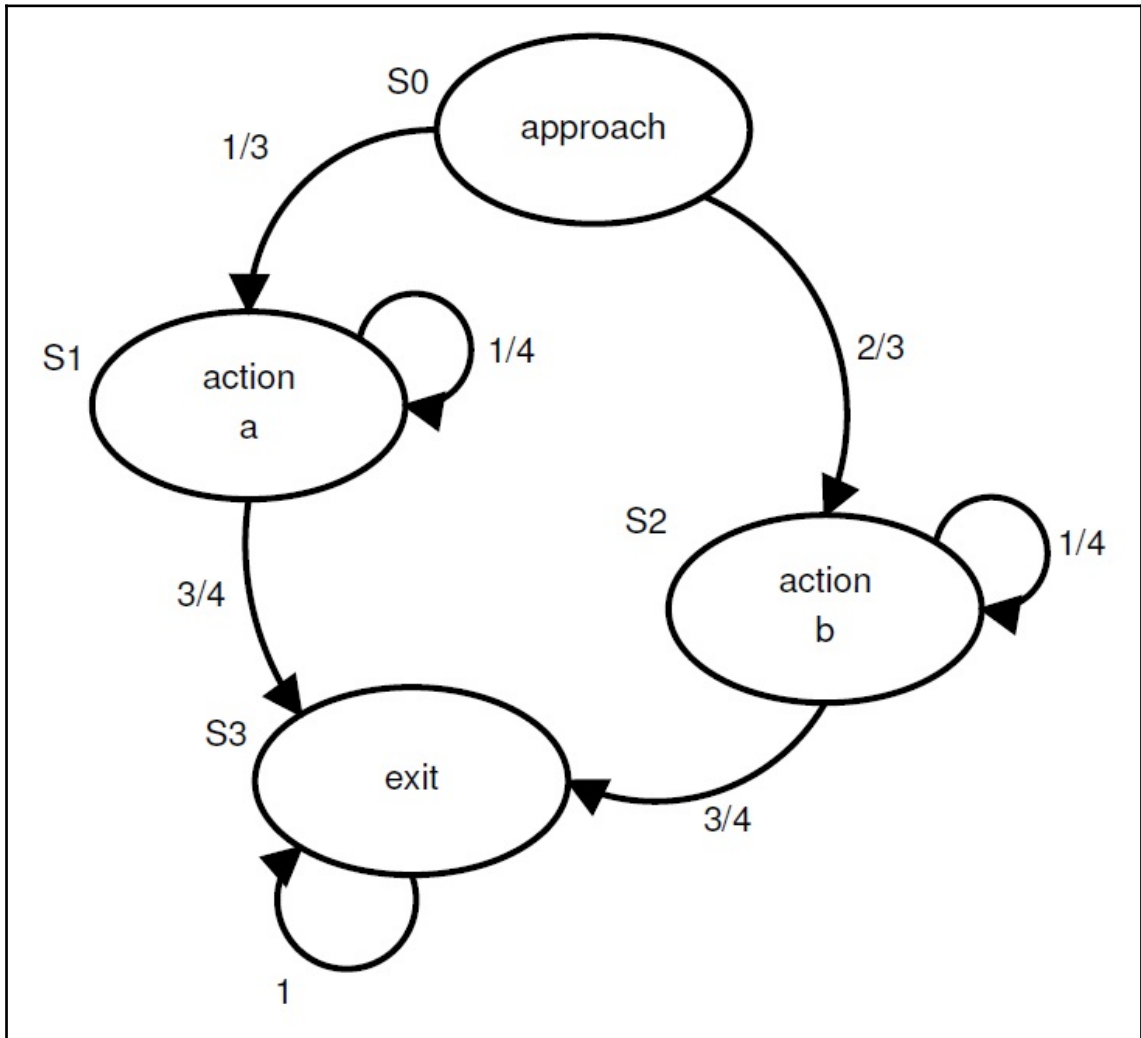
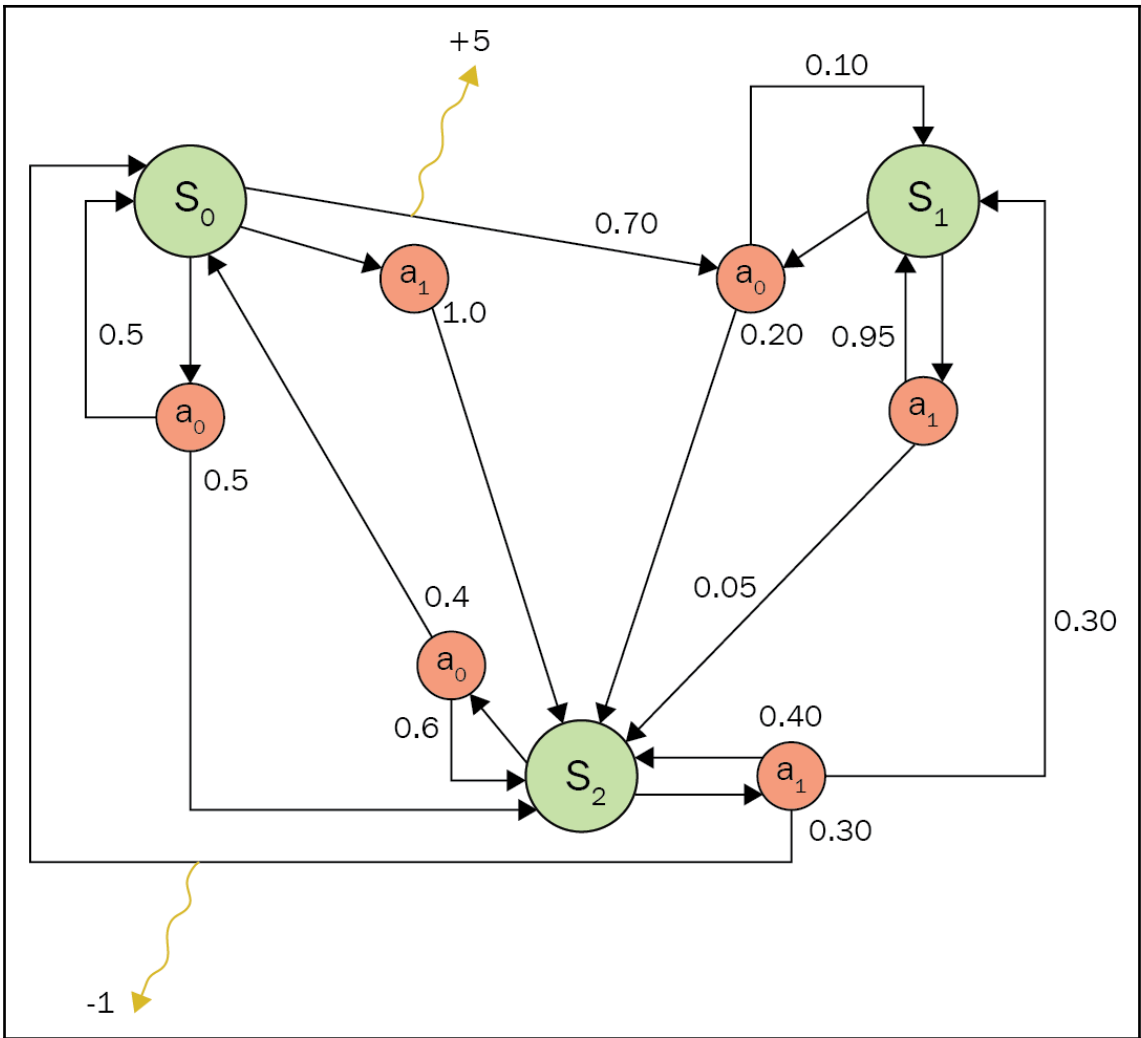


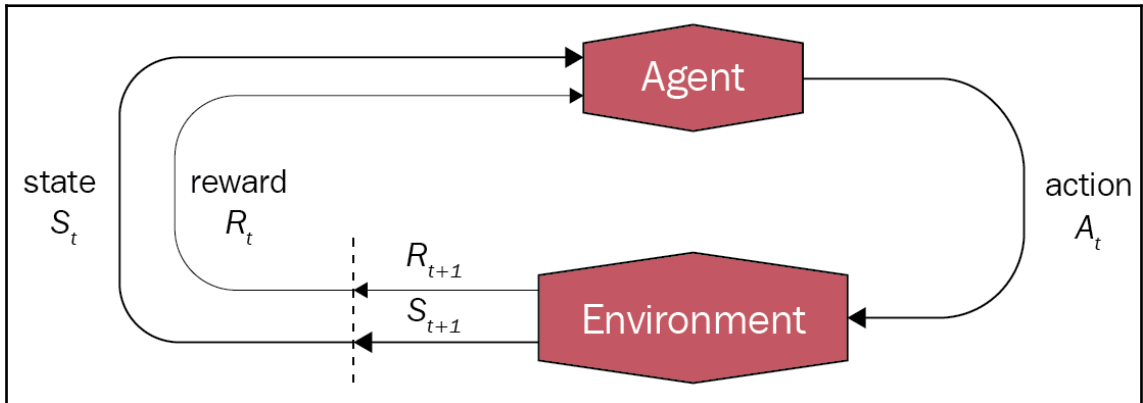
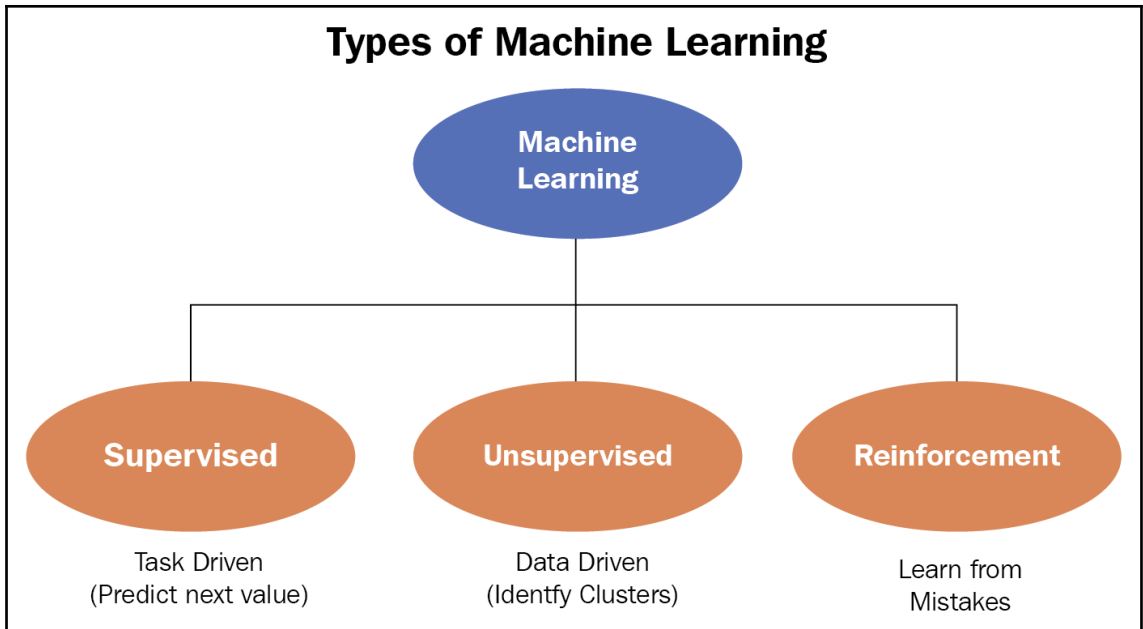
# Chapter 1, Brushing Up on Reinforcement Learning Concepts





---

## Types of Machine Learning





```
import gym
env = gym.make("Taxi-v2")
env.render()

+-----+
|R: | : :G|
| : : : :|
| : | : : :| |
| | : | : |
|Y| : |B: |
+-----+
```

```
print("Action Space {}".format(env.action_space))
print("State Space {}".format(env.observation_space))

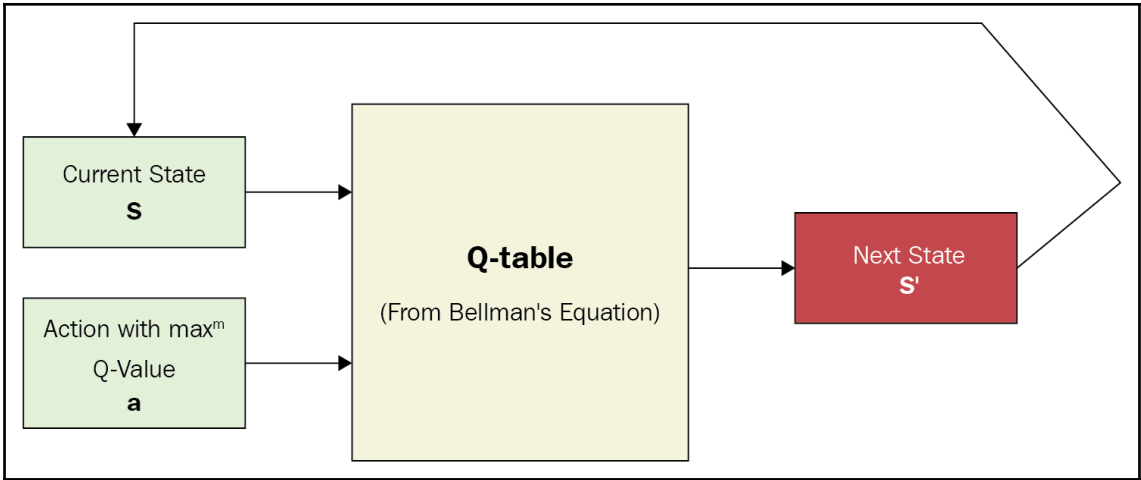
Action Space Discrete(6)
State Space Discrete(500)
```

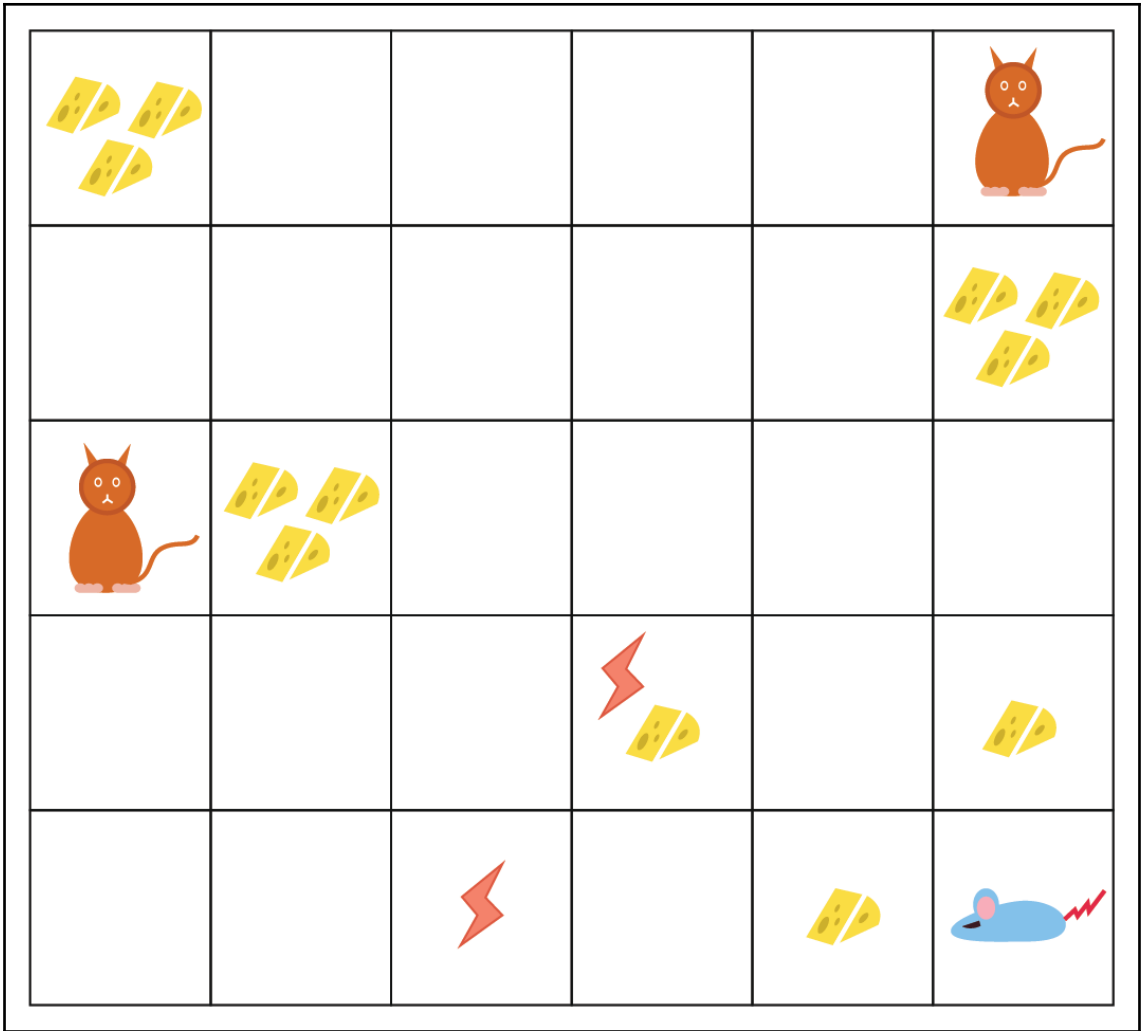
Initialized

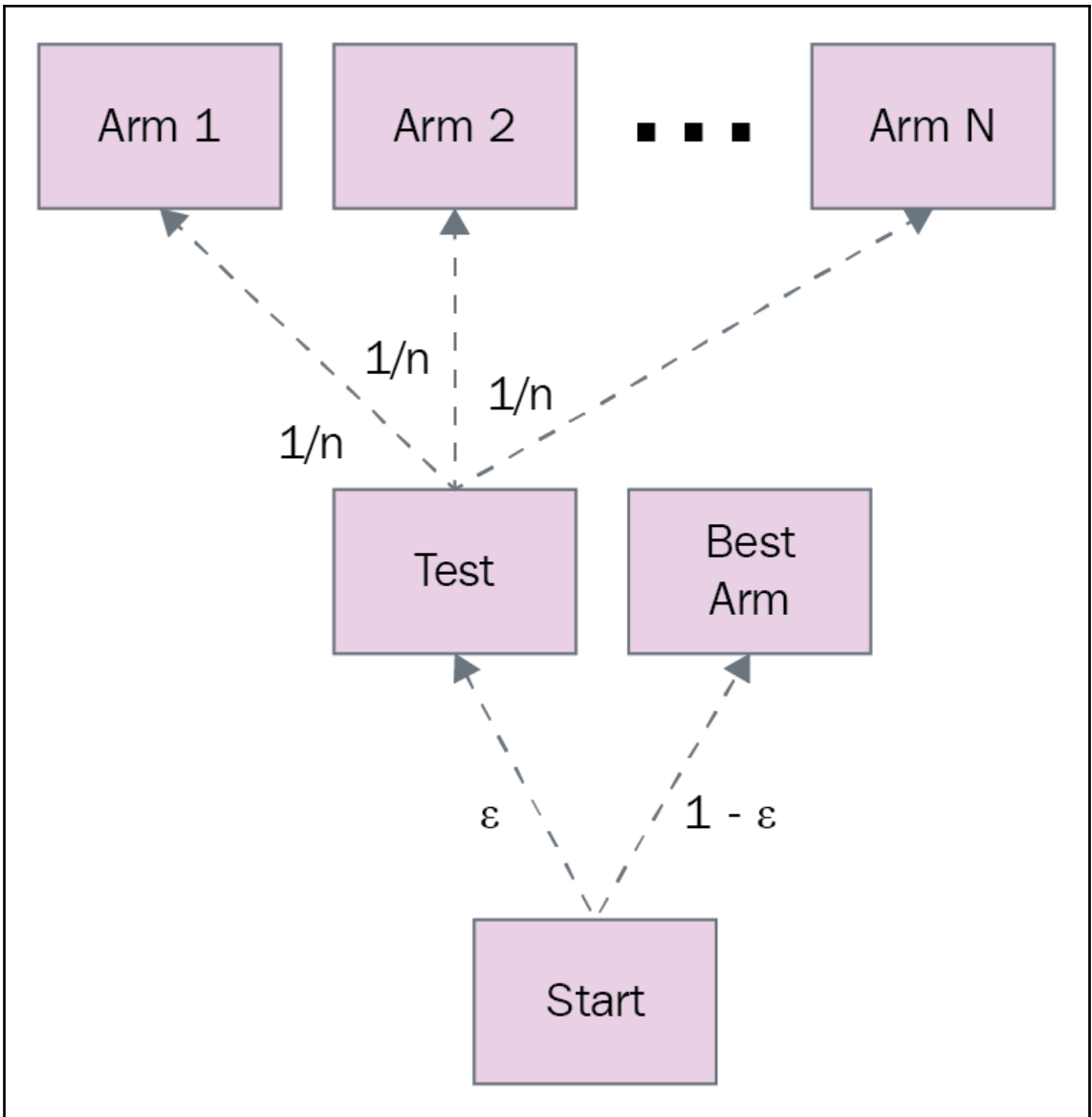
Q-Table		Actions					
		South (0)	North (1)	East (2)	West (3)	Pickup (4)	Dropoff (5)
States	0	0	0	0	0	0	0
	.	.	.	.	.	.	.
	.	.	.	.	.	.	.
	.	.	.	.	.	.	.
	327	0	0	0	0	0	0
	.	.	.	.	.	.	.
.	.	.	.	.	.	.	
.	.	.	.	.	.	.	
499	0	0	0	0	0	0	

Training

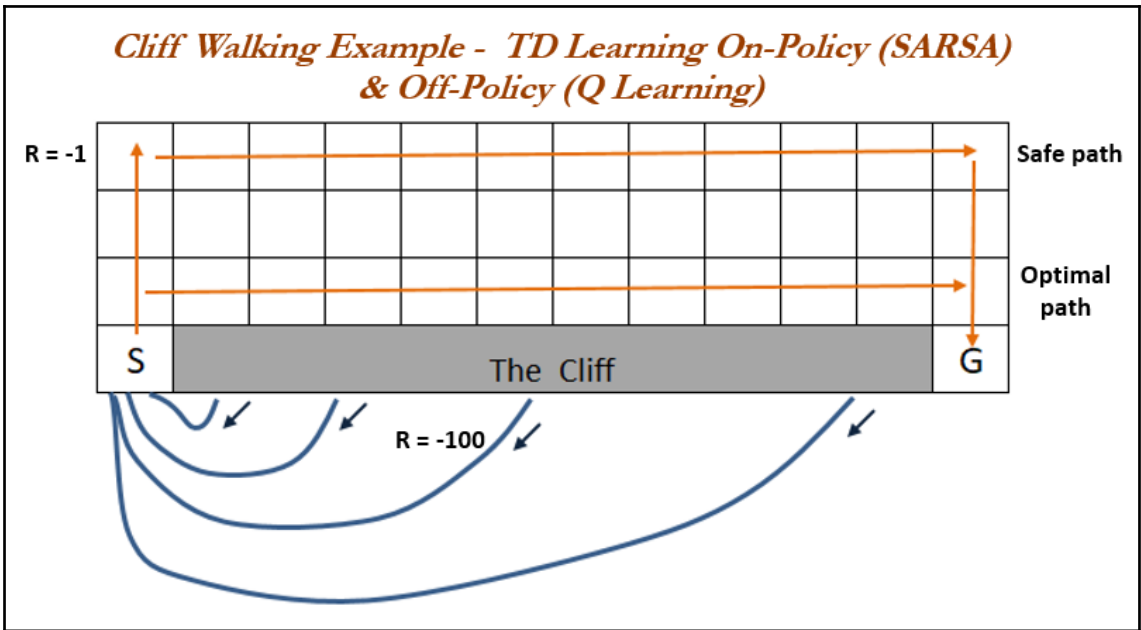
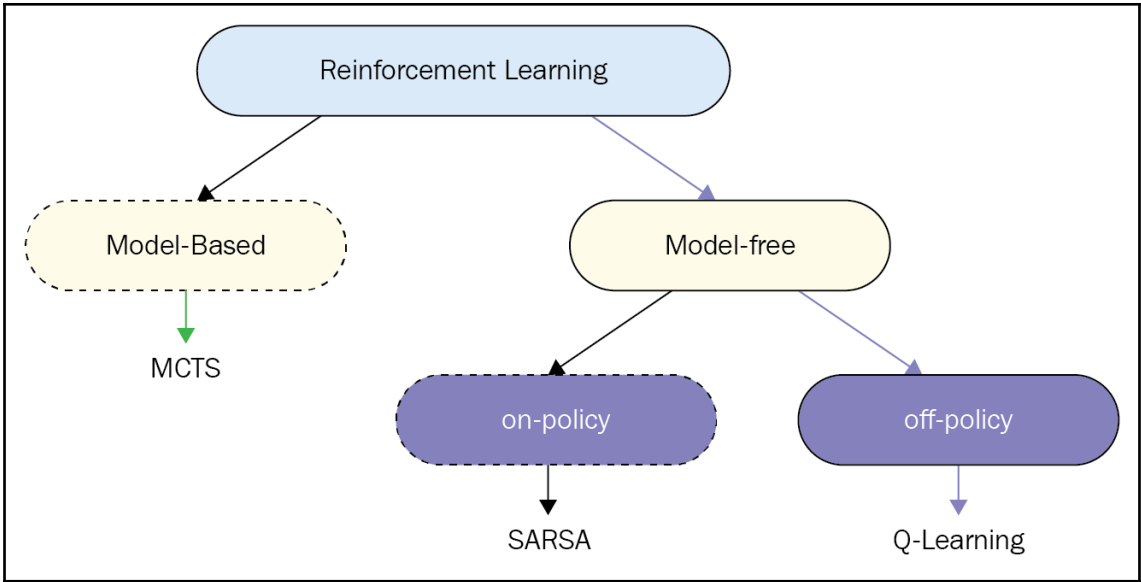
Q-Table		Actions					
		South (0)	North (1)	East (2)	West (3)	Pickup (4)	Dropoff (5)
States	0	0	0	0	0	0	0
	.	.	.	.	.	.	.
	.	.	.	.	.	.	.
	.	.	.	.	.	.	.
	328	-2.30108105	-1.97092096	-2.30357004	-2.20591839	-10.3607344	-8.5583017
	.	.	.	.	.	.	.
.	.	.	.	.	.	.	
.	.	.	.	.	.	.	
499	9.96984239	4.02706992	12.96022777	29	3.32877873	3.38230603	





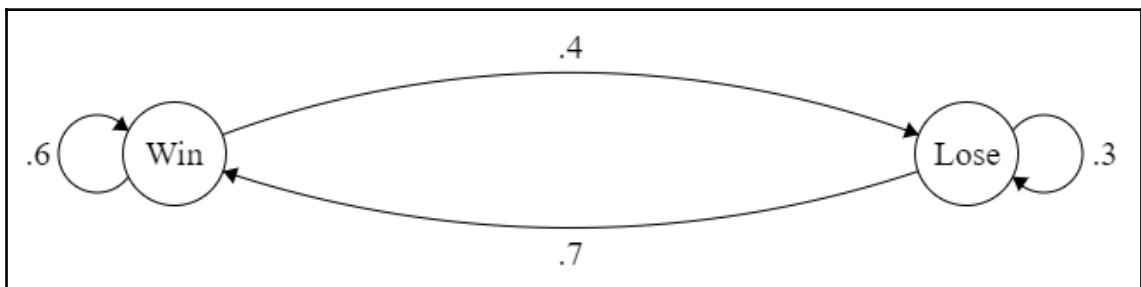
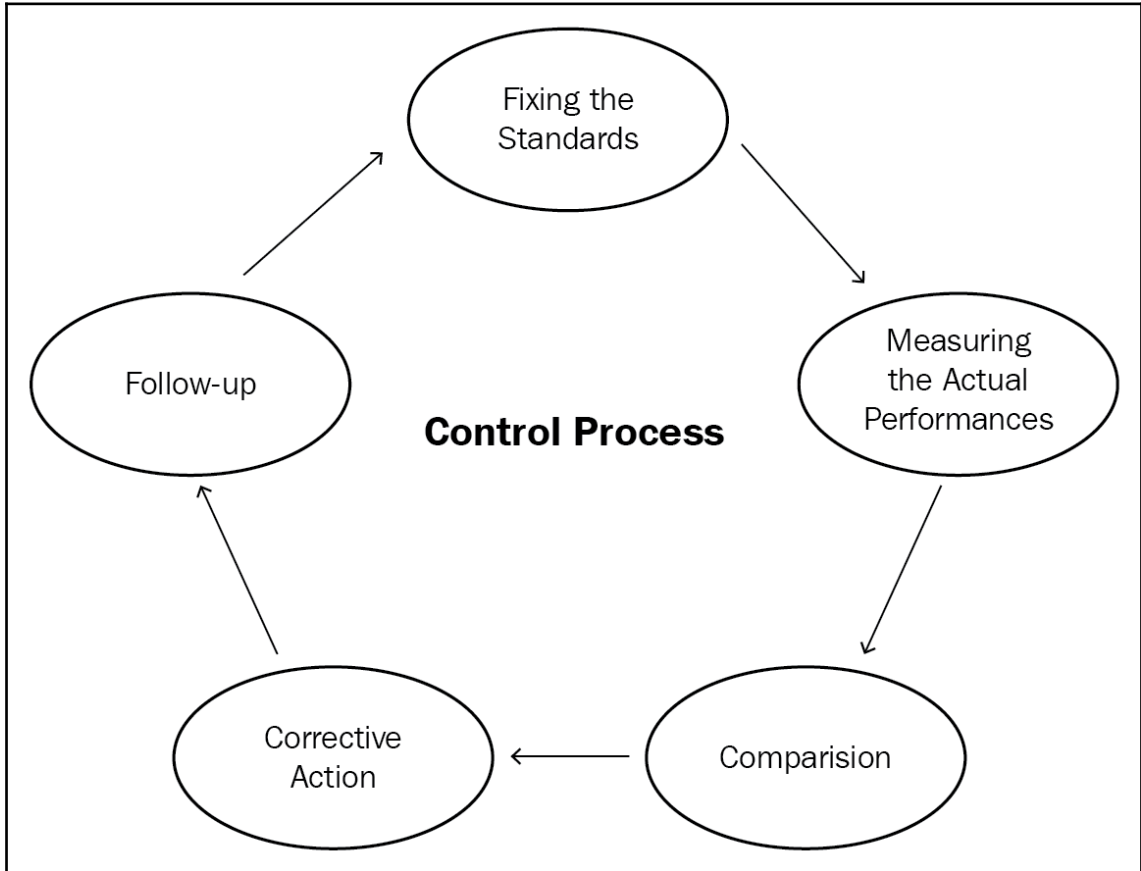


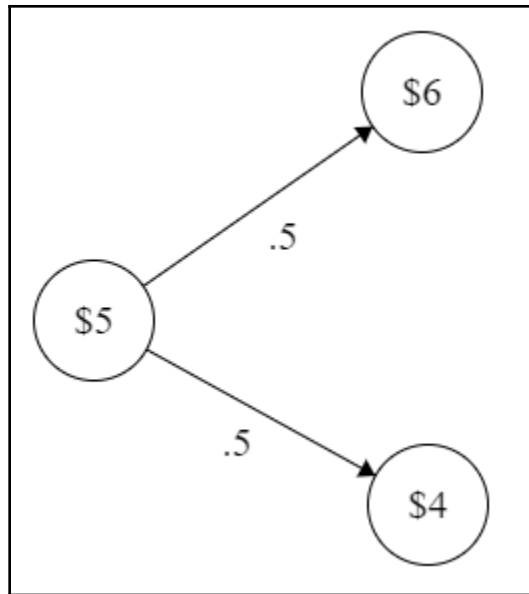


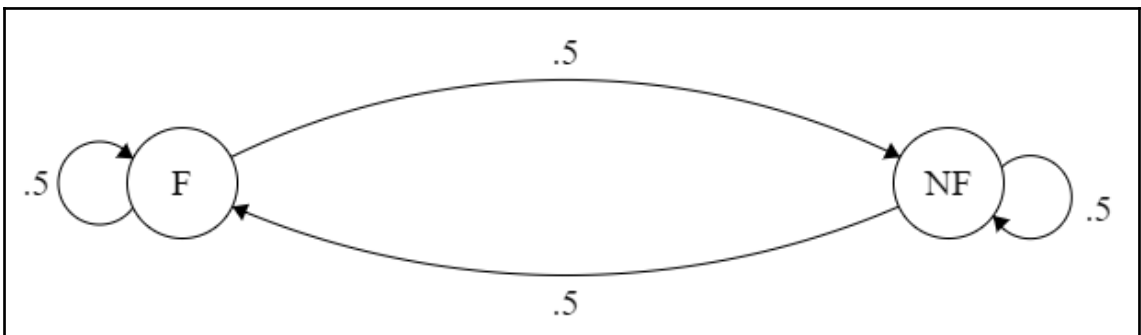
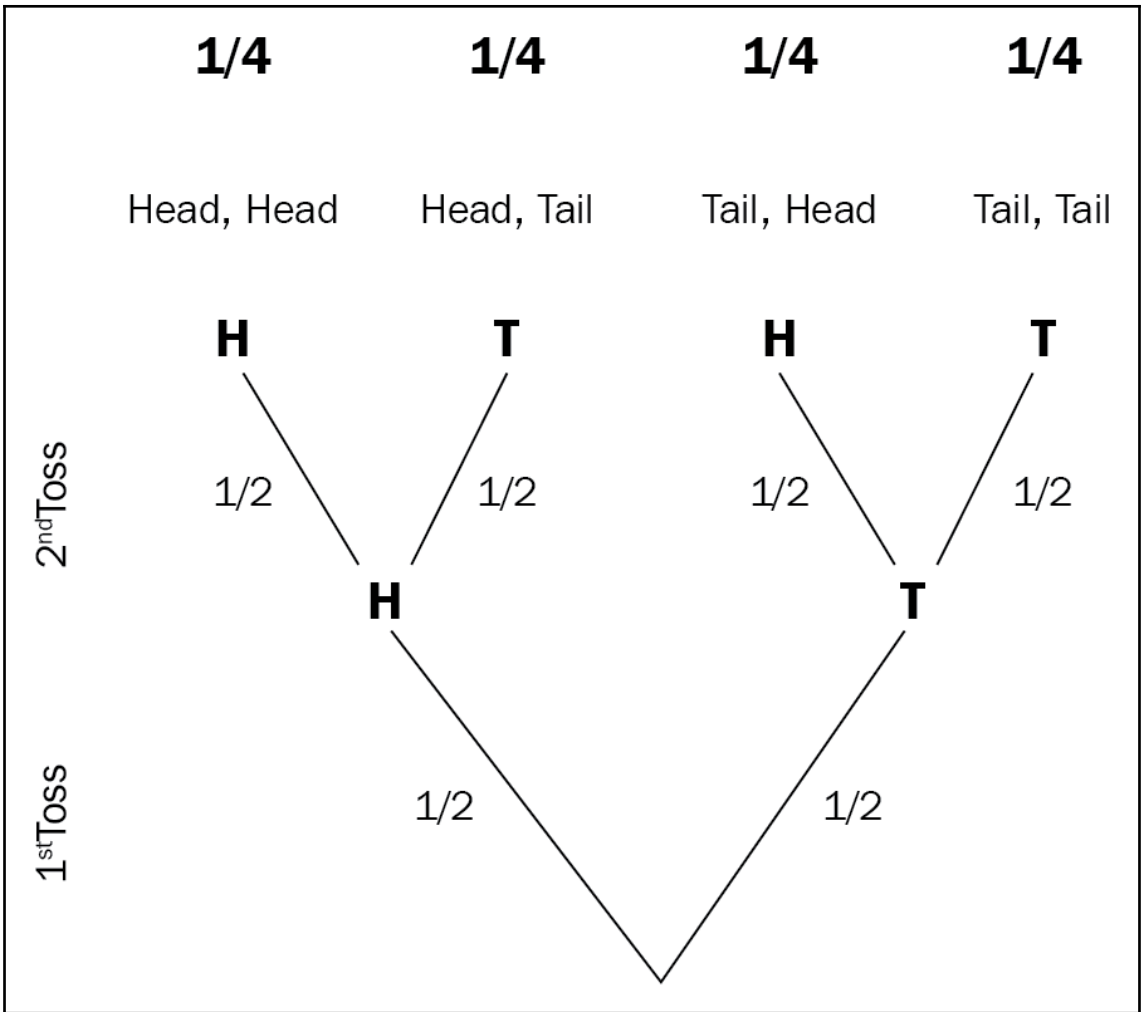


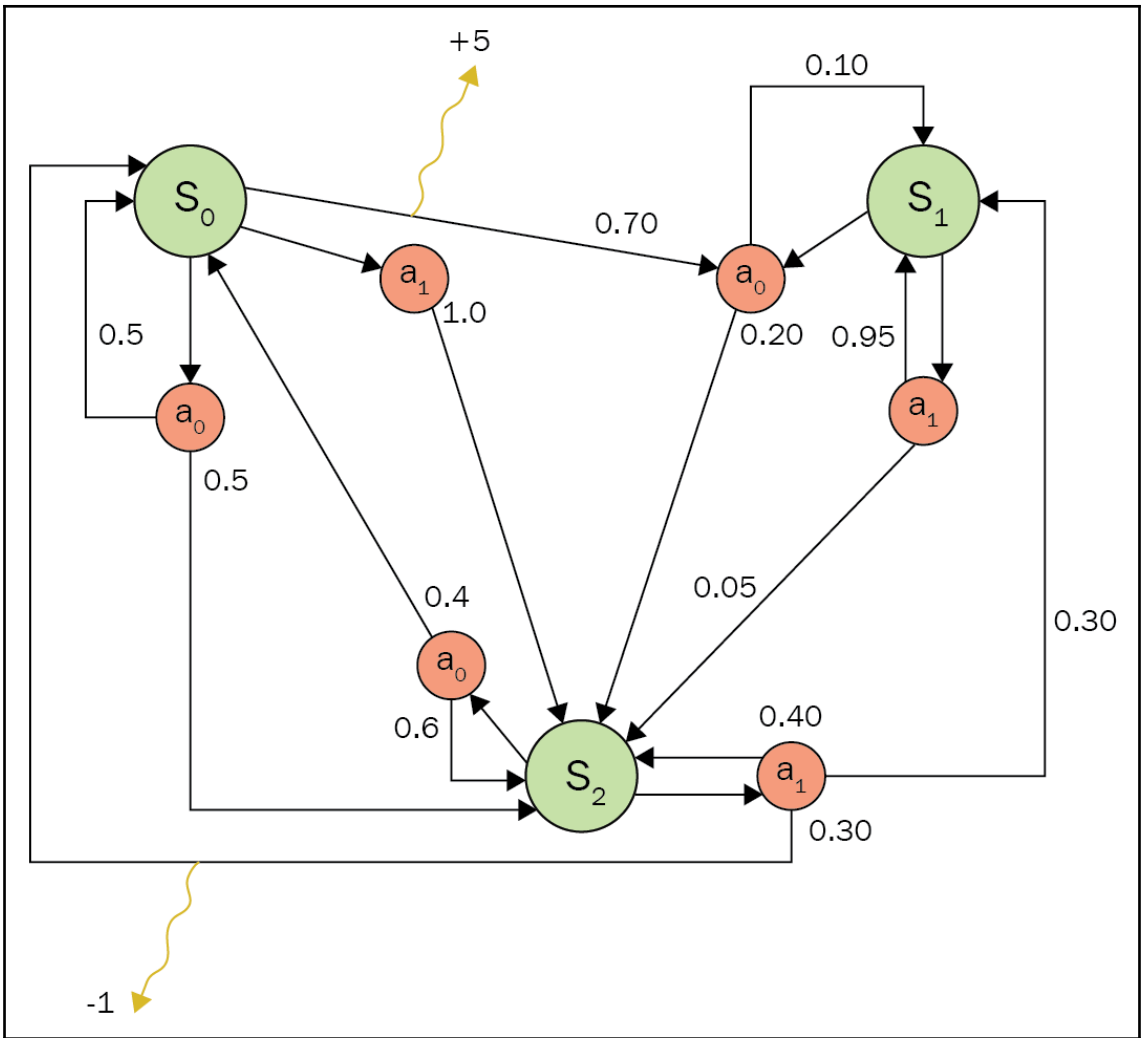
---

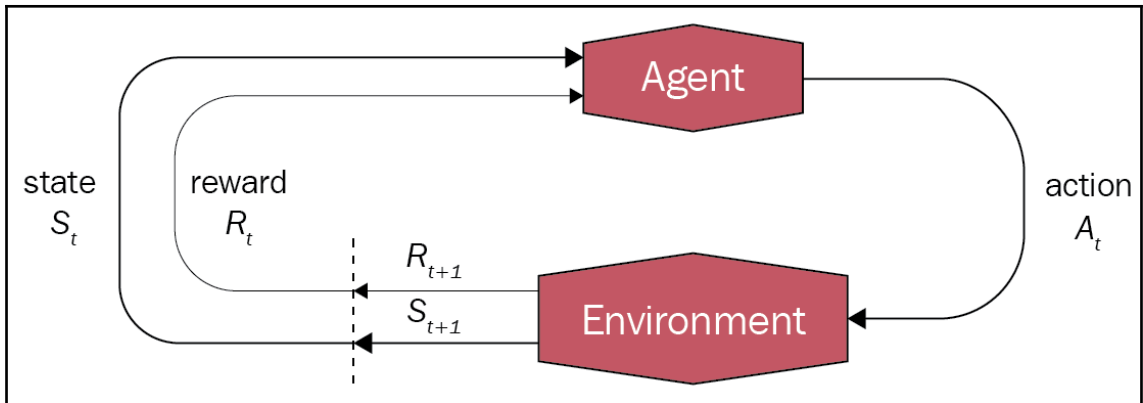
## Chapter 2, Getting Started with the Q-Learning Algorithm









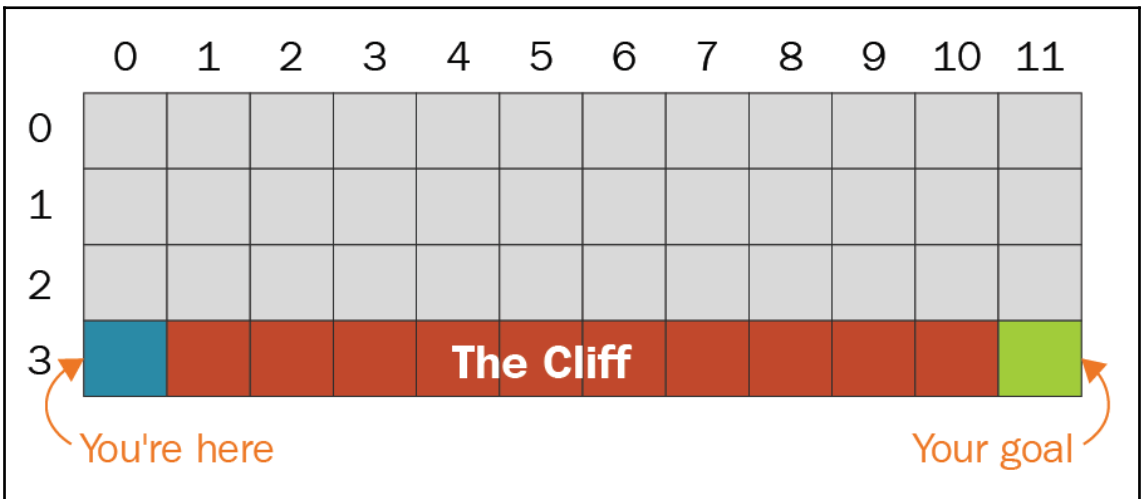
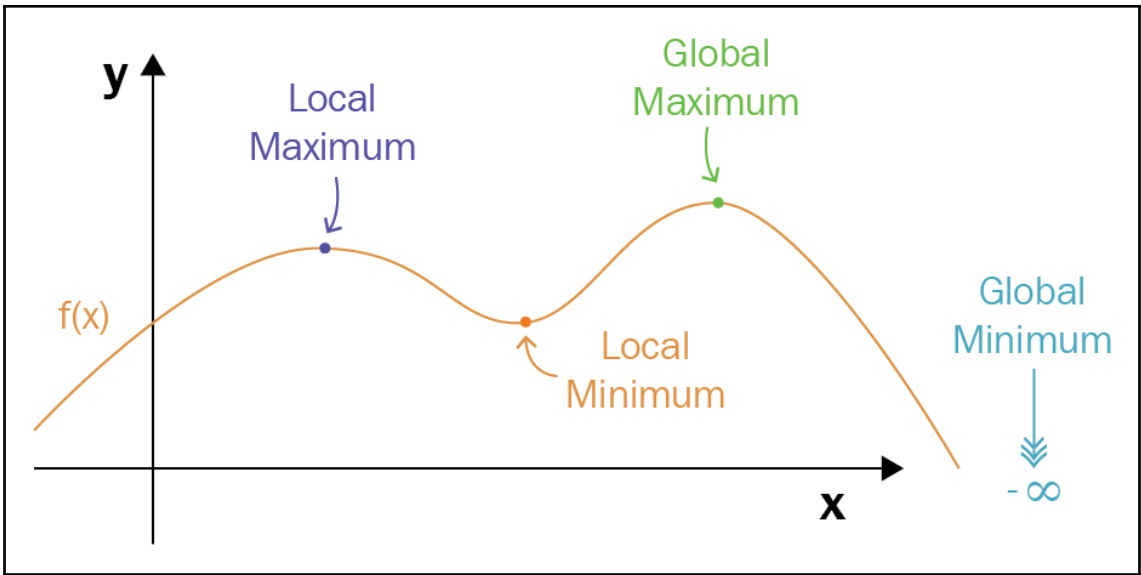


```
import gym
env = gym.make("Taxi-v2")
env.render()
```

```
+-----+
|R: | : :G|
| : : : :|
| : | : :|
|Y| : |B: |
+-----+
```

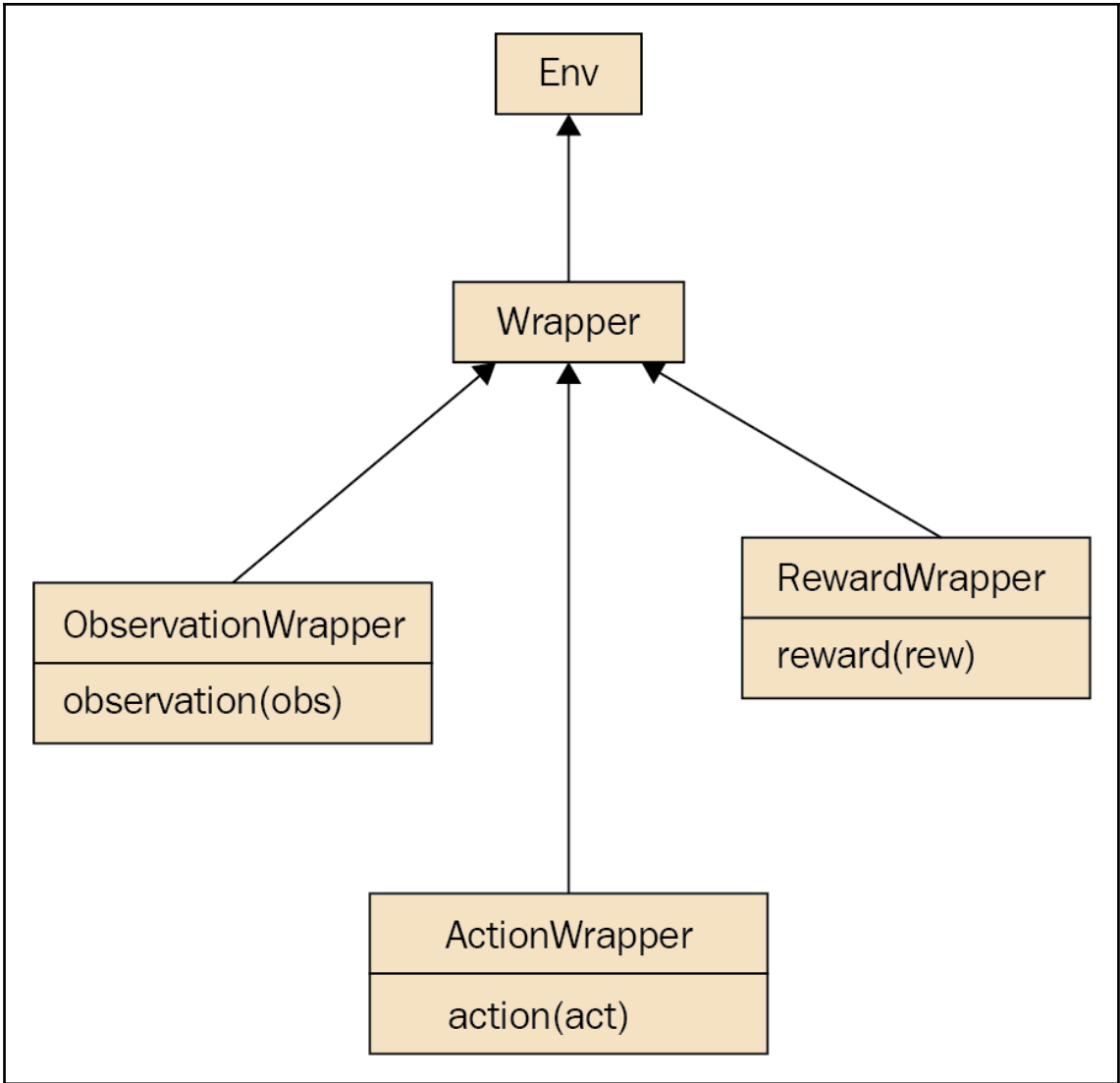
---

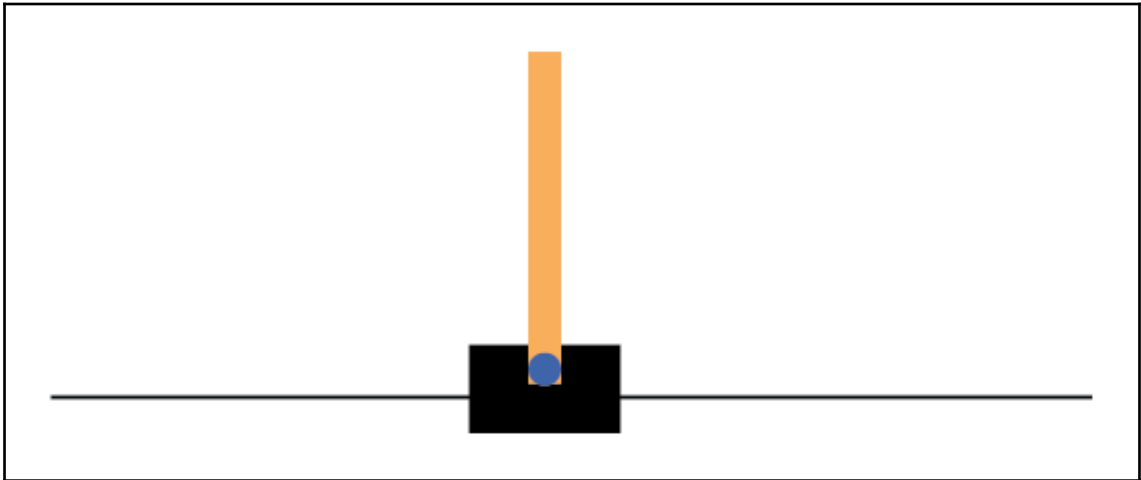
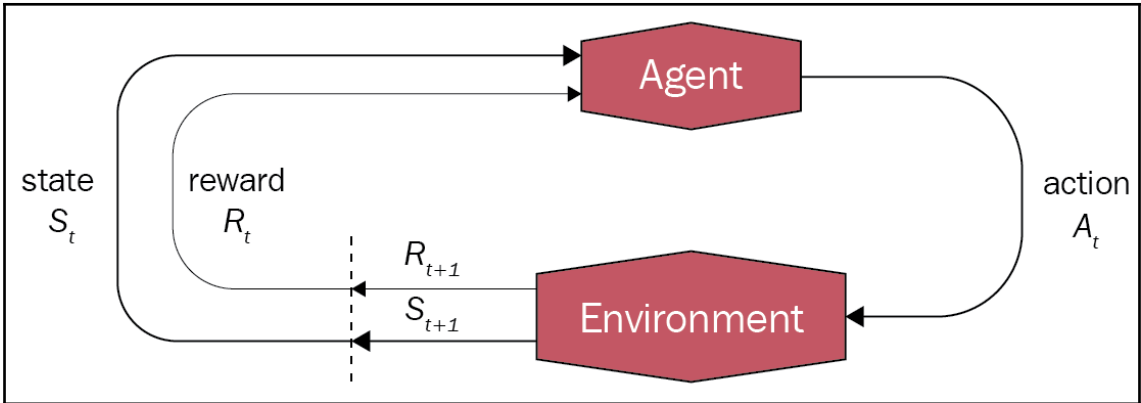
Arm	Reward
1	0
2	1
3	0
4	1
2	0
4	1
4	1
2	0
1	1

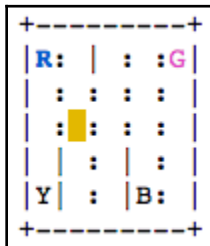
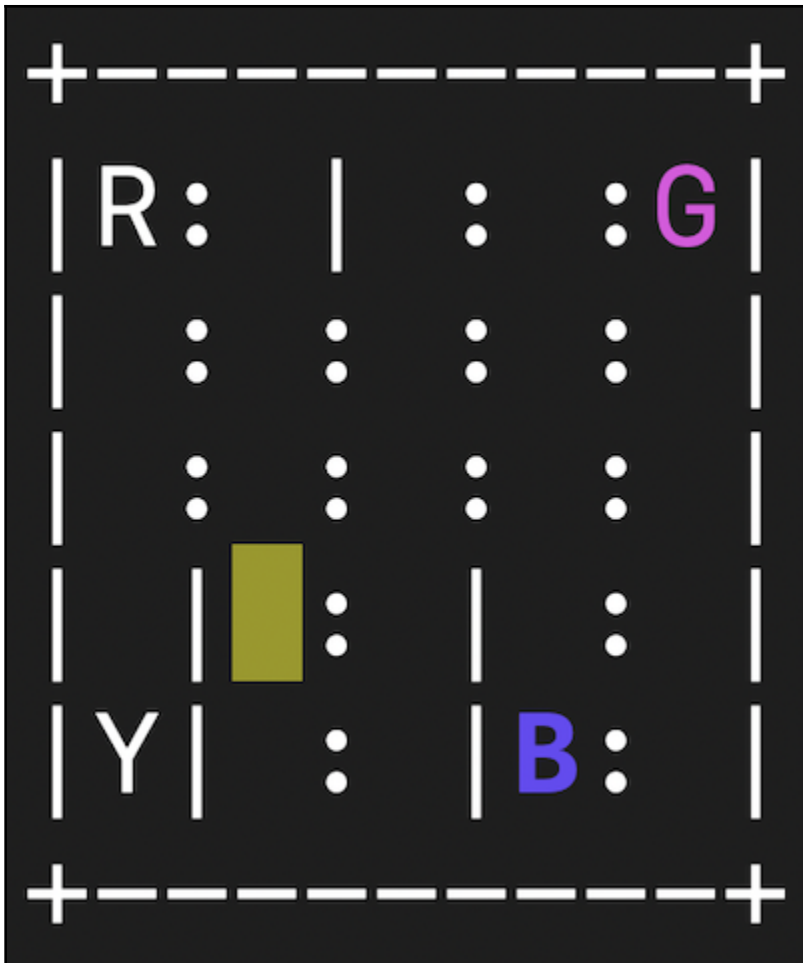


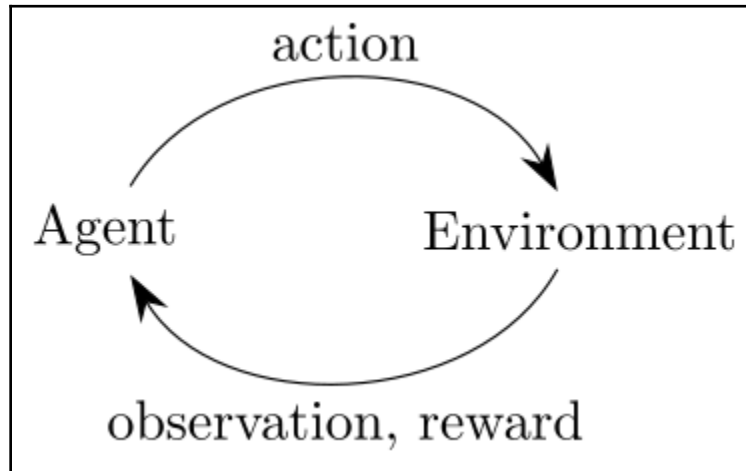
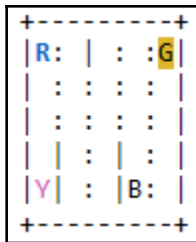
## Chapter 3, Setting Up Your First Environment with OpenAI Gym











```
print("Action Space {}".format(env.action_space))  
print("State Space {}".format(env.observation_space))
```

```
Action Space Discrete(6)  
State Space Discrete(500)
```

```
env.step(1)
(382, -1, False, {'prob': 1.0})

env.render()

+-----+
|R: | : :G|
| : : : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(North)
```

```
env.env.s = 20
env.render()

+-----+
|R: | : :G|
| : : : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(North)

env.env.s = 50
env.render()

+-----+
|R: | : :G|
| : : : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(North)
```

---

```
env.env.s = 50  
env.render()
```

```
+-----+  
|R: |█: :G|  
| : : : : |  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+  
(South)
```

```
env.step(0)  
env.render()
```

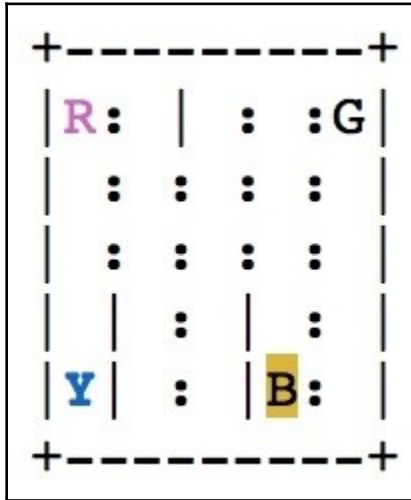
```
+-----+  
|R: | : :G|  
| : :█: : |  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+  
(South)
```

```
+-----+  
|R: | : :G|  
| : : : : |  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+  
(Dropoff)
```

---

```
+-----+
|R: | : :G|
| : : : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
      (Dropoff)
      2124
```

# Chapter 4, Teaching a Smartcab to Drive Using Q-Learning



Q-table initialised at zero					After few episodes					Eventually				
	UP	DOWN	LEFT	RIGHT		UP	DOWN	LEFT	RIGHT		UP	DOWN	LEFT	RIGHT
0	0	0	0	0	0	0	0	0	0	0	0	0	0.45	0
1	0	0	0	0	1	0	0	0	0	1	0	1.01	0	0
2	0	0	0	0	2	0	2.25	2.25	0	2	0	2.25	2.25	0
3	0	0	0	0	3	0	0	5	0	3	0	0	5	0
4	0	0	0	0	4	0	0	0	0	4	0	0	0	0
5	0	0	0	0	5	0	0	0	0	5	0	0	0	0
6	0	0	0	0	6	0	5	0	0	6	0	5	0	0
7	0	0	0	0	7	0	0	2.25	0	7	0	0	2.25	0
8	0	0	0	0	8	0	0	0	0	8	0	0	0	0



---

```
+-----+
|R: | : :G|
| : : : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(Dropoff)
```

```
+-----+
|R: | : :G|
| : : : : |
| : : : : |
| | : | : |
|Y| : |G: |
+-----+
(Dropoff)
1998
Counter: 1998
```

```
+-----+
|R: | : :G|
| : : : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(Dropoff)
524
Counter: 524
```

---

```
total_epochs = 0
episodes = 100

for episode in range(0, episodes):
    epochs = 0
    reward = 0
    state = env.reset()
    while reward != 20:
        action = env.action_space.sample()
        state, reward, done, info = env.step(action)
        epochs += 1
    total_epochs += epochs
print("Average timesteps taken: {}".format(total_epochs/episodes))
```

Average timesteps taken: 2464.7

```

Q = np.zeros([env.observation_space.n, env.action_space.n])

gamma = 0.1
alpha = 0.1
epsilon = 0.1
total_epochs = 0
episodes = 100

for episode in range(episodes):
    epochs = 0
    reward = 0
    state = env.reset()

    while reward != 20:
        if np.random.rand() < epsilon:
            action = env.action_space.sample()
        else:
            action = np.argmax(Q[state])
        next_state, reward, done, info = env.step(action)
        Q[state, action] = Q[state, action] + alpha * (reward + gamma * \
            np.max(Q[next_state]) - Q[state, action])

        state = next_state
        epochs += 1
        total_epochs += epochs

print("Average timesteps taken: {}".format(total_epochs/episodes))

```

Average timesteps taken: 663.97

---

```
total_epochs = 0
episodes = 10000

for episode in range(episodes):
    epochs = 0
    reward = 0
    state = env.reset()
    while reward != 20:
        action = env.action_space.sample()
        state, reward, done, info = env.step(action)
        epochs += 1
    total_epochs += epochs
print("Average timesteps taken: {}".format(total_epochs/episodes))
```

Average timesteps taken: 2255.7757

```
Q = np.zeros([env.observation_space.n, env.action_space.n])

gamma = 0.1
alpha = 0.1
epsilon = 0.1
total_epochs = 0
episodes = 10000

for episode in range(episodes):
    epochs = 0
    reward = 0
    state = env.reset()

    while reward != 20:
        if np.random.rand() < epsilon:
            action = env.action_space.sample()
        else:
            action = np.argmax(Q[state])
        next_state, reward, done, info = env.step(action)
        Q[state, action] = Q[state, action] + alpha * (reward + gamma * \
            np.max(Q[next_state]) - Q[state, action])

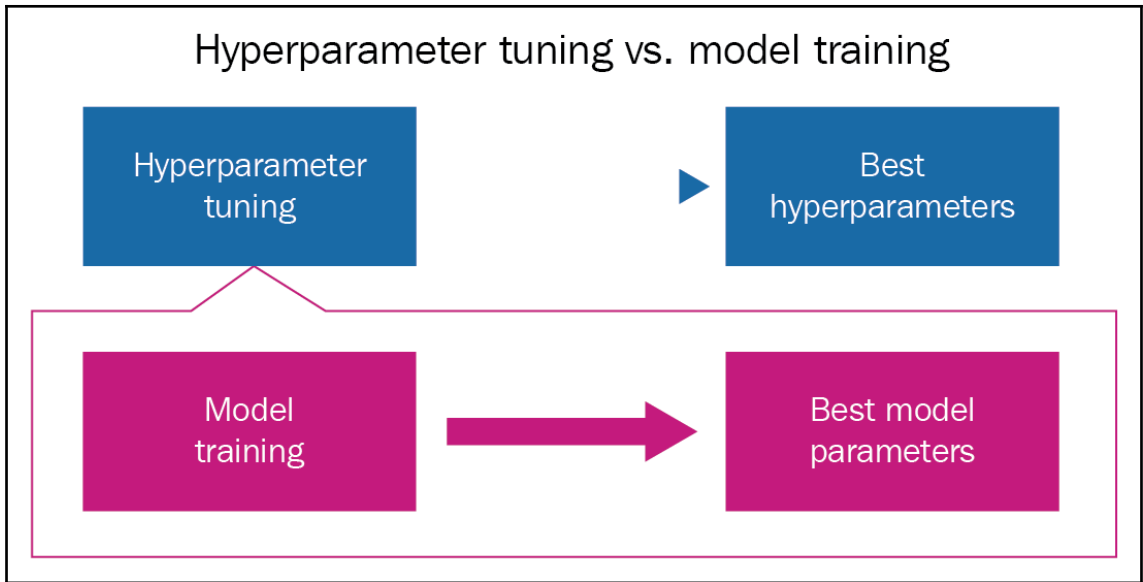
        state = next_state
        epochs += 1
    total_epochs += epochs

print("Average timesteps taken: {}".format(total_epochs/episodes))
```

Average timesteps taken: 46.0174

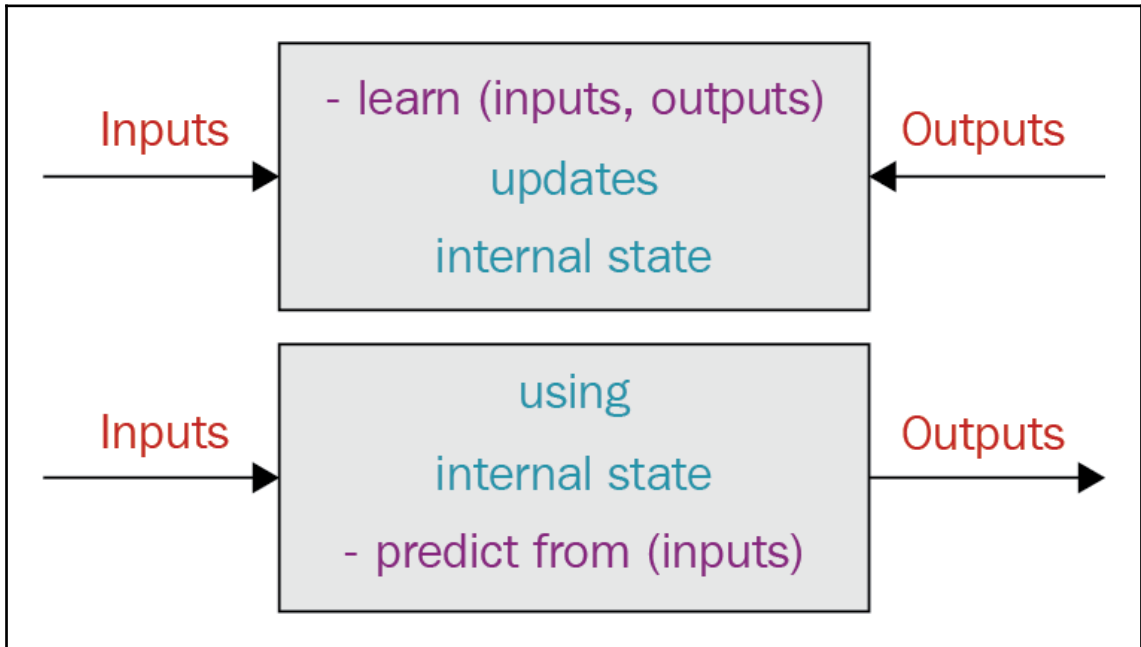
---

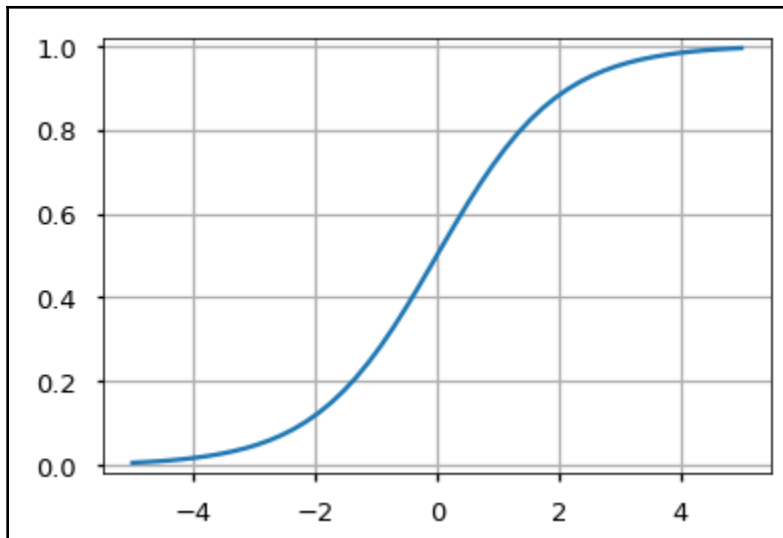
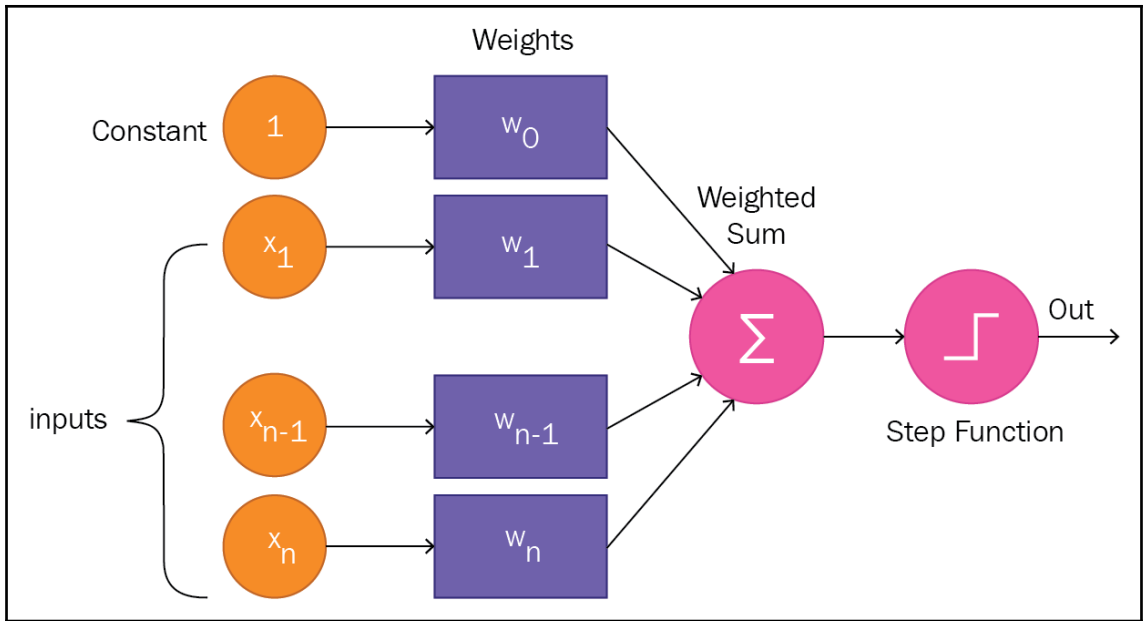
## Hyperparameter tuning vs. model training



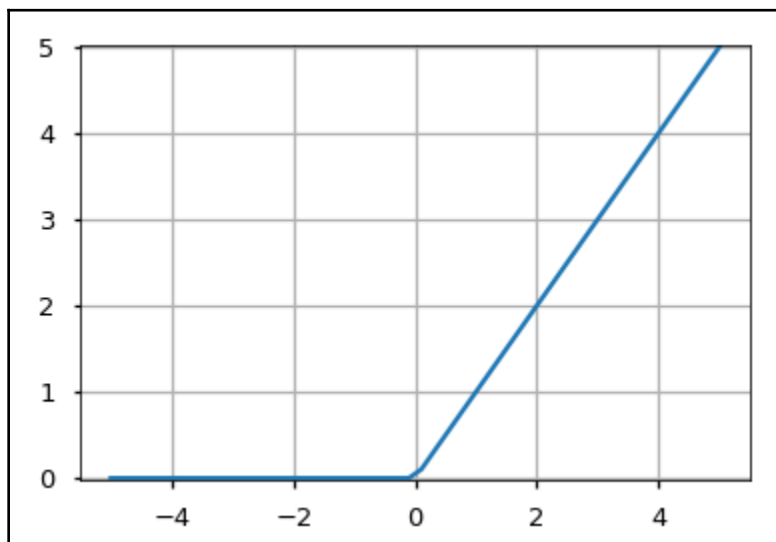
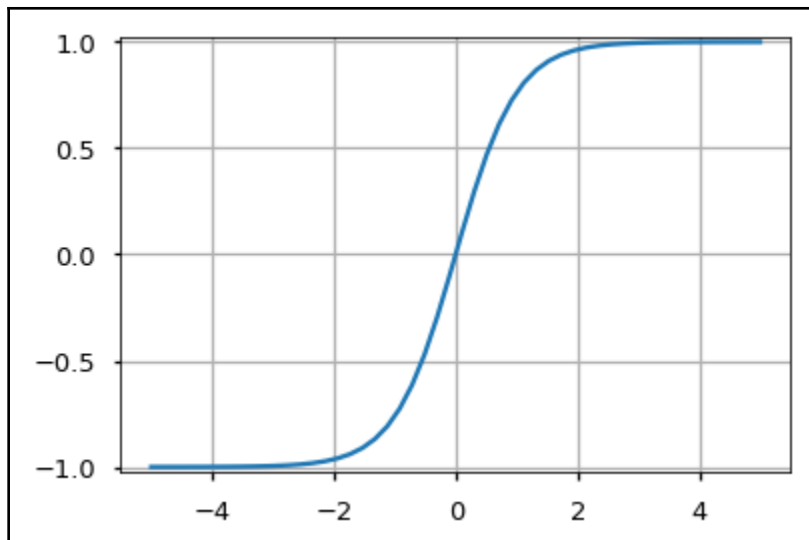
---

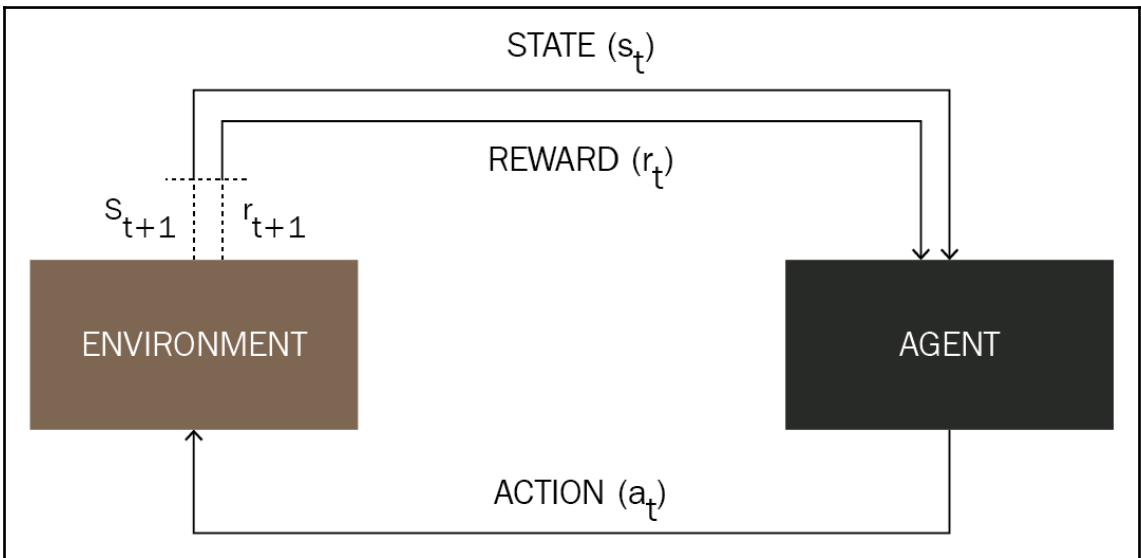
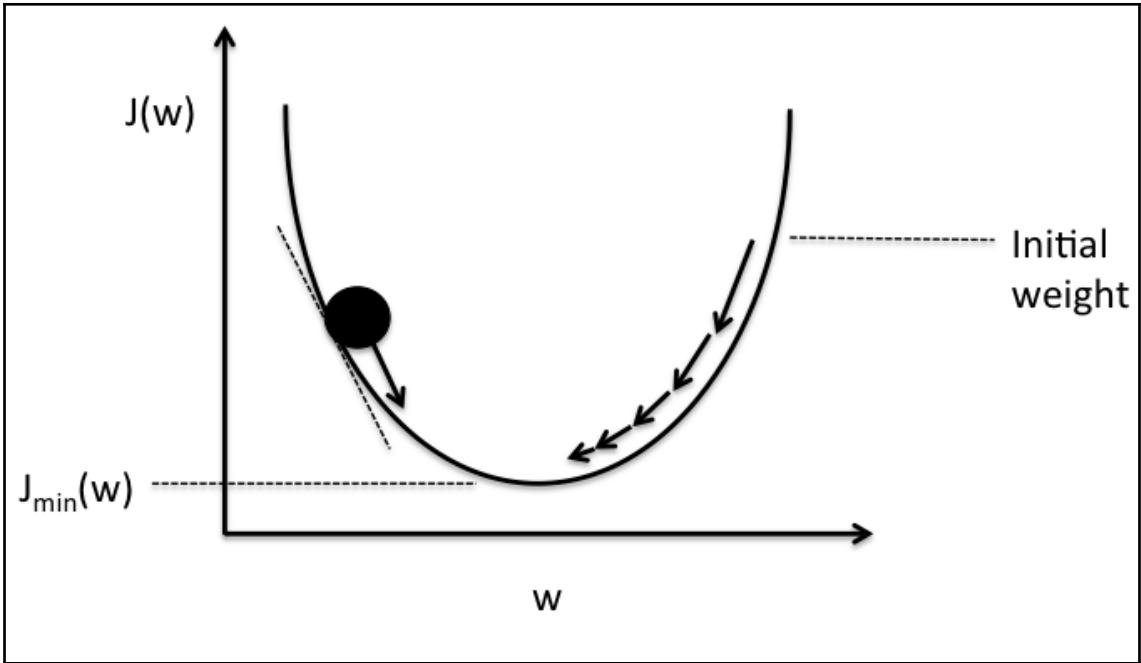
## Chapter 5, Building Q-Networks with TensorFlow





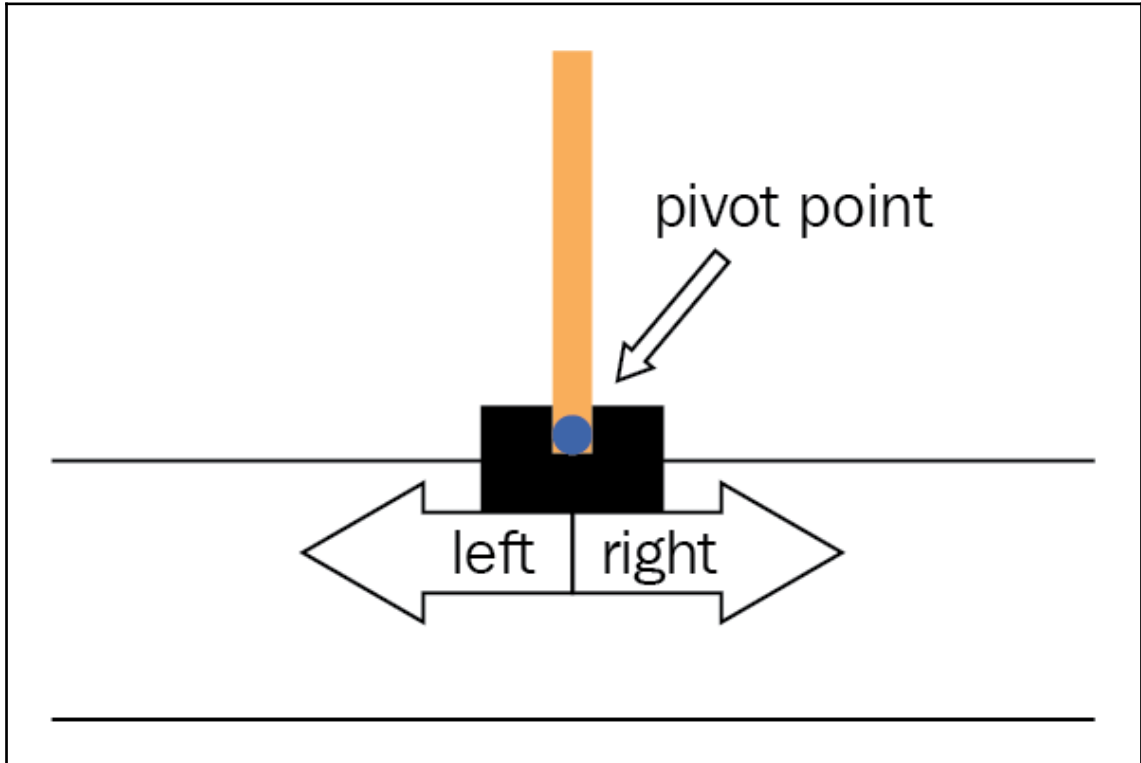


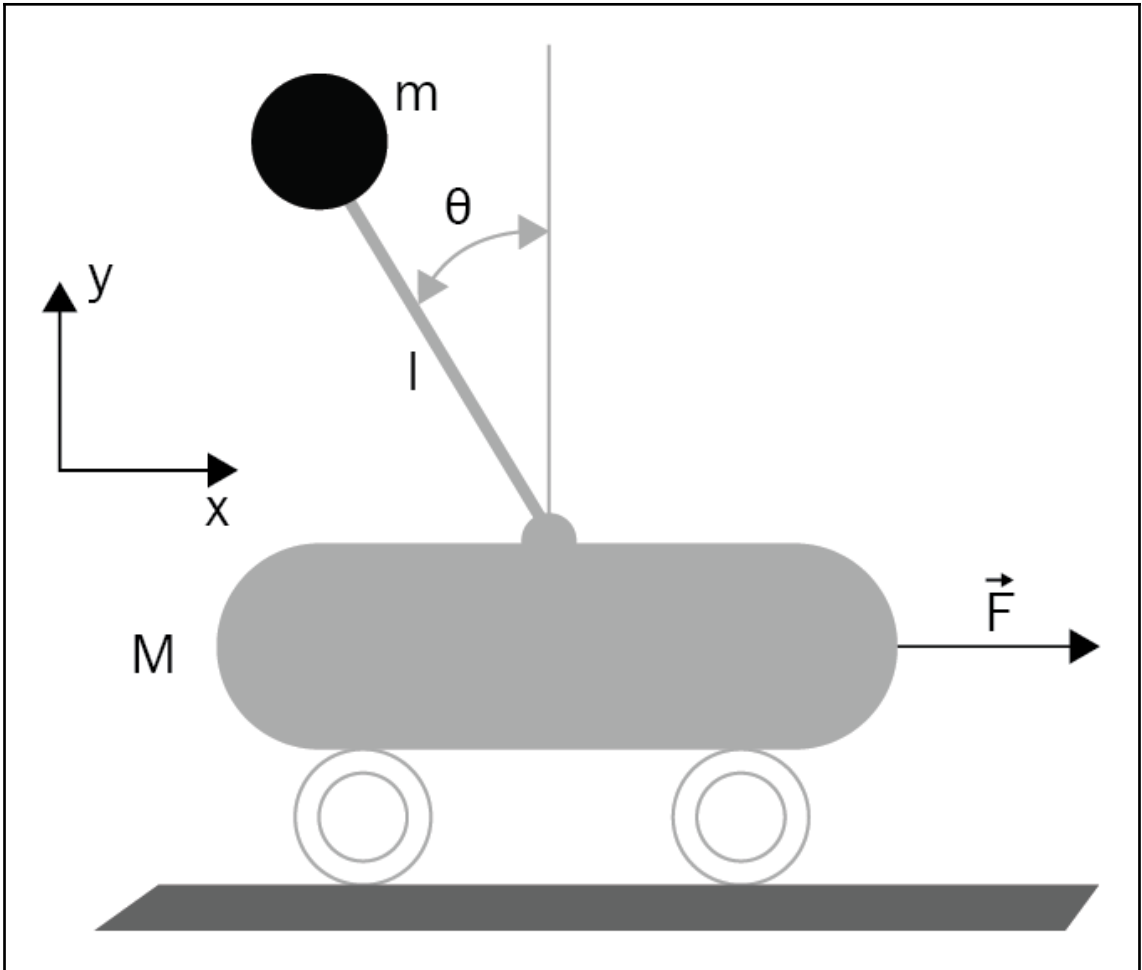


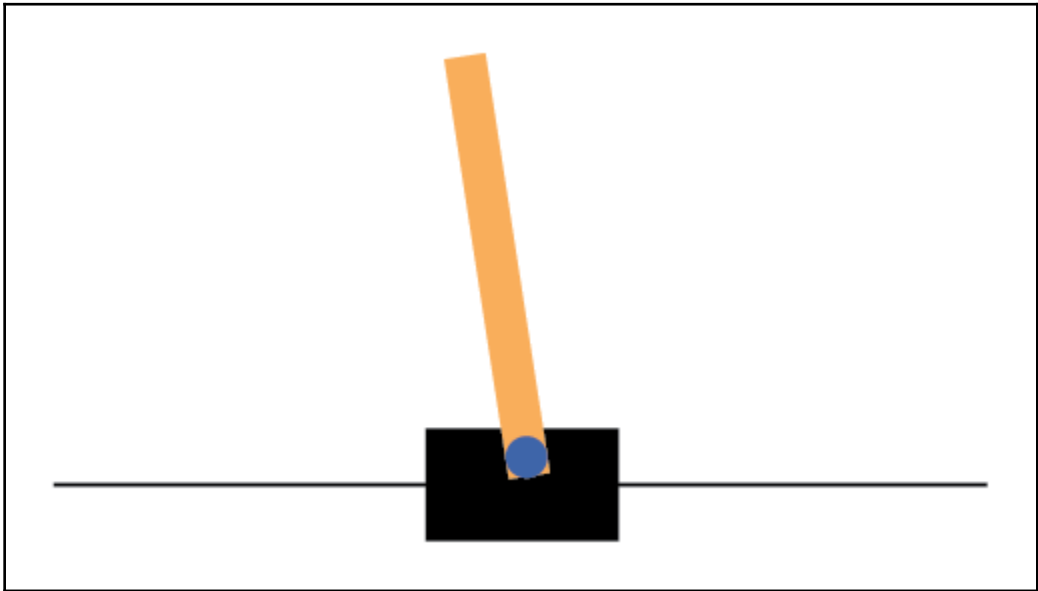
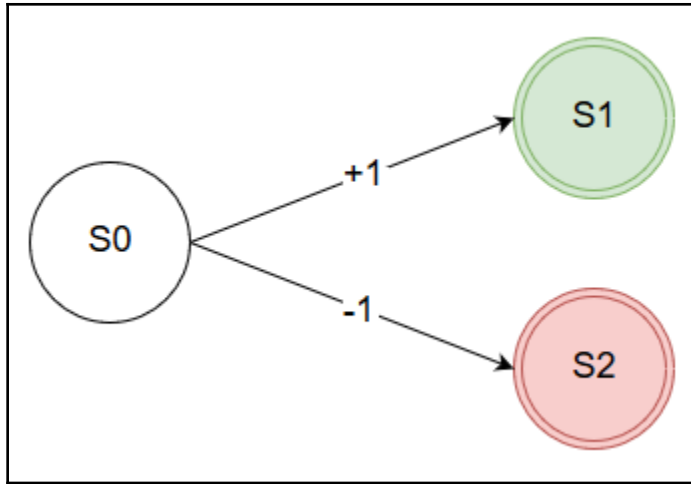


---

## Chapter 6, Digging Deeper into Deep Q-Networks with Keras and TensorFlow







---

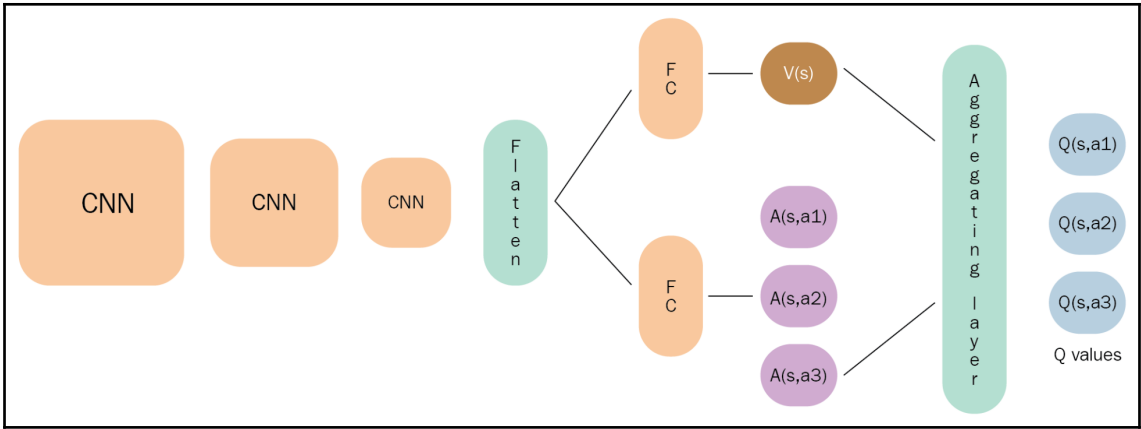
```
if __name__ == "__main__":  
    cartpole()
```

```
Epoch: 450 Score: 12  
Epoch: 451 Score: 24  
Epoch: 452 Score: 19  
Epoch: 453 Score: 15  
Epoch: 454 Score: 16  
Epoch: 455 Score: 19  
Epoch: 456 Score: 38  
Epoch: 457 Score: 18  
Epoch: 458 Score: 22  
Epoch: 459 Score: 13  
Epoch: 460 Score: 19  
Epoch: 461 Score: 16  
Epoch: 462 Score: 28  
Epoch: 463 Score: 10  
Epoch: 464 Score: 12  
Epoch: 465 Score: 21  
Epoch: 466 Score: 14  
Epoch: 467 Score: 26  
Epoch: 468 Score: 21  
Epoch: 469 Score: 22
```

---

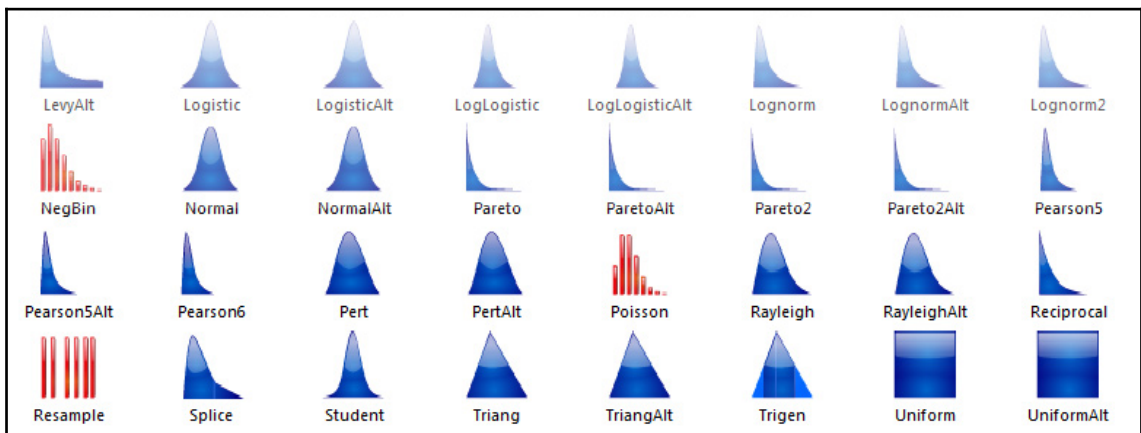
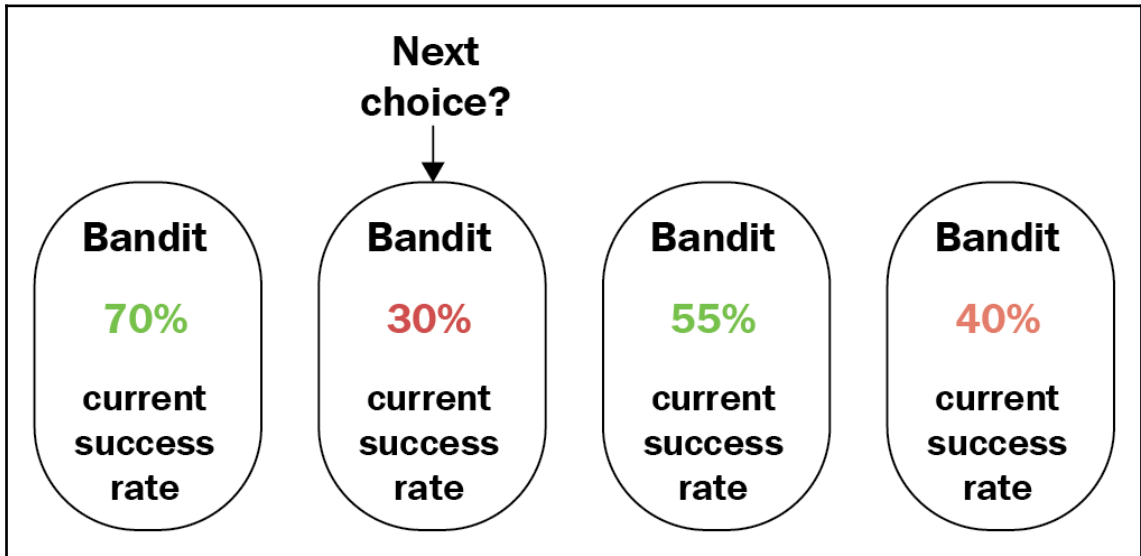
```
if __name__ == "__main__":  
    cartpole()
```

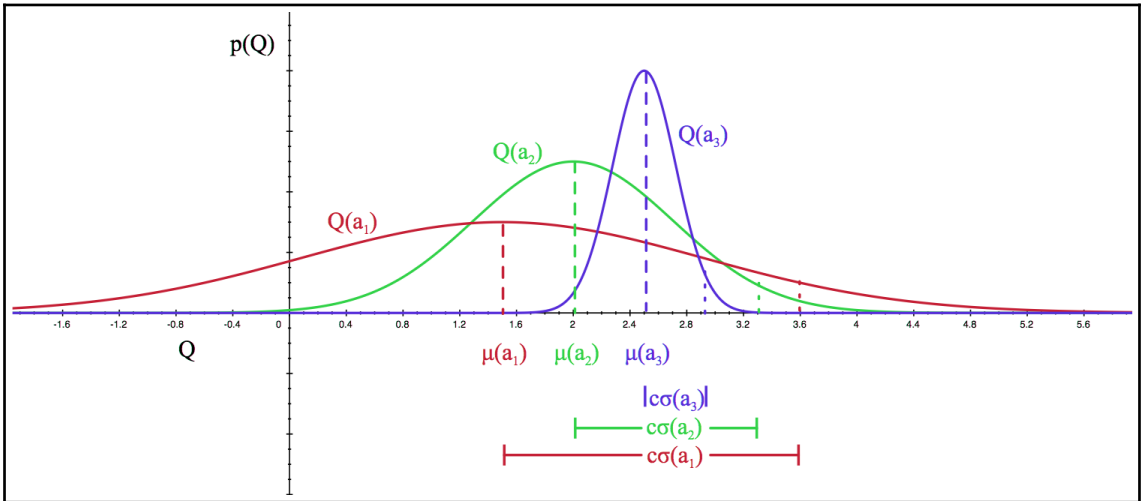
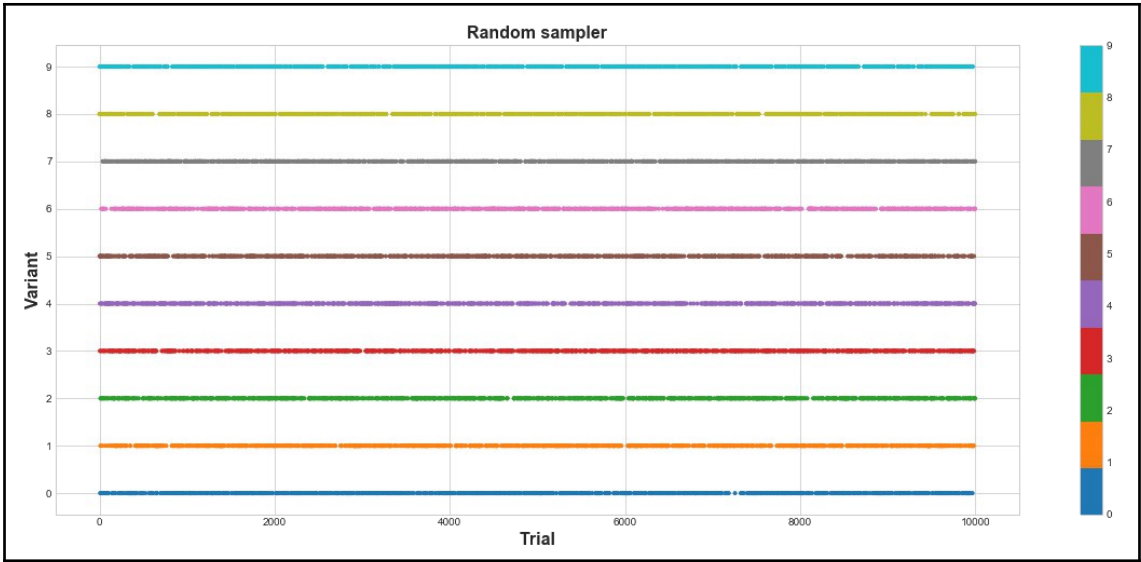
```
Epoch: 1 Score: 30  
Epoch: 2 Score: 13  
Epoch: 3 Score: 11  
Epoch: 4 Score: 10  
Epoch: 5 Score: 16  
Epoch: 6 Score: 12  
Epoch: 7 Score: 13  
Epoch: 8 Score: 23  
Epoch: 9 Score: 27  
Epoch: 10 Score: 9  
Epoch: 11 Score: 41  
Epoch: 12 Score: 8  
Epoch: 13 Score: 12  
Epoch: 14 Score: 13  
Epoch: 15 Score: 11  
Epoch: 16 Score: 24  
Epoch: 17 Score: 47  
Epoch: 18 Score: 62  
Epoch: 19 Score: 62  
Epoch: 20 Score: 100  
Epoch: 21 Score: 79  
Epoch: 22 Score: 124  
Epoch: 23 Score: 110  
Epoch: 24 Score: 129  
Epoch: 25 Score: 126  
Epoch: 26 Score: 176  
Epoch: 27 Score: 211  
Epoch: 28 Score: 201
```

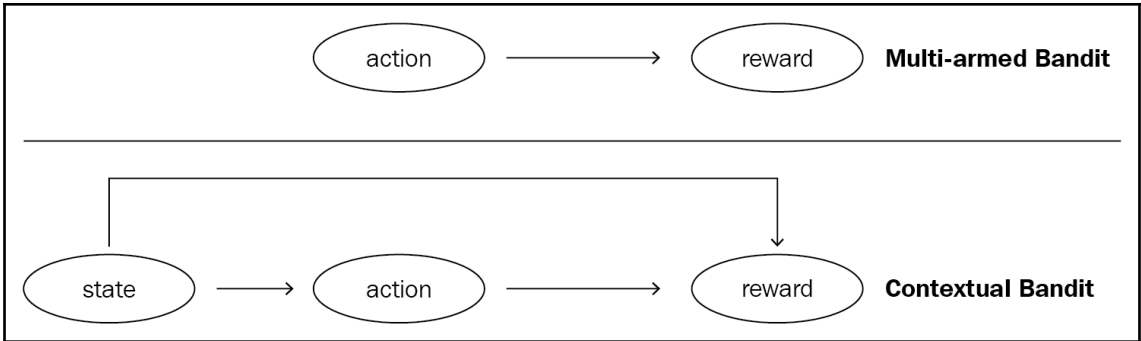


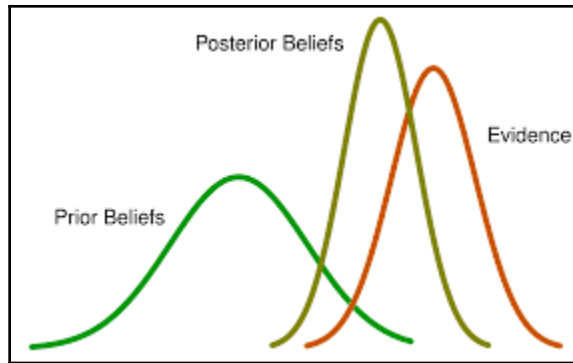


# Chapter 7, Decoupling Exploration and Exploitation in Multi-Armed Bandits









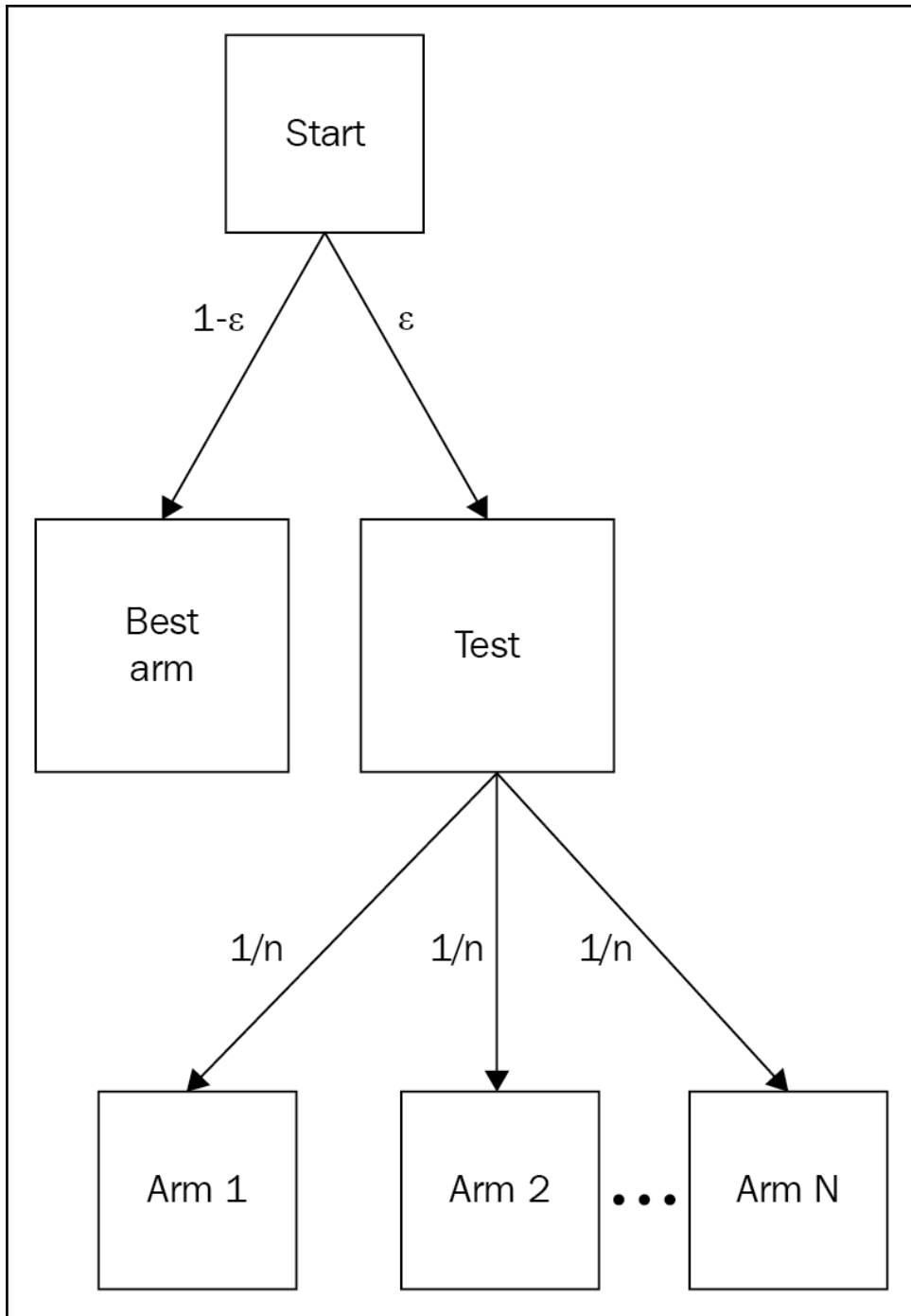
```
In [28]: df
Out[28]:
```

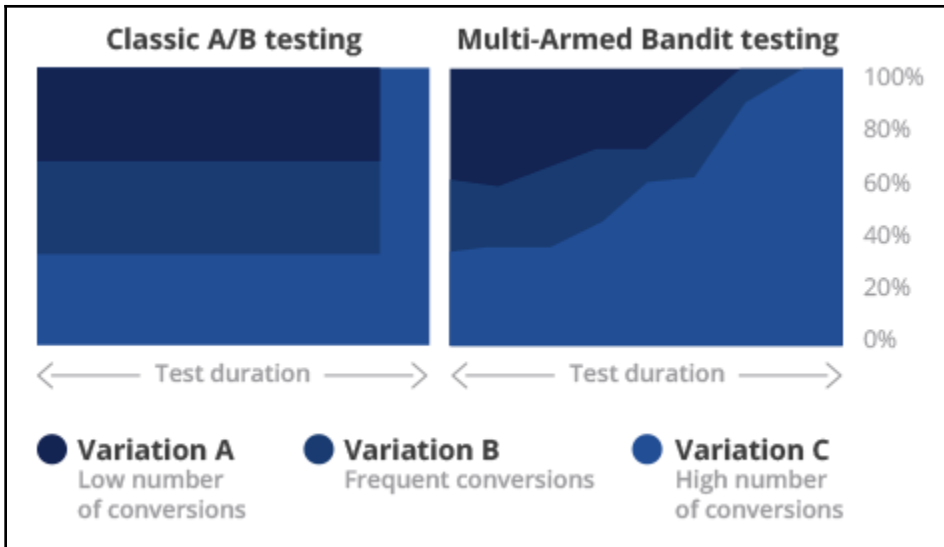
	A	B	C	D	E
0	0	0	0	0	0
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0
5	0	1	0	0	1
6	0	0	0	0	0
7	0	0	0	0	0
8	0	1	0	0	0
9	0	0	0	0	0
10	1	0	0	0	0
11	0	0	0	0	0
12	0	0	0	0	0

---

```
In [33]: pd.Series(ad_list).value_counts(normalize=True)
```

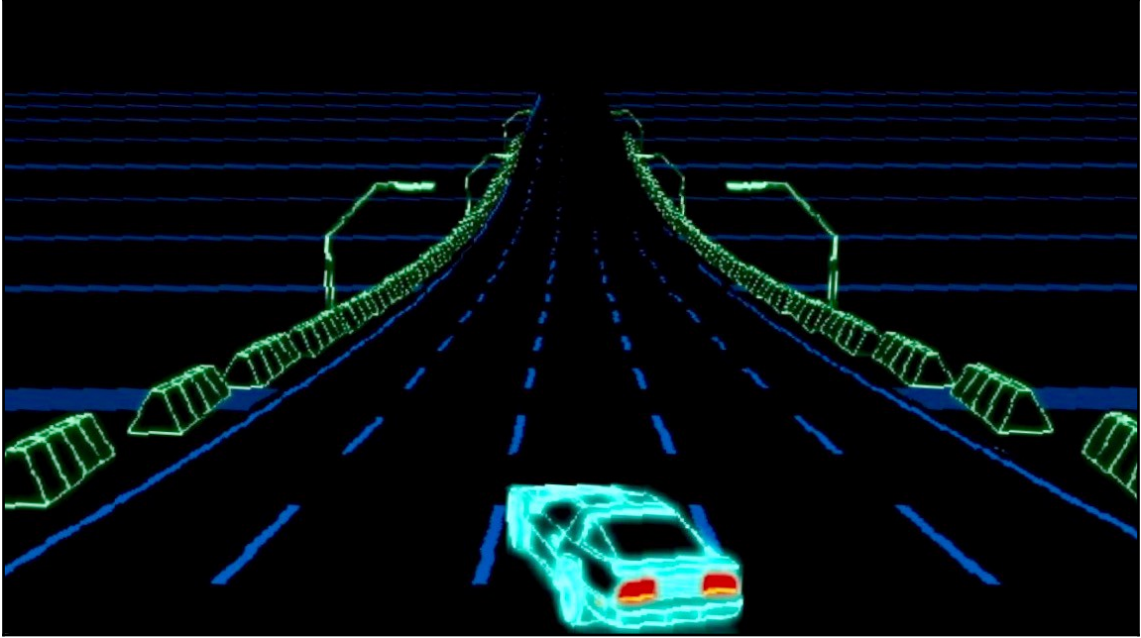
```
Out[33]: 0    0.2057  
         2    0.2006  
         3    0.2000  
         4    0.1991  
         1    0.1946  
         dtype: float64
```



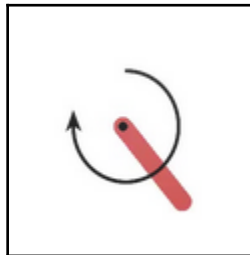
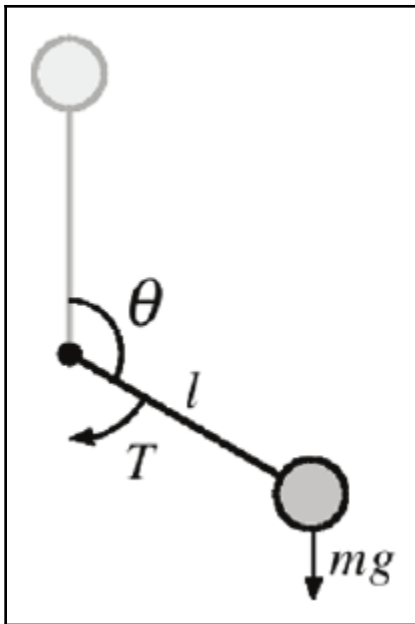


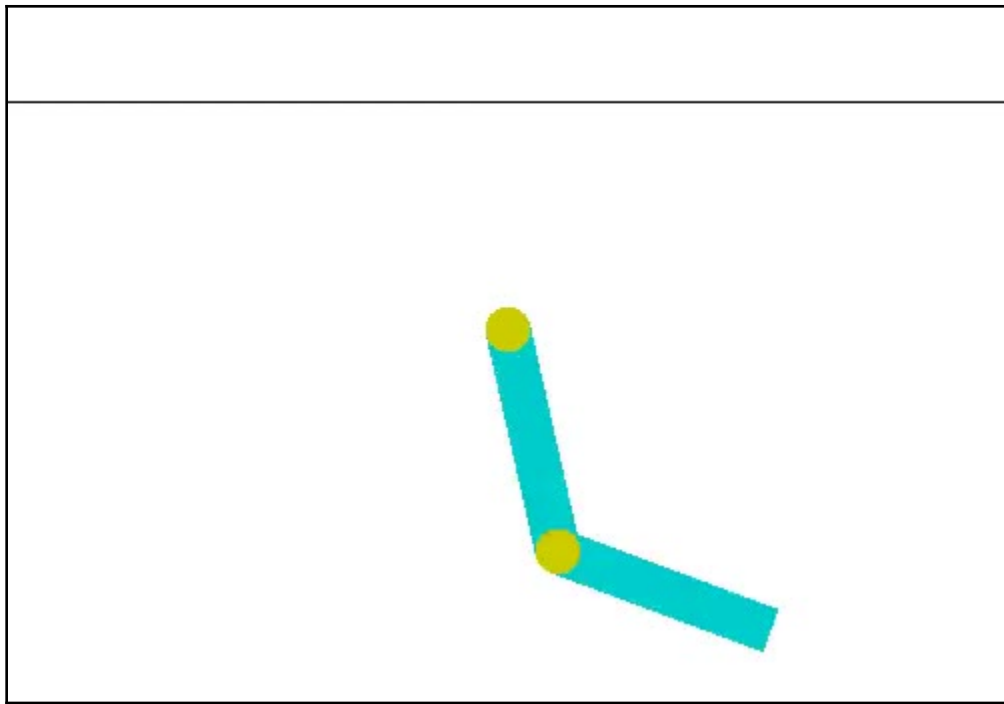
---

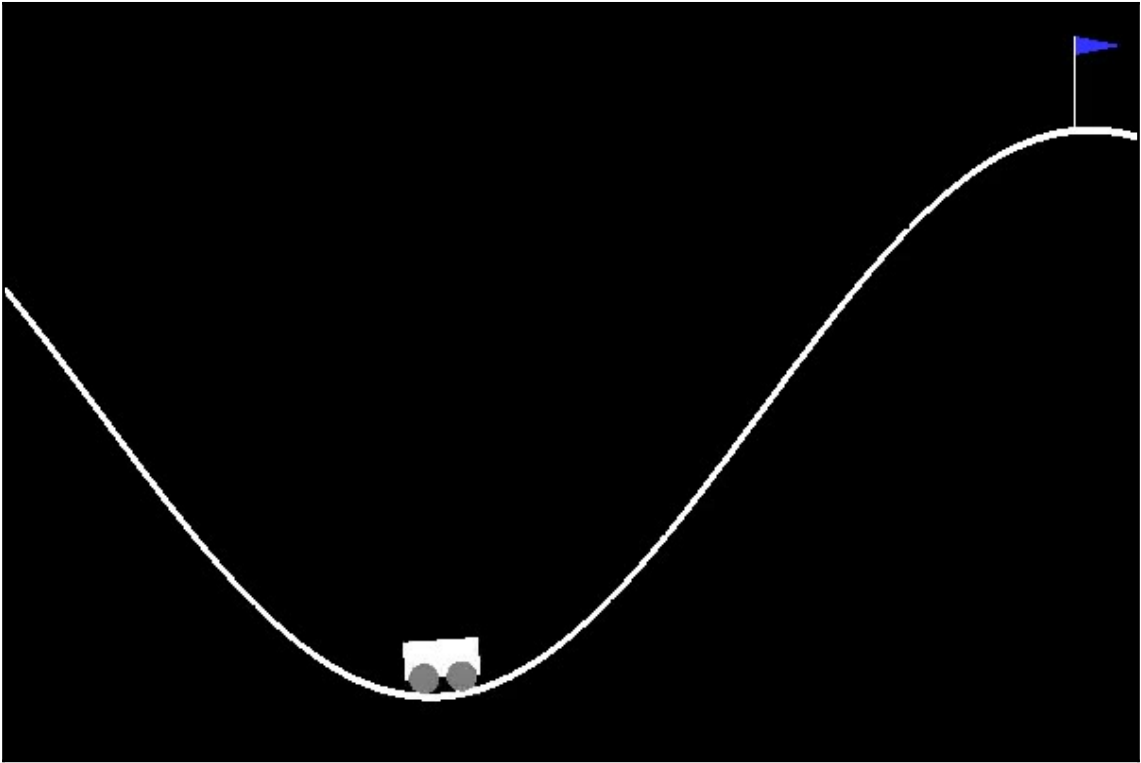
## Chapter 8, Further Q-Learning Research and Future Projects

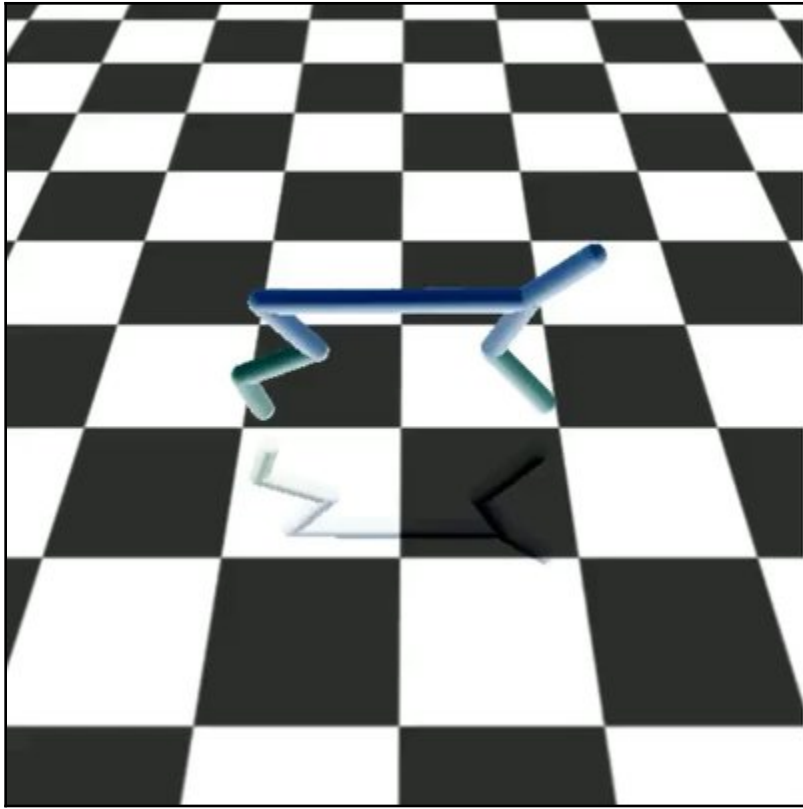


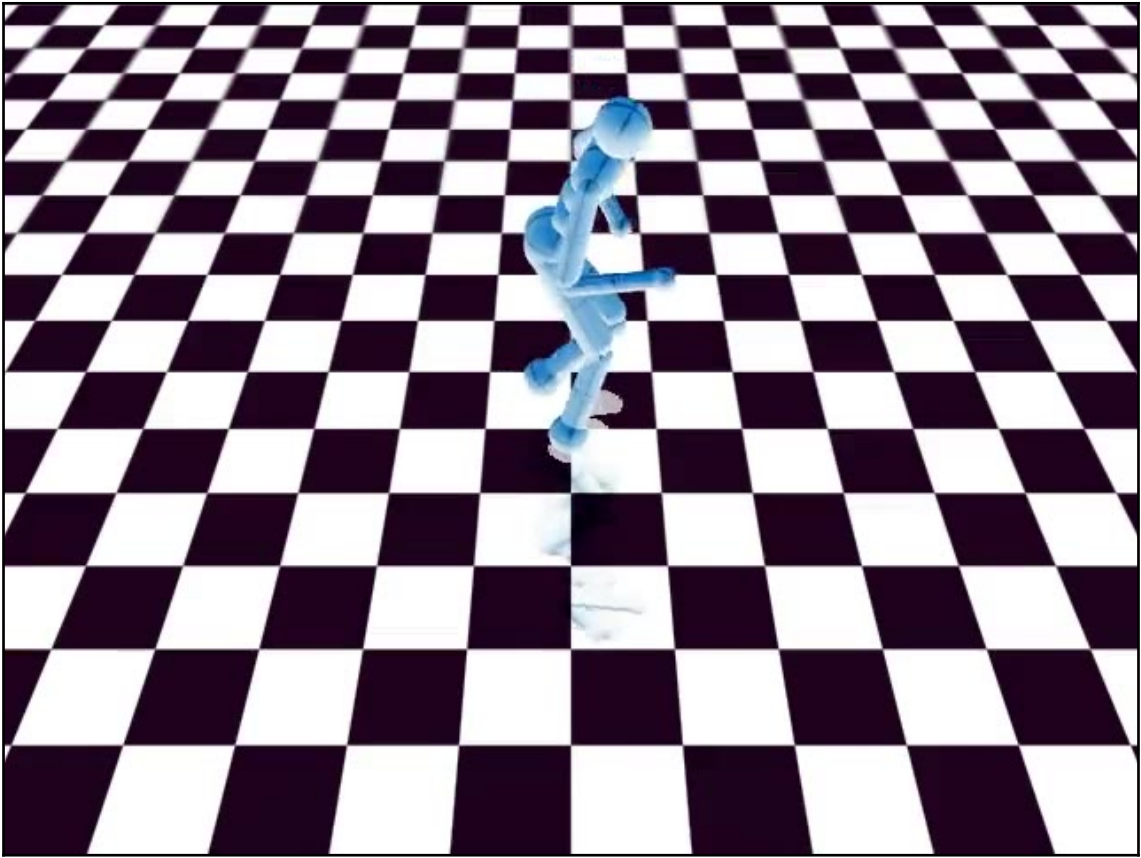


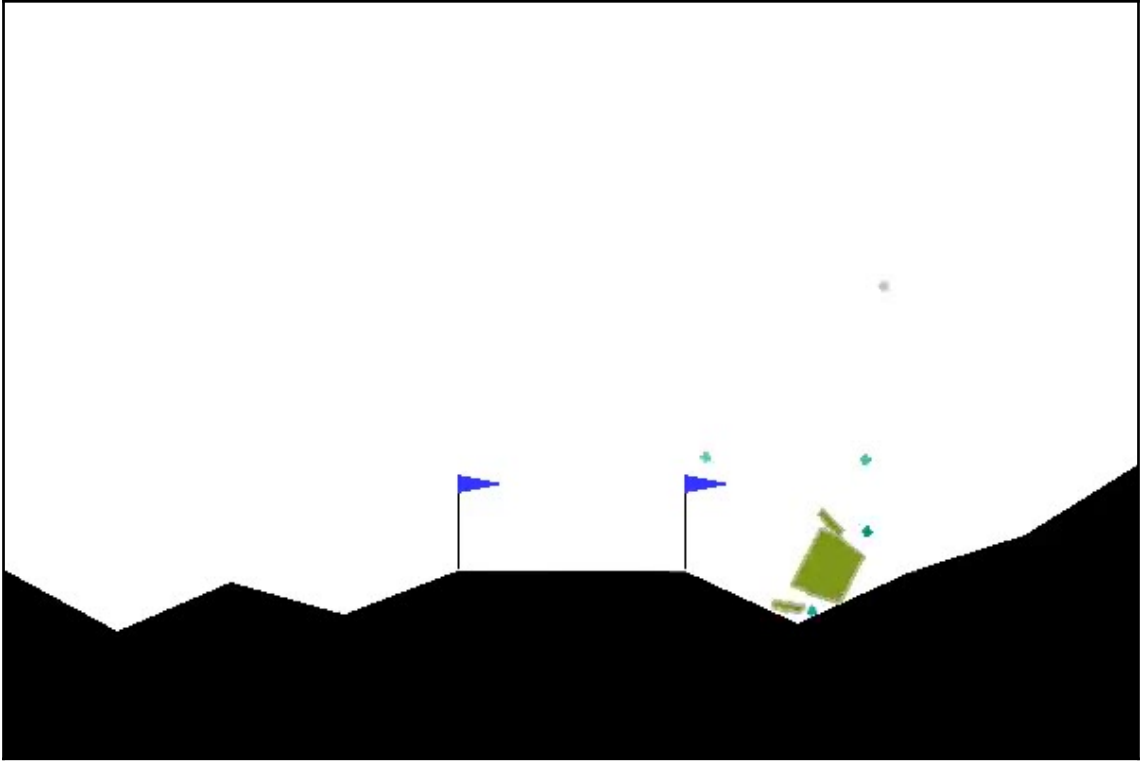


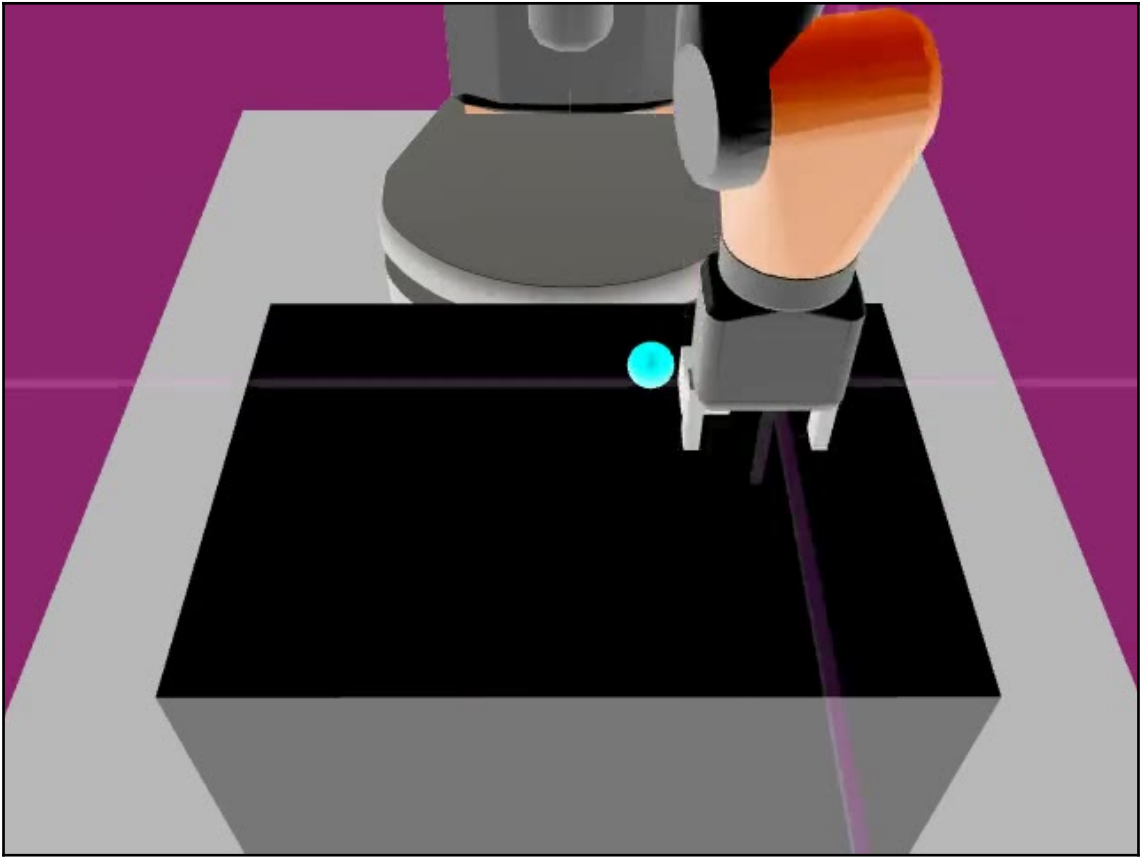














Total length of input instance: 4, step: 3

```
=====
Observation Tape   :  DADC
Output Tape       :  DAD
Targets           :  DADC

Current reward     :  1.000
Cumulative reward  :  3.000
Action            :  Tuple(move over input: right,
                          write to the output tape: True
                          prediction: D)
```



00:43





