

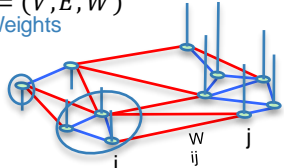
# Cheap Bandits

Manjesh K Hanawal, Venkatesh Saligrama (BU, Boston, USA)  
Michal Valko, Remi Munos (INRIA, Lille, France)



## Problem Statement

- Undirected Graph:  $G = (V, E, W)$ 
  - $N$  Nodes,  $W = \{w_{ij}\}$ : Weights



- Signal on Graph
  - Reward Function  $f: V \rightarrow R$

- Locate maxima  $u^* = \operatorname{argmax}_{v \in V} f(v)$

- Actions:
  - Noisy Cluster Averages; Differentiated Costs

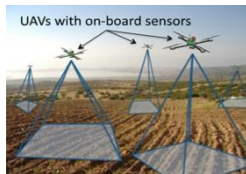
Goal: Locate  $u^*$  with min Cost?

## Motivation and Application

Surveillance/Geography

Forest Cover Dataset: Labeled samples on 30m<sup>2</sup> region  
**Nodes:** Forest Regions; **Edge weights:** feature similarity;  
**Rewards:** Density of species. Locate highest density.  
**Actions:** Zoom-in (high cost); Zoom-out (low cost).

Other Examples:  
 Sensor Networks  
 Radar Search  
 Online advertisements



Large Number of nodes  
 Few samples to observe.  $T \ll N$

## Reward Function

- Linear Reward:  $f = Q\alpha^*$ 
  - $Q$  is the eigenvectors of the graph Laplacian
  - Linear bandits, with parameter  $\alpha^*$

- Smooth Reward
  - Neighboring nodes have similar reward

$$(u, v) \in E \Rightarrow f(u) \sim f(v):$$

$$|Lf|_2^2 = \sum_{(u,v)} w_{uv} (f(u) - f(v))^2$$

- Assume  $|Lf|_2^2 \leq c$  [Valko et al. ICML'15]

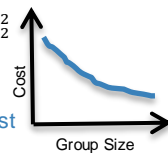
## Actions and Costs

- Actions set is subset of simplex:  $S \subset \Delta^N$ 
  - Node actions: Sample nodes  
For a node  $u \in V$  use  $s(u) = \delta(u - v)$
  - Group actions: Sample Average of a subset  
For a subset  $A \in V$  use  $s(u) = \frac{1}{|A|} \sum_{v \in V} \delta(u - v)$

- Cost of actions

- For all  $s \in S$

$$C(s) = \sum_{u,v} (s(u) - s(v))^2 = |Ls|_2^2$$



- Why this cost function
  - Larger the group size smaller the cost
  - Probing Nodes has high cost
  - In Fourier domain: Energy of  $s$

## Learning Setting and Objective

- Policy ( $\pi$ ): In each round  $t$ , select action  $s_t \in S$

- Observed Reward at rounds  $t$ :  
 $r(s_t) = \sum_{u \in V} s_t(u) f(u) + \epsilon_t$  (noise)

- Regret of policy  $\pi$ :

$$R_T(\pi) = Tf(u^*) - E[\sum r(s_t)]$$

- Cost of policy  $\pi$ :

$$C_T(\pi) = \sum_t C(s_t)$$

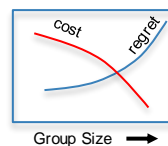
- Goal:

$$\min_{(\pi, S)} C_T(\pi)$$

$$s.t. R_T(\pi) \leq R_T^*$$

- Conflicting goals:
  - Node actions give better estimates, but costly
  - Group actions give poor estimates, but cheaper

What is the best Regret Constraint  $R_T^*$ ?



## Lower Bounds

- No smoothness constraints ( $c \rightarrow \infty$ )

- Arbitrary bounded set of action

$$R_T = \Omega(N\sqrt{T})$$
 [Dani et al. COLT'08]

- Finite set of actions

$$R_T = \Omega(\sqrt{NT})$$
 [Chu et al. AISTATS'08]

- For smooth reward

$$R_T = \Omega(\sqrt{dT}): d \ll N$$

$$d = \max\{i: \lambda_i(i-1) \leq \frac{1}{\log T}\}$$
 [Valko et al. ICML'14]

- $d$ -sparsely connected clusters
- Need at least one sample from each cluster

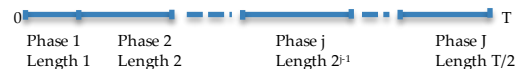
We aim for  $R_T^* = O(d\sqrt{T})$

## CheapUCB Algorithm

- Main idea: use similarity of neighborhood

- Group actions provide good node information  
 $u \in A: f(u) \sim \frac{1}{|A|} \sum_{v \in A} f(v) + \text{const}$

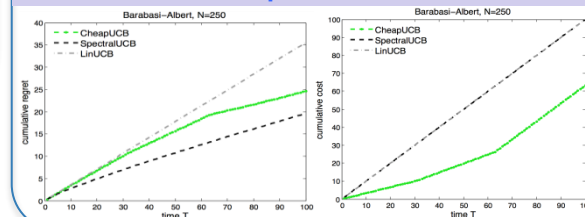
- CheapUCB: UCB based algorithm. Selects arms optimistically like SpectralUCB [Valko et al. ICML'15]



- Phases: Split the  $T$  into  $J = \lceil \log T \rceil$  phases
  - Length: Phase  $j=1, 2, \dots, J$  is of  $2^{j-1}$  rounds
  - Select action: In phase  $j$  select groups of size  $J+1$
- Zoom-in slowly using progressively costly actions

Algorithm	Regret bound	Cost
SpectralUCB (ICML'14)	$O(d\sqrt{T})$	$T$
CheapUCB (This paper)	$O(d\sqrt{T})$	$3T/4$

## Experiments



CheapUCB achieves at least 25% cost savings