

RL@Coursera

week 5

Policy gradient methods



Yandex
Data Factory

LAMBDA A stylized yellow arrow pointing to the right, with a small red dot at the tip.



Small experiment

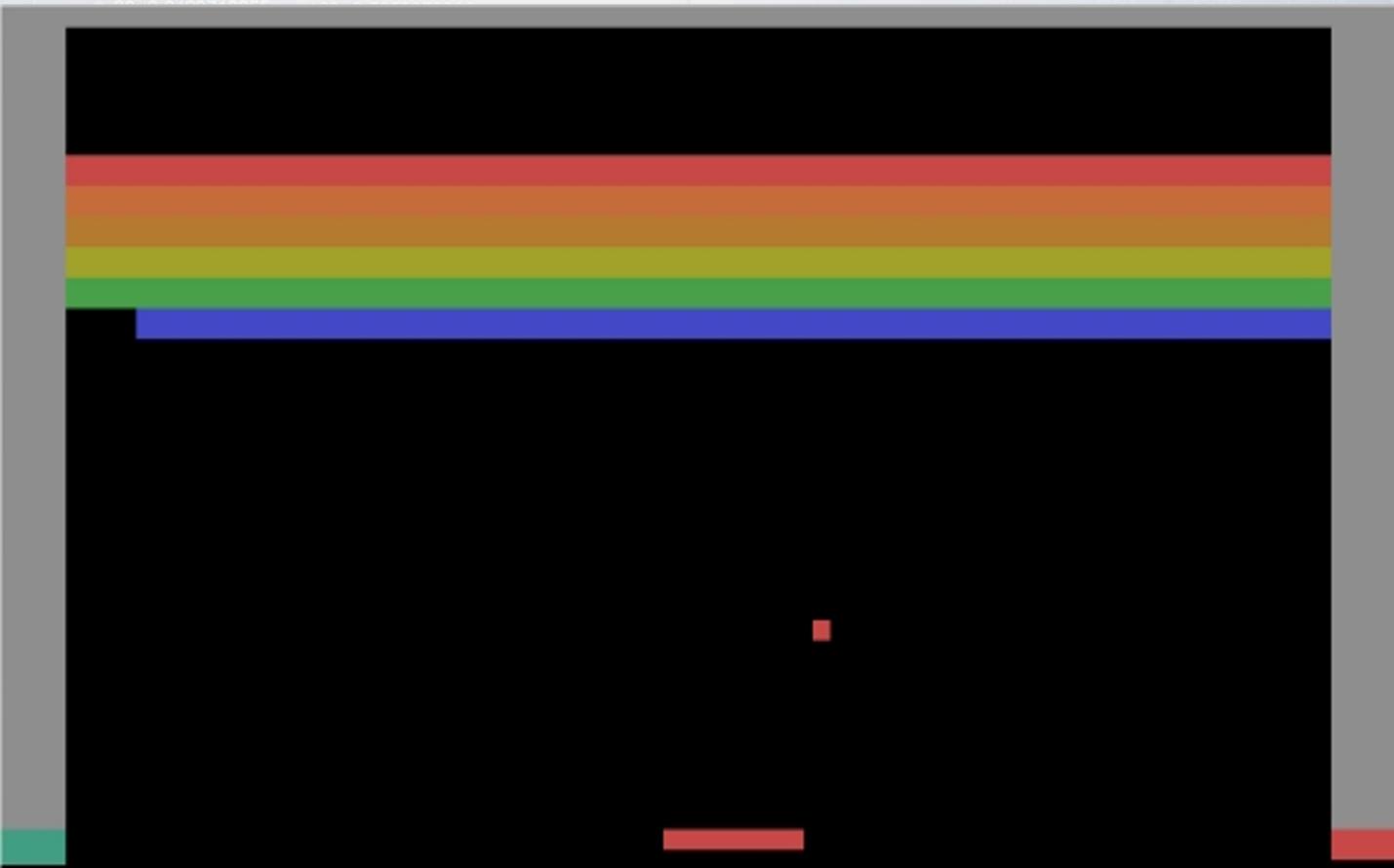
The next slide contains a question

Please respond as fast as you can!



Small experiment

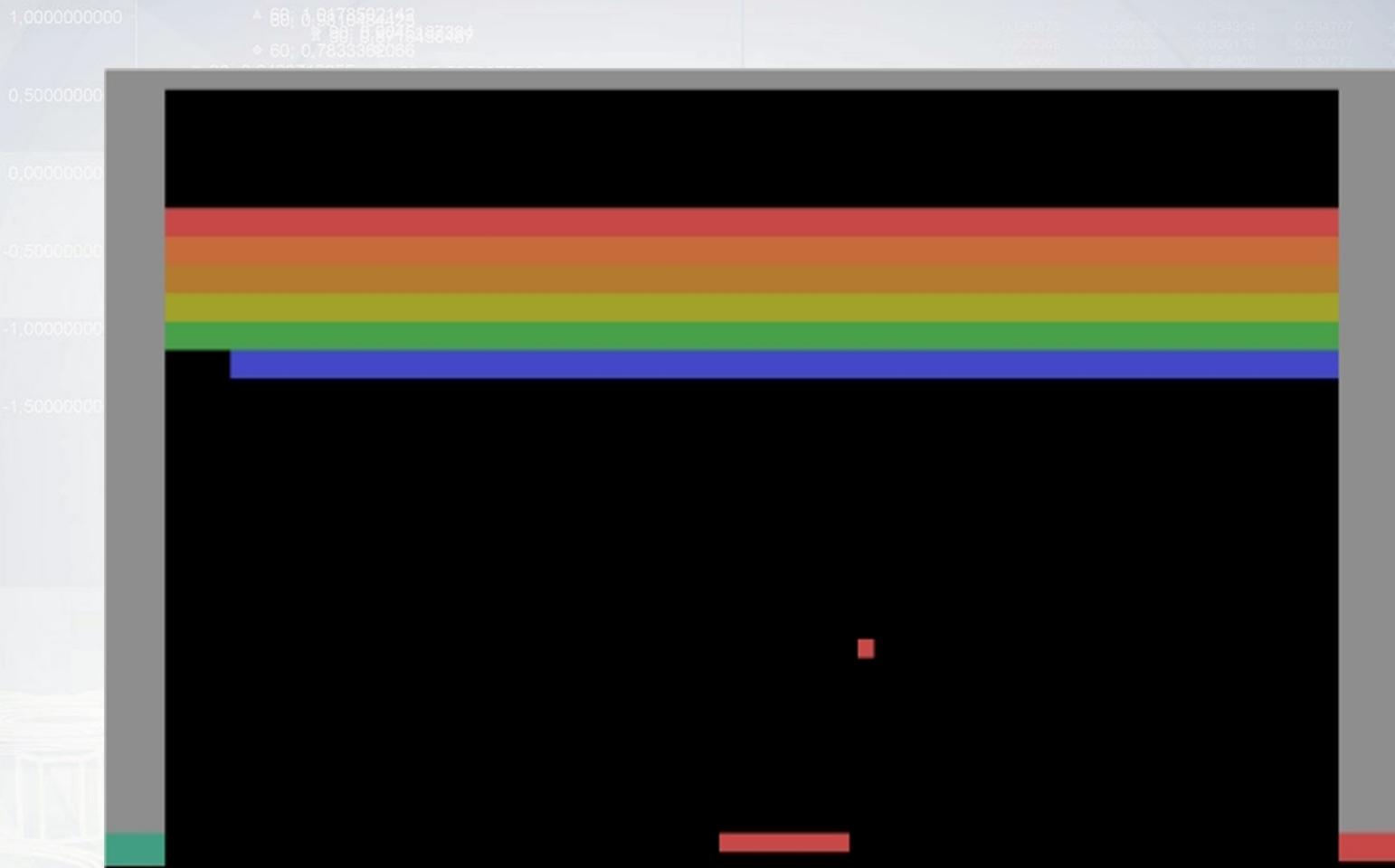
▲ 60; 1,01178592114
● 60; 0,801785186389
◆ 60; 0,7833362066



Left or Right?



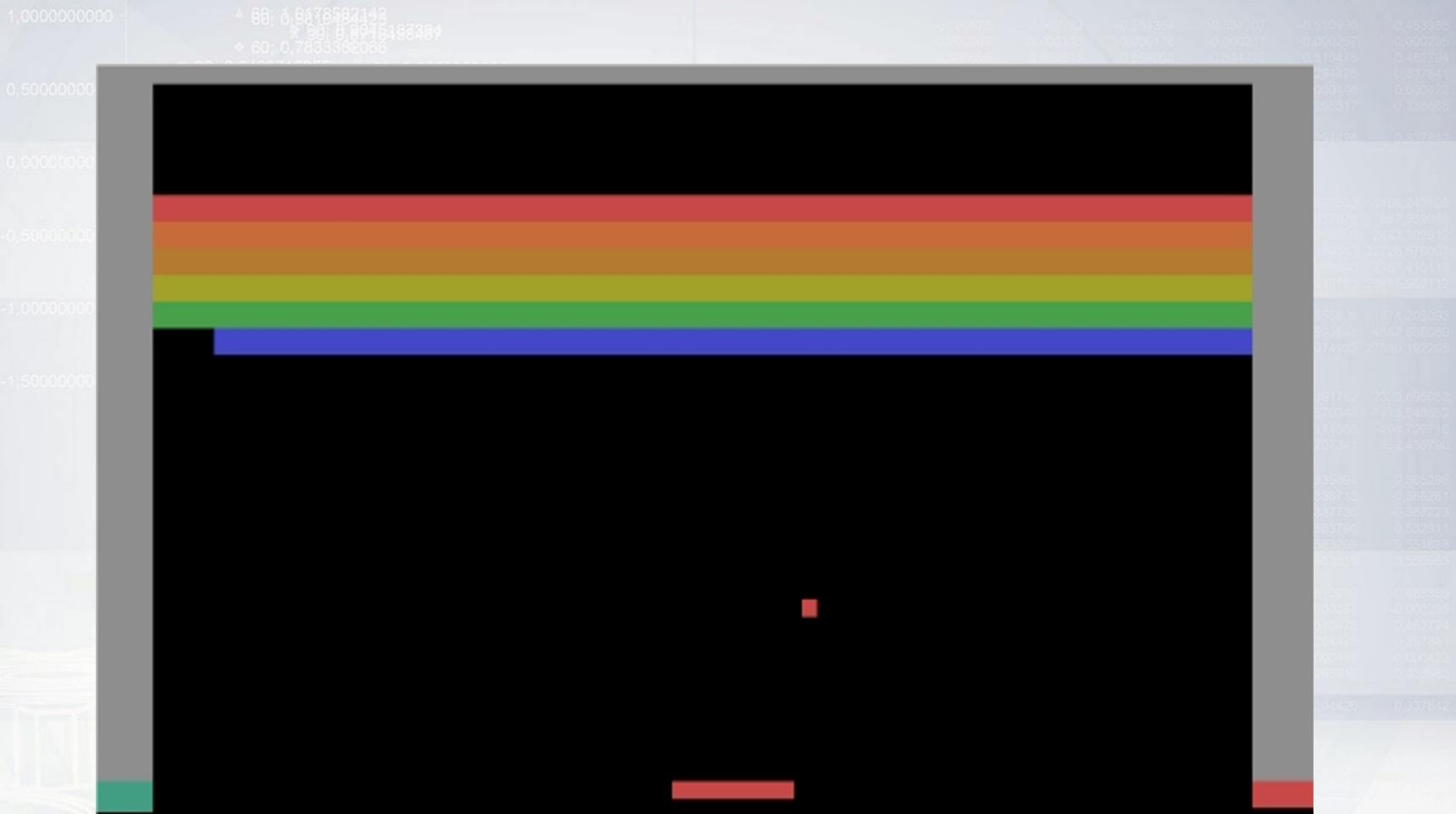
Small experiment



Right! Ready for next one?



Small experiment

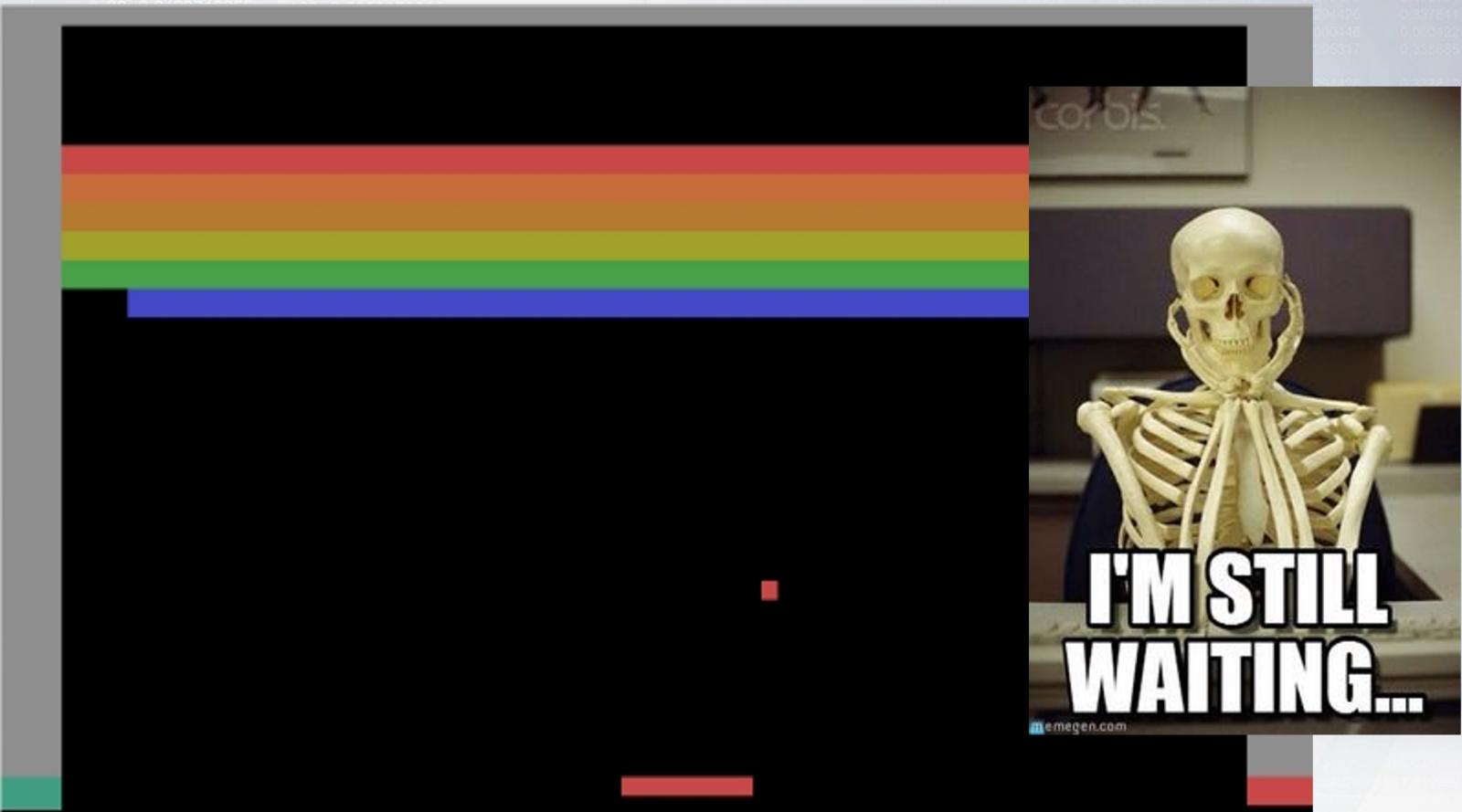


What's $Q(s,\text{right})$ under $\gamma=0.99$?



Small experiment

1.0000000000
▲ 60; 0.98178592142
■ 60; 0.99751486384
◆ 60; 0.7833362066



What's $Q(s,\text{right})$ under $\gamma=0.99$?



Approximation error

$$L \approx E[Q(s_t, a_t) - (r_t + \gamma \cdot \max_{a'} Q(s_t, a'))]^2$$

Simple 2-state

	True world	(A)	(B)
Q(s0,a0)	1	1	2
Q(s0,a1)	2	2	1
Q(s1,a0)	3	3	3
Q(s1,a1)	100	50	100

Trivia: Which prediction is better
(A/B)?



Approximation error

$$L \approx E[Q(s_t, a_t) - (r_t + \gamma \cdot \max_{a'} Q(s_t, a'))]^2$$

Simple 2-state

	True world	(A)	(B)
Q(s0,a0)	1	1	2
Q(s0,a1)	2	2	1
Q(s1,a0)	3	3	3
Q(s1,a1)	100	50	100

better
policy

less
MSE



Approximation error

$$L \approx E[Q(s_t, a_t) - (r_t + \gamma \cdot \max_{a'} Q(s_t, a'))]^2$$

Simple 2-state

	True world	(A)	(B)
Q(s0,a0)	1	1	2
Q(s0,a1)	2	2	1
Q(s1,a0)	3	3	3
Q(s1,a1)	100	50	100

Q-learning will prefer
worse policy (B)!

better
policy

less
MSE



Conclusion

- Often computing q-values is harder than picking optimal actions!
- We could avoid learning value functions by directly learning agent's policy $\pi_\theta(s|a)$

Q: what algorithm works that way?



Conclusion

- Often computing q-values is harder than picking optimal actions!
- We could avoid learning value functions by directly learning agent's policy $\pi_\theta(s|a)$

Q: what algorithm works that way?

(of those we
studied)



Conclusion

- Often computing q-values is harder than picking optimal actions!
- We could avoid learning value functions by directly learning agent's policy $\pi_\theta(s|a)$

Q: what algorithm works that way?

e.g. crossentropy method

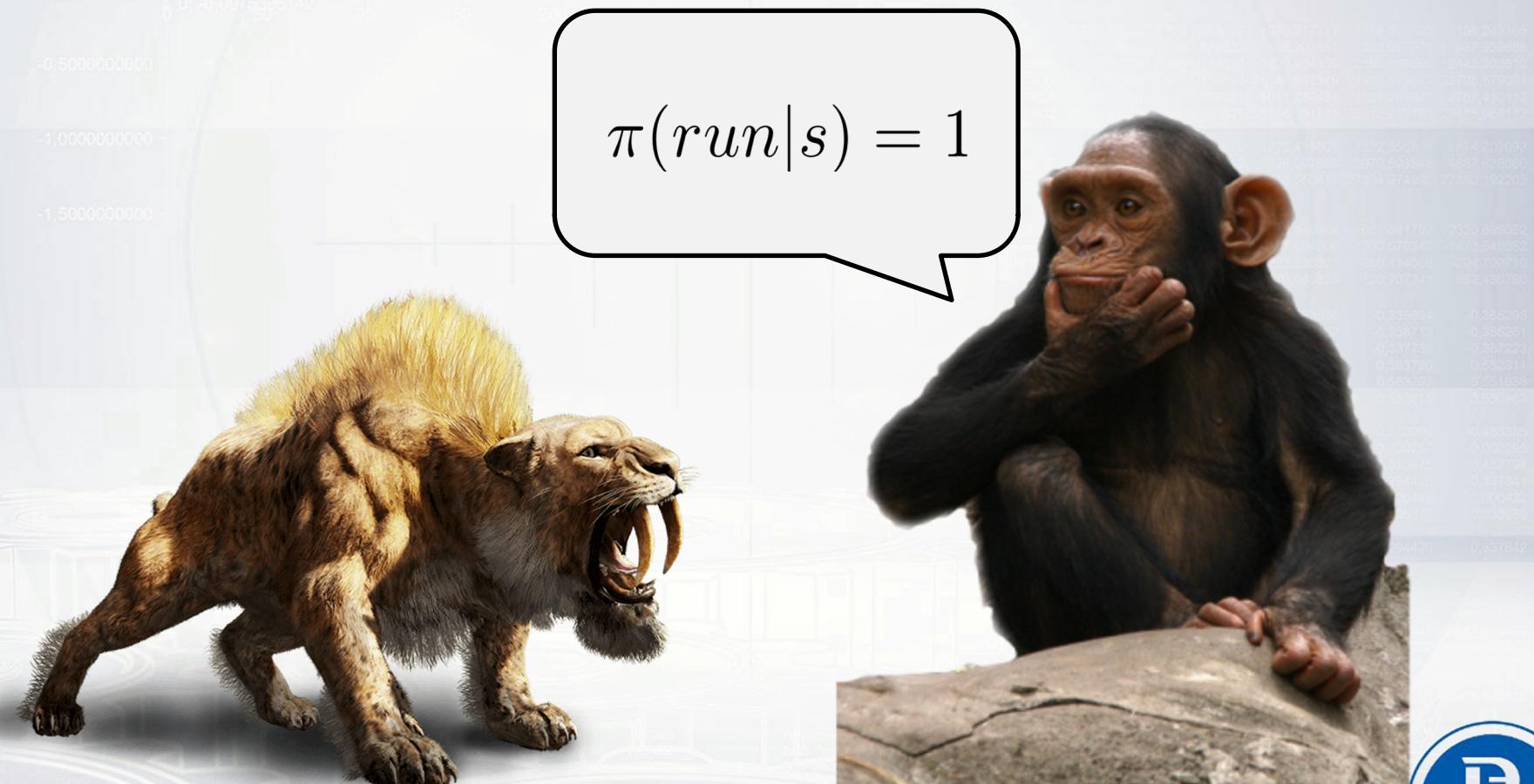


NOT how humans survived

argmax
[$Q(s, \text{pet the tiger})$
 $Q(s, \text{run from tiger})$
 $Q(s, \text{provoke tiger})$
 $Q(s, \text{ignore tiger})]$



how humans survived



Policies

In general, two kinds

Deterministic policy

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

Stochastic policy

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



Policies

In general, two kinds

Deterministic policy

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

Stochastic policy

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

Q: Any case where stochastic is better?



Policies

In general, two kinds

Deterministic policy

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

Stochastic policy

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

e.g. rock-paper-scissors

Q: Any case where stochastic is better?



Policies

In general, two kinds

Deterministic policy

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

same action
each time

Stochastic policy

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

sampling takes
care of
exploration



Policies

In general, two kinds

Deterministic policy

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

same action
each time

Stochastic policy

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

sampling takes
care of
exploration

Q: how to represent policy in continuous action space?



Policies

In general, two kinds

Deterministic policy

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

same action
each time

Stochastic policy

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

sampling takes
care of
exploration

categorical, normal, mixture of normal, whatever



Two approaches

- **Value based:**

Learn value function $Q_\theta(s, a)$ or $V_\theta(s)$

Infer policy $\pi(a|s) = \# \left[a = \underset{a}{\operatorname{argmax}} Q_\theta(s, a) \right]$

- **Policy based:**

Explicitly learn policy $\pi_\theta(s, a)$ $\pi_\theta(s) \rightarrow a$

Implicitly maximize reward over policy



Recap: crossentropy method

- Initialize policy params $\theta_0 \leftarrow \text{random}$

- Loop:

- Sample N sessions
- elite = take M best sessions and concatenate

$$\theta_{i+1} = \theta_i + \alpha \bigtriangledown \sum_i \log \pi_{\theta_i}(a_i|s_i) \cdot [s_i, a_i \in \text{Elite}]$$



Policy gradient main idea

Why so complicated?

We'd rather simply maximize G over pi!



Objective

- Expected reward:

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a, s', a', \dots)$$

...

- Expected discounted reward:

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} G(s, a)$$



Objective

- Expected reward:

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a, s', a', \dots)$$

...

- Expected discounted reward:

$R(z)$ settin

$$G(s, a) = r + \gamma \cdot G(s', a')$$

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} G(s, a)$$



Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a)$$



Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a) = \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds$$



Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a) = \int p(s) \int \pi_\theta(a|s) R(s, a) da ds$$



state visitation frequency
(may depend on policy)

Reward for
1-step session



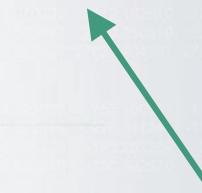
Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a) = \int p(s) \int \pi_\theta(a|s) R(s, a) da ds$$



state visitation frequency
(may depend on policy)



Reward for
1-step session

Q: how do we compute
that?



Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a) = \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds$$

Reward for
agent's action

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in z_i} R(s, a)$$

sample N sessions
by following $\pi_\theta(a|s)$



Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a) = \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds$$

Reward for
agent's action

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in z_i} R(s, a)$$

sample N sessions

Can we optimize policy now?



Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} R(s, a) = \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds$$

parameters “sit” here

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in z_i} R(s, a)$$

We don't know how to compute $dJ/d\theta$



Optimization

- Finite differences

- Change policy a little, evaluate

$$\nabla J \approx \frac{J_{\theta+\epsilon} - J_\theta}{\epsilon}$$

- Stochastic optimization
 - Good old crossentropy method
 - Maximize probability of “elite” actions



Optimization

- Finite differences

- Change policy a little, evaluate

$$\nabla J \approx \frac{J_{\theta+\epsilon} - J_\theta}{\epsilon}$$

- Stochastic optimization
 - Good old crossentropy method
 - Maximize probability of “elite” actions



Optimization

- Finite differences

- Change policy a little, evaluate

$$\nabla J \approx \frac{J_{\theta+\epsilon} - J_\theta}{\epsilon}$$

VERY noisy, especially
if both Js are sampled

- Stochastic optimization
 - Good old crossentropy method
 - Maximize probability of “elite” actions

stochastic MDPs;
Throws a most sessions away



Objective

- Consider an 1-step process for simplicity

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(a|s)}}{E} R(s, a) = \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds$$

Wish list:

- Analytical gradient
- Easy/stable approximations



Logderivative trick

Simple math

$$\nabla \log \pi(z) = ???$$

(try chain rule)



Logderivative trick

Simple math

$$\nabla \log \pi(z) = \frac{1}{\pi(z)} \cdot \pi(z)$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$



Policy gradient

Analytical inference

$$J = \int p(s) \int_{s} \pi_{\theta}(a|s) R(s, a) da ds$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$



Policy gradient

Analytical inference

$$J = \int p(s) \int_{s} \pi_{\theta}(a|s) R(s, a) da ds$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$

$$J = \int p(s) \int_{a} \pi_{\theta}(a|s) \nabla \log \pi_{\theta}(a|s) R(s, a) da ds$$

Trivia: Anything curious about that formula?



Policy gradient

Analytical inference

$$J = \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$

$$J = \int_s p(s) \int_a \nabla \log \pi_\theta(a|s) R(s, a) da ds$$

that's expectation
:)



Discounted reward case

Replace R with Q :

True action

$$E(G[s, a]) \text{ value}$$

a.k.a.

$$J = \int p(s) \int_{s} \pi_{\theta}(a|s) Q(s, a) da ds$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$

$$J = \int p(s) \int_{s} \pi_{\theta}(a|s) \nabla \log \pi_{\theta}(a|s) Q(s, a) da ds$$

that's expectation
:)



Policy gradient (REINFORCE)

- Policy gradient

$$\nabla J = \underset{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}}{E} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$

- Approximate with sampling

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in z_i} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$



REINFORCE algorithm

- Initialize NN weights $\theta_0 \leftarrow \text{random}$

- Loop:

- Sample N sessions \mathbf{z} under current $\pi_\theta(a|s)$
- Evaluate policy gradient

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$

Ascend

$$\theta_{i+1} \leftarrow \theta_i + \alpha \cdot \nabla J$$



REINFORCE algorithm

- Initialize NN weights $\theta_0 \leftarrow \text{random}$

Q: is it on- or off-policy?

- Loop:
 - Sample N sessions \mathbf{z} under current $\pi_\theta(a|s)$
 - Evaluate policy gradient

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$

Ascend

$$\theta_{i+1} \leftarrow \theta_i + \alpha \cdot \nabla J$$



REINFORCE algorithm

- Initialize NN weights $\theta_0 \leftarrow \text{random}$

actions under current policy
= on-policy

- Loop:

- Sample N sessions \mathbf{z} under current $\pi_\theta(a|s)$
- Evaluate policy gradient

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$

Ascend

$$\theta_{i+1} \leftarrow \theta_i + \alpha \cdot \nabla J$$



REINFORCE algorithm

- Initialize NN weights $\theta_0 \leftarrow \text{random}$

- Loop:

- Sample N sessions \mathbf{z} under current $\pi_\theta(a|s)$
- Evaluate policy gradient

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$

Ascend

$$\theta_{i+1} \leftarrow \theta_i + \alpha \cdot \nabla J$$

What is better for learning:
random action in good state or
great action in bad state?



REINFORCE baseline

- Initialize NN weights $\theta_0 \leftarrow \text{random}$

- Loop:

- Sample N sessions \mathbf{z} under current $\pi_\theta(a|s)$
- Evaluate policy gradient

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$

Ascend

$$\theta_{i+1} \leftarrow \theta_i + \alpha \cdot \nabla J$$

$$Q(s, a) = V(s) + A(s, a)$$

Actions influence $A(s, a)$ only, so $V(s)$

is irrelevant



REINFORCE baseline

- Initialize NN weights $\theta_0 \leftarrow \text{random}$

- Loop:

- Sample N sessions \mathbf{z} under current $\pi_\theta(a|s)$
- Evaluate policy gradient

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in z_i} \nabla \log \pi_\theta(a|s) \cdot (Q(s, a) - b(s))$$

Ascend

$$\theta_{i+1} \leftarrow \theta_i + \alpha \cdot \nabla J$$

Anything that doesn't depend on action

ideally, $b(s) = V(s)$



Actor-critic

- Learn both $V(s)$ and $\pi_\theta(a|s)$
- Hope for best of both worlds :)



image: youtube: mr tivikov



Advantage actor-critic

- Idea: learn both $\pi_\theta(a|s)$ and $V_\theta(s)$
- Use $V_\theta(s)$ to learn $\pi_\theta(a|s)$ faster!

Non-trivia: how can we estimate $A(s, a)$ from (s, a, r, s') and V -function?



Advantage actor-critic

- Idea: learn both $\pi_\theta(a|s)$ and $V_\theta(s)$
- Use $V_\theta(s)$ to learn $\pi_\theta(a|s)$ faster!

Non-trivia: how can we estimate $A(s, a)$ from (s, a, r, s') and V -function?



Advantage actor-critic

- Idea: learn both $\pi_\theta(a|s)$ and $V_\theta(s)$
- Use $V_\theta(s)$ to learn $\pi_\theta(a|s)$ faster!

$$A(s, a) = Q(s, a) - V(s)$$

$$Q(s, a) = r + \gamma \cdot V(s')$$

$$A(s, a) = r + \gamma \cdot V(s') - V(s)$$



Advantage actor-critic

- Idea: learn both $\pi_\theta(a|s)$ and $V_\theta(s)$
- Use $V_\theta(s)$ to learn $\pi_\theta(a|s)$ faster!

$$A(s, a) = Q(s, a) - V(s)$$

$$Q(s, a) = r + \gamma \cdot V(s')$$

$$A(s, a) = r + \gamma \cdot V(s') - V(s)$$

Also: n-step version



Advantage actor-critic

- Idea: learn both $\pi_\theta(a|s)$ and $V_\theta(s)$
- Use $V_\theta(s)$ to learn $\pi_\theta(a|s)$ faster!

$$A(s, a) = r + \gamma \cdot V(s') - V(s)$$

$$Q(s, a) = r + \gamma \cdot V(s')$$

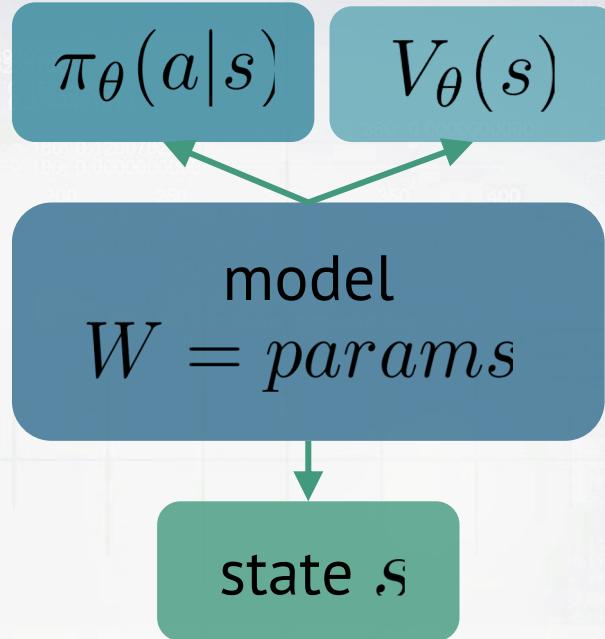
$$\nabla J_{actor} \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in z_i} \nabla \log \pi_\theta(a|s) \cdot A(s, a)$$

consider
const

Trivia: How do we train V then?



Advantage actor-critic



$$\nabla J_{actor} \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(a|s) \cdot A(s, a)$$

$$L_{critic} \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} (V_\theta(s) - [r + \gamma \cdot V(s')])^2$$



value-based Vs policy-based

Value-based

Q-learning, SARSA,value-iteration

Policy-based + Actor-critic

REINFORCE, Advantage Actor-Critic, Crossentropy Method



value-based Vs policy-based



Q-learning, SARSA,value-iteration

Policy-based + Actor-critic

360: 0.0000000000

REINFORCE, Advantage
Actor-Critic, Crossentropy
Method

Your
guess?



Casey Study: A3C

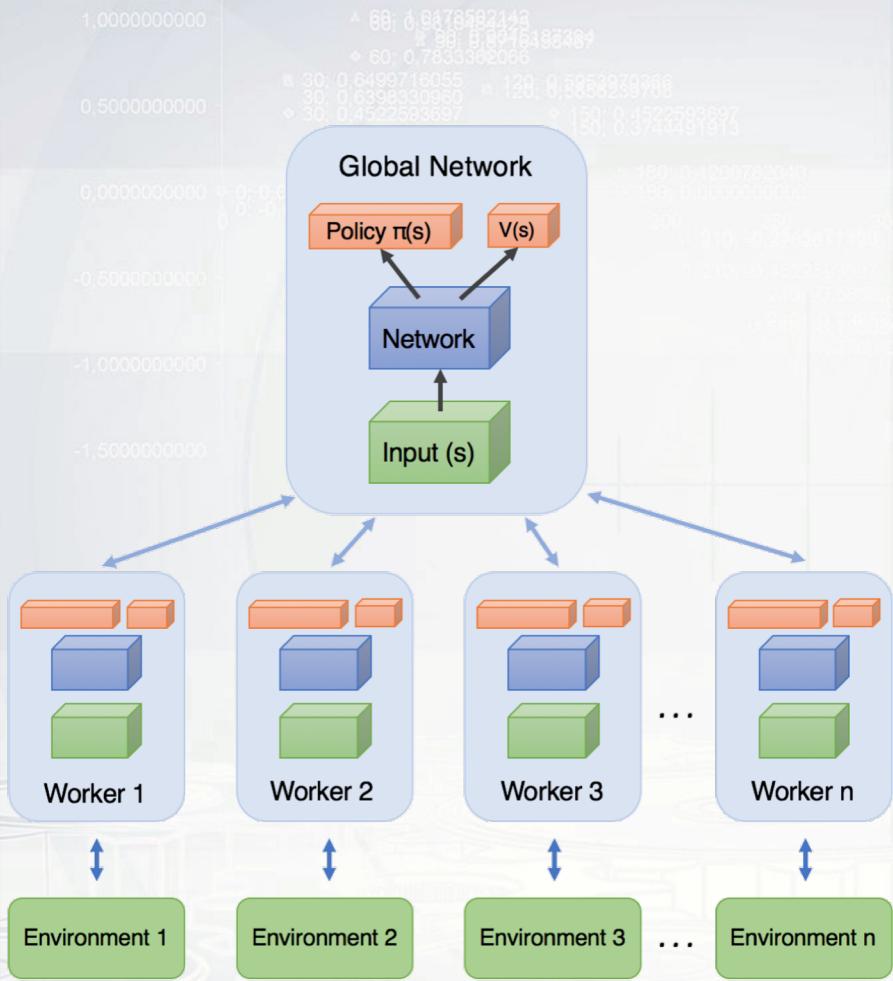


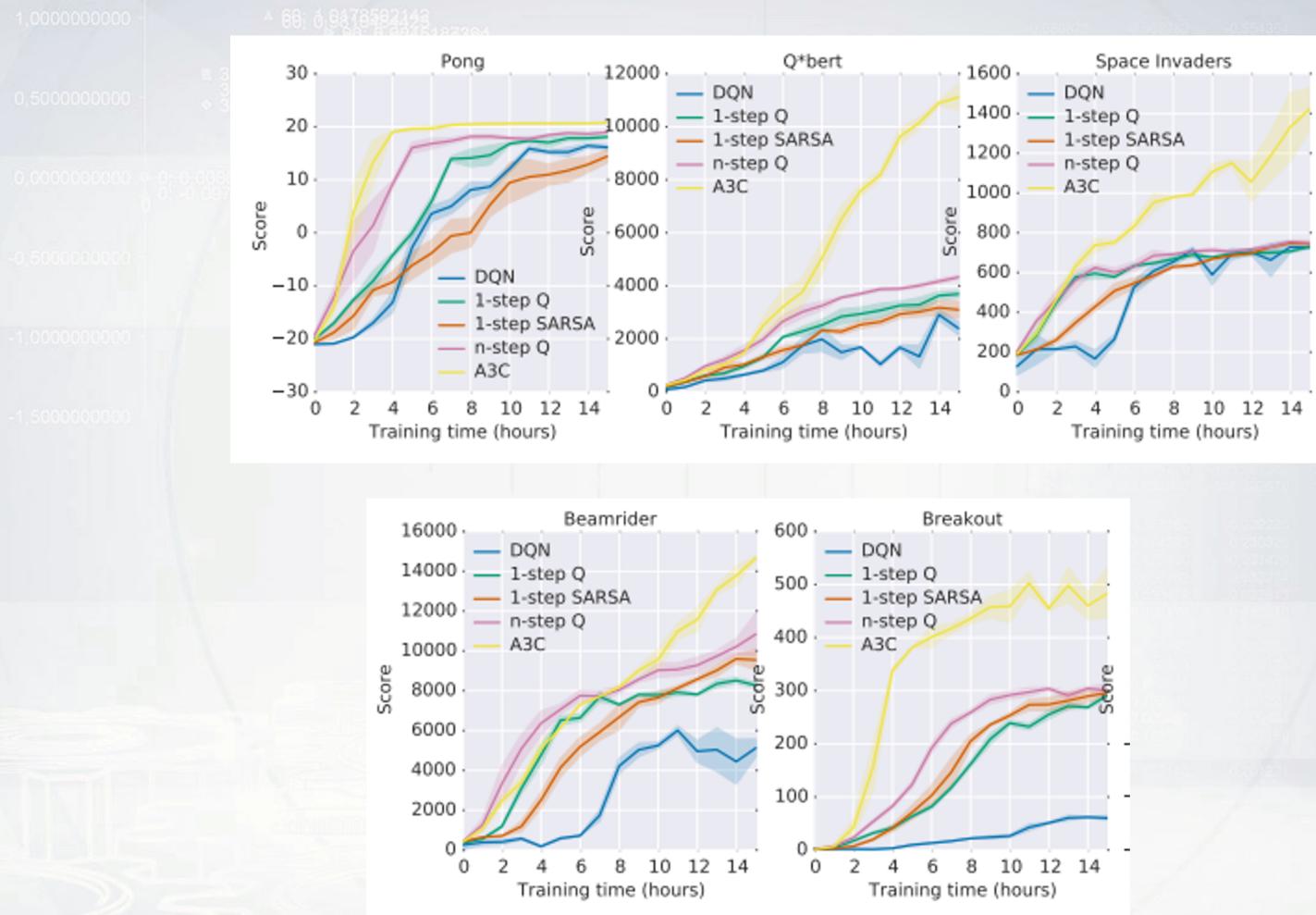
image: <http://bit.ly/2z0kBNU>

A popular implementation
of advantage actor-critic
using neural net agent

- No experience replay
- Many parallel sessions
- Asynchronous updates

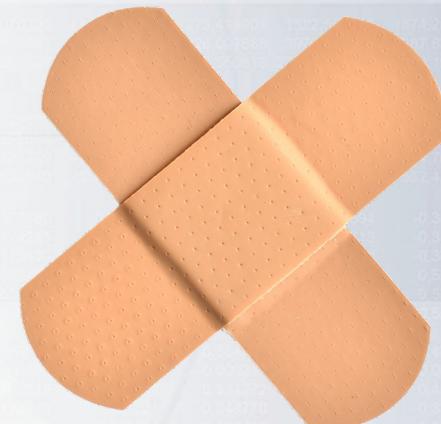


Casey Study: A3C results



Duct tape zone

- $V(s)$ errors less important than in Q-learning
 - actor still learns even if critic is random, just slower
- Regularize with entropy
 - to prevent premature convergence
- Learn on parallel sessions
 - Or super-small experience replay
- Use logsoftmax for numerical stability



value-based Vs policy-based

Value-based

Policy-based + Actor-critic

Q-learning, SARSA,
value-iteration

REINFORCE, Advantage
Actor-Critic, Crossentropy
Method

Solves harder problem

Solves easier problem

Explicit exploration

Innate exploration &
stochasticity

Evaluates states & actions

Easier for continuous
actions

Easier to train off-policy

Compatible with supervised
learning



Supervised pre-training

Reinforcement learning usually takes long to find optimal policy completely from scratch.

We can use existing knowledge to help it!

- Human experience
- Known heuristic
- Previous system



Supervised pre-training

- Supervised learning:

$$\nabla llh = \underset{x, y_{opt} \sim D}{E} \nabla \log P_\theta(y_{opt}|x)$$

- Policy gradient:

$$\nabla J = \underset{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}}{E} \nabla \log \pi(a|s) \cdot Q(s, a)$$



Supervised pre-training

- Supervised learning:

$$\nabla llh = \underset{s, a_{opt} \sim D}{E} \nabla \log \pi_\theta(a_{opt}|s)$$

- Policy gradient:

$$\nabla J = \underset{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}}{E} \nabla \log \pi(a|s) \cdot Q(s, a)$$



Supervised pre-training

- Supervised learning:

$$\nabla llh = \underset{s, a_{opt} \sim D}{E} \nabla \log \pi_\theta(a_{opt}|s)$$

- Policy gradient:

$$\nabla J = \underset{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}}{E} \nabla \log \pi(a|s) \cdot Q(s, a)$$

Q: what's different? (apart from $Q(s, a)$)



Supervised pre-training

- Supervised learning:

$$\nabla llh = \underset{s, a_{opt} \sim D}{E} \nabla \log \pi_\theta(a_{opt}|s)$$

reference

- Policy gradient:

$$\nabla J = \underset{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}}{E} \nabla \log \pi(a|s) \cdot Q(s, a)$$

generated



More out there

This domain is huge!

- Trust Region Policy optimization
 - policy gradient on steroids
- Deterministic policy gradients
 - Off-policy & less variance
- Many many more!
 - Research happens as you read this line



Let's code!

