

Exploration vs. Exploitation

Krishna Devkota

Bielefeld University
5th May 2017

Reinforcement Learning in Autonomous Social agents

Outline

Overview

Random exploration (naive approach)

ϵ -greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Exploration vs. Exploitation (In Practice)

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

► Restaurant Selection

Exploitation: Go to a Restaurant that you know well

Exploration: Try a new restaurant

► Playing a game

Exploitation: Play a move that you are confident of

Exploration: Try a new move that you haven't played much

► Advertisement

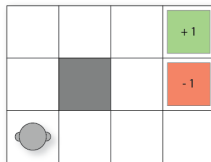
Exploitation: Play an advertisement that has been received well

Exploration: Try a new Ad that has not been played yet

Quick Recap of Markov Decision Process (MDP)

Markov decision processes formally describe an environment for reinforcement learning (decision making)

- ▶ A MDP can be represented as a 4-tuple $\langle S, A, P, R \rangle$ where:



- ▶ S is a set of states
- ▶ A is set of all the actions that the agent can take
- ▶ $P(s'|s, a)$ is a function that defines the transition probability (*Markovian*)
- ▶ $R(s, a)$ is the reward function, which gives the probability of receiving reward r after choosing a in state s

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Three different classes of approach to the problem

- ▶ Random exploration
 - ▶ Explore random actions (e.g. ϵ – *greedy*, softmax)
- ▶ Optimism in the face of uncertainty
 - ▶ estimate uncertainty on value
 - ▶ Prefer to explore states/actions with highest uncertainty
- ▶ Information state space
 - ▶ Consider agent's information as part of its state
 - ▶ Look ahead to see how information helps reward

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Random exploration (naive approach)

Multi-armed Bandit Problem (single-state MDP)

One of the simplest way to model a exploration/ exploitation dilemma is using a multi-armed Bandit Problem

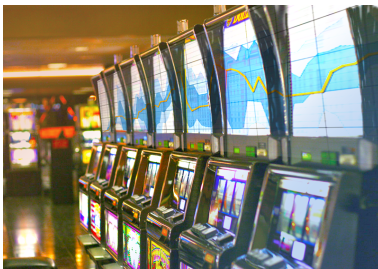


Figure: Bandit Machines

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Multi-armed bandit problem

Krishna Devkota

Random exploration (naive approach)

Initialization

Acting

Bayes- adaptive MDPs and Gittins indices

Summary

Formalizing Multi-armed Bandit problem

- ▶ Can be represented as a Tuple $\langle A, R \rangle$
- ▶ A is known set of Arms that can be Pulled i.e. actions that can be taken (m)
- ▶ $R^a(r) = \mathcal{P}[r|a]$ is the unknown Probability distribution over rewards
- ▶ Action taken at each step t by the agent : $a_t \in A$
- ▶ Reward generated by the environment: $r_t \sim R^{a_t}$

- Goal : maximize cumulative reward

$$\sum_{\tau=1}^t r_{\tau}$$

Exploration vs. Exploitation

Krishna Devkota

Overview

Random exploration (naive approach)

- €greedy
- Softmax
- Optimistic Initialization

Acting Optimistically in Uncertain situations

Bayes- adaptive
MDPs and Gittins
indices

Summary

Regret

► Action-value

mean reward for action a , $Q(a) = E[r|a]$

► Optimal value (V^*)

$$V^* = Q(a^*) = \max_{a \in A} Q(a)$$

► Regret

Opportunity lost for each step

$$l_t = E[V^* - Q(a_t)]$$

► Total Regret

Opportunity lost over all the steps

$$L_t = E\left[\sum_{\tau=1}^t (V^* - Q(a_\tau))\right]$$

■ Goal : maximize cumulative reward, hence minimize total regret

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

- ▶ Count ($N_t(a)$)
Expected number of times an action is taken,
- ▶ Gap (Δ_a)
Difference in value between action (a) and optimal action (a^*) $\Delta_a = V^* - Q(a)$
- ▶ Regret as a function of count and gap

$$L_t = \sum_{a \in A} E[N_t(a)] \Delta_a$$

■ Goal : Find a good algorithm, so we visit the bad state the least amount of times i.e. small counts for large gaps

■ Problem : Optimal value (V^*) unknown, and hence the gaps

Regret

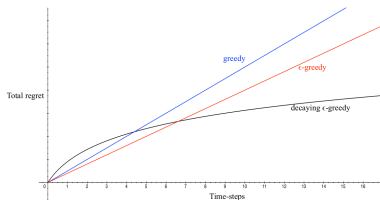


Figure: Comparing Regret values for different naive-algorithms

- ▶ If an algorithm **forever explores**, it will have a linear regret (ϵ - greedy)
- ▶ If an algorithm **never explores**, it will have a linear regret (greedy)
- ▶ So, how do we achieve sub-linear total regret?

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

A quick look at Greedy algorithm



Figure: Greedy algorithm based on local optimum

- Consider algorithm that estimates $\hat{Q}_t(a)$ which is closest to $Q_t(a)$ i.e. the **MC evaluation**:

$$\hat{Q}_t(a) = \frac{1}{N_t(a)} \sum_{t=1}^T r_t 1(a_t = a)$$

- Using **greedy algorithm** gives us: $a_t^* = \operatorname{argmax}_{a \in A} \hat{Q}_t(a)$
- **Problem**: We might get stuck onto a suboptimal action again and again

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Outline

Overview

Random exploration (naive approach)

ϵ greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

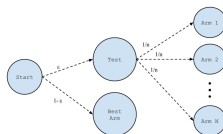
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

The ϵ – greedy algorithm



- ▶ Using the ϵ – greedy algorithm, we want to introduce some randomness into our greedy approach
- ▶ How do we do that??
 - ▶ with probability $1 - \epsilon$ act greedily, i.e. select $a_t^* = \operatorname{argmax}_{a \in A} \hat{Q}_t(a)$
 - ▶ with probability ϵ , select a random action

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Advantages of ϵ – greedy exploration:

- ▶ Simplest idea for ensuring continual exploration
- ▶ All m-actions are tried with non-zero probability

Drawback of ϵ – greedy exploration:

- ▶ Random actions selected uniformly. The worst possible action is just as likely to be selected as the second best action

Outline

Overview

Random exploration (naive approach)

ϵ greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Softmax

ϵ – *greedy* selected the random actions uniformly, giving equal weights to the good and the bad

Softmax remedies this by assigning a rank or weight to each of the actions, according to their *action-value* estimate

- ▶ Grade action probabilities by estimated values
- ▶ weight actions using linear combination of features $\phi(s, a)^T \theta$
- ▶ Probability of action is proportional to exponentiated weight

$$\pi_{\theta}(s, a) \propto e^{\phi(s, a)^T \theta}$$

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

- ▶ Bias exploration towards promising actions
- ▶ The most common softmax uses a Gibbs (or Boltzmann) distribution

Advantages:

- ▶ As appropriate weight associated with each action, the worst actions are unlikely to be chosen
- ▶ Good in scenarios where the worst actions are very unfavourable

Outline

Overview

Random exploration (naive approach)

ϵ greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Optimistic Initialization

- ▶ initialize $Q(a)$ to high value (i.e. assume all of our actions pay the best possible

$$Q(a) = r_{max}$$

- ▶ Use MC evaluation to incrementally update action value
- ▶ $\hat{Q}_t(a_t) = \hat{Q}_{t-1} + \frac{1}{N_t(a_t)}(r_t - \hat{Q}_{t-1})$

Advantage:

- ▶ Encourages exploration of unknown values

Drawback:

- ▶ We need to know the maximum possible reward r_{max}
- ▶ Can still get caught in suboptimal action

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Comparing greedy, ϵ -greedy, and Optimistic Initialization

For a 10-armed testbed, $N = 10$ possible actions, 1000 plays $Q(a)$ are chosen randomly from a Normal distribution $N(0, 1)$

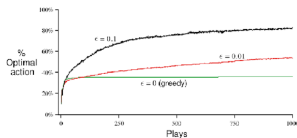


Figure: Greedy vs. ϵ -greedy

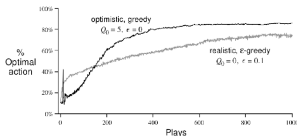


Figure: Normal case vs. Optimistically initialized case

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Methods based on Upper Confidence Bounds (UCBs)

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

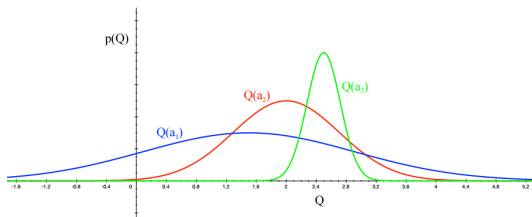
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Acting Optimistically in Uncertain situations



- ▶ The more uncertain we are about an action value
- ▶ The more important it is to explore that action
- ▶ It could turn out to be the best action

Outline

Overview

Random exploration (naive approach)

ϵ -greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Krishna Devkota

- ## Summary

Methods based on Upper Confidence Bounds (UCBs)

- ▶ **Confidence interval**
a range of values within which we are sure the mean lies with a certain probability
- ▶ For an action which has been tried less often, our estimated reward is less accurate so the confidence interval is larger
- ▶ It shrinks as we get more information (i.e. try the action more often)
- ▶ So, instead of trying the action with the highest mean, we can try the action with the highest upper bound on its confidence interval
- ▶ This is known as an optimistic policy

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Methods based on Upper Confidence Bounds (UCBs)

Steps:

- ▶ For each action value, estimate an upper confidence $\hat{U}_t(a)$ such that:
 $Q(a) \leq \hat{Q}_t(a) + \hat{U}_t(a)$ with high probability
- ▶ determined by the number of times $N(a)$ has been selected
 - ▶ Small $N_t(a) \Rightarrow$ *large* $\hat{U}_t(a)$
estimated value is uncertain
 - ▶ Large $N_t(a) \Rightarrow$ *small* $\hat{U}_t(a)$
estimated value is accurate
- ▶ Select action that maximizes the UCB
 $a_t = \operatorname{argmax}_{a \in A} \hat{Q}_t(a) + \hat{U}_t(a)$

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

- ▶ To solve for the bounds, we can turn to:
Chernoff-Hoeffding bound
- ▶ Then, pick a probability p for the true value to exceed the UCB
- ▶ Solve for $U_t(a)$, and reducing probability p , as more rewards are observed
- ▶ As $t \rightarrow \infty$, we select optimal action as given by:

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

This gives us the UCB1 algorithm:

$$a_t = \operatorname{argmax}_{a \in A} Q(a) + \sqrt{\frac{2 \log t}{N_t(a)}}$$

Quick recap of methods based on UCB

- ▶ Each arm is assigned an UCB for its mean reward
- ▶ Arm with the largest bound to be played
- ▶ Bound is not conventional upper limit for a confidence interval, hence difficult to compute
- ▶ However, making some basic assumptions, the expected number of times suboptimal arm a would be played by time t is:

$$E(n_{at}) \leq \left(\frac{1}{K(a, a^*)} + o(1) \right) \log t$$

where $K(a, a^*)$ is the Kullback-Leibler divergence between the reward distributions for arm a and optimal arm a^*

- ▶ This bound essentially says that the optimal arm will be played exponentially more often than any of the suboptimal arms, for large t

Outline

Overview

Random exploration (naive approach)

ϵ -greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Bayesian Bandits

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

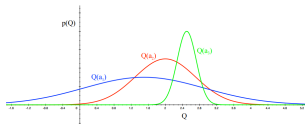
Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

How do we exploit prior knowledge about rewards?

- ▶ Recall, the distribution over our action-value function (Q) was unknown



- ▶ Instead, say we start with some prior distribution over the action value function
- ▶ Let $p[Q|w]$ be some distribution over action-value function (Q), where w is the parameter

Bayesian Bandits

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

- ▶ The parameters \mathbf{w} could be (say) the *mean* (μ) and the variances (σ) of each of our arms
- ▶ We could then compute posterior distribution over \mathbf{w} by using the Bayesian methods

$$p[\mathbf{w}|R_1, \dots, R_t]$$

- ▶ Use this posterior to guide exploration i.e. Probability matching
- ▶ Better performance for accurate prior

Random Probability Matching

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

- ▶ Randomized probability matching combines many positive aspects of the heuristic strategies mentioned above
- ▶ Probability matching selects action a according to probability that a is the optimal action

$$\pi(a|h_t) = P[Q(a) > Q(a'), \forall a' \neq a | h_t]$$

- ▶ Uncertain actions have higher probability of being max
- ▶ Can be difficult to compute analytically from posterior

Outline

Overview

Random exploration (naive approach)

ϵ -greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Thompson sampling

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

- way to implement probability matching

$$\pi(a|h_t) = P[Q(a) > Q(a'), \forall a' \neq a|h_t]$$

$$= E_{R|h_t}[1(a = \operatorname{argmax}_{a \in A} Q(a))]$$

Steps:

- ▶ Use Bayes law to compute posterior distribution $p[R|h_t]$, where $h_t = a_1, r_1, \dots, a_{t-1}, r_{t-1}$ is the history
- ▶ Sample a reward distribution R from posterior
- ▶ Compute action-value function $Q(a) = E[R_a]$
- ▶ Select action maximising value on sample, $a_t = \operatorname{argmax}_{a \in A} Q(a)$

Advantages of Probability matching techniques:

- ▶ the tuning parameters, and the decay schedule evolves in a principled, data-determined way
- ▶ In other methods, the parameters are arbitrarily set by analyst, and incorrect values bear huge costs
- ▶ Thompson sampling achieves Lai and Robbins lower bound

Disadvantages of Probability matching techniques:

- ▶ There is a need to sample from the posterior distribution
- ▶ This can require substantially more computing than other heuristics

Value of Information

- ▶ Why is exploration useful?
- ▶ Because we gain information
- ▶ Sometimes, sacrificing immediate rewards will be beneficial in the long run
- ▶ Other times, when on extremely limited budget, immediate reward might be beneficial
- ▶ Information gain is higher in uncertain situations, hence exploration is important
- ▶ What if we could quantify all those information, and use them to make informed decisions?

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Information State Space

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

- ▶ Bandits as a single-step decision problems
- ▶ Can be expanded to sequential decision-making problems
- ▶ Add a new information state \tilde{s}
 - ▶ \tilde{s} summarizes history in certain statistical way i.e.
 $\tilde{s}_t = f(h_t)$
- ▶ With each action, transition to a new information state \tilde{s}' with a certain probability \tilde{P}
- ▶ If we augment this info into our state space, we'll get a MDP in information state space

$$\tilde{M} = \langle \tilde{S}, A, \tilde{P}, R, \gamma \rangle$$

Outline

Overview

Random exploration (naive approach)

ϵ -greedy

Softmax

Optimistic Initialization

Acting Optimistically in Uncertain situations

Upper Confidence Bounds

Bayesian Bandits

Thompson sampling

Information State Space

Bayes- adaptive MDPs and Gittins indices

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

$$\tilde{M} = \langle \tilde{S}, A, \tilde{P}, R, \gamma \rangle$$

This can then be solved using reinforcement learning:

- ▶ Model-free reinforcement learning (Q-learning)
- ▶ Bayesian model-based reinforcement learning (Gittins indices)

Contextual Bandits

- ▶ Using similar idea, if we now add State Information to our Multi-armed bandit Tuple, we will get a Contextual Bandit

$$\langle A, S, R \rangle$$

- ▶ A is the set of actions, S is an unknown distribution over States, R is an unknown distribution over rewards
- ▶ Using the same set of Principles, this can then be solved

Summary

■ We saw how the problem of exploration/ exploitation can be tricky sometimes

■ We looked at some of the heuristic strategies to handle the dilemma

- ▶ Equal allocation
- ▶ Play-the-winner
- ▶ Deterministic greedy strategies
- ▶ Hybrid strategies such as ϵ - greedy, and Softmax

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ greedy
Softmax
Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds
Bayesian Bandits
Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Summary

■ We looked at some strategies based on bounds in both Frequentist and Bayesian approach

- ▶ UCB1
- ▶ Random Probability matching (Thompson sampling)

■ We introduced an information theoretic criteria to quantify Information value

■ Finally, we saw how this Bandit problem can be expanded into a full Markov Decision Problem

Exploration vs.
Exploitation

Krishna Devkota

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Overview

Random
exploration (naive
approach)

ϵ -greedy

Softmax

Optimistic
Initialization

Acting
Optimistically in
Uncertain
situations

Upper Confidence
Bounds

Bayesian Bandits

Thompson sampling

Information State
Space

Bayes- adaptive
MDPs and Gittins
indices

Summary

Thank you!



D. Silver.

<http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>

Scott, Steven L. "A modern Bayesian look at the multiarmed bandit." *Applied Stochastic Models in Business and Industry* 26.6 (2010): 639-658.



Steven L. Scott

A modern Bayesian look at the multiarmed bandit
Applied Stochastic Models in Business and Industry 26.6
, (2010): 639-658.