

Generalized Linear Contextual Bandits

Information Theory, Statistics, and Learning Course Project

Amirreza Velae
Amirabbas Afzali

Sharif University of Technology

May 12, 2025

1 Introduction and Background

- Why Bandit Algorithms?
- Online optimization
 - A Gentle Introduction to Online Optimization
 - Expert Advice
- Bandit Algorithms
 - Connection to Bandit Algorithms
- Know your Enemy
 - Adaptive vs. Oblivious Adversaries

2 Future Work

AlphaGo: RL and Bandit Algorithms Showdown

What is AlphaGo?

- Developed by DeepMind to play the board game Go.
- Utilizes deep RL and MCTS.

Significance:

- Defeated world champion, Lee Sedol, in 2016.
- History-making move 37.



Figure: AlphaGo playing against Lee Sedol

Definition: Online optimization involves making a sequence of decisions based on incoming data, where each decision is made without knowledge of future data.

Fundamental Assumption:

- 1 **Bunded Losses:** The losses determined by an adversary should not be allowed to be unbounded.
- 2 **Bounded Decision Set:** The decision set must be somehow bounded and/or structured, though not necessarily finite.

Prediction from Expert Advice

- **Prediction from Expert Advice** is a framework in online learning and decision-making.
- Involves making sequential predictions by aggregating advice from a set of experts.
- The goal is to perform nearly as well as the best expert in hindsight.

Key Components

- ① **Experts:** Sources that provide predictions or advice.
- ② **Learner (Main Character):** Aggregates expert advice to make decisions.
- ③ **Feedback:** Information about the actual outcome to update future predictions.
- ④ **Objective:** Minimize the difference between the learner's performance and the best expert's performance.

Family Members as Experts

- Main character decides whether to take an umbrella each day.
- Three family members act as **experts** providing daily weather predictions, the father, mother, and brother.



Figure: Family Members as Experts

Example Scenario

Table: Daily Weather Predictions and Outcomes

Day	Father	Mother	Brother	Actual Weather
1	No Rain	No Rain	Rain	?

Example Scenario

Table: Daily Weather Predictions and Outcomes

Day	Father	Mother	Brother	Actual Weather
1	No Rain	No Rain	Rain	Rain
2	Rain	No Rain	No Rain	?

Example Scenario

Table: Daily Weather Predictions and Outcomes

Day	Father	Mother	Brother	Actual Weather
1	No Rain	No Rain	Rain	Rain
2	Rain	No Rain	No Rain	No Rain
3	Rain	Rain	No Rain	?

Example Scenario

Table: Daily Weather Predictions and Outcomes

Day	Father	Mother	Brother	Actual Weather
1	No Rain	No Rain	Rain	Rain
2	Rain	No Rain	No Rain	No Rain
3	Rain	Rain	No Rain	Rain
4	No Rain	Rain	Rain	?

Example Scenario

Table: Daily Weather Predictions and Outcomes

Day	Father	Mother	Brother	Actual Weather
1	No Rain	No Rain	Rain	Rain
2	Rain	No Rain	No Rain	No Rain
3	Rain	Rain	No Rain	Rain
4	No Rain	Rain	Rain	No Rain
5	Rain	No Rain	No Rain	?

Example Scenario

Table: Daily Weather Predictions and Outcomes

Day	Father	Mother	Brother	Actual Weather
1	No Rain	No Rain	Rain	Rain
2	Rain	No Rain	No Rain	No Rain
3	Rain	Rain	No Rain	Rain
4	No Rain	Rain	Rain	No Rain
5	Rain	No Rain	No Rain	Rain
Cost	2	3	3	-

Weighted Majority Algorithm

- **Mechanism:**

- Assigns a weight to each expert based on their past performance.
- Aggregates predictions by considering these weights.
- Updates weights multiplicatively based on the correctness of experts' predictions.

- **Applications:**

- Ensemble learning in machine learning.
- Financial decision-making.

How Weighted Majority Works

Algorithm Steps:

- ➊ **Initialization:** Assign equal weights to all experts.
- ➋ **Prediction**

$$\text{Prediction} = \text{sign} \left(\sum_{i=1}^N w_t(i) \cdot \text{Prediction}_t(i) \right)$$

- ➌ **Update Weights:** After observing the outcome, update weights:

$$w_{t+1}(i) = \begin{cases} w_t(i) & \text{if expert } i \text{ was correct} \\ w_t(i) \cdot \varepsilon & \text{if expert } i \text{ was incorrect} \end{cases}$$

- ➍ **Iteration:** Repeat the prediction and update steps for each round.

Parameters:

- N : Number of experts.
- $\varepsilon \in (0, 1)$: Penalty factor for incorrect experts.

Regret Bound for Weighted Majority

Lemma : Denote by M_t the number of mistakes the algorithm makes until time t , and by $M_t(i)$ the number of mistakes made by expert i until time t . Then, for any expert $i \in [N]$ we have

$$M_T \leq 2(1 + \varepsilon)M_T(i) + \frac{2 \log N}{\varepsilon}$$

Corollary : The regret of the WM algorithm is bounded by

$$M_T \leq 2M_T(i^*) + O(\sqrt{M_T(i^*) \log N})$$

where i^* is the best expert.

Proof: Just let $\varepsilon^* = \sqrt{\frac{\log N}{M_T(i^*)}}$.

Randomized Weighted Majority Algorithm

- Algorithm Steps:

- 1 Initialization:** Set $w_1(i) = 1$ for all experts.

- 2 Probability Assignment:**

$$P_t(i) = \frac{w_t(i)}{\sum_{j=1}^N w_t(j)}$$

- 3 Expert Selection:** Choose expert i with probability $P_i(t)$.

- 4 Prediction and Update:** Make prediction based on selected expert and update weights as in Weighted Majority.

- Better Regret Bound:

$$\mathbb{E}[M_T] \leq (1 + \varepsilon) M_T(i^*) + \frac{\log N}{\varepsilon}$$

Introduction to Multi-Armed Bandit Algorithms

- **Definition:** A framework for making a sequence of decisions under uncertainty, aiming to maximize cumulative rewards.
- **Key Concepts:**
 - **Exploration:** Trying different actions to gather more information.
 - **Exploitation:** Selecting the best-known action to maximize immediate reward.
- **Types of Bandit Problems:**
 - **Stochastic Bandits:** Rewards are drawn from fixed probability distributions.
 - **Contextual Bandits:** Incorporates contextual information to make more informed decisions.
- **Applications:**
 - Recommendation Systems
 - Clinical Trials
 - Online Advertising

Bandit Algorithms in Online Advertising

- **Use Case:** Optimizing Ad Selection to Maximize CTR
- **How It Works:**
 - Each ad variant is considered an arm of the bandit.
 - The algorithm dynamically selects which ad to display based on past performance.
 - Balances exploration (trying new ads) with exploitation (showing top-performing ads).

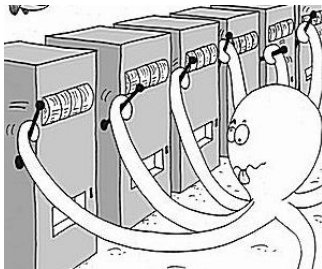


Figure: Bandit Setting

Online Optimization vs. Bandit Algorithms

Both online optimization and bandit algorithms involve sequential decisions, but differ in feedback and objectives.

Online Optimization

- **Full Feedback:** Receives complete information about all possible actions after each decision.
- **Objective:** Minimize cumulative loss compared to the best fixed decision in hindsight.

Bandit Algorithms

- **Partial Feedback:** Only receives feedback for the action actually taken, not for all possible actions.
- **Objective:** Balance exploration and exploitation to maximize cumulative rewards.

Exploration and Exploitation

Now that we have a way to convert an oblivious adversary to an adaptive one, we can focus on using online optimization algorithms to solve bandit problems.

- **Exploration:** Use the perturbation or other methods to explore different experts and their predictions.
- **Exploitation:** Follow the leader to exploit the best-performing expert.

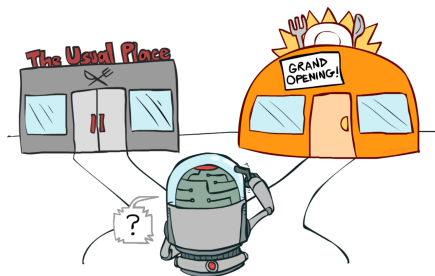


Figure: Exploration vs. Exploitation in Bandit Algorithms

Adaptive vs. Oblivious Adversaries

Oblivious Adversary

- **Definition:** Fixes the sequence of events beforehand.
- **Traits:**
 - Non-responsive.
 - Simpler to analyze.



Adaptive Adversary

- **Definition:** Adjusts based on algorithm's past actions.
- **Traits:**
 - Responsive and dynamic.
 - Harder to counter.



Use your Maximum Power

In general, oblivious adversaries are easier to handle than adaptive adversaries. However, we can convert an oblivious adversary to an adaptive by being pessimistic.



Figure: An evil adversary in machine learning using its maximum power

Regret Bound for Any Adversary

Lemma: Fix T , let H^* be the set of decision histories of length 0 to $T-1$, and let K^* be the set of all cost histories of length 0 to $T-1$. Then, fix a decision algorithm $\mathcal{A} : K^* \rightarrow \Delta(S)$, where $\Delta(S)$ is the set of probability distributions on the set S of possible decisions. Define

$$R(\mathcal{A}, \mathcal{V}) = \mathbb{E}_{\mathcal{A}, \mathcal{V}} \left[\sum_{t=1}^T \mathbf{c}^t \cdot \mathbf{x}^t - \min_{\mathbf{x} \in S} \sum_{t=1}^T \mathbf{c}^t \cdot \mathbf{x} \right].$$

Let \mathcal{V} be an arbitrary adversary. Then, there exists an oblivious adversary \mathcal{V}' such that

$$R(\mathcal{A}, \mathcal{V}') \geq R(\mathcal{A}, \mathcal{V}).$$

- ① **Exploring while Exploiting:** I think the exploitation phase has some information that can be used to explore better.
- ② **Hyperparameter Tuning:** The parameters γ can be tuned to improve the performance of the algorithm.
- ③ **Game Theory:** The algorithm can be viewed as an Stackelberg game. Thus there is a possibility of using game theory.
- ④ **Experiments:** The algorithm hasn't tested on real-world data. It would be interesting to see how it performs in practice.
- ⑤ **Bound Tightening:** Investigate whether the current bounds of $O(T^{3/4}\sqrt{\ln T})$ against adaptive adversaries and $O(T^{2/3})$ against oblivious adversaries can be improved to $O(\sqrt{T})$.

References I



Tor Lattimore and Csaba Szepesvári.

Bandit Algorithms.

Cambridge University Press, 2020.



Elad Hazan.

Introduction to Online Convex Optimization.

Cambridge University Press, 2016.



H. B. McMahan and A. Blum.

Online Geometric Optimization in the Bandit Setting Against an Adaptive Adversary.

In J. Shawe-Taylor and Y. Singer, editors, *Learning Theory*, Springer Berlin Heidelberg, 2004, pp. 109–123.



P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire.

The Nonstochastic Multiarmed Bandit Problem.

SIAM Journal on Computing, 32(1):48–77, 2002.



A. Kalai and S. Vempala.

Efficient algorithms for online decision problems.

Journal of Computer and System Sciences, 71(3):291–307, 2005.