

Dimensionality Reduction: CUR Decomposition

Mining of Massive Datasets
Leskovec, Rajaraman, and Ullman
Stanford University



CUR Decomposition

Frobenius norm:

$$\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2}$$

- **Goal: Express A as a product of matrices C,U,R**
Make $\|A - C \cdot U \cdot R\|_F$ small
- **“Constraints” on C and R:**

$$\begin{pmatrix} \text{red bar} & \text{blue bar} & \text{dark red bar} \end{pmatrix} \approx \begin{pmatrix} \text{red bar} & \text{red bar} & \text{red bar} & \text{blue bar} & \text{dark red bar} & \text{dark red bar} \end{pmatrix} \cdot \begin{pmatrix} U \end{pmatrix} \cdot \begin{pmatrix} R \end{pmatrix}$$

$A \qquad C \qquad U \qquad R$

CUR Decomposition

Frobenius norm:

$$\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2}$$

- **Goal: Express A as a product of matrices C,U,R**
Make $\|A - C \cdot U \cdot R\|_F$ small
- **“Constraints” on C and R:**

$$\begin{pmatrix} \text{red bar} \\ \text{brown bar} \\ \text{blue bar} \end{pmatrix} \approx \begin{pmatrix} \text{red bar} \\ \text{brown bar} \\ \text{blue bar} \end{pmatrix} \cdot \begin{pmatrix} \text{gray square} \end{pmatrix} \cdot \begin{pmatrix} \text{red bar} \\ \text{brown bar} \\ \text{blue bar} \end{pmatrix}$$

$A \qquad C \qquad U \qquad R$

Pseudo-inverse of the intersection of C and R

CUR: Provably good approx. to SVD

- **Let:**

\mathbf{A}_k be the “best” rank k approximation to \mathbf{A} (that is, \mathbf{A}_k is SVD of \mathbf{A})

Theorem [Mahoney & Drineas]

CUR in $O(m \cdot n)$ time achieves

- $\|\mathbf{A} - \mathbf{CUR}\|_F \leq \|\mathbf{A} - \mathbf{A}_k\|_F + \epsilon \|\mathbf{A}\|_F$

with probability at least $1 - \delta$, by picking

- $O(k \log(1/\delta)/\epsilon^2)$ columns, and

- $O(k^2 \log^3(1/\delta)/\epsilon^6)$ rows

In practice:
Pick $4k$ cols/rows

CUR: How it Works

- Sampling columns (similarly for rows):
- **ColumnSelect** algorithm:

Input: matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, sample size c

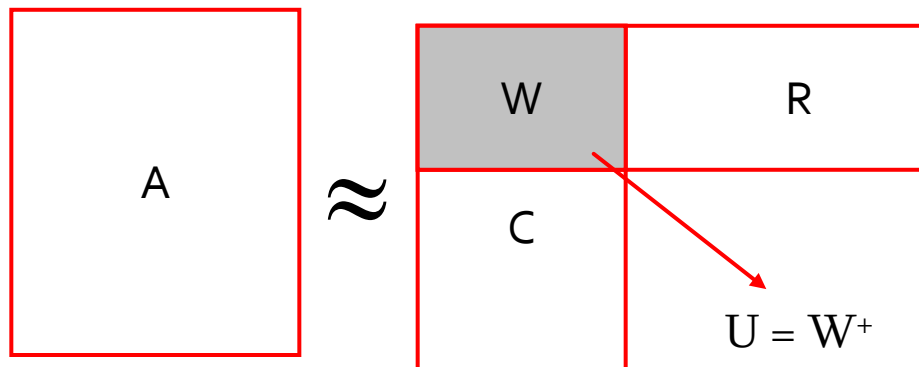
Output: $\mathbf{C}_d \in \mathbb{R}^{m \times c}$

1. for $x = 1 : n$ [column distribution]
2. $P(x) = \sum_i \mathbf{A}(i, x)^2 / \sum_{i,j} \mathbf{A}(i, j)^2$
3. for $i = 1 : c$ [sample columns]
4. Pick $j \in 1 : n$ based on distribution $P(x)$
5. Compute $\mathbf{C}_d(:, i) = \mathbf{A}(:, j) / \sqrt{cP(j)}$

Note this is a randomized algorithm, same column can be sampled more than once

Computing U

- Let \mathbf{W} be the “intersection” of sampled columns \mathbf{C} and rows \mathbf{R}
 - Let SVD of $\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}^T$
- Then: $\mathbf{U} = \mathbf{W}^+ = \mathbf{Y} \mathbf{Z}^+ \mathbf{X}^T$
 - \mathbf{Z}^+ : reciprocals of non-zero singular values: $Z_{ii}^+ = 1/Z_{ii}$
 - \mathbf{W}^+ is the “pseudoinverse”



Why pseudoinverse works?
 $\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}$ then $\mathbf{W}^{-1} = \mathbf{X}^{-1} \mathbf{Z}^{-1} \mathbf{Y}^{-1}$
Due to orthonormality
 $\mathbf{X}^{-1} = \mathbf{X}^T$ and $\mathbf{Y}^{-1} = \mathbf{Y}^T$
Since \mathbf{Z} is diagonal $\mathbf{Z}^{-1} = 1/Z_{ii}$
Thus, if \mathbf{W} is non-singular,
pseudoinverse is the true
inverse

CUR: Provably good approx. to SVD

- For example:
 - Select $c = O\left(\frac{k \log k}{\epsilon^2}\right)$ columns of A using **ColumnSelect** algorithm
 - Select $r = O\left(\frac{k \log k}{\epsilon^2}\right)$ rows of A using **ColumnSelect** algorithm
 - Set $U = W^+$
- **Then:** $\|A - CUR\|_F \leq (2 + \epsilon) \|A - A_k\|_F$
with probability 98%