```python
1   class Agent():
2    def __init__(self, ff_inputsize):
3     self.ff_inputsize = ff_inputsize
4     self.model = DDDQN_model
5     self.memory = Memory(10000, self)  #for definition see code
6     self.action_repeat = 4
7     self.update_frequency = 4
8     self.batch_size = 32
9     self.replaystartsize = 1000
10    self.epsilon = 0.05
11    self.last_action = None
12    self.repeated_action_for = self.action_repeat


14   def runInference(self, gameState, pastState):
15    self.addToMemory(gameState, pastState)
16    ff_inputs = self.getAgentState(*gameState)
17    self.repeated_action_for += 1
18    if self.repeated_action_for < self.action_repeat:
19     toUse, toSave = self.last_action
20    else:
21     self.repeated_action_for = 0
22     if self.canLearn() and np.random.random() > self.epsilon:
23      action, _ = self.model.inference(ff_inputs)
24      toSave = self.dediscretize(action[0])
25      toUse = "["+str(throttle)+", "+str(brake)+", "+str(steer)+"]"
26     else:
27      toUse, toSave = self.randomAction() #for definition see code
28     self.last_action = toUse, toSave
29    self.containers.outputval.update(toUse, toSave)
30    self.numsteps += 1
31    if self.numsteps % self.update_frequency == 0 and len(self.memory) > self.replaystartsize):
32     self.learnStep()


34   def learnStep(self):
35    QLearnInputs = self.memory.sample(self.batch_size)
36    self.model.q_learn(QLearnInputs)


38   def addToMemory(self, gameState, pastState):
39     s = self.getAgentState(*pastState)  #for definition see code
40     a = self.getAction(*pastState)   #for definition see code
41     r = self.calculateReward(*gameState)#for definition see code
42     s2= self.getAgentState(*gameState)  #for definition see code
43     t = False #will be updated if episode did end
44     self.memory.append([s,a,r,s2,t])



47   class DuelDQN():
48    def __init__(self, name, ff_inputsize, num_actions):
49     with tf.variable_scope(name, initializer = tf.random_normal_initializer(0, 1e-3)):
50      #for the inference
51      self.ff_inputs = tf.placeholder(tf.float32, shape=[None, ff_inputsize], name="ff_inputs")
52      self.fc1 = tf.layers.dense(self.ff_inputs, 400, activation=tf.nn.relu)
53      #modifications from the Dueling DQN architecture
54      self.streamA, self.streamV = tf.split(self.fc1,2,1)
55      xavier_init = tf.contrib.layers.xavier_initializer()
```

```python
56      neutral_init = tf.random_normal_initializer(0, 1e-50)
57      self.AW = tf.Variable(xavier_init([200,self.num_actions]))
58      self.VW = tf.Variable(neutral_init([200,1]))
59      self.Advantage = tf.matmul(self.streamA,self.AW)
60      self.Value = tf.matmul(self.streamV,self.VW)
61      self.Qout = self.Value + tf.subtract(self.Advantage,tf.reduce_mean(self.Advantage,axis=1,keep_dims=True))
62      self.Qmax = tf.reduce_max(self.Qout, axis=1)
63      self.predict = tf.argmax(self.Qout,1)
64      #for the learning
65      self.targetQ = tf.placeholder(shape=[None],dtype=tf.float32)
66      self.targetA = tf.placeholder(shape=[None],dtype=tf.int32)
67      self.targetA_OH = tf.one_hot(self.targetA, self.num_actions, dtype=tf.float32)
68      self.compareQ = tf.reduce_sum(tf.multiply(self.Qout, self.targetA_OH), axis=1)
69      self.td_error = tf.square(self.targetQ - self.compareQ)
70      self.q_loss = tf.reduce_mean(self.td_error)
71      q_trainer = tf.train.AdamOptimizer(learning_rate=0.00025)
72      q_OP = q_trainer.minimize(self.q_loss)
73    self.trainables = tf.get_collection(tf.GraphKeys.TRAINABLE_VARIABLES, scope=name)


76  def _netCopyOps(fromNet, toNet, tau = 1):
77   op_holder = []
78   for idx,var in enumerate(fromNet.trainables[:]):
79    op_holder.append(toNet.trainables[idx].assign((var.value()*tau) + ((1-tau)*toNet.trainables[idx].value())))
80   return op_holder
```

# see agent
```python
1   #see agent
2   class DDDQN_model():
3       def __init__(self, session, ff_inputsize, num_action):
4           self.session = session
5           self.onlineQN = DuelDQN("onlineNet", ff_inputsize, num_action)
6           self.targetQN = DuelDQN("targetNet", ff_inputsize, num_action)
7           self.session.run(tf.global_variables_initializer())
8           self.session.run(_netCopyOps(self.targetQN, self.onlineQN))

10      #see agent
11      def inference(self, statesBatch): #called for every step t

13          return self.session.run([self.onlineQN.predict, self.onlineQN.Qout], feed_dict={self.onlineQN.
                ↳ ff_inputs: statesBatch})
14      #see agent
15      #see agent
16      #see agent
17      def q_learn(self, batch): #also called for every step t
18          oldstates, actions, rewards, newstates, terminals = batch
19          action = self.session.run(self.onlineQN.predict,feed_dict={self.onlineQN.ff_inputs:newstates})
20          folgeQ = self.session.run(self.targetQN.Qout,feed_dict={self.targetQN.ff_inputs:newstates})
21          doubleQ = folgeQ[range(len(terminals)),action]
22          consider_stateval = -(terminals - 1)
23          targetQ = rewards + (0.99 * doubleQ * consider_stateval)
24          self.session.run(self.onlineQN.q_OP, feed_dict={self.onlineQN.ff_inputs:oldstates, self.onlineQN.
                ↳ targetQ:targetQ, self.onlineQN.targetA:actions})
25          self.session.run(_netCopyOps(self.onlineQN, self.targetQN, 0.001))
26          return
```

1. Initialize replay memory $D$ to capacity N

5. Initialize action-value function $Q(s,a;\theta)$ with random weights $\theta$

8. Initialize target action-value function $Q(s,a;\theta^-)$ with weights $\theta^- = \theta$

9. **For** episode = 1,M **do**

10. Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

11. **For** 1 = 1,T **do**

12. With probability $\epsilon$ select random action $a_t$

13. **otherwise** select $a_t = argmax_a Q(\phi(s_t),a;\theta)$

15. Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

16. Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

17. Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$

19. Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$

20. Define $a^{max}(\phi_{j+1};\theta) = argmax_{a'} Q(\phi_{j+1},a';\theta)$

21. Define $Q^{j+1} = Q(\phi_{j+1}, a^{max}(\phi_{t+1};\theta);\theta^-)$

23. **If** episode terminates at step $j+1$ **then** set $y_j = r_j$,
    ↳ **otherwise** set $y_j = r_j + \gamma * Q^{j+1}$

24. Perform a gradient descent step on $(y_j - Q(\phi_j,a_j;\theta))^2$ with respect to
    ↳ the network parameters $\theta$

25. Update target network: $\theta^- \leftarrow \tau * \theta + (1-\tau)\theta^-$

26. **End For**

27. **End For**