

CMPUT 365: RL

K-Armed Bandits

Rupam Mahmood

Jan 10, 2022

Admin

- Due dates:
 - C1M1 practice quiz: Tomorrow
 - C1M1 graded notebook: This saturday
- Each week:
 - Focus on one Coursera module
 - Read the chapter
 - Watch the Coursera videos
 - Work on practice quiz (important for exams and written assignments)
 - Solve worksheet questions
 - Finish graded notebook

Coursera videos: Video 1

- First introduction to problem of decision making
- Objects: ^{Sets A} actions, ^{Sets R} rewards, time steps or trials, values

they are associated with some uncertainty

only takes a single value input, not a set value input. \Rightarrow

$$q_{\star}(a) = \underbrace{E}_{\text{expectation}}[\underbrace{R}_{\text{can't appear as an input at the LHS 因为它不是一个随机变量}} | \underbrace{A=a}_{\text{choosing a particular arm "a" / 即是一个RV.}}]$$

the reward of choosing a particular action "a".

Random variable A takes a particular outcome.

- Check [this extra video](#) connecting bandits with probabilities we just learned

$$\begin{aligned} y &= E[R | A=a] \\ &= \sum_r P(R=r | A=a) r \end{aligned}$$

- 我们要 decide which action is best. 就要 find the value of taking each action. **RP Action-Values** (Policy Action-Values Function)

Define the value of expecting an action as the expected reward: (goal: maximize the expected reward)

is defined as

expected: $q_*(a) \triangleq \mathbb{E}[R_t | A_t = a] \quad \forall a \in \{1, \dots, k\}$

reward + given selected action a

(Value of an action, not known)

for each possible action $1 \leq k$.

If the agent (like) chose the action has the highest value, it achieve the goal. This procedure is called

$\text{Argmax}_a q_*(a)$

the mean of the distribution of each action.

$= \sum_r p(r|a) r \Rightarrow$ the sum of all possible rewards

(multiply possible reward with the probability of observing that reward).

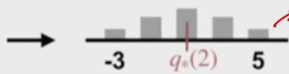
Calculating $q_*(a)$



低风险的分布, 11 与 9 各 0.5

$\rightarrow q_*(a) = .5 \times -11 + .5 \times 9 = -1$

the expected value of (低风险的) distribution



中风险的分布

$\rightarrow q_*(a) = 1$



$\rightarrow q_*(a) = 3$

不同的 action 都有 不同的 distribution.

Video 2

- Estimating action values

The agent doesn't know q_* . But the \Rightarrow agent knows Q . The agent wants to estimate q_* .

$q_*(a) = E[R|A=a] \approx Q_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$

$q_*(a)$ is the true quantity. $Q_t(a)$ is the estimator (the sample-average method to estimating q_*). a is constant. Sample average.

each time you get the different values, estimator $Q(a)$ will get closer and closer to the true value $q_*(a)$.

(The estimator usually has the big letter. \therefore It depends on R , giving a different value.)

- Why can't we choose actions greedily while estimating action values?

$$A_t \doteq \underset{a}{\operatorname{argmax}} Q_t(a)$$

A_t is a random value.
 \therefore The greedy choice is bias in a sense.

shway to Everyone

1:37 PM

5

It might be suboptimal but the agent does not know and it keeps choosing that suboptimal action greedily.

- If we had perfect estimates of action values, can we choose greedily?

Video 3

- Estimating action values incrementally

$$Q_n \doteq \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

Cameron Jen to Everyone

1:40 PM



So Q_n = Estimated value before choosing the action for the n -th time right? *yes!!*

$$Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$$

R_n has been taken, the number of times that each action has been taken.

finding the recursive update rule for sample average.

- Incremental learning is more generally applicable to both stationary and nonstationary problems

$$\text{NewEstimate} \leftarrow \text{OldEstimate} + \underbrace{\text{StepSize}}_{\text{can be sth else.}} [\text{Target} - \text{OldEstimate}]$$

$$Q_{n+1} \doteq Q_n + \underbrace{\alpha}_{\text{constant stepsize. (不变的)}} [R_n - Q_n]$$

*constant stepsize. (不变的) \Rightarrow we can solve the bandit problem change over time.
常设步长的 stepsize, 则步长为 $\frac{1}{n-1}$.*

Video 4 & 5

- (4) The exploration-exploitation tradeoff
- epsilon-greedy, as a simple method to balance exploration and exploitation

$$A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \varepsilon \quad (\text{breaking ties randomly}) \\ \text{a random action} & \text{with probability } \varepsilon \end{cases}$$

- (5) Optimistic initial values encourage early exploration
- Limitation of optimistic initialization: not well suited for nonstationary problems

greedy 方法 - 一定会得到最优解!! 因为有可能 agent 会特别喜欢某个 meal 进而反复选它, 导致他根本没有发现其他 meal 的价值.

Exploration versus Exploitation

∴ 该选择 explore
(尝试点一个新菜)

还是 exploit?
(选以前点过的菜).

- **Exploration** - **improve** knowledge for **long-term** benefit \Rightarrow get more accurate estimate to our value
- **Exploitation** - **exploit** knowledge for **short-term** benefit \Rightarrow might get more reward

- How do we choose when to **explore** and when to **exploit**?

ϵ -greedy: 继续 explore 的概率:

Life of a bandit agent

A simple bandit algorithm

Initialize, for $a = 1$ to k :

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

tail to what
action we take
and what reward
we receive.

$$A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases} \quad \text{(breaking ties randomly)}$$
$$R \leftarrow \text{bandit}(A)$$
$$N(A) \leftarrow N(A) + 1$$
$$Q(A) \leftarrow Q(A) + \underbrace{\frac{1}{N(A)}}_{\text{constant}} [R - Q(A)]$$

Video 6

- Upper confidence bound action-selection uses uncertainty in the estimates to drive exploration

Mohamed Ahmed to Everyone 1:50 PM

MA For example issues with using bigger e is once you explored a lot and found the better actions, bigger e would still prefer to explore more

exploration, 对已经 explore 的概率. 是为什么还要去探索呢?

$$A_t \doteq \underset{a}{\operatorname{argmax}} \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

if only Q_t .
you choose the greedy action.

探索的概率加了个因子. (if you select a random action).

- $N_t(a)$ being in the denominator decreases a 's uncertainty estimate after being chosen
- $\ln t$ in the numerator increases the uncertainty estimates of actions that are not chosen
- Over time, increases become smaller but remains unbounded

Slido questions

- Action value: 1,
- Step size: 1,
- Epsilon-greedy: 1,

The action with the highest estimated value is called

☐ optimal action

☒ greedy action \rightarrow the largest true value.

Send

Voting as Anonymous

Acceptable Use - Privacy Policy

Talking:



From Sajjad Kavyani Baghbaderani...
greedy

A good step size for a non-stationary task is

- ☐ $1/t$ \Rightarrow 若越久, 则 step size will be smaller and smaller. Earlier values get more weight. 越久的值必越多是 the sample average $\bar{x} = \frac{x_1 + \dots + x_{n-1}}{n-1}$
- ☒ constant \Rightarrow you will estimate the change from the beginning to the end
be able to

Send

Voting as Anonymous

Acceptable Use - Privacy Policy

Talking: Rupam Mahmood

Recording

You are viewing Rupam Mahmood's screen

View Options

View

365w22-bandits - Google Slides X

365w22-bandits

365w22-bandits

365w22-bandits

Course: CMPUT 365 LEC B1 - X

+

← → ↺

https://app.sli.do/event/7BdTJeEnLpUdpZsnu9B4Ff/embed/polls/85a2e82b-e83f-4803-8a4d-f6dbc2e57395

☆

Inbox - ashique@ua...

University of Albert...

evernote - Goog...

Rupam's Website

CMPUT 340 Winter ...

CMPUT 365 Winter ...

Dictionary.com | Me...

Edmonton, Alberta ...

describing the expe...

diary - Google Docs

Fellows Intranet

Other Bookmarks

Is it possible that in an epsilon-greedy method, the action taken by the agent randomly could be the greedy action?

✓

Yes

⇒ 有 $1-\epsilon$ 的概率是 $\text{argmax}_a Q(a)$. ∴ 当有 ϵ 的概率遇到 a random action 时, 该 action 依然有可能是停留在 $\text{argmax}_a Q(a)$ 里的.

○

No

∴ probability of greedy action = $(1-\epsilon) + \frac{\epsilon}{\# \text{ of actions}}$

Send

Voting as Anonymous

Acceptable Use · Privacy Policy

Talking: Rupam Mahmood

Unmute

Start Video

Participants 49

Chat

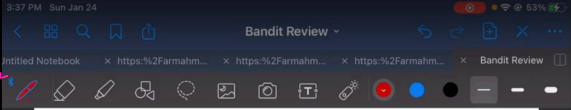
Share Screen

Record

Reactions

Leave

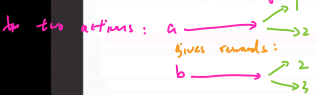
Bandit problems are a sequence of two experiments.



Bandits review

- ✓ What is the experiment? \Rightarrow the process the agent choosing actions.
- ✓ What are the outcomes?
- ✓ What are the random variables involved?

In Bandit Problem, we have an agent choosing actions and an environment providing rewards.



It is an experiment of {a, b}. It is a Random Variable: A for actions.

\therefore $P(A=a)$ & $P(A=b)$.

outcomes of {1, 2, 3}, the probability is depend on the action that was chosen. It is a random variable R. \therefore $P(R=1 | A=a)$

probability of $R=1$ base on what action was taken. If pick action b, $\text{if } A=b$, $\text{h} \therefore P(R=1 | A=b)=0$

\therefore the experiment is outcomes of {a, 1}, {a, 2}, {b, 2}, {b, 3}.

when we talk about probabilities here, we talk about probabilities with

joint event: $P(A=a, R=1)$
 $= P(R=1 | A=a) P(A=a)$


comes from the reward process.

4:05 PM Sun Jan 24 Contextual bandits

tebook x https:%2Farmahm... x https:%2Farmahm... x https:%2Farmahm... x Bandit Review x Marko

Contextual Bandits

In Bandit problems, the agent choose from n actions, and one of them gives the highest reward:

Actions: 0  1 \Rightarrow The action that gives the highest reward doesn't depend on the context.

The best action depends on the context!!

eg. 苹果, the best action is the first one: 0 0 0

香蕉, the best action is the second one: 0 1 0

The context 可用一个 R.V. 来表示, 记 S as state. Depending on the state, the agent choose an action A and receive the reward R .

\therefore The whole experiment is: $P(S=s, A=a, R=r)$ \Rightarrow a composition of 3 different events.

\leftarrow hit state being presented

\leftarrow can be computed using the condition probability chosen by the agent

definition, P :

$$P(R=r | S=s, A=a) \times P(A=a | S=s) \times P(S=s)$$

the reward depends on the state and the action

The action depends on the state

probability of the environment choosing this state.