

RL@Coursera week 5

Policy gradient methods



Yandex
Data Factory

LAMBDA 



Small experiment

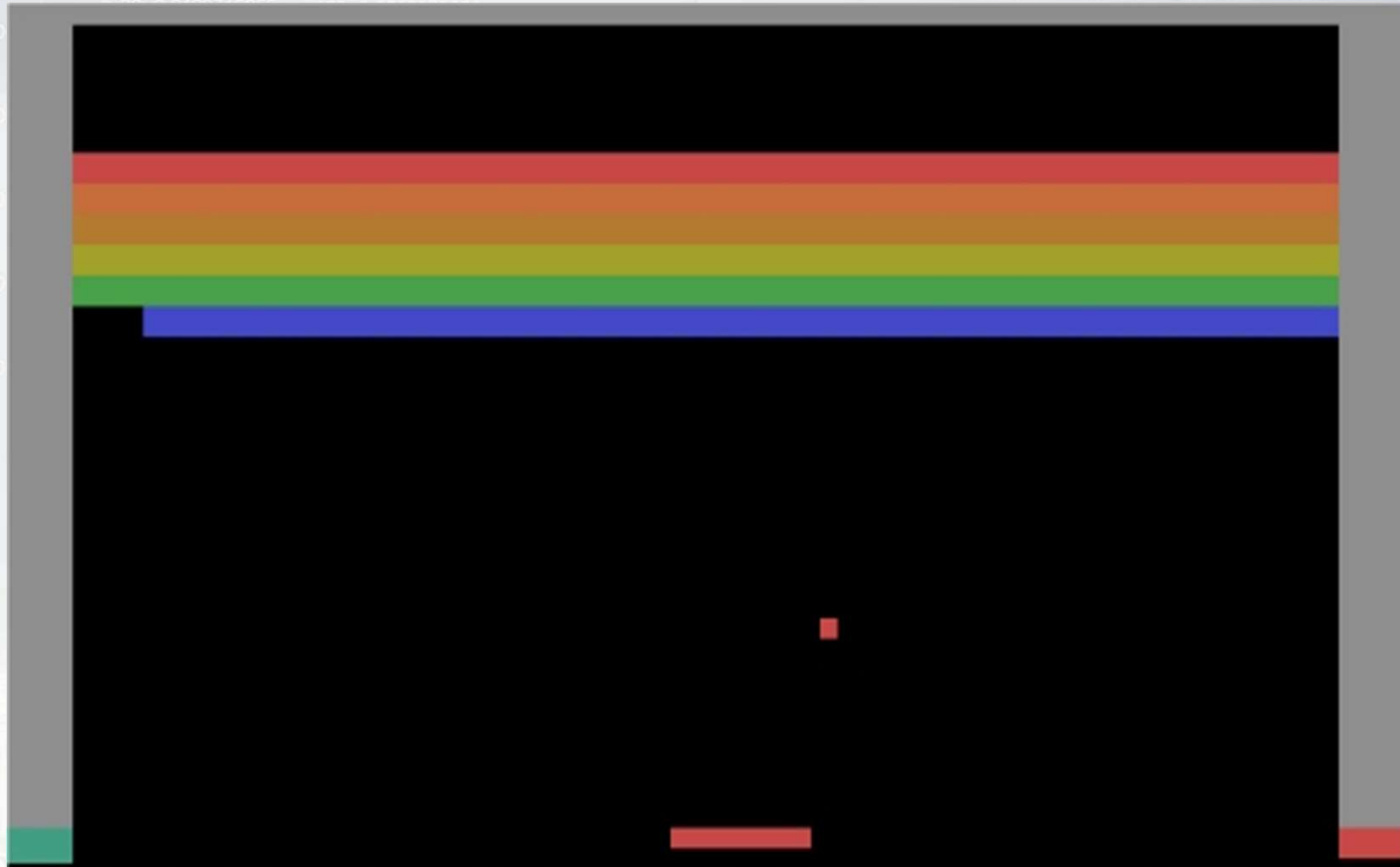
The next slide contains a question

Please respond as fast as you can!



Small experiment

▲ 68: 0.9476582143
■ 99: 0.9048488384
◆ 60: 0.7833382066

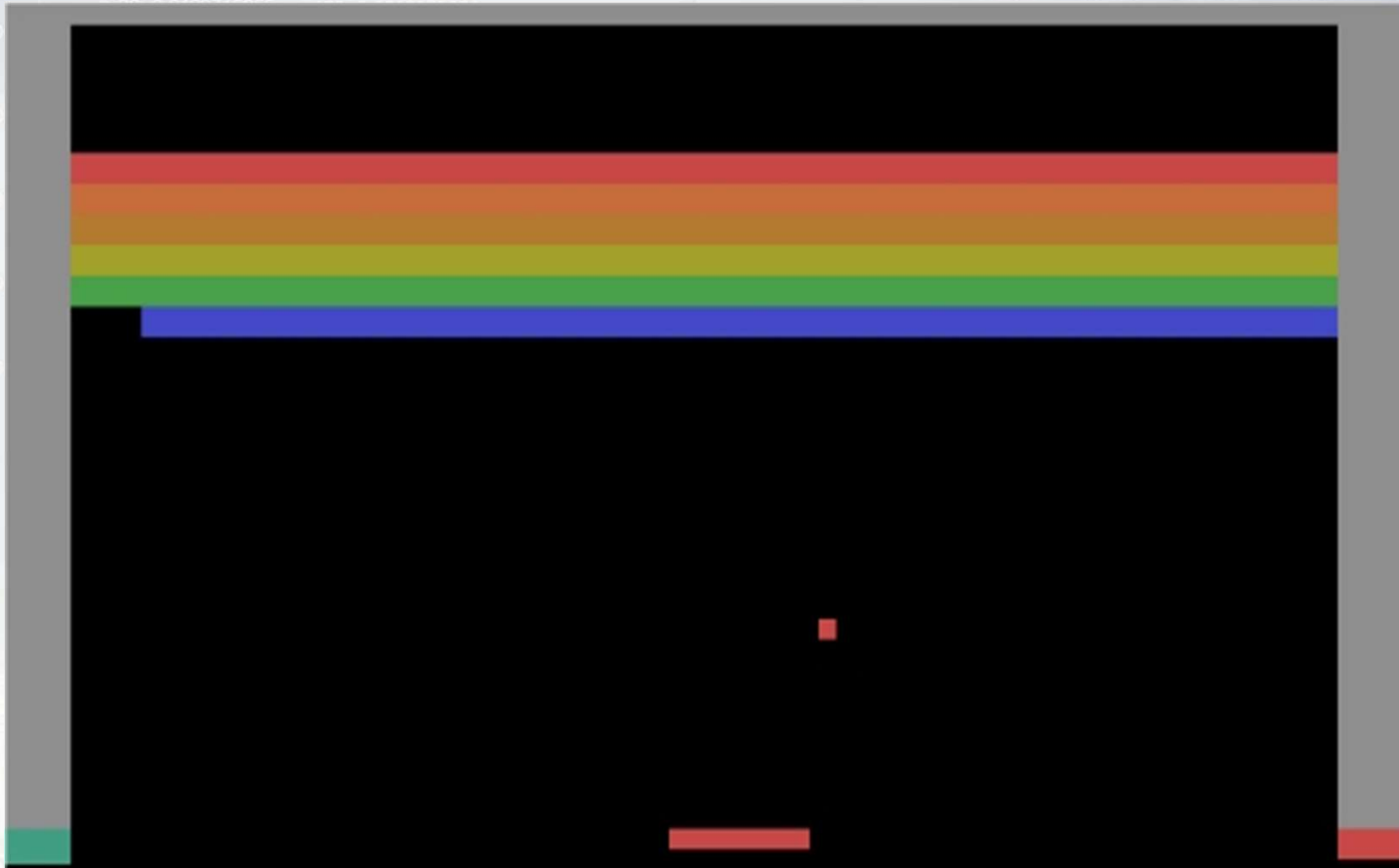


Left or Right?



Small experiment

▲ 68: 0.9476582143
■ 99: 0.9048488384
◆ 60: 0.7833382066

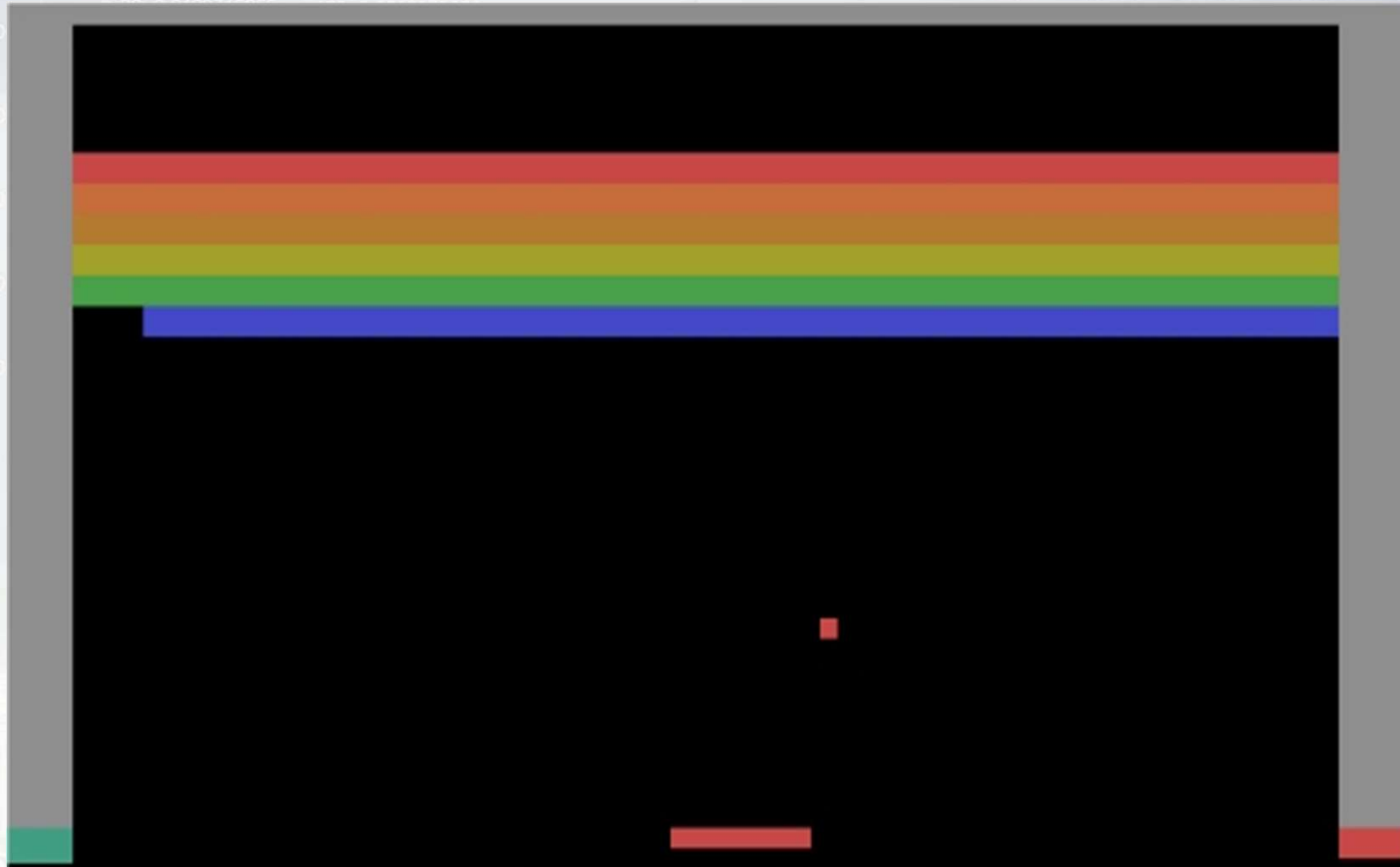


Right! Ready for next one?



Small experiment

▲ 68: 0.9476582143
■ 99: 0.9976488384
◆ 60: 0.7833382066



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



Small experiment

1.0000000000

▲ 68: 0.9478582143
■ 99: 0.9948488384
◆ 60: 0.7833382066

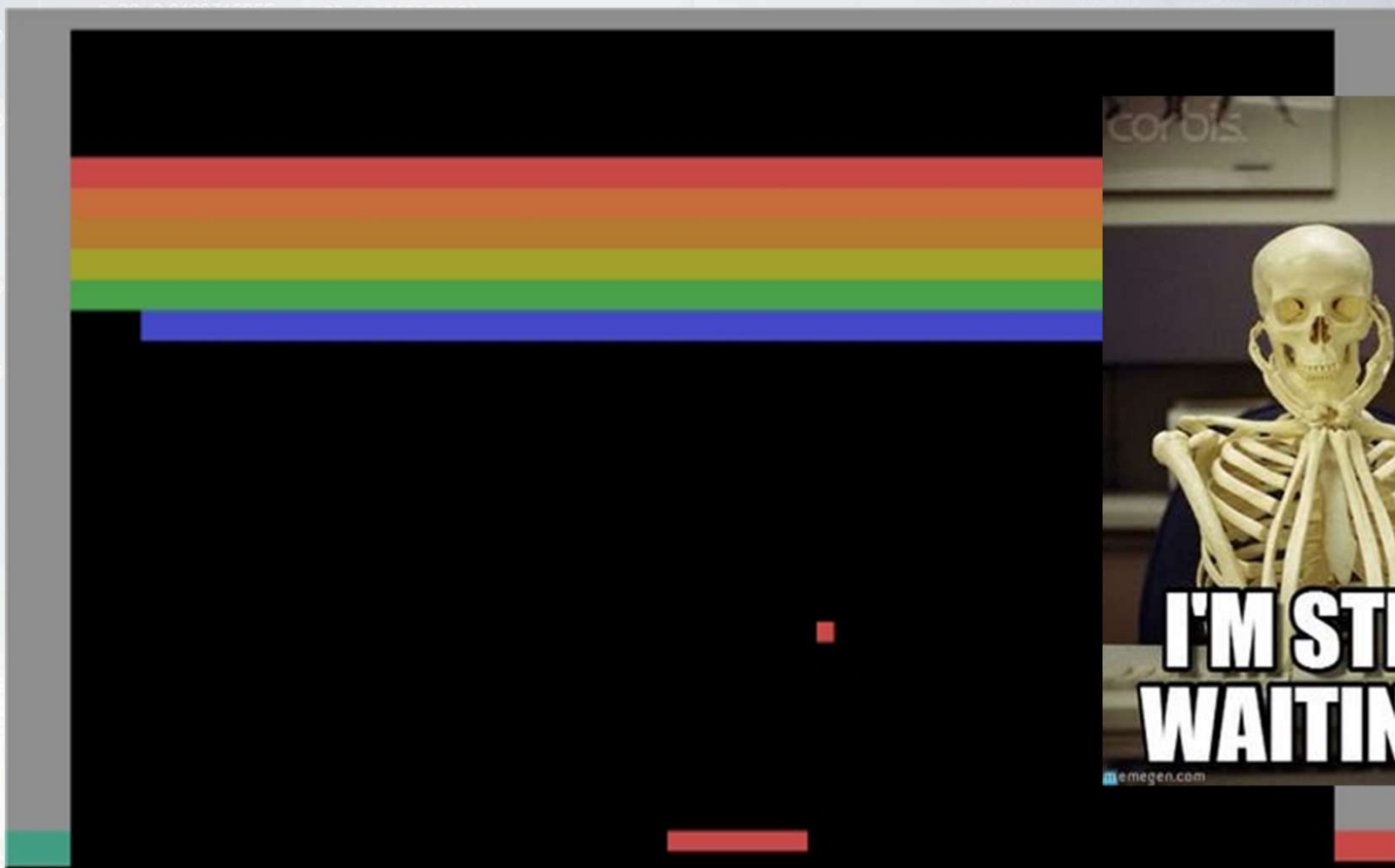
0.5000000000

0.0000000000

-0.5000000000

-1.0000000000

-1.5000000000



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



Approximation error

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

Assuming
2 nets

	Simple 2-state		
	<u>True</u> 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

const, not
dependent on
Q-fn parameters

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

while A does not capture Q-values identically, picking $\arg\max_a Q(s,a)$ will still yield a^* everywhere

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

not always representative

- 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**

- learning prob of actions
- requires lots of sampling
- very efficient in complex env



NOT how humans survived

Q-learning

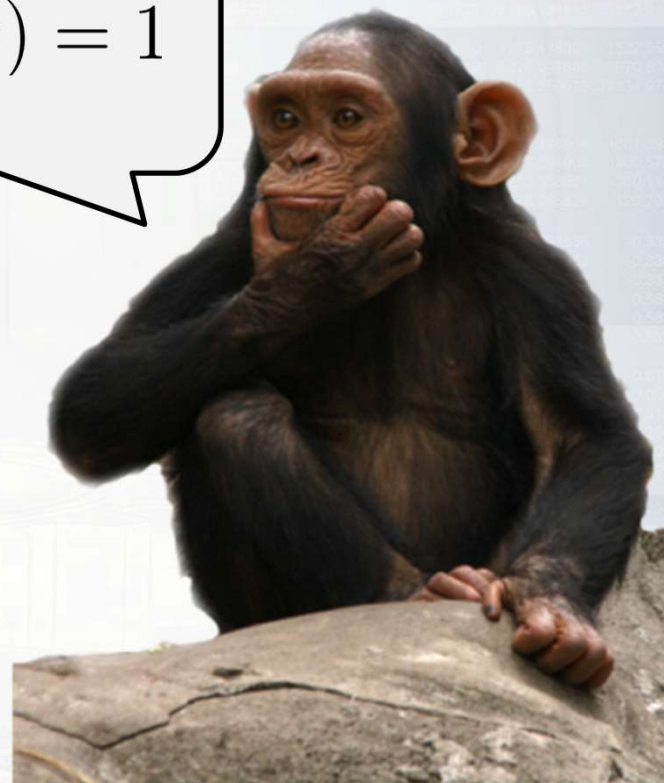
Approach

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

$$\pi(run|s) = 1$$



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

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

/
cross entropy method

* learning prob dist

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?

* Having a deterministic policy makes it easier
for an adversary to adapt



Policies

In general, two kinds

Deterministic policy

$$a = \pi_{\theta}(a|s) \leftarrow \text{same action each time}$$

Stochastic policy

$$a \sim \pi_{\theta}(a|s) \leftarrow \text{sampling takes care of exploration}$$

with Q-learning, it switch bet
random & optimal according to
some prob!



Policies

In general, two kinds

Deterministic policy

$$a = \pi_{\theta}(a|s) \leftarrow \text{same action each time}$$

Stochastic policy

$$a \sim \pi_{\theta}(a|s) \leftarrow \text{sampling takes care of exploration}$$

Q: how to represent policy in continuous action space?

→ trying some distributions



Policies

In general, two kinds

Deterministic policy

$$a = \pi_{\theta}(a|s) \leftarrow \text{same action each time}$$

Stochastic policy

$$a \sim \pi_{\theta}(a|s) \leftarrow \text{sampling takes care of exploration}$$

categorical, normal, mixture of normal, whatever



Two approaches

- **Value based:**

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

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

- **Policy based:**

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

Implicitly maximize reward over policy

Having suboptimal values yield
suboptimal policy



Recap: crossentropy method

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

- Loop:

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

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



Policy gradient main idea

Why so complicated?

We'd rather simply maximize G over π !



Objective

- Expected reward:

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

- Expected discounted reward:

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



Objective

- Expected reward:

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

$R(z)$ setting

"Undiscounted
Version"

- Expected discounted reward:

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

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



Objective

- Consider an 1-step process for simplicity

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



Objective

- Consider an 1-step process for simplicity

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



Objective

- Consider an 1-step process for simplicity

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

state visitation frequency
(may depend on policy)

Reward for
1-step session

indirectly agent will not affect,
a more complicated agent will
start to affect the state

what agent directly
influence
e.g NN prediction



Objective

- Consider an 1-step process for simplicity

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

state visitation frequency
(may depend on policy)

**Reward for
1-step session**

Q: how do we compute that?

Impractical !!



Objective

- Consider an 1-step process for simplicity

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

Reward for
agent's action

SAMPLING
using MC

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

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



Objective

- Consider an 1-step process for simplicity

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

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

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

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/dtheta

not shown when sampling



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

problems with such approaches

- 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

* We will use cross entropy with some tweaks;
to prevent it from favoring lucky outcomes



Objective

- Consider an 1-step process for simplicity

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

Wish list:

- Analytical gradient
- Easy/stable approximations



Logderivative trick

Simple math

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

(try chain rule)



Logderivative trick

Simple math

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

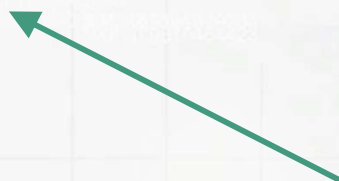
$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$



Policy gradient

Analytical inference

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

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




Policy gradient

Analytical inference

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

Doing so, we're
approx J via
sampling over s & a

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

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

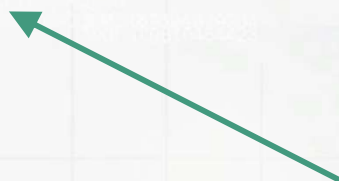
Trivia: Anything curious about that formula?



Policy gradient

Analytical inference

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

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


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

that's expectation

.)



Discounted reward case

To generalize for more than 1-step

Replace R with Q :)

True action value
aka $E(G[s,a])$

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

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

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

that's expectation
:)



Policy gradient (REINFORCE)

- Policy gradient

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

- Approximate with sampling

[MC version]

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



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

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

Ascend

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

* G obtained from discounted rewards
somehow rep Q -fn



REINFORCE algorithm

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

Q: is it on- or off-policy? — *Q-learning*
SARSA

- Loop:

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

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

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

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

Ascend

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

* We sample traj under current policy



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

Advantage



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

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

baseline, depends
on s not a

Anything that doesn't depend on action

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

* optimal policy will not change doing so



Actor-critic

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



image: youtube: mr tivikov



Advantage actor-critic

- Idea: learn both $\pi_{\theta}(a|s)$ & $V_{\theta}(s)$

~~✗~~ 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!

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

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

Also: n-step
version

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



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!

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

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

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

consider
const

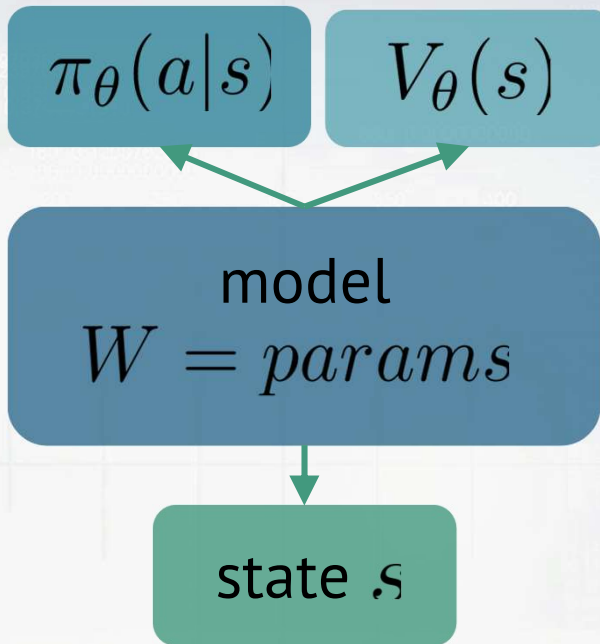
Trivia: How do we train V then?



Advantage actor-critic

Consider const

$$\min E [r + \gamma V_{\pi}(s') - V_{\pi}(s)]^2$$



picking action

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

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

est how good
particular state

→ improving convergence as well



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

Value-based

Q-learning, SARSA, value-iteration

Policy-based + Actor-critic

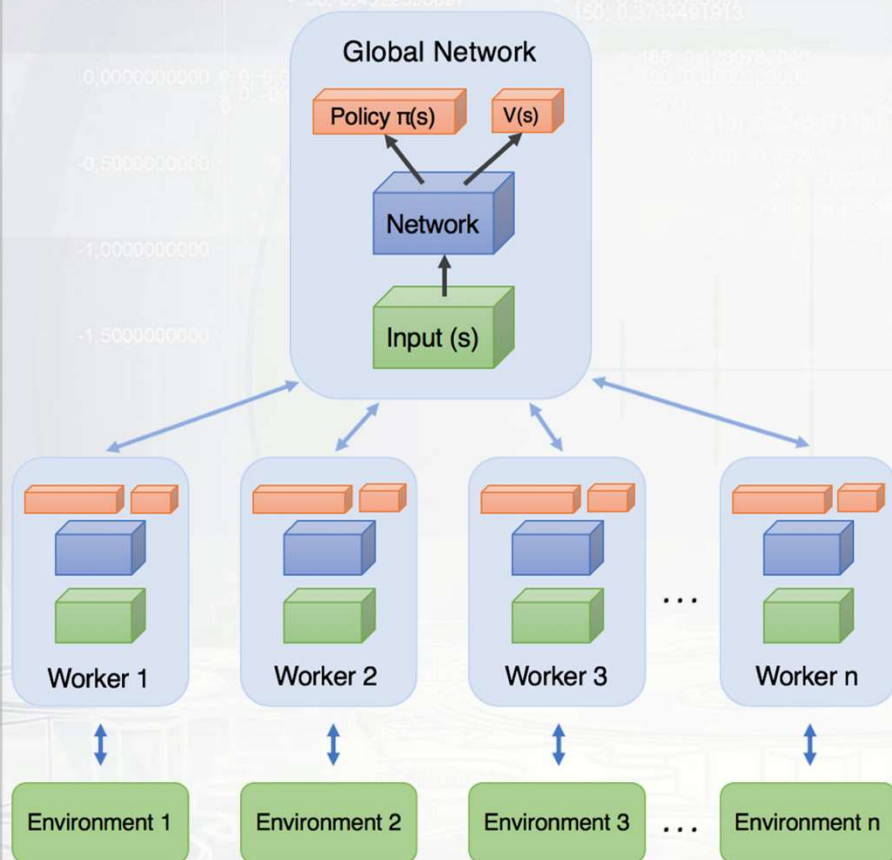
REINFORCE, Advantage Actor-Critic, Crossentropy Method

Your guess?



Casy Study: A3C

Async Actor Critic



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

- No experience replay
- Many parallel sessions
- Asynchronous updates

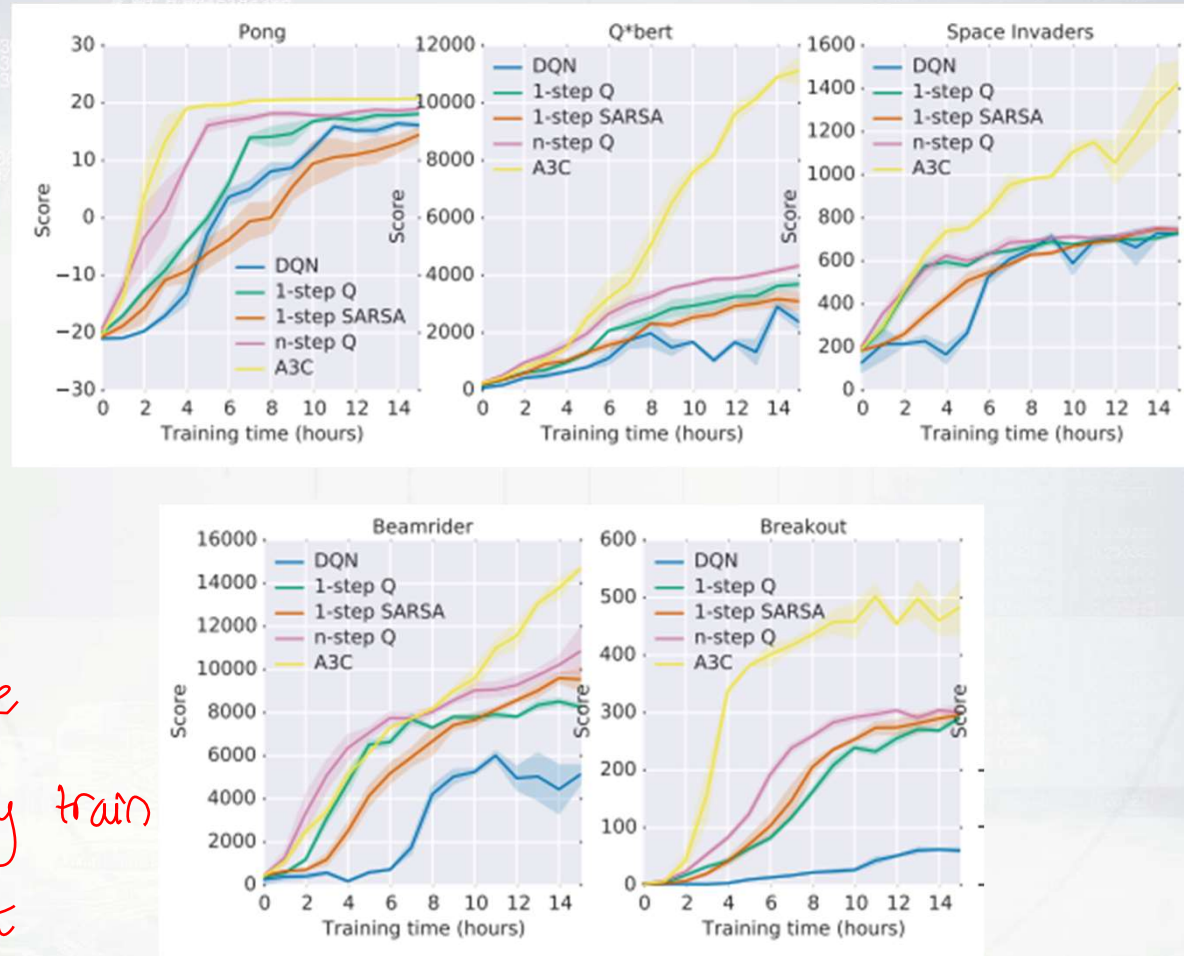
image: <http://bit.ly/2z0kBNu>



Casy Study: A3C results

A3C + LSTM

at least several seq
transitions $[s_1, a_1, r_1, s'_1, a'_1, \dots, s_k]$ will be
required to effectively train
RNN inside the agent



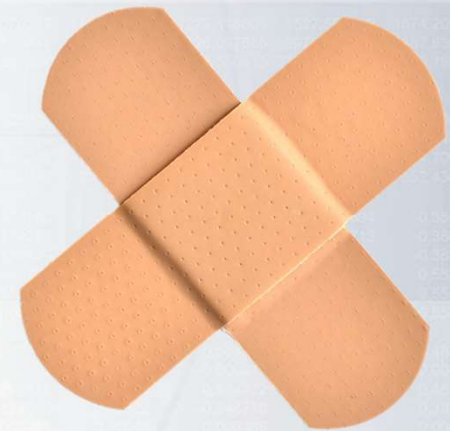
Duct tape zone

LOSS \nearrow policy-based ; Policy grad
 \searrow TD, value-based (less important)

* 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

learning
value

learning policy

Value-based

Q-learning, SARSA,
value-iteration

Solves harder problem

* Explicit exploration

Evaluates states & actions

Easier to train off-policy ;

increasing sample efficiency

Policy-based + Actor-critic

REINFORCE, Advantage Actor-Critic, Crossentropy Method

Solves easier problem

Innate exploration & stochasticity

* Easier for continuous actions

* Compatible with supervised learning

on-policy



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:

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

- Policy gradient:

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



Supervised pre-training

- Supervised learning:

initializing from some prior

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

- Policy gradient:

improving

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



Supervised pre-training

- Supervised learning:

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

- Policy gradient:

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

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



Supervised pre-training

- Supervised learning:

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

reference

- Policy gradient:

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

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

narrowing the policy space



Let's code!

