

Using the concept of Google's PageRank to enhance student motivation in mathematics

Francisco Pedroche Sánchez
Department of Applied Mathematics,
Politechnic University of Valencia,
Camí de Vera s/n 46022, Valencia, Spain
pedroche@mat.upv.es

Abstract - The fact that the search engine Google is a popular tool known by our students can be used to catch their attention in classroom. Google recognizes that the PageRank™ algorithm is still the basic tool of the engine. In fact, this algorithm has already caught the interest of the researchers in mathematics and it is an active field of investigation. In this paper it is shown how to use the key-concept of the search engine Google, namely the PageRank, to illustrate some topics of mathematics. In order to explain how this algorithm works one has to deal with the following items: random process, random walk, conditional probability, total probability theorem, Markov chain, incidence matrix, stochastic matrix, eigenvalues, eigenvectors, vector space, spectral radius, matrix product, directed graph, stationary state, numerical methods and iterative methods for linear systems of equations. This communication describes how to introduce these concepts in a very suggestive way and connected with the topic of Google's PageRank.

Index Terms – Iterative methods, matrix analysis, PageRank, probability.

INTRODUCTION

The concept of crawling the internet is a common term for our students. For this reason, it is normal that when we use this term in the classroom, the students feel that the teacher is going to explain something familiar for them; something connected with their world made up with internet-related terms like myspace.com, messenger.com, secondlife.com or popomundo.com. When you talk about Google, they use to think: this person is going to say something interesting!. In fact, I have first noticed the existence of Google in the late 90's when one of my students told me that there was a fantastic searcher called that way. Some years later, in the course of a research project about Markov chains, I became in contact with the theoretical foundations of Google and I began to incorporate some of the concepts of Google as a part of the material that I used as examples in the classrooms. In this communication I try to put in order some of the ideas that I have been using in the classroom for some years. The aim of this paper is to show how these ideas can be used to illustrate real applications of mathematics to students of mathematics, computer science or engineering programs. Since the purpose of this paper is to focus on Google's concepts we have omitted rigorous definitions of the terms

involved letting this task to the teachers who would use these notes in their classrooms. The paper is structured as follows. Each section is named after the main concept that it is treated therein. The aim in each section is to show the connection of each topic with the main algorithm of Google: the PageRank algorithm [1].

SYSTEMS OF LINEAR EQUATIONS

Let us consider a system of n linear equations with n unknowns of the form

$$\left. \begin{array}{cccc} a_{11}x_1 + & a_{12}x_2 + & \cdots & + a_{1n}x_n = b_1 \\ a_{21}x_1 + & a_{22}x_2 + & \cdots & + a_{2n}x_n = b_2 \\ \vdots & \vdots & & \vdots \\ a_{n1}x_1 + & a_{n2}x_2 + & \cdots & + a_{nn}x_n = b_n \end{array} \right\} \quad (1)$$

where a_{ij} , x_i and b_i are real numbers for all the values of the indices i and j . The engineers of Google have to use a system of linear equations which size is that of the web that is crawled by Google. That is to say they have to handle systems of the order of 20.000 millions of equations with 20.000 millions of unknowns. How they manage to handle this amazing number of equations? Indeed this problem has been addressed as the largest matrix computation ever made [2]. Obviously, the first thing that we need to do is to use an easy notation. Try to imagine that we have to write explicitly (1) with millions of equations! That is a nonsense. With the aid of the concepts of matrix, column vector (a particular type of matrix) and product of two matrices we can write (1) in the very short form

$$Ax = b \quad (2)$$

Where we have denoted by A the matrix of the system- an ordered array of $n \times n$ elements with a_{ij} laying on the intersection of the i th row and the j th column of A - and by x and b the vector column of unknowns and the vector column of the right hand side of (1), respectively. Therefore we have now a useful notation to treat, at least theoretically, this enormous system of equations. In the following sections we see how a system like (2) can be associated with Google's PageRank and we also will show how to solve this problem.

CONNECTIVITY MATRIX

In order to illustrate how to construct a matrix associated with a set of web pages crawled by Google, let us consider the Figure 1. In this figure four web pages with their outgoing links (outlinks) are shown. The destination of the links is also indicated. For example, in page 3 there are outlinks that go to pages 1 and 2.

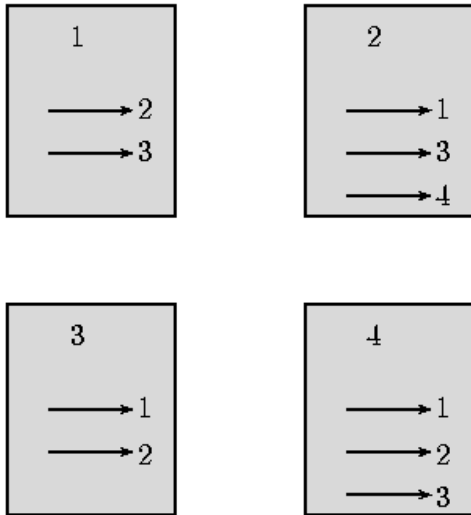


FIGURE 1

A WEB CONSISTING OF FOUR PAGES SHOWING THEIR OUTLINKS.

How can we describe the structure of the links between these four pages? One way is defining what is called the directed graph associated with the set of connections. In fact, this graph is just a schematically way of depicting the same information that is exhibited in Figure 1. In the present case, the directed graph associated with the connections of the web shown in Figure 1 is shown in Figure 2. The pages have been substituted by numbers called vertexes.

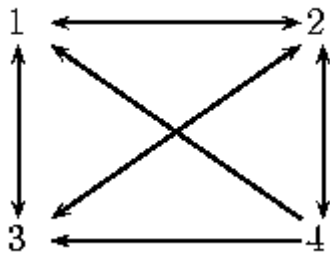


FIGURE 2

DIRECTED GRAPH ASSOCIATED WITH THE SET OF PAGES OF FIGURE 1.

As you can see we have gained more clarity in our treatment since it is easier to draw Figure 2 than Figure 1. Now we can define a matrix, called connectivity matrix, which summarizes all the information in Figure 2. We define the connectivity matrix associated with a directed graph as the square matrix with elements g_{ij} such that they have the value 1 if there is a connection from page j to page i , with

$i \neq j$, and 0 otherwise. Therefore the connectivity matrix for the directed graph in Figure 2 is

$$G = \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{pmatrix} \quad (3)$$

Now we have a matrix to illustrate the connections between the web pages. The connectivity matrices are called nonnegative matrices since their elements are all nonnegative. These matrices are also called sparse matrices since they use to have a great percentage of zeros. There are special procedures to work with sparse matrices; see [3] for some methods for storing the data of sparse matrices.

In order to show how Google constructs a matrix A to form a system of equations like (2) we need first to make some comments about random walks. We do this in the next section.

RANDOM WALK

To illustrate the characteristics of a Markov process we present in this section the classic model of a random walk, also called "drunkard's walk". Let us assume that a person is walking on a line with n discrete positions according to the following rule: The walker tosses a coin and if it is head he will go a step to the right, and otherwise he goes a step to the left; see Figure 3.

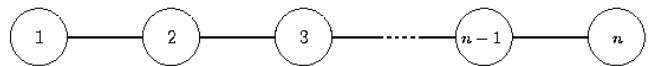


FIGURE 3

A LINE SHOWING DISCRETE POSITIONS.

We also admit that when the person reaches an extreme of the line then in the next movement he turns back, independently of the result of the flip. We are interested in describing the movement of this walker. In particular, when the person begins at position marked as '1', what can we say about the position (or state) of the walker when he has taken some steps? As you can observe we are in front of a random problem. We can not give a deterministic answer to this problem. To fix ideas, let us describe the time by the expression

$$t = k \Delta t, \quad k = 0, 1, 2, \dots \quad (4)$$

and let us denote by $p_i(k)$ the probability of finding the walker in the position i on the line of Figure 3, at instant time given by k . In order to describe the state of the walker in a certain time k we construct the probability distribution vector v_k - which is a positive vector such that the sum of

their components is the unity- and it is given by the expression

$$v_k = \begin{bmatrix} p_1(k) \\ p_2(k) \\ \vdots \\ p_n(k) \end{bmatrix} \quad (5)$$

For example, if we start the experiment in $t = 0$, with the walker in the position '1' of the line depicted in Figure 3 we have that the initial state of the walker is given by

$$v_0 = \begin{bmatrix} p_1(0) \\ p_2(0) \\ \vdots \\ p_n(0) \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (6)$$

since we know that the walker is in the position '1' and no uncertainty is assumed. In the following throwing of the coin we know, since we have assumed this condition, that the walker will go to the position marked as '2'. Therefore it is clear that

$$v_1 = \begin{bmatrix} p_1(1) \\ p_2(1) \\ \vdots \\ p_n(1) \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad (7)$$

In the next movement we have uncertainty since the walker can go to the left or to the right. Therefore it is clear that

$$v_2 = \begin{bmatrix} p_1(2) \\ p_2(2) \\ p_3(2) \\ \vdots \\ p_n(2) \end{bmatrix} = \begin{bmatrix} 1/2 \\ 0 \\ 1/2 \\ \vdots \\ 0 \end{bmatrix} \quad (8)$$

We remark that the position of the walker in any time described by k is given by a state vector v_k . The vector v_0 is called the initial state vector. Observe that in each consecutive movement we are distributing the probabilities and this allows us to compute the state vector v_k at any time. But if we want to describe the state given $k = 1000$ it would be very tedious to compute all the intermediate states to reach this state! How can we manage to solve this situation? We show this in the following section in where we introduce Markov chains and we will treat again this

example of the random walk with steps on the line depicted on Figure 3.

MARKOV CHAINS

The concepts of matrices and the product of a matrix by a vector will help us to write some computations easily. First, let us define a matrix P , called a transition matrix, associated with the random walk in the following form. Let $P = (p_{ij})$ where p_{ij} is the probability that the walker goes to the position i at time $k + 1$ given that he was in position j at time k . In our example it is easy to show that

$$p_{ij} = \begin{cases} 1/2 & \text{if } i = j - 1, j = 2, 3, \dots, n - 1 \\ 1/2 & \text{if } i = j + 1, j = 2, 3, \dots, n - 1 \\ 1 & \text{if } i = 2, j = 1 \\ 1 & \text{if } i = n - 1, j = n \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

And this can be written in matrix form as follows

$$P = \begin{bmatrix} 0 & 1/2 & & & & & \\ 1 & 0 & 1/2 & & & & \\ & 1/2 & 0 & 1/2 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & & 1/2 & 0 & 1 \\ & & & & & 1/2 & 0 \end{bmatrix} \quad (10)$$

Where p_{ij} lies in the position given by row i and column j . Note that this matrix is sparse. Since it is nonnegative and each column adds the unity this matrix is called a stochastic matrix. Now it is easy to see -and here we have a good opportunity to motivate the somewhat strange form of the concept of product of matrices- that we have

$$v_1 = Pv_0, \quad v_2 = Pv_1, \quad v_3 = Pv_2. \quad (11)$$

We note now that

$$v_3 = Pv_2 = PPv_1 = P^2v_1 = P^3v_0 \quad (12)$$

Here the teacher may explain the notation for the product of matrices and their properties. In our example, it is easy to show that for any instant time given by k the following equations hold

$$v_k = Pv_{k-1} \quad k = 1, 2, \dots \quad (13)$$

$$v_k = P^k v_0 \quad k = 0, 1, 2, \dots \quad (14)$$

And (14) implies that the state vector in the instant time k can be computed just using the initial state and the rules of our walker which are given by matrix P . For example, let us consider a random walk on a line with only four permitted positions, i.e., $n = 4$, with the notation of Figure 3. In this case the transition matrix is, according to (10)

$$P = \begin{bmatrix} 0 & 1/2 & 0 & 0 \\ 1 & 0 & 1/2 & 0 \\ 0 & 1/2 & 0 & 1 \\ 0 & 0 & 1/2 & 0 \end{bmatrix} \quad (15)$$

Let us assume that the walker begins its walk at the position "1", i.e., v_0 has the form given by (6). We are interested, for example, in the instant time given by $k = 10$. What can be said about the walker?. An easy computation (here we recommend to use some software such as Matlab™, Derive™ or Mathematica™) shows that

$$v_{10} = P^{10} v_0 = \begin{bmatrix} 1/3 \\ 0 \\ 2/3 \\ 0 \end{bmatrix}, \quad (16)$$

which means that, at the time given by $k = 10$, the probability to find the walker at position "1" is $1/3$ and the probability to find him at position "3" is $2/3$, while the rest of positions allocate the walker with zero probability. In conclusion, the location most likely to find the walker is position "3", when $k = 10$.

We remark here that matrix P , given by (10), does not depend on the time, i.e., it does not depend on k . This means, in words, that our walker has always the same rules. In these cases it is said that the Markov chain is homogeneous, since it is not a function of time. We are now interested in the following topics:

1. Whether there exists a state vector, called stationary state vector, such that

$$v_{est} = P v_{est} \quad (17)$$

2. Whether there exists a limit state vector v_∞ , such that

$$v_\infty = \lim_{k \rightarrow \infty} P^k v_0 \quad (18)$$

3. Whether v_{est} and v_∞ are equal.

In the example above it is easy to see that

$$v_{est} = \begin{bmatrix} 1/6 \\ 1/3 \\ 1/3 \\ 1/6 \end{bmatrix}, \quad (19)$$

and this is an stationary state vector since it satisfies (17) where P is given by (15). Here the teacher can suggest the following homework for the students: to compute that this system oscillates, when k is large enough (taking $k=15$ the phenomenon is already observed) between the states

$$v_m = \begin{bmatrix} 1/3 \\ 0 \\ 2/3 \\ 0 \end{bmatrix}, \quad v_{m+1} = \begin{bmatrix} 0 \\ 2/3 \\ 0 \\ 1/3 \end{bmatrix} \quad (20)$$

We remark that (19) means that in our example we have that the vector v_∞ does not exist when v_0 is given by (6). Therefore, this is an example showing that v_∞ and v_{est} are not equal. However, the PageRank vector is a vector that satisfies $v_{est} = v_\infty$. Indeed, the matrix P for the Google problem has another interesting property: the equality $v_{est} = v_\infty$ is satisfied for any choice of the initial state vector v_0 . In the next section we focus on how the Google matrix is constructed.

THE GOOGLE MATRIX

We have already seen how to associate a connectivity matrix G to a set of pages. Now, we show how to modify G to obtain a stochastic matrix P . Given a matrix $G = (g_{ij})$ we define the quantities

$$c_j = \sum_{i=1}^n g_{ij}, \quad 1 \leq j \leq n \quad (21)$$

and the stochastic matrix $P = (p_{ij})$ given by

$$p_{ij} = \begin{cases} g_{ij} / c_j & \text{if } c_j \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad 1 \leq i, j \leq n \quad (22)$$

In our example, where G is given by (3), we have that P results to be

$$P = \begin{pmatrix} 0 & 1/3 & 1/2 & 1/3 \\ 1/2 & 0 & 1/2 & 1/3 \\ 1/2 & 1/3 & 0 & 1/3 \\ 0 & 1/3 & 0 & 0 \end{pmatrix} \quad (23)$$

Matrix P given by (23) is called the Google matrix of the web considered.

GOOGLE'S PAGERANK

We have said that the PageRank vector is a vector that satisfies $v_{est} = v_{\infty}$ and this is true for any choice of the initial state vector v_0 . Therefore, if we denote by x the PageRank vector we have, from (17), that

$$Px = x \quad (24)$$

And this is a particular case of the eigenvector problem that we study in a first course of university. We recall that given a square matrix A , the nonzero vectors x that satisfy the equation

$$Ax = \lambda x, \quad (25)$$

with λ real (or complex) numbers, are called eigenvectors of A associated with λ . Therefore, the computation of the PageRank vector is an eigenvector problem! In fact, the PageRank vector is a positive eigenvector corresponding to the eigenvalue $\lambda = 1$ and such that their components sum the unity. In general, the Google matrix P given by (22) will not have the desired properties and some minor modifications have to be made in order to ensure that the matrix admits an eigenvector with the properties mentioned above. The Perron-Frobenius theorem is the base to perform these minor modifications. In more detail, when P is stochastic and primitive (nonnegative, irreducible¹ and $\lambda = 1$ is the only eigenvalue which norm is the spectral radius²) then the PageRank vector exists. In this case it is said that the Markov chain is ergodic. See, for example, [4] for details.

A LINEAR SYSTEM FOR GOOGLE

Some authors [5] [6] showed that the PageRank problem can be formulated as the following linear problem

$$(I - \alpha P)y = v \quad (26)$$

¹ A matrix is irreducible if its associated directed graph is such that one can go from any vertex to any other vertex, perhaps in several steps.

² The spectral radius is the maximum of the absolute values of the eigenvalues

Where v is a probability distribution vector and $0 < \alpha < 1$. Once the unknown y has been solved the PageRank vector x is then computed using the expression

$$x = \frac{y}{e^T y} \quad (27)$$

where $e^T = [1 \ 1 \ \dots \ 1]$. The linear system (26) is formally a system like (1) but with the coefficient matrix $A = I - \alpha P$. This matrix results to be a nonsingular M-matrix [6] which means that some known iterative methods for linear systems can be applied for solving (26). In the next section we introduce these kinds of methods.

ITERATIVE METHODS

We have shown that the PageRank problem can be seen either as (24) if we consider the problem as an eigenvalue problem or as a linear system like (26). When we compute eigenvectors in the classroom we use matrices of order 3 or 4 and we do the computations by hand. Nevertheless when we have a matrix of the size of billions we can not use the direct methods that we show in the classroom. Not even with a computer a direct methods would be useful! We have to deal with iterative methods. An easy way to introduce the nature of the iterative methods is the following. Let us consider the following equation with one unknown

$$3x = 6 \quad (28)$$

Clearly, the solution is $x = \frac{1}{3}6 = 2$. Let us imagine

that we don't know how to compute the fraction $\frac{1}{3}$ and thus we have to invent a method to solve (28) by inspection, i.e., we may give arbitrary values to x in the hope that we chose the exact value that satisfies (28). This method has the disadvantage that we have infinite values of x to prove. Obviously, in the case of having a system of n linear equations the possibility of finding a solution by inspection is very small. Let us suppose, instead, that we do know how to compute the fraction $\frac{1}{5}$ and we do the following: we split the number 3 as $3 = 5 - 2$ and therefore (28) can be written as

$$(5 - 2)x = 6 \quad (29)$$

and this allows us to write

$$x = \frac{1}{5}2x + \frac{1}{5}6 \quad (30)$$

and now we apply the following: a random value of x would not probably solve this equation and therefore the right hand side of (30) would be different than its left hand side. Let us

denote by x_0 the initial guess that we prove on the right hand side of (30) and by x_1 the value that we obtain in the left hand side, i.e.,

$$x_1 = \frac{1}{5}2x_0 + \frac{1}{5}6 \quad (31)$$

And we want to fulfill (31) with $x_1 = x_0$. Since we shall prove a random value of x_0 it is not likely that $x_1 = x_0$. Therefore we will have to repeat this method but taking now as a guess the value of x_1 . In conclusion we formally have to deal with the successive approximations

$$x_{k+1} = \frac{2}{5}x_k + \frac{6}{5}, \quad k = 0,1,2,\dots \quad (32)$$

This expression is called an iterative scheme. In the Table 1 we show the values for the successive x_k when we begin with $x_0=1.0$.

TABLE I
VALUES FOR x_k GIVEN BY (32)

k	x_k
0	1.000
1	1.600
2	1.840
3	1.936
4	1.974
5	1.990
6	1.996
7	1.998
8	1.999

In Table 1 it is shown that x_k converges to 2.0 as k becomes greater. In this point the teacher may explain the difference between convergent and divergent methods and the theorems related. In the case of linear systems of the form $Ax = b$ the convergence of the iterative methods can be studied in terms of the splitting of matrix A in the form $A = M - N$ where M is a nonsingular matrix, and considering the iterative scheme

$$x_{k+1} = M^{-1}Nx_k + M^{-1}b, \quad k = 0,1,2,\dots \quad (33)$$

It is known [4] that when A is nonsingular the convergence of this scheme it is guaranteed when the spectral radius of the iteration matrix $M^{-1}N$ is less than 1.0. When the coefficient matrix A is a singular matrix (as in the case of Google) the convergence conditions are more complicated and require, among other conditions, that the spectral radius of $M^{-1}N$ be less or equal than 1.0 ; see [7], [8].

The classical method of solving the eigenvalue problem (23) in an iterative way is using the scheme

$$x_{k+1} = Px_k \quad k = 0,1,2,\dots \quad (34)$$

which is called *Power method* and actually is the method that Google uses to compute the PageRank vector; see [1] and [6] for details.

CONCLUSIONS

Some concepts related with Google's PageRank has been shown. In each section we have motivated the terms introduced mainly relating them with the performing of the inner algorithm of the Google searcher. In our experience, these sorts of explanations use to catch the attention of our students. They also allow to illustrate that linear algebra is a broadly applicable branch of mathematics and reveal some connections between separate parts of mathematics.

REFERENCES

- [1] Page, L., Brin, S., Motwani, R. and Winograd, T. "The PageRank Citation Ranking: Bringing Order to the Web". *Stanford Digital Library Technologies Project*. 1999.
- [2] Moler, C. B. "The World's Largest Matrix Computation. Google's PageRank is an eigenvector of a matrix of order 2.7 billion". *MATLAB News & Notes*. 2002.
- [3] Barrett, R., Berry, M. W., Chan, T. F., Demmel, J., Donato, J. et al, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*. SIAM, Philadelphia, Pennsylvania, USA. 1994.
- [4] Berman, A. and Plemmons, R. J. *Nonnegative matrices in the mathematical sciences*, SIAM, Philadelphia, Pennsylvania, USA. 1994.
- [5] Del Corso, G. M., Gulli, A. and Romani, F. "Fast PageRank computation via a sparse linear system", *Lecture notes in computer science*, No 3243, 2004, pp. 118-130.
- [6] Langville, A. N. and Meyer, C. D. "Deeper inside PageRank". *Internet Mathematics*, Vol 1, No. 3, 2005, pp. 335-380.
- [7] Bru, R., Pedroche, F. and Szyld, D. B. "Cálculo del vector PageRank de Google mediante el método aditivo de Schwarz". *Congreso de Métodos Numéricos en Ingeniería* 2005. J.L. Pérez Aparicio et al (ed.), Granada, España. 2005. p. 263.
- [8] Pedroche, F. "Métodos de cálculo del vector PageRank". *Bol. Soc. Esp. Mat. Apl.* 39, pp. 7-30. 2007.
- [9] Nelson, R. *Probability, stochastic processes, and queueing theory*, Springer-Verlag, Ney York, USA. 1995.