

---

# Q-Learning: Deep Deterministic Policy Gradient

Zander Zemliak  
CMSI 5350  
December 6, 2022

---

# Overview of Topics

- I. Introduction to Reinforcement Learning (RL)
  - II. RL algorithm: Q-Learning
  - III. Deep Q-Learning
  - IV. DQN Algorithm: Deep Deterministic Policy Gradient (DDPG)
  - V. Stock Trading Agent built Using the DDPG Algorithm
-

---

# Reinforcement Learning

---

---

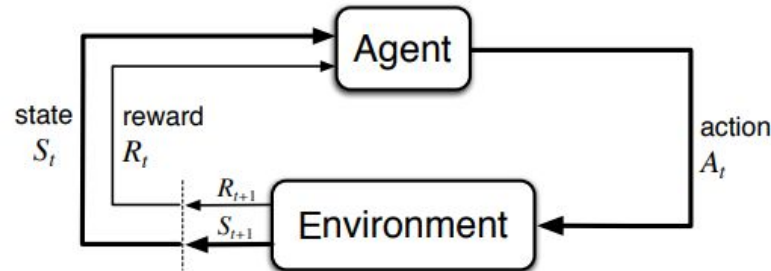
# Introduction to RL

- RL is one of the main disciplines of machine learning.
  - Differences between the three main ML disciplines:
    - Supervised learning is concerned with learning an approximate mapping between an input and an output using a labeled dataset.
    - Unsupervised learning is concerned with finding a pattern in unlabeled data.
    - RL agents learn by interacting with their environment.
  - RL agents main objective is to take actions that maximize reward.
  - A dog learns that when it sits on command it gets a treat.
-

---

# Markov Decision Process (MDP)

- MDP is an essential framework for RL, it is a mathematical model of decision making where the outcomes can be partly arbitrary or under control.
- Four main components for MDP: states, actions, effects of actions on future states, effects of actions on future rewards.



- Agriculture ex: How much to plant given the current state of water and soil?

---

## Brief MDP Math

- An MDP can be described as:  $MDP = (S, A, T, R, \gamma)$ 
  - $S$  = set of possible states
  - $A$  = set of possible actions
  - $T$  = set of transition probabilities where the probability of going from  $s$  to  $s'$  is described as  $P(s'|s, a)$ .
  - $R$  = set of rewards
  - $\gamma$  = discount factor (determines how much agent cares about future rewards and is necessary for algorithm convergence)

- We will use the Bellman equation to evaluate the max current and future reward:

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

- The value for our current state,  $V(s)$ , is equal to the action which maximizes current reward in state  $s$  with action  $a$  plus the value of the next state  $s'$ .
-

---

# Q-Learning

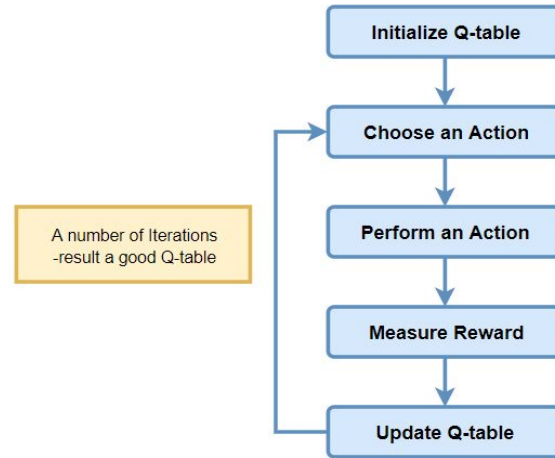
---

# Q-Learning

- Q-learning is a *value based* RL algorithm that attempts to learn the quality (Q) of actions in particular states that are defined by future value.
- The mapping of action to state is known as the policy ( $\pi$ ) and these values are stored in a Q-table with two columns: state-action pair and value.
- Agents choose actions in certain states that result in the highest reward  $Q^*(s,a)$ .
- Q-table values are updated using TD:

$$Q_{\text{new}}(s, a) = Q(s, a) + \alpha [R(s,a) + \gamma \max_a Q(s', a=a) - Q(s, a)]$$

- TD target looks very familiar to Bellman equation...
- This will converge towards  $\pi^*$





---

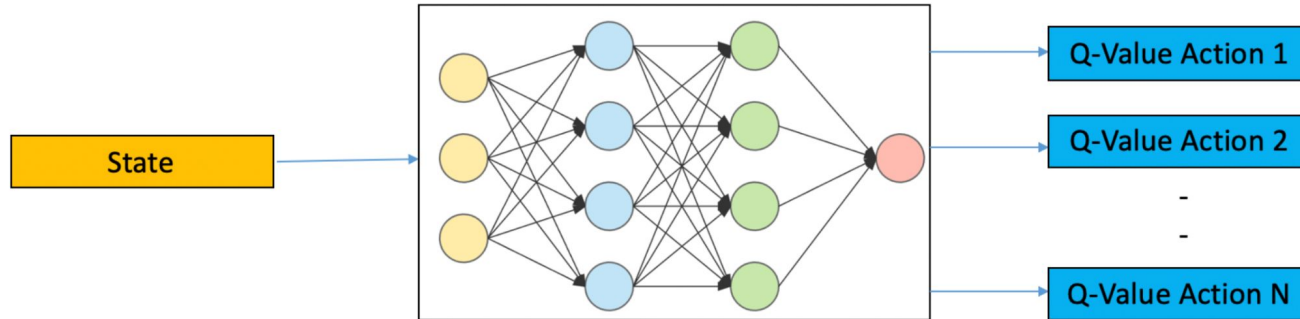
# Deep Q-Learning

---

---

# Deep Q-Learning

- Very similar idea to Q-Learning but instead of using a Q-table with (state, action) pairs mapped to a Q-value for choosing actions, we have a neural network that maps input states to (action, Q-value) pairs.
- Since in real life we could have many (state, action) pairs for a given environment, it is better to use a NN for a large state space because of its generalizability.



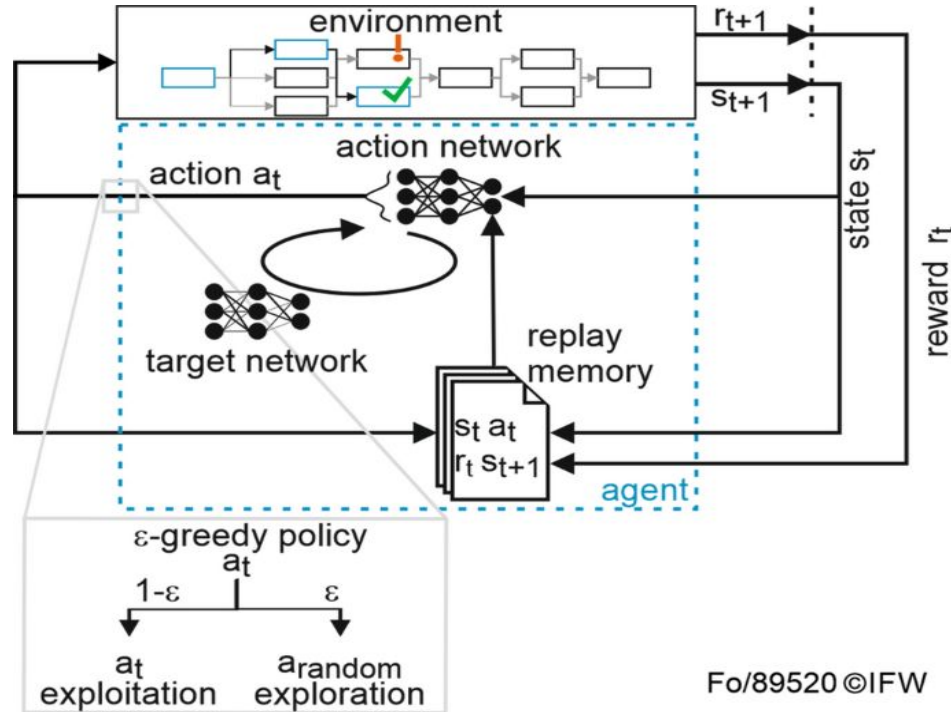
---

# Training a DQN

- DQN training works by taking an action via epsilon-greedy, observing the episode of that action, storing it in memory, and then sampling from the replay memory to update main prediction NN.
- We use two networks: a main and target network
  - Main NN: Where our Q-values are derived
  - Target NN: Same parameters as the main NN however, this network is not trained. It is synchronized with the main NN after N steps. This allows us to train the main NN with a single target thus making training stable.
- Weights are updated in the main network via the Q-Learning Loss function:

$$L(\theta) = \left( \overbrace{(r + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}; \theta^{target}))}^{\text{Target Q value}} - \underbrace{Q(s, a; \theta^{pred})}_{\text{Predicted Q value}} \right)^2$$

# Training a DQN



Fo/89520 ©IFW

---

# **DQN Algorithm: Deep Deterministic Policy Gradient (DDPG)**

---

---

# DDPG Overview

- Q-Learning and Deep Q-Learning are great for environments that have discrete action-state spaces. But, what if our action-state space is *continuous*?
- What if the amount of actions we can take in any given state is not finite?
- Remember, the ultimate goal of our agent is to choose the best action for any state:

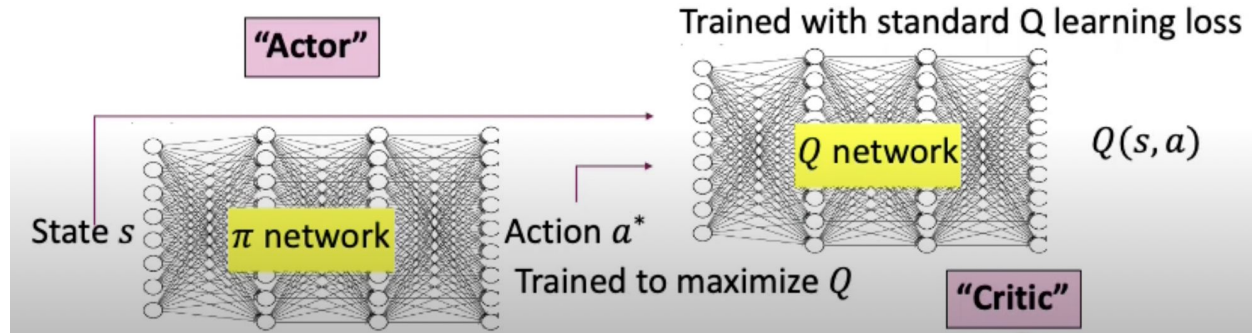
$$\pi(s) = a^*(s)$$

$$\text{where; } a^*(s) = \operatorname{argmax}_a Q(s, a)$$

- To solve this issue of having a continuous action-state space, DDPG uses an actor network and critic network.
    - Actor network: Responsible for mapping state to  $a^*$
    - Critic network: Responsible for mapping  $(s, a^*)$  tuple to  $Q(s, a)$
-

# DDPG Training

- The actor network is trained to maximize the Q-value from the critic network so its update rule uses gradient **ascent**.
- The critic network is updated using the traditional Q-Learning Loss function with.
- Similar to DQN, each actor and critic network both have target networks to help with stabilizing learning.



---

# DDPG Algorithm Steps

1. Initialize actor  $\theta$  params ( $\pi$  net), critic  $\phi$  params (Q net), and empty replay buffer B
2. Set target params equal to their respective main networks (just like DQ Learning)
3. Repeat until convergence:
  - a. Observe state  $s$  and get  $a^*$  from critic and execute  $a^*$
  - b. Observe state  $s'$ , get reward  $r$ , and set  $d$  = is in terminal state (0=no, 1=yes)
  - c. Store episode  $(s, a, r, s', d)$  in replay buffer
  - d. If  $d = 1$  reset environment parameters otherwise continue
  - e. If time to update:
    - i. For however many updates:
      1. Get sample batch from B
      2. Compute target via Bellman:  $y(r, s', d) = r + \gamma(1-d)Q_{\phi_{\text{target}}}(s', \pi_{\theta_{\text{target}}})$
      3. Update critic net via one step of GD using  $y(r, s', d)$
      4. Update actor net via one step of GA using:  $Q_{\phi}(s, \pi_{\theta}(s))$
      5. Update target networks with their respective main network



---

# DDPG Example: Stock Trading Agent

---

---

# Problem Justification

- Managing a portfolio of stocks to create a comfortable return can be very difficult:
    - Emotions can impact one's due diligence on a planned strategy.
    - Trading in live time can be stressful.
    - You don't know what you're doing.
  - Because the stock market changes every second of everyday, the state space is **extremely** large. DDPG is intended for large state spaces.
  - Using DDPG, we can potentially find a viable trading strategy within the vast complexities of the financial markets.
-

---

# Problem Overview

## Some Stuff for Q-Learning:

- Actions: buy, sell, or hold a stock.
- State space: N stocks we are interested in that have features such as closing, high, low, volume, day of week, rsi, ema, sma, etc..
- Reward: Return of stock of trade.
- Therefore, overall goal of agent is to maximize trading returns

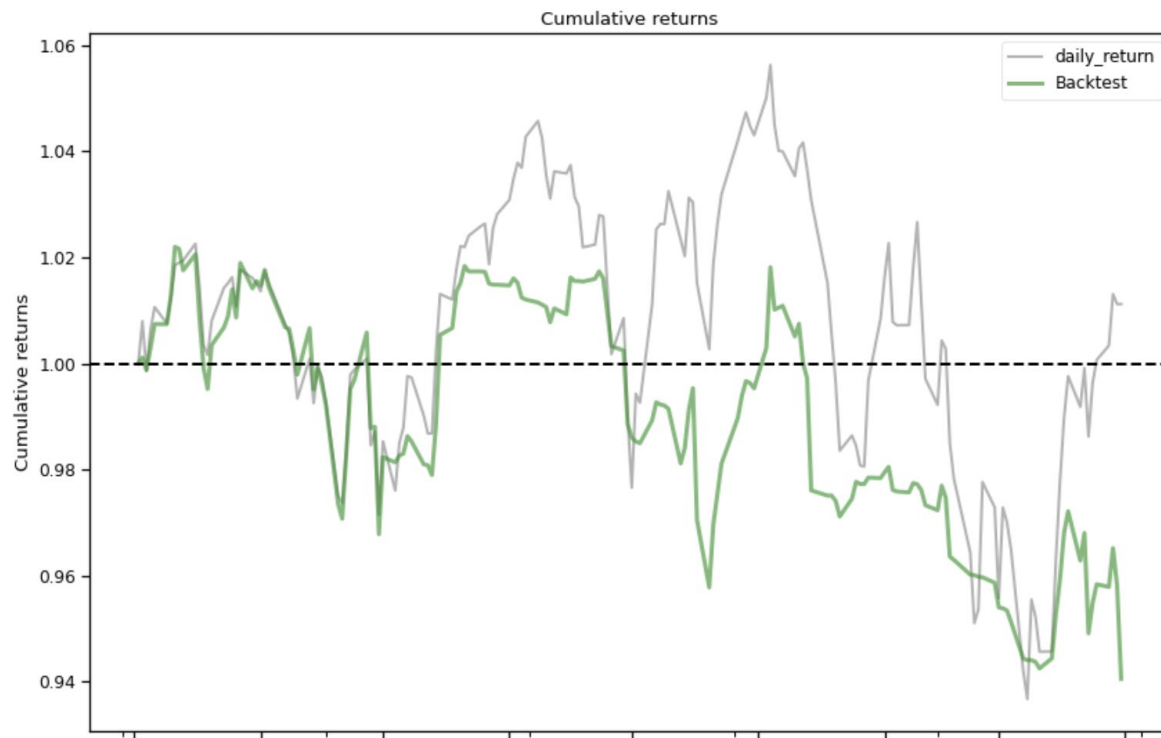
## Data:

- We will use the 30 stocks from the Dow 30 and use the overall return of the Dow 30 in a given year as the baseline for our agent.
  - For training, we will use 10 years of historical data from each of the 30 stocks.
  - For validation and backtesting, we will use 1 year of historical data from the 30 stocks.
-

# Dataset

date	open	high	low	close	adjcl	volume	tic	day	macd	boll_ub	boll_lb	rsi_30	cci_30	dx_30	close_30_sma	close_60_sma
2009-01-02	3.067143	3.261429	3.041429	3.241071	2.762747	746015200	AAPL	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	3.241071	3.241071
2009-01-02	58.590000	59.080002	57.750000	58.990002	44.219173	6547900	AMCN	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	58.990002	58.990002
2009-01-02	18.570000	19.520000	18.400000	19.330000	15.418568	18955700	AXP	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	19.330000	19.330000
2009-01-02	42.799999	45.560001	42.779999	45.250000	33.941097	7010200	BA	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	45.250000	45.250000
2009-01-02	44.910000	46.980000	44.709999	46.910000	31.729940	7117200	CAT	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	46.910000	46.910000
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
2022-03-31	517.099976	521.890015	509.670013	509.970001	506.568329	3979700	UNH	3	9.082670	525.617835	475.733175	57.179088	102.187053	34.560233	490.984004	481.827002
2022-03-31	223.910004	225.919998	220.440002	221.770004	220.463120	10759500	V	3	2.785993	233.259210	187.598787	53.368487	85.681562	12.340181	212.350666	216.093500
2022-03-31	51.660000	51.750000	50.930000	50.939999	48.878143	31027600	VZ	3	-0.600788	54.738811	49.623188	42.542530	-93.798458	24.604509	52.704333	52.972667
2022-03-31	45.290001	45.740002	44.169998	44.770000	43.221767	23284700	WBA	3	-0.360896	48.792783	45.911217	39.979683	-177.475807	26.487251	46.905000	49.042667
2022-03-31	148.789993	150.539993	148.179993	148.919998	147.737213	9054600	WMT	3	2.213517	149.642543	138.487459	58.973003	144.878746	45.338909	141.572668	140.756834
[97481 rows x 16 columns]																
date	open	high	low	close	adjcl	volume	tic	day	macd	boll_ub	boll_lb	rsi_30	cci_30	dx_30	close_30_sma	close_60_sma
2009-01-02	3.067143	3.261429	3.041429	3.241071	2.762747	746015200	AAPL	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	3.241071	3.241071
2009-01-02	58.590000	59.080002	57.750000	58.990002	44.219173	6547900	AMCN	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	58.990002	58.990002
2009-01-02	18.570000	19.520000	18.400000	19.330000	15.418568	18955700	AXP	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	19.330000	19.330000
2009-01-02	42.799999	45.560001	42.779999	45.250000	33.941097	7010200	BA	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	45.250000	45.250000
2009-01-02	44.910000	46.980000	44.709999	46.910000	31.729940	7117200	CAT	4	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	46.910000	46.910000
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
2022-03-31	517.099976	521.890015	509.670013	509.970001	506.568329	3979700	UNH	3	9.082670	525.617835	475.733175	57.179088	102.187053	34.560233	490.984004	481.827002
2022-03-31	223.910004	225.919998	220.440002	221.770004	220.463120	10759500	V	3	2.785993	233.259210	187.598787	53.368487	85.681562	12.340181	212.350666	216.093500
2022-03-31	51.660000	51.750000	50.930000	50.939999	48.878143	31027600	VZ	3	-0.600788	54.738811	49.623188	42.542530	-93.798458	24.604509	52.704333	52.972667
2022-03-31	45.290001	45.740002	44.169998	44.770000	43.221767	23284700	WBA	3	-0.360896	48.792783	45.911217	39.979683	-177.475807	26.487251	46.905000	49.042667
2022-03-31	148.789993	150.539993	148.179993	148.919998	147.737213	9054600	WMT	3	2.213517	149.642543	138.487459	58.973003	144.878746	45.338909	141.572668	140.756834

# Results



---

# Results

- Model performed poorly when compared to the Dow Index.
  - Model appears to have adopted a FOMO strategy.
  - Annual Returns (backtested on August 2021 to October 2022):
    - DOW Index Return: +1.68%
    - DDPG Backtest: -8.74%
  - DDPG model did not perform as poorly as the S&P 500 though: -17.09%
-

---

# Future Improvements

- Adjusting hyperparameters:
    - Changing learning rate
    - Lowering batch sample size
    - Lowering the number of steps.
  - Using different features for training (deeper research into technical indicators).
  - Using a different set of stocks that have a lot more volatility.
  - Try using one stock?
-

---

**Thank you!**  
**Questions?**

---