1) Large scale matrices.   ← features

$$1 \quad 2 \quad \cdots j \cdots n$$

objects $\begin{matrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ m \end{matrix}$

$m \times n$

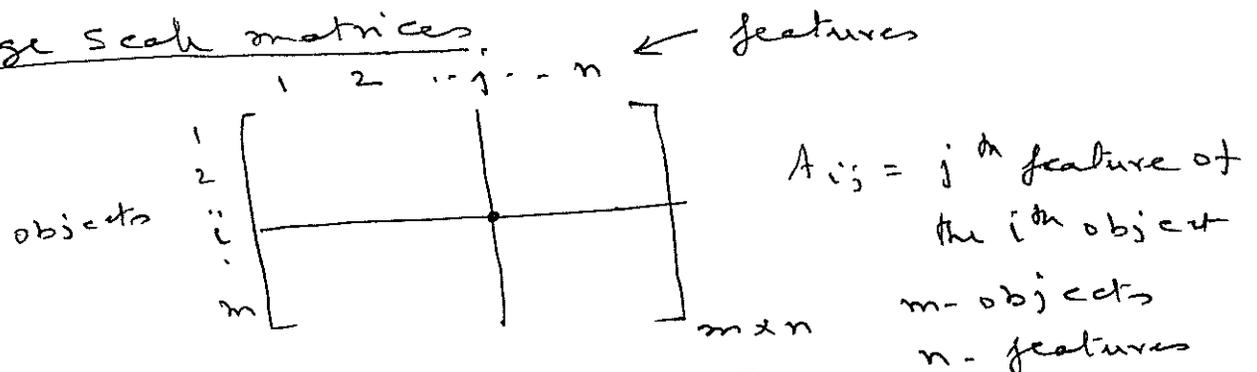$A_{ij} = j^{th}$ feature of the $i^{th}$ object

m - objects
n - features

. Geometrically: m objects in $R^n$.   $m < n$
                  points

1) Astronomy:
   objects: Different angular regions in space
   Features: Elements from frequency bands

2) Genetics
   objects: list of genes
   Features: Level $A_{ij}$ of the $i^{th}$ gene in $j^{th}$ individual

3) Document / Text
   objects: list of documents
   Features: $A_{ij}$ = frequency of $j^{th}$ term in $i^{th}$ document

4) Discretization of ODE / PDE

5) Kernel matrices when describing pairwise relations.

Properties of these matrices

. Singular values decay quickly suggesting the potential to compress or reduce dimension

## Finding Structures with randomness

. Two step process

1) construct a low-dimensional subspace that captures the action of a matrix

2) Project/restrict the matrix to this subspace and process that matrix using standard factorization- QR, SVD, CUR of the resulting reduced matrix.

Stage 1 : Compute an approximate basis for the range of the input matrix $A|_{m \times n}$ consisting of columns of $Q$ which are orthonormal : $Q_{m \times n}$.

$$A \approx Q Q^T A \longrightarrow (1)$$

where $n$ is small but large enough to capture the action of $A$

Stage 2 : Given $Q$ that satisfies (1), use $Q$ to get standard factorization of $A$.

. Stage 1 is implemented by randomized approach - random sampling and random projection

. Stage 2 uses std. deterministic approach

. This separation is key to the success of this approach.

An example: SVD computation

1) From (1): $A \approx Q Q^T A$. SET $B = Q^T A \rightarrow (2)$

$Q|_{m \times n}$ — $n$ is the rank of Approx. $\quad n \times n$

$Q^T Q = I_n, \quad Q Q^T$ - Projection

This gives a _low-rank factorization_

$$A \approx Q B$$

2) Compute SVD of $B = \tilde{U} \Sigma V^T$

3) Set $U = Q \tilde{U}$ and $A \approx Q B$

$$= Q \tilde{U} \Sigma V^T$$
$$= U \Sigma V^T$$

---

_Two ways_

1) _Fixed Precision Approximation_

**Statement**
- Given $A$ and $\varepsilon > 0$, seek a $Q$ with ~~orthogonal~~ orthogonal columns and $k = k(\varepsilon)$:
$$\| A - Q Q^T A \| \leq \varepsilon \longrightarrow (3)$$
- Dimension $k$ of the range of $Q$ captures the action of $A$.

**Pathway to soln**
- SVD furnishes an optimal answer to this fixed precision problem. For each $j$
$$\min_{\text{rank}(x) \leq j} \| A - x \| \leq \sigma_{j+1} \longrightarrow (4)$$
- choosing $x = Q Q^T A$ where $Q_{m \times k}$ has $k$ orthogonal left singular vectors of $A$ which guarantees (3)

~~Which guarantees~~

2) <u>Fixed rank Approx</u>: Given A and k and an over sampling parameter $p > 0$, Construct $Q$ with $(k+p)$ orthonormal columns

~~$BA - QRQ^T A$~~ $\| A - Q Q^T A \| = \min_{\text{rank}(x) \leq k} \| A - x \|$

$\rightarrow$ <mark>(5)</mark>

<u>Remark</u> ~~Note~~

<u>How randomness helps to solve fixed rank problems?</u>

· <u>Seek a basis for Range (A) with rank k:</u>

· Let $w$ be a random vector. ~~and $y = Aw$~~ Then $y = Aw$ is a random sample from Range (A)

· Build $Y = \{y_1, y_2 \cdots y_k\}$ : $y_i = A w_i$

· This set Y is in <u>general position</u>.

· Then <u>orthogonalize</u> Y.

<u>Note</u>: Let $A = B + E$ $\begin{cases} B - \text{rank } k \\ E - \text{Perturbation} \end{cases}$

<u>oversampling</u> $\boxed{W_i \in \mathbb{R}^n}$  $y_i = A w_i = B w_i + E w_i$  $1 \leq i \leq k+p \rightarrow$ <mark>(6)</mark>

· Goal is to find a basis that covers B

· E shifts the direction sample vectors outside of the range of B which can prevent the Span (Y) to cover the Range (B)

· The enriched set in (b) can reduce the leak outside of Range (B)

- Small $p$ helps $p = 5$ to $10$.
- This is the basis for randomized algorithm.

## Fixed Rank Algorithm: Prototype

Given: $A \in \mathbb{R}^{m \times n}$ Target Rank $= k$,
over sampling parameter $p$.

output: $Q \big|_{m \times (k+p)}$

1) Draw $n \times (k+p)$ random matrix $\Omega \big|_{n \times (k+p)}$

2) Compute $Y \big|_{m \times (k+p)} = A \big|_{m \times n} \Omega \big|_{n \times (k+p)} = [y_1, y_2 \cdots y_{k+p}]$

3) orthonormalize the columns of $Y = QR$
and get $Q \big|_{m \times (k+p)}$. This is the basis
for Range $(A)$

Gram-Schmidt
Householder
Givens

## Performance Analysis:

Theorem 1.1 $A \in \mathbb{R}^{m \times n}$, select a target rank $k$
$\geq 2$ and an oversize parameter $p \geq 2$ where
$(k+p) < \min\{m, n\}$. Execute the prototype
algorithm with $\Omega = [\Omega_{ij}]$, $\Omega_{ij} \sim iid$
$N(0,1)$ to obtain $Q \big|_{m \times (k+p)}$ with
orthonormal columns.

Then
$$\mathbb{E} \| A - QQ^T A \| \leq \left\{ \left[ 1 + \frac{4\sqrt{k+p}}{p-1} \right] (m \wedge n)^{1/2} \right\} \sigma_{k+1}$$
$\downarrow$
(w.r.to random $\Omega$)
$\rightarrow (7)$

- known as a sharper than deterministic analog based on RRQR Algorithm of Gu and Eisenstat (1996)

- Thanks to measure concentration: the following holds:

$$\| A - QQ^\top A \| \leq \left[ 1 + 11 \sqrt{k+p} \cdot (m \wedge n)^{1/2} \right] \sigma_{k+1}$$

holds with ~~at least~~ probability $(1 - 6 p^{-p})$. $p = 5$ gives good results.

## Historical facts. (Existence)

**Existence**

1) **Column selection**: Every $A \in \mathbb{R}^{m \times n}$

Contains a ~~submatrix~~ C with $k$ ~~columns~~ Columns s.t.

$$\| A - CC^+ A \| \leq \sqrt{1 + k(n-k)} \, \| A - A_{(k)} \| \qquad \rightarrow (8)$$

where $A_{(k)}$ is the ~~best~~ $k$-rank approximation to A.

[ A. F. Ruston (1964) "Auerbach's Theorem", Math. Proc. Cambridge Phil. Society Vol 56 (1964) pp 476-480

Note: . Column Selection is NP-hard.
. Efficient deterministic RRQR method will achieve the above bound in (8)

Column Selection

2) Using this there exists randomized algorithm for fixed rank approximation:

$$\| A - QQ^\top A \| \leq \min_{\text{rank}(x) \leq k} \| A - x \| \quad \rightarrow (9)$$
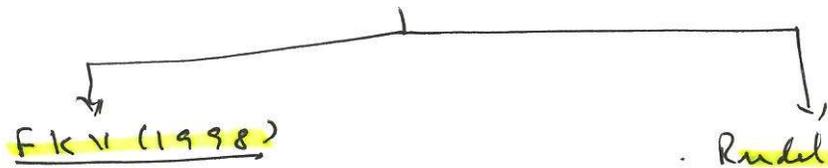
· This approach use sampling probabilities
based on
    (1)  Squared 2-norms of rows / columns
    (2)  Leverage Score that reflects the
         relative importance of columns

· Columns are selected using this distribution

· Earliest method for randomized Column Selection
    is by Frieze, kannan, Vempala (1998)

**FKV (1998)**

· Given A and $l(\varepsilon, k)$
select B :
$$\| A - B \|_F \leq \| A - A_{(k)} \|_F + \varepsilon \| A \|_F$$

**Rudelson + Vershynin (2007)**

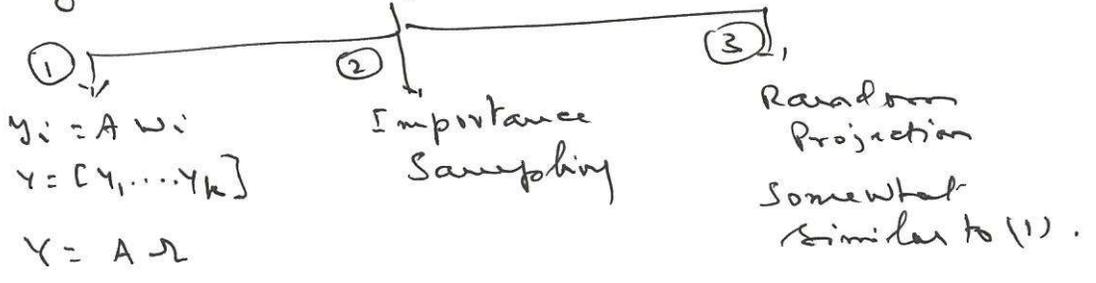· using the same sampling, proved $\| \cdot \|_2$ bounds

· **Deshpande et al (2006)**
$$\| A - B \|_F \leq (1+\varepsilon) \| A - A_{(k)} \|_F$$

**3 Dimension reduction** using random Projection

    · Started in 1984 - Johnson - Lindenstrauss

    · $l_2$ embedding

Summary :-
    · Stage 1 is realized in one of 3-ways

    ① $y_i = A w_i$
    $Y = [y_1, ..., y_k]$
    $Y = A \Omega$

    ② Importance Sampling

    ③ Random Projection Somewhat similar to (1).

# Sampling based randomized Algorithms

Fact    $X \in R^n$ a random Vector           $\rightarrow$ (1)

- $E(x) \in R^n$ _____

- $Var(x) = E\left[ \| x - E(x) \|^2 \right]$

$$= E\left[ x - E(x) \right]^T \left[ x - E(x) \right]$$

$$= \sum_i E\left( x_i - E(x_i) \right)^2 - \text{Total Variance}$$

$E\left[ xx^T - x^T E(x) - E(x)^T x + E^T(x) E(x) \right]$

$= E[xx^T] - E^T(x) E(x)$

$= E[\| x\|^2] - \| E(x)\|^2$

(2a)

- $Cov(x) = E\left[ x - E(x)\right]\left[ x - E(x) \right]^T \rightarrow$ (2)

- $Var(x) = tr\left[ Cov(x) \right] \longrightarrow$ (3)

_____

| Problem 1 | Matrix - Vector product

- $A \in R^{m \times n}$, $V \in R^n$, $AV = \sum_j A_{*j} V_j \rightarrow$ (4)

$$= \text{Sum of } n \text{ vectors}$$

- $(AV)$ can be estimated using a sample of column vectors of A

- Error in the estimate is measured by the variance of the estimate.

- Two ways to sample the Columns:

UNIFORM            Importance
DISTRIBUTION        distribution

## Motivating
## Exercise 1   Scaled random Variable

• Let $x$ be a r.v. and $\boxed{x = \dfrac{a_i}{p_i}}$ with prob $= p_i$

with $\sum\limits_{i=1}^{n} p_i = 1$.   Let $n = 2$ for simplicity

• $E(x) = \dfrac{a_1}{p_1} \cdot p_1 + \dfrac{a_2}{p_2} \cdot p_2 = a_1 + a_2 \overset{\Delta}{=} \bar{a}$

• $Var(x) = p_1 \left[ \dfrac{a_1}{p_1} - \bar{a} \right]^2 + p_2 \left[ \dfrac{a_2}{p_2} - \bar{a} \right]^2$

$= p_1 \left[ \dfrac{a_1^2}{p_1^2} - \dfrac{2 a_1 \bar{a}}{p_1} + \bar{a}^2 \right] + p_2 \left[ \dfrac{a_2^2}{p_2^2} + \bar{a}^2 - \dfrac{2 a_2 \bar{a}}{p_2} \right]$

$= \dfrac{a_1^2}{p_1} + \dfrac{a_2^2}{p_2} + (p_1 + p_2) \bar{a}^2 - (2 a_1 \bar{a} + 2 a \bar{a})$

$= \dfrac{a_1^2}{p_1} + \dfrac{a_2^2}{p_2} + \bar{a}^2 - 2 \bar{a}(a_1 + a_2)$
$\begin{bmatrix} (a_1+a_2)^2 \\ -2(a_1+a_2)(a_1+a_2) \\ = -(a_1+a_2)^2 \end{bmatrix}$

$= \dfrac{a_1^2}{p_1} + \dfrac{a_2^2}{p_2} - (a_1 + a_2)^2$

---

**Special Case**   Set $p_1 = p_2 = \frac{1}{2}$ — __uniform Sampling__

$E(x) = a_1 + a_2$

$Var(x) = 2 a_1^2 + 2 a_2^2 - (a_1^2 + 2 a_1 a_2 + a_2^2)$

$= a_1^2 - 2 a_1 a_2 + a_2^2 = (a_1 - a_2)^2 \geqslant 0$

• This is higher than the sampling below.

---

**optimal $p_1$ & $p_2$**   What choice of $p_i$'s will min. $Var(x)$?

$Min \left( \dfrac{a_1^2}{p_1} + \dfrac{a_2^2}{p_2} \right)$ when $p_1 + p_2 = 1$

$L(p, \lambda) = \dfrac{a_1^2}{p_1} + \dfrac{a_2^2}{p_2} + \lambda(p_1 + p_2 - 1)$

$\dfrac{\partial L}{\partial p_1} = -\dfrac{a_1^2}{p_1^2} + \lambda = 0 \implies \lambda = \dfrac{a_1^2}{p_1^2} \implies \boxed{p_1^2 = \dfrac{a_1^2}{\lambda}}$   $\cdots$

$$\frac{\partial L}{\partial p_2} = -\frac{a_2^2}{p_2^2} + \lambda = 0 \implies \lambda = \frac{a_2^2}{p_2^2} \implies \boxed{p_2^2 = \frac{a_2^2}{\lambda}}$$

· But $p_1 + p_2 = 1 \implies \frac{a_1}{\sqrt{\lambda}} + \frac{a_2}{\sqrt{\lambda}} = 1 \implies a_1 + a_2 = \sqrt{\lambda}$

∴ optimal $\quad p_1^* = \frac{a_1}{a_1 + a_2} \qquad p_2^* = \frac{a_2}{a_1 + a_2}$

$$\text{Var}(x) = \frac{a_1^2}{p_1} + \frac{a_2^2}{p_2} - (a_1 + a_2)^2$$

$$= \frac{a_1^2 (a_1 + a_2)}{a_1} + \frac{a_2^2 (a_1 + a_2)}{a_2} - (a_1 + a_2)^2$$

$$= a_1 (a_1 + a_2) + a_2 (a_1 + a_2) - (a_1 + a_2)^2$$

$$= (a_1 + a_2)^2 - (a_1 + a_2)^2 = 0$$

∴ Variance reduction is ~~possible~~ using a data dependent sampling distribution

---

· <u>Return to Matrix- Vector multiply</u>

· What is a good sampling distribution for the columns of $A$?

· { Define a random vector $x = \frac{A_{*j} \, v_j}{p_j}$

with probability $p_j$, $1 \le j \le n$

· scaled random vector scaled by $p_j$

· $E(x) = \sum_j p_j \frac{A_{*j} \, v_j}{p_j} = \sum A_{*j} v_j = Av$

$\longrightarrow (4)$

(ii) $x$ is an <u>unbiased random vector</u>

$$Var(x) = E\|x\|^2 - \|E(x)\|^2 \qquad \boxed{\text{from (2a) in page (1)}}$$

$$= \sum_j \frac{\|A_{*j} V_j\|^2}{P_j^2} P_j - \|AV\|^2 \qquad \boxed{\text{From (4)}}$$

$$(\because \|ax\| = |a|\|x\|)$$

$$\boxed{\text{Similar to the example}}$$

$$= \sum_j \frac{N_j^2 \|A_{*j}\|^2}{P_j} - \underbrace{\|AV\|^2}_{\substack{\text{constant} \\ \text{w.r. to } P_j}} \longrightarrow \boxed{(5)}$$

. choice of $P_j$ depends on the first term on the r. h. s. of (5).

. $LS_{col}(A)$ — $\boxed{\begin{array}{l}\text{Distribution that depends on the} \\ \text{Squared length of columns of } A\end{array}}$

$\boxed{\begin{array}{c}\text{choice of} \\ P_j\end{array}}$ . column $j$ is picked with prob. $\left.\right\} = P_j = \dfrac{\|A_{*j}\|^2}{\sum_j \|A_{*j}\|^2}$

$$= \frac{\|A_{*j}\|^2}{\|A\|_F^2} \longrightarrow \boxed{(6)}$$

. Substitute (6) in (5)

$$Var(x) = \sum_j N_j^2 \frac{\|A_{*j}\|^2}{\|A_{*j}\|^2} \cdot \|A\|_F^2 - \|AV\|^2$$

$$= \|A\|_F^2 \sum_j N_j^2 - \|AV\|^2$$

$$= \|A\|_F^2 \|N\|_2^2 - \|AV\|^2$$

$\boxed{\begin{array}{l}\text{useful when} \\ \|AV\| \text{ is comparable} \\ \text{to } \|A\|_F \|N\|_2\end{array}}$

$$\leq \|A\|_F^2 \|N\|_2^2 \longrightarrow \boxed{(7)}$$

$$LS_{col}(A) = \left\{ P_1, P_2, \cdots P_m \mid \sum P_j = 1, P_j \text{ is in (6)} \right\}$$

. To approximate $AV$ :

. Perform $s$ trials (with replacement) and take the average

$$Y = \frac{1}{s} \sum_{t=1}^{s} \left[ \frac{A_{* j_t} N_{j_t}}{P_{j_t}} \right] \approx AV$$

$$= \frac{1}{s} \sum_{t=1}^{s} \qquad X_t = \sum_{t} \left( \frac{X_t}{s} \right)$$

$\therefore E\{Y\} = \frac{1}{s} \sum_{t=1}^{s} E(X_t) = \frac{s (AV)}{(s)} = AV$

. $\text{Var}(Y) = \frac{1}{s^2} \sum_{t} \text{Var}(X_t) \qquad$ [using (7)]

$$\leq \frac{1}{s^2} \cdot \left[ \sum_{j=1}^{n} \| A \|_F \| N \|_2^2 \right]$$

$$= \frac{1}{s} \| A \|_F^2 \| N \|_2^2 \longrightarrow \text{(8)}$$

<u>Note :-</u>

(1) We could row sampling as well using a similar argument to get $RS(A)$ distribution.

Problem 2    MATRIX- MATRIX Product

$$C = A \begin{vmatrix} B \end{vmatrix}_{n \times k}$$
$$m \times n$$

— Three ways of multiplying Matrices

$$C = AB = (AB_{*1}, AB_{*2}, \cdots AB_{*k})$$

. Apply Matrix-Vector product routine $k$ times and collate the result.

claim : $\begin{bmatrix} n \text{ is the common \# of columns of } A \text{ and} \\ \text{rows of } B. \end{bmatrix}$

. Define $p_i$ : $\sum_{i=1}^{n} p_i = 1$

. Define a random matrix $\overbrace{x_j =}^{\text{with } x=} \dfrac{A_{*j} B_{j*}}{p_j} \rightarrow (1)$

with probability $p_j$ , $1 \leq j \leq n$

. $E(x) = \sum_{j=1}^{n} x_j p_j = \sum_{j=1}^{n} A_{*j} B_{j*} = AB$
$$\rightarrow (2)$$

. $p_j = \dfrac{\|A_{*j}\|^2 \|B_{j*}\|^2}{\sum_{j=1}^{n} \|A_{*j}\|^2 \|B_{j*}\|^2} \qquad 1 \leq j \leq n \rightarrow (3)$

. Define $Y = \dfrac{1}{s} \sum_{t=1}^{s} x_t = \dfrac{1}{s} \sum_{t=1}^{s} \left( \dfrac{A_{*j_t} B_{j_t *}}{p_{j_t}} \right)$
$$\rightarrow (4)$$

. $E(Y) = AB \qquad \Rightarrow \text{unbiased} \rightarrow (4)$

. $Var(Y) = \dfrac{1}{s} \sum_{t=1}^{s} \dfrac{\|A_{*j_t}\|^2 \|B_{j_t *}\|^2}{p_j} \rightarrow (5)$

$\boxed{\text{using } p_j \text{ as in (3)}}$

$$- \|AB\|_F^2$$

$$\leq \dfrac{1}{s} \|A\|_F^2 \|B\|_F^2 \rightarrow (6)$$

# Implementation

- Let $k_1, k_2, \ldots k_s$ be $s$ integers chosen in $s$ trials.

- Then,

$$\frac{1}{s} \sum x_i = \frac{1}{s} \left[ \frac{A_{*k_1} B_{k_1 *}}{p_{k_1}} + \frac{A_{*k_2} B_{k_2 *}}{p_{k_2}} + \cdots \frac{A_{*k_s} B_{k_s *}}{p_{k_s}} \right]$$

$$= \left[ A_{*k_1}, A_{*k_2}, \ldots A_{*k_s} \right] \begin{bmatrix} B_{k_1 *} / s\, p_{k_1} \\ B_{k_2 *} / s\, p_{k_2} \\ \vdots \\ B_{k_s *} / s\, p_{k_s} \end{bmatrix}$$

$$= C \, \tilde{B}$$

$\downarrow$
chosen
columns
of $A$

$\searrow$ chosen rows of $\tilde{B}$
scaled by $s\, p_k$

$$\begin{bmatrix} A \end{bmatrix} \begin{bmatrix} B \end{bmatrix} \approx \begin{bmatrix} C \end{bmatrix} \begin{bmatrix} \tilde{B} \end{bmatrix} \quad \rightarrow \quad \textcircled{7}$$

$\underset{m \times n}{} \quad \underset{n \times p}{} \quad \underset{m \times s}{} \quad \underset{s \times p}{}$

- Error bound is given by (6) above.

Special cases:

1) If $B = A^T$, we can approximate the Gramian $A A^T$ and hence the singular values of $A$.

2) When $B = A^T$: Start with

$$\text{Var}(Y) \leq \frac{1}{s} \sum \frac{\|A\|_F^2 \, \|A^T\|_F^2}{p_j} - \|A A^T\|_F^2$$

$$= \frac{1}{s} \sum \underbrace{\frac{\|A\|_F^4}{p_j} - \|A A^T\|_F^2}_{\text{minimize w.r. to } p_j} \longrightarrow (8)$$

Exercise 2 ⊛ let $a_1, a_2, \ldots a_n$ be $+$ real numbers. Prove that $\sum_{j=1}^{n} \frac{a_j}{x_j}$ when $\sum x_j = 1$ attains the minimum when $x_j = \dfrac{\sqrt{a_j}}{\sum \sqrt{a_j}}$. (See Exercise 1)

Problem 3.    Low rank approximation for A is called Sketch of A, where the columns and rows each picked in ~~direct~~ individual trials based on importance distribution based on squared lengths of columns/rows is a good sketch.    Let $A \in \mathbb{R}^{m \times n}$

1) Pick $s$ columns of A and form $C|_{m \times s}$ based on $LS_{col}(A)$

2) Pick $r$ rows of A from $LS_{row}(A)$ and form $R_{s \times n}$

3) Using C and R express: $A \approx CUR$ by computing U

Note:  · You may think of using SVD to get a good sketch. While it is true, it takes a long time:   $A = U \wedge V^T$

· Besides the columns of U and rows of V are ~~combinations of those of A~~ not directly ~~related to those of A~~ from A but are linear combinations.

· In here, we directly pick the columns and row of A by sampling and is called interpolative approximation

· SVD gives best approximation, the CUR is not optimal in that sense.

How to choose $v$? $\quad A \in \mathbb{R}^{m \times n}$

1) An Intuition: $\quad A = A I$

- Sample $s$ columns of $A$ and get $C|_{m \times s}$
- Let $W$ be the corresponding $s$ rows of $I$ scaled as in Matrix-Matrix product.

  - $A = A I = C W \longrightarrow ①$

  - From (b): $\quad \mathbb{E} \|A I - C W\|^2 \le \dfrac{\|A\|_F^2 \|I\|_F^2}{s}$

  $$= \|A\|_F^2 \left(\dfrac{n}{s}\right) \longrightarrow ②$$

- Since we need $s > n$, this is not viable, but the idea stands.

2) Modification: $\quad A \in \mathbb{R}^{m \times n}$

- Consider $C|_{m \times s}$ and $R_{\wedge \times n}$ that we selected earlier, where $s < m$, $r < n$.

- If $R$ is of full-rank, then

  $R R^T \big|_{n \times n}$ is SPD

  $R^+ \big|_{n \times n} = R^T (R R^T)^{-1}$ is the Moore-Penrose inverse

  $P \big|_{n \times n} = R^T (R R^T)^{-1} R$ is the projection onto the row space of $R$.

- Recall $P$ acts as identity on the ~~row sp~~ Row space of $R$.

**Illustration** eg: $\quad R = (1, 1) \quad R R^T = (1,1)\begin{pmatrix}1\\1\end{pmatrix} = 2$

$$R^+ = R^T (R R^T)^{-1} = \frac{1}{2}\begin{pmatrix}1\\1\end{pmatrix}$$

$$P = R^T (R R^T)^{-1} R = \frac{1}{2}\begin{pmatrix}1\\1\end{pmatrix}(1,1) = \frac{1}{2}\begin{pmatrix}1&1\\1&1\end{pmatrix}$$

$\boxed{\text{Row space of } R}$ $\quad y = R^T a \Rightarrow P y = R^T (R R^T)^{-1} R R^T a$

$a \neq \begin{pmatrix}a_1\\a_2\end{pmatrix}$ $\qquad\qquad = R^T a = y$

· So we can replace $I$ by $P$ in the above analysis.

· $$AP = A \underline{R^T(RR^T)^{-1}} R \approx CUR \qquad \boxed{U = R^T(RR^T)^{-1}}$$

∴ By the inequality (6):

$$E\left[\|AP - CUR\|_2^2\right] \leq E\left[\|AP - CUR\|_F^2\right]$$

$\boxed{\|A\|_2^2 \leq \|A\|_F^2}$

$$\leq \frac{\|A\|_F^2 \, \|P\|_F^2}{\wedge} = \|A\|_F^2 \left(\frac{\wedge}{\partial}\right)$$

$$\longrightarrow (3)$$

where $\|P\|_F^2 = \wedge$. (why?)

· Recall $\|A\|_F^2 = tr(AA^T)$

$R_{n \times m}$
$R^+_{m \times n}$
$P_{m \times m}$
$(RR^T)^{-1}_{n \times n}$

$$\Rightarrow \|P\|_F^2 = tr[P] = tr[R^T(RR^T)^{-1}R]$$

$\boxed{tr(AB) = tr(BA)}$ $\quad = tr[RR^T(RR^T)^{-1}]$

$$= tr[I_n] = n \longrightarrow (4)$$

Combining:

$$A - CUR = A - AP + AP - CUR$$

Triangle inequality

$$\|A - CUR\|_2^2 \leq \|A - AP\|_2^2 + \|AP - CUR\|_2^2$$

Recall: $\qquad x \leq y + z$

$$\Rightarrow x^2 \leq (y + z)^2 = y^2 + 2yz + z^2 < 2(y^2 + z^2)$$

$$\Rightarrow 0 \leq (y^2 - 2yz + z^2)$$
$$= (y - z)^2$$

∴ $\|A - CUR\|_2^2 \leq 2\|A - AP\|_2^2 + 2\|AP - CUR\|_2^2$