

Modeling web-crawlers on the Internet with random walks on graphs

Doron Shahar

December 11, 2014

Motivation

- The state of the Internet at any time can be modeled by a graph where websites are the vertices and the links between websites are the edges of the graph.

Motivation

- The state of the Internet at any time can be modeled by a graph where websites are the vertices and the links between websites are the edges of the graph.
- The growth of the Internet can be modeled by a random graph.

Motivation

- The state of the Internet at any time can be modeled by a graph where websites are the vertices and the links between websites are the edges of the graph.
- The growth of the Internet can be modeled by a random graph.
- Our key interest shall be in web crawlers. A web crawler browses the Internet for the purpose of indexing websites. This is one the way in which search engines find new sites to list. We can model the way in which the web crawler browses by a random walk on a graph.

Motivation

- The state of the Internet at any time can be modeled by a graph where websites are the vertices and the links between websites are the edges of the graph.
- The growth of the Internet can be modeled by a random graph.
- Our key interest shall be in web crawlers. A web crawler browses the Internet for the purpose of indexing websites. This is one the way in which search engines find new sites to list. We can model the way in which the web crawler browses by a random walk on a graph.
- My goal then is to describe properties of random walks on the random graphs that model the growth of the Internet.

The Model $G_{d,n}$

The random graph $G_{d,n}$ is formed by introducing what are called mini-vertices. Some of these mini-vertices are then identified to form one vertex in the final graph.

The Model $G_{d,n}$

The random graph $G_{d,n}$ is formed by introducing what are called mini-vertices. Some of these mini-vertices are then identified to form one vertex in the final graph.

The first mini-vertex is introduced with a single self-loop. Each additional mini-vertex added to the graph attaches to one of the previous mini-vertices with a probability proportional to the degree of the existing mini-vertices. This process continues until there are dn mini-vertices. The first d mini-vertices are then identified to form the first vertex; the next set of d mini-vertices are identified to form the second vertex; and so on.

The Model $G_{d,n}$

The random graph $G_{d,n}$ is formed by introducing what are called mini-vertices. Some of these mini-vertices are then identified to form one vertex in the final graph.

The first mini-vertex is introduced with a single self-loop. Each additional mini-vertex added to the graph attaches to one of the previous mini-vertices with a probability proportional to the degree of the existing mini-vertices. This process continues until there are dn mini-vertices. The first d mini-vertices are then identified to form the first vertex; the next set of d mini-vertices are identified to form the second vertex; and so on.

Alternatively, you may view the graph as being built one vertex at a time with each vertex attaching with d edges to the other vertices with the possibility of self-loops.

Convergence to a Stationary State

The random walk on $G_{d,n}$ is irreducible and aperiodic, because the graph is connected and the first vertex has a selfloop. By the ergodic theorem, P^t ($t \in \mathbb{N}$) converges to the matrix Π whose rows are all π . Here P is the transition matrix, and π is the stationary distribution.

Convergence to a Stationary State

The random walk on $G_{d,n}$ is irreducible and aperiodic, because the graph is connected and the first vertex has a selfloop. By the ergodic theorem, P^t ($t \in \mathbb{N}$) converges to the matrix Π whose rows are all π . Here P is the transition matrix, and π is the stationary distribution.

Physically, this means that for large time t the initial state of the system does not influence the final state of the system. In regards to the Internet, this means that eventually the probability that the web crawler starts visits a particular website will be the same regardless of which website it started at. This “eventually” is a time t when $P^t \approx \Pi$. It is related to the time when it is likely that the web-crawler has indexed every website.

Reversible Markov Chains

Definition

A Markov Chain with transition matrix P is reversible if there is probability distribution π such that $\pi(i)P(i,j) = \pi(j)P(j,i)$.

Reversible Markov Chains

Definition

A Markov Chain with transition matrix P is reversible if there is probability distribution π such that $\pi(i)P(i,j) = \pi(j)P(j,i)$.

Theorem

Every random walk on a graph is a reversible Markov chain.

Reversible Markov Chains

Definition

A Markov Chain with transition matrix P is reversible if there is probability distribution π such that $\pi(i)P(i,j) = \pi(j)P(j,i)$.

Theorem

Every random walk on a graph is a reversible Markov chain.

The following matrices will be useful when proving facts about reversible Markov chains.

Definition

Let D be the diagonal matrix with $D(i,i) = \pi(i)$. Let $Q = DP$, and $A = D^{1/2}PD^{-1/2}$.

Reversible Markov Chains

Definition

A Markov Chain with transition matrix P is reversible if there is probability distribution π such that $\pi(i)P(i,j) = \pi(j)P(j,i)$.

Theorem

Every random walk on a graph is a reversible Markov chain.

The following matrices will be useful when proving facts about reversible Markov chains.

Definition

Let D be the diagonal matrix with $D(i,i) = \pi(i)$. Let $Q = DP$, and $A = D^{1/2}PD^{-1/2}$.

Theorem

Q and A are symmetric matrices.

Key Bounds

Definition

The relative pointwise distance is $\Delta(t) = \max_{i,j} \left| \frac{P^t(i,j)}{\pi(j)} - 1 \right| = \max_{i,j} \left| \frac{P^t(i,j)}{\Pi(i,j)} - 1 \right|$.

Definition

$\lambda_{max} = \max\{|\lambda_2|, |\lambda_n|\} = \max\{|\lambda_2|, |\lambda_3|, \dots, |\lambda_n|\}$

Key Bounds

Definition

The relative pointwise distance is $\Delta(t) = \max_{i,j} \left| \frac{P^t(i,j)}{\pi(j)} - 1 \right| = \max_{i,j} \left| \frac{P^t(i,j)}{\Pi(i,j)} - 1 \right|$.

Definition

$\lambda_{max} = \max\{|\lambda_2|, |\lambda_n|\} = \max\{|\lambda_2|, |\lambda_3|, \dots, |\lambda_n|\}$

Theorem

$$\Delta(t) \leq \frac{1}{\min_k \pi(k)} \lambda_{max}^t$$

Key Bounds

Definition

The relative pointwise distance is $\Delta(t) = \max_{i,j} \left| \frac{P^t(i,j)}{\pi(j)} - 1 \right| = \max_{i,j} \left| \frac{P^t(i,j)}{\Pi(i,j)} - 1 \right|$.

Definition

$\lambda_{max} = \max\{|\lambda_2|, |\lambda_n|\} = \max\{|\lambda_2|, |\lambda_3|, \dots, |\lambda_n|\}$

Theorem

$$\Delta(t) \leq \frac{1}{\min_k \pi(k)} \lambda_{max}^t$$

λ_2 can be bounded in terms of conductance. Therefore, we will be able to find bounds for λ_{max} in terms of conductance when $\lambda_{max} = \lambda_2$.

Theorem (Cheeger's inequality)

$$\lambda_2 < 1 - \frac{(\Phi_G)^2}{2}$$

Random walks with $\lambda_{max} = \lambda_2$

Changes to the random walk can be made to shift all the eigenvalues to be non-negative. Then $1 = \lambda'_1 \geq \lambda'_2 \geq \dots \geq \lambda'_n \geq 0$. So $\lambda'_{max} = \lambda'_2$.

Random walks with $\lambda_{max} = \lambda_2$

Changes to the random walk can be made to shift all the eigenvalues to be non-negative. Then $1 = \lambda'_1 \geq \lambda'_2 \geq \dots \geq \lambda'_n \geq 0$. So $\lambda'_{max} = \lambda'_2$.

Example

A lazy random walk is one in which there is a probability of $1/2$ that you remain at the same vertex. Then new transition matrix is $P' = \frac{I+P}{2}$. All of the eigenvalues of P' are between 1 and 0. $\lambda'_2 = \lambda'_{max}$, and $\lambda'_2 = \frac{1+\lambda_2}{2}$.

Random walks with $\lambda_{max} = \lambda_2$

Changes to the random walk can be made to shift all the eigenvalues to be non-negative. Then $1 = \lambda'_1 \geq \lambda'_2 \geq \dots \geq \lambda'_n \geq 0$. So $\lambda'_{max} = \lambda'_2$.

Example

A lazy random walk is one in which there is a probability of $1/2$ that you remain at the same vertex. Then new transition matrix is $P' = \frac{I+P}{2}$. All of the eigenvalues of P' are between 1 and 0. $\lambda'_1 = \lambda'_{max}$, and $\lambda'_2 = \frac{1+\lambda_2}{2}$.

Example

A random jump time with a Poisson distribution can be introduced. $H_t(x, y) = \sum_{k=0}^{\infty} e^{-t} \frac{t^k}{k!} P^k(x, y)$ is the new transition matrix from time 0 to t . Its eigenvalues are between 1 and 0. $\lambda'_1 = \lambda'_{max}$, and $\lambda'_2 = e^{-t(1-\lambda_2)}$.

Random walks with $\lambda_{max} = \lambda_2$

Changes to the random walk can be made to shift all the eigenvalues to be non-negative. Then $1 = \lambda'_1 \geq \lambda'_2 \geq \dots \geq \lambda'_n \geq 0$. So $\lambda'_{max} = \lambda'_2$.

Example

A lazy random walk is one in which there is a probability of $1/2$ that you remain at the same vertex. Then new transition matrix is $P' = \frac{I+P}{2}$. All of the eigenvalues of P' are between 1 and 0. $\lambda'_2 = \lambda'_{max}$, and $\lambda'_2 = \frac{1+\lambda_2}{2}$.

Example

A random jump time with a Poisson distribution can be introduced. $H_t(x, y) = \sum_{k=0}^{\infty} e^{-t} \frac{t^k}{k!} P^k(x, y)$ is the new transition matrix from time 0 to t . Its eigenvalues are between 1 and 0. $\lambda'_2 = \lambda'_{max}$, and $\lambda'_2 = e^{-t(1-\lambda_2)}$.

The latter random walk is a continuous-time Markov chain. We still have the same bound as before for $\Delta(t)$. I found a simple proof of this result using the definition of H_t and the result for P^t .

Definitions concerning conductance

Definition

Let G be a undirected multigraph with self-loops with vertex set V and edge set E .

- 1 $d_G(v)$ is the degree of a vertex, where self loops contribute 2 to the degree.
- 2 The volume of a set $S \subseteq V$ is $\text{Vol}_G(S) = \sum_{v \in S} d_G(v)$
- 3 The cutset of $S \subseteq V$ is
 $C_G(S, S^C) = \{e \in E : e \text{ has one endpoint in } S \text{ and one in } S^C\}$
- 4 The conductance Φ_G of G is
$$\Phi_G = \min \left\{ \frac{|C_G(S, S^C)|}{\text{Vol}_G(S)} : \emptyset \neq S \subseteq V, \text{Vol}_G(S) \leq \frac{\text{Vol}_G(V)}{2} \right\}.$$

Definitions concerning conductance

Definition

Let G be a undirected multigraph with self-loops with vertex set V and edge set E .

- 1 $d_G(v)$ is the degree of a vertex, where self loops contribute 2 to the degree.
- 2 The volume of a set $S \subseteq V$ is $\text{Vol}_G(S) = \sum_{v \in S} d_G(v)$
- 3 The cutset of $S \subseteq V$ is
 $C_G(S, S^C) = \{e \in E : e \text{ has one endpoint in } S \text{ and one in } S^C\}$
- 4 The conductance Φ_G of G is
$$\Phi_G = \min \left\{ \frac{|C_G(S, S^C)|}{\text{Vol}_G(S)} : \emptyset \neq S \subseteq V, \text{Vol}_G(S) \leq \frac{\text{Vol}_G(V)}{2} \right\}.$$

I have generalized this definition to all graphs.

Definition

$$\Phi_G = \min \left\{ 1, \frac{|C_G(S, S^C)|}{\text{Vol}_G(S)} : S \subseteq V, \emptyset \neq S \neq V, 0 \leq \text{Vol}_G(S) \leq \frac{\text{Vol}_G(V)}{2} \right\}.$$

Theorem

- $0 \leq \Phi_G \leq 1$.
- $\Phi_G = 0$ iff G is disconnected.
- Hubs are the only graphs with four or more vertices that have $\Phi_G = 1$.

Results about Conductance

Theorem

- $0 \leq \Phi_G \leq 1$.
- $\Phi_G = 0$ iff G is disconnected.
- Hubs are the only graphs with four or more vertices that have $\Phi_G = 1$.

Example

Let G be a multigraph with an odd number of vertices and one edge connecting every vertex to every other vertex, with a self loop from every vertex to itself. $\Phi_G = 1/2$.

Results about Conductance

Theorem

- $0 \leq \Phi_G \leq 1$.
- $\Phi_G = 0$ iff G is disconnected.
- Hubs are the only graphs with four or more vertices that have $\Phi_G = 1$.

Example

Let G be a multigraph with an odd number of vertices and one edge connecting every vertex to every other vertex, with a self loop from every vertex to itself. $\Phi_G = 1/2$.

Theorem

$\Phi_G = \min\{1, \frac{Q(S, S^c)}{P(S)} : S \subseteq V, \emptyset \neq S \neq V, P(S) \leq 1/2\}$, where
 $Q(S, S^c) = \sum_{i \in S, j \in S^c} Q(i, j)$ and $P(S) = \sum_{i \in S} \pi(i)$.

Theorem Bounding Conductance

Theorem

For all integers $d \geq 2$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} < \frac{\alpha}{2d+\alpha}) = 0$ for all $\alpha < \min\{\frac{2d-3}{4}, \alpha_0\}$ where α_0 is the unique positive solution less than d of the equation $\frac{3}{2}\alpha(1 + \ln(d) - \ln(\alpha)) = (d - 1 - \frac{5}{2}\alpha) \ln 2$.

Theorem Bounding Conductance

Theorem

For all integers $d \geq 2$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} < \frac{\alpha}{2d+\alpha}) = 0$ for all $\alpha < \min\{\frac{2d-3}{4}, \alpha_0\}$ where α_0 is the unique positive solution less than d of the equation $\frac{3}{2}\alpha(1 + \ln(d) - \ln(\alpha)) = (d - 1 - \frac{5}{2}\alpha) \ln 2$.

Idea behind Proof:

Bound $\Phi_{G_{d,n}}$ below in terms of edge expansion $\rho_{G_{d,n}} = \min\{\frac{|C_G(S, S^c)|}{|S|}\}$

Theorem Bounding Conductance

Theorem

For all integers $d \geq 2$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} < \frac{\alpha}{2d+\alpha}) = 0$ for all $\alpha < \min\{\frac{2d-3}{4}, \alpha_0\}$ where α_0 is the unique positive solution less than d of the equation $\frac{3}{2}\alpha(1 + \ln(d) - \ln(\alpha)) = (d - 1 - \frac{5}{2}\alpha) \ln 2$.

Idea behind Proof:

Bound $\Phi_{G_{d,n}}$ below in terms of edge expansion $\rho_{G_{d,n}} = \min\{\frac{|C_G(S, S^c)|}{|S|}\}$
 $P(\rho_{G_{d,n}} < \alpha) \leq \sum_{k=2}^{n/2} \binom{n}{k} \alpha k \binom{dn}{\alpha k} f(k)$ where $f(k) \geq P(A = \text{Good}(S))$ for all S with $|S| = k \leq n/2$ and for all A with $|A| = j$ for any $j \leq \alpha k - 1$.

Theorem Bounding Conductance

Theorem

For all integers $d \geq 2$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} < \frac{\alpha}{2d+\alpha}) = 0$ for all $\alpha < \min\{\frac{2d-3}{4}, \alpha_0\}$ where α_0 is the unique positive solution less than d of the equation $\frac{3}{2}\alpha(1 + \ln(d) - \ln(\alpha)) = (d - 1 - \frac{5}{2}\alpha) \ln 2$.

Idea behind Proof:

Bound $\Phi_{G_{d,n}}$ below in terms of edge expansion $\rho_{G_{d,n}} = \min\{\frac{|C_G(S, S^c)|}{|S|}\}$

$P(\rho_{G_{d,n}} < \alpha) \leq \sum_{k=2}^{n/2} \binom{n}{k} \alpha k \binom{dn}{\alpha k} f(k)$ where $f(k) \geq P(A = \text{Good}(S))$ for all S with $|S| = k \leq n/2$ and for all A with $|A| = j$ for any $j \leq \alpha k - 1$.

Lastly, we want to find a suitable candidate for $f(k)$.

$P(A = \text{Good}(S)) \leq \binom{dn}{|A|/2} \times \binom{dn-|A|/2}{dk-|A|/2}^{-1}$. This is a bound that I found, which is slightly better than the one given in the paper.

Theorem Bounding Conductance

Theorem

For all integers $d \geq 2$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} < \frac{\alpha}{2d+\alpha}) = 0$ for all $\alpha < \min\{\frac{2d-3}{4}, \alpha_0\}$ where α_0 is the unique positive solution less than d of the equation $\frac{3}{2}\alpha(1 + \ln(d) - \ln(\alpha)) = (d - 1 - \frac{5}{2}\alpha) \ln 2$.

Idea behind Proof:

Bound $\Phi_{G_{d,n}}$ below in terms of edge expansion $\rho_{G_{d,n}} = \min\{\frac{|C_G(S, S^c)|}{|S|}\}$

$P(\rho_{G_{d,n}} < \alpha) \leq \sum_{k=2}^{n/2} \binom{n}{k} \alpha k \binom{dn}{\alpha k} f(k)$ where $f(k) \geq P(A = \text{Good}(S))$ for all S with $|S| = k \leq n/2$ and for all A with $|A| = j$ for any $j \leq \alpha k - 1$.

Lastly, we want to find a suitable candidate for $f(k)$.

$P(A = \text{Good}(S)) \leq \binom{dn}{|A|/2} \times \binom{dn-|A|/2}{dk-|A|/2}^{-1}$. This is a bound that I found, which is slightly better than the one given in the paper.

$P(\rho_{G_{d,n}} < \alpha) \leq \sum_{k=2}^{n/2} \alpha k (2)^{\frac{\alpha}{2}k} (\frac{ed}{\alpha})^{\frac{3}{2}\alpha k} (\frac{k}{n})^{(d-1-2\alpha)k}$ after some trickery.

The first or last term in the sum is the largest. So $n/2$ times the first or last term bounds $P(\rho_{G_{d,n}} < \alpha)$. The result goes to zero for specified α .

Good and Bad Mini-vertices

Definition

A mini-vertex is said to be associated with a set if its corresponding vertex is in the set. A mini-vertex is good with respect to a set $S \subseteq V$ if it is associated with S and its father is associated with S^c or if it is associated with S^c and its father is associated with S . A mini-vertex is said to be bad with respect to S if it is not good with respect to S . By convention, mini-vertex 1 is bad. $Good(S)$ is the set of all mini-vertices that are good with respect to S .

Good and Bad Mini-vertices

Definition

A mini-vertex is said to be associated with a set if its corresponding vertex is in the set. A mini-vertex is good with respect to a set $S \subseteq V$ if it is associated with S and its father is associated with S^c or if it is associated with S^c and its father is associated with S . A mini-vertex is said to be bad with respect to S if it is not good with respect to S . By convention, mini-vertex 1 is bad. $Good(S)$ is the set of all mini-vertices that are good with respect to S .

Note that $|Good(S)| = |C_G(S, S^c)|$. Also note that every good mini-vertex contributes 1 to the volume of S ; every bad mini-vertex associated with S contributes 2 to the volume of S ; and every bad mini-vertex associated with S^c contributes 0 to the volume of S . Using these facts and some simple but neat computations, I found the bound $P(A = Good(S)) \leq \binom{dn}{|A|/2} \times \binom{dn-|A|/2}{dk-|A|/2}^{-1}$ in the same way a different bound was found in the paper.

Results

In the following table are the values that conductance is bounded above.

$\Phi_G \geq$	d=2	d=3	d=4	d=5
Their bound	0.00257956	0.02168342	0.02439024	0.01960784
My bound	0.02139017	0.03031210	0.03499005	0.03785976

Results

In the following table are the values that conductance is bounded above.

$\Phi_G \geq$	d=2	d=3	d=4	d=5
Their bound	0.00257956	0.02168342	0.02439024	0.01960784
My bound	0.02139017	0.03031210	0.03499005	0.03785976

Using Cheeger's inequality, I found the values that λ_2 is bounded below.

$\lambda_2 \leq$	d=2	d=3	d=4	d=5
Their bound	0.99999667	0.99976491	0.99970256	0.99980777
My bound	0.99977123	0.99954059	0.99938785	0.99928332

Results

In the following table are the values that conductance is bounded above.

$\Phi_G \geq$	d=2	d=3	d=4	d=5
Their bound	0.00257956	0.02168342	0.02439024	0.01960784
My bound	0.02139017	0.03031210	0.03499005	0.03785976

Using Cheeger's inequality, I found the values that λ_2 is bounded below.

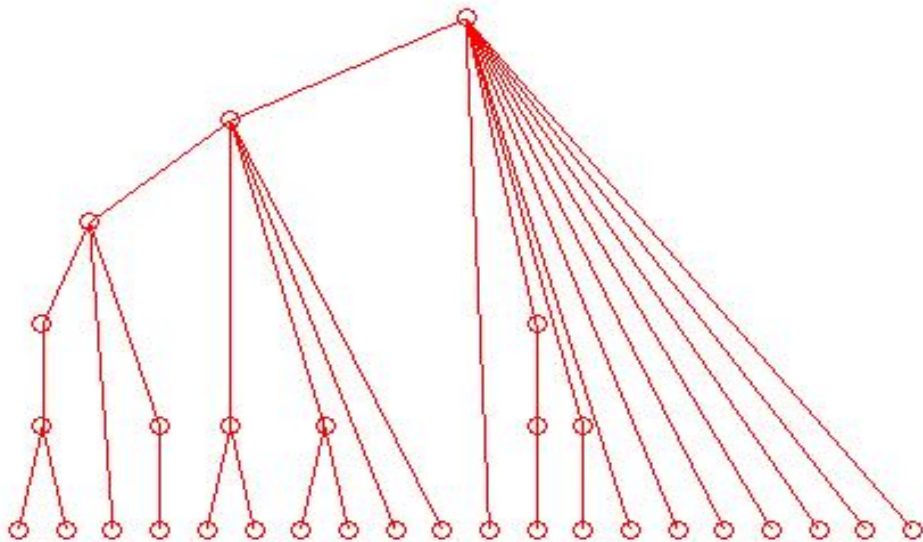
$\lambda_2 \leq$	d=2	d=3	d=4	d=5
Their bound	0.99999667	0.99976491	0.99970256	0.99980777
My bound	0.99977123	0.99954059	0.99938785	0.99928332

λ_2 seems to approach the following constants based on a program I wrote.

λ_2	d=2	d=3	d=4	d=5
	0.8354	0.7212	0.6428	0.5852

I conjecture that the constant λ_2 decreases with d and approaches 0.5.

A randomly generated graph with 30 vertices ($d=1$)



Theorem for $d = 1$

Theorem

For $d = 1$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} > \beta) = 0$ for all $\beta > 0$.

Theorem for $d = 1$

Theorem

For $d = 1$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} > \beta) = 0$ for all $\beta > 0$.

Idea behind Proof:

Let S_n be the set of the second vertex added to the graph and all of its descendants at time n . Then S_n^c is the set of the first vertex and all of its descendants except for the second vertex and its descendants.

$|C_G(S_n, S_n^c)| = 1$ for all n because the only edge in $C_G(S_n, S_n^c)$ is the edge connecting the first and second vertices.

Theorem for $d = 1$

Theorem

For $d = 1$, $\lim_{n \rightarrow \infty} P(\Phi_{G_{d,n}} > \beta) = 0$ for all $\beta > 0$.

Idea behind Proof:

Let S_n be the set of the second vertex added to the graph and all of its descendants at time n . Then S_n^c is the set of the first vertex and all of its descendants except for the second vertex and its descendants.

$|C_G(S_n, S_n^c)| = 1$ for all n because the only edge in $C_G(S_n, S_n^c)$ is the edge connecting the first and second vertices.

Both $\text{Vol}_{G_{d,n}}(S_n)$ and $\text{Vol}_{G_{d,n}}(S_n^c)$ approach infinity. Using these sets to estimate conductance gives us the desired result, since

$$\Phi_{G_{d,n}} \leq \frac{|C_G(S_n, S_n^c)|}{\min\{\text{Vol}_{G_{d,n}}(S_n), \text{Vol}_{G_{d,n}}(S_n^c)\}}.$$

Summary

For $d \geq 2$, λ_2 is very likely to be less than $1 - c$ for a constant c . So P^t most likely approaches Π faster than $(1 - c)^t$ when the graph is large. When $d = 1$, this is not the case.

If this is a good model of the Internet with $d \geq 2$, then a web-crawler will be able to index new websites well even as the Internet continues to grow. If this is a good model with $d = 1$, then the web-crawler will perform more and more poorly as the Internet grows.