

Lecture 11

Lecture 11-13 Fast Reinforcement Learning

Idea

Move towards algorithms that can be more computational-wise and sample-wise efficient?

Also mean that the optimal decision can be calculated quickly (require less time to react) and need less data to compute on.

Key terms

1. Bandits
2. Multi-armed bandits
3. Regret
4. Upper confidence bound (UCB)
5. Hoeffding's Inequality

Bandits

Actions have no influence on
next observation & reward

Slot machine



Bandits

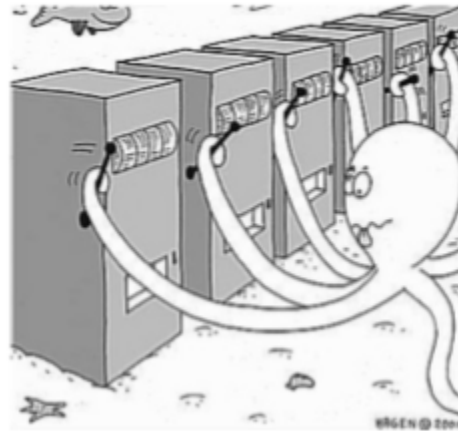
Actions have no influence on
next observation & reward

Slot machine = one-armed bandit
= each action is independent



Multi-armed Bandits

Multiple slot machine with multiple arms



Can choose which arm(action) to pull
and the slot machine will return the
rewards

Regret

Opportunity loss

Just an indicator for
algorithm evaluation

Regret

Opportunity loss

$$l_t = \mathbb{E}[V^* - Q(a_t)]$$

Difference of reward from current pull of the arm with the optimal value

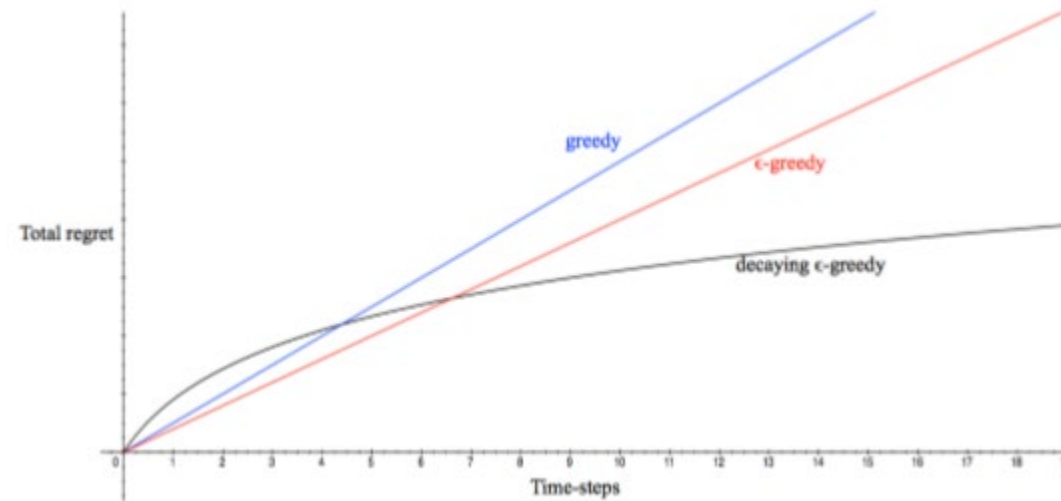
Or in other words this part is the Gap $\Delta_i = V^* - Q(a_i)$

Then calculate total regret =
regret for all timesteps

$$\begin{aligned} L_t &= \mathbb{E} \left[\sum_{\tau=1}^t V^* - Q(a_\tau) \right] \\ &= \sum_{a \in \mathcal{A}} \mathbb{E}[N_t(a)] (V^* - Q(a)) \\ &= \sum_{a \in \mathcal{A}} \mathbb{E}[N_t(a)] \Delta_a \end{aligned}$$

Count of the given action (arm pulled)

Upper confidence bound



Aim is to get sublinear regret

Because linear regret just means getting worse over time (continuously not picking the optimal action)

But first, In order to define the upper confidence bound...

Hoeffding's Inequality

$$\mathbb{P} [\mathbb{E}[X] > \bar{X}_n + u] \leq \exp(-2nu^2)$$



$$\bar{X}_n + u \geq \mathbb{E}[X] \text{ w.p. prob } \geq 1 - \delta/t^2$$

Upper bound can deviate from the expected value by more than a certain amount

Hoeffding's Inequality

$$\mathbb{P} [\mathbb{E}[X] > \bar{X}_n + u] \leq \exp(-2nu^2)$$



$$\bar{X}_n + u \geq \mathbb{E}[X] \text{ w. prob } \geq 1 - \delta/t^2$$

Where...

$$u = \sqrt{\frac{1}{2n} \log(t^2/\delta)}$$

Upper confidence bound

$$U_t(a_t) = \hat{Q}(a_t) + \sqrt{\frac{1}{2n(a_t)} \log(t^2/\delta)}$$

Upper confidence bound

$$U_t(a_t) = \hat{Q}(a_t) + \sqrt{\frac{2 \log t}{N_t(a_t) \log(1/\delta)}}$$

UCB1 algorithm is choosing the action (arm) with the best upper bound value

$$a_t = \arg \max_{a \in \mathcal{A}} \left[\hat{Q}(a) + \sqrt{\frac{2 \log t}{N_t(a)}} \right]$$

Upper confidence bound

$$U_t(a_t) = \hat{Q}(a_t) + \sqrt{\frac{2 \log t}{N_t(a_t) \log(1/\delta)}}$$

UCB1 algorithm is choosing the action (arm) with the best upper bound value

$$a_t = \arg \max_{a \in \mathcal{A}} \left[\hat{Q}(a) + \sqrt{\frac{2 \log t}{N_t(a)}} \right]$$

Key terms of lecture 12

1. Bayesians bandit
2. Thompson Sampling
3. Probably Approximately Correct (PAC)

Bayesian

$$\begin{array}{c} \text{Posterior} \\ \downarrow \\ P(A|B) \end{array} = \frac{\begin{array}{c} \text{Likelihood} \\ \downarrow \\ P(B|A) \end{array} * \begin{array}{c} \text{Prior} \\ \downarrow \\ P(A) \end{array}}{\begin{array}{c} P(B) \\ \uparrow \\ \text{Evidence} \end{array}}$$

Bayesian bandit

$$\begin{array}{c} \text{Posterior} \\ \downarrow \\ P(A|B) \end{array} = \frac{\begin{array}{c} \text{Likelihood} \\ \downarrow \\ P(B|A) \end{array} * \begin{array}{c} \text{Prior} \\ \downarrow \\ P(A) \end{array}}{\begin{array}{c} \uparrow \\ P(B) \\ \text{Evidence} \end{array}}$$

$$p(\phi_i | r_{i1}) = \frac{p(r_{i1} | \phi_i) p(\phi_i)}{\int_{\phi_i} p(r_{i1} | \phi_i) p(\phi_i) d\phi_i}$$

Bayesian bandit

$$\begin{array}{c}
 \text{Posterior} \downarrow \\
 P(A|B) = \frac{
 \begin{array}{c}
 \text{Likelihood} \downarrow \\
 P(B|A) * \text{Prior} \downarrow \\
 P(A)
 \end{array}
 }{
 \begin{array}{c}
 \uparrow \\
 P(B) \\
 \text{Evidence}
 \end{array}
 }
 \end{array}$$

Data likelihood

$$p(\phi_i | r_{i1}) = \frac{p(r_{i1} | \phi_i) p(\phi_i)}{\int_{\phi_i} p(r_{i1} | \phi_i) p(\phi_i) d\phi_i}$$

Reward of particular action depends on this parameter

Distribution of rewards

Bayesian bandit

$$\begin{array}{c}
 \text{Posterior} \\
 \downarrow \\
 P(A|B) = \frac{
 \begin{array}{c}
 \text{Likelihood} \\
 \downarrow \\
 P(B|A) * P(A) \\
 \uparrow \\
 \text{Prior}
 \end{array}
 }{
 \begin{array}{c}
 P(B) \\
 \uparrow \\
 \text{Evidence}
 \end{array}
 }
 \end{array}$$

$$p(\phi_i | r_{i1}) = \frac{p(r_{i1} | \phi_i) p(\phi_i)}{\int_{\phi_i} p(r_{i1} | \phi_i) p(\phi_i) d\phi_i}$$

Conjugate....

Bayesian bandit

$$\text{Regret}(\mathcal{A}, T; \theta) = \sum_{t=1}^T \mathbb{E} [Q(a^*) - Q(a_t)]$$

$$\text{BayesRegret}(\mathcal{A}, T; \theta) = \mathbb{E}_{\theta \sim p_\theta} \left[\sum_{t=1}^T \mathbb{E} [Q(a^*) - Q(a_t) | \theta] \right]$$

Thompson Sampling

Because computing optimal action from posterior can be difficult, a simpler approach would be as following

Thompson Sampling

-
- 1: Initialize prior over each arm a , $p(\mathcal{R}_a)$
 - 2: **loop**
 - 3: For each arm a **sample** a reward distribution \mathcal{R}_a from posterior
 - 4: Compute action-value function $Q(a) = \mathbb{E}[\mathcal{R}_a]$
 - 5: $a_t = \arg \max_{a \in \mathcal{A}} Q(a) \leftarrow$
 - 6: Observe reward r
 - 7: Update posterior $p(\mathcal{R}_a|r)$ using Bayes law
 - 8: **end loop**
-

Easier understood through Bernoulli toy example

In attempt to understand...

PAC Probably Approximately Correct

Input epsilon and delta in all but N steps

The algorithm selects action which its true Q value will be greater than the best possible value of that state subtracted by epsilon

With probability of at least $1 - \delta$

N is a polynomial function of size of S, A and gamma, epsilon and delta

Probably Approximately Correct RL
 input ϵ, δ on all but N steps
 our alg will select an action
 $Q^*(s_t, a_t) \geq V^*(s_t) - \epsilon$
 w/prob at least $1 - \delta$
 where N poly func ($|S|, |A|, \gamma, \epsilon, \delta$)
 Yeah.

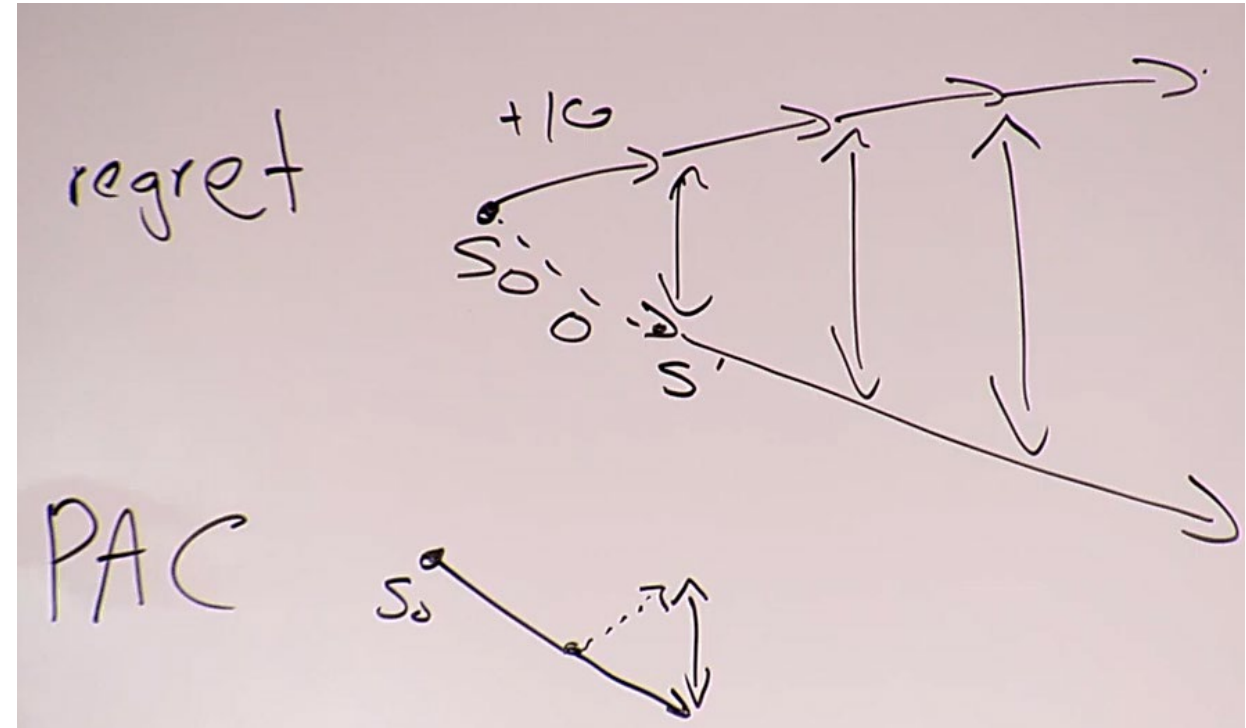
lecture 13

In attempt to understand...

PAC Probably Approximately Correct

PAC can be thought of as making the most out of its circumstances

While regret judges on whether good decisions have been made from the beginning or not.



In attempt to understand...

PAC Probably Approximately Correct

Sufficient condition for PAC

1. Optimism $\underline{Q}_t(s,a) \geq Q^*(s,a) - \epsilon \quad \forall s,a,t$
 computed value should be optimistic with respect to
 real Q values

2. Accuracy $V_t(s) - V^{\pi_t}(s) \leq \epsilon$
 μ will define further. MDP related to true MDP
 \S MDP defined in MBLE-EB

3) Bounded learning complexity:
 - total # of updates to Q
 - # times visit an "unknown" (s,a) pair
 bounded by $\S(\epsilon, \delta)$

In attempt to understand...

MBIE-EB Model Based Interval Estimation with Exploration Bonus

```

1: Given  $\epsilon, \delta, m$ 
2:  $\beta = \frac{1}{1-\gamma} \sqrt{0.5 \ln(2|S||A|m/\delta)}$ 
3:  $n_{sas}(s, a, s') = 0 \quad s \in S, a \in A, s' \in S$ 
4:  $rc(s, a) = 0, n_{sa}(s, a) = 0, \tilde{Q}(s, a) = 1/(1-\gamma) \quad \forall s \in S, a \in A$ 
5:  $t = 0, s_t = s_{init}$ 
6: loop
7:    $a_t = \arg \max_{a \in A} Q(s_t, a)$ 
8:   Observe reward  $r_t$  and state  $s_{t+1}$ 
9:    $n_{sa}(s_t, a_t) = n(s_t, a_t) + 1, n_{sas}(s_t, a_t, s_{t+1}) = n_{sas}(s_t, a_t, s_{t+1}) + 1$ 
10:   $rc(s_t, a_t) = \frac{rc(s_t, a_t)n_{sa}(s_t, a_t) + r_t}{(n_{sa}(s_t, a_t) + 1)}$ 
11:   $\hat{R}(s, a) = \frac{rc(s_t, a_t)}{n(s_t, a_t)}$  and  $\hat{T}(s'|s, a) = \frac{n_{sas}(s_t, a_t, s')}{n_{sa}(s_t, a_t)} \quad \forall s' \in S$ 
12:  while not converged do
13:     $\tilde{Q}(s, a) = \hat{R}(s, a) + \gamma \sum_{s'} \hat{T}(s'|s, a) \max_{a'} \tilde{Q}(s', a') + \underbrace{\frac{\beta}{\sqrt{n_{sa}(s, a)}}}_{\text{Bonus}} \quad \forall s \in S, a \in A$ 
14:  end while
15: end loop

```

Diagram annotations:

- Backup**: Points to the term $\gamma \sum_{s'} \hat{T}(s'|s, a) \max_{a'} \tilde{Q}(s', a')$.
- Bonus**: Points to the term $\frac{\beta}{\sqrt{n_{sa}(s, a)}}$.
- Transition model**: Points to the term $\hat{T}(s'|s, a)$.