

Multi-armed Bandits

Rupam Mahmood



January 15, 2020



Expectations review

- Expectation: $E[X] = \sum x P(X=x)$ \checkmark
- summarizes the outcomes of an experiments \checkmark
- Conditional expectation w.r.t. an event: $E[X | Y=y] = \sum x P(X=x| Y=y)$ \checkmark
- Conditional expectation w.r.t a random variable: E[X | Y] \checkmark
- ... is a random variable that maps outcomes to numbers: \checkmark
- E[X | Y](y) = E[X | Y=y] \checkmark

... is a function of Y that "best approximates" X \checkmark

Worksheet question 1

Action values

$q_*(a) \doteq E[R | A = a]$

✓ What is the sample space?

✓ What is the optimal behavior?

Estimating action values

 $Q_{n+1} = \frac{1}{n} \sum_{i=1}^{n} R_i = \frac{1}{n} \sum_{i=1}^{n-1} R_i + \frac{1}{n} R_n$



Life of a bandit agent

Initialize N and Q

Loop forever:

take an action A based on Q and an action-selection strategy

Observe reward R

Update estimates N and Q

Pseudocode

A simple bandit algorithm

Initialize, for
$$a = 1$$
 to k :
 $Q(a) \leftarrow 0$
 $N(a) \leftarrow 0$

Loop forever:

 $A \leftarrow \begin{cases} \operatorname{arg\,max}_a Q(a) & \text{with probability } 1 - \varepsilon & \text{(breaking ties randomly)} \\ \text{a random action} & \text{with probability } \varepsilon \end{cases}$ $R \leftarrow bandit(A)$ $N(A) \leftarrow N(A) + 1$ $Q(A) \leftarrow Q(A) + \frac{1}{N(A)} \left[R - Q(A) \right]$

Life of an RL-glue bandit agent

```
for i in range(num_steps):
```

. . .

reward, _, action, _ = rl_glue.rl_step()

A simple bandit algorithm

Initialize, for
$$a = 1$$
 to k :
 $Q(a) \leftarrow 0$
 $N(a) \leftarrow 0$

Loop forever:agent step $A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases}$ (breaking ties randomly) $R \leftarrow bandit(A)$ environment step $N(A) \leftarrow N(A) + 1$ $Q(A) \leftarrow Q(A) + \frac{1}{N(A)} \left[R - Q(A) \right]$

Review & clarifications

- \checkmark action?
- \checkmark why it could be constant?
- \checkmark agent randomly could be the greedy action again?

Is the action which has highest expected reward(value) the optimal

What exactly the stepsize is? The meaning and influence of it? and

Is it possible that in e-greedy, with probability e, the action taken by