# Frequent \*

STREAMING ALGORITHMS FOR APPROXIMATIONS

### Reminder

- Please start testing your code on Elephant
- Start thinking about your final project
  - ► Teams
  - Topics
  - Proposals due before fall break
- Review material
  - Mid-term
  - Review class



Obtain the frequency f(i) of each item in the stream



CS 5965/6965 - Big Data Systems - Fall 2014

### Think EXACT VS. APPROXIMATIONS



Lets maintain less than l < d counters (Misra-Gries) With at least 1 empty counter



If an item has a counter, add 1 to that counter



Otherwise, create a new counter and set it to 1



### But now we don't have less than *l* counters





Let  $\delta$  be the median counter value at time t





Decrease all counters by  $\delta$ 



and continue ...

CS 5965/6965 - Big Data Systems - Fall 2014

### The approximated counts are f'

# $f'(\square) = 1$ $f'(\square) = 0$

We increase the count by only 1 for each item appearance

Because we decrease each counter by at most  $\delta_t$  at time t

Calculating the total approximate frequencies:

• Setting  $l = 2/\epsilon$  gives us

We increase the count by only 1 for each item appearance

### $f'(i) \le f(i)$

Because we decrease each counter by at most  $\delta_t$  at time t

$$f'(i) \ge f(i) - \sum_{t} \delta_t$$

Calculating the total approximate frequencies;

$$0 \leq \sum_{i} f'(i) \leq \sum_{t} 1 - \frac{l}{2} \delta_{t} = n - \frac{l}{2} \sum_{t} \delta_{t}$$
$$\sum_{t} \delta_{t} \leq \frac{2n}{l}$$

• Setting  $l = 2/\epsilon$  gives us

 $|f(i) - f'(i)| \le \epsilon n$ 



### What about higher dimensions?



# Matrix Sketching

- Data is usually represented as a matrix
- For most Big Data applications, this matrix is too large for one machine
- In many cases, the matrix is too large to even fit in distributed memory
- Need to optimize for data access
- Streaming algorithm
  - Generate approximation by accessing (streaming) data once



# Matrix Sketching

Efficiently compute a concisely representable matrix B such that

 $B \approx A$  or  $BB^T \approx AA^T$ 

- Working with B is often "good enough"
  - Dimension reduction
  - Classification
  - Regression
  - Matrix Multiplication (approximate)
  - Recommendation systems



### Frequent Directions (Edo Liberty '13)

• Efficiently maintain a matrix B with only  $l = 2/\epsilon$  columns such that,

 $\|AA^T - BB^T\|_2 \le \epsilon \|A\|_f^2$ 

- $\blacktriangleright$  Intuitive approach  $\rightarrow$  extend frequent items
  - How to estimate the frequency of (frequent) items in a streaming fashion?



# Frequent Directions (Liberty 2013)

maintain a sketch of at most *l* columns

maintain the invariant that some columns are empty (zero-valued)

1	2		
3	7		
1	4		
9	2		
3	5		
4	8		
8	1		

d



2

5

6

9

2

3

Stream in matrix *A* one column at a time

Input vectors are simply stored In empty columns



CS 5965/6965 - Big Data Systems - Fall 2014



When the sketch is 'full' We need to zero out some columns

1	2	3	8	6	
3	7	5	7	0	
1	4	4	3	3	
9	2	2	3	7	
3	5	6	6	1	
4	8	1	5	4	
8	1	2	6	2	

d



### Singular Value Decomposition SWISS ARMY KNIFE OF LINEAR ALGEBRA



**Goal**: Given a  $m \times n$  matrix A,

$$A = U \Sigma V^* = \sum_{j=1}^n \sigma_j \boldsymbol{u}_j \boldsymbol{v}_j^*$$

 $m \times n$   $m \times n$   $n \times n$   $n \times n$ 

 $\sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_n \ge 0$  are the singular values of *A*  $u_1, u_2, \dots, u_n$  are orthonormal, the left singular vectors of *A*, and  $v_1, v_2, \dots, v_n$  are orthonormal, the right singular vectors of *A*.

### Singular Value Decomposition

- Closely related problems:
  - Eigenvalue decomposition  $A \approx V \Lambda V^*$
  - Spanning columns or rows  $A \approx C U R$

### ► Applications:

- Principal Component Analysis: Form an empirical covariance matrix from some collection of statistical data. By computing the singular value decomposition of the matrix, you find the directions of maximal variance
- Finding spanning columns or rows: Collect statistical data in a large matrix. By finding a set of spanning columns, you can identify some variables that "explain" the data. (Say a small collection of genes among a set of recorded genomes, or a small number of stocks in a portfolio)
- Relaxed solutions to k-means clustering: Relaxed solutions can be found via the singular value decomposition
- **PageRank**: primary eigenvector

# Singular values, intuition



- Blue circles are m data points in 2D
  The SVD of the m × 2 matrix
  - V<sub>1</sub>: 1st (right) singular vector: direction of maximal variance,
  - σ<sub>1</sub>: how much of data variance is explained by the first singular vector
  - V<sub>2</sub>: 2nd (right) singular vector: direction of maximal variance, after removing projection of the data along first singular vector.
  - σ<sub>2</sub>: measures how much of the data variance is explained by the second singular vector



### SVD - Interpretation

 $M = U\Sigma V^*$  - example:





### SVD - Interpretation

 $M = U\Sigma V^*$  - example:





### SVD - Interpretation

 $M = U\Sigma V^*$  - example:

•  $U\Sigma$  gives the coordinates of the points in the projection axis





### Dimensionality reduction

set the smallest eigenvalues to zero:





### Dimensionality reduction





### Dimensionality reduction





# Frequent Directions (Liberty 2013)

maintain a sketch of at most *l* columns

maintain the invariant that some columns are empty (zero-valued)

1	2		
3	7		
1	4		
9	2		
3	5		
4	8		
8	1		

d



2

5

6

9

2

3

Stream in matrix *A* one column at a time

Input vectors are simply stored In empty columns



CS 5965/6965 - Big Data Systems - Fall 2014



When the sketch is 'full' We need to zero out some columns

1	2	3	8	6	
3	7	5	7	0	
1	4	4	3	3	
9	2	2	3	7	
3	5	6	6	1	
4	8	1	5	4	
8	1	2	6	2	

d



 $B = U\Sigma V^T$ 

 $V^T$ 



$$B_{new} = U\Sigma$$

CS 5965/6965 - Big Data Systems - Fall 2014



 $BB^{T} = B_{new}B_{new}^{T}$ 

The columns of *B* are now orthogonal and in decreasing magnitude

$B_{new} = U\Sigma$						
-9.2	-0.4	-5.1	1.5	0.9		
-10.4	-4.4	2.0	0.5	0.3		
-6.5	-1.9	-1.1	-0.9	-1.9		
-9.9	6.7	0.1	-1.1	-1.4		
-9.6	-3.2	1.1	1.5	-1.4		
-10.1	-0.8	0.3	-4.0	1.6		
-9.0	3.8	2.1	2.7	1.3		

d

CS 5965/6965 - Big Data Systems - Fall 2014



Let  $\delta = \left\| B_{l/2} \right\|^2$ 

				$\rightarrow \delta$	$= \left\  B_{l/2} \right\ ^2$
-9.2	-0.4	-5.1	1.5	0.9	
-10.4	-4.4	2.0	0.5	0.3	
-6.5	-1.9	-1.1	-0.9	-1.9	
-9.9	6.7	0.1	-1.1	-1.4	$\rangle d$
-9.6	-3.2	1.1	1.5	-1.4	
-10.1	-0.8	0.3	-4.0	1.6	
-9.0	3.8	2.1	2.7	1.3	
				,	

d



CS 5965/6965 - Big Data Systems - Fall 2014

### Reduce column $l_2^2$ -norms by $\delta$

-8.9	-0.3		
-10.1	-3.5		
-6.3	-1.5		
-9.6	5.2		
-9.3	-2.5		
-9.8	-0.7		
-8.7	2.9		ر

d



2

5

6

9

2

3

-0.3 -8.9 -3.5 -10.1 -1.5 -6.3 d5.2 -9.6 -9.3 -2.5 -9.8 -0.7 2.9 -8.7

CS 5965/6965 - Big Data Systems - Fall 2014

Start aggregating columns again

 $\langle \rangle_{0}^{\circ}$ 

Input:  $l, A \in \mathbb{R}^{d \times n}$  $B \leftarrow \text{all zeros matrix} \in \mathbb{R}^{d \times l}$ for  $i \in [n]$  do Insert  $A_i$  into a zero valued column of B if B has no zero valued columns then  $[U, \Sigma, V] \leftarrow SVD(B)$  $\delta \leftarrow \sigma_{l/2}^2$  $\hat{\Sigma} \leftarrow \sqrt{\max(\Sigma^2 - I_l \delta, 0)}$  $B \leftarrow U\hat{\Sigma}$ Return: B



### Bounding the error

We first bound  $||AA^T - BB^T||$ 

$$\sup_{\|x\|=1} \|xA\|^{2} - \|xB\|^{2} = \sup_{\|x\|=1} \sum_{t=1}^{n} [\langle x, A_{t} \rangle^{2} + \|xB^{t-1}\|^{2} - \|xB^{t}\|^{2}]$$
$$= \sup_{\|x\|=1} \sum_{t=1}^{n} [\|xC^{t}\|^{2} - \|xB^{t}\|^{2}]$$
$$\leq \sum_{t=1}^{n} \|C^{t^{T}}C^{t} - B^{t^{T}}B^{t}\| \cdot \|x\|^{2}$$
$$= \sum_{t=1}^{n} \delta_{t}$$

Which gives,

$$\|AA^T - BB^T\| \le \sum_{t=1}^n \delta_t$$



### Bounding the error

We compute the Frobenius norm of the final sketch,

$$0 \leq ||B||_{f}^{2} = \sum_{t=1}^{n} [||B^{t}||_{f}^{2} - ||B^{t-1}||_{f}^{2}]$$
  
= 
$$\sum_{t=1}^{n} [(||C^{t}||_{f}^{2} - ||B^{t-1}||_{f}^{2}) - (||C^{t}||_{f}^{2} - ||B^{t}||_{f}^{2})]$$
  
= 
$$\sum_{t=1}^{n} ||A_{t}||^{2} - tr (C^{t^{T}}C^{t} - B^{t^{T}}B^{t})]$$
  
$$\leq ||A||_{f}^{2} - \frac{l}{2} \sum_{t=1}^{n} \delta_{t}$$

Which gives,

$$\sum_{t=1}^n \delta_t \le \frac{2\|A\|_f^2}{l}$$



### Bounding the error

We saw that:

 $\|AA^T - BB^T\| \le \sum \delta_t$ 

and,

 $\overline{\sum \delta_t} \le \frac{2\|A\|_f^2}{l}$ 

setting  $l = \frac{2}{\epsilon}$  gives us,

 $\|AA^T - BB^T\| \le \epsilon \|A\|_f^2$ 



# Divide & Conquer

- Sketching can be implemented in a divide & conquer fashion as well
- Let  $A = [A_1; A_2]$
- Compute the sketches  $B_1, B_2$  of matrices  $A_1, A_2$
- Compute the sketch C of the matrix  $[B_1; B_2]$
- It can be shown that

$$||AA^{T} - CC^{T}|| \le \frac{2||A||_{f}^{2}}{l}$$

