

10 Section (viii): PageRank Algorithm

1

We are now ready to explore an algorithm that you, yourself, have probably already used today (as well as perhaps a billion other people - estimates are that there are around 3.5 billion PageRank queries run per day)! The PageRank algorithm is the engine behind Google Search, and today it is commonly ranked as one of the most influential & important algorithms in the field of data science (source: Springer).

Put simply, PageRank is an algorithm used to rank websites in search engine results, by applying a numerical weighting/measure to a semantically-relevant subgraph of the web.

As we shall see, this ranking procedure bears a strong resemblance to several of the centrality computations - particularly Katz centrality - that we have encountered.

The key difference between PageRank and several of these other centrality measures concerns what might be called "prestige dilution." Recall that with Katz centrality, nodes/webpages of high prestige or importance tend to ^{freely} propagate this prestige to their neighbors, and so on with their neighbors' neighbors. Depending on the particular real-world network of

interest, This "dilution" phenomenon might not be malum in se.

However, with respect to ranking the importance of linked documents/webpages, such centrality dilution is potentially parlous.

For instance, the wikipedia webpage should, naturally, carry a high centrality measure (and indeed it does according to PageRank's criteria, as any empirical Google search will attest).

But suppose that wikipedia links directly a great many auxiliary webpages (that is, webpages with an intrinsically small centrality/importance) - indeed it does link to many of these sorts of webpages.

If one used Katz centrality to rank pages in this network, then the prestige of wikipedia would spread to a multitude of less important pages. Certainly this kind of dilution is generally undesirable for webpage rankings. PageRank proposes a simple solution.

We define centrality now as a variation of Katz centrality in which the centrality a vertex derives from its neighbors is proportional to their centrality divided by their out-degree.

(Note that our graph is a digraph here).

Mathematically, then, we define the centrality of node x_i as:

$$x_i = \alpha \sum_j A_{ij} \left(\frac{x_j}{k_{out}^j} \right) + \beta$$

Where, in the previous formula, A_{ij} are elements of the adjacency matrix for the network, α is an "attenuation factor" (see below), and β is a "free" parameter essentially enabling vertices with zero in-degree to collect centrality - as seen with Katz centrality measures; k_j^{out} denotes the out-degree of vertex j in the graph.

Put succinctly, this variant of "PageRank centrality" (note that the algorithm/computations presented here represent a simplified version of the True PageRank algorithm - some of the finer details are, of course, well-kept industry trade secrets), is basically a "normalized" version of Katz centrality.

Looking at the formula on the previous page, the astute reader should question: what happens if $k_j^{out} = 0$? In this case our formula contains the indeterminate expression $(\frac{0}{0})$, which we would like to avoid. We eschew this problem by artificially setting $k_j^{out} = 1$ for all such vertices (note that this "inflationary" term does not affect the exact solution / iterative convergence of the PR algorithm).

As before, we prefer to write the previous formula in matrix form:

$$\vec{x} = \alpha A D^{-1} \vec{x} + \beta \vec{1}$$

matrix D here is the Diagonal Matrix defined as $D_{ii} = \max(1, k_i^{out})$.

Where does D^{-1} come from? Recall that for a diagonal matrix,

4

$$D = \begin{bmatrix} a_{11} & & \phi \\ & a_{22} & \\ \phi & & \ddots \\ & & & a_{nn} \end{bmatrix}, \quad D^{-1} \text{ is simply the new diagonal matrix consisting of the reciprocal values of } D.$$

$$\text{Thus } D^{-1} = \begin{bmatrix} \frac{1}{a_{11}} & & \phi \\ & \frac{1}{a_{22}} & \\ \phi & & \ddots \\ & & & \frac{1}{a_{nn}} \end{bmatrix}, \quad \text{noting that } a_{ii} \neq 0 \text{ for all } 1 \leq i \leq n.$$

Reading the formula $\vec{x} = \alpha A D^{-1} \vec{x} + \beta \vec{1}$, we can now see that the left-term ($\alpha A D^{-1} \vec{x}$) propagates centrality from a vertex to its neighbors, "normalized" by vertex-out-degree & scaled, additionally, by the attenuation factor α . The right-term ($\beta \vec{1}$) effectively scales this centrality to avoid penalizing the zero in-degree vertices.

What is α ? Google Search sets $\alpha = 0.85$ for PageRank computations.

In one sense (recall convergence of the eigenvector centrality measures)

This damping/attenuation < 1 guarantees convergence of the iterated calculations of PR measures. Secondly, we can think of $\alpha = 0.85$ as an "empirically-discovered" (i.e. it "works")

hyperparameter value, reflecting the fact that we believe in general

that there is an $1 - \alpha = 15\%$ chance a given user won't follow any of the ranked pages suggested by Google Search.

As is typical, we set $\beta=1$ & solve for \vec{x} :

$$\vec{x} = \alpha AD^{-1} \vec{x} = \vec{1}$$
$$\underbrace{(-\alpha AD^{-1} \vec{x})}_{(-\alpha AD^{-1} \vec{x})}$$

$$\rightarrow \vec{x} - \alpha AD^{-1} \vec{x} = \vec{1}$$

$$\rightarrow \underbrace{(I - \alpha AD^{-1})}_{\text{invertible}} \vec{x} = \vec{1}$$

$$\rightarrow \underbrace{(I - \alpha AD^{-1})^{-1}}_{=I} \cdot (I - \alpha AD^{-1}) \vec{x} = (I - \alpha AD^{-1})^{-1} \vec{1}$$

$$\rightarrow \vec{x} = (I - \alpha AD^{-1})^{-1} \vec{1}$$

Notice the two instances of inverse operations. This is computationally inefficient, so we clean up.

$$\rightarrow \vec{x} = \underbrace{((D - \alpha A) D^{-1})^{-1}}_{\text{pull out } D^{-1}} \vec{1}$$

$$\text{(Recall: } (AB)^{-1} = B^{-1}A^{-1} \text{)}$$

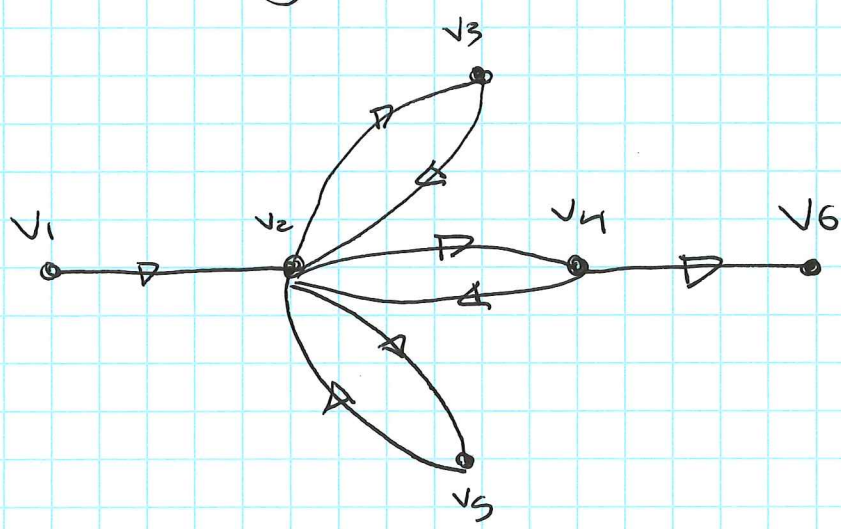
$$\rightarrow \boxed{\vec{x} = D (D - \alpha A)^{-1} \vec{1}}$$

key formula! \rightarrow

PageRank works well ~~very~~ web-ranking for the reason that having links to your page from important pages elsewhere is a strong indication that your page may also be important. The added ingredient of normalizing centrality by out-degrees of pages that merely point to a great many pages do not freely bestow their centrality to a multitude of other pages.

Ex.

We compute the PageRank score for each vertex in the following network.



$$D = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

Recall: $D_{ii} = \max(1, k_i^{out})$

Adjacency Matrix

We compute the solution to the PR algorithm system:

$$\vec{x} = D(D - \alpha A)^{-1} \vec{1}, \text{ with } \alpha = .85. \text{ This yields:}$$

$$\vec{x} = \begin{bmatrix} 3.706 \\ 11.23 \\ 3.706 \\ 4.88 \\ 3.706 \\ 2.763 \end{bmatrix}$$

Consequently, the ranking (from high to low) of vertices by the PR algorithm is as follows:

$$\{v_2, v_4, v_1, v_3, v_5, v_6\}$$