Wrap-up of Bandits Sequential decision making Markov Decision Processes

Gradient-Bandit Algorithms

• Let $H_t(a)$ be a learned preference for taking action a

$$\Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \pi_t(a)$$

Note that this allows us to work with unnormalized preferences and turn them into probabilities!

Same idea as using potentials in graphical models

Softmax (Boltzmann) Exploration

• Let $H_t(a)$ be a learned preference for taking action a

$$\Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \pi_t(a)$$

Consider
$$H_t(a) = Q_t(a)/T$$

This is Boltzmann or softmax exploration!

If the temperature T is very large (towards infinity) - same as uniform

If temperature T goes to 0, same as greedy

Derivation of gradient-bandit algorithm

In exact gradient ascent:

$$H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E} [R_t]}{\partial H_t(a)}, \qquad (1)$$

where:

$$\mathbb{E}[R_t] \doteq \sum_b \pi_t(b) q_*(b),$$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[\sum_b \pi_t(b) q_*(b) \right]$$
$$= \sum_b q_*(b) \frac{\partial \pi_t(b)}{\partial H_t(a)}$$
$$= \sum_b \left(q_*(b) - X_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)},$$

where X_t does not depend on b, because $\sum_b \frac{\partial \pi_t(b)}{\partial H_t(a)} = 0$.

$$\begin{split} \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} &= \sum_b \left(q_*(b) - X_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)} \\ &= \sum_b \pi_t(b) \left(q_*(b) - X_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)} / \pi_t(b) \\ &= \mathbb{E} \left[\left(q_*(A_t) - X_t \right) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t) \right] \\ &= \mathbb{E} \left[\left(R_t - \bar{R}_t \right) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t) \right], \end{split}$$

where here we have chosen $X_t = \overline{R}_t$ and substituted R_t for $q_*(A_t)$, which is permitted because $\mathbb{E}[R_t|A_t] = q_*(A_t)$. For now assume: $\frac{\partial \pi_t(b)}{\partial H_t(a)} = \pi_t(b)(\mathbf{1}_{a=b} - \pi_t(a))$. Then:

$$= \mathbb{E}\left[\left(R_t - \bar{R}_t\right)\pi_t(A_t)\left(\mathbf{1}_{a=A_t} - \pi_t(a)\right)/\pi_t(A_t)\right] \\= \mathbb{E}\left[\left(R_t - \bar{R}_t\right)\left(\mathbf{1}_{a=A_t} - \pi_t(a)\right)\right].$$

 $H_{t+1}(a) = H_t(a) + \alpha (R_t - \bar{R}_t) (\mathbf{1}_{a=A_t} - \pi_t(a)), \text{ (from (1), QED)}$

Thus it remains only to show that

$$\frac{\partial \pi_t(b)}{\partial H_t(a)} = \pi_t(b) \big(\mathbf{1}_{a=b} - \pi_t(a) \big).$$

Recall the standard quotient rule for derivatives:

$$\frac{\partial}{\partial x} \left[\frac{f(x)}{g(x)} \right] = \frac{\frac{\partial f(x)}{\partial x} g(x) - f(x) \frac{\partial g(x)}{\partial x}}{g(x)^2}.$$

Using this, we can write...

Quotient Rule: $\frac{\partial}{\partial x} \left[\frac{f(x)}{g(x)} \right] = \frac{\frac{\partial f(x)}{\partial x} g(x) - f(x) \frac{\partial g(x)}{\partial x}}{g(x)^2}$

$$\begin{split} \frac{\partial \pi_t(b)}{\partial H_t(a)} &= \frac{\partial}{\partial H_t(a)} \pi_t(b) \\ &= \frac{\partial}{\partial H_t(a)} \left[\frac{e^{H_t(b)}}{\sum_{c=1}^k e^{H_t(c)}} \right] \\ &= \frac{\frac{\partial e^{H_t(b)}}{\partial H_t(a)} \sum_{c=1}^k e^{H_t(c)} - e^{H_t(b)} \frac{\partial \sum_{c=1}^k e^{H_t(c)}}{\partial H_t(a)}}{\left(\sum_{c=1}^k e^{H_t(c)}\right)^2} \qquad (Q.R.) \\ &= \frac{\mathbf{1}_{a=b} e^{H_t(a)} \sum_{c=1}^k e^{H_t(c)} - e^{H_t(b)} e^{H_t(a)}}{\left(\sum_{c=1}^k e^{H_t(c)}\right)^2} \qquad (\frac{\partial e^x}{\partial x} = e^x) \\ &= \frac{\mathbf{1}_{a=b} e^{H_t(b)}}{\sum_{c=1}^k e^{H_t(c)}} - \frac{e^{H_t(b)} e^{H_t(a)}}{\left(\sum_{c=1}^k e^{H_t(c)}\right)^2} \\ &= \mathbf{1}_{a=b} \pi_t(b) - \pi_t(b) \pi_t(a) \\ &= \pi_t(b) (\mathbf{1}_{a=b} - \pi_t(a)). \qquad (Q.E.D.) \end{split}$$

Summary Comparison of Bandit Algorithms



Lecture 9: Exploration and Exploitation

Multi-Armed Bandits

Bayesian Bandits

Probability Matching

 Probability matching selects action a according to probability that a is the optimal action

$$\pi(a \mid h_t) = \mathbb{P}\left[Q(a) > Q(a'), \forall a' \neq a \mid h_t\right]$$

Probability matching is optimistic in the face of uncertainty
 Uncertain actions have higher probability of being max
 Can be difficult to compute analytically from posterior

Lecture 9: Exploration and Exploitation

Multi-Armed Bandits

Bayesian Bandits

Thompson Sampling

Thompson sampling implements probability matching

$$\pi(a \mid h_t) = \mathbb{P}\left[Q(a) > Q(a'), \forall a' \neq a \mid h_t
ight]$$

= $\mathbb{E}_{\mathcal{R} \mid h_t}\left[\mathbf{1}(a = \operatorname*{argmax}_{a \in \mathcal{A}} Q(a))
ight]$

- Use Bayes law to compute posterior distribution $p[\mathcal{R} \mid h_t]$
- Sample a reward distribution *R* from posterior
- Compute action-value function $Q(a) = \mathbb{E}[\mathcal{R}_a]$
- Select action maximising value on sample, $a_t = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(a)$
- Thompson sampling achieves Lai and Robbins lower bound!

Thompson Sampling



- A simple natural Bayesian heuristic
 - Maintain a belief(distribution) for the unknown parameters
 - Each time, pull arm *a* and observe a reward *r*
- Initialize priors using belief distribution
 - For t=1:T:
 - Sample random variable X from each arm's belief distribution
 - Select the arm with largest X
 - Observe the result of selected arm
 - Update prior belief distribution for selected arm

Simple Example







Initially assumes arm *i* with prior Beta(1,1) on μ_i

FLORIDA INTERNATIONAL UNIVERSITY Miami's public research university

• $S_i = #$ "Success", $F_i = #$ "Failure"

Algorithm 1: Thompson Sampling for Bernoulli bandits

```
\begin{array}{l} S_i = 0, F_i = 0.\\ \text{foreach } t = 1, 2, \dots, \text{do}\\ & | \quad \text{For each arm } i = 1, \dots, N, \text{ sample } \theta_i(t) \text{ from the } \text{Beta}(S_i + 1, F_i + 1) \text{ distribution.}\\ & \text{Play arm } i(t) := \arg \max_i \theta_i(t) \text{ and observe reward } r_t.\\ & \text{If } r = 1, \text{ then } S_i = S_i + 1, \text{ else } F_i = F_i + 1.\\ & \text{end} \end{array}
```



Initialization





- For each round:
 - Sample random variable X from each arm's Beta Distribution





- For each round:
 - Sample random variable X from each arm's Beta Distribution
 - Select the arm with largest X





- For each round:
 - Sample random variable X from each arm's Beta Distribution
 - Select the arm with largest X
 - Observe the result of selected arm





- For each round:
 - Sample random variable X from each arm's Beta Distribution
 - Select the arm with largest X
 - Observe the result of selected arm
 - Update prior Beta distribution for selected arm



	Single State	Associative
Instructive feedback		
Evaluative feedback		

	Single State	Associative
Instructive feedback		
Evaluative feedback	Bandits (Function optimization)	

	Single State	Associative
Instructive feedback		Supervised learning
Evaluative feedback	Bandits (Function optimization)	

	Single State	Associative
Instructive feedback	Averaging (Imitiation)	Supervised learning
Evaluative feedback	Bandits (Function optimization)	

	Single State	Associative
Instructive	Averaging	Supervised
feedback	(Imitiation)	learning
Evaluative	Bandits	Contextual
feedback	(Function optimization)	bandits



Agent and environment interact at discrete time steps: t = 0, 1, 2, 3, ...Agent observes state at step t: $S_t \in S$ produces action at step t: $A_t \in \mathcal{A}(S_t)$ gets resulting reward: $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$ and resulting next state: $S_{t+1} \in S^+$

Trajectory and History

□ A sequence of states, actions and rewards

$$\cdots \underbrace{S_t}_{A_t} \underbrace{R_{t+1}}_{A_{t+1}} \underbrace{S_{t+1}}_{A_{t+1}} \underbrace{R_{t+2}}_{A_{t+2}} \underbrace{S_{t+3}}_{A_{t+3}} \underbrace{S_{t+3}}_{A_{t+3}} \cdots$$

 \Box We will sometimes use the notation $\tau_{ij} = S_i A_i R_{i+1} \dots S_j$

- We will use the term *history* to refer to the trajectory prior to the current time step
- In general, next states and rewards can depend on the history since the beginning of time

Markov Property

- □ An assumption about the environment
- Next state and reward depend only on the previous state and action, and noting else that happened in the past

$$p(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a) = p(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a, \tau_t), \forall \tau_t$$

- The assumption is useful to develop, analyze and understand algorithms
- □ It does NOT mean it has to always hold

Markov Decision Processes

- ☐ If a reinforcement learning task has the Markov Property, it is basically a Markov Decision Process (MDP).
- □ If state and action sets are finite, it is a **finite MDP**.
- **T** To define a finite MDP, you need to give:
 - state and action sets
 - one-step "dynamics"

$$p(s', r | s, a) = \mathbf{Pr}\{S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a\}$$

$$p(s'|s,a) \doteq \Pr\{S_{t+1} = s' \mid S_t = s, A_t = a\} = \sum_{r \in \mathcal{R}} p(s',r|s,a)$$
$$r(s,a) \doteq \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s',r|s,a)$$

Policy at step $t = \pi_t =$

a mapping from states to action probabilities $\pi_t(a \mid s) =$ probability that $A_t = a$ when $S_t = s$

Special case - *deterministic policies*: $\pi_t(s)$ = the action taken with prob=1 when $S_t = s$

- Reinforcement learning methods specify how the agent changes its policy as a result of experience.
- Roughly, the agent's goal is to get as much reward as it can over the long run.

An Example Finite MDP

Recycling Robot

- At each step, robot has to decide whether it should (1) actively search for a can, (2) wait for someone to bring it a can, or (3) go to home base and recharge.
- Searching is better but runs down the battery; if runs out of power while searching, has to be rescued (which is bad).
- Decisions made on basis of current energy level: high, low.
- Reward = number of cans collected

Recycling Robot MDP

$$\begin{split} & S = \left\{ \texttt{high}, \texttt{low} \right\} \\ & \mathcal{A}(\texttt{high}) = \left\{ \texttt{search}, \texttt{wait} \right\} \\ & \mathcal{A}(\texttt{low}) = \left\{ \texttt{search}, \texttt{wait}, \texttt{recharge} \right\} \end{split}$$

 r_{search} = expected no. of cans while searching r_{wait} = expected no. of cans while waiting

 $r_{\rm search} > r_{\rm wait}$



The Markov Property

- By "the state" at step t, the book means whatever information is available to the agent at step t about its environment.
- The state can include immediate "sensations," highly processed sensations, and structures built up over time from sequences of sensations.
- Ideally, a state should summarize past sensations so as to retain all "essential" information, i.e., it should have the Markov Property:

$$\mathbf{Pr}\{R_{t+1} = r, S_{t+1} = s' \mid S_0, A_0, R_1, \dots, S_{t-1}, A_{t-1}, R_t, S_t, A_t\} = p(s', r \mid s, a) = \mathbf{Pr}\{R_{t+1} = r, S_{t+1} = s' \mid S_t, A_t\}$$

 \Box for all $s' \in S^+$, $r \in \mathcal{R}$, and all histories $S_0, A_0, R_1, ..., S_{t-1}, A_{t-1}, R_t, S_t, A_t$.

The Meaning of Life

(goals, rewards, and returns)

Goals and Rewards

- □ Is a scalar reward signal an adequate notion of a goal? maybe not, but it is surprisingly flexible.
- A goal should specify **what** we want to achieve, not **how** we want to achieve it.
- A goal must be outside the agent's direct control—thus outside the agent.
- □ The agent must be able to measure success:
 - explicitly;
 - frequently during its lifespan.

The reward hypothesis

- That all of what we mean by goals and purposes can be well thought of as the maximization of the cumulative sum of a received scalar signal (reward)
- □ A sort of *null hypothesis*.
 - Probably ultimately wrong, but so simple we have to disprove it before considering anything more complicated

Rewards and returns

- The objective in RL is to maximize long-term future reward
- That is, to choose A_t so as to maximize $R_{t+1}, R_{t+2}, R_{t+3}, \ldots$
- But what exactly should be maximized?
- The discounted return at time t: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \qquad \gamma \in [0, 1)$ the discount rate

γ	Reward sequence	Return
0.5(or any)	1000	
0.5	002000	
0.9	002000	
0.5	-12632000	

4 value functions

	state values	action values
prediction	v_{π}	q_{π}
control	v_*	q_*

- All theoretical objects, mathematical ideals (expected values)
- Distinct from their estimates:

 $V_t(s) = Q_t(s,a)$

Values are *expected* returns

• The value of a state, given a policy:

 $v_{\pi}(s) = \mathbb{E}\{G_t \mid S_t = s, A_{t:\infty} \sim \pi\} \qquad v_{\pi} : S \to \Re$

- The value of a state-action pair, given a policy: $q_{\pi}(s, a) = \mathbb{E}\{G_t \mid S_t = s, A_t = a, A_{t+1:\infty} \sim \pi\}$ $q_{\pi}: S \times \mathcal{A} \to \Re$
- The optimal value of a state:

$$v_*(s) = \max_{\pi} v_{\pi}(s) \qquad v_* : \mathcal{S} \to \Re$$

• The optimal value of a state-action pair:

$$q_*(s,a) = \max_{\pi} q_{\pi}(s,a) \qquad q_* : \mathcal{S} \times \mathcal{A} \to \Re$$

- Optimal policy: π_* is an optimal policy if and only if $\pi_*(a|s) > 0$ only where $q_*(s, a) = \max_b q_*(s, b) \quad \forall s \in S$
 - in other words, π_* is optimal iff it is greedy wrt q_*



What policy is optimal? A: left B: Right C: Other IF 8=0? IF X=.99 It &= 2?

Return

Suppose the sequence of rewards after step *t* is:

$$R_{t+1}, R_{t+2}, R_{t+3}, \dots$$

What do we want to maximize?

At least three cases, but in all of them, we seek to maximize the **expected return**, $E\{G_t\}$, on each step *t*.

- <u>Total reward</u>, G_t = sum of all future reward in the episode
- <u>Discounted reward</u>, G_t = sum of all future *discounted* reward
- <u>Average reward</u>, G_t = average reward per time step

Episodic tasks: interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze

In episodic tasks, we almost always use simple *total reward*:

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T,$$

where *T* is a final time step at which a **terminal state** is reached, ending an episode.

Continuing Tasks

Continuing tasks: interaction does not have natural episodes, but just goes on and on...

In this class, for continuing tasks we will always use *discounted return*:

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^{k} R_{t+k+1},$$

where $\gamma, 0 \le \gamma \le 1$, is the **discount rate**.

shortsighted $0 \leftarrow \gamma \rightarrow 1$ farsighted

Typically, $\gamma = 0.9$

An Example: Pole Balancing



As an **episodic task** where episode ends upon failure:

reward = +1 for each step before failure

 \Rightarrow return = number of steps before failure

As a continuing task with discounted return:

reward = -1 upon failure; 0 otherwise \Rightarrow return = $-\gamma^{k}$, for k steps before failure

In either case, return is maximized by avoiding failure for as long as possible.

Another Example: Mountain Car



Get to the top of the hill as quickly as possible.

reward = -1 for each step where **not** at top of hill \Rightarrow return = - number of steps before reaching top of hill

Return is maximized by minimizing number of steps to reach the top of the hill.

A Trick to Unify Notation for Returns

- □ In episodic tasks, we number the time steps of each episode starting from zero.
- □ We usually do not have to distinguish between episodes, so instead of writing $S_{t,j}$ for states in episode *j*, we write just S_t
- Think of each episode as ending in an absorbing state that always produces reward of zero:

$$(S_0 \xrightarrow{R_1 = +1} S_1 \xrightarrow{R_2 = +1} S_2 \xrightarrow{R_3 = +1} x_{k_3 = +1} \xrightarrow{R_4 = 0} R_5 = 0$$

$$\vdots$$
We can cover all cases by writing $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$

where γ can be 1 only if a zero reward absorbing state is always reached.

Value Functions

□ The value of a state is the expected return starting from that state; depends on the agent's policy:

State - value function for policy
$$\pi$$
:
 $v_{\pi}(s) = E_{\pi} \left\{ G_t \mid S_t = s \right\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right\}$

The value of an action (in a state) is the expected return starting after taking that action from that state; depends on the agent's policy:

Action - value function for policy
$$\pi$$
:
 $q_{\pi}(s,a) = E_{\pi} \left\{ G_t \mid S_t = s, A_t = a \right\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right\}$

Gridworld

- Actions: north, south, east, west; deterministic.
- □ If would take agent off the grid: no move but reward = -1
- Other actions produce reward = 0, except actions that move agent out of special states A and B as shown.



State-value function for equiprobable random policy; $\gamma = 0.9$

Bellman Equation for a Policy $\boldsymbol{\pi}$

The basic idea:

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

= $R_{t+1} + \gamma \left(R_{t+2} + \gamma R_{t+3} + \gamma^{2} R_{t+4} + \cdots \right)$
= $R_{t+1} + \gamma G_{t+1}$

So:

$$v_{\pi}(s) = E_{\pi} \{ G_t | S_t = s \}$$

$$= E_{\pi} \{ R_{t+1} + \gamma v_{\pi} (S_{t+1}) | S_t = s \}$$

Or, without the expectation operator:

$$v_{\pi}(s) = \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[r + \gamma v_{\pi}(s')\right]$$

More on the Bellman Equation

$$v_{\pi}(s) = \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[r + \gamma v_{\pi}(s')\right]$$

This is a set of equations (in fact, linear), one for each state. The value function for π is its unique solution.

Backup diagrams:





Input π , the policy to be evaluated Initialize an array V(s) = 0, for all $s \in S^+$ Repeat

$$\begin{array}{l} \Delta \leftarrow 0\\ \text{For each } s \in \mathbb{S}:\\ v \leftarrow V(s)\\ V(s) \leftarrow \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[r + \gamma V(s')\right]\\ \Delta \leftarrow \max(\Delta, |v - V(s)|)\\ \text{until } \Delta < \theta \text{ (a small positive number)}\\ \text{Output } V \approx v_{\pi} \end{array}$$

Gridworld

- Actions: north, south, east, west; deterministic.
- □ If would take agent off the grid: no move but reward = -1
- Other actions produce reward = 0, except actions that move agent out of special states A and B as shown.



State-value function for equiprobable random policy; $\gamma = 0.9$