

Target: figure out a set of strategies to use for the entropy case.

Differences between the usual RL paradigm and us:

- we are looking for entropy, not reward

Usually, high-reward states stay high-reward. For us, high-entropy states will hopefully turn into low-entropy states as we sample more.

Normally, as $t \rightarrow \infty$, $\hat{p}_n \rightarrow p^*$: p^* mean of samples

For us, as $t \rightarrow \infty$, $h(p_{\text{actual}}) \rightarrow 0$.

CUT gives us:

$$\text{Let } \hat{X}_n = \{1, n\}^3 \text{ be a sequence of iid RVs.}$$

$$S_n = \sum_i X_i, \forall n \geq 1 \text{ with } E[S_n] = \frac{n}{3} E[X_1] \text{ and } \text{Var}[S_n] = \frac{n}{3} \text{Var}[X_1].$$

$$S_n = \frac{S_n - np}{\sqrt{np}} \left(n - \frac{np}{\sqrt{np}} \right)^2 \text{ for } p: E[X_1] \text{ and } \sigma^2 = \text{Var}[X_1]$$

If $S_n \sim F(x)$, then $F_n(t) \rightarrow F(t)$ as $n \rightarrow \infty$.

In other words,

$$\hat{p} \sim G_n(t) \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right)$$

The more samples we get, the distribution of \hat{p} becomes normal & narrower $\sim \frac{1}{\sqrt{n}}$.

So our estimate of $Q(s,a)$ gets closer to the true value as $\frac{1}{\sqrt{n}}$, which means the probability approaches what?

$$P(Q(s,a) > 0) = 1 - \Phi\left(\frac{\mu}{\sqrt{n}}\right) \text{ where } \Phi \text{ is the cdf of the standard normal.}$$

$$= 1 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{\mu}{\sqrt{n}}} e^{-\frac{t^2}{2}} dt = 1 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{\mu}{\sqrt{n}}} e^{-\frac{t^2}{2}} dt$$

As $n \rightarrow \infty$, the term $\frac{\mu}{\sqrt{n}}$ becomes $\pm\infty$ as $n \rightarrow \infty$ and $\hat{p} \rightarrow p^*$ depending on the sign of p^* . When $p^* > 0$, $\lim P(Q(s,a) > 0) \rightarrow 1$ and when $p^* < 0$, $\lim P(Q(s,a) > 0) \rightarrow 0$.

This is all kind of obvious. What am I trying to do? For true $Q(s,a) \neq 0$,

$$\lim_{n \rightarrow \infty} h(Q(s,a)) = \lim_{n \rightarrow \infty} -p(Q(s,a) > 0) \log(p(Q(s,a) > 0)) - (1-p(Q(s,a) > 0)) \log(1-p(Q(s,a) > 0))$$

$$= 0 \text{ (using } \lim_{n \rightarrow \infty} x \log x = 0)$$

Usually RL wants estimates of value $Q(s,a)$ for state s, a to $g(s)$; true value a is p^* .

We want $h(Q(s,a)) \rightarrow 0$ as well.

$$g(s) \rightarrow g(s) \text{ as } \frac{1}{\sqrt{n}}$$

Can I get a similar expression for h ?

$$h(Q(s,a)) \rightarrow 0 \text{ as } ?$$

David Silver exploration/exploitation lecture (9 in RL course)

- 1 intro
- 2 multiarmed bandits
- 3 confidence bounds
- 4 MCTS

Principles:

Naive exploitation

Add noise to greedy policy, $(E[\text{greedy}])$

Optimistic init

Optimism in the face of uncertainty

Naive estimator of uncertainty on value

Information-theoretic search

Consider agent's info as part of state

Look where to seek info helps reward

Correct but computationally difficult

Probability matching

Select action according to probability, they are best

Multiarmed bandit case.

$\langle t, p \rangle > 0$ no more.

Regret definition:

$$g(a) = E[R|A=a]$$

$$V_a = g(a) + \frac{\partial g}{\partial a}$$

$$\text{regret } L = E[V^* - g(A)]$$

$$\text{total regret } L_t = E[\sum V^* - g(A_t)]$$

Maximizing cumulative reward = minimizing total regret

What about for our case? The optimal action is the one with the most information.

They compare to best-case scenario, where optimal action is taken all the time.

What's our best-case? Best states to sample would still have max information.

The most uncertainty, it's not quite the same as choosing actions, is it? It's choosing states. Or rather, the action is to sample a state. Some states you're not about and some you've not seen yet.

What is the MAB case? The goal there is most reward. The goal here is most information. I guess the equivalent would be knowing exactly what the distributions look like. The estimate is $\hat{h}(s)$ compared to some "true" $h(s)$. Then I'd only be sampling around the border. So maybe later I don't actually get the information from a state I chose to sample at. I'm getting an improvement on the estimate of the actual information.

$$\text{One option is } L = E[V^* - g(a)] + E[g(a^*) - g(a)].$$

$$\text{Our regret step is } L = E[V^* - h(p(a))]$$

$$L = \frac{1}{n} \sum_{s \in S} g(s) \ln(p(s))$$

↑ this value has weight 1/n
↑ this value has weight 1/n
↑ this value has weight 1/n
↑ this value has weight 1/n

In the MAB case, say I have two bins that are visited equally and I can choose when to sample if my sampling is limited. One is certain, another's not.

I can choose between two actions, which gives me the most information? And what's the best case? The dice are Δ so at least I want to be only picking state 1. I approach the generating distribution $g(s)$.

$$L = E[h(s^*) - h(s)] \text{ only relevant need to switch reward case.}$$

The count $N(s)$ is expected # selections for a .

Gap Δ is $V^* - g(a)$.

[Real time & motivation]

Regret includes both:

$$L_a = E\left[\sum_{s \in S} v_s - g(A_s)\right]$$

$$:= E[N(s)] (V^* - g(s))$$

$$+ E[N(s)] \Delta_s.$$

$$L_v = E\left[\frac{1}{n} \sum_{s \in S} v_s - h(s)\right]$$

↑ the vector v has weight 1/n
↑ the vector v has weight 1/n
↑ the vector v has weight 1/n
↑ the vector v has weight 1/n

$$= E[N(s)] (h(s) - g(s))$$

$$+ E[N(s)] \Delta_s$$

Problem is we don't know the gap. Normally, we only don't know h^* . Here, we don't really know $h(s)$ either. But I suppose that's true for the first case, where $g(s)$ is just a sample.

Need to recall the relationship between a sample (theta-pacing) and the V^* reward.

The sample becomes part of a pool $\sum_{s \in S} p(s) v_s + (1-p(s)) h(s)$ and we go $f(D) \rightarrow f(s)$.

Then $h(p(s)) = p(s) v_s + (1-p(s)) h(s)$

$$= -\left(1 - \frac{p(s)}{p(s) + (1-p(s))}\right) \ln\left(\frac{p(s)}{p(s) + (1-p(s))}\right) + \left(\frac{p(s)}{p(s) + (1-p(s))}\right) \ln\left(\frac{p(s)}{p(s) + (1-p(s))}\right)$$

$$= -\frac{p(s)}{p(s) + (1-p(s))} \ln\left(\frac{p(s)}{p(s) + (1-p(s))}\right) + \frac{(1-p(s))}{p(s) + (1-p(s))} \ln\left(\frac{p(s)}{p(s) + (1-p(s))}\right)$$

So every point $g(s)$ is same $\Delta_h(s)$.

Normally:

$$Q(a) = \max_{s \in S} \sum_{a \in A} R_a$$

We estimate with the mean.

Instead, we have

$$h(s) = g(f(D)).$$

Greedy has linear total regret.

One sol: optimistic initialization. Initialize values to maximum possible, then act greedily.

MAB has linear total regret. Can get linear out of actions forever.

Compare to Greedy, the opposite of great sort of: continue-to-explore forever.

Ensures minimum regret $L \geq \sum_{s \in S} \Delta_s$ over time.

Softmax exploration also has linear total regret.

Can get sublinear regret, however, with a simple fix.

Pick a decay schedule:

$$\min_{a \in A} \Delta_a, C > 0$$

decay so that Δ is scaled by the smallest regret.
What? Gap between best & second best.
Does other touches the threshold work.

$$C = \max_{a \in A} \frac{1}{\Delta_a}, \frac{1}{\Delta_a} \leq C$$

This requires knowledge of the gaps, however when gaps are small, explore more. When gaps are large, explore less.

Greedy logarithmic regret.

Now we want to do this without knowing R .

First note that there's a lower bound on regret (it is logarithmic). What doesn't depend on? Similarity between best arm & other arms.

Hard problems have similar options - lots of noise but one is still better.

Formally, we describe this as Δ and the similarity in distributions

$D(p_1 || p_2)$.

What is our hardest case? Lots of noise in the HGS. In the QGS. In the Q's, the Q distribution itself is very wide. In the HGS case?

What distributions must be the same? The HGS vs. the QGS. The Q's vs. the Q's.

estimate, which means the distributions of both p and g , which in turn means the distribution of $Q(s)$ versus $Q^*(s)$.

The lower bound depends on similarity between the most informative (boundary) states and the other states.

$(Q(s)) \approx Q^*(s)$

Information-theoretic bound:

$$D(p_1 || p_2) = D(p_1) + D(p_2) - D(p_1, p_2)$$

What does this tell us? The lower bound is the average of the two distributions.

Greedy has linear total regret.

One sol: optimistic initialization. Initialize values to maximum possible, then act greedily.

MAB has linear total regret. Can get linear out of actions forever.

Compare to Greedy, the opposite of great sort of: continue-to-explore forever.

Ensures minimum regret $L \geq \sum_{s \in S} \Delta_s$ over time.

Softmax exploration also has linear total regret.

These things are not necessarily true.

Greedy has linear total regret.

Greedy has linear total regret.