

# Fast Mixing for Discrete Point Processes



— YALE INSTITUTE FOR —  
NETWORK SCIENCE

Patrick Rebeschini<sup>1</sup> and Amin Karbasi<sup>2</sup>

<sup>1</sup>patrick.rebeschini@yale.edu <sup>2</sup>amin.karbasi@yale.edu  
Yale Institute for Network Science, Yale University

## 1. Role of submodularity in probability?

- **Combinatorial optimization:** Submodularity extensively studied.

Let  $V$  be a finite set and  $f: S \in 2^V \rightarrow f(S) \in \mathbb{R}$  be a set function.

$$\Delta_i f(S) := f(S \cup \{i\}) - f(S) \quad (\text{gradient of } f)$$

Function  $f$  is **submodular** if for each  $i, j \notin S, i \neq j$

$$\Delta_i \Delta_j f(S) \equiv \Delta_j f(S \cup \{i\}) - \Delta_j f(S) \leq 0 \quad (\text{Hessian of } f)$$

- **Probability:** Submodularity recently investigated to compute  $\mathbb{P}(S \ni i)$  in

$$\mathbb{P}(S = S) := \frac{e^{-\beta f(S)}}{Z}, \quad \beta > 0.$$

**ISSUE:** (Djolonga and Krause, 2014) **bounds exp. bad in  $|V| := \text{card } V$ .**

**Q:** Can we get dimension-free bounds?

**Q:** Is submodularity right notion?

## 2. Fast mixing MCMC: control on Hessian

**GOAL:** Investigate **fundamental** property of  $f$  to get fast mixing MCMC.

Consider local-update (Glauber dynamics type) Markov chains:  
(systematic-scan, Metropolis-Hasting algorithm also considered in the paper)

**ALGORITHM 1.** Random-scan Gibbs sampler

Sample  $S_0 \in 2^V$  from a given distribution (e.g., uniform); Set  $S \leftarrow S_0$ ;

**for**  $s = 1, \dots, t$  **do**

**for**  $|V|$  **times do**

    Sample  $i \in V$  uniformly. Draw  $C \in \{0, 1\}$  with  $\mathbb{P}(C=0) = \frac{1}{1 + \exp \Delta_i f(S \setminus \{i\})}$ ;

    If  $C = 0$  then set  $S \leftarrow S \setminus \{i\}$ , else set  $S \leftarrow S \cup \{i\}$ ;

$S_s \leftarrow S$ ;

**Output:** Markov chain  $S_0, S_1, \dots, S_t$ .

**MAIN RESULT:** For a **generic** set function  $f$ , if

$$\beta \|M\|_\infty \equiv \beta \max_{i \in V} \sum_{j \in V} M_{ij} \leq \gamma < 1 \quad \text{where} \quad M_{ij} \propto \max_{S \in 2^V: S \not\ni i, j} |\Delta_i \Delta_j f(S)|$$

then  $S_0, S_1, \dots, S_t$  is fast mixing (mixing time  $\tau(\epsilon) \leq \left\lceil \frac{\log(|V|\epsilon^{-1})}{1-\gamma} \right\rceil$ ) and

$$\left\| \frac{1}{N} \sum_{k=1}^N \mathbf{1}(S_t^{[k]} \ni i) - \mathbb{P}(S \ni i) \right\|_2 \leq \gamma^t + \frac{1}{\sqrt{N}},$$

where  $S^{[1]}, \dots, S^{[N]}$  are  $N$  independent copies of the Markov chain.

- Proof relies on theory of Dobrushin uniqueness for Gibbs measures.
- **Key result:** Bound does not depend on dimension  $|V|$ .
- **Key property:** **Dimension-free uniform control on Hessian.**
- **Submodularity not enough:** Phase transition as a function of  $\beta$  for convergence rate of Glauber dynamics for Ising model.

**NOTE:** No previous literature on Hessian of set functions.

## 3. Hessian and decay of correlation

**Hessian captures decay of correlations in probability.**

Examples in metric space ( $d$  is metric):

- Exponential decay of correlations:  $\max_{S \in 2^V: S \not\ni i, j} |\Delta_i \Delta_j f(S)| \leq \alpha e^{-\alpha' d(i, j)}$ .
- Finite-range correlations:  $\max_{S \in 2^V: S \not\ni i, j} |\Delta_i \Delta_j f(S)| \leq \begin{cases} c & \text{if } d(i, j) \leq r, \\ 0 & \text{if } d(i, j) > r. \end{cases}$

## 4. Hessian and curvature

- Many results in submodular optimization for *monotone* functions (i.e.,  $\Delta_i f(S) \geq 0$  for each  $i, S$ ) rely on notion of **curvature** (based on **gradient**):

$$c := 1 - \min_{i \in V} \frac{\min_{S \in 2^V: S \not\ni i} \Delta_i f(S)}{f(\{i\})} = 1 - \min_{i \in V} \frac{\Delta_i f(V \setminus \{i\})}{f(\{i\})} \in [0, 1].$$

We have  $c = 0$  if and only if function is *modular*, i.e.,  $f(S) = \sum_{i \in S} w_i$ .

Curvature is convenient as easy to compute (minimum is over  $|V|$  terms).

- **Hessian** is a **more natural concept to characterize "curvature"**.
- **Hessian** also **captures locality**.

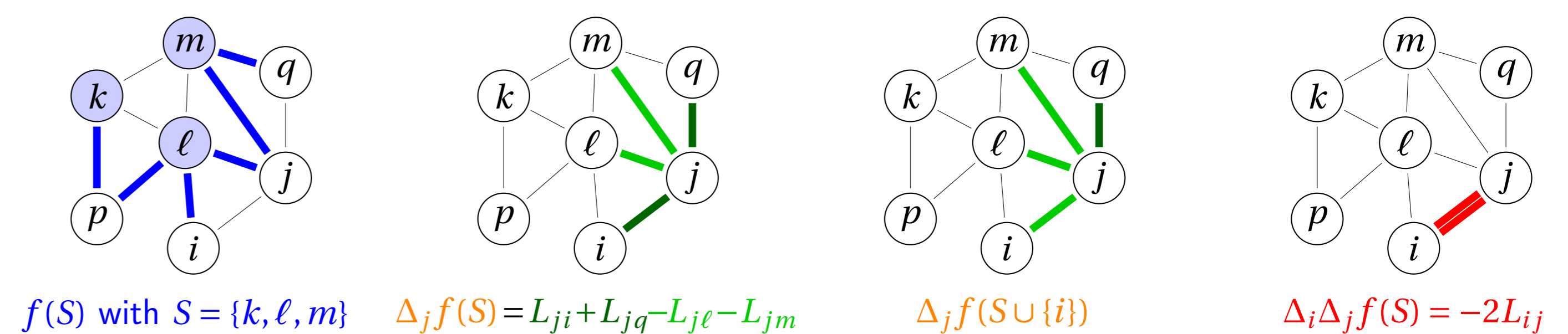
$$\begin{pmatrix} -1 & -c & -c \\ -c & \dots & -c \\ -c & -c & -1 \end{pmatrix} \leq \frac{\Delta_i \Delta_j f(S)}{f(\{i\}) \wedge f(\{j\})} \leq \begin{pmatrix} c-1 & 0 & 0 \\ 0 & \dots & 0 \\ 0 & 0 & c-1 \end{pmatrix}$$

- In general  $\max_{S \in 2^V: S \not\ni i, j} |\Delta_i \Delta_j f(S)|$  is not easy to compute (maximum is over  $2^{|V|-2}$  terms). However:

- In many canonical applications (cut function, coverage function, etc.) Hessian is **sparse** and can be **easily computed or uniformly bounded**.
- In other applications (e.g., determinantal point processes) more assumptions are needed to uniformly control Hessian.

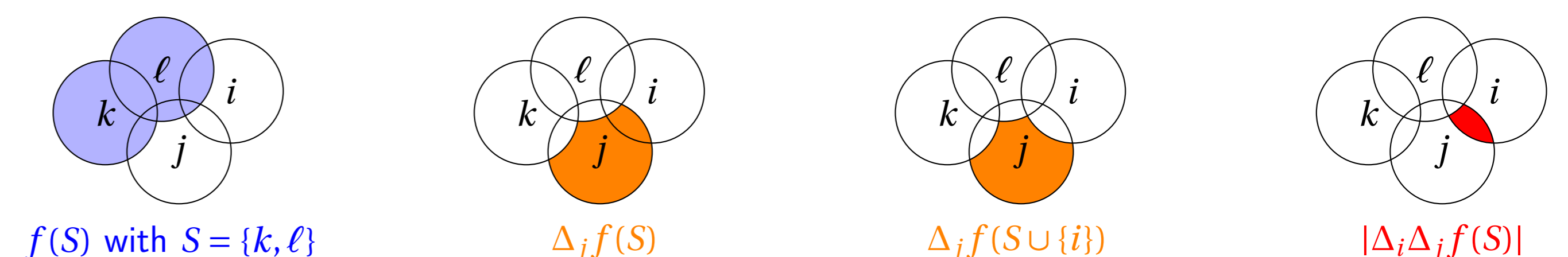
## 5. Cut function

- $(V, E)$  complete graph.  $L_{ij} = L_{ji} \geq 0$  weight associated to edge  $\{i, j\} \in E$ .
- $f(S) = f(V \setminus S) := \sum_{k \in S} \sum_{\ell \in V \setminus S} L_{k\ell}$  and  $f(\emptyset) = f(V) := 0$ .
- $\Delta_i \Delta_j f(S) = -2L_{ij}$  for any  $S \in 2^V$ .



## 6. Coverage function

- $V$  set of points in  $\mathbb{R}^2$ .  $B_i$  ball centered at  $i \in V$ .  $f(S) := \text{Vol}(\cup_{i \in S} B_i)$ .
- $|\Delta_i \Delta_j f(S)| = \text{Vol}(B_i \cap B_j \setminus \cup_{k \in S} B_k) \leq \text{Vol}(B_i \cap B_j)$



## 7. Determinantal point process

- $V = \{1, \dots, n\}$ .  $L \in \mathbb{R}^{n \times n}$  pos. definite.  $(X_v)_{v \in V}$  Gaussian r.v.'s covariance  $L$ .
- $f(S) := \log \det L_S$  where  $L_S := (L_{ij})_{i, j \in S}$  and  $f(\emptyset) := 0$ .
- $\Delta_i \Delta_j f(S) = -2I(X_i; X_j | X_S)$
- **CAVEAT:** Conditional mutual information not monotone in  $S$  ( $\neq$  entropy).

## 8. Back to optimization!

**NEXT STEP:** Use Hessian in **combinatorial optimization**.

- Dimension-free uniform control on Hessian can be exploited to **get fastest convergence rates for ordinary greedy-type algorithms**.
- Work to be posted soon on **arXiv**.