

436 **A Basic Facts about Matrix Norms**

437 In this section, we list some basic facts about matrix norms that will be helpful in comprehending the
 438 subsequent proofs.

439 **A.1 Matrix norms induced by vector norms**

440 Suppose a vector norm $\|\cdot\|_\alpha$ on \mathbb{R}^n and a vector norm $\|\cdot\|_\beta$ on \mathbb{R}^m are given. Any matrix $M \in \mathbb{R}^{m \times n}$
 441 induces a linear operator from \mathbb{R}^n to \mathbb{R}^m with respect to the standard basis, and one defines the
 442 corresponding *induced norm* or *operator norm* by

$$\|M\|_{\alpha,\beta} = \sup \left\{ \frac{\|Mv\|_\beta}{\|v\|_\alpha}, v \in \mathbb{R}^n, v \neq \mathbf{0} \right\}.$$

443 If the p -norm for vectors ($1 \leq p \leq \infty$) is used for both spaces \mathbb{R}^n and \mathbb{R}^m , then the corresponding
 444 operator norm is

$$\|M\|_p = \sup_{v \neq \mathbf{0}} \frac{\|Mv\|_p}{\|v\|_p}.$$

445 The matrix 1-norm and ∞ -norm can be computed by

$$\|M\|_1 = \max_{1 \leq j \leq m} \sum_{i=1}^n |M_{ij}|,$$

446 that is, the maximum absolute column sum of the matrix M ;

$$\|M\|_\infty = \max_{1 \leq m \leq n} \sum_{j=1}^m |M_{ij}|,$$

447 that is, the maximum absolute row sum of the matrix M .

448 **Remark** In the special case of $p = 2$, the induced matrix norm $\|\cdot\|_2$ is called the *spectral norm*,
 449 and is equal to the largest singular value of the matrix.

450 For square matrices, we note that the name “spectral norm” does not imply the quantity is directly
 451 related to the spectrum of a matrix, unless the matrix is symmetric.

Example We give the following example of a stochastic matrix P , whose spectral radius is 1, but
 its spectral norm is greater than 1.

$$P = \begin{bmatrix} 0.9 & 0.1 \\ 0.25 & 0.75 \end{bmatrix} \quad \|P\|_2 \approx 1.0188$$

452 **A.2 Matrix (p, q) -norms**

453 The Frobenius norm of a matrix $M \in \mathbb{R}^{m \times n}$ is defined as

$$\|M\|_F = \sqrt{\sum_{j=1}^n \sum_{i=1}^m |M_{ij}|^2},$$

454 and it belongs to a family of entry-wise matrix norms: for $1 \leq p, q \leq \infty$, the matrix (p, q) -norm is
 455 defined as

$$\|M\|_{p,q} = \left(\sum_{j=1}^n \left(\sum_{i=1}^m |M_{ij}|^p \right)^{q/p} \right)^{1/q}.$$

456 The special case $p = q = 2$ is the Frobenius norm $\|\cdot\|_F$, and $p = q = \infty$ yields the max norm
 457 $\|\cdot\|_{\max}$.

458 **A.3 Equivalence of norms**

459 For any two matrix norms $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$, we have that for all matrices $M \in \mathbb{R}^{m \times n}$,

$$r\|M\|_\alpha \leq \|M\|_\beta \leq s\|M\|_\alpha$$

460 for some positive numbers r and s . In particular, the following inequality holds for the 2-norm $\|\cdot\|_2$
461 and the ∞ -norm $\|\cdot\|_\infty$:

$$\frac{1}{\sqrt{n}}\|M\|_\infty \leq \|M\|_2 \leq \sqrt{m}\|M\|_\infty.$$

462 **B Proof of Proposition 1**

463 It is straightforward to check that $\|X - \mathbf{1}\gamma_X\|_F$ satisfies the two axioms of a node similarity measure:

464 1. $\|X - \mathbf{1}\gamma_X\|_F = 0 \iff X = \mathbf{1}\gamma_X \iff X_i = \gamma_X$ for all node i .

465 2. Let $\gamma_X = \frac{\mathbf{1}^\top X}{N}$ and $\gamma_Y = \frac{\mathbf{1}^\top Y}{N}$, then $\gamma_X + \gamma_Y = \frac{\mathbf{1}^\top (X+Y)}{N} = \gamma_{X+Y}$. So

$$\begin{aligned} \mu(X+Y) &= \|(X+Y) - \mathbf{1}(\gamma_X + \gamma_Y)\|_F = \|X - \mathbf{1}\gamma_X + Y - \mathbf{1}\gamma_Y\|_F \\ &\leq \|X - \mathbf{1}\gamma_X\|_F + \|Y - \mathbf{1}\gamma_Y\|_F \\ &= \mu(X) + \mu(Y). \end{aligned}$$

466 **C Proof of Lemma 1**

467 According to the formulation (6):

$$X_{\cdot i}^{(t+1)} = \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t W_{j_{k+1}j_k}^{(k)} \right) D_{j_{t+1}}^{(t)} P^{(t)} \dots D_{j_1}^{(0)} P^{(0)} X_{j_0}^{(0)},$$

468 we thus obtain that

$$\begin{aligned} \|X_{\cdot i}^{(t+1)}\|_\infty &= \left\| \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t W_{j_{k+1}j_k}^{(k)} \right) D_{j_{t+1}}^{(t)} P^{(t)} \dots D_{j_1}^{(0)} P^{(0)} X_{j_0}^{(0)} \right\|_\infty \\ &\leq \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t |W_{j_{k+1}j_k}^{(k)}| \right) \|D_{j_{t+1}}^{(t)} P^{(t)} \dots D_{j_1}^{(0)} P^{(0)}\|_\infty \|X_{j_0}^{(0)}\|_\infty \\ &\leq \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t |W_{j_{k+1}j_k}^{(k)}| \right) \|X_{j_0}^{(0)}\|_\infty \\ &\leq C_0 \left(\sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t |W_{j_{k+1}j_k}^{(k)}| \right) \right) \\ &= C_0 \|(|W^{(0)}| \dots |W^{(t)}|)_{\cdot i}\|_1, \end{aligned}$$

469 where C_0 equals the maximal entry in $|X^{(0)}|$.

The assumption **A3** implies that there exists $C' > 0$ such that for all $t \in \mathbb{N}_{\geq 0}$ and $i \in [d]$,

$$\|(|W^{(0)}| \dots |W^{(t)}|)_{\cdot i}\|_1 \leq C' N.$$

Hence there exists $C'' > 0$ such that for all $t \in \mathbb{N}_{\geq 0}$ and $i \in [d]$, we have

$$\|\bar{X}_{\cdot i}^{(t)}\|_\infty \leq C'',$$

470 proving the existence of $C > 0$ such that $\|X^{(t)}\|_{\max} \leq C$ for all $t \in \mathbb{N}_{\geq 0}$.

471 **D Proof of Lemma 2**

472 Lemma 2 is a direct corollary of Lemma 1 and the assumption that $\Psi(\cdot, \cdot)$ assigns bounded attention
473 scores to bounded inputs.

474 **E Proof of Lemma 3**

475 **E.1 Auxiliary results**

476 We make use of the following sufficient condition for the ergodicity of the infinite products of
477 row-stochastic matrices.

478 **Lemma 7** (Corollary 5.1 [2]). *Consider a sequence of row-stochastic matrices $\{S^{(t)}\}_{t=0}^\infty$. Let a_t
479 and b_t be the smallest and largest entries in $S^{(t)}$, respectively. If $\sum_{t=0}^\infty \frac{a_t}{b_t} = \infty$, then $\{S^{(t)}\}_{t=0}^\infty$ is
480 ergodic.*

In order to make use of the above result, we first show that long products of $P^{(t)}$'s from $\mathcal{P}_{\mathcal{G}, \epsilon}$ will eventually become strictly positive. For $t_0 \leq t_1$, we denote

$$P^{(t_1:t_0)} = P^{(t_1)} \dots P^{(t_0)}.$$

481 **Lemma 8.** *Under the assumption **A1**, there exist $T \in \mathbb{N}$ and $c > 0$ such that for all $t_0 \geq 0$,*

$$c \leq P_{ij}^{(t_0+T:t_0)} \leq 1, \forall 1 \leq i, j \leq N.$$

482

483 *Proof.* Fix any $T \in \mathbb{N}_{\geq 0}$. Since $\|P^{(t)}\|_\infty \leq 1$ for any $P^{(t)} \in \mathcal{P}_{\mathcal{G}, \epsilon}$, it follows that $\|P^{(t_0+T:t_0)}\|_\infty \leq$
484 1 and hence $P_{ij}^{(t_0+T:t_0)} \leq 1$, for all $1 \leq i, j \leq N$.

485 To show the lower bound, without loss of generality, we will show that there exist $T \in \mathbb{N}$ and $c > 0$
486 such that

$$P_{ij}^{(T:0)} \geq c, \forall 1 \leq i, j \leq N.$$

487 Since each $P^{(t)}$ has the same connectivity pattern as the original graph \mathcal{G} , it follows from the
488 assumption **A1** that there exists $T \in \mathbb{N}$ such that $P^{(T:0)}$ is a positive matrix, following a similar
489 argument as the one for Proposition 1.7 in [4]: For each pair of nodes i, j , since we assume that the
490 graph \mathcal{G} is connected, there exists $r(i, j)$ such that $P_{ij}^{(r(i,j):0)} > 0$. on the other hand, since we also
491 assume each node has a self-loop, $P_{ii}^{(t:0)} > 0$ for all $t \geq 0$ and hence for $t \geq r(i, j)$,

$$P_{ij}^{(t:0)} \geq P_{ii}^{(t-r(i,j))} P_{ij}^{(r(i,j):0)} > 0.$$

492 For $t \geq t(i) := \max_{j \in \mathcal{G}} r(i, j)$, we have $P_{ij}^{(t:0)} > 0$ for all node j in \mathcal{G} . Finally, if $t \geq T := \max_{i \in \mathcal{G}} t(i)$,
493 then $P_{ij}^{(t:0)} > 0$ for all pairs of nodes i, j in \mathcal{G} . Notice that $P_{ij}^{(T:0)}$ is a weighted sum of walks of
494 length T between nodes i and j , and hence $P_{ij}^{(T:0)} > 0$ if and only if there exists a walk of length
495 T between nodes i and j . Since for all $t \in \mathbb{N}_{\geq 0}$, $P_{ij}^{(t)} \geq \epsilon$ if $(i, j) \in E(\mathcal{G})$, we conclude that
496 $P_{ij}^{(T:0)} \geq \epsilon^T := c$. \square

497 **E.2 Proof of Lemma 3**

498 Given the sequence $\{P^{(t)}\}_{t=0}^\infty$, we use $T \in \mathbb{N}$ from Lemma 8 and define

$$\bar{P}^{(k)} := P^{((k+1)T:kT)}.$$

499 Then $\{P^{(t)}\}_{t=0}^\infty$ is ergodic if and only if $\{\bar{P}^{(k)}\}_{k=0}^\infty$ is ergodic. Notice that by Lemma 8, for all
500 $k \in \mathbb{N}_{\geq 0}$, there exists $c > 0$ such that $c \leq \bar{P}_{ij}^{(k)} \leq 1, \forall 1 \leq i, j \leq N$. Then Lemma 3 is a direct
501 consequence of Lemma 7.

502 **F Proof of Lemma 5**

503 **F.1 Notations and auxiliary results**

504 Consider a sequence $\{D^{(t)}P^{(t)}\}_{t=0}^{\infty}$ in $\mathcal{M}_{\mathcal{G},\epsilon}$. For $t_0 \leq t_1$, define

$$Q_{t_0,t_1} := D^{(t_1)}P^{(t_1)} \dots D^{(t_0)}P^{(t_0)}$$

505 and

$$\delta_t = \|D^{(t)} - I_N\|_{\infty},$$

506 where I_N denotes the $N \times N$ identity matrix. It is also useful to define

$$\begin{aligned} \hat{Q}_{t_0,t_1} &:= P^{(t_1)}Q_{t_0,t_1-1} \\ &:= P^{(t_1)}D^{(t_1-1)}P^{(t_1-1)} \dots D^{(t_0)}P^{(t_0)}. \end{aligned}$$

507 We start by proving the following key lemma, which states that long products of matrices in $\mathcal{M}_{\mathcal{G},\epsilon}$
508 eventually become a contraction in ∞ -norm.

509 **Lemma 9.** *There exist $0 < c < 1$ and $T \in \mathbb{N}$ such that for all $t_0 \leq t_1$,*

$$\|\hat{Q}_{t_0,t_1+T}\|_{\infty} \leq (1 - c\delta_{t_1})\|\hat{Q}_{t_0,t_1}\|_{\infty}.$$

510

511 *Proof.* First observe that for every $T \geq 0$,

$$\begin{aligned} \|\hat{Q}_{t_0,t_1+T}\|_{\infty} &\leq \|P^{(t_1+T)}D^{(t_1+T-1)}P^{(t_1+T-1)} \dots D^{(t_1+1)}P^{(t_1+1)}D^{(t_1)}\|_{\infty}\|\hat{Q}_{t_0,t_1}\|_{\infty} \\ &\leq \|P^{(t_1+T)}P^{(t_1+T-1)} \dots P^{(t_1+1)}D^{(t_1)}\|_{\infty}\|\hat{Q}_{t_0,t_1}\|_{\infty}, \end{aligned}$$

512 where the second inequality is based on the following element-wise inequality:

$$P^{(t_1+T)}P^{(t_1+T-1)} \dots D^{(t_1+1)}P^{(t_1+1)} \leq_{\text{ew}} P^{(t_1+T)}P^{(t_1+T-1)} \dots P^{(t_1+1)}.$$

513 By Lemma 8, there exist $T \in \mathbb{N}$ and $0 < c < 1$ such that

$$(P^{(t_1+T)} \dots P^{(t_1+1)})_{ij} \geq c, \forall 1 \leq i, j \leq N.$$

514 Since the matrix product $P^{(t_1+T)}P^{(t_1+T-1)} \dots P^{(t_1+1)}$ is row-stochastic, multiplying it with the
515 diagonal matrix $D^{(t_1)}$ from right decreases the row sums by at least $c(1 - D_{\min}^{(t_1)}) = c\delta_{t_1}$, where
516 $D_{\min}^{(t_1)}$ here denotes the smallest diagonal entry of the diagonal matrix $D^{(t_1)}$. Hence,

$$\|P^{(t_1+T)}P^{(t_1+T-1)} \dots P^{(t_1+1)}D^{(t_1)}\|_{\infty} \leq 1 - c\delta_{t_1}.$$

517

□

518 **F.2 Proof of Lemma 4**

519 Now define $\beta_k := \prod_{t=0}^k (1 - c\delta_t)$ and let $\beta := \lim_{k \rightarrow \infty} \beta_k$. Note that β is well-defined because the
520 partial product is non-increasing and bounded from below. Then we present the following result,
521 which is stated as Lemma 4 in the main paper and from which the ergodicity of any sequence in
522 $\mathcal{M}_{\mathcal{G},\epsilon}$ is an immediate result.

523 **Lemma 4.** Let $\beta_k := \prod_{t=0}^k (1 - c\delta_t)$ and $\beta := \lim_{k \rightarrow \infty} \beta_k$.

- 524 1. If $\beta = 0$, then $\lim_{k \rightarrow \infty} Q_{0,k} = 0$;
- 525 2. If $\beta > 0$, then $\lim_{k \rightarrow \infty} BQ_{0,k} = 0$.

526 *Proof.* We will prove the two cases separately.

527 **[Case $\beta = 0$]** We will show that $\beta = 0$ implies $\lim_{k \rightarrow \infty} \|\hat{Q}_{0,k}\|_\infty = 0$, and as a result,
 528 $\lim_{k \rightarrow \infty} \|Q_{0,k}\|_\infty = 0$. For $0 \leq j \leq T - 1$, let us define

$$\beta^j := \prod_{k=0}^{\infty} (1 - \delta_{j+kT}).$$

529 Then by Lemma 9, we get that

$$\lim_{k \rightarrow \infty} \|\hat{Q}_{0,kT}\|_\infty \leq \beta^j \|\hat{Q}_{0,j}\|_\infty.$$

530 By construction, $\beta = \prod_{j=0}^{T-1} \beta^j$. Hence, if $\beta = 0$ then $\beta^{j_0} = 0$ for some $0 \leq j_0 \leq T - 1$, which
 531 yields $\lim_{k \rightarrow \infty} \|\hat{Q}_{0,k}\|_\infty = 0$. Consequently, $\lim_{k \rightarrow \infty} \|Q_{0,k}\|_\infty = 0$ implies that $\lim_{k \rightarrow \infty} Q_{0,k} = 0$.

532 **[Case $\beta > 0$]** First observe that if $\beta > 0$, then $\forall 0 < \eta < 1$, there exist $m \in \mathbb{N}_{\geq 0}$ such that

$$\prod_{t=m}^{\infty} (1 - c\delta_t) > 1 - \eta. \quad (8)$$

533 Using $1 - x \leq e^{-x}$ for all $x \in \mathbb{R}$, we deduce

$$\prod_{t=m}^{\infty} e^{-c\delta_t} > 1 - \eta.$$

534 It also follows from (8) that $1 - c\delta_t > 1 - \eta$, or equivalently $\delta_t < \frac{\eta}{c}$ for $t \geq m$. Choosing $\eta < \frac{c}{2}$
 535 thus ensures that $\delta_t < \frac{1}{2}$ for $t \geq m$. Putting this together with the fact that, there exists² $b > 0$ such
 536 that $1 - x \geq e^{-bx}$ for all $x \in [0, \frac{1}{2}]$, we obtain

$$\prod_{t=m}^{\infty} (1 - \delta_t) \geq \prod_{t=m}^{\infty} e^{-b\delta_t} > (1 - \eta)^{\frac{b}{c}} := 1 - \eta'. \quad (9)$$

537 Define the product of row-stochastic matrices $P^{(M:m)} := P^{(M)} \dots P^{(m)}$. It is easy to verify the
 538 following element-wise inequality:

$$\left(\prod_{t=m}^M (1 - c\delta_t) \right) P^{(M:m)} \leq_{\text{ew}} Q_{m,M} \leq_{\text{ew}} P^{(M:m)},$$

539 which together with (9) leads to

$$(1 - \eta') P^{(M:m)} \leq_{\text{ew}} Q_{m,M} \leq_{\text{ew}} P^{(M:m)}. \quad (10)$$

540 Therefore,

$$\begin{aligned} \|BQ_{m,M}\|_\infty &= \|B(Q_{m,M} - P^{(M:m)}) + BP^{(M:m)}\|_\infty \\ &\leq \|B(Q_{m,M} - P^{(M:m)})\|_\infty + \|BP^{(M:m)}\|_\infty \\ &= \|B(Q_{m,M} - P^{(M:m)})\|_\infty \\ &\leq \|B\|_\infty \|Q_{m,M} - P^{(M:m)}\|_\infty \\ &\leq \eta' \|B\|_\infty \\ &\leq \eta' \sqrt{N}, \end{aligned}$$

541 where the last inequality is due to the fact that $\|B\|_2 = 1$. By definition, $Q_{0,M} = Q_{m,M} Q_{0,m-1}$,
 542 and hence

$$\|BQ_{0,M}\|_\infty \leq \|BQ_{m,M}\|_\infty \|Q_{0,m-1}\|_\infty \leq \|BQ_{m,M}\|_\infty \leq \eta' \sqrt{N}. \quad (11)$$

²Choose, e.g., $b = 2 \log 2$.

543 The above inequality (11) holds when taking $M \rightarrow \infty$. Then taking $\eta \rightarrow 0$ implies $\eta' \rightarrow 0$ and
 544 together with (11), we conclude that

$$\lim_{M \rightarrow \infty} \|BQ_{0,M}\|_\infty = 0,$$

545 and therefore,

$$\lim_{M \rightarrow \infty} BQ_{0,M} = 0.$$

546

□

547 F.3 Proof of Lemma 5

548 Notice that both cases $\beta = 0$ and $\beta > 0$ in Lemma 4 imply the ergodicity of $\{D^{(t)}P^{(t)}\}_{t=0}^\infty$. Hence
 549 the statement is a direct corollary of Lemma 4.

550 G Proof of Lemma 6

551 In order to show that $\text{JSR}(\tilde{\mathcal{M}}_{\mathcal{G},\epsilon}) < 1$, we start by making the following observation.

552 **Lemma 10.** *A sequence $\{M^{(n)}\}_{n=0}^\infty$ is ergodic if and only if $\prod_{n=0}^t \tilde{M}^{(n)}$ converges to the zero
 553 matrix.*

554 *Proof.* For any $t \in \mathbb{N}_{\geq 0}$, it follows from the third property of the orthogonal projection B (see, Page
 555 6 of the main paper) that

$$B \prod_{n=0}^t M^{(n)} = \prod_{n=0}^t \tilde{M}^{(n)} B.$$

556 Hence

$$\begin{aligned} \{M^{(n)}\}_{n=0}^\infty \text{ is ergodic} &\iff \lim_{t \rightarrow \infty} B \prod_{n=0}^t M^{(n)} = 0 \\ &\iff \lim_{t \rightarrow \infty} \prod_{n=0}^t \tilde{M}^{(n)} B = 0 \\ &\iff \lim_{t \rightarrow \infty} \prod_{n=0}^t \tilde{M}^{(n)} = 0. \end{aligned}$$

557

□

558 Next, we utilize the following result, as a means to ensure a joint spectral radius strictly less than 1
 559 for a bounded set of matrices.

560 **Lemma 11** (Proposition 3.2 in [6]). *For any bounded set of matrices \mathcal{M} , $\text{JSR}(\mathcal{M}) < 1$ if and only
 561 if for any sequence $\{M^{(n)}\}_{n=0}^\infty$ in \mathcal{M} , $\prod_{n=0}^t M^{(n)}$ converges to the zero matrix.*

562 Here, ‘‘bounded’’ means that there exists an upper bound on the norms of the matrices in the set. Note
 563 that $\mathcal{M}_{\mathcal{G},\epsilon}$ is bounded because $\|DP\|_\infty \leq 1$, $DP \in \mathcal{M}_{\mathcal{G},\epsilon}$. To show that $\tilde{\mathcal{M}}_{\mathcal{G},\epsilon}$ is also bounded, let
 564 $\tilde{M} \in \tilde{\mathcal{M}}_{\mathcal{G},\epsilon}$, then by definition, we have

$$\tilde{M}B = BM, M \in \mathcal{M}_{\mathcal{G},\epsilon} \Rightarrow \tilde{M} = BMB^T,$$

565 since $BB^T = I_{N-1}$. As a result,

$$\|\tilde{M}\|_2 = \|BMB^T\|_2 \leq \|M\|_2 \leq \sqrt{N},$$

566 where the first inequality is due to $\|B\|_2 = \|B^T\|_2 = 1$, and the second inequality follows from
 567 $\|M\|_\infty \leq 1$.

568 Combining Lemma 5, Lemma 10 and Lemma 11, we conclude that $\text{JSR}(\tilde{\mathcal{M}}_{\mathcal{G},\epsilon}) < 1$.

569 H Proof of Theorem 1

570 Recall the formulation of $X_{\cdot i}^{(t+1)}$ in (6):

$$X_{\cdot i}^{(t+1)} = \sigma(P^{(t)}(X^{(t)}W^{(t)})_{\cdot i}) = \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t W_{j_{k+1}j_k}^{(k)} \right) D_{j_{t+1}}^{(t)} P^{(t)} \dots D_{j_1}^{(0)} P^{(0)} X_{j_0}^{(0)}.$$

571 Then it follows that

$$\begin{aligned} \|BX_{\cdot i}^{(t+1)}\|_2 &= \left\| \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t W_{j_{k+1}j_k}^{(k)} \right) BD_{j_{t+1}}^{(t)} P^{(t)} \dots D_{j_1}^{(0)} P^{(0)} X_{j_0}^{(0)} \right\|_2 \\ &\leq \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t |W_{j_{k+1}j_k}^{(k)}| \right) \left\| BD_{j_{t+1}}^{(t)} P^{(t)} \dots D_{j_1}^{(0)} P^{(0)} X_{j_0}^{(0)} \right\|_2 \\ &= \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t |W_{j_{k+1}j_k}^{(k)}| \right) \left\| \tilde{D}_{j_{t+1}}^{(t)} \tilde{P}^{(t)} \dots \tilde{D}_{j_1}^{(0)} \tilde{P}^{(0)} BX_{j_0}^{(0)} \right\|_2 \\ &\leq \sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t |W_{j_{k+1}j_k}^{(k)}| \right) Cq^{t+1} \|BX_{j_0}^{(0)}\|_2 \\ &\leq C'q^{t+1} \left(\sum_{j_{t+1}=i, (j_t, \dots, j_0) \in [d]^{t+1}} \left(\prod_{k=0}^t |W_{j_{k+1}j_k}^{(k)}| \right) \right) \\ &= C'q^{t+1} \|(|W^{(0)}| \dots |W^{(t)}|)_{\cdot i}\|_1, \end{aligned}$$

572 where $C' = C \max_{j \in [d]} \|BX_j^{(0)}\|_2$ and $\|\cdot\|_1$ denotes the 1-norm. Specifically, the first inequality follows
573 from the triangle inequality, and the second inequality is due to the property of the joint spectral
574 radius in (7), where $\text{JSR}(\mathcal{M}_{G, \epsilon}) < q < 1$.

Since $\|Bx\|_2 = \|x\|_2$ if $x^\top \mathbf{1} = 0$ for $x \in \mathbb{R}^N$, we also have that if $X^\top \mathbf{1} = 0$ for $X \in \mathbb{R}^{N \times d}$, then

$$\|BX\|_F = \|X\|_F,$$

575 using which we obtain that

$$\begin{aligned} \mu(X^{(t+1)}) &= \|X^{(t+1)} - \mathbf{1}\gamma_{X^{(t+1)}}\|_F = \|BX^{(t+1)}\|_F = \sqrt{\sum_{i=1}^d \|BX_{\cdot i}^{(t+1)}\|_2^2} \\ &\leq C'q^{t+1} \sqrt{\sum_{i=1}^d \|(|W^{(0)}| \dots |W^{(t)}|)_{\cdot i}\|_1^2} \\ &\leq C'q^{t+1} \sqrt{\left(\sum_{i=1}^d \|(|W^{(0)}| \dots |W^{(t)}|)_{\cdot i}\|_1 \right)^2} \\ &= C'q^{t+1} \|(|W^{(0)}| \dots |W^{(t)}|)_{1,1}\|, \end{aligned}$$

where $\|\cdot\|_{1,1}$ denotes the matrix $(1, 1)$ -norm (recall from Section A.2 that for a matrix $M \in \mathbb{R}^{m \times n}$, we have $\|M\|_{1,1} = \sum_{i=1}^m \sum_{j=1}^n |M_{ij}|$). The assumption **A3** implies that there exists C'' such that for all $t \in \mathbb{N}_{\geq 0}$,

$$\|(|W^{(0)}| \dots |W^{(t)}|)_{1,1}\| \leq C'' d^2.$$

Thus we conclude that there exists C_1 such that for all $t \in \mathbb{N}_{\geq 0}$,

$$\mu(X^{(t)}) \leq C_1 q^t.$$

576 **I Proof of Proposition 2**

577 Since $D_{\text{deg}}^{-1}A$ is similar to $D_{\text{deg}}^{-1/2}AD_{\text{deg}}^{-1/2}$, they have the same spectrum. For $D_{\text{deg}}^{-1}A$, the smallest
 578 nonzero entry has value $1/d_{\text{max}}$, where d_{max} is the maximum node degree in \mathcal{G} . On the other hand,
 579 it follows from the definition of $\mathcal{P}_{\mathcal{G},\epsilon}$ that

$$\epsilon d_{\text{max}} \leq 1.$$

580 Therefore, $\epsilon \leq 1/d_{\text{max}}$ and thus $D_{\text{deg}}^{-1}A \in \mathcal{P}_{\mathcal{G},\epsilon}$.

581 We proceed by proving the following result.

582 **Lemma 12.** *For any M in \mathcal{M} , the spectral radius of M denoted by $\rho(M)$, satisfies*

$$\rho(M) \leq \text{JSR}(\mathcal{M}).$$

583

584 *Proof.* Gelfand’s formula states that $\rho(M) = \lim_{k \rightarrow \infty} \|M^k\|^{1/k}$, where the quantity is independent of the
 585 norm used [3]. Then comparing with the definition of the joint spectral radius, we can immediately
 586 conclude the statement. \square

587 Let $B(D_{\text{deg}}^{-1}A) = \tilde{P}B$. By definition, $\tilde{P} \in \tilde{\mathcal{M}}_{\mathcal{G},\epsilon}$ since $D_{\text{deg}}^{-1}A \in \mathcal{P}_{\mathcal{G},\epsilon}$ as shown before the lemma.
 588 Moreover, the spectrum of \tilde{P} is the spectrum of $D_{\text{deg}}^{-1}A$ after reducing the multiplicity of eigenvalue 1
 589 by one. Under the assumption **A1**, the eigenvalue 1 of $D_{\text{deg}}^{-1}A$ has multiplicity 1, and hence $\rho(\tilde{P}) = \lambda$,
 590 where λ is the second largest eigenvalue of $D_{\text{deg}}^{-1}A$. Putting this together with Lemma 12, we conclude
 591 that

$$\lambda \leq \text{JSR}(\tilde{\mathcal{M}}_{\mathcal{G},\epsilon})$$

592 as desired.

593 **J Numerical Experiments**

594 Here we provide more details on the numerical experiments presented in Section 5. All models were
 595 implemented with PyTorch [5] and PyTorch Geometric [1].

596 **Datasets** We used `torch_geometric.datasets.planetoid` provided in PyTorch Geometric
 597 for all the three datasets: Cora, CiteSeer, and PubMed with their default training and test splits.

598 **Model details**

- 599 • For GAT, we consider the architecture proposed in Veličković et al. [7] with each attentional
 600 layer sharing the parameter a in $\text{LeakyReLU}(a^\top [W^\top X_i || W^\top X_j])$, $a \in \mathbb{R}^{2d'}$ to compute the
 601 attention scores.
- 602 • For GCN, we consider the standard random walk graph convolution $D_{\text{deg}}^{-1}A$. That is, the update
 603 rule of each graph convolutional layer can be written as

$$X' = D_{\text{deg}}^{-1}AXW,$$

604 where X and X' are the input and output node representations, respectively, and W is the shared
 605 learnable weight matrix in the layer.

606 **Compute** We trained all of our models on a Telsa V100 GPU.

607 **Training details** In all experiments, we used the Adam optimizer using a learning rate of 0.00001
 608 and 0.0005 weight decay and trained for 1000 epoch.

609 **References**

- 610 [1] Matthias Fey and Jan E. Lenssen. Fast Graph Representation Learning with PyTorch Geometric.
611 In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 612 [2] Darald J. Hartfiel. *Nonhomogeneous Matrix Products*. 2002.
- 613 [3] Peter D. Lax. *Functional Analysis*. 2002.
- 614 [4] David A. Levin and Yuval Peres and Elizabeth L. Wilmer. *Markov Chains and Mixing Times*.
615 2008.
- 616 [5] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
617 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
618 Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
619 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style,
620 High-Performance Deep Learning Library. In *NeurIPS*, 2019.
- 621 [6] Jacques Theys. Joint Spectral Radius: theory and approximations. *Ph. D. dissertation*, 2005.
- 622 [7] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
623 Bengio. Graph Attention Networks. In *ICLR*, 2018.