

# Subset Selection

## Deterministic vs. Randomized

**Ilse Ipsen**

North Carolina State University

**Joint work with: Stan Eisenstat, Yale**

**Mary Beth Broadbent, Martin Brown, Kevin Penner**

# Subset Selection

Given: real or complex matrix  $A$   
integer  $k$

Determine permutation matrix  $P$  so that

$$AP = \left( \underbrace{A_1}_k \quad A_2 \right)$$

- **Important columns  $A_1$**   
Columns of  $A_1$  are 'very' linearly independent
- **Redundant columns  $A_2$**   
Columns of  $A_2$  are 'well' represented by  $A_1$

# Subset Selection Requirements

- Important columns  $A_1$

Smallest singular value  $\sigma_k(A_1)$  should be 'large'

- Redundant columns  $A_2$

$\min_Z \|A_1 Z - A_2\|$  should be 'small' (two norm)

# Subset Selection Requirements

- Important columns  $A_1$

Smallest singular value  $\sigma_k(A_1)$  should be 'large'

$$\sigma_k(A)/\gamma \leq \sigma_k(A_1) \leq \sigma_k(A)$$

for some  $\gamma$

- Redundant columns  $A_2$

$\min_Z \|A_1 Z - A_2\|$  should be 'small' (two norm)

$$\sigma_{k+1}(A) \leq \min_Z \|A_1 Z - A_2\| \leq \gamma \sigma_{k+1}(A)$$

for some  $\gamma$

# Outline

- **Deterministic algorithms: strong RRQR, SVD**
- **Randomized 2-phase algorithm**
- **Perturbation analysis of randomized algorithm**
- **Numerical experiments: randomized vs. deterministic**
- **New deterministic algorithm**

# Deterministic Subset Selection

Businger & Golub (1965) **QR with column pivoting**

Faddev, Kublanovskaya & Faddeeva (1968)

Golub, Klema & Stewart (1976)

Gragg & Stewart (1976)

Stewart (1984)

Foster (1986)

T. Chan (1987)

Hong & Pan (1992)

Chandrasekaran & Ipsen (1994)

Gu & Eisenstat (1996) **Strong RRQR**

# First Deterministic Algorithm

Rank revealing QR decomposition

$$AP = Q \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{R}_{22} \end{pmatrix} \quad \text{where} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I}$$

- Important columns

$$Q \begin{pmatrix} \mathbf{R}_{11} \\ \mathbf{0} \end{pmatrix} = \mathbf{A}_1 \quad \text{and} \quad \sigma_i(\mathbf{A}_1) = \sigma_i(\mathbf{R}_{11}) \quad 1 \leq i \leq k$$

- Redundant columns

$$Q \begin{pmatrix} \mathbf{R}_{12} \\ \mathbf{R}_{22} \end{pmatrix} = \mathbf{A}_2 \quad \min_{\mathbf{Z}} \|\mathbf{A}_1 \mathbf{Z} - \mathbf{A}_2\| = \|\mathbf{R}_{22}\|$$

## Strong RRQR (Gu & Eisenstat 1996)

Input:  $m \times n$  matrix  $\mathbf{A}$ ,  $m \geq n$ , integer  $k$

Output:  $\mathbf{AP} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{R}_{22} \end{pmatrix}$   $\mathbf{R}_{11}$  is  $k \times k$

- $\mathbf{R}_{11}$  is well conditioned

$$\frac{\sigma_i(\mathbf{A})}{\sqrt{1 + k(n - k)}} \leq \sigma_i(\mathbf{R}_{11}) \leq \sigma_i(\mathbf{A}) \quad 1 \leq i \leq k$$

- $\mathbf{R}_{22}$  is small

$$\sigma_{k+j}(\mathbf{A}) \leq \sigma_j(\mathbf{R}_{22}) \leq \sqrt{1 + k(n - k)} \sigma_{k+j}(\mathbf{A})$$

- Offdiagonal block not too large

$$\left| \left( \mathbf{R}_{11}^{-1} \mathbf{R}_{12} \right)_{ij} \right| \leq 1$$



# Strong RRQR Algorithm

- 1 Compute some QR decomposition with column pivoting

$$AP_{\text{initial}} = Q \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{R}_{22} \end{pmatrix}$$

- 2 Repeat

Exchange a column of  $\begin{pmatrix} \mathbf{R}_{11} \\ \mathbf{0} \end{pmatrix}$  with a column of  $\begin{pmatrix} \mathbf{R}_{12} \\ \mathbf{R}_{22} \end{pmatrix}$   
Update permutations P, retriangularize

until  $|\det(\mathbf{R}_{11})|$  stops increasing

- 3 Output:  $AP_{\text{final}} = \underbrace{(\mathbf{A}_1)}_k \quad \underbrace{(\mathbf{A}_2)}_{n-k}$

## Second Deterministic Algorithm

Singular value decomposition

$$AP = \underbrace{\begin{pmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{pmatrix}}_{\substack{k \\ n-k}} = \mathbf{U} \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} \mathbf{V}_{11} & \mathbf{V}_{12} \\ \mathbf{V}_{21} & \mathbf{V}_{22} \end{pmatrix}$$

- Important columns  $\mathbf{A}_1$

$$\frac{\sigma_i(\mathbf{A})}{\|\mathbf{V}_{11}^{-1}\|} \leq \sigma_i(\mathbf{A}_1) \leq \sigma_i(\mathbf{A}) \quad \text{for all } 1 \leq i \leq k$$

- Redundant columns  $\mathbf{A}_2$

$$\sigma_{k+1}(\mathbf{A}) \leq \min_{\mathbf{Z}} \|\mathbf{A}_1 \mathbf{Z} - \mathbf{A}_2\| \leq \|\mathbf{V}_{11}^{-1}\| \sigma_{k+1}(\mathbf{A})$$

[Hong & Pan 1992]

# Almost Strong RRQR Algorithm

In the spirit of Golub, Klema and Stewart (1976)

- 1 Compute SVD

$$A = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1 \\ V_2 \end{pmatrix}$$

- 2 Apply strong RRQR to  $V_1$ :  $V_1 P = \underbrace{(V_{11})}_k \underbrace{(V_{12})}_{n-k}$

$$\frac{1}{\sqrt{1 + k(n - k)}} \leq \sigma_i(V_{11}) \leq 1 \quad 1 \leq i \leq k$$

- 3 Output:  $AP = \underbrace{(A_1)}_k \underbrace{(A_2)}_{n-k}$

# Almost Strong RRQR

Produces permutation  $P$  so that

$$AP = \left( \underbrace{A_1}_k \quad \underbrace{A_2}_{n-k} \right)$$

where

- Important columns  $A_1$

$$\frac{\sigma_i(A)}{\sqrt{1 + k(n - k)}} \leq \sigma_i(A_1) \leq \sigma_i(A) \quad 1 \leq i \leq k$$

- Redundant columns  $A_2$

$$\sigma_{k+1}(A) \leq \min_Z \|A_1 Z - A_2\| \leq \sqrt{1 + k(n - k)} \sigma_{k+1}(A)$$

# Deterministic Subset Selection

- Algorithms: strong RRQR, SVD
- **Permuting columns of  $A$  corresponds to permuting right singular vector matrix  $V$**
- Perturbation bounds in terms of  $V$

$$\sigma_{k+1}(A) \leq \min_Z \|A_1 Z - A_2\| \leq \|V_{11}^{-1}\| \sigma_{k+1}(A)$$

- Operation count for  $m \times n$  matrix,  $m \geq n$

$$\min_Z \|A_1 Z - A_2\| \leq \sqrt{1 + f^2 k(n - k)} \sigma_{k+1}(A)$$

in  $\mathcal{O}(mn^2 + n^3 \log_f n)$  flops

# Randomized Algorithms

**Frieze, Kannan & Vempala 1998, 2004**

**Drineas, Kannan & Mahoney 2006**

**Deshpande, Rademacher, Vempala & Wang 2006**

**Rudelson & Vershynin 2007**

**Liberty, Woolfe, Martinsson, Rokhlin & Tygert 2007**

**Drineas, Mahoney & Muthukrishnan 2006, 2008**

**Boutsidis, Mahoney & Drineas 2008, 2009**

**Civril & Magdon-Ismail 2009**

## Survey paper:

**Halko, Martinsson & Tropp 2009**

# 2-Phase Randomized Algorithm

Boutsidis, Mahoney & Drineas 2009

- 1 **Randomized Phase:**  
Sample small number ( $\approx k \log k$ ) of columns
- 2 **Deterministic Phase:**  
Apply rank revealing QR to sampled columns

With 70% probability:

- Two norm

$$\min_{\mathbf{Z}} \|\mathbf{A}_1 \mathbf{Z} - \mathbf{A}_2\|_2 \leq \mathcal{O} \left( k^{3/4} \log^{1/2} k (n - k)^{1/4} \right) \|\Sigma_2\|_2$$

- Frobenius norm

$$\min_{\mathbf{Z}} \|\mathbf{A}_1 \mathbf{Z} - \mathbf{A}_2\|_F \leq \mathcal{O} \left( k \log^{1/2} k \right) \|\Sigma_2\|_F$$

# Deterministic vs. Randomized Algorithms

- Want: permutation  $P$  so that  $AP = \left( \underbrace{A_1}_k \quad \underbrace{A_2}_{n-k} \right)$

- Compute SVD

$$A = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1 \\ V_2 \end{pmatrix}$$

- Obtain  $P$  from  $k$  dominant right singular vectors  $V_1$



# Deterministic vs. Randomized Algorithms

- Want: permutation  $P$  so that  $AP = \left( \underbrace{A_1}_k \quad \underbrace{A_2}_{n-k} \right)$

- Compute SVD

$$A = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1 \\ V_2 \end{pmatrix}$$

- Obtain  $P$  from  $k$  dominant right singular vectors  $V_1$
- **Deterministic:**  
Apply RRQR to **all** columns of matrix  $V_1$
- **Randomized:**  
Apply RRQR to **subset** of columns of **scaled** matrix  $V_1 D$

## 2-Phase Randomized Algorithm

- 1 Compute SVD

$$A = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1 \\ V_2 \end{pmatrix}$$

- 2 **Randomized phase:**

Scale:  $V_1 \rightarrow V_1 D$

Sample  $c$  columns:  $(V_1 D) P_s = \underbrace{(V_{1s} D_s)}_{\hat{c}} *$

- 3 **Deterministic phase:**

Apply RRQR to  $V_{1s} D_s$ :  $\underbrace{(V_{1s} D_s)}_k P_d = \underbrace{(V_{11} D_1)}_k *$

- 4 **Output:**  $AP_s P_d = \underbrace{(A_1)}_k \underbrace{(A_2)}_{n-k}$

## Perturbation Bounds

SVD : 
$$AP = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_{11} & V_{12} \\ V_{21} & V_{22} \end{pmatrix} \quad V_{11} \text{ is } k \times k$$

- For deterministic algorithms

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| / \sigma_k(V_{11})$$

or

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| + \|\Sigma_2 V_{21}\| / \sigma_k(V_{11})$$

## Perturbation Bounds

SVD : 
$$AP = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_{11} & V_{12} \\ V_{21} & V_{22} \end{pmatrix} \quad V_{11} \text{ is } k \times k$$

- For deterministic algorithms

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| / \sigma_k(V_{11})$$

or

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| + \|\Sigma_2 V_{21}\| / \sigma_k(V_{11})$$

- For randomized 2-phase algorithm

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| + \|\Sigma_2 V_{21} D_1\| / \sigma_k(V_{11} D_1)$$

for any nonsingular matrix  $D_1$

# Perturbation Bounds for Randomized Algorithm

- $D_1$  is scaling matrix

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| + \|\Sigma_2 V_{21} D_1\| / \sigma_k(V_{11} D_1)$$

- $V_{11} D_1$  comes from RRQR:  $\underbrace{(V_{1s} D_s)}_{\hat{c}} P_d = \underbrace{(V_{11} D_1)}_k *$

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| + \sqrt{1 + k(\hat{c} - k)} \|\Sigma_2 V_{21} D_1\| / \sigma_k(V_{1s} D_s)$$

# Perturbation Bounds for Randomized Algorithm

- $D_1$  is scaling matrix

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| + \|\Sigma_2 V_{21} D_1\| / \sigma_k(V_{11} D_1)$$

- $V_{11} D_1$  comes from RRQR:  $\underbrace{(V_{1s} D_s)}_{\hat{c}} P_d = \underbrace{(V_{11} D_1)}_k *$

$$\min_Z \|A_1 Z - A_2\| \leq \|\Sigma_2\| + \sqrt{1 + k(\hat{c} - k)} \|\Sigma_2 V_{21} D_1\| / \sigma_k(V_{1s} D_s)$$

- We need **with high probability**:

$$\|\Sigma_2 V_{21} D_1\| \approx \|\Sigma_2\| \quad \sigma_k(V_{1s} D_s) \gg 0$$

# Probabilistic Bounds: Frobenius Norm

Column  $i$  of  $X$  sampled with probability  $p_i$

Scaling matrix  $D_{ii} = 1/\sqrt{p_i}$  with probability  $p_i$

# Probabilistic Bounds: Frobenius Norm

Column  $i$  of  $\mathbf{X}$  sampled with probability  $p_i$

Scaling matrix  $\mathbf{D}_{ii} = 1/\sqrt{p_i}$  with probability  $p_i$

- **Frobenius norm**  $\|\mathbf{X} \mathbf{D}\|_F^2 = \text{trace}(\mathbf{X} \mathbf{D}^2 \mathbf{X}^T)$
- **Linearity**  $\mathbb{E} [\|\mathbf{X} \mathbf{D}\|_F^2] = \text{trace}(\mathbf{X} \underbrace{\mathbb{E} [\mathbf{D}^2]}_I \mathbf{X}^T)$
- **Scaling**  $\mathbb{E} [\mathbf{D}_{ii}^2] = p_i * \frac{1}{p_i} + (1 - p_i) * 0 = 1$
- **Expected value**  $\mathbb{E} [\|\mathbf{X} \mathbf{D}\|_F^2] = \|\mathbf{X}\|_F^2$



# Probabilistic Bounds: Frobenius Norm

Column  $i$  of  $\mathbf{X}$  sampled with probability  $p_i$

Scaling matrix  $\mathbf{D}_{ii} = 1/\sqrt{p_i}$  with probability  $p_i$

- **Frobenius norm**  $\|\mathbf{X} \mathbf{D}\|_F^2 = \text{trace}(\mathbf{X} \mathbf{D}^2 \mathbf{X}^T)$
- **Linearity**  $\mathbb{E} [\|\mathbf{X} \mathbf{D}\|_F^2] = \text{trace}(\mathbf{X} \underbrace{\mathbb{E} [\mathbf{D}^2]}_I \mathbf{X}^T)$
- **Scaling**  $\mathbb{E} [\mathbf{D}_{ii}^2] = p_i * \frac{1}{p_i} + (1 - p_i) * 0 = 1$
- **Expected value**  $\mathbb{E} [\|\mathbf{X} \mathbf{D}\|_F^2] = \|\mathbf{X}\|_F^2$
- **Markov's inequality**

$$\text{Prob} \left[ \|\mathbf{X} \mathbf{D}\|_F^2 \leq \alpha \|\mathbf{X}\|_F^2 \right] \geq 1 - \frac{1}{\alpha}$$

# Randomized Subset Selection: Frobenius Norm

- Perturbation bound

$$\min_{\mathbf{Z}} \|\mathbf{A}_1 \mathbf{Z} - \mathbf{A}_2\|_F \leq \|\boldsymbol{\Sigma}_2\|_F + \sqrt{1 + k(\hat{c} - k)} \|\boldsymbol{\Sigma}_2 \mathbf{V}_{21} \mathbf{D}_1\|_F / \sigma_k(\mathbf{V}_{1s} \mathbf{D}_s)$$

- With probability  $1 - \frac{1}{\alpha}$

$$\|\boldsymbol{\Sigma}_2 \mathbf{V}_{21} \mathbf{D}_1\|_F \leq \sqrt{\alpha} \|\boldsymbol{\Sigma}_2 \mathbf{V}_2\|_F = \sqrt{\alpha} \|\boldsymbol{\Sigma}_2\|_F$$

- Holds for **any** probability distribution

# Randomized Subset Selection: Frobenius Norm

- Perturbation bound

$$\min_{\mathbf{Z}} \|\mathbf{A}_1 \mathbf{Z} - \mathbf{A}_2\|_F \leq \|\boldsymbol{\Sigma}_2\|_F + \sqrt{1 + k(\hat{c} - k)} \|\boldsymbol{\Sigma}_2 \mathbf{V}_{21} \mathbf{D}_1\|_F / \sigma_k(\mathbf{V}_{1s} \mathbf{D}_s)$$

- With probability  $1 - \frac{1}{\alpha}$

$$\|\boldsymbol{\Sigma}_2 \mathbf{V}_{21} \mathbf{D}_1\|_F \leq \sqrt{\alpha} \|\boldsymbol{\Sigma}_2 \mathbf{V}_2\|_F = \sqrt{\alpha} \|\boldsymbol{\Sigma}_2\|_F$$

- Holds for **any** probability distribution
- **Still to show:**  $\sigma_k(\mathbf{V}_{1s} \mathbf{D}_s) \gg 0$  with high probability

## Sampling: Frobenius Norm

When is  $\sigma_k(\mathbf{V}_{1s}\mathbf{D}_s) \gg 0$  with high probability?

- **Expected** number of sampled columns:  $c$
- Column  $i$  of  $\mathbf{V}_1$  sampled with “probability”

$$p_i = \min\{1, c q_i\} \quad \text{where} \quad q_i = \|(\mathbf{V}_1)_i\|_2^2/k$$

Try to sample columns with large norm

- Scaling matrix  $\mathbf{D} = (1/\sqrt{p_1} \quad \dots \quad 1/\sqrt{p_n})$
- If  $c = \Theta(k \log k)$  then with **probability**  $\geq .9$

$$\sigma_k(\mathbf{V}_{1s}\mathbf{D}_s) \geq 1/2$$

[Boutsidis, Mahoney & Drineas 2009]

# Sampling: Frobenius Norm

How many columns should  $V_{1s}$  actually have?

- **Expected** number of sampled columns from  $V_1$ :  $c$
- **Actual** number of columns in  $V_{1s}$ :  $\hat{c}$

$$E[\hat{c}] \leq c$$

# Sampling: Frobenius Norm

How many columns should  $V_{1s}$  actually have?

- **Expected** number of sampled columns from  $V_1$ :  $c$
- **Actual** number of columns in  $V_{1s}$ :  $\hat{c}$

$$E[\hat{c}] \leq c$$

- If  $c q_i \leq 1$  for all  $i$ , and  $\hat{c} \geq k$  then

$$\sigma_k(V_{1s}D_s) \leq \sqrt{\frac{\hat{c}}{c}}$$

- If  $\hat{c} < c/4$  then  $\sigma_k(V_{1s}D_s) < 1/2$

Make sure that **enough** columns are **actually** sampled

## Randomized Subset Selection: Two Norm

$$\min_{\mathbf{Z}} \|\mathbf{A}_1 \mathbf{Z} - \mathbf{A}_2\|_2 \leq \|\boldsymbol{\Sigma}_2\|_2 + \sqrt{1 + k(\hat{c} - k)} \|\boldsymbol{\Sigma}_2 \mathbf{V}_{21} \mathbf{D}_1\|_2 / \sigma_k(\mathbf{V}_{1s} \mathbf{D}_s)$$

- Probability distribution  $\mathbf{p}_i = \min\{1, c \mathbf{q}_i\}$

$$\mathbf{q}_i = \frac{1}{2} \frac{\|(\mathbf{V}_1)_i\|_2^2}{k} + \frac{1}{2} \left( \frac{\|(\boldsymbol{\Sigma}_2 \mathbf{V}_2)_i\|_2}{\|\boldsymbol{\Sigma}_2 \mathbf{V}_2\|_F} \right)^2$$

- With probability  $\geq .9$

$$\|\boldsymbol{\Sigma}_2 \mathbf{V}_{21} \mathbf{D}_1\|_2 \leq \gamma \left( \frac{(n - k + 1) \log c}{c} \right)^{1/4} \|\boldsymbol{\Sigma}_2\|_2$$

[Boutsidis, Mahoney & Drineas 2009]

# Numerical Experiments

## Compare strong RRQR and randomized 2-phase algorithm

- Subset selection for 2 norm
- Matrix orders  $n \leq 500, 2000$
- $0 \leq k \leq 240$
- Matrices:
  - Kahan, random, scaled random, triangular
  - numerical rank  $k$
- Randomized algorithm: **run 40 times**
- **Iterative determination of  $c$** 
  - $c = 2k$
  - while  $\sigma_k(\mathbf{V}_{1s}\mathbf{D}_s) < 1/2$  do  $c = 2 * c$



# Accuracy

Residuals  $\min_z \|A_1 Z - A_2\|_2$  for  $n = 500$  and  $k = 20$

Matrix	RRQR	Randomized		
		max	min	mean
Kahan	$7 \times 10^0$	$2 \times 10^0$	$5 \times 10^{-1}$	$5 \times 10^{-1}$
num. rank k	$7 \times 10^0$	$2 \times 10^1$	$8 \times 10^0$	$1 \times 10^1$
triangular	$3 \times 10^0$	$1 \times 10^0$	$1 \times 10^0$	$1 \times 10^0$
random	$3 \times 10^1$	$3 \times 10^1$	$3 \times 10^1$	$3 \times 10^1$
scaled rand	$4 \times 10^1$	$5 \times 10^1$	$4 \times 10^1$	$5 \times 10^1$

**No significant difference in accuracy between deterministic and randomized algorithms**

## Number of Sampled Columns

Values of  $c$  for  $n = 500$  and  $k = 20$

Matrix	$c$ values tried	most frequent	mean $\hat{c}$
Kahan	40, 80, 160, 320, 640	$2k = 40$	56
num. rank $k$	40, 80, 160	$4k = 80$	83
triangular	40, 80, 160	$4k = 80$	97
random	40, 80, 160	$4k = 80$	93
scaled rand	40, 80, 160	$4k = 80$	97

$c = 4k$  seems to be a good value

## Different Probability Distributions

- Two norm

$$q_i = \frac{1}{2} \frac{\|(\mathbf{V}_1)_i\|_2^2}{k} + \frac{1}{2} \left( \frac{\|(\boldsymbol{\Sigma}_2 \mathbf{V}_2)_i\|_2}{\|\boldsymbol{\Sigma}_2 \mathbf{V}_2\|_F} \right)^2$$

expensive

$$q_i = \frac{1}{2} \frac{\|(\mathbf{V}_1)_i\|_2^2}{k} + \frac{1}{2} \frac{\|(\mathbf{A})_i\|_2^2 - \|(\mathbf{A} \mathbf{V}_1^T \mathbf{V}_1)_i\|_2^2}{\|\mathbf{A}\|_F^2 - \|\mathbf{A} \mathbf{V}_1^T \mathbf{V}_1\|_F^2}$$

numerically unstable (can be negative)

- Frobenius norm

$$q_i = \frac{\|(\mathbf{V}_1)_i\|_2^2}{k}$$

numerically stable and cheap

## Different Probability Distributions

Residuals  $\min_Z \|A_1 Z - A_2\|_2$  for  $n = 500$  and  $k = 20$

Matrix	P	max	min	mean	$\hat{c}$
num. rank k	2	$2 \times 10^1$	$8 \times 10^0$	$1 \times 10^1$	83
	F	$1 \times 10^1$	$7 \times 10^0$	$1 \times 10^1$	88
triangular	2	$1 \times 10^0$	$1 \times 10^0$	$1 \times 10^0$	97
	F	$1 \times 10^0$	$1 \times 10^0$	$8 \times 10^0$	91
random	2	$3 \times 10^1$	$3 \times 10^1$	$3 \times 10^1$	93
	F	$3 \times 10^1$	$3 \times 10^1$	$3 \times 10^1$	87
scaled rand	2	$5 \times 10^1$	$4 \times 10^1$	$5 \times 10^1$	97
	F	$5 \times 10^1$	$4 \times 10^1$	$5 \times 10^1$	93
Kahan	2	$2 \times 10^0$	$5 \times 10^{-1}$	$5 \times 10^{-1}$	56
	F	$7 \times 10^{-1}$	$5 \times 10^{-1}$	$5 \times 10^{-1}$	42

No difference in accuracy between 2 norm and Frobenius norm probability distributions

## Ideas from Randomized Algorithm

$$AP = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1 \\ \mathbf{V}_2 \end{pmatrix} \quad \mathbf{V}_1 = \underbrace{(\mathbf{V}_{11})}_k \quad \underbrace{(\mathbf{V}_{12})}_{n-k}$$

- Order columns of  $\mathbf{V}_1$  in order of decreasing norms
- Apply strong RRQR to columns of largest norm

### Intuition

$$\|\mathbf{V}_{11}\|_F^2 = k - \|\mathbf{V}_{12}\|_F^2$$

This means:  $\|\mathbf{V}_{11}\|_F$  large implies  $\|\mathbf{V}_{12}\|_F$  small

$$\sigma_k(\mathbf{V}_{11})^2 = 1 - \|\mathbf{V}_{12}\|_2^2$$

This means:  $\|\mathbf{V}_{12}\|_2$  small implies  $\sigma_k(\mathbf{V}_{11})$  large

# New 2-Phase Deterministic Algorithm

- 1 Compute SVD

$$A = U \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1 \\ V_2 \end{pmatrix}$$

- 2 **Ordering phase:**

Order according to decreasing column norms

$$V_1 \rightarrow V_1 P_O$$

Select leading  $4k$  columns:  $A P_O = \underbrace{(A_S)}_{4k} *$

- 3 **Rank revealing phase:**

Apply strong RRQR to  $A_S$ :  $A_S P_R = \underbrace{(A_1)}_k *$

## Results for New Deterministic Algorithm

- $n = 2000$  and  $k = 40$
- Residuals  $\min_z \|A_1 Z - A_2\|_2$
- Time ratio  $TR = \text{time}(\text{new algorithm})/\text{time}(\text{RRQR})$

Matrix	Residuals		$\sigma_k(A_1)$		TR
	RRQR	new	RRQR	new	
Kahan	$4 \times 10^0$	$4 \times 10^0$	$8 \times 10^{-2}$	$8 \times 10^{-2}$	0.11
random	$9 \times 10^1$	$9 \times 10^1$	$1 \times 10^1$	$1 \times 10^1$	0.03
s. rand	$1 \times 10^2$	$1 \times 10^2$	$2 \times 10^1$	$2 \times 10^1$	0.07
triang	$4 \times 10^0$	$3 \times 10^1$	$3 \times 10^{-1}$	$3 \times 10^{-1}$	0.02

New algorithm appears to be as accurate as RRQR and possibly faster

# Summary

## Subset selection

Given: real or complex matrix  $A$ , integer  $k$

Want:  $AP = \left( \underbrace{A_1}_{k} \quad A_2 \right)$  with

$$\sigma_k(A_1) \approx \sigma_k(A) \quad \min_Z \|A_1 Z - A_2\| \approx \sigma_{k+1}(A)$$

- Deterministic algorithms: strong RRQR, SVD
- Randomized 2-phase algorithm
- Randomized algorithm: no more accurate than strong RRQR for matrices of order  $\leq 2000$
- Numerical issues with randomized algorithm
- New deterministic algorithm:  
As accurate as strong RRQR and perhaps faster