

# Reinforcement Learning Basics Any% Speedrun

David Quarel

ARENA

Thursday, 8th June 2023

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - ① We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees “real world” once it's already performing well

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - 1 We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees “real world” once it's already performing well
  - 2 Data is i.i.d between batches

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - 1 We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees “real world” once it's already performing well
  - 2 Data is i.i.d between batches
  - 3 No planning required, future predictions don't depend on past predictions

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - ① We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees “real world” once it's already performing well
  - ② Data is i.i.d between batches
  - ③ No planning required, future predictions don't depend on past predictions
- RL is vastly different: Agent takes actions in an interactive environment, receive scalar reward as feedback. This lends itself to several problems:

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - ① We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees “real world” once it's already performing well
  - ② Data is i.i.d between batches
  - ③ No planning required, future predictions don't depend on past predictions
- RL is vastly different: Agent takes actions in an interactive environment, receive scalar reward as feedback. This lends itself to several problems:
  - ① **Sparse reward:** Very little feedback during learning

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - 1 We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees "real world" once it's already performing well
  - 2 Data is i.i.d between batches
  - 3 No planning required, future predictions don't depend on past predictions
- RL is vastly different: Agent takes actions in an interactive environment, receive scalar reward as feedback. This lends itself to several problems:
  - 1 **Sparse reward:** Very little feedback during learning
  - 2 **Reward attribution:** Hard to tell which action was the one that caused the good reward



# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - 1 We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees “real world” once it's already performing well
  - 2 Data is i.i.d between batches
  - 3 No planning required, future predictions don't depend on past predictions
- RL is vastly different: Agent takes actions in an interactive environment, receive scalar reward as feedback. This lends itself to several problems:
  - 1 **Sparse reward:** Very little feedback during learning
  - 2 **Reward attribution:** Hard to tell which action was the one that caused the good reward
  - 3 **No ground truth** Optimal or even good policies may be unknown, (in pure RL settings) no data from good players to compare against

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - 1 We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees "real world" once it's already performing well
  - 2 Data is i.i.d between batches
  - 3 No planning required, future predictions don't depend on past predictions
- RL is vastly different: Agent takes actions in an interactive environment, receive scalar reward as feedback. This lends itself to several problems:
  - 1 **Sparse reward:** Very little feedback during learning
  - 2 **Reward attribution:** Hard to tell which action was the one that caused the good reward
  - 3 **No ground truth** Optimal or even good policies may be unknown, (in pure RL settings) no data from good players to compare against
  - 4 **Explore vs. Exploit tradeoff:**
    - 1 **Exploration:** Taking actions to learn how the world works (and improve the policy).
    - 2 **Exploitation:** Taking actions that maximise the expected sum of reward given current policy.

# What is Reinforcement Learning?

- Up to this point we've been mostly within the regime of supervised learning: Given some labelled data, train a model to minimise loss, then deploy to classify new data.
  - ① We have access to labelled training data, and only deploy the agent after we get good performance. Agent only sees “real world” once it's already performing well
  - ② Data is i.i.d between batches
  - ③ No planning required, future predictions don't depend on past predictions
- RL is vastly different: Agent takes actions in an interactive environment, receive scalar reward as feedback. This lends itself to several problems:
  - ① **Sparse reward:** Very little feedback during learning
  - ② **Reward attribution:** Hard to tell which action was the one that caused the good reward
  - ③ **No ground truth** Optimal or even good policies may be unknown, (in pure RL settings) no data from good players to compare against
  - ④ **Explore vs. Exploit tradeoff:**
    - ① **Exploration:** Taking actions to learn how the world works (and improve the policy).
    - ② **Exploitation:** Taking actions that maximise the expected sum of reward given current policy.
  - ⑤ **Online only:** No clear distinction between training and testing. Agent gets dumped in the environment and must learn on the fly.

# How Much Information is the Machine Given during Learning?

## ▶ “Pure” Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

## ▶ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

## ▶ Self-Supervised Learning (cake génoise)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



# How Much **(cherry)** is the Machine Given during Learning?

- ▶ “Pure” **(cherry)** Learning **(cherry)**
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ Millions of **(cherry)** samples

- ▶ Supervised Learning (fruit)
  - ▶ The machine predicts a category of a few numbers
  - ▶ Learning rate
  - ▶ 10,000 samples per example

- ▶ Self-supervised Learning (cake exercise)
  - ▶ The machine predicts a few numbers for any
  - ▶ Learning rate
  - ▶ 10,000 samples per example



# Multi-Armed Bandits

- The simplest type of RL environment with interaction: (equivalent to MDP with 1-state)

# Multi-Armed Bandits

- The simplest type of RL environment with interaction: (equivalent to MDP with 1-state)
- Agent has a set of “arms” (actions)  $\mathcal{A}$ . Environment has a family of reward distributions  $\{p_a\}_{a \in \mathcal{A}}$  for each action.

# Multi-Armed Bandits

- The simplest type of RL environment with interaction: (equivalent to MDP with 1-state)
- Agent has a set of “arms” (actions)  $\mathcal{A}$ . Environment has a family of reward distributions  $\{p_a\}_{a \in \mathcal{A}}$  for each action.
- On timestep  $t$ , agent chooses action  $a_t$  and receives reward  $r_t \sim p_{a_t}(\cdot)$ . Distributions  $p_i$  are unknown to agent.



# Multi-Armed Bandits

- The simplest type of RL environment with interaction: (equivalent to MDP with 1-state)
- Agent has a set of “arms” (actions)  $\mathcal{A}$ . Environment has a family of reward distributions  $\{p_a\}_{a \in \mathcal{A}}$  for each action.
- On timestep  $t$ , agent chooses action  $a_t$  and receives reward  $r_t \sim p_{a_t}(\cdot)$ . Distributions  $p_i$  are unknown to agent.
- Want to always choose the arm with the highest expected payout:

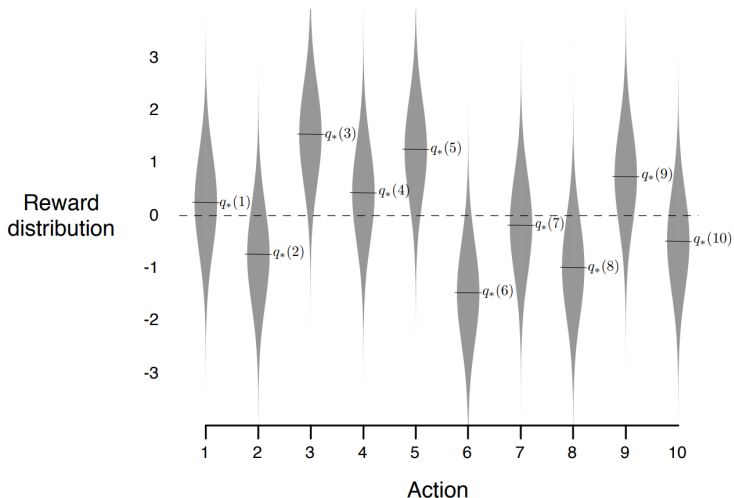
$$q_*(a) = \mathbb{E}[r_t | a_t = a]$$

# Multi-Armed Bandits

- The simplest type of RL environment with interaction: (equivalent to MDP with 1-state)
- Agent has a set of “arms” (actions)  $\mathcal{A}$ . Environment has a family of reward distributions  $\{p_a\}_{a \in \mathcal{A}}$  for each action.
- On timestep  $t$ , agent chooses action  $a_t$  and receives reward  $r_t \sim p_{a_t}(\cdot)$ . Distributions  $p_i$  are unknown to agent.
- Want to always choose the arm with the highest expected payout:

$$q_*(a) = \mathbb{E}[r_t | a_t = a]$$

- Need to balance trying all the arms to get a good estimate of the value of each arm, v.s. always trying to pull the best arm.



**Figure 2.1:** An example bandit problem from the 10-armed testbed. The true value  $q_*(a)$  of each of the ten actions was selected according to a normal distribution with mean zero and unit variance, and then the actual rewards were selected according to a mean  $q_*(a)$ , unit-variance normal distribution, as suggested by these gray distributions.

# Bandito Algorithm

- Keep track of  $\hat{Q}(a)$ , the estimated value of each arm after  $t$  arm-pulls

$$\hat{Q}_t(a) = \frac{\text{sum of rewards when } a \text{ taken up to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1, a_i=a}^{t-1} r_t}{\sum_{i=1, a_i=a}^{t-1} 1}$$

# Bandito Algorithm

- Keep track of  $\hat{Q}(a)$ , the estimated value of each arm after  $t$  arm-pulls

$$\hat{Q}_t(a) = \frac{\text{sum of rewards when } a \text{ taken up to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1, a_i=a}^{t-1} r_t}{\sum_{i=1, a_i=a}^{t-1} 1}$$

- $\hat{Q}_t(a)$  represents the empirical average reward obtained from arm  $a$  up to time  $t$ .

# Bandito Algorithm

- Keep track of  $\hat{Q}(a)$ , the estimated value of each arm after  $t$  arm-pulls

$$\hat{Q}_t(a) = \frac{\text{sum of rewards when } a \text{ taken up to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1, a_i=a}^{t-1} r_t}{\sum_{i=1, a_i=a}^{t-1} 1}$$

- $\hat{Q}_t(a)$  represents the empirical average reward obtained from arm  $a$  up to time  $t$ .
- In practice, easier to init  $\hat{Q}_1(a) = \hat{R}_1(a) = \hat{N}_1(a) = 0$  and

$$\begin{aligned}\hat{R}_{t+1}(a) &\leftarrow \hat{R}_t(a) + r_t \llbracket a_t = a \rrbracket & \hat{N}_{t+1}(a) &\leftarrow \hat{N}_t(a) + \llbracket a_t = a \rrbracket \\ \hat{Q}_{t+1}(a) &\leftarrow \frac{\hat{R}_{t+1}(a)}{\hat{N}_{t+1}(a)}\end{aligned}$$

where  $\llbracket P \rrbracket = 1$  if  $P$  evaluates to True, else  $\llbracket P \rrbracket = 0$ .

# Bandito Algorithm

- Keep track of  $\hat{Q}(a)$ , the estimated value of each arm after  $t$  arm-pulls

$$\hat{Q}_t(a) = \frac{\text{sum of rewards when } a \text{ taken up to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1, a_i=a}^{t-1} r_t}{\sum_{i=1, a_i=a}^{t-1} 1}$$

- $\hat{Q}_t(a)$  represents the empirical average reward obtained from arm  $a$  up to time  $t$ .
- In practice, easier to init  $\hat{Q}_1(a) = \hat{R}_1(a) = \hat{N}_1(a) = 0$  and

$$\begin{aligned}\hat{R}_{t+1}(a) &\leftarrow \hat{R}_t(a) + r_t \llbracket a_t = a \rrbracket & \hat{N}_{t+1}(a) &\leftarrow \hat{N}_t(a) + \llbracket a_t = a \rrbracket \\ \hat{Q}_{t+1}(a) &\leftarrow \frac{\hat{R}_{t+1}(a)}{\hat{N}_{t+1}(a)}\end{aligned}$$

where  $\llbracket P \rrbracket = 1$  if  $P$  evaluates to True, else  $\llbracket P \rrbracket = 0$ .

- Choose arm with highest estimated payout:  $a_t := \arg \max \hat{Q}_t(a)$ .

# Bandito Algorithm

- Keep track of  $\hat{Q}(a)$ , the estimated value of each arm after  $t$  arm-pulls

$$\hat{Q}_t(a) = \frac{\text{sum of rewards when } a \text{ taken up to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1, a_i=a}^{t-1} r_t}{\sum_{i=1, a_i=a}^{t-1} 1}$$

- $\hat{Q}_t(a)$  represents the empirical average reward obtained from arm  $a$  up to time  $t$ .
- In practice, easier to init  $\hat{Q}_1(a) = \hat{R}_1(a) = \hat{N}_1(a) = 0$  and

$$\begin{aligned}\hat{R}_{t+1}(a) &\leftarrow \hat{R}_t(a) + r_t \llbracket a_t = a \rrbracket & \hat{N}_{t+1}(a) &\leftarrow \hat{N}_t(a) + \llbracket a_t = a \rrbracket \\ \hat{Q}_{t+1}(a) &\leftarrow \frac{\hat{R}_{t+1}(a)}{\hat{N}_{t+1}(a)}\end{aligned}$$

where  $\llbracket P \rrbracket = 1$  if  $P$  evaluates to True, else  $\llbracket P \rrbracket = 0$ .

- Choose arm with highest estimated payout:  $a_t := \arg \max \hat{Q}_t(a)$ .
- Problem:** Can get stuck with a suboptimal arm.



# Encouraging Exploration

- 1 **First approach:** Just do random stuff every now and again, hope for the best

$$a_t^{\epsilon\text{-greedy}} = \begin{cases} \text{Do random action} & \text{Prob } \epsilon \\ \arg \max_{a'} Q_t(a') & \text{Prob } 1 - \epsilon \end{cases}$$

# Encouraging Exploration

- 1 **First approach:** Just do random stuff every now and again, hope for the best

$$a_t^{\epsilon\text{-greedy}} = \begin{cases} \text{Do random action} & \text{Prob } \epsilon \\ \arg \max_{a'} Q_t(a') & \text{Prob } 1 - \epsilon \end{cases}$$

- 2 **Better approach:** Give a bonus to actions seldom taken

$$a_t^{UCB} = \arg \max_{a'} \left( Q_t(a') + c \sqrt{\frac{\ln t}{N_t(a')}} \right)$$

# Encouraging Exploration

- 1 **First approach:** Just do random stuff every now and again, hope for the best

$$a_t^{\epsilon\text{-greedy}} = \begin{cases} \text{Do random action} & \text{Prob } \epsilon \\ \arg \max_{a'} Q_t(a') & \text{Prob } 1 - \epsilon \end{cases}$$

- 2 **Better approach:** Give a bonus to actions seldom taken

$$a_t^{UCB} = \arg \max_{a'} \left( Q_t(a') + c \sqrt{\frac{\ln t}{N_t(a')}} \right)$$

- 3 **Intuition:** Error of  $Q_t(a)$  is  $\propto \frac{1}{\sqrt{N_t(a)}}$ . Add a bonus proportional to variance, so actions with high variance  $\equiv$  few samples get explored

# Encouraging Exploration

- 1 **First approach:** Just do random stuff every now and again, hope for the best

$$a_t^{\epsilon\text{-greedy}} = \begin{cases} \text{Do random action} & \text{Prob } \epsilon \\ \arg \max_{a'} Q_t(a') & \text{Prob } 1 - \epsilon \end{cases}$$

- 2 **Better approach:** Give a bonus to actions seldom taken

$$a_t^{UCB} = \arg \max_{a'} \left( Q_t(a') + c \sqrt{\frac{\ln t}{N_t(a')}} \right)$$

- 3 **Intuition:** Error of  $Q_t(a)$  is  $\propto \frac{1}{\sqrt{N_t(a)}}$ . Add a bonus proportional to variance, so actions with high variance  $\equiv$  few samples get explored
- 4 Add  $\ln t$  to numerator to ensure every action is sampled infinitely often (in case you get an unlucky run).  $\ln t$  is optimal because math.  $c = 2$  works good in practice.

# Agent-Environment Interaction Loop (MDPs)

- Environment has *states*  $\mathcal{S}$ , *actions*  $\mathcal{A}$ , *rewards*  $\mathcal{R}$ , *environment distribution*  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{R} \rightarrow [0, 1]$ .

# Agent-Environment Interaction Loop (MDPs)

- Environment has *states*  $\mathcal{S}$ , *actions*  $\mathcal{A}$ , *rewards*  $\mathcal{R}$ , *environment distribution*  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{R} \rightarrow [0, 1]$ .
  - Think of  $p(s, a, s', r)$  as  $\Pr(s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a)$ . We write  $p(s', r | s, a)$  for clarity.

# Agent-Environment Interaction Loop (MDPs)

- Environment has *states*  $\mathcal{S}$ , *actions*  $\mathcal{A}$ , *rewards*  $\mathcal{R}$ , *environment distribution*  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{R} \rightarrow [0, 1]$ .
  - Think of  $p(s, a, s', r)$  as  $\Pr(s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a)$ . We write  $p(s', r | s, a)$  for clarity.
- In timestep  $t$ , agent samples  $a_t \sim \pi(s_t)$  from *policy*  $\pi_t$ . Environment samples  $(s_{t+1}, r_{t+1}) \sim p(\cdot | s_t, a_t)$ .

# Agent-Environment Interaction Loop (MDPs)

- Environment has *states*  $\mathcal{S}$ , *actions*  $\mathcal{A}$ , *rewards*  $\mathcal{R}$ , *environment distribution*  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{R} \rightarrow [0, 1]$ .
  - Think of  $p(s, a, s', r)$  as  $\Pr(s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a)$ . We write  $p(s', r | s, a)$  for clarity.
- In timestep  $t$ , agent samples  $a_t \sim \pi(s_t)$  from *policy*  $\pi_t$ . Environment samples  $(s_{t+1}, r_{t+1}) \sim p(\cdot | s_t, a_t)$ .
- Generates an interaction history, or *trajectory*

$$s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, r_3, s_3, \dots$$



# Agent-Environment Interaction Loop (MDPs)

- Environment has *states*  $\mathcal{S}$ , *actions*  $\mathcal{A}$ , *rewards*  $\mathcal{R}$ , *environment distribution*  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{R} \rightarrow [0, 1]$ .
  - Think of  $p(s, a, s', r)$  as  $\Pr(s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a)$ . We write  $p(s', r | s, a)$  for clarity.
- In timestep  $t$ , agent samples  $a_t \sim \pi(s_t)$  from *policy*  $\pi_t$ . Environment samples  $(s_{t+1}, r_{t+1}) \sim p(\cdot | s_t, a_t)$ .
- Generates an interaction history, or *trajectory*

$$s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, r_3, s_3, \dots$$

- Agent may choose to update choice of policy at any timestep. Most RL algorithms focus on the mechanism that does this.

# Agent-Environment Interaction Loop (MDPs)

- Environment has *states*  $\mathcal{S}$ , *actions*  $\mathcal{A}$ , *rewards*  $\mathcal{R}$ , *environment distribution*  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{R} \rightarrow [0, 1]$ .
  - Think of  $p(s, a, s', r)$  as  $\Pr(s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a)$ . We write  $p(s', r | s, a)$  for clarity.
- In timestep  $t$ , agent samples  $a_t \sim \pi(s_t)$  from *policy*  $\pi_t$ . Environment samples  $(s_{t+1}, r_{t+1}) \sim p(\cdot | s_t, a_t)$ .
- Generates an interaction history, or *trajectory*

$$s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, r_3, s_3, \dots$$

- Agent may choose to update choice of policy at any timestep. Most RL algorithms focus on the mechanism that does this.

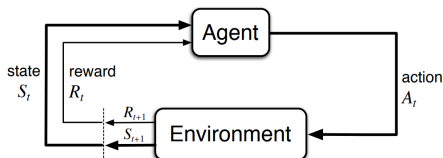


Figure: Agent-Environment interaction loop

# Objective of the Agent

- At timestep  $t$ , the *return*  $G_t$  is the sum of all future rewards:

$$G_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots$$

# Objective of the Agent

- At timestep  $t$ , the *return*  $G_t$  is the sum of all future rewards:

$$G_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots$$

- **Goal:** Maximise the return.

# Objective of the Agent

- At timestep  $t$ , the *return*  $G_t$  is the sum of all future rewards:

$$G_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots$$

- **Goal:** Maximise the return.
  - For episodic (finite length interaction) environments of maximum duration  $T$ , return  $G_t = r_{t+1} + r_{t+2} + \dots + r_T$  well defined.

# Objective of the Agent

- At timestep  $t$ , the *return*  $G_t$  is the sum of all future rewards:

$$G_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots$$

- **Goal:** Maximise the return.
  - For episodic (finite length interaction) environments of maximum duration  $T$ , return  $G_t = r_{t+1} + r_{t+2} + \dots + r_T$  well defined.
- **Problems:** (for continuing environments)
  - The return may diverge or be undefined (compare 2, 2, 2, 2, ... with 1, 1, 1, 1, ...).
  - The agent might be lazy (compare 1, 1, 1 ... with 0, 0, ..., 0, 1, 1, 1, ...).
  - The environment is stochastic, and the rewards are often up to chance. How to trade-off unlikely big rewards with likely small rewards?
  - May desire rewards now to be more valuable than rewards later: \$100 now? Or \$110 in a year?

# Objective of the Agent

- At timestep  $t$ , the *return*  $G_t$  is the sum of all future rewards:

$$G_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots$$

- **Goal:** Maximise the return.
  - For episodic (finite length interaction) environments of maximum duration  $T$ , return  $G_t = r_{t+1} + r_{t+2} + \dots + r_T$  well defined.
- **Problems:** (for continuing environments)
  - The return may diverge or be undefined (compare 2, 2, 2, 2, ... with 1, 1, 1, 1, ...).
  - The agent might be lazy (compare 1, 1, 1 ... with 0, 0, ..., 0, 1, 1, 1, ...).
  - The environment is stochastic, and the rewards are often up to chance. How to trade-off unlikely big rewards with likely small rewards?
  - May desire rewards now to be more valuable than rewards later: \$100 now? Or \$110 in a year?
- **Solutions:**
  - Add a discount factor  $\gamma \in [0, 1)$  so rewards more imminent are worth more, and the return is always well defined.

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$$

- Want agent to choose actions to maximise the *expected* return. 

# Core Assumptions

These kind of environments are called *Markov Decision Processes* (MDPs), and have the following “nice” properties



# Core Assumptions

These kind of environments are called *Markov Decision Processes* (MDPs), and have the following “nice” properties

- 1 **Stationary:** The environmental distribution  $p$  is fixed and does not change over time
  - *Old data is as useful as new data*

# Core Assumptions

These kind of environments are called *Markov Decision Processes* (MDPs), and have the following “nice” properties

- 1 **Stationary:** The environmental distribution  $p$  is fixed and does not change over time
  - *Old data is as useful as new data*
- 2 **Markovian:** The behaviour of the environment at timestep  $t$  depends only on the current state  $s_t$  and action  $a_t$ .
  - *Only need to consider the current state to act optimally, the past is irrelevant*

# Core Assumptions

These kind of environments are called *Markov Decision Processes* (MDPs), and have the following “nice” properties

- 1 **Stationary:** The environmental distribution  $p$  is fixed and does not change over time
  - *Old data is as useful as new data*
- 2 **Markovian:** The behaviour of the environment at timestep  $t$  depends only on the current state  $s_t$  and action  $a_t$ .
  - *Only need to consider the current state to act optimally, the past is irrelevant*
- 3 **Fully Observable:** The state is a full description of the world
  - *Agent always has access to sufficient information to choose the optimal action*

# Core Assumptions

These kind of environments are called *Markov Decision Processes* (MDPs), and have the following “nice” properties

- 1 **Stationary:** The environmental distribution  $p$  is fixed and does not change over time
  - *Old data is as useful as new data*
- 2 **Markovian:** The behaviour of the environment at timestep  $t$  depends only on the current state  $s_t$  and action  $a_t$ .
  - *Only need to consider the current state to act optimally, the past is irrelevant*
- 3 **Fully Observable:** The state is a full description of the world
  - *Agent always has access to sufficient information to choose the optimal action*
- 4 **Reward Hypothesis:**

*“That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward).” -Rich Sutton*

  - *Reward alone is sufficient to communicate any possible goal or desired behaviour*

# Value Function

- Want to define the “goodness” (*value*) of a state, so the agent can take actions to move towards “good” states, and away from “bad” states.

# Value Function

- Want to define the “goodness” (*value*) of a state, so the agent can take actions to move towards “good” states, and away from “bad” states.
- The value of a state depends also on how the agent chooses actions, called a *policy*  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ . Actions are sampled  $a \sim \pi(\cdot|s)$ .

# Value Function

- Want to define the “goodness” (*value*) of a state, so the agent can take actions to move towards “good” states, and away from “bad” states.
- The value of a state depends also on how the agent chooses actions, called a *policy*  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ . Actions are sampled  $a \sim \pi(\cdot|s)$ .

## Value Function

$$\begin{aligned}V_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | s_t = s] \\ &= \mathbb{E}_{\pi}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s]\end{aligned}$$

(Expectation is also with respect to the environment  $p$ .)

# Bellman Equation

We note that since

$$\begin{aligned}G_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\ &= r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) \\ &= r_{t+1} + \gamma G_{t+1}\end{aligned}$$



# Bellman Equation

We note that since

$$\begin{aligned}G_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\ &= r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) \\ &= r_{t+1} + \gamma G_{t+1}\end{aligned}$$

we can then define the value function recursively,

$$V_\pi(s) = \mathbb{E}_\pi[G_t \mid s_t = s]$$

# Bellman Equation

We note that since

$$\begin{aligned}G_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\ &= r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) \\ &= r_{t+1} + \gamma G_{t+1}\end{aligned}$$

we can then define the value function recursively,

$$\begin{aligned}V_\pi(s) &= \mathbb{E}_\pi[G_t \mid s_t = s] \\ &= \mathbb{E}_\pi[r_{t+1} + \gamma G_{t+1} \mid s_t = s]\end{aligned}$$

# Bellman Equation

We note that since

$$\begin{aligned}G_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\&= r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) \\&= r_{t+1} + \gamma G_{t+1}\end{aligned}$$

we can then define the value function recursively,

$$\begin{aligned}V_\pi(s) &= \mathbb{E}_\pi[G_t \mid s_t = s] \\&= \mathbb{E}_\pi[r_{t+1} + \gamma G_{t+1} \mid s_t = s] \\&= \mathbb{E}_\pi[r_{t+1} \mid s_t = s] + \gamma \mathbb{E}_\pi[G_{t+1} \mid s_t = s]\end{aligned}$$

# Bellman Equation

We note that since

$$\begin{aligned}G_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\&= r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) \\&= r_{t+1} + \gamma G_{t+1}\end{aligned}$$

we can then define the value function recursively,

$$\begin{aligned}V_\pi(s) &= \mathbb{E}_\pi[G_t \mid s_t = s] \\&= \mathbb{E}_\pi[r_{t+1} + \gamma G_{t+1} \mid s_t = s] \\&= \mathbb{E}_\pi[r_{t+1} \mid s_t = s] + \gamma \mathbb{E}_\pi[G_{t+1} \mid s_t = s] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) r \\&\quad + \gamma \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) \mathbb{E}_\pi[G_{t+1} \mid s_{t+1} = s']\end{aligned}$$

# Bellman Equation

We note that since

$$\begin{aligned}G_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\&= r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) \\&= r_{t+1} + \gamma G_{t+1}\end{aligned}$$

we can then define the value function recursively,

$$\begin{aligned}V_\pi(s) &= \mathbb{E}_\pi[G_t \mid s_t = s] \\&= \mathbb{E}_\pi[r_{t+1} + \gamma G_{t+1} \mid s_t = s] \\&= \mathbb{E}_\pi[r_{t+1} \mid s_t = s] + \gamma \mathbb{E}_\pi[G_{t+1} \mid s_t = s] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) r \\&\quad + \gamma \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) \mathbb{E}_\pi[G_{t+1} \mid s_{t+1} = s'] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) (r + \gamma V_\pi(s'))\end{aligned}$$

This gives the **Bellman equation**

## Bellman Equation

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

This gives the **Bellman equation**

## Bellman Equation

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

- Equation is linear in  $V_{\pi}(\cdot)$ , giving a set of **linear** simultaneous equations.

This gives the **Bellman equation**

## Bellman Equation

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

- Equation is linear in  $V_{\pi}(\cdot)$ , giving a set of **linear** simultaneous equations.
- Given policy  $\pi$ , can now easily solve for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$



This gives the **Bellman equation**

## Bellman Equation

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

- Equation is linear in  $V_{\pi}(\cdot)$ , giving a set of **linear** simultaneous equations.
- Given policy  $\pi$ , can now easily solve for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$
- Computing  $V_{\pi}$  from  $\pi$  is called **policy evaluation**.

# Value Function (simplified)

Assume policy  $\pi : S \rightarrow A$  is deterministic, define *transition probability*  $T(s' | s, a) := \sum_{r \in \mathcal{R}} p(s', r | s, a)$  and assume reward  $r_{t+1} := R(s_t, a_t, s_{t+1})$  is deterministic function of  $s_t, a_t, s_{t+1}$ .

## Bellman Equation

$$V_{\pi}(s) = \sum_{s'} T(s' | s, a) (R(s, a, s') + \gamma V_{\pi}(s'))$$

where  $a = \pi(s)$ .

# Value Function (simplified)

Assume policy  $\pi : S \rightarrow A$  is deterministic, define *transition probability*  $T(s' | s, a) := \sum_{r \in \mathcal{R}} p(s', r | s, a)$  and assume reward  $r_{t+1} := R(s_t, a_t, s_{t+1})$  is deterministic function of  $s_t, a_t, s_{t+1}$ .

## Bellman Equation

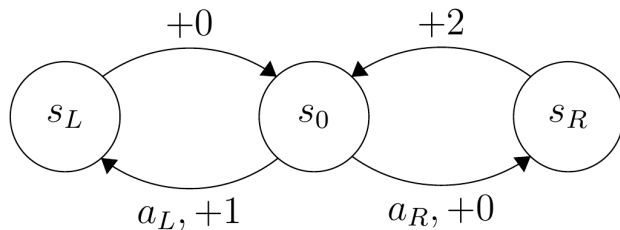
$$V_{\pi}(s) = \sum_{s'} T(s' | s, a) (R(s, a, s') + \gamma V_{\pi}(s'))$$

where  $a = \pi(s)$ .

Only need to sum over all states to find  $V_{\pi}(s)$  in terms of  $\{V_{\pi}(s_1), \dots, V_{\pi}(s_n)\}$ .

# Example Environment

- States  $\mathcal{S} = s_0, s_L, s_R$ , actions  $\mathcal{A} = \{a_L, a_R\}$ , rewards  $\mathcal{R} = \{0, 1, 2\}$ .
- Each transition indicates if an action is taken, the reward returned and which state to transition to
- What is the best action from state  $s_0$ ?



# Optimal Bellman

- Policy  $\pi_1$  is **better** than  $\pi_2$  ( $\pi_1 \geq \pi_2$ ) if  $\forall s. V_{\pi_1}(s) \geq V_{\pi_2}(s)$ . A policy is **optimal** if it is better than all other policies.

# Optimal Bellman

- Policy  $\pi_1$  is **better** than  $\pi_2$  ( $\pi_1 \geq \pi_2$ ) if  $\forall s. V_{\pi_1}(s) \geq V_{\pi_2}(s)$ . A policy is **optimal** if it is better than all other policies.
- **Theorem:** An optimal policy  $\pi^*$  always exists. Define optimal value function as

$$V_*(s) := V_{\pi^*}(s) \equiv \max_{\pi} V_{\pi}(s)$$

# Optimal Bellman

- Policy  $\pi_1$  is **better** than  $\pi_2$  ( $\pi_1 \geq \pi_2$ ) if  $\forall s. V_{\pi_1}(s) \geq V_{\pi_2}(s)$ . A policy is **optimal** if it is better than all other policies.
- **Theorem:** An optimal policy  $\pi^*$  always exists. Define optimal value function as

$$V_*(s) := V_{\pi^*}(s) \equiv \max_{\pi} V_{\pi}(s)$$

## Optimal Bellman Equation

$$V_*(s) = \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s'))$$

# Optimal Bellman

- Policy  $\pi_1$  is **better** than  $\pi_2$  ( $\pi_1 \geq \pi_2$ ) if  $\forall s. V_{\pi_1}(s) \geq V_{\pi_2}(s)$ . A policy is **optimal** if it is better than all other policies.
- **Theorem:** An optimal policy  $\pi^*$  always exists. Define optimal value function as

$$V_*(s) := V_{\pi^*}(s) \equiv \max_{\pi} V_{\pi}(s)$$

## Optimal Bellman Equation

$$V_*(s) = \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s'))$$

Gives a set of **non-linear** simultaneous equations with variables  $V_*(s_1), V_*(s_2), \dots$



# Optimal Bellman

- Policy  $\pi_1$  is **better** than  $\pi_2$  ( $\pi_1 \geq \pi_2$ ) if  $\forall s. V_{\pi_1}(s) \geq V_{\pi_2}(s)$ . A policy is **optimal** if it is better than all other policies.
- **Theorem:** An optimal policy  $\pi^*$  always exists. Define optimal value function as

$$V_*(s) := V_{\pi^*}(s) \equiv \max_{\pi} V_{\pi}(s)$$

## Optimal Bellman Equation

$$V_*(s) = \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s'))$$

Gives a set of **non-linear** simultaneous equations with variables  $V_*(s_1), V_*(s_2), \dots$

**Problem:** No clear way to solve for  $V_*(\cdot)$

# Optimal Bellman

- Policy  $\pi_1$  is **better** than  $\pi_2$  ( $\pi_1 \geq \pi_2$ ) if  $\forall s. V_{\pi_1}(s) \geq V_{\pi_2}(s)$ . A policy is **optimal** if it is better than all other policies.
- **Theorem:** An optimal policy  $\pi^*$  always exists. Define optimal value function as

$$V_*(s) := V_{\pi^*}(s) \equiv \max_{\pi} V_{\pi}(s)$$

## Optimal Bellman Equation

$$V_*(s) = \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s'))$$

Gives a set of **non-linear** simultaneous equations with variables  $V_*(s_1), V_*(s_2), \dots$

**Problem:** No clear way to solve for  $V_*(\cdot)$

Can't just compute  $V_{\pi}$  using policy evaluation for all  $\pi$ , as there are  $|\mathcal{A}|^{|\mathcal{S}|}$  many to choose from.

- Obviously we have that

$$\begin{aligned} V_*(s) &= \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s')) \\ &\geq \sum_{s'} T(s'|s, \pi(s)) (R(s, \pi(s), s') + \gamma V_*(s')) = V_\pi(s) \end{aligned}$$

# Policy Improvement

- Obviously we have that

$$\begin{aligned} V_*(s) &= \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s')) \\ &\geq \sum_{s'} T(s'|s, \pi(s)) (R(s, \pi(s), s') + \gamma V_*(s')) = V_\pi(s) \end{aligned}$$

- Given a policy  $\pi_n$ , can feed it through the optimal Bellman equation to get a better policy  $\pi_{n+1}$

# Policy Improvement

- Obviously we have that

$$\begin{aligned} V_*(s) &= \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s')) \\ &\geq \sum_{s'} T(s'|s, \pi(s)) (R(s, \pi(s), s') + \gamma V_*(s')) = V_\pi(s) \end{aligned}$$

- Given a policy  $\pi_n$ , can feed it through the optimal Bellman equation to get a better policy  $\pi_{n+1}$

## Policy Improvement

$$\pi_{n+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi_n}(s'))$$

# Policy Iteration

## Policy Improvement (I)

$$\pi_{n+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi_n}(s'))$$

## Policy Evaluation (E)

Solve

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$

# Policy Iteration

## Policy Improvement (I)

$$\pi_{n+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi_n}(s'))$$

## Policy Evaluation (E)

Solve

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$

- Start with arbitrary policy  $\pi_0$ .

# Policy Iteration

## Policy Improvement (I)

$$\pi_{n+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi_n}(s'))$$

## Policy Evaluation (E)

Solve

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$

- Start with arbitrary policy  $\pi_0$ .
- Note that  $\pi_*$  is fixed point of policy improvement.



# Policy Iteration

## Policy Improvement (I)

$$\pi_{n+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi_n}(s'))$$

## Policy Evaluation (E)

Solve

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$

- Start with arbitrary policy  $\pi_0$ .
- Note that  $\pi_*$  is fixed point of policy improvement.
- Alternate until policy is stable

# Policy Iteration

## Policy Improvement (I)

$$\pi_{n+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi_n}(s'))$$

## Policy Evaluation (E)

Solve

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$

- Start with arbitrary policy  $\pi_0$ .
- Note that  $\pi_*$  is fixed point of policy improvement.
- Alternate until policy is stable

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{E} \pi_2 \xrightarrow{I} V_{\pi_2} \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} V_{\pi_*} \xrightarrow{I} \pi_*$$

# Policy Iteration

## Policy Improvement (I)

$$\pi_{n+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi_n}(s'))$$

## Policy Evaluation (E)

Solve

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) (r + \gamma V_{\pi}(s'))$$

for  $V_{\pi}(s_1), V_{\pi}(s_2), \dots$

- Start with arbitrary policy  $\pi_0$ .
- Note that  $\pi_*$  is fixed point of policy improvement.
- Alternate until policy is stable

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{E} \pi_2 \xrightarrow{I} V_{\pi_2} \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} V_{\pi_*} \xrightarrow{I} \pi_*$$

**Theorem:** Policy iteration converges to optimal policy in finitely many steps!

# Problems with Policy Iteration

- Requires white-box access to the environmental distribution  $T$  and reward function  $R$ .
- Only works for environments with few enough states and actions to sweep through.

# Problems with Policy Iteration

- Requires white-box access to the environmental distribution  $T$  and reward function  $R$ .
- Only works for environments with few enough states and actions to sweep through.

For the moment, we weaken only the first assumption, and assume the environment is now a black box, from which state-reward pairs  $(s', r)$  can be sampled given state-action pairs  $(s, a)$  as input.

# Temporal Difference Learning

- **Goal:** Perform policy evaluation without access to environmental distribution.

# Temporal Difference Learning

- **Goal:** Perform policy evaluation without access to environmental distribution.
- **Motivation:** Consider once again the value function:

$$V_{\pi}(s) = \mathbb{E}_{\substack{a=\pi(s) \\ s' \sim T(\cdot|s,a)}} [R(s, a, s') + \gamma V_{\pi}(s')]$$

# Temporal Difference Learning

- **Goal:** Perform policy evaluation without access to environmental distribution.
- **Motivation:** Consider once again the value function:

$$V_{\pi}(s) = \mathbb{E}_{\substack{a=\pi(s) \\ s' \sim T(\cdot|s,a)}} [R(s, a, s') + \gamma V_{\pi}(s')]$$

On timestep  $t$ , this is “*on average*”, equal to the actual reward  $r_{t+1}$ , plus the discounted value of the actual next state  $s_{t+1}$ .

$$V_{\pi}(s_t) \approx r_{t+1} + \gamma V_{\pi}(s_{t+1})$$



# Temporal Difference Learning

- **Goal:** Perform policy evaluation without access to environmental distribution.
- **Motivation:** Consider once again the value function:

$$V_{\pi}(s) = \mathbb{E}_{\substack{a=\pi(s) \\ s' \sim T(\cdot|s,a)}} [R(s, a, s') + \gamma V_{\pi}(s')]$$

On timestep  $t$ , this is “*on average*”, equal to the actual reward  $r_{t+1}$ , plus the discounted value of the actual next state  $s_{t+1}$ .

$$V_{\pi}(s_t) \approx r_{t+1} + \gamma V_{\pi}(s_{t+1})$$

We define the **TD-Error** as the difference

$$\delta_t := r_{t+1} + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t)$$

# Temporal Difference Learning

- **Goal:** Perform policy evaluation without access to environmental distribution.
- **Motivation:** Consider once again the value function:

$$V_{\pi}(s) = \mathbb{E}_{\substack{a=\pi(s) \\ s' \sim T(\cdot|s,a)}} [R(s, a, s') + \gamma V_{\pi}(s')]$$

On timestep  $t$ , this is “on average”, equal to the actual reward  $r_{t+1}$ , plus the discounted value of the actual next state  $s_{t+1}$ .

$$V_{\pi}(s_t) \approx r_{t+1} + \gamma V_{\pi}(s_{t+1})$$

We define the **TD-Error** as the difference

$$\delta_t := r_{t+1} + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t)$$

This then gives us an update rule to improve on our estimate  $\hat{V}_{\pi}$  of  $V_{\pi}$ , similar to SGD, called  $TD(0)$ .

$$\begin{aligned} \hat{V}_{\pi}(s_t) &\leftarrow \hat{V}_{\pi}(s_t) + \alpha \delta_t \\ &\equiv \hat{V}_{\pi}(s_t) + \alpha \left( r_{t+1} + \gamma \hat{V}_{\pi}(s_{t+1}) - \hat{V}_{\pi}(s_t) \right) \end{aligned}$$

where  $\alpha \in (0, 1]$  is the **learning rate**.

# Q-Value

- **Q-value** is the expected return from state  $s$ , taking action  $a$ , and thereafter following policy  $\pi$ .

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | s_t = s, a_t = a]$$

# Q-Value

- **Q-value** is the expected return from state  $s$ , taking action  $a$ , and thereafter following policy  $\pi$ .

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | s_t = s, a_t = a]$$

Contrast with the value function

$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | s_t = s]$$

# Q-Value

- **Q-value** is the expected return from state  $s$ , taking action  $a$ , and thereafter following policy  $\pi$ .

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | s_t = s, a_t = a]$$

Contrast with the value function

$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | s_t = s]$$

## Q-value Bellman

$$Q_{\pi}(s, a) = \sum_{s'} T(s' | s, a) (R(s, a, s') + \gamma Q_{\pi}(s', a'))$$

where  $a' = \pi(s')$

## Optimal Q-value Bellman

$$Q_{*}(s, a) = \sum_{s'} T(s' | s, a) \left( R(s, a, s') + \max_{a'} Q_{*}(s', a') \right)$$

# Q-value vs. Value

Can state  $Q$  in terms of  $V$ , and vice-versa.

$$Q_{\pi}(s, a) = \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi}(s'))$$

$$V_{\pi}(s) = \sum_{s'} T(s'|s, \pi(s)) (R(s, a, s') + \gamma Q_{\pi}(s', \pi(s')))$$

$$Q_{*}(s, a) = \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{*}(s'))$$

$$V_{*}(s) = \max_a \sum_{s'} T(s'|s, a) \left( R(s, a, s') + \gamma \max_{a'} Q_{*}(s', a') \right)$$

# Q-value vs. Value

Can state  $Q$  in terms of  $V$ , and vice-versa.

$$Q_{\pi}(s, a) = \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{\pi}(s'))$$

$$V_{\pi}(s) = \sum_{s'} T(s'|s, \pi(s)) (R(s, a, s') + \gamma Q_{\pi}(s', \pi(s')))$$

$$Q_{*}(s, a) = \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_{*}(s'))$$

$$V_{*}(s) = \max_a \sum_{s'} T(s'|s, a) \left( R(s, a, s') + \gamma \max_{a'} Q_{*}(s', a') \right)$$

(exercise to the reader...)

# Motivation for Q-Value

- So far, we have been learning a policy  $\pi$ , and using  $\pi$  to compute  $V_\pi$ .



# Motivation for Q-Value

- So far, we have been learning a policy  $\pi$ , and using  $\pi$  to compute  $V_\pi$ .
- Even if we were given  $V_*$  directly, can't recover  $\pi_*$  without white-box access to  $T$  and  $R$  (environment).

$$\pi_*(s) = \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s'))$$

# Motivation for Q-Value

- So far, we have been learning a policy  $\pi$ , and using  $\pi$  to compute  $V_\pi$ .
- Even if we were given  $V_*$  directly, can't recover  $\pi_*$  without white-box access to  $T$  and  $R$  (environment).

$$\pi_*(s) = \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s'))$$

- However, given  $Q_*$ , we can directly recover  $\pi_*$

$$\pi_*(s) = \arg \max_a Q_*(s, a)$$

# Motivation for Q-Value

- So far, we have been learning a policy  $\pi$ , and using  $\pi$  to compute  $V_\pi$ .
- Even if we were given  $V_*$  directly, can't recover  $\pi_*$  without white-box access to  $T$  and  $R$  (environment).

$$\pi_*(s) = \arg \max_a \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V_*(s'))$$

- However, given  $Q_*$ , we can directly recover  $\pi_*$

$$\pi_*(s) = \arg \max_a Q_*(s, a)$$

- **Idea:** Learn  $Q_*$  instead, recover policy  $\pi_*$

# SARSA: On-Policy TD Control

Apply same argument as TD(0) to the Q-Value

$$Q_*(s, a) = \mathbb{E}_{s' \sim T(\cdot | s, a)} [R(s, a, s') + \gamma Q_*(s', \pi_*(s'))]$$

# SARSA: On-Policy TD Control

Apply same argument as TD(0) to the Q-Value

$$Q_*(s, a) = \mathbb{E}_{s' \sim T(\cdot | s, a)} [R(s, a, s') + \gamma Q_*(s', \pi_*(s'))]$$

On timestep  $t$ , this is “*on average*”, equal to the actual reward  $r_{t+1}$ , plus the discounted Q-value of the actual next state-action pair  $s_{t+1}, a_{t+1}$ .

$$Q_*(s_t, a_t) \approx r_{t+1} + \gamma Q_*(s_{t+1}, a_{t+1})$$

# SARSA: On-Policy TD Control

Apply same argument as TD(0) to the Q-Value

$$Q_*(s, a) = \mathbb{E}_{s' \sim T(\cdot | s, a)} [R(s, a, s') + \gamma Q_*(s', \pi_*(s'))]$$

On timestep  $t$ , this is “on average”, equal to the actual reward  $r_{t+1}$ , plus the discounted Q-value of the actual next state-action pair  $s_{t+1}, a_{t+1}$ .

$$Q_*(s_t, a_t) \approx r_{t+1} + \gamma Q_*(s_{t+1}, a_{t+1})$$

## SARSA Update Rule

$$\hat{Q}_*(s_t, a_t) \leftarrow \hat{Q}_*(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \hat{Q}_*(s_{t+1}, a_{t+1}) - \hat{Q}_*(s_t, a_t) \right)$$

where  $\alpha \in (0, 1]$  is the **learning rate**.

# SARSA: On-Policy TD Control

Apply same argument as TD(0) to the Q-Value

$$Q_*(s, a) = \mathbb{E}_{s' \sim T(\cdot|s,a)} [R(s, a, s') + \gamma Q_*(s', \pi_*(s'))]$$

On timestep  $t$ , this is “on average”, equal to the actual reward  $r_{t+1}$ , plus the discounted Q-value of the actual next state-action pair  $s_{t+1}, a_{t+1}$ .

$$Q_*(s_t, a_t) \approx r_{t+1} + \gamma Q_*(s_{t+1}, a_{t+1})$$

## SARSA Update Rule

$$\hat{Q}_*(s_t, a_t) \leftarrow \hat{Q}_*(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \hat{Q}_*(s_{t+1}, a_{t+1}) - \hat{Q}_*(s_t, a_t) \right)$$

where  $\alpha \in (0, 1]$  is the **learning rate**.

Actions drawn from  $\varepsilon$ -greedy strategy

$$\pi^{\varepsilon\text{-greedy}}(s) = \begin{cases} \text{do random stuff} & \text{prob } \varepsilon \\ \arg \max_a \hat{Q}_*(s, a) & \text{prob } 1 - \varepsilon \end{cases}$$

# SARSA: On-Policy TD Control

Apply same argument as TD(0) to the Q-Value

$$Q_*(s, a) = \mathbb{E}_{s' \sim T(\cdot|s, a)} [R(s, a, s') + \gamma Q_*(s', \pi_*(s'))]$$

On timestep  $t$ , this is “on average”, equal to the actual reward  $r_{t+1}$ , plus the discounted Q-value of the actual next state-action pair  $s_{t+1}, a_{t+1}$ .

$$Q_*(s_t, a_t) \approx r_{t+1} + \gamma Q_*(s_{t+1}, a_{t+1})$$

## SARSA Update Rule

$$\hat{Q}_*(s_t, a_t) \leftarrow \hat{Q}_*(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \hat{Q}_*(s_{t+1}, a_{t+1}) - \hat{Q}_*(s_t, a_t) \right)$$

where  $\alpha \in (0, 1]$  is the **learning rate**.

Actions drawn from  $\varepsilon$ -greedy strategy

$$\pi^{\varepsilon\text{-greedy}}(s) = \begin{cases} \text{do random stuff} & \text{prob } \varepsilon \\ \arg \max_a \hat{Q}_*(s, a) & \text{prob } 1 - \varepsilon \end{cases}$$

**Theorem:** Under “niceness” conditions SARSA guaranteed to converge to  $Q_*$ .



# Q-Learning: Off-Policy TD Control

- Why learn from  $a_{t+1}$  when it was a random exploration action? Why not instead learn from the action  $\arg \max_{a'} Q(s_{t+1}, a')$  that should have been taken?

# Q-Learning: Off-Policy TD Control

- Why learn from  $a_{t+1}$  when it was a random exploration action? Why not instead learn from the action  $\arg \max_{a'} Q(s_{t+1}, a')$  that should have been taken?

## Q-Learning Update Rule

$$\hat{Q}_*(s_t, a_t) \leftarrow \hat{Q}_*(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_{a'} \hat{Q}(s_{t+1}, a') - \hat{Q}(s_t, a_t) \right)$$

Actions taken via  $\varepsilon$ -greedy strategy over  $\hat{Q}_*(s, a)$ .

# Q-Learning: Off-Policy TD Control

- Why learn from  $a_{t+1}$  when it was a random exploration action? Why not instead learn from the action  $\arg \max_{a'} Q(s_{t+1}, a')$  that should have been taken?

## Q-Learning Update Rule

$$\hat{Q}_*(s_t, a_t) \leftarrow \hat{Q}_*(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_{a'} \hat{Q}(s_{t+1}, a') - \hat{Q}(s_t, a_t) \right)$$

Actions taken via  $\epsilon$ -greedy strategy over  $\hat{Q}_*(s, a)$ .

**Theorem:** Under “niceness” conditions Q-learning guaranteed to converge to  $Q_*$ .

# SARSA v.s. Q-Learning

## SARSA Update Rule

$$\hat{Q}_*(s_t, a_t) \leftarrow \hat{Q}_*(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \hat{Q}_*(s_{t+1}, a_{t+1}) - \hat{Q}_*(s_t, a_t) \right)$$

## Q-Learning Update Rule

$$\hat{Q}_*(s_t, a_t) \leftarrow \hat{Q}_*(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_{a'} \hat{Q}(s_{t+1}, a') - \hat{Q}(s_t, a_t) \right)$$

- Q-Learning (usually) tends to converge faster than SARSA, and chooses more aggressive/risky moves
- SARSA learns from the moves that were actually taken, including any exploration
- In “risky” environments, SARSA will learn to avoid getting near dangerous situations (to avoid accidentally taking a very bad exploratory move). Q-Learning will not.

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

Loop for each step of episode:

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A'$ ;

until  $S$  is terminal

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

Initialize  $S$

Loop for each step of episode:

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

until  $S$  is terminal

# Beyond Tabular Learning

- All the methods up to this point assume sweeping through all state-action pairs is tractable

# Beyond Tabular Learning

- All the methods up to this point assume sweeping through all state-action pairs is tractable
- What about large/continuous state spaces?

# Beyond Tabular Learning

- All the methods up to this point assume sweeping through all state-action pairs is tractable
- What about large/continuous state spaces?
  - State aggregation?



# Beyond Tabular Learning

- All the methods up to this point assume sweeping through all state-action pairs is tractable
- What about large/continuous state spaces?
  - State aggregation?
  - Parameterised policy  $\pi_\theta$ , learn best  $\theta$ ?

# Beyond Tabular Learning

- All the methods up to this point assume sweeping through all state-action pairs is tractable
- What about large/continuous state spaces?
  - State aggregation?
  - Parameterised policy  $\pi_\theta$ , learn best  $\theta$ ?
  - Craft a heuristic by hand?

# Beyond Tabular Learning

- All the methods up to this point assume sweeping through all state-action pairs is tractable
- What about large/continuous state spaces?
  - State aggregation?
  - Parameterised policy  $\pi_\theta$ , learn best  $\theta$ ?
  - Craft a heuristic by hand?
- In general, would like the agent to learn useful features for us

# Beyond Tabular Learning

- All the methods up to this point assume sweeping through all state-action pairs is tractable
- What about large/continuous state spaces?
  - State aggregation?
  - Parameterised policy  $\pi_\theta$ , learn best  $\theta$ ?
  - Craft a heuristic by hand?
- In general, would like the agent to learn useful features for us
  - Something deep learning excels at!

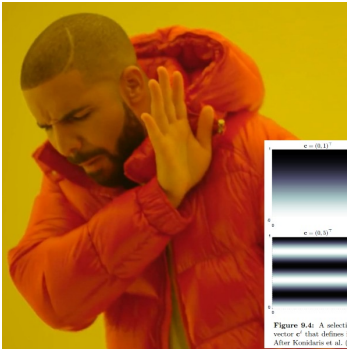


Figure 9.4: A selection of six vector  $e^i$  that defines it ( $s_1$  is 1). After Konradis et al. (2011).



9.5 Feature Construction for Linear Methods

- 9.5.1 Polynomials . . . . .
- 9.5.2 Fourier Basis . . . . .
- 9.5.3 Coarse Coding . . . . .
- 9.5.4 Tile Coding . . . . .
- 9.5.5 Radial Basis Functions . . . . .

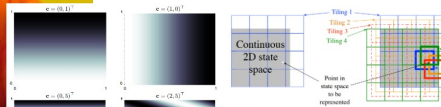


Figure 9.9: Multiple, overlapping grid-tilings on a limited two-dimension are offset from one another by a uniform amount in each dimension.

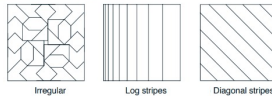
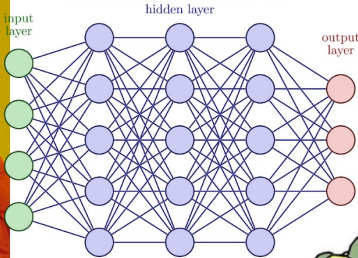


Figure 9.12: Tilings need not be grids. They can be arbitrarily shaped and non-still in many cases being computationally efficient to compute.



# Difficulties with using Neural Networks for RL

Neural networks expect to be trained in a supervised learning fashion, with batches of data fed in, loss computed, and gradients backpropagated.

**Idea:** Reduce the reinforcement learning problem to a supervised learning problem?

- Interaction with environment is *NOT* i.i.d

# Difficulties with using Neural Networks for RL

Neural networks expect to be trained in a supervised learning fashion, with batches of data fed in, loss computed, and gradients backpropagated.

**Idea:** Reduce the reinforcement learning problem to a supervised learning problem?

- Interaction with environment is *NOT* i.i.d
  - Collect many trajectories, dump into a buffer and shuffle

# Difficulties with using Neural Networks for RL

Neural networks expect to be trained in a supervised learning fashion, with batches of data fed in, loss computed, and gradients backpropagated.

**Idea:** Reduce the reinforcement learning problem to a supervised learning problem?

- Interaction with environment is *NOT* i.i.d
  - Collect many trajectories, dump into a buffer and shuffle
- Rewards are sparse



# Difficulties with using Neural Networks for RL

Neural networks expect to be trained in a supervised learning fashion, with batches of data fed in, loss computed, and gradients backpropagated.

**Idea:** Reduce the reinforcement learning problem to a supervised learning problem?

- Interaction with environment is *NOT* i.i.d
  - Collect many trajectories, dump into a buffer and shuffle
- Rewards are sparse
  - $\epsilon$ -greedy explore, hope for the best

# Difficulties with using Neural Networks for RL

Neural networks expect to be trained in a supervised learning fashion, with batches of data fed in, loss computed, and gradients backpropagated.

**Idea:** Reduce the reinforcement learning problem to a supervised learning problem?

- Interaction with environment is *NOT* i.i.d
  - Collect many trajectories, dump into a buffer and shuffle
- Rewards are sparse
  - $\epsilon$ -greedy explore, hope for the best
- No ground truth to compare against

# Difficulties with using Neural Networks for RL

Neural networks expect to be trained in a supervised learning fashion, with batches of data fed in, loss computed, and gradients backpropagated.

**Idea:** Reduce the reinforcement learning problem to a supervised learning problem?

- Interaction with environment is *NOT* i.i.d
  - Collect many trajectories, dump into a buffer and shuffle
- Rewards are sparse
  - $\epsilon$ -greedy explore, hope for the best
- No ground truth to compare against
  - Bootstrap from current estimates (i.e. Q-Learning)

# Deep Q-Networks (DQN)

- The Q-Value estimate  $\hat{Q}_*(s, a; \theta)$  is now stored as a network, with parameters  $\theta$ . Recall the TD-error for Q-Learning

$$\delta_t = r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta) - Q_*(s_t, a_t; \theta)$$

# Deep Q-Networks (DQN)

- The Q-Value estimate  $\hat{Q}_*(s, a; \theta)$  is now stored as a network, with parameters  $\theta$ . Recall the TD-error for Q-Learning

$$\delta_t = r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta) - Q_*(s_t, a_t; \theta)$$

- **Idea:** Accumulate experience  $(s^i, a^i, r^i, s_{\text{new}}^i)$  via interaction, optimise  $\theta$  to minimise loss  $L(\theta)$

# Deep Q-Networks (DQN)

- The Q-Value estimate  $\hat{Q}_*(s, a; \theta)$  is now stored as a network, with parameters  $\theta$ . Recall the TD-error for Q-Learning

$$\delta_t = r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta) - Q_*(s_t, a_t; \theta)$$

- **Idea:** Accumulate experience  $(s^i, a^i, r^i, s_{\text{new}}^i)$  via interaction, optimise  $\theta$  to minimise loss  $L(\theta)$ 
  - In practice, experience is accumulated in a buffer, and batches are sampled at random to make data “more i.i.d”

# Deep Q-Networks (DQN)

- The Q-Value estimate  $\hat{Q}_*(s, a; \theta)$  is now stored as a network, with parameters  $\theta$ . Recall the TD-error for Q-Learning

$$\delta_t = r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta) - Q_*(s_t, a_t; \theta)$$

- **Idea:** Accumulate experience  $(s^i, a^i, r^i, s_{\text{new}}^i)$  via interaction, optimise  $\theta$  to minimise loss  $L(\theta)$ 
  - In practice, experience is accumulated in a buffer, and batches are sampled at random to make data “more i.i.d”
  - Also use separate set of parameters  $\theta_{\text{target}}$  for the target network, copy weights every so often for stability

# Deep Q-Networks (DQN)

- The Q-Value estimate  $\hat{Q}_*(s, a; \theta)$  is now stored as a network, with parameters  $\theta$ . Recall the TD-error for Q-Learning

$$\delta_t = r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta) - Q_*(s_t, a_t; \theta)$$

- **Idea:** Accumulate experience  $(s^i, a^i, r^i, s_{\text{new}}^i)$  via interaction, optimise  $\theta$  to minimise loss  $L(\theta)$ 
  - In practice, experience is accumulated in a buffer, and batches are sampled at random to make data “more i.i.d”
  - Also use separate set of parameters  $\theta_{\text{target}}$  for the target network, copy weights every so often for stability

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \left( r^i + \gamma \max_{a'} Q_*(s_{\text{new}}, a'; \theta_{\text{target}}) - Q_*(s_t, a_t; \theta) \right)^2$$



# Deep Q-Networks (DQN)

- The Q-Value estimate  $\hat{Q}_*(s, a; \theta)$  is now stored as a network, with parameters  $\theta$ . Recall the TD-error for Q-Learning

$$\delta_t = r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta) - Q_*(s_t, a_t; \theta)$$

- **Idea:** Accumulate experience  $(s^i, a^i, r^i, s_{\text{new}}^i)$  via interaction, optimise  $\theta$  to minimise loss  $L(\theta)$ 
  - In practice, experience is accumulated in a buffer, and batches are sampled at random to make data “more i.i.d”
  - Also use separate set of parameters  $\theta_{\text{target}}$  for the target network, copy weights every so often for stability

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \left( r^i + \gamma \max_{a'} Q_*(s_{\text{new}}, a'; \theta_{\text{target}}) - Q_*(s_t, a_t; \theta) \right)^2$$

Then, perform gradient update step over parameters

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} L(\theta)$$

# Deep Q-Networks (DQN) for Episodic Environments

- Slightly modify the TD-error, depending if  $s_{t+1}$  is a terminal state.
- Assume environment returns  $(s_{t+1}, r_{t+1}, d_{t+1}) \sim p(\cdot | s_t, a_t)$ , where  $d_{t+1}$  (done) indicates if the episode ended on timestep  $t + 1$ .

$$\delta_t = y_t - Q_*(s_t, a_t; \theta)$$
$$y_t = \begin{cases} r_{t+1} & d_{t+1} = \text{True} \\ r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta_{\text{target}}) & d_{t+1} = \text{False} \end{cases}$$

# Deep Q-Networks (DQN) for Episodic Environments

- Slightly modify the TD-error, depending if  $s_{t+1}$  is a terminal state.
- Assume environment returns  $(s_{t+1}, r_{t+1}, d_{t+1}) \sim p(\cdot | s_t, a_t)$ , where  $d_{t+1}$  (done) indicates if the episode ended on timestep  $t + 1$ .

$$\delta_t = y_t - Q_*(s_t, a_t; \theta)$$
$$y_t = \begin{cases} r_{t+1} & d_{t+1} = \text{True} \\ r_{t+1} + \gamma \max_{a'} Q_*(s_{t+1}, a'; \theta_{\text{target}}) & d_{t+1} = \text{False} \end{cases}$$

Loss function is now

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N (y_t - Q_*(s_t, a_t; \theta))^2$$

---

**Algorithm 1** Deep Q-Learning with Replay Buffer

---

**Input:** Environment  $p$ , Number of episodes  $M$ , replay buffer size  $N$

```
1: Initialise replay buffer  $\mathcal{D}$  to capacity  $N$ 
2: for episode = 1 to  $M$  do
3:   Sample initial state  $s$  from environment
4:    $d \leftarrow \text{False}$ 
5:   Initialize target parameters  $\theta_{\text{target}} \leftarrow \theta$ 
6:   while  $d = \text{False}$  do
7:      $a \leftarrow \begin{cases} \text{random action} & \text{prob } \varepsilon \\ \arg \max_{a'} Q(s, a'; \theta) & \text{prob } 1 - \varepsilon \end{cases}$ 
8:     Sample  $(s_{\text{new}}, r, d) \sim p(\cdot | s, a)$ 
9:     Store experience  $(s, a, r, s_{\text{new}}, d)$  in  $\mathcal{D}$ 
10:     $s_{\text{new}} \leftarrow s$ 
11:    if Learning on this step then
12:      Sample minibatch  $B \leftarrow \{(s^i, a^i, r^i, s_{\text{new}}^i, d^i)\}_{i=1}^{|B|}$  from  $\mathcal{D}$ 
13:      for  $j = 1$  to  $|B|$  do
14:         $y^j \leftarrow \begin{cases} r^j & d^j = \text{True} \\ r^j + \gamma \max_{a'} Q(s_{\text{new}}^j, a'; \theta_{\text{target}}) & d^j = \text{False} \end{cases}$ 
15:      end for
16:      Define loss  $L(\theta) = \frac{1}{|B|} \sum_{i=1}^{|B|} (y^i - Q(s^i, a^i; \theta))^2$ 
17:      Gradient descent step  $\theta \leftarrow \theta - \eta \nabla_{\theta} L(\theta)$ 
18:    end if
19:    if Update target this step then
20:       $\theta_{\text{target}} \leftarrow \theta$ 
21:    end if
22:  end while
23: end for
```

- State space  $(x, v, \theta, \omega) \subseteq \mathbb{R}^4$ , representing
  - $-4.8 \leq x \leq 4.8$ , position of the cart (meters)
  - $-\infty \leq v \leq \infty$ , velocity of the cart (meters/second)
  - $-28^\circ \leq \theta \leq 28^\circ$ , angle of the pole (measured from vertical) (degrees)
  - $-\infty \leq \omega \leq \infty$ , angular velocity of the pole (degrees/second)
- Actions:  $\{L, R\}$  Apply a force of 10 newtons to the left/right of the cart
- Environment: Takes old state  $s_t = (x_t, v_t, \theta_t, \omega_t)$  and force  $a_t \in L, R$ , simulates the physics of the cartpole system using Euler's method in a 20ms timestep, returns the new state space  $s_{t+1} = (x_{t+1}, v_{t+1}, \theta_{t+1}, \omega_{t+1})$  and reward  $r_{t+1} = 1$
- Episode terminates if  $|x| \geq 2.4$  (the cart rolls off the track) or  $|\theta| \geq 12^\circ$  (the pole moves too far off vertical) or 500 timesteps (= 10 seconds) elapse.
- Initial state sampled uniformly from  $[-0.05, 0.05]^4$  (to avoid agent memorising a sequence of actions).
- Agent knows nothing about poles, or carts, or the laws of physics. Has to infer all of this from a vector of 4 numbers, and then determine a strategy to keep the cart centred and the pole upright

# SpinnyPole

- State space  $(x, v, \theta, \omega) \subseteq \mathbb{R}^4$ , representing
  - $-4.8 \leq x \leq 4.8$ , position of the cart (meters)
  - $-\infty \leq v \leq \infty$ , velocity of the cart (meters/second)
  - $-28^\circ \leq \theta \leq 28^\circ$ , angle of the pole (measured from vertical) (degrees)
  - $-\infty \leq \omega \leq \infty$ , angular velocity of the pole (degrees/second)
- Actions:  $\{L, R\}$  Apply a force of 10 newtons to the left/right of the cart
- Environment: Takes old state  $s_t = (x_t, v_t, \theta_t, \omega_t)$  and force  $a_t \in L, R$ , simulates the physics of the cartpole system using Euler's method in a 20ms timestep, returns the new state space  $s_{t+1} = (x_{t+1}, v_{t+1}, \theta_{t+1}, \omega_{t+1})$  and reward  $r_{t+1} = ???$
- Episode terminates if  $|x| \geq 2.4$  (the cart rolls off the track) or  ~~$|\theta| \geq 12^\circ$  (the pole moves too far off vertical)~~ or 1000 timesteps (=20 seconds) elapse.
- Initial state sampled uniformly from  $[-0.05, 0.05]^4$  (to avoid agent memorising a sequence of actions).
- Agent knows nothing about poles, or carts, or the laws of physics. Has to infer all of this from a vector of 4 numbers, and then determine a strategy to keep the cart centred and the pole upright **spin as fast as possible without moving off track**

# Vanilla Policy Gradient (VPG)

- Learn  $\pi$  directly.  $\pi$  is stochastic, push up (down) probability  $\pi(a|s)$  of good (bad) actions, converge to  $\pi^*$ .

# Vanilla Policy Gradient (VPG)

- Learn  $\pi$  directly.  $\pi$  is stochastic, push up (down) probability  $\pi(a|s)$  of good (bad) actions, converge to  $\pi^*$ .
- Policy  $\pi_{\theta}$  is parameterised by  $\theta$ , such that  $\nabla_{\theta}\pi_{\theta}$  exists



# Vanilla Policy Gradient (VPG)

- Learn  $\pi$  directly.  $\pi$  is stochastic, push up (down) probability  $\pi(a|s)$  of good (bad) actions, converge to  $\pi^*$ .
- Policy  $\pi_{\theta}$  is parameterised by  $\theta$ , such that  $\nabla_{\theta}\pi_{\theta}$  exists
- Measure of performance  $J(\theta)$  (gain)

# Vanilla Policy Gradient (VPG)

- Learn  $\pi$  directly.  $\pi$  is stochastic, push up (down) probability  $\pi(a|s)$  of good (bad) actions, converge to  $\pi^*$ .
- Policy  $\pi_{\theta}$  is parameterised by  $\theta$ , such that  $\nabla_{\theta}\pi_{\theta}$  exists
- Measure of performance  $J(\theta)$  (gain)
- Update step  $\theta \leftarrow \theta + \eta \widehat{\nabla_{\theta} J(\theta)}$

# Vanilla Policy Gradient (VPG)

- Learn  $\pi$  directly.  $\pi$  is stochastic, push up (down) probability  $\pi(a|s)$  of good (bad) actions, converge to  $\pi^*$ .
- Policy  $\pi_\theta$  is parameterised by  $\theta$ , such that  $\nabla_\theta \pi_\theta$  exists
- Measure of performance  $J(\theta)$  (gain)
- Update step  $\theta \leftarrow \theta + \eta \widehat{\nabla_\theta J(\theta)}$
- Learn preferences  $h(s, a, \theta)$ , and (assuming  $|\mathcal{A}|$  “small”) define softmax policy

$$\pi_\theta^{\text{softmax}}(a|s) = \frac{\exp(h(s, a, \theta)/T)}{\sum_{a'} \exp(h(s, a', \theta)/T)}$$

where  $T$  is temperature (hyperparameter).

- Use neural network to learn  $h(s, a, \theta)$

## • Advantages

- $\pi_{\epsilon\text{-greedy}}$  always does uniformly random moves when exploring.  $\pi_{\theta}^{\text{softmax}}$  is still stochastic, but biased towards good moves
- $\pi_{\theta}^{\text{softmax}}$  is continuous w.r.t preferences  $h(s, a, \theta)$ .  $\pi_{\epsilon\text{-greedy}}$  might dramatically change behaviour in response to small perturbations in  $\hat{Q}_* \equiv$  better convergence

## • Disadvantages

- More computationally expensive/more complex
- $\pi_{\theta}^{\text{softmax}}$  will play near uniform for two states with similar values.  $\pi_{\epsilon\text{-greedy}}$  will choose the best
- $\pi_{\theta}^{\text{softmax}}$  will only converge to deterministic policy with a temperature schedule (especially for states with similar value), hard to choose temperature scale a priori/requires domain knowledge

# Log-derivative trick

Note that

$$\frac{d}{dx} \log f(x) = \frac{1}{f(x)} \cdot \frac{d}{dx} f(x)$$

Hence, multiplying by  $f(x)$ ,

$$\frac{d}{dx} f(x) = f(x) \frac{d}{dx} \log f(x)$$

Or, in the form we will use it

$$\nabla_{\theta} P_{\theta}(x) = P_{\theta}(x) \nabla_{\theta} \log P_{\theta}(x)$$

# Policy Gradient Framework

- Assume episodic environment, length  $t'$ , no discount  $\gamma = 1$ . Assume fixed starting state  $s_0 = s_{\text{start}}$ .
- Define  $J(\theta) = V_{\pi_\theta}(s_{\text{start}})$ .
- Let  $\tau = s_{\text{start}}, a_0, r_1, s_1, \dots, s_{t'}$  denote a trajectory
- $G(\tau) = \sum_{t=0}^{t'} r_t$  is the undiscounted return for trajectory  $\tau$ .
- $\Pr(\tau|\theta) = \prod_{k=0}^{t'} \pi_\theta(a_k|s_k) T(s_{k+1}|s_k, a_k)$  is the probability of sampling  $\tau$  from environment given  $\theta$ .

# Policy Gradient Framework

- Assume episodic environment, length  $t'$ , no discount  $\gamma = 1$ . Assume fixed starting state  $s_0 = s_{\text{start}}$ .
- Define  $J(\theta) = V_{\pi_\theta}(s_{\text{start}})$ .
- Let  $\tau = s_{\text{start}}, a_0, r_1, s_1, \dots, s_{t'}$  denote a trajectory
- $G(\tau) = \sum_{t=0}^{t'} r_t$  is the undiscounted return for trajectory  $\tau$ .
- $\Pr(\tau|\theta) = \prod_{k=0}^{t'} \pi_\theta(a_k|s_k) T(s_{k+1}|s_k, a_k)$  is the probability of sampling  $\tau$  from environment given  $\theta$ .

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)]$$

# Policy Gradient Framework

- Assume episodic environment, length  $t'$ , no discount  $\gamma = 1$ . Assume fixed starting state  $s_0 = s_{\text{start}}$ .
- Define  $J(\theta) = V_{\pi_\theta}(s_{\text{start}})$ .
- Let  $\tau = s_{\text{start}}, a_0, r_1, s_1, \dots, s_{t'}$  denote a trajectory
- $G(\tau) = \sum_{t=0}^{t'} r_t$  is the undiscounted return for trajectory  $\tau$ .
- $\Pr(\tau|\theta) = \prod_{k=0}^{t'} \pi_\theta(a_k|s_k) T(s_{k+1}|s_k, a_k)$  is the probability of sampling  $\tau$  from environment given  $\theta$ .

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)] \\ &= \nabla_{\theta} \sum_{\tau} \Pr(\tau|\theta) G(\tau)\end{aligned}$$



# Policy Gradient Framework

- Assume episodic environment, length  $t'$ , no discount  $\gamma = 1$ . Assume fixed starting state  $s_0 = s_{\text{start}}$ .
- Define  $J(\theta) = V_{\pi_\theta}(s_{\text{start}})$ .
- Let  $\tau = s_{\text{start}}, a_0, r_1, s_1, \dots, s_{t'}$  denote a trajectory
- $G(\tau) = \sum_{t=0}^{t'} r_t$  is the undiscounted return for trajectory  $\tau$ .
- $\Pr(\tau|\theta) = \prod_{k=0}^{t'} \pi_\theta(a_k|s_k) T(s_{k+1}|s_k, a_k)$  is the probability of sampling  $\tau$  from environment given  $\theta$ .

$$\begin{aligned}\nabla_\theta J(\theta) &= \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \\ &= \nabla_\theta \sum_{\tau} \Pr(\tau|\theta) G(\tau) \\ &= \sum_{\tau} \nabla_\theta \Pr(\tau|\theta) G(\tau)\end{aligned}$$

# Policy Gradient Framework

- Assume episodic environment, length  $t'$ , no discount  $\gamma = 1$ . Assume fixed starting state  $s_0 = s_{\text{start}}$ .
- Define  $J(\theta) = V_{\pi_\theta}(s_{\text{start}})$ .
- Let  $\tau = s_{\text{start}}, a_0, r_1, s_1, \dots, s_{t'}$  denote a trajectory
- $G(\tau) = \sum_{t=0}^{t'} r_t$  is the undiscounted return for trajectory  $\tau$ .
- $\Pr(\tau|\theta) = \prod_{k=t}^{t'} \pi_\theta(a_k|s_k) T(s_{k+t}|s_k, a_k)$  is the probability of sampling  $\tau$  from environment given  $\theta$ .

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)] \\ &= \nabla_{\theta} \sum_{\tau} \Pr(\tau|\theta) G(\tau) \\ &= \sum_{\tau} \nabla_{\theta} \Pr(\tau|\theta) G(\tau) \\ &= \sum_{\tau} \Pr(\tau|\theta) (\nabla_{\theta} \log \Pr(\tau|\theta) G(\tau)) \quad (\text{Log Derivative trick})\end{aligned}$$

# Policy Gradient Framework

- Assume episodic environment, length  $t'$ , no discount  $\gamma = 1$ . Assume fixed starting state  $s_0 = s_{\text{start}}$ .
- Define  $J(\theta) = V_{\pi_\theta}(s_{\text{start}})$ .
- Let  $\tau = s_{\text{start}}, a_0, r_1, s_1, \dots, s_{t'}$  denote a trajectory
- $G(\tau) = \sum_{t=0}^{t'} r_t$  is the undiscounted return for trajectory  $\tau$ .
- $\Pr(\tau|\theta) = \prod_{k=t}^{t'} \pi_\theta(a_k|s_k) T(s_{k+t}|s_k, a_k)$  is the probability of sampling  $\tau$  from environment given  $\theta$ .

$$\begin{aligned}\nabla_\theta J(\theta) &= \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \\ &= \nabla_\theta \sum_{\tau} \Pr(\tau|\theta) G(\tau) \\ &= \sum_{\tau} \nabla_\theta \Pr(\tau|\theta) G(\tau) \\ &= \sum_{\tau} \Pr(\tau|\theta) (\nabla_\theta \log \Pr(\tau|\theta) G(\tau)) \quad (\text{Log Derivative trick}) \\ &= \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \Pr(\tau|\theta) G(\tau)]\end{aligned}$$

Note that

$$\nabla_{\theta} \log \Pr(\tau|\theta) = \nabla_{\theta} \log \prod_{k=t}^{t'} \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k)$$

Note that

$$\begin{aligned}\nabla_{\theta} \log \Pr(\tau|\theta) &= \nabla_{\theta} \log \prod_{k=t}^{t'} \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k)\end{aligned}$$

Note that

$$\begin{aligned}\nabla_{\theta} \log \Pr(\tau|\theta) &= \nabla_{\theta} \log \prod_{k=t}^{t'} \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) + \log T(s_{k+t}|s_k, a_k)\end{aligned}$$

Note that

$$\begin{aligned}\nabla_{\theta} \log \Pr(\tau|\theta) &= \nabla_{\theta} \log \prod_{k=t}^{t'} \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) + \log T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) + \cancel{\nabla_{\theta} \sum_{k=t}^{t'} \log T(s_{k+t}|s_k, a_k)}\end{aligned}$$

Note that

$$\begin{aligned}\nabla_{\theta} \log \Pr(\tau|\theta) &= \nabla_{\theta} \log \prod_{k=t}^{t'} \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) + \log T(s_{k+t}|s_k, a_k) \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) + \cancel{\nabla_{\theta} \sum_{k=t}^{t'} \log T(s_{k+t}|s_k, a_k)} \\ &= \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k)\end{aligned}$$



# Vanilla Policy Gradient (VPG)

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) G(\tau) \right]$$

## Clever trick 1: The future cannot affect the past

- $\pi_{\theta}(a_j | s_j)$  gets bumped by the full return  $G(\tau)$ . Obviously  $a_t$  has no effect on  $r_0, r_1, \dots, r_{t-1}$
- At timestep  $k$ , swap full return  $G(\tau)$  with partial return  $\sum_{j=k}^{t'} r_j$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) \sum_{j=k}^{t'} R(s_j, a_j, s_{j+1}) \right] \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \nabla_{\theta} \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) Q_{\pi_{\theta}}(s_k, a_k) \right] \end{aligned}$$

---

## Algorithm 1 Vanilla Policy Gradient ( $\gamma = 1$ )

---

**Input:** Environment  $p$ , Number of episodes  $M$

- 1: **for** episode = 1 to  $M$  **do**
  - 2:     Generate episode  $s_0, a_0, r_1, s_1, \dots, s_{T-1}, a_{T-1}, r_T$
  - 3:     Define  $G_t = \sum_{i=t+1}^T r_i$  for  $0 \leq t \leq T - 1$
  - 4:     Define gain  $J(\boldsymbol{\theta}) = \sum_{t=0}^{T-1} G_t \log \pi_{\boldsymbol{\theta}}(a_t | s_t)$
  - 5:     Gradient ascent step  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
  - 6: **end for**
-

---

## Algorithm 2 Efficient Vanilla Policy Gradient ( $\gamma = 1$ )

---

**Input:** Environment  $p$ , Number of episodes  $M$

- 1: **for** episode = 1 to  $M$  **do**
  - 2:     Generate episode  $s_0, a_0, r_1, s_1, \dots, s_{T-1}, a_{T-1}, r_T$
  - 3:     Initialise array  $G = \{G_0, G_1, \dots, G_{T-1}\}$
  - 4:      $G_{T-1} \leftarrow r_T$
  - 5:     **for** timestep in episode  $t = T - 2, T - 1$  to 0 **do**
  - 6:          $G_t \leftarrow r_{t+1} + G_{t+1}$
  - 7:     **end for**
  - 8:     Define gain  $J(\theta) = \sum_{t=0}^{T-1} G_t \log \pi_{\theta}(a_t | s_t)$
  - 9:     Gradient ascent step  $\theta \leftarrow \theta + \eta \nabla_{\theta} J(\theta)$
  - 10: **end for**
-

# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

Proof:

$$\sum_x P_\theta(x) = 1$$

# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

Proof:

$$\sum_x P_\theta(x) = 1$$

$$\nabla_\theta \sum_x P_\theta(x) = \nabla_\theta 1 = 0$$

# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

Proof:

$$\sum_x P_\theta(x) = 1$$

$$\nabla_\theta \sum_x P_\theta(x) = \nabla_\theta 1 = 0$$

$$\nabla_\theta \sum_x P_\theta(x) = 0$$

# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

Proof:

$$\sum_x P_\theta(x) = 1$$

$$\nabla_\theta \sum_x P_\theta(x) = \nabla_\theta 1 = 0$$

$$\nabla_\theta \sum_x P_\theta(x) = 0$$

$$\sum_x \nabla_\theta P_\theta(x) = 0$$



# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

Proof:

$$\sum_x P_\theta(x) = 1$$

$$\nabla_\theta \sum_x P_\theta(x) = \nabla_\theta 1 = 0$$

$$\nabla_\theta \sum_x P_\theta(x) = 0$$

$$\sum_x \nabla_\theta P_\theta(x) = 0$$

Apply log-derivative trick

# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

Proof:

$$\sum_x P_\theta(x) = 1$$

$$\nabla_\theta \sum_x P_\theta(x) = \nabla_\theta 1 = 0$$

$$\nabla_\theta \sum_x P_\theta(x) = 0$$

$$\sum_x \nabla_\theta P_\theta(x) = 0$$

Apply log-derivative trick

$$\sum_x P_\theta(x) \nabla_\theta P_\theta(x) = 0$$

# Expected Grad-Log-Prob (EGLP) Lemma

Let  $P_\theta$  be a parameterised probability distribution over random variable  $x$ . Then

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

Proof:

$$\sum_x P_\theta(x) = 1$$

$$\nabla_\theta \sum_x P_\theta(x) = \nabla_\theta 1 = 0$$

$$\nabla_\theta \sum_x P_\theta(x) = 0$$

$$\sum_x \nabla_\theta P_\theta(x) = 0$$

Apply log-derivative trick

$$\sum_x P_\theta(x) \nabla_\theta P_\theta(x) = 0$$

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$$

**Corollary of EGLP:** For any function  $b$  that depends only on state  $s_t$ ,

$$\mathbb{E}_{a_t \sim \pi_\theta} [\nabla_{\theta} \log \pi_\theta(a_t | s_t) b(s_t)] = 0$$

**Corollary of EGLP:** For any function  $b$  that depends only on state  $s_t$ ,

$$\mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0$$

So, can add/subtract any such **baseline function**  $b$  into VPG without changing the result (in expectation),

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) \left( Q_{\pi_\theta}(s_k, a_k) - b(s_k) \right) \right]$$

**Corollary of EGLP:** For any function  $b$  that depends only on state  $s_t$ ,

$$\mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0$$

So, can add/subtract any such **baseline function**  $b$  into VPG without changing the result (in expectation),

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) \left( Q_{\pi_\theta}(s_k, a_k) - b(s_k) \right) \right]$$

**Clever trick 2:** Choose  $b(s_t) = V_{\pi_\theta}(s_t)$ , the on-policy value function

$$\begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) \left( Q_{\pi_\theta}(s_k, a_k) - V_{\pi_\theta}(s_k) \right) \right] \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) A_{\pi_\theta}(s_t, a_t) \right] \end{aligned}$$

where  $A_\pi(s, a) := Q_\pi(s, a) - V_\pi(s)$  is the **advantage** function

**Corollary of EGLP:** For any function  $b$  that depends only on state  $s_t$ ,

$$\mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0$$

So, can add/subtract any such **baseline function**  $b$  into VPG without changing the result (in expectation),

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) \left( Q_{\pi_\theta}(s_k, a_k) - b(s_k) \right) \right]$$

**Clever trick 2:** Choose  $b(s_t) = V_{\pi_\theta}(s_t)$ , the on-policy value function

$$\begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) \left( Q_{\pi_\theta}(s_k, a_k) - V_{\pi_\theta}(s_k) \right) \right] \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) A_{\pi_\theta}(s_t, a_t) \right] \end{aligned}$$

where  $A_\pi(s, a) := Q_\pi(s, a) - V_\pi(s)$  is the **advantage** function

- $V_{\pi_\theta}$  learned by separate critic network.

**Corollary of EGLP:** For any function  $b$  that depends only on state  $s_t$ ,

$$\mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0$$

So, can add/subtract any such **baseline function**  $b$  into VPG without changing the result (in expectation),

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) \left( Q_{\pi_\theta}(s_k, a_k) - b(s_k) \right) \right]$$

**Clever trick 2:** Choose  $b(s_t) = V_{\pi_\theta}(s_t)$ , the on-policy value function

$$\begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) \left( Q_{\pi_\theta}(s_k, a_k) - V_{\pi_\theta}(s_k) \right) \right] \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \sum_{k=t}^{t'} \log \pi_\theta(a_k | s_k) A_{\pi_\theta}(s_t, a_t) \right] \end{aligned}$$

where  $A_\pi(s, a) := Q_\pi(s, a) - V_\pi(s)$  is the **advantage** function

- $V_{\pi_\theta}$  learned by separate critic network.
- Reduces variance, only update policy when critic disagrees



# WARNING

Everything beyond this point, I am less certain about. Where I make my best guess, or am uncertain, I mark it with  $\neg\_(\_)\_/$ .

# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

- Don't blame me, the PPO paper use this convention too!

# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

- Don't blame me, the PPO paper use this convention too!

Define policy gradient "loss" (gain?)

$$L^{PG}(\theta) = \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) A_{\pi_{\theta}}(s_t, a_t) \right]$$

where  $\hat{\mathbb{E}}$  indicates the expectation is approximated by a batch of samples, and  $\hat{A}(s_t, a_t) = \hat{Q}(s_t, a_t) - \hat{V}_{\phi}(s_t)$ , where

# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

- Don't blame me, the PPO paper use this convention too!

Define policy gradient "loss" (gain?)

$$L^{PG}(\theta) = \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) A_{\pi_{\theta}}(s_t, a_t) \right]$$

where  $\hat{\mathbb{E}}$  indicates the expectation is approximated by a batch of samples, and  $\hat{A}(s_t, a_t) = \hat{Q}(s_t, a_t) - \hat{V}_{\phi}(s_t)$ , where

- $Q(s_t, a_t) = \sum_{k=t}^{t'} R(s_t, a_t, s_{t+1})$  Q-value computed using empirical return

# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

- Don't blame me, the PPO paper use this convention too!

Define policy gradient "loss" (gain?)

$$L^{PG}(\theta) = \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) A_{\pi_{\theta}}(s_t, a_t) \right]$$

where  $\hat{\mathbb{E}}$  indicates the expectation is approximated by a batch of samples, and  $\hat{A}(s_t, a_t) = \hat{Q}(s_t, a_t) - \hat{V}_{\phi}(s_t)$ , where

- $Q(s_t, a_t) = \sum_{k=t}^{t'} R(s_t, a_t, s_{t+1})$  Q-value computed using empirical return
- $\hat{V}_{\phi}(s_t)$  computed using critic network

# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

- Don't blame me, the PPO paper use this convention too!

Define policy gradient "loss" (gain?)

$$L^{PG}(\theta) = \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) A_{\pi_{\theta}}(s_t, a_t) \right]$$

where  $\hat{\mathbb{E}}$  indicates the expectation is approximated by a batch of samples, and  $\hat{A}(s_t, a_t) = \hat{Q}(s_t, a_t) - \hat{V}_{\phi}(s_t)$ , where

- $Q(s_t, a_t) = \sum_{k=t}^{t'} R(s_t, a_t, s_{t+1})$  Q-value computed using empirical return
- $\hat{V}_{\phi}(s_t)$  computed using critic network
- Note that  $\hat{A}(s_t, a_t)$  has no dependence on  $\theta$ .

# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

- Don't blame me, the PPO paper use this convention too!

Define policy gradient "loss" (gain?)

$$L^{PG}(\theta) = \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) A_{\pi_{\theta}}(s_t, a_t) \right]$$

where  $\hat{\mathbb{E}}$  indicates the expectation is approximated by a batch of samples, and  $\hat{A}(s_t, a_t) = \hat{Q}(s_t, a_t) - \hat{V}_{\phi}(s_t)$ , where

- $Q(s_t, a_t) = \sum_{k=t}^{t'} R(s_t, a_t, s_{t+1})$  Q-value computed using empirical return
- $\hat{V}_{\phi}(s_t)$  computed using critic network
- Note that  $\hat{A}(s_t, a_t)$  has no dependence on  $\theta$ .
- However, this leads to destructively large policy updates



# Empirical Policy Gradient

**Note:** Policy gradient uses gradient **ascent**, so we actually **maximise** loss!

- Don't blame me, the PPO paper use this convention too!

Define policy gradient "loss" (gain?)

$$L^{PG}(\theta) = \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k | s_k) A_{\pi_{\theta}}(s_t, a_t) \right]$$

where  $\hat{\mathbb{E}}$  indicates the expectation is approximated by a batch of samples, and  $\hat{A}(s_t, a_t) = \hat{Q}(s_t, a_t) - \hat{V}_{\phi}(s_t)$ , where

- $Q(s_t, a_t) = \sum_{k=t}^{t'} R(s_t, a_t, s_{t+1})$  Q-value computed using empirical return
- $\hat{V}_{\phi}(s_t)$  computed using critic network
- Note that  $\hat{A}(s_t, a_t)$  has no dependence on  $\theta$ .
- However, this leads to destructively large policy updates

# Importance Sampling

**Main Idea:** Use samples from one distribution to estimate the expected value of a function under a different distribution.

In RL, policy  $\pi$  being learned about is **target policy** (usually  $\pi_*$ ), policy generating behaviour  $\beta$  is **behaviour policy**.

# Importance Sampling

**Main Idea:** Use samples from one distribution to estimate the expected value of a function under a different distribution.

In RL, policy  $\pi$  being learned about is **target policy** (usually  $\pi_*$ ), policy generating behaviour  $\beta$  is **behaviour policy**.

**On-Policy:** target=behaviour

# Importance Sampling

**Main Idea:** Use samples from one distribution to estimate the expected value of a function under a different distribution.

In RL, policy  $\pi$  being learned about is **target policy** (usually  $\pi_*$ ), policy generating behaviour  $\beta$  is **behaviour policy**.

**On-Policy:** target=behaviour

- SARSA: Target policy  $\pi_*^\epsilon$ , behaviour policy  $\pi_*^\epsilon$  (on-policy)

# Importance Sampling

**Main Idea:** Use samples from one distribution to estimate the expected value of a function under a different distribution.

In RL, policy  $\pi$  being learned about is **target policy** (usually  $\pi_*$ ), policy generating behaviour  $\beta$  is **behaviour policy**.

**On-Policy:** target=behaviour

- SARSA: Target policy  $\pi_*^\epsilon$ , behaviour policy  $\pi_*^\epsilon$  (on-policy)
- Q-Learning: Target policy  $\pi_*$ , behaviour policy  $\pi_*^\epsilon$  (off-policy)

# Importance Sampling

**Main Idea:** Use samples from one distribution to estimate the expected value of a function under a different distribution.

In RL, policy  $\pi$  being learned about is **target policy** (usually  $\pi_*$ ), policy generating behaviour  $\beta$  is **behaviour policy**.

**On-Policy:** target=behaviour

- SARSA: Target policy  $\pi_*^\epsilon$ , behaviour policy  $\pi_*^\epsilon$  (on-policy)
- Q-Learning: Target policy  $\pi_*$ , behaviour policy  $\pi_*^\epsilon$  (off-policy)

If  $\pi$  is very different from  $\beta$ , high variance, bad learning.

# Importance Sampling

Given starting state  $s_t$ , the probability of a particular state-action trajectory from timestep  $t$  to  $t'$

$$\tau = a_t, s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots, a_{t'-1}, s_{t'}$$

is

$$\begin{aligned}\Pr(\tau | s_t, a_{t:t'-1} \sim \pi) &= \pi(a_t | s_t) T(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1}) \dots T(s_{t'} | s_{t'-1}, a_{t'-1}) \\ &= \prod_{k=t}^{t'-1} \pi(a_k | s_k) T(s_{k+1} | s_k, a_k)\end{aligned}$$

# Importance Sampling

Given starting state  $s_t$ , the probability of a particular state-action trajectory from timestep  $t$  to  $t'$

$$\tau = a_t, s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots, a_{t'-1}, s_{t'}$$

is

$$\begin{aligned}\Pr(\tau | s_t, a_{t:t'-1} \sim \pi) &= \pi(a_t | s_t) T(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1}) \dots T(s_{t'} | s_{t'-1}, a_{t'-1}) \\ &= \prod_{k=t}^{t'-1} \pi(a_k | s_k) T(s_{k+1} | s_k, a_k)\end{aligned}$$

**Importance-sampling ratio:**  $\rho_{t:t'-1}$  The ratio of the likelihood of the trajectory under target and behaviour policies.

$$\rho_{t:t'-1} = \frac{\prod_{k=t} \pi(a_k | s_k) T(s_{k+1} | s_k, a_k)}{\prod_{k=t} \beta(a_k | s_k) T(s_{k+1} | s_k, a_k)} = \frac{\prod_{k=t} \pi(a_k | s_k)}{\prod_{k=t} \beta(a_k | s_k)}$$

- No dependency on environment distribution  $T$ !



# Importance Sampling

Want to estimate  $V_\pi$ , but only have returns  $G_t^\beta$  obtained from  $\beta$ .  $G_t^\beta$  has the wrong expectation

$$\mathbb{E}[G_t^\beta | s_t = s] = V_\beta(s)$$

# Importance Sampling

Want to estimate  $V_\pi$ , but only have returns  $G_t^\beta$  obtained from  $\beta$ .  $G_t^\beta$  has the wrong expectation

$$\mathbb{E}[G_t^\beta | s_t = s] = V_\beta(s)$$

Transform with the importance sampling ratio!

$$\mathbb{E}[\rho_{t:t'-1} G_t^\beta | s_t = s] = V_\pi(s)$$

# Importance Sampling

Want to estimate  $V_\pi$ , but only have returns  $G_t^\beta$  obtained from  $\beta$ .  $G_t^\beta$  has the wrong expectation

$$\mathbb{E}[G_t^\beta | s_t = s] = V_\beta(s)$$

Transform with the importance sampling ratio!

$$\mathbb{E}[\rho_{t:t'-1} G_t^\beta | s_t = s] = V_\pi(s)$$

( $\backslash$ \_ $\wedge$ ) $\_$  $\wedge$ ) During policy gradient, data is sampled from  $\pi_{\theta_{old}}$ , but target is  $\pi_\theta$ , so we would rather optimise

$$\hat{\mathbb{E}} \left[ \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}(s_t, a_t) \right]$$

called the **surrogate** objective.

# Justifying the surrogate objective

- Take  $L^{PG}(\theta)$ , and subtract out  $\log \pi_{\theta_{old}}(a_t|s_t)\hat{A}_t(s_t, a_t)$  (no dependence on  $\theta$ , maximising  $\theta$  is unchanged)

# Justifying the surrogate objective

- Take  $L^{PG}(\theta)$ , and subtract out  $\log \pi_{\theta_{old}}(a_t|s_t)\hat{A}_t(s_t, a_t)$  (no dependence on  $\theta$ , maximising  $\theta$  is unchanged)

$$\begin{aligned} & \arg \max_{\theta} L^{PG}(\theta) \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) A_{\pi_{\theta}}(s_t, a_t) - \log \pi_{\theta_{old}}(a_t|s_t) \hat{A}_t(s_t, a_t) \right] \end{aligned}$$

# Justifying the surrogate objective

- Take  $L^{PG}(\theta)$ , and subtract out  $\log \pi_{\theta_{old}}(a_t|s_t)\hat{A}_t(s_t, a_t)$  (no dependence on  $\theta$ , maximising  $\theta$  is unchanged)

$$\begin{aligned} & \arg \max_{\theta} L^{PG}(\theta) \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) A_{\pi_{\theta}}(s_t, a_t) - \log \pi_{\theta_{old}}(a_t|s_t) \hat{A}_t(s_t, a_t) \right] \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \frac{\pi_{\theta}(a_k|s_k)}{\pi_{\theta_{old}}(a_k|s_k)} \hat{A}_{\pi_{\theta}}(s_t, a_t) \right] \end{aligned}$$

# Justifying the surrogate objective

- Take  $L^{PG}(\theta)$ , and subtract out  $\log \pi_{\theta_{old}}(a_t|s_t)\hat{A}_t(s_t, a_t)$  (no dependence on  $\theta$ , maximising  $\theta$  is unchanged)

$$\begin{aligned} & \arg \max_{\theta} L^{PG}(\theta) \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) A_{\pi_{\theta}}(s_t, a_t) - \log \pi_{\theta_{old}}(a_t|s_t) \hat{A}_t(s_t, a_t) \right] \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \frac{\pi_{\theta}(a_k|s_k)}{\pi_{\theta_{old}}(a_k|s_k)} \hat{A}_{\pi_{\theta}}(s_t, a_t) \right] \end{aligned}$$

log is monotonic/Jensens theorem/idk  $\neg \setminus (\sphericalangle) \_ / \sphericalangle$

# Justifying the surrogate objective

- Take  $L^{PG}(\theta)$ , and subtract out  $\log \pi_{\theta_{old}}(a_t|s_t)\hat{A}_t(s_t, a_t)$  (no dependence on  $\theta$ , maximising  $\theta$  is unchanged)

$$\begin{aligned} & \arg \max_{\theta} L^{PG}(\theta) \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) A_{\pi_{\theta}}(s_t, a_t) - \log \pi_{\theta_{old}}(a_t|s_t) \hat{A}_t(s_t, a_t) \right] \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \frac{\pi_{\theta}(a_k|s_k)}{\pi_{\theta_{old}}(a_k|s_k)} \hat{A}_{\pi_{\theta}}(s_t, a_t) \right] \end{aligned}$$

log is monotonic/Jensens theorem/idk  $\neg \setminus (\setminus) \setminus /$

$$= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \frac{\pi_{\theta}(a_k|s_k)}{\pi_{\theta_{old}}(a_k|s_k)} \hat{A}_{\pi_{\theta}}(s_t, a_t) \right]$$



# Justifying the surrogate objective

- Take  $L^{PG}(\theta)$ , and subtract out  $\log \pi_{\theta_{old}}(a_t|s_t)\hat{A}_t(s_t, a_t)$  (no dependence on  $\theta$ , maximising  $\theta$  is unchanged)

$$\begin{aligned} & \arg \max_{\theta} L^{PG}(\theta) \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \pi_{\theta}(a_k|s_k) A_{\pi_{\theta}}(s_t, a_t) - \log \pi_{\theta_{old}}(a_t|s_t) \hat{A}_t(s_t, a_t) \right] \\ &= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \log \frac{\pi_{\theta}(a_k|s_k)}{\pi_{\theta_{old}}(a_k|s_k)} \hat{A}_{\pi_{\theta}}(s_t, a_t) \right] \end{aligned}$$

log is monotonic / Jensens theorem / idk  $\setminus \setminus (\setminus) \setminus /$

$$= \arg \max_{\theta} \hat{\mathbb{E}} \left[ \sum_{k=t}^{t'} \frac{\pi_{\theta}(a_k|s_k)}{\pi_{\theta_{old}}(a_k|s_k)} \hat{A}_{\pi_{\theta}}(s_t, a_t) \right]$$

- Learn a policy  $\pi_{\theta}$  (actor) and a value  $V_{\phi}(s)$  (critic). Actor acts, critic critiques.

# Trust Region Policy Optimisation (TRPO)

(Drop summations, write  $\hat{A}_t \equiv \hat{A}(s_t, a_t)$  for brevity).

The goal is maximisation of  $L^{CPI}(\theta)$  w.r.t  $\theta$  defined as

$$L^{CPI}(\theta) = \mathbb{E} \left[ \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \hat{A}_t \right]$$

subject to the constraint

$$\hat{\mathbb{E}} \left[ \text{KL}[\pi_{\theta_{\text{old}}}(\cdot|s_t) \parallel \pi_{\theta}(\cdot|s_t)] \right] \leq \delta$$

to avoid the two distributions changing too much.

# Trust Region Policy Optimisation (TRPO)

(Drop summations, write  $\hat{A}_t \equiv \hat{A}(s_t, a_t)$  for brevity).

The goal is maximisation of  $L^{CPI}(\theta)$  w.r.t  $\theta$  defined as

$$L^{CPI}(\theta) = \mathbb{E} \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right]$$

subject to the constraint

$$\hat{\mathbb{E}} \left[ \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t) \parallel \pi_{\theta}(\cdot | s_t)] \right] \leq \delta$$

to avoid the two distributions changing too much.

Here,  $\text{KL}(p \parallel q) := \sum_x p(x) \log \frac{p(x)}{q(x)}$  is the **Kullback-Liebler divergence**, or KL-divergence, that measures the “distance” between two probability distributions.

# Trust Region Policy Optimisation (TRPO)

(Drop summations, write  $\hat{A}_t \equiv \hat{A}(s_t, a_t)$  for brevity).

The goal is maximisation of  $L^{CPI}(\theta)$  w.r.t  $\theta$  defined as

$$L^{CPI}(\theta) = \mathbb{E} \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right]$$

subject to the constraint

$$\hat{\mathbb{E}} \left[ \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t) \parallel \pi_{\theta}(\cdot | s_t)] \right] \leq \delta$$

to avoid the two distributions changing too much.

Here,  $\text{KL}(p \parallel q) := \sum_x p(x) \log \frac{p(x)}{q(x)}$  is the **Kullback-Liebler divergence**, or KL-divergence, that measures the “distance” between two probability distributions. Constrained optimisation is problematic to deal with, but unconstrained optimisation with a KL-penalty

$$\underset{\theta}{\text{maximise}} \mathbb{E} \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t - \beta \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t) \parallel \pi_{\theta}(\cdot | s_t)] \right]$$

requires an additional hyperparameter  $\beta$ . Via experimentation, could not find a single  $\beta$  suitable for many different environments.

Instead, allow for unconstrained optimisation, but “clip” the result, so the policy can't drift too far.

Instead, allow for unconstrained optimisation, but “clip” the result, so the policy can't drift too far. Letting  $r_t(\boldsymbol{\theta}) = \frac{\pi_{\boldsymbol{\theta}}(a_t|s_t)}{\pi_{\boldsymbol{\theta}_{\text{old}}}(a_t|s_t)}$  denote the **probability ratio**, TRPO maximises

$$L^{CPI}(\boldsymbol{\theta}) = \mathbb{E} \left[ \frac{\pi_{\boldsymbol{\theta}}(a_t|s_t)}{\pi_{\boldsymbol{\theta}_{\text{old}}}(a_t|s_t)} \hat{A}(s_t, a_t) \right] = \mathbb{E} \left[ r_t(\boldsymbol{\theta}) \hat{A}(s_t, a_t) \right]$$

Instead, allow for unconstrained optimisation, but “clip” the result, so the policy can't drift too far. Letting  $r_t(\boldsymbol{\theta}) = \frac{\pi_{\boldsymbol{\theta}}(a_t|s_t)}{\pi_{\boldsymbol{\theta}_{\text{old}}}(a_t|s_t)}$  denote the **probability ratio**, TRPO maximises

$$L^{CPI}(\boldsymbol{\theta}) = \mathbb{E} \left[ \frac{\pi_{\boldsymbol{\theta}}(a_t|s_t)}{\pi_{\boldsymbol{\theta}_{\text{old}}}(a_t|s_t)} \hat{A}(s_t, a_t) \right] = \mathbb{E} \left[ r_t(\boldsymbol{\theta}) \hat{A}(s_t, a_t) \right]$$

We define the clip “loss” as

$$L^{CLIP}(\boldsymbol{\theta}) = \hat{E} \left[ \min(r_t(\boldsymbol{\theta}) \hat{A}(s), \text{clip}(r_t(\boldsymbol{\theta}), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

for hyperparameter  $\epsilon = 0.2$ .



Instead, allow for unconstrained optimisation, but “clip” the result, so the policy can’t drift too far. Letting  $r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$  denote the **probability ratio**, TRPO maximises

$$L^{CPI}(\theta) = \mathbb{E} \left[ \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \hat{A}(s_t, a_t) \right] = \mathbb{E} \left[ r_t(\theta) \hat{A}(s_t, a_t) \right]$$

We define the clip “loss” as

$$L^{CLIP}(\theta) = \hat{E} \left[ \min(r_t(\theta) \hat{A}(s), \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

for hyperparameter  $\epsilon = 0.2$ .

**Intuition:** Clip the ratio  $r_t(\theta)$  inside  $[1 - \epsilon, 1 + \epsilon]$ , then take the min of the clipped and unclipped to get a lower bound (pessimistic) on the unclipped objective.

The critic is simply trained against the returns from the environment

$$L_t^{VF}(\theta) = \frac{1}{2}(V_\theta(s_t) - G_t)^2$$

The critic is simply trained against the returns from the environment

$$L_t^{VF}(\theta) = \frac{1}{2}(V_\theta(s_t) - G_t)^2$$

We also add an **entropy bonus** to incentivise exploration by increasing the **entropy** of the distribution. The entropy  $H_\pi(s)$  of a policy  $\pi$  in state  $s$  is defined as

$$H_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{1}{\pi(a|s)}$$

Entropy can be thought of as a measure of how “random” the distribution is.

The critic is simply trained against the returns from the environment

$$L_t^{VF}(\theta) = \frac{1}{2}(V_\theta(s_t) - G_t)^2$$

We also add an **entropy bonus** to incentivise exploration by increasing the **entropy** of the distribution. The entropy  $H_\pi(s)$  of a policy  $\pi$  in state  $s$  is defined as

$$H_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{1}{\pi(a|s)}$$

Entropy can be thought of as a measure of how “random” the distribution is. Combine them all, with hyperparameters  $c_1, c_2$ .

$$L_t(\theta)^{CLIP+VF+S} = \hat{\mathbb{E}}_t[L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 H_{\pi_\theta}(s_t)]$$

The critic is simply trained against the returns from the environment

$$L_t^{VF}(\theta) = \frac{1}{2}(V_\theta(s_t) - G_t)^2$$

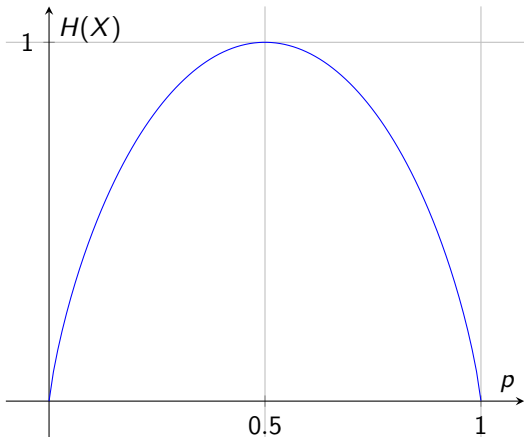
We also add an **entropy bonus** to incentivise exploration by increasing the **entropy** of the distribution. The entropy  $H_\pi(s)$  of a policy  $\pi$  in state  $s$  is defined as

$$H_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{1}{\pi(a|s)}$$

Entropy can be thought of as a measure of how “random” the distribution is. Combine them all, with hyperparameters  $c_1, c_2$ .

$$L_t(\theta)^{CLIP+VF+S}(\theta) = \hat{\mathbb{E}}_t[L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 H_{\pi_\theta}(s_t)]$$

Maximise  $L_t(\theta)^{CLIP+VF+S}(\theta)$  w.r.t  $\theta$



$$H(X) = -p \log_2 p - (1-p) \log_2 (1-p)$$

Figure: The entropy of a policy over two actions with  $\pi(a|s) = p$

# TD( $\lambda$ ) using Eligibility Traces

## TD(0) update

$$\hat{V}_\pi(s_t) \leftarrow \hat{V}_\pi(s_t) + \alpha \left( r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t) \right)$$

- Only provides an update based on the most recent state

# TD( $\lambda$ ) using Eligibility Traces

## TD(0) update

$$\hat{V}_\pi(s_t) \leftarrow \hat{V}_\pi(s_t) + \alpha \left( r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t) \right)$$

- Only provides an update based on the most recent state
- What if the pivotal action was taken far in the past, that lead to this desirable state?



# TD( $\lambda$ ) using Eligibility Traces

## TD(0) update

$$\hat{V}_\pi(s_t) \leftarrow \hat{V}_\pi(s_t) + \alpha \left( r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t) \right)$$

- Only provides an update based on the most recent state
- What if the pivotal action was taken far in the past, that lead to this desirable state?

One solution is to keep track of the **Eligibility Trace**, the number of times a state has been visited, discounted geometrically via a parameter  $\lambda$ , called the **trace decay**, and by  $\gamma$ , the **discount rate**.

$$E^0(s) := 0$$

$$E^t(s) := \gamma \lambda E^{t-1}(s) + \delta_{s,s_t}$$

# TD( $\lambda$ ) using Eligibility Traces

## TD(0) update

$$\hat{V}_\pi(s_t) \leftarrow \hat{V}_\pi(s_t) + \alpha \left( r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t) \right)$$

- Only provides an update based on the most recent state
- What if the pivotal action was taken far in the past, that lead to this desirable state?

One solution is to keep track of the **Eligibility Trace**, the number of times a state has been visited, discounted geometrically via a parameter  $\lambda$ , called the **trace decay**, and by  $\gamma$ , the **discount rate**.

$$E^0(s) := 0$$

$$E^t(s) := \gamma \lambda E^{t-1}(s) + \delta_{s,s_t}$$

**Motivation:** States that are more recent/have been

The discounting allows for more recent visits to contribute more to the count than past visits (which may be valuable for non-stationary environments.)

# TD( $\lambda$ ) using Eligibility Traces

## TD(0) update

$$\hat{V}_\pi(s_t) \leftarrow \hat{V}_\pi(s_t) + \alpha \left( r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t) \right)$$

- Only provides an update based on the most recent state
- What if the pivotal action was taken far in the past, that lead to this desirable state?

One solution is to keep track of the **Eligibility Trace**, the number of times a state has been visited, discounted geometrically via a parameter  $\lambda$ , called the **trace decay**, and by  $\gamma$ , the **discount rate**.

$$E^0(s) := 0$$

$$E^t(s) := \gamma \lambda E^{t-1}(s) + \delta_{s,s_t}$$

**Motivation:** States that are more recent/have been

The discounting allows for more recent visits to contribute more to the count than past visits (which may be valuable for non-stationary environments.)

# Generalised Advantage Estimation

Penalising TD updates using the eligibility trace, this gives us the update rule for TD( $\lambda$ ). On timestep  $t$ , perform update

$$\forall s \in \mathcal{S}, \hat{V}_\pi(s) := \hat{V}_\pi(s) + \alpha E^t(s) \left( r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t) \right)$$

Above expression can be unrolled for the advantage function (exercise to reader.)

$$\hat{A}_t = \delta_t + (\gamma\lambda)\delta_{t+1} + \dots + (\gamma\lambda)^{T-t+1}\delta_{T-1}$$

where  $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ .

- Sutton and Barto, Reinforcement Learning: An Introduction  
<http://incompleteideas.net/book/the-book-2nd.html>
- OpenAI Intro to Policy Optimisation [https://spinningup.openai.com/en/latest/spinningup/rl\\_intro3.html](https://spinningup.openai.com/en/latest/spinningup/rl_intro3.html)
- PPO Paper <https://arxiv.org/pdf/1707.06347.pdf>
- Generalised Advantage Estimation  
<https://arxiv.org/pdf/1506.02438.pdf>
- $\backslash\_(\u201c)\_/\_$

- Sutton and Barto, Reinforcement Learning: An Introduction  
<http://incompleteideas.net/book/the-book-2nd.html>
- Yan LeCun's cake analogy, NeurIPS 2016,  
<https://www.youtube.com/watch?v=0unt2Y4qxQo&t=1072s>
- Jay Bailey's DQN Distillation <https://www.lesswrong.com/posts/kyvCNgx9oAwJCuevo/deep-q-networks-explained>
- Playing Atari with Deep Reinforcement Learning  
<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>
- Human-level control through deep reinforcement learning  
<https://www.nature.com/articles/nature14236>
- Rainbow DQN <https://arxiv.org/abs/1710.02298>

# Robins-Monro Convergence conditions

The Robins-Monro convergence conditions are properties of the learning rate that are usually required in most proofs to ensure convergence.

Let  $\alpha_t$  denote the learning rate at time  $t$ . Then the conditions are

$$\sum_{t=1}^{\infty} \alpha_t = \infty \text{ and } \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Intuitively, the first condition ensures that the steps are large enough to eventually overcome any initial conditions/random fluctuations, and the second condition ensures that eventually the steps become small enough to ensure convergence.

# Robins-Monro Convergence conditions

The Robins-Monro convergence conditions are properties of the learning rate that are usually required in most proofs to ensure convergence.

Let  $\alpha_t$  denote the learning rate at time  $t$ . Then the conditions are

$$\sum_{t=1}^{\infty} \alpha_t = \infty \text{ and } \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Intuitively, the first condition ensures that the steps are large enough to eventually overcome any initial conditions/random fluctuations, and the second condition ensures that eventually the steps become small enough to ensure convergence.

An example of such a learning rate would be  $\alpha_t = 1/t$ .

Note that the usual method of choosing  $\forall t, \alpha_t = \alpha \in \mathbb{R}$  fails the RM conditions.