# Actor-Critic to Policy iteration

$$\sum_{\substack{\text{minibatch} \\ \text{f }(s,a)}} \nabla_\theta \log \pi_\theta(a/s) \, A^\pi(a,s)$$

Let's consider 'one' particular state & see what Actor-critic __does__ !
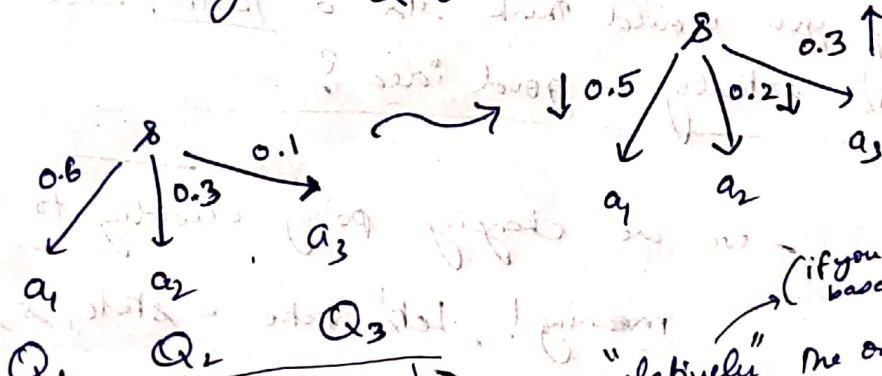
Say, 3 actions possible &

$$P_i = \pi_\theta(a_i/s)$$



Since $A^\pi(a_i, s) = Q^\pi(a_i, s) \boxed{- V^\pi(s)}$

If we remove the baseline, (since it's (same) here as we consider only one state)

Then, $P_i$ is increased (or) decreased according to $Q^\pi(a_i, s)$
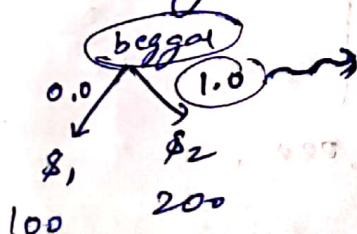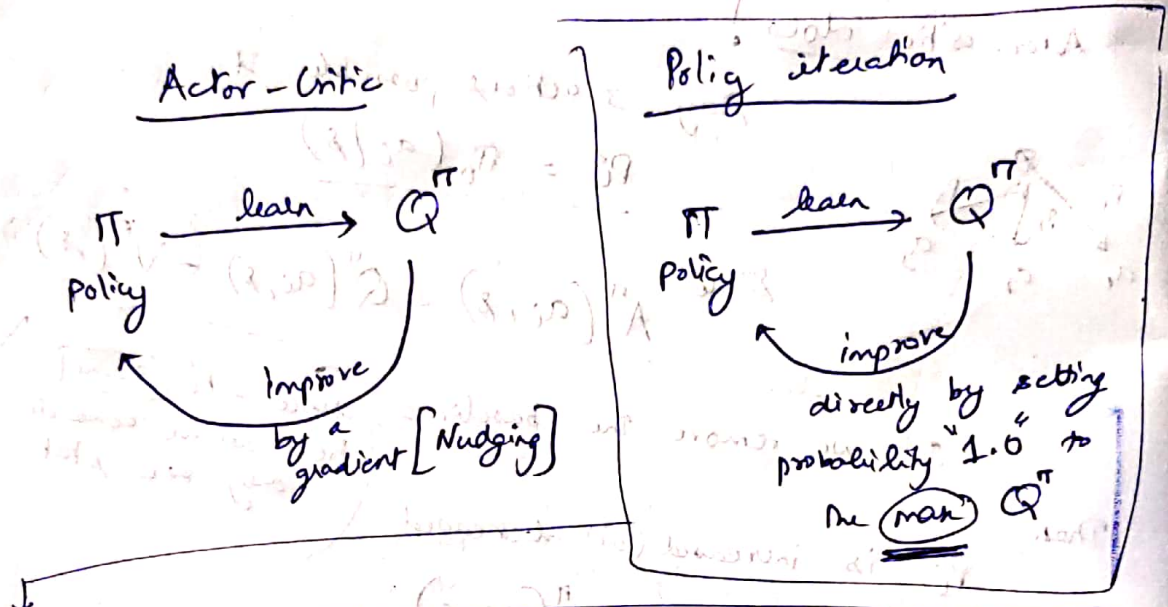
That is,



Say $\boxed{100 \quad , \quad 200 \quad \boxed{300}} \longrightarrow$ then "relatively", the one $P_1$ the highest $Q$'s probability is increased, & lowest $Q$'s probability decreased

$Q_1 \qquad Q_2 \qquad Q_3$

(if you include baseline (or) otherwise)

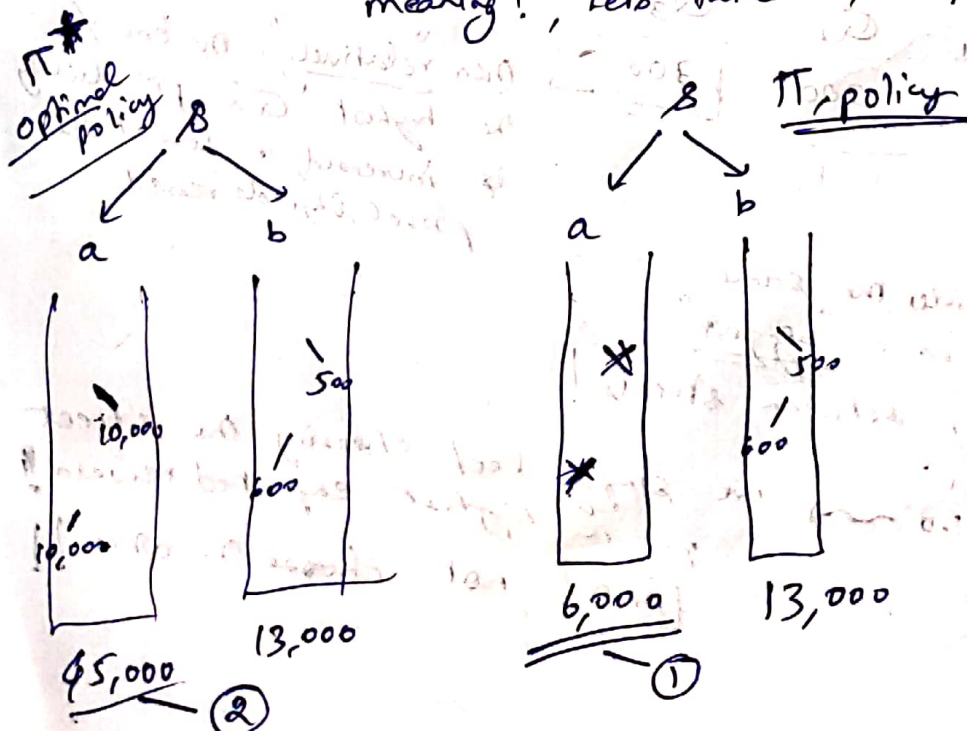If you consider the same scenario as a "~~beggar~~ beggas" choosing between streets !



we simply keep choosing the street of the highest expected reward! $\boxed{AND}$ not choose the other !!

So, if that's the case, then its instead of
[nudging] the probabilities (only) a little bit,
we can completely change it [like] in policy
iteration! which is the (only) difference between
the algorithms!

| Actor - Critic | Policy iteration |
|---|---|

$\Pi \xrightarrow{\text{learn}} Q^{\pi}$
Policy

Improve
by a
gradient [Nudging]

$\Pi \xrightarrow{\text{learn}} Q^{\pi}$
Policy

improve
directly by setting
probability "1.0" to
the (max) $Q^{\pi}$

* why you would think its a [BAD] idea when
its actually a good idea?

- we are changing policy [according to $Q^{\pi}$, not $Q^*$]
meaning!, lets take a state, s, & 2 actions

$\Pi^*$
optimal
policy

s

a          b

$\Pi$, policy

s

a          b

a: 10,000
10,000
b: 500
600
13,000
$5,000      ②

a: ✗
✗
6,000
b: 500
600
13,000
①

So,
under optimal policy all, actions under all states "are"
optimal! So, $Q^*(a,s) = 45,000$

But
since $\pi$, is
some policy,

So, $Q^\pi(a,s) = 6,000$

it has been
configured if
actions sub-optimal
under `a` [but]
the same as optimal
policy under `b`

So, Because $\pi$ is a [BAD]
policy, it's Q-value
estimate from `a` looks bad!

As, a result, update to $\pi$ [only in one state]
under policy iteration
changed the policy to $\pi'$

⊛

$\pi'$

$a$ $s$ $\to b$  which [prefers] `b` instead of a!

So, $\pi^*$   $\pi$    $\pi'$
best    bad    worse
policy  policy  policy

Since it's the
same $\pi$, if the
policy changed for only
one state, where
only the sub-optimal
was selected!

So, if you "rate" the policies
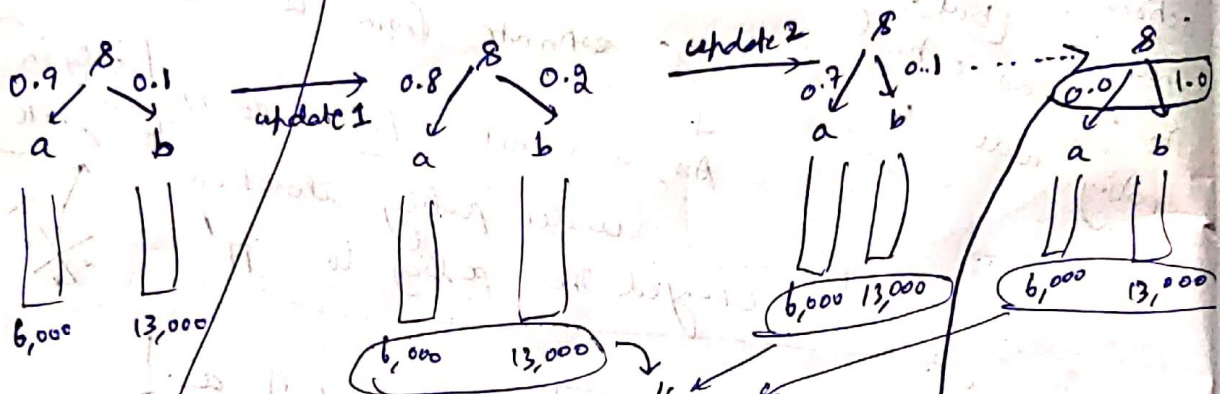by the [number] of sub-optimal
choices the policy makes,

Then $\pi'$ does not look like an improvement to $\pi$
but "like" a defmovement!

Disclaimer! All this reasoning is assuming, we
are applying policy iteration update
to [only one] state, (not) everyone

But policy iteration is [proved] to improve policy. (So, what is happening?)

Proof (via)
Actor-Critic

Say, you apply actor-critic update (also) to [only one] state, then,



0.9 / 8 \ 0.1
a      b
6,000   13,000

—update 1→

0.8 / 8 \ 0.2
a      b
6,000   13,000

—update 2→

0.7 / 8 \ 0.1
a      b
6,000  13,000

· · · · ·

0.0 / 8 \ 1.0
a      b
6,000   13,000

remains the "same",
since in the new policy,
only (8) is changed & others
are all same!

[ Assuming, there are (no) loops
& V(s) is not used in
the computation of V(s) ]

(but in reality
loops exist)

So, Actor-Critic @ [only one] state "also" produces same result

So, how is our intuition that it (won't) improve is wrong?

This is because we are evaluating the
policies π*, π, (π') on the metrics of [number] of
optimal decisions?
in which case, π' ranks [LAST] !

But in terms of expected reward, as our metric,
π' is (not) LAST !

This is because, not all suboptimal choices

are EQUAL

meaning -

(not) all mistakes [Cost] EQUAL

You can make a 100 small mistakes ← still receive 1000 points

(but) by making only a few mistakes say 5, you can receive 100 points

Example → death ↳ 1 mistake, but costliest !

↳ (but) skipping lunch oneday → small cost !!

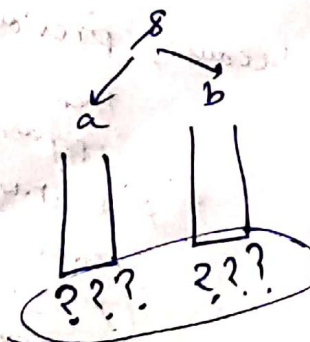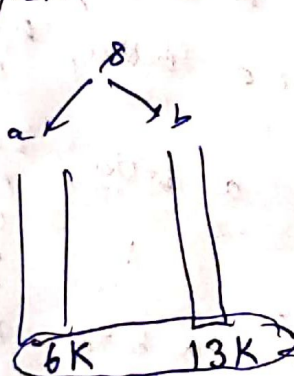1 mistake but

So, in metrics of expected reward

$\Pi^*$    $\Pi'$    $\Pi''$
best policy   (improved) policy   [BAD] policy

} This is the same !

** why you would think it's a bad idea when it's actually a good idea ?

↳ Say, you are doing policy iteration (on) [all states !]
Then → since the policy is changed everywhere

a ↙ $s$ ↘ b

a ↙ $s$ ↘ b

??? ???

6K    13K

we don't know what the estimate will be in this (new) policy
But we update(s) using the invalid old policy (!)

Not only are we using in correct values to update ⑧!

↳ There will also be [OSCILLATION]

↳ Since we have proved that the policy improves,

So,



decisions here are improved

→ decisions here are improved

so, new returns

Obviously this wud improve, in fact it will improve, till

45K !

6K    13K

19K    13K

So, now a seems better,

previously b seemed better,

So, the policy keeps (on) oscillating !

But point is it oscillates only once!

↳ But point is it oscillates, only a [Fixed] number of times!

— The oscillations are bounded and [NOT] infinite!
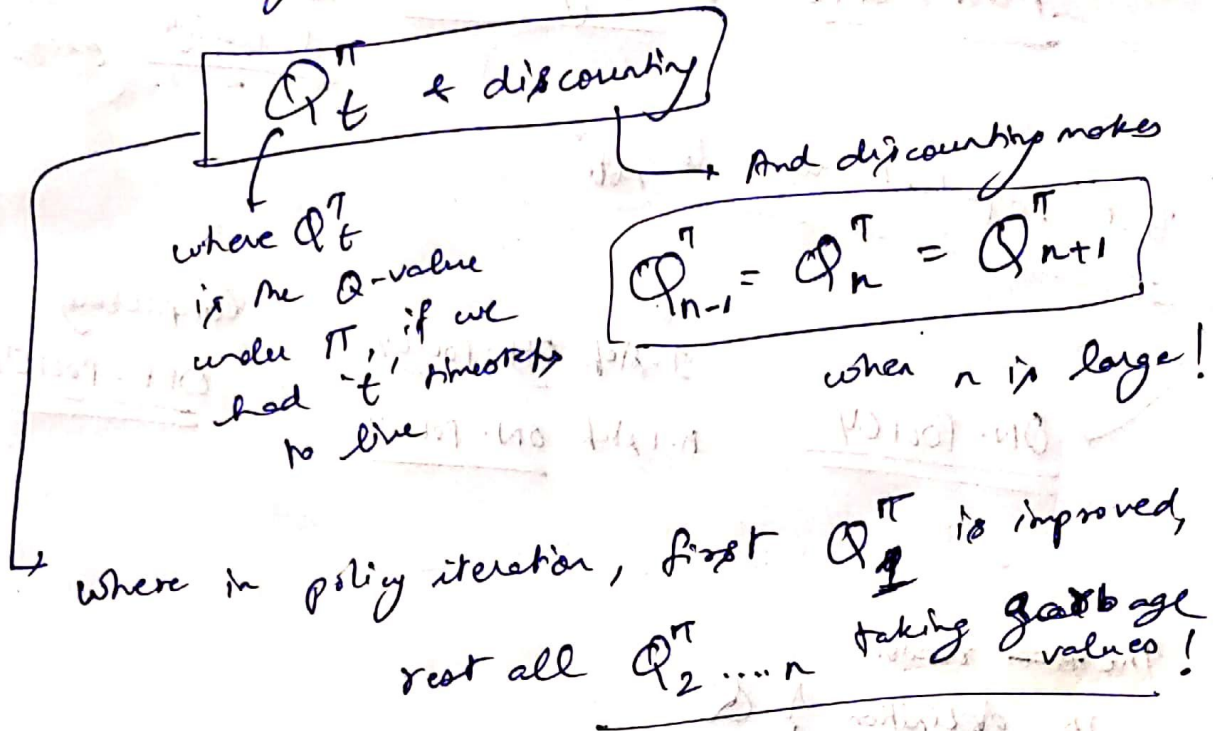
But why does the oscillation [STOP] ??

↳ This is because, previously   a worse, b better

then improvement   a better, b worse

And it will remain the   a better, b worse

(same) here onwards   a better b worse

<u>In other words,</u>

The way, we calculate $\underline{Q^{\pi}}$ is using

$$\boxed{Q_t^{\pi} \text{ & discounting}}$$

→ And discounting makes

where $Q_t^{\pi}$ is the Q-value under $\pi$, if we had "t" timesteps to live

$$\boxed{Q_{n-1}^{\pi} = Q_n^{\pi} = Q_{n+1}^{\pi}}$$

when $n$ is large!

↳ where in policy iteration, first $Q_1^{\pi}$ is improved, rest all $Q_{2\ldots n}^{\pi}$ <u>taking garbage values</u>!

then

if $Q_1^{\pi}$ settled $Q_2^{\pi}$ settles

$\underline{Q_{3 \ldots n}^{\pi}}$ oscillates!

So, slowly

$Q_{1 \ldots k}^{\pi}$ are settled & $Q_{k+1 \ldots n}^{\pi}$ oscillate

↳ until $\boxed{Q_n^{\pi} \text{ settles}}$

which you could do <u>directly</u>

Evaluate <u>only</u>

$Q_0^{\pi}$, then $Q_1^{\pi}$ then $Q_2^{\pi} \ldots$ until you get

$\underline{Q_n^{\pi}}$

↳ which is called as <u>Q-learning</u>!

Here the transition from

Actor - Critic  to  Policy Iteration  to  Q-learing (ends)

But important thing to (note)

→ ON - POLICY

1-step OFF - POLICY,
n-1 step ON - POLICY

Completely
OFF - POLICY

The reason
is, definition of $\underline{\underline{Q}}$



A policy can be anything, (even)
'0' probability for $\underline{a}$,
But to calculate

$\underline{Q(a,s)}$, we still
need to sample 'a' action,
1,000 times to find
expected return !

In Actor-Critic,
there is always,
some probability
for all actions !
So, all $Q(a_i, s_j)$
are calculated !

In policy iteration,
it is [deterministic]
policy !

So, for 'one' step, i.e for only the
1st step $(a,s)$, we need to sample
off-policy 1,000 times, but
Then on , we need to only sample
according to policy to calculate
$Q^{\pi}(a,s)$

## Last of all

$Q^*(a,s)$ is completely off-policy

## To summarize,

Actor-critic's policy has off-policy built in, by having "some" probability for all actions in the policy, so, (purely) on-policy!

Policy iteration evaluates $Q^\pi(a,s)$, where,

a is off-policy
→ Rest all, on-policy!

$Q^*(a,s)$ is strictly off-policy!

off-policy → ↓

$$Q^*(a,s) = \mathop{E}_{s'} \left( \max_{a'} \right) Q^*(a',s')$$

max takes care of using the optimal policy [always]!

---

DONE