# Last Time

.

# Last Time

- Bandits

Explore - Exploit
Greedy
→ ε - Greedy
softmax
→ UCB
Thompson
Interval

# Guiding Questions

# Guiding Questions

- What is Policy Gradient?

# Guiding Questions

- What is Policy Gradient?
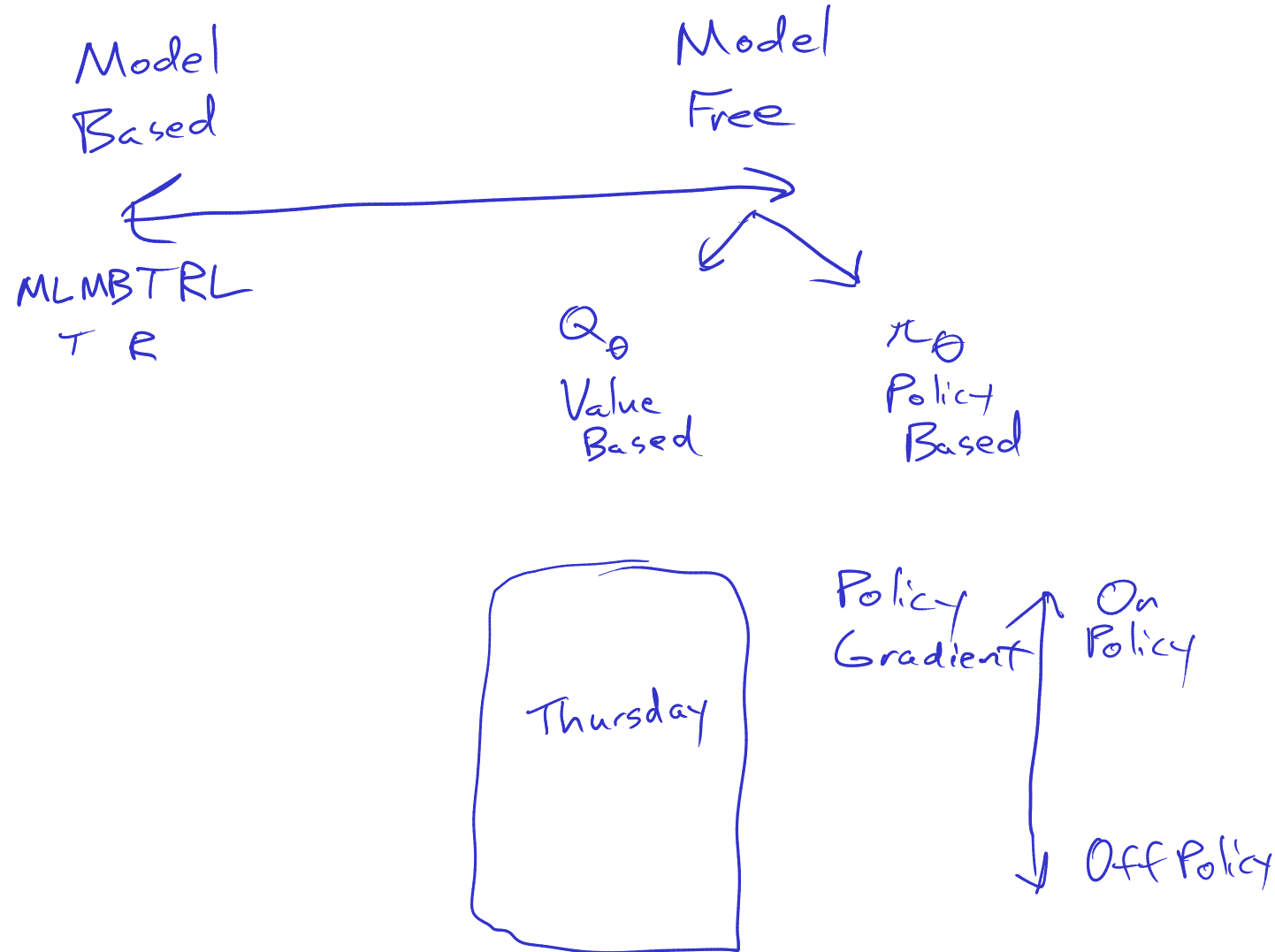- What tricks are needed for it to work effectively?

# Map

# Map

Challenges in RL

- Exploration and Exploitation
- Credit Assignment
- Generalization

# Map

Challenges in RL

- Exploration and Exploitation
- Credit Assignment ←
- Generalization

Model
Based

Model
Free

MLMBTRL
T R

$Q_\theta$
Value
Based

$\pi_\theta$
Policy
Based

Thursday

Policy
Gradient

On
Policy

Off Policy

# Review: Gradient Ascent

# Review: Gradient Ascent

$$U(\theta)$$

- Definition of Gradient

# Review: Gradient Ascent

- Definition of Gradient
- Gradient Ascent

# Review: Gradient Ascent

$$U(\theta)$$

$$\nabla U(\theta) = \left[ \frac{\partial U}{\partial \theta_1}, \quad \ldots \quad \frac{\partial U}{\partial \theta_n} \right]$$

loop

$$\theta \leftarrow \theta + \alpha \widehat{\nabla U}(\theta)$$

$\theta_2$

Stoch Gradient Ascent

$\nabla U(\theta)$

Gradient Ascent

$\theta_1$

$$\nabla U(\theta) = E\left[ \widehat{\nabla U}(\theta) \right]$$

- Definition of Gradient
- Gradient Ascent
- Stochastic Gradient Ascent

$$\pi_\theta$$
$$U(\theta) = E\left[ \sum_{k=0}^{d} \gamma^k r_k \right]$$

4

# Additional Notation

# Additional Notation

- Probabilistic parameterized policies

# Additional Notation

- Probabilistic parameterized policies
- initial state distribution

# Additional Notation

- Probabilistic parameterized policies
- initial state distribution
- trajectory: $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$

# Additional Notation

$$\pi_\theta(a|s)$$

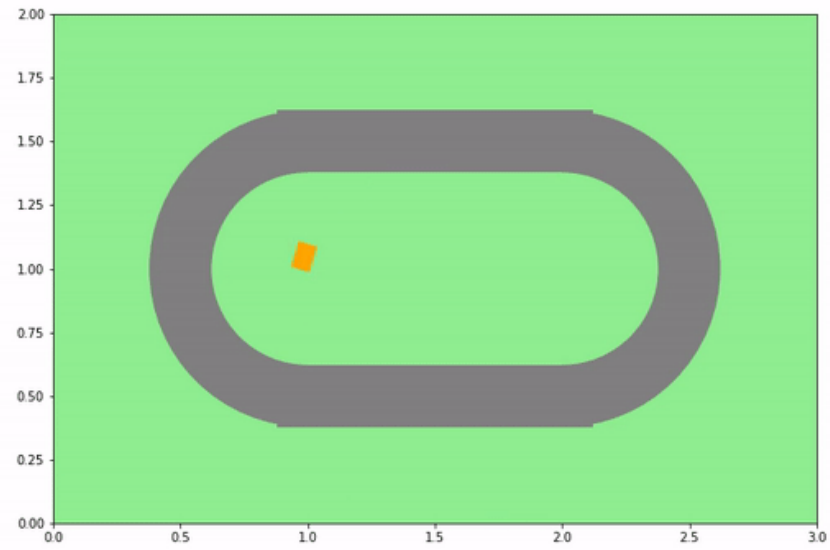$$(S, A, R, T, \gamma) \longrightarrow (S, A, R, T, \gamma, P(s_0))$$

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$
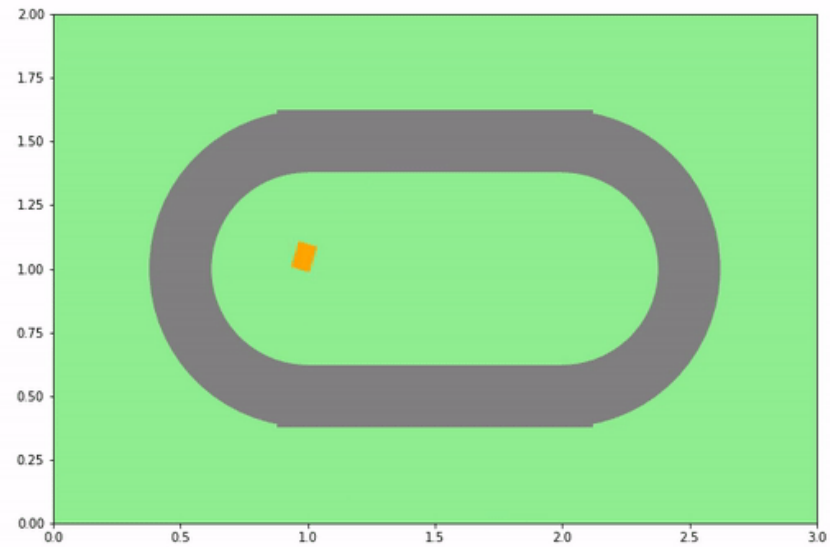
$$A(s,a) = Q(s,a) - V(s)$$

- Probabilistic parameterized policies
- initial state distribution
- trajectory: $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$
- advantage function

# Tricks

# Tricks

# Tricks



$$\nabla_\theta U(\theta)$$

For policy gradient, 3 tricks

- Likelihood Ratio
- Reward to go
- Baseline Subtraction

# Log Derivative

$$\sum \gamma^t r_t = R(\tau)$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau)\, d\tau$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau) \, d\tau$$

$$\nabla U(\theta) = \nabla_\theta \int p_\theta(\tau) R(\tau) \, d\tau$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla U(\theta) = \nabla_\theta \int p_\theta(\tau) R(\tau)\, d\tau$$

$$= \int \nabla_\theta\, p_\theta(\tau) R(\tau)\, d\tau$$

$\frac{\partial \log}{\partial x} \quad \frac{1}{x}$

$\nabla \log \theta =$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla U(\theta) = \nabla_\theta \int p_\theta(\tau) R(\tau)\, d\tau$$

$$= \int \nabla_\theta\, p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla_\theta \log p_\theta(\tau) = \nabla_\theta\, p_\theta(\tau) / p_\theta(\tau)$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla U(\theta) = \nabla_\theta \int p_\theta(\tau) R(\tau)\, d\tau$$

$$= \int \nabla_\theta\, p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla_\theta \log p_\theta(\tau) = \nabla_\theta\, p_\theta(\tau) / p_\theta(\tau)$$

$$\therefore \quad \nabla_\theta\, p_\theta(\tau) = p_\theta(\tau)\, \nabla_\theta \log p_\theta(\tau)$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau) \, d\tau$$

$$\nabla U(\theta) = \nabla_\theta \int p_\theta(\tau) R(\tau) \, d\tau$$

$$= \int \nabla_\theta \, p_\theta(\tau) R(\tau) \, d\tau$$

$$\nabla_\theta \log p_\theta(\tau) = \nabla_\theta \, p_\theta(\tau) / p_\theta(\tau)$$

$$\therefore \; \nabla_\theta \, p_\theta(\tau) = p_\theta(\tau) \, \nabla_\theta \log p_\theta(\tau)$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla U(\theta) = \nabla_\theta \int p_\theta(\tau) R(\tau)\, d\tau$$

$$= \int \nabla_\theta\, p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla_\theta \log p_\theta(\tau) = \nabla_\theta\, p_\theta(\tau) / p_\theta(\tau)$$

$$\therefore \quad \nabla_\theta\, p_\theta(\tau) = p_\theta(\tau)\, \nabla_\theta \log p_\theta(\tau)$$

$$\nabla U = E\left[\widehat{\nabla U}\right)$$

$$= \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) R(\tau)\, d\tau$$

# Log Derivative

$$U(\theta) = \mathrm{E}[R(\tau)]$$

$$= \int p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla U(\theta) = \nabla_\theta \int p_\theta(\tau) R(\tau)\, d\tau$$

$$= \int \nabla_\theta\, p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla_\theta \log p_\theta(\tau) = \nabla_\theta\, p_\theta(\tau) / p_\theta(\tau)$$

$$\therefore \quad \nabla_\theta\, p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

$$= \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) R(\tau)\, d\tau$$

$$\nabla_\theta \log \pi_\theta$$

$$= \mathrm{E}\left[\nabla_\theta \log p_\theta(\tau) R(\tau)\right]$$

$$\nabla U(\theta)$$

# Trajectory Probability Gradient

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau)$$

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau)$$

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau) \qquad\qquad \tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$

$$p_\theta(\tau) = p(s_0) \prod_{k=0}^{d} T(s_{k+1} \mid s_k, a_k) \, \pi_\theta(a_k \mid s_k)$$

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau) \qquad\qquad \tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$

$$p_\theta(\tau) = p(s_0) \prod_{k=0}^{d} T(s_{k+1} \mid s_k, a_k)\, \pi_\theta(a_k \mid s_k)$$

$$\log p_\theta(\tau)$$

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau) \qquad\qquad \tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$

$$\log(ab) = \log a + \log b$$

$$p_\theta(\tau) = p(s_0) \prod_{k=0}^{d} T(s_{k+1} \mid s_k, a_k)\, \pi_\theta(a_k \mid s_k)$$

$$\nabla \log p_\theta(\tau) = \log p(s_0) + \sum_{k=0}^{d} \log T(s_{k+1} \mid s_k, a_k) + \sum_{k=0}^{d} \log \pi_\theta(a_k \mid s_k)$$

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau) \qquad\qquad \tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$

$$p_\theta(\tau) = p(s_0) \prod_{k=0}^{d} T(s_{k+1} \mid s_k, a_k) \, \pi_\theta(a_k \mid s_k)$$

$$\log p_\theta(\tau) = \log p(s_0) + \sum_{k=0}^{d} \log T(s_{k+1} \mid s_k, a_k) + \sum_{k=0}^{d} \log \pi_\theta(a_k \mid s_k)$$

$$\nabla_\theta \log p_\theta(\tau)$$

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau) \qquad\qquad \tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$

$$p_\theta(\tau) = p(s_0) \prod_{k=0}^{d} T(s_{k+1} \mid s_k, a_k)\, \pi_\theta(a_k \mid s_k)$$

$$\log p_\theta(\tau) = \log p(s_0) + \sum_{k=0}^{d} \log T(s_{k+1} \mid s_k, a_k) + \sum_{k=0}^{d} \log \pi_\theta(a_k \mid s_k)$$

$$\nabla_\theta \log p_\theta(\tau) = \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)$$

Grid World

$\theta \in \mathbb{R}^{|S| \times |A|}$

$\pi_\theta(a \mid s) = \dfrac{\theta_{s,a}}{\sum\limits_{a} \theta_{s,a}}$

# Trajectory Probability Gradient

$$\nabla_\theta \log p_\theta(\tau) \qquad\qquad \tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_d, a_d, r_d)$$

$$p_\theta(\tau) = p(s_0) \prod_{k=0}^{d} T(s_{k+1} \mid s_k, a_k) \, \pi_\theta(a_k \mid s_k)$$

$$\log p_\theta(\tau) = \log p(s_0) + \sum_{k=0}^{d} \log T(s_{k+1} \mid s_k, a_k) + \sum_{k=0}^{d} \log \pi_\theta(a_k \mid s_k)$$

$$\nabla_\theta \log p_\theta(\tau) = \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)$$

$$\nabla U(\theta) = E\left[ \; \nabla_\theta \log p_\theta(\tau) \, R(\tau) \right]$$

$$\nabla U(\theta) = \mathrm{E}\left[ \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau) \right]$$

# Policy Gradient

# Policy Gradient

loop

$\quad \tau \leftarrow \text{simulate}(\pi_\theta)$

$\quad \theta \leftarrow \theta + \alpha \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)$

# Policy Gradient

On Policy — only use
data from current policy

Off Policy

loop

$\tau \leftarrow \text{simulate}(\pi_\theta)$

$\theta \leftarrow \theta + \alpha \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)$

On Policy!

# Causality

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)\right]$$

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)\right]$$

$$= \mathrm{E}\left[\left(\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\right)\left(\sum_{k=0}^{d} \gamma^k r_k\right)\right]$$

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)\right]$$

$$= \mathrm{E}\left[\left(\sum_{k=0}^{d} \underbrace{\nabla_\theta \log \pi_\theta(a_k \mid s_k)}_{f_k}\right)\left(\sum_{k=0}^{d} \gamma^k r_k\right)\right]$$

# Causality

$$\nabla U(\theta) = \mathrm{E} \left[ \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau) \right]$$

$$= \mathrm{E} \left[ \left( \sum_{k=0}^{d} \underbrace{\nabla_\theta \log \pi_\theta(a_k \mid s_k)}_{\textcolor{blue}{f_k}} \right) \left( \sum_{k=0}^{d} \gamma^k r_k \right) \right]$$

$$= \mathrm{E} \left[ (f_0 + \ldots + f_d) \left( \gamma^0 r_0 + \ldots \gamma^d r_d \right) \right]$$

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)\right]$$

$$= \mathrm{E}\left[\left(\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(\underbrace{a_k \mid s_k}_{f_k})\right)\left(\sum_{k=0}^{d} \gamma^k r_k\right)\right]$$

$$= \mathrm{E}\left[(f_0 + \ldots + f_d)\left(\gamma^0 r_0 + \ldots \gamma^d r_d\right)\right]$$

$$= \mathrm{E}\left[\begin{array}{l} f_0 \gamma^0 r_0 + f_0 \gamma^1 r_1 + f_0 \gamma^2 r_2 + \ldots + f_0 \gamma^d r_d \\ + f_1 \gamma^0 r_0 + f_1 \gamma^1 r_1 + f_1 \gamma^2 r_2 + \ldots f_1 \gamma^d r_d \\ \vdots \\ + f_d \gamma^0 r_0 + f_d \gamma^1 r_1 + f_d \gamma^2 r_2 + \ldots f_d \gamma^d r_d \end{array}\right]$$

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)R(\tau)\right]$$

$$= \mathrm{E}\left[\left(\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\right)\left(\sum_{k=0}^{d} \gamma^k r_k\right)\right]$$

$$\underbrace{\phantom{\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)}}_{\color{blue}{f_k}}$$

$$= \mathrm{E}\left[(f_0 + \ldots + f_d)\left(\gamma^0 r_0 + \ldots \gamma^d r_d\right)\right]$$

$$= \mathrm{E}\left[\begin{array}{l} f_0\gamma^0 r_0 + f_0\gamma^1 r_1 + f_0\gamma^2 r_2 + \ldots + f_0\gamma^d r_d \\ + f_1\gamma^0 r_0 + f_1\gamma^1 r_1 + f_1\gamma^2 r_2 + \ldots f_1\gamma^d r_d \\ \vdots \\ + f_d\gamma^0 r_0 + f_d\gamma^1 r_1 + f_d\gamma^2 r_2 + \ldots f_d\gamma^d r_d \end{array}\right]$$

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)\right]$$

$$= \mathrm{E}\left[\left(\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\right)\left(\sum_{k=0}^{d} \gamma^k r_k\right)\right]$$

$$f_k$$

$$= \mathrm{E}\left[(f_0 + \ldots + f_d)\left(\gamma^0 r_0 + \ldots \gamma^d r_d\right)\right]$$

$$= \mathrm{E}\left[\begin{array}{l} f_0\gamma^0 r_0 + f_0\gamma^1 r_1 + f_0\gamma^2 r_2 + \ldots + f_0\gamma^d r_d \\ + f_1\gamma^0 r_0 + f_1\gamma^1 r_1 + f_1\gamma^2 r_2 + \ldots f_1\gamma^d r_d \\ \vdots \\ + f_d\gamma^0 r_0 + f_d\gamma^1 r_1 + f_d\gamma^2 r_2 + \ldots f_d\gamma^d r_d \end{array}\right]$$

$$= \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\left(\sum_{l=k}^{d} \gamma^l r_l\right)\right]$$

10

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)\right]$$

$$= \mathrm{E}\left[\left(\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\right)\left(\sum_{k=0}^{d} \gamma^k r_k\right)\right]$$

$$\underbrace{\phantom{\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)}}_{\textstyle f_k}$$

$$= \mathrm{E}\left[(f_0 + \ldots + f_d)\left(\gamma^0 r_0 + \ldots \gamma^d r_d\right)\right]$$

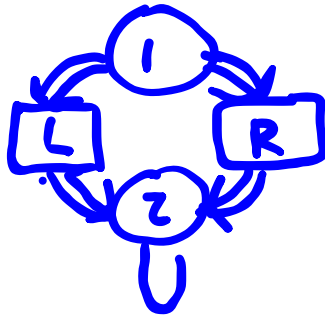$$= \mathrm{E}\left[\begin{array}{l} f_0\gamma^0 r_0 + f_0\gamma^1 r_1 + f_0\gamma^2 r_2 + \ldots + f_0\gamma^d r_d \\ + f_1\gamma^0 r_0 + f_1\gamma^1 r_1 + f_1\gamma^2 r_2 + \ldots f_1\gamma^d r_d \\ \vdots \\ + f_d\gamma^0 r_0 + f_d\gamma^1 r_1 + f_d\gamma^2 r_2 + \ldots f_d\gamma^d r_d \end{array}\right]$$

$$= \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\left(\sum_{l=k}^{d} \gamma^l r_l\right)\right] \qquad = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k r_{k,\text{to-go}}\right]$$

10

# Causality

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) R(\tau)\right]$$

$$= \mathrm{E}\left[\left(\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\right)\left(\sum_{k=0}^{d} \gamma^k r_k\right)\right]$$

$$f_k$$

$$= \mathrm{E}\left[(f_0 + \ldots + f_d)\left(\gamma^0 r_0 + \ldots \gamma^d r_d\right)\right]$$

$$= \mathrm{E}\left[ \begin{array}{l} f_0\gamma^0 r_0 + f_0\gamma^1 r_1 + f_0\gamma^2 r_2 + \ldots + f_0\gamma^d r_d \\ + f_1\gamma^0 r_0 + f_1\gamma^1 r_1 + f_1\gamma^2 r_2 + \ldots f_1\gamma^d r_d \\ \vdots \\ + f_d\gamma^0 r_0 + f_d\gamma^1 r_1 + f_d\gamma^2 r_2 + \ldots f_d\gamma^d r_d \end{array} \right]$$

$$= \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\left(\sum_{l=k}^{d} \gamma^l r_l\right)\right] \qquad = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k r_{k,\text{to-go}}\right] \; Q^\theta(s_k, a_k)$$

# Discuss

$\cup$

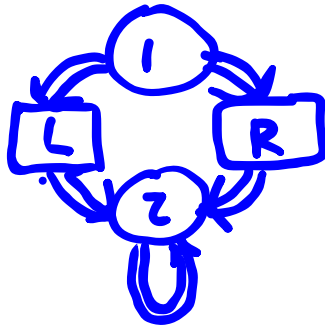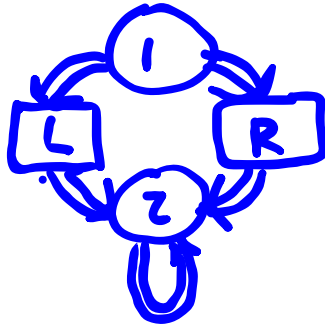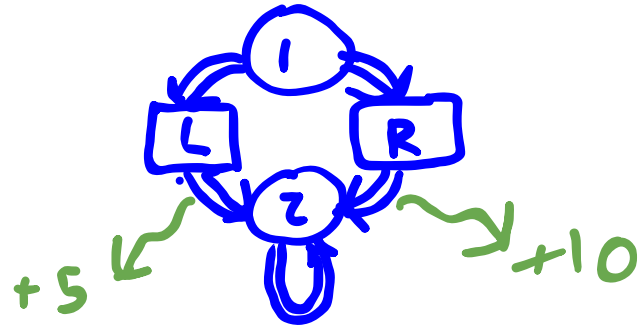# Discuss
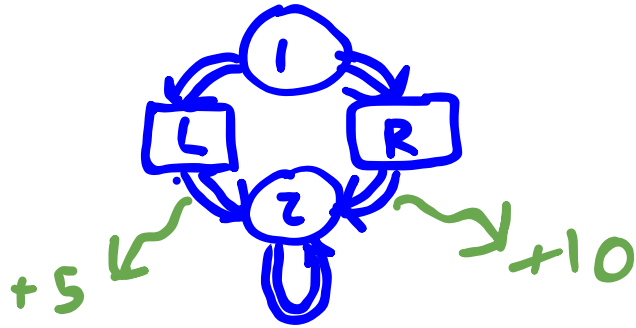
# Discuss

# Discuss

$A = \{L, R\}$

# Discuss

$A = \{L, R\}$

# Discuss

$A = \{L, R\}$



loop

$\qquad \tau \leftarrow \text{simulate}(\pi_\theta)$

$\qquad \theta \leftarrow \theta + \alpha \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$
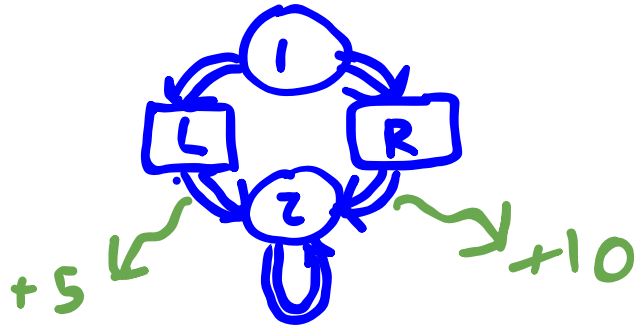
# Discuss

A = {L, R}



loop

$\quad \tau \leftarrow \mathrm{simulate}(\pi_\theta)$

$\quad \theta \leftarrow \theta + \alpha \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$

1. 1. Given $\theta = (0.2, 0.8)$ calculate $\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$ for two cases, (a) where $a_0 = R$ and (b) where $a_0 = L$

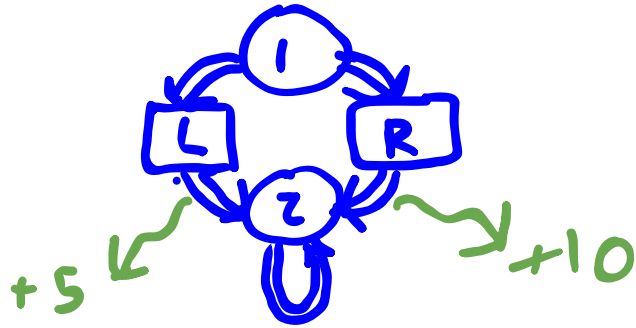# Discuss

$A = \{L, R\}$



loop

$$\tau \leftarrow \text{simulate}(\pi_\theta)$$

$$\theta \leftarrow \theta + \alpha \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$$

1. 1. Given $\theta = (0.2, 0.8)$ calculate $\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$ for two cases, (a) where $a_0 = R$ and (b) where $a_0 = L$
2. What happens if $\theta_1 \to 0$

# Discuss

$A = \{L, R\}$



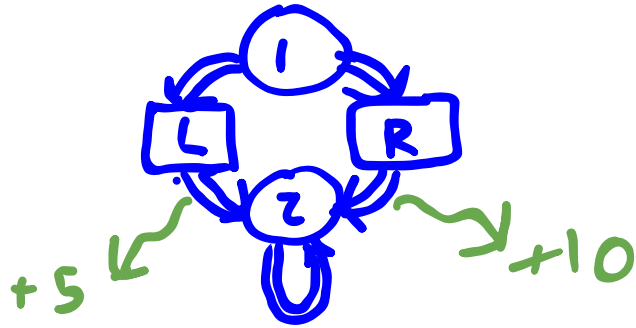$$\pi_\theta(a = L \mid s) = \frac{\theta_1}{\theta_1 + \theta_2}$$

loop

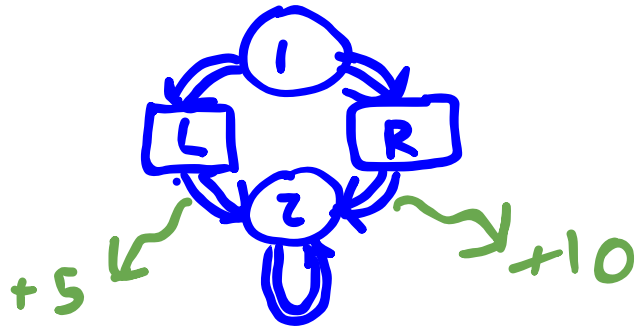    $\tau \leftarrow \mathrm{simulate}(\pi_\theta)$

    $\theta \leftarrow \theta + \alpha \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$

1. 1. Given $\theta = (0.2, 0.8)$ calculate $\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$ for two cases, (a) where $a_0 = R$ and (b) where $a_0 = L$
2. What happens if $\theta_1 \to 0$

# Discuss

$A = \{L, R\}$

$\theta = (\theta_1, \theta_2)$



$$\pi_\theta(a = L \mid s) = \frac{\theta_1}{\theta_1 + \theta_2}$$

$$\pi_\theta(a = R \mid s) = \frac{\theta_2}{\theta_1 + \theta_2}$$

$\log \pi_\theta(a = L \mid s) = \log \theta_1 - \log(\theta_1 + \theta_2)$

$\dfrac{\partial \log \pi_\theta(a = L \mid s)}{\partial \theta_1} = \dfrac{1}{\theta_1} - \dfrac{1}{\theta_1 + \theta_2}$

$\dfrac{\partial \log \pi_\theta(a = L \mid s)}{\partial \theta_2} = -\dfrac{1}{\theta_1 + \theta_2}$

$\nabla_\theta \log \pi_\theta(a_0 = L \mid s_0 = 1) = \begin{bmatrix} \frac{1}{0.2} - \frac{1}{1} \\ -1 \end{bmatrix} = \begin{bmatrix} 4 \\ -1 \end{bmatrix}$

$\nabla_\theta \log \pi_\theta(a_0 = L \mid s_0 = 1) \, r_{\text{togo}} = \begin{bmatrix} 4 \\ -1 \end{bmatrix}(5 \quad -7.5)$

$\nabla_\theta \log \pi_\theta(a_0 = R \mid s_0 = 1) \, r_{\text{togo}} = \begin{bmatrix} -1 \\ 0.25 \end{bmatrix}(10 \quad -7.5)$

loop

$\quad \tau \leftarrow \text{simulate}(\pi_\theta)$

$\quad \theta \leftarrow \theta + \alpha \sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$

1. 1. Given $\theta = (0.2, 0.8)$ calculate $\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k r_{k,\text{to go}}$ for two cases, (a) where $a_0 = R$ and (b) where $a_0 = L$
2. What happens if $\theta_1 \to 0$

# Baseline Subtraction

# Baseline Subtraction

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k r_{k,\text{to-go}}\right]$$

# Baseline Subtraction

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k r_{k,\text{to-go}}\right]$$

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k \left(r_{k,\text{to-go}} - r_{\text{base}}(s_k)\right)\right]$$

# Baseline Subtraction

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k r_{k,\text{to-go}}\right]$$

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k \left(r_{k,\text{to-go}} - r_{\text{base}}(s_k)\right)\right]$$

does not bias
(proof in book)

# Baseline Subtraction

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k r_{k,\text{to-go}}\right]$$

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k \left(r_{k,\text{to-go}} - r_{\text{base}}(s_k)\right)\right]$$

*does not bias*
*(proof in book)*

$$r_{\text{base},i} = \frac{\mathbb{E}_{a,s,r_{\text{to-go}},k}\left[\ell_i(a,s,k)^2 r_{\text{to-go}}\right]}{\mathbb{E}_{a,s,k}\left[\ell_i(a,s,k)^2\right]}$$

# Baseline Subtraction

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k r_{k,\text{to-go}}\right]$$

$$\nabla U(\theta) = \mathrm{E}\left[\sum_{k=0}^{d} \nabla_\theta \log \pi_\theta(a_k \mid s_k)\, \gamma^k \left(r_{k,\text{to-go}} - r_{\text{base}}(s_k)\right)\right]$$

*does not bias (proof in book)*

*Optimal*

$$r_{\text{base},i} = \frac{\mathbb{E}_{a,s,r_{\text{to-go}},k}\left[\ell_i(a,s,k)^2 r_{\text{to-go}}\right]}{\mathbb{E}_{a,s,k}\left[\ell_i(a,s,k)^2\right]}$$

$$\ell_i(a,s,k) = \gamma^{k-1}\frac{\partial}{\partial \theta_i} \log \pi_\theta(a \mid s)$$

*practical*

$$\hat{V}(s)$$

# Guiding Questions

$\pi_\theta$

S.G.D.

$\nabla U(\theta)$

- What is Policy Gradient?
- What tricks are needed for it to work effectively?

— Log Derivative

— Causality

— Baseline Subtraction