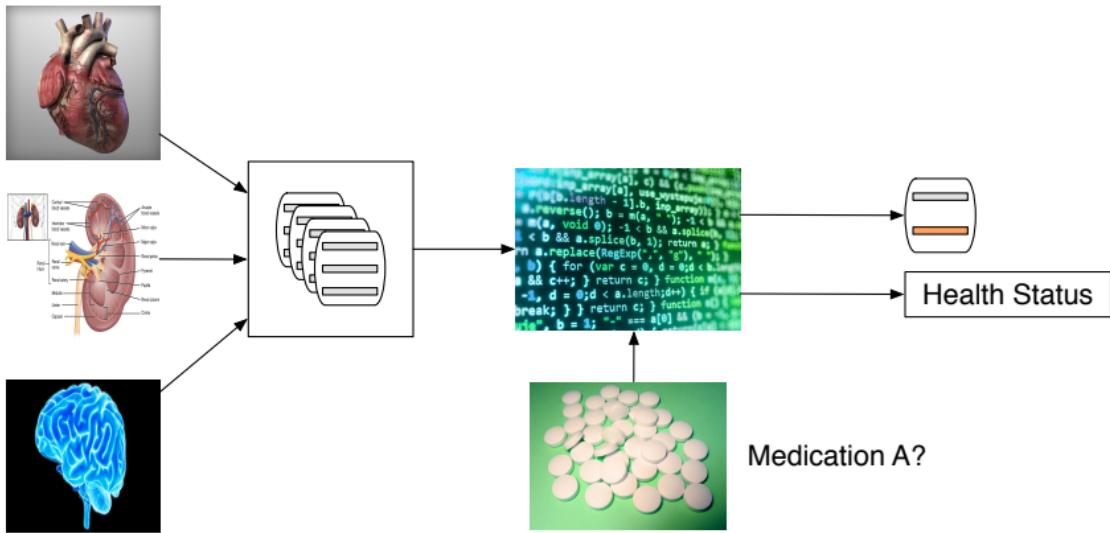


Machine Learning Reinforcement Learning

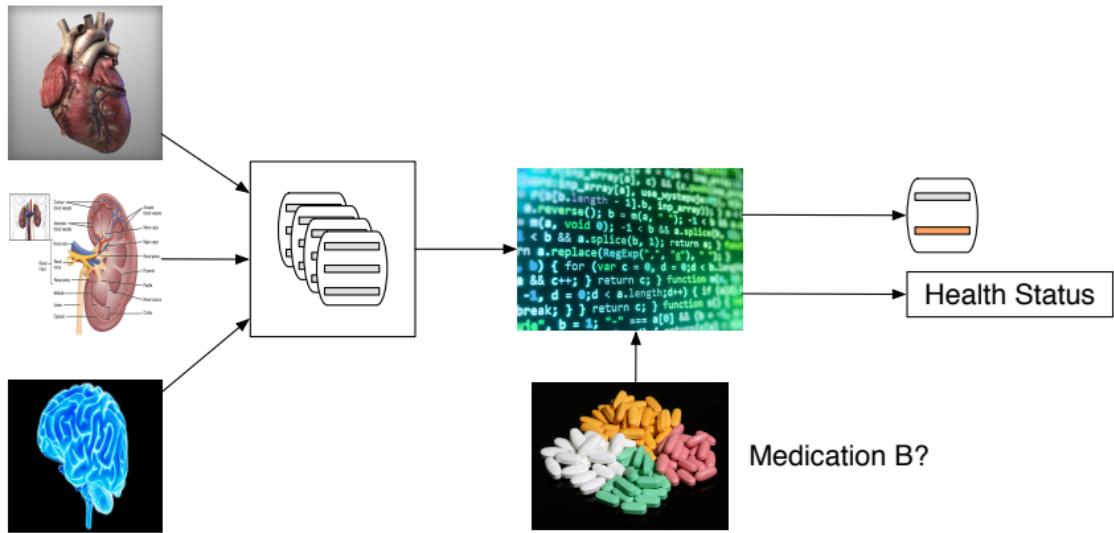
Rajesh Ranganath

- We have a goal
- We want to take actions to reach that goal
- The actions depend on the **state**

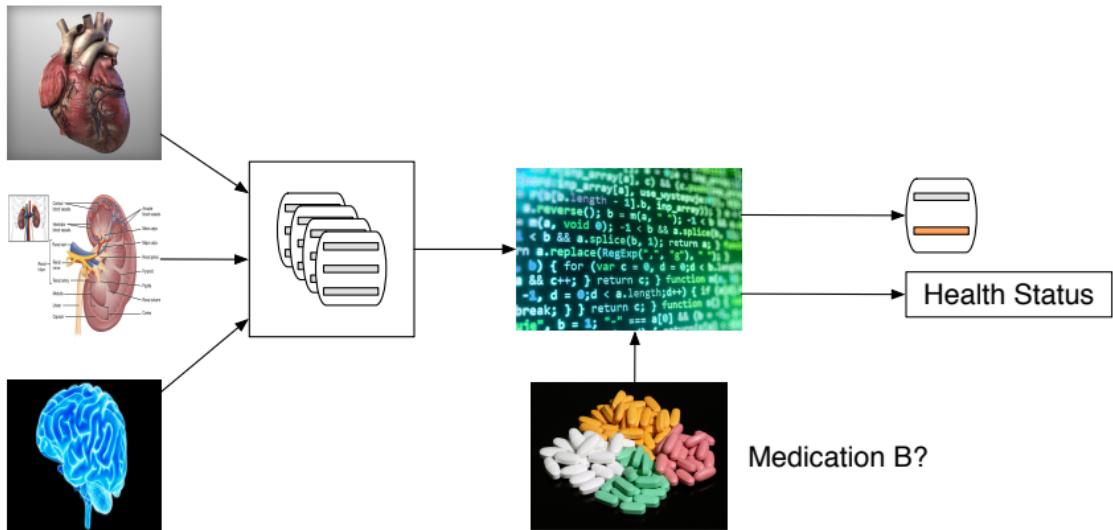
Interventions in Healthcare



Interventions in Healthcare



Medication A or Medication B?



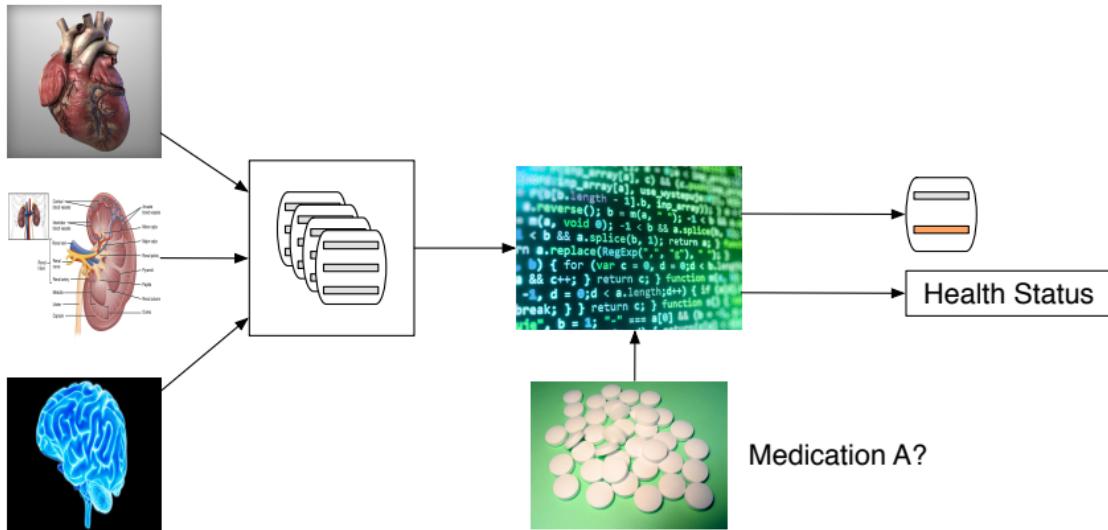
Question: Should we give medication A or medication B?

Causal Inference

Causal inference seeks to estimate the effect of an intervention

- Which statin to give to a patient with hyperlipidemia?
- Which medication to give to a depressed person?
- Which patients to give hospice care to?

Medication A Dosing?



Question: How much of medication A should we give every week?

Solutions?

- Build a time series model and make prediction?
Why should this work?
- Do causal inference from a fixed time point?
Sequential Decision Making?

Markov Decision Processes

A Markov Decision Process

- S : A collection of states \mathbf{s}
- A : A collection of actions \mathbf{a} to take in each state
- $p(\mathbf{s}_{t+1} | \mathbf{a}_t, \mathbf{s}_t)$: State transition
- $r(r_{t+1} | \mathbf{a}_t, \mathbf{s}_t)$: Reward distribution
Don't really need to be random

Markov Decision Processes

Mapping the medical example

- S : The physiological state of a person
- A : Drug A or B along with dosage
- $p(s_{t+1} | a_t, s_t)$: How physiology changes with the drug/dosage given
- $r(r_{t+1} | a_t, s_t)$: Quality of life for day after medication

Why Markov?

The Markov assumptions says

Things only depend on the previous time step

Is this okay?

Why Markov?

The Markov assumptions says

Things only depend on the previous time step

Is this okay?

Can append previous states into one big one. All Markov

What about states that are hidden?

Markov Decision Processes

What's the goal of working with MDPs?

Markov Decision Processes

What's the goal of working with MDPs?

Find a policy: π

$$\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$$

that maximizes total γ -weighted reward

$$\sum_t \gamma_t r_t$$

How do we do this?

First question: What do we know?

Let's assume

- We can simulate a policy
- Observe rewards for the that policy
- A finite time horizon T

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

What's random?

- The actions sampled from the policy
- The state transition $p(s_{t+1} | s_t, a_t)$
- The rewards $r(r_{t+1} | s_t, a_t)$

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

How can we optimize this?

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

How can we optimize this?

Score function gradients

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

$$\nabla_\theta \mathcal{L}(\theta) = \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau] = \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) r_\tau]$$

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

$$\nabla_\theta \mathcal{L}(\theta) = \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau] = \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) r_\tau]$$

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

$$\hat{\nabla}_\theta \mathcal{L}(\theta) = \frac{1}{S} \sum_{s=1}^S \nabla_\theta \log \pi_\theta(\tau_s) r_{\tau^s}$$

Algorithm

1. Sample from current policy
2. Observe state transition and rewards
3. Compute Monte Carlo unbiased gradient
4. Maximize via stochastic optimization

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

$$\nabla_\theta \mathcal{L}(\theta) = \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau] = \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) r_\tau]$$

Problem?

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

$$\nabla_\theta \mathcal{L}(\theta) = \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau] = \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) r_\tau]$$

Problem?

Variance!

Policy Gradients

Goal maximize

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r_\tau]$$

A trajectory τ of actions, states and rewards

$$r_\tau = \sum_{t=1}^T \gamma_t r_t$$

- Can Rao-Blackwellize
- Use control variates or baseline using $\mathbb{E}[\nabla_\theta \log \pi_\theta(\tau)] = 0$

$$\begin{aligned}\nabla_\theta \mathcal{L}(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) r_\tau] \\ &= \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) (r_\tau - b)]\end{aligned}$$

Actor-Critic Methods

Assume $\gamma_t = 1$ for simplicity

Actor-Critic Methods

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(\tau) r_{\tau}] \\&= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[r_{\tau} \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right] \\&= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t'=1}^T r_{t'} \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right] \\&= \mathbb{E} \left[\mathbb{E} \left[\sum_{t'=1}^T r_{t'} \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) | \mathbf{s}_{\leq t}, \mathbf{a}_{<t} \right] \right] \\&= \mathbb{E} \left[\mathbb{E} \left[\sum_{t=1}^T \sum_{t'=1}^t r_{t'} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) | \mathbf{s}_{\leq t}, \mathbf{a}_{<t} \right] \right] \\&\quad + \mathbb{E} \left[\mathbb{E} \left[\sum_{t=1}^T \sum_{t'=t+1}^T r_{t'} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) | \mathbf{s}_{\leq t}, \mathbf{a}_{<t} \right] \right] \\&= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=1}^T \sum_{t'=t+1}^T r_{t'} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right]\end{aligned}$$

Actor-Critic Methods

$$\sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_\theta} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) \nabla_\theta \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) \right]$$

Define future reward estimate

$$\sum_{t'=t+1}^T r_{t'}$$

Policy gradient uses sample based future reward estimate.
Other ideas?

Actor-Critic Methods

$$\sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_\theta} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) \nabla_\theta \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) \right]$$

Define future reward estimate

$$\sum_{t'=t+1}^T r_{t'}$$

Policy gradient uses sample based future reward estimate.
Other ideas?

What if the reward was estimated directly?

$$\sum_{t=1}^T \mathbb{E} \left[\mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right] \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right]$$

Define the *Q-function*:

$$Q_t^{\pi}(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right]$$

$$\sum_{t=1}^T \mathbb{E} \left[\mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right] \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right]$$

Take

$$Q_t^{\pi}(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right]$$

Substitute

$$\nabla_{\theta} \mathcal{L} = \sum_{t=1}^T \mathbb{E} \left[Q_t^{\pi}(\mathbf{a}_t, \mathbf{s}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right]$$

What about variance?

Define the *Value function*

$$V_t^\pi(\mathbf{s}_t) = \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} \mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right]$$

Then

$$\sum_{t=1}^T \mathbb{E} [V_t^\pi(\mathbf{s}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)] = 0$$

Can use the value function to reduce variance!

Define the *Advantage function*

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = Q_t^\pi(\mathbf{a}_t, \mathbf{s}_t) - V_t^\pi(\mathbf{s}_t)$$

- Measures action \mathbf{a}_t better than the average policy action
- It's the advantage of action \mathbf{a}_t over average under π

$$\nabla_{\theta} \mathcal{L} = \sum_{t=1}^T \mathbb{E} [A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)]$$

How do we make practical use of this decomposition?

Rewriting

$$Q_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right]$$

with

$$V_t^\pi(\mathbf{s}_t) = \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} \mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right],$$

yields

$$\begin{aligned} Q_t^\pi(\mathbf{a}_t, \mathbf{s}_t) &= \mathbb{E} \left[\left(\sum_{t'=t+1}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right] = \mathbb{E} \left[r_{t+1} + \left(\sum_{t'=t+2}^T r_{t'} \right) | \mathbf{a}_t, \mathbf{s}_t \right] \\ &= \mathbb{E} \left[r_{t+1} + \mathbb{E} \left[\left(\sum_{t'=t+2}^T r_{t'} \right) | \mathbf{a}_{t+1}, \mathbf{s}_{t+1} \right] | \mathbf{a}_t, \mathbf{s}_t \right] \\ &= \mathbb{E} \left[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t \right] \end{aligned}$$

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

Monte Carlo

$$V_t^\pi(\mathbf{s}_t) \approx \frac{1}{J} \sum_{j=1}^J \sum_{t'=t+1}^T r_{t'}^j$$

Needs to sample sub-trajectories

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

Single Monte Carlo Estimate

$$V^\pi(\mathbf{s}_t) \approx \sum_{t'=t+1}^T r_{t'}$$

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

What about regression?

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

Data

$$\left(s_t^j, \sum_{t'=t+1}^T r_{t'}^j \right)_{j=1}^J$$

Take a parametric function f_ϕ

$$\min_\phi \sum_{j=1}^J \left(\sum_{t'=t+1}^T r_{t'}^j - f_\phi(s_t^j) \right)^2$$

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

Can we reuse the previous estimate?

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

Split

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_{t+1}) + r_{t+1}$$

Use current reward along with previous estimate of future rewards

Advantage Actor-Critic

Simplify

$$A_t^\pi(\mathbf{a}_t, \mathbf{s}_t) = \mathbb{E}[r_{t+1} + V_{t+1}^\pi(\mathbf{s}_{t+1}) | \mathbf{a}_t, \mathbf{s}_t] - V_t^\pi(\mathbf{s}_t)$$

Only need to estimate $V_t^\pi(\mathbf{s}_t)$. How?

$$\min_{\phi} \sum_{j=1}^J (f_\phi^*(s_{t+1}^j) + r_{t+1}^j - f_\phi(s_t^j))^2$$

- Trades bias for variance
- Can put together into an algorithm

Start with initial policy π_θ and value function f_ϕ

Run until convergence

1. Sample from policy π_θ
2. Observe states and rewards
3. Update policy with policy gradients

$$\nabla_\theta \mathcal{L} = \sum_{t=1}^T \mathbb{E}[A^\pi(\mathbf{a}_t, \mathbf{s}_t) \nabla_\theta \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)]$$

4. Update value function with stochastic gradients

$$\nabla_\phi = \sum_{j=1}^J (f_\phi^*(s_{t+1}^j) + r_{t+1}^j - f_\phi(s_t^j)) \nabla_\phi f_\phi(s_t^j)$$

Actor-Critic with Advantage function

What's the actor? What's the critic?

- Actor: Policy
- Critic: Advantage, Q, Value functions

Can we learn a policy with just critics?

- Actor: Policy
- Critic: Advantage, Q, Value functions

Can we learn a policy with just critics?

- $Q^\pi(s, a)$: Reward under policy π for a when in s
- Use Q^π to find a better π^*

$$\pi^*(a | s) = \delta_{a=\operatorname{argmax}_{a'} Q(a', s)}$$

Pick the best action

Fitted Value Iteration

Assume reward is $\sum_{t=0}^{\infty} \gamma^t r_t$, then

$$\begin{aligned} Q^\pi(\mathbf{s}, \mathbf{a}) &= \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{s}, \mathbf{a} \right] \\ &= \mathbb{E}_\pi \left[r(\mathbf{s}, \mathbf{a}) + \sum_{t=1}^{\infty} \gamma^t r_t \mid \mathbf{s}, \mathbf{a} \right] \\ &= r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})} \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{s}' \right] \\ &= r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})} [V(\mathbf{s}')] \end{aligned}$$

Fitted Value Iteration

Assume reward is $\sum_{t=0}^{\infty} \gamma^t r_t$, then

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t | s, a \right] \\ &= \mathbb{E}_\pi \left[r(s, a) + \sum_{t=1}^{\infty} \gamma^t r_t | s, a \right] \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim p(s' | s, a)} \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t | s' \right] \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim p(s' | s, a)} [V(s')] \end{aligned}$$

Can use this to build an algorithm

- $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s' | s, a)} [V(s')]$
- $V(s) = \underbrace{\max_a Q(s, a)}_{\text{Choose Best Policy}}$

Fitted Value Iteration

Assume reward is $\sum_{t=0}^{\infty} \gamma^t r_t$, then

$$\begin{aligned} Q^\pi(\mathbf{s}, \mathbf{a}) &= \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{s}, \mathbf{a} \right] \\ &= \mathbb{E}_\pi \left[r(\mathbf{s}, \mathbf{a}) + \sum_{t=1}^{\infty} \gamma^t r_t \mid \mathbf{s}, \mathbf{a} \right] \\ &= r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})} \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{s}' \right] \\ &= r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})} [V(\mathbf{s}')] \end{aligned}$$

Storing all of V can be hard. Use a function f_ϕ

- $y_j = \max_{a'} r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' \mid \mathbf{s}_j, \mathbf{a}')} [f_\phi(\mathbf{s}')]$

- $\min_\phi \sum_{j=1}^J (y_j - f_\phi(\mathbf{s}_j))^2$

Needs knowledge of dynamics

Fitted Q Iteration

Data comes in tuples $\mathbf{s}', \mathbf{s}, \mathbf{a}, r$

Idea: Use Q -function to approximate the Value function

$$V(\mathbf{s}') = \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}')$$

Use function to learn Q

- Step one:

$$\begin{aligned}y_j &= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}_j, \mathbf{a}_j)} [f_\phi(\mathbf{s}')] \\&= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}'_j, \mathbf{a}')$$

- $\min_{\phi} \sum_{j=1}^J (y_j - Q_\phi(\mathbf{s}_j, \mathbf{a}_j))^2$

Doesn't require samples from a particular policy!

Fitted Q Iteration

Data comes in tuples $\mathbf{s}', \mathbf{s}, \mathbf{a}, r$

Idea: Use Q -function to approximate the Value function

$$V(\mathbf{s}') = \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}')$$

Use function to learn Q

- Step one:

$$\begin{aligned}y_j &= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}_j, \mathbf{a}_j)} [f_\phi(\mathbf{s}')] \\&= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}'_j, \mathbf{a}')$$

- $\min_{\phi} \sum_{j=1}^J (y_j - Q_\phi(\mathbf{s}_j, \mathbf{a}_j))^2$

Doesn't require samples from a particular policy!

- How many y_j 's to collect?
- What if the collection is too big?

Online Q Learning

Collection a single tuple $\mathbf{s}', \mathbf{s}, \mathbf{a}, r$ and update Q with it

- Take some action \mathbf{a}_j in \mathbf{s}_j , observe $\mathbf{s}'_j, r(\mathbf{s}_j, \mathbf{a}_j)$
- Step one:

$$\begin{aligned}y_j &= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}_j, \mathbf{a}_j)} [f_{\phi}(\mathbf{s}')] \\&= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}'_j, \mathbf{a}')\end{aligned}$$

- $\phi = \phi + \alpha \nabla_{\phi} Q_{\phi}(\mathbf{s}_j, \mathbf{a}_j)(y_j - Q_{\phi}(\mathbf{s}_j, \mathbf{a}_j))$

Only needs a single sample at a time

Online Q Learning

Collection a single tuple s', s, a, r and update Q with it

- Take some action a_j in s_j , observe $s'_j, r(s_j, a_j)$
- Step one:

$$\begin{aligned}y_j &= r(s_j, a_j) + \gamma \mathbb{E}_{s' \sim p(s' | s_j, a_j)} [f_\phi(s')] \\&= r(s_j, a_j) + \gamma \max_{a'} Q(s'_j, a')\end{aligned}$$

- $\phi = \phi + \alpha \nabla_\phi Q_\phi(s_j, a_j)(y_j - Q_\phi(s_j, a_j))$

Only needs a single sample at a time

Can something go wrong?

**What if the first time in state s , we take an action
that rarely gives a bad future state/reward**

Exploration vs Exploitation

Exploration:

- Take new actions to understand if they are good

Exploitation:

- Take actions that are known to be good

Tradeoff when trying to learn, but still be good while learning

Online q-iteration exploits too much

Epsilon Greedy and Entropy

Two ways:

- Epsilon Greedy: Sample a random action with probability ϵ otherwise choose maximum
- Entropy Sampling: $\pi(\mathbf{a}_t | \mathbf{s}_t) \propto \exp\left(\frac{1}{T}Q_\phi(\mathbf{a}_t, \mathbf{s}_t)\right)$

Helps trade exploration and exploitation

More trouble?

Why does it work?

What is online Q learning?

What is online Q learning?

Not a gradient method

$$y_j = r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}'_j, \mathbf{a}')$$

Objective

$$\min(y_j - Q_\phi(\mathbf{a}_j, \mathbf{s}_j))^2$$

Same as

$$\min(r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q_\phi(\mathbf{s}'_j, \mathbf{a}') - Q_\phi(\mathbf{a}_j, \mathbf{s}_j))^2$$

Online Q-iteration

$$\boldsymbol{\phi} = \boldsymbol{\phi} + \alpha \nabla_{\boldsymbol{\phi}} Q_\phi(\mathbf{s}_j, \mathbf{a}_j)(y_j - Q_\phi(\mathbf{s}_j, \mathbf{a}_j))$$

is not gradient!

More trouble?

Collection a single tuple $\mathbf{s}', \mathbf{s}, \mathbf{a}, r$ and update Q with it

- Take some action \mathbf{a}_j in \mathbf{s}_j , observe $\mathbf{s}'_j, r(\mathbf{s}_j, \mathbf{a}_j)$
- Step one:

$$\begin{aligned}y_j &= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}_j, \mathbf{a}_j)} [f_{\phi}(\mathbf{s}')] \\&= r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}'_j, \mathbf{a}')\end{aligned}$$

- $\phi = \phi + \alpha \nabla_{\phi} (y_j - Q(\mathbf{a}_j, \mathbf{s}_j))$

How are the tuples collected?

- Start in state \mathbf{s}
- Take action \mathbf{a}
- Go to state \mathbf{s}'
- Set $\mathbf{s} = \mathbf{s}'$

Samples are not iid, standard regression unhappy!

Fixes

Correlated Samples:

- Store old samples in a buffer
- Update Q function with that buffer
- Can be parallelized

Not a proper gradient:

- Use a second Q function

$$\min(r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q_{\hat{\phi}}(\mathbf{s}'_j, \mathbf{a}') - Q_{\phi}(\mathbf{a}_j, \mathbf{s}_j))^2$$

Guarantees?

Deep Q Learning

Deep Q Learning

- Save ϕ as $\hat{\phi}$
- Collect dataset $(\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}'_j, r_j)_{j=1}^J$
- For K iterations:
 - Sample a batch of data from the data set
 - Update parameters

$$\phi = \phi + \alpha \nabla_\phi Q(\mathbf{s}_j, \mathbf{a}_j)(r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q_{\hat{\phi}}(\mathbf{s}'_j, \mathbf{a}') - Q_\phi(\mathbf{a}_j, \mathbf{s}_j))$$

- Use Deep Networks for Q_ϕ

Deep Q Learning

Deep Q Learning

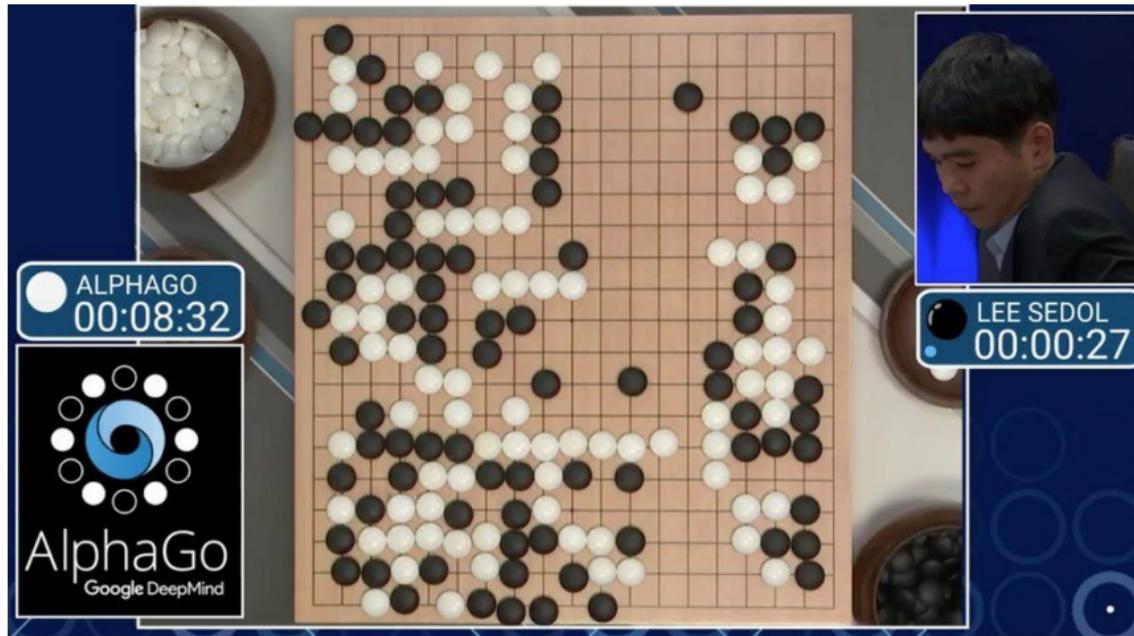
- Save ϕ as $\hat{\phi}$
- Collect dataset $(\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}'_j, r_j)_{j=1}^J$
- For K iterations:
 - Sample a batch of data from the data set
 - Update parameters

$$\phi = \phi + \alpha \nabla_\phi Q(\mathbf{s}_j, \mathbf{a}_j)(r(\mathbf{s}_j, \mathbf{a}_j) + \gamma \max_{\mathbf{a}'} Q_{\hat{\phi}}(\mathbf{s}'_j, \mathbf{a}') - Q_\phi(\mathbf{a}_j, \mathbf{s}_j))$$

- Use Deep Networks for Q_ϕ

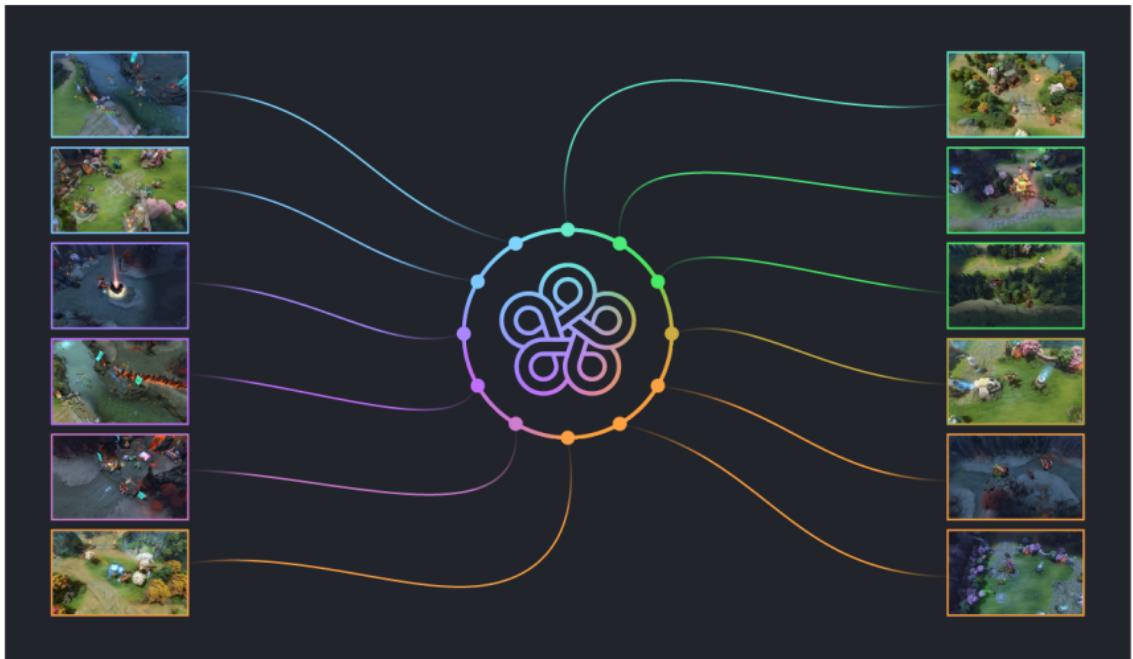
Guarantees?

This has led to lots of fun



[BBC]

This has led to lots of fun



[OpenAI Five]

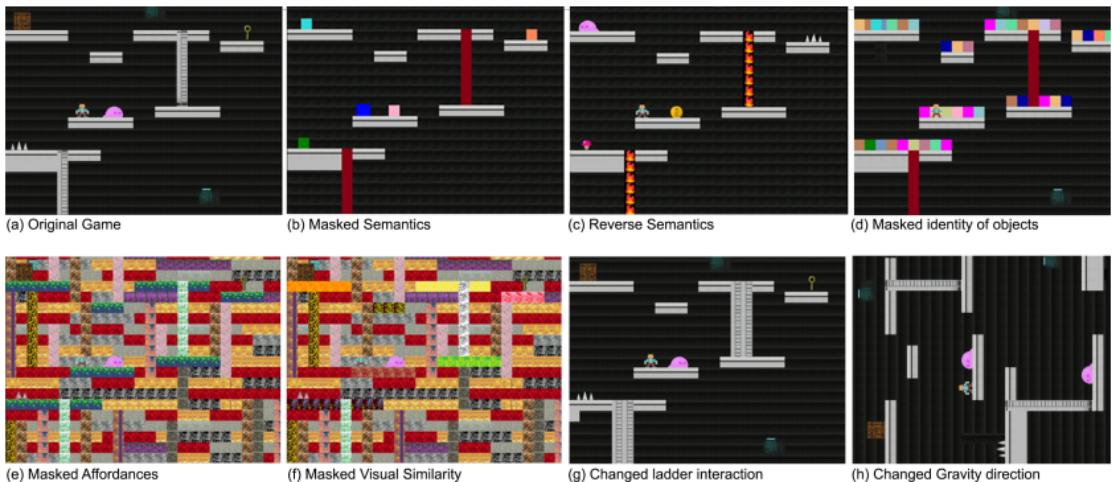
We were expecting to need sophisticated algorithmic ideas, such as hierarchical reinforcement learning, but we were surprised by what we found: the fundamental improvement we needed for this problem was scale. Achieving and utilizing that scale wasn't easy and was the bulk of our research effort!

[OpenAI Five]

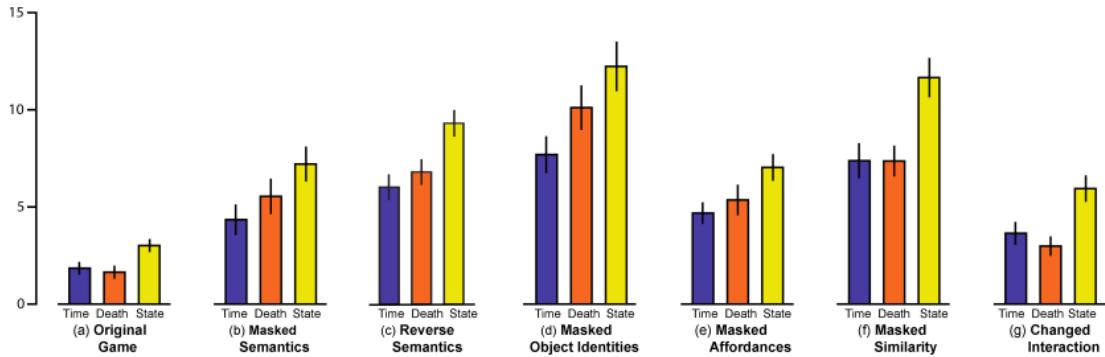
- Today's class went from foundations to deep q learning quickly
- Variants plus massive scale led to some “super-human” performance

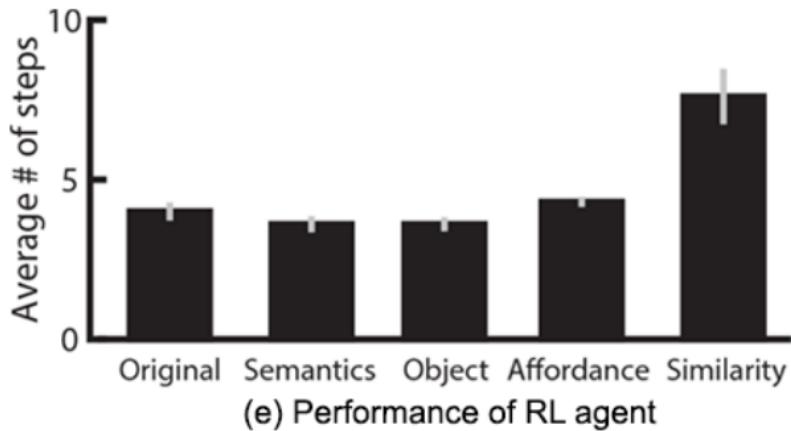
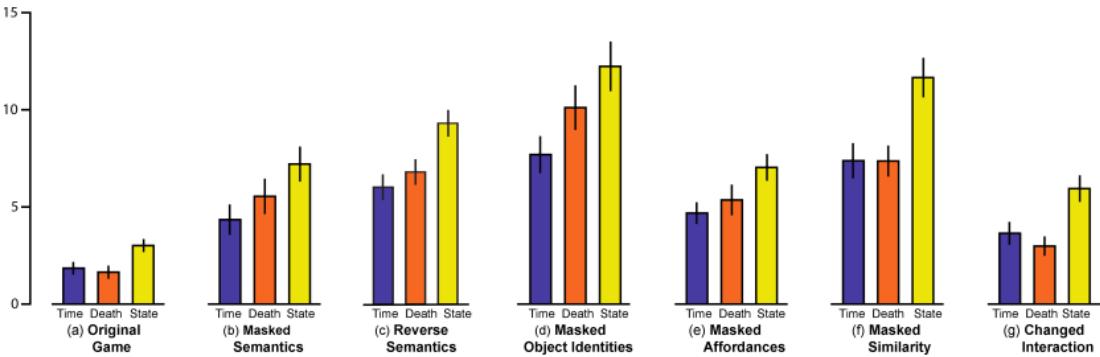
Is everything solved?

Investigating Human Priors for Playing Video Games



[Dubey+ 2018]





[Dubey+ 2018]

A Large Question

How to integrate, retain, and refine knowledge?

Learned about Sequential Decision Making: Markov Decision Processes

Policy Gradients: Directly update a policy of the world

Actor-Critic: Use estimates of future value of state to update policy

Q-learning: Estimate future value of state to define implicit policy

Deep Q-learning: Estimate future value of state to define an implicit policy with deep learning

**What's the underlying principles?
Feels like a collection of bandaids**

Where does the reward come from?