

8.recurrent-q-learning-agent

September 29, 2021

```
[1]: import numpy as np
import pandas as pd
import tensorflow as tf
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
```

```
[2]: df = pd.read_csv('../dataset/G00G-year.csv')
df.head()
```

```
[2]:
```

	Date	Open	High	Low	Close	Adj Close	\
0	2016-11-02	778.200012	781.650024	763.450012	768.700012	768.700012	
1	2016-11-03	767.250000	769.950012	759.030029	762.130005	762.130005	
2	2016-11-04	750.659973	770.359985	750.560974	762.020020	762.020020	
3	2016-11-07	774.500000	785.190002	772.549988	782.520020	782.520020	
4	2016-11-08	783.400024	795.632996	780.190002	790.510010	790.510010	

	Volume
0	1872400
1	1943200
2	2134800
3	1585100
4	1350800

```
[3]: from collections import deque
import random

class Agent:

    LEARNING_RATE = 0.003
    BATCH_SIZE = 32
    LAYER_SIZE = 256
    OUTPUT_SIZE = 3
    EPSILON = 0.5
    DECAY_RATE = 0.005
    MIN_EPSILON = 0.1
    GAMMA = 0.99
```

```

MEMORIES = deque()
MEMORY_SIZE = 300

def __init__(self, state_size, window_size, trend, skip):
    self.state_size = state_size
    self.window_size = window_size
    self.half_window = window_size // 2
    self.trend = trend
    self.skip = skip
    tf.reset_default_graph()
    self.INITIAL_FEATURES = np.zeros((4, self.state_size))
    self.X = tf.placeholder(tf.float32, (None, None, self.state_size))
    self.Y = tf.placeholder(tf.float32, (None, self.OUTPUT_SIZE))
    cell = tf.nn.rnn_cell.LSTMCell(self.LAYER_SIZE, state_is_tuple = False)
    self.hidden_layer = tf.placeholder(tf.float32, (None, 2 * self.
↪LAYER_SIZE))
    self.rnn, self.last_state = tf.nn.dynamic_rnn(inputs=self.X, cell=cell,
                                                    dtype=tf.float32,
                                                    initial_state=self.
↪hidden_layer)
    self.logits = tf.layers.dense(self.rnn[:, -1], self.OUTPUT_SIZE)
    self.cost = tf.reduce_sum(tf.square(self.Y - self.logits))
    self.optimizer = tf.train.AdamOptimizer(learning_rate = self.
↪LEARNING_RATE).minimize(self.cost)
    self.sess = tf.InteractiveSession()
    self.sess.run(tf.global_variables_initializer())

    def _memorize(self, state, action, reward, new_state, dead, rnn_state):
        self.MEMORIES.append((state, action, reward, new_state, dead,
↪rnn_state))
        if len(self.MEMORIES) > self.MEMORY_SIZE:
            self.MEMORIES.popleft()

    def _construct_memories(self, replay):
        states = np.array([a[0] for a in replay])
        new_states = np.array([a[3] for a in replay])
        init_values = np.array([a[-1] for a in replay])
        Q = self.sess.run(self.logits, feed_dict={self.X:states, self.
↪hidden_layer:init_values})
        Q_new = self.sess.run(self.logits, feed_dict={self.X:new_states, self.
↪hidden_layer:init_values})
        replay_size = len(replay)
        X = np.empty((replay_size, 4, self.state_size))
        Y = np.empty((replay_size, self.OUTPUT_SIZE))
        INIT_VAL = np.empty((replay_size, 2 * self.LAYER_SIZE))
        for i in range(replay_size):

```

```

        state_r, action_r, reward_r, new_state_r, dead_r, rnn_memory =
→replay[i]
        target = Q[i]
        target[action_r] = reward_r
        if not dead_r:
            target[action_r] += self.GAMMA * np.amax(Q_new[i])
        X[i] = state_r
        Y[i] = target
        INIT_VAL[i] = rnn_memory
    return X, Y, INIT_VAL

    def get_state(self, t):
        window_size = self.window_size + 1
        d = t - window_size + 1
        block = self.trend[d : t + 1] if d >= 0 else -d * [self.trend[0]] +
→self.trend[0 : t + 1]
        res = []
        for i in range(window_size - 1):
            res.append(block[i + 1] - block[i])
        return np.array(res)

    def buy(self, initial_money):
        starting_money = initial_money
        states_sell = []
        states_buy = []
        inventory = []
        state = self.get_state(0)
        init_value = np.zeros((1, 2 * self.LAYER_SIZE))
        for k in range(self.INITIAL_FEATURES.shape[0]):
            self.INITIAL_FEATURES[k,:] = state
        for t in range(0, len(self.trend) - 1, self.skip):
            action, last_state = self.sess.run([self.logits, self.last_state],
                                                feed_dict={self.X: [self.
→INITIAL_FEATURES],
                                                         self.hidden_layer:
→init_value})
            action, init_value = np.argmax(action[0]), last_state
            next_state = self.get_state(t + 1)

            if action == 1 and initial_money >= self.trend[t]:
                inventory.append(self.trend[t])
                initial_money -= self.trend[t]
                states_buy.append(t)
                print('day %d: buy 1 unit at price %f, total balance %f' % (t,
→self.trend[t], initial_money))

            elif action == 2 and len(inventory):

```

```

        bought_price = inventory.pop(0)
        initial_money += self.trend[t]
        states_sell.append(t)
        try:
            invest = ((close[t] - bought_price) / bought_price) * 100
        except:
            invest = 0
        print(
            'day %d, sell 1 unit at price %f, investment %f %%, total_
↪balance %f,'
            % (t, close[t], invest, initial_money)
        )

        new_state = np.append([self.get_state(t + 1)], self.
↪INITIAL_FEATURES[:3, :], axis = 0)
        self.INITIAL_FEATURES = new_state
        invest = ((initial_money - starting_money) / starting_money) * 100
        total_gains = initial_money - starting_money
        return states_buy, states_sell, total_gains, invest

    def train(self, iterations, checkpoint, initial_money):
        for i in range(iterations):
            total_profit = 0
            inventory = []
            state = self.get_state(0)
            starting_money = initial_money
            init_value = np.zeros((1, 2 * self.LAYER_SIZE))
            for k in range(self.INITIAL_FEATURES.shape[0]):
                self.INITIAL_FEATURES[k, :] = state
            for t in range(0, len(self.trend) - 1, self.skip):

                if np.random.rand() < self.EPSILON:
                    action = np.random.randint(self.OUTPUT_SIZE)
                else:
                    action, last_state = self.sess.run([self.logits,
                                                         self.last_state],
                                                         feed_dict={self.X: [self.
↪INITIAL_FEATURES],
                                                         self.hidden_layer:
↪init_value})

                    action, init_value = np.argmax(action[0]), last_state

                next_state = self.get_state(t + 1)

                if action == 1 and starting_money >= self.trend[t]:
                    inventory.append(self.trend[t])

```

```

        starting_money -= self.trend[t]

    elif action == 2 and len(inventory) > 0:
        bought_price = inventory.pop(0)
        total_profit += self.trend[t] - bought_price
        starting_money += self.trend[t]

    invest = ((starting_money - initial_money) / initial_money)
    new_state = np.append([self.get_state(t + 1)], self.
→ INITIAL_FEATURES[:3, :], axis = 0)
    self._memorize(self.INITIAL_FEATURES, action, invest, new_state,
                    starting_money < initial_money, init_value[0])
    self.INITIAL_FEATURES = new_state
    batch_size = min(len(self.MEMORIES), self.BATCH_SIZE)
    replay = random.sample(self.MEMORIES, batch_size)
    X, Y, INIT_VAL = self._construct_memories(replay)

    cost, _ = self.sess.run([self.cost, self.optimizer],
                             feed_dict={self.X: X, self.Y: Y,
                                           self.hidden_layer: INIT_VAL})
    self.EPSILON = self.MIN_EPSILON + (1.0 - self.MIN_EPSILON) * np.
→ exp(-self.DECAY_RATE * i)

    if (i+1) % checkpoint == 0:
        print('epoch: %d, total rewards: %f.3, cost: %f, total money: %
→ %f'%(i + 1, total_profit, cost,

→ starting_money))

```

```

[4]: close = df.Close.values.tolist()
initial_money = 10000
window_size = 30
skip = 1
batch_size = 32
agent = Agent(state_size = window_size,
               window_size = window_size,
               trend = close,
               skip = skip)
agent.train(iterations = 200, checkpoint = 10, initial_money = initial_money)

```

WARNING:tensorflow:<tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x7fef003b2d30>: Using a concatenated state is slower and will soon be deprecated. Use state_is_tuple=True.

```

epoch: 10, total rewards: 449.400388.3, cost: 0.117951, total money: 7420.680355
epoch: 20, total rewards: 513.109983.3, cost: 0.187314, total money: 7552.130003
epoch: 30, total rewards: 1755.114813.3, cost: 0.337607, total money:
6759.834784

```

```

epoch: 40, total rewards: 545.719909.3, cost: 0.555657, total money: 9529.079894
epoch: 50, total rewards: 593.435182.3, cost: 0.399239, total money: 6611.165162
epoch: 60, total rewards: 285.174678.3, cost: 0.071772, total money: 6314.564631
epoch: 70, total rewards: 169.200014.3, cost: 0.796504, total money: 4264.030030
epoch: 80, total rewards: 520.019840.3, cost: 0.567794, total money: 6501.959842
epoch: 90, total rewards: 498.320189.3, cost: 0.245750, total money: 9481.210204
epoch: 100, total rewards: 1572.605044.3, cost: 1.142984, total money:
11572.605044
epoch: 110, total rewards: 297.584960.3, cost: 0.973414, total money:
10297.584960
epoch: 120, total rewards: 912.394901.3, cost: 2.032860, total money:
6987.034854
epoch: 130, total rewards: 22.109988.3, cost: 0.097879, total money:
10022.109988
epoch: 140, total rewards: 471.779909.3, cost: 0.532008, total money:
10471.779909
epoch: 150, total rewards: 215.255126.3, cost: 0.236825, total money:
10215.255126
epoch: 160, total rewards: 147.780093.3, cost: 0.432537, total money:
9174.450076
epoch: 170, total rewards: 203.309817.3, cost: 0.413111, total money:
10203.309817
epoch: 180, total rewards: 76.350403.3, cost: 0.132205, total money: 8084.520385
epoch: 190, total rewards: 173.749880.3, cost: 1.325852, total money:
10173.749880
epoch: 200, total rewards: 4.325196.3, cost: 0.500293, total money: 8987.685181

```

```

[5]: states_buy, states_sell, total_gains, invest = agent.buy(initial_money =
    ↪initial_money)

```

```

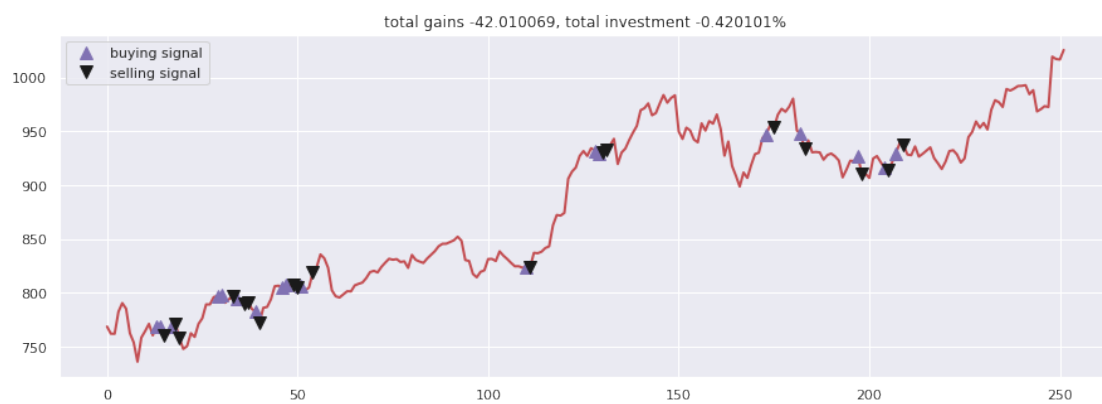
day 13: buy 1 unit at price 769.200012, total balance 9230.799988
day 14: buy 1 unit at price 768.270020, total balance 8462.529968
day 15, sell 1 unit at price 760.989990, investment -1.067346 %, total balance
9223.519958,
day 17: buy 1 unit at price 768.239990, total balance 8455.279968
day 18, sell 1 unit at price 770.840027, investment 0.334519 %, total balance
9226.119995,
day 19, sell 1 unit at price 758.039978, investment -1.327712 %, total balance
9984.159973,
day 29: buy 1 unit at price 797.070007, total balance 9187.089966
day 30: buy 1 unit at price 797.849976, total balance 8389.239990
day 33, sell 1 unit at price 796.419983, investment -0.081552 %, total balance
9185.659973,
day 34: buy 1 unit at price 794.559998, total balance 8391.099975
day 36, sell 1 unit at price 789.909973, investment -0.995175 %, total balance
9181.009948,
day 37, sell 1 unit at price 791.549988, investment -0.378827 %, total balance
9972.559936,

```

day 39: buy 1 unit at price 782.789978, total balance 9189.769958
 day 40, sell 1 unit at price 771.820007, investment -1.401394 %, total balance 9961.589965,
 day 46: buy 1 unit at price 804.789978, total balance 9156.799987
 day 47: buy 1 unit at price 807.909973, total balance 8348.890014
 day 49, sell 1 unit at price 807.880005, investment 0.383954 %, total balance 9156.770019,
 day 50, sell 1 unit at price 804.609985, investment -0.408460 %, total balance 9961.380004,
 day 51: buy 1 unit at price 806.070007, total balance 9155.309997
 day 54, sell 1 unit at price 819.309998, investment 1.642536 %, total balance 9974.619995,
 day 110: buy 1 unit at price 824.320007, total balance 9150.299988
 day 111, sell 1 unit at price 823.559998, investment -0.092198 %, total balance 9973.859986,
 day 128: buy 1 unit at price 932.169983, total balance 9041.690003
 day 129: buy 1 unit at price 928.780029, total balance 8112.909974
 day 130, sell 1 unit at price 930.599976, investment -0.168425 %, total balance 9043.509950,
 day 131, sell 1 unit at price 932.219971, investment 0.370372 %, total balance 9975.729921,
 day 173: buy 1 unit at price 947.159973, total balance 9028.569948
 day 175, sell 1 unit at price 953.419983, investment 0.660924 %, total balance 9981.989931,
 day 182: buy 1 unit at price 947.799988, total balance 9034.189943
 day 183, sell 1 unit at price 934.090027, investment -1.446504 %, total balance 9968.279970,
 day 197: buy 1 unit at price 926.960022, total balance 9041.319948
 day 198, sell 1 unit at price 910.979980, investment -1.723919 %, total balance 9952.299928,
 day 204: buy 1 unit at price 915.890015, total balance 9036.409913
 day 205, sell 1 unit at price 913.809998, investment -0.227103 %, total balance 9950.219911,
 day 207: buy 1 unit at price 929.570007, total balance 9020.649904
 day 209, sell 1 unit at price 937.340027, investment 0.835872 %, total balance 9957.989931,

```

[6]: fig = plt.figure(figsize = (15,5))
plt.plot(close, color='r', lw=2.)
plt.plot(close, '^', markersize=10, color='m', label = 'buying signal',
    ↳markevery = states_buy)
plt.plot(close, 'v', markersize=10, color='k', label = 'selling signal',
    ↳markevery = states_sell)
plt.title('total gains %f, total investment %f%%'%(total_gains, invest))
plt.legend()
plt.show()
  
```



[]: