# Mining Massive Datasets

## Lecture 12

**Artur Andrzejak**

[http://pvs.ifi.uni-heidelberg.de](http://pvs.ifi.uni-heidelberg.de)

# Note on Slides

A substantial part of these slides come (either verbatim or in a modified form) from the book
Mining of Massive Datasets
by Jure Leskovec, Anand Rajaraman, Jeff Ullman (Stanford University).
For more information, see the website accompanying the book: http://www.mmds.org.

# Infinite Data

| High dim. data | Graph data | Infinite data | Machine learning | Apps |
|---|---|---|---|---|
| Locality sensitive hashing | PageRank, SimRank | Filtering data streams | SVM | Recommender systems |
| Clustering | Community Detection | Web advertising | Decision Trees | Association Rules |
| Dimensionality reduction | Spam Detection | Queries on streams | Perceptron, kNN | Duplicate document detection |

**Programming in Spark & MapReduce**

# Theory of Stream Processing

- **Last lecture:**

  - **Stream filtering**

  - **Sampling a fixed proportion of a stream**

    - Sample size grows as the stream grows

- **This lecture:**

  - **Sampling a fixed-size sample**

    - Reservoir sampling

# Sampling from a Data Stream: Sampling a fixed-size sample

# Maintaining a fixed-size sample

- **Problem 2: Fixed-size sample**
- Suppose we need to maintain a random sample S of **size exactly *s* tuples**

  - Why? Don't know length of stream in advance
- Suppose at time n we have seen n items

  - Each item is in the sample S with equal prob. s/n

**How to think about the problem: say s = 2**
**Stream:** a x c y z k c d e g…
At **n= 5,** each of the first 5 tuples is included in the sample **S** with equal prob.
At **n= 7,** each of the first 7 tuples is included in the sample **S** with equal prob.
**Impractical solution would be to store all the *n* tuples seen so far and out of them pick *s* at random**

6

# Solution: Fixed Size Sample

- **Algorithm (a.k.a. Reservoir Sampling)**

  - Store all the first $s$ elements of the stream to $S$
  - Suppose we have seen $n-1$ elements, and now the $n^{th}$ element arrives ($n > s$)
    - With probability $s/n$, keep the $n^{th}$ element, else discard it
    - If we picked the $n^{th}$ element, then it replaces one of the $s$ elements in the sample $S$, picked uniformly at random

- **Claim:** This algorithm maintains a sample $S$ with the desired property:

  - After $n$ elements, the sample contains each element seen so far with probability $s/n$

# Proof: By Induction

- **We prove this by induction:**
  - Assume that after *n* elements, the sample contains each element seen so far with probability *s/n*
  - We need to show that after seeing element *n+1* the sample maintains the property
    - Sample contains each element seen so far with probability *s/(n+1)*
- **Base case:**
  - After we see **n=s** elements the sample **S** has the desired property
    - Each out of **n=s** elements is in the sample with probability *s/s = 1*

# Proof: By Induction

- **Inductive hypothesis:** After ***n*** elements, the sample ***S*** contains each element seen so far with prob. ***s/n***
- **Now element *n+1* arrives**
- **Inductive step:** For elements already in ***S***, probability that the algorithm keeps it in ***S*** is:

$$\left(1 - \frac{s}{n+1}\right) + \left(\frac{s}{n+1}\right)\left(\frac{s-1}{s}\right) = \frac{n}{n+1}$$

Element **n+1** discarded    Element **n+1** not discarded    Element in the sample not picked

- So, at time ***n,*** tuples in ***S*** were there with prob. **s/n**
- Time ***n→n+1,*** tuple stayed in ***S*** with prob. **n/(n+1)**
- So prob. tuple is in ***S*** at time ***n+1*** $= \dfrac{s}{n} \cdot \dfrac{n}{n+1} = \dfrac{s}{n+1}$
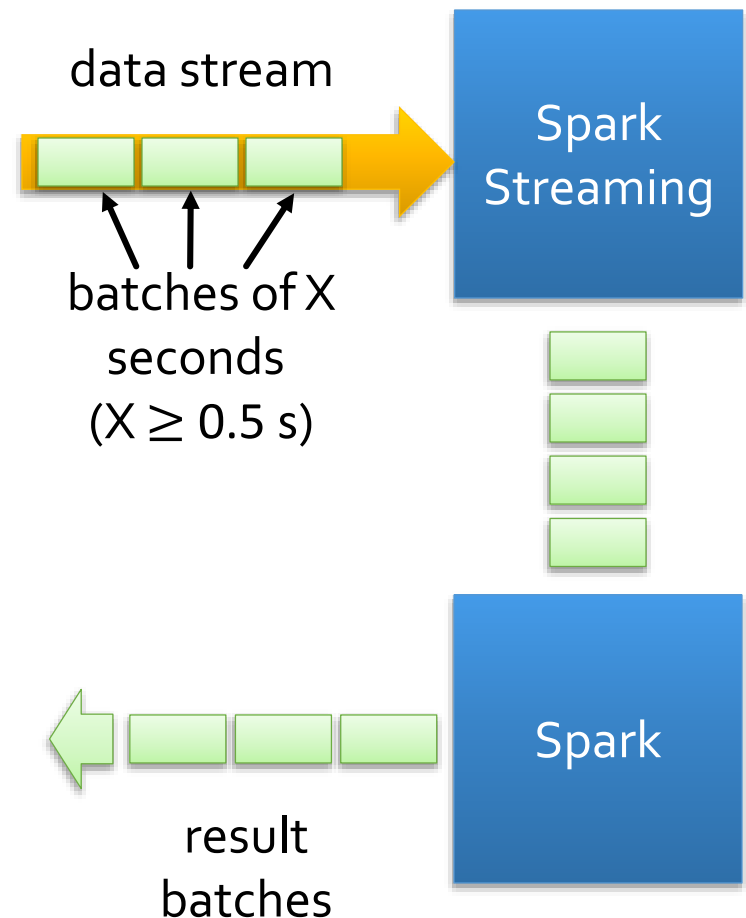
# Spark Streaming

# Advanced Programming

# Spark Streaming: Concept

- Process stream as a **series of small batch jobs**
  - Chop up the live stream into batches of X seconds
  - Spark treats each batch of data as an RDD and processes them using (normal) RDD operations
  - The results of the RDD operations are returned in batches

data stream

batches of X seconds
(X ≥ 0.5 s)

Spark Streaming

Spark

result batches

# Transformations on DStreams

- Many "normal" Spark transformations are available, and some additional ones

| map | flatMap | filter |
|---|---|---|
| repartition | union | count |
| reduce | count | countByValue |
| reduceByKey | join | cogroup |
| **transform** | **updateStateByKey** | |

# Operation "transform"

- The **transform** operation applies an arbitrary RDD-to-RDD function (i.e. a transformation) on each RDD in a DStream
  - Allows using any RDD operation not in the DStream API
- Example: join each RDD in a DStream with additional (precomputed) information

```
# "Normal" RDD containing spam information
spamInfoRDD = sc.pickleFile(<loadPath>)
# join data stream with spamInfoRDD
cleanDStream = wordCounts.transform( lambda
                rdd: rdd.join(spamInfoRDD).filter(...) )
```

# Operation "updateStateByKey"

- Allows to maintain arbitrary state and update it with new data from a stream
  - Input DStream contains (key, value)-pairs
- Two components:
  - **State**: any datatype (e.g. primitive, object, list,…)
  - **Update function f** of the form:
    newState = **f** (newValues, oldState)
- **f** will be called for each key **k** <u>separately</u>;
  newValues  is a sequence of new values (for **k**)
- Usage:
  - statesDStream = inputDStream.updateStateByKey(**f**)

> The updateStateByKey method returns a <u>new DStream</u>
> which contains RDDs that has the state data of each key

# Example for updateStateByKey

- We want to count number of occurrences of each word <u>since the start of the stream</u>
- The input DStream contains pairs (<word>, 1)
- We need a separate state for each encountered word w, namely a count of w (= integer)
- The update function in Python:

> Different for each <word>!

```
def updateFunction (newValues, oldCount):
        if oldCount is None:
                oldCount = 0
        return sum (newValues, oldCount)
```

Python's sum(iterable[, start])
(https://docs.python.org/2/library/functions.html#sum)

This is a sequence of 1's for a current word w since last call of the updateFunction (for w)

15

# Excerpt from stateful_network_wordcount.py ([link](link))

Call:    stateful_network_wordcount.py localhost 9999

```
…
def updateFunc(new_values, last_sum):
        return sum(new_values) + (last_sum or 0)

lines = ssc.socketTextStream(sys.argv[1], int(sys.argv[2]))
running_counts = lines.flatMap(lambda line: line.split(" "))\
                .map(lambda word: (word, 1))\
                .updateStateByKey(updateFunc)

running_counts.pprint()

# Start the processing pipeline
ssc.start()
ssc.awaitTermination()
```
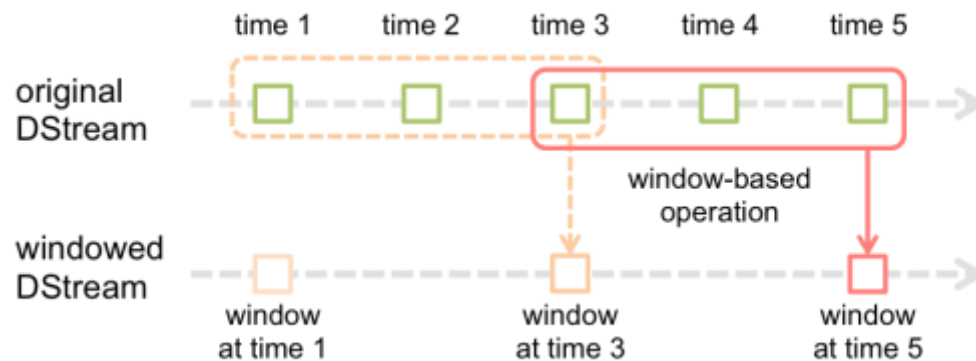
# Spark Streaming

# Window Operations

# Window Operations

- Transformations over a **sliding window** of data, i.e. most recent stream fragment of fixed length



Spark Streaming Programming Guide v1.2,
http://spark.apache.org/docs/latest/streaming-programming-guide.html

- Here, a **windowed operation** is performed every 2 sec. (= **sliding interval**) over the last 3 sec. of data (= **window length**)
- Example: create a new DStream via
            window (windowLength, slideInterval)

# Overview: Window Operations

- **window** (*windowLength*, *slideInterval*)
  - Construct a new Dstream (Example 1)
- **reduceByKeyAndWindow** (*func[, invFunc], windowLength, slideInterval, [numTasks]*)
  - Example 2
- **countByWindow** (*windowLength, slideInterval*)
  - Sliding window count of elements in the stream
- **reduceByWindow** (*func, windowLength, slideInterval*)
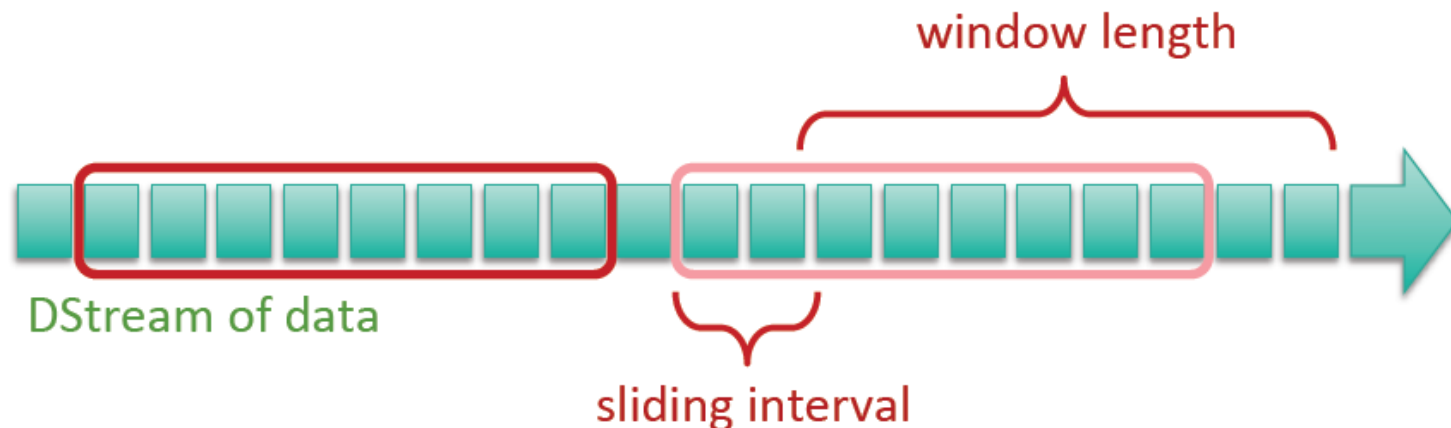  - A new single-element stream from aggregating elements over a sliding interval using func

# Window Op Example 1 (Scala)

val ssc = new StreamingContext (sparkContext, Seconds(1))
val hashTags = tweets.flatMap(status => getTags(status))
val tagCounts = hashTags.**window**(Minutes(1), Seconds(5)).**countByValue**()

Sliding window operation
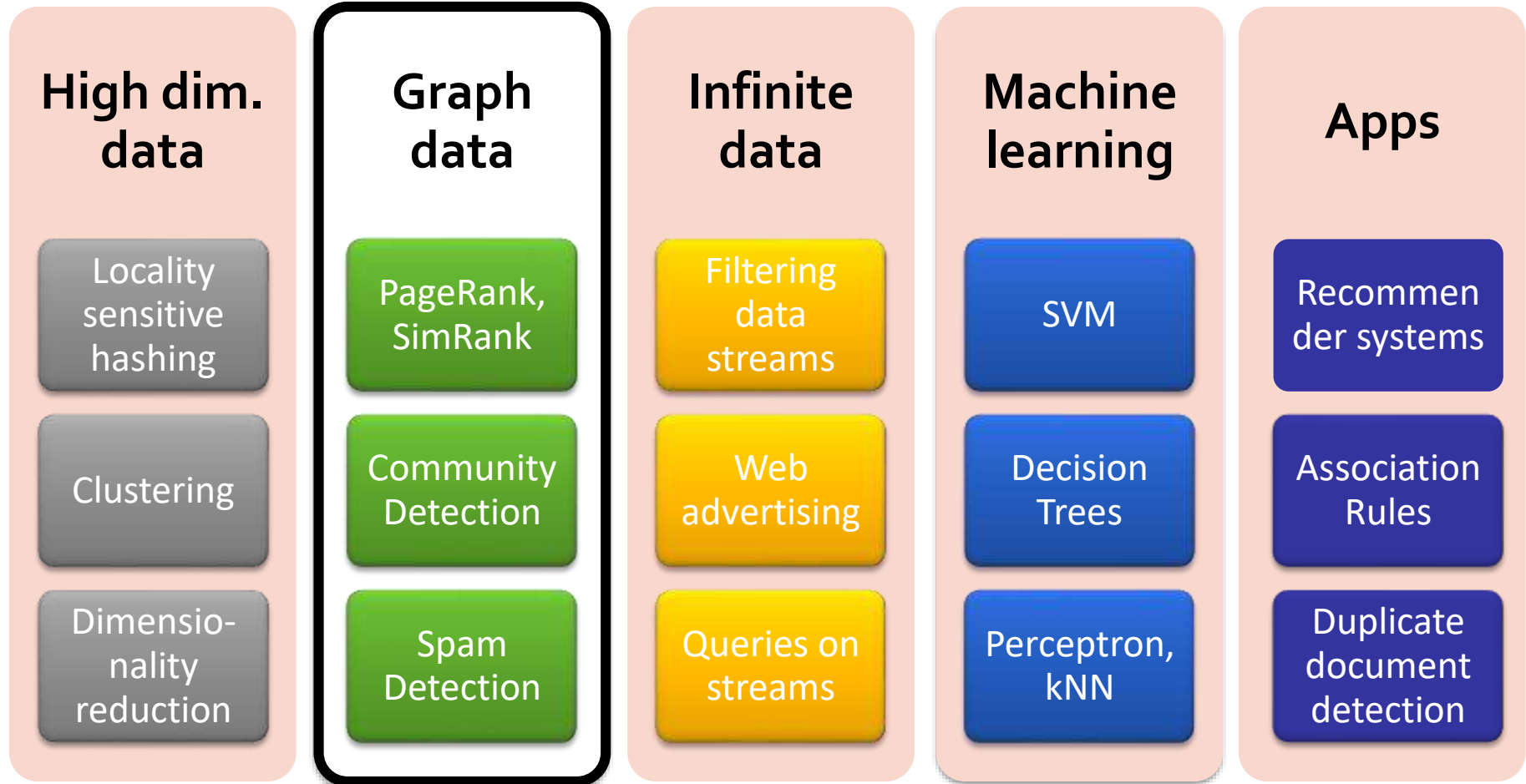
Window **length**

Sliding **interval**



window length

DStream of data

sliding interval

From: Tathagata Das, Spark Streaming, Spark Summit 2014

# Word Count in a Window (Ex. 2)

- We do count of word occurrences every 10 sec. over the last 30 sec. of stream data
  - Again, input stream contains pairs (<word>, 1)

> windowedWordCounts =
> pairs.reduceByKeyAndWindow( lambda x, y: x + y,
> lambda x, y: x - y, 30, 10)

- Here: reduceByKeyAndWindow (func, invFunc, windowLength, slideInterval, [numTasks])
- func is clear (= adding up counts), but why invFunc?

- invFunc "substracts" values which leave the window
  - => Efficient handling by reusing the result of previous window

# Link Analysis

# Infinite Data

| High dim. data | Graph data | Infinite data | Machine learning | Apps |
|---|---|---|---|---|
| Locality sensitive hashing | PageRank, SimRank | Filtering data streams | SVM | Recommender systems |
| Clustering | Community Detection | Web advertising | Decision Trees | Association Rules |
| Dimensionality reduction | Spam Detection | Queries on streams | Perceptron, kNN | Duplicate document detection |

## Programming in Spark & MapReduce

# Web as a Graph

- **Web as a directed graph:**
  - **Nodes: Webpages**
  - **Edges: Hyperlinks**

I teach a class on Networks.

CS224W: Classes are in the Gates building

Computer Science Department at Stanford

Stanford University

# Web as a Graph

- **Web as a directed graph:**
  - **Nodes: Webpages**
  - **Edges: Hyperlinks**

I teach a class on Networks.

CS224W: Classes are in the Gates building

Computer Science Department at Stanford

Stanford University

# Web as a Directed Graph

# Broad Question



- **How to organize the Web?**
- **First try:** Human curated **Web directories**

  - Yahoo, DMOZ, LookSmart

- **Second try: Web Search**

  - **Information Retrieval** investigates: Find relevant docs in a small and trusted set

    - Newspaper articles, Patents, etc.

  - **But:** Web is **huge**, full of untrusted documents, random things, web spam, etc.

# Web Search: 2 Challenges

**Two challenges of web search:**

- **(1) Web contains many sources of information Who to "trust"?**
  - **Trick:** Trustworthy pages may point to each other!

- **(2) What is the "best" answer to query "newspaper"?**
  - No single right answer
  - **Trick:** Pages that actually know about newspapers might all be pointing to many newspapers

# Ranking Nodes on the Graph

- **All web pages are not equally "important"**

  www.joe-schmoe.com vs. www.stanford.edu

- There is large diversity in the web-graph node connectivity
  **Let's rank the pages by the link structure!**

# PageRank

# The "Flow" Formulation

# Links as Votes

- **Idea: Links as votes**

  - **Page is more important if it has more links**

    - In-coming links? Out-going links?

- **Think of in-links as votes:**

  - www.stanford.edu has 23,400 in-links

  - www.joe-schmoe.com has 1 in-link

- **Are all in-links are equal?**

  - **Links from important pages count more**

  - Recursive question!

# Example: PageRank Scores



33

# Simple Recursive Formulation

- Each link's vote is proportional to the **importance** of its source page

- If page $j$ with importance $r_j$ has $n$ out-links, each link gets $r_j / n$ votes

- Page $j$'s own importance is the sum of the votes on its in-links

$$r_j = r_i/3 + r_k/4$$

# PageRank: The "Flow" Model

- A "vote" from an important page is worth more
- A page is important if it is pointed to by other important pages
- Define a "rank" $r_j$ for page $j$

$$r_j = \sum_{i \to j} \frac{r_i}{d_i}$$

$d_i$ ... **out-degree of node** $i$



**"Flow" equations:**

$$r_y = r_y/2 + r_a/2$$
$$r_a = r_y/2 + r_m$$
$$r_m = r_a/2$$

# Solving the Flow Equations

**Flow equations:**
$$r_y = r_y/2 + r_a/2$$
$$r_a = r_y/2 + r_m$$
$$r_m = r_a/2$$

- **3 equations, 3 unknowns, no constants**
  - No unique solution
  - All solutions equivalent modulo the scale factor
- **Additional constraint forces uniqueness:**
  - $r_y + r_a + r_m = 1$
  - **Solution:** $r_y = \frac{2}{5},\ r_a = \frac{2}{5},\ r_m = \frac{1}{5}$
- **Gaussian elimination method works for small examples, but we need a better method for large web-size graphs**
- **We need a new formulation!**

# PageRank: Matrix Formulation

- **Stochastic adjacency matrix $M$**
  - Let page $i$ has $d_i$ out-links
  - If $i \rightarrow j$, then $M_{ji} = \dfrac{1}{d_i}$ else $M_{ji} = 0$
    - $M$ is a **column stochastic matrix** (i.e. columns sum to 1)

- **Rank vector $r$:** vector with an entry per page
  - $r_i$ is the importance score of page $i$
  - $\sum_i r_i = 1$

- **The flow equations can be written as**

$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

$$\boldsymbol{r} = \boldsymbol{M} \cdot \boldsymbol{r}$$

# Example

- Remember the flow equation: $r_j = \sum_{i \to j} \dfrac{r_i}{d_i}$
- Flow equation in the matrix form

$$M \cdot r = r$$

- Suppose page $i$ links to 3 pages, including $j$

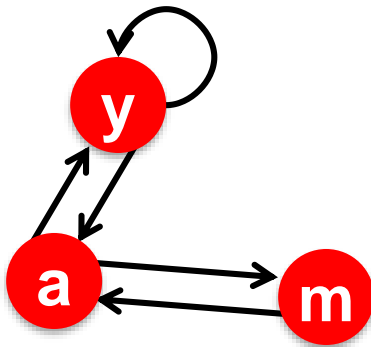# Eigenvector Formulation

- The flow equations can be written
$$r = M \cdot r$$

- So the **rank vector** *r* is an **eigenvector** of the stochastic web matrix *M*

  - In fact, it is its first or **principal** eigenvector, with corresp. eigenvalue *1*

    - Largest eigenvalue of *M* is **1** since *M* is column stochastic (with non-negative entries)

      - *We know **r** is unit length and each column of **M** sums to one, so $Mr \leq 1$*

- **We can now efficiently solve for *r*!**
**The method is called Power iteration**

> **NOTE:** *x* is an **eigenvector** with the corresponding **eigenvalue λ** if:
> $$Ax = \lambda x$$

# Example: Flow Equations & M

|   | y | a | m |
|---|---|---|---|
| **y** | ½ | ½ | 0 |
| **a** | ½ | 0 | 1 |
| **m** | 0 | ½ | 0 |

$$r = M \cdot r$$

$$r_y = r_y/2 + r_a/2$$
$$r_a = r_y/2 + r_m$$
$$r_m = r_a/2$$

$$\begin{bmatrix} y \\ a \\ m \end{bmatrix} = \begin{bmatrix} ½ & ½ & 0 \\ ½ & 0 & 1 \\ 0 & ½ & 0 \end{bmatrix} \begin{bmatrix} y \\ a \\ m \end{bmatrix}$$

# Power Iteration Method

- Given a web graph with *n* nodes, where the nodes are pages and edges are hyperlinks
- **Power iteration**: a simple iterative scheme
    - Suppose there are *N* web pages
    - Initialize: $\mathbf{r}^{(0)} = [1/N,....,1/N]^T$
    - Iterate: $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$
    - Stop when $|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}|_1 < \varepsilon$

        $|\mathbf{x}|_1 = \sum_{1 \le i \le N} |x_i|$ is the **L₁** norm
        Can use any other vector norm, e.g., Euclidean

$$r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i}$$

$d_i$ .... out-degree of node i

# PageRank: How to solve?

- **Power Iteration:**
  - Set $r_j = 1/N$
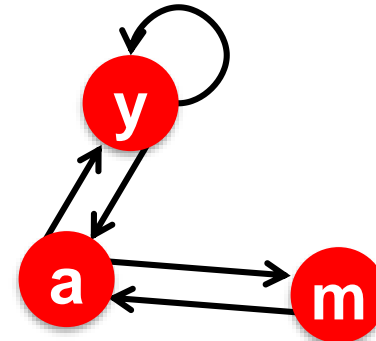  - **1:** $r'_j = \sum_{i \to j} \frac{r_i}{d_i}$
  - **2:** $r = r'$
  - Goto **1**

- **Example:**

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 \\ 1/3 \\ 1/3 \end{matrix}$$

Iteration 0, 1, 2, …

|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 0 |
| a | ½ | 0 | 1 |
| m | 0 | ½ | 0 |

$r_y = r_y/2 + r_a/2$
$r_a = r_y/2 + r_m$
$r_m = r_a/2$

# PageRank: How to solve?

- **Power Iteration:**
  - Set $r_j = 1/N$
  - **1:** $r'_j = \sum_{i \to j} \frac{r_i}{d_i}$
  - **2:** $r = r'$
  - Goto **1**

- **Example:**



|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 0 |
| a | ½ | 0 | 1 |
| m | 0 | ½ | 0 |

$$r_y = r_y/2 + r_a/2$$
$$r_a = r_y/2 + r_m$$
$$r_m = r_a/2$$

$$
\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} =
\begin{array}{cccccc}
1/3 & 1/3 & 5/12 & 9/24 & & 6/15 \\
1/3 & 3/6 & 1/3 & 11/24 & \ldots & 6/15 \\
1/3 & 1/6 & 3/12 & 1/6 & & 3/15 \\
\end{array}
$$

Iteration 0, 1, 2, …

43

# Random Walk Interpretation

- **Imagine a random web surfer:**
  - At any time $t$, surfer is on some page $i$
  - At time $t + 1$, the surfer follows an out-link from $i$ uniformly at random
  - Ends up on some page $j$ linked from $i$
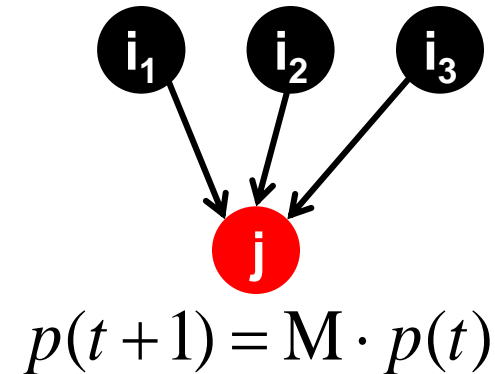  - Process repeats indefinitely
- **Let:**
  - $p(t)$ … vector whose $i^{\text{th}}$ coordinate is the prob. that the surfer is at page $i$ at time $t$
  - So, $p(t)$ is a probability distribution over pages

$$r_j = \sum_{i \to j} \frac{r_i}{d_{\text{out}}(i)}$$

# The Stationary Distribution

- **Where is the surfer at time $t+1$?**

  - Follows a link uniformly at random
    $$p(t+1) = M \cdot p(t)$$

- Suppose the random walk reaches a state
  $$p(t+1) = M \cdot p(t) = p(t)$$

  then $p(t)$ is **stationary distribution** of a random walk

- **Our original rank vector $r$ satisfies $r = M \cdot r$**

  - **So, $r$ is a stationary distribution for the random walk**

$p(t+1) = M \cdot p(t)$

# Existence and Uniqueness

- A central result from the theory of random walks (a.k.a. Markov processes):

> For graphs that satisfy **certain conditions**, the **stationary distribution is unique** and eventually will be reached no matter what the initial probability distribution at time **t = 0**

# Thank you.

Questions?

# Additional Slides

# Why Power Iteration works? (1)

- **Power iteration:**

  A method for finding dominant eigenvector (the vector corresponding to the largest eigenvalue)

  - $r^{(1)} = M \cdot r^{(0)}$

  - $r^{(2)} = M \cdot r^{(1)} = M(Mr^{(1)}) = M^2 \cdot r^{(0)}$

  - $r^{(3)} = M \cdot r^{(2)} = M(M^2 r^{(0)}) = M^3 \cdot r^{(0)}$

- **Claim:**

  Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \ldots M^k \cdot r^{(0)}, \ldots$ approaches the dominant eigenvector of $M$

# Why Power Iteration works? (2)

- **Claim:** Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \dots M^k \cdot r^{(0)}, \dots$ approaches the dominant eigenvector of $M$
- **Proof:**
  - Assume $M$ has $n$ linearly independent eigenvectors, $x_1, x_2, \dots, x_n$ with corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$, where $\lambda_1 > \lambda_2 > \cdots > \lambda_n$
  - Vectors $x_1, x_2, \dots, x_n$ form a basis and thus we can write: $r^{(0)} = c_1\, x_1 + c_2\, x_2 + \cdots + c_n\, x_n$
  - $M r^{(0)} = M(c_1\, x_1 + c_2\, x_2 + \cdots + c_n\, x_n)$
    $$= c_1(Mx_1) + c_2(Mx_2) + \cdots + c_n(Mx_n)$$
    $$= c_1(\lambda_1 x_1) + c_2(\lambda_2 x_2) + \cdots + c_n(\lambda_n x_n)$$
  - **Repeated multiplication on both sides produces** $M^k r^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \cdots + c_n(\lambda_n^k x_n)$

# Why Power Iteration works? (3)

- **Claim:** Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \dots M^k \cdot r^{(0)}, \dots$ approaches the dominant eigenvector of $M$
- **Proof (continued):**
  - Repeated multiplication on both sides produces
    $$M^k r^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \cdots + c_n(\lambda_n^k x_n)$$

  - $M^k r^{(0)} = \lambda_1^k \left[ c_1 x_1 + c_2 \left( \frac{\lambda_2}{\lambda_1} \right)^k x_2 + \cdots + c_n \left( \frac{\lambda_2}{\lambda_1} \right)^k x_n \right]$

  - Since $\lambda_1 > \lambda_2$ then fractions $\frac{\lambda_2}{\lambda_1}, \frac{\lambda_3}{\lambda_1} \dots < 1$

    and so $\left( \frac{\lambda_i}{\lambda_1} \right)^k = 0$ as $k \to \infty$ (for all $i = 2 \dots n$).

  - **Thus: $M^k r^{(0)} \approx c_1 \left( \lambda_1^k x_1 \right)$**
    - Note if $c_1 = 0$ then the method won't converge

51