# Reinforcement Learning: Policy-based Methods

## Benjamin Rosman

School of Computer Science and Applied Mathematics
University of the Witwatersrand
South Africa

African Masters of Machine Intelligence (Kigali, Rwanda)

January 13-17 2020

R A I L
L A B

Benjamin Rosman

# Further reading…

- Reinforcement Learning: An Introduction (Sutton and Barto)
  - Chapter 13

- CS 294-122 (Sergey Levine)
  - Lectures 4 and 5
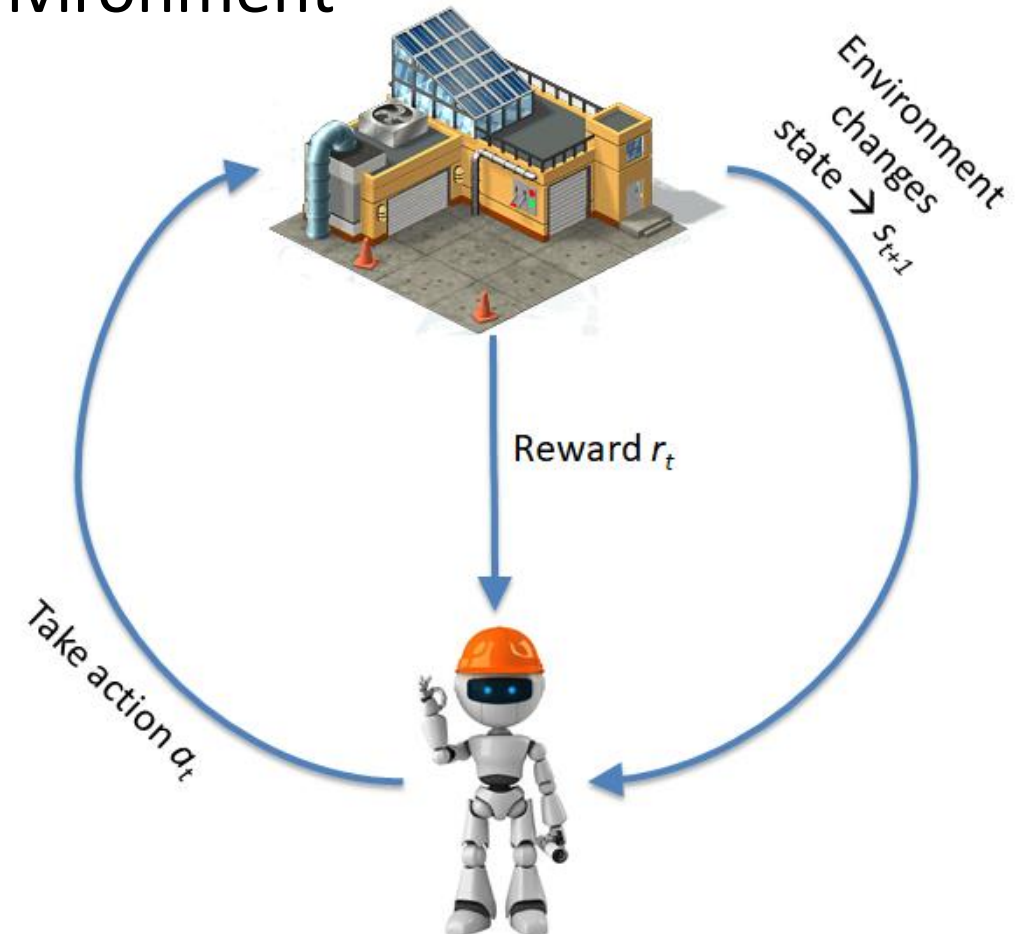
- COMPM050/COMPGI13 (David Silver)
  - Lecture 7

# A quick recap

Decision maker (agent) exists within an environment

Agent takes **actions** based on the environment **state**

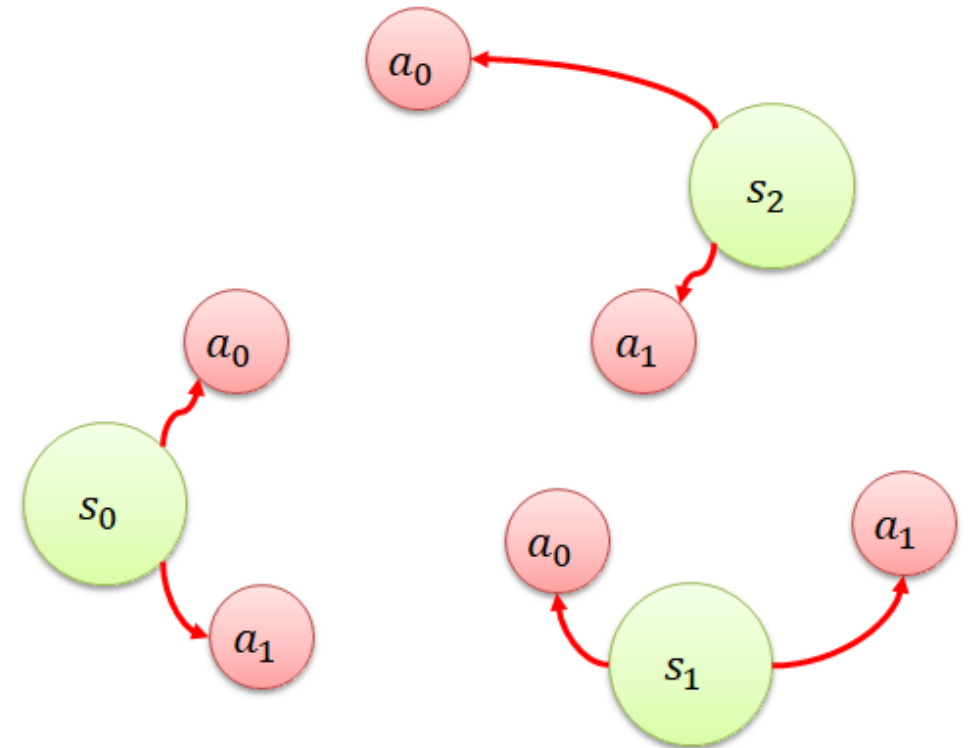Environment **state** updates

Agent receives feedback as **rewards**



Environment changes state $\rightarrow s_{t+1}$

Reward $r_t$

Take action $a_t$

Benjamin Rosman

# A model for decision making

Markov Decision Process (MDP)

M = $\langle S, A, T, R, \gamma \rangle$

- States: encode world configurations

- Actions: choices made by agent
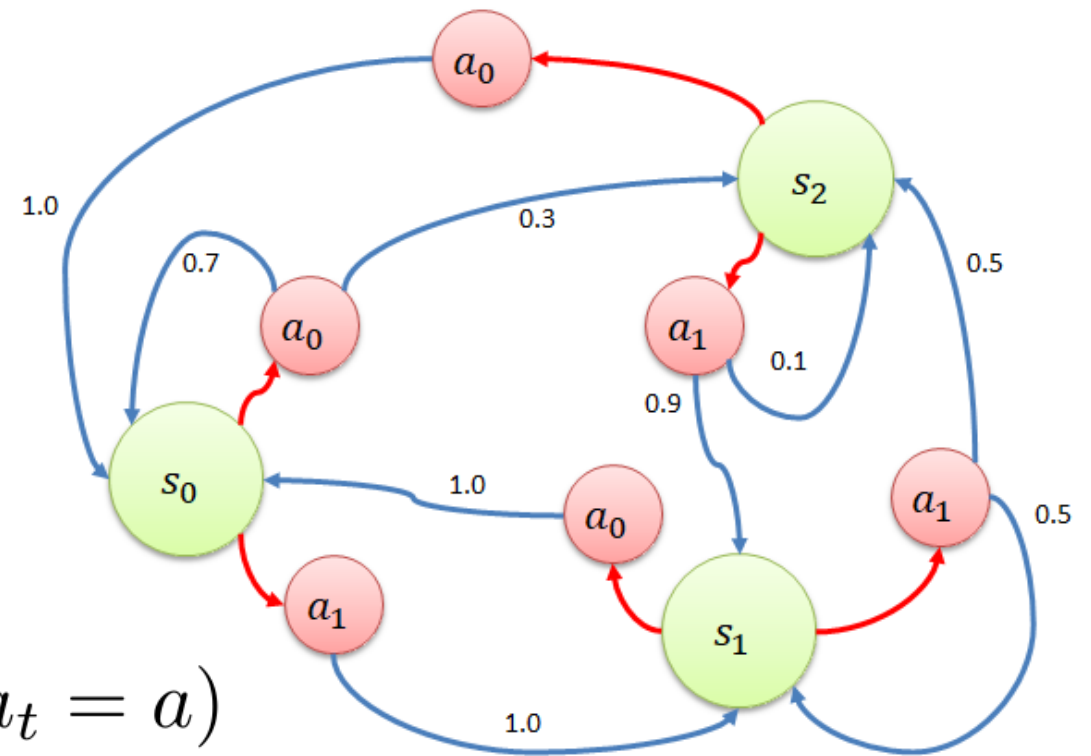
# A model for decision making

Markov Decision Process (MDP)

M = $\langle S, A, T, R, \gamma \rangle$

- Transition function: how the world evolves under actions

  Stochastic!

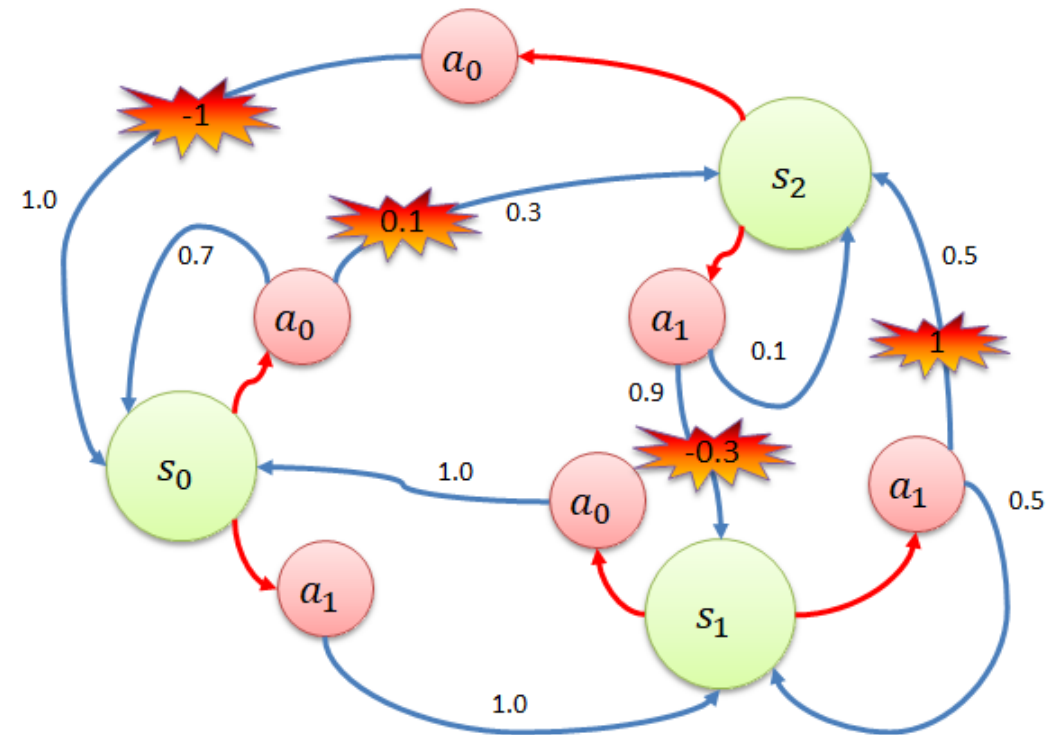$$T(s, a, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$$

# A model for decision making

Markov Decision Process (MDP)

M = $\langle S, A, T, R, \gamma \rangle$

- Rewards: feedback signal to agent

$$R(s,a) = E[r_t | s_t = s, a_t = a]$$
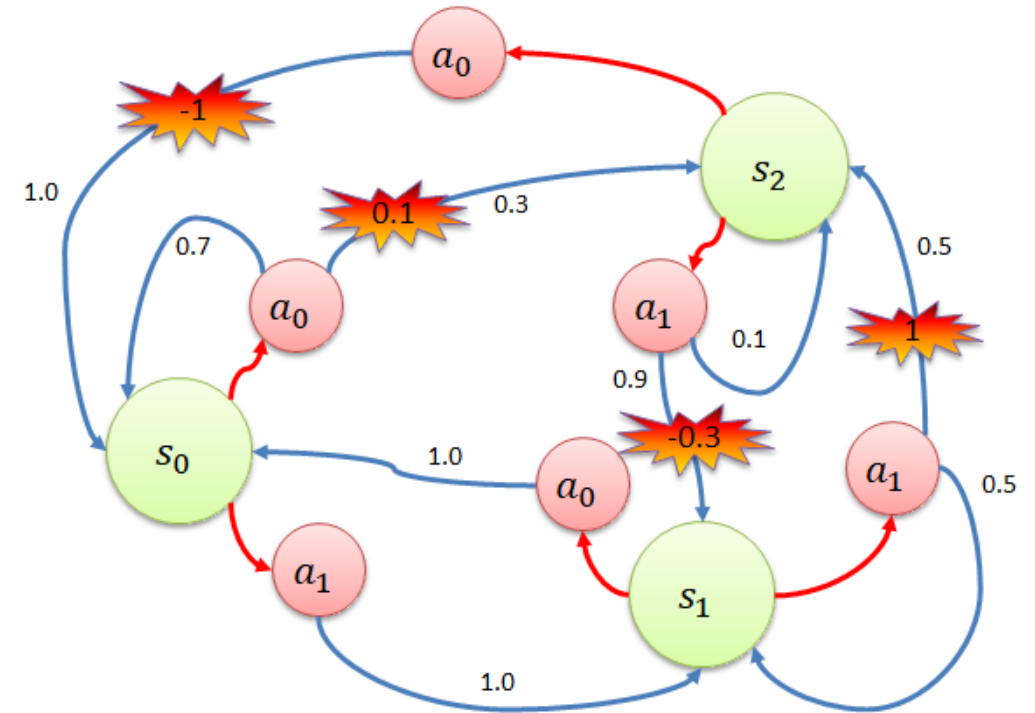


Benjamin Rosman

# A model for decision making

Markov Decision Process (MDP)

$M = \langle S, A, T, R, \gamma \rangle$

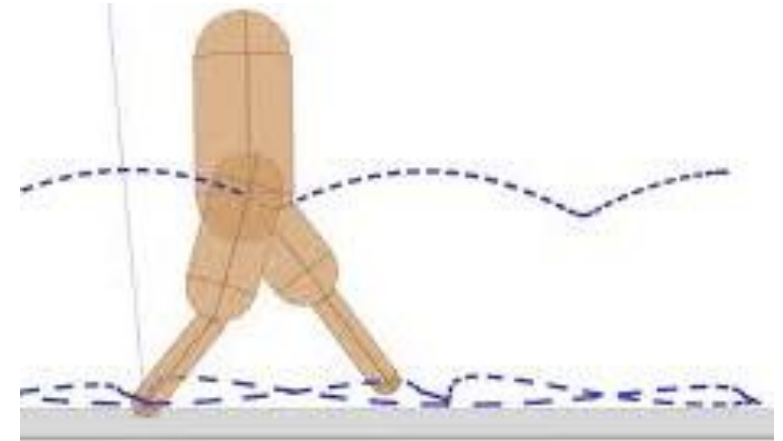- $\gamma \in [0,1]$ discounting for future rewards

Particularly important for **continuing** tasks

The value of something now is usually greater than in the future.

# Example

- I want to get this robot to walk as far as possible:
  - What is $S, A, T, R$?



- I want to play this game as well as possible:
  - What is $S, A, T, R$?

# Why can't we act greedily?

Cannot just rely on the **instantaneous** reward function
Tradeoff: <span style="color:red">don't just act myopically (short term)</span>



1 step                    5 steps

Notion of **value** to codify the goodness of a state, considering a policy running into the future

- Represented as a **value function**

# Warning: reward functions

- The reward tells the agent **what** we want it to achieve
    - Learning is about figuring out **how** to achieve it

- Be careful: <span style="color:red">the agent will literally do exactly what you ask it</span>
    - It doesn't "know what you mean"

- Where might this be a problem?

Benjamin Rosman

# Value functions

Allow us to act greedily by considering estimates of the future

**Optimal** policy:



a) gridworld   b) $v_*$   c) $\pi_*$

NB: Equivalence between optimal policy and optimal value function

# Value functions: recursion

$V(s) \Rightarrow$ expected return starting at $s$ and following $\pi$
Suggests dependence on value of next state $s'$

Bellman Equation:

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V^\pi(s')$$

value of s     immediate reward     for all possible next states     the probability of reaching that state with $\pi$     value of s'

# Solution techniques

- If you know $T$ and $R$:
  - Iteratively compute solution
  - Dynamic programming
  - Policy iteration and value iteration

- If you don't know $T$ and $R$: (more common case)
  - Explore the environment and collect samples
  - Use samples to learn $T$ and $R$
    - Model-based RL (more next time)
  - Use samples to learn $V$
    - Model-free RL
    - E.g. TD-learning, SARSA, Q-learning, …

# Q-Learning

- Initialise $Q(s, a)$ for all $s, a$
- Repeat (for each episode):
  - Initialise $s$
  - Repeat (for each step of episode):
    - Choose $a$ from $s$ using policy from $Q$ ($\epsilon$-greedy)
    - Take action $a$, observe $r, s'$ ——————————→ act
    - $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$  ————→ learn

    take best next action
    (so far)

    - $s \leftarrow s'$
  - Until $s$ is terminal

# Function approximation

- This works well if you can tabulate the value function

| State | $a_1$ | ... | $a_m$ |
|---|---|---|---|
| $s_1$ | 1.7 | ... | 2.8 |
| $s_2$ | -3.9 | ... | -3.1 |
| ... | ... | ... | ... |
| $s_n$ | 0.2 | ... | -0.1 |

- But if the state space is large/continuous
  - Approximate this with a learned function
  - Why?
    - Infinity? May not visit every state (many times)? Need to generalise?

Benjamin Rosman

# Deep Q-Networks



[Mnih et al., 2015]

# DQN

- Take action $a_t$ according to $\epsilon$-greedy policy
  - Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory $D$
- <span style="color:red">Sample mini-batch of transitions $(s, a, r, s')$ from $D$</span>
  - Experience replay: decorrelate samples
- <span style="color:green">Compute Q-learning targets w.r.t old fixed parameters $w^-$</span>
  - Fixed Q-targets: avoid oscillations
- Optimise MSE between Q-network and Q-learning targets
  - $L_i(w_i) = \mathbb{E}_{s,a,r,s' \sim D_i} \left[ \left( r + \gamma \max_{a'} Q(s', a', w_i^-) - Q(s, a, w_i) \right)^2 \right]$
  - Using stochastic gradient descent on $\nabla_{w_i} L_i(w_i)$

# RL approaches



Value
Function
Based

Model
Based

π ← Q ← T, R

# Questions?

# RL approaches

Policy
Search

Value
Function
Based

Model
Based

π ← Q ← T, R

# Policy Search

- Learn policy directly:

$$\pi_\theta(s, a) = \pi(s, a; \theta)$$

- Parameterise policy: learn parameters of policy

- Why?
  - When might it be easier to learn a policy than a value function?
  - Learning a Q-function can be complicated
    - Policy may be much simpler than learning a value for each state-action
  - Injecting information?
  - Stochasticity?

Benjamin Rosman

# Policy Search

Trajectory $\tau$

$$p(\tau|\theta) = p_\theta(s_1, a_1, s_2, a_2, \ldots, s_T, a_T)$$

$$= p(s_1) \prod_{t=1}^{T} \pi_\theta(a_t|s_t) p(s_{t+1}|s_t, a_t)$$

- Objective function?

  - Maximise return given $\theta$:

    - $J(\theta) = \mathbb{E}[\sum_t \gamma^t r_t | \pi_\theta]$

    - $\theta^* = argmax_\theta J(\theta)$

    - Always more efficient to follow the gradient!

Note:
The return $R$, cost $J$, utility $U$ are often used interchangeably here

Benjamin Rosman

# Hill Climbing

- What if you can't differentiate $\pi$?


- Sample-based optimisation:
  - Sample some $\theta$ values near your current best $\theta$
  - Compute return
  - Approximate a gradient
    - Finite differences
  - Adjust your current best $\theta$ to give the highest value


- Other approaches, e.g. genetic algorithms

# Aibo gait optimization

## from Kohl and Stone, ICRA 2004



Fig. 2. The elliptical locus of the Aibo's foot. The half-ellipse is defined by length, height, and position in the $x$-$y$ plane.

All told, the following set of 12 parameters define the Aibo's gait [10]:

- The front locus (3 parameters: height, $x$-pos., $y$-pos.)
- The rear locus (3 parameters)
- Locus length
- Locus skew multiplier in the $x$-$y$ plane (for turning)
- The height of the front of the body
- The height of the rear of the body
- The time each foot takes to move through its locus
- The fraction of time each foot spends on the ground

Benjamin Rosman

# Using gradients

- If we can differentiate $\pi$
  - Compute and ascend $\partial R / \partial \theta$

- This is the gradient of return w.r.t policy parameters

- $\theta_{t+1} = \theta_t + \Delta\theta_t$
  $= \theta_t + \alpha \nabla J(\theta_t)$

- These are called **policy gradient methods**

# Policy Gradient Theorem

- Why does this work?

- Relate the **gradient of performance with respect to the policy parameter** to the **gradient of the policy**

- Policy gradient theorem:
  - $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) \, Q^{\pi_\theta}(s, a)]$

- No explicit dependence on distribution of states (or model)

# What is a good form for a parameterised function $f$?

- Simplest thing you can do?
  - *Linear value function approximation*
  - Use set of basis functions $\phi_1, \dots, \phi_n$
  - $f$ is a linear function of them:

  - $\hat{f} = \boldsymbol{w} \cdot \Phi(s, a) = \sum_{i=1}^{n} w_i \phi_i(s, a)$
  - We'll want to learn parameters $\boldsymbol{w}$

- Neural network:
  - $f = f(s, a; \boldsymbol{w})$

Basis functions $\phi(x)$:
- Could be polynomial in state vars:
  - 1st order: $[1, x, y]$
  - 2nd order: $[1, x, y, x^2, y^2, xy]$
  - This is a Taylor expansion
- Others:
  - Fourier basis
  - Wavelet basis
  - …

# REINFORCE (Monte-Carlo Policy Gradient)

- REINFORCE: one particularly popular sample-based estimate of the gradient
  - Based on the policy gradient theorem, but approximate the $\mathbb{E}$ with sampled trajectories (Monte-Carlo samples)

The return $R_t$ acts as an estimate of $Q^{\pi_\theta}(s, a)$

$$\Delta\theta_t = \alpha R_t \frac{\nabla_\theta \pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} = \alpha R_t \nabla_\theta \log \pi_\theta(a_t|s_t)$$

Log derivative trick: $\frac{\nabla_\theta x}{x} = \nabla_\theta \log x$

Benjamin Rosman

# REINFORCE algorithm

- Initialise $\theta$

- For each episode
  - Choose actions according to $\pi_\theta$: $a \sim \pi_\theta(a|s)$
  - Gather samples $\{s_1, a_1, r_1, \ldots, s_T, a_T, r_T\}$
  - For $t = 1 \; to \; T$
    - $\theta \leftarrow \theta + \alpha R_t \nabla_\theta \log \pi_\theta(a_t|s_t)$

(that's it)

# Deriving REINFORCE

- Cost:
$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)} [r(\tau)]$$
$$= \int_\tau r(\tau) p(\tau;\theta) \mathrm{d}\tau$$

$$\tau = (s_0, a_0, r_0, s_1, \ldots)$$

- Derivative:
$$\nabla_\theta J(\theta) = \int_\tau r(\tau) \nabla_\theta p(\tau;\theta) \mathrm{d}\tau$$

Log derivative trick: $\dfrac{\nabla_\theta x}{x} = \nabla_\theta \log x$

- Transformation:
$$\nabla_\theta p(\tau;\theta) = \boxed{p(\tau;\theta)} \frac{\nabla_\theta p(\tau;\theta)}{\boxed{p(\tau;\theta)}} = p(\tau;\theta) \nabla_\theta \log p(\tau;\theta)$$

- Substitute:
$$\nabla_\theta J(\theta) = \int_\tau \left( r(\tau) \nabla_\theta \log p(\tau;\theta) \right) p(\tau;\theta) \mathrm{d}\tau$$
$$= \mathbb{E}_{\tau \sim p(\tau;\theta)} [r(\tau) \nabla_\theta \log p(\tau;\theta)]$$
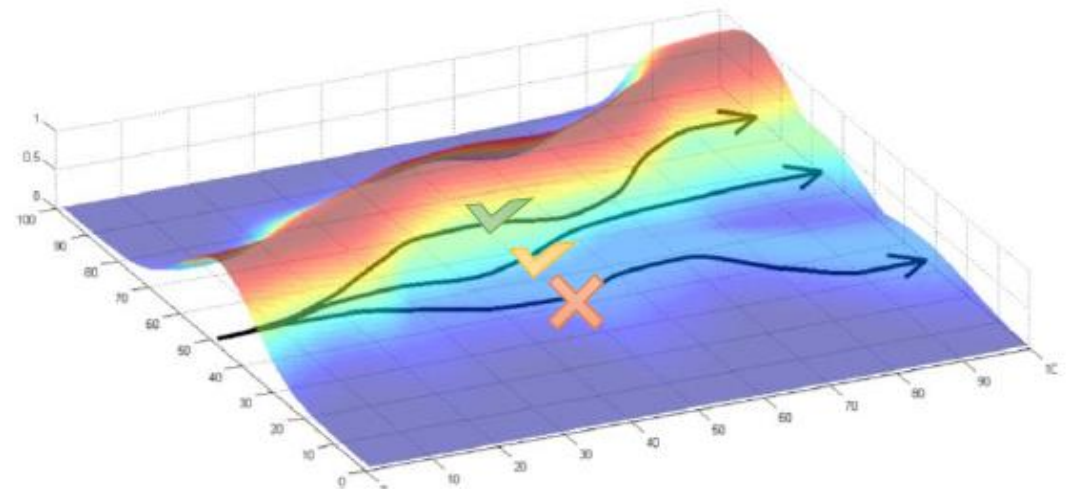
# Deriving REINFORCE

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)}\left[r(\tau)\nabla_\theta \log p(\tau;\theta)\right]$$

- Recall:
$$p(\tau;\theta) = \prod_{t \geq 0} p(s_{t+1}|s_t, a_t)\pi_\theta(a_t|s_t)$$

- So:
$$\log p(\tau;\theta) = \sum_{t \geq 0} \log p(s_{t+1}|s_t, a_t) + \log \pi_\theta(a_t|s_t)$$

- Derivative:
$$\nabla_\theta \log p(\tau;\theta) = \sum_{t \geq 0} \nabla_\theta \log \pi_\theta(a_t|s_t)$$
<span style="color:red">Note we lose dependence on dynamics</span>

- And so estimate:
$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} r(\tau)\nabla_\theta \log \pi_\theta(a_t|s_t)$$

Benjamin Rosman

# Interpretation

- Gradient estimator:  $\nabla_\theta J(\theta) \approx \sum_{t \geq 0} \boxed{r(\tau)} \boxed{\nabla_\theta \log \pi_\theta(a_t | s_t)}$

- Think of this as saying:
  - If $r(\tau)$ is high: push up probabilities of seen actions
  - If $r(\tau)$ is low: push down probabilities of seen actions

- Simple version of credit assignment

# Variance

- <span style="color:red">This gradient estimator (MC) turns out to have a high variance</span>

$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_\theta \log \pi_\theta(a_t|s_t)$$

  - Why?
    - These are all samples of a run of a policy!
  - Slow convergence

- Correct with a baseline:

- $\nabla_\theta J(\theta) \approx \sum_{t \geq 0} (r(\tau) - b(s_t)) \nabla_\theta \log \pi_\theta(a_t|s_t)$

- Because $b$ could be 0, this is a generalisation of REINFORCE
  - Will converge asymptotically to a local minimum

Benjamin Rosman

# Why can we use a baseline?

- $\nabla_\theta J(\theta) \approx \sum_{t \geq 0} (r(\tau) - b(s_t)) \nabla_\theta \log \pi_\theta(a_t|s_t)$

- Intuition: can add or subtract $b$ from $r$ without biasing algorithm
    - As long as $b$ not a function of $a_t$

- Mathematically (with some notational abuse):

Keeps the gradient unbiased, i.e. doesn't change the expected value, but can change the variance!

$$\mathbb{E}_{\pi_\theta} \left[ \sum_t b \nabla_\theta \log \pi_\theta(a_t|s_t) \right]$$

$$= \int \left[ \sum_t b \pi_\theta(a_t|s_t) \nabla_\theta \log \pi_\theta(a_t|s_t) \right] d\tau$$

$$= \int b \nabla_\theta \pi_\theta(\tau) \, d\tau \qquad \text{Log derivative trick: } \frac{\nabla_\theta x}{x} = \nabla_\theta \log x$$

$$= b \nabla_\theta \int \pi_\theta(\tau) \, d\tau = b \nabla_\theta 1 = 0$$

Benjamin Rosman

# Choice of baseline

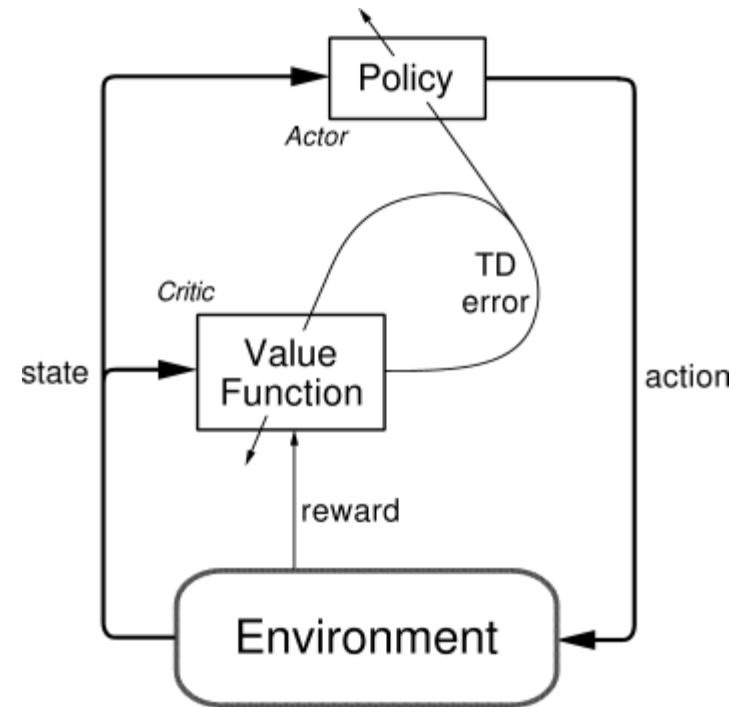- What baseline to use?

- $b(s_t) = V(s_t)$
  - Change based on whether or not reward was better than expected

  - Term $r(\tau) - b(s_t)$ resembles advantage function:
    - $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$
    - Measures how much better $a$ is than whatever $\pi$ would have done

- Suggests we should be learning $\pi$ and $V$!

# Learning with a baseline

# Actor-Critic

- Combine ideas from <span style="color:red">policy</span> and <span style="color:blue">value function</span> methods
  - Approximate both the <span style="color:red">policy</span> and the <span style="color:blue">value function</span>
- <span style="color:red">Actor improvement</span>
  - Policy parameterised by $\theta$
- <span style="color:blue">Critic evaluation</span>
  - Value function parameterised by $\omega$
  - Either $V(s; \omega)$ or $Q(s, a; \omega)$

- Keep track of two sets of parameters

# Actor-Critic pseudocode

- Input: parameterised forms for $\pi_\theta(s|a)$ and $V_\omega(s)$ <span style="color:red">What forms could we use?</span>

- Input: learning rates $\alpha_\omega > 0$ and $\alpha_\theta > 0$

- For each episode:
  - Initialise $s$
  - For each time step:
    - Choose $a \sim \pi_\theta(s|a)$     <span style="color:red">Policy is stochastic: this is a random draw</span>
    - Take $a$, observe $s', r$
    - $\delta \leftarrow r + \gamma V_\omega(s') - V_\omega(s)$     <span style="color:red">Compute TD error</span>
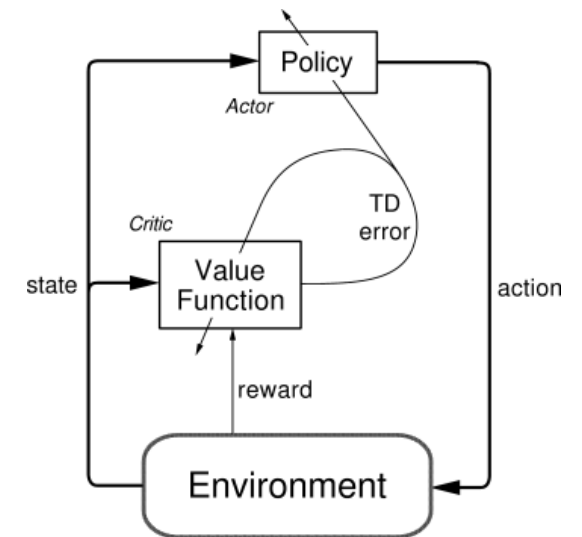    - $\omega \leftarrow \omega + \alpha_\omega \delta \nabla_\omega V_\omega(s)$
    - $\theta \leftarrow \theta + \alpha_\theta \delta \nabla_\theta \log \pi_\theta(a|s)$     <span style="color:red">Update parameters by gradient ascent</span>
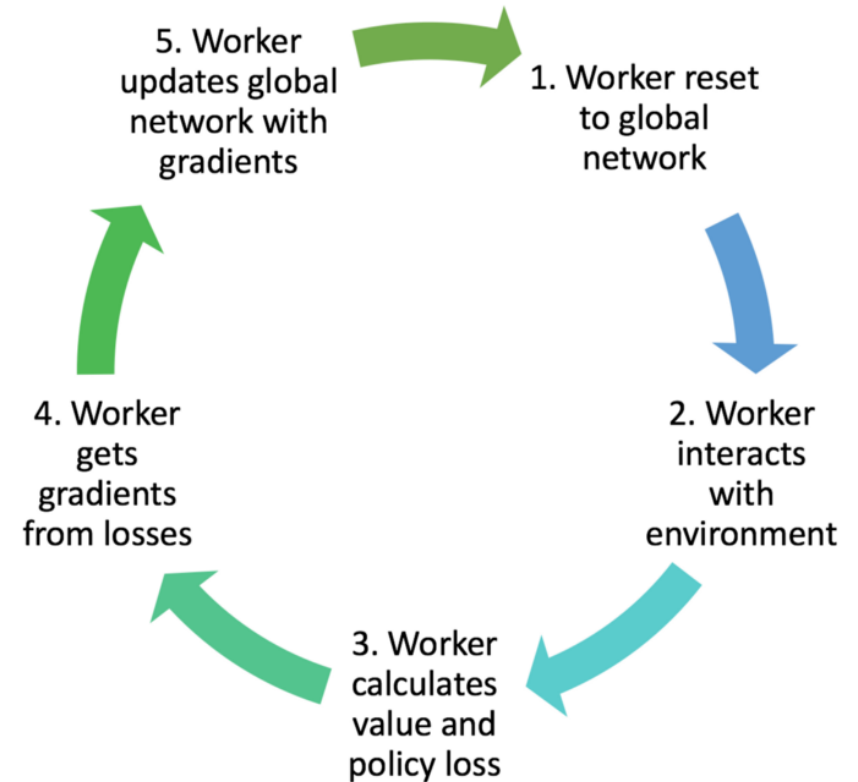    - $s \leftarrow s'$

# Many different ways to do the updates

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a)\ v_t\right] \qquad \text{REINFORCE}$$

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a)\ Q^w(s, a)\right] \qquad \text{Q Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a)\ A^w(s, a)\right] \qquad \text{Advantage Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a)\ \delta\right] \qquad \text{TD Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a)\ \delta e\right] \qquad \text{TD}(\lambda)\ \text{Actor-Critic}$$

$$G_\theta^{-1}\nabla_\theta J(\theta) = w \qquad \text{Natural Actor-Critic}$$

Benjamin Rosman

# Asynchronous Advantage Actor-Critic (A3C)

[Mnih et al., 2016]

- Actor-Critic can be easily parallelised

- Why is this useful?
  - Speed up exploration of state space

- Have multiple agents training with shared parameters



5. Worker updates global network with gradients

1. Worker reset to global network

2. Worker interacts with environment

3. Worker calculates value and policy loss

4. Worker gets gradients from losses

**Algorithm S3** Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

// Assume global shared parameter vectors $\theta$ and $\theta_v$ and global shared counter $T = 0$
// Assume thread-specific parameter vectors $\theta'$ and $\theta'_v$

Spin up a new agent/thread

Initialize thread step counter $t \leftarrow 1$
**repeat**
    Reset gradients: $d\theta \leftarrow 0$ and $d\theta_v \leftarrow 0$.
    Synchronize thread-specific parameters $\theta' = \theta$ and $\theta'_v = \theta_v$
    $t_{start} = t$

Use global parameters

    Get state $s_t$
    **repeat**
        Perform $a_t$ according to policy $\pi(a_t|s_t; \theta')$
        Receive reward $r_t$ and new state $s_{t+1}$

Act

        $t \leftarrow t + 1$
        $T \leftarrow T + 1$
    **until** terminal $s_t$ **or** $t - t_{start} == t_{max}$

$$R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t \text{// Bootstrap from last state} \end{cases}$$

    **for** $i \in \{t - 1, \dots, t_{start}\}$ **do**
        $R \leftarrow r_i + \gamma R$

Update local parameters using advantage functions

        Accumulate gradients wrt $\theta'$: $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$
        Accumulate gradients wrt $\theta'_v$: $d\theta_v \leftarrow d\theta_v + \partial(R - V(s_i; \theta'_v))^2/\partial\theta'_v$
    **end for**
    Perform asynchronous update of $\theta$ using $d\theta$ and of $\theta_v$ using $d\theta_v$.

Update global parameters

**until** $T > T_{max}$
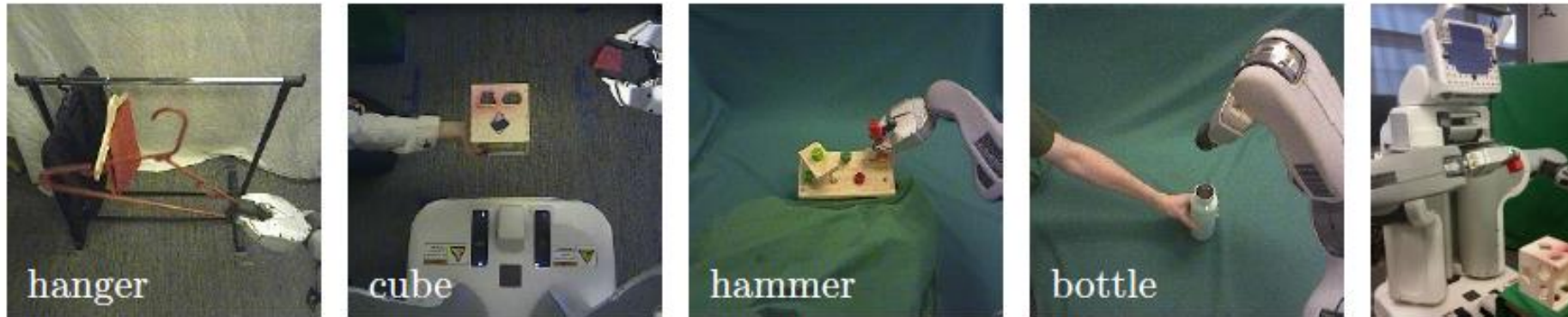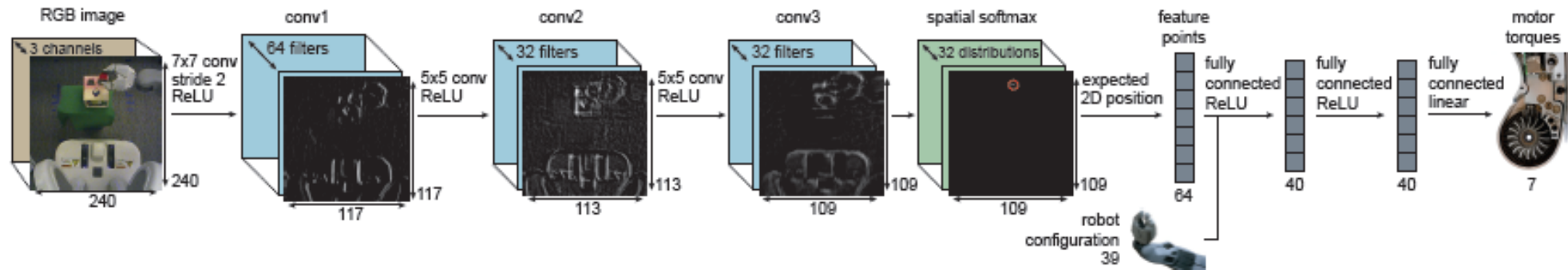
[Mnih et al., 2016]
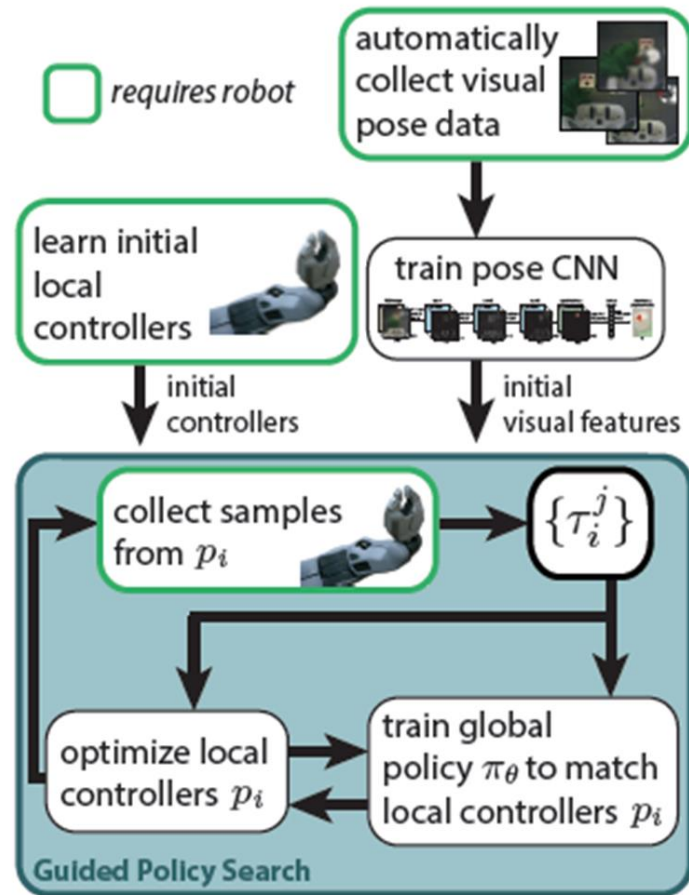
# Deep policy search



Figure 1: Our method learns visuomotor policies that directly use camera image observations (left) to set motor torques on a PR2 robot (right).



[Levine et al., 2016]

# Deep policy search



[Levine et al., 2016]

# Robotics



Learned Visuomotor Policy: Shape sorting cube

[Levine et al., 2016]

# Conclusion

- Recap:
  - The RL setting, MDPs
  - Rewards and value functions
  - Q-learning
  - Function approximation → DQN
- Policy-based methods
  - Gradient free
    - Hill climbing
  - Gradient based
    - Policy Gradient Theorem
    - REINFORCE
    - Baselines
    - Actor-Critic (A3C)
  - Incorporate ideas from supervised learning, deep learning, etc.