# Multi-Armed Bandits
## Efficient Algorithms for Learning with Semi-Bandit Feedback

Abdul, Adwait, Anirudh

$30^{th}$ May 2021

# Overview

# Motivation and Introduction

### Goals of this section

- Motivate using Online Least Cost Path problem.
- Understanding the Semi-Bandits setting.
- Understanding the Combinatorial problem.
- Understanding Linear Generalization.
- Formally introducing the problem.

- We have a (directed or undirected) graph $G$ with
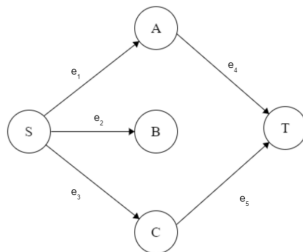    - $L$ = number of edge,
    - $w_i$ = cost associate with the edge $i$.
    - $S$ = start node and,
    - $T$ = target node.
- The goal: Find the cheapest path from $S$ to $T$ w.r.to the cost associated with the path.

# Motivation
## The Online Least Cost Path problem - An example



- In the above graph $L = 5$.
- A path is denoted as a $5 - D$ vector.
  - Path SAT, $A_1 =$

    | 1 | 0 | 0 | 1 | 0 |
    |---|---|---|---|---|

  - Path SCT, $A_2 =$

    | 0 | 0 | 1 | 0 | 1 |
    |---|---|---|---|---|

- Set of valid paths, $\mathcal{A} = \{A_1, A_2\} \subseteq \{0, 1\}^L$

- Let us say we have decided to go with the path SAT.

| 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|

- To compute the cost, we may get one of the following:

  1. Full information setting:

  $$w = \begin{array}{|c|c|c|c|c|} \hline 0.3 & 0.2 & 0.7 & 0.5 & 0.4 \\ \hline \end{array}$$

  2. Semi-Bandits setting:

  $$w = \begin{array}{|c|c|c|c|c|} \hline 0.3 & x & x & 0.5 & x \\ \hline \end{array}$$

  3. Full-Bandits setting:

  $$l = \begin{array}{|c|} \hline 0.8 \\ \hline \end{array}$$

- For the Full information and Semi-Bandits setting, the loss is
  $l = A \cdot w$

# Motivation
## How is this different from the familiar MAB setting?

- In the MAB setting, we choose $i\epsilon[L]$.
  - $i$ can take one of $L$ values.
- However, here we choose $A\epsilon\mathcal{A} \subseteq \{0,1\}^L$.
  - $A$ can take one of potentially $2^L$ values.

- That is, we want to choose one or more arms in each round.
- In the Online Least Cost path problem,
  Edge $\equiv$ Arm
  Path $\equiv$ multiple Arms.
- This is an example of Combinatorial Optimization problem.

# Motivation

- An immediate issue with this setting is that applying algorithms such as *EXP* − 3 may not have a good regret bound, since both *L* and $|\mathcal{A}|$ can be very large in real life settings.
- Can we make it independent of *L*?
  - Here's a sneak peak...

# Motivation
## Linear Generalization

- In this setting (and many other real life settings), we have access to features of edges, such as,
    1. Length,
    2. Traffic, and
    3. Road Quality.

|       | Length | Traffic | Road Quality |
|-------|--------|---------|--------------|
| $e_1$ | 10     | 3       | 7            |
| $e_2$ | 4      | 6       | 3            |
| $e_3$ | 8      | 1       | 4            |
| $e_4$ | 6      | 2       | 5            |
| $e_5$ | 3      | 7       | 2            |

# Motivation
## Linear Generalization

- Hence, we know a (possibly imperfect) generalization matrix $\Phi$.
  In this case,

$$\Phi = \begin{bmatrix} 10 & 3 & 7 \\ 4 & 6 & 3 \\ 8 & 1 & 4 \\ 6 & 2 & 5 \\ 3 & 7 & 2 \end{bmatrix}$$

- We assume that the cost of an edge is a linear combination of its feature values.

$$cost(e) = \theta_1 \ length(e) + \theta_2 \ traffic(e) + \theta_3 \ quality(e) = \vec{\theta} \cdot \phi_e$$

- Now, the task is to estimate the entries in $\vec{\theta}$.
- Notice the advantage:
  The values to be estimated equal the number of features, and not $L$.

Let's formalize what all we have seen till now.

# Introduction
## Formalism

### A Combinatorial Optimization Problem

Can be represented as a triple $(E, \mathcal{A}, w)$ where,

- $E$ = Set of $L$ arms.
- $\mathcal{A}$ = Set of (some or all) subsets $A$ of $E$ with $|A| \leq K$.
- $w : E \to \mathbb{R}$ is the weight function.

- The total weight of $A \in \mathcal{A}$ is $\sum_{e \in A} w_e$.

## Semi-Bandit Problem

The losses associated to only the chosen arms are seen.

## Linear Generalization

For $e \epsilon [L]$, $w(e) = \Phi_e \theta$,

where $\Phi_e$ is the feature vector of the arm $e$, and $\theta$ indicates the weights of the features.

# Introduction
## Formalism

### The Combinatorial semi-bandits

It is an online learning problem where at each step the learning agent chooses a subset of arms $A$ subject to combinatorial constraints, and then observes weights of the selected arms, $\{w(e) : e \in A\}$, and gets their sum $f(A, w)$ as a payoff.

### The Goal

- The learner chooses a subset of arms $A^t$ at round $t$. The loss suffered at this round $R_t = f(A^*, w_t) - f(A^t, w)$ (In the Reward setup).

- We want the cumulative Regret to be *good*,
  And the algorithm to be efficient when applied in Large scale settings.

# Efficient Algorithm for Learning with Semi-Bandit Feedback (Paper 1)

# Learning with Semi-Bandit Feedback

- The paper considers the problem of online Combinatorial Optimization under semi-bandit feedback. (No linear Generalization)
- The paper assumes:
    - Finite Decision Set (potentially very large)
    - Efficient offline combinatorial optimization is possible
    - Elements of the decision set can be described with $d$-dimensional arrays with at most $m$ non-zero entries.
    - For example a decision vector in which 3 arms out of possible 7 arms are played will look like $\begin{bmatrix} 1 & 0 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$

-

**Parameters:** set of decision vectors $\mathcal{S} = \{v(1), v(2), \ldots, v(N)\} \subseteq \{0, 1\}^d$, number of rounds $T$;

For all $t = 1, 2, \ldots, T$, repeat

1. The learner chooses a probability distribution $p_t$ over $\{1, 2, \ldots, N\}$.
2. The learner draws an action $I_t$ randomly according to $p_t$. Consequently, the learner plays decision vector $V_t = v(I_t)$.
3. The environment chooses loss vector $\ell_t$.
4. The learner suffers loss $V_t^\top \ell_t$.
5. The learner observes some feedback based on $\ell_t$ and $V_t$.

- Generic framework to accommodate a number of interesting problem instances such as path planning, ranking and matching problems , finding minimum weight spanning trees and cut sets etc.

# Normal Bandit Setting

- Total $N$ arms and we chose only one arm.
- A distribution $p_t$ over the arms where $p_{t,i} = \mathcal{P}[l_t = i | \mathcal{F}_{t-1}]$ where $\mathcal{F}_{t-1}$ is the history of the learners observation and choice made upto step $t-1$.
- Most bandit algorithms rely on feeding some loss estimates to black box algorithm like Hedge etc. A common loss estimate used is
  $\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i}} \mathbb{I}[l_t = i]$          (An unbiased estimate)
- Almost all existing algorithms use some form of the above loss estimate. But $p_t$ is not readily available for all algorithms like FPL.
- Following loss estimate is proposed

## Geometric Resampling

- In round $t$ the learner draws $I_t \sim p_t$
- For n = 1, 2 , ....
    - Let $n \leftarrow n + 1$
    - Draw $I'(n) \sim p_t$
    - If $I'(n) = I_t$, break
- Let $K_t = n$

- $\hat{\ell_{t,i}} = \ell_{t,i} \mathbb{I}[I_t = i] K_t$ (Also unbiased)
- The number samples can be unbounded, so we cap the number of samples by $M$ and use $\hat{K_t} = min(K_t, M)$
- Introduces some bias, but if $M$ is chosen appropriately the performance does not hurt much.

## Generalizing Geometric Sampling to Semi-Bandit Feedback

- In round $t$ again the learner draws $l_t \sim p_t$
- For each arm draw $M$ samples (because anyway we are capping by $M$). We can effectively draw only $M$ samples and use them for all the arms being played. So, draw $M$ additional indices $I'_t(1), I'_t(2) \dots I'_t(M) \sim p_t$
- Define $K_{t,j} = min\{1 \leq s \leq M : v_j(I'_t(s)) = v_j(I_t)\}$ for all arms.
- Use Loss estimate for each arm as $\hat{\ell}_{t,j} = K_{t,j} V_{t,j} \ell_{t,j}$

- Since $V_{t,j}$ is non-zero only for the arms being played the estimates are well defined.

# FPL with Geometric Sampling

- 

**Input:** $\mathcal{S} = \{v(1), v(2), \ldots, v(N)\} \subseteq \{0,1\}^d$, $\eta \in \mathbb{R}^+$, $M \in \mathbb{Z}^+$;
**Initialization:** $\widehat{L}(1) = \cdots = \widehat{L}(d) = 0$;
**for** $t=1,\ldots,T$ **do**

    Draw $Z(1), \ldots, Z(d)$ independently from distribution $\text{Exp}(\eta)$;

    Choose action $I = \underset{i \in \{1,2,\ldots,N\}}{\arg\min} \left\{ v(i)^\top \left( \widehat{L} - Z \right) \right\}$;

    $K(1) = \cdots = K(d) = M$;

    $k = 0$;                /* Counter for reoccurred indices */

    **for** $n=1,\ldots,M\text{-}1$ **do**             /* Geometric Resamplig */

        Draw $Z'(1), \ldots, Z'(d)$ independently from distribution $\text{Exp}(\eta)$;

        $I(n) = \underset{i \in \{1,2,\ldots,N\}}{\arg\min} \left\{ v(i)^\top \left( \widehat{L} - Z' \right) \right\}$;

        **for** $j=1,\ldots,d$ **do**

                **if** $v(I(n))(j) = 1$ & $K(j) = M$ **then**

                        $K(j) = n$;

                        $k = k + 1$;

                        **if** $k = \|v(I)\|_1$ **then** break;    /* All indices reoccurred */

                **end**

        **end**

    **end**

    **for** $j=1,\ldots,d$ **do** $\widehat{L}(j) = \widehat{L}(j) + K(j)v(I)(j)\ell(j)$ ;        /* Update */

**end**

# Computational Efficiency under Semi Bandit Feedback

- The expected number of times the algorithm draws an action up to time step $T$ can be upper bounded by $dT$. This can be shown as

- For each co-ordinate that the original arm had 1, we take sampling until we get 1 in the same co-ordinate again. (Here we do not assume cutoff)
- Let $q_{t,k} = \mathbb{E}[V_{t,k}|\mathcal{F}_{t-1}]$
- At time step $t$ for a given co-ordinate $k$ with 1, the expected number of samples is $\frac{1}{q_{t,k}}$ while the probability of co-ordinate $k$ being 1 is $q_{t,k}$
- So, expected number of samples is $\sum_{k=1}^{d} q_{t,k} \frac{1}{q_{t,k}} = d$. ( For one round)
- For $T$ rounds it will be $dT$.
- The expected running time is comforting.

# Regret Bound

- The total expected regret of FPL with Geometric Resampling satisfies $R_n \leq \frac{m(\log d + 1)}{\eta} + \eta m dT + \frac{dT}{eM}$ under semi-bandit information.
- For the full information setting if $C_T = \sum_{t=1}^{T} \mathbb{E}[V_t^{\mathsf{T}} \ell_t]$ then the expected regret of FPL satisfies $R_n \leq \frac{m(\log d + 1)}{\eta} + \eta m C_t$ under full information.

---

- For $\eta = \sqrt{\frac{\log d + 1}{dT}}$ and $M \geq \frac{\sqrt{dT}}{eM(\sqrt{\log d + 1})}$, $R_n = \mathcal{O}(m\sqrt{dT \log d})$
- For $\eta = \sqrt{\frac{(\log d + 1)}{mT}}$ $R_n = \mathcal{O}(m^{\frac{3}{2}}\sqrt{T(\log d + 1)})$.

- Full information setting $\rightarrow$ Optimal Regret $\rightarrow \mathcal{O}(m\sqrt{T\log d})$
  FPL-with-GR achives $\mathcal{O}(m^{\frac{3}{2}}\sqrt{T(\log d + 1)})$ off by a factor of $\sqrt{m}$
- Semi-bandit setting $\rightarrow$ Optimal Regret $\mathcal{O}(\sqrt{mdT})$.
  FPL-with-GR achives $\mathcal{O}(m\sqrt{dT\log d})$ off by a factor of $\sqrt{m\log d}$

# Conclusion and Future work

- The paper has introduced the first general efficient algorithm for online combinatorial optimization under semi-bandit feedback.
- It remains as an open problem whether the gaps we have shown can be closed for FPL-style algorithms.
- The most important open problem is the development efficient online linear optimization with full bandit feedback.

# Efficient Learning in Large Scale Combinatorial Semi-Bandits (Paper 2)

# Large Scale Combinatorial Semi-Bandits

- This paper looks at efficient algorithms for stochastic combinatorial semi-bandits with **linear generalization**
- We look at algorithms that have regret bound **INDEPENDENT OF L** (hence large scale)

# Setup Used for Combinatorial optimization

## Combinatorial setup

$$(E, \mathcal{A}, P)$$

- $E = \{1, \ldots, L\} \equiv$ arms (the ground set)
- $\mathcal{A} \subseteq \{A \subseteq E : |A| \leq K\} \equiv$ allowed combinations of arms
- $P \equiv$ probability distribution over weights $w \in \mathbb{R}^L$ on $E$. ($\bar{w} = \mathbb{E}[w]$)
- After the arms $A$ are pulled, we observe the individual return of each arm $\{w(e) : e \in A\}$
- $f(A, w) \equiv$ loss on pulling $A \in \mathcal{A}$

We assume linear generalization(possibly imperfect) and that the agent knows the generalization matrix $\Phi \in \mathbb{R}^{L \times d}$. If $\bar{w} \in span[\Phi]$ we call it as coherent learning case otherwise agnostic. WLOG rank$[\Phi] = d$.

# Performance Metrics

- $R_t = f(A^*, w_t) - f(A^t, w)$ where $A^* = ORACLE(E, \mathcal{A}, \bar{w})$
- For fixed $\bar{w}$,
  $R(T) = \sum_{t=1}^{T} \mathbb{E}[R_t | \bar{w}] \equiv$ **EXPECTED CUMULATIVE REGRET**
  where expectation is over random weights and possible randomization in the algorithm.
- for randomly generated $\bar{w}$,
  $R_{bayes}(T) = \sum_{t=1}^{T} \mathbb{E}[R_t] \equiv$ **BAYES CUMULATIVE REGRET** where expectation is over random weights, possible randomization in the algorithm and also over $\bar{w}$.
- $\theta^* = argmin_\theta ||\bar{w} - \Phi\theta||$. Since $rank[\Phi] = d$, $\theta^*$ is uniquely defined. Also in coherent learning case $\bar{w} = \Phi\theta^*$.

# ALGORITHMS

# KALMAN FILTERING for updating parameters

## Algorithm 1

**Input:** $\bar{\theta}_t, \Sigma_t, \sigma$, and feature-observation pairs $\{(\phi_e, \mathbf{w}_t(e)) : e \in A^t\}$

Initialize $\bar{\theta}_{t+1} \leftarrow \bar{\theta}_t$ and $\Sigma_{t+1} \leftarrow \Sigma_t$

**for** $k = 1, \ldots, |A^t|$ **do**

  Update $\bar{\theta}_{t+1}$ and $\Sigma_{t+1}$ as follows, where $a_k^t$ is the $k$th element in $A^t$

$$\bar{\theta}_{t+1} \leftarrow \left[ I - \frac{\Sigma_{t+1}\phi_{a_k^t}\phi_{a_k^t}^T}{\phi_{a_k^t}^T\Sigma_{t+1}\phi_{a_k^t} + \sigma^2} \right] \bar{\theta}_{t+1} + \left[ \frac{\Sigma_{t+1}\phi_{a_k^t}}{\phi_{a_k^t}^T\Sigma_{t+1}\phi_{a_k^t} + \sigma^2} \right] \mathbf{w}_t\left(a_k^t\right) \quad \text{and} \quad \Sigma_{t+1} \leftarrow \Sigma_{t+1} - \frac{\Sigma_{t+1}\phi_{a_k^t}\phi_{a_k^t}^T\Sigma_{t+1}}{\phi_{a_k^t}^T\Sigma_{t+1}\phi_{a_k^t} + \sigma^2},$$

**end for**

**Output:** $\bar{\theta}_{t+1}$ and $\Sigma_{t+1}$

# Combinatorial Linear Thompson Sampling

## CombLinTS

**Input:** Combinatorial structure $(E, \mathcal{A})$, generalization matrix $\Phi \in \mathbb{R}^{L \times d}$, algorithm parameters $\lambda, \sigma > 0$, oracle ORACLE

Initialize $\Sigma_1 \leftarrow \lambda^2 I \in \mathbb{R}^{d \times d}$ and $\bar{\theta}_1 = 0 \in \mathbb{R}^d$

**for all** $t = 1, 2, \ldots, n$ **do**

    Sample $\theta_t \sim N\left(\bar{\theta}_t, \Sigma_t\right)$

    Compute $A^t \leftarrow \texttt{ORACLE}(E, \mathcal{A}, \Phi\theta_t)$

    Choose set $A^t$, and observe $\mathrm{w}_t(e), \forall e \in A^t$

    Compute $\bar{\theta}_{t+1}$ and $\Sigma_{t+1}$ based on Algorithm 1

**end for**

# CombLinTS: Key Points

- $\lambda \rightarrow$ inverse regularization parameter.
- smaller $\lambda$ makes the covariance matrix $\sum_t$ closer to 0.
- smaller $\lambda \Rightarrow$ narrower prior $\Rightarrow$ insufficient exploration $\Rightarrow$ degraded performance of CombLinTS.
- $\sigma$ controls the decrease rate of covariance matrix $\Sigma_t$.
- Large $\sigma$ will lead to slow learning and a smaller $\sigma$ will make the algorithm quickly converge to some sub-optimal coefficient vector.

## Regret Bound

- If $\bar{w} = \Phi\theta^*$, the prior on $\theta^*$ is $\mathcal{N}(0, \lambda^2 I)$, the noises $(w_t(e) - \bar{w}(e))$ are i.i.d sampled from $\mathcal{N}(0, \sigma^2)$ then CombLinTS guarantees:
  $R_{Bayes}(T) = \tilde{O}(K\lambda\sqrt{dT\min\{\ln(L), d\}})$

- The conditions ensure it is a coherent gaussian case.
- The $\tilde{O}$ notation hides the logarithmic factors.
- The regret bound is a minimum of two bounds. The first bound is $L$-dependent as $O(\sqrt{\ln(L)})$. The second bound is $L$-independent but is $\tilde{O}(d)$ instead of $\tilde{O}(\sqrt{d})$.

# Combinatorial Linear UCB

## CombLinUCB

**Input:** Combinatorial structure $(E, \mathcal{A})$, generalization matrix $\Phi \in \mathbb{R}^{L \times d}$, algorithm parameters $\lambda, \sigma, c > 0$, oracle ORACLE

Initialize $\Sigma_1 \leftarrow \lambda^2 I \in \mathbb{R}^{d \times d}$ and $\bar{\theta}_1 = 0 \in \mathbb{R}^d$

**for all** $t = 1, 2, \ldots, n$ **do**

Define the UCB weight vector $\hat{\mathbf{w}}_t$ as

$$\hat{\mathbf{w}}_t(e) = \langle \phi_e, \bar{\theta}_t \rangle + c\sqrt{\phi_e^T \Sigma_t \phi_e} \quad \forall e \in E$$

Compute $A^t \leftarrow \text{ORACLE}(E, \mathcal{A}, \hat{\mathbf{w}}_t)$

Choose set $A^t$, and observe $\mathbf{w}_t(e), \forall e \in A^t$

Compute $\bar{\theta}_{t+1}$ and $\Sigma_{t+1}$ based on Algorithm 1

**end for**

# CombLinUCB: Key Points

- $\lambda \rightarrow$ inverse regularization parameter.
- $\sigma$ controls the decrease rate of covariance matrix $\Sigma_t$..
- The constant $c$ controls the degree of optimism.
- Small $c \Rightarrow$ algorithm might converge to some sub-optimal coefficient vector due to insufficient exploration.
- Large $c \Rightarrow$ excessive exploration and slow learning.

# CombLinUCB: Regret Bound

Assuming $P$ is a subset of $[0,1]^L$, the stochastic item weights are statistically independent under $P$. Then for the coherent learning case i.e. $\bar{w} = \Phi\theta$ we have: For any $\lambda$, $\sigma$ and $\delta \in (0,1)$ and any $c$ satisfying

$$c \geq \frac{1}{\sigma}\sqrt{d\ln(1 + \frac{Tk\lambda^2}{d\sigma^2}) + 2\ln(\frac{1}{\delta})} + \frac{||\theta^*||_2}{\lambda}$$ we have:

$$R(T) \leq 2cK\lambda\sqrt{\frac{dT\ln(1+\frac{Tk\lambda^2}{d\sigma^2})}{\ln(1+\frac{\lambda^2}{\sigma^2})}} + TK\delta$$

Specifically if we chose $\lambda = \sigma = 1$, $\delta = \frac{1}{nK}$ and $c$ is the lower bound on its above condition then:

## Regret Bound

$$R(T) = \tilde{O}(Kd\sqrt{T})$$

# Comparison with state of art

## Standard Result

Standard results have been set for the no generalisation case.

No generalization $\Rightarrow$ lower bound of $\Omega(\sqrt{LKT})$

$\Rightarrow \Phi = I \Rightarrow L = d \Rightarrow$ No generalization lower bound $= \Omega(\sqrt{KdT})$

## CombLinTS

$\tilde{O}(\sqrt{K \min\{\ln(L), d\}})$ larger than the $\Omega(\sqrt{KdT})$

## CombLinUCB

$\tilde{O}(\sqrt{Kd})$ larger than the $\Omega(\sqrt{KdT})$

## Note

The $\Omega(\sqrt{d})$ and $\Omega(\sqrt{K})$ factors are due to linear generalization. Full tightness analysis has been left to future work.

# Conclusion and Future Work

- This paper has introduced two learning algorithms CombLinTS and CombLinUCB for stochastic combinatorial semi bandits with linear genrelariation.
- The main contribution has been that the paper has successfully introduced $L$ independent bounds which is highly useful for real life problems e.g Online advertisements(millions of users and products).
- This paper has left it open to derive bounds for the agnostic learning case.
- Another open problem is how to extend the results to combinatorial semi bandits with non-linear generalization.

# TAKEAWAYS

# What have we seen so far

- **Geometric Resampling:**
  - Don't need $p_t$ explicitly;
  - Computationally efficient;
  - Fits into any semi bandit algorithm

- **Geometric Resampling:**
    - Don't need $p_t$ explicitly;
    - Computationally efficient;
    - Fits into any semi bandit algorithm
- **Linear Generalization:**
    - Makes regret bound L independent

THANK YOU