

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Introduction to  
Assignment Project Exam Help  
**Informa** |  
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`  
Lecture 18: Li

# Assignment Project Exam Help

## Today's lecture

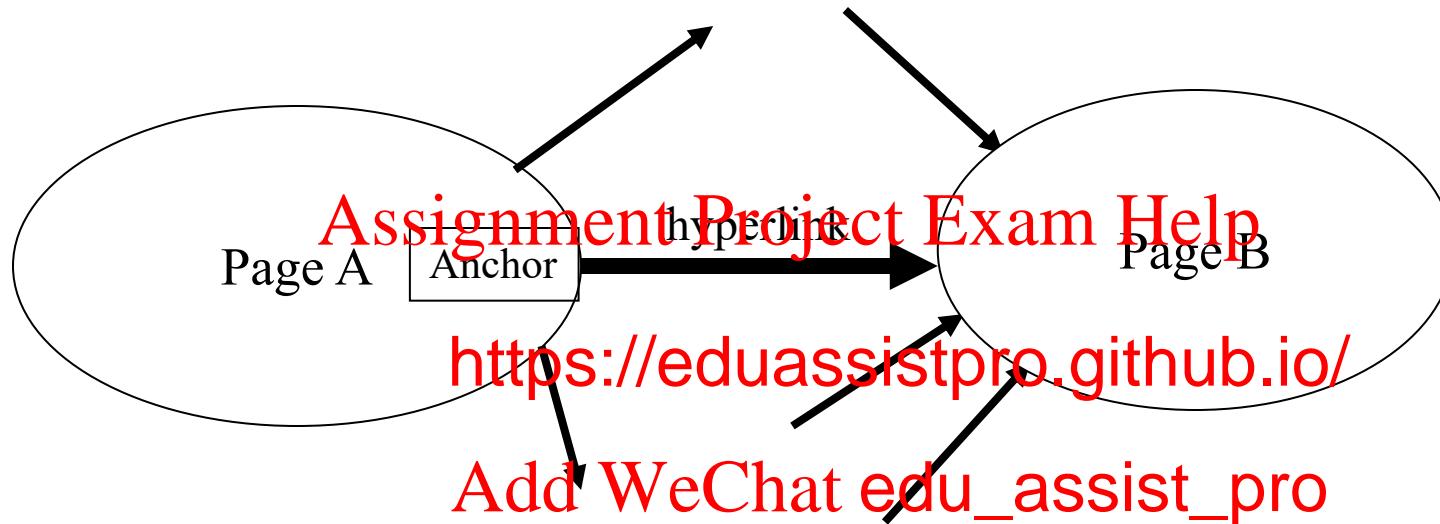
- Anchor text
- Link analysis for ranking
  - Pagerank and variants

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## The Web as a Directed Graph



**Assumption 1:** A hyperlink between pages denotes author perceived relevance (quality signal)

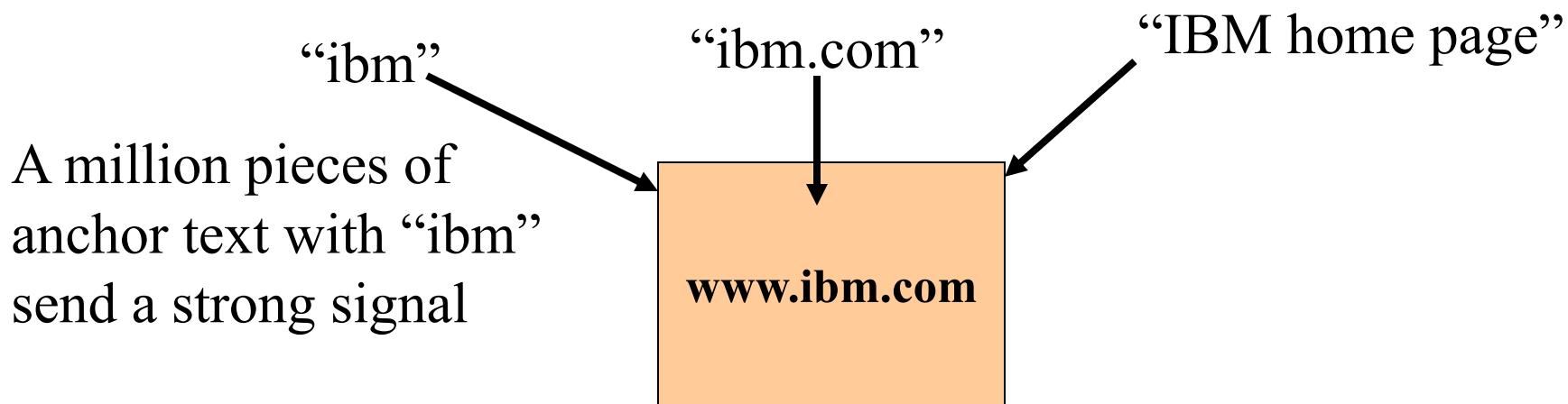
**Assumption 2:** The text in the anchor of the hyperlink describes the target page (textual context)

# Anchor Text

www Worm Add WeChat edu\_assist\_pro

- For **ibm** how to distinguish between:
  - IBM's home page (mostly graphical)
  - IBM's copyright page (high term freq. for 'ibm')
  - Rival's spa <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

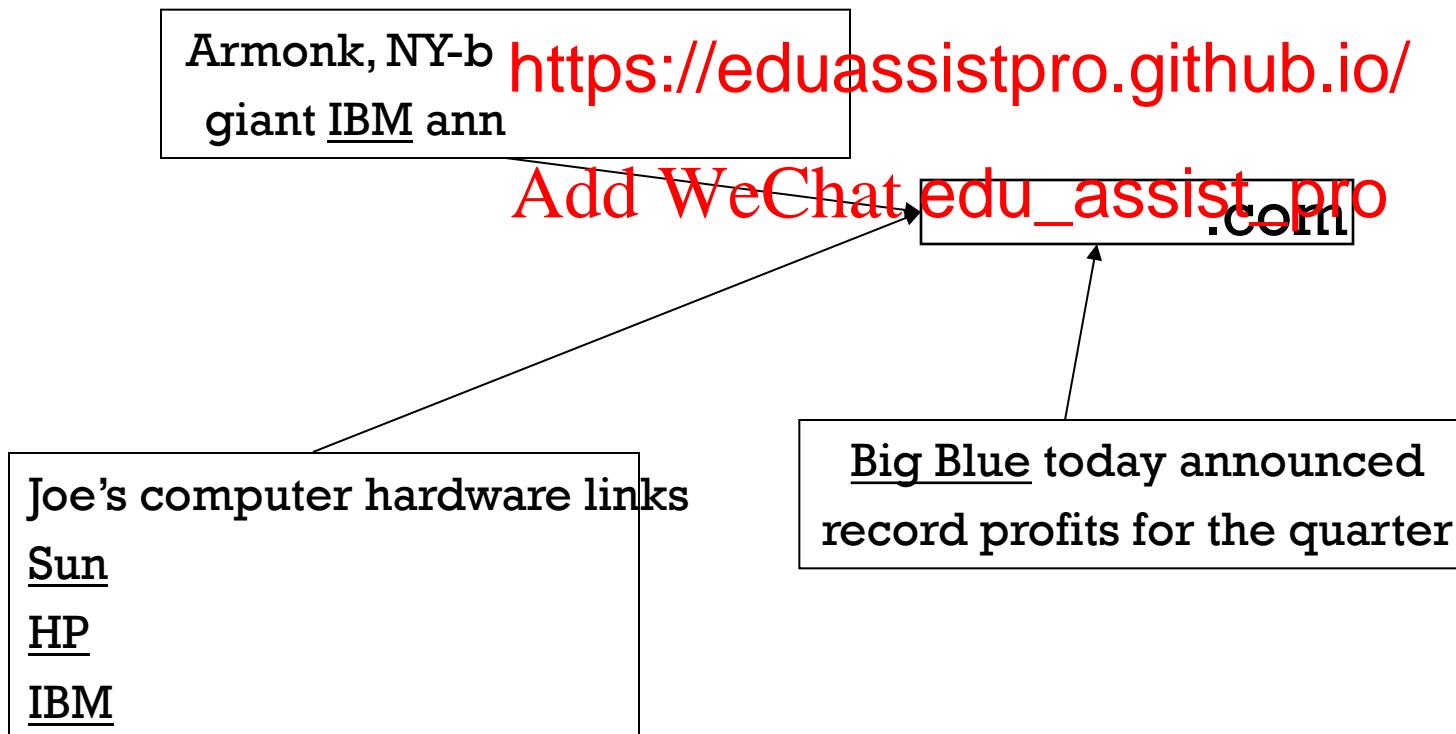


# Assignment Project Exam Help

## Indexing anchor text

- When indexing a document  $D$ , include anchor text from links pointing to  $D$ .

Assignment Project Exam Help



# Assignment Project Exam Help

## Indexing anchor te

- Can sometimes have unexpected side effects - e.g., ***evil empire***.
- Can score anchor text with weight depending on the authority of the site
  - E.g., if we were to assume that yahoo.com is authoritative, the anchor text from them

# Assignment Project Exam Help

## Anchor Text

Add WeChat edu\_assist\_pro

---

- Other applications
    - Weighting/filtering links in the graph
    - Generating scores from anchor text
- Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

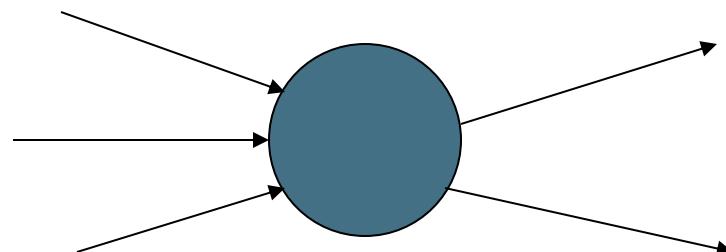
## Citation Analysis

- Citation frequency
- Co-citation coupling frequency
  - Cocitations with a given author measures “impact”
  - Cocitation <https://eduassistpro.github.io/>
- Bibliographic coupling
- Articles that co-cite the same articles are related
- Citation indexing
  - Who is this author cited by? (Garfield 1972)
- Pagerank preview: Pinski and Narin ’60s

# Assignment Project Exam Help

## Query-independent

- First generation: using link counts as simple measures of popularity.
- Two basic suggestions:
  - Undirected <https://eduassistpro.github.io/>
    - Each page gets a score = the number of out-links ( $3+2=5$ ).
  - Directed popularity:
    - Score of a page = number of its in-links (3).



# Assignment Project Exam Help

## Query processing

- First retrieve all pages meeting the text query (say ***venture capital***).
- Order these by their link popularity (either variant on the previous s <https://eduassistpro.github.io/>)
- More nuanced – use link co-easure of static goodness (Lecture 7), with text match score

# Assignment Project Exam Help

## Spamming simple Add WeChat edu\_assist\_pro

---

- *Exercise:* How do you spam each of the following heuristics so your page gets a high score?
  1. Each page gets a static score – the number of in-links plus the <https://eduassistpro.github.io/>
  2. Static score of a page =  $\frac{1}{n} \sum_{i=1}^n s_i$  in-links.

# Assignment Project Exam Help

## Ideas of PageRank

---

- Inlinks as votes
  - [www.stanford.edu](http://www.stanford.edu) has 23,400 inlinks
  - [www.joe-schmiede.com](http://www.joe-schmiede.com) has 1 inlink
- Web pages are <https://eduassistpro.github.io/>
  - [www.joe-schmiede.com](http://www.joe-schmiede.com)
  - vs. [www.stanford.edu](http://www.stanford.edu) ➡ p2
- Are all inlinks equal?
  - Recursive question!

# Assignment Project Exam Help

## Pagerank Add WeChat edu\_assist\_pro

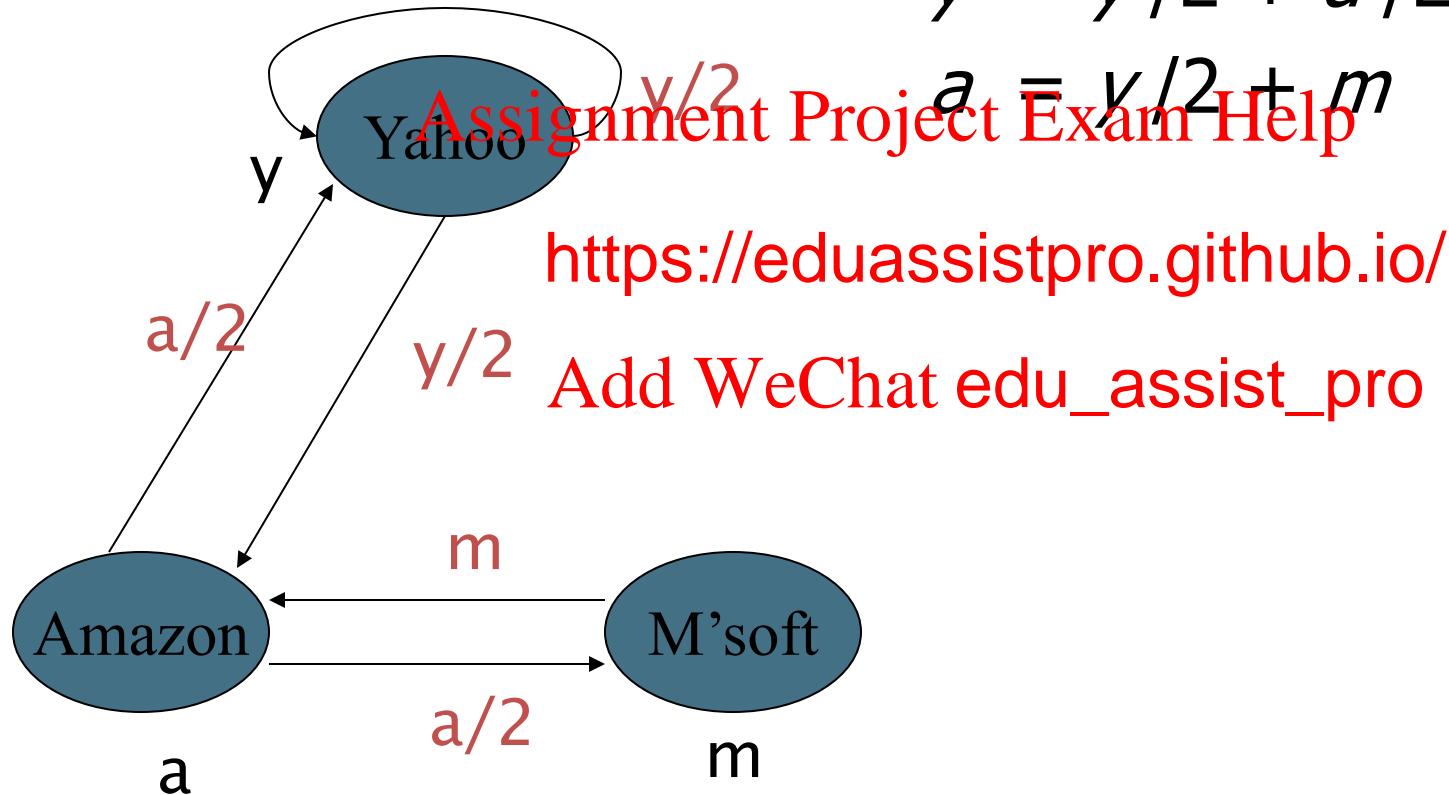
- Imagine a browser doing a *random walk* on web pages:
  - Start at a random page
  - At each step, e along one of the links on that page
- “In the steady state” each page has a long-term visit rate - use this as the page’s score.



# Assignment Project Exam Help

## Example – the Sim w' Model

The web in 1839



$$y = y/2 + a/2$$

$$a = y/2 + m$$

# Assignment Project Exam Help

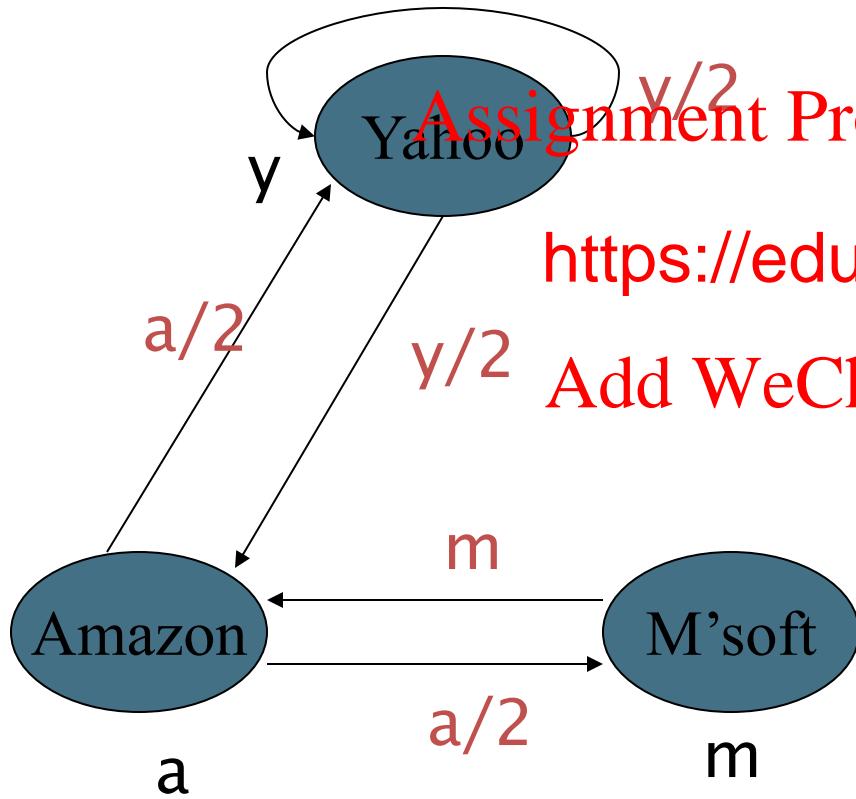
## Solving the flow eq

- 3 equations, 3 unknowns, no constants
  - No unique solution
  - All solutions equivalent modulo scale factor
- Additional co <https://eduassistpro.github.io/>
  - $y+a+m = 1$
  - $y = 2/5, a = 2/5, m = 1/5$
- Gaussian elimination method works for small examples, but we need a better method for large graphs

# Assignment Project Exam Help

## Example – the Sim~~WeChat~~<sup>edu\_assist\_pro</sup>w Model

The web in 1839



$$y_{new} = y_{old}/2 + a_{old}/2$$

$$a_{new} = y_{old}/2 + m_{old}$$

$$m = a_{old}/2$$

<https://eduassistpro.github.io/>

$y_{new}$	$= y_{old}/2 + a_{old}/2$	$1/3$	$5/12 \dots 2/5$
$a_{new}$	$= y_{old}/2 + m_{old}$	$1/2$	$1/3 \dots 2/5$
$m_{new}$	$= a_{old}/2$	$1/3$	$1/6 \dots 1/5$

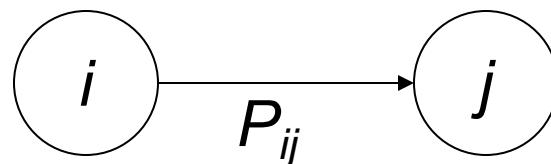
Matrix-based characterization of the computation is simpler and more useful for the general case.

# Assignment Project Exam Help

## Markov chains

- A Markov chain consists of  $n$  states, plus an  $n \times n$  transition probability matrix  $\mathbf{P}$ .
- At each step, we are in exactly one of the states.
- For  $1 \leq i, j \leq n$ , <https://eduassistpro.github.io/> is the probability of  $j$  being the next state given we are currently in state  $i$ .

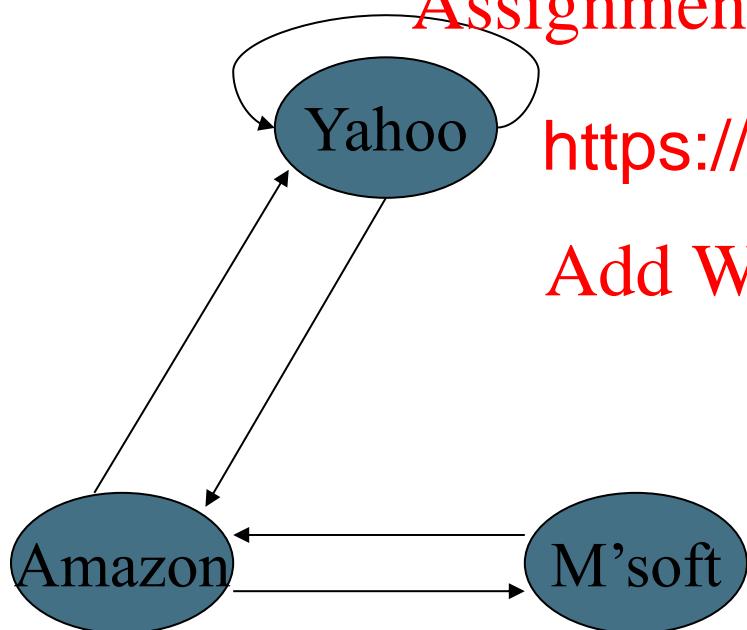
$P_{ii} > 0$   
is OK.



# Assignment Project Exam Help

## Markov chains

- Clearly, for all i,  $\sum_{j=1}^n P_{ij} = 1.$
- Markov chains are abstractions of random walks.



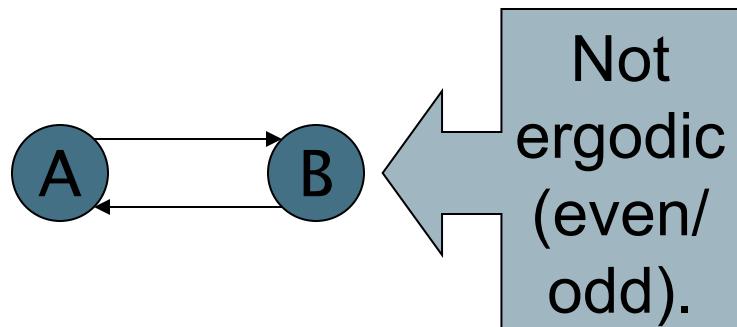
https://eduassistpro.github.io/

a	m
0	1/2
m	0 1 0

# Assignment Project Exam Help

## Ergodic Markov ch

- A Markov chain is ergodic if
  - you have a path from any state to any other
  - For any state  $s_0$ , there exists a transient time  $T_0$ , such that at a fixed time  $T > T_0$ , the probability of being at state  $s$  is nonzero.



# Assignment Project Exam Help

Ergodic Markov chain  
Add WeChat edu\_assist\_pro

---

- For any ergodic Markov chain, there is a unique long-term visit rate for each state.
  - *Steady-state distribution.* <https://eduassistpro.github.io/>
- Over a long time period, visit each state in proportion to this rate.
- It doesn't matter where we start.

# Assignment Project Exam Help

## Probability vectors

- A probability (row) vector  $\mathbf{x} = (x_1, \dots x_n)$  tells us where the walk is at any point.
- E.g.,  $(000\dots 1\dots 000)$  means we're in state  $i$ .

$1 \quad i \quad \text{https://eduassistpro.github.io/}$

Add WeChat edu\_assist\_pro

More generally, the vector  $\mathbf{x} = (x_1, \dots x_n)$  means the walk is in state  $i$  with probability  $x_i$ .

$$\sum_{i=1}^n x_i = 1.$$

# Assignment Project Exam Help

## Change in probability

- If the probability vector is  $\mathbf{x} = (x_1, \dots, x_n)$  at this step, what is it at the next step?
- Recall that  $\mathbf{P}$  is the transition probability matrix  $\mathbf{P}$  from state  $i$ . So from  $\mathbf{x}$ , our next state is distributed as  $\mathbf{x}\mathbf{P}$ .

# Assignment Project Exam Help

## How do we compute the steady-state vector?

- Let  $\mathbf{a} = (a_1, \dots, a_n)$  denote the row vector of steady-state probabilities.
- If our current position is described by  $\mathbf{a}$ , then the next step is directed by <https://eduassistpro.github.io/>
- But  $\mathbf{a}$  is the steady state, so
- Solving this matrix equation.
  - So  $\mathbf{a}$  is the (left) eigenvector for  $\mathbf{P}$ .
  - (Corresponds to the “principal” eigenvector of  $\mathbf{P}$  with the largest eigenvalue.)
  - Transition probability matrices always have largest eigenvalue 1.

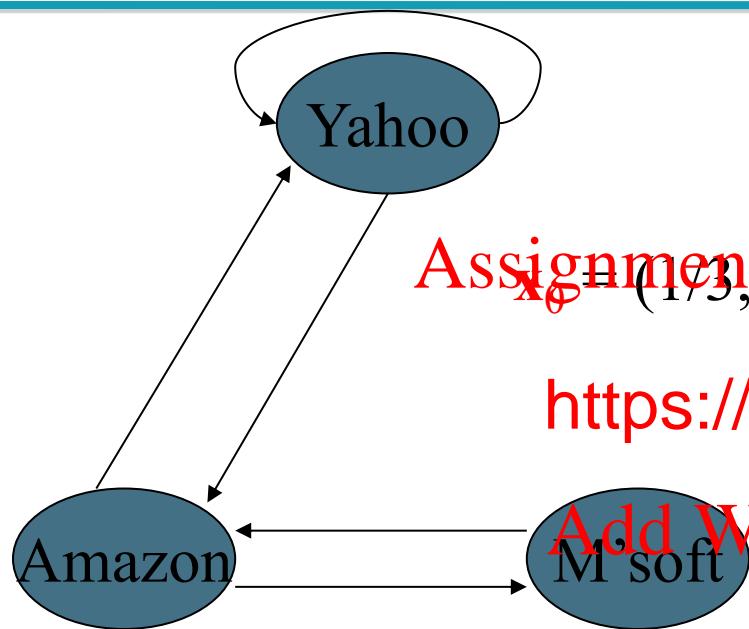
# Assignment Project Exam Help

## One way of compu

- Recall, regardless of where we start, we eventually reach the steady state  $\mathbf{a}$ .
- Start with any distribution (say  $\mathbf{x} = (1/n, 1/n, \dots, 1/n)$ ).
- After one step <https://eduassistpro.github.io/>
- after two steps at  $\mathbf{x} \mathbf{P}^2$ , then  $\mathbf{x} \mathbf{P}^3$  on
- “Eventually” means for “large  $n$ ”  $\mathbf{a}$ .
- Algorithm: multiply  $\mathbf{x}$  by increasing powers of  $\mathbf{P}$  until the product looks stable.

# Assignment Project Exam Help

## Power Iteration Ex



$$\mathbf{P}$$

	y	a	m
y	1/2	1/2	0
a	1/2	0	1/2
m	1	0	

Assignment Project Exam Help  
 $\mathbf{x}_0 = (1/3, 1/3, 1/3)$   
<https://eduassistpro.github.io/>

$$\begin{aligned}
 a_{new} &= \frac{a_{old}}{2} + \frac{m_{old}}{2} \\
 m_{new} &= \frac{a_{old}}{2}
 \end{aligned}$$

	$\mathbf{x}_0$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_t$
y	1/3	1/3	5/12	3/8	2/5
a	=	1/3	1/2	1/3	11/24
m		1/3	1/6	1/4	1/6

# Assignment Project Exam Help

## Spider traps

Add WeChat edu\_assist\_pro

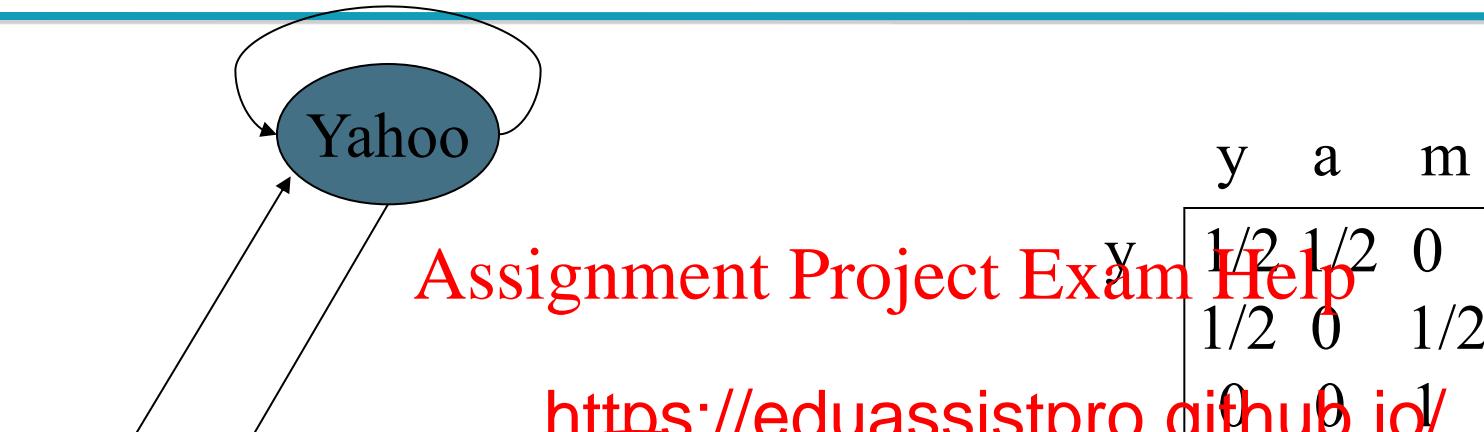
---

- A group of pages is a **spider trap** if there are no links from within the group to outside the group
  - Random surfer gets trapped
- Spider traps v <https://eduassistpro.github.io/> needed for the random walk t

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

Microsoft ~~Add WeChat edu\_assist\_pro~~ trap



y a m

y	$1/2$	$1/2$	0
a	$1/2$	0	$1/2$
m	0	0	1

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

y	$1/3$	$1/3$	$1/4$	$5/24$	...	0
a	$1/3$	$1/6$	$1/6$	$1/8$	...	0
m	$1/3$	$1/2$	$7/12$	$2/3$		1

# Assignment Project Exam Help

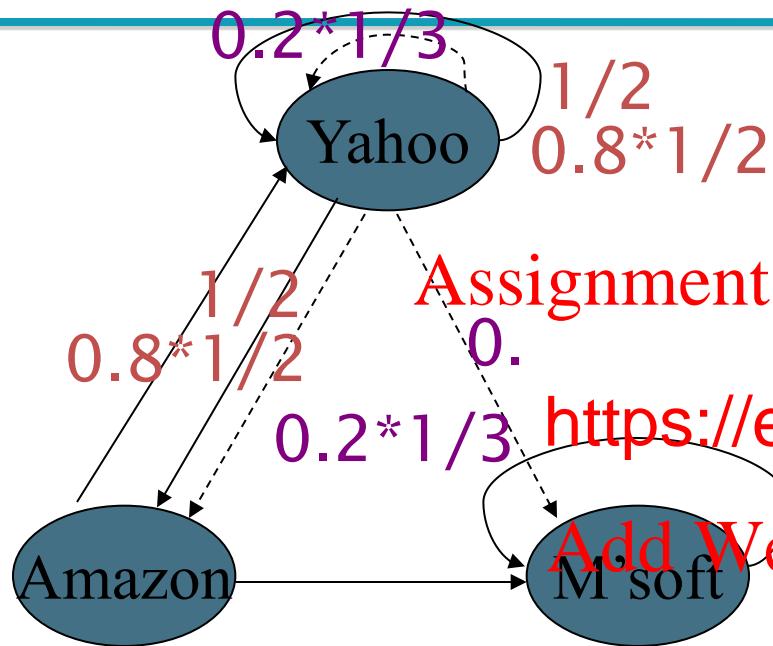
## Random teleports

---

- The Google solution for spider traps
- At each time step, the random surfer has two options: Assignment Project Exam Help
  - With probability  $\beta$ , jump to <https://eduassistpro.github.io/>
  - With probability  $1-\beta$ , jump to a random page uniformly at random
- Common values for  $\beta$  are in the range 0.8 to 0.9
- Surfer will teleport out of spider trap within a few time steps

# Assignment Project Exam Help

Random teleports ~~Add WeChat edu\_assist\_pro~~



~~Assignment Project Exam Help~~

~~https://eduassistpro.github.io/~~

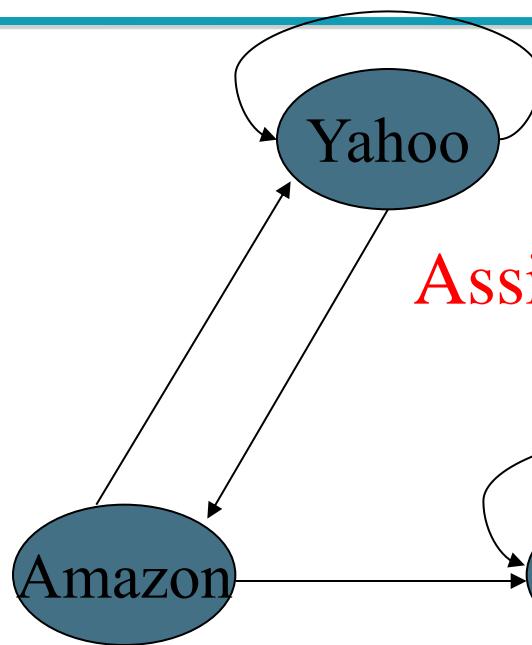
0	0	1	0.2
1/3	1/3	1/3	
1/3	1/3	1/3	

1/3	1/3	1/3
1/3	1/3	1/3
1/3	1/3	1/3

y	7/15	7/15	1/15
a	7/15	1/15	7/15
m	1/15	1/15	13/15

# Assignment Project Exam Help

Random teleports ~~Add WeChat edu\_assist\_pro~~)



0.8

1/2	1/2	0
1/2	0	0
0	1/2	1

+ 0.2

1/3	1/3	1/3
1/3	1/3	1/3
1/3	1/3	1/3

Assignment Project Exam Help

<https://eduassistpro.github.io/>

5 7/15 1/15  
5

1/15 13/15

y            0.33  0.33  0.28  0.26            7/33

a   =       0.33  0.20  0.20  0.18   ...      5/33

m            0.33  0.47  0.52  0.56            7/11

# Assignment Project Exam Help

## Dead ends

Add WeChat edu\_assist\_pro

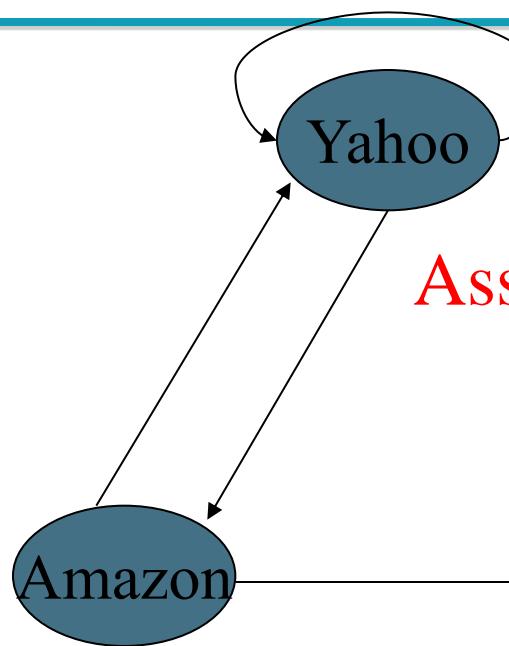
---

- Pages with no outlinks are “dead ends” for the random surfer
  - Nowhere to go on next step
- Especially co <https://eduassistpro.github.io/>
  - URLs that have not yet been

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

Microsoft ~~Add WeChat~~ become ~~edu\_assist\_pro~~ end



0.8

1/2	1/2	0
1/2	0	1/2
0	0	0

+ 0.2

1/3	1/3	1/3
1/3	1/3	1/3
1/3	1/3	1/3

Assignment Project Exam Help

<https://eduassistpro.github.io/>

5	7/15	1/15
5	7/15	1/15

Non-stochastic!

y	0.33	0.33	0.262	0.216	...	0
a =	0.33	0.20	0.182	0.143	...	0
m	0.33	0.20	0.129	0.111		0

# Assignment Project Exam Help

## Dealing with dead ends

- Teleport
  - Follow random teleport links **with probability 1.0** from dead-ends
  - Adjust matrix
    - How? <https://eduassistpro.github.io/>

- (Suggested by **Google**)
  - Preprocess the graph to eliminate dead-ends
  - Might require multiple passes
  - Compute page rank on reduced graph
  - **Approximate** values for deadends by propagating values from reduced graph

Q: Why approximate values and why errors are insignificant?

# Assignment Project Exam Help

## Pagerank Add WeChat edu\_assist\_pro

---

- Preprocessing:
  - Given graph of links, build matrix  $\tilde{P}$ .
  - From it compute  $a$ .
    - $a$  is the principal eigenvector of  $\tilde{P}$ :  
$$\tilde{P} = \frac{1}{n}$$
<https://eduassistpro.github.io/>
    - The entry  $a_i$  is a number between 0 and 1: the pagerank of page  $i$ .
- Query processing:
  - Retrieve pages meeting query.
  - Rank them by their pagerank.
  - Order is query-*independent*.

# Assignment Project Exam Help

## The reality

- Pagerank is used in google, but is hardly the full story of ranking
  - Many sophisticated features are used
  - Some address <https://eduassistpro.github.io/>
  - Machine learning heavily used
- Pagerank still very useful for the crawl policy

# Assignment Project Exam Help

## Pagerank

- How realistic is the random surfer model?
  - (Does it matter?)
  - What if we modeled the back button?
  - Surfer behavior makes short paths
  - Search engines make jumps non-random.
- Biased Surfer Models
  - Weight edge traversal probabilities based on match with topic/query (non-uniform edge selection)
  - Bias jumps to pages on topic (e.g., based on personal bookmarks & categories of interest)

# Assignment Project Exam Help

## Topic Specific Pager

- Goal – pagerank values that depend on query *topic*
- Conceptually, we use a random surfer who teleports, with say 10% probability, using the following rule: <https://eduassistpro.github.io/>
  - Selects a topic (say, one category) based on a  $q$ -specific distribution over the categories
  - Teleport to a page uniformly at random within the chosen topic
- Sounds hard to implement: can't compute PageRank at query time!

only  
randomly  
teleport  
to a  
subset of  
pages

# Assignment Project Exam Help

## Topic Specific Pagerank

Add WeChat edu\_assist\_pro

- **Offline:** Compute pagerank for *individual* topics
  - Query independent as before
  - Each page has multiple pagerank scores - one for each ODP category, with weight at category
- **Online:** Query <https://eduassistpro.github.io/> (distribution of weights over) topics
  - Generate a dynamic pagerank score
$$\text{e} = \text{weighted sum of topic-specific pageranks}$$

Add WeChat edu\_assist\_pro

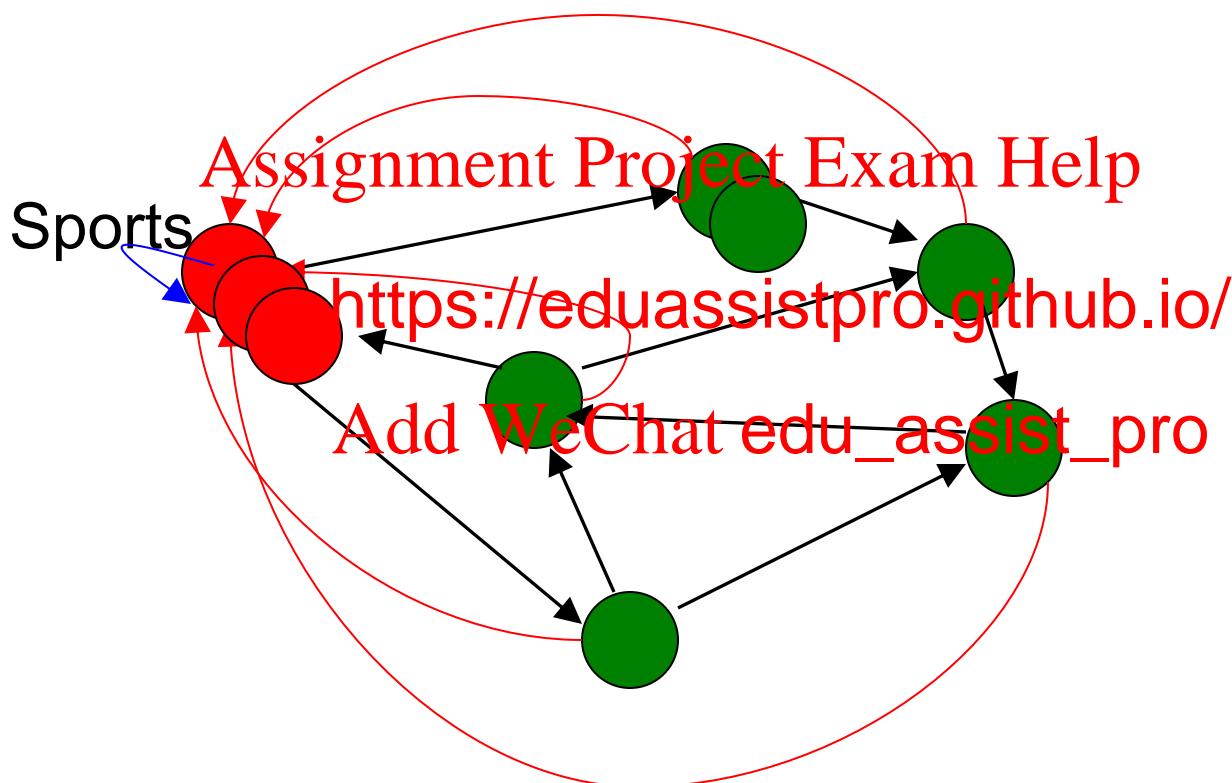
# Assignment Project Exam Help

## Influencing PageRank

- Input:
  - Web graph  $W$
  - Influence vector  $v$ : (page  $\rightarrow$  topic)
- Output:
  - Rank vector  $r$ : (page  $\rightarrow$  page importance wrt  $v$ )
  - $r = PR(W, v)$

# Assignment Project Exam Help

## Non-uniform Teleportation



Teleport with 10% probability to a Sports page

# Assignment Project Exam Help

## Interpretation of Co<sub>ij</sub> Add WeChat edu\_assist\_pro Score

---

- Given a set of personalization vectors  $\{v_j\}$

$$PR(W, \sum_j [w_j \cdot v_j]) = \sum_j [w_j \cdot PR(W, v_j)]$$

Assignment Project Exam Help

Given a user's p <https://eduassistpro.github.io/>, express as a combination of the "basis" v<sub>j</sub>

# Assignment Project Exam Help

## PageRank as a Line

[Jeh & Widom, KDD 2003] Add WeChat edu\_assist\_pro [Optional]

- Preference vectors  $\mathbf{u} \rightarrow$  specifies the random teleport probability distribution
  - A column vector of  $n$  dimensions ( $n$ : # of vertices in  $G$ )
  - sum up to 1. <https://eduassistpro.github.io/>
- Personalized distribution over the vertices
  - Also a column vector of  $n$  dimensions
- $v$  is determined by  $u$  via a linear system

$$\mathbf{v} = (1 - \beta)\mathbf{P}^T \mathbf{v} + \beta \mathbf{u}$$


# Linearity of Pagerank

- For two preference vectors  $\mathbf{u}_1$  and  $\mathbf{u}_2$ , let  $\mathbf{v}_1$  and  $\mathbf{v}_2$  be the corresponding personalized page rank vectors, then we can prove that

$$\lambda \mathbf{v}_1 + (1 - \lambda) \mathbf{v}_2 = (\text{https://eduassistpro.github.io}) + \beta(\lambda \mathbf{u}_1 + (1 - \lambda) \mathbf{u}_2)$$

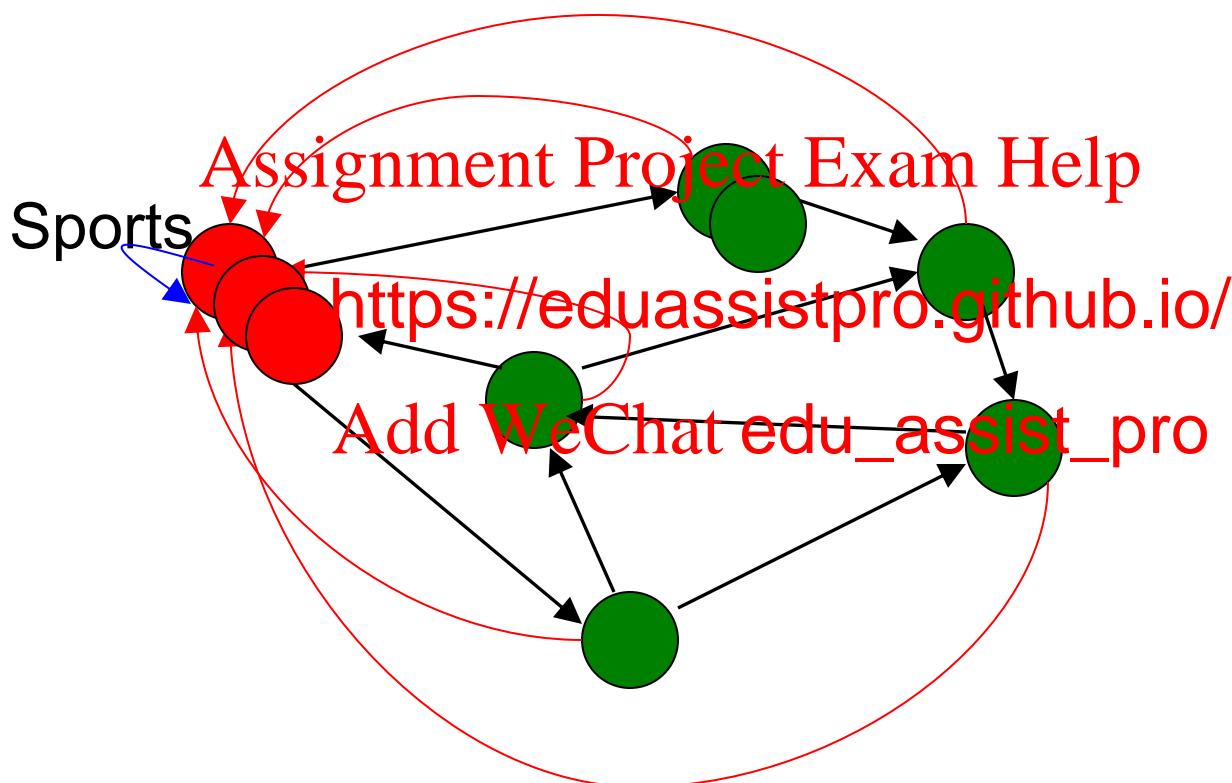
Add WeChat edu\_assist\_pro

- Implication:
  - Personalized pagerank vectors induced by a linear combination of preference vectors can be computed as the same linear combination of corresponding personalized pagerank vectors.

# Assignment Project Exam Help

## Interpretation

Add WeChat edu\_assist\_pro

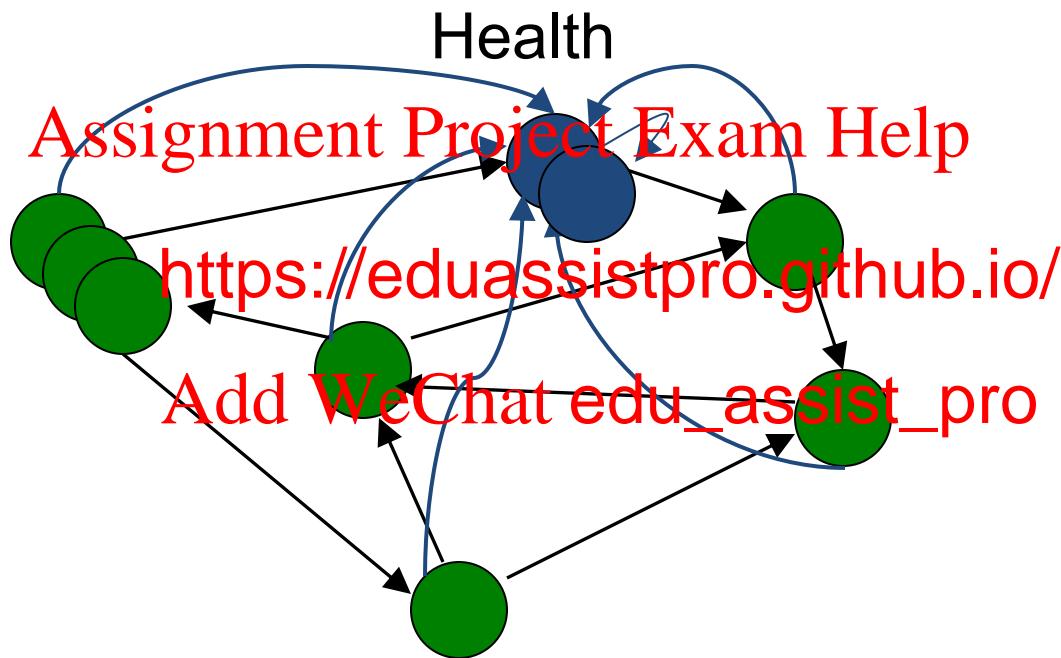


10% Sports teleportation

# Assignment Project Exam Help

## Interpretation

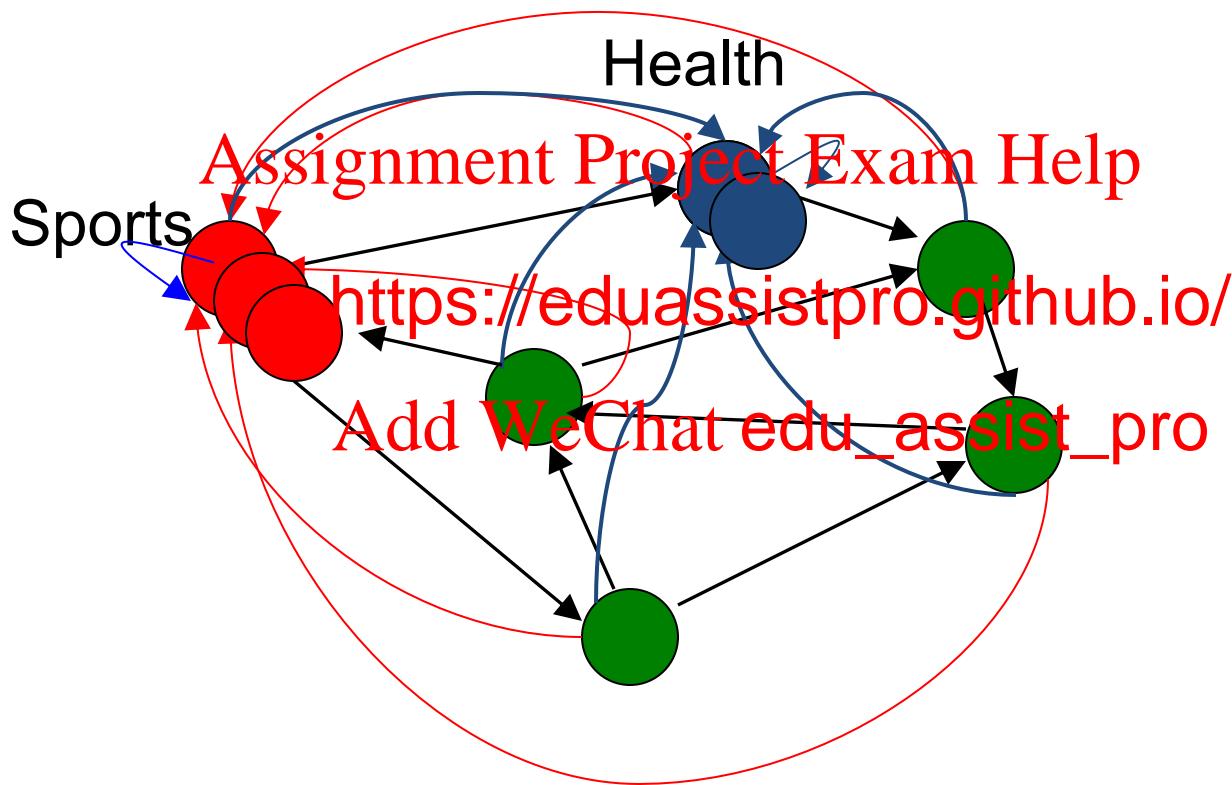
Add WeChat edu\_assist\_pro



10% Health teleportation

# Assignment Project Exam Help

## Interpretation



$pr = (0.9 PR_{\text{sports}} + 0.1 PR_{\text{health}})$  gives you:  
9% sports teleportation, 1% health teleportation

# Assignment Project Exam Help

## Resources

- IIR Chap 21
- <http://www2004.org/proceedings/docs/1p309.pdf>
- <http://www2.cs.berkeley.edu/pubs/2004/1p595.pdf>
- <http://www2.cs.berkeley.edu/pubs/2004/1p595.pdf>
- <https://eduassistpro.github.io/refered/p270/kamvar-270-xhtml/index.html>
- <http://www2003.org/cdrom/refered/p641/xhtml/p641-mccurley.html>
- Glen Jeh and Jennifer Widom: Scaling Personalized Web Search. KDD 2003.