Lecture 2: Learning from Evaluative Feedback

> or "Bandit Problems"

Edward L. Thorndike (1874-1949)

Learning by Trial-and-Error

Law of Effect:

"Of several responses to the same situation, those which are accompanied or closely followed by satisfaction to the animal will, other things being equal, be more firmly connected with the situation, so that, when it recurs, they will be more likely to recur; those which are accompanied or closely followed by discomfort to the animal will, other things being equal, have their connections with that situation weakened, so that when it recurs, they will be less likely to recur."

Thorndike, 1911

Some Distinctions

• Evaluating actions vs. instructing by giving correct actions

Puzzle Box

- Pure evaluative feedback depends totally on the action taken. Pure instructive feedback depends not at all on the action taken.
- Supervised learning is instructive; optimization is evaluative
- Associative vs. Non-associative:
 - Associative: inputs mapped to outputs; learn the best output for each input
 - Non-associative: "learn" (find) one best output
- *n*-armed bandit (at least how we treat it) is:
 - Non-associative
 - Evaluative feedback

The *n*-Armed Bandit Problem • Choose repeatedly from one of *n* actions; each choice is called a **play** • After each play a_t , you get a reward r_t , where $E\langle r_t | a_t \rangle = Q^*(a_t)$ These are unknown **action values** Distribution of r_t depends only on a_t • Objective is to Maximize the reward in the long term, e.g., over 1000 plays

To solve the *n*-armed bandit problem, you must **explore** a variety of actions and the **exploit** the best of them

5

5

The Exploration/Exploitation Dilemma

Suppose you form estimates
Q_t(a) ≈ Q^{*}(a) action value estimates

The greedy action at t is:
a_t^{*} = arg max _a Q_t(a)
a_t = a_t^{*} ⇒ exploitation
a_t ≠ a_t^{*} ⇒ exploitation

You can't exploit all the time; you can't explore all the time
You can never stop exploring; but you should always reduce exploring

6

Observation of the end of the en

ε-Greedy Action Selection

• Greedy action selection:

$$a_t = a_t^* = \arg\max_a Q_t(a)$$

• ε-Greedy:

$$a_t = \begin{cases} a_t^* \text{ with probability } 1 - \varepsilon \\ \text{random action with probability } \varepsilon \end{cases}$$

... the simplest way to try to balance exploration and exploitation

8

10-Armed Testbed

- n = 10 possible actions
- Each $Q^*(a)$ is chosen randomly from a normal distribution: $\eta(0,1)$
- each r_t is also normal: $\eta(Q^*(a_t), 1)$
- 1000 plays
- repeat the whole thing 2000 times and average the results

ε-Greedy Methods on the 10-Armed Testbed



Softmax Action Selection

- Softmax action selection methods grade action probs. by estimated values.
- The most common softmax uses a Gibbs, or Boltzmann, distribution:

Choose action a on play t with probability

$$\frac{e^{Q_t(a)/\tau}}{\sum_{b=1}^n e^{Q_t(b)/\tau}},$$
 where τ is the

"computational temperature"

II

Binary Bandit Tasks

Suppose you have just **two** actions: $a_t = 1$ or $a_t = 2$ and just **two** rewards: $r_t = success$ or $r_t = failure$

Then you might infer a target or desired action:

 $d_t = \begin{cases} a_t & \text{if success} \\ \text{the other action} & \text{if failure} \end{cases}$

and then always play the action that was most often the target

Call this the **supervised algorithm** It works fine on deterministic tasks...

12







Incremental Implementation

Recall the sample average estimation method:

The average of the first k rewards is (dropping the dependence on a):

 $Q_k = \frac{r_1 + r_2 + \cdots + r_k}{k}$

Can we do this incrementally (without storing all the rewards)?

We could keep a running sum and count, or, equivalently:

$$Q_{k+1} = Q_k + \frac{1}{k+1} [r_{k+1} - Q_k]$$

This is a common form for update rules:

NewEstimate = *OldEstimate* + *StepSize*[*Target* – *OldEstimate*]

16



Optimistic Initial Values

- All methods so far depend on $Q_0(a)$, i.e., they are **biased**.
- Suppose instead we initialize the action values **optimistically**,

i.e., on the 10-armed testbed, use $Q_0(a) = 5$ for all a



Reinforcement Comparison

- Compare rewards to a reference reward, $\overline{r_t}$, e.g., an average of observed rewards
- Strengthen or weaken the action taken depending on $r_t \overline{r_t}$
- Let $p_t(a)$ denote the **preference** for action a
- Preferences determine action probabilities, e.g., by Gibbs distribution:

$$\pi_{t}(a) = \Pr\{a_{t} = a\} = \frac{e^{p_{t}(a)}}{\sum_{b=1}^{n} e^{p_{t}(b)}}$$

• Then: $p_{t+1}(a_t) = p_t(a) + [r_t - \overline{r}_t]$ and $\overline{r}_{t+1} = \overline{r}_t + \alpha [r_t - \overline{r}_t]$

Performance of a Reinforcement Comparison Method



19

Pursuit Methods

- Maintain both action-value estimates and action preferences
- Always "pursue" the greedy action, i.e., make the greedy action more likely to be selected
- After the *t*-th play, update the action values to get Q_{t+1}
- The new greedy action is $a_{t+1}^* = \arg \max_{a} Q_{t+1}(a)$
- Then:

$$\pi_{t+1}(a_{t+1}^*) = \pi_t(a_{t+1}^*) + \beta \Big[1 - \pi_t(a_{t+1}^*) \Big]$$

21

and the probs. of the other actions decremented to maintain the sum of 1



Performance of a Pursuit Method





Next Class

• Read Chapter 3