

# Deep Reinforcement Learning on Stock Data

## Reinforcement Learning

Reinforcement Learning (RL) is a type of machine learning paradigm that involves an agent learning to make decisions by interacting with an environment.

The agent learns to achieve a goal or maximize a reward by taking suitable actions in different situations.

It's based on the idea of learning from feedback, wherein the agent receives rewards or penalties based on its actions.

## Import Libraries

```
In [2]: import time
import copy
import numpy as np
import pandas as pd
import chainer
import chainer.functions as F
import chainer.links as L
from plotly import tools
from plotly.graph_objs import *
from plotly.offline import init_notebook_mode, iplot, iplot_mpl
init_notebook_mode()
```

## Read the data

```
In [3]: data = pd.read_csv('../input/Data/Stocks/goog.us.txt')
data['Date'] = pd.to_datetime(data['Date'])
data = data.set_index('Date')
print(data.index.min(), data.index.max())
data.head()
```

2014-03-27 00:00:00 2017-11-10 00:00:00

```
Out[3]:      Open   High    Low   Close  Volume  OpenInt
               Date
2014-03-27  568.00  568.00  552.92  558.46   13052      0
2014-03-28  561.20  566.43  558.67  559.99   41003      0
2014-03-31  566.89  567.00  556.93  556.97   10772      0
2014-04-01  558.71  568.45  558.71  567.16    7932      0
2014-04-02  599.99  604.83  562.19  567.00  146697      0
```

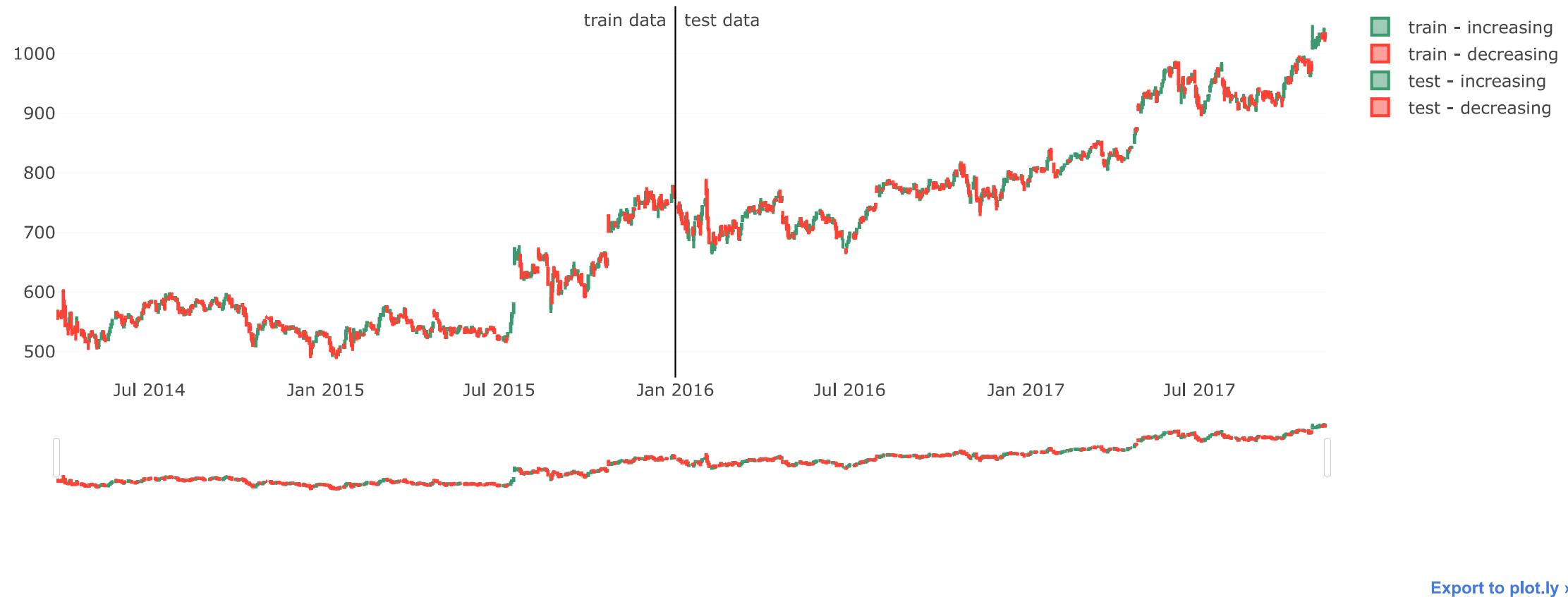
## Train and Test Datasets

```
In [4]: date_split = '2016-01-01'  
train = data[:date_split]  
test = data[date_split:]  
len(train), len(test)
```

```
Out[4]: (446, 470)
```

```
In [5]: def plot_train_test(train, test, date_split):  
  
    data = [  
        Candlestick(x=train.index, open=train['Open'], high=train['High'], low=train['Low'], close=train['Close'], name='train'),  
        Candlestick(x=test.index, open=test['Open'], high=test['High'], low=test['Low'], close=test['Close'], name='test')  
    ]  
    layout = {  
        'shapes': [  
            {'x0': date_split, 'x1': date_split, 'y0': 0, 'y1': 1, 'xref': 'x', 'yref': 'paper', 'line': {'color': 'rgb(0,0,0)', 'width': 1}}  
        ],  
        'annotations': [  
            {'x': date_split, 'y': 1.0, 'xref': 'x', 'yref': 'paper', 'showarrow': False, 'xanchor': 'left', 'text': ' test data'},  
            {'x': date_split, 'y': 1.0, 'xref': 'x', 'yref': 'paper', 'showarrow': False, 'xanchor': 'right', 'text': 'train data '}  
        ]  
    }  
    figure = Figure(data=data, layout=layout)  
    iplot(figure)
```

```
In [6]: plot_train_test(train, test, date_split)
```



In [7]:

```
class Environment1:

    def __init__(self, data, history_t=90):
        self.data = data
        self.history_t = history_t
        self.reset()

    def reset(self):
        self.t = 0
        self.done = False
        self.profits = 0
        self.positions = []
        self.position_value = 0
        self.history = [0 for _ in range(self.history_t)]
        return [self.position_value] + self.history # obs

    def step(self, act):
        reward = 0

        # act = 0: stay, 1: buy, 2: sell
        if act == 1:
            self.positions.append(self.data.iloc[self.t, :]['Close'])
        elif act == 2: # sell
            if len(self.positions) == 0:
                reward = -1
            else:
                profits = 0
```

```
for p in self.positions:
    profits += (self.data.iloc[self.t, :]['Close'] - p)
reward += profits
self.profits += profits
self.positions = []

# set next time
self.t += 1
self.position_value = 0
for p in self.positions:
    self.position_value += (self.data.iloc[self.t, :]['Close'] - p)
self.history.pop(0)
self.history.append(self.data.iloc[self.t, :]['Close'] - self.data.iloc[(self.t-1), :]['Close'])

# clipping reward
if reward > 0:
    reward = 1
elif reward < 0:
    reward = -1

return [self.position_value] + self.history, reward, self.done # obs, reward, done
```

```
In [9]: # DQN

def train_dqn(env):

    class Q_Network(chainer.Chain):

        def __init__(self, input_size, hidden_size, output_size):
            super(Q_Network, self).__init__()
            fc1 = L.Linear(input_size, hidden_size),
            fc2 = L.Linear(hidden_size, hidden_size),
            fc3 = L.Linear(hidden_size, output_size)
        )

        def __call__(self, x):
            h = F.relu(self.fc1(x))
            h = F.relu(self.fc2(h))
            y = self.fc3(h)
            return y

        def reset(self):
            self.zerograds()

    Q = Q_Network(input_size=env.history_t+1, hidden_size=100, output_size=3)
```

```

Q_ast = copy.deepcopy(Q)
optimizer = chainer.optimizers.Adam()
optimizer.setup(Q)

epoch_num = 50
step_max = len(env.data)-1
memory_size = 200
batch_size = 20
epsilon = 1.0
epsilon_decrease = 1e-3
epsilon_min = 0.1
start_reduce_epsilon = 200
train_freq = 10
update_q_freq = 20
gamma = 0.97
show_log_freq = 5

memory = []
total_step = 0
total_rewards = []
total_losses = []

start = time.time()
for epoch in range(epoch_num):

    pobs = env.reset()
    step = 0
    done = False
    total_reward = 0
    total_loss = 0

    while not done and step < step_max:

        # select act
        pact = np.random.randint(3)
        if np.random.rand() > epsilon:
            pact = Q(np.array(pobs, dtype=np.float32).reshape(1, -1))
            pact = np.argmax(pact.data)

        # act
        obs, reward, done = env.step(pact)

        # add memory
        memory.append((pobs, pact, reward, obs, done))
        if len(memory) > memory_size:
            memory.pop(0)

        # train or update q
        if len(memory) == memory_size:
            if total_step % train_freq == 0:
                shuffled_memory = np.random.permutation(memory)
                memory_idx = range(len(shuffled_memory))
                for i in memory_idx[::batch_size]:
                    batch = np.array(shuffled_memory[i:i+batch_size])
                    b_pobs = np.array(batch[:, 0].tolist(), dtype=np.float32).reshape(batch_size, -1)
                    b_pact = np.array(batch[:, 1].tolist(), dtype=np.int32)
                    b_reward = np.array(batch[:, 2].tolist(), dtype=np.int32)
                    b_obs = np.array(batch[:, 3].tolist(), dtype=np.float32).reshape(batch_size, -1)
                    b_done = np.array(batch[:, 4].tolist(), dtype=np.bool)

```

```

q = Q(b_pobs)
maxq = np.max(Q.ast(b_obs).data, axis=1)
target = copy.deepcopy(q.data)
for j in range(batch_size):
    target[j, b_pact[j]] = b_reward[j]+gamma*maxq[j]*(not b_done[j])
Q.reset()
loss = F.mean_squared_error(q, target)
total_loss += loss.data
loss.backward()
optimizer.update()

if total_step % update_q_freq == 0:
    Q.ast = copy.deepcopy(Q)

# epsilon
if epsilon > epsilon_min and total_step > start_reduce_epsilon:
    epsilon -= epsilon_decrease

# next step
total_reward += reward
pobs = obs
step += 1
total_step += 1

total_rewards.append(total_reward)
total_losses.append(total_loss)

if (epoch+1) % show_log_freq == 0:
    log_reward = sum(total_rewards[((epoch+1)-show_log_freq):])/show_log_freq
    log_loss = sum(total_losses[((epoch+1)-show_log_freq):])/show_log_freq
    elapsed_time = time.time()-start
    print('\t'.join(map(str, [epoch+1, epsilon, total_step, log_reward, log_loss, elapsed_time])))
    start = time.time()

return Q, total_losses, total_rewards

```

In [10]: `Q, total_losses, total_rewards = train_dqn(Environment1(train))`

5	0.0999999999999992	2225	-62.0	38994.727623	9.185798168182373
10	0.0999999999999992	4450	-69.6	5872.36153181	10.086313247680664
15	0.0999999999999992	6675	-71.6	490.885859407	9.929555654525757
20	0.0999999999999992	8900	-55.0	228.797397427	9.453218221664429
25	0.0999999999999992	11125	-25.6	182.087579884	9.757601499557495
30	0.0999999999999992	13350	-5.8	337.967664644	10.504354476928711
35	0.0999999999999992	15575	-15.4	1159.82199649	12.502948999404907
40	0.0999999999999992	17800	-9.4	57.2816223043	10.434857845306396
45	0.0999999999999992	20025	8.6	28.5512007653	10.884730577468872
50	0.0999999999999992	22250	7.2	22.5572660003	10.777376413345337

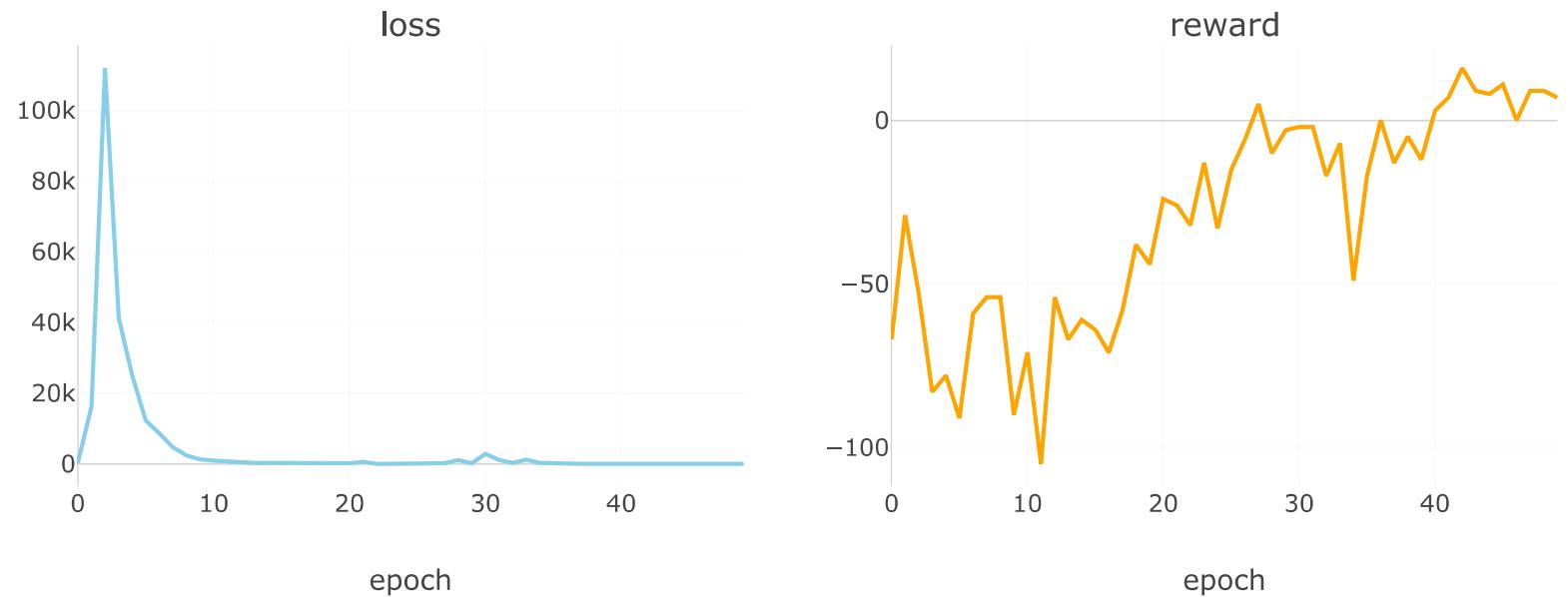
In [11]: `def plot_loss_reward(total_losses, total_rewards):`

```

figure = tools.make_subplots(rows=1, cols=2, subplot_titles=('loss', 'reward'), print_grid=False)
figure.append_trace(Scatter(y=total_losses, mode='lines', line=dict(color='skyblue')), 1, 1)
figure.append_trace(Scatter(y=total_rewards, mode='lines', line=dict(color='orange')), 1, 2)
figure['layout']['xaxis1'].update(title='epoch')
figure['layout']['xaxis2'].update(title='epoch')
figure['layout'].update(height=400, width=900, showlegend=False)
iplot(figure)

```

```
In [12]: plot_loss_reward(total_losses, total_rewards)
```



[Export to plot.ly »](#)

```
In [13]: def plot_train_test_by_q(train_env, test_env, Q, algorithm_name):
```

```
# train
pobs = train_env.reset()
train_acts = []
train_rewards = []

for _ in range(len(train_env.data)-1):

    pact = Q(np.array(pobs, dtype=np.float32).reshape(1, -1))
    pact = np.argmax(pact.data)
    train_acts.append(pact)

    obs, reward, done = train_env.step(pact)
    train_rewards.append(reward)

    pobs = obs

train_profits = train_env.profits

# test
pobs = test_env.reset()
test_acts = []
test_rewards = []

for _ in range(len(test_env.data)-1):

    pact = Q(np.array(pobs, dtype=np.float32).reshape(1, -1))
    pact = np.argmax(pact.data)
    test_acts.append(pact)
```

```

obs, reward, done = test_env.step(pact)
test_rewards.append(reward)

pobs = obs

test_profits = test_env.profits

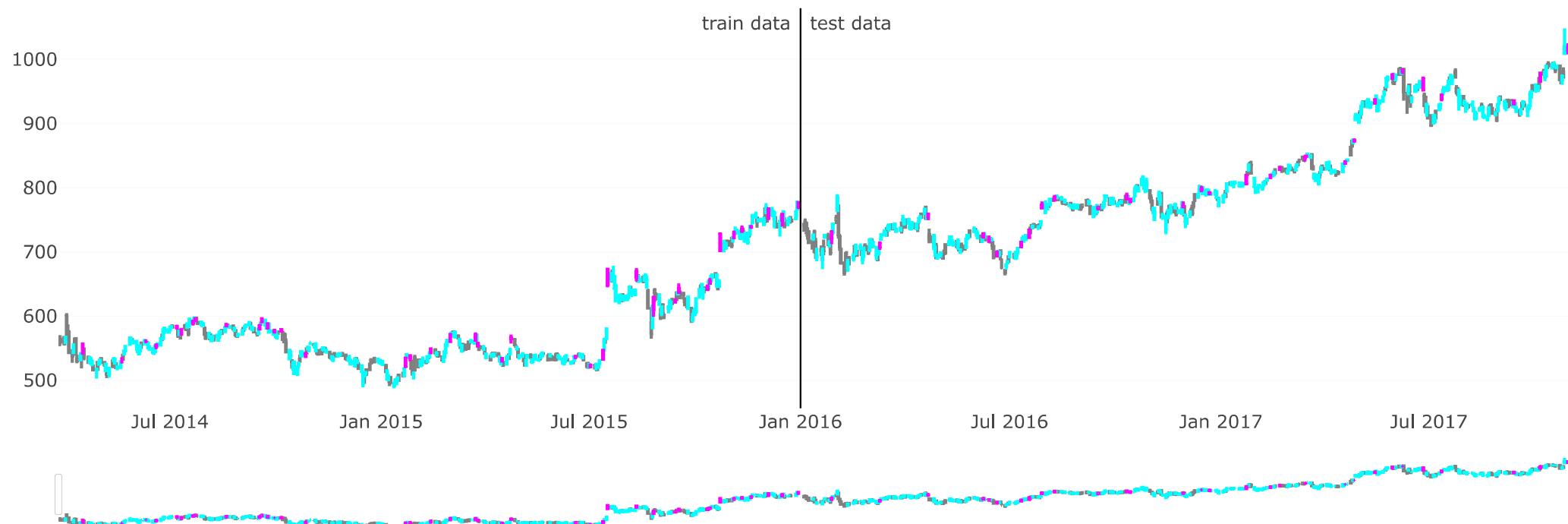
# plot
train_copy = train_env.data.copy()
test_copy = test_env.data.copy()
train_copy['act'] = train_acts + [np.nan]
train_copy['reward'] = train_rewards + [np.nan]
test_copy['act'] = test_acts + [np.nan]
test_copy['reward'] = test_rewards + [np.nan]
train0 = train_copy[train_copy['act'] == 0]
train1 = train_copy[train_copy['act'] == 1]
train2 = train_copy[train_copy['act'] == 2]
test0 = test_copy[test_copy['act'] == 0]
test1 = test_copy[test_copy['act'] == 1]
test2 = test_copy[test_copy['act'] == 2]
act_color0, act_color1, act_color2 = 'gray', 'cyan', 'magenta'

data = [
    Candlestick(x=train0.index, open=train0['Open'], high=train0['High'], low=train0['Low'], close=train0['Close'], increasing=dict(line=dict(color=act_color0)), decreasing=dict(line=dict(color=act_color0))),
    Candlestick(x=train1.index, open=train1['Open'], high=train1['High'], low=train1['Low'], close=train1['Close'], increasing=dict(line=dict(color=act_color1)), decreasing=dict(line=dict(color=act_color1))),
    Candlestick(x=train2.index, open=train2['Open'], high=train2['High'], low=train2['Low'], close=train2['Close'], increasing=dict(line=dict(color=act_color2)), decreasing=dict(line=dict(color=act_color2))),
    Candlestick(x=test0.index, open=test0['Open'], high=test0['High'], low=test0['Low'], close=test0['Close'], increasing=dict(line=dict(color=act_color0)), decreasing=dict(line=dict(color=act_color0))),
    Candlestick(x=test1.index, open=test1['Open'], high=test1['High'], low=test1['Low'], close=test1['Close'], increasing=dict(line=dict(color=act_color1)), decreasing=dict(line=dict(color=act_color1))),
    Candlestick(x=test2.index, open=test2['Open'], high=test2['High'], low=test2['Low'], close=test2['Close'], increasing=dict(line=dict(color=act_color2)), decreasing=dict(line=dict(color=act_color2)))
]
title = '{}: train s-reward {}, profits {}, test s-reward {}, profits {}'.format(
    algorithm_name,
    int(sum(train_rewards)),
    int(train_profits),
    int(sum(test_rewards)),
    int(test_profits)
)
layout = {
    'title': title,
    'showlegend': False,
    'shapes': [
        {'x0': date_split, 'x1': date_split, 'y0': 0, 'y1': 1, 'xref': 'x', 'yref': 'paper', 'line': {'color': 'rgb(0,0,0)', 'width': 1}}
    ],
    'annotations': [
        {'x': date_split, 'y': 1.0, 'xref': 'x', 'yref': 'paper', 'showarrow': False, 'xanchor': 'left', 'text': 'test data'},
        {'x': date_split, 'y': 1.0, 'xref': 'x', 'yref': 'paper', 'showarrow': False, 'xanchor': 'right', 'text': 'train data'}
    ]
}
figure = Figure(data=data, layout=layout)
iplot(figure)

```

In [14]: `plot_train_test_by_q(Environment1(train), Environment1(test), Q, 'DQN')`

DQN: train s-reward 28, profits 2248, test s-reward 14, profits 2831



[Export to plot.ly »](#)

In [15]: # Double DQN

```
def train_ddqn(env):

    class Q_Network(chainer.Chain):

        def __init__(self, input_size, hidden_size, output_size):
            super(Q_Network, self).__init__(
                fc1 = L.Linear(input_size, hidden_size),
                fc2 = L.Linear(hidden_size, hidden_size),
                fc3 = L.Linear(hidden_size, output_size)
            )

        def __call__(self, x):
            h = F.relu(self.fc1(x))
            h = F.relu(self.fc2(h))
            y = self.fc3(h)
            return y

        def reset(self):
            self.zerograds()

    Q = Q_Network(input_size=env.history_t+1, hidden_size=100, output_size=3)
    Q_ast = copy.deepcopy(Q)
    optimizer = chainer.optimizers.Adam()
    optimizer.setup(Q)
```

```

epoch_num = 50
step_max = len(env.data)-1
memory_size = 200
batch_size = 50
epsilon = 1.0
epsilon_decrease = 1e-3
epsilon_min = 0.1
start_reduce_epsilon = 200
train_freq = 10
update_q_freq = 20
gamma = 0.97
show_log_freq = 5

memory = []
total_step = 0
total_rewards = []
total_losses = []

start = time.time()
for epoch in range(epoch_num):

    pobs = env.reset()
    step = 0
    done = False
    total_reward = 0
    total_loss = 0

    while not done and step < step_max:

        # select act
        pact = np.random.randint(3)
        if np.random.rand() > epsilon:
            pact = Q(np.array(pobs, dtype=np.float32).reshape(1, -1))
            pact = np.argmax(pact.data)

        # act
        obs, reward, done = env.step(pact)

        # add memory
        memory.append((pobs, pact, reward, obs, done))
        if len(memory) > memory_size:
            memory.pop(0)

        # train or update q
        if len(memory) == memory_size:
            if total_step % train_freq == 0:
                shuffled_memory = np.random.permutation(memory)
                memory_idx = range(len(shuffled_memory))
                for i in memory_idx[::batch_size]:
                    batch = np.array(shuffled_memory[i:i+batch_size])
                    b_pobs = np.array(batch[:, 0].tolist(), dtype=np.float32).reshape(batch_size, -1)
                    b_pact = np.array(batch[:, 1].tolist(), dtype=np.int32)
                    b_reward = np.array(batch[:, 2].tolist(), dtype=np.int32)
                    b_obs = np.array(batch[:, 3].tolist(), dtype=np.float32).reshape(batch_size, -1)
                    b_done = np.array(batch[:, 4].tolist(), dtype=np.bool)

                    q = Q(b_pobs)
                    """ <<< DQN -> Double DQN
                    maxq = np.max(Q.ast(b_obs).data, axis=1)
                    === """

```

```

        indices = np.argmax(q.data, axis=1)
        maxqs = Q.ast(b_obs).data
        """ >>> """
        target = copy.deepcopy(q.data)
        for j in range(batch_size):
            """ <<< DQN -> Double DQN
            target[j, b_pact[j]] = b_reward[j]+gamma*maxq[j]*(not b_done[j])
            === """
            target[j, b_pact[j]] = b_reward[j]+gamma*maxqs[j, indices[j]]*(not b_done[j])
            """ >>> """
        Q.reset()
        loss = F.mean_squared_error(q, target)
        total_loss += loss.data
        loss.backward()
        optimizer.update()

    if total_step % update_q_freq == 0:
        Q.ast = copy.deepcopy(Q)

    # epsilon
    if epsilon > epsilon_min and total_step > start_reduce_epsilon:
        epsilon -= epsilon_decrease

    # next step
    total_reward += reward
    pobs = obs
    step += 1
    total_step += 1

    total_rewards.append(total_reward)
    total_losses.append(total_loss)

    if (epoch+1) % show_log_freq == 0:
        log_reward = sum(total_rewards[((epoch+1)-show_log_freq):])/show_log_freq
        log_loss = sum(total_losses[((epoch+1)-show_log_freq):])/show_log_freq
        elapsed_time = time.time()-start
        print('\t'.join(map(str, [epoch+1, epsilon, total_step, log_reward, log_loss, elapsed_time])))
        start = time.time()

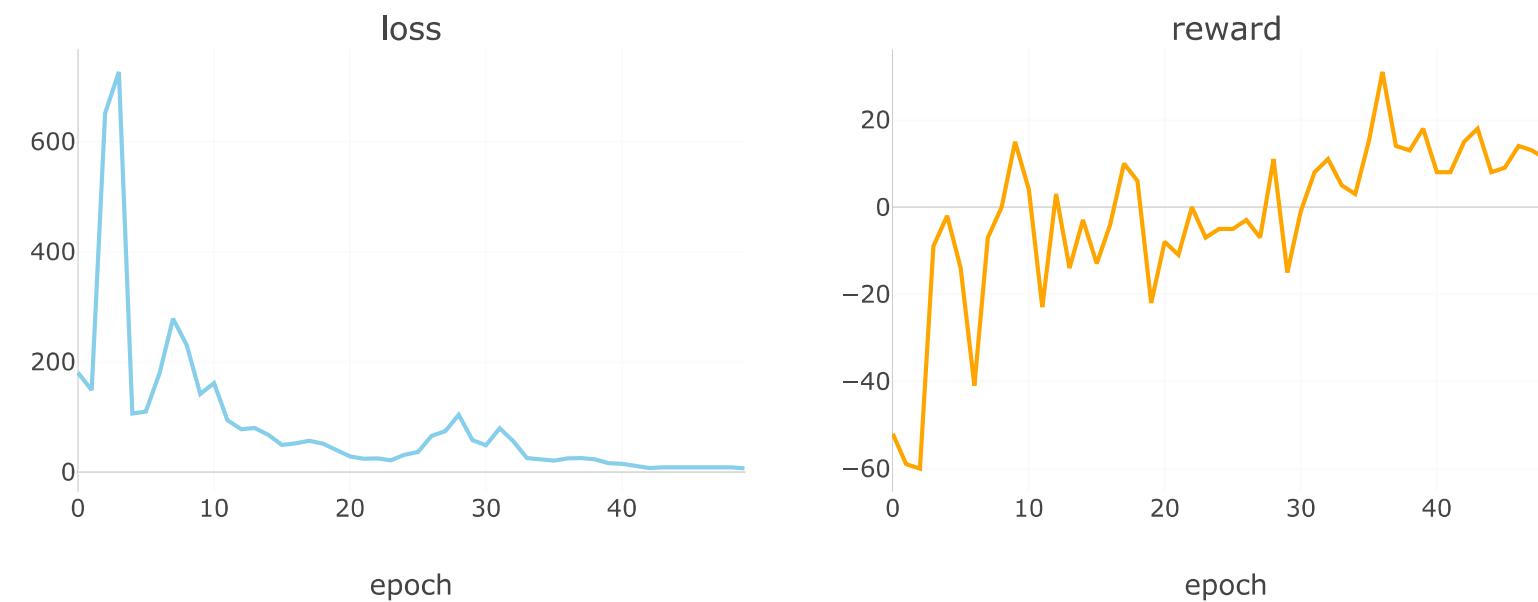
return Q, total_losses, total_rewards

```

In [16]: `Q, total_losses, total_rewards = train_ddqn(Environment1(train))`

5	0.0999999999999992	2225	-36.4	362.532034132	5.976006746292114
10	0.0999999999999992	4450	-9.4	188.051366159	7.1547322273254395
15	0.0999999999999992	6675	-6.6	96.137116237	6.8685243129730225
20	0.0999999999999992	8900	-4.6	50.2297730058	6.529151439666748
25	0.0999999999999992	11125	-6.2	25.7622961797	6.623997926712036
30	0.0999999999999992	13350	-3.8	67.7237939358	7.270005941390991
35	0.0999999999999992	15575	5.2	46.4835077588	7.701663970947266
40	0.0999999999999992	17800	18.2	21.8936307523	7.602718114852905
45	0.0999999999999992	20025	11.4	9.97410044037	7.157979726791382
50	0.0999999999999992	22250	10.2	7.85483987899	7.09658408164978

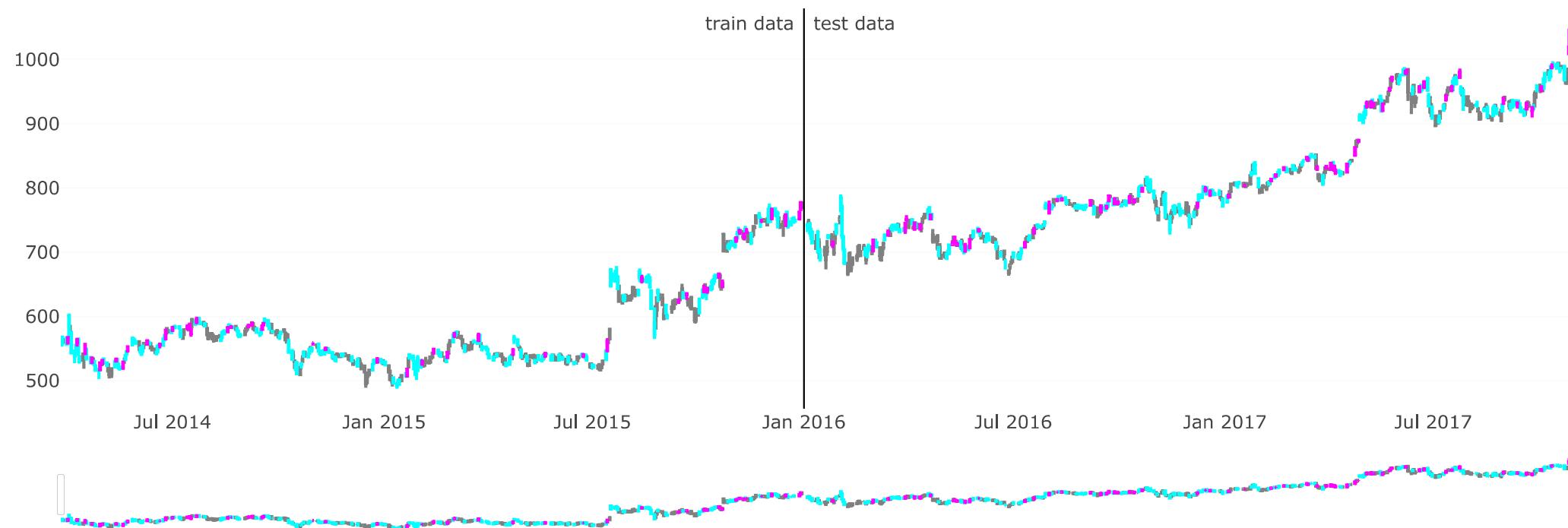
In [17]: `plot_loss_reward(total_losses, total_rewards)`



[Export to plot.ly »](#)

In [18]: `plot_train_test_by_q(Environment1(train), Environment1(test), Q, 'Double DQN')`

Double DQN: train s-reward 24, profits 1084, test s-reward 12, profits 1297



[Export to plot.ly »](#)

In [19]: # Dueling Double DQN

```
def train_ddqn(env):

    """ <<< Double DQN -> Dueling Double DQN
    class Q_Network(chainer.Chain):

        def __init__(self, input_size, hidden_size, output_size):
            super(Q_Network, self).__init__(
                fc1 = L.Linear(input_size, hidden_size),
                fc2 = L.Linear(hidden_size, hidden_size),
                fc3 = L.Linear(hidden_size, output_size)
            )

        def __call__(self, x):
            h = F.relu(self.fc1(x))
            h = F.relu(self.fc2(h))
            y = self.fc3(h)
            return y

        def reset(self):
            self.zerograds()
    === """
    class Q_Network(chainer.Chain):

        def __init__(self, input_size, hidden_size, output_size):
            super(Q_Network, self).__init__(
                fc1 = L.Linear(input_size, hidden_size),
                fc2 = L.Linear(hidden_size, hidden_size),
                fc3 = L.Linear(hidden_size, hidden_size//2),
                fc4 = L.Linear(hidden_size, hidden_size//2),
                state_value = L.Linear(hidden_size//2, 1),
                advantage_value = L.Linear(hidden_size//2, output_size)
            )
            self.input_size = input_size
            self.hidden_size = hidden_size
            self.output_size = output_size

        def __call__(self, x):
            h = F.relu(self.fc1(x))
            h = F.relu(self.fc2(h))
            hs = F.relu(self.fc3(h))
            ha = F.relu(self.fc4(h))
            state_value = self.state_value(hs)
            advantage_value = self.advantage_value(ha)
            advantage_mean = (F.sum(advantage_value, axis=1)/float(self.output_size)).reshape(-1, 1)
            q_value = F.concat([state_value for _ in range(self.output_size)], axis=1) + (advantage_value - F.concat([advantage_mean for _ in range(self.output_size)], axis=1))
            return q_value

        def reset(self):
            self.zerograds()
    === >>> """

Q = Q_Network(input_size=env.history_t+1, hidden_size=100, output_size=3)
Q_ast = copy.deepcopy(Q)
optimizer = chainer.optimizers.Adam()
optimizer.setup(Q)

epoch_num = 50
```



```

maxqs = Q.ast(b_obs).data
"""
>>> """
target = copy.deepcopy(q.data)
for j in range(batch_size):
    """ <<< DQN -> Double DQN
    target[j, b_pact[j]] = b_reward[j]+gamma*maxq[j]*(not b_done[j])
    === """
    target[j, b_pact[j]] = b_reward[j]+gamma*maxqs[j, indices[j]]*(not b_done[j])
    """
>>> """
Q.reset()
loss = F.mean_squared_error(q, target)
total_loss += loss.data
loss.backward()
optimizer.update()

if total_step % update_q_freq == 0:
    Q_ast = copy.deepcopy(Q)

# epsilon
if epsilon > epsilon_min and total_step > start_reduce_epsilon:
    epsilon -= epsilon_decrease

# next step
total_reward += reward
pobs = obs
step += 1
total_step += 1

total_rewards.append(total_reward)
total_losses.append(total_loss)

if (epoch+1) % show_log_freq == 0:
    log_reward = sum(total_rewards[((epoch+1)-show_log_freq):])/show_log_freq
    log_loss = sum(total_losses[((epoch+1)-show_log_freq):])/show_log_freq
    elapsed_time = time.time()-start
    print('\t'.join(map(str, [epoch+1, epsilon, total_step, log_reward, log_loss, elapsed_time])))
    start = time.time()

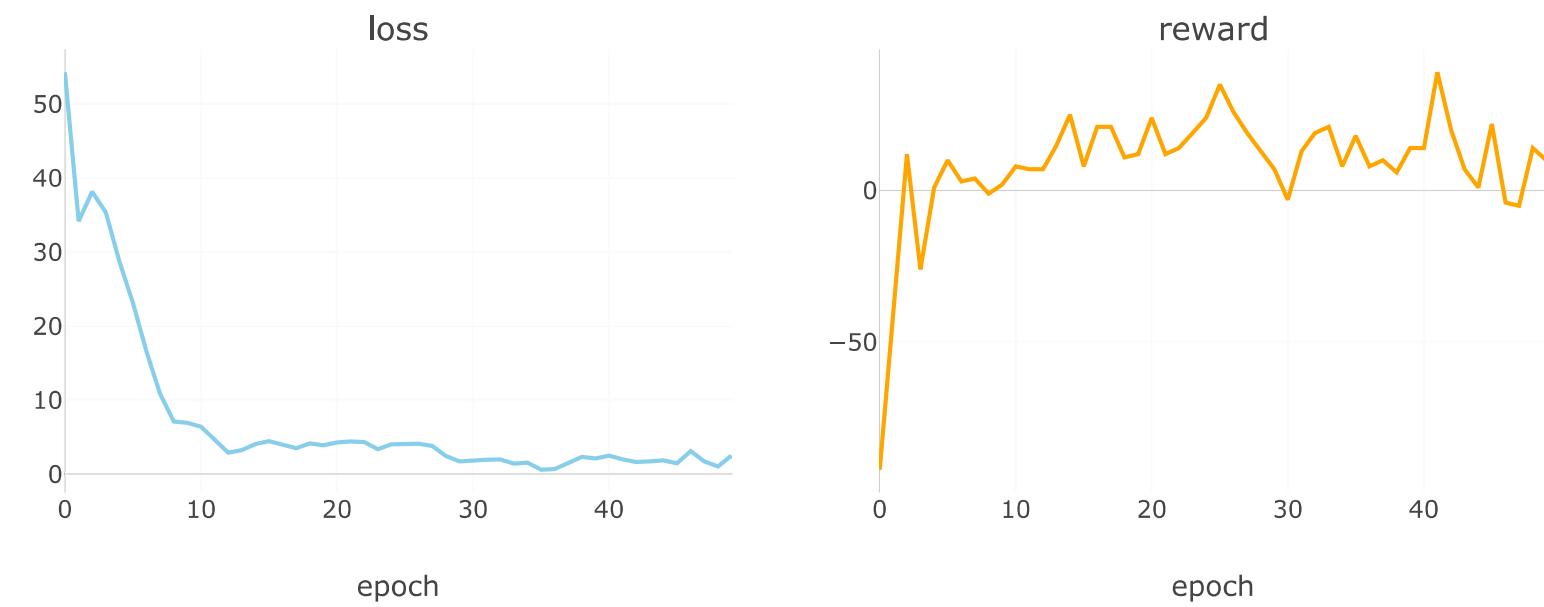
return Q, total_losses, total_rewards

```

In [20]: `Q, total_losses, total_rewards = train_ddqn(Environment1(train))`

5	0.0999999999999992	2225	-29.4	38.1641087133	9.779456615447998
10	0.0999999999999992	4450	3.6	12.950515821	10.579248428344727
15	0.0999999999999992	6675	12.4	4.27871555081	10.769827604293823
20	0.0999999999999992	8900	14.6	4.02654048405	11.711466312408447
25	0.0999999999999992	11125	18.6	4.1195472979	10.950128078460693
30	0.0999999999999992	13350	20.0	3.25838643042	11.056329011917114
35	0.0999999999999992	15575	11.6	1.76036339309	10.970870733261108
40	0.0999999999999992	17800	11.2	1.46561827543	10.56775975227356
45	0.0999999999999992	20025	16.2	1.97831109846	11.008494853973389
50	0.0999999999999992	22250	7.4	1.99162399785	10.179644584655762

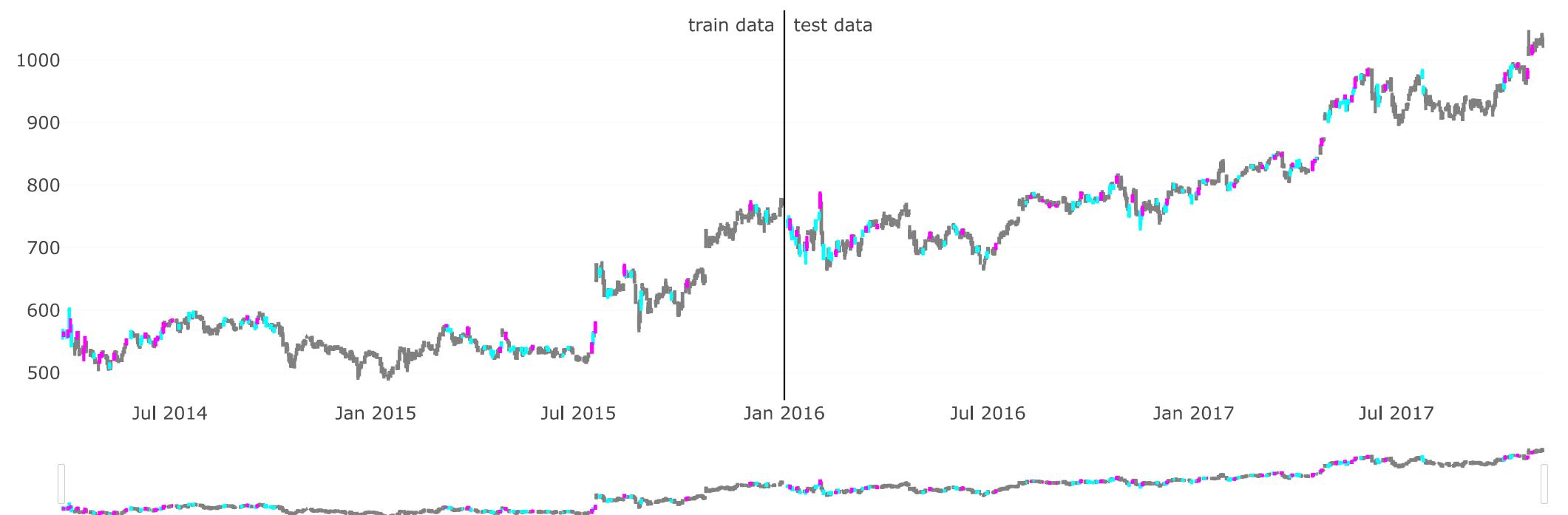
In [21]: `plot_loss_reward(total_losses, total_rewards)`



[Export to plot.ly »](#)

```
In [22]: plot_train_test_by_q(Environment1(train), Environment1(test), Q, 'Dueling Double DQN')
```

Dueling Double DQN: train s-reward 14, profits 331, test s-reward 27, profits 785



[Export to plot.ly »](#)