

4/6/2023

Supervised Learning

Dataset = $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$

Unsupervised Learning

Dataset $\{x^{(i)}\}_{i=1}^n$

Someone gives us dataset
Learning algorithm has no influence
on what the dataset is
"passive"

Bandits & Reinforcement Learning

Learning algorithm takes actions:

- ① Influence what info is observed
↑
ie data
- ② Influence state of agent / world

Example: Students selecting classes

- Action: choose a class to take
- ✓ - Information: You learn whether you like the class you take
don't learn about other classes
- X - State: taking into class unlocks more advanced classes

Bandits

- Taking action gives info only about that action
- No state changing between actions

Medicine

K treatment options for a condition
w/ unknown effectiveness

Patients come in one at a time

Actions: Prescribe one of K treatments

Information: Outcome of that patient

Early on, try all the treatments

Eventually learn which is best & mostly prescribe that

News Headlines K headlines for article, want people to click

Action: For each visitor to site, show 1 of the K headlines
Info: Did they click on article

Stochastic Multi-Armed Bandit Problems
one-armed bandit = slot machine
 K different slot machines, learn which gives most \$\$
Rewards are random

Set of actions $\{1, \dots, K\}$ ("arms")

Each action a has reward distribution $P_a(r)$ ← Unknown ahead of time
Fixed for all t

- What payoff did you receive?
- Health of patient after treatment
 - Did person click or not (1 or 0)

(In bandits/RL we maximize reward and not minimizing loss)

Player play this game for T rounds

At each time $t=1, \dots, T$:

- Player chooses action $A_t \in \{1, \dots, K\}$
- Player receive reward $R_t \sim P_{A_t}(R)$

Depends on observed rewards R_1, \dots, R_{t-1} .
So A_t is also a random variable

Goals: maximize total reward $\sum_{t=1}^T R_t$
 $\in \mathbb{R}$

Example
 $K=2$

| t | A_t | R_t |
|-----|-------|-------|
| 1 | 1 | 1 |
| 2 | 2 | 0 |
| 3 | 1 | 0 |
| 4 | 1 | 0 |
| 5 | 2 | 1 |
| 6 | 2 | 1 |
| 7 | 2 | 1 |

// Action 1 better, lets try it more

// maybe action 1 isnt that good?

If $T=7$ then total reward = 4

Regret: How well did your strategy do compared to the optimal strategy

Define $N(a) =$ expected reward when choosing action $a = \mathbb{E} [R]$
 $R \sim P_a(R)$

Optimal action $a^* = \operatorname{argmax}_a N(a)$

Expected Regret: Expected difference between optimal strategy & player's strategy

$$\underbrace{N(a^*) \cdot T}_{\text{Expected reward of optimal}} - \underbrace{\mathbb{E} \left[\sum_{t=1}^T R_t \right]}_{\text{Expected reward for player}}$$

$$= N(a^*) \cdot T - \sum_{t=1}^T N(A_t)$$

Regret close to 0 good
 large regret is bad

Exploration vs Exploitation

Want to try all the actions enough times to learn which is better
ie gain knowledge that's useful later

Use current knowledge to do what currently seems best

"I like math classes
I will keep taking math"

"I want to try lots of subjects before specializing"

Algorithm: Upper Confidence Bound (UCB) Algorithm

Idea:

- Player is estimating $N(a)$ for every a
- Estimates are uncertain
↳ we will represent this as a confidence interval
"I think $N(a)$ is between 0.6 and 0.8"
- At each t , choose action with largest upper bound

↑
lower bound

↑
upper bound

Why? : Optimism in the face of uncertainty

If a is action with largest upper bound:

Either ① it's actually good = very good news

② It's not so good \Rightarrow can update our estimates & try something else

UCB Algorithm: Assume $0 \leq R_t \leq 1$

At time t , Let $n_t(a)$ denote # of times we tried a up until time t

Size of dataset we collected about a

| t | A_t | R_t |
|-----|-------|-------|
| 1 | 1 | 1 |
| 2 | 2 | 0 |
| 3 | 1 | 0 |
| 4 | 1 | 0 |
| 5 | 2 | 1 |
| 6 | 2 | 1 |
| 7 | 2 | 1 |

$n_8(1) = 3, n_8(2) = 4$

Let $\hat{N}_t(a)$ be sample mean

of rewards when taking action a in the data before time t

$\hat{N}_8(1) = 1/3, \hat{N}_8(2) = 3/4$

How much uncertainty is in a sample mean?
 If n examples
 Variance of sample mean
 $= \frac{\sigma^2}{n}$ ← Variance of one sample
 ⇒ Standard deviation of sample mean is σ/\sqrt{n}

For UCB:

For action a , use confidence interval of \pm at time t

$\pm \sqrt{\frac{2 \log t}{n_t(a)}}$

so it's $O\left(\frac{1}{\sqrt{n_t(a)}}\right)$

ie $N(a) \in \left[\hat{N}_t(a) - \sqrt{\frac{2 \log t}{n_t(a)}}, \hat{N}_t(a) + \sqrt{\frac{2 \log t}{n_t(a)}} \right]$

$= UCB_t(a)$

$$UCB_t(a) = \hat{N}_t(a) + \sqrt{\frac{2 \log t}{N_t(a)}}$$

Exploitation term
 "a is good if we think its reward is high"

Exploration
 "a is useful if we haven't tried it very much yet"

Gets smaller as $N_t(a)$ gets larger

Full algorithm:

1. For $t=1, \dots, K$: try each action once
2. For $t=K+1, \dots, T$: choose $A_t = \operatorname{argmax}_a UCB_t(a)$

What happens to $UCB_t(a)$ over time?

$$\sqrt{\frac{2 \log t}{N_t(a)}}$$

Gets bigger over time very slowly

\Rightarrow we never completely rule out an action.

Gets bigger over time
 \Rightarrow UCB gets closer to \hat{N}
 \Rightarrow do explorations

If $N_t(a)$ constant eventually this UCB gets large

Can prove a bound on Regret of UCB

In particular, Regret is $O\left(\sqrt{KT \log T}\right)$

This is good because it's sublinear

If we average across timesteps, average regret is

$$O\left(\frac{\sqrt{KT \log T}}{T}\right) \rightarrow 0 \text{ as } T \rightarrow \infty$$

After enough time, gap w/ optimal strategy is negligible.