

Information Directed Sampling for Stochastic Bandits with Graph Feedback

Fang Liu¹, Swapna Buccapatnam² and Ness Shroff¹

¹The Ohio State University

²AT&T Labs Research



THE OHIO STATE UNIVERSITY



Introduction

We study stochastic multi-armed bandits with graph feedback:

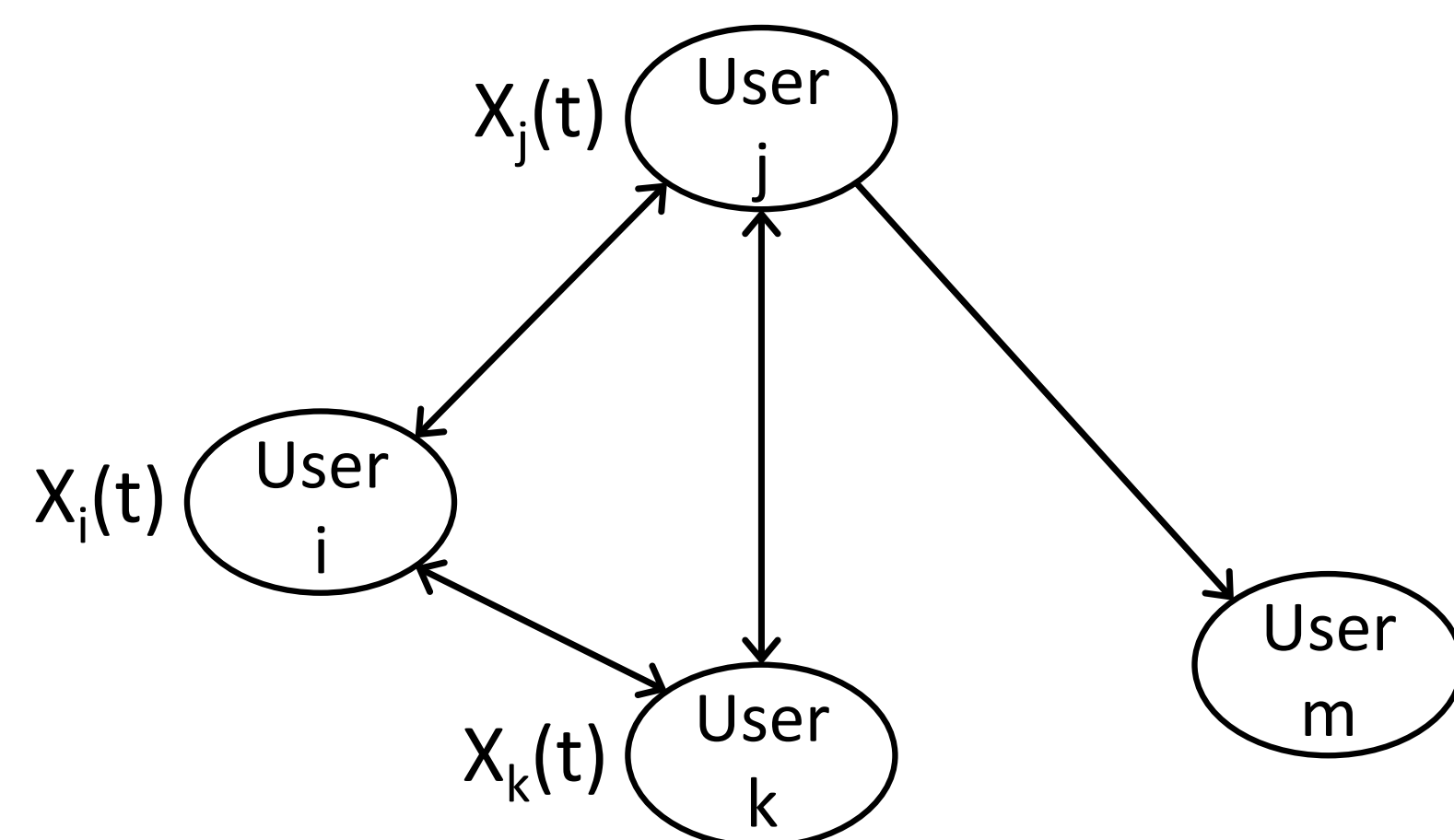
- Bandit: only obtain the reward of the chosen arm;
- Graph feedback: playing one arm may reveal other arms;
- Goal: online learning to minimize the regret due to uncertainty.

Motivation:

- Side observations are available
- Can reduce the regret
- How to use the graph information?

Applications:

- Online advertising on social media
- Recommendation systems with social connections like Yelp, Tripadvisor, ...



Model

Basic setting:

- Discrete time model of horizon T , there are K arms;
- At each time t , choosing an arm A_t returns a reward Y_{t,A_t} drawn from the distribution of arm A_t ;
- A^* : arm with the highest expected reward.
- Regret: expected loss compared to the oracle that plays arm A^* each time.

$$\mathbb{E}[R(T)] = \mathbb{E} \left[\sum_{t=1}^T Y_{t,A^*} - Y_{t,A_t} \right]$$

Graph feedback $G_t = (K, E_t)$ may change over time.

- Deterministic graph: G_t is known.
- Random graph: G_t is unknown.
- Erdos-Renyi graph

Randomized policy π_t .

- Update the posterior distribution α_t of A^*
- Δ_t : instantaneous expected regret
- h_t : information gain

Algorithm

Algorithm 1 Meta-algorithm for Information Directed Sampling with Graph Feedback

Input: Time horizon T and feedback graph model $(G_t)_{t \leq T}$
for t **from** 1 **to** T **do**
 Updating statistics: compute α_t , Δ_t and h_t accordingly.
 Generating policy: generate π_t as a function of $(\alpha_t, \Delta_t, h_t, G_t)$. (To be determined)
 Sampling: sample A_t according to π_t , play action A_t and receive reward Y_{t,A_t} .
 Observations: observe $Y_{t,a}$ if $(A_t, a) \in E_t$, where $G_t = (K, E_t)$ is the graph generated by G_t .
end for

TS-N policy: $\pi_t^{\text{TS-N}} = \alpha_t$. Unaware of graph, probability matching.

IDS-N policy: $\min_{\pi_t \in \mathcal{S}} (\pi_t^T \Delta_t)^2 / (\pi_t^T G_t h_t)$. Minimize the information ratio.

IDSN-LP policy: $\min_{\pi_t \in \mathcal{S}} \pi_t^T \Delta_t$ s.t. $\pi_t^T G_t h_t \geq \alpha_t^T G_t h_t$. Minimize the regret.

IDS-LP policy: $\min_{\pi_t \in \mathcal{S}} \pi_t^T \Delta_t$ s.t. $\pi_t^T G_t h_t \geq \alpha_t^T h_t$. Minimize the regret.

Analysis

Theorem 1. For any (deterministic or random) graph feedback, the Bayesian regret of IDS-LP is

$$\mathbb{E}[R(T, \pi^{\text{IDS-LP}})] \leq \sqrt{\frac{K}{2} T H(\alpha_1)}$$

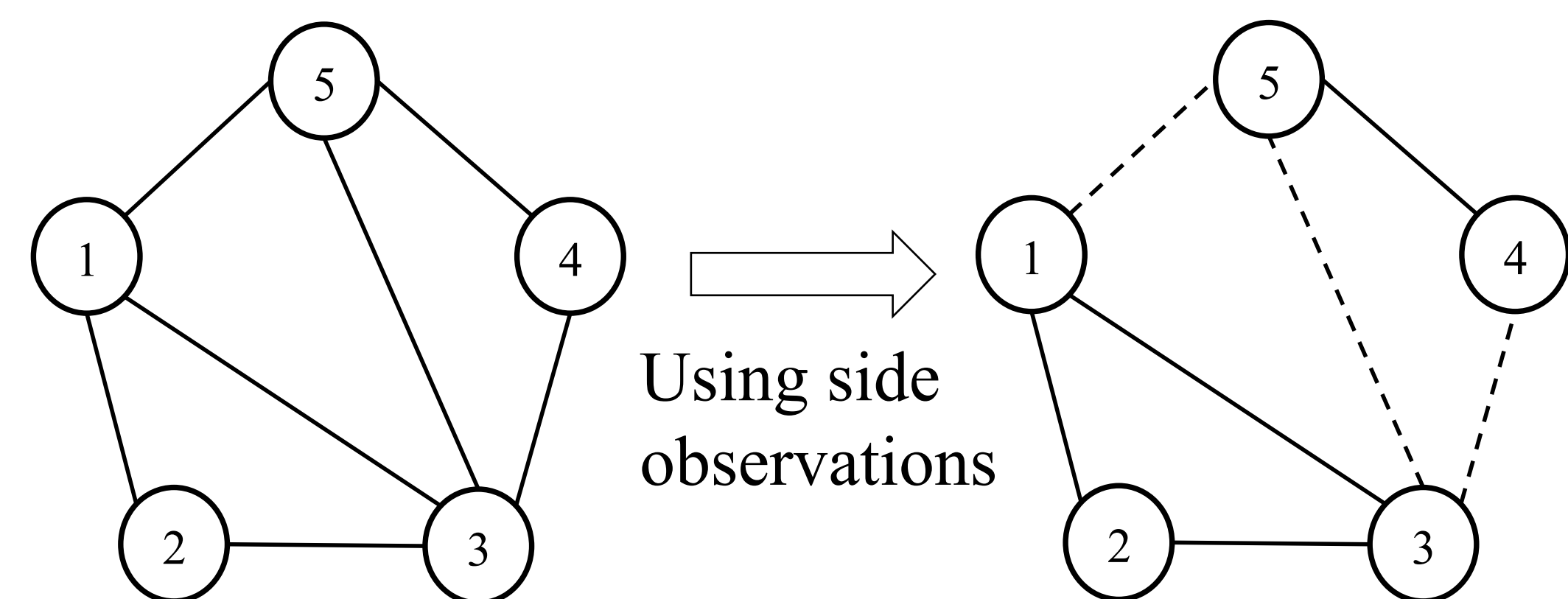
Theorem 2. For any deterministic graph feedback (G_1, G_2, G_3, \dots) , the Bayesian regrets of TS-N, IDS-N and IDSN-LP are upper-bounded by

$$\sqrt{\sum_{t=1}^T \frac{\chi(G_t)}{2} H(\alpha_1)}$$

Theorem 3. For any random graph feedback (r_1, r_2, r_3, \dots) , the Bayesian regrets of TS-N, IDS-N and IDSN-LP are upper-bounded by

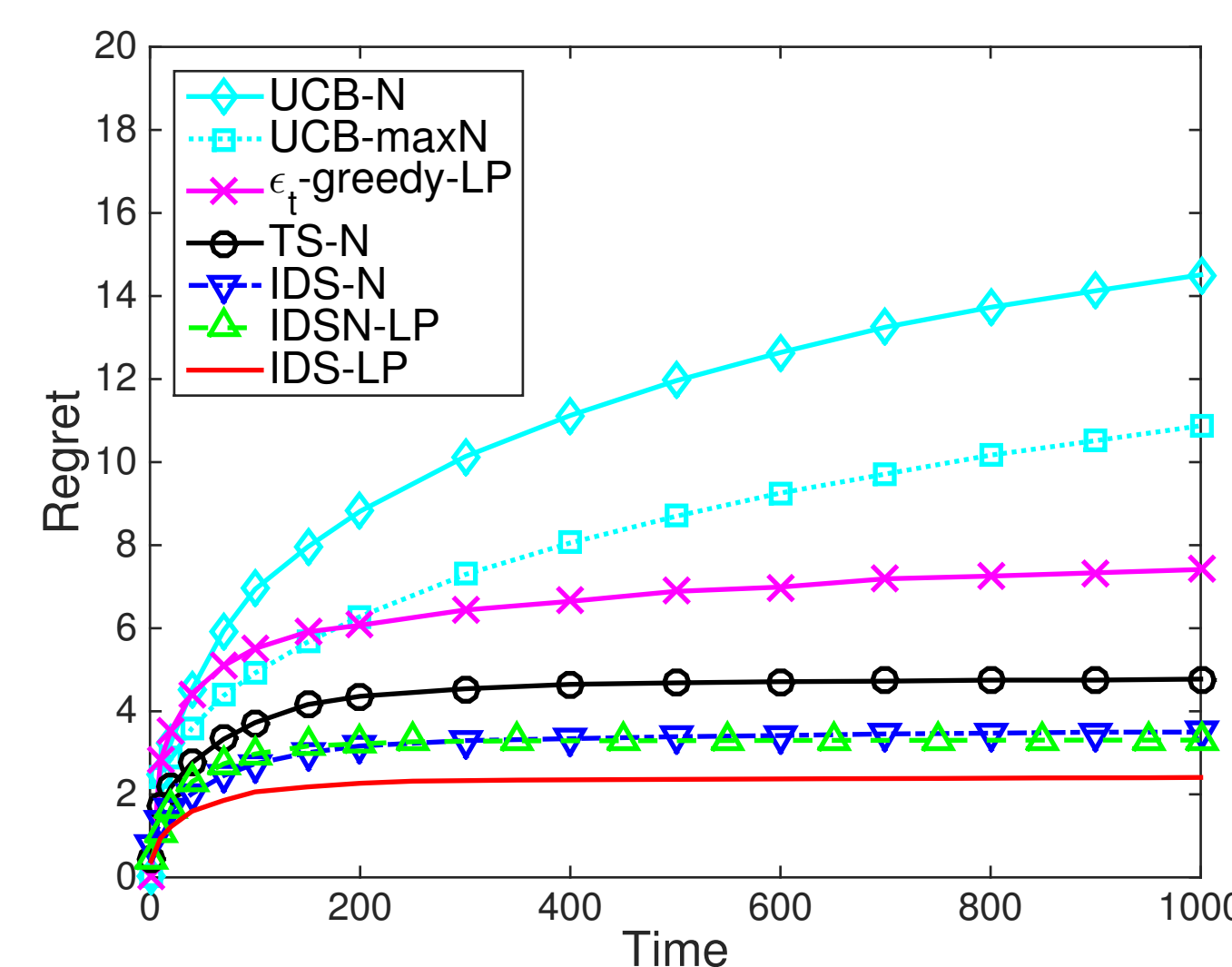
$$\sqrt{\sum_{t=1}^T \frac{K}{2(Kr_t + 1 - r_t)} H(\alpha_1)}$$

Clique cover number, $\chi(G)$, is the smallest cardinality of clique partition.

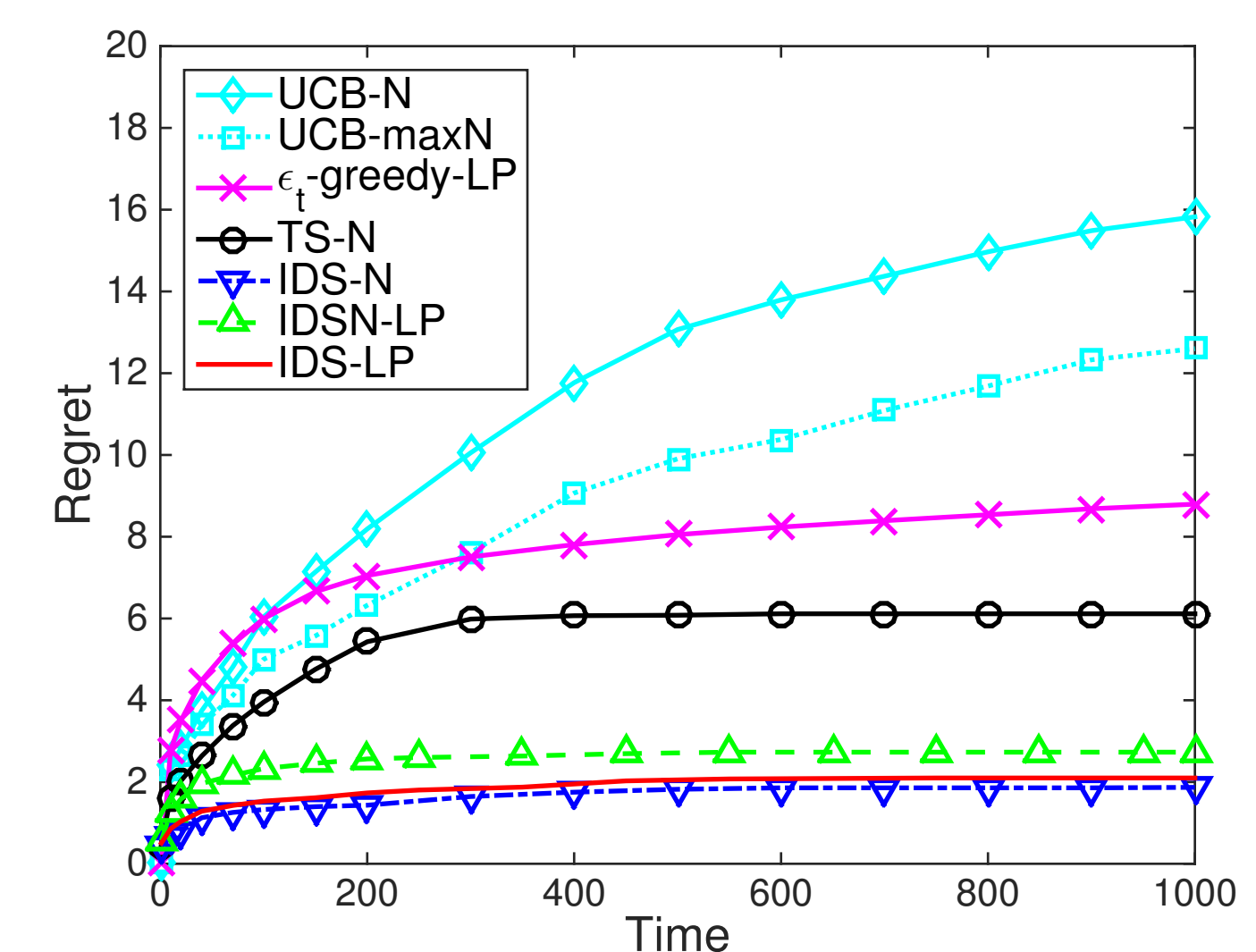


Evaluation

Deterministic graph

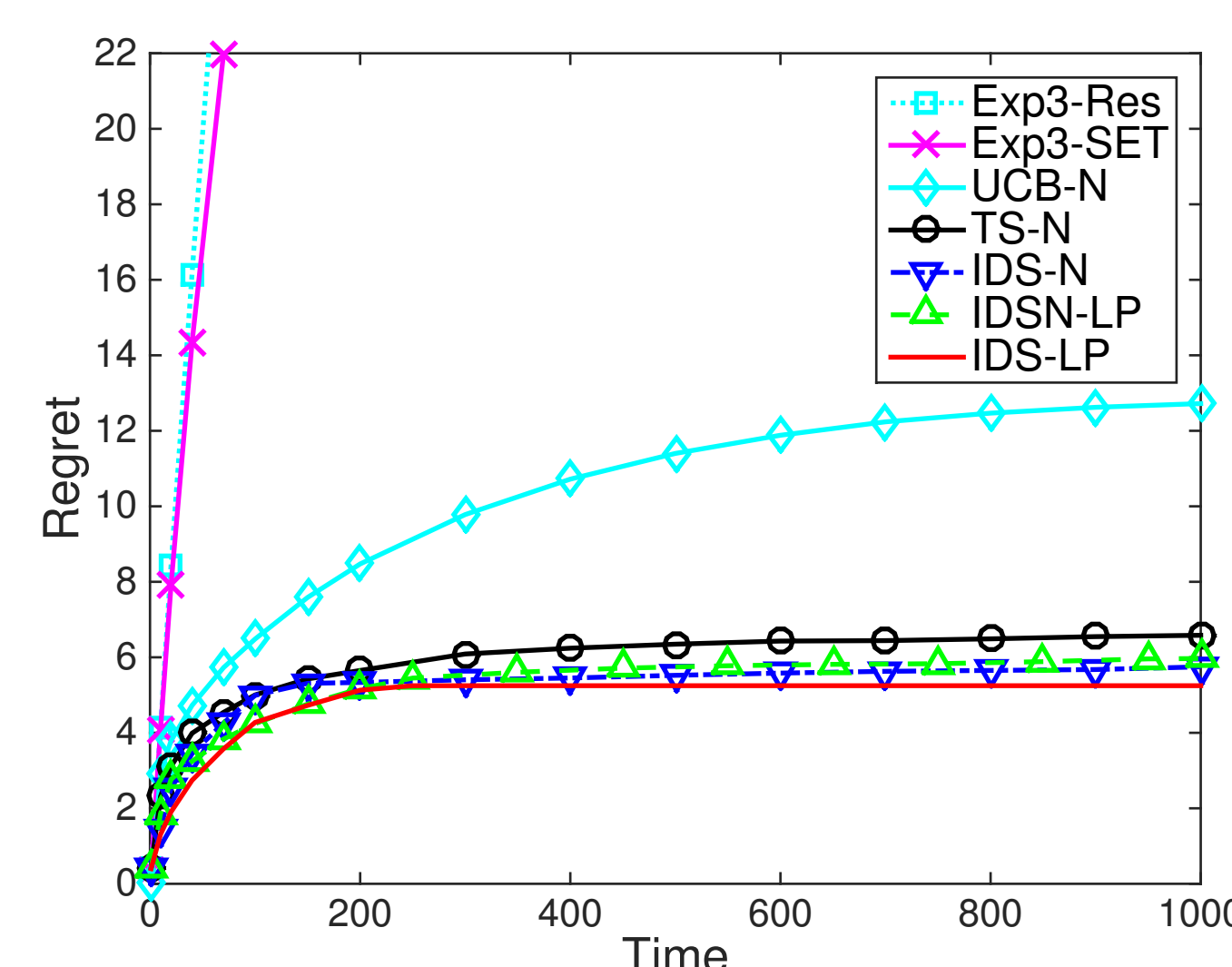


Time-invariant graphs

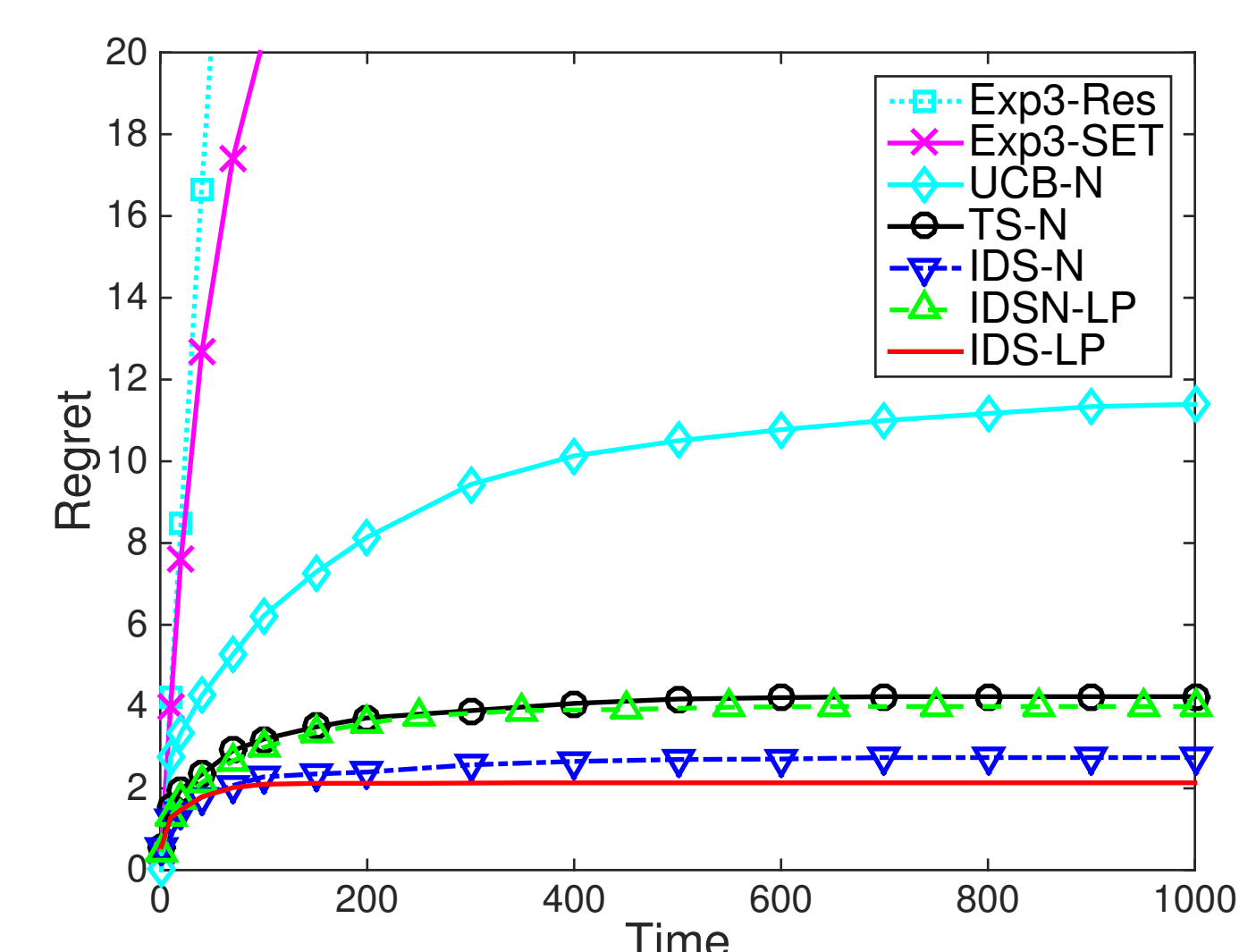


Time-variant graphs

Erdos Renyi random graph



Time-invariant $r=0.25$



Time-variant r_t