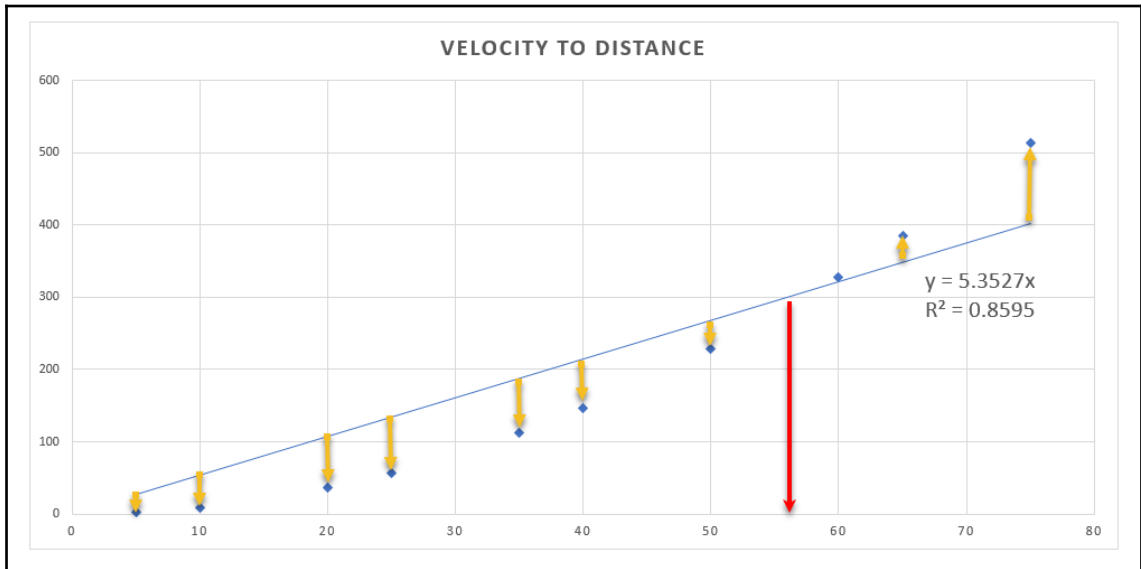
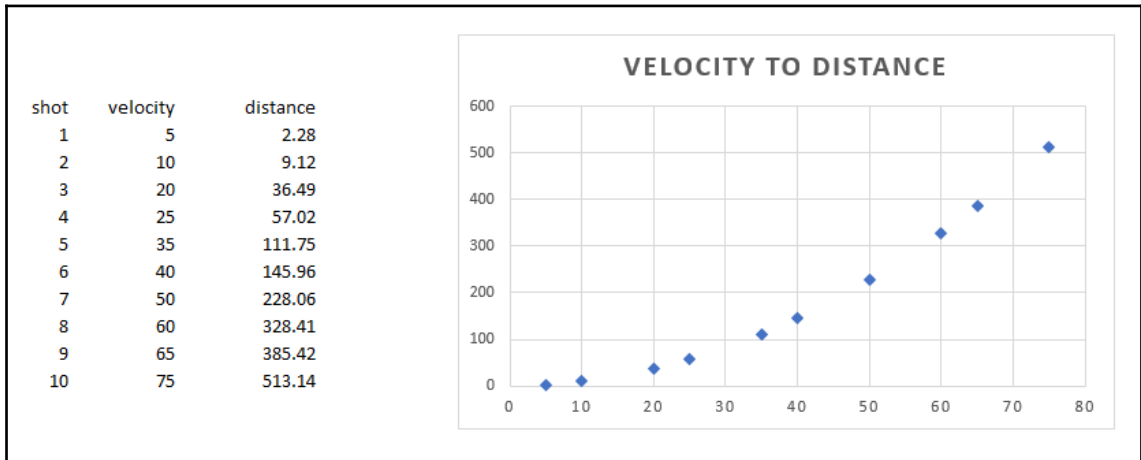


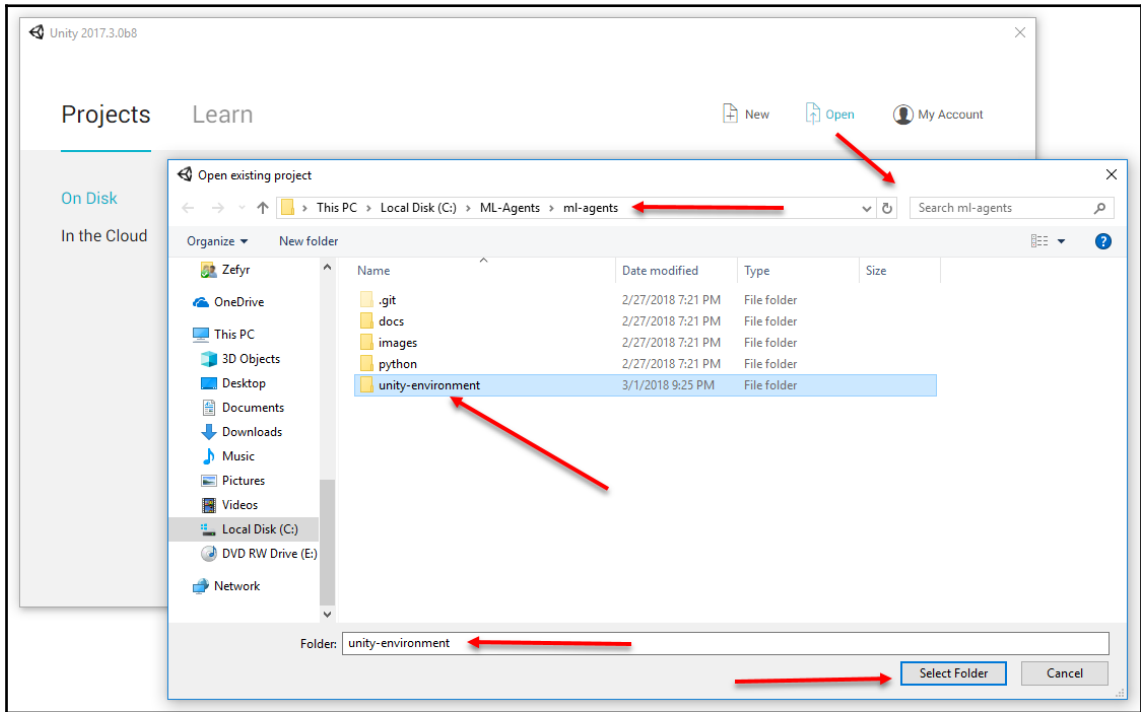
Chapter 1: Introducing Machine Learning and ML- Agents

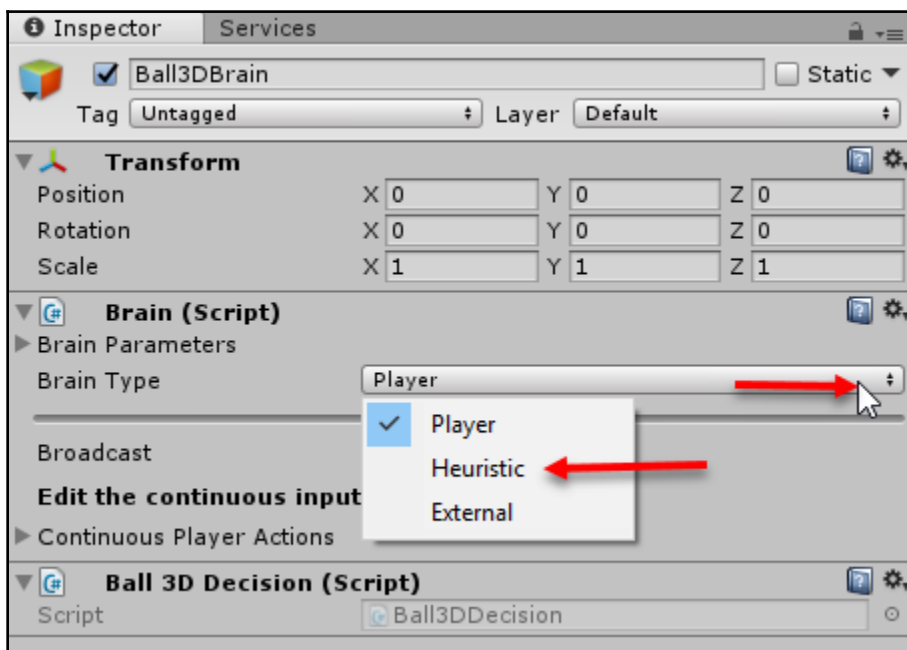
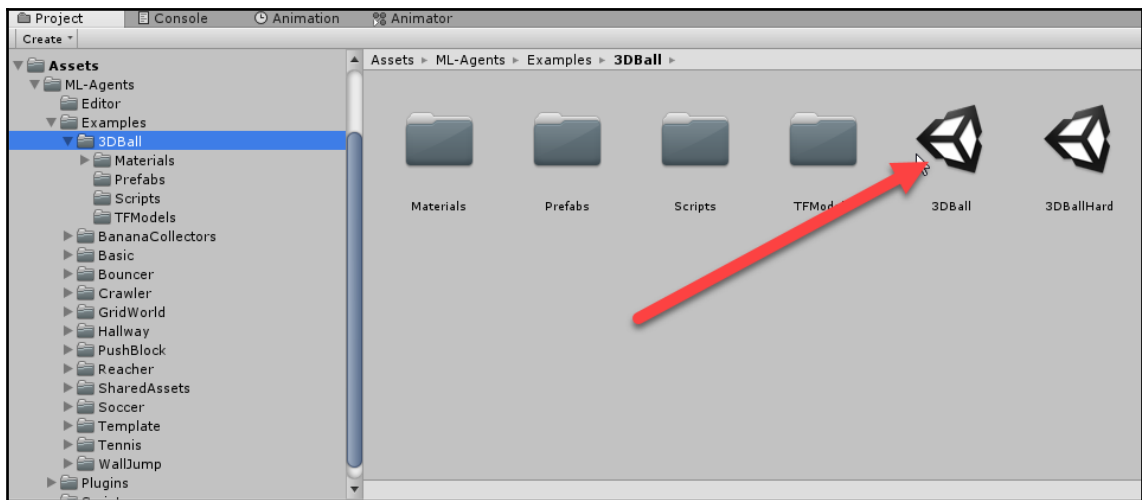


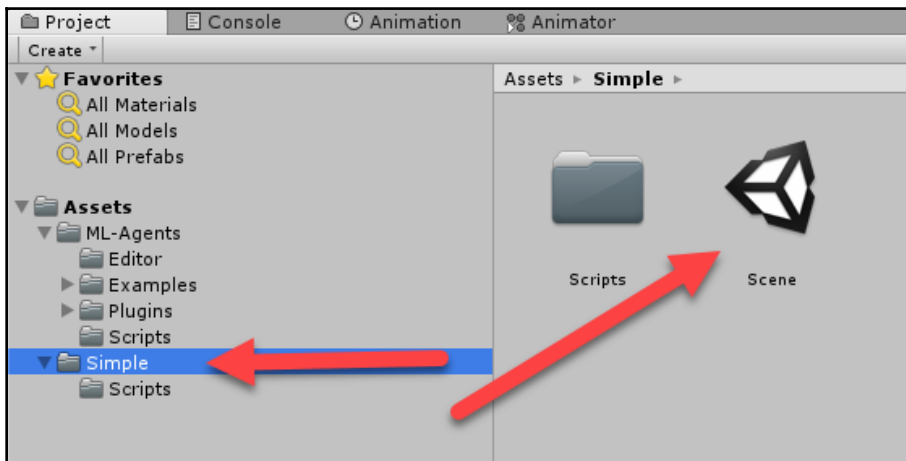
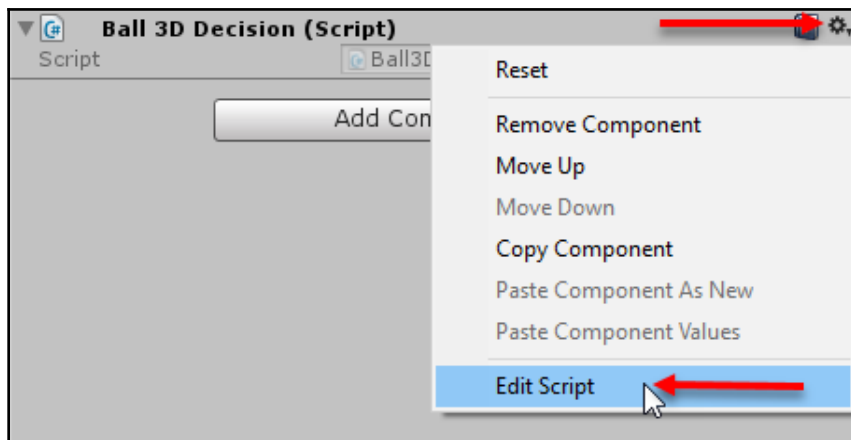
$$d = 5.3527v$$

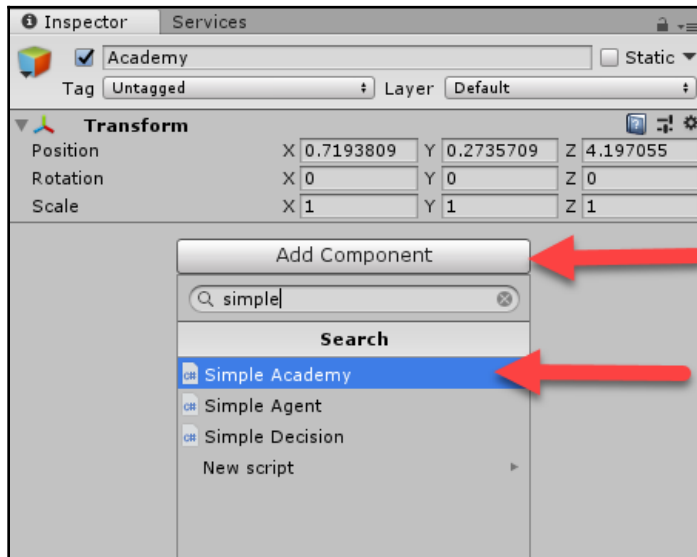
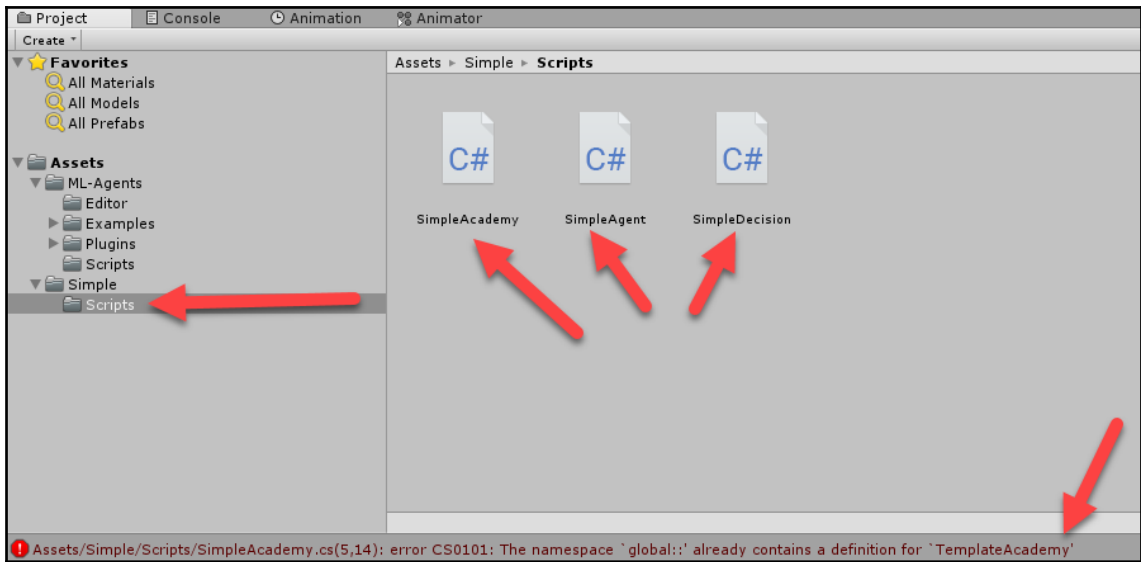
$$v = d/5.3527$$

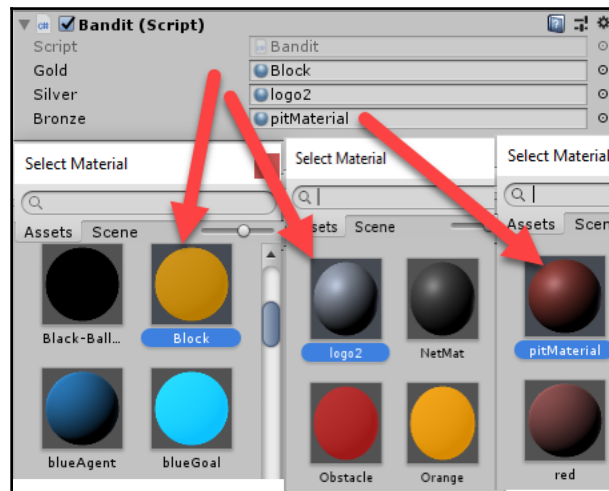
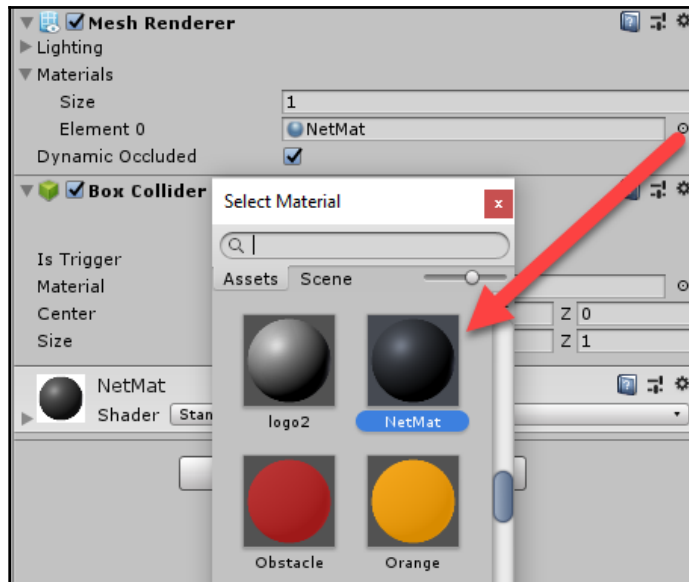
$$v = 300/5.3527 = 56.05$$

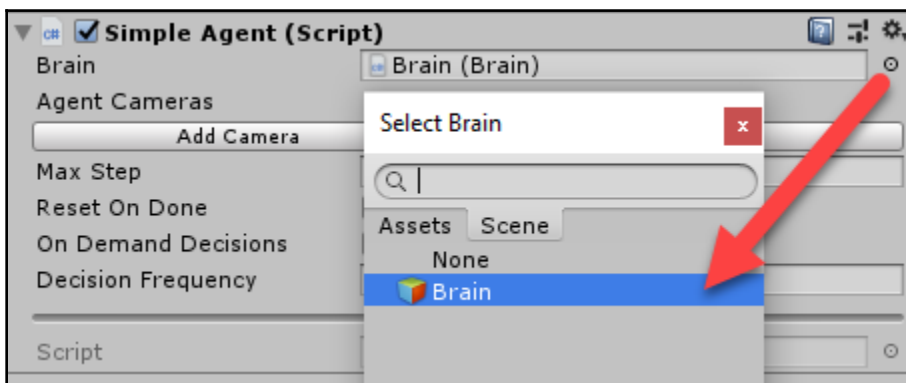
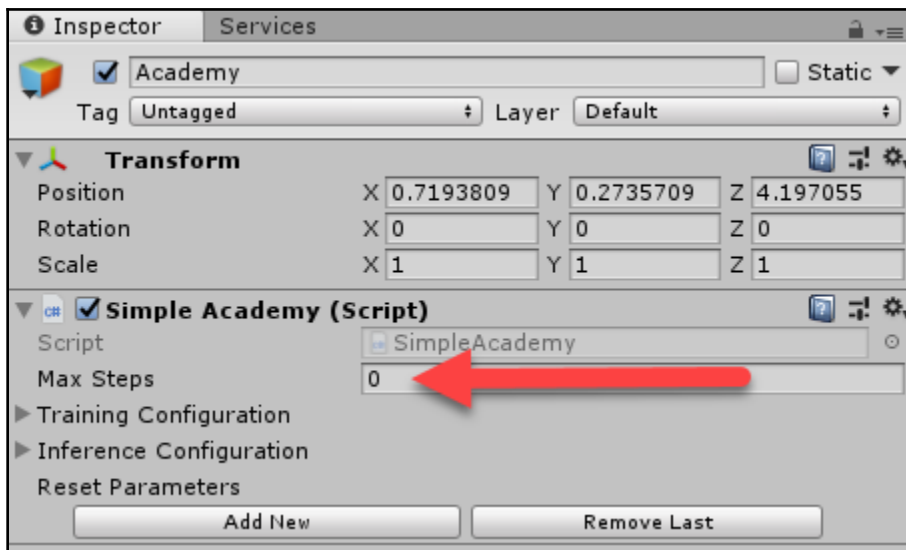












Simple Agent (Script)

Brain: Brain (Brain)

Agent Cameras: Add Camera Remove Camera

Max Step: 0

Reset On Done: ☒

On Demand Decisions: ☐

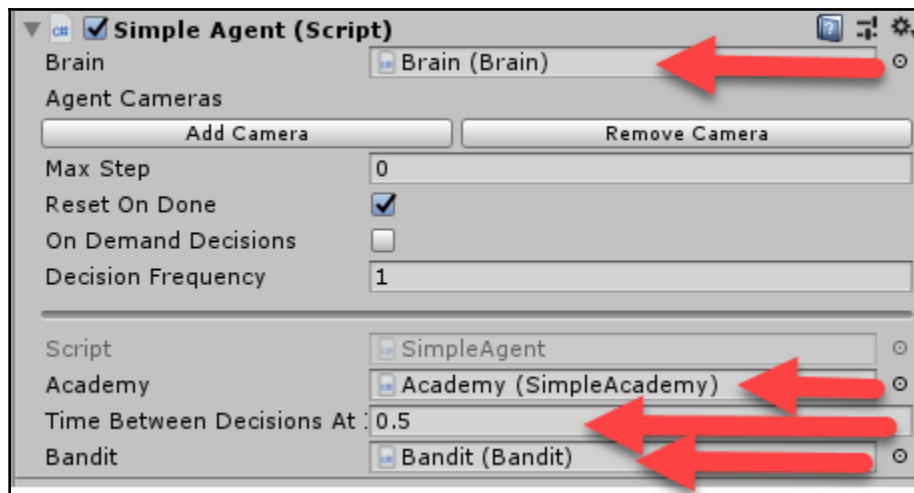
Decision Frequency: 1

Script: SimpleAgent

Academy: Academy (SimpleAcademy)

Time Between Decisions At: 0.5

Bandit: Bandit (Bandit)

A screenshot of the 'Simple Agent (Script)' configuration window. The window has a title bar with a dropdown arrow, a 'CR' icon, and a checked checkbox next to the title. Below the title bar, there are several sections. The first section is 'Brain', with a dropdown menu showing 'Brain (Brain)'. A red arrow points to this dropdown. The second section is 'Agent Cameras', with two buttons: 'Add Camera' and 'Remove Camera'. The third section contains four settings: 'Max Step' with a value of 0, 'Reset On Done' with a checked checkbox, 'On Demand Decisions' with an unchecked checkbox, and 'Decision Frequency' with a value of 1. A horizontal line separates this section from the next. The fourth section contains three settings: 'Script' with a dropdown showing 'SimpleAgent', 'Academy' with a dropdown showing 'Academy (SimpleAcademy)', and 'Time Between Decisions At' with a value of 0.5. Red arrows point to the 'Academy' dropdown and the 'Time Between Decisions At' field. The fifth section contains 'Bandit' with a dropdown showing 'Bandit (Bandit)'. A red arrow points to this dropdown. The window has a standard macOS-style title bar with a close button, a maximize button, and a settings gear icon.

Inspector

Services

Brain

Static

Tag Untagged

Layer Default

Transform

PositionX 0Y 0Z 0

RotationX 0Y 0Z 0

ScaleX 1Y 1Z 1

Brain (Script)

Brain Parameters

Vector Observation

Space TypeDiscrete

Space Size1

Stacked Vectors1

Visual Observation

Size0

Vector Action

Space TypeDiscrete

Space Size4

Action Descriptions

Size4

Element 01

Element 12

Element 23

Element 34

Brain TypePlayer

Broadcast

Edit the discrete inputs for your actions

Default Action0

Discrete Player Actions

Size4

Element 0

KeyA

Value1

Element 1

KeyS

Value2

Element 2

KeyD

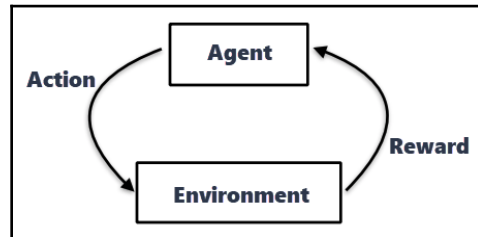
Value3

Element 3

KeyF

Value4

Chapter 2: The Bandit and Reinforcement Learning



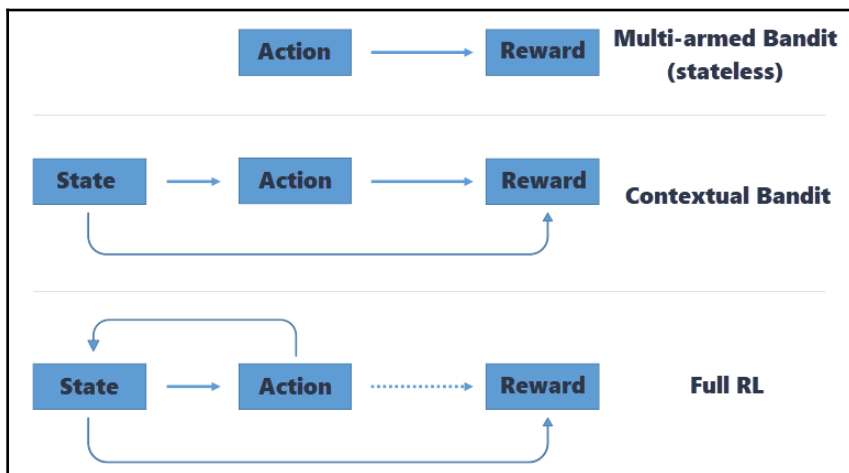
$$V(a) = V(a) + \alpha \times (r - V(a))$$

$$V(a) =$$

Simple Decision (Script)		Simple Decision (Script)	
Script	SimpleDecision	Script	SimpleDecision
Learning Rate	0.9	Learning Rate	0.9
Values		Values	
Size	4	Size	4
Element 0	0	Element 0	3
Element 1	0	Element 1	1
Element 2	0	Element 2	1
Element 3	0	Element 3	2

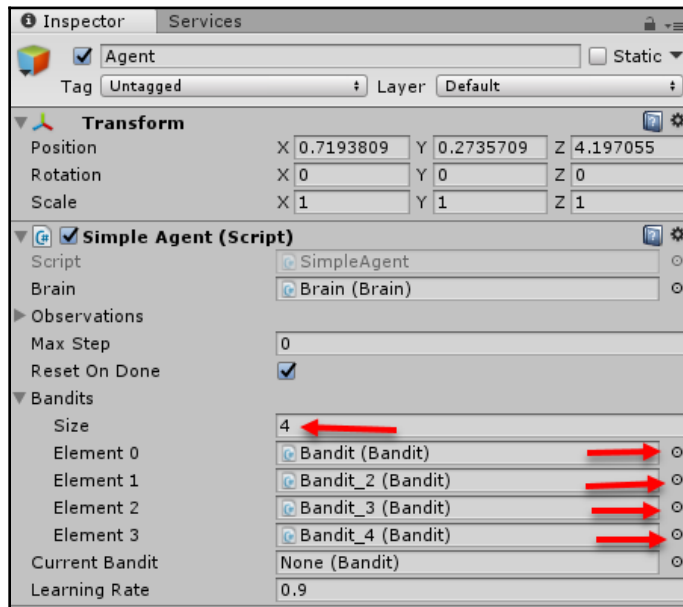
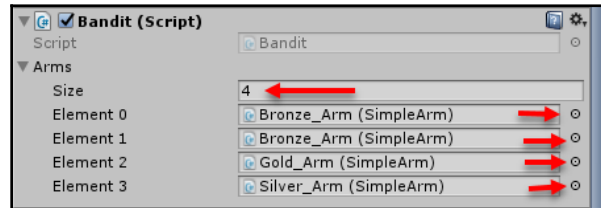
Configuration

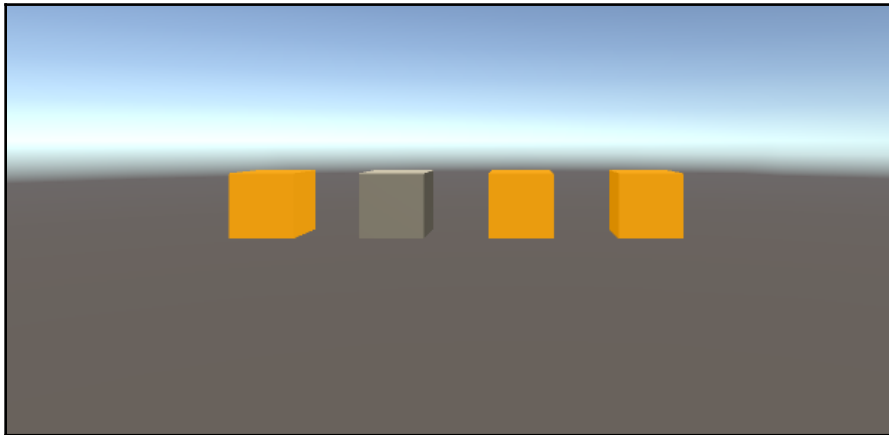
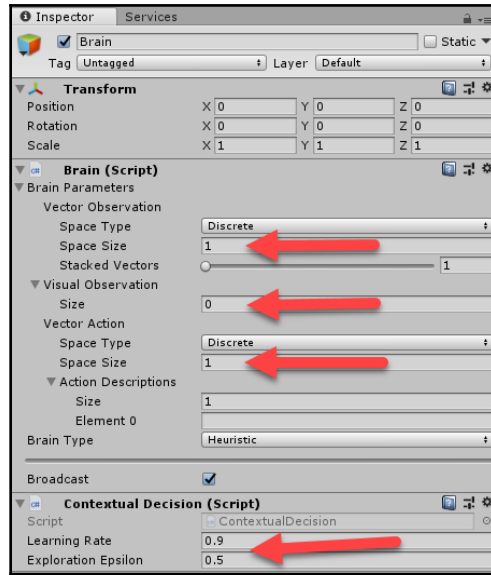
Completed



$$Q[s, a] = Q[s, a] + \alpha \times (r - Q[s, a])$$

$$Q[s, a] =$$





$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

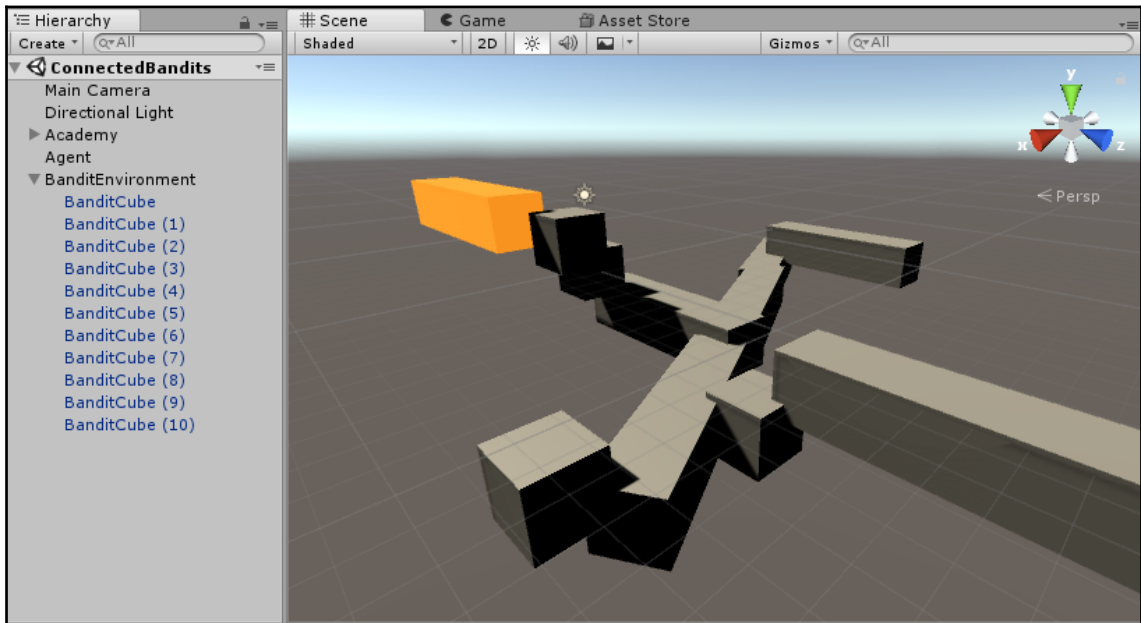
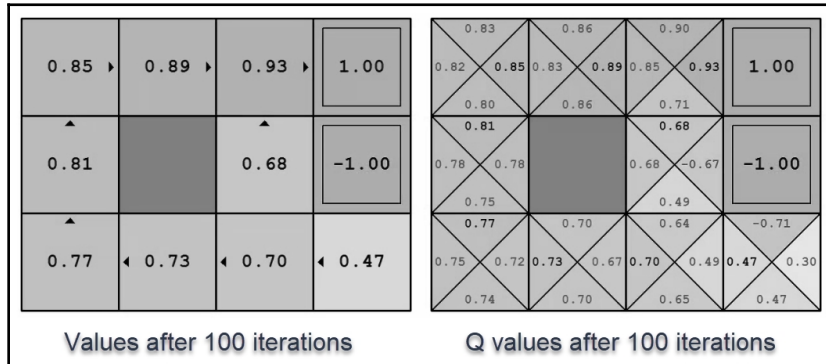
$$r = \text{reward}$$

$$\gamma = \text{gamma (reward discount factor 0 - 1.0)}$$

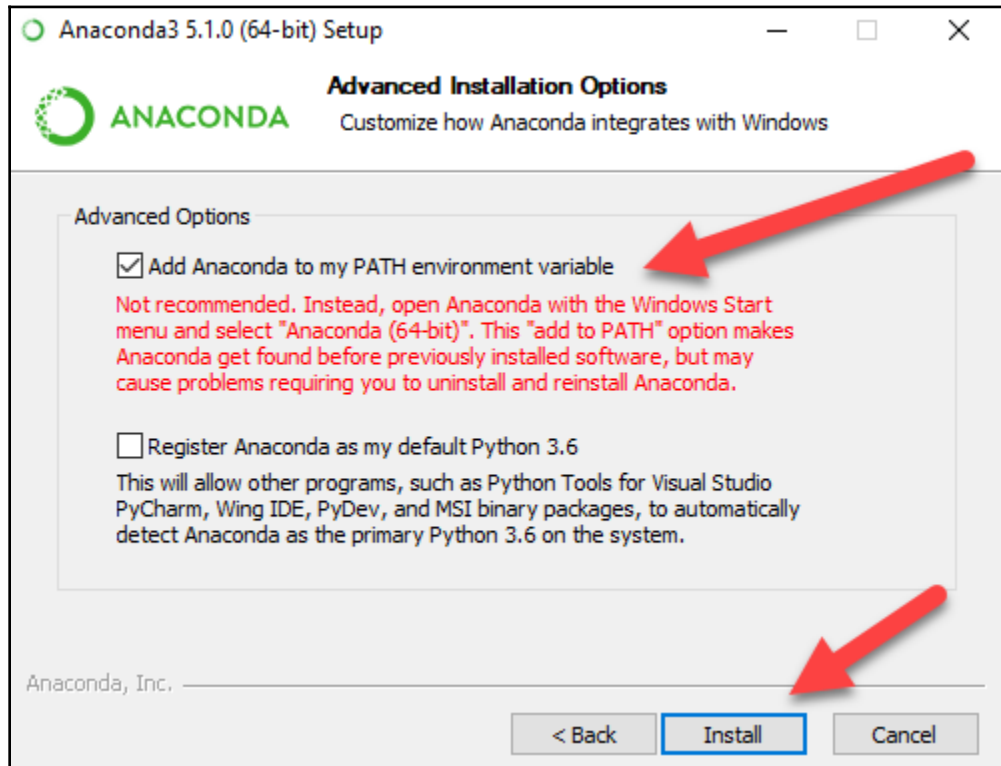
$$\max_{a'} = \text{maximum of all actions for state } a'$$

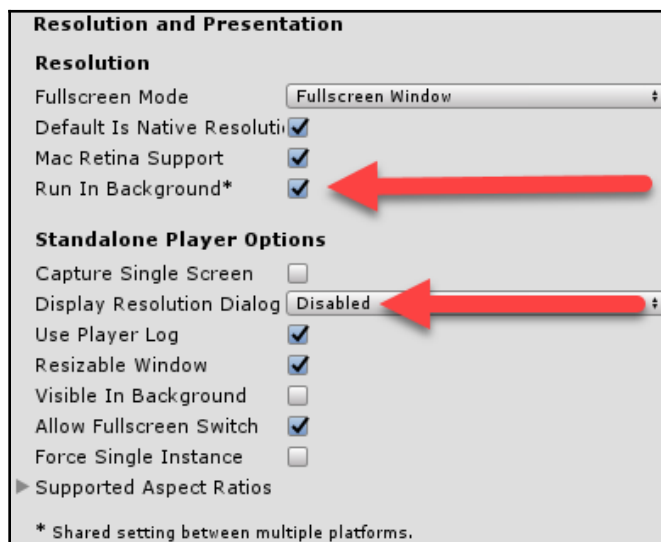
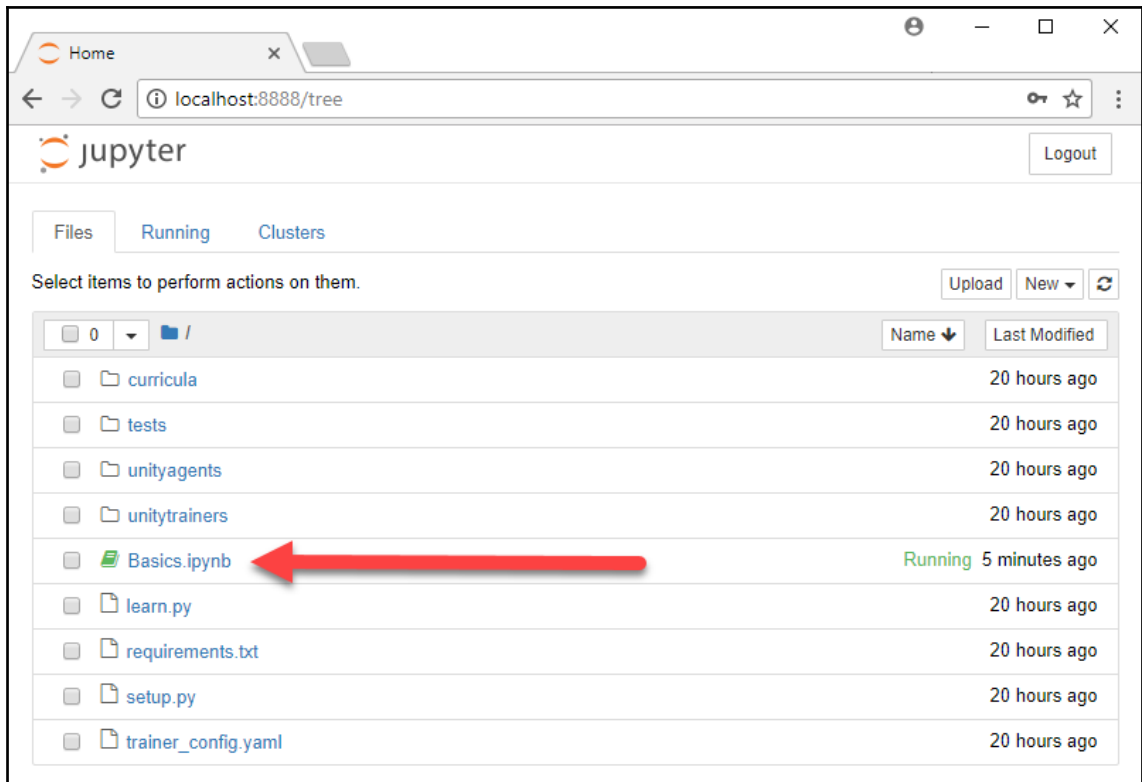
$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_a Q_t(s_{t+1}, a) - Q_t(s_{t+1}, a_t))$$

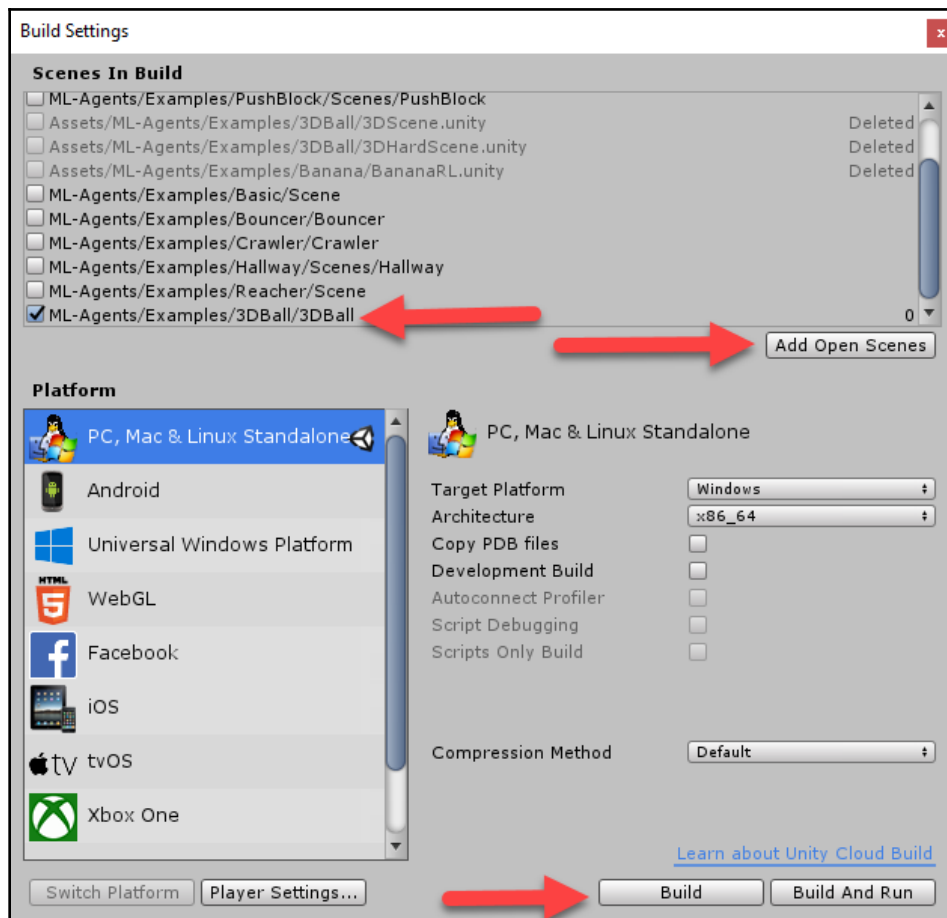
$\alpha = \text{learning rate}$

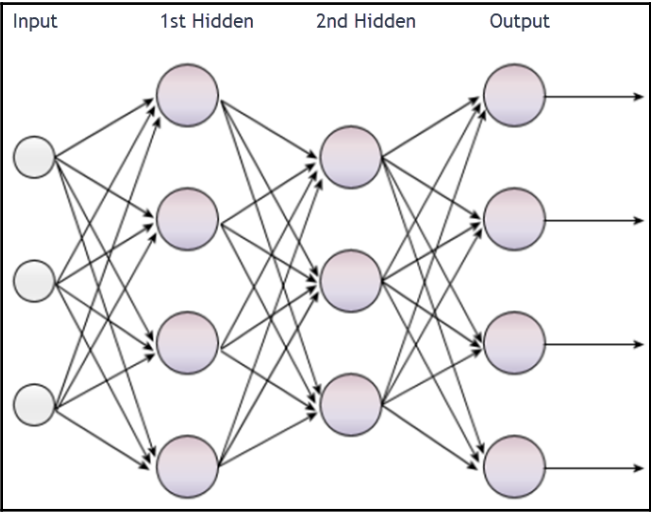
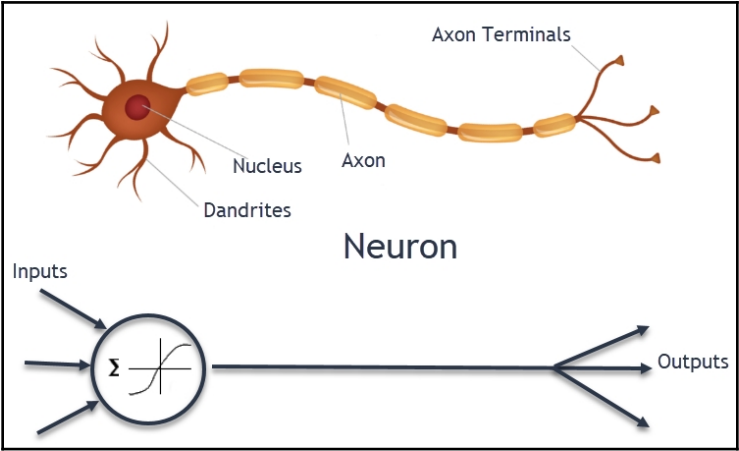


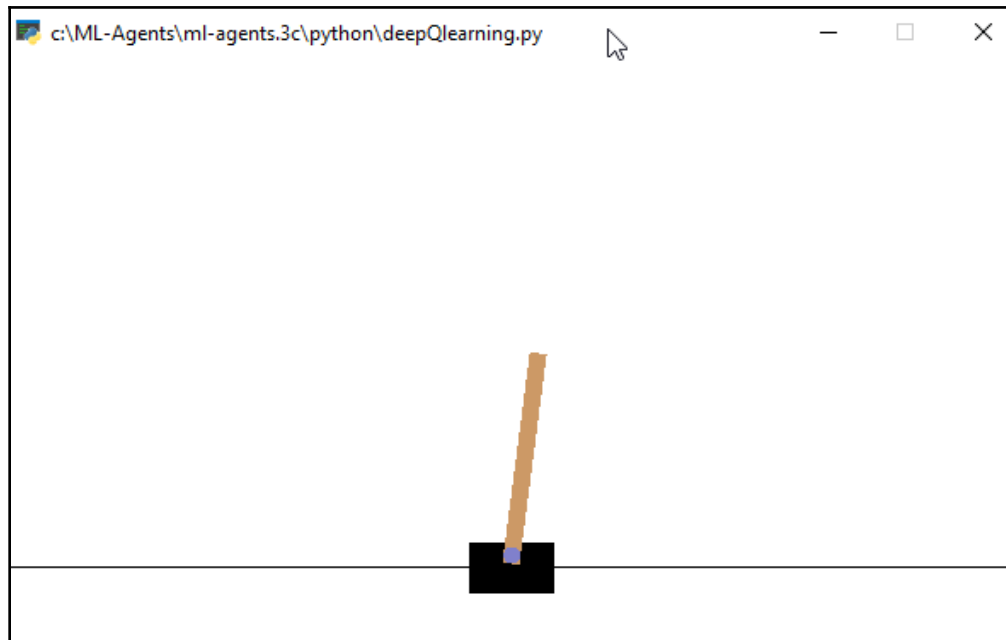
Chapter 3: Deep Reinforcement Learning with Python











```
20 model = Sequential()
21 model.add(Flatten(input_shape=(1,) + env.observa
22 model.add(Dense(16))
23 model.add(Activation('relu'))
24 model.add(Dense(nb_actions))
25 model.add(Activation('linear'))
26 print(model.summary())
27
28 policy = EpsGreedyQPolicy()
29 memory = SequentialMemory(limit=50000, window_le
30 dqn = DQNAgent(model=model, nb_actions=nb_action
31 target_model_update=1e-2, policy=policy)
32 dqn.compile(Adam(lr=1e-3), metrics=['mae'])
```

```
7
8 from rl.
9 from rl.
10 from rl.
11
12 ENV_NAME
13
14 # Get th
15 env = gy
16 np.rando
17 env.seed
18 nb_actio
19
20 model =
21 model.ad
22 model.ad
23 model.ad
24 model.ad
25 model.ad
26 print(model.summary())
27
28 policy = EpsGreedyQPolicy()
```

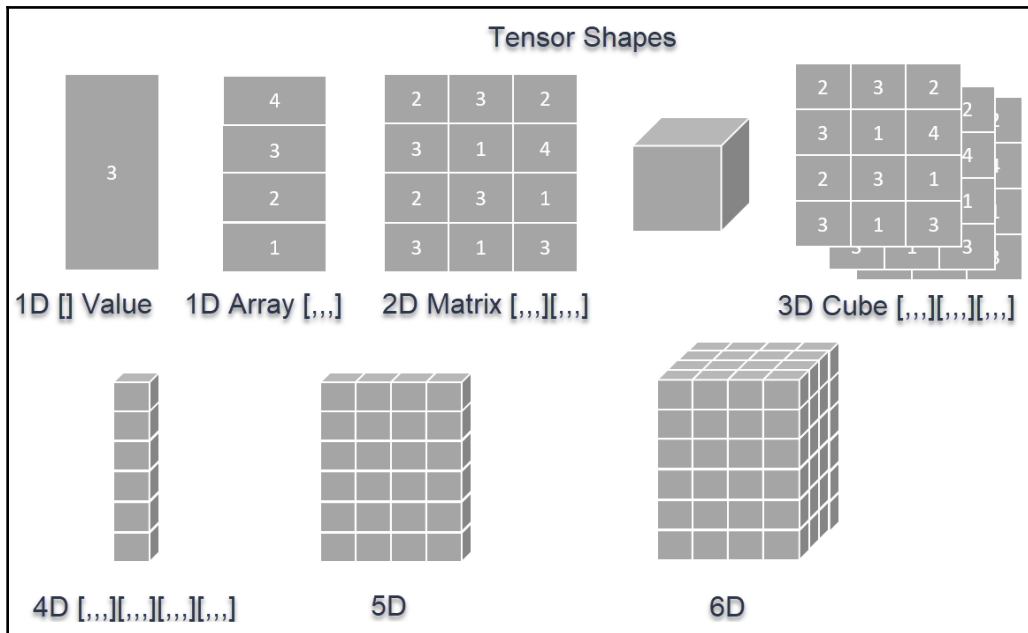


`<keras.models.Sequential object at`

- `input: <Tensor>`
- `OVERLOADABLE_OPERATORS: {'__truediv__', '_`
- `_consumers: [<tf.Operation 'flatt...type=Sh`
- `_dtype: tf.float32`
- `_handle_data: None`
- `_id: 0`
- `_keras_history: (<keras.engine.topolo...29`
- `_keras_shape: (None, 1, 4)`
- `_op: <tf.Operation 'flatten_1_input' type=`
- `_shape: TensorShape([`
- `_dims: [Dimension(None), Dimension(1), Di`
- `dims: [Dimension(None), Dimension(1), Dim`
- `ndims: 3`
- `[0]: Dimension(None)`
- `[1]: Dimension(1)`
- `[2]: Dimension(4)`
- `_uses_learning_phase: False`
- `_value_index: 0`

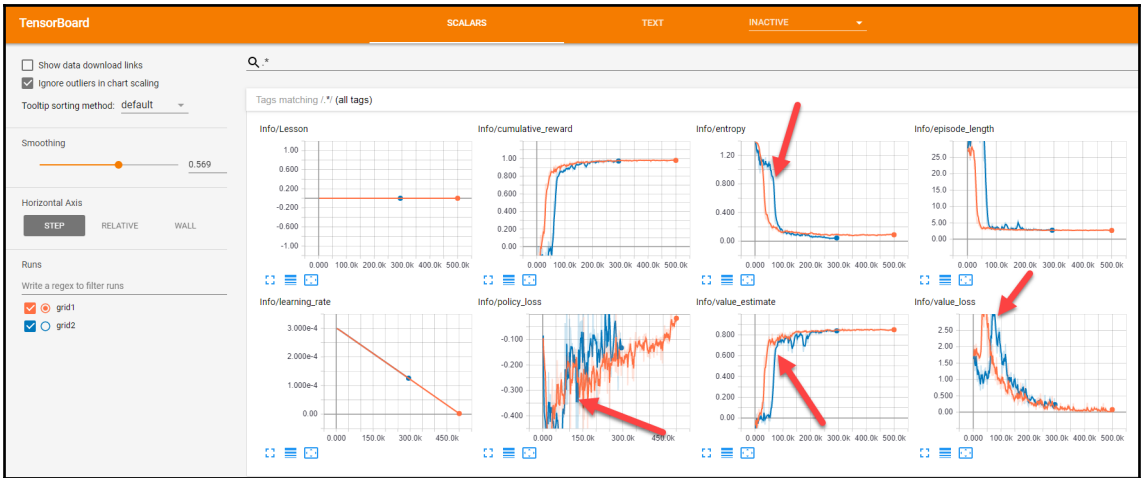
`print(model.summary())`

`policy = EpsGreedyQPolicy()`

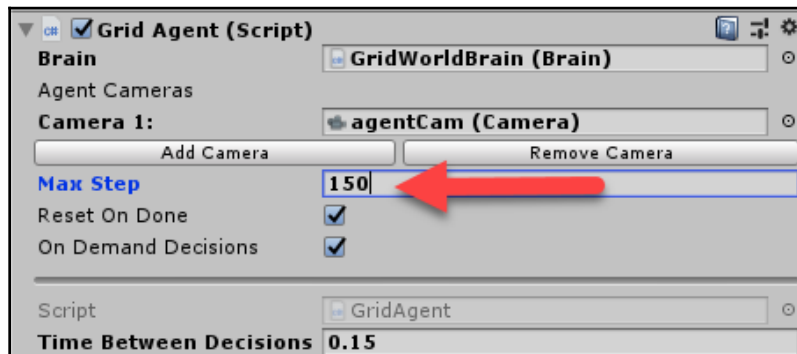
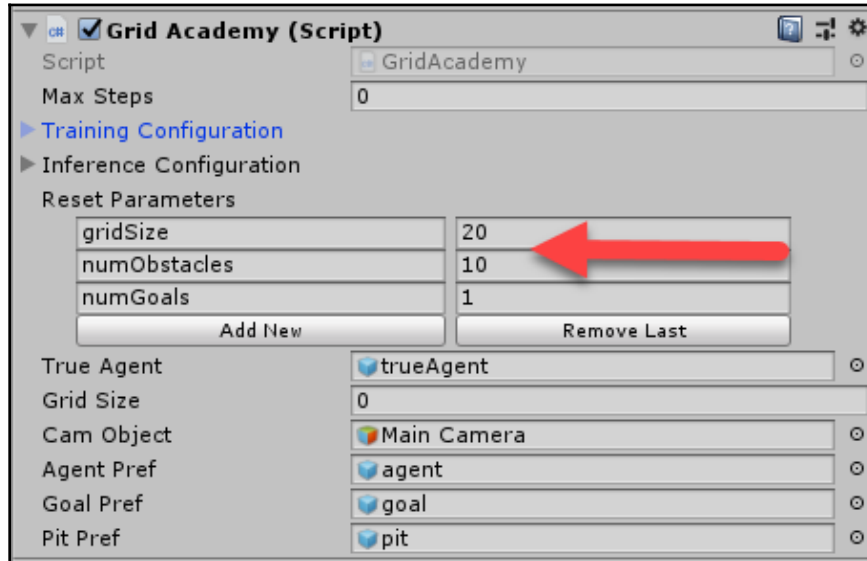


```

Anaconda Prompt
INFO:unityagents: GridWorldBrain: Step: 464000. Mean Reward: 0.982. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 466000. Mean Reward: 0.983. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 468000. Mean Reward: 0.967. Std of Reward: 0.175.
INFO:unityagents: GridWorldBrain: Step: 470000. Mean Reward: 0.983. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 472000. Mean Reward: 0.981. Std of Reward: 0.015.
INFO:unityagents: GridWorldBrain: Step: 474000. Mean Reward: 0.979. Std of Reward: 0.090.
INFO:unityagents: GridWorldBrain: Step: 476000. Mean Reward: 0.979. Std of Reward: 0.089.
INFO:unityagents: GridWorldBrain: Step: 478000. Mean Reward: 0.982. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 480000. Mean Reward: 0.982. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 482000. Mean Reward: 0.983. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 484000. Mean Reward: 0.982. Std of Reward: 0.015.
INFO:unityagents: GridWorldBrain: Step: 486000. Mean Reward: 0.983. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 488000. Mean Reward: 0.984. Std of Reward: 0.015.
INFO:unityagents: GridWorldBrain: Step: 490000. Mean Reward: 0.983. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 492000. Mean Reward: 0.984. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 494000. Mean Reward: 0.983. Std of Reward: 0.014.
INFO:unityagents: GridWorldBrain: Step: 496000. Mean Reward: 0.976. Std of Reward: 0.123.
INFO:unityagents: GridWorldBrain: Step: 498000. Mean Reward: 0.984. Std of Reward: 0.014.
INFO:unityagents: Saved Model
INFO:unityagents: GridWorldBrain: Step: 500000. Mean Reward: 0.979. Std of Reward: 0.089.
INFO:unityagents: Saved Model
INFO:unityagents: Saved Model
INFO:unityagents: List of nodes to export :
INFO:unityagents:      action
INFO:unityagents:      value_estimate
INFO:unityagents:      action_probs
INFO:tensorflow: Restoring parameters from ./models/grid1\model-500000.cptk
INFO:tensorflow: Restoring parameters from ./models/grid1\model-500000.cptk
INFO:tensorflow: Froze 8 variables.
INFO:tensorflow: Froze 8 variables.
  
```



Chapter 4: Going Deeper with Deep Learning



```
INFO:unityagents: GridWorldBrain: Step: 2000. Mean Reward: -1.394. Std of Reward: 0.468.  
INFO:unityagents: GridWorldBrain: Step: 4000. Mean Reward: -1.390. Std of Reward: 0.733.  
INFO:unityagents: GridWorldBrain: Step: 6000. Mean Reward: -1.319. Std of Reward: 0.490.  
INFO:unityagents: GridWorldBrain: Step: 8000. Mean Reward: -1.273. Std of Reward: 0.683.  
INFO:unityagents: GridWorldBrain: Step: 10000. Mean Reward: -1.158. Std of Reward: 0.647.  
INFO:unityagents: GridWorldBrain: Step: 12000. Mean Reward: -1.413. Std of Reward: 0.274.  
INFO:unityagents: GridWorldBrain: Step: 14000. Mean Reward: -1.109. Std of Reward: 0.883.  
INFO:unityagents: GridWorldBrain: Step: 16000. Mean Reward: -1.251. Std of Reward: 0.472.  
INFO:unityagents: GridWorldBrain: Step: 18000. Mean Reward: -1.203. Std of Reward: 0.593.  
INFO:unityagents: GridWorldBrain: Step: 20000. Mean Reward: -1.199. Std of Reward: 0.672.
```

Brain (Script)

▼ Brain Parameters

Vector Observation

Space Type: Continuous

Space Size: 0

Stacked Vectors: 1

▼ Visual Observation

Size: 1

▼ Element 0

Width: 128

Height: 128

Black And White: ☒

Vector Action

Space Type: Discrete

Space Size: 4

▼ Action Descriptions

Size: 4

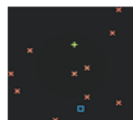
Element 0:

Element 1:

Element 2:

Element 3:

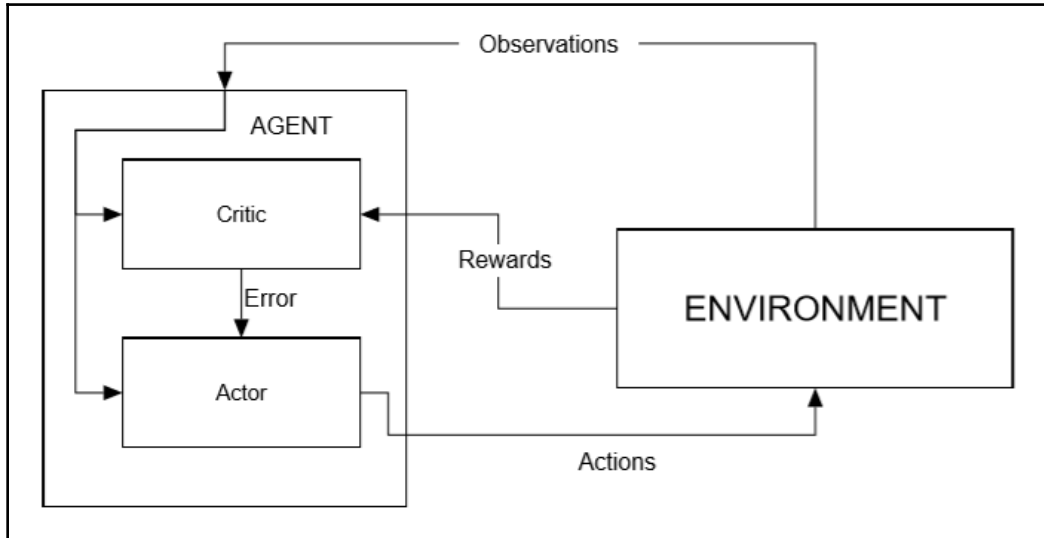
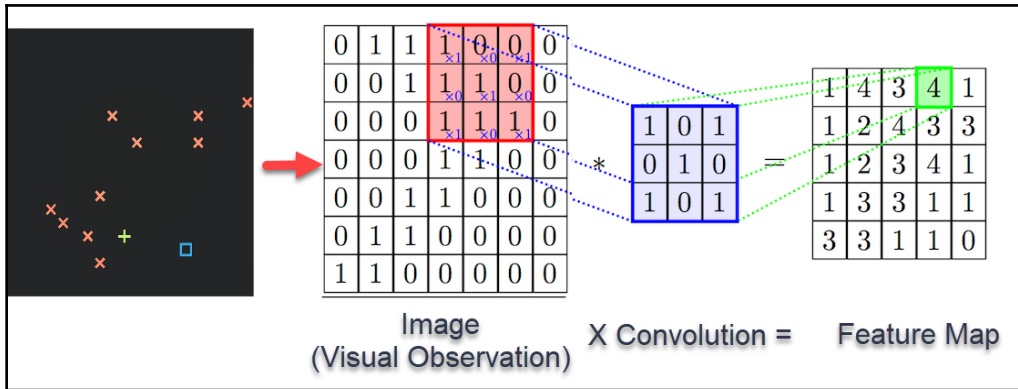
Brain Type: External

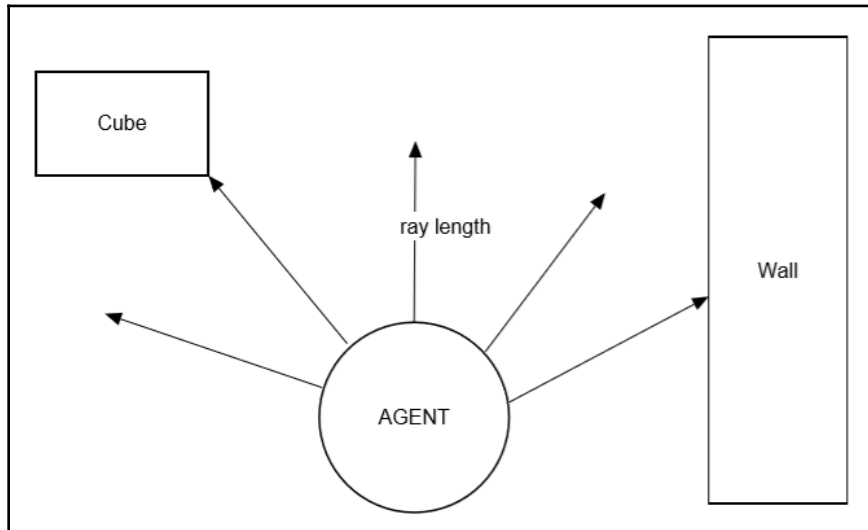
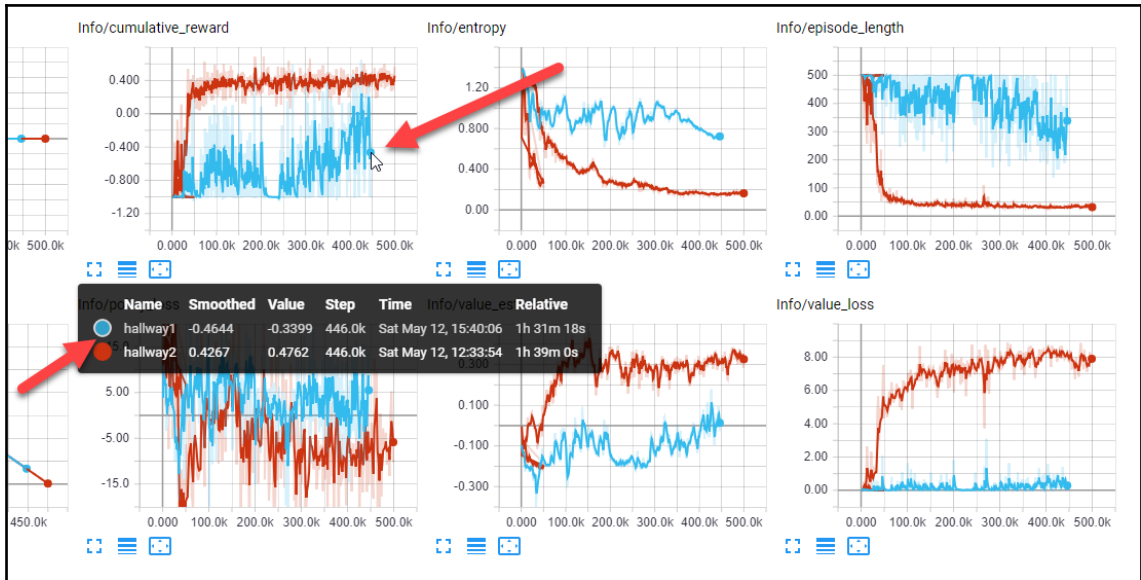


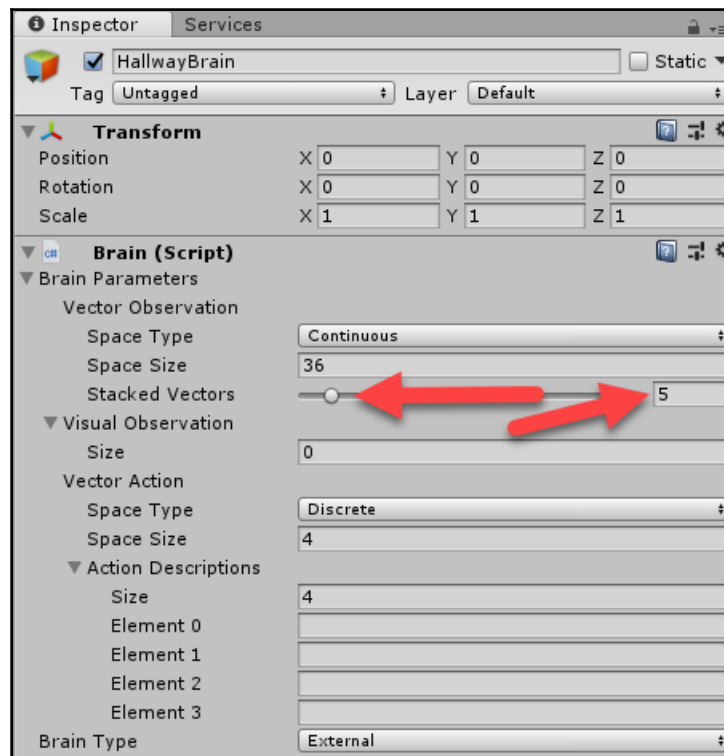
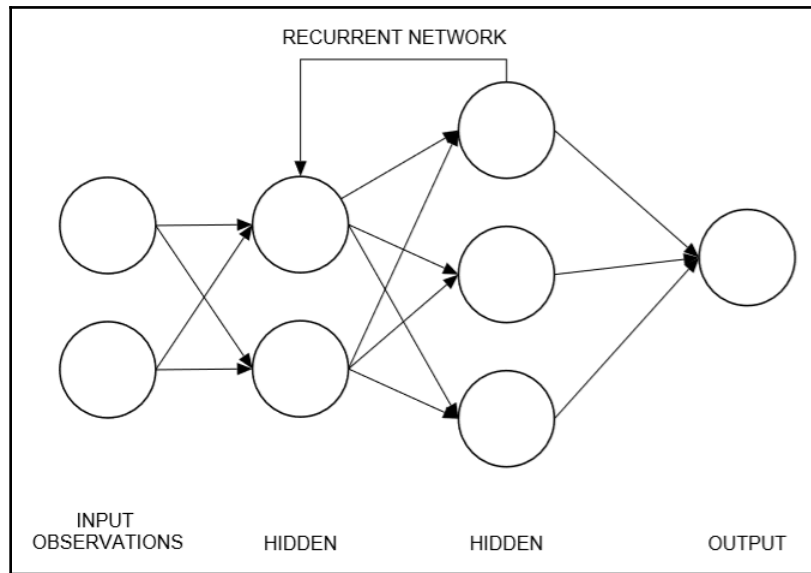
84x84
color



128X128
B&W

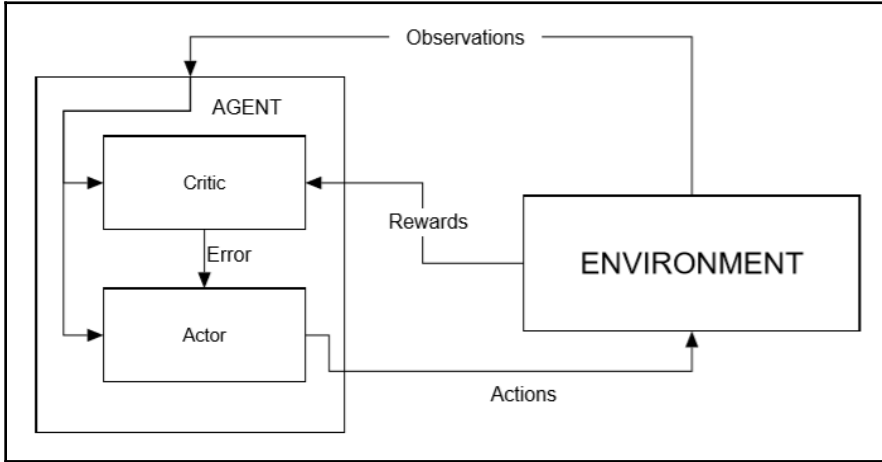






$$Advantage : A = Q(s, a) - V(s)$$

$$Advantage : A = R - V(s)$$



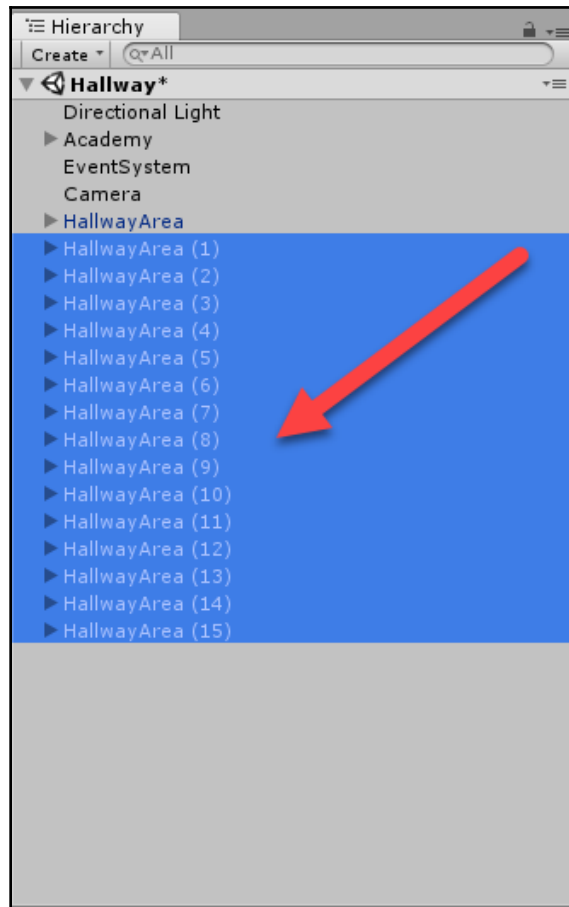
$$ValueLoss : L = \Sigma(R - V(s))^2 (SumSquaredError)$$

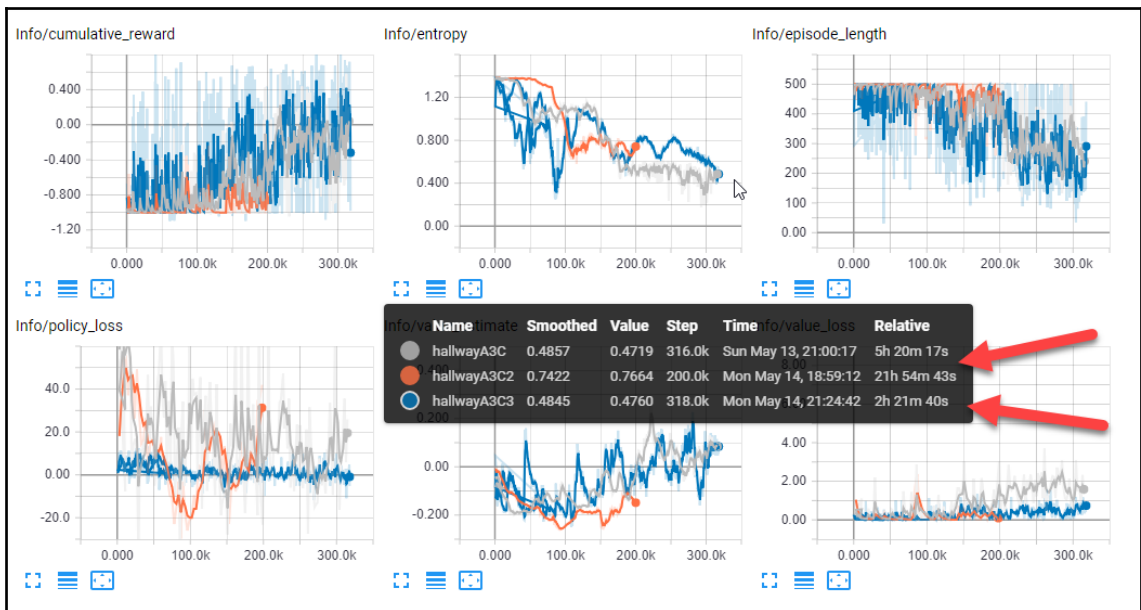
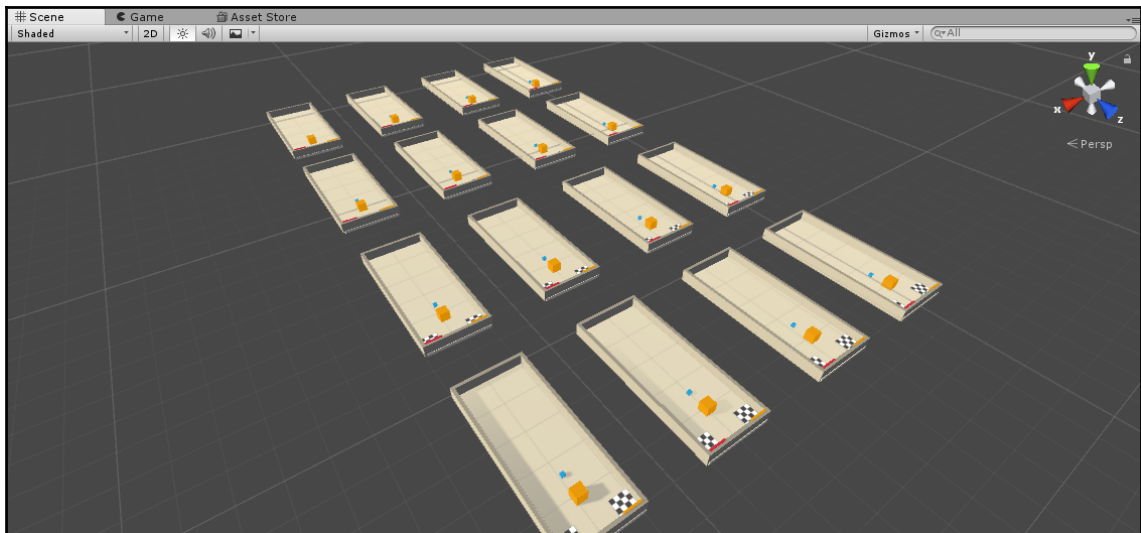
$$PolicyLoss : L = -\log(\pi(a|s)) * A(s)$$

$$H(\pi) = -\Sigma(P(x)\log(P(x)))$$

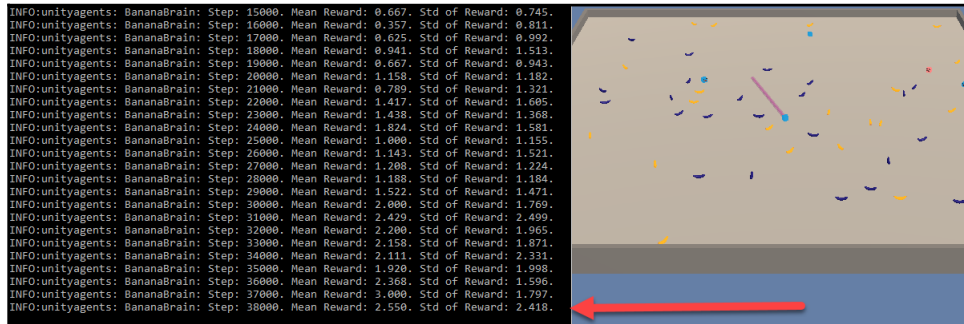
$$PolicyLoss : L = -\log(\pi(a|s)) * A(s) - \beta * H(\pi)$$

$$L = 0.5 * \Sigma(R - V(s))^2 - \log(\pi(a|s)) * A(s) - \beta * H(\pi)$$



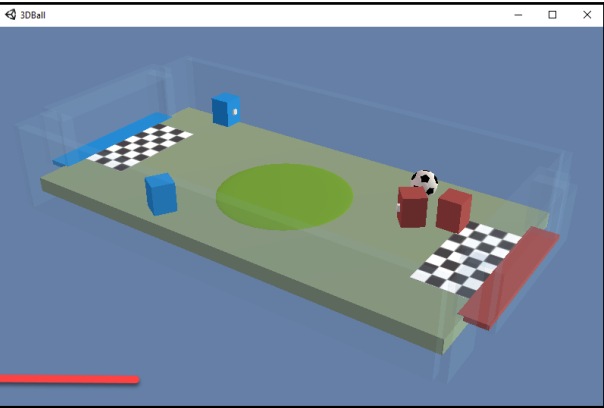


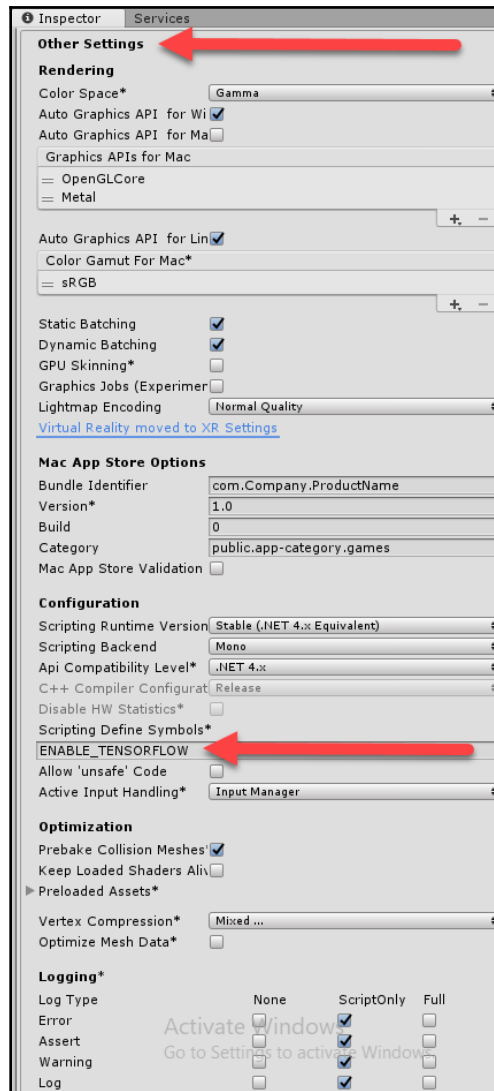
Chapter 5: Playing the Game

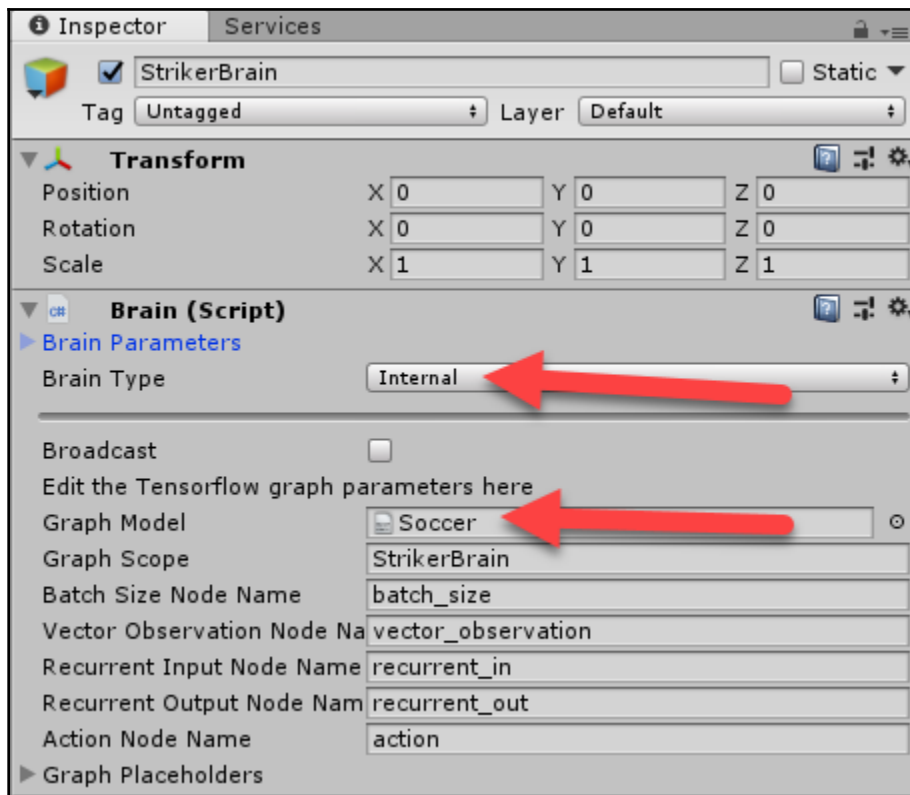


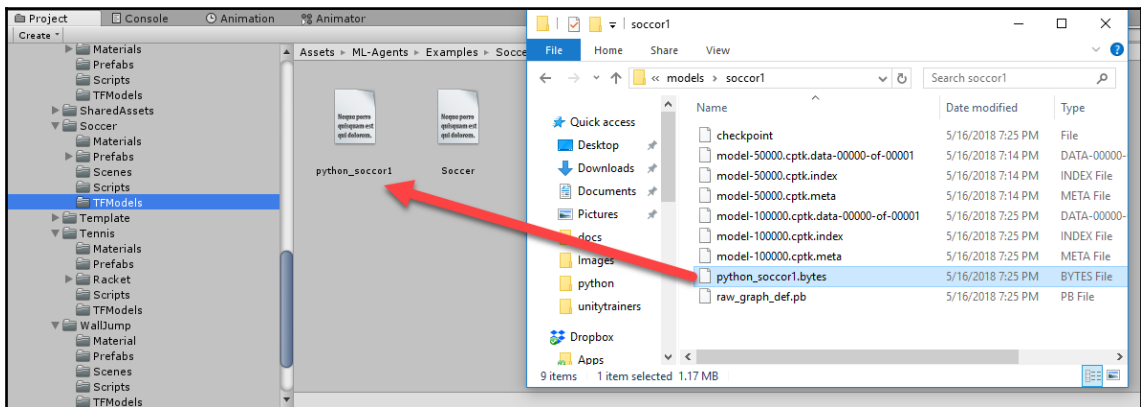
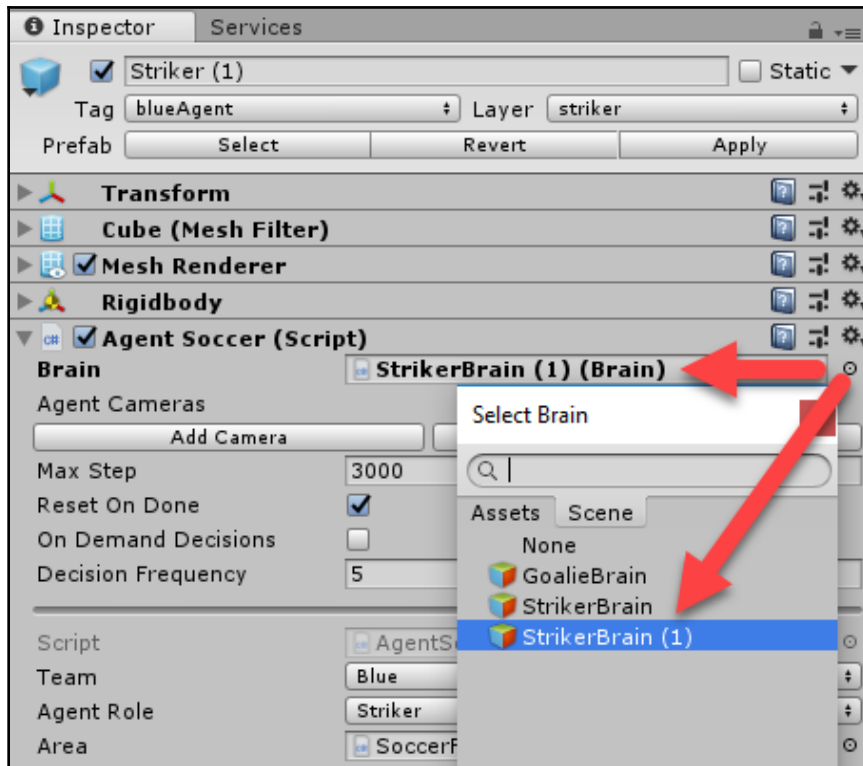
```
Anaconda Prompt - python python/learn.py python/python.exe --run-id=banana1 --train

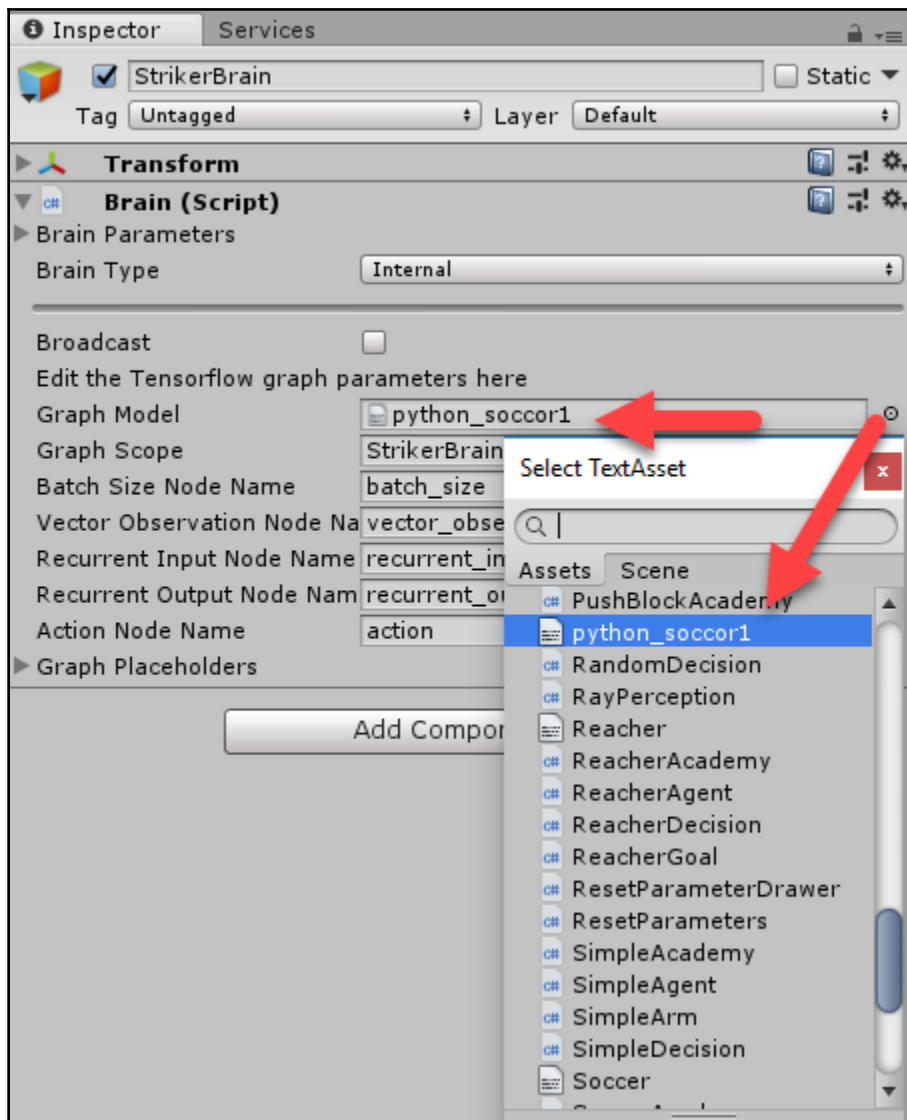
use_recurrent: False
graph_scope: GoalieBrain
summary_path: ./summaries/banana1_GoalieBrain
memory_size: 256
(INFO:unityagents:Hyperparameters for the PPO Trainer of brain StrikerBrain:
batch_size: 128
beta: 0.005
buffer_size: 2048
epsilon: 0.2
gamma: 0.99
hidden_units: 256
lambda: 0.95
learning_rate: 0.0003
max_steps: 1.0e5
normalize: False
num_epoch: 3
num_layers: 2
time_horizon: 128
sequence_length: 64
summary_freq: 2000
use_recurrent: False
graph_scope: StrikerBrain
summary_path: ./summaries/banana1_StrikerBrain
memory_size: 256
(INFO:unityagents: GoalieBrain: Step: 2000. Mean Reward: 0.198. Std of Reward: 0.782.
(INFO:unityagents: StrikerBrain: Step: 2000. Mean Reward: -0.198. Std of Reward: 0.782.
(INFO:unityagents: GoalieBrain: Step: 4000. Mean Reward: 0.736. Std of Reward: 0.583.
(INFO:unityagents: StrikerBrain: Step: 4000. Mean Reward: -0.736. Std of Reward: 0.583.
(INFO:unityagents: GoalieBrain: Step: 6000. Mean Reward: 0.313. Std of Reward: 0.629.
(INFO:unityagents: StrikerBrain: Step: 6000. Mean Reward: -0.313. Std of Reward: 0.629.
```

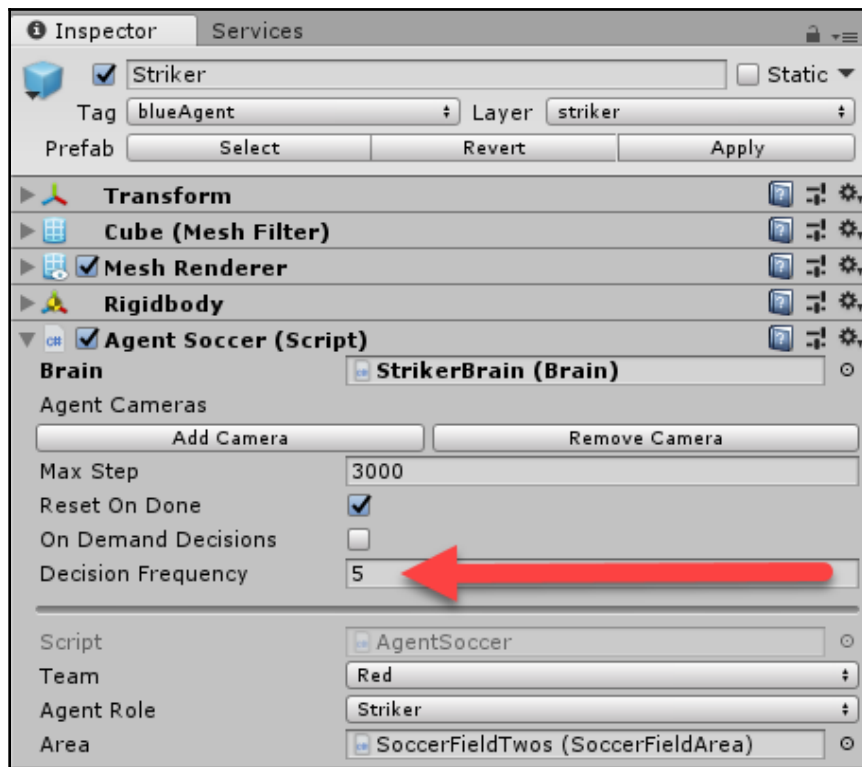


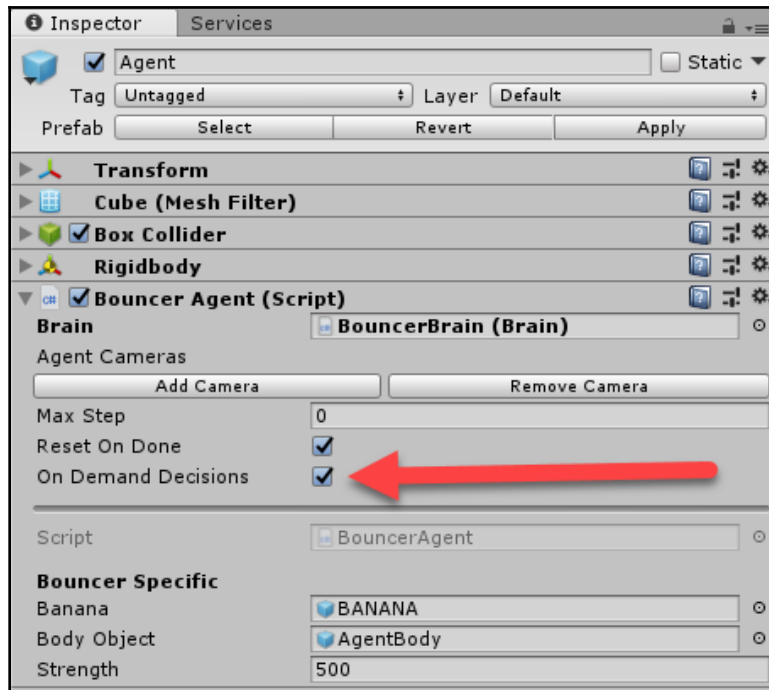


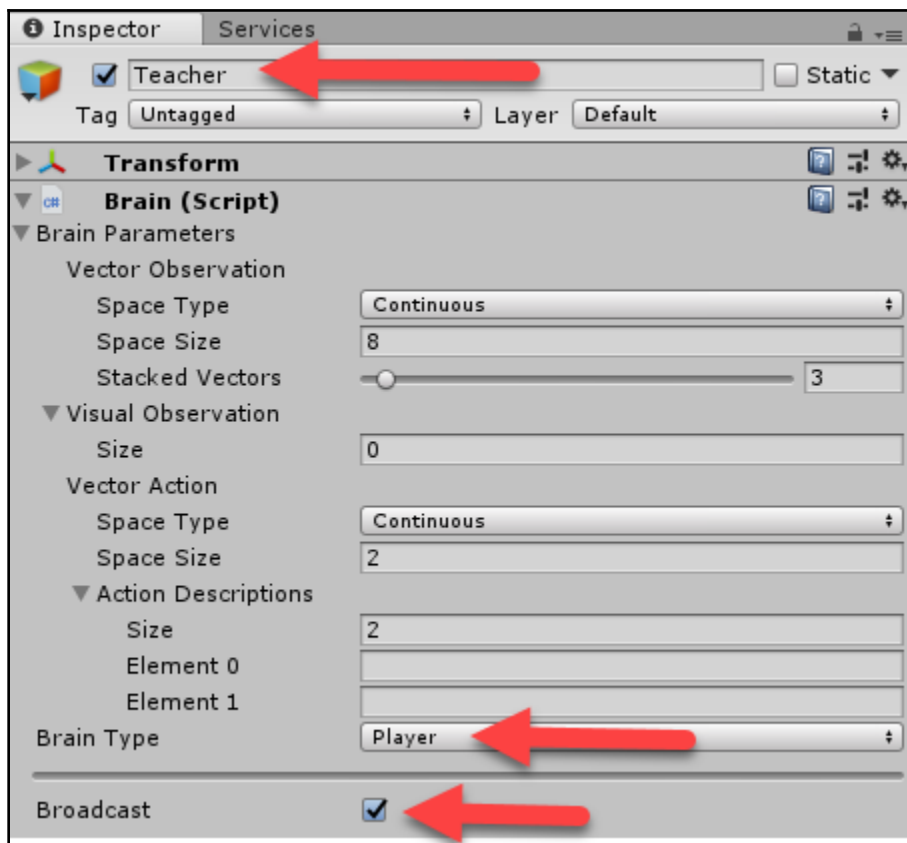


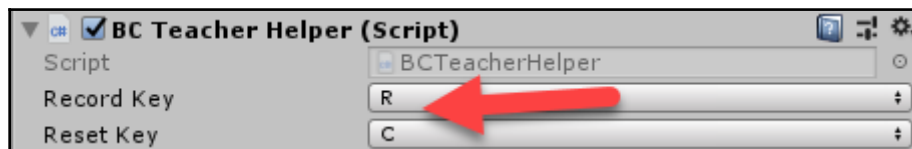
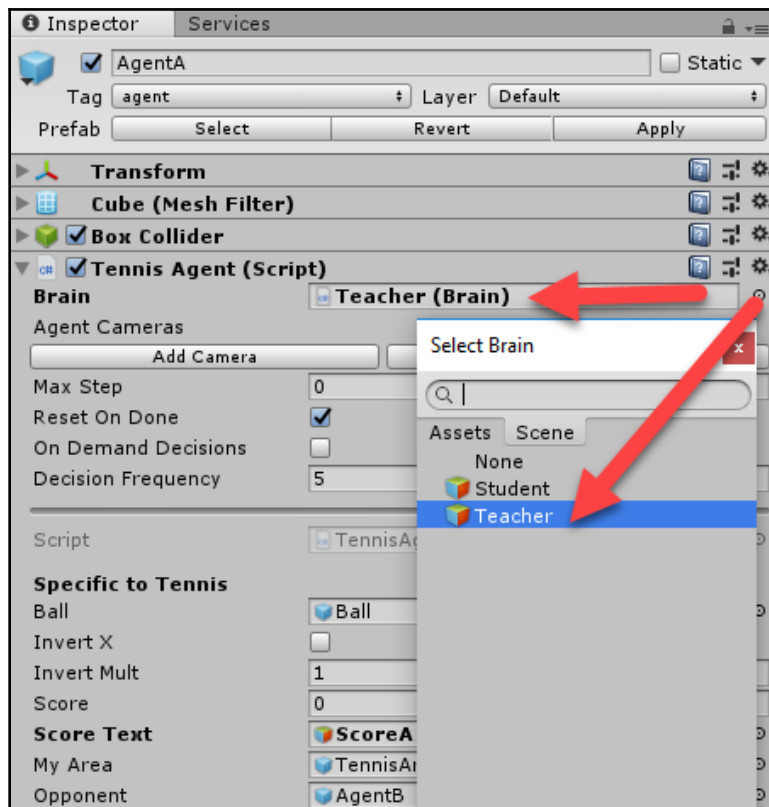


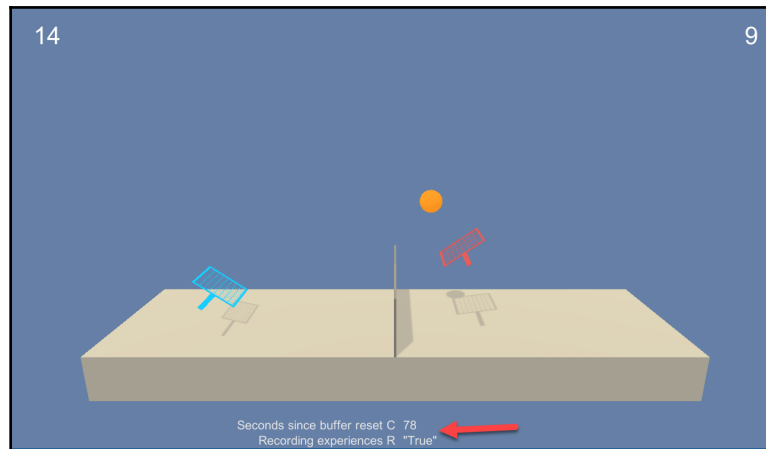












Inspector **Services**

☒ Academy ☐ Static
Tag: Untagged Layer: Default

Transform

☒ **Wall Jump Academy (Script)**

Script: WallJumpAcademy

Max Steps: 10000

► Training Configuration

► Inference Configuration

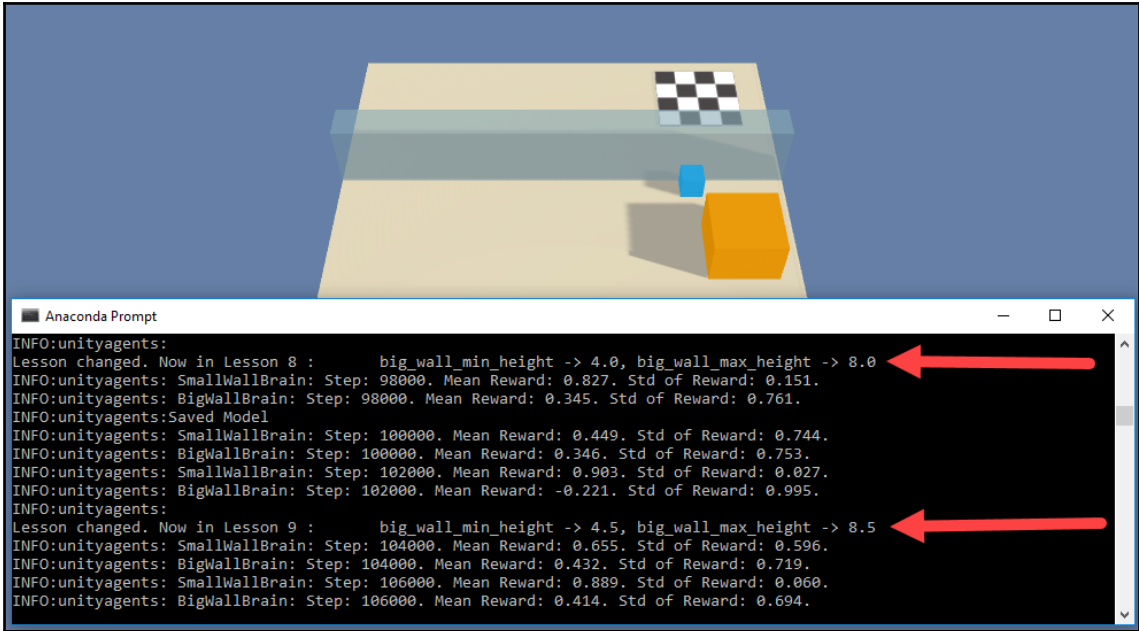
Reset Parameters

big_wall_min_height	8
small_wall_height	4
no_wall_height	0
big_wall_max_height	8

Add New Remove Last

Specific to WallJump

Agent Run Speed	1.5
Agent Jump Height	2.75
Goal Scored Material	SuccessGround
Fail Material	FailGround



Chapter 6: Terrarium Revisited – A Multi-Agent Ecosystem



