Bandits: Thompson Sampling and Gittins Index

Lucas Janson and Sham Kakade

CS/Stat 184: Introduction to Reinforcement Learning Fall 2022

1

Today

- Feedback from last lecture
- Recap
- Thompson sampling
- Gittins index

Feedback from feedback forms

- 1. Thank you to everyone who filled out the forms!
- 2. Bit of confusion about Bayesian bandits
 - I'll review briefly, but section this week and HW 1 will give examples

Today



- Recap
- Thompson sampling
- Gittins index

Recap

Recap

• Algorithms we've seen so far: pure greedy, pure exploration, ETC, ε -greedy, and UCB

Recap

- Algorithms we've seen so far: pure greedy, pure exploration, ETC, ε -greedy, and UCB
- Bayesian bandit augments bandit environment with a prior distribution, allowing arm means to be treated as random with known distribution

Today



- Thompson sampling
- Gittins index

A Bayesian bandit augments the bandit environment we've been working in so far with a prior distribution on the unknown reward distributions: $\pi(\nu^{(1)}, ..., \nu^{(K)})$

A Bayesian bandit augments the bandit environment we've been working in so far with a prior distribution on the unknown reward distributions: $\pi(\nu^{(1)}, ..., \nu^{(K)})$

Bayes rule at time step *t* gives us a distribution (called the posterior distribution) $\mathbb{P}(\boldsymbol{\mu} \mid r_0, a_0, r_1, a_1, \dots, r_{t-1}, a_{t-1})$

that exactly characterizes our uncertainty about μ .

A Bayesian bandit augments the bandit environment we've been working in so far with a prior distribution on the unknown reward distributions: $\pi(\nu^{(1)}, ..., \nu^{(K)})$

Bayes rule at time step *t* gives us a distribution (called the posterior distribution) $\mathbb{P}(\boldsymbol{\mu} \mid r_0, a_0, r_1, a_1, \dots, r_{t-1}, a_{t-1})$

that exactly characterizes our uncertainty about μ .

Note that although we are now treating μ as random, we still assume its value is only drawn once (from the prior) and then stays the same throughout *t*

A Bayesian bandit augments the bandit environment we've been working in so far with a prior distribution on the unknown reward distributions: $\pi(\nu^{(1)}, ..., \nu^{(K)})$

Bayes rule at time step *t* gives us a distribution (called the posterior distribution) $\mathbb{P}(\boldsymbol{\mu} \mid r_0, a_0, r_1, a_1, \dots, r_{t-1}, a_{t-1})$

that exactly characterizes our uncertainty about μ .

Note that although we are now treating μ as random, we still assume its value is only drawn once (from the prior) and then stays the same throughout *t*

What changes with *t* is our information about μ , i.e., the posterior distribution, as we collect more and more data by pulling arms via a bandit algorithm

<u>Bayesian Bernoulli bandit</u> with <u>uniform prior</u> on μ gives a running posterior on the mean of each arm k that is Beta $(1 + \#\{\text{arm } k \text{ successes}\}, 1 + \#\{\text{arm } k \text{ failures}\})$

Bayesian Bernoulli bandit with uniform prior on μ gives a running posterior on the mean of each arm k that is Beta $(1 + #\{arm k successes\}, 1 + #\{arm k failures\})$ (derived by Bayes rule and some algebra; see this week's section or HW1 for more details)

Bayesian Bernoulli bandit with uniform prior on μ gives a running posterior on the mean of each arm k that is Beta $(1 + \#\{\operatorname{arm} k \operatorname{successes}\}, 1 + \#\{\operatorname{arm} k \operatorname{failures}\})$ (derived by Bayes rule and some algebra; see this week's section or HW1 for more details)

Beta (α_k, β_k) has mean (posterior mean = what we expect $\mu^{(k)}$ to be): $\frac{\alpha_k}{\alpha_k + \beta_k} = \frac{1 + \#\{\operatorname{arm} k \text{ successes}\}}{2 + \#\{\operatorname{arm} k \text{ pulls}\}}$

which starts at 1/2 and approaches the sample mean of arm k with more pulls.

Bayesian Bernoulli bandit with uniform prior on μ gives a running posterior on the mean of each arm k that is Beta $(1 + \#\{\operatorname{arm} k \operatorname{successes}\}, 1 + \#\{\operatorname{arm} k \operatorname{failures}\})$ (derived by Bayes rule and some algebra; see this week's section or HW1 for more details)

Beta (α_k, β_k) has mean (posterior mean = what we expect $\mu^{(k)}$ to be): $\frac{\alpha_k}{\alpha_k + \beta_k} = \frac{1 + \#\{\operatorname{arm} k \text{ successes}\}}{2 + \#\{\operatorname{arm} k \text{ pulls}\}}$

which starts at 1/2 and approaches the sample mean of arm k with more pulls.

Beta(α_k, β_k) has variance (posterior variance \approx how uncertain we are about $\mu^{(k)}$):

$$\frac{\alpha_k}{\alpha_k + \beta_k} \times \frac{\beta_k}{\alpha_k + \beta_k} \times \frac{1}{\alpha_k + \beta_k + 1}$$

which decreases at a rate of roughly $1/\#\{\operatorname{arm} k \text{ pulls}\}$

Bayesian bandit environment means that at every time step, we know the distribution of the arm reward distributions conditioned on everything we've seen so far

Bayesian bandit environment means that at every time step, we know the distribution of the arm reward distributions conditioned on everything we've seen so far In particular, we know the exact probability, given everything we've seen so far, that each arm is the true optimal arm, i.e.,

 $\forall k$, we know $\mathbb{P}(k = k^* \mid r_0, a_0, ..., r_{t-1}, a_{t-1})$

Bayesian bandit environment means that at every time step, we know the distribution of the arm reward distributions conditioned on everything we've seen so far In particular, we know the exact probability, given everything we've seen so far, that each arm is the true optimal arm, i.e.,

 $\forall k$, we know $\mathbb{P}(k = k^* \mid r_0, a_0, ..., r_{t-1}, a_{t-1})$

Thompson sampling: sample from this distribution to determine next arm to pull

Bayesian bandit environment means that at every time step, we know the distribution of the arm reward distributions conditioned on everything we've seen so far In particular, we know the exact probability, given everything we've seen so far, that each arm is the true optimal arm, i.e.,

 $\forall k$, we know $\mathbb{P}(k = k^* \mid r_0, a_0, ..., r_{t-1}, a_{t-1})$

Thompson sampling: sample from this distribution to determine next arm to pull

For t = 0, ..., T - 1:

 $a_t \sim \text{distribution of } k^{\star} \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Bayesian bandit environment means that at every time step, we know the distribution of the arm reward distributions conditioned on everything we've seen so far In particular, we know the exact probability, given everything we've seen so far, that each arm is the true optimal arm, i.e.,

 $\forall k$, we know $\mathbb{P}(k = k^* \mid r_0, a_0, ..., r_{t-1}, a_{t-1})$

Thompson sampling: sample from this distribution to determine next arm to pull

For t = 0, ..., T - 1:

$$a_t \sim \text{distribution of } k^{\star} \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$$

(In practice, usually draw a sample $\mu_t \sim \text{distribution of } \mu \mid r_0, a_0, \dots, r_{t-1}, a_{t-1} \text{ and then compute}$ $a_t = \arg \max_k \mu_t^{(k)}$, which is the same thing as $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$)

Bayesian bandit environment means that at every time step, we know the distribution of the arm reward distributions conditioned on everything we've seen so far In particular, we know the exact probability, given everything we've seen so far, that each arm is the true optimal arm, i.e.,

 $\forall k$, we know $\mathbb{P}(k = k^* \mid r_0, a_0, ..., r_{t-1}, a_{t-1})$

Thompson sampling: sample from this distribution to determine next arm to pull

For t = 0, ..., T - 1:

$$a_t \sim \text{distribution of } k^{\star} \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$$

(In practice, usually draw a sample $\mu_t \sim \text{distribution of } \mu \mid r_0, a_0, \dots, r_{t-1}, a_{t-1} \text{ and then compute}$ $a_t = \arg \max_k \mu_t^{(k)}$, which is the same thing as $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$)

That's it! Statistically, this is a super simple and elegant algorithm

Bayesian bandit environment means that at every time step, we know the distribution of the arm reward distributions conditioned on everything we've seen so far In particular, we know the exact probability, given everything we've seen so far, that each arm is the true optimal arm, i.e.,

 $\forall k$, we know $\mathbb{P}(k = k^* \mid r_0, a_0, ..., r_{t-1}, a_{t-1})$

Thompson sampling: sample from this distribution to determine next arm to pull

For t = 0, ..., T - 1:

$$a_t \sim \text{distribution of } k^{\star} \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$$

(In practice, usually draw a sample $\mu_t \sim \text{distribution of } \mu \mid r_0, a_0, \dots, r_{t-1}, a_{t-1} \text{ and then compute}$ $a_t = \arg \max_k \mu_t^{(k)}$, which is the same thing as $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$)

That's it! Statistically, this is a super simple and elegant algorithm (though computationally, it may not be easy to update the posterior at each time step)

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:a) Sample the optimal arm as much as possible (duh)

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:

- a) Sample the optimal arm as much as possible (duh)
- b) Ensure arms that might still be optimal aren't overlooked

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:

- a) Sample the optimal arm as much as possible (duh)
- b) Ensure arms that might still be optimal aren't overlooked
- c) Not waste undue time on less promising arms

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:

- a) Sample the optimal arm as much as possible (duh)
- b) Ensure arms that might still be optimal aren't overlooked
- c) Not waste undue time on less promising arms

Intuitively: want to sample arms proportionally to how promising they are

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:

- a) Sample the optimal arm as much as possible (duh)
- b) Ensure arms that might still be optimal aren't overlooked
- c) Not waste undue time on less promising arms
- Intuitively: want to sample arms proportionally to how promising they are

This is exactly what Thompson sampling does, where "promising" is encoded very naturally as: "the probability that the arm is the optimal arm, given all the data so far"

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:

- a) Sample the optimal arm as much as possible (duh)
- b) Ensure arms that might still be optimal aren't overlooked
- c) Not waste undue time on less promising arms
- Intuitively: want to sample arms proportionally to how promising they are

This is exactly what Thompson sampling does, where "promising" is encoded very naturally as: "the probability that the arm is the optimal arm, given all the data so far"

No arbitrary δ tuning parameter, but do have to choose prior π

<u>Thompson sampling</u>: $a_t \sim \text{distribution of } k^* \mid r_0, a_0, \dots, r_{t-1}, a_{t-1}$

Why is this a good idea?

A good tradeoff of exploration vs exploitation should:

- a) Sample the optimal arm as much as possible (duh)
- b) Ensure arms that might still be optimal aren't overlooked
- c) Not waste undue time on less promising arms
- Intuitively: want to sample arms proportionally to how promising they are

This is exactly what Thompson sampling does, where "promising" is encoded very naturally as: "the probability that the arm is the optimal arm, given all the data so far"

No arbitrary δ tuning parameter, but do have to choose prior π π can often be chosen "uninformatively" to a default prior such as the uniform, or can encode nuanced prior information/belief about the arms' reward distributions

Thompson sampling vs other algorithms

Thompson sampling samples arms proportionally to how promising they are

Thompson sampling vs other algorithms

Thompson sampling samples arms proportionally to how promising they are Note this sampling is much more sophisticated than, say, ε -greedy, which really just samples according to 2 categories: "most promising" and "other"

Thompson sampling vs other algorithms

Thompson sampling samples arms proportionally to how promising they are Note this sampling is much more sophisticated than, say, ε -greedy, which really just samples according to 2 categories: "most promising" and "other" But it's also quite different from UCB, whose OFU approach doesn't really involve "sampling" at all, i.e., every a_t for UCB is a *deterministic* function of the previous data
Thompson sampling vs other algorithms

Thompson sampling samples arms proportionally to how promising they are Note this sampling is much more sophisticated than, say, ε -greedy, which really just samples according to 2 categories: "most promising" and "other" But it's also quite different from UCB, whose OFU approach doesn't really involve

"sampling" at all, i.e., every a_t for UCB is a *deterministic* function of the previous data

My interpretation: OFU provides a simple heuristic to accomplish what Thompson sampling does by design, namely, sample arms according to how promising they are

Thompson sampling vs other algorithms

Thompson sampling samples arms proportionally to how promising they are Note this sampling is much more sophisticated than, say, ε -greedy, which really just samples according to 2 categories: "most promising" and "other"

But it's also quite different from UCB, whose OFU approach doesn't really involve "sampling" at all, i.e., every a_t for UCB is a *deterministic* function of the previous data

My interpretation: OFU provides a simple heuristic to accomplish what Thompson sampling does by design, namely, sample arms according to how promising they are

Thompson sampling can do this because of the Bayesian bandit: assuming a prior on the reward distributions makes the arm means random, otherwise it wouldn't even make sense to talk about "the probability that an arm is the best arm"

Thompson sampling vs other algorithms

Thompson sampling samples arms proportionally to how promising they are Note this sampling is much more sophisticated than, say, ε -greedy, which really just samples according to 2 categories: "most promising" and "other"

But it's also quite different from UCB, whose OFU approach doesn't really involve "sampling" at all, i.e., every a_t for UCB is a *deterministic* function of the previous data

My interpretation: OFU provides a simple heuristic to accomplish what Thompson sampling does by design, namely, sample arms according to how promising they are

Thompson sampling can do this because of the Bayesian bandit: assuming a prior on the reward distributions makes the arm means random, otherwise it wouldn't even make sense to talk about "the probability that an arm is the best arm"

Although derived from the Bayesian bandit, Thompson sampling has excellent practical performance across bandit problems, whether or not they are Bayesian!

Thompson sampling has excellent performance in practice, but is still just a heuristic

Thompson sampling has excellent performance in practice, but is still just a heuristic However, asymptotically, i.e., as $T \rightarrow \infty$, it actually is optimal in a certain sense:

Thompson sampling has excellent performance in practice, but is still just a heuristic However, asymptotically, i.e., as $T \to \infty$, it actually is optimal in a certain sense: Recall our instance-dependent UCB regret bound proved that with high probability, $N_t^{(k)} \leq \frac{2\ln(2KT/\delta)}{g_k}$ or, equivalently, $\frac{N_t^{(k)}}{2\ln(2KT/\delta)} \leq \frac{1}{g_k}$

Thompson sampling has excellent performance in practice, but is still just a heuristic However, asymptotically, i.e., as $T \to \infty$, it actually is optimal in a certain sense: Recall our instance-dependent UCB regret bound proved that with high probability, $N_t^{(k)} \le \frac{2\ln(2KT/\delta)}{g_k} \text{ or, equivalently, } \frac{N_t^{(k)}}{2\ln(2KT/\delta)} \le \frac{1}{g_k}$ There is actually a lower-bound result that says that for <u>any</u> bandit algorithm: $\liminf_{T \to \infty} \frac{\mathbb{E}[N_T^{(k)}]}{\ln(T)} \ge \frac{1}{d(\nu^{(k^*)}, \nu^{(k)})},$

where d is a distance between distributions called the Kullback-Leibler divergence

Thompson sampling has excellent performance in practice, but is still just a heuristic However, asymptotically, i.e., as $T \to \infty$, it actually is optimal in a certain sense: Recall our instance-dependent UCB regret bound proved that with high probability, $N_t^{(k)} \le \frac{2\ln(2KT/\delta)}{g_k} \text{ or, equivalently, } \frac{N_t^{(k)}}{2\ln(2KT/\delta)} \le \frac{1}{g_k}$ There is actually a lower-bound result that says that for <u>any</u> bandit algorithm: $\liminf_{T \to \infty} \frac{\mathbb{E}[N_T^{(k)}]}{\ln(T)} \ge \frac{1}{d(\nu^{(k^*)}, \nu^{(k)})},$

where d is a distance between distributions called the Kullback-Leibler divergence

It turns out that Thompson sampling satisfies this lower-bound with equality!

Thompson sampling has excellent performance in practice, but is still just a heuristic However, asymptotically, i.e., as $T \to \infty$, it actually is optimal in a certain sense: Recall our instance-dependent UCB regret bound proved that with high probability, $N_t^{(k)} \le \frac{2\ln(2KT/\delta)}{g_k} \text{ or, equivalently, } \frac{N_t^{(k)}}{2\ln(2KT/\delta)} \le \frac{1}{g_k}$ There is actually a lower-bound result that says that for <u>any</u> bandit algorithm: $\liminf_{T \to \infty} \frac{\mathbb{E}[N_T^{(k)}]}{\ln(T)} \ge \frac{1}{d(\nu^{(k^*)}, \nu^{(k)})},$

where d is a distance between distributions called the Kullback-Leibler divergence

It turns out that Thompson sampling satisfies this lower-bound with equality! So it is asymptotically optimal, not just in its rate, but its constant too!

Thompson sampling has excellent performance in practice, but is still just a heuristic However, asymptotically, i.e., as $T \to \infty$, it actually is optimal in a certain sense: Recall our instance-dependent UCB regret bound proved that with high probability, $N_t^{(k)} \le \frac{2\ln(2KT/\delta)}{g_k} \text{ or, equivalently, } \frac{N_t^{(k)}}{2\ln(2KT/\delta)} \le \frac{1}{g_k}$ There is actually a lower-bound result that says that for <u>any</u> bandit algorithm: $\liminf_{T \to \infty} \frac{\mathbb{E}[N_T^{(k)}]}{\ln(T)} \ge \frac{1}{d(\nu^{(k^*)}, \nu^{(k)})},$

where d is a distance between distributions called the Kullback-Leibler divergence

It turns out that Thompson sampling satisfies this lower-bound with equality! So it is asymptotically optimal, not just in its rate, but its constant too! (UCB is not, but there are more complicated versions of it that are)

So Thompson sampling is basically exactly optimal for large T

So Thompson sampling is basically exactly optimal for large T

What could go wrong for smaller T? Suppose K = 2 and T = 3, and:

So Thompson sampling is basically exactly optimal for large T

What could go wrong for smaller T? Suppose K = 2 and T = 3, and:

• t = 0: $a_0 = 1$, $r_0 = 1$

So Thompson sampling is basically exactly optimal for large T

What could go wrong for smaller T? Suppose K = 2 and T = 3, and:

- t = 0: $a_0 = 1$, $r_0 = 1$
- $t = 1: a_1 = 2, r_1 = 0$

So Thompson sampling is basically exactly optimal for large T

What could go wrong for smaller T? Suppose K = 2 and T = 3, and:

- t = 0: $a_0 = 1$, $r_0 = 1$
- $t = 1: a_1 = 2, r_1 = 0$
- t = 2 (last time step, with $\hat{\mu}_2^{(1)} = 1$ and $\hat{\mu}_2^{(2)} = 0$): $a_2 = ?$

So Thompson sampling is basically exactly optimal for large T

What could go wrong for smaller T? Suppose K = 2 and T = 3, and:

- t = 0: $a_0 = 1$, $r_0 = 1$
- $t = 1: a_1 = 2, r_1 = 0$
- t = 2 (last time step, with $\hat{\mu}_2^{(1)} = 1$ and $\hat{\mu}_2^{(2)} = 0$): $a_2 = ?$

Thompson sampling has a decent probability of choosing $a_2 = 2$, since with just one sample from each arm, Thompson sampling isn't sure which arm is best.

So Thompson sampling is basically exactly optimal for large T

What could go wrong for smaller T? Suppose K = 2 and T = 3, and:

- t = 0: $a_0 = 1$, $r_0 = 1$
- $t = 1: a_1 = 2, r_1 = 0$
- t = 2 (last time step, with $\hat{\mu}_2^{(1)} = 1$ and $\hat{\mu}_2^{(2)} = 0$): $a_2 = ?$

Thompson sampling has a decent probability of choosing $a_2 = 2$, since with just one sample from each arm, Thompson sampling isn't sure which arm is best.

But $a_2 = 1$ is clear right choice here: there is no future value to learning more, i.e., no reason to explore rather than exploit.

So Thompson sampling is basically exactly optimal for large T

What could go wrong for smaller T? Suppose K = 2 and T = 3, and:

- t = 0: $a_0 = 1$, $r_0 = 1$
- $t = 1: a_1 = 2, r_1 = 0$
- t = 2 (last time step, with $\hat{\mu}_2^{(1)} = 1$ and $\hat{\mu}_2^{(2)} = 0$): $a_2 = ?$

Thompson sampling has a decent probability of choosing $a_2 = 2$, since with just one sample from each arm, Thompson sampling isn't sure which arm is best.

But $a_2 = 1$ is clear right choice here: there is no future value to learning more, i.e., no reason to explore rather than exploit.

Thompson sampling doesn't know this, and neither does UCB (although UCB wouldn't happen to make the same mistake in this case).

For small T, Thompson sampling is not greedy enough

For small *T*, Thompson sampling is not greedy enough

Fix: add a tuning parameter to make it more greedy. Some possibilities:

For small *T*, Thompson sampling is not greedy enough

Fix: add a tuning parameter to make it more greedy. Some possibilities:

• Update the Beta parameters by $1+\epsilon$ instead of just 1 each time

For small *T*, Thompson sampling is not greedy enough

Fix: add a tuning parameter to make it more greedy. Some possibilities:

- Update the Beta parameters by $1 + \epsilon$ instead of just 1 each time
- Instead of just taking one sample of µ and computing the greedy action with respect to it, take n samples, compute the greedy action with respect to each, and pick the mode of those greedy actions

For small *T*, Thompson sampling is not greedy enough

Fix: add a tuning parameter to make it more greedy. Some possibilities:

- Update the Beta parameters by $1 + \epsilon$ instead of just 1 each time
- Instead of just taking one sample of µ and computing the greedy action with respect to it, take n samples, compute the greedy action with respect to each, and pick the mode of those greedy actions

All of these favor arms that the algorithm has more confidence are good (i.e., arms that have worked well so far), as opposed to arms that *may* be good

For small *T*, Thompson sampling is not greedy enough

Fix: add a tuning parameter to make it more greedy. Some possibilities:

- Update the Beta parameters by $1 + \epsilon$ instead of just 1 each time
- Instead of just taking one sample of µ and computing the greedy action with respect to it, take n samples, compute the greedy action with respect to each, and pick the mode of those greedy actions

All of these favor arms that the algorithm has more confidence are good (i.e., arms that have worked well so far), as opposed to arms that *may* be good

Such tuning can greatly improve Thompson sampling's performance even for reasonably large T (the asymptotic optimality of vanilla TS is *very* asymptotic)

Today



• Gittins index

So far, we have always taken the time horizon T to be fixed

So far, we have always taken the time horizon T to be fixed Another model we might consider is for the T to be *random*, and a simple yet intuitive distribution for it is to imagine that the bandit ends at each time step with a fixed probability $1 - \gamma$, given that it has reached that time step:

So far, we have always taken the time horizon *T* to be fixed Another model we might consider is for the *T* to be *random*, and a simple yet intuitive distribution for it is to imagine that the bandit ends at each time step with a fixed probability $1 - \gamma$, given that it has reached that time step: 1. $\mathbb{P}(T = 1) = 1 - \gamma$

So far, we have always taken the time horizon *T* to be fixed Another model we might consider is for the *T* to be *random*, and a simple yet intuitive distribution for it is to imagine that the bandit ends at each time step with a fixed probability $1 - \gamma$, given that it has reached that time step: 1. $\mathbb{P}(T = 1) = 1 - \gamma$ 2. $\mathbb{P}(T = 2) = \mathbb{P}(T = 2 \mid T > 1)\mathbb{P}(T > 1) = (1 - \gamma)\gamma$

So far, we have always taken the time horizon *T* to be fixed Another model we might consider is for the *T* to be *random*, and a simple yet intuitive distribution for it is to imagine that the bandit ends at each time step with a fixed probability $1 - \gamma$, given that it has reached that time step: 1. $\mathbb{P}(T = 1) = 1 - \gamma$

- 2. $\mathbb{P}(T=2) = \mathbb{P}(T=2 \mid T > 1)\mathbb{P}(T > 1) = (1 \gamma)\gamma$
- 3. $\mathbb{P}(T=3) = \mathbb{P}(T=3 \mid T > 2)\mathbb{P}(T > 2) = (1 \gamma)(1 (1 \gamma) (1 \gamma)\gamma)$

So far, we have always taken the time horizon *T* to be fixed Another model we might consider is for the *T* to be *random*, and a simple yet intuitive distribution for it is to imagine that the bandit ends at each time step with a fixed probability $1 - \gamma$, given that it has reached that time step: 1. $\mathbb{P}(T = 1) = 1 - \gamma$ 2. $\mathbb{P}(T = 2) = \mathbb{P}(T = 2 \mid T > 1)\mathbb{P}(T > 1) = (1 - \gamma)\gamma$ 3. $\mathbb{P}(T = 3) = \mathbb{P}(T = 3 \mid T > 2)\mathbb{P}(T > 2) = (1 - \gamma)(1 - (1 - \gamma) - (1 - \gamma)\gamma)$ $= (1 - \gamma)\gamma^2$

So far, we have always taken the time horizon T to be fixed Another model we might consider is for the T to be random, and a simple yet intuitive distribution for it is to imagine that the bandit ends at each time step with a fixed probability $1 - \gamma$, given that it has reached that time step: 1. $\mathbb{P}(T = 1) = 1 - \gamma$ 2. $\mathbb{P}(T=2) = \mathbb{P}(T=2 \mid T > 1)\mathbb{P}(T > 1) = (1 - \gamma)\gamma$ 3. $\mathbb{P}(T=3) = \mathbb{P}(T=3 \mid T > 2)\mathbb{P}(T > 2) = (1 - \gamma)(1 - (1 - \gamma) - (1 - \gamma)\gamma)$ $= (1 - \gamma)\gamma^{2}$... *T* is geometric, i.e., $\mathbb{P}(T = n) = (1 - \gamma)\gamma^{n-1}$

So far, we have always taken the time horizon T to be fixed Another model we might consider is for the T to be random, and a simple yet intuitive distribution for it is to imagine that the bandit ends at each time step with a fixed probability $1 - \gamma$, given that it has reached that time step: 1. $\mathbb{P}(T = 1) = 1 - \gamma$ 2. $\mathbb{P}(T=2) = \mathbb{P}(T=2 \mid T > 1)\mathbb{P}(T > 1) = (1 - \gamma)\gamma$ 3. $\mathbb{P}(T=3) = \mathbb{P}(T=3 \mid T > 2)\mathbb{P}(T > 2) = (1 - \gamma)(1 - (1 - \gamma) - (1 - \gamma)\gamma)$ $= (1 - \gamma)\gamma^{2}$... *T* is geometric, i.e., $\mathbb{P}(T = n) = (1 - \gamma)\gamma^{n-1}$ Thus, assuming T independent of the data, then we'll get to r_t w/p: $P(T > t) = \sum_{n=1}^{\infty} (1-x)y^{n-1} = 1 - (1-x)\left|\sum_{n=1}^{t} y^{n-1}\right| = 1 - (1-x)\frac{1-x^{t}}{1-x} = y^{t}$

Exactly optimizing the expected reward

So we see any given r_t w/p γ^t
Exactly optimizing the expected reward

t=0

So we see any given r_t w/p $\gamma^t \Rightarrow$ total reward, in expectation only over *T*, is $\sum_{t} \gamma^t r_t$

Exactly optimizing the expected reward

So we see any given r_t w/p $\gamma^t \Rightarrow$ total reward, in expectation only over *T*, is $\sum \gamma^t r_t$

In the <u>Bayesian bandit</u>, there is an algorithm that <u>exactly</u> optimizes $\mathbb{E}\left[\sum_{i=1}^{\infty} \gamma^{t} r_{t}\right]$



Today's summary:

Thompson sampling

- Operates in Bayesian bandit environment
- Samples optimal arm from its (posterior) distribution
- Achieves strong performance in practice

Gittins index

- Operates in Bayesian bandit with random horizon
- Exactly optimal in terms of expected (discounted) reward
- Some computational details to work out on HW1

Next time:

- Contextual bandits
- Other flavors of bandits

1-minute feedback form: https://bit.ly/3RHtlxy

