

A Dynamical System on Bipartite Graphs

Kishore Papineni Google inc. papineni@google.com

ABSTRACT

This paper poses a non-linear dynamical system on bipartite graphs and shows its stability under certain conditions. The dynamical system changes the weights on the nodes of the graph in each time step. The underlying weight transformation is non-linear, motivated by information gain in a document retrieval setting. Stability analysis of this problem is therefore more involved than that of PageRank-like algorithms. We show convergence using methods from Lyapunov theory and also provide some examples of how the algorithm performs when ranking keywords and sentences in a set of documents.

KEYWORDS

Ranking sentences, Information gain, Lyupanov theory

ACM Reference Format:

Kishore Papineni and Pratik Worah. 2018. A Dynamical System on Bipartite Graphs. In *The 27th ACM International Conference on Information and Knowledge Management (CIKM '18), October 22–26, 2018, Torino, Italy.* ACM, New York, NY, USA, 4 pages. https://doi.org/10.1145/3269206.3269271

1 INTRODUCTION

1.1 Motivation

Suppose we are to summarize a corpus of documents that do not have explicit hyperlinks among them. Arguably, a good summary would contain important concepts from the corpus and arise from important sentences in these documents. Given a collection of words and sentences from a corpus of related documents, our goal then is to associate a weight with each sentence and each word (or phrase) which can be interpreted as their importance in the corpus.

From the viewpoint of natural language understanding, words that occur too frequently are usually not important. For example, words such as 'the' are usually not important and typically will not be the keywords of a document. This is the informal intuition behind Inverse Document Frequency (IDF) [11]. However, IDF is the highest for rarest words; but just as very frequent words in a document are unimportant, very rare words are also not expected to be keywords of a document. Therefore, neither too frequent nor very rare words are going to be keywords in a document, and the sweet-spot lies somewhere in the middle.

The above informal intuition was first described in [5], and was captured in [10] in a mathematical setting, which defined the *information gain* of a word as the Kullback-Leibler divergence of the



This work is licensed under a Creative Commons Attribution International 4.0 License.

CIKM '18, October 22–26, 2018, Torino, Italy © 2018 Copyright held by the owner/author(s). ACM ISBN 978-1-4503-6014-2/18/10. https://doi.org/10.1145/3269206.3269271 Pratik Worah Google inc. pworah@google.com

word's relative frequency of occurence in a document. In other words, the information gain from a given word is

$$g_{w} := \mathbf{f}_{w} \cdot (\mathbf{f}_{w} - 1 - \ln(\mathbf{f}_{w})),$$

where f_w denotes the fraction of sentences containing the word w. Note that the gain, as a function of f, is a positve and inverted U-shaped function in [0, 1] (the range of f), with a maximum near 0.2. Of course, there is nothing magical about the constant 0.2 i.e., ideal keywords are not necessarily those that occur with a frequency of 20%. This particular fraction is merely an artifact of using equal weights on all sentences to measure the Kullback-Leibler divergence. Therefore, assigning the correct weights to sentences is important before we can identify keywords. In particular, long sentences with many keywords should intuitively be more important from a natural language understanding and summarization point of view than short sentences which contain mostly filler words. Hence, we are now facing a circular problem: We need to identify important sentences before we can identify important words, but to identify important sentences we need to identify important words.

The circular nature of our problem suggests the following fixedpoint type algorithm to compute a weighted version of information gain: Assign a (reasonable) intial weight to all words and sentences. For each word w and each sentence s, compute sum_w and sum_s – the sum of weights of sentences that contain w and the sum of weights of words in s. Update the weight of each w and s as follows:

$$\begin{split} \text{Weight}_{w} &\leftarrow \text{sum}_{w} \cdot (\text{sum}_{w} - 1 - \ln(\text{sum}_{w})), \\ \text{Weight}_{s} &\leftarrow \text{sum}_{s} \cdot (\text{sum}_{s} - 1 - \ln(\text{sum}_{s})). \end{split}$$

Repeat the updates until the weights of all vertices do not change significantly between two iterations.

Mathematically, it is not clear under what conditions i.e., for what kinds of corpus and what initial choice of weights, does the above algorithm converge to a "optimal" weight distribution. Intuitively, it is not clear how the important sentences and keywords generated by using such a weighted information gain qualitatively compare with some other conventional approach, or even in some absolute sense.

This paper applies basic linear algebra and Lyapunov theory to prove convergence of the above iterated update algorithm when the underlying word-sentence graph obeys certain eigenvalue constraints. Moreover, we have also implemented the above fixed-point algorithm and ran it on a small corpus of documents to see if there was a natural qualitative difference between the sentences produced with weights and using a conventional method.

Finally, we observe that there is no reason why the algorithm can't be applied to more general settings. For instance, instead of a collection of words and sentences, we can apply it to any container which is a collection of tokens – a corpus of images and the entities within those images; a corpus of videos and the entities within those videos, and so on.

1.2 Related Work

The fixed-point approach described in the introduction implicitly assumes a graph structure between sentences and words in a document. In the underlying graph, every word and sentence corresponds to a vertex, and a word is linked to a sentence by an undirected edge, if it occurs in the sentence. There is already a significant amount of work which applies fixed-point algorithms over an underlying graph structure to compute a rank or score which has some relevance in the real world. For example, the initial papers [9], [4] in the areas of citation analysis, social networks and analyzing the link structure of the world-wide web, use similar ideas. More importantly, in the area of lexical analysis and ranking of documents there has been significant work using such graph based models, see for example [6], [8] and the book [7].

In [9], the authors work on a graph where the vertices are webpages and the edges reflect the link structure of the web. If $\delta_+(v)$ and $\delta_-(v)$ denote the out-degree and in-degree neighbours of a vertex v, then the score of v (on which its Page-rank is based) is defined as:

$$S(v) := (1 - c) + c \cdot \sum_{u \in \delta_{-}(v)} \frac{S(u)}{|\delta_{+}(u)|},$$
(1)

where the constant c is around 0.8. The idea is to iterate until the score stabilizes and then use the scores to rank web-pages. Note that the initial scores do not matter too much as the scores converge rapidly [9].

[4] introduced a similar algorithm (HITS), which ranked webpages by implicitly ascribing a bipartitie structure to the graph – vertices correspond to "authority" pages (corresponding to a large δ_{-}) or hub pages (corresponding to a large δ_{+}). Therefore, we have two scores for each vertex: an authority score and a hub score – both calculated as a weighted linear sum over their neighbours.

More relevant to us, [8] introduced TextRank an unsupervised procedure to extract and rank keywords and sentences, or more generally text units, from a lexical corpus. Their approach is similar to Page's approach, but now the graph vertices correspond to text units (for example, sentences or phrases) and edges reflect some semantic or syntactic connection between text units. While the algorithm and score used is similar to that in [9], one important difference was that the edges can be weighed and so their counterpart of Equation 1 would be:

$$S(\upsilon) := (1-\varepsilon) + \varepsilon \cdot \sum_{u \in \delta_{-}(\upsilon)} \frac{W_{u\upsilon} \cdot S(u)}{\sum_{w \in \delta_{+}(u)} W_{uw}},$$
(2)

where W is the weight matrix of the underlying graph and the constant c may be chosen depending upon the exact scenario at hand.

While the viewpoints in each case may be different, the relevant theme, at least for us, remains the same i.e., there is an iterated local computation on a graph that leads to a score for each vertex, and this score is interesting from the perspective of some real world problem. However, in all the above cases the score is computed using a linear function, as in Equations 1 and 2, for example. Linearity together with high connectivity in the graph structure leads to rapid mixing, thus convergence is ensured and is typically very fast. This linearity is in contrast to our situation, where we iteratively update based on a Kullback-Leibler divergence type of function– a non-linear function, and so convergence does not follow from earlier ideas. Therefore, we go back to first principles and analyze convergence of the underlying dynamical system using Lyapunov theory.

2 RESULTS

In this section, we recap the mathematical model for our problem and then provide a proof of Theorem 1 and 3, our main mathematical results. We provide experimental verification in the next section.

We are given an undirected bipartite graph $G \equiv (X, Y, W)$, where X and Y denote the vertex sets and W denotes the edge weights i.e., $W_{uv} \in \mathbb{R}$ is the weight of the edge connecting vertices u and v in G. With |X| = n and |Y| = m, W is a $(n+m) \times (n+m)$ block offdiagonal symmetric matrix¹. Let $x_u(t) \in \mathbb{R}^n$ and $y_v(t) \in \mathbb{R}^m$, where $u \in X$ and $v \in Y$, denote the values of the weight distribution at time t on the vertices u and v respectively. Let

$$z(t) := \begin{pmatrix} x(t) \\ y(t) \end{pmatrix}, \tag{3}$$

denote the column vector of all weights, which is the state of the dynamical system that we will induce on the graph.

Following [10], we define the time-invariant gain function $g : \mathbb{R} \to \mathbb{R}$ as follows:

$$g(x) := x \cdot (x - 1 - \ln |x|), x \in \mathbb{R}, \qquad (4)$$

define the vector version $\bar{g} : \mathbb{R}^{n+m} \to \mathbb{R}$, as the component-wise application of g to each element in the argument.

At time $t + \Delta t$, we update the weight vectors x and y based on the following update rule:

$$z(t + \Delta t) = \bar{g}(Wz(t)).$$
(5)

Note that edge weights are constant, only the vertex weights change over time. We thus have a nonlinear time-invariant dynamical system on the graph. We make a mild assumption that *W* is full-rank. In our motivating example, this assumption is satisfied when no two tokens appear in exactly the same set of containers. If they do, we can replace by a joint token. When *W* is non-binary, we can perturb *W* to make it full-rank. This concludes the description of our underlying model. We can state our main theorem as follows.

THEOREM 1. Assume that there exists $z_e := (x_e, y_e) \in \mathbb{R}^{n+m}$, a positive solution to the system $\overline{g}(Wz_e) = z_e$. Suppose that λ denotes the maximum eigenvalue of the graph G^2 obtained by squaring G i.e., squaring the adjacency matrix of G, and then multiplying the weight of each edge $uv \in G^2$ by $g'(z_e(u)) \cdot g'(z_e(v))$. The algorithm for updating vertex weights given above converges to some stable point z_e if $\lambda < 1$, and the initial value z_0 lies in the basin of attraction of z_e .

We first sketch the main ideas behind the proof in the scalar case before presenting the actual proof. Consider the following scalar ODE:

$$\frac{\mathrm{d}}{\mathrm{d}t}x(t) = \hat{g}(x) := \bar{g}(x) - x(t) \tag{6}$$

$$= x(t)(x(t) - 1 - \ln x(t)) - x(t).$$
(7)

The rate of change in weight x(t) in time Δt is equal to the change in x(t) in one update step, and as $\Delta t \rightarrow 0$, we get the ODE above.

¹In our motivating application, W is the adjacency matrix with $w_{uv} \in \{0, 1\}$.

The ODE has two fixed points: (1) x = 0.158..., and (2) x = 3.146..., with \hat{g} .

Our goal is to analyze whether either of the fixed points is locally asymptotically stable and by Lyapunov's criteria such a one dimensional update rule for x(t) would lead to convergence with $x^* \simeq 0.158...$, as long as the initial value x(0) was close to x^* .

We need the following indirect version of Lyapunov's stability criterion for higher dimensions:

THEOREM 2 ([3]). Lyapunov stability criterion: Let x = 0 be an equilibrium point for $\dot{x} = f(x)$, where $f : D \to \mathbb{R}^n$ is continuous and differentiable and D is a neighborhood of the origin. Let $J \equiv \frac{\partial}{\partial x} f(x)|_{x=0}$ denote the Jacobian matrix, then

- The origin is asymptotically stable if the real part of all eigenvalues of J are negative.
- The origin is unstable if the real part for at least one of the eigenvalues is positive.

Consider the following n + m dimensional ODE:

$$\frac{\mathrm{d}}{\mathrm{d}t}z(t) = \hat{g}(z(t)) := \bar{g}(Wz(t)) - z(t) \,. \tag{8}$$

The above dynamical system captures our update rule. Note that any $z_e \in \mathbb{R}^{n+m}$ that satisfies $\bar{g}(Wz_e) = z_e$ will be a candidate for a stable fixed point.

Our next task is to characterize the eigenvalues of the Jacobian at fixed point z_e . The uv^{th} entry of the Jacobian matrix at z_e is given by:

$$\frac{\partial}{\partial z_{\upsilon}} \left(g \left(\sum_{w \in N(u)} W_{u\upsilon} z_w \right) - z_u \right)_{|z_e}, \tag{9}$$

where N(u) denotes the set of neighbours of vertex u. By the chain rule, it equals:

$$\left(\frac{\partial}{\partial z_{\upsilon}}s(u)\right)\left(\frac{\partial}{\partial s}g(s)\right)_{|z_{e}} - \delta_{u}(\upsilon), \tag{10}$$

where we have used $s(u) := \sum_{w \in N(u)} W_{uv} z_w$ and δ is the kronecker delta function. But, the last expression is equivalent to:

$$W_{uv} \cdot g'(s)|_{z_e} - \delta_u(v). \tag{11}$$

Therefore, the Jacobian about z_e has the form:

$$\mathbf{J}_{|z_e} = \mathbf{D} \cdot \left(\frac{\mathbf{0} \quad | \mathbf{M} |}{\mathbf{M}^T \mid \mathbf{0}} \right) - \mathbf{I}, \tag{12}$$

where M is the $n \times m$ block matrix which is the incidence matrix of the bipartitie graph *G*, I is the n + m dimensional identity matrix, and D is the diagonal matrix with it's uth entry given by:

$$D_u = g'\left(\sum_{w \in N(u)} W_{uv} z_e(w)\right).$$
(13)

Let, H := J + I. If we can place an upper-bound on the real part of the eigenvalues of H so that

$$\operatorname{Re}\left(\lambda_{\max}(\mathrm{H})\right) < 1$$
 (14)

then the eigenvalues of J are simply obtained by translation of the eigenvalues of H, so that $Re(\lambda_{\min}(J)) < 0$ – the fixed point corresponding to z_e is then stable. The condition which ensures Equation 14 will be our condition in Theorem 1.

Let us rewrite D as:

$$\mathbf{D} = \begin{pmatrix} \mathbf{D}_1 & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{D}_2 \end{pmatrix},\tag{15}$$

where D_1 and D_2 are $n \times n$ and $m \times m$ diagonal matrices, respectively.

Proposition 1.

$$\mathbf{H}\mathbf{H}^{T} = \begin{pmatrix} \mathbf{D}_{1}\mathbf{M}\mathbf{M}^{T}\mathbf{D}_{1} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{D}_{2}\mathbf{M}^{T}\mathbf{M}\mathbf{D}_{2} \end{pmatrix}.$$
 (16)

The eigenvalues of HH^T are bounded by the maximum of the eigenvalues of $D_1 \text{MM}^T D_1$ and $D_2 \text{M}^T \text{MD}_2$. It follows from the definition of spectral norm that the absolute value of the eigenvalues of H is upper-bounded by the square-root of the maximum eigenvalue of HH^T . Therefore, the condition that both matrices $D_1 \text{MM}^T D_1$ and $D_2 \text{M}^T \text{MD}_2$ have all eigenvalues upper-bounded by 1 is sufficient to ensure a stable fixed point.

However, the matrices $D_1 MM^T D_1$ and $D_2 M^T MD_2$ are simply the adjacency matrix of the graphs on the partite sets *X* and *Y* of *G*, obtained by squaring the weighted bipartite graph *G* and reweighing every edge uv in the squared graph by $g'(z_e(u))g'(z_e(v))$. Hence, if the eigenvalues of these graphs are all upper-bounded by 1, the statement of Theorem 1 follows.

2.1 Basin of Attraction

The next task is to determine the starting values of z for which we can be assured of convergence. In this case, the conceptually easiest solution is to construct a Lyapunov function, as in the one dimensional case above, and show that V(z) is locally positive definite and $\dot{V}(z)$ is locally negative definite. However, things are more complicated in higher dimensions and such explicit constructions are not easy for non-linear dynamical systems. Therefore, the standard approach is to construct an ellipsoidal Lyapunov function for the linearized system, which is expected to be globally asymptotically stable, and then show it is locally asymptotically stable in a large enough region. Typically, even this problem is difficult but the symmetries of the Jacobian matrix makes it simpler.

THEOREM 3. Given matrix M which satisfies the conditions of Theorem 1, the radius of convergence about z_e is non-decreasing in λ , where λ is as in Theorem 1.

Consider a linear system $\dot{z} = Az$, it is asymptotically stable if the Jacobian J has all eigenvalues with negative real parts. A more direct characterization is that if the linear system is asymptotically stable then for any positive definite matrix Q, there exists a positive definite matrix P such that

$$A^T P + PA = -Q \tag{17}$$

which follows from using $V(z) = z^T P z$ as the Lyapunov function [2].

The same kind of analysis can be done for non-linear systems i.e., linearize the system about the fixed point (Taylor expansion with higher order terms dropped), compute a Lyapunov function which shows global asymptotic stability, then use it as a candidate for the non-linear system. *V* already has all the proprties we need except for one: $\dot{V}(z)$ may not be negative definite, but we can show it is locally negative definite and prove a lower bound on the convergence radius.

We know,

$$\frac{\mathrm{d}}{\mathrm{d}t}z(t) = \mathbf{J}_{|z_e}z(t) + \mathbf{h}(z(t)), \tag{18}$$

where h is what remains after substracting $J_{|z_e}$ from the RHS of Equation 8 i.e. $\bar{g}(z) - z$. In our notation, $A = J_{|z_e}$, so that Equation 17 reads:

$$\mathbf{J}^T P + P \mathbf{J} = -Q \tag{19}$$

Choosing Q = I gives the matrix equation:

$$\mathbf{J}^T P + P \mathbf{J} = -\mathbf{I}.$$
 (20)

Note that, the matrix *P* in Equation 20 is assumed symmetric, so that the left and right eigenvectors coincide (after transposition). Multiplying the LHS and RHS of Equation 20 by v_{max} and v_{max}^T , the unit eigenvector corresponding to the maximum eigenvalue $\lambda_{\text{max}}(P)$, we get,

$$\lambda_{\max}(P)\left(v_{\max}^T \mathbf{J}^T v_{\max} + v_{\max}^T \mathbf{J} v_{\max}\right) = -1.$$
(21)

By the definition of spectral norm,

$$\max\left(|v_{\max}^{T} J v_{\max}|, |v_{\max}^{T} J^{T} v_{\max}|\right) \le \sqrt{\lambda_{\max}(J J^{T})}, \qquad (22)$$

and so,

$$\sqrt{\lambda_{\max}(\mathbf{J}\mathbf{J}^T)} \ge \frac{1}{2\lambda_{\max}(P)}.$$
(23)

Now, it is also known (see [3] or [2]) that

$$\begin{split} \dot{V}(z) &= z^{T} (\mathbf{J}^{T} P + P \mathbf{J}) z + 2 z^{T} P \mathbf{h} \\ &\leq - \| z \|^{2} + 2 \lambda_{\max}(P) \| z \| \cdot \| \mathbf{h}(z) \| \\ &\leq - \left(1 - 2 \lambda_{\max}(P) \frac{\| \mathbf{h}(z) \|}{\| z \|} \right) \| z \|^{2}, \end{split}$$

where $\|\cdot\|$ denotes ℓ_2 norm. The RHS is negative as long as

$$2\lambda_{\max}(P)\frac{\|\mathbf{h}(z)\|}{\|z\|} \leq 1,$$
 (24)

which is equivalent, by Equation 23, to

$$\frac{\|\mathbf{h}(z)\|}{\|z\|} \leq \frac{1}{2\lambda_{\max}(P)}$$
(25)

$$\leq \sqrt{\lambda_{\max}(JJ^T)}.$$
 (26)

As $\lambda_{\max}(JJ^T)$ increases \dot{V} is at least going to remain locally negative definite around z_e . Recall from Proposition 1 and the discussion that followed that $\lambda_{\max}(JJ^T)$ is upper-bounded by the largest eigenvalue of the reweighed square of G i.e., $\lambda (\equiv \lambda_{\max}(DW^2D^T))$. Therfore, as λ increases, the time derivative of the candidate Lyapunov function either remains negative or may become negative. Hence, the statement of Theorem 3 follows.

3 EXPERIMENTAL RESULTS

We present some qualitative comparisons with the MMR algorithm [1] below. We believe the top-ranked sentences for the reflection algorithm have more information content. Further experimental results will appear in a full paper. Our comparison in the table simply displays the top-ranked sentence, for each algorithm, when both algorithms are run on the concatenated text in the set of top twenty documents obtained using google web-search for some arbitrary queries related to monarch butterflies.

Query	Weighted Gain	MMR
monarch butterfly parasites	Protozoan parasites such as Ophryocystis	Monarchs have many natural enemies -
	elektroscirrha and a microsporidian Nosema	predators, parasitoids, and parasites can
	species have also been identified in wild and	harm monarch eggs, larvae, pupae, and
	captive monarchs (McLaughlin and Myers	adults.
	1970, Leong et al. 1992;1997, Altizer and Ober-	
	hauser 1999, O. Taylor, personal communica-	
	tion).	
monarch butterfly migration	While the practice of transferring monarchs	The Monarch butterfly migrates for 2 rea-
	from place to place is generally not condoned	sons.
	by scientists, some reciprocal transfers of	
	tagged monarchs have demonstrated that	
	monarchs from east of the Rocky Mountains	
	will migrate south if transferred west, in the	
	range of the western population (rather than	
	SW).	
monarch butterfly adult	The four stages of the monarch butterfly life	The King of Butterflies - The Monarch But-
	cycle are the egg, the larvae (caterpillar), the	terfly
	pupa (chrysalis), and the adult butterfly.	
monarch butterfly climate	Aside from the ecological significance of	Given that monarchs largely depend on the
	these migrations - monarch butterflies are	genus Asclepias as larval host plants, the ef-
	the only insects known to migrate to warmer	fects of climate change on monarch north-
	climates more than 2,500 miles away - the	ward migrations will most likely be mediated
	butterflies' five-month layover in Mexico be-	by climate change effects on Asclepias.
	fore returning to the United States has be-	
	come one of the region's main tourist attrac-	
	tions and economic drivers.	
monarch butterfly captive rearing	Moreover, these subtle variation appear to	Chrysoperla rufilabris (Green Lacewing) lar-
	have biological significance; monarchs with	vae is a common inhabitant of milkweed
	darker shades of orange (approaching red)	species and a voracious predator of monarch
	show higher flight ability in captive settings	eggs.
	[13], and a recent study provided evidence	
	that the degree of black pigment is related to	
	migration distance in wild-caught monarchs	
	[14].	

Table 1: Some example sentences

REFERENCES

- Jaime Carbonell and Jade Goldstein. 1998. The Use of MMR, Diversity-based Reranking for Reordering Documents and Producing Summaries. 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (1998).
- [2] Emilio Frazzoli and Munther Dahleh. 2011. MIT Open Courseware. https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-241j-dynamic-systems-and-control-spring-2011 (2011).
- [3] Hassan Khalil. 2002. Nonlinear Systems. Prentice Hall, 3rd edition (2002).
- [4] Jon Kleinberg. 1999. Authoritative sources in a hyperlinked environment. J. ACM (1999).
- [5] Hans Peter Luhn. 1958. The Automatic Creation of Literature Abstracts. IBM journal of Research and Development (1958).
- [6] Rada Mihalcea. 2004. Graph based Ranking Algorithms for Sentence Extraction, Applied to Text Summarization. (2004).
- [7] Rada Mihalcea and Dragomir Radev. 2011. Graph-based Natural Language Processing and Information Retrieval. Cambridge University Press (2011).
- [8] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing Order into Texts. Conference on Empirical Methods in Natural Language Processing (2004).
- [9] Larry Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1998. The PageRank Citation Ranking: Bringing Order to the Web. Stanford University Technical Report (1998).
- [10] Kishore Papineni. 2001. Why Inverse Document Frequency? Proceedings of the second meeting of the North American Chapter of the Association for Computational Linguistics on Language Technologies (2001).
- [11] Karen Sparck-Jones. 1973. Index Term Weighting. Information Storage and Retrieval, 9 (1973).