

Carnegie Mellon University

95-885 Data Science and Big Data

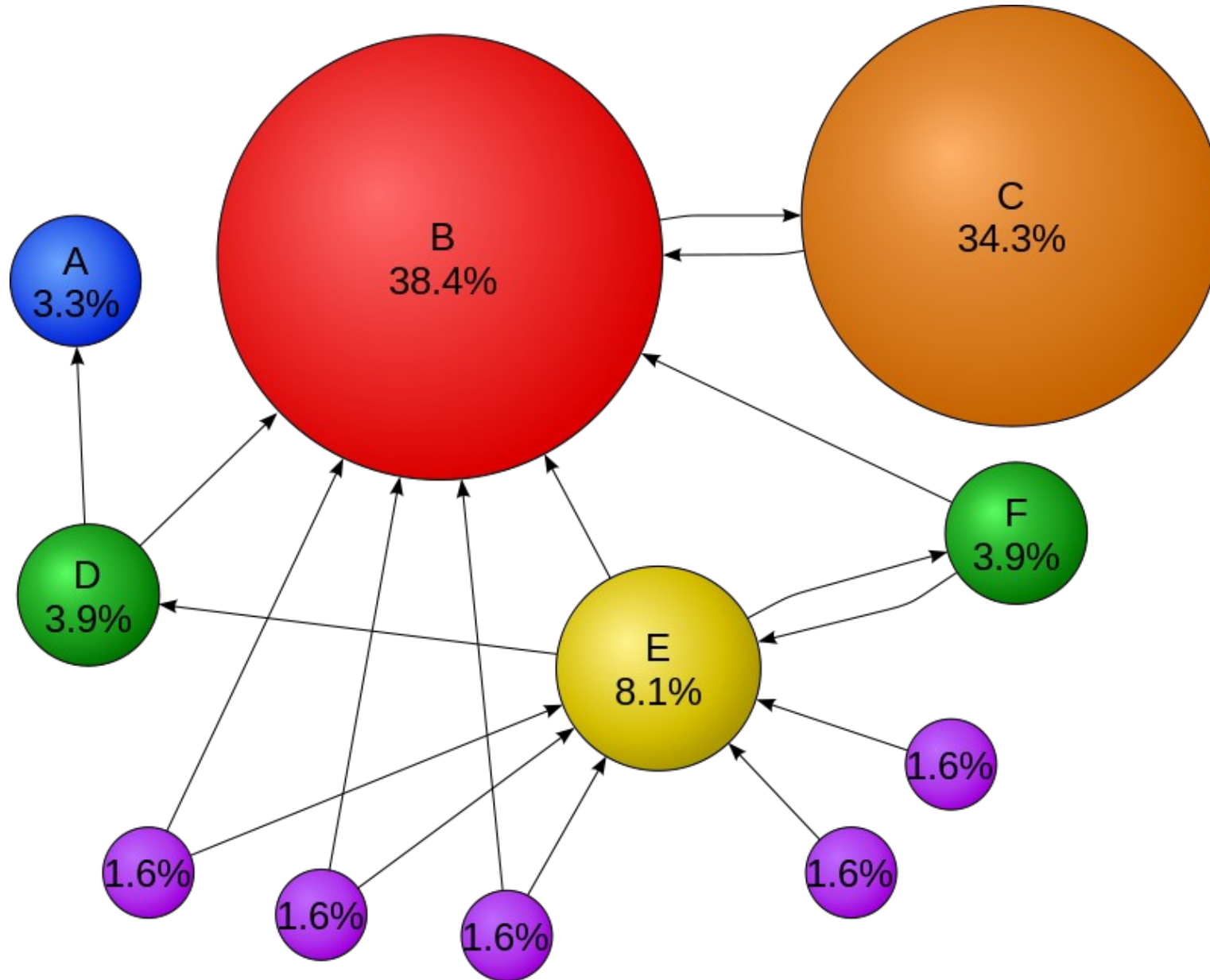
PageRank with MapReduce

Agenda

- The essence of the algorithm
- Simulating with a spreadsheet
- Expressing with a mapper and reducer

THE INTUITION BEHIND THE PAGERANK ALGORITHM

A Sample Set of Interconnected Web Pages



Pagerank (page)
==
*Probability of navigating
to that page*

*If you sum all the
probabilities on the adjacent
figure you will get 100*
 $3.3 + 38.4 + \dots + 1.6$

The Intuition behind the PageRank algorithm

- $A \rightarrow B$ means A thinks B is worth something
- A page is important if it is pointed to by other important pages
 - e.g., a link to a page from Wikipedia etc
 - the “slashdot effect”
- Note that this measure of importance is calculated *independent* of the *content* of the page
- “wisdom of the crowds”

Quantifying the Intuition

- We get to a page in one of two ways

1. By randomly jumping to it
2. By navigating (clicking) to it from the page we currently are on

- Hence the probability of getting to a page x is:

- $P(x) = P(\text{randomly jumping to the page}) + P(\text{navigating to it by clicking a link})$

- Suppose the probability of deciding to navigate from a page i.e., click a link on the page is β , then the probability of not clicking a link i.e., randomly jumping to a page is $1-\beta$

- If there are a total of N pages then,

- $P(\text{randomly jumping to a page}) = \frac{1-\beta}{N}$

Probability of Navigating to a Page

- Suppose page y_1 links to page x
- Then the probability of navigating to x from y_1 is:

- $PR(y_1) * 1/out(y_1)$

- where $PR(y_1)$ is the probability of getting to y_1

and $out(y_1)$ is the number of outlinks from y_1

- We may have several 'y's linking to x so

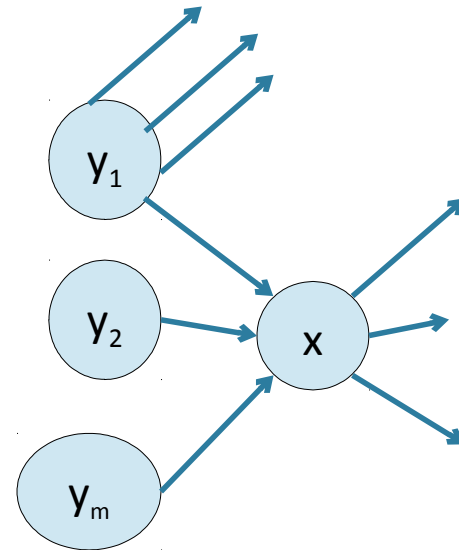
$PR(x)$ is

- $PR(y_1)/out(y_1) + \dots PR(y_m)/out(y_m) +$

- The probability of deciding to navigate in the first place is

any of the y s (pages that link to x) is:

$$\beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$$



Combining the two pieces together we have

$$PR(x) = \frac{1 - \beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$$

Probability of clicking on a link

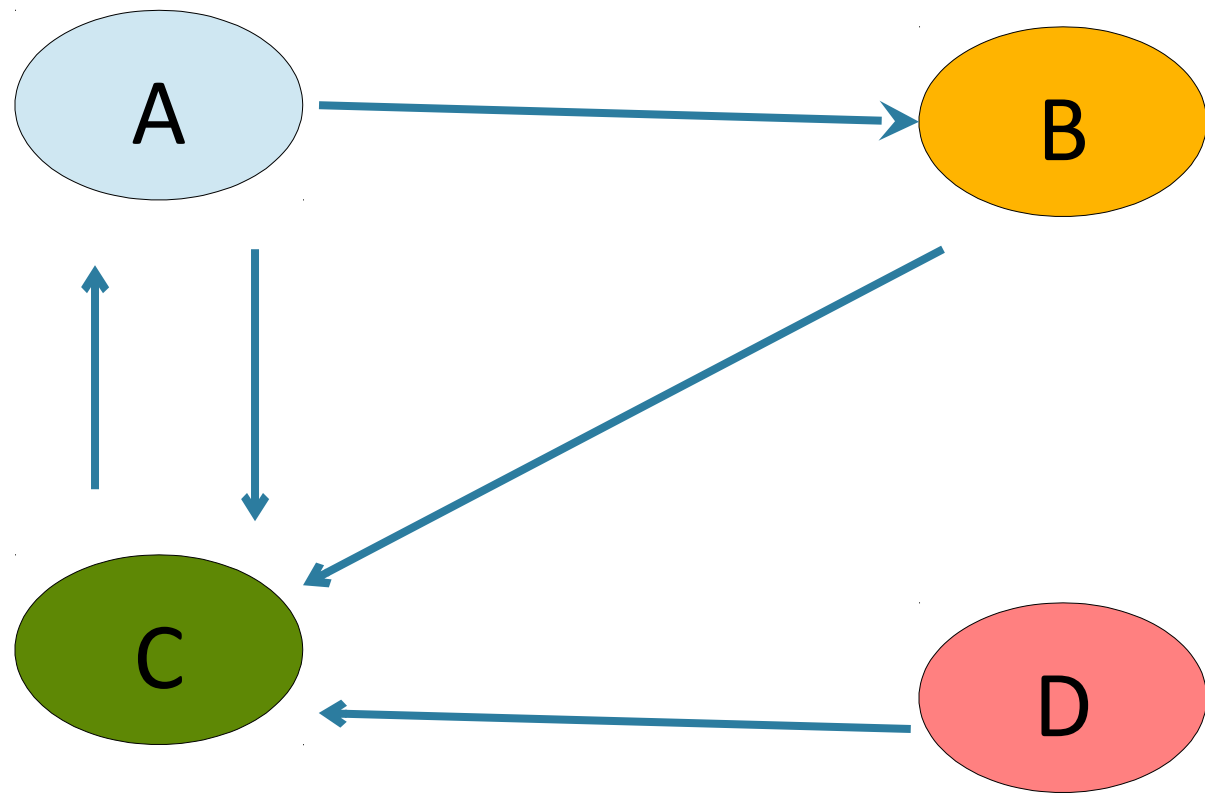
Total number of pages

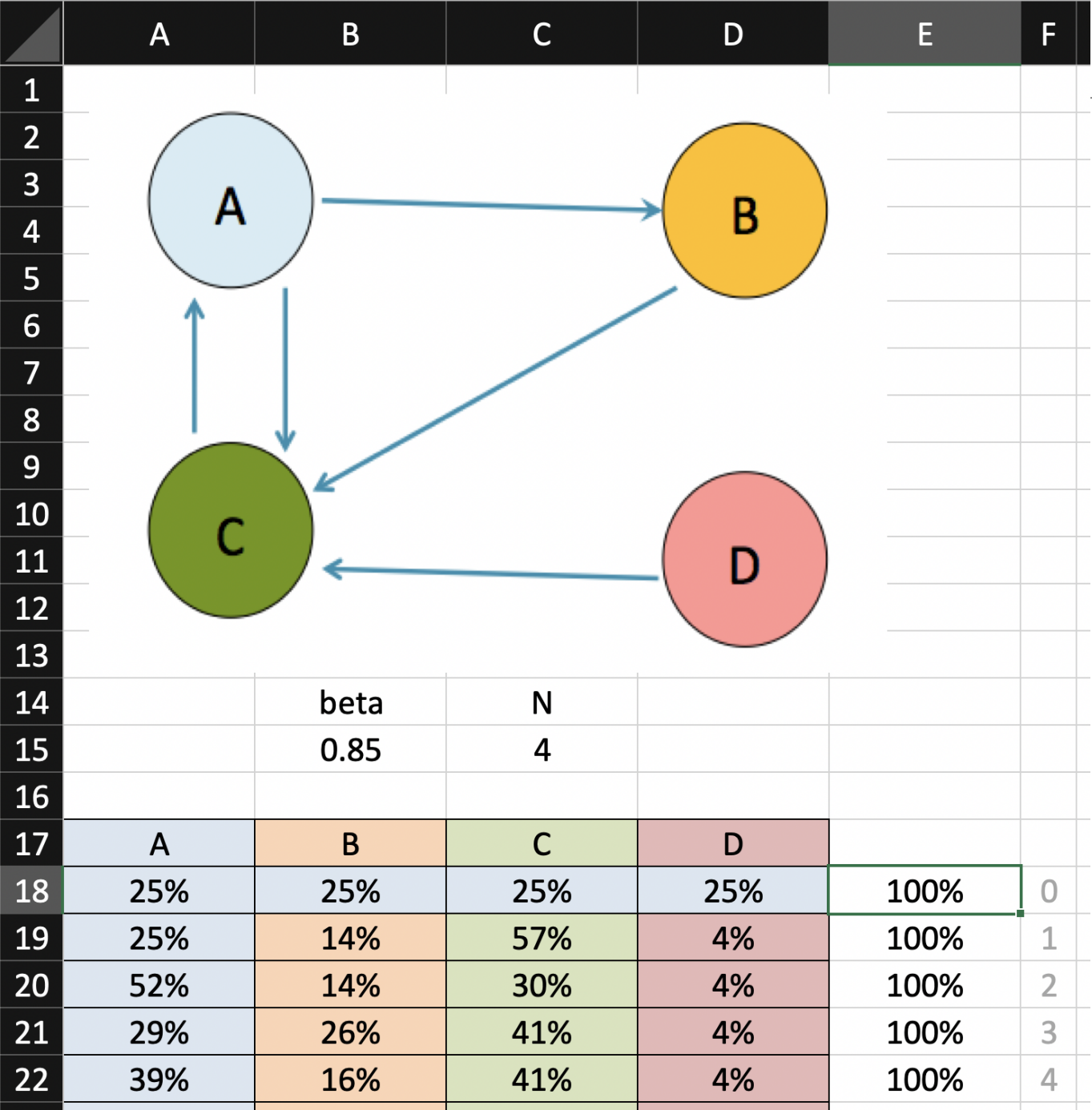
$P(x) = P(\text{randomly jumping to the page}) + P(\text{navigating to it by clicking})$

PAGERANK ON A SPREADSHEET

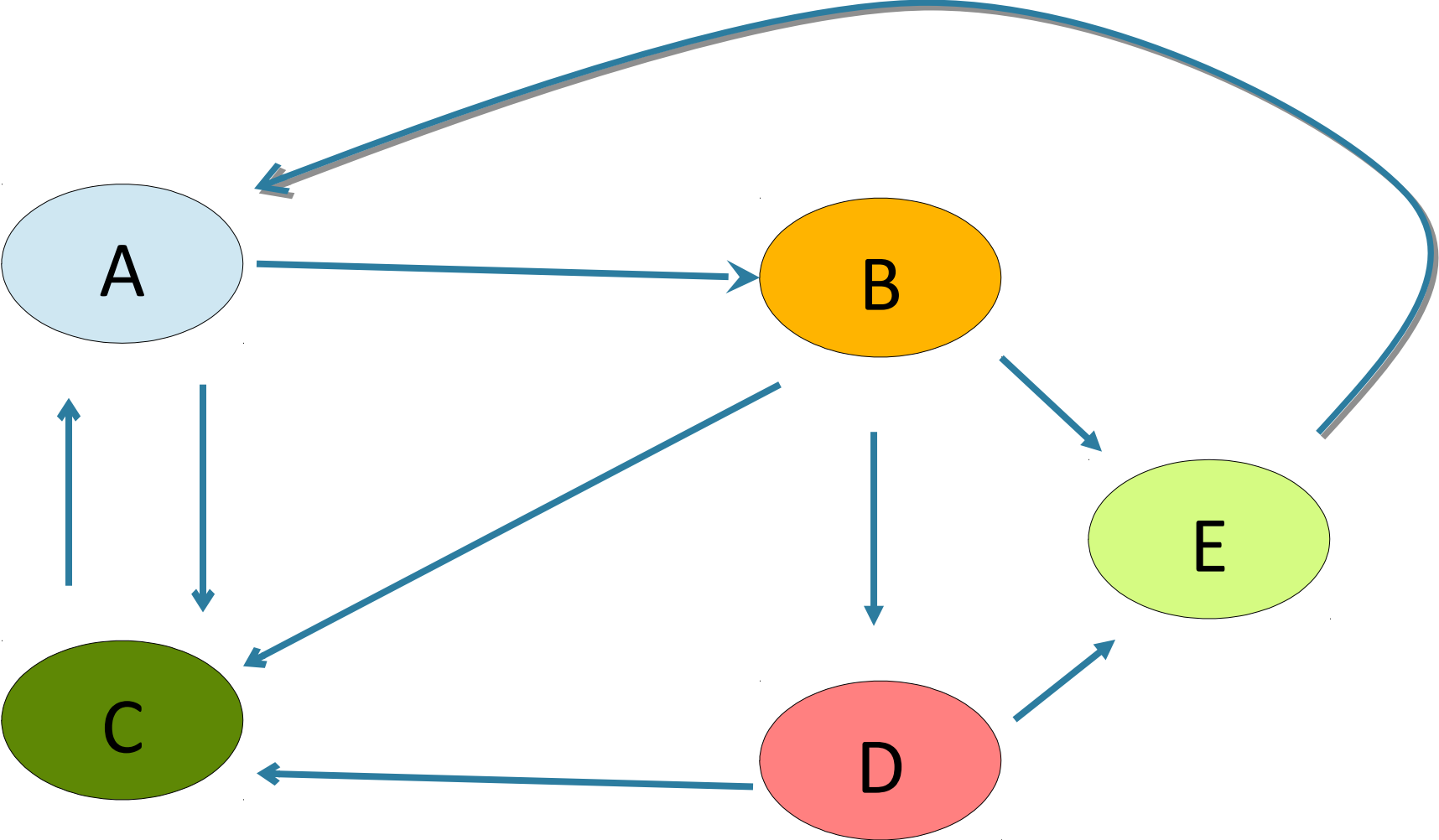
Spreadsheet available on Canvas schedule

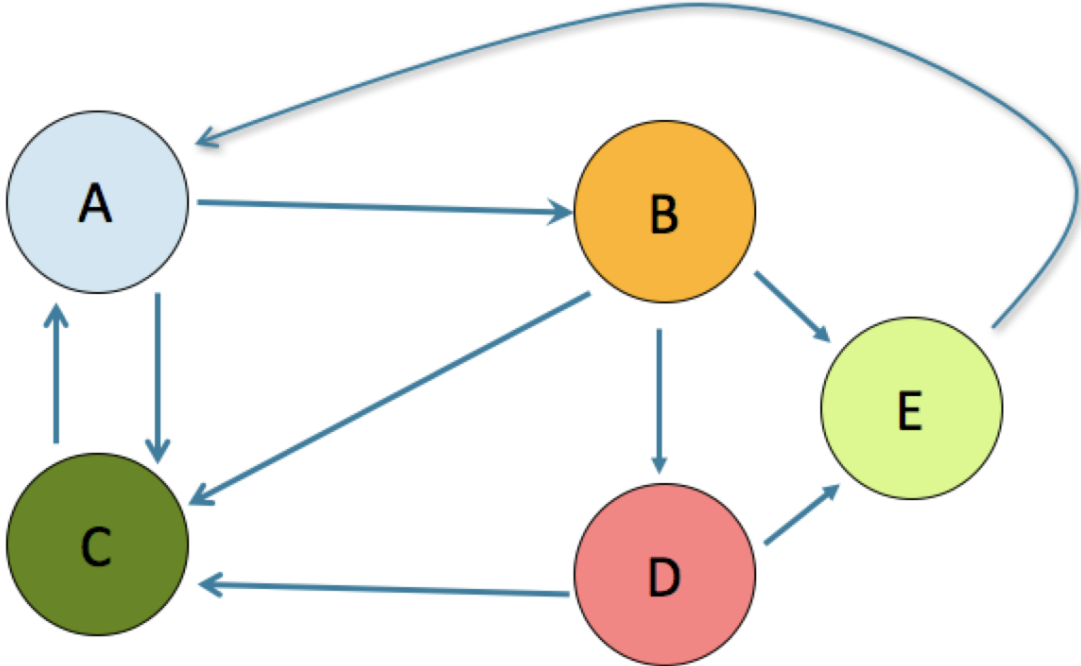
Spreadsheet: Working through the PageRank Algorithm





Exercise



	A	B	C	D	E	F	G
1							
2							
3							
4							
5							
6							
7							
8							
9							
10							
11							
12							
13							
14		beta	N				
15		0.85	5				
16							
17	A	B	C	D	E		
18	20%	20%	20%	20%	20%	100%	0
19	37%	12%	26%	9%	17%	100%	1
20	39%	19%	26%	6%	10%	100%	2
21	33%	20%	28%	8%	11%	100%	3
22	36%	17%	26%	9%	12%	100%	4

PAGERANK: MAPPER AND REDUCER

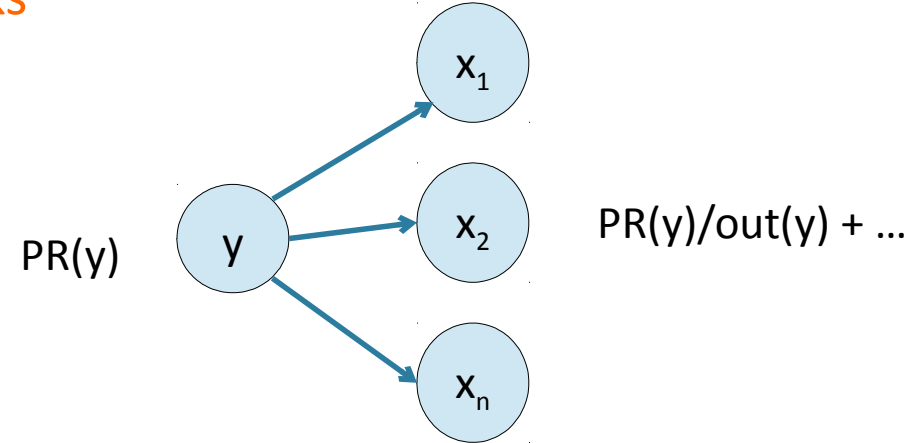
Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

—for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

—**yield** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

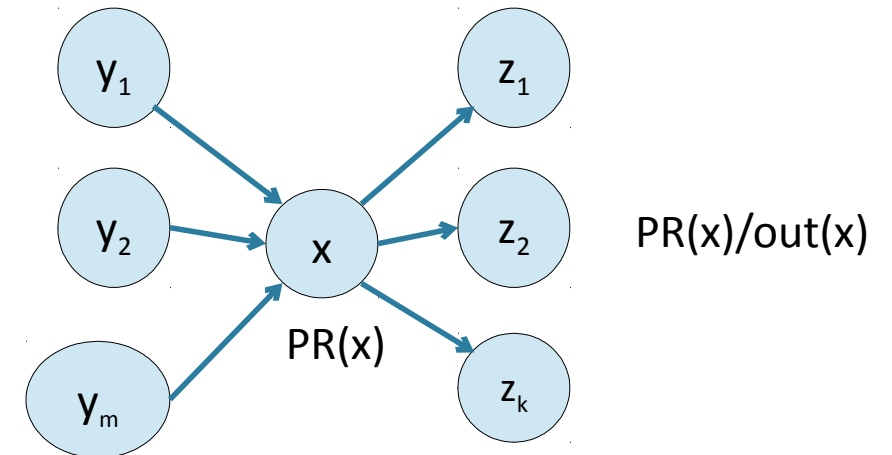
node, out-links



■ **Reducer** $\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\}, \{z_1, \dots, z_k\} \rangle$ node, ΔPR from in-links

—**compute** $PR(x) = \frac{1 - \beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

—for $j = 1..k$, **yield** $\langle z_j, \frac{PR(x)}{out(x)} \rangle$
 $\langle x, \{z_1, z_2, \dots, z_k\} \rangle$



$$PR(x) = [\dots] + PR(y_1)/out(y_1) + \dots$$

Expressing PageRank with MapReduce

- Mapper

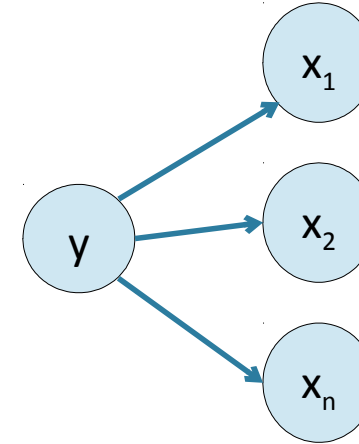
- Reducer

Expressing PageRank with MapReduce

■ Mapper

node, out-links

■ Reducer

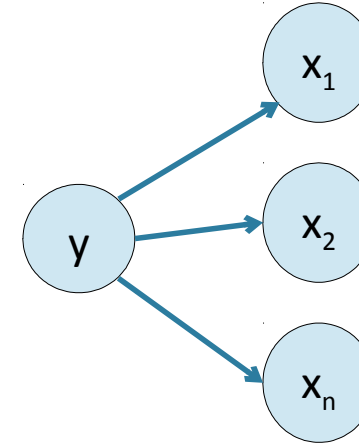


Expressing PageRank with MapReduce

■ Mapper $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

node, out-links

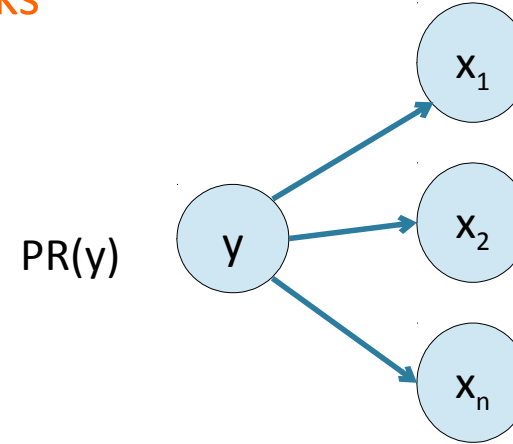
■ Reducer



Expressing PageRank with MapReduce

■ Mapper $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

node, out-links

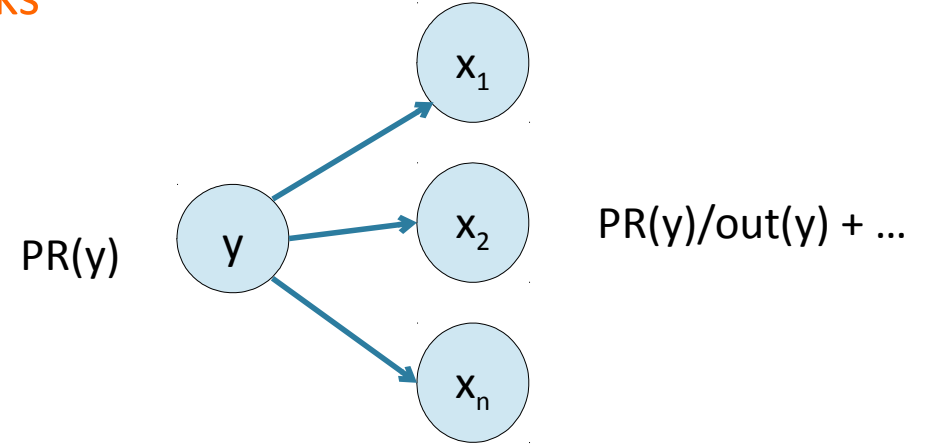


■ Reducer

Expressing PageRank with MapReduce

■ Mapper $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

node, out-links



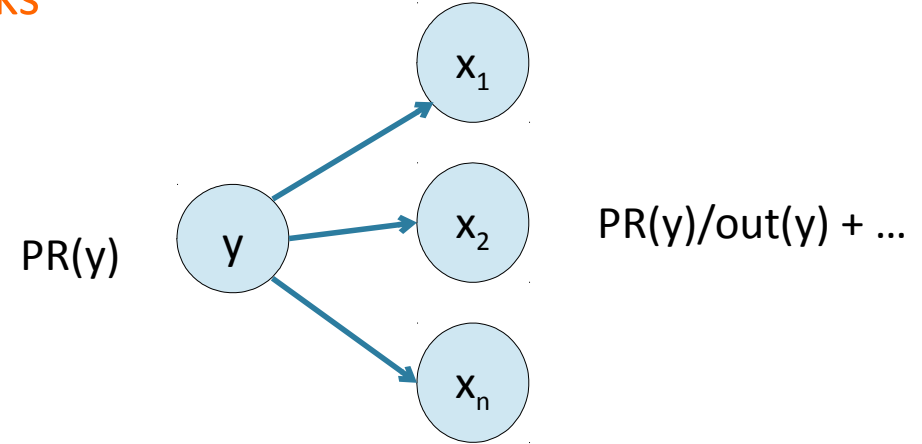
■ Reducer

Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

—for $i = 1..n$, **yield** $\left\langle x_i, \frac{PR(y)}{out(y)} \right\rangle$

node, out-links



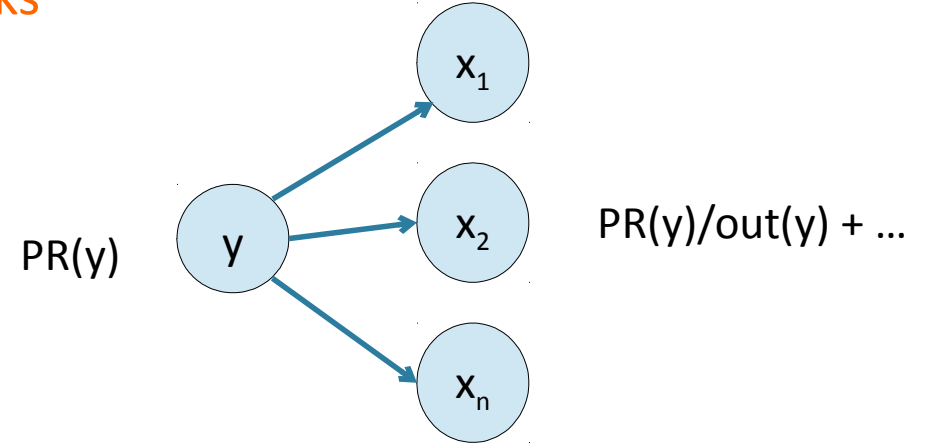
■ **Reducer**

Expressing PageRank with MapReduce

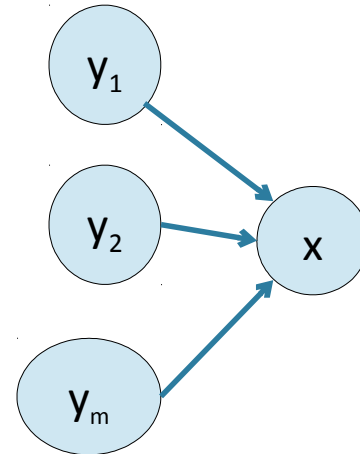
■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

—for $i = 1..n$, **yield** $\left\langle x_i, \frac{PR(y)}{out(y)} \right\rangle$

node, out-links



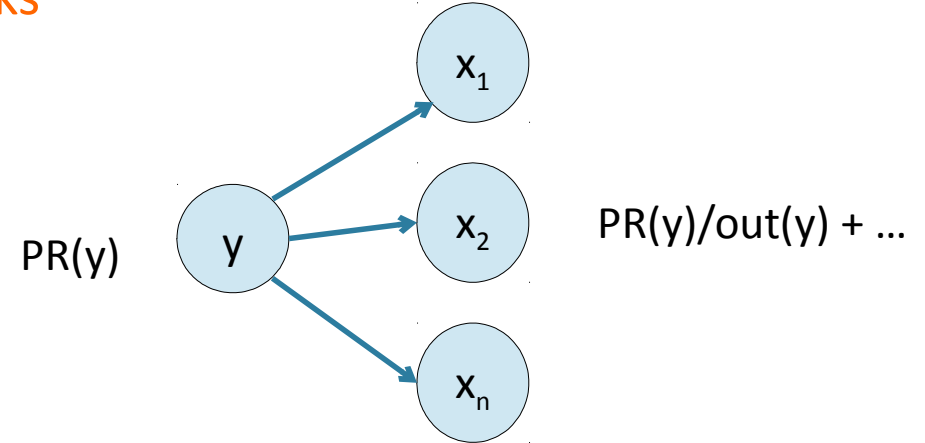
■ **Reducer**



Expressing PageRank with MapReduce

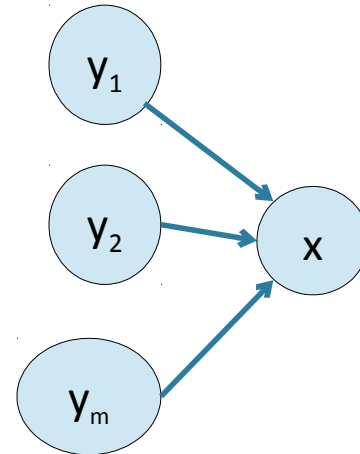
■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$
—for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

node, out-links



■ **Reducer**

node, ΔPR from in-links

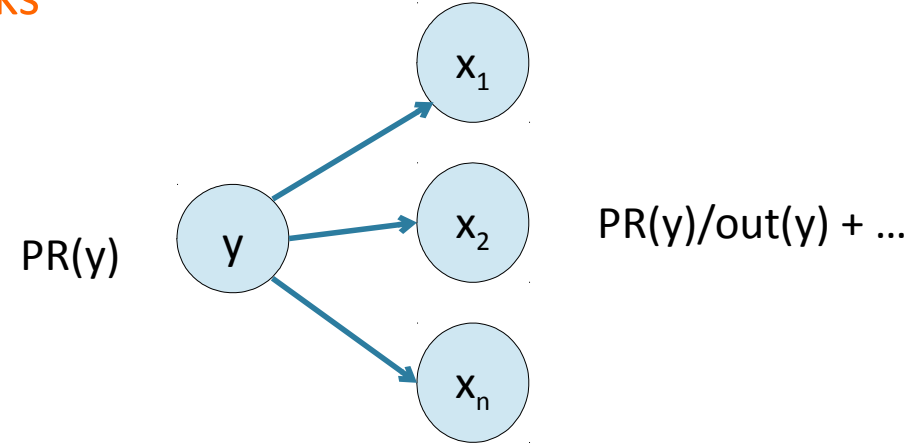


Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

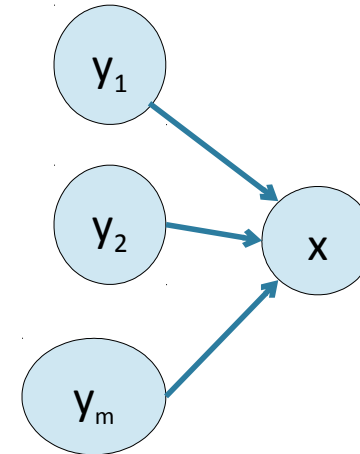
—for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

node, out-links



■ **Reducer** $\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\} \rangle$

node, ΔPR from in-links

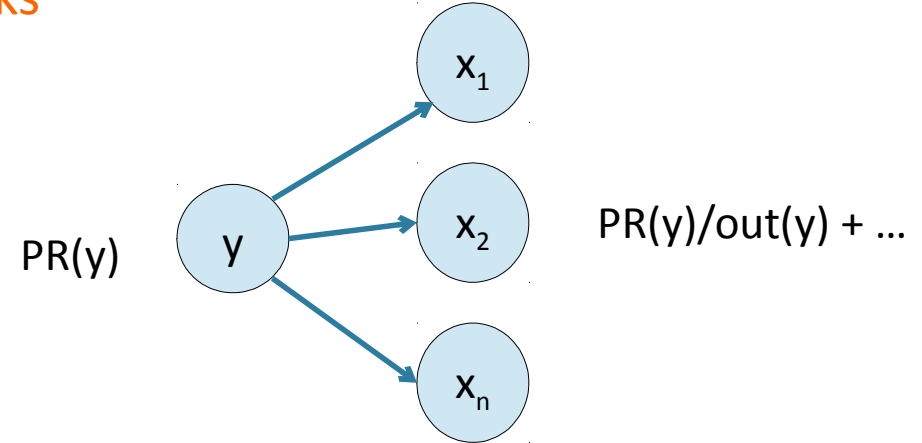


Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

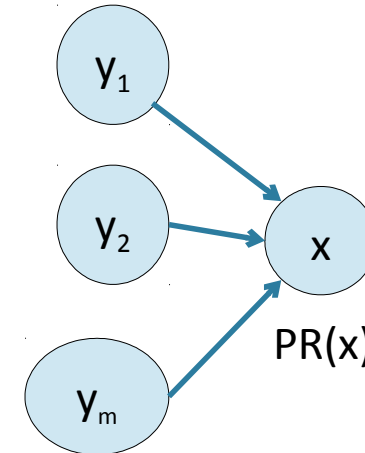
—for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

node, out-links



■ **Reducer** $\left\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\} \right\rangle$

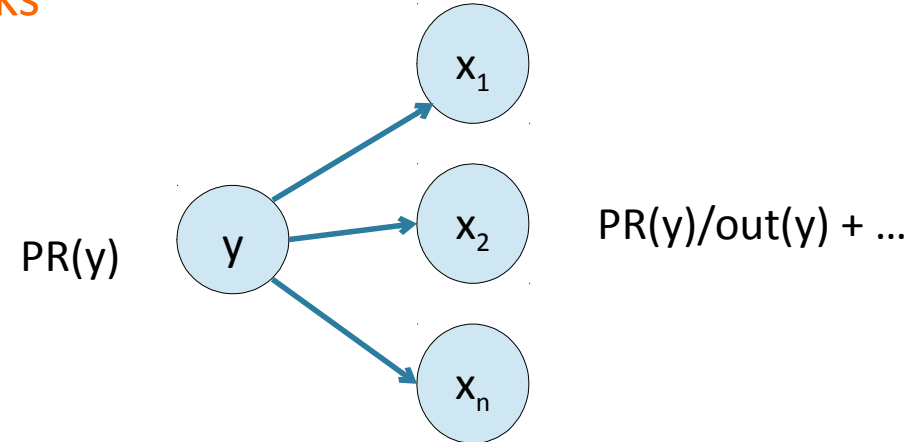
node, ΔPR from in-links



Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$
–for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

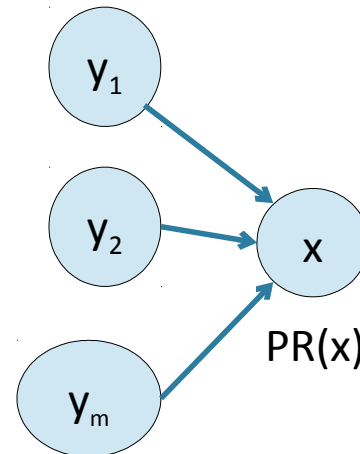
node, out-links



■ **Reducer** $\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\} \rangle$

node, ΔPR from in-links

–**Compute** $PR(x) = \frac{1-\beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

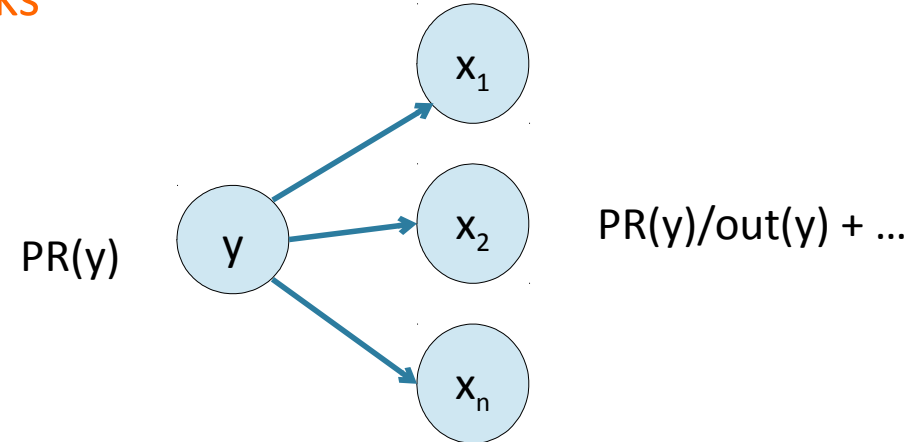


$PR(x) = [\dots] + PR(y_1)/out(y_1) + \dots$

Expressing PageRank with MapReduce

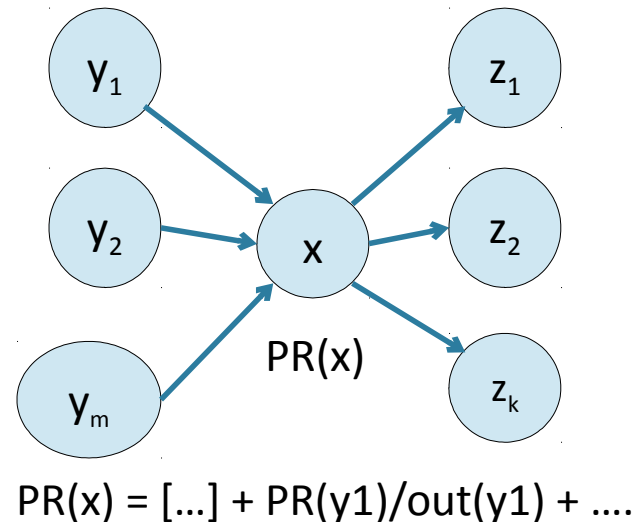
■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$
– for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

node, out-links



■ **Reducer** $\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\} \rangle$
– **Compute** $PR(x) = \frac{1-\beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

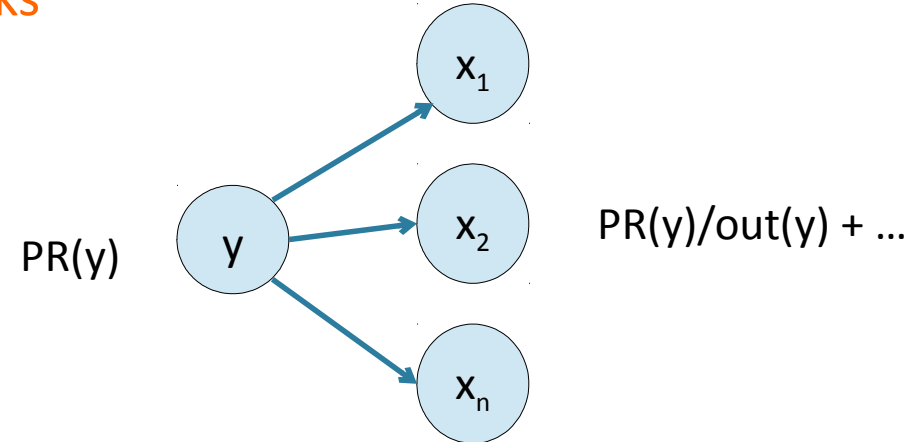
node, ΔPR from in-links



Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$
—for $i = 1..n$, **yield** $\left\langle x_i, \frac{PR(y)}{out(y)} \right\rangle$

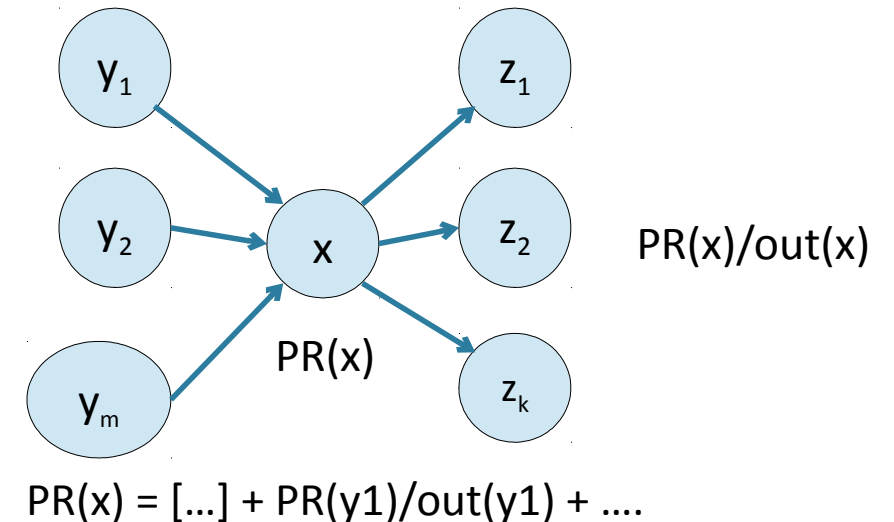
node, out-links



■ **Reducer** $\left\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\} \right\rangle$

node, ΔPR from in-links

—**Compute** $PR(x) = \frac{1-\beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

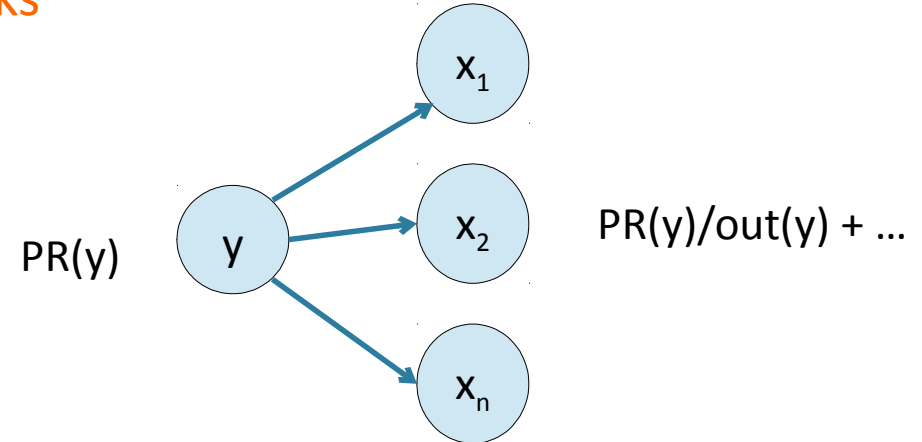


Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

—for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

node, out-links

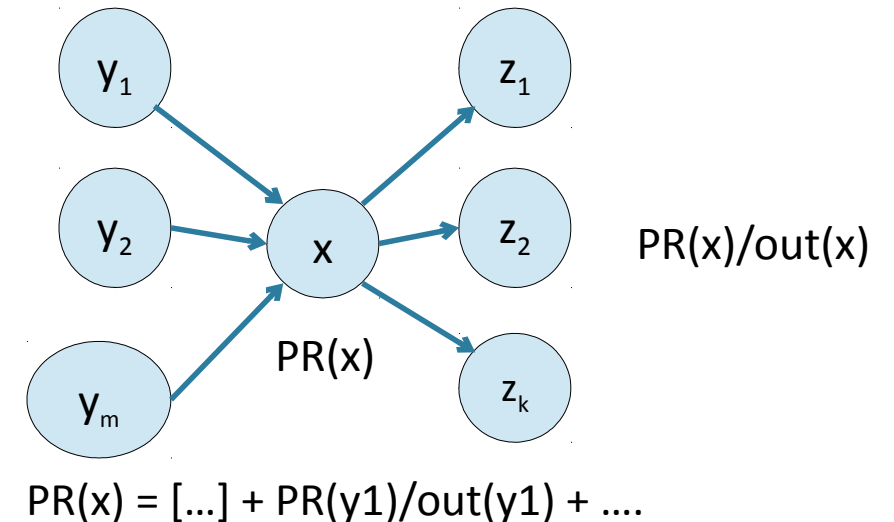


■ **Reducer** $\left\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\} \right\rangle$

—**Compute** $PR(x) = \frac{1-\beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

node, ΔPR from in-links

—for $j = 1, k$ **yield** $\langle z_j, \frac{PR(x)}{out(x)} \rangle$



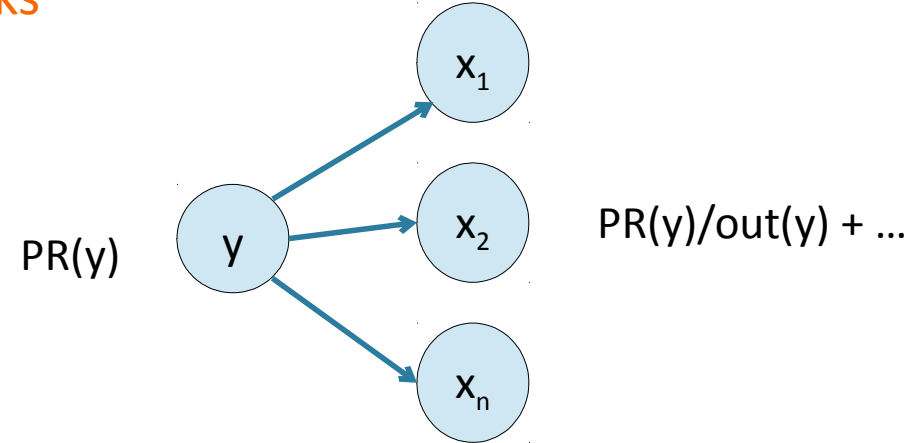
Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

node, out-links

– for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

– **yield** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

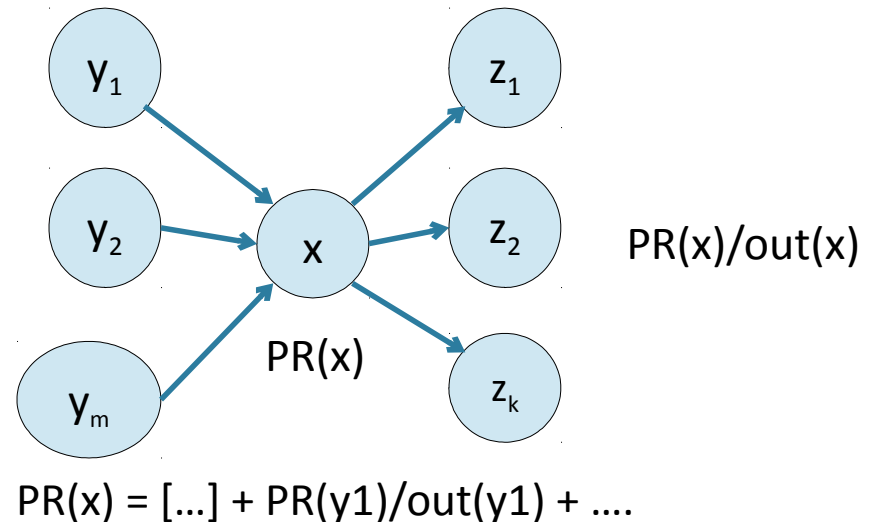


■ **Reducer** $\left\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)} \right\} \right\rangle$

node, ΔPR from in-links

– **Compute** $PR(x) = \frac{1-\beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

– for $j = 1, k$ **yield** $\langle z_j, \frac{PR(x)}{out(x)} \rangle$



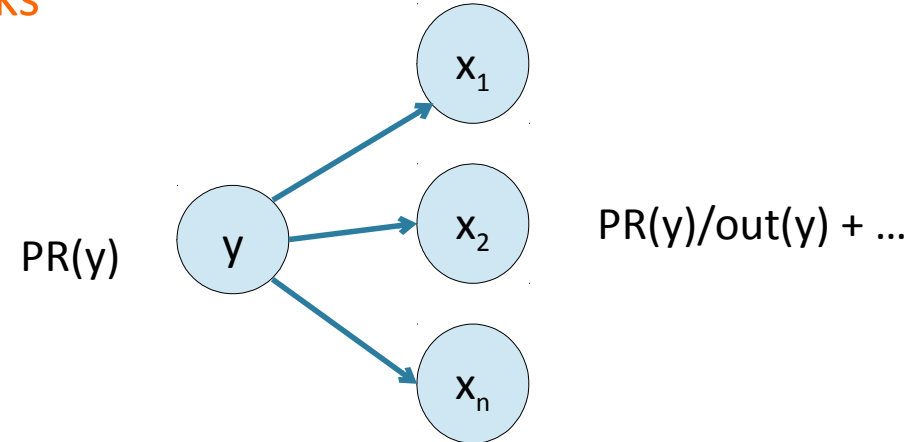
Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

node, out-links

—for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

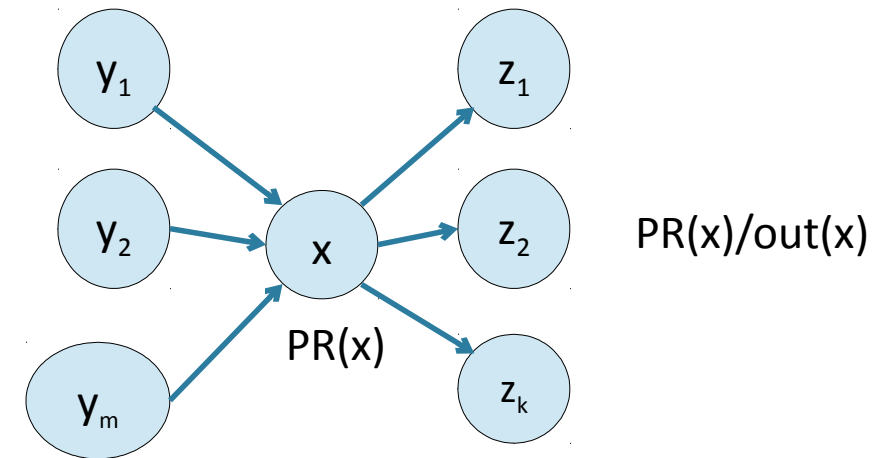
—**yield** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$



■ **Reducer** $\left\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)}, \{z_1, \dots, z_k\} \right\} \right\rangle$ node, ΔPR from in-links

—**compute** $PR(x) = \frac{1 - \beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

—for $j = 1..k$, **yield** $\langle z_j, \frac{PR(x)}{out(x)} \rangle$



$$PR(x) = [\dots] + PR(y_1)/out(y_1) + \dots$$

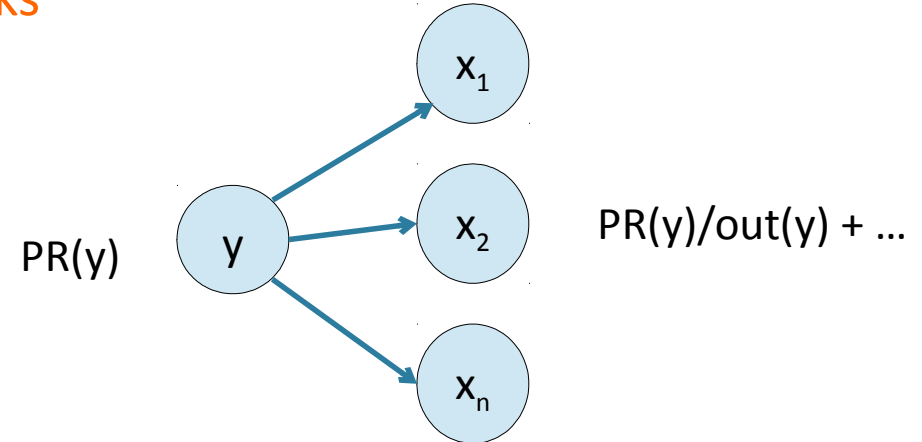
Expressing PageRank with MapReduce

■ **Mapper** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$

node, out-links

–for $i = 1..n$, **yield** $\langle x_i, \frac{PR(y)}{out(y)} \rangle$

–**yield** $\langle y, \{x_1, x_2, \dots, x_n\} \rangle$



■ **Reducer** $\langle x, \left\{ \frac{PR(y_1)}{out(y_1)}, \dots, \frac{PR(y_m)}{out(y_m)}, \{z_1, \dots, z_k\} \right\} \rangle$ node, + ΔPR from in-links

–**compute** $PR(x) = \frac{1 - \beta}{N} + \beta * \sum_{y \rightarrow x} \frac{PR(y)}{out(y)}$

–for $j = 1..k$, **yield** $\langle z_j, \frac{PR(x)}{out(x)} \rangle$
 $\langle x, \{z_1, z_2, \dots, z_k\} \rangle$

