## A MODIFICATION OF THE COLLEY MATRIX METHOD

### COLE R. MANSCHOT

ABSTRACT. The Colley Matrix Method avoids the issue of singular matrices to provide a simple method for ranking teams without bias or ad hoc adjustments. When modifying the Colley Method to use a pure win percentage, the resulting equivalent matrix equation is no longer full rank causing a complication of the method. There is a simple solution that still exists through a closer analysis of the methodology.

### 1. INTRODUCTION

Before the Bowl Championship Series, the National Collegiate Athletic Association had multiple football teams claiming to be the national champions depending on what newspaper or coaches polls had which school ranked number one at the end of the season. In 2003 this discrepancy between rankings and polls resulted in the selection of LSU and USC as conational champions.

The question of how to rank head-to-head competition has resulted in a variety of methods to determine who is "best." Whether determining which sports team is the best in the league, which page is the best fit in a Google search, or what movie suggestions to make on Netflix, there is a definitive need for a methodology of pairwise comparisons.

The challenge in determining rankings arises over arguments in what factors should be considered. With respect to sports (which will be the focus of this paper), many argue that win margin, strength of schedule, home field advantage, and a myriad of other factors should be accounted for to give an accurate picture and representation of the current standings. Ideally, a ranking methodology should be simple and easy to understand, easily computed, bias-free, and objective.

In 2002, Wes Colley published his methodology for how to rank collegiate football team [1]. The goal was a simple methodology founded in two basic principles:

- Rankings should reflect wins and losses only, not the margin of victory
- Strength of schedule should be accounted for by the rankings in some methodology

This paper will discuss the original methodology offered and used by Wes Colley, the challenges that arise from using a pure win percentage rather than the modified percentage used by Colley, a proposed solution to that challenge, and an example illustrating the slight differences between these models.

#### COLE R. MANSCHOT

## 2. Iterative Scheme

Colley begins his ranking through an iterative process that updates the rank of the team after each step. Each team begins with an initial ranking of one plus their number of wins over two plus their number of total games played, or

(1) 
$$r^{i} = \frac{1 + n_{w,i}}{2 + n_{tot,i}}.$$

Hence before any games have been played, all teams have an unbiased ranking of  $\frac{1}{2}$ . This equality prior to play removes conference or historic bias allowing the rankings to solely update based on current season play. To motivate the need for iterative updating, we simply acknowledge that in NCAA DI Football or Basketball the strength of schedule can vary greatly even within conference play and there is no guarantee that all teams have played teams of the same skill level.

First we define the *effective number of wins* as Colley states in order to adjust for strength of schedule. The purpose is essentially a weighting on the value of the wins (or losses) against teams based on their opponent's rating. The effective number of wins is

$$n_{w,i}^{eff} = \frac{n_{w,i} - n_{l,i}}{2} + \sum_{j} G_{i,j} \cdot r^{j},$$

where  $r^{j}$  is the ranking of the  $j^{th}$  team,  $n_{w,i}$  and  $n_{l,i}$  are the number of wins and loses of the  $i^{th}$  team (respectively), and  $G_{i,j}$  is the number of games played between teams i and j. Note that if all opponents have a ranking of  $\frac{1}{2}$  then  $n_{w,i}^{eff}$  will equal the actual number of wins for team i.

The iterative method works by calculating r for each team, then calculating  $n_{w,i}^{eff}$ , recalculating r with the new number effective wins, then recalculating the effective number of wins with the new ranking, so on and so forth.

In order to obtain the final Colley ranking, one simply takes the limit of the rankings given by the iterative scheme. The existence of this limit and a simple method for computation of the limit can be obtained using linear algebra.

# 3. As a Matrix Equation

For small groups of teams, these calculations may be relatively simple. However, as the number of teams involved grows, it becomes worthwhile to look at this iterative method as a system of equations.

In order to write these rankings in matrix form, it is necessary to define a matrix C and a column vector **b** where

$$C_{i,j} = \frac{G_{i,j}}{2 + n_{tot,i}},$$
$$b_i = \frac{1 + \frac{n_{w,i} - n_{l,i}}{2}}{2 + n_{tot,i}}.$$

If there are N teams, then C is an  $N \times N$  matrix and **b** is a column vector of length N. The matrix C and the vector **b** are related to the iterative Colley scheme in that if the original ranking is  $\mathbf{r}_0$  given by the modified win percentage in (1), then the modified ranking after one iteration is given by

$$\mathbf{r_1} = \mathbf{b} + C\mathbf{r_0}$$

because the entry (1, j) of each side is equivalent.

$$\begin{aligned} r_1^j &= \frac{1 + n_{w,j}^{eff}}{2 + n_{tot,j}} = \frac{1 + \frac{n_{w,j} - n_{l,j}}{2} + \sum_j (G_{i,j} \cdot r_0^j)}{2 + n_{tot,j}} \\ &= \frac{1 + \frac{n_{w,j} - n_{l,j}}{2}}{2 + n_{tot,j}} + \frac{\sum_j (G_{i,j} \cdot r_0^j)}{2 + n_{tot,j}} = \frac{1 + \frac{n_{w,j} - n_{l,j}}{2}}{2 + n_{tot,j}} + \sum_j (\frac{G_{i,j}}{2 + n_{tot,j}} \cdot r_0^j) \\ &= (b + Cr_0)_j. \end{aligned}$$

Continuing to iterate using the Colley method gives a sequence of rankings  $r_1, r_2, r_3, \ldots$  with formulas given by

$$\mathbf{r_1} = \mathbf{b} + C\mathbf{r_0}$$
  

$$\mathbf{r_2} = \mathbf{b} + C\mathbf{r_1} = \mathbf{b} + C\mathbf{b} + C\mathbf{r_0}$$
  
:  

$$\mathbf{r_n} = \sum_{k=0}^{n-1} C^k \mathbf{b} + C^n \mathbf{r_0}.$$

# 4. Solution to the Matrix Equation

In order to calculate the final ranking, we would use an infinite number of iterations, or  $\mathbf{r}_{\infty}$ . Following from  $\mathbf{r}_{\mathbf{n}} = \mathbf{b} + C\mathbf{r}_{\mathbf{n-1}}$ , then

$$\lim_{n \to \infty} \mathbf{r_n} = \lim_{n \to \infty} \sum_{k=0}^{n-1} C^k \mathbf{b} + C^n \mathbf{r_0}.$$

When evaluating this limit and to ensure the existence of a solution, we note that C is non-negative, primitive, and irreducible (since a reducible matrix implies comparing two sets of teams where teams in one set do not play any teams in the other set). We also note that by construction, all row sums  $\sum_{j=1}^{n} C_{i,j}$  are strictly less than one. This implies that the spectral radius  $\rho(C) < 1$ . Then by the Perron-Frobenius theorem, all the eigenvalues  $\lambda$  of C are such that  $|\lambda| < 1$ . Therefore it can be shown that

$$\lim_{n \to \infty} C^n = 0_{N \times N},$$

Where  $0_{N \times N}$  is the  $N \times N$  zero matrix. Using this, the limit can now be evaluated in the following way:

$$\lim_{n \to \infty} \mathbf{r_n} = \sum_{k=0}^{\infty} C^k \mathbf{b} + \mathbf{0}_{n \times n} = (I - C)^{-1} \mathbf{b}.$$

Observing that (I - C) has a non-zero determinant implies the existence of this solution showing that

(2)  $\mathbf{r}_{\infty} = (I - C)^{-1} \mathbf{b}.$ 

### 5. The Modification

In this section, we will examine modifying Colley's iterative procedure by using the pure win percentage  $\frac{n_{win,i}}{n_{tot,i}}$  rather than the modified win percentage used by Colley,  $\frac{1+n_{win,i}}{2+n_{tot,i}}$ . In his paper, Colley states that using a pure win percentage would result in a singular matrix causing insufficient conditions to find a solution to the ranking. Here we will propose a solution to this problem.

This change in computation results in the need to redefine variables from the Colley Matrix Method Background section. We will call the similar variable of C matrix  $\tilde{C}$  and of **b** the column vector  $\tilde{\mathbf{b}}$  such that

$$\tilde{C}_{i,j} = \frac{G_{i,j}}{n_{tot,i}}$$
$$\tilde{b}_j = \frac{n_{w,i} - n_{l,i}}{2} \cdot \frac{1}{n_{tot,i}}$$

where  $n_{i,j}$  is the number of times team *i* played team *j*. It follows from the same iterative procedure that we get a new sequence

(3)  

$$\begin{aligned}
\tilde{\mathbf{r}}_{1} &= \tilde{\mathbf{b}} + \tilde{C}\mathbf{r}_{0} \\
\tilde{\mathbf{r}}_{2} &= \tilde{\mathbf{b}} + \tilde{C}\mathbf{r}_{1} = \tilde{\mathbf{b}} + \tilde{C}\mathbf{b} + \tilde{C}\mathbf{r}_{0} \\
&\vdots \\
\tilde{\mathbf{r}}_{n} &= \sum_{k=0}^{n-1} \tilde{C}^{k}\mathbf{b} + \tilde{C}^{n}\mathbf{r}_{0}.
\end{aligned}$$

In trying to compute  $\lim_{n\to\infty} \tilde{\mathbf{r}}_n$ , there arise two problems:

- The matrix  $(I \tilde{C})$  is no longer non-singular, therefore  $\sum_{k=0}^{\infty} \tilde{C}^k$  no longer converges.
- The eigenvalues  $\lambda$  of  $\tilde{C}$  are no longer strictly less than one, therefore  $\lim_{n\to\infty} \tilde{C}^n$  does not go to 0.

#### MODIFIED COLLEY MATRIX

Both of these complications come from the fact that the matrix  $\tilde{C}$  now has an eigenvalue  $\lambda = 1$ . This follows from the Perron-Frobenius Theorem since the matrix  $\tilde{C}$  has all row sums equal to 1. In fact, the Perron-Frobenius Theorem tells us that there are two possible cases.

- There exist eigenvalues  $\lambda_1 = 1, \lambda_2 = -1$ . However, this only occurs when the schedule graph of games played is bipartite. Since this rarely happens in practice, we will not consider this case here (despite the existence of a solution to this problem).
- There is one eigenvalue  $\lambda_1 = 1$  with all other eigenvalues  $|\lambda_i| < 1, i > 1$ . We will assume this is the case for the rest of the paper.

## 6. The Solution

In the case in which there is one eigenvalue  $\lambda_1 = 1$  and the other eigenvalues are strictly less than one, we note that the row sums of  $\tilde{C}$  are now one and the entries  $\tilde{C}_{i,j} \in [0,1]$ . Therefore, the matrix  $\tilde{C}$  can be thought of as a transition matrix for a stochastic process, specifically a random walk on a graph.

It becomes further necessary to introduce one last vector in order to show the existence of the solution for a non-singular matrix with an eigenvalue of  $\lambda_1 = 1$  with algebraic multiplicity 1. Let  $\pi$  be the unique stationary distribution of  $\tilde{C}$  with eigenvalue of  $\lambda_1 = 1$  such that  $\pi \tilde{C} = \pi$ . Then because  $\tilde{C}$  is the transition matrix for a random walk on a graph, we know that

$$\pi(j) = \frac{n_{tot,j}}{\sum_{i} n_{tot,i}}$$

Recalling the ranking as defined in the iterative scheme,  $\tilde{\mathbf{r}}_{\mathbf{n}} = \sum_{k=0}^{n-1} \tilde{C}^k \tilde{\mathbf{b}} + \tilde{C}^n \tilde{\mathbf{r}}_{\mathbf{0}}$ , then one must show that this limit converges. First, since  $\tilde{C}$  is the transition matrix for an aperiodic random walk on a graph with stationary distribution  $\pi$ , then  $\lim_{n\to\infty} \tilde{C}^n \tilde{\mathbf{r}}_{\mathbf{0}}$  goes to a column vector with all entries equal to  $\pi \cdot \tilde{\mathbf{r}}_{\mathbf{0}} = \frac{1}{2}$ . That is

(4) 
$$\lim_{n \to \infty} \tilde{C}^n \tilde{\mathbf{r}}_0 = \frac{1}{2} \mathbf{1}_{N \times 1}$$

because the rows of a transition matrix of a finite Markov Process approaches the stationary distribution as  $n \to \infty$ .

A little more complicated is the convergence of

$$\lim_{n\to\infty}\sum_{k=0}^{n-1}\tilde{C}^k\tilde{\mathbf{b}}.$$

This sum does not obviously converge since the matrix  $\tilde{C}$  has an eigenvalue of  $\lambda_1 = 1$ . To show this does converge, decompose  $\mathbb{R}^N = W_1 \oplus W_2$  where

$$W_1 = \operatorname{span}\{\mathbf{1}_{N \times 1}\}$$

is the eigenspace for  $\lambda_1 = 1$  and

$$W_2 = \{ v \in \mathbb{R}^N : \pi \cdot v = 0 \}$$

#### COLE R. MANSCHOT

is an invariant subspace for the transformation  $T(v) = \tilde{C}v$ .  $W_2$  is an invariant subspace under  $(\tilde{C})$  because

$$\pi \cdot (\tilde{C}v) = (\pi \tilde{C}) \cdot v = \pi \cdot v = 0, \quad \text{if } v \in W_2$$

Moreover, we claim  $\tilde{\mathbf{b}} \in W_2$ . This is easily seen because

$$\pi \cdot \tilde{\mathbf{b}} = \sum_{k=1}^{N} \frac{n_{tot,k}}{\sum_{i} n_{tot,i}} \cdot \frac{n_{w,k} - n_{l,k}}{2} \cdot \frac{1}{n_{tot,k}} = \frac{1}{2\sum_{i} n_{tot,i}} \sum_{k=1}^{N} (n_{w,k} - n_{l,k}) = 0.$$

From here we define  $T_2: W_2 \to W_2$  as the transformation  $T(v) = \tilde{C}v$  restricted on the subspace  $W_2$ . Then,

$$\lim_{n \to \infty} \sum_{k=0}^{n-1} \tilde{C}^k \tilde{\mathbf{b}} = \lim_{n \to \infty} \sum_{k=0}^{n-1} T_2^k \tilde{\mathbf{b}}.$$

Moreover, T has eigenvalues  $\lambda_1 = 1, |\lambda_i| < 1$ , for i > 1 from the construction above, so  $T_2$  has eigenvalues of strictly less than 1 and  $\mathbf{\tilde{b}} \in W_2$  as shown before. Therefore the limit must exist.

### 7. Computing the Limit

Having proved the limit of (3) as  $n \to \infty$  exists, it remains to show how to calculate this limit in a simple and efficient manner similar to Colley's solution given by (2). We claim that the limit solution to the system using the pure win percentage

$$\tilde{\mathbf{r}}_{\infty} = \sum_{k=0}^{\infty} \tilde{C}^k \tilde{\mathbf{b}} + \frac{1}{2} \mathbf{1}_{N \times 1}$$

is equivalent to the solution from the set of equations

(5) 
$$\begin{cases} (I - \tilde{C})\tilde{\mathbf{r}} = \tilde{\mathbf{b}} \\ \pi \cdot \tilde{\mathbf{r}} = \frac{1}{2} \end{cases}$$

The rank of  $(I - \tilde{C}) = n - \dim(E_{\lambda_1}) = n - 1$ . Showing that  $\pi$  is linearly independent of  $(I - \tilde{C})$  will prove that there exists a unique r satisfying the second solution method. Since  $(I - \tilde{C})$  has row sums of 0, and the row sum of  $\pi$  is non-zero, then it is clear that the rank of

$$\left[\begin{array}{c} (I-\tilde{C})\\ \pi \end{array}\right]$$

is n by direct application of the definition of linear independence and that the rank of the column space equals the rank of the row space.

Therefore,  $\pi$  is indeed row independent of  $(I - \hat{C})$  and there does exist a solution  $\tilde{\mathbf{r}}$  to the second set of equations. Now one must show these are equivalent rankings  $\tilde{\mathbf{r}}$  which will be done through substitution.

$$(I - \tilde{C}) \cdot \tilde{\mathbf{r}}_{\infty} = (I - \tilde{C}) \cdot \left(\sum_{k=0}^{\infty} \tilde{C}^{k} \tilde{\mathbf{b}} + \frac{1}{2} \mathbf{1}_{N \times 1}\right)$$
$$= \left(\sum_{k=0}^{\infty} \tilde{C}^{k} \tilde{\mathbf{b}} + \frac{1}{2} \mathbf{1}_{N \times 1}\right) - \left(\sum_{k=0}^{\infty} \tilde{C}^{k+1} \tilde{\mathbf{b}} + \tilde{C} \frac{1}{2} \mathbf{1}_{N \times 1}\right)$$
$$= \left(\sum_{k=0}^{\infty} \tilde{C}^{k} \tilde{\mathbf{b}} - \sum_{k=0}^{\infty} \tilde{C}^{k+1} \tilde{\mathbf{b}}\right) + \left(\frac{1}{2} \mathbf{1}_{N \times 1} - \frac{1}{2} \mathbf{1}_{N \times 1}\right)$$
$$= \tilde{\mathbf{b}}$$

and

$$\pi \cdot \tilde{\mathbf{r}}_{\infty} = \pi \cdot \left( \sum_{k=0}^{\infty} \tilde{C}^k \tilde{b} + \frac{1}{2} \mathbf{1}_{N \times 1} \right) = \sum_{k=0}^{\infty} \pi \tilde{\mathbf{b}} + \pi \frac{1}{2} \mathbf{1}_{N \times 1}$$
$$= 0 + \sum_{j=0}^n \frac{n_{tot,j}}{2 \sum_i n_{tot,i}}$$
$$= \frac{\sum_{j=0}^n n_{tot,j}}{2 \sum_i n_{tot,i}}$$
$$= \frac{1}{2}.$$

Therefore,  $\tilde{\mathbf{r}}_{\infty}$  is in fact the unique solution to the system (5).

# 8. B1G BASKETBALL 2015-16 SEASON

To illustrate that both ranking methodologies yield a similar solution, we will use the B1G Basketball (pre-tournament) inter-conference play for the 2015 to 2016 season. The reason that non-conference play will be excluded is to maintain the small size of teams necessary to complete the ranking.

A side-by-side comparison of the rank value r, the rank, and the conference ranking is shown in the table below.

University	Colley $r$	Rank	Modified Colley $\tilde{r}$	Rank	Conference Rank
Indiana	0.7686	1	0.7950	1	1
Michigan State	0.6873	2	0.7068	2	2
Iowa	0.6706	3	0.6919	3	5
Wisconsin	0.6704	4	0.6916	4	6
Purdue	0.6621	5	0.6819	5	3
Maryland	0.6619	6	0.6813	6	4
Ohio State	0.5801	7	0.5866	7	7
Michigan	0.5504	8	0.5564	8	8
Northwestern	0.4361	9	0.4274	9	9
Penn State	0.4141	10	0.4062	10	10
Nebraska	0.3514	11	0.3349	11	11
Illinois	0.2985	12	0.2757	12	12
Minnesota	0.1501	13	0.1109	13	13
Rutgers	0.0984	14	0.0534	14	14

It is clear from the table that the values of r differ slightly between the methods, yet they both preserve the average ranking of 1/2.

The Big Ten determines the Conference Rank strictly based on inter-conference play and the win percentage. The differences between the Colley rank, Modified Colley rank, and the Conference rank come from the tie-breaking rules used by the Big Ten when teams have the same in-conference win percentage, which in this case was a four way tie between Iowa, Maryland, Purdue, and Wisconsin.

### 9. CONCLUSION

When computing the ranking of head-to-head competition, ultimately one can maintain the same principles of the Colley Matrix Method of no bias or ad hoc adjustments while using an unmodified win percentage. This simple solution is the one that satisfies the system (5).

### References

[1] Wesley N. Colley. Colley's Bias Free College Football Ranking Method: The Colley Matrix Explained. www.colleymatrix.com. May 30, 2002.