

✓ Frozen Lake with Q-Learning

✓ Understanding the state and action space

- Link to study the environment: https://gymnasium.farama.org/environments/toy_text/frozen_lake/

```

1 import gymnasium as gym
2
3 env = gym.make("FrozenLake-v1",
4                 map_name = "4x4",
5                 is_slippery = False,
6                 render_mode = 'rgb_array')
7
8 print(f"The possible state space: {env.observation_space.n}")
9 print(f"The possible action space: {env.action_space.n}")

```

↗ The possible state space: 16
The possible action space: 4

✓ Create a Q-Learning Table

```

1 import numpy as np
2
3 # A function that creates Q Table based on the size of state space and action space size.
4 def initialize_q_table(state_space, action_space):
5     """
6     A function that creates Q Table based on the size of state space and action space size.
7
8     Args:
9         state_space (int): Number of available unique state in the state space.
10        action_space (int): Number of available actions in the state space.
11
12    Returns:
13        q_table (float array): A Q Table with the shape of (state_space, action_space).
14
15    """
16    q_table = np.zeros((state_space, action_space))
17    return q_table
18

```

✓ Create and test the Q-Learning Table

```

1 state_space = env.observation_space.n
2 action_space = env.action_space.n
3
4 QTable = initialize_q_table(state_space, action_space)
5 print(f"The shape of the Q Table is: {QTable.shape}")

```

↗ The shape of the Q Table is: (16, 4)

✓ Define Greedy Policy

- Always take the action that has the highest value.

```

1 def greedy_policy(q_table, state):
2     """
3     Greedy Policy - always take the action that has the highest q value within a state.
4
5     Args:
6         q_table (float array): A Q Table that has the size of (state_space, action_space).
7         state (int): The current state that the agent is in.
8
9     Returns:
10        action (int): The action to be taken by the agent under the greed_policy.
11
12    """
13    action = np.argmax(q_table[state])
14    return action

```

✓ Define Epsilon-Greedy Policy

- The agent has a probability of ϵ in taking a random action (Exploration) and a probability of $1-\epsilon$ in following Greedy Policy (Exploitation).

```

1 def epsilon_greedy_policy(q_table, state, epsilon, env):
2
3     """
4     Epsilon Greedy Policy - The agent has a probability of  $\epsilon$  in taking a random action (Exploration) and a probability of  $1-\epsilon$  in following Greedy Policy (Exploitation).
5
6     Args:
7         q_table (float array): A Q Table that has the size of (state_space, action_space).
8         state (int): The current state that the agent is in.
9         epsilon (float): A number that decides if exploration or exploitation.
10        env (gymnasium.env): The environment that the the agent is in.
11
12    Returns:
13        action (int): The action to be taken by the agent under the greed_policy.
14    """
15
16    # Sample a random number
17    probability = np.random.uniform(0, 1)
18
19    # Exploitation - Follow Greedy Policy
20    if probability > epsilon:
21        action = greedy_policy(q_table, state)
22    # Exploration - Take random action
23    else:
24        action = env.action_space.sample()
25
26    return action
27

```

✓ Defining Hyperparameters

```

1 # Training parameters
2 n_training_episodes = 10000 # Total training episodes
3 learning_rate = 0.7         # Learning rate
4
5 # Evaluation parameters
6 n_eval_episodes = 100       # Total number of test episodes
7
8 # Environment parameters
9 env_id = "FrozenLake-v1"    # Name of the environment
10 max_steps = 99              # Max steps per episode
11 gamma = 0.95                # Discounting rate
12 eval_seed = []              # The evaluation seed of the environment
13
14 # Exploration parameters
15 max_epsilon = 1.0           # Exploration probability at start
16 min_epsilon = 0.05          # Minimum exploration probability
17 decay_rate = 0.0005         # Exponential decay rate for exploration prob

```

✓ Create the training loop

Q-Learning

```

1 from tqdm.auto import tqdm
2
3 def train(n_training_episodes, max_steps, env, q_table, epsilon_min, epsilon_max, decay_rate, lr, gamma):
4
5     """
6     Training Loop for a Q-Learning Agent.
7
8     Args:
9         n_training_episodes (int): Number of episodes to be trained.
10        max_steps (int): Maximum number of steps per episodes.
11        env (gymnasium.env): The environment that the agent is in.
12        q_table (float array): A Q Table that has the size of (state_space, action_space).
13        epsilon_min (float): Lower bound for the epsilon value, expected between range of 0 to 1.
14        epsilon_max (float): Upper bound for the epsilon value, expected between range of 0 to 1.
15        decay_rate (float): Decay rate for the epsilon value.
16        lr (float): Learning rate for the agent.
17        gamma (float): Discount factor for the reward.
18
19    Returns:
20        q_table (float array): A trained Q Table that has the size of (state_space, action_space).
21    """
22
23
24    for episode in tqdm(range(n_training_episodes)):
25        # Reset status of termination or truncation
26        terminated = False
27        truncated = False
28
29        # Reset state
30        state, info = env.reset()
31
32        # Adjust epsilon
33        epsilon = epsilon_min + (epsilon_max - epsilon_min) * np.exp(-decay_rate * episode)
34
35        # Begin a new episode step by step
36        for step in range(max_steps):
37
38            # Sample an action based on epsilon_greedy_policy
39            action = epsilon_greedy_policy(q_table, state, epsilon, env)
40
41            # Perform the sampled action and observe the new state and received reward
42            new_state, reward, terminated, truncated, info = env.step(action)
43
44            # Update Q(s,a) := Q(s,a) + lr [R(s,a) + gamma * max Q(s',a') - Q(s,a)]
45            q_table[state][action] = q_table[state][action] + lr * (reward + gamma * (np.max(q_table[new_state])) - q_table[state][action])
46
47            # Check if reached end of episodes
48            if terminated or truncated:
49                break
50
51            # Update the state
52            state = new_state
53
54    return q_table
55
56

```

✓ Perform training

```

1 QTable_frozen_lake = train(n_training_episodes = n_training_episodes,
2                             max_steps = max_steps,
3                             env = env,
4                             q_table = QTable,
5                             epsilon_min = min_epsilon,
6                             epsilon_max = max_epsilon,
7                             decay_rate = decay_rate,
8                             lr = learning_rate,
9                             gamma = gamma)

```



100%

10000/10000 [00:02<00:00, 5055.18it/s]

✓ Create evaluation code

```

1 def evaluate_agent(env, max_steps, n_eval_episodes, q_table, seed):
2
3     """
4     Evaluation for a Q-Learning Agent.
5
6     Args:
7         n_eval_episodes (int): Number of episodes to be evaluated.
8         max_steps (int): Maximum number of steps per episodes.
9         env (gymnasium.env): The environment that the agent is in.
10        q_table (float array): A Q Table that has the size of (state_space, action_space).
11        seed (int list): A list that consists of the initial state of the environment.
12
13    Returns:
14        mean_reward (float): Mean reward received over the evaluated episodes.
15        std_reward (float): Standard deviation reward received over the evaluated episodes.
16    """
17
18    # Initialize a list to store reward from each episode
19    episode_rewards = []
20
21    # Begin each episode
22    for episode in tqdm(range(n_eval_episodes)):
23
24        # Reset the environment parameters
25        if seed:
26            state, info = env.reset(seed = seed[episode])
27        else:
28            state, info = env.reset()
29        terminated = False
30        truncated = False
31
32        # Reset the reward per episode
33        reward_per_ep = 0
34
35        # Repeat steps
36        for step in range(max_steps):
37
38            # Sample an action using greedy_policy and step with that
39            action = greedy_policy(q_table, state)
40            state, reward, terminated, truncated, info = env.step(action)
41
42            # Add the reward
43            reward_per_ep += reward
44
45            # Check if terminated or truncated
46            if terminated or truncated:
47                break
48
49        # Append the episode for this episode to the list
50        episode_rewards.append(reward_per_ep)
51
52    # End of all episodes - Calculate mean and standard deviation of reward
53    mean_reward = np.mean(episode_rewards).item()
54    std_reward = np.std(episode_rewards).item()
55
56    return mean_reward, std_reward
57


```

✓ Perform evaluation

```

1 mean_reward, std_reward = evaluate_agent(n_eval_episodes = n_eval_episodes,
2                                           max_steps = max_steps,
3                                           env = env,
4                                           q_table = QTable_frozen_lake,
5                                           seed = None)
6
7 print(f"The trained agent's performance: {mean_reward:.2f} +/- {std_reward:.2f}")

```

 100% 100/100 [00:00<00:00, 4137.25it/s]

The trained agent's performance: 1.00 +/- 0.00

✓ Functions made to push to Hugging Face Hub

- This code is provided by the tutorial: <https://colab.research.google.com/github/huggingface/deep-rl-class/blob/master/notebooks/unit2/unit2.ipynb#scrollTo=paOynXy3aoJW>
- **package_to_hub** cannot be used to push models to hub because the architecture used here is custom made and not from stable-baselines3 as in Unit 1.

```

1 # Import Library
2
3 from huggingface_hub import HfApi, snapshot_download
4 from huggingface_hub.repocard import metadata_eval_result, metadata_save
5
6 from pathlib import Path
7 import datetime
8 import json
9
10 import imageio
11 import random
12 import pickle

```

```

1 # A function made to record an episode to be showcased on the model card.
2
3 def record_video(env, Qtable, out_directory, fps=1):
4     """
5     Generate a replay video of the agent
6     :param env
7     :param Qtable: Qtable of our agent
8     :param out_directory
9     :param fps: how many frame per seconds (with taxi-v3 and frozenlake-v1 we use 1)
10    """
11    images = []
12    terminated = False
13    truncated = False
14    state, info = env.reset(seed=random.randint(0,500))
15    img = env.render()
16    images.append(img)
17    while not terminated or truncated:
18        # Take the action (index) that have the maximum expected future reward given that state
19        action = np.argmax(Qtable[state][:])
20        state, reward, terminated, truncated, info = env.step(action) # We directly put next_state = state for recording logic
21        img = env.render()
22        images.append(img)
23    imageio.mimsave(out_directory, [np.array(img) for i, img in enumerate(images)], fps=fps)

```

```

1 # Push to Hub
2
3 def push_to_hub(
4     repo_id, model, env, video_fps=1, local_repo_path="hub"
5 ):
6     """
7     Evaluate, Generate a video and Upload a model to Hugging Face Hub.
8     This method does the complete pipeline:
9     - It evaluates the model
10    - It generates the model card
11    - It generates a replay video of the agent
12    - It pushes everything to the Hub
13
14    :param repo_id: repo_id: id of the model repository from the Hugging Face Hub
15    :param env
16    :param video_fps: how many frame per seconds to record our video replay
17    (with taxi-v3 and frozenlake-v1 we use 1)
18    :param local_repo_path: where the local repository is
19    """
20    _, repo_name = repo_id.split("/")
21
22    eval_env = env
23    api = HfApi()
24
25    # Step 1: Create the repo
26    repo_url = api.create_repo(
27        repo_id=repo_id,
28        exist_ok=True,
29    )
30
31    # Step 2: Download files
32    repo_local_path = Path(snapshot_download(repo_id=repo_id))
33
34    # Step 3: Save the model
35    if env.spec.kwarg.get("map_name"):
36        model["map_name"] = env.spec.kwarg.get("map_name")
37    if env.spec.kwarg.get("is_slippery", "") == False:

```

```

38     model["slippery"] = False
39
40     # Pickle the model
41     with open((repo_local_path) / "q-learning.pkl", "wb") as f:
42         pickle.dump(model, f)
43
44     # Step 4: Evaluate the model and build JSON with evaluation metrics
45     mean_reward, std_reward = evaluate_agent(
46         eval_env, model["max_steps"], model["n_eval_episodes"], model["qtable"], model["eval_seed"]
47     )
48
49     evaluate_data = {
50         "env_id": model["env_id"],
51         "mean_reward": mean_reward,
52         "n_eval_episodes": model["n_eval_episodes"],
53         "eval_datetime": datetime.datetime.now().isoformat()
54     }
55
56     # Write a JSON file called "results.json" that will contain the
57     # evaluation results
58     with open(repo_local_path / "results.json", "w") as outfile:
59         json.dump(evaluate_data, outfile)
60
61     # Step 5: Create the model card
62     env_name = model["env_id"]
63     if env.spec.kwarg.get("map_name"):
64         env_name += "-" + env.spec.kwarg.get("map_name")
65
66     if env.spec.kwarg.get("is_slippery", "") == False:
67         env_name += "-" + "no_slippery"
68
69     metadata = {}
70     metadata["tags"] = [env_name, "q-learning", "reinforcement-learning", "custom-implementation"]
71
72     # Add metrics
73     eval = metadata_eval_result(
74         model_pretty_name=repo_name,
75         task_pretty_name="reinforcement-learning",
76         task_id="reinforcement-learning",
77         metrics_pretty_name="mean_reward",
78         metrics_id="mean_reward",
79         metrics_value=f"{mean_reward:.2f} +/- {std_reward:.2f}",
80         dataset_pretty_name=env_name,
81         dataset_id=env_name,
82     )
83
84     # Merges both dictionaries
85     metadata = {**metadata, **eval}
86
87     model_card = f"""
88     # **Q-Learning** Agent playing **{env_id}**
89     This is a trained model of a **Q-Learning** agent playing **{env_id}** .
90
91     ## Usage
92
93     ```python
94
95     model = load_from_hub(repo_id="{repo_id}", filename="q-learning.pkl")
96
97     env = gym.make(model["env_id"])
98     ```
99     """
100
101     evaluate_agent(env, model["max_steps"], model["n_eval_episodes"], model["qtable"], model["eval_seed"])
102
103     readme_path = repo_local_path / "README.md"
104     readme = ""
105     print(readme_path.exists())
106     if readme_path.exists():
107         with readme_path.open("r", encoding="utf8") as f:
108             readme = f.read()
109     else:
110         readme = model_card
111
112     with readme_path.open("w", encoding="utf-8") as f:
113         f.write(readme)
114
115     # Save our metrics to Readme metadata
116     metadata_save(readme_path, metadata)
117
118     # Step 6: Record a video
119     video_path = repo_local_path / "replay.mp4"

```

```

120 record_video(env, model["qtable"], video_path, video_fps)
121
122 # Step 7. Push everything to the Hub
123 api.upload_folder(
124     repo_id=repo_id,
125     folder_path=repo_local_path,
126     path_in_repo=".",
127 )
128
129 print("Your model is pushed to the Hub. You can view your model here: ", repo_url)

```

```

1 # Create the dictionary that describe the model.
2
3 model = {
4     "env_id": env_id,
5     "max_steps": max_steps,
6     "n_training_episodes": n_training_episodes,
7     "n_eval_episodes": n_eval_episodes,
8     "eval_seed": eval_seed,
9
10    "learning_rate": learning_rate,
11    "gamma": gamma,
12
13    "max_epsilon": max_epsilon,
14    "min_epsilon": min_epsilon,
15    "decay_rate": decay_rate,
16
17    "qtable": QTable_frozen_lake
18 }

```

✓ Push to Hugging Face Hub

- Create a new token with write role here: <https://huggingface.co/settings/tokens>

```

1 from huggingface_hub import notebook_login
2
3 notebook_login()

```



```

1 push_to_hub(repo_id = "wengti0608/q-FrozenLake-v1-4x4-noSlippery",
2             model = model,
3             env = env)

```



```

Fetching 5 files: 100%                               5/5 [00:00<00:00, 3.37it/s]
q-learning.pkl: 100%                                915/915 [00:00<00:00, 1.37kB/s]
results.json: 100%                                  118/118 [00:00<00:00, 6.93kB/s]
100%                                                  100/100 [00:00<00:00, 3321.25it/s]
100%                                                  100/100 [00:00<00:00, 2850.90it/s]
True
Uploading...: 100%                                   915/915 [00:00<00:00, 2.54kB/s]
Your model is pushed to the Hub. You can view your model here: https://huggingface.co/wengti0608/q-FrozenLake-v1-4x4-noSlippery

```

✓ Taxi with Q-Learning

- Link to study the environment: https://gymnasium.farama.org/environments/toy_text/taxi/

✓ Understanding the state and action space

```

1 import gymnasium as gym
2
3 env = gym.make("Taxi-v3",
4               render_mode = 'rgb_array')
5
6 print(f"The possible state space: {env.observation_space.n}")
7 print(f"The possible action space: {env.action_space.n}")

```




```

The possible state space: 500
The possible action space: 6

```

✓ Create a Q-Learning Table

```
1 state_space = env.observation_space.n
2 action_space = env.action_space.n
3
4 q_table = initialize_q_table(state_space, action_space)
5
6 print(f"The shape of the q_table: {q_table.shape}")
```

 The shape of the q_table: (500, 6)

✓ Define the hyperparameter

```
1 # Training parameters
2 n_training_episodes = 25000 # Total training episodes
3 learning_rate = 0.7 # Learning rate
4
5 # Evaluation parameters
6 n_eval_episodes = 100 # Total number of test episodes
7
8 # DO NOT MODIFY EVAL_SEED
9 eval_seed = [16,54,165,177,191,191,120,80,149,178,48,38,6,125,174,73,50,172,100,148,146,6,25,40,68,148,49,167,9,97,164,176,61,7,54,161,131,184,51,170,12,120,113,95,126,51,98,36,135,54,82,45,95,89,59,95,124,9,113,58,85,51,134,121,169,105,21,30,11,50,65,12,43,82,112,102,168,123,97,21,83,158,26,80,63,5,81,32,11,28,148] # Evaluation seed, this ensures that all classmates agents are trained on
12 # Each seed has a specific starting state
13
14 # Environment parameters
15 env_id = "Taxi-v3" # Name of the environment
16 max_steps = 99 # 99 # Max steps per episode
17 gamma = 0.95 # Discounting rate
18
19 # Exploration parameters
20 max_epsilon = 1.0 # Exploration probability at start
21 min_epsilon = 0.05 # Minimum exploration probability
22 decay_rate = 0.005 # Exponential decay rate for exploration prob
```

✓ Perform Training

```
1 QTable_taxi = train(n_training_episodes = n_training_episodes,
2                     max_steps = max_steps,
3                     env = env,
4                     q_table = q_table,
5                     epsilon_min = min_epsilon,
6                     epsilon_max = max_epsilon,
7                     decay_rate = decay_rate,
8                     lr = learning_rate,
9                     gamma = gamma)
10
```

 100% 25000/25000 [01:41<00:00, 243.04it/s]

✓ Perform evaluation

```
1 mean_reward, std_reward = evaluate_agent(env = env,
2                                         max_steps = max_steps,
3                                         n_eval_episodes = n_eval_episodes,
4                                         q_table = QTable_taxi,
5                                         seed = eval_seed)
6
7 print(f"The agent has a performance of: {mean_reward:.2f} +/- {std_reward:.2f}")
```

 100% 100/100 [00:00<00:00, 1637.52it/s]

The agent has a performance of: 7.56 +/- 2.71

✓ Push to Hugging Face Hub

- Create a new token with write role here: <https://huggingface.co/settings/tokens>


```
1 from huggingface_hub import notebook_login
2 notebook_login()
```



Invalid user token.

```
1 # Create the dictionary that describe the model.
2
3 model = {
4     "env_id": env_id,
5     "max_steps": max_steps,
6     "n_training_episodes": n_training_episodes,
7     "n_eval_episodes": n_eval_episodes,
8     "eval_seed": eval_seed,
9
10    "learning_rate": learning_rate,
11    "gamma": gamma,
12
13    "max_epsilon": max_epsilon,
14    "min_epsilon": min_epsilon,
15    "decay_rate": decay_rate,
16
17    "qtable": QTable_taxi
18 }
```

```
1 push_to_hub(repo_id = "wengti0608/q-Taxi-v3_note",
2             model = model,
3             env = env)
```



Show hidden output

✓ Load models downloaded from Hugging Face Hub

- Cannot use built-in **load_from_hub** because this is a custom made model.
- The following code is provided by the tutorial: https://colab.research.google.com/github/huggingface/deep-rl-class/blob/master/notebooks/unit2/unit2.ipynb#scrollTo=AB6n_hhg7YS

```
1 from huggingface_hub import hf_hub_download
2
3 def load_from_hub(repo_id: str, filename: str) -> str:
4     """
5     Download a model from Hugging Face Hub.
6     :param repo_id: id of the model repository from the Hugging Face Hub
7     :param filename: name of the model zip file from the repository
8     """
9     # Get the model from the Hub, download and cache the model on your local disk
10    pickle_model = hf_hub_download(
11        repo_id=repo_id,
12        filename=filename
13    )
14
15    with open(pickle_model, 'rb') as f:
16        downloaded_model_file = pickle.load(f)
17
18    return downloaded_model_file
```

```
1 model = load_from_hub(repo_id = "wengti0608/q-Taxi-v3",
2                       filename = "q-learning.pkl")
3
4 mean_reward, std_reward = evaluate_agent(env = gym.make(model['env_id']),
5                                         max_steps = model['max_steps'],
6                                         n_eval_episodes = model['n_eval_episodes'],
7                                         q_table = model['qtable'],
8                                         seed = model['eval_seed'])
9
10 print(f"The model has a performance of {mean_reward:.2f} +/- {std_reward:.2f}")
```



q-learning.pkl: 100%

24.6k/24.6k [00:00<00:00, 25.9kB/s]

100%

100/100 [00:00<00:00, 3158.08it/s]

The model has a performance of 7.56 +/- 2.71