

Google produced search results based on the connectivity of webpages.

---

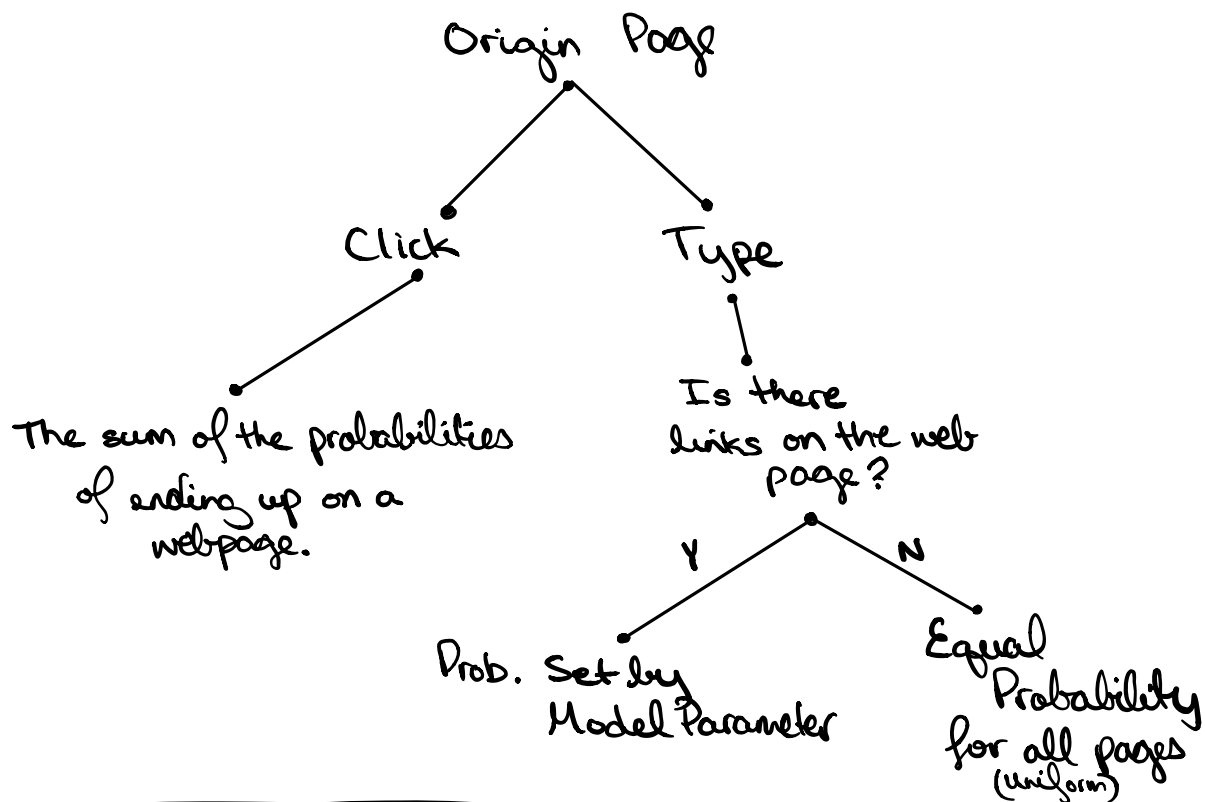
A web-surfer navigates the web through these actions,

① Clicking on webpage links

② Navigating to a new page by typing.

---

A workflow would look like,

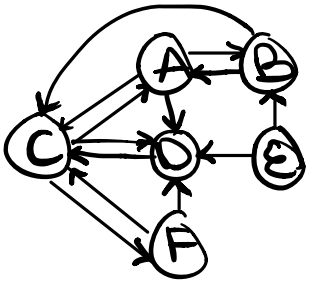


---

Inherent in the design of PageRank algorithm is,

- the underlying connectivity is reflective a page's search rank.
  - second is that important websites link to other important websites
-

# A Graph of Page Connectivity



Each node represents a page and each **directed edge is a link** that connects the source webpage to the target.

**Page surfers** are modelled as **making random walks**

choosing edges travel with a uniform probability dictated by the number of outbound edges of the current node.

Given the random walks, the graph, and a random initial state of users in the network - the traffic through those cycles will converge. Users can then be found predictably on certain pages with certain ratios/densities.

**Page surfers** can also **teleport** (by typing).

The number of page surfers on each page once the network activity converges is expressed by the vector  $\underline{r}$ .

$$\underline{r} = \begin{bmatrix} r_A \\ r_B \\ r_C \\ r_D \\ r_E \\ r_F \end{bmatrix}$$

The number of users on each website can be generalized with the expression

$$r^{i+1} = L r^i$$

This models the number of users per page at minute  $i$ . It holds true before and after convergence.

After convergence  $r^{(i+1)} = r^i$ , can instead write

$$L r = r$$

This is an eigenequation for matrix  $L$  with eigenvalue 1, determined by the probabilistic structure of the matrix  $L$ . (i.e. network has no loss of <sup>normalized pr.</sup> const. traffic)

The matrix  $L$ .

$$\begin{bmatrix} L_{A \rightarrow A} & L_{B \rightarrow A} & L_{C \rightarrow A} & L_{D \rightarrow A} & L_{E \rightarrow A} & L_{F \rightarrow A} \\ L_{A \rightarrow B} & L_{B \rightarrow B} & L_{C \rightarrow B} & L_{D \rightarrow B} & L_{E \rightarrow B} & L_{F \rightarrow B} \\ L_{A \rightarrow C} & L_{B \rightarrow C} & L_{C \rightarrow C} & L_{D \rightarrow C} & L_{E \rightarrow C} & L_{F \rightarrow C} \\ L_{A \rightarrow D} & L_{B \rightarrow D} & L_{C \rightarrow D} & L_{D \rightarrow D} & L_{E \rightarrow D} & L_{F \rightarrow D} \\ L_{A \rightarrow E} & L_{B \rightarrow E} & L_{C \rightarrow E} & L_{D \rightarrow E} & L_{E \rightarrow E} & L_{F \rightarrow E} \\ L_{A \rightarrow F} & L_{B \rightarrow F} & L_{C \rightarrow F} & L_{D \rightarrow F} & L_{E \rightarrow F} & L_{F \rightarrow F} \end{bmatrix}$$

Each column represent the (sum=1) probability leaving a page for another page. Each row represents how likely you are to enter a website from another (doesn't need to sum=1).

Can try to calculate each of the eigenvectors of the network graph; however this isn't scalable for high dimension systems. Instead power iteration can be used - the power series.

To add the opportunity for walks to jump/teleport to another vertice. The model can be adjusted to include additional variables

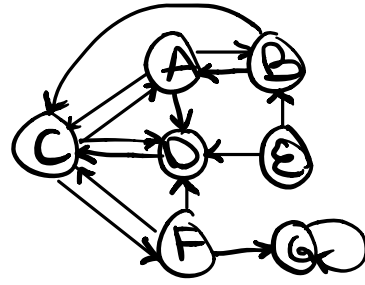
$$M = dL + \frac{1-d}{n} J, \quad d \text{ is the prob of clicking a link}$$

Typically  $1-d$  is small.

$1-d$  is the prob of typing

$J_{n \times n}$  is a matrix of 1 elements computing to equal probabilities of leaving any page for any other page.

## THE DAMPING PARAM



The power iteration applied to the network graph will yield  $\textcircled{G}$  as the highest rank page. This is due to an oversimplified model which only allows for walks - once a path enters the reflexive  $\textcircled{G}$  node the walk only has the opportunity to revisit the  $\textcircled{G}$  node.