

---

# **actor-critic Documentation**

*Release 0.1*

**Jan Robine**

**Aug 08, 2018**



---

## Contents

---

<b>1</b>	<b>API Documentation</b>	<b>3</b>
<b>2</b>	<b>Quickstart Guide</b>	<b>29</b>
	<b>Python Module Index</b>	<b>33</b>



Documentation of the *actor-critic* repository on [GitHub](#).



---

<i>actorcritic</i>	The root package.
--------------------	-------------------

---

## 1.1 actorcritic

The root package.

<i>agents</i>	Contains <i>agents</i> , which are an abstraction from environments.
<i>baselines</i>	Contains <i>baselines</i> , which are used to compute the <i>advantage</i> .
<i>kfac_utils</i>	Contains utilities that concern K-FAC.
<i>model</i>	Contains the base class of actor-critic models.
<i>multi_env</i>	Contains classes that provide the ability to run multiple environments in subprocesses.
<i>nn</i>	Contains utilities that concern TensorFlow and neural networks.
<i>objectives</i>	Contains <i>objectives</i> that are used to optimize actor-critic models.
<i>policies</i>	Contains <i>policies</i> that determine the behavior of an <i>agent</i> .
<i>envs</i>	Contains functions that are dedicated to certain environments.
<i>examples</i>	Contains examples of how to use this project.

### 1.1.1 actorcritic.agents

Contains *agents*, which are an abstraction from environments.

## Functions

---

<code>transpose_list(values)</code>	Transposes a list of lists.
-------------------------------------	-----------------------------

---

## Classes

---

<code>Agent</code>	Takes environments and a model (containing a policy) and provides <code>interact()</code> , which manages operations such as selecting actions from the model and stepping in the environments.
<code>MultiEnvAgent(multi_env, model, num_steps)</code>	An agent that maintains multiple environments (via <code>MultiEnv</code> ) and samples multiple steps.
<code>SingleEnvAgent(env, model, num_steps)</code>	An agent that maintains a single environment and samples multiple steps.

---

**class** actorcritic.agents.**Agent**

Bases: `object`

Takes environments and a model (containing a policy) and provides `interact()`, which manages operations such as selecting actions from the model and stepping in the environments.

**See also:**

This allows to create multi-step agents, like `SingleEnvAgent` and `MultiEnvAgent`.

**interact** (*session*)

Samples actions from the model, and steps in the environments.

**Parameters** `session` (`tf.Session`) – A session that will be used to compute the actions.

**Returns**

`tuple` – A tuple of (*observations*, *actions*, *rewards*, *terminals*, *next\_observations*, *infos*).

All values are in *batch-major* format, meaning that the rows determine the batch and the columns determine the time: [*batch*, *time*]. In our case the rows correspond to the environments and the columns correspond to the steps: [*environment*, *step*]. The opposite is the *time-major* format: [*time*, *batch*] or [*step*, *environment*].

Example:

If the agent maintains 3 environments and samples for 5 steps, the result would consist of a matrix (`list of list`) with shape [3, 5]:

```
[ [step 1, step 2, step 3, step 4, step 5], # environment 1
  [step 1, step 2, step 3, step 4, step 5], # environment 2
  [step 1, step 2, step 3, step 4, step 5] ] # environment 3
```

*observations*, *actions*, *rewards*, *terminals*, and *infos* are collected during sampling and have the shape [*environments*, *steps*].

*next\_observations* contains the observations that the agent received at last, but did not use for selecting actions yet. These e.g. can be used to bootstrap the remaining returns. Has the shape [*environments*, 1].

**class** actorcritic.agents.**MultiEnvAgent** (*multi\_env*, *model*, *num\_steps*)

Bases: `actorcritic.agents.Agent`

An agent that maintains multiple environments (via *MultiEnv*) and samples multiple steps.

`__init__(multi_env, model, num_steps)`

#### Parameters

- **multi\_env** (*MultiEnv*) – Multiple environments.
- **model** (*ActorCriticModel*) – A model to sample actions.
- **num\_steps** (*int*) – The number of steps to take in *interact()*.

**interact** (*session*)

Samples actions from the model, and steps in the environments.

**Parameters** **session** (*tf.Session*) – A session that will be used to compute the actions.

#### Returns

*tuple* – A tuple of (*observations, actions, rewards, terminals, next\_observations, infos*).

All values are in *batch-major* format, meaning that the rows determine the batch and the columns determine the time: [*batch, time*]. In our case the rows correspond to the environments and the columns correspond to the steps: [*environment, step*]. The opposite is the *time-major* format: [*time, batch*] or [*step, environment*].

Example:

If the agent maintains 3 environments and samples for 5 steps, the result would consist of a matrix (*list of list*) with shape [3, 5]:

```
[ [step 1, step 2, step 3, step 4, step 5], # environment 1
  [step 1, step 2, step 3, step 4, step 5], # environment 2
  [step 1, step 2, step 3, step 4, step 5] ] # environment 3
```

*observations, actions, rewards, terminals, and infos* are collected during sampling and have the shape [*environments, steps*].

*next\_observations* contains the observations that the agent received at last, but did not use for selecting actions yet. These e.g. can be used to bootstrap the remaining returns. Has the shape [*environments, 1*].

**class** `actorcritic.agents.SingleEnvAgent` (*env, model, num\_steps*)

Bases: *actorcritic.agents.Agent*

An agent that maintains a single environment and samples multiple steps.

`__init__(env, model, num_steps)`

#### Parameters

- **env** (*gym.Env*) – An environment.
- **model** (*ActorCriticModel*) – A model to sample actions.
- **num\_steps** (*int*) – The number of steps to take in *interact()*.

**interact** (*session*)

Samples actions from the model and steps in the environment.

**Parameters** **session** (*tf.Session*) – A session that will be used to compute the actions.

#### Returns

*tuple* – A tuple (*observations, actions, rewards, terminals, next\_observations, infos*).

All values are in *batch-major* format, meaning that the rows determine the batch and the columns determine the time:  $[batch, time]$ . In our case we have *one* environment so the row corresponds to the environment and the columns correspond to the steps:  $[I, step]$ . The opposite is the *time-major* format:  $[time, batch]$  or  $[step, I]$ .

*observations*, *actions*, *rewards*, *terminals*, and *infos* are collected during sampling and have the shape  $[I, steps]$ .

*next\_observations* contains the observation that the agent received at last, but did not use for selecting an action yet. This e.g. can be used to bootstrap the remaining return. Has the shape  $[I, I]$ .

`actorcritic.agents.transpose_list` (*values*)

Transposes a list of lists. Can be used to convert from *time-major* format to *batch-major* format and vice versa.

### Example

Input:

```
[1, 2, 3, 4],
[5, 6, 7, 8],
[9, 10, 11, 12]]
```

Output:

```
[1, 5, 9],
[2, 6, 10],
[3, 7, 11],
[4, 8, 12]]
```

**Parameters** *values* (*list of list*) – Values to transpose.

**Returns** *list of list* – The transposed values.

## 1.1.2 actorcritic.baselines

Contains *baselines*, which are used to compute the *advantage*.

### Classes

---

<i>Baseline</i>	A wrapper class for the baseline that is subtracted from the target values to get the <i>advantage</i> .
<i>StateValueFunction</i> ( <i>value</i> )	A baseline defined by a state-value function.

---

**class** `actorcritic.baselines.Baseline`

Bases: `object`

A wrapper class for the baseline that is subtracted from the target values to get the *advantage*.

**register\_predictive\_distribution** (*layer\_collection*, *random\_seed=None*)

Registers the predictive distribution of this baseline in the specified `kfac.LayerCollection` (required for K-FAC).

**Parameters**

- **layer\_collection** (`kfac.LayerCollection`) – A layer collection used by the `KfacOptimizer`.
- **random\_seed** (`int`, optional) – A random seed for sampling from the predictive distribution.

**Raises** `NotImplementedError` – If this baseline does not support K-FAC.

**value**

`tf.Tensor` – The output values of this baseline.

**class** `actorcritic.baselines.StateValueFunction` (*value*)

Bases: `actorcritic.baselines.Baseline`

A baseline defined by a state-value function.

**\_\_init\_\_** (*value*)

**Parameters** **value** (`tf.Tensor`) – The output values of this state-value function.

**register\_predictive\_distribution** (*layer\_collection*, *random\_seed=None*)

Registers the predictive distribution (normal distribution) of this state-value function in the specified `kfac.LayerCollection` (required for K-FAC).

**Parameters**

- **layer\_collection** (`kfac.LayerCollection`) – A layer collection used by the `KfacOptimizer`.
- **random\_seed** (`int`, optional) – A random seed for sampling from the predictive distribution.

**value**

`tf.Tensor` – The output values of this state-value function.

### 1.1.3 actorcritic.kfac\_utils

Contains utilities that concern K-FAC.

#### Classes

---

<code>ColdStartPeriodicInvUpdateKfacOpt(...)</code>	A modified <code>KfacOptimizer</code> that runs the inverse operation periodically and uses a standard SGD optimizer for a few updates in the beginning, called <i>cold updates</i> and <i>cold optimizer</i> .
---	---

---

**class** `actorcritic.kfac_utils.ColdStartPeriodicInvUpdateKfacOpt` (*num\_cold\_updates*, *cold\_optimizer*, *invert\_every*, *\*\*kwargs*)

Bases: `kfac.python.ops.optimizer.KfacOptimizer`

A modified `KfacOptimizer` that runs the inverse operation periodically and uses a standard SGD optimizer for a few updates in the beginning, called *cold updates* and *cold optimizer*.

This can be used to slowly initialize the parameters in the beginning before using the heavy K-FAC optimizer. The covariances get updated every step (after the *cold updates*).

**See also:**

- `kfac.PeriodicInvCovUpdateKfacOpt`
- The idea is taken from the [original ACKTR implementation](#).

`__init__` (*num\_cold\_updates*, *cold\_optimizer*, *invert\_every*, *\*\*kwargs*)

#### Parameters

- **num\_cold\_updates** (*int*) – The number of *cold updates* in the beginning before using the actual K-FAC optimizer.
- **cold\_optimizer** (*tf.train.Optimizer*) – An optimizer that is used for the *cold updates*.
- **invert\_every** (*int*) – The inverse operation gets called every *invert\_every* steps (after the *cold updates* have finished).

**apply\_gradients** (*grads\_and\_vars*, *global\_step=None*, *name=None*)

Applies gradients to variables.

#### Parameters

- **grads\_and\_vars** – List of (gradient, variable) pairs.
- **\*args** – Additional arguments for `super.apply_gradients`.
- **\*\*kwargs** – Additional keyword arguments for `super.apply_gradients`.

**Returns** An *Operation* that applies the specified gradients.

## 1.1.4 actorcritic.model

Contains the base class of actor-critic models.

### Classes

---

`ActorCriticModel`(*observation\_space*, *ac-* *tion\_space*) Represents a model (e.g.

---

**class** `actorcritic.model.ActorCriticModel` (*observation\_space*, *action\_space*)

Bases: `object`

Represents a model (e.g. a neural net) that provides the functionalities required for actor-critic algorithms. Provides a policy, a baseline (that is subtracted from the target values to compute the advantage) and the values used for bootstrapping from next observations (ideally the values of the baseline), and the placeholders.

`__init__` (*observation\_space*, *action\_space*)

#### Parameters

- **observation\_space** (*gym.spaces.Space*) – A space that determines the shape of the *observations\_placeholder* and the *bootstrap\_observations\_placeholder*.
- **action\_space** (*gym.spaces.Space*) – A space that determines the shape of the *actions\_placeholder*.

**actions\_placeholder**

*tf.Tensor* – The placeholder for the sampled actions.

**baseline**

*Baseline* – The baseline used by this model.

**bootstrap\_observations\_placeholder**

`tf.Tensor` – The placeholder for the sampled next observations. These are used to compute the *bootstrap\_values*.

**bootstrap\_values**

`tf.Tensor` – The bootstrapped values that are computed based on the observations passed to the *bootstrap\_observations\_placeholder*.

**observations\_placeholder**

`tf.Tensor` – The placeholder for the sampled observations.

**policy**

*Policy* – The policy used by this model.

**register\_layers** (*layer\_collection*)

Registers the layers of this model (neural net) in the specified `kfac.LayerCollection` (required for K-FAC).

**Parameters** `layer_collection` (`kfac.LayerCollection`) – A layer collection used by the `KfacOptimizer`.

**Raises** `NotImplementedError` – If this model does not support K-FAC.

**register\_predictive\_distributions** (*layer\_collection*, *random\_seed=None*)

Registers the predictive distributions of the policy and the baseline in the specified `kfac.LayerCollection` (required for K-FAC).

**Parameters**

- **layer\_collection** (`kfac.LayerCollection`) – A layer collection used by the `KfacOptimizer`.
- **random\_seed** (`int`, optional) – A random seed used for sampling from the predictive distributions.

**rewards\_placeholder**

`tf.Tensor` – The placeholder for the sampled rewards (scalars).

**sample\_actions** (*observations*, *session*)

Samples actions from the policy based on the specified observations.

**Parameters**

- **observations** – The observations that will be passed to the *observations\_placeholder*.
- **session** (`tf.Session`) – A session that will be used to compute the values.

**Returns** `list of list` – A list of lists of actions. The shape equals the shape of *observations*.

**select\_max\_actions** (*observations*, *session*)

Selects actions from the policy that have the highest probability (mode) based on the specified observations.

**Parameters**

- **observations** – The observations that will be passed to the *observations\_placeholder*.
- **session** (`tf.Session`) – A session that will be used to compute the values.

**Returns** `list of list` – A list of lists of actions. The shape equals the shape of *observations*.

**terminals\_placeholder**

`tf.Tensor` – The placeholder for the sampled terminals (booleans).

### 1.1.5 actorcritic.multi\_env

Contains classes that provide the ability to run multiple environments in subprocesses.

#### Functions

---

<code>create_subprocess_envs(env_fns)</code>	Utility function that creates environments by calling the functions in <code>env_fns</code> and wrapping the returned environments in <code>SubprocessEnvs</code> .
--	---

---

#### Classes

---

<code>MultiEnv(envs)</code>	An environment that maintains multiple <code>SubprocessEnvs</code> and executes them in parallel.
<code>SubprocessEnv(env_fn)</code>	Maintains a <code>gym.Env</code> inside a subprocess, so it can run concurrently.

---

**class** `actorcritic.multi_env.MultiEnv(envs)`

Bases: `object`

An environment that maintains multiple `SubprocessEnvs` and executes them in parallel.

The environments will be reset automatically when a terminal state is reached. That means that `reset()` actually only has to be called once in the beginning.

`__init__(envs)`

**Parameters** `envs` (`list` of `SubprocessEnv`) – The environments. The observation and action spaces must be equal across the environments.

**action\_space**

`gym.spaces.Space` – The action space used by all environments.

**close()**

Closes all environments.

**envs**

`list` of `gym.Env` – The environments.

**observation\_space**

`gym.spaces.Space` – The observation space used by all environments.

**reset()**

Resets all environments.

**Returns** `list` – A list of observations received from each environment.

**step(actions)**

Proceeds one step in all environments.

**Parameters** `actions` (`list`) – A list of actions to be executed in the environments.

**Returns** `tuple` – A tuple of (*observations*, *rewards*, *terminals*, *infos*). Each element is a list containing the values received from the environments.

**class** `actorcritic.multi_env.SubprocessEnv` (*env\_fn*)

Bases: `gym.core.Env`

Maintains a `gym.Env` inside a subprocess, so it can run concurrently. If the subprocess ends unexpectedly, it will be recreated automatically without interrupting the execution.

To use the subprocess `start()` has to be called first. After that `initialize()` has to be called to retrieve the observation space and the action space from the underlying environment. The purpose of these methods is that multiple `SubprocessEnvs` can be created and started in parallel without blocking the execution, which creates the underlying `gym.Env` already. Afterwards `start()`, which blocks the execution, can be called in parallel. See `create_subprocess_envs()` which implements this idea.

`__init__` (*env\_fn*)

**Parameters** `env_fn` (`callable`) – A function that returns a `gym.Env`. It will be called inside the subprocess, so watch out for referencing variables on the main process or the like. It possibly will be called multiple times, since the subprocess will be recreated when it unexpectedly ends.

**action\_space**

`gym.spaces.Space` – The action space of the underlying environment. Does not block the execution. `start()` and `initialize()` must have been called.

**close** ()

Closes the subprocess.

**initialize** ()

Retrieves the observation space and the action space from the environment in the subprocess. This method blocks until the execution is finished. `start()` must have been called.

**observation\_space**

`gym.spaces.Space` – The observation space of the underlying environment. Does not block the execution. `start()` and `initialize()` must have been called.

**render** (*mode='human'*)

Remotely calls `gym.Env.render()` in the subprocess. This method blocks until execution is finished. `start()` and `initialize()` must have been called.

**Parameters** `mode` (`str`) – The *mode* argument passed to `gym.Env.render()`.

**Returns** The value returned by `gym.Env.render()`.

**reset** (\*\**kwargs*)

Remotely calls `gym.Env.reset()` in the underlying environment. This method blocks until execution is finished. `start()` and `initialize()` must have been called.

**Parameters** `kwargs` (`dict`) – Keyword arguments passed to `gym.Env.reset()`.

**Returns** The value returned by `gym.Env.reset()`.

**start** ()

Starts the subprocess. Does not block. You should call `initialize()` afterwards.

**step** (*action*)

Remotely calls `gym.Env.step()` in the underlying environment. This method blocks until execution is finished. `start()` and `initialize()` must have been called.

**Parameters** `action` – The *action* argument passed to `gym.Env.step()`.

**Returns** `tuple` – A tuple of (*observation*, *reward*, *terminal*, *info*). The values returned by `gym.Env.step()`.

**class** `actorcritic.multi_env._AutoResetWrapper` (*env*)

Bases: `gym.core.Wrapper`

**reset** (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns:** **observation (object): the initial observation of the** space.

**step** (*action*)

Run one timestep of the environment's dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment's state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** **action (object)** – an action provided by the environment

**Returns** *observation (object)* – agent's observation of the current environment reward (float) : amount of reward returned after previous action done (boolean): whether the episode has ended, in which case further step() calls will return undefined results info (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

`actorcritic.multi_env.create_subprocess_envs` (*env\_fns*)

Utility function that creates environments by calling the functions in *env\_fns* and wrapping the returned environments in `SubprocessEnvs`. They will be started and initialized in parallel.

**Parameters** **env\_fns (list of callable)** – A list of functions that return a `gym.Env`. They should not be instances of `SubprocessEnv`.

**Returns** *list of SubprocessEnv* – A list of the created environments.

## 1.1.6 actorcritic.nn

Contains utilities that concern TensorFlow and neural networks.

### Functions

<code>conv2d(input, params, stride, padding)</code>	Creates a 2D convolutional layer with bias (without activation).
<code>conv2d_params(num_input_channels, ...)</code>	Creates weights and bias variables for a 2D convolutional layer.
<code>flatten(input)</code>	Flattens inputs but keeps the batch size.
<code>fully_connected(input, params)</code>	Creates a fully connected layer with bias (without activation).
<code>fully_connected_params(input_size, ...)</code>	Creates weights and bias variables for a fully connected layer.
<code>linear_decay(start_value, end_value, step, ...)</code>	Applies linear decay from <i>start_value</i> to <i>end_value</i> .

### Classes

---

<code>ClipGlobalNormOptimizer(optimizer, clip_norm)</code>	A <code>tf.train.Optimizer</code> that wraps around another optimizer and minimizes the loss by clipping gradients using the global norm ( <code>tf.clip_by_global_norm()</code> ).
--	---

---

**class** `actorcritic.nn.ClipGlobalNormOptimizer` (*optimizer, clip\_norm, name=None*)

Bases: `tensorflow.python.training.optimizer.Optimizer`

A `tf.train.Optimizer` that wraps around another optimizer and minimizes the loss by clipping gradients using the global norm (`tf.clip_by_global_norm()`).

**See also:**

- [https://www.tensorflow.org/versions/r1.2/api\\_docs/python/tf/clip\\_by\\_global\\_norm](https://www.tensorflow.org/versions/r1.2/api_docs/python/tf/clip_by_global_norm)
- <https://stackoverflow.com/questions/36498127/how-to-apply-gradient-clipping-in-tensorflow/43486487#43486487>

`__init__` (*optimizer, clip\_norm, name=None*)

#### Parameters

- **optimizer** (`tf.train.Optimizer`) – An optimizer whose gradients will be clipped.
- **clip\_norm** (`tf.Tensor` or `float`) – Value for the global norm (passed to `tf.clip_by_global_norm()`).
- **name** (`string`, optional) – A name for this optimizer.

**apply\_gradients** (*grads\_and\_vars, global\_step=None, name=None*)

Apply gradients to variables.

This is the second part of `minimize()`. It returns an *Operation* that applies gradients.

#### Parameters

- **grads\_and\_vars** – List of (gradient, variable) pairs as returned by `compute_gradients()`.
- **global\_step** – Optional *Variable* to increment by one after the variables have been updated.
- **name** – Optional name for the returned operation. Default to the name passed to the *Optimizer* constructor.

**Returns** An *Operation* that applies the specified gradients. If `global_step` was not `None`, that operation also increments `global_step`.

#### Raises

- `TypeError` – If `grads_and_vars` is malformed.
- `ValueError` – If none of the variables have gradients.

`actorcritic.nn.conv2d` (*input, params, stride, padding*)

Creates a 2D convolutional layer with bias (without activation).

#### Parameters

- **input** (`tf.Tensor`) – The input values.

- **params** (tuple of (tf.Variable, tf.Variable)) – A tuple of (*weights*, *bias*). Probably obtained by `conv2d_params()`.
- **stride** (int) – The stride of the convolution.
- **padding** (string) – The padding of the convolution. One of 'VALID', 'SAME'.

**Returns** tf.Tensor – The output values.

`actorcritic.nn.conv2d_params` (*num\_input\_channels*, *num\_filters*, *filter\_extent*, *dtype*, *weights\_initializer*, *bias\_initializer*)

Creates weights and bias variables for a 2D convolutional layer. These can be used in `conv2d()`.

#### Parameters

- **num\_input\_channels** (int) – The size of the input layer.
- **num\_filters** (int) – The output size. Number of filters to apply.
- **filter\_extent** (int) – The spatial extent of the filters. Determines the size of the weights.
- **dtype** (tf.DType) – The data type of the variables.
- **weights\_initializer** (tf.keras.initializers.Initializer) – An initializer for the weights.
- **bias\_initializer** (tf.keras.initializers.Initializer) – An initializer for the bias.

**Returns** tuple of (tf.Variable, tf.Variable) – A tuple of (*weights*, *bias*).

`actorcritic.nn.flatten` (*input*)

Flattens inputs but keeps the batch size.

**Parameters** *input* (tf.Tensor) – Input values of shape [*batch\_size*, *d<sub>1</sub>*, ..., *d<sub>n</sub>*].

**Returns** tf.Tensor – Flattened input values of shape [*batch\_size*, *d<sub>1</sub>* \* ... \* *d<sub>n</sub>*].

`actorcritic.nn.fully_connected` (*input*, *params*)

Creates a fully connected layer with bias (without activation).

#### Parameters

- **input** (tf.Tensor) – The input values.
- **params** (tuple of (tf.Variable, tf.Variable)) – A tuple of (*weights*, *bias*). Probably obtained by `fully_connected_params()`.

**Returns** tf.Tensor – The output values.

`actorcritic.nn.fully_connected_params` (*input\_size*, *output\_size*, *dtype*, *weights\_initializer*, *bias\_initializer*)

Creates weights and bias variables for a fully connected layer. These can be used in `fully_connected()`.

#### Parameters

- **input\_size** (int) – The size of the input layer.
- **output\_size** (int) – The output size. Number of units.
- **dtype** (tf.DType) – The data type of the variables.
- **weights\_initializer** (tf.keras.initializers.Initializer) – An initializer for the weights.
- **bias\_initializer** (tf.keras.initializers.Initializer) – An initializer for the bias.

**Returns** tuple of (tf.Variable, tf.Variable) – A tuple of (*weights*, *bias*).

`actorcritic.nn.linear_decay` (*start\_value*, *end\_value*, *step*, *total\_steps*, *name=None*)

Applies linear decay from *start\_value* to *end\_value*. The value at a specific step is computed as:

```
value = (start_value - end_value) * (1 - step / total_steps) + end_value
```

#### Parameters

- **start\_value** (tf.Tensor or float) – The start value.
- **end\_value** (tf.Tensor or float) – The end value.
- **step** (tf.Tensor) – The current step (e.g. `global_step`).
- **total\_step** (int or tf.Tensor) – The total number of steps. Steps to reach `end_value`.
- **name** (string, optional) – A name for the operation.

**Returns** tf.Tensor – The linear decayed value.

## 1.1.7 actorcritic.objectives

Contains *objectives* that are used to optimize actor-critic models.

### Classes

<code>A2CObjective</code> ( <i>model</i> [, <i>discount_factor</i> , ...])	An objective that defines the loss of the policy and the baseline according to the A3C and A2C/ACKTR papers.
<code>ActorCriticObjective</code>	An objective takes an <code>ActorCriticModel</code> and determines how it is optimized.

```
class actorcritic.objectives.A2CObjective (model, discount_factor=0.99, entropy_regularization_strength=0.01, name=None)
```

Bases: `actorcritic.objectives.ActorCriticObjective`

An objective that defines the loss of the policy and the baseline according to the A3C and A2C/ACKTR papers.

The rewards are discounted and the policy loss uses entropy regularization. The baseline is optimized using a squared error loss.

The policy objective uses entropy regularization:

```
J(theta) = log(policy(state, action | theta)) * (target_values - baseline) + beta_
↳* entropy(policy)
```

where *beta* determines the strength of the entropy regularization.

**See also:**

- <https://arxiv.org/pdf/1602.01783.pdf> (A3C)
- <https://arxiv.org/pdf/1708.05144.pdf> (A2C/ACKTR)

`__init__` (*model*, *discount\_factor*=0.99, *entropy\_regularization\_strength*=0.01, *name*=None)

#### Parameters

- **model** (*ActorCriticModel*) – A model that provides the policy and the baseline that will be optimized.
- **discount\_factor** (*float*) – Used for discounting the rewards. Should be between [0, 1].
- **entropy\_regularization\_strength** (*float* or *tf.Tensor*) – Determining the strength of the entropy regularization. Corresponds to the *beta* parameter in A3C.
- **name** (*string*, optional) – A name for this objective.

#### **baseline\_loss**

*tf.Tensor* – The current loss of the baseline of the model.

#### **mean\_entropy**

*tf.Tensor* – The current mean entropy of the policy of the model.

#### **policy\_loss**

*tf.Tensor* – The current loss of the policy of the model.

**class** `actorcritic.objectives.ActorCriticObjective`

Bases: `object`

An objective takes an *ActorCriticModel* and determines how it is optimized. It defines the loss of the policy and the loss of the baseline, and can create train operations based on these losses.

#### **baseline\_loss**

*tf.Tensor* – The current loss of the baseline of the model.

**optimize\_separate** (*policy\_optimizer*, *baseline\_optimizer*, *policy\_kwargs*=None, *baseline\_kwargs*=None)

Creates an operation that minimizes the policy loss and the baseline loss separately. This means that it minimizes the losses using two different optimizers.

#### Parameters

- **policy\_optimizer** (*tf.train.Optimizer*) – An optimizer that is used for the policy loss.
- **baseline\_optimizer** (*tf.train.Optimizer*) – An optimizer that is used for the baseline loss.
- **policy\_kwargs** (*dict*, optional) – Keyword arguments passed to the `minimize()` method of the *policy\_optimizer*.
- **baseline\_kwargs** (*dict*, optional) – Keyword arguments passed to the `minimize()` method of the *baseline\_optimizer*.

**Returns** *tf.Operation* – An operation that updates both the policy and the baseline.

**optimize\_shared** (*optimizer*, *baseline\_loss\_weight*=0.5, *\*\*kwargs*)

Creates an operation that minimizes both the policy loss and the baseline loss using the same optimizer. This is used for models that share parameters between the policy and the baseline. The shared loss is defined as:

```
shared_loss = policy_loss + baseline_loss_weight * baseline_loss
```

where *baseline\_loss\_weight* determines the ‘learning rate’ relative to the policy loss.

#### Parameters

- **optimizer** (`tf.train.Optimizer`) – An optimizer that is used for both the policy loss and the baseline loss.
- **baseline\_loss\_weight** (`float` or `tf.Tensor`) – Determines the relative ‘learning rate’.
- **kwargs** (`dict`, optional) – Keyword arguments passed to the `minimize()` method of the optimizer.

**Returns** `tf.Operation` – An operation that updates both the policy and the baseline.

**policy\_loss**

`tf.Tensor` – The current loss of the policy of the model.

## 1.1.8 actorcritic.policies

Contains *policies* that determine the behavior of an *agent*.

### Classes

<code>DistributionPolicy(distribution, actions[, ...])</code>	Base class for stochastic policies that follow a concrete <code>tf.distributions.Distribution</code> .
<code>Policy</code>	Base class for stochastic policies.
<code>SoftmaxPolicy(logits, actions[, ...])</code>	A stochastic policy that follows a categorical distribution.

**class** `actorcritic.policies.DistributionPolicy` (*distribution*, *actions*, *random\_seed=None*)

Bases: `actorcritic.policies.Policy`

Base class for stochastic policies that follow a concrete `tf.distributions.Distribution`. Implements the required methods based on this distribution.

`__init__` (*distribution*, *actions*, *random\_seed=None*)

**Parameters**

- **distribution** (`tf.distributions.Distribution`) – The distribution.
- **actions** (`tf.Tensor`) – The input actions used to compute the log-probabilities. Must have the same shape as the inputs.
- **random\_seed** (`int`, optional) – A random seed used for sampling.

**entropy**

`tf.Tensor` – Computes the entropy of this policy based on the inputs that are provided for computing the probabilities. The shape equals the shape of the inputs.

**log\_prob**

`tf.Tensor` – Computes the log-probability of the given actions based on the inputs that are provided for computing the probabilities. The shape equals the shape of the actions and the inputs.

**mode**

`tf.Tensor` – Selects actions from this policy which have the highest probability (mode) based on the inputs that are provided for computing the probabilities. The shape equals the shape of the inputs.

**sample**

`tf.Tensor` – Samples actions from this policy based on the inputs that are provided for computing the

probabilities. The shape equals the shape of the inputs.

**class** actorcritic.policies.Policy

Bases: `object`

Base class for stochastic policies.

**entropy**

`tf.Tensor` – Computes the entropy of this policy based on the inputs that are provided for computing the probabilities. The shape equals the shape of the inputs.

**log\_prob**

`tf.Tensor` – Computes the log-probability of the given actions based on the inputs that are provided for computing the probabilities. The shape equals the shape of the actions and the inputs.

**mode**

`tf.Tensor` – Selects actions from this policy which have the highest probability (mode) based on the inputs that are provided for computing the probabilities. The shape equals the shape of the inputs.

**register\_predictive\_distribution** (*layer\_collection*, *random\_seed=None*)

Registers the predictive distribution of this policy in the specified `kfac.LayerCollection` (required for K-FAC).

**Parameters**

- **layer\_collection** (`kfac.LayerCollection`) – A layer collection used by the `KfacOptimizer`.
- **random\_seed** (`int`, optional) – A random seed for sampling from the predictive distribution.

**Raises** `NotImplementedError` – If this policy does not support K-FAC.

**sample**

`tf.Tensor` – Samples actions from this policy based on the inputs that are provided for computing the probabilities. The shape equals the shape of the inputs.

**class** actorcritic.policies.SoftmaxPolicy (*logits*, *actions*, *random\_seed=None*,  
*name=None*)

Bases: `actorcritic.policies.DistributionPolicy`

A stochastic policy that follows a categorical distribution.

**\_\_init\_\_** (*logits*, *actions*, *random\_seed=None*, *name=None*)

**Parameters**

- **logits** (`tf.Tensor`) – The input logits (or ‘scores’) used to compute the probabilities.
- **actions** (`tf.Tensor`) – The input actions used to compute the log-probabilities. Must have the same shape as *logits*.
- **random\_seed** (`int`, optional) – A random seed used for sampling.
- **name** (`string`, optional) – A name for this policy.

**register\_predictive\_distribution** (*layer\_collection*, *random\_seed=None*)

Registers the predictive distribution of this policy in the specified `kfac.LayerCollection` (required for K-FAC).

**Parameters**

- **layer\_collection** (`kfac.LayerCollection`) – A layer collection used by the `KfacOptimizer`.

- **random\_seed** (*int*, optional) – A random seed for sampling from the predictive distribution.

### 1.1.9 actorcritic.envs

Contains functions that are dedicated to certain environments.

<i>atari</i>	Contains functions that are dedicated to Atari environments.
--------------	--

#### actorcritic.envs.atari

Contains functions that are dedicated to Atari environments.

<i>model</i>	An implementation of an actor-critic model that is aimed at Atari games.
<i>wrappers</i>	Contains <i>wrappers</i> that can wrap around environments to modify their functionality.

#### actorcritic.envs.atari.model

An implementation of an actor-critic model that is aimed at Atari games.

#### Classes

<i>AtariModel</i> ( <i>observation_space</i> , <i>action_space</i> )	An <i>ActorCriticModel</i> that follows the A3C and ACKTR paper.
--	--

```
class actorcritic.envs.atari.model.AtariModel (observation_space, action_space,
                                             conv3_num_filters=64, random_seed=None, name=None)
```

Bases: *actorcritic.model.ActorCriticModel*

An *ActorCriticModel* that follows the A3C and ACKTR paper.

The observations are sent to three convolutional layers followed by a fully connected layer, each using rectifier activation functions (ReLU). The policy and the baseline use fully connected layers built on top of the last hidden fully connected layer separately. The policy layer has one unit for each action and its outputs are used as logits for a categorical distribution (softmax). The baseline layer has only one unit which represents its value.

The weights of the layers are orthogonally initialized.

Detailed network architecture:

- Conv2D: 32 filters 8x8, stride 4
- ReLU
- Conv2D: 64 filters 4x4, stride 2
- ReLU
- Conv2D: 64 filters 3x3, stride 1 (number of filters based on argument *conv3\_num\_filters*)

- Flatten
- Fully connected: 512 units
- ReLU
- Fully connected (policy): units = number of actions / Fully connected (baseline): 1 unit

A2C uses 64 filters in the third convolutional layer. ACKTR uses 32.

The policy is a *SoftmaxPolicy*. The baseline is a *StateValueFunction*.

**See also:**

This network architecture was originally used in: <https://www.nature.com/articles/nature14236>

`__init__(observation_space, action_space, conv3_num_filters=64, random_seed=None, name=None)`

**Parameters**

- **observation\_space** (`gym.spaces.Space`) – A space that determines the shape of the `observations_placeholder` and the `bootstrap_observations_placeholder`.
- **action\_space** (`gym.spaces.Space`) – A space that determines the shape of the `actions_placeholder`.
- **conv3\_num\_filters** (`int`, optional) – Number of filters used for the third convolutional layer, defaults to 64. ACKTR uses 32.
- **random\_seed** (`int`, optional) – A random seed used for sampling from the *actor-critic.policies.SoftmaxPolicy*.
- **name** (`string`, optional) – A name for this model.

**register\_layers** (`layer_collection`)

Registers the layers of this model (neural net) in the specified `kfac.LayerCollection` (required for K-FAC).

**Parameters** `layer_collection` (`kfac.LayerCollection`) – A layer collection used by the `KfacOptimizer`.

## actorncritic.envs.atari.wrappers

Contains *wrappers* that can wrap around environments to modify their functionality.

The implementations of these wrappers are adopted from OpenAI.

## Classes

---

<code>AtariClipRewardWrapper(env)</code>	A wrapper that clips the rewards between -1 and 1.
<code>AtariEpisodicLifeWrapper(env)</code>	A wrapper that ends episodes (returns <i>terminal</i> = True) after a life in the Atari game has been lost.
<code>AtariFireResetWrapper(env)</code>	A wrapper that executes the 'FIRE' action after the environment has been reset.
<code>AtariFrameskipWrapper(env, frameskip)</code>	A wrapper that skips frames.
<code>AtariInfoClearWrapper(env)</code>	A wrapper that removes unnecessary data in the <i>info</i> returned by <code>gym.Env.step()</code> .

---

Continued on next page

Table 17 – continued from previous page

<code>AtariNoopResetWrapper(env, noop_max)</code>	A wrapper that executes a random number of ‘ <i>NOOP</i> ’ actions.
<code>AtariPreprocessFrameWrapper(env)</code>	A wrapper that scales the observations from 210x160 down to 84x84 and converts from RGB to grayscale by extracting the luminance.
<code>EpisodeInfoWrapper(env)</code>	A wrapper that stores episode information in the <i>info</i> returned by <code>gym.Env.step()</code> at the end of an episode.
<code>FrameStackWrapper(env, num_stacked_frames)</code>	A wrapper that stacks the last observations.
<code>RenderWrapper(env[, fps])</code>	A wrapper that calls <code>gym.Env.render()</code> every step.

**class** `actorcritic.envs.atari.wrappers.AtariClipRewardWrapper` (*env*)

Bases: `gym.core.RewardWrapper`

A wrapper that clips the rewards between -1 and 1.

`__init__` (*env*)

**Parameters** *env* (`gym.Env`) – An environment that will be wrapped.

**class** `actorcritic.envs.atari.wrappers.AtariEpisodicLifeWrapper` (*env*)

Bases: `gym.core.Wrapper`

A wrapper that ends episodes (returns *terminal* = True) after a life in the Atari game has been lost.

`__init__` (*env*)

**Parameters** *env* (`gym.Env`) – An environment that will be wrapped.

**reset** (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns: observation (object): the initial observation of the** space.

**step** (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling *reset()* to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** *action* (*object*) – an action provided by the environment

**Returns** *observation* (*object*) – agent’s observation of the current environment  
*reward* (float) : amount of reward returned after previous action done  
*done* (boolean): whether the episode has ended, in which case further step() calls will return undefined results  
*info* (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**class** `actorcritic.envs.atari.wrappers.AtariFireResetWrapper` (*env*)

Bases: `gym.core.Wrapper`

A wrapper that executes the ‘*FIRE*’ action after the environment has been reset.

`__init__` (*env*)

**Parameters** *env* (`gym.Env`) – An environment that will be wrapped.

**reset** (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns: observation (object): the initial observation of the** space.

**step** (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling *reset()* to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** `action` (*object*) – an action provided by the environment

**Returns** `observation` (*object*) – agent’s observation of the current environment  
`reward` (float) : amount of reward returned after previous action done  
`done` (boolean): whether the episode has ended, in which case further `step()` calls will return undefined results  
`info` (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**class** `actorcritic.envs.atari.wrappers.AtariFrameskipWrapper` (*env*, *frameskip*)

Bases: `gym.core.Wrapper`

A wrapper that skips frames.

`__init__` (*env*, *frameskip*)

**Parameters**

- `env` (`gym.Env`) – An environment that will be wrapped.
- `frameskip` (`int`) – Every *frameskip*-th frame is used. The remaining frames are skipped.

`reset` (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns:** `observation` (*object*): the initial observation of the space.

`step` (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** `action` (*object*) – an action provided by the environment

**Returns** `observation` (*object*) – agent’s observation of the current environment  
`reward` (float) : amount of reward returned after previous action done  
`done` (boolean): whether the episode has ended, in which case further `step()` calls will return undefined results  
`info` (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**class** `actorcritic.envs.atari.wrappers.AtariInfoClearWrapper` (*env*)

Bases: `gym.core.Wrapper`

A wrapper that removes unnecessary data in the *info* returned by `gym.Env.step()`. This reduces the amount of inter-process data.

**Warning:** `AtariEpisodicLifeWrapper` does not work afterwards, so it should be used *before*.

`__init__` (*env*)

**Parameters** `env` (`gym.Env`) – An environment that will be wrapped.

`reset` (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns:** `observation` (*object*): the initial observation of the space.

`step` (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** *action* (*object*) – an action provided by the environment

**Returns** *observation* (*object*) – agent’s observation of the current environment reward (float) : amount of reward returned after previous action done (boolean): whether the episode has ended, in which case further `step()` calls will return undefined results *info* (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**class** `actorcritic.envs.atari.wrappers.AtariNoopResetWrapper` (*env*, *noop\_max*)  
Bases: `gym.core.Wrapper`

A wrapper that executes a random number of ‘*NOOP*’ actions.

`__init__` (*env*, *noop\_max*)

**Parameters**

- **env** (`gym.Env`) – An environment that will be wrapped.
- **noop\_max** (`int`) – The maximum number of ‘*NOOP*’ actions. The number is selected randomly between 1 and *noop\_max*.

**reset** (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns:** **observation (object): the initial observation of the** space.

**step** (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** *action* (*object*) – an action provided by the environment

**Returns** *observation* (*object*) – agent’s observation of the current environment reward (float) : amount of reward returned after previous action done (boolean): whether the episode has ended, in which case further `step()` calls will return undefined results *info* (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**class** `actorcritic.envs.atari.wrappers.AtariPreprocessFrameWrapper` (*env*)  
Bases: `gym.core.ObservationWrapper`

A wrapper that scales the observations from 210x160 down to 84x84 and converts from RGB to grayscale by extracting the luminance.

`__init__` (*env*)

**Parameters** **env** (`gym.Env`) – An environment that will be wrapped.

**class** `actorcritic.envs.atari.wrappers.EpisodeInfoWrapper` (*env*)  
Bases: `gym.core.Wrapper`

A wrapper that stores episode information in the *info* returned by `gym.Env.step()` at the end of an episode. More specifically, if an episode is terminal, *info* will contain the key ‘*episode*’ which has a `dict` value containing the ‘*total\_reward*’, which is the cumulative reward of the episode.

---

**Note:** If you want to get the cumulative reward of the entire episode, `AtariEpisodicLifeWrapper` should be used *after* this wrapper.

---

`__init__` (*env*)

**Parameters** **env** (`gym.Env`) – An environment that will be wrapped.

**static** `get_episode_rewards_from_info_batch` (*infos*)

Utility function that extracts the episode rewards, that are inserted by the `EpisodeInfoWrapper`, out of the *infos*.

**Parameters** *infos* (`list` of `list`) – A batch-major list of *infos* as returned by `interact()`.

**Returns** `numpy.ndarray` – A batch-major array with the same shape as *infos*. It contains the episode reward of an *info* at the corresponding position. If no episode reward was in an *info*, the result will contain `numpy.nan` respectively.

**reset** (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns:** **observation (object): the initial observation of the** `space`.

**step** (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** *action* (*object*) – an action provided by the environment

**Returns** *observation* (*object*) – agent’s observation of the current environment reward (float) : amount of reward returned after previous action done (boolean): whether the episode has ended, in which case further `step()` calls will return undefined results *info* (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**class** `actorcritic.envs.atari.wrappers.FrameStackWrapper` (*env*, *num\_stacked\_frames*)

Bases: `gym.core.Wrapper`

A wrapper that stacks the last observations. The observations returned by this wrapper consist of the last frames.

**\_\_init\_\_** (*env*, *num\_stacked\_frames*)

**Parameters**

- **env** (`gym.Env`) – An environment that will be wrapped.
- **num\_stacked\_frames** (`int`) – The number of frames that will be stacked.

**reset** (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns:** **observation (object): the initial observation of the** `space`.

**step** (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** *action* (*object*) – an action provided by the environment

**Returns** *observation* (*object*) – agent’s observation of the current environment reward (float) : amount of reward returned after previous action done (boolean): whether the episode has ended, in which case further `step()` calls will return undefined results *info* (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**class** `actorcritic.envs.atari.wrappers.RenderWrapper` (*env*, *fps=None*)

Bases: `gym.core.Wrapper`

A wrapper that calls `gym.Env.render()` every step.

`__init__(env, fps=None)`

#### Parameters

- **env** (`gym.Env`) – An environment that will be wrapped.
- **fps** (`int`, `float`, optional) – If it is not `None`, the steps will be slowed down to run at the specified frames per second by waiting  $1.0/fps$  seconds every step.

**reset** (*\*\*kwargs*)

Resets the state of the environment and returns an initial observation.

**Returns:** **observation (object): the initial observation of the** space.

**step** (*action*)

Run one timestep of the environment’s dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment’s state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** **action (object)** – an action provided by the environment

**Returns** *observation (object)* – agent’s observation of the current environment reward (float) : amount of reward returned after previous action done (boolean): whether the episode has ended, in which case further step() calls will return undefined results info (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

### 1.1.10 actorcritic.examples

Contains examples of how to use this project.

---

<code>atari</code>	Contains examples that deal with Atari environments.
--------------------	--

---

#### actorcritic.examples.atari

Contains examples that deal with Atari environments.

---

<code>a2c_acktr</code>	An example of how to use <i>A2C</i> and <i>ACKTR</i> to learn to play an Atari game.
------------------------	--

---

#### actorcritic.examples.atari.a2c\_acktr

An example of how to use *A2C* and *ACKTR* to learn to play an Atari game.

### Functions

---

<code>create_environments(env_id, num_envs)</code>	Creates multiple Atari environments that run in subprocesses.
<code>create_optimizer(acktr, model, learning_rate)</code>	Creates an optimizer based on whether <i>ACKTR</i> or <i>A2C</i> is used.
<code>load_model(saver, checkpoint_path, session)</code>	Loads the latest model checkpoint (with the neural network parameters) from a directory.

---

Continued on next page

Table 20 – continued from previous page

<code>make_atari_env(env_id, render)</code>	Creates a <code>gym.Env</code> and wraps it with all Atari wrappers in <code>actorcritic.envs.atari.wrappers</code> .
<code>save_model(saver, checkpoint_path, ...)</code>	Saves a model checkpoint to a directory.
<code>train_a2c_acktr(acktr, env_id, num_envs, ...)</code>	Trains an Atari model using <i>A2C</i> or <i>ACKTR</i> .

`actorcritic.examples.atari.a2c_acktr.create_environments(env_id, num_envs)`  
Creates multiple Atari environments that run in subprocesses.

#### Parameters

- **env\_id** (`string`) – An id passed to `gym.make()` to create the environments.
- **num\_envs** (`int`) – The number of environments (and subprocesses) that will be created.

**Returns** `list` of `gym Wrapper` – The environments.

`actorcritic.examples.atari.a2c_acktr.create_optimizer(acktr, model, learning_rate)`  
Creates an optimizer based on whether *ACKTR* or *A2C* is used. *A2C* uses the RMSProp optimizer, *ACKTR* uses the K-FAC optimizer. This function is not restricted to Atari models and can be used generally.

#### Parameters

- **acktr** (`bool`) – Whether to use the optimizer of *ACKTR* or *A2C*.
- **model** (`ActorCriticModel`) – A model that is needed for K-FAC to register the neural network layers and the predictive distributions.
- **learning\_rate** (`float` or `tf.Tensor`) – A learning rate for the optimizer.

`actorcritic.examples.atari.a2c_acktr.load_model(saver, checkpoint_path, session)`  
Loads the latest model checkpoint (with the neural network parameters) from a directory.

#### Parameters

- **saver** (`tf.train.Saver`) – A saver object to restore the model.
- **checkpoint\_path** (`string`) – A directory where the checkpoint is loaded from.
- **session** (`tf.Session`) – A session which will contain the loaded variable values.

`actorcritic.examples.atari.a2c_acktr.make_atari_env(env_id, render)`  
Creates a `gym.Env` and wraps it with all Atari wrappers in `actorcritic.envs.atari.wrappers`.

#### Parameters

- **env\_id** (`string`) – An id passed to `gym.make()`.
- **render** (`bool`) – Whether this environment should be rendered.

**Returns** `gym.Env` – The environment.

`actorcritic.examples.atari.a2c_acktr.save_model(saver, checkpoint_path, model_name, step, session)`  
Saves a model checkpoint to a directory.

#### Parameters

- **saver** (`tf.train.Saver`) – A saver object to save the model.
- **checkpoint\_path** (`string`) – A directory where the model checkpoint will be saved.
- **model\_name** (`string`) – A name of the model. The checkpoint file in the `checkpoint_path` directory will have this name.
- **step** (`int` or `tf.Tensor`) – A number that is appended to the checkpoint file name.

- **session** (`tf.Session`) – A session whose variables will be saved.

```
actorcritic.examples.atari.a2c_acktr.train_a2c_acktr(acktr, env_id, num_envs,  
                                                    num_steps, checkpoint_path,  
                                                    model_name,          sum-  
                                                    mary_path=None)
```

Trains an Atari model using *A2C* or *ACKTR*. Automatically saves and loads the trained model.

### Parameters

- **acktr** (`bool`) – Whether the *ACKTR* or the *A2C* algorithm should be used. *A2C* uses the RMSProp optimizer and 64 filters in the third convolutional layer of the neural network. *ACKTR* uses the K-FAC optimizer and 32 filters.
- **env\_id** (`string`) – An id passed to `gym.make()` to create the environments.
- **num\_envs** (`int`) – The number of environments that will be used (so *num\_envs* subprocesses will be created). *A2C* normally uses 16. *ACKTR* normally uses 32.
- **num\_steps** (`int`) – The number of steps to take in each iteration. *A2C* normally uses 5. *ACKTR* normally uses 20.
- **checkpoint\_path** (`string`) – A directory where the model's checkpoints will be loaded and saved.
- **model\_name** (`string`) – A name of the model. The files in the *checkpoint\_path* directory will have this name.
- **summary\_path** (`string`, optional) – A directory where the TensorBoard summaries will be saved. If not specified, no summaries will be saved.



The basic idea of reinforcement learning is to find a behavior for an *agent* inside an *environment* that leads to a maximal *reward*. Such a behavior is called a *policy* and it decides what *action* to take based on the current *observation* (also called *state*).

For example, the environment can be an Atari game. In this case the reward is the score, the actions are the controller actions, and the current frame/image of the game is an observation.

The `gym` library ([GitHub](#)) by OpenAI provides several types of environments. A basic reinforcement learning setup to learn a policy for the *Breakout* environment could look like this:

```
import gym

# create the environment
env = gym.make('BreakoutNoFrameskip-v4')

# receive an initial observation (frame) to select the first action
observation = env.reset()

while True:
    # let the current policy select an action
    action = policy(observation)

    # execute the action and take one step in the environment (go to next frame)
    next_observation, reward, terminal, info = env.step(action)

    # improve the policy based on this experience
    improve_policy(observation, action, reward, terminal, next_observation)

    observation = next_observation

    if terminal:
        observation = env.reset()
```

`terminal` indicates whether the game ended, so the game has to be reset. `reward` is just a number that represents the points that were achieved in this step. `info` contains debug information (the current number of lives).

*A2C* and *ACKTR* actually use *multiple environments* at once by running them in multiple subprocesses. This means that we can improve the policy faster, since we simply have more observations and rewards available. For that reason there is *MultiEnv*:

```
from actorcritic.multi_env import MultiEnv

envs = create_environments() # create multiple environments
multi_env = MultiEnv(envs)
```

Yet the crucial parts are `policy(observation)` and `improve_policy(observation, action, reward, next_observation)`. We need to know how to define a policy and especially how to improve it.

*Actor-critic* methods define the policy as a probability distribution, such that it computes the probability of every action based on the current observation. Then these probabilities are used to sample one of the actions. For example, if the ball approaches the bottom in *Breakout*, the probability to move the paddle towards the ball should be high.

We typically use a neural network to compute these probabilities. Then the observations (frames) are sent into the network, which produces a score for every action. These scores can be passed in the softmax function to obtain probabilities. *AtariModel* provides a neural network and a policy made for Atari environments:

```
from actorcritic.envs.atari.model import AtariModel

# observation_space and action_space define the type and shape of the observations,
↳ and actions
# e.g. the size of the frames
model = AtariModel(multi_env.observation_space, multi_env.action_space)
```

Additionally *A2C* and *ACKTR* do not take one step only and improve the policy immediately. Instead they take multiple steps and use all the experienced observations and rewards to improve the policy. A *MultiEnvAgent* simplifies this process. It takes the neural network and the policy (the ‘*model*’), and the environments. Then we just have to call `interact()` and it uses the policy to take multiple steps:

```
from actorcritic.agents import MultiEnvAgent

agent = MultiEnvAgent(multi_env, model, num_steps=5)

while True:
    # take 5 steps in all environments
    # session is a tf.Session used to compute the values of the neural network
    observations, actions, rewards, terminals, next_observations, infos = agent.
↳ interact(session)

    # improve the policy based on this experience
    improve_policy(observations, actions, rewards, terminals, next_observations)
```

In *actor-critic* methods we do not define a loss function directly, but a *policy objective* function to optimize the neural network. It needs the observations, the actions, and the rewards that the agent experienced. Then we can learn through the policy objective, which looks at the rewards in order to decide whether the actions were good or not.

Furthermore we need a *baseline* function that enhances the policy objective. It should express how much reward we can expect if we would follow our policy proceeding from the observations we just have seen. This helps the policy to decide whether the actions it has taken actually were better or worse than expected. This *baseline* function is the ‘*critic*’ of *actor-critic* (the policy is the ‘*actor*’). It distinguishes actor-critic methods from *policy gradient* methods which just have an ‘*actor*’.

Unfortunately we do not have such a *baseline* function. That is why we will learn the *baseline*, too, at the same time as the policy. Therefore an *ActorCriticModel* like the *AtariModel* has to provide a baseline. *A2C* and *ACKTR* use the *state-value function* which indeed tells us how much reward we can expect from a given observation.

It can be beneficial to use the same neural network as the policy for the baseline. *AtariModel* does exactly this.

In summary we need a *ActorCriticObjective*. The policy objective of A2C and ACKTR is implemented in *A2CObjective*. It *discounts* the rewards and uses *entropy regularization* (see *A2CObjective*).

```
from actorcritic.objectives import A2CObjective

objective = A2CObjective(model, discount_factor=0.99, entropy_regularization_
↳strength=0.01)
```

Next we need an optimizer for our neural network:

```
import tensorflow as tf

# A2C uses the RMSProp optimizer
optimizer = tf.train.RMSPropOptimizer(learning_rate=0.0007)

# create an 'optimize' operation that we can call
# use optimize_shared() since we share the network between the policy and the baseline
optimize_op = objective.optimize_shared(optimizer)
```

That is all. We can use all variables defined above to run the A2C algorithm:

```
while True:
    # take multiple steps in all environments
    observations, actions, rewards, terminals, next_observations, infos = agent.
↳interact(session)

    # improve the policy and the baseline
    session.run(optimize_op, feed_dict={
        model.observations_placeholder: observations,
        model.bootstrap_observations_placeholder: next_observations,
        model.actions_placeholder: actions,
        model.rewards_placeholder: rewards,
        model.terminals_placeholder: terminals
    })
```

*bootstrap\_observations\_placeholder* is needed to compute the *bootstrap\_values*, which are used in the policy objective.

In order to use *ACKTR* we just have to change the optimizer to a *kfac.KfacOptimizer*.

See *a2c\_acktr.py* for a full implementation, especially how to implement *create\_environments()* and how to use the K-FAC optimizer.



**a**

actorcritic, 3  
actorcritic.agents, 3  
actorcritic.baselines, 6  
actorcritic.envs, 19  
actorcritic.envs.atari, 19  
actorcritic.envs.atari.model, 19  
actorcritic.envs.atari.wrappers, 20  
actorcritic.examples, 25  
actorcritic.examples.atari, 25  
actorcritic.examples.atari.a2c\_acktr,  
25  
actorcritic.kfac\_utils, 7  
actorcritic.model, 8  
actorcritic.multi\_env, 10  
actorcritic.nn, 12  
actorcritic.objectives, 15  
actorcritic.policies, 17



## Symbols

- `__AutoResetWrapper` (class in `actorcritic.multi_env`), 12
  - `__init__()` (`actorcritic.agents.MultiEnvAgent` method), 5
  - `__init__()` (`actorcritic.agents.SingleEnvAgent` method), 5
  - `__init__()` (`actorcritic.baselines.StateValueFunction` method), 7
  - `__init__()` (`actorcritic.envs.atari.model.AtariModel` method), 20
  - `__init__()` (`actorcritic.envs.atari.wrappers.AtariClipRewardWrapper` method), 21
  - `__init__()` (`actorcritic.envs.atari.wrappers.AtariEpisodicLifeWrapper` method), 21
  - `__init__()` (`actorcritic.envs.atari.wrappers.AtariFireResetWrapper` method), 21
  - `__init__()` (`actorcritic.envs.atari.wrappers.AtariFrameskipWrapper` method), 22
  - `__init__()` (`actorcritic.envs.atari.wrappers.AtariInfoClearWrapper` method), 22
  - `__init__()` (`actorcritic.envs.atari.wrappers.AtariNoopResetWrapper` method), 23
  - `__init__()` (`actorcritic.envs.atari.wrappers.AtariPreprocessFrameWrapper` method), 23
  - `__init__()` (`actorcritic.kfac_utils.ColdStartPeriodicInvUpdateKfacOpt` method), 8
  - `__init__()` (`actorcritic.model.ActorCriticModel` method), 8
  - `__init__()` (`actorcritic.model.ActorCriticModel` method), 8
  - `__init__()` (`actorcritic.multi_env.MultiEnv` method), 10
  - `__init__()` (`actorcritic.multi_env.SubprocessEnv` method), 11
  - `__init__()` (`actorcritic.nn.ClipGlobalNormOptimizer` method), 13
  - `__init__()` (`actorcritic.objectives.A2CObjective` method), 15
  - `__init__()` (`actorcritic.policies.DistributionPolicy` method), 17
  - `__init__()` (`actorcritic.policies.SoftmaxPolicy` method), 18
- ## A
- `A2CObjective` (class in `actorcritic.objectives`), 15
  - `action_space` (`actorcritic.multi_env.MultiEnv` attribute), 10
  - `action_space` (`actorcritic.multi_env.SubprocessEnv` attribute), 11
  - `actions_placeholder` (`actorcritic.model.ActorCriticModel` attribute), 8
  - `actorcritic` (module), 3
  - `actorcritic.agents` (module), 3
  - `actorcritic.baselines` (module), 6
  - `actorcritic.envs` (module), 19
  - `actorcritic.envs.atari` (module), 19
  - `actorcritic.envs.atari.model` (module), 19
  - `actorcritic.envs.atari.wrappers` (module), 20
  - `actorcritic.examples` (module), 25
  - `actorcritic.examples.atari` (module), 25
  - `actorcritic.examples.atari.a2c_acktr` (module), 25
  - `actorcritic.kfac_utils` (module), 7
  - `actorcritic.model` (module), 8
  - `actorcritic.multi_env` (module), 10
  - `actorcritic.nn` (module), 12
  - `actorcritic.objectives` (module), 15
  - `actorcritic.policies` (module), 17
  - `ActorCriticModel` (class in `actorcritic.model`), 8
  - `ActorCriticObjective` (class in `actorcritic.objectives`), 16
  - `Agent` (class in `actorcritic.agents`), 4
  - `apply_gradients()` (`actorcritic.kfac_utils.ColdStartPeriodicInvUpdateKfacOpt` method), 8
  - `apply_gradients()` (`actorcritic.nn.ClipGlobalNormOptimizer` method), 13
  - `AtariClipRewardWrapper` (class in `actorcritic.envs.atari.wrappers`), 21

AtariEpisodicLifeWrapper (class in actor-critic.envs.atari.wrappers), 21  
 AtariFireResetWrapper (class in actor-critic.envs.atari.wrappers), 21  
 AtariFrameskipWrapper (class in actor-critic.envs.atari.wrappers), 22  
 AtariInfoClearWrapper (class in actor-critic.envs.atari.wrappers), 22  
 AtariModel (class in actorcritic.envs.atari.model), 19  
 AtariNoopResetWrapper (class in actor-critic.envs.atari.wrappers), 23  
 AtariPreprocessFrameWrapper (class in actor-critic.envs.atari.wrappers), 23

## B

baseline (actorcritic.model.ActorCriticModel attribute), 8  
 Baseline (class in actorcritic.baselines), 6  
 baseline\_loss (actorcritic.objectives.A2CObjective attribute), 16  
 baseline\_loss (actorcritic.objectives.ActorCriticObjective attribute), 16  
 bootstrap\_observations\_placeholder (actor-critic.model.ActorCriticModel attribute), 9  
 bootstrap\_values (actorcritic.model.ActorCriticModel attribute), 9

## C

ClipGlobalNormOptimizer (class in actorcritic.nn), 13  
 close() (actorcritic.multi\_env.MultiEnv method), 10  
 close() (actorcritic.multi\_env.SubprocessEnv method), 11  
 ColdStartPeriodicInvUpdateKfacOpt (class in actor-critic.kfac\_utils), 7  
 conv2d() (in module actorcritic.nn), 13  
 conv2d\_params() (in module actorcritic.nn), 14  
 create\_environments() (in module actor-critic.examples.atari.a2c\_acktr), 26  
 create\_optimizer() (in module actor-critic.examples.atari.a2c\_acktr), 26  
 create\_subprocess\_envs() (in module actor-critic.multi\_env), 12

## D

DistributionPolicy (class in actorcritic.policies), 17

## E

entropy (actorcritic.policies.DistributionPolicy attribute), 17  
 entropy (actorcritic.policies.Policy attribute), 18  
 envs (actorcritic.multi\_env.MultiEnv attribute), 10  
 EpisodeInfoWrapper (class in actor-critic.envs.atari.wrappers), 23

## F

flatten() (in module actorcritic.nn), 14  
 FrameStackWrapper (class in actor-critic.envs.atari.wrappers), 24  
 fully\_connected() (in module actorcritic.nn), 14  
 fully\_connected\_params() (in module actorcritic.nn), 14

## G

get\_episode\_rewards\_from\_info\_batch() (actor-critic.envs.atari.wrappers.EpisodeInfoWrapper static method), 23

## I

initialize() (actorcritic.multi\_env.SubprocessEnv method), 11  
 interact() (actorcritic.agents.Agent method), 4  
 interact() (actorcritic.agents.MultiEnvAgent method), 5  
 interact() (actorcritic.agents.SingleEnvAgent method), 5

## L

linear\_decay() (in module actorcritic.nn), 15  
 load\_model() (in module actor-critic.examples.atari.a2c\_acktr), 26  
 log\_prob (actorcritic.policies.DistributionPolicy attribute), 17  
 log\_prob (actorcritic.policies.Policy attribute), 18

## M

make\_atari\_env() (in module actor-critic.examples.atari.a2c\_acktr), 26  
 mean\_entropy (actorcritic.objectives.A2CObjective attribute), 16  
 mode (actorcritic.policies.DistributionPolicy attribute), 17  
 mode (actorcritic.policies.Policy attribute), 18  
 MultiEnv (class in actorcritic.multi\_env), 10  
 MultiEnvAgent (class in actorcritic.agents), 4

## O

observation\_space (actorcritic.multi\_env.MultiEnv attribute), 10  
 observation\_space (actorcritic.multi\_env.SubprocessEnv attribute), 11  
 observations\_placeholder (actor-critic.model.ActorCriticModel attribute), 9  
 optimize\_separate() (actor-critic.objectives.ActorCriticObjective method), 16  
 optimize\_shared() (actor-critic.objectives.ActorCriticObjective method), 16

## P

policy (actorcritic.model.ActorCriticModel attribute), 9  
 Policy (class in actorcritic.policies), 18  
 policy\_loss (actorcritic.objectives.A2CObjective attribute), 16  
 policy\_loss (actorcritic.objectives.ActorCriticObjective attribute), 17

## R

register\_layers() (actorcritic.envs.atari.model.AtariModel method), 20  
 register\_layers() (actorcritic.model.ActorCriticModel method), 9  
 register\_predictive\_distribution() (actorcritic.baselines.Baseline method), 6  
 register\_predictive\_distribution() (actorcritic.baselines.StateValueFunction method), 7  
 register\_predictive\_distribution() (actorcritic.policies.Policy method), 18  
 register\_predictive\_distribution() (actorcritic.policies.SoftmaxPolicy method), 18  
 register\_predictive\_distributions() (actorcritic.model.ActorCriticModel method), 9  
 render() (actorcritic.multi\_env.SubprocessEnv method), 11  
 RenderWrapper (class in actorcritic.envs.atari.wrappers), 24  
 reset() (actorcritic.envs.atari.wrappers.AtariEpisodicLifeWrapper method), 21  
 reset() (actorcritic.envs.atari.wrappers.AtariFireResetWrapper method), 21  
 reset() (actorcritic.envs.atari.wrappers.AtariFrameskipWrapper method), 22  
 reset() (actorcritic.envs.atari.wrappers.AtariInfoClearWrapper method), 22  
 reset() (actorcritic.envs.atari.wrappers.AtariNoopResetWrapper method), 23  
 reset() (actorcritic.envs.atari.wrappers.EpisodeInfoWrapper method), 24  
 reset() (actorcritic.envs.atari.wrappers.FrameStackWrapper method), 24  
 reset() (actorcritic.envs.atari.wrappers.RenderWrapper method), 25  
 reset() (actorcritic.multi\_env.\_AutoResetWrapper method), 12  
 reset() (actorcritic.multi\_env.MultiEnv method), 10  
 reset() (actorcritic.multi\_env.SubprocessEnv method), 11  
 rewards\_placeholder (actorcritic.model.ActorCriticModel attribute), 9

## S

sample (actorcritic.policies.DistributionPolicy attribute), 17  
 sample (actorcritic.policies.Policy attribute), 18  
 sample\_actions() (actorcritic.model.ActorCriticModel method), 9  
 save\_model() (in module actorcritic.examples.atari.a2c\_acktr), 26  
 select\_max\_actions() (actorcritic.model.ActorCriticModel method), 9  
 SingleEnvAgent (class in actorcritic.agents), 5  
 SoftmaxPolicy (class in actorcritic.policies), 18  
 start() (actorcritic.multi\_env.SubprocessEnv method), 11  
 StateValueFunction (class in actorcritic.baselines), 7  
 step() (actorcritic.envs.atari.wrappers.AtariEpisodicLifeWrapper method), 21  
 step() (actorcritic.envs.atari.wrappers.AtariFireResetWrapper method), 21  
 step() (actorcritic.envs.atari.wrappers.AtariFrameskipWrapper method), 22  
 step() (actorcritic.envs.atari.wrappers.AtariInfoClearWrapper method), 22  
 step() (actorcritic.envs.atari.wrappers.AtariNoopResetWrapper method), 23  
 step() (actorcritic.envs.atari.wrappers.EpisodeInfoWrapper method), 24  
 step() (actorcritic.envs.atari.wrappers.FrameStackWrapper method), 24  
 step() (actorcritic.envs.atari.wrappers.RenderWrapper method), 25  
 step() (actorcritic.multi\_env.\_AutoResetWrapper method), 12  
 step() (actorcritic.multi\_env.MultiEnv method), 10  
 step() (actorcritic.multi\_env.SubprocessEnv method), 11  
 SubprocessEnv (class in actorcritic.multi\_env), 11

## T

terminals\_placeholder (actorcritic.model.ActorCriticModel attribute), 9  
 train\_a2c\_acktr() (in module actorcritic.examples.atari.a2c\_acktr), 27  
 transpose\_list() (in module actorcritic.agents), 6

## V

value (actorcritic.baselines.Baseline attribute), 7  
 value (actorcritic.baselines.StateValueFunction attribute), 7