# CartPole-v1 with TensorFlow2

Problem [gym.openai.com/envs/CartPole-v1/](gym.openai.com/envs/CartPole-v1/)

## Setup

```
1 import numpy as np
2 import os
3 import pandas as pd
4 import random
5 from google.colab import drive
```

```
1 try:
2     import gym
3 except:
4     !pip install gym
5 import gym
6 print(gym.__version__)
```

```
    0.10.11
```

```
1 try:
2     import tensorflow as tf
3 except:
4     !pip install tensorflow==2.0.0-beta1
5     import tensorflow as tf
6 if tf.__version__[0] == "1":
7     !pip install tensorflow==2.0.0-beta1
8     import tensorflow as tf
9 print(tf.__version__)
```

```
    2.0.0-beta1
```

```
1 drive.mount('/content/gdrive', force_remount=False)
```

```
    Drive already mounted at /content/gdrive; to attempt to forcib
```

## Visualize gym environment - random action

```
1 import gym
2 env = gym.make("CartPole-v1")
3 no_actions = env.action_space.n
```

```
4 no_observations = env.observation_space.shape[0]
5 print(no_actions)
6 print(no_observations)
```

⊡→   2
     4

Action and observation interpretations: [github.com/openai/gym/blob/master/gym/envs/cl](github.com/openai/gym/blob/master/gym/envs/cl)

```
1 observation = env.reset()
2 for ep in range(10):
3     state = env.reset()
4     for t in range(100):
5 #         env.render()
6         action = env.action_space.sample()
7         observation, reward, done, info = env.step(action)
8         if done:
9             print("Episode finished after {} timesteps".format(t+1))
10            break
11 env.close()
```

⊡→   Episode finished after 32 timesteps
     Episode finished after 12 timesteps
     Episode finished after 25 timesteps
     Episode finished after 13 timesteps
     Episode finished after 12 timesteps
     Episode finished after 14 timesteps
     Episode finished after 33 timesteps
     Episode finished after 16 timesteps
     Episode finished after 11 timesteps
     Episode finished after 22 timesteps

## Train with TF2, Keras

```
1 def generate_model():
2     model = tf.keras.Sequential()
3     model.add(tf.keras.layers.Dense(
4         10,
5         activation="relu",
6         input_shape=(no_observations,)
7         )
8     )
9     model.add(tf.keras.layers.Dense(10, activation="relu"))
10    model.add(tf.keras.layers.Dense(no_actions, activation="linear"))
11    model.compile(loss="mse", optimizer=tf.keras.optimizers.Adam(lr=0.001))
12    return model
```

```
1 # ALPHA = 1
2 # GAMMA = 0.95
3 # model = generate_model()
4 # def fit_model(state, action, reward, next_state, done):
```

```
 5 #      q = model.predict(np.array([state]))
 6 #      q_next = model.predict(np.array([next_state]))
 7 #      q_update = reward
 8 #      if not done:
 9 #          q_update = ((1-ALPHA)*q[0, action] +
10 #                      ALPHA*(reward + GAMMA * q_next.max(axis=1)))
11 #      q[0, action] = q_update
12 #      model.fit(np.array([state]), q, epochs=1, batch_size=1)
```

```
 1 ALPHA = 1
 2 GAMMA = 0.95
 3 state_history = []
 4 action_history = []
 5 reward_history = []
 6 next_state_history = []
 7 model = generate_model()
 8 def fit_model(state, action, reward, next_state, done):
 9     state_history = []
10     action_history = []
11     next_state_history = []
12     reward_history = []
13     state_history.append(state)
14     action_history.append(action)
15     next_state_history.append(next_state)
16     reward_history.append(reward)
17     q = model.predict(np.array(state_history))
18     # print("q")
19     # print(q)
20     q_next = model.predict(np.array(next_state_history))
21     # print("q next")
22     # print(q_next)
23     q_update = reward_history
24     # print("q update")
25     # print(q_update)
26     if not done:
27         q_update = ((1-ALPHA)*q[np.arange(len(q)), action_history] +
28                     ALPHA*(reward_history + GAMMA * q_next.max(axis=1)))
29     q[np.arange(len(q)), action_history] = q_update
30     # print(q_update)
31     # print(next_best_actions)
32     # print(q)
33     model.fit(np.array(state_history), q, epochs=1, batch_size=32)
```

```
 1 %%capture
 2 for ep in range(1500):
 3     state = env.reset()
 4     for t in range(1000):
 5 #         env.render()
 6         if random.random() < 0.7:
 7             action = env.action_space.sample()
 8         else:
 9             action = model.predict(np.array([state])).argmax(axis=1)[0]
10         next_state, reward, done, info = env.step(action)
11         fit_model(state, action, reward, next_state, done)
```

```
11        IIt_model(state, action, reward, next_state, done)
12        state = next_state
13        if done:
14            print(f"Episode {ep} finished after {t} timesteps")
15            break
```

## Test with TF2, Keras

```
1 for ep in range(20):
2     state = env.reset()
3 #     env.render()
4     for t in range(1000):
5         action = model.predict(np.array([state])).argmax(axis=1)[0]
6         next_state, reward, done, info = env.step(action)
7         if done:
8             print(f"Episode {ep} finished after {t} timesteps")
9             break
10        state = next_state
11 env.close()
```

```
Episode 0 finished after 250 timesteps
Episode 1 finished after 292 timesteps
Episode 2 finished after 435 timesteps
Episode 3 finished after 255 timesteps
Episode 4 finished after 268 timesteps
Episode 5 finished after 239 timesteps
Episode 6 finished after 323 timesteps
Episode 7 finished after 227 timesteps
Episode 8 finished after 224 timesteps
Episode 9 finished after 245 timesteps
Episode 10 finished after 276 timesteps
Episode 11 finished after 262 timesteps
Episode 12 finished after 499 timesteps
Episode 13 finished after 228 timesteps
Episode 14 finished after 277 timesteps
Episode 15 finished after 294 timesteps
Episode 16 finished after 402 timesteps
Episode 17 finished after 240 timesteps
Episode 18 finished after 294 timesteps
Episode 19 finished after 499 timesteps
```

```
1 loc = "/content/gdrive/My Drive/Colab Notebooks/gym-openai-cartpole/model.h
2 model.save(loc)
3 loaded_model = tf.keras.models.load_model(loc)
4 loaded_model.summary()
```

```
Model: "sequential_1"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| dense_3 (Dense) | (None, 10) | 50 |
| dense_4 (Dense) | (None, 10) | 110 |

```
                            ─
_____
dense_5 (Dense)                (None, 2)                       22
===================================================================
Total params: 182
Trainable params: 182
Non-trainable params: 0
_____
```