptimizer(learning_rate=lr).minimize(pi_loss, var_list=get_vars('main/pi'), name='train_pi'
       )
225.
226.    # Value train op
227.    with tf.control_dependencies([train_pi_op]):
228.        train_q_op = MpiAdamOptimizer(learning_rate=lr).minimize(q_loss, var_list=get_vars('main/q'), name='train_q'
           )
229.
230.    if fixed_entropy_bonus is None:
231.        entreg_optimizer = MpiAdamOptimizer(learning_rate=lr)
232.        with tf.control_dependencies([train_q_op]):
233.            train_entreg_op = entreg_optimizer.minimize(alpha_loss, var_list=get_vars('entreg'))
234.
235.    if use_costs and fixed_cost_penalty is None:
236.        costpen_optimizer = MpiAdamOptimizer(learning_rate=lr * lr_scale)
237.
238.        if fixed_entropy_bonus is None:
239.            with tf.control_dependencies([train_entreg_op]):
240.                train_costpen1_op = costpen_optimizer.minimize(beta1_loss, var_list=get_vars('costpen1'))
241.        else:
242.            with tf.control_dependencies([train_q_op]):
243.                train_costpen1_op = costpen_optimizer.minimize(beta1_loss, var_list=get_vars('costpen1'))
244.
245.        if fixed_entropy_bonus is None:
246.            with tf.control_dependencies([train_entreg_op]):
247.                train_costpen2_op = costpen_optimizer.minimize(beta2_loss, var_list=get_vars('costpen2'))
248.        else:
249.            with tf.control_dependencies([train_q_op]):
250.                train_costpen2_op = costpen_optimizer.minimize(beta2_loss, var_list=get_vars('costpen2'))
251.
252.    target_update = get_target_update('main', 'target', polyak)
253.
254.    with tf.control_dependencies([train_pi_op]):
255.        with tf.control_dependencies([train_q_op]):
256.            grouped_update = tf.group([target_update])
257.
258.    if fixed_entropy_bonus is None:
259.        grouped_update = tf.group([grouped_update, train_entreg_op])
260.    if use_costs and fixed_cost_penalty is None:
```

```python
261.        grouped_update = tf.group([grouped_update, train_costpen1_op])
262.        grouped_update = tf.group([grouped_update, train_costpen2_op])
263.
264.    def get_action(o, deterministic=False):
265.        act_op = mu if deterministic else pi
266.        return sess.run(act_op, feed_dict={x_ph: o.reshape(1, -1)})[0]
267.
268.    config = tf.ConfigProto()
269.    config.gpu_options.per_process_gpu_memory_fraction = 0.8
270.    sess = tf.Session(config=config)
271.
272.    writer = tf.summary.FileWriter(".logs/",sess.graph)
273.    # Initializing targets to match main variables
274.    target_init = get_target_update('main', 'target', 0.0)
275.    sess.run(tf.global_variables_initializer())
276.    sess.run(target_init)
277.    sess.run(sync_all_params())
278.    # Setup model saving
279.    logger.setup_tf_saver(sess, inputs={'x': x_ph, 'a': a_ph},
280.                    outputs={'mu': mu, 'pi': pi, 'qr1': qr1, 'qr2': qr2, 'qc1': qc1, 'qc2': qc2})
281.    start_time = time.time()
282.    o, r, d, ep_ret, ep_cost1, ep_cost2, ep_len= env.reset(), 0, False, 0, 0, 0, 0
283.    total_steps = steps_per_epoch * epochs * max_ep_len
284.
285.    # variables to measure in an update
286.    vars_to_get = dict(LossPi=pi_loss, LossQR1=qr1_loss, LossQR2=qr2_loss, LossQC1=qc1_loss, LossQC2=qc2_lo
    ss,
287.                    QR1Vals=qr1, QR2Vals=qr2, QC1Vals=qc1, QC2Vals=qc2, LogPi=logp_pi, PiEntropy=pi_entropy,
288.                    Alpha=alpha, LogAlpha=log_alpha, LossAlpha=alpha_loss)
289.    if use_costs:
290.        vars_to_get.update(dict(Beta1=beta1, LogBeta1=log_beta1, LossBeta1=beta1_loss,
291.                        Beta2=beta2, LogBeta2=log_beta2, LossBeta2=beta2_loss,))
292.
293.    print('starting training', proc_id())
294.    cum_cost1 = 0
295.    cum_cost2 = 0
296.    local_steps = 0
297.    local_steps_per_epoch = steps_per_epoch // num_procs()
298.    local_batch_size = batch_size // num_procs()
299.    epoch_start_time = time.time()
300.    for t in range(total_steps // num_procs()):
301.        if t > local_start_steps:
302.            a = get_action(o)
303.        else:
```

```python
304.            a = env.action_sample()
305.        # Step the env
306.        o2, r, c1, c2, d = env.step(a,o,ep_len)
307.        r *= reward_scale  # yee-haw
308.        ep_ret += r
309.        ep_cost1 += c1
310.        ep_cost2 += c2
311.        ep_len += 1
312.        local_steps += 1
313.        cum_cost1 += c1
314.        cum_cost2 += c2
315.        replay_buffer.store(o, a, r, o2, d, c1, c2)
316.        o = o2
317.        if d or (ep_len == max_ep_len):
318.            logger.store(EpRet=ep_ret, EpCost1=ep_cost1, EpCost2=ep_cost2)
319.            o, r, d, ep_ret, ep_cost1, ep_cost2, ep_len = env.reset(n = np.random.randint(n_wind_data)),\
320.                0, False, 0, 0, 0, 0
321.
322.        if t > 0 and t % update_freq == 0:
323.            for j in range(200):
324.                batch = replay_buffer.sample_batch(local_batch_size)
325.                feed_dict = {x_ph: batch['obs1'],
326.                             x2_ph: batch['obs2'],
327.                             a_ph: batch['acts'],
328.                             r_ph: batch['rews'],
329.                             c1_ph: batch['costs1'],
330.                             c2_ph: batch['costs2'],
331.                             d_ph: batch['done'],
332.                             }
333.                if t < local_update_after:
334.                    logger.store(**sess.run(vars_to_get, feed_dict))
335.                else:
336.                    values, _ = sess.run([vars_to_get, grouped_update], feed_dict)
337.                    logger.store(**values)
338.            ETA = (time.time()-epoch_start_time)*(total_steps//update_freq-t//update_freq)
339.            print(('Training: [epoch:%d|%d] [batch:%d|%d] ETA %d:%d:%d')\
340.              % (t // (local_steps_per_epoch * max_ep_len),epochs-1,\
341.                (t%(local_steps_per_epoch * max_ep_len)//max_ep_len),steps_per_epoch-1,\
342.                ETA//3600, ETA%3600//60, ETA%60))
343.            epoch_start_time = time.time()
344.
345.        # End of epoch wrap-up
346.        if t > 0 and t % (local_steps_per_epoch * max_ep_len) == 0:
347.            epoch = t // (local_steps_per_epoch * max_ep_len)
```

```python
348.         cumulative_cost1 = mpi_sum(cum_cost1)
349.         cumulative_cost2 = mpi_sum(cum_cost2)
350.         cost_rate1 = cumulative_cost1 / ((epoch + 1) * steps_per_epoch)
351.         cost_rate2 = cumulative_cost2 / ((epoch + 1) * steps_per_epoch)
352.
353.         # Save model
354.         if (epoch % save_freq == 0) or (epoch == epochs - 1):
355.             saver = tf.train.Saver()
356.             saver.save(sess,'./saved_models/Case118_%d/Model'%epoch)
357.         logger.log_tabular('Epoch', epoch)
358.         logger.log_tabular('EpRet', with_min_and_max=True)
359.         logger.log_tabular('EpCost1', with_min_and_max=True)
360.         logger.log_tabular('EpCost2', with_min_and_max=True)
361.         logger.log_tabular('CumulativeCost1', cumulative_cost1)
362.         logger.log_tabular('CumulativeCost2', cumulative_cost2)
363.         logger.log_tabular('CostRate1', cost_rate1)
364.         logger.log_tabular('CostRate2', cost_rate2)
365.         logger.log_tabular('LossPi', with_min_and_max=True)
366.         logger.log_tabular('LossQR1', with_min_and_max=True)
367.         logger.log_tabular('LossQC1', with_min_and_max=True)
368.         logger.log_tabular('LossQC2', with_min_and_max=True)
369.         logger.log_tabular('PiEntropy', with_min_and_max=True)
370.         logger.log_tabular('TotalTime', time.time() - start_time)
371.         logger.dump_tabular()
372.     writer.close()
```