

---

# **cartpole-tf-dqn**

**Tianhao Zhou**

**Aug 24, 2020**



## CONTENTS

<b>1</b>	<b>Replay Buffer</b>	<b>3</b>
<b>2</b>	<b>Visualizer</b>	<b>5</b>
<b>3</b>	<b>CLI Entrypoint</b>	<b>7</b>
<b>4</b>	<b>Configuration</b>	<b>9</b>
<b>5</b>	<b>Train</b>	<b>11</b>
<b>6</b>	<b>Utilities</b>	<b>13</b>
<b>7</b>	<b>Tests</b>	<b>15</b>
	<b>Python Module Index</b>	<b>17</b>
	<b>Index</b>	<b>19</b>



---

```

class dqn_agent.DqnAgent (state_space, action_space, gamma, lr, verbose, checkpoint_location,
                           model_location, persist_progress_option, mode, epsilon)
    DQN agent with production policy and benchmark

collect_policy (state)
    The policy for collecting data points which can contain some randomness to encourage exploration.

    Returns action

load_checkpoint ()
    Loads training checkpoint into the underlying model

    Returns None

load_model ()
    Loads previously saved model :return: None

policy (state)
    Outputs a action based on model

    Parameters state – current state

    Returns action

random_policy (state)
    Outputs a random action

    Parameters state – current state

    Returns action

save_checkpoint ()
    Saves training checkpoint

    Returns None

save_model ()
    Saves model to file system

    Returns None

train (state_batch, next_state_batch, action_batch, reward_batch, done_batch, batch_size)
    Train the model on a batch

    Parameters
        • state_batch – batch of states
        • next_state_batch – batch of next states
        • action_batch – batch of actions
        • reward_batch – batch of rewards
        • done_batch – batch of done status
        • batch_size – the size of the batch

    Returns loss history

update_target_network ()
    Updates the target Q network with the parameters from the currently trained Q network.

    Returns None

```



## REPLAY BUFFER

```
class replay_buffer.DqnReplayBuffer(max_size)
    DQN replay buffer to keep track of game play records

    can_sample_batch(batch_size)
        Returns if a batch can be sampled

            Parameters batch_size – the size of the batch to be sampled
            Returns (bool) if can sample

    get_volume()
        Gets the current length of the records

            Returns (int) the length of the records

    record(state, reward, next_state, action, done)
        Puts a game play state into records

            Parameters
                • state – current game state
                • reward – reward after taking action
                • next_state – state after taking action
                • action – action taken
                • done – if the episode is finished

            Returns None

    sample_batch(batch_size)
        Samples a batch from the records

            Parameters batch_size – the size of the batch to be sampled
            Returns sample batch
```



---

CHAPTER  
TWO

---

## VISUALIZER

Training progress visualizer

**class** visualizer.DummyTrainingVisualizer

Used when no logging is required

**get\_ui\_feedback()**

A dummy logger that does nothing

**Returns** None

**log\_loss**(*loss*)

A dummy logger that does nothing

**Parameters** **loss** – a list of loss history

**Returns** None

**log\_reward**(*reward*)

A dummy logger that does nothing

**Parameters** **reward** – a list of reward history

**Returns** None

**class** visualizer.StreamlitTrainingVisualizer

Used when runs with stream lit

**get\_ui\_feedback()**

Gets the user defined config from the UI

**Returns** config

**log\_loss**(*loss*)

Adds a loss history to the chart

**Parameters** **loss** – a list of loss history

**Returns** None

**log\_reward**(*reward*)

Adds a reward history to the chart

**Parameters** **reward** – a list of reward history

**Returns**

**class** visualizer.TrainingVisualizer

Base training visualizer

**abstract** **get\_ui\_feedback()**

Gets the configuration from UI

**Returns** None

**abstract log\_loss**(*loss*)

Logs a loss history to the desired visualization

**Parameters** **loss** – a list of loss history

**Returns** None

**abstract log\_reward**(*reward*)

Logs a reward history to the desired visualization

**Parameters** **reward** – a list of reward history

**Returns** None

**visualizer.get\_training\_visualizer**(*visualizer\_type*)

A factory wrapper to generate training progress visualizers.

**Parameters** **visualizer\_type** – (str) the type of the visualizer to create

**Returns** TrainingVisualizer

---

CHAPTER  
THREE

---

## CLI ENTRYPPOINT

`entrypoint.main()`  
The CLI entrypoint to the APIs

**Returns** None



---

CHAPTER  
FOUR

---

## CONFIGURATION

### Config

```
config.DEFAULT_BATCH_SIZE = 128
    The default batch size the model should be trained on

config.DEFAULT_CHECKPOINT_LOCATION = './checkpoints'
    The default location to store the training checkpoints

config.DEFAULT_ENV_NAME = 'CartPole-v0'
    The OpenAI environment name to be used

config.DEFAULT_EPSILON = 0.05
    The default value for epsilon

config.DEFAULT_EVAL_EPS = 10
    The default number of episode the model should be evaluated with

config.DEFAULT_GAMMA = 0.95
    The default discount rate for the Q learning

config.DEFAULT_LEARNING_RATE = 0.001
    The default learning rate

config.DEFAULT_MAX_REPLAY_HISTORY = 1000000
    The default max length of the replay buffer

config.DEFAULT_MIN_STEPS = 10
    The minimum number of steps the evaluation should run per episode so that the tester can better visualize how the agent is doing.

config.DEFAULT_MODE = 'train'
    The default mode the program should run in

config.DEFAULT_MODEL_LOCATION = './model'
    The default location to store the best performing models

config.DEFAULT_NUM_ITERATIONS = 50000
    The default number of iteration to train the model

config.DEFAULT_PAUSE_TIME = 0
    The default value for pausing before execution starts to make time for screen recording. It's only available in testing mode since it's pointless to do so while training.

config.DEFAULT_RENDER_OPTION = 'none'
    The default value for rendering option

config.DEFAULT_TARGET_NETWORK_UPDATE_FREQUENCY = 120
    How often the target Q network should get parameter update from the training Q network.
```

```
config.DEFAULT_VERBOSITY_OPTION = 'progress'  
    The default verbosity option  
  
config.DEFAULT_VISUALIZER_TYPE = 'none'  
    The default visualizer type  
  
config.MODE_OPTIONS = ['train', 'test']  
    The supported modes  
  
config.RENDER_OPTIONS = ['none', 'collect']  
    The available render options:  
        • none: don't render anything  
        • collect: render the game play while collecting data  
  
config.VERBOSITY_OPTIONS = ['progress', 'loss', 'policy', 'init']  
    The available verbosity options:  
        • progress: show the training progress  
        • loss: show the logging information from loss calculation  
        • policy: show the logging information from policy generation  
        • init: show the logging information from initialization
```

## TRAIN

```
train.train_model(num_iterations=50000,      batch_size=128,      max_replay_history=1000000,
                  gamma=0.95,          eval_eps=10,         learning_rate=0.001,        tar-
                  get_network_update_frequency=120,      checkpoint_location='./checkpoints',
                  model_location='./model',    verbose='progress',   visualizer_type='none',    ren-
                  der_option='none',    persist_progress_option='all', epsilon=0.05)
```

Trains a DQN agent by playing episodes of the Cart Pole game

### Parameters

- **epsilon** – epsilon is the probability that a random action is chosen
- **target\_network\_update\_frequency** – how frequent target Q network gets updates
- **num\_iterations** – the number of episodes the agent will play
- **batch\_size** – the training batch size
- **max\_replay\_history** – the limit of the replay buffer length
- **gamma** – discount rate
- **eval\_eps** – the number of episode per evaluation
- **learning\_rate** – the learning rate of the back propagation
- **checkpoint\_location** – the location to save the training checkpoints
- **model\_location** – the location to save the pre-trained models
- **verbose** – the verbosity level which can be progress, loss, policy and init
- **visualizer\_type** – the type of visualization to be used
- **render\_option** – if the game play should be rendered
- **persist\_progress\_option** – if the training progress should be saved

**Returns** (maximum average reward, baseline average reward)



## UTILITIES

### Utilities

`utils.collect_episode(env, policy, buffer, render_option)`

Collect steps from a single episode play and record with replay buffer

#### Parameters

- `env` – OpenAI gym environment
- `policy` – DQN agent policy
- `buffer` – reinforcement learning replay buffer
- `render_option` – (bool) if should render the game play

#### Returns

None

`utils.collect_steps(env, policy, buffer, render_option, current_state, n_steps)`

Collects a single step from the game environment with policy specified. It is currently not used in favor of collect\_episode API.

#### Parameters

- `n_steps` – the number of steps to collect
- `current_state` – the current state of the environment
- `env` – OpenAI gym environment
- `policy` – DQN agent policy
- `buffer` – reinforcement learning replay buffer
- `render_option` – (bool) if should render the game play

#### Returns

None

`utils.compute_avg_reward(env, policy, num_episodes)`

Compute the average reward across num\_episodes under policy

#### Parameters

- `env` – OpenAI gym environment
- `policy` – DQN agent policy
- `num_episodes` – the number of episode to take average from

#### Returns

(int) average reward

`utils.play_episode(env, policy, render_option, min_steps)`

Play an episode with the given policy.

**Parameters**

- **min\_steps** – the minimum steps the game should be played
- **env** – the OpenAI gym environment
- **policy** – the policy that should be used to generate actions
- **render\_option** – how the game play should be rendered

**Returns** episode reward`utils.play_episodes (env, policy, render_option, num_eps, pause_time, min_steps)`

Play episodes with the given policy

**Parameters**

- **min\_steps** – the minimum steps the game should be played
- **pause\_time** – the time that should pause to prepare for screen recording
- **env** – the OpenAI gym environment
- **policy** – the policy that should be used to generate actions
- **render\_option** – how the game play should be rendered
- **num\_eps** – how many episodes should be played

**Returns** average episode reward

---

CHAPTER  
SEVEN

---

TESTS

Tests for model training

**class** train\_test.**TestTrain**(methodName='runTest')

Test suite for model training

**test\_sanity\_check()**

Tests if the model training finishes without crashing

**Returns** None

**test\_training\_effectiveness()**

Test if the model training can achieve a performance better than a random policy

**Returns** None



## PYTHON MODULE INDEX

### C

config, 9

### t

train\_test, 15

### u

utils, 13

### v

visualizer, 5



# INDEX

## C

can\_sample\_batch () (replay\_buffer.DqnReplayBuffer method), 3  
collect\_episode () (in module utils), 13  
collect\_policy () (dqn\_agent.DqnAgent method), 1  
collect\_steps () (in module utils), 13  
compute\_avg\_reward () (in module utils), 13  
config module, 9

## D

DEFAULT\_BATCH\_SIZE (in module config), 9  
DEFAULT\_CHECKPOINT\_LOCATION (in module config), 9  
DEFAULT\_ENV\_NAME (in module config), 9  
DEFAULT\_EPSILON (in module config), 9  
DEFAULT\_EVAL\_EPS (in module config), 9  
DEFAULT\_GAMMA (in module config), 9  
DEFAULT\_LEARNING\_RATE (in module config), 9  
DEFAULT\_MAX\_REPLAY\_HISTORY (in module config), 9  
DEFAULT\_MIN\_STEPS (in module config), 9  
DEFAULT\_MODE (in module config), 9  
DEFAULT\_MODEL\_LOCATION (in module config), 9  
DEFAULT\_NUM\_ITERATIONS (in module config), 9  
DEFAULT\_PAUSE\_TIME (in module config), 9  
DEFAULT\_RENDER\_OPTION (in module config), 9  
DEFAULT\_TARGET\_NETWORK\_UPDATE\_FREQUENCY (in module config), 9  
DEFAULT\_VERBOSITY\_OPTION (in module config), 9  
DEFAULT\_VISUALIZER\_TYPE (in module config), 10  
DqnAgent (class in dqn\_agent), 1  
DqnReplayBuffer (class in replay\_buffer), 3  
DummyTrainingVisualizer (class in visualizer), 5

## G

get\_training\_visualizer () (in module visualizer), 6  
get\_ui\_feedback () (visualizer.DummyTrainingVisualizer method), 5

get\_ui\_feedback () (visualizer.StreamlitTrainingVisualizer method), 5  
get\_ui\_feedback () (visualizer.TrainingVisualizer method), 5  
get\_volume () (replay\_buffer.DqnReplayBuffer method), 3

## L

load\_checkpoint () (dqn\_agent.DqnAgent method), 1  
load\_model () (dqn\_agent.DqnAgent method), 1  
log\_loss () (visualizer.DummyTrainingVisualizer method), 5  
log\_loss () (visualizer.StreamlitTrainingVisualizer method), 5  
log\_loss () (visualizer.TrainingVisualizer method), 6  
log\_reward () (visualizer.DummyTrainingVisualizer method), 5  
log\_reward () (visualizer.StreamlitTrainingVisualizer method), 5  
log\_reward () (visualizer.TrainingVisualizer method), 6

## M

main () (in module entrypoint), 7  
MODE\_OPTIONS (in module config), 10  
module config, 9  
train\_test, 15  
utils, 13  
visualizer, 5

## P

play\_episode () (in module utils), 13  
play\_episodes () (in module utils), 14  
policy () (dqn\_agent.DqnAgent method), 1

## R

random\_policy () (dqn\_agent.DqnAgent method), 1  
record () (replay\_buffer.DqnReplayBuffer method), 3  
RENDER\_OPTIONS (in module config), 10

## S

sample\_batch() (*replay\_buffer.DqnReplayBuffer method*), 3  
save\_checkpoint() (*dqn\_agent.DqnAgent method*), 1  
save\_model() (*dqn\_agent.DqnAgent method*), 1  
StreamlitTrainingVisualizer (*class in visualizer*), 5

## T

test\_sanity\_check() (*train\_test.TestTrain method*), 15  
test\_training\_effectiveness() (*train\_test.TestTrain method*), 15  
TestTrain (*class in train\_test*), 15  
train() (*dqn\_agent.DqnAgent method*), 1  
train\_model() (*in module train*), 11  
train\_test  
    module, 15  
TrainingVisualizer (*class in visualizer*), 5

## U

update\_target\_network() (*dqn\_agent.DqnAgent method*), 1  
utils  
    module, 13

## V

VERBOSITY\_OPTIONS (*in module config*), 10  
visualizer  
    module, 5