# Data Mining and Machine Learning

# Page Ran

Peter Jančovič

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Objectives

- To understand the basic idea of the PageRank of a document in a corpus

- To understand how to calculate PageRank

- To understa https://eduassistpro.github.io/ underlies PageRank

Assignment Project Exam Help

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Not all documents are equal

- So far, whether or not a document *d* is retrieved in response to a query *q* depends only on *sim*(*q,d*)

- Assumption is that all documents are equal - relevance of ~~~~ ry depends only on the similarity score

- This is clearly not true (co ~~~~ *ipedia* with my home page)

- Prior importance of a document is its <u>Page rank</u>

- Probabilistic interpretation of Page rank

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# The *prior* probability of a document

- Suppose that we could assign a probability $P(d)$ to each document $d$ in our corpus

- Think of $P(d)$ as the probability that $d$ is a relevant document b _____ a query $q$

- $P(d)$ is the prior (or *a prio* _____ ility of $d$

- In this case, whether $d$ is returned in response to a query $q$ depends on $sim(q,d)$ and $P(d)$

- We will treat $P(d)$ as the Page rank of $d$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Retrieval using prior probabilities

- Retrieval based only on *sim*(*q,d*) assumes that *P*(*d*) is the same for all documents

- This case is called equal priors

- Intuitively w https://eduassistpro.github.io/ stimate more meaningful priors

- Assumption: the *prior* relevance of a document to any query is related to how often that document is accessed

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Citation indices

- Similar idea used to measure quality of academic papers

- If a paper $p$ contains important results or ideas, then lots of papers will refer to it

- The citations i ⎯⎯⎯⎯⎯ any papers refer to a given paper $p$

- Citations index is a standard q ⎯⎯⎯ re in research assessment

- But, quality of a paper depends not only on the <u>quantity</u> of papers that cite it but on their <u>quality</u> – their citation indices

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Basics of Page Rank

- For a document, or a page, *d* on the web, we could defined the <u>Page Rank</u> *pr*(*d*) to be the number of documents that have a hyperlink to *d*

- This relies o ~~cracy of~~ the web – users 'vote with their m ns'

- The ranking of a documen onse to *q* depends on both *sim*(*q,d*) <u>and</u> *pr*(*d*)

- But not all links are equal

Data Mining and Machine Learning

UNIVERSITY<sup>OF</sup> BIRMINGHAM

# The "Random Surfer Model"

- The solution is to allocate a <u>weight</u> of $w_{de}$ to the hyperlink from document $d$ to document $e$

- $w_{de}$ can be thought of as the <u>probability</u> of following the link to p age $d$

- If $l(d)$ denotes the number nks from $d$, setting $w_{de} = 1/l(d)$ corres e <u>random surfer model</u>:  on any page any of the available links are chosen with equal probability

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# The "Intentional Surfer Model"

- In reality all links on a page are not clicked with equal probability

- A better alternative is to estimate the $w_{de}$s using actual statist                                    surfers

- This is the in_____

- Organizations like Google            store this kind of information (I assume!)

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Simplified Page Rank Calculation

- Once $pr(d)$ is accepted as a measure of the importance of $d$ there is a natural consequence

- In the calculation of $pr(d)$, a hyperlink from a page $d_1$ to $d$ shoul <span style="color:red">Assignment Project Exam Help</span> a hyperlink from page $d_2$ to $d$ <span style="color:red">https://eduassistpro.github.io/</span>

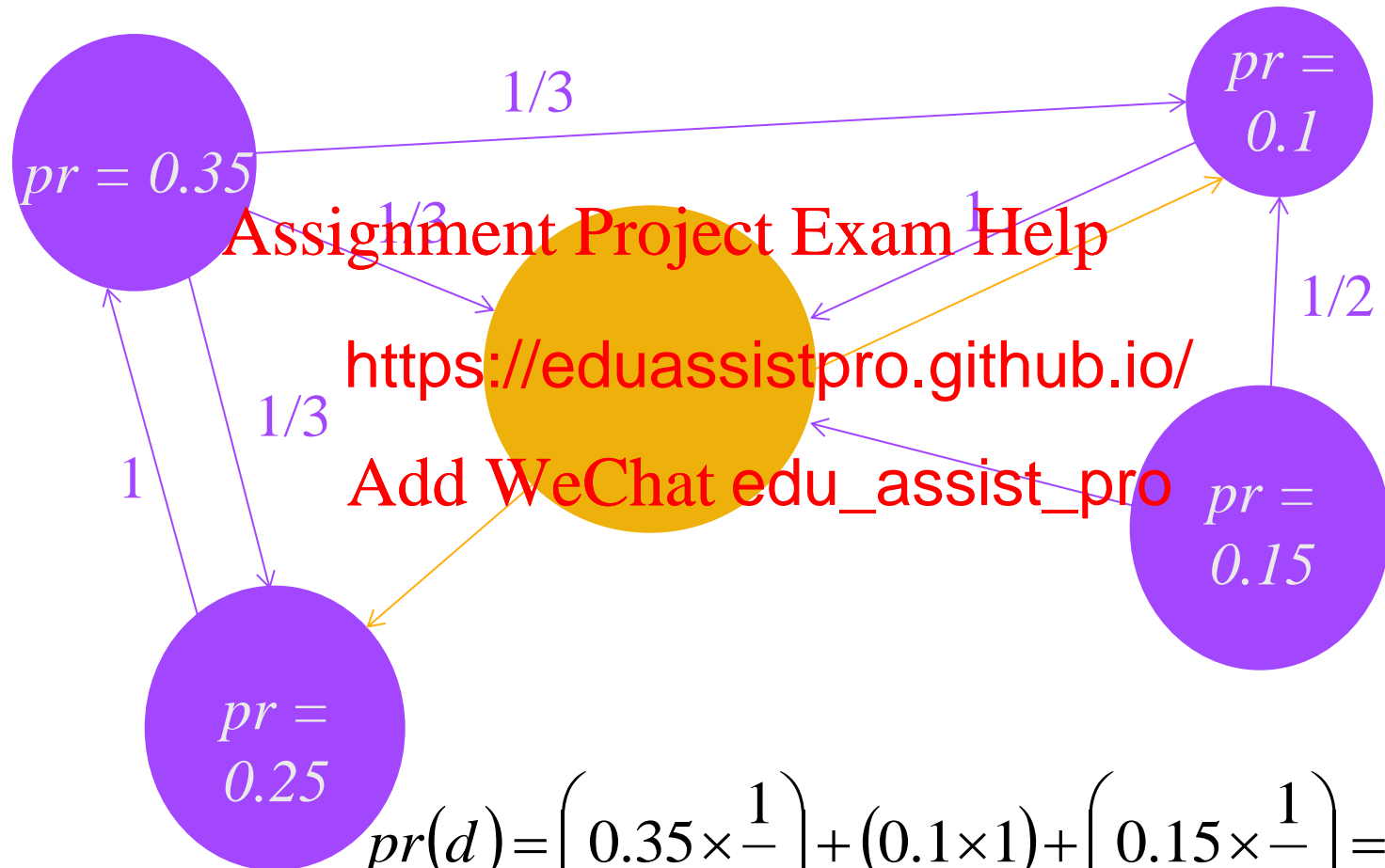- This motivates: <span style="color:red">Add WeChat edu_assist_pro</span>

$$pr(d) = \sum_{e \in L(d)} pr(e) w_{ed}$$

  where $L(d)$ is the set of pages which link to page $d$

- This is the <u>simplified Page rank</u> calculation

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Simplified Page Rank Calculation



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$$pr(d) = \left(0.35 \times \frac{1}{3}\right) + \left(0.1 \times 1\right) + \left(0.15 \times \frac{1}{2}\right) = 0.292$$

UNIVERSITY$^{OF}$
BIRMINGHAM

# Example

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Taken from wikipedia: see http://en.wikipedia.org/wiki/PageRank

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Simplified Page Rank Calculation

- Of course, changing *pr*(*d*) will change the Page Ranks of the other pages, which in turn will change *pr*(*d*)....  Assignment Project Exam Help

- Hence the d  https://eduassistpro.github.io/  is recursive, and *pr*(*d*) is calculated iterativ

  Add WeChat edu_assist_pro

$$pr_{n+1}(d) = \sum_{e \in L(d)} pr_n(e) w_{ed}$$

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Markov Chain interpretation

- Let

$$W = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1D} \\ w_{21} & w_{22} & \cdots & w_{2D} \\ \vdots & \vdots & & \vdots \\ w_{d1} & w & & \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & \cdots & \vdots \\ w_{D1} & w_{D2} & \cdots & w_{DD} \end{bmatrix}$$

where $w_{ij}$ is the probabilit               following a hyperlink between the $i^{th}$ and $j^{th}$ pages and $D$ is the number of pages – this is the <u>page transition probability matrix</u>

- Notice that each row of $W$ sums to 1

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Markov Chain interpretation

- Let $pr_n^T = [pr_n(1), pr_n(2), ..., pr_n(D)]$ - $pr_n(i)$ is the Page Rank of the $i^{th}$ page after $n$ iterations

- Then $pr_{(n+1)} = W^T pr_n$, or $pr_n = (W^T)^n pr_0$

- In Markov C $\qquad$ is the <u>transition</u> <u>probability</u> f <span style="color:red">https://eduassistpro.github.io/</span>

- Can think of $w_{de}$ as the <span style="color:red">Add WeChat edu_assistpro</span> page $e$ at time $t+1$ given page $d$ at time $t$: $P(e @ t+1 \mid d @ t)$

- $pr_n$ is an estimate of the probability distribution over all of the pages after the $n^{th}$ iteration

- In this case $\displaystyle\sum_d pr_n(d) = 1$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Markov chain interpretation

$$
\begin{bmatrix}
pr_{n+1}(1) \\
pr_{n+1}(2) \\
. \\
pr_{n+1}(d) \\
. \\
. \\
pr_{n+1}(D)
\end{bmatrix}
=
\begin{bmatrix}
w_{11} & & & \\
w_{12} & w_{22} & & w_{D2} \\
. & . & . & . \\
w_{1d} & w_{2d} & \cdots & w_{Dd} \\
. & . & \cdots & . \\
. & . & \cdots & . \\
w_{1D} & w_{2D} & \cdots & w_{DD}
\end{bmatrix}
\begin{bmatrix}
pr_n(1) \\
. \\
. \\
. \\
. \\
pr_n(D)
\end{bmatrix}
$$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Markov Chain interpretation

- If this system converges, then

$$
\begin{bmatrix} pr(1) \\ pr(2) \\ . \\ pr(d) \\ . \\ . \\ pr(D) \end{bmatrix} = \begin{bmatrix} w_{11} & w & & \\ w_{12} & w_{22} & & w_{D2} \\ . & . & . & . \\ w_{1d} & w_{2d} & w_{bd} & . \\ . & . & . & . \\ . & . & . & . \\ w_{1D} & w_{2D} & ... & w_{DD} \end{bmatrix} \begin{bmatrix} pr(1) \\ . \\ . \\ . \\ . \\ pr(D) \end{bmatrix}
$$

- $pr = W^T pr$

- In other words $pr$ is an <u>eigenvector</u> of $W^T$ with eigenvalue 1

Data Mining and Machine Learning

UNIVERSITY$^{OF}$
BIRMINGHAM

# Damping Factor

- The model we have used to develop Page Rank is a "random surfer" model with 'proper' hyperlink probabilities

- The random <span style="color:red">stop clicking</span>

- The probability that the ra er continues clicking when he arrives a called the <u>damping factor</u> and denoted by $\delta$

- A typical value of $\delta$ is 0.85

<span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">https://eduassistpro.github.io/</span>

<span style="color:red">Add WeChat edu_assist_pro</span>

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Page Rank

- Taking into account the damping factor,

$$pr(d) = \left(\frac{1-\delta}{N}\right) + \delta\left(\sum_{e \in L(d)} pr(e) \times w_{ed}\right)$$

where *N* is the number of documents

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Notes

- Assuming that *p*(*e*) is the probability of the page d, then this formula preserves $\sum_{d} pr(d) = 1$

- The formula assigns a "floor" value of $\frac{1-\delta}{N}$ to a page that has                                    ks (so that it has non-zero pag

- In addition, the damping fa                    es the effect of past estimates of PageRank on the present estimate

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Notes

- This lecture presents a probabilistic approach to Page rank

- "PageRank" is a trademark of Google

- It was devel https://eduassistpro.github.io/ 1995 and 1998

- Larry Page is one of the founders of Google Inc.

- A high PageRank is a valuable asset for a www page, for example to attract advertising

- Hence the precise details of the Google PageRank algorithm are secret!

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM