

Regret Minimization and Posterior Thompson Sampling

Monday Morning Session, CDT Summer School

Outline

- Regret
- Intuition behind various learning algorithms
 - Multi-arm bandits
 - Contextual bandits
 - Markov decision process
- Some details of Bayesian posterior sampling
- Coding

Regret of a learning algorithm

$$\text{Regret} = \max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}]$$

Sum of rewards that
the optimal policy π can get

Sum of rewards that
the learning algorithm gets

So, regret is the excess reward that the algorithm misses compared to the best possible policy in a given class

Question: Does the optimal quantity depend on the learning algorithm? What other things does it depend on?

$$\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}]$$

Question: Does the optimal quantity depend on the learning algorithm? What other things does it depend on?

$$\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}]$$

The nature of underlying environment – this object does not depend on the learning algorithm.

Regret, optimal policy, and learning algorithms for two-arm bandits

Optimal policy: Multi-arm bandits (2-arms)

- Two potential outcomes $\{R_{t+1}(1), R_{t+1}(0)\}$ drawn from some distribution \mathcal{P} .
- At time t , take action A_t and receive outcome R_{t+1} .
- Consistency in observed outcomes, i.e., $R_{t+1} = R_{t+1}(A_t)$.
- Observed data: $(A_1, R_2), (A_2, R_3), (A_3, R_4), \dots, (A_T, R_{T+1})$

We assume stationarity across time



Optimal policy: Multi-arm bandits (2-arms)

- Two potential outcomes $\{R_{t+1}(1), R_{t+1}(0)\}$ drawn from some distribution \mathcal{P} .
- At time t , take action A_t and receive outcome R_{t+1} .
- Consistency in observed outcomes, i.e., $R_{t+1} = R_{t+1}(A_t)$.
- Observed data: $(A_1, R_2), (A_2, R_3), (A_3, R_4), \dots, (A_T, R_{T+1})$
- Regret minimization: Suffices to search over deterministic stationary policies
- Let $r(a) \triangleq \mathbb{E}[R_{j+1}(a)]$ be the stationary mean potential outcome, then

We assume stationarity across time



$$\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] = \sum_{j=1}^T r(a^*) \text{ where } a^* = \arg \max_a r(a)$$

Exercise: Prove!

Learning optimal policy: 2-arm bandits

- Regret = $\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}] = T \max(r(1), r(0)) - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}]$
- Regret would be minimized, if we knew the maximum reward or optimal policy

Learning optimal policy: 2-arm bandits

- Regret = $\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}] = T \max(r(1), r(0)) - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}]$
- Regret would be minimized, if we knew the maximum reward or optimal policy
- **Question:** Knowledge of what quantity suffices to learn the optimal policy?
 - $r(1)$ and $r(0)$?

Learning optimal policy: 2-arm bandits

- Regret = $\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}] = T \max(r(1), r(0)) - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}]$
- Regret would be minimized, if we knew the maximum reward or optimal policy
- **Question:** Knowledge of what quantity suffices to learn the optimal policy?
 - $r(1)$ and $r(0)$?
 - $r(1) - r(0)$?

Learning optimal policy: 2-arm bandits

- Regret = $\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}] = T \max(r(1), r(0)) - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}]$
- Regret would be minimized, if we knew the maximum reward or optimal policy
- **Question:** Knowledge of what quantity suffices to learn the optimal policy?
 - $r(1)$ and $r(0)$?
 - $r(1) - r(0)$?
 - $\text{sign}(r(1) - r(0))$?

Learning optimal policy: 2-arm bandits

- Regret = $\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}] = T \max(r(1), r(0)) - \mathbb{E}_{\pi^L}[\sum_{j=1}^T R_{j+1}]$
- Regret would be minimized, if we knew the maximum reward or optimal policy
- **Question:** Knowledge of what quantity suffices to learn the optimal policy?
 - $r(1)$ and $r(0)$?
 - $r(1) - r(0)$?
 - $\text{sign}(r(1) - r(0))$?
- Is there an advantage of one term over the other?

Typical bandit algorithm has two components (2-arm):

These components may be explicit or implicit

1. **Learning algorithm:**

- Uses the data so far to produce estimates $\hat{r}(1)$ and $\hat{r}(0)$ with uncertainty $\hat{\sigma}(1)$ and $\hat{\sigma}(0)$ for these quantities

2. **Optimization algorithm:**

- Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

Typical bandit algorithm has two components (2-arm):

These components may be explicit or implicit

1. **Learning algorithm:**

- Uses the data so far to produce estimates $\hat{r}(1)$ and $\hat{r}(0)$ with uncertainty $\hat{\sigma}(1)$ and $\hat{\sigma}(0)$ for these quantities
- where, *typically* $\hat{r}(a) \rightarrow r(a)$ and $\hat{\sigma}(a) \rightarrow 0$ as the number of pulls of arm $a \rightarrow \infty$

2. **Optimization algorithm:**

- Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

Typical bandit algorithm has two components (2-arm):

These components may be explicit or implicit

1. **Learning algorithm:**

- Uses the data so far to produce estimates $\hat{r}(1)$ and $\hat{r}(0)$ with uncertainty $\hat{\sigma}(1)$ and $\hat{\sigma}(0)$ for these quantities
- where, *typically* $\hat{r}(a) \rightarrow r(a)$ and $\hat{\sigma}(a) \rightarrow 0$ as the number of pulls of arm $a \rightarrow \infty$

2. **Optimization algorithm:**

- Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration
- Typically, the policy is a “good guess” of the optimal policy based on the data so far

Typical bandit algorithm has two components (2-arm):

These components may be explicit or implicit

1. **Learning algorithm:**

- Uses the data so far to produce estimates $\hat{r}(1)$ and $\hat{r}(0)$ with uncertainty $\hat{\sigma}(1)$ and $\hat{\sigma}(0)$ for these quantities
- where, *typically* $\hat{r}(a) \rightarrow r(a)$ and $\hat{\sigma}(a) \rightarrow 0$ as the number of pulls of arm $a \rightarrow \infty$

2. **Optimization algorithm:**

- Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration
- Typically, the policy is a “good guess” of the optimal policy based on the data so far
- The data collected by executing the learning policy is then fed into the learning algorithm to produce next set of estimates

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a)$$

Is this a good choice?

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

Is this a good choice?

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a)$$

This is the optimal policy if $\hat{r} = r$

What can go wrong when $\hat{r} \neq r$

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a)$$

This is the optimal policy if $\hat{r} = r$

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a) + c\hat{o}(a)$$

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a)$$

This is the optimal policy if $\hat{r} = r$

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a) + c\hat{\sigma}(a)$$

Rationale behind Upper Confidence Bound

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

➡ $\pi^L = \arg \max_a \hat{r}(a)$

This is the optimal policy if $\hat{r} = r$

➡ $\pi^L = \arg \max_a \hat{r}(a) + c\hat{\sigma}(a)$

Rationale behind Upper Confidence Bound

➡ $\pi^L = \begin{cases} \arg \max_a \hat{r}(a) & \text{with high probability} \\ \text{the other action} & \text{with low probability} \end{cases}$

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a)$$

This is the optimal policy if $\hat{r} = r$

$$\Rightarrow \pi^L = \arg \max_a \hat{r}(a) + c\hat{\sigma}(a)$$

Rationale behind Upper Confidence Bound

$$\Rightarrow \pi^L = \begin{cases} \arg \max_a \hat{r}(a) & \text{with high probability} \\ \text{the other action} & \text{with low probability} \end{cases}$$

Rationale behind ϵ -greedy

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

- ➡ $\pi^L = \arg \max_a \hat{r}(a)$ This is the optimal policy if $\hat{r} = r$
- ➡ $\pi^L = \arg \max_a \hat{r}(a) + c\hat{\sigma}(a)$ Rationale behind Upper Confidence Bound
- ➡ $\pi^L = \begin{cases} \arg \max_a \hat{r}(a) & \text{with high probability} \\ \text{the other action} & \text{with low probability} \end{cases}$ Rationale behind ϵ -greedy
- ➡ If $Z(a) \sim \mathcal{N}(\hat{r}(a), \hat{\sigma}(a))$ denotes a posterior on $r(a)$,
 $\pi^L = \begin{cases} 1 & \text{with probability } \mathbb{P}(Z(1) > Z(0)) \\ 0 & \text{otherwise.} \end{cases}$

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

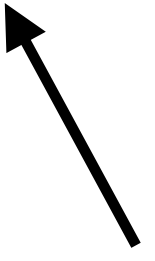
- ➡ $\pi^L = \arg \max_a \hat{r}(a)$ This is the optimal policy if $\hat{r} = r$
- ➡ $\pi^L = \arg \max_a \hat{r}(a) + c\hat{\sigma}(a)$ Rationale behind Upper Confidence Bound
- ➡ $\pi^L = \begin{cases} \arg \max_a \hat{r}(a) & \text{with high probability} \\ \text{the other action} & \text{with low probability} \end{cases}$ Rationale behind ϵ -greedy
- ➡ If $Z(a) \sim \mathcal{N}(\hat{r}(a), \hat{\sigma}(a))$ denotes a posterior on $r(a)$,
 $\pi^L = \begin{cases} 1 & \text{with probability } \mathbb{P}(Z(1) > Z(0)) \\ 0 & \text{otherwise.} \end{cases}$ Rationale behind Posterior/Thompson sampling

Discussion questions

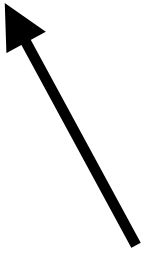
- What would happen to the learning policies in the bandit algorithms if one arm has significantly higher mean reward than the other? Which variant of algorithms would you prefer if the two arms have similar means?
- There are two possible ways to implement posterior sampling:
 - Sample action 1 with probability $p \triangleq \mathbb{P} \left[\mathcal{N}(\hat{r}(1) - \hat{r}(0), \hat{\sigma}_1^2 + \hat{\sigma}_2^2) > 0 \right]$ and action 0 with $1 - p$.
 - Sample two random variables $Z(a) \sim \mathcal{N}(\hat{r}(a), \hat{\sigma}^2(a))$ and assign $\arg \max_a Z(a)$
 - Are both approaches correct? Is there a benefit of one approach over the other?
 - **Hint:** In which case do we have explicit probabilities of randomization? What happens when we don't have Gaussian posterior?

Same ideas applied to
contextual bandits

Contextual bandits

- State, two potential outcomes $\{S_t, R_{t+1}(1), R_{t+1}(0)\}$ drawn from some distribution \mathcal{P}
 - At time t , take action A_t and receive outcome R_{t+1}
 - Consistency in observed outcomes, i.e., $R_{t+1} = R_{t+1}(A_t)$
 - Observed data: $(S_1, A_1, R_2), (S_2, A_2, R_3), (S_3, A_3, R_4), \dots, (S_T, A_T, R_{T+1})$
- At each time, a state/context variable is observed
(still assuming stationarity across time points)**
- 

Contextual bandits

- State, two potential outcomes $\{S_t, R_{t+1}(1), R_{t+1}(0)\}$ drawn from some distribution \mathcal{P}
 - At time t , take action A_t and receive outcome R_{t+1}
 - Consistency in observed outcomes, i.e., $R_{t+1} = R_{t+1}(A_t)$
 - Observed data: $(S_1, A_1, R_2), (S_2, A_2, R_3), (S_3, A_3, R_4), \dots, (S_T, A_T, R_{T+1})$
 - Regret minimization: Would like to use the states, so the policy class is $\{\pi : \mathcal{S} \rightarrow \Delta^{|A|-1}\}$
 - Let $r(s, a) \triangleq \mathbb{E}[R_{j+1}(a) | S_j = s]$ be the stationary mean reward outcome, then
- At each time, a state/context variable is observed (still assuming stationarity across time points)**
- 

$$\max_{\pi} \mathbb{E}_{\pi}[\sum_{j=1}^T R_{j+1}] = \mathbb{E}_{\pi^*}[\sum_{j=1}^T r(S_j, \pi^*(S_j))] \text{ where } \pi^*(s) = \arg \max_a r(s, a)$$

Exercise: Prove!

Typical “contextual” bandit algorithm has two components (2-arm): These components may be explicit or implicit

1. **Learning algorithm:**

- Uses the data so far to produce estimates $\hat{r}(s,1)$ and $\hat{r}(s,0)$ with uncertainty $\hat{\sigma}(s,1)$ and $\hat{\sigma}(s,0)$ for these quantities
- where, typically $\hat{r}(s, a) \rightarrow r(s, a)$ and $\hat{\sigma}(s, a) \rightarrow 0$ as the number of pulls of arm a under state $s \rightarrow \infty$

2. **Optimization algorithm:**

- Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration
- Typically, the policy is a “good guess” of the optimal policy based on the data so far
- The data collected by executing the learning policy is then fed into the learning algorithm to produce next set of estimates

Typical “contextual” bandit algorithm has two components (2-arm): These components may be explicit or implicit

1. Learning algorithm:

Typical approach: Model
 $r(s, a) = \beta_a^\top \phi(s, a)$ and estimate β_a

- Uses the data so far to produce estimates $\hat{r}(s, 1)$ and $\hat{r}(s, 0)$ with uncertainty $\hat{\sigma}(s, 1)$ and $\hat{\sigma}(s, 0)$ for these quantities
- where, *typically* $\hat{r}(s, a) \rightarrow r(s, a)$ and $\hat{\sigma}(s, a) \rightarrow 0$ as the number of pulls of arm a under state $s \rightarrow \infty$

2. Optimization algorithm:

- Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration
- Typically, the policy is a “good guess” of the optimal policy based on the data so far
- The data collected by executing the learning policy is then fed into the learning algorithm to produce next set of estimates

Different principled heuristics in the optimization algorithm lead to different algorithms

Optimization algorithm:

Uses the estimates to construct a “learning” policy that pulls the arm in the next iteration

- ➡ $\pi^L(s) = \arg \max_a \hat{r}(s, a)$ This is the optimal policy if $\hat{r} = r$
- ➡ $\pi^L(s) = \arg \max_a \hat{r}(s, a) + c\hat{\sigma}(s, a)$ Rationale behind Upper Confidence Bound
- ➡ $\pi^L = \begin{cases} \arg \max_a \hat{r}(s, a) & \text{with high probability} \\ \text{the other action} & \text{with low probability} \end{cases}$ Rationale behind ϵ -greedy
- ➡ If $Z(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(a)$,
 $\pi^L = \begin{cases} 1 & \text{with probability } \mathbb{P}(Z(s, 1) > Z(s, 0)) \\ 0 & \text{otherwise.} \end{cases}$ Rationale behind Posterior/Thompson sampling

Moving on to MDPs

Markov Decision process

- Sequence of states $\{S_0, S_1(1), S_1(0), S_2(1), S_2(0), \dots\}$ drawn from some distribution \mathcal{P} , satisfying Markovian property

Markov Decision process

- Sequence of states $\{S_0, S_1(1), S_1(0), S_2(1), S_2(0), \dots\}$ drawn from some distribution \mathcal{P} , satisfying Markovian property **Assuming stationarity in the transitions between states**

Markov Decision process

- Sequence of states $\{S_0, S_1(1), S_1(0), S_2(1), S_2(0), \dots\}$ drawn from some distribution \mathcal{P} , satisfying Markovian property **Assuming stationarity in the transitions between states**
- At time $t = 0, 1, \dots$, take action A_t , observe S_{t+1} , and receive outcome R_{t+1} .
- Consistency in observed states, i.e., $S_{t+1} = S_{t+1}(A_t)$ and reward $R_{t+1} = h(S_{t+1}) = h(S_{t+1}(A_t))$.
- Observed data: $(S_1, A_1, R_2), (S_2, A_2, R_3), (S_3, A_3, R_4), \dots$,

Markov Decision process

- Sequence of states $\{S_0, S_1(1), S_1(0), S_2(1), S_2(0), \dots\}$ drawn from some distribution \mathcal{P} , satisfying Markovian property

Assuming stationarity in the transitions between states

- At time $t = 0, 1, \dots$, take action A_t , observe S_{t+1} , and receive outcome R_{t+1} .
- Consistency in observed states, i.e., $S_{t+1} = S_{t+1}(A_t)$ and reward $R_{t+1} = h(S_{t+1}) = h(S_{t+1}(A_t))$.
- Observed data: $(S_1, A_1, R_2), (S_2, A_2, R_3), (S_3, A_3, R_4), \dots$,
- Now the benchmark is discounted sum of rewards $\mathbb{E}_\pi[\sum_{j=0}^{\infty} \gamma^j R_{j+1}]$
- Would like to use the states, so the policy class is $\{\pi : \mathcal{S} \rightarrow \Delta^{|A|-1}\}$.
 - Here the policy class is restricted based on the fact that under stationary transitions, the optimal policy is stationary

Unpacking the optimal policy for Markov decision process

- Let $r(s, a) \triangleq \mathbb{E}[R_{j+1} | S_j = s, A_j = a] = \mathbb{E}[h(S_{j+1}(a)) | S_j = s]$

- Suppose initial state is deterministically s , then

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi[\sum_{j=0}^{\infty} \gamma^j R_{j+1} | S_0 = s] = \sum_a \pi(a | s) \left[r(s, a) + \underbrace{\gamma \sum_{j=1}^{\infty} \mathbb{E}[\gamma^{j-1} R_{j+1} | S_0 = s, A_0 = a]}_{\triangleq H^\pi(s, a)} \right] \\ &= \sum_a \pi(a | s) \left[r(s, a) + \gamma H^\pi(s, a) \right] \\ &= \sum_a \pi(a | s) \left[Q^\pi(s, a) \right] \end{aligned}$$

- Hence the optimal policy $\pi^\star = \arg \max V^\pi$ also satisfies

$$\pi^\star(a | s) = \arg \max_a Q^{\pi^\star}(s, a) = \arg \max_a r(s, a) + H^{\pi^\star}(s, a)$$

Exercise: Prove!

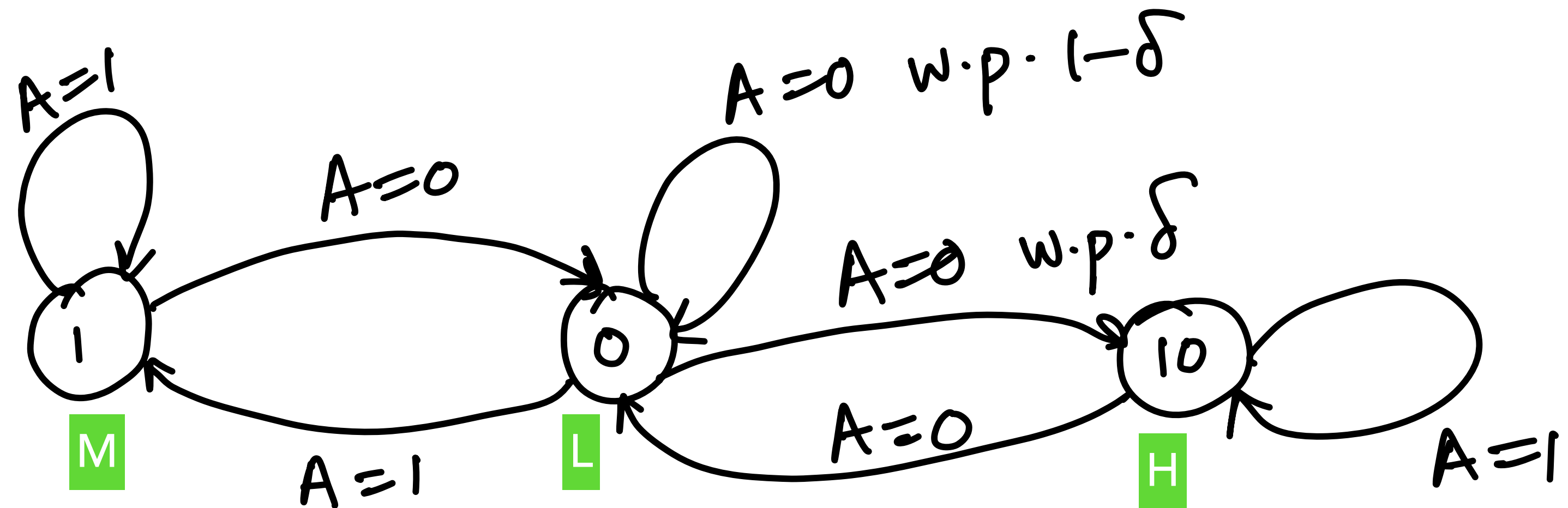
- **Typical goal in MDP:** Find this optimal policy!

Think what goes wrong if we just solve $\arg \max_a r(s, a)$?

Thinking about delayed effects:

What if we just maximize immediate rewards?

- Consider 3 state (state space L, M, H) MDP with 2 actions (0, 1)
- The reward in the states is deterministic and equal to 1, 0, and 10 in the M, L, and H state.
- The transitions are marked by probability (some of them are deterministic)



- What happens for this problem:

$$\arg \max_a r(s, a) \quad \text{vs} \quad \arg \max_a r(s, a) + H^*(s, a)$$

Exercises

- What is the mean reward function $r(s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a] = \mathbb{E}[R_{t+1}(a) | S_t = s]$?
- Derive the policy $\arg \max_a r(s, a)$ and show that it is sub-optimal.
- Let $\delta = 0.5$. Derive the optimal value function, optimal Q function, H function, and the optimal policy (which is take action 0 in state L, M and action 1 in state H).
- You can use Bellman value iterations to obtain optimal value function as well. The code is included in `value_iteration.py`

Key ideas behind Q-learning

- **Question:** Suppose someone gave us the values of Q^{π^*} but not π^* . How do we recover π^* ?

Key ideas behind Q-learning

- **Question:** Suppose someone gave us the values of Q^{π^*} but not π^* . How do we recover π^* ?
 - $\pi^*(a | s) = \arg \max_a Q^{\pi^*}(s, a)$

Key ideas behind Q-learning

- **Question:** Suppose someone gave us the values of Q^{π^*} but not π^* . How do we recover π^* ?
 - $\pi^*(a | s) = \arg \max_a Q^{\pi^*}(s, a)$
- **Question:** What if we had an estimate \hat{Q}^{π^*} with uncertainty $\hat{\sigma}$?

One class of RL algorithms for MDPs

1. Learning algorithm:

- Uses the data so far to produce estimates $\hat{Q}^{\pi^*}r(s,1)$ and $\hat{Q}^{\pi^*}r(s,0)$ with uncertainty $\hat{\sigma}(s,1)$ and $\hat{\sigma}(s,0)$ for these quantities
- where, *typically* $\hat{Q}^{\pi^*}(s,a) \rightarrow Q^{\pi^*}(s,a)$ and $\hat{\sigma}(s,a) \rightarrow 0$ as the number of pulls of action a under state $s \rightarrow \infty$

2. Optimization algorithm:

- Uses the estimates to construct a “learning” policy that samples the action for the next iteration
- $\pi^L(s) = \arg \max_a \hat{Q}^{\pi^*}(s,a) + c\hat{\sigma}(s,a)$
- If $Z(s,a) \sim \mathcal{N}(\hat{Q}^{\pi^*}(s,a), \hat{\sigma}(s,a))$ denotes a posterior on $Q^{\pi^*}(s,a)$,
$$\pi^L(s) = \begin{cases} 1 & \text{with probability } \mathbb{P}(Z(s,1) > Z(s,0)) \\ 0 & \text{otherwise.} \end{cases}$$

A simpler RL algorithm class with a different target

- Recall $\pi^\star(a | s) = \arg \max_a Q^{\pi^\star}(s, a) = \arg \max_a r(s, a) + H^{\pi^\star}(s, a)$. Suppose someone gave us H^{π^\star} but not r . How do we estimate π^\star ?

A simpler RL algorithm class with a different target

- Recall $\pi^\star(a | s) = \arg \max_a Q^{\pi^\star}(s, a) = \arg \max_a r(s, a) + H^{\pi^\star}(s, a)$. Suppose someone gave us H^{π^\star} but not r . How do we estimate π^\star ?
- **Learning algorithm:** Produce estimates $\hat{r}(s,1)$ and $\hat{r}(s,0)$ with uncertainty $\hat{\sigma}(s,1)$ and $\hat{\sigma}(s,0)$.

A simpler RL algorithm class with a different target

- Recall $\pi^\star(a | s) = \arg \max_a Q^{\pi^\star}(s, a) = \arg \max_a r(s, a) + H^{\pi^\star}(s, a)$. Suppose someone gave us H^{π^\star} but not r . How do we estimate π^\star ?
- **Learning algorithm:** Produce estimates $\hat{r}(s,1)$ and $\hat{r}(s,0)$ with uncertainty $\hat{\sigma}(s,1)$ and $\hat{\sigma}(s,0)$.
- **Optimization algorithm:** Uses the estimates to construct a policy.
 - $\pi^L = \arg \max_a \hat{r}(s, a) + c\hat{\sigma}(s, a) + \gamma H^{\pi^\star}(s, a)$
 - If $Z'(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,
$$\pi^L = \begin{cases} 1 & \text{with probability } \mathbb{P} \left[Z'(s,1) - Z'(s,0) > -\gamma [H^{\pi^\star}(s,1) - H^{\pi^\star}(s,0)] \right] \\ 0 & \text{otherwise.} \end{cases}$$

A simpler RL algorithm class with a different target

- Recall $\pi^\star(a | s) = \arg \max_a Q^{\pi^\star}(s, a) = \arg \max_a r(s, a) + H^{\pi^\star}(s, a)$. Suppose someone gave us H^{π^\star} but not r . How do we estimate π^\star ?
- **Learning algorithm:** Produce estimates $\hat{r}(s,1)$ and $\hat{r}(s,0)$ with uncertainty $\hat{\sigma}(s,1)$ and $\hat{\sigma}(s,0)$.
- **Optimization algorithm:** Uses the estimates to construct a policy.
 - $\pi^L = \arg \max_a \hat{r}(s, a) + c\hat{\sigma}(s, a) + \gamma H^{\pi^\star}(s, a)$
 - If $Z'(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,
$$\pi^L = \begin{cases} 1 & \text{with probability } \mathbb{P} \left[Z'(s,1) - Z'(s,0) > -\gamma [H^{\pi^\star}(s,1) - H^{\pi^\star}(s,0)] \right] \\ 0 & \text{otherwise.} \end{cases}$$

Discuss: Why is this a simpler algorithm class?

Steps towards the HeartSteps RL algorithm

Action 1 = **Send** a notification vs Action 0 = **Do nothing**

If $Z'(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,

$$\pi^L(s) = \begin{cases} 1 & \text{with probability } \mathbb{P} \left[Z'(s, 1) - Z'(s, 0) > -\gamma [H^{\pi^*}(s, 1) - H^{\pi^*}(s, 0)] \right] \\ 0 & \text{otherwise.} \end{cases}$$

HeartSteps further clips the probabilities between 0.2 and 0.8

Discuss:

Why does this make intuitive sense?
What are the quantities on the LHS and RHS?
Why is this a simpler algorithm class compared to estimating Q^{π^*} ?

This is the HeartSteps RL algorithm

Action 1 = **Send** a notification vs Action 0 = **Do nothing**

If $Z'(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,

$$\pi^L(s) = \begin{cases} 1 & \text{with probability } \mathbb{P} \left[Z'(s, 1) - Z'(s, 0) > -\gamma [H^{\pi^*}(s, 1) - H^{\pi^*}(s, 0)] \right] \\ 0 & \text{otherwise.} \end{cases}$$

LHS is like the immediate treatment effect of sending a notification

RHS is like the delayed effects (denoted as $\eta(s)$)

This is the HeartSteps RL algorithm

Action 1 = **Send** a notification vs Action 0 = **Do nothing**

If $Z'(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,

$$\pi^L(s) = \begin{cases} 1 & \text{with probability } \mathbb{P} \left[Z'(s, 1) - Z'(s, 0) > -\gamma [H^{\pi^*}(s, 1) - H^{\pi^*}(s, 0)] \right] \\ 0 & \text{otherwise.} \end{cases}$$

LHS is like the immediate treatment effect of sending a notification

RHS is like the delayed effects (denoted as $\eta(s)$)

- HeartSteps V2/V3 RL algorithm does posterior sampling for LHS
- Estimates the RHS from HeartSteps V1 data
- Then clips the probabilities between 0.2 and 0.8

Bayesian Thompson Sampling with Gaussian linear model

Bayesian Thompson sampling for 2-arm bandits

- Parameters unknown, put a Gaussian prior on them
- Gaussian prior $\mathcal{N}(\bar{\alpha}_a, \bar{\sigma}_a^2)$ on $r(a)$ + Gaussian likelihood on $\mathcal{N}(r(a), \sigma^2)$ on $R_t(a)$, σ known \rightarrow Gaussian posterior $\mathcal{N}(\tilde{\alpha}_a, \tilde{\sigma}_a^2)$ on $r(a)$
- When we observe (A_1, R_2) . What are the posterior parameters?

$$\bullet \tilde{\sigma}_a^2 = \left(\frac{1}{\bar{\sigma}_a^2} + \frac{\mathbf{1}(A_1 = a)}{\sigma^2} \right)^{-1} \text{ and } \tilde{\alpha}_a = \tilde{\sigma}_a^2 \left(\frac{\alpha_a}{\bar{\sigma}_a^2} + \frac{\mathbf{1}(A_1 = a)R_2}{\sigma^2} \right)$$

- Now what happens after observing T tuples? Can use induction with – posterior being the new prior

$$\bullet \tilde{\sigma}_a^2 = \left(\frac{1}{\bar{\sigma}_a^2} + \frac{\sum_{j=1}^T \mathbf{1}(A_j = a)}{\sigma^2} \right)^{-1} \text{ and } \tilde{\alpha}_a = \tilde{\sigma}_a^2 \left(\frac{\alpha_a}{\bar{\sigma}_a^2} + \frac{\sum_{j=1}^T \mathbf{1}(A_j = a)R_{j+1}}{\sigma^2} \right)$$

Bayesian Thompson sampling with linear model

- We didn't discuss so far, how the estimate of \hat{r} is generated
- Model: $r(s, a) = \alpha_a^\top s$ and $R_{t+1}(a) | S_t = s \sim \mathcal{N}(r(s, a), \sigma^2)$ with σ known
 - can choose non-linear model, like $\alpha_a^\top \phi(s, a)$,
 - can also do posterior updates for σ
 - or non-Gaussian noise in reward
- Or we can also model: $r(s, a) = \alpha^\top g(s) + a\beta^\top f(s) \rightarrow$ HeartSteps used this

Bayesian Thompson sampling with linear model for contextual bandits

- Suppose we model $r(s, a) = \alpha_a^\top \phi(s, a)$
- We put a prior $\mathcal{N}(\bar{\alpha}_a, \bar{\Sigma}_a)$ on $\alpha_a \rightarrow \alpha_a^\top \phi(s, a) \sim \mathcal{N}(\bar{\alpha}_a^\top \phi(s, a), \phi(s, a)^\top \bar{\Sigma}_a \phi(s, a))$
- Observe (A_1, S_1, R_2) . What are the posteriors?

$$\bullet \tilde{\Sigma}_a = \left(\bar{\Sigma}_a^{-1} + \frac{\mathbf{1}(A_1 = a) \phi(S_1, A_1) \phi(S_1, A_1)^\top}{\sigma^2} \right)^{-1} \text{ and } \tilde{\alpha}_a = \tilde{\Sigma}_a \left(\bar{\Sigma}_a^{-1} \alpha_a + \frac{\mathbf{1}(A_1 = a) \phi(S_1, A_1) R_2}{\sigma^2} \right)$$

- Now what happens after observing T tuples? Can use induction with – posterior being the new prior

$$\bullet \tilde{\Sigma}_a = \left(\bar{\Sigma}_a^{-1} + \frac{\sum_{j=1}^T \mathbf{1}(A_j = a) \phi(S_j, A_j) \phi(S_j, A_j)^\top}{\sigma^2} \right)^{-1} \text{ and } \tilde{\alpha}_a = \tilde{\Sigma}_a \left(\bar{\Sigma}_a^{-1} \alpha_a + \frac{\sum_{j=1}^T \mathbf{1}(A_j = a) \phi(S_j, A_j) R_{j+1}}{\sigma^2} \right)$$

Posterior sampling

- If $Z(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,
$$\pi^L(s) = \begin{cases} 1 & \text{with probability } \mathbb{P}(Z(s,1) > Z(s,0)) \\ 0 & \text{otherwise.} \end{cases}$$
- Sample action 1 with probability $\mathbb{P} \left[\mathcal{N}((\tilde{\alpha}_1 - \tilde{\alpha}_0)^\top \phi(s, a), \tilde{\Sigma}_1 + \tilde{\Sigma}_2) > 0 \right]$ – can compute this probability exactly – closed form for Gaussian case

Posterior sampling

- If $Z(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,
$$\pi^L(s) = \begin{cases} 1 & \text{with probability } \mathbb{P}(Z(s, 1) > Z(s, 0)) \\ 0 & \text{otherwise.} \end{cases}$$
- Sample action 1 with probability $\mathbb{P} \left[\mathcal{N}((\tilde{\alpha}_1 - \tilde{\alpha}_0)^\top \phi(s, a), \tilde{\Sigma}_1 + \tilde{\Sigma}_2) > 0 \right]$ – can compute this probability exactly – closed form for Gaussian case
 - Or if we model the treatment effect as $\beta^\top f(s)$, then
Sample action 1 with probability $\mathbb{P} \left[\mathcal{N}(\tilde{\beta}^\top f(s), \tilde{\Sigma}_\beta) > 0 \right]$

Posterior sampling v1

- If $Z(s, a) \sim \mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,
$$\pi^L(s) = \begin{cases} 1 & \text{with probability } \mathbb{P}(Z(s, 1) > Z(s, 0)) \\ 0 & \text{otherwise.} \end{cases}$$
- Sample action 1 with probability $\mathbb{P} \left[\mathcal{N}((\tilde{\alpha}_1 - \tilde{\alpha}_0)^\top \phi(s, a), \tilde{\Sigma}_1 + \tilde{\Sigma}_2) > 0 \right]$ – can compute this probability exactly – closed form for Gaussian case
 - Or if we model the treatment effect as $\beta^\top f(s)$, then
Sample action 1 with probability $\mathbb{P} \left[\mathcal{N}(\tilde{\beta}^\top f(s), \tilde{\Sigma}_\beta) > 0 \right]$
 - When we also have delayed effects proxy $\eta(s)$:
Sample action 1 with probability $\mathbb{P} \left[\mathcal{N}(\tilde{\beta}^\top f(s), \tilde{\Sigma}_\beta) > \eta(s) \right]$

Posterior sampling v2

- If $\mathcal{N}(\hat{r}(s, a), \hat{\sigma}(s, a))$ denotes a posterior on $r(s, a)$,
$$\pi^L = \begin{cases} 1 & \text{with probability } \mathbb{P}(r(s, 1) > r(s, 0)) \\ 0 & \text{otherwise.} \end{cases}$$
- Sample α_1, α_0 from the posterior and assign $\text{action} = \arg \max_a \alpha_a^\top \phi(s, a)$
 - Or if we model the treatment effect: Sample β from the posterior and assign $a = \mathbf{1}(\beta^\top f(s) > 0)$.
 - Or if we model the treatment effect with delayed effects proxy as $\eta(s)$: Sample β from the posterior and assign $a = \mathbf{1}(\beta^\top f(s) > \eta(s))$.

Discuss

- Sample action 1 with probability $\mathbb{P} \left[\mathcal{N}((\tilde{\alpha}_1 - \tilde{\alpha}_0)^\top \phi(s, a), \tilde{\Sigma}_1 + \tilde{\Sigma}_2) > 0 \right]$
- Sample α_1, α_0 from the posterior and assign action $= \arg \max_a \alpha_a^\top \phi(s, a)$
 - Which of these two approaches are correct?
 - What are the pros and cons of the two approaches?
 - **Hint:** What happens when we don't have Gaussian posterior?

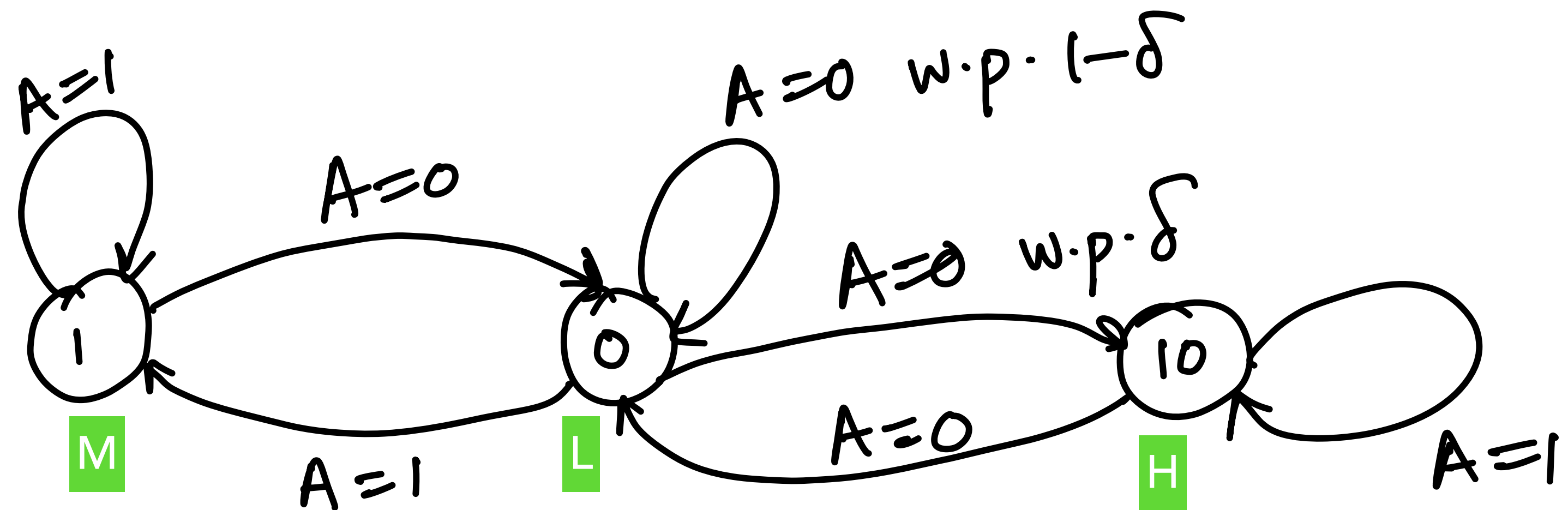
Coding exercises overview

- **Exercise 1:** 2-armed bandits
 - four algorithms: UCB, eps-greedy, BTS (2 variants)
 - Code up the two variants of BTS
- **Exercise 2:** 2-armed linear contextual bandits
 - Code up the two variants of BTS
- **Exercise 3:** 2-armed tabular MDP [builds on our example from earlier; see next slide]
 - Code up two variants of BTS with no delayed effect estimate
 - Code up two variants of BTS with some delayed effect estimate

Thinking about delayed effects:

What if we just maximize immediate rewards?

- Consider 3 state (state space L, M, H) MDP with 2 actions (0, 1)
- The reward in the states is deterministic and equal to 1, 0, and 10 in the M, L, and H state.
- The transitions are marked by probability (some of them are deterministic)



• What happens for this problem:

$$\arg \max_a r(s, a) \quad \text{vs} \quad \arg \max_a r(s, a) + H(s, a)^*$$

We have given you the optimal V and H function in the code for this exercise and we implement the two policies above (with a Bayesian Thompson sampling algorithm for the mean rewards)