

Origins of Actor-Critic

Policy Gradient

$$\leftarrow E_{\pi \sim \pi_0(\pi)} [\nabla \log \pi_0(\pi) \lambda(\pi)]$$

$$= E_{\pi \sim \pi_0(\pi)} \left[\sum_{t=1}^T \nabla \log \pi_0(a_t | s_t) \lambda(\pi) \right]$$

$$= \sum_{t=1}^T E_{\substack{a, s \\ \sim \pi_0}} \nabla \log \pi_0(a_t | s_t) \lambda(\pi)$$

Now, if this was a CONSTANT,

Then $\nabla \log \pi_0(a_t | s_t)$ can be taken

common,

across all the different trajectories!

$$\sum_{t=1}^T \underbrace{\sum_{a, s} \pi_0(a, s)}_{P_0} \nabla \log \pi_0(a | s) (K)$$

And this would be the
marginal distribution
probability!

But, Now, Because " K " is not same for
all the trajectories, the grad. can't be
GROUPED, taken common to get the

marginal probability in front!

Instead,

$$\text{say } \pi_0(a/s) = \frac{l}{N}$$

Total number

of trajectories,

Subclass
of trajectories

That all
through
different
sources

landed on the
same
(a/s)

$\pi_0(a/s)$ will be distributed,
the following
way

$$\sum_{t=1}^T \sum_{\text{all } a/s} \nabla \log \pi_0(a/s) [l_1 \tilde{Q}_1 + l_2 \tilde{Q}_2 + \dots + l_n \tilde{Q}_n]$$

Since we are estimating the policy gradient through Monte-Carlo!

$$l = i_1 + i_2 + \dots + i_n$$

$$\delta_g = \sum_{t=1}^T \sum_{\text{all } a/s} \nabla \log \pi_0(a/s) \left[\frac{l_1 \tilde{Q}_1 + l_2 \tilde{Q}_2 + \dots + l_n \tilde{Q}_n}{N} \right]$$

is the policy gradient,
where. $\tilde{Q}_1, \tilde{Q}_2, \dots$ are the different sources via
which it's carried in, starting from a_t, p_t &
 i_1, i_2, \dots are the frequencies for each source!

Now,

Instead of dividing by N ,

If you divide by " l ",

Then it becomes (an) average of \cancel{N} values

$$\text{So } \sum_{t=1}^T \text{ all } a_t s_t$$

$$\nabla \log \pi_\theta(a_t | s_t) \left[\frac{a_1 s_1 + a_2 s_2 + \dots + a_T s_T}{l} \right]$$

This is averaging of \cancel{N} estimates

$$\nabla \sum_{t=1}^T \nabla \log \pi_\theta(a_t | s_t) \left[Q(a_t, s_t) \right]$$

Actor-Critic ($\hat{\pi}$) born!

So, doing this, changes the gradient, it introduces a "bias"!

Also,

$$l < N$$

so, the gradients are now being divided by a smaller no. so, the gradients will be

BIGGER!

But, In actor-critic,

There is **NO** expectation
! there does not interact with anymore!

Since Actor-Critic is the RESULT
of an expectation!

"biased"

So in Actor-Critic, there is ONLY
a sum!

So if we know the Q-values perfectly
somehow!

Then, we don't need (loss) of sampling where
 $N \rightarrow \infty$

We just sample enough to get a
wide variety of states & actions, to sum
over gradients to θ
from all those places
& update our
network!

A minibatch

Smaller the batch,

Higher the threshold &
oscillations!

So, ① Expectation in Policy Gradients gives
rise to some sort of an
averaging over Q estimates

We replace it by Q value taking
in a little bit (bias) just like
"causality"

Since the expectation is consumed in
Q-values, there is NO expectation in
actor, ONLY summation, just generate a nice
variety, a minibatch

Now, The CRITIC

Only for the 1st step,

we generate LOT of samples ($N \rightarrow \infty$) like policy gradients & train Hill convergence, To make the critic ~~learn~~ Q-function!

From next step on, we just generate a mini-batch, (for variety), for updating the actor! ~~and use every step~~!

And since the policy is just changing a little bit @ each gradient step, we use the same minibatch to take one gradient step in "our" critic also, to update it for the next policy! Since the Q-function is also not changing @ every step!

So, @ every step

- Sample a minibatch
- Update Critic to reflect current ~~next~~ policy
- Update current policy actor using critic

~~to~~ next policy

Training Critic

Monte-Carlo estimates

We can use the estimates we use in Policy Gradients as our training targets for our critic!

$$V(s) = \sum_{t=t}^T r_t,$$

This is

Now since we have many samples, so many estimates, a single sample estimate of the actual expectation for V .

Since all of them are targets together for our network, you might think gradient descent, would settle at the average

2 these targets to minimize loss which is nothing BUT the expectation of the estimates!

But this anyway Policy Gradients also does for us! There is nothing new!

So, what is the advantage of introducing an entire network to do the same averaging as that of PG?

Critic learns the V function and NOT the Q function.

This is because V fn is a function of the state $[V(s)]$ while Q fn is a function of both the state and the action $[Q(s, a)]$

So, if action is several dimensions, then the network grows! & there are a lot of parameters to train

So, $V(s)$ is learnt by critic to reduce network size!

And Q is approximated using V !

$$Q(s_t, a) = \mathbb{E}_{s_{t+1}} V(s_{t+1})$$

we instead take a single $V(s_{t+1})$ which is a single sample estimate of the expectation & approximate Q by it

$$Q(s_t, a) = V(s_{t+1})$$

found in the same

The point of introducing a network for the same averaging is,

Practically when you sample, Only the most probable states, will repeat a lot! Let's call these Type A states

And type B states get sampled less,

\hat{s}^b

So, in PCA, expectation of \hat{s}^a is accurate, while for \hat{s}^b , it gets overfit!

So of a network you get generalization!

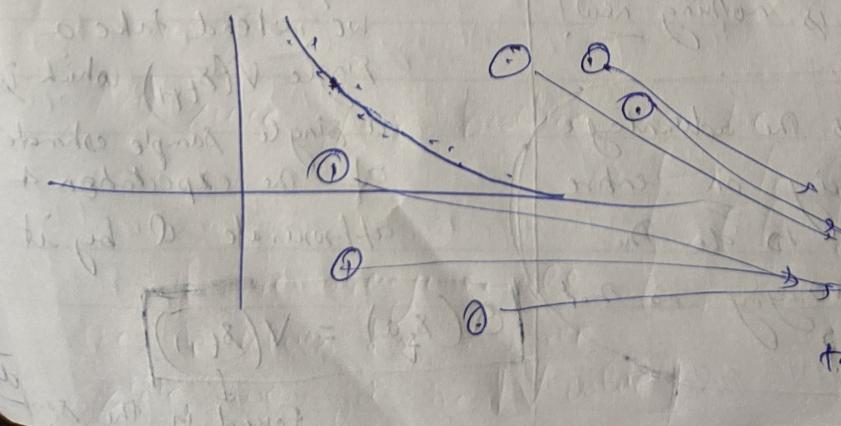
all the noise are fitted!
in fact the noise is taken as the exact value!

overfitting to noise would change the network so badly that it will error a lot for \hat{s}^a the majority, so loss will be more!

So, instead, noise is ignored,

\hat{s}^a are properly fitted, & \hat{s}^b values follow from the

[GENERALIZATION]



All the noise are IGNORED, Instead fitted values are taken from the generalization!

Bootstrapped Estimates

We can do TD-learning, Temporal difference learning!

So,

$$\hat{V}(s_t) = \mathbb{E} V(s_{t+1})$$

$$V(s_t) = \mathbb{E}_{a_t} Q(s_t, a_t) \quad \text{take a sample}$$

$$\approx Q(s_t, a_t) \quad \text{from the sampled trajectory}$$

$$\approx \delta_t + \mathbb{E}_{s_{t+1}} V(s_{t+1}) \quad \text{take a sample}$$

$$\approx \delta_t + V(s_{t+1}) \quad \text{from the sampled trajectory!}$$

So, take $\boxed{\delta_t + V(s_{t+1})}$ as the target for $V(s_t)$

Funny

thing is \downarrow This is from the network's output!

So, this is the V value

for the "previous" policy!

Since in one gradient step, the policies don't change that much!

we can approximate previous policy's V !

V as current policy's V !

So consider bootstrapped estimates,
if you observe the critics training!

It is

$$V(\beta_t) \leftarrow \hat{r}_t + V(\beta_{t+1})$$

~~So each time the network is expected to predict a value~~

~~Suppose all rewards are positive in a given problem, no negative rewards!~~

Then, each time the network is expected to predict a value bigger than its previous prediction!

→ This would result in explosion, where soon the V function would be outputting ∞

So, we need to introduce DISCOUNTING

$$V(\beta_t) \leftarrow \hat{r}_t + \gamma V(\beta_{t+1})$$

in the rewards of the problem

Now, Baseline

if your eqn uses Q values,

$$\sum_{t=1}^T \sum_{a \in \mathcal{A}} \frac{\log \pi(a|s)}{n \pi_0} Q(a, s)$$

And since your baseline can be

fn of time
fn of state
constant

(paved in advanced & basic baselines)

IF we average all the Q values for a given state, we get the V value

$$\delta_0, \sum_{t=1}^T \sum_{a_t, s_t} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left[Q(s_t, a_t) - V(s_t) \right]$$

And this is called the
ADVANTAGE
And it's evaluated
of approximation of

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

This is next state
in the sample

Now, You have options,
you can use
the critic to evaluate Advantage

OR in "Policy Gradients", you can use
critic to evaluate (only) the baseline
 $V(s)$

OR You can take the n^{th} advantage!

$$V(s_t) = \delta_t + \gamma \delta_{t+1} + \dots + \gamma^{n-1} \delta_{t+n} + \gamma^n V(s_{t+n})$$

where for n steps, you can use
values from the sample & use @ the
 n^{th} step decide to use the critic

OR Generalized Advantage Estimator, where you
take all n -step advantages ($n=1$ to T) & discount
(λ) them!