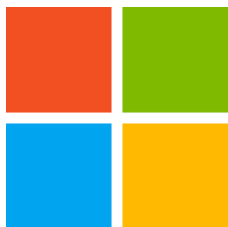


A Matrix Chernoff Bound for Markov Chains and Its Application to Co-occurrence Matrices

Jiezhong Qiu, Chi Wang, Ben Liao, Richard Peng, Jie Tang



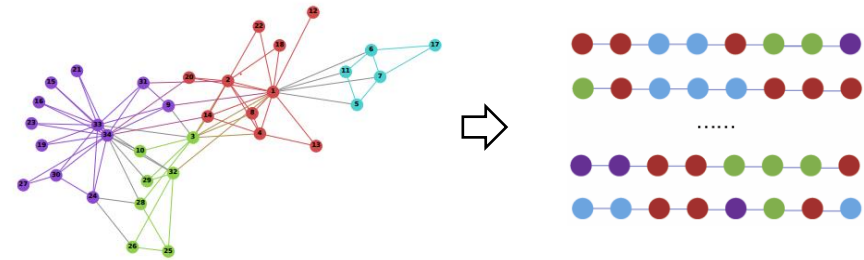
Tencent Quantum Lab
腾讯量子实验室



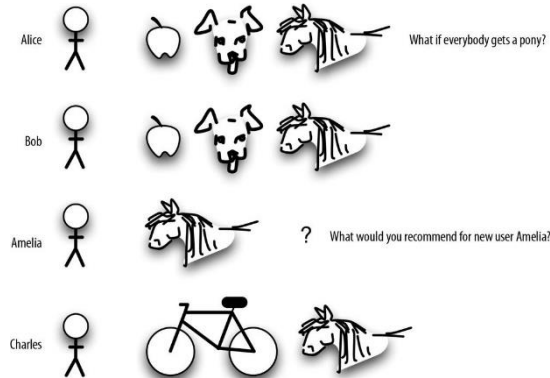
The Application to Co-occurrence Matrices

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

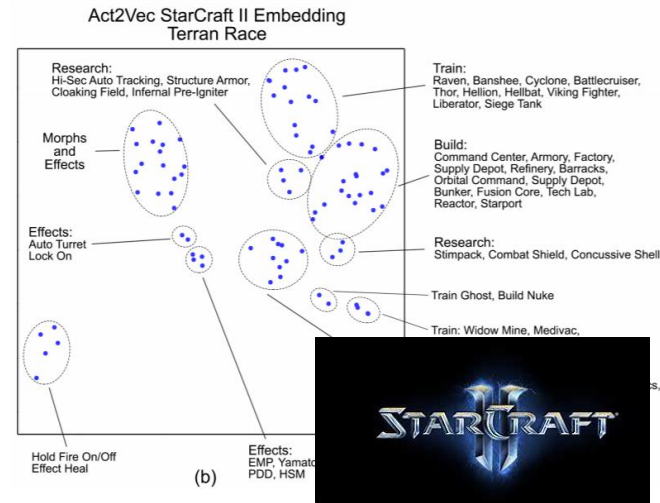
NLP
(LDA, Word2vec, Glove)



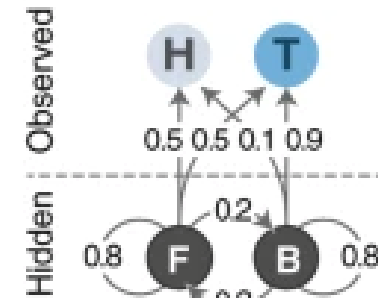
Graph Learning
(DeepWalk, node2vec, metapath2vec)



Recommendation System
(Pin2Vec, Item2vec)



Reinforcement Learning
(Act2Vec)

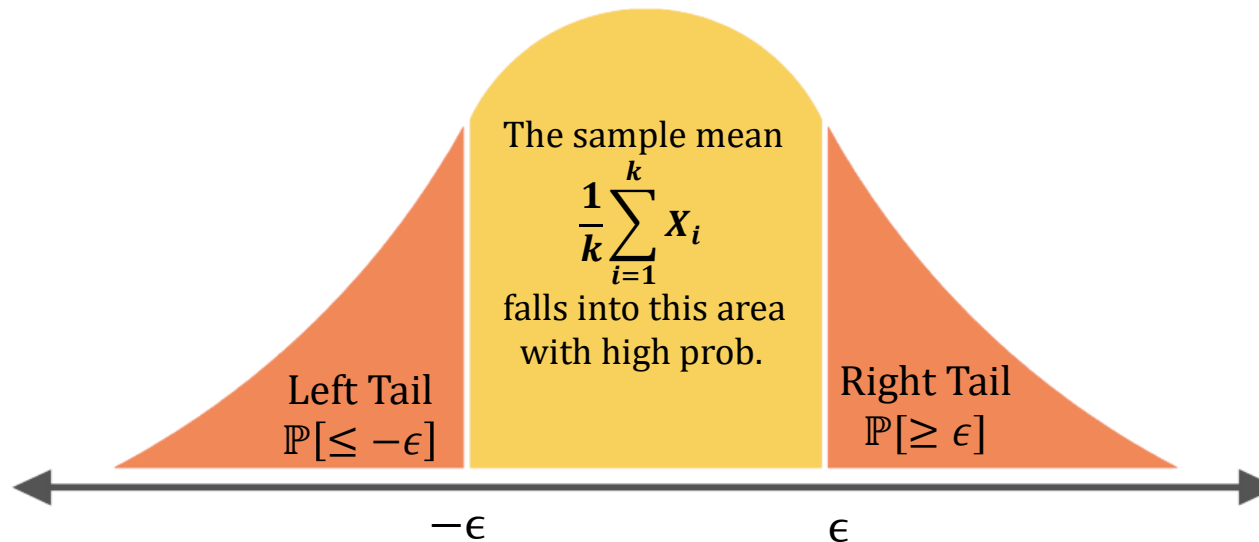


Hidden Markov Models
(Emission Co-occurrence)

Chernoff Bounds

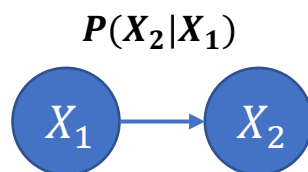
Theorem (Chernoff Bound, 1952): If X_1, X_2, \dots, X_k are **independent** zero-mean **scalar-valued** random variables with $|X_i| \leq 1$. Then for $\epsilon \in (0, 1)$

$$\mathbb{P}\left(\left|\frac{1}{k}\sum_{i=1}^k X_i\right| \geq \epsilon\right) \leq 2\exp(-k\epsilon^2/4)$$

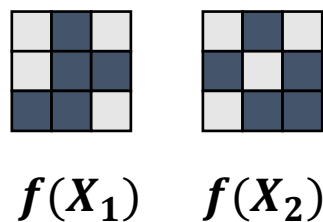


A Matrix Chernoff Bound for Markov Chains

Independence
Markov Dependence



Scalar-valued
Random Variables
Matrix-valued
Random Variables

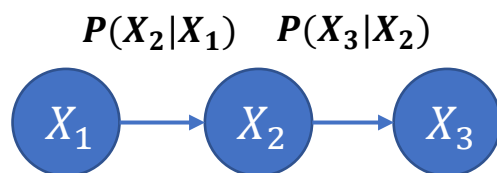


Sample Mean Matrix

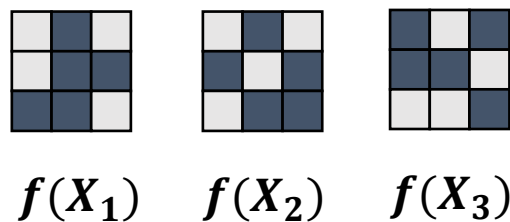
$$\frac{1}{2}(f(X_1) + f(X_2))$$

A Matrix Chernoff Bound for Markov Chains

Independence
Markov Dependence



Scalar-valued
Random Variables
Matrix-valued
Random Variables

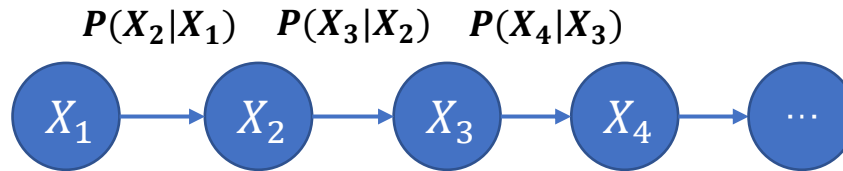


Sample Mean Matrix

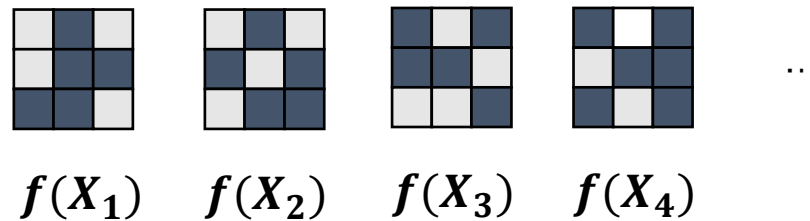
$$\frac{1}{3} (f(X_1) + f(X_2) + f(X_3))$$

A Matrix Chernoff Bound for Markov Chains

Independence
Markov Dependence

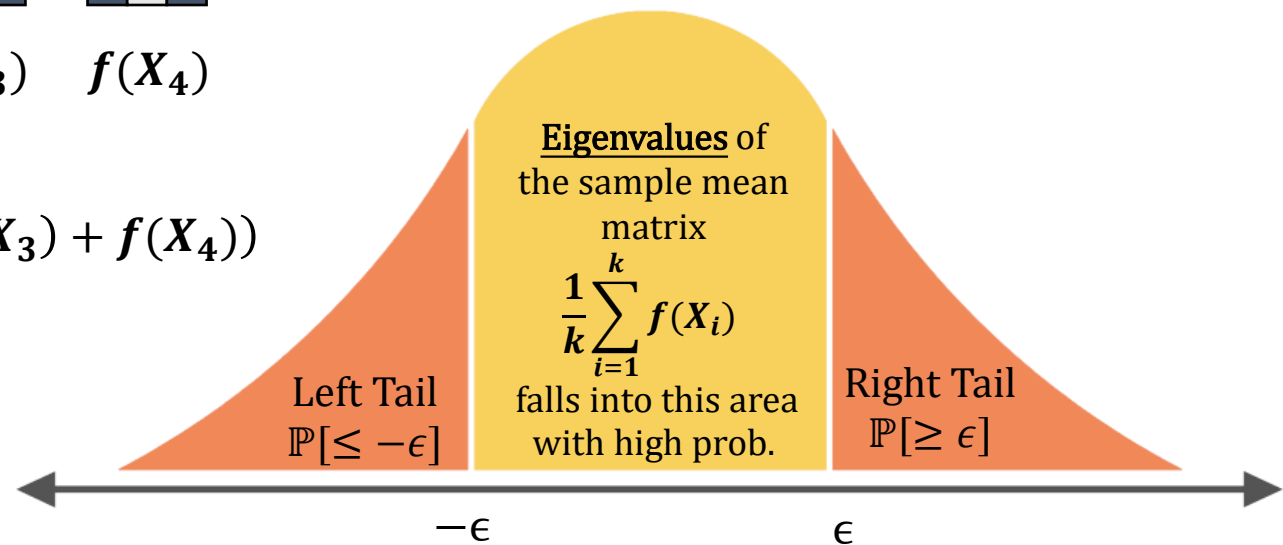


Scalar-valued
Random Variables
Matrix-valued
Random Variables



Sample Mean Matrix

$$\frac{1}{4}(f(X_1) + f(X_2) + f(X_3) + f(X_4))$$



A Matrix Chernoff Bound for Markov Chains

$$\mathbb{P} \left[\lambda_{\min} \left(\frac{1}{k} \sum_{i=1}^k f(X_i) \right) \leq -\epsilon \right] \quad \text{and} \quad \mathbb{P} \left[\lambda_{\max} \left(\frac{1}{k} \sum_{i=1}^k f(X_i) \right) \geq \epsilon \right]$$

Comparison	X	f	Tail Probability
Chernoff '52	i.i.d scalars	Identity	$\exp(-\Omega(k\epsilon^2))$
Tropp'12	i.i.d matrices	identity	$d\exp(-\Omega(k\epsilon^2))$
CLLM'12	Random walk on a regular Markov with spectral expansion λ	$[N] \rightarrow \mathbb{C}$	$\exp(-\Omega(k(1-\lambda)\epsilon^2))$
GLSS'18	Random walk on an undirected regular graph with second eigenvalue λ	$[N] \rightarrow \mathbb{C}^{d \times d}$	$d\exp(-\Omega(k(1-\lambda)\epsilon^2))$
Ours	Random walk on a regular Markov chain with spectral expansion λ	$[N] \rightarrow \mathbb{C}^{d \times d}$	$d\exp(-\Omega(k(1-\lambda)\epsilon^2))$

A Matrix Chernoff Bound for Markov Chains

Theorem: Let P be an regular Markov chain with state space $[N]$, stationary distribution π and spectral expansion λ . Let $f: [N] \rightarrow \mathbb{C}^{d \times d}$ be a matrix-valued function such that

1. $\forall X \in [N], f(X)$ is Hermitian and $\|f(X)\|_2 \leq 1$;
2. $\sum_{X \in [N]} \pi_X f(X) = 0$.

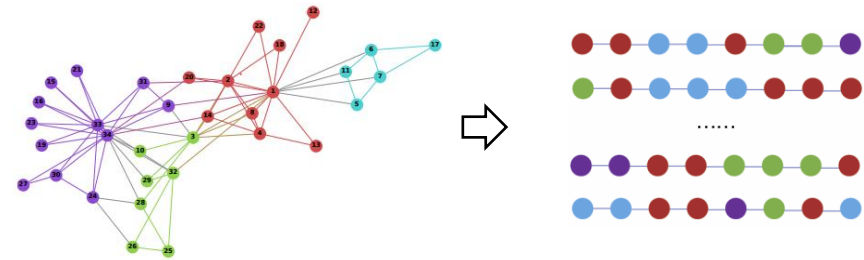
Let (X_1, X_2, \dots, X_k) denote a k -step random walk on P starting from an initial distribution ϕ . Then for $\epsilon \in (0, 1)$:

$$\mathbb{P} \left[\lambda_{\min} \left(\frac{1}{k} \sum_{i=1}^k f(X_i) \right) \leq -\epsilon \right] \leq \|\phi\|_{\pi} d^2 \exp(-k(1 - \lambda)\epsilon^2/72)$$
$$\mathbb{P} \left[\lambda_{\max} \left(\frac{1}{k} \sum_{i=1}^k f(X_i) \right) \geq \epsilon \right] \leq \|\phi\|_{\pi} d^2 \exp(-k(1 - \lambda)\epsilon^2/72)$$

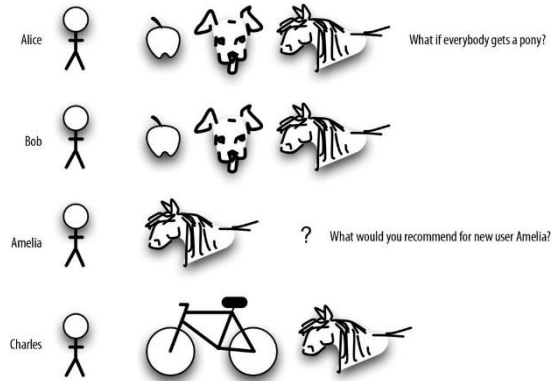
The Application to Co-occurrence Matrices

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

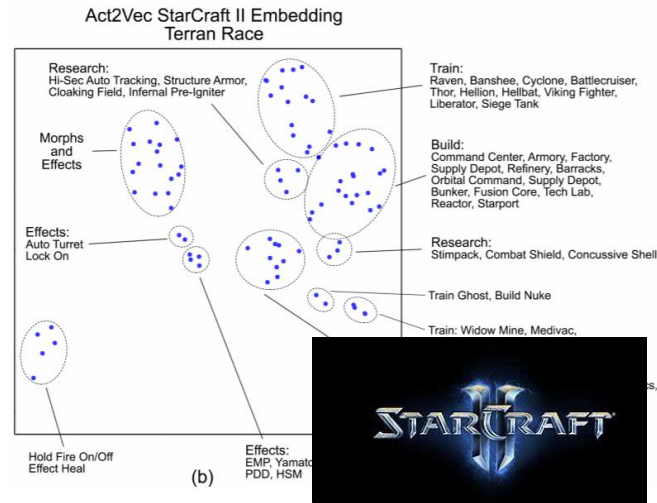
NLP
(LDA, Word2vec, Glove)



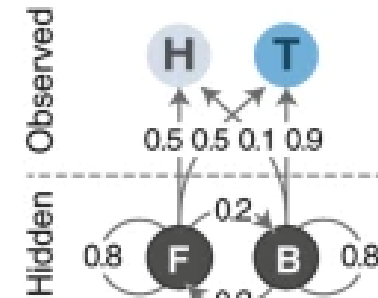
Graph Representation Learning
(DeepWalk, node2vec, metapath2vec)



Recommendation System
(Pin2Vec, Item2vec)



Reinforcement Learning
(Act2Vec)

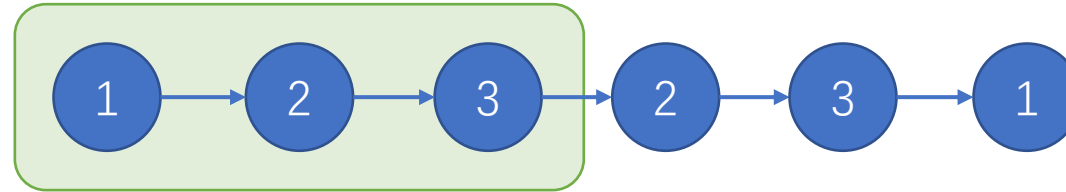


Hidden Markov Models
(Emission Co-occurrence)

Co-occurrence Matrix of Sequential Data

Sliding Window 1

$X_1 = (1,2,3)$

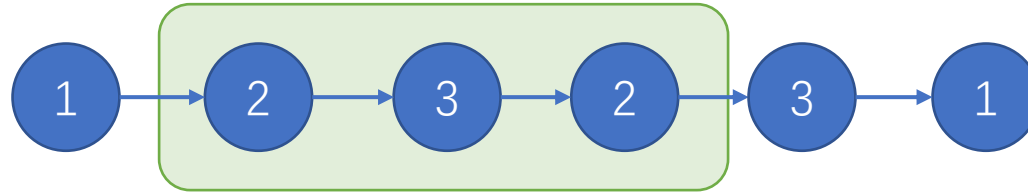


$$\mathbf{C} = \frac{1}{4} \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}$$

Co-occurrence Matrix of Sequential Data

Sliding Window 2

$$X_2 = (2, 3, 2)$$

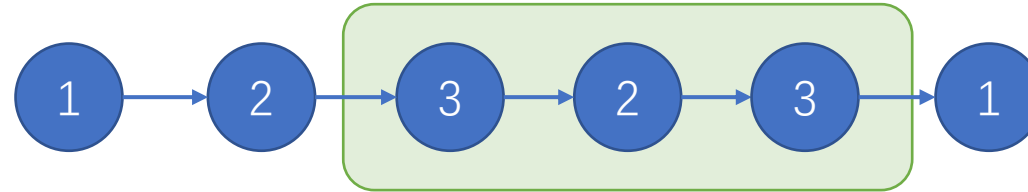


$$\mathbf{C} = \frac{1}{2} \left[\frac{1}{4} \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} + \frac{1}{4} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 0 \end{pmatrix} \right]$$

Co-occurrence Matrix of Sequential Data

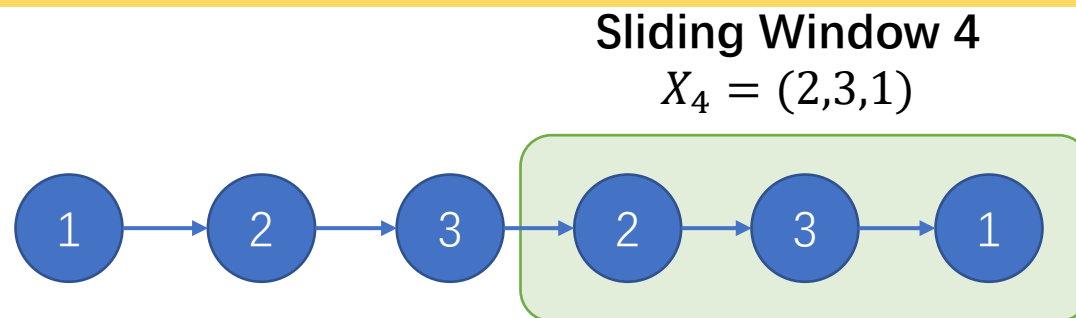
Sliding Window 3

$$X_3 = (3, 2, 3)$$



$$\mathbf{C} = \frac{1}{3} \left[\frac{1}{4} \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} + \frac{1}{4} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 0 \end{pmatrix} + \frac{1}{4} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 2 \end{pmatrix} \right]$$

Markov chain Matrix Chernoff Bound!



$$\begin{aligned} \mathbf{C} &= \frac{1}{4} \left[\frac{1}{4} \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} + \frac{1}{4} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 0 \end{pmatrix} + \frac{1}{4} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 2 \end{pmatrix} + \frac{1}{4} \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \right] \\ &= \frac{1}{4} (\mathbf{f}(X_1) + \mathbf{f}(X_2) + \mathbf{f}(X_3) + \mathbf{f}(X_4)) \end{aligned}$$

Observation 1:

Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_{L-T}$ be the sequence of sliding windows, and \mathbf{f} maps a sliding window to the co-occurrence matrix within this window. The co-occurrence matrix \mathbf{C} can be written as the **sample mean** of $\mathbf{f}(\mathbf{X}_1), \mathbf{f}(\mathbf{X}_2), \dots, \mathbf{f}(\mathbf{X}_{L-T})$:

$$\mathbf{C} = \frac{1}{L-T} \sum_{k=1}^{L-T} \mathbf{f}(\mathbf{X}_k)$$

Observation 2: If the input sequence v_1, v_2, \dots is a Markov Chain, then $\mathbf{X}_1, \mathbf{X}_2, \dots$ is a Markov Chain, too.

Convergence Rate of Co-occurrence Matrices

- The co-occurrence matrix:

$$\mathcal{C} = \frac{1}{L-T} \sum_{k=1}^{L-T} f(X_k)$$

- The asymptotic expectation of \mathcal{C} (denote $\Pi = \text{diag}(\pi)$):

$$\mathbb{A}\mathbb{E}[\mathcal{C}] = \lim_{L \rightarrow +\infty} \mathbb{E}[\mathcal{C}] = \sum_{r=1}^T \frac{1}{2T} (\Pi P^r + (\Pi P^r)^\top)$$

Theorem: Let P be a regular Markov chain with state space $[n]$, stationary distribution π and mixing time τ . Let (v_1, \dots, v_L) be a L -step random walk on P starting from a distribution ϕ . Given $\epsilon \in (0, 1)$, the probability that the co-occurrence matrix \mathcal{C} deviates from its asymptotic expectation $\mathbb{A}\mathbb{E}[\mathcal{C}]$ (in 2-norm) is bounded by:

$$\mathbb{P}(\|\mathcal{C} - \mathbb{A}\mathbb{E}[\mathcal{C}]\|_2 \geq \epsilon) \leq 2(\tau + T) \|\phi\|_\pi n^2 \exp\left(-\frac{\epsilon^2(L-T)}{576(\tau + T)}\right)$$

Roughly, one needs $L = O(\tau(\log n + \log \tau)/\epsilon^2)$ samples to guarantee good estimation to the co-occurrence matrix.

Experiments

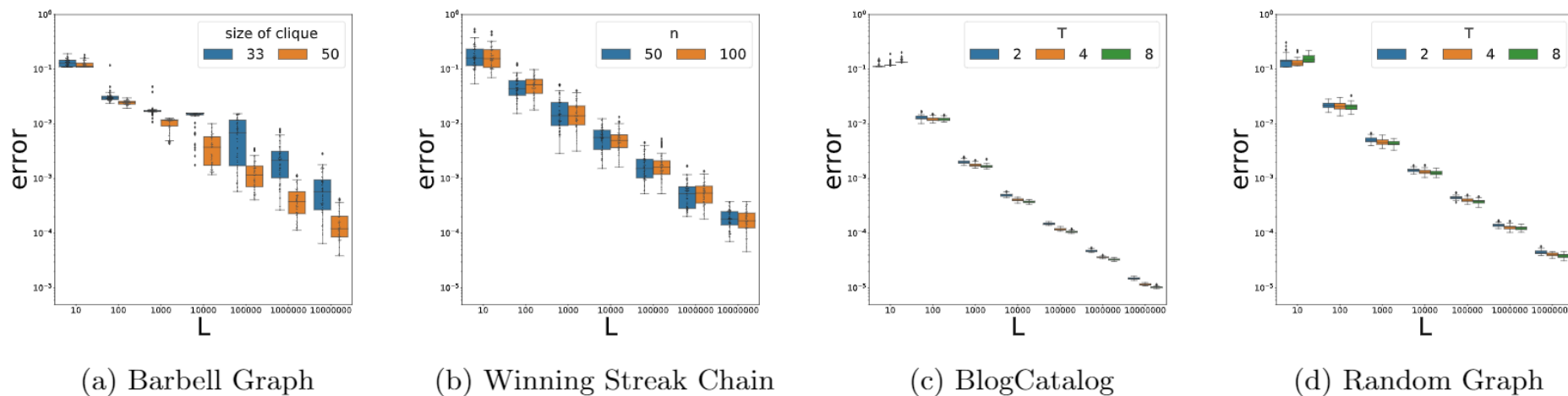


Figure 1: The convergence rate of co-occurrence matrices on Barbell graph, winning streak chain, BlogCatalog graph, and random graph (in log-log scale). The x -axis is the trajectory length L and the y -axis is the approximation error $\|\mathbf{C} - \mathbb{A}\mathbb{E}[\mathbf{C}]\|_2$. Each experiment contains 64 trials, and the error bar is presented.

Thanks!

<https://arxiv.org/abs/2008.02464>