



Correlated Multi-armed Bandits

CS 6780 Advanced Machine Learning

Zhengdi Shen
Bangrui Chen
Saul Toscano Palmerin

April 24, 2015



1 Motivation

2 Problem Setting

3 Approach

Bandits

- Pull arms sequentially so as to maximize the total expected reward, e.g, gambler faces at a row of slot machines.
Robbins in 1952, realizing the importance of the problem.

Bandits

- Pull arms sequentially so as to maximize the total expected reward, e.g, gambler faces at a row of slot machines. Robbins in 1952, realizing the importance of the problem.
- The crucial tradeoff the gambler faces at each trial is between "exploitation" of the machine that has the highest expected payoff and "exploration" to get more information about the expected payoffs of the other machines.

Independent Arms

- Rewards for each arm are generally assumed to be independent of each other.

Independent Arms

- Rewards for each arm are generally assumed to be independent of each other.
- A theorem, the Gittins Index published first by John C. Gittins gives an optimal policy in the Markov setting for maximizing the expected discounted reward.

Independent Arms

- Rewards for each arm are generally assumed to be independent of each other.
- A theorem, the Gittins Index published first by John C. Gittins gives an optimal policy in the Markov setting for maximizing the expected discounted reward.
- This assumption enables us to consider each arm separately, but Lai and Robbins proved the regret under an arbitrary policy increases linearly with number of arms. Most policies that assume independence require each arm to be tried at least once, and are impractical in settings involving many arms.

Dependent Arms

- What if they are dependent? E.g., ads on similar topics, using similar text/phrases, should have similar rewards.

Dependent Arms

- What if they are dependent? E.g., ads on similar topics, using similar text/phrases, should have similar rewards.
- In such a setting, the information obtained from pulling one arm can change our understanding of other arms. Here, we want a policy whose regret is independent of the number of arms.

Linearly Parameterized Bandits

- Mersereau proposed a simple model that the reward of each arm depends on a single random variable z , with a known prior distribution.

Linearly Parameterized Bandits

- Mersereau proposed a simple model that the reward of each arm depends on a single random variable z , with a known prior distribution.
- Since the reward of each arm depends on a single random variable, the mean rewards are perfectly correlated.

Linearly Parameterized Bandits

- Mersereau proposed a simple model that the reward of each arm depends on a single random variable z , with a known prior distribution.
- Since the reward of each arm depends on a single random variable, the mean rewards are perfectly correlated.
- Under certain assumptions, the cumulative Bayes risk over T periods under a greedy policy admits an $O(\log(T))$ upper bound, independent of the number of arms.

Linearly Parameterized Bandits

- Mersereau proposed a simple model that the reward of each arm depends on a single random variable z , with a known prior distribution.
- Since the reward of each arm depends on a single random variable, the mean rewards are perfectly correlated.
- Under certain assumptions, the cumulative Bayes risk over T periods under a greedy policy admits an $O(\log(T))$ upper bound, independent of the number of arms.
- In this project, we are going to consider a model that the reward of each arm depends linearly on a multivariate random variable, with a known prior distribution.

A More Practical Motivation

For a newly registered user on Yelp, how should yelp select the forwarding restaurants at each time so that it can maximize the expected average rating of this new user?

① Motivation

② Problem Setting

③ Approach

Problem Setting

- We represent each restaurant with a 20 dimensional binary vector X , and the 20 features are Pizza, Sandwiches, Food, Nightlife, American(new), Bars, American(traditional), Mexican, Chinese, Italian, Japanese, Fast Food, Burgers, Breakfast and Brunch, Coffee and Tea, Delis, Indian, Thai, Sushi Bars, Mediterranean and Asian Fusion. Denote the set of arms as U^{20} .

Problem Setting

- We represent each restaurant with a 20 dimensional binary vector X , and the 20 features are Pizza, Sandwiches, Food, Nightlife, American(new), Bars, American(traditional), Mexican, Chinese, Italian, Japanese, Fast Food, Burgers, Breakfast and Brunch, Coffee and Tea, Delis, Indian, Thai, Sushi Bars, Mediterranean and Asian Fusion. Denote the set of arms as U^{20} .
- We assume each user has a user preference vector θ corresponding to the 20 different features listed above, with $\theta \sim N(\mu_0, \Sigma_0)$ with μ_0 and Σ_0 known (This can be calculated from the historical data).

Problem Setting

- Y_t : the reward of playing arm $X \in U^{20}$ in period t , which is given by

$$Y_t = X \cdot \theta + W_t$$

where the measurement error term $\{W_t : t \geq 1\}$ is i.i.d distributed with $N(0, \sigma^2)$.

- Denote the posterior distribution after t observations as μ_t and Σ_t respectively.

Problem Setting

- Y_t : the reward of playing arm $X \in U^{20}$ in period t , which is given by

$$Y_t = X \cdot \theta + W_t$$

where the measurement error term $\{W_t : t \geq 1\}$ is i.i.d distributed with $N(0, \sigma^2)$.

- \mathcal{H}_t : the set of possible histories until the end of period t .
- Denote the posterior distribution after t observations as μ_t and Σ_t respectively.

The Model

- $\psi = (\psi_1, \psi_2, \dots)$: Policy ψ is a sequence of functions such that $\psi_t : \mathcal{H}_{t-1} \rightarrow \mathbb{U}^{20}$ selects an arm in period t based on the history until the end of period $t-1$.

The Model

- $\psi = (\psi_1, \psi_2, \dots)$: Policy ψ is a sequence of functions such that $\psi_t : \mathcal{H}_{t-1} \rightarrow \mathbb{U}^{20}$ selects an arm in period t based on the history until the end of period $t-1$.
- For any policy ψ and $\theta_0 \in \mathbb{R}^{20}$, the T -period cumulative regret under ψ given $\theta = \theta_0$, denoted by $\text{Regret}(\theta_0, T, \psi)$, is defined by

$$\text{Regret}(\theta_0, T, \psi) = \sum_{t=1}^T E \left[\max_{X \in \mathbb{U}^{20}} X \cdot \theta_0 - X_t \cdot \theta_0 \mid \theta = \theta_0 \right]$$

where for any $t \geq 1$, $X_t \in \mathbb{U}^{20}$ is the arm chosen under ψ in period t .

① Motivation

② Problem Setting

③ Approach

"Optimal" Algorithm

Lower bound for regret

For an arbitrary policy, the regret is at least $\Omega(r\sqrt{T})$ under some regularity conditions, where the set of arms is compact in R^r .

Regret for PEPE algorithm

The Phased Exploration and Greedy Exploitation algorithm has regret $\Omega(r\sqrt{T})$ under some regularity conditions.

Question: Is the PEPE algorithm really optimal? The big O notation might hide a large constant!

PEGE

Phased Exploration and Greedy Exploitation

Description: For each cycle $c \geq 1$, complete the following two phases.

- (1) **Exploration (r periods):** For $k = 1, 2, \dots, r$, play arm $\mathbf{b}_k \in \mathbb{U}_r$ given in Assumption 1(b), and observe the reward $Y^{b_k}(c)$. Compute the OLS estimate $\hat{\theta}(c) \in \mathbb{R}^r$, given by

$$\begin{aligned}\hat{\theta}(c) &= \frac{1}{c} \left(\sum_{k=1}^r \mathbf{b}_k \mathbf{b}_k' \right)^{-1} \sum_{s=1}^c \sum_{k=1}^r \mathbf{b}_k Y^{b_k}(s) \\ &= \theta + \frac{1}{c} \left(\sum_{k=1}^r \mathbf{b}_k \mathbf{b}_k' \right)^{-1} \sum_{s=1}^c \sum_{k=1}^r \mathbf{b}_k W^{b_k}(s)\end{aligned}$$

where for any k , $Y^{b_k}(s)$, and $W^{b_k}(s)$ denote the observed reward and the error random variable associated with playing arm \mathbf{b}_k in cycle s .

PEGE

Phased Exploration and Greedy Exploitation

- **Exploitation (c periods):** Play the greedy arm $\mathbf{G}(c) = \arg \max_{X \in \mathbb{U}^r} X \cdot \hat{\theta}(c)$ for c periods.

UCB

Upper Confidence Bound

Given $\theta \sim N(\mu_0, \Sigma_0)$, for t from 1 to T :

- Play arm $X_{i_t} = \arg \max \{ \mu_{t-1} \cdot X_i + 1.96 X_i' \Sigma_{t-1} X_i \}$
- Calculate μ_t and Σ_t based on arm Y_t , reward X_{i_t} and μ_{t-1} and Σ_{t-1} .

Question?

Thanks for your time!