📖 tensorflow / **agents**

---

Branch: master ▾     **agents** / agents / ppo / **algorithm.py**               Find file    Copy path

**vincentvanhoucke** Fix obesrv typo          104a68b 7 days ago

**2** contributors

---

559 lines (493 sloc)    23.1 KB

```python
1   # Copyright 2017 The TensorFlow Agents Authors.
2   #
3   # Licensed under the Apache License, Version 2.0 (the "License");
4   # you may not use this file except in compliance with the License.
5   # You may obtain a copy of the License at
6   #
7   #       http://www.apache.org/licenses/LICENSE-2.0
8   #
9   # Unless required by applicable law or agreed to in writing, software
10  # distributed under the License is distributed on an "AS IS" BASIS,
11  # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
12  # See the License for the specific language governing permissions and
13  # limitations under the License.
14
15  """Proximal Policy Optimization algorithm.
16
17  Based on John Schulman's implementation in Python and Theano:
18  https://github.com/joschu/modular_rl/blob/master/modular_rl/ppo.py
19  """
20
21  from __future__ import absolute_import
22  from __future__ import division
```

```python
23   from __future__ import print_function
24
25   import collections
26
27   import tensorflow as tf
28
29   from agents.ppo import memory
30   from agents.ppo import normalize
31   from agents.ppo import utility
32
33
34   _NetworkOutput = collections.namedtuple(
35       'NetworkOutput', 'policy, mean, logstd, value, state')
36
37
38   class PPOAlgorithm(object):
39     """A vectorized implementation of the PPO algorithm by John Schulman."""
40
41     def __init__(self, batch_env, step, is_training, should_log, config):
42       """Create an instance of the PPO algorithm.
43
44       Args:
45         batch_env: In-graph batch environment.
46         step: Integer tensor holding the current training step.
47         is_training: Boolean tensor for whether the algorithm should train.
48         should_log: Boolean tensor for whether summaries should be returned.
49         config: Object containing the agent configuration as attributes.
50       """
51       self._batch_env = batch_env
52       self._step = step
53       self._is_training = is_training
54       self._should_log = should_log
55       self._config = config
56       self._observ_filter = normalize.StreamingNormalize(
57           self._batch_env.observ[0], center=True, scale=True, clip=5,
```

```python
                name='normalize_observ')
        self._reward_filter = normalize.StreamingNormalize(
                self._batch_env.reward[0], center=False, scale=True, clip=10,
                name='normalize_reward')
        # Memory stores tuple of observ, action, mean, logstd, reward.
        template = (
                self._batch_env.observ[0], self._batch_env.action[0],
                self._batch_env.action[0], self._batch_env.action[0],
                self._batch_env.reward[0])
        self._memory = memory.EpisodeMemory(
                template, config.update_every, config.max_length, 'memory')
        self._memory_index = tf.Variable(0, False)
        use_gpu = self._config.use_gpu and utility.available_gpus()
        with tf.device('/gpu:0' if use_gpu else '/cpu:0'):
            # Create network variables for later calls to reuse.
            self._network(
                    tf.zeros_like(self._batch_env.observ)[:, None],
                    tf.ones(len(self._batch_env)), reuse=None)
            cell = self._config.network(self._batch_env.action.shape[1].value)
            with tf.variable_scope('ppo_temporary'):
                self._episodes = memory.EpisodeMemory(
                        template, len(batch_env), config.max_length, 'episodes')
                self._last_state = utility.create_nested_vars(
                        cell.zero_state(len(batch_env), tf.float32))
                self._last_action = tf.Variable(
                        tf.zeros_like(self._batch_env.action), False, name='last_action')
                self._last_mean = tf.Variable(
                        tf.zeros_like(self._batch_env.action), False, name='last_mean')
                self._last_logstd = tf.Variable(
                        tf.zeros_like(self._batch_env.action), False, name='last_logstd')
        self._penalty = tf.Variable(
                self._config.kl_init_penalty, False, dtype=tf.float32)
        self._policy_optimizer = self._config.policy_optimizer(
                self._config.policy_lr, name='policy_optimizer')
        self._value_optimizer = self._config.value_optimizer(
```

```python
 93            self._config.value_lr, name='value_optimizer')
 94
 95      def begin_episode(self, agent_indices):
 96        """Reset the recurrent states and stored episode.
 97
 98        Args:
 99          agent_indices: 1D tensor of batch indices for agents starting an episode.
100
101        Returns:
102          Summary tensor.
103        """
104        with tf.name_scope('begin_episode/'):
105          reset_state = utility.reinit_nested_vars(self._last_state, agent_indices)
106          reset_buffer = self._episodes.clear(agent_indices)
107          with tf.control_dependencies([reset_state, reset_buffer]):
108            return tf.constant('')
109
110      def perform(self, observ):
111        """Compute batch of actions and a summary for a batch of observation.
112
113        Args:
114          observ: Tensor of a batch of observations for all agents.
115
116        Returns:
117          Tuple of action batch tensor and summary tensor.
118        """
119        with tf.name_scope('perform/'):
120          observ = self._observ_filter.transform(observ)
121          network = self._network(
122              observ[:, None], tf.ones(observ.shape[0]), self._last_state)
123          action = tf.cond(
124              self._is_training, network.policy.sample, lambda: network.mean)
125          logprob = network.policy.log_prob(action)[:, 0]
126          # pylint: disable=g-long-lambda
127          summary = tf.cond(self._should_log, lambda: tf.summary.merge([
```

```python
128            tf.summary.histogram('mean', network.mean[:, 0]),
129            tf.summary.histogram('std', tf.exp(network.logstd[:, 0])),
130            tf.summary.histogram('action', action[:, 0]),
131            tf.summary.histogram('logprob', logprob)]), str)
132        # Remember current policy to append to memory in the experience callback.
133        with tf.control_dependencies([
134            utility.assign_nested_vars(self._last_state, network.state),
135            self._last_action.assign(action[:, 0]),
136            self._last_mean.assign(network.mean[:, 0]),
137            self._last_logstd.assign(network.logstd[:, 0])]):
138          return tf.check_numerics(action[:, 0], 'action'), tf.identity(summary)
139
140    def experience(self, observ, action, reward, unused_done, unused_nextob):
141      """Process the transition tuple of the current step.
142
143      When training, add the current transition tuple to the memory and update
144      the streaming statistics for observations and rewards. A summary string is
145      returned if requested at this step.
146
147      Args:
148        observ: Batch tensor of observations.
149        action: Batch tensor of actions.
150        reward: Batch tensor of rewards.
151        unused_done: Batch tensor of done flags.
152        unused_nextob: Batch tensor of successor observations.
153
154      Returns:
155        Summary tensor.
156      """
157      with tf.name_scope('experience/'):
158        return tf.cond(
159            self._is_training,
160            lambda: self._define_experience(observ, action, reward), str)
161
162    def _define_experience(self, observ, action, reward):
```

```python
163        """Implement the branch of experience() entered during training."""
164        update_filters = tf.summary.merge([
165            self._observ_filter.update(observ),
166            self._reward_filter.update(reward)])
167      with tf.control_dependencies([update_filters]):
168        if self._config.train_on_agent_action:
169          # NOTE: Doesn't seem to change much.
170          action = self._last_action
171        batch = observ, action, self._last_mean, self._last_logstd, reward
172        append = self._episodes.append(batch, tf.range(len(self._batch_env)))
173      with tf.control_dependencies([append]):
174        norm_observ = self._observ_filter.transform(observ)
175        norm_reward = tf.reduce_mean(self._reward_filter.transform(reward))
176        # pylint: disable=g-long-lambda
177        summary = tf.cond(self._should_log, lambda: tf.summary.merge([
178            update_filters,
179            self._observ_filter.summary(),
180            self._reward_filter.summary(),
181            tf.summary.scalar('memory_size', self._memory_index),
182            tf.summary.histogram('normalized_observ', norm_observ),
183            tf.summary.histogram('action', self._last_action),
184            tf.summary.scalar('normalized_reward', norm_reward)]), str)
185        return summary
186
187    def end_episode(self, agent_indices):
188      """Add episodes to the memory and perform update steps if memory is full.
189
190      During training, add the collected episodes of the batch indices that
191      finished their episode to the memory. If the memory is full, train on it,
192      and then clear the memory. A summary string is returned if requested at
193      this step.
194
195      Args:
196        agent_indices: 1D tensor of batch indices for agents starting an episode.
197
```

```python
198      Returns:
199        Summary tensor.
200      """
201    with tf.name_scope('end_episode/'):
202      return tf.cond(
203          self._is_training,
204          lambda: self._define_end_episode(agent_indices), str)

206  def _define_end_episode(self, agent_indices):
207    """Implement the branch of end_episode() entered during training."""
208    episodes, length = self._episodes.data(agent_indices)
209    space_left = self._config.update_every - self._memory_index
210    use_episodes = tf.range(tf.minimum(
211        tf.shape(agent_indices)[0], space_left))
212    episodes = [tf.gather(elem, use_episodes) for elem in episodes]
213    append = self._memory.replace(
214        episodes, tf.gather(length, use_episodes),
215        use_episodes + self._memory_index)
216    with tf.control_dependencies([append]):
217      inc_index = self._memory_index.assign_add(tf.shape(use_episodes)[0])
218    with tf.control_dependencies([inc_index]):
219      memory_full = self._memory_index >= self._config.update_every
220      return tf.cond(memory_full, self._training, str)

222  def _training(self):
223    """Perform multiple training iterations of both policy and value baseline.
224
225    Training on the episodes collected in the memory. Reset the memory
226    afterwards. Always returns a summary string.
227
228    Returns:
229      Summary tensor.
230    """
231    with tf.name_scope('training'):
232      assert_full = tf.assert_equal(
```

```python
233              self._memory_index, self._config.update_every)
234          with tf.control_dependencies([assert_full]):
235            data = self._memory.data()
236          (observ, action, old_mean, old_logstd, reward), length = data
237          with tf.control_dependencies([tf.assert_greater(length, 0)]):
238            length = tf.identity(length)
239          observ = self._observ_filter.transform(observ)
240          reward = self._reward_filter.transform(reward)
241          policy_summary = self._update_policy(
242              observ, action, old_mean, old_logstd, reward, length)
243          with tf.control_dependencies([policy_summary]):
244            value_summary = self._update_value(observ, reward, length)
245          with tf.control_dependencies([value_summary]):
246            penalty_summary = self._adjust_penalty(
247                observ, old_mean, old_logstd, length)
248          with tf.control_dependencies([penalty_summary]):
249            clear_memory = tf.group(
250                self._memory.clear(), self._memory_index.assign(0))
251          with tf.control_dependencies([clear_memory]):
252            weight_summary = utility.variable_summaries(
253                tf.trainable_variables(), self._config.weight_summaries)
254            return tf.summary.merge([
255                policy_summary, value_summary, penalty_summary, weight_summary])
256
257    def _update_value(self, observ, reward, length):
258      """Perform multiple update steps of the value baseline.
259
260      We need to decide for the summary of one iteration, and thus choose the one
261      after half of the iterations.
262
263      Args:
264        observ: Sequences of observations.
265        reward: Sequences of reward.
266        length: Batch of sequence lengths.
267
```

```python
268        Returns:
269          Summary tensor.
270        """
271      with tf.name_scope('update_value'):
272        loss, summary = tf.scan(
273            lambda _1, _2: self._update_value_step(observ, reward, length),
274            tf.range(self._config.update_epochs_value),
275            [0., ''], parallel_iterations=1)
276        print_loss = tf.Print(0, [tf.reduce_mean(loss)], 'value loss: ')
277        with tf.control_dependencies([loss, print_loss]):
278          return summary[self._config.update_epochs_value // 2]
279
280    def _update_value_step(self, observ, reward, length):
281      """Compute the current value loss and perform a gradient update step.
282
283      Args:
284        observ: Sequences of observations.
285        reward: Sequences of reward.
286        length: Batch of sequence lengths.
287
288      Returns:
289        Tuple of loss tensor and summary tensor.
290      """
291      loss, summary = self._value_loss(observ, reward, length)
292      gradients, variables = (
293          zip(*self._value_optimizer.compute_gradients(loss)))
294      optimize = self._value_optimizer.apply_gradients(
295          zip(gradients, variables))
296      summary = tf.summary.merge([
297          summary,
298          tf.summary.scalar('gradient_norm', tf.global_norm(gradients)),
299          utility.gradient_summaries(
300              zip(gradients, variables), dict(value=r'.*'))])
301      with tf.control_dependencies([optimize]):
302        return [tf.identity(loss), tf.identity(summary)]
```

```python
303
304    def _value_loss(self, observ, reward, length):
305      """Compute the loss function for the value baseline.
306
307      The value loss is the difference between empirical and approximated returns
308      over the collected episodes. Returns the loss tensor and a summary strin.
309
310      Args:
311        observ: Sequences of observations.
312        reward: Sequences of reward.
313        length: Batch of sequence lengths.
314
315      Returns:
316        Tuple of loss tensor and summary tensor.
317      """
318      with tf.name_scope('value_loss'):
319        value = self._network(observ, length).value
320        return_ = utility.discounted_return(
321            reward, length, self._config.discount)
322        advantage = return_ - value
323        value_loss = 0.5 * self._mask(advantage ** 2, length)
324        summary = tf.summary.merge([
325            tf.summary.histogram('value_loss', value_loss),
326            tf.summary.scalar('avg_value_loss', tf.reduce_mean(value_loss))])
327        value_loss = tf.reduce_mean(value_loss)
328        return tf.check_numerics(value_loss, 'value_loss'), summary
329
330    def _update_policy(
331        self, observ, action, old_mean, old_logstd, reward, length):
332      """Perform multiple update steps of the policy.
333
334      The advantage is computed once at the beginning and shared across
335      iterations. We need to decide for the summary of one iteration, and thus
336      choose the one after half of the iterations.
337
```

```
338        Args:
339          observ: Sequences of observations.
340          action: Sequences of actions.
341          old_mean: Sequences of action means of the behavioral policy.
342          old_logstd: Sequences of action log stddevs of the behavioral policy.
343          reward: Sequences of rewards.
344          length: Batch of sequence lengths.
345
346        Returns:
347          Summary tensor.
348        """
349        with tf.name_scope('update_policy'):
350          return_ = utility.discounted_return(
351              reward, length, self._config.discount)
352          value = self._network(observ, length).value
353          if self._config.gae_lambda:
354            advantage = utility.lambda_return(
355                reward, value, length, self._config.discount,
356                self._config.gae_lambda)
357          else:
358            advantage = return_ - value
359          mean, variance = tf.nn.moments(advantage, axes=[0, 1], keep_dims=True)
360          advantage = (advantage - mean) / (tf.sqrt(variance) + 1e-8)
361          advantage = tf.Print(
362              advantage, [tf.reduce_mean(return_), tf.reduce_mean(value)],
363              'return and value: ')
364          advantage = tf.Print(
365              advantage, [tf.reduce_mean(advantage)],
366              'normalized advantage: ')
367          # pylint: disable=g-long-lambda
368          loss, summary = tf.scan(
369              lambda _1, _2: self._update_policy_step(
370                  observ, action, old_mean, old_logstd, advantage, length),
371              tf.range(self._config.update_epochs_policy),
372              [0., ''], parallel_iterations=1)
```

```python
373            print_loss = tf.Print(0, [tf.reduce_mean(loss)], 'policy loss: ')
374          with tf.control_dependencies([loss, print_loss]):
375            return summary[self._config.update_epochs_policy // 2]
376
377      def _update_policy_step(
378          self, observ, action, old_mean, old_logstd, advantage, length):
379        """Compute the current policy loss and perform a gradient update step.
380
381        Args:
382          observ: Sequences of observations.
383          action: Sequences of actions.
384          old_mean: Sequences of action means of the behavioral policy.
385          old_logstd: Sequences of action log stddevs of the behavioral policy.
386          advantage: Sequences of advantages.
387          length: Batch of sequence lengths.
388
389        Returns:
390          Tuple of loss tensor and summary tensor.
391        """
392        network = self._network(observ, length)
393        loss, summary = self._policy_loss(
394            network.mean, network.logstd, old_mean, old_logstd, action,
395            advantage, length)
396        gradients, variables = (
397            zip(*self._policy_optimizer.compute_gradients(loss)))
398        optimize = self._policy_optimizer.apply_gradients(
399            zip(gradients, variables))
400        summary = tf.summary.merge([
401            summary,
402            tf.summary.scalar('gradient_norm', tf.global_norm(gradients)),
403            utility.gradient_summaries(
404                zip(gradients, variables), dict(policy=r'.*'))])
405        with tf.control_dependencies([optimize]):
406          return [tf.identity(loss), tf.identity(summary)]
407
```

```python
408    def _policy_loss(
409        self, mean, logstd, old_mean, old_logstd, action, advantage, length):
410      """Compute the policy loss composed of multiple components.
411
412      1. The policy gradient loss is importance sampled from the data-collecting
413         policy at the beginning of training.
414      2. The second term is a KL penalty between the policy at the beginning of
415         training and the current policy.
416      3. Additionally, if this KL already changed more than twice the target
417         amount, we activate a strong penalty discouraging further divergence.
418
419      Args:
420        mean: Sequences of action means of the current policy.
421        logstd: Sequences of action log stddevs of the current policy.
422        old_mean: Sequences of action means of the behavioral policy.
423        old_logstd: Sequences of action log stddevs of the behavioral policy.
424        action: Sequences of actions.
425        advantage: Sequences of advantages.
426        length: Batch of sequence lengths.
427
428      Returns:
429        Tuple of loss tensor and summary tensor.
430      """
431      with tf.name_scope('policy_loss'):
432        entropy = utility.diag_normal_entropy(mean, logstd)
433        kl = tf.reduce_mean(self._mask(utility.diag_normal_kl(
434            old_mean, old_logstd, mean, logstd), length), 1)
435        policy_gradient = tf.exp(
436            utility.diag_normal_logpdf(mean, logstd, action) -
437            utility.diag_normal_logpdf(old_mean, old_logstd, action))
438        surrogate_loss = -tf.reduce_mean(self._mask(
439            policy_gradient * tf.stop_gradient(advantage), length), 1)
440        kl_penalty = self._penalty * kl
441        cutoff_threshold = self._config.kl_target * self._config.kl_cutoff_factor
442        cutoff_count = tf.reduce_sum(
```

```python
443                 tf.cast(kl > cutoff_threshold, tf.int32))
444           with tf.control_dependencies([tf.cond(
445               cutoff_count > 0,
446               lambda: tf.Print(0, [cutoff_count], 'kl cutoff! '), int)]):
447             kl_cutoff = (
448                 self._config.kl_cutoff_coef *
449                 tf.cast(kl > cutoff_threshold, tf.float32) *
450                 (kl - cutoff_threshold) ** 2)
451         policy_loss = surrogate_loss + kl_penalty + kl_cutoff
452         summary = tf.summary.merge([
453             tf.summary.histogram('entropy', entropy),
454             tf.summary.histogram('kl', kl),
455             tf.summary.histogram('surrogate_loss', surrogate_loss),
456             tf.summary.histogram('kl_penalty', kl_penalty),
457             tf.summary.histogram('kl_cutoff', kl_cutoff),
458             tf.summary.histogram('kl_penalty_combined', kl_penalty + kl_cutoff),
459             tf.summary.histogram('policy_loss', policy_loss),
460             tf.summary.scalar('avg_surr_loss', tf.reduce_mean(surrogate_loss)),
461             tf.summary.scalar('avg_kl_penalty', tf.reduce_mean(kl_penalty)),
462             tf.summary.scalar('avg_policy_loss', tf.reduce_mean(policy_loss))])
463         policy_loss = tf.reduce_mean(policy_loss, 0)
464         return tf.check_numerics(policy_loss, 'policy_loss'), summary
465
466     def _adjust_penalty(self, observ, old_mean, old_logstd, length):
467       """Adjust the KL policy between the behavioral and current policy.
468
469       Compute how much the policy actually changed during the multiple
470       update steps. Adjust the penalty strength for the next training phase if we
471       overshot or undershot the target divergence too much.
472
473       Args:
474         observ: Sequences of observations.
475         old_mean: Sequences of action means of the behavioral policy.
476         old_logstd: Sequences of action log stddevs of the behavioral policy.
477         length: Batch of sequence lengths.
```

```python
478
479      Returns:
480        Summary tensor.
481      """
482    with tf.name_scope('adjust_penalty'):
483      network = self._network(observ, length)
484      assert_change = tf.assert_equal(
485          tf.reduce_all(tf.equal(network.mean, old_mean)), False,
486          message='policy should change')
487      print_penalty = tf.Print(0, [self._penalty], 'current penalty: ')
488      with tf.control_dependencies([assert_change, print_penalty]):
489        kl_change = tf.reduce_mean(self._mask(utility.diag_normal_kl(
490            old_mean, old_logstd, network.mean, network.logstd), length))
491        kl_change = tf.Print(kl_change, [kl_change], 'kl change: ')
492        maybe_increase = tf.cond(
493            kl_change > 1.3 * self._config.kl_target,
494            # pylint: disable=g-long-lambda
495            lambda: tf.Print(self._penalty.assign(
496                self._penalty * 1.5), [0], 'increase penalty '),
497            float)
498        maybe_decrease = tf.cond(
499            kl_change < 0.7 * self._config.kl_target,
500            # pylint: disable=g-long-lambda
501            lambda: tf.Print(self._penalty.assign(
502                self._penalty / 1.5), [0], 'decrease penalty '),
503            float)
504      with tf.control_dependencies([maybe_increase, maybe_decrease]):
505        return tf.summary.merge([
506            tf.summary.scalar('kl_change', kl_change),
507            tf.summary.scalar('penalty', self._penalty)])
508
509  def _mask(self, tensor, length):
510    """Set padding elements of a batch of sequences to zero.
511
512    Useful to then safely sum along the time dimension.
```

```
513
514        Args:
515          tensor: Tensor of sequences.
516          length: Batch of sequence lengths.
517
518        Returns:
519          Masked sequences.
520        """
521        with tf.name_scope('mask'):
522          range_ = tf.range(tensor.shape[1].value)
523          mask = tf.cast(range_[None, :] < length[:, None], tf.float32)
524          masked = tensor * mask
525          return tf.check_numerics(masked, 'masked')
526
527      def _network(self, observ, length=None, state=None, reuse=True):
528        """Compute the network output for a batched sequence of observations.
529
530        Optionally, the initial state can be specified. The weights should be
531        reused for all calls, except for the first one. Output is a named tuple
532        containing the policy as a TensorFlow distribution, the policy mean and log
533        standard deviation, the approximated state value, and the new recurrent
534        state.
535
536        Args:
537          observ: Sequences of observations.
538          length: Batch of sequence lengths.
539          state: Batch of initial recurrent states.
540          reuse: Python boolean whether to reuse previous variables.
541
542        Returns:
543          NetworkOutput tuple.
544        """
545        with tf.variable_scope('network', reuse=reuse):
546          observ = tf.convert_to_tensor(observ)
547          use_gpu = self._config.use_gpu and utility.available_gpus()
```

```python
548          with tf.device('/gpu:0' if use_gpu else '/cpu:0'):
549            observ = tf.check_numerics(observ, 'observ')
550            cell = self._config.network(self._batch_env.action.shape[1].value)
551            (mean, logstd, value), state = tf.nn.dynamic_rnn(
552                cell, observ, length, state, tf.float32, swap_memory=True)
553          mean = tf.check_numerics(mean, 'mean')
554          logstd = tf.check_numerics(logstd, 'logstd')
555          value = tf.check_numerics(value, 'value')
556          policy = tf.contrib.distributions.MultivariateNormalDiag(
557              mean, tf.exp(logstd))
558          return _NetworkOutput(policy, mean, logstd, value, state)
```