

e-PGPathshala

Subject : Computer Science

Paper: Data Analytics

Module No 30: CS/DA/30 - Dimensionality

Reduction: SVD case study & CUR

Quadrant 1 – e-text

1.1 REVIEW OF THE KEY IDEAS of SVD

- The SVD factors A into $U\Sigma V^T$, with r singular values $\sigma_1 \geq \dots \geq \sigma_r > 0$.
- The numbers $\sigma_1^2, \dots, \sigma_r^2$ are the nonzero eigenvalues of AA^T and $A^T A$.
- The orthonormal columns of U and V are eigenvectors of AA^T and $A^T A$.
- Those columns hold orthonormal bases for the four fundamental subspaces of A .
- Those bases diagonalize the matrix: $Av_i = \sigma_i u_i$ for $i \leq r$. This is $AV = U\Sigma$.
 $A = \sigma_1 u_1 v_1^T + \dots + \sigma_r u_r v_r^T$ and σ_1 is the maximum of the ratio $\|Ax\| / \|x\|$.

1.2 Learning Objectives

- Understand SVD using a Case Study
- Learn the process of CUR approximation technique
- Know the performance of SVD vs CUR

1.3 Case Study:

the query is *gold silver truck* and the "collection" consists of just three "documents":

- d1: *Shipment of gold damaged in a fire.*
- d2: *Delivery of silver arrived in a silver truck.*
- d3: *Shipment of gold arrived in a truck.*

Terms	d1	d2	d3	q
↓	↓	↓	↓	↓
a	1	1	1	0
arrived	0	1	1	0
damaged	1	0	0	0
delivery	0	1	0	0
fire	1	0	0	0
gold	1	0	1	1
in	1	1	1	0
of	1	1	1	0
shipment	1	0	1	0
silver	0	2	0	1
truck	0	1	1	1

$A =$

$q =$

1.3.1 Dimensionality Reduction: Computing U_k , S_k , V_k and V_k^T

$$U = \begin{bmatrix} -0.4201 & 0.0748 & -0.0460 \\ -0.2995 & -0.2001 & 0.4078 \\ -0.1206 & 0.2749 & -0.4538 \\ -0.1576 & -0.3046 & -0.2006 \\ -0.1206 & 0.2749 & -0.4538 \\ -0.2626 & 0.3794 & 0.1547 \\ -0.4201 & 0.0748 & -0.0460 \\ -0.4201 & 0.0748 & -0.0460 \\ -0.2626 & 0.3794 & 0.1547 \\ -0.3151 & -0.6093 & -0.4013 \\ -0.2995 & -0.2001 & 0.4078 \end{bmatrix} \quad S = \begin{bmatrix} 4.0989 & 0.0000 & 0.0000 \\ 0.0000 & 2.3616 & 0.0000 \\ 0.0000 & 0.0000 & 1.2737 \end{bmatrix}$$

$$V = \begin{bmatrix} -0.4945 & 0.6492 & -0.5780 \\ -0.6458 & -0.7194 & -0.2556 \\ -0.5817 & 0.2469 & 0.7750 \end{bmatrix} \quad V^T = \begin{bmatrix} -0.4945 & -0.6458 & -0.5817 \\ 0.6492 & -0.7194 & 0.2469 \\ -0.5780 & -0.2556 & 0.7750 \end{bmatrix}$$

$$U \approx U_k = \begin{bmatrix} -0.4201 & 0.0748 \\ -0.2995 & -0.2001 \\ -0.1206 & 0.2749 \\ -0.1576 & -0.3046 \\ -0.1206 & 0.2749 \\ -0.2626 & 0.3794 \\ -0.4201 & 0.0748 \\ -0.4201 & 0.0748 \\ -0.2626 & 0.3794 \\ -0.3151 & -0.6093 \\ -0.2995 & -0.2001 \end{bmatrix} \quad k = 2 \quad S \approx S_k = \begin{bmatrix} 4.0989 & 0.0000 \\ 0.0000 & 2.3616 \end{bmatrix}$$

$$V \approx V_k = \begin{bmatrix} -0.4945 & 0.6492 \\ -0.6458 & -0.7194 \\ -0.5817 & 0.2469 \end{bmatrix} \quad V^T \approx V_k^T = \begin{bmatrix} -0.4945 & 0.6458 & -0.5817 \\ 0.6492 & -0.7194 & 0.2469 \end{bmatrix}$$

1.3.2 Incorporating the Query and Ranking the Documents:

$$d = d^T U_k S_k^{-1}$$

$$q = q^T U_k S_k^{-1}$$

$$\text{sim}(q, d) = \text{sim}(q^T U_k S_k^{-1}, d^T U_k S_k^{-1})$$

$$q = q^T U_k S_k^{-1}$$

$$q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$k = 2$$

$$\begin{bmatrix} -0.4201 & 0.0748 \\ -0.2995 & -0.2001 \\ -0.1206 & 0.2749 \\ -0.1576 & -0.3046 \\ -0.1206 & 0.2749 \\ -0.2626 & 0.3794 \\ -0.4201 & 0.0748 \\ -0.4201 & 0.0748 \\ -0.2626 & 0.3794 \\ -0.3151 & -0.6093 \\ -0.2995 & -0.2001 \end{bmatrix}$$

$$\begin{bmatrix} 1 & \\ 4.0989 & 0.0000 \\ & 1 \\ 0.0000 & 2.3616 \end{bmatrix}$$

$$q = \begin{bmatrix} -0.2140 & -0.1821 \end{bmatrix}$$

$$d1(-0.4945, 0.6492)$$

$$d2(-0.6458, -0.7194)$$

$$d3(-0.5817, 0.2469)$$

1.3.3 Cosine similarities

Cosine similarity is a measure of similarity between two non-zero vectors of an inner product space that measures the cosine of the angle between them. The cosine of 0° is 1, and it is less than 1 for any other angle.

$$\text{sim}(q, d) = \frac{q \bullet d}{\|q\| \|d\|}$$

$$\text{sim}(q, d_1) = \frac{(-0.2140)(-0.4945) + (-0.1821)(0.6492)}{\sqrt{(-0.2140)^2 + (-0.1821)^2} \sqrt{(-0.4945)^2 + (0.6492)^2}} = -0.0541$$

$$\text{sim}(q, d_2) = \frac{(-0.2140)(-0.6458) + (-0.1821)(-0.7194)}{\sqrt{(-0.2140)^2 + (-0.1821)^2} \sqrt{(-0.6458)^2 + (-0.7194)^2}} = 0.9910$$

$$\text{sim}(q, d_3) = \frac{(-0.2140)(-0.5817) + (-0.1821)(0.2469)}{\sqrt{(-0.2140)^2 + (-0.1821)^2} \sqrt{(-0.5817)^2 + (0.2469)^2}} = 0.4478$$

Ranking documents in descending order

$$d_2 > d_3 > d_1$$

1.4 CUR Decomposition

In large-data applications, it is normal for the matrix M being decomposed to be very sparse; that is, most entries are 0. For example, a matrix representing many documents (as rows) and the words they contain (as columns) will be sparse, because most words are not present in most documents. Similarly, a matrix of customers and products will be sparse because most people do not buy most products. We cannot deal with dense matrices that have millions or billions of rows and/or columns. However, with SVD, even if M is sparse, U and V will be dense.⁴ Since Σ is diagonal, it will be sparse, but Σ is usually much smaller than U and V , so its sparseness does not help.

Consider another approach to decomposition, called CUR-decomposition. The merit of this approach lies in the fact that if M is sparse, then the two large matrices (called C and R for “columns” and “rows”) analogous to U and V in SVD are also sparse. Only the matrix in the middle (analogous to Σ in SVD) is dense, but this matrix is small so the density does not hurt too much. Unlike SVD, which gives an exact decomposition as long as the parameter r is taken to be at least as great as the rank of the matrix M , CUR-decomposition is an approximation no matter how large we make r . There is a theory that guarantees convergence to M as r gets larger, but typically you have to make r so large to get, say within 1% that the method becomes impractical. Nevertheless, decomposition with a relatively small value of r has a good probability of being a useful and accurate decomposition.

Definition of CUR

Let M be a matrix of m rows and n columns. Pick a target number of “concepts” r to be used in the decomposition. A CUR-decomposition of M is a randomly chosen set of r columns of M , which form the $m \times r$ matrix C , and a randomly chosen set of r rows of M , which form the $r \times n$ matrix R . There is also an $r \times r$ matrix U that is constructed from C and R as follows:

1. Let W be the $r \times r$ matrix that is the intersection of the chosen columns of C and the chosen rows of R . That is, the element in row i and column j of W is the element of M whose column is the j^{th} column of C and whose row is the i^{th} row of R .
2. Compute the SVD of W ; say $W = X \Sigma Y^T$.
3. Compute Σ^+ , the Moore-Penrose pseudo inverse of the diagonal matrix Σ . That is, if the i^{th} diagonal element of Σ is $\sigma \neq 0$, then replace it by $1/\sigma$. But if the i^{th} element is 0, leave it as 0.
4. Let $U = Y (\Sigma^+)^2 X^T$.

Choosing rows and columns properly recall that the choice of rows and columns is random. However, this choice must be biased so that the more important rows and columns have a better chance of being picked. The measure of importance we must use is the square of the Frobenius norm, that is, the sum of the squares of the elements of the row or column. Let $f = \sum_{i,j} m_{ij}^2$, the square of the Frobenius norm of M . Then each time we select a row, the probability p_i with which we select row i is $\sum_j m_{ij}^2 / f$. Each time we select a column, the probability q_j with which we select column j is $\sum_i m_{ij}^2 / f$.

- 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 :

$$\begin{pmatrix} \left(\begin{array}{|c|c|c|} \hline \text{red} & \text{blue} & \text{brown} \\ \hline \end{array} \right) & A & \approx & \left(\begin{array}{|c|c|c|c|} \hline \text{red} & \text{red} & \text{red} & \text{red} \\ \hline \end{array} \right) & C & \cdot & \left(\begin{array}{|c|} \hline U \\ \hline \end{array} \right) & \cdot & \left(\begin{array}{|c|c|c|} \hline \text{red} & \text{brown} & \text{blue} \\ \hline \end{array} \right) & R \end{pmatrix}$$

$$\begin{pmatrix} \left(\begin{array}{|c|c|c|} \hline \text{red} & \text{brown} & \text{blue} \\ \hline \end{array} \right) & A & \approx & \left(\begin{array}{|c|} \hline C \\ \hline \end{array} \right) & \cdot & \left(\begin{array}{|c|} \hline U \\ \hline \end{array} \right) & \cdot & \left(\begin{array}{|c|c|c|} \hline \text{red} & \text{red} & \text{red} \\ \hline \end{array} \right) & R \end{pmatrix}$$

A C U R A C U R

Pseudo-inverse of the intersection of C and R

How it Works:

- Sampling columns (similarly for rows):

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

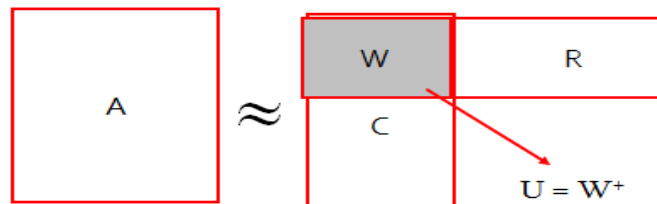
Output: $\mathbf{C}_d \in \mathbb{R}^{m \times c}$

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?

A pseudoinverse is a matrix inverse-like object that may be defined for a complex matrix, even if it is not necessarily square. For any given complex matrix, it is possible to define many possible pseudoinverses.

- $\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}$ 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 Z is diagonal $Z^{-1} = 1/Z_{ii}$
- Thus, if W is nonsingular, pseudoinverse is the true inverse

1.5 CUR DECOMPOSITION

Having selected each of the columns of M , we scale each column by dividing its elements by the square root of the expected number of times this column would be picked. That is, we divide the elements of the j th column of M , if it is selected, by $\sqrt{r_j}$. The scaled column of M becomes a column of C . Rows of M are selected for R in the analogous way. For each row of R we select from the rows of M , choosing row i with probability p_i . Recall p_i is the sum of the squares of the elements of the i th row divided by the sum of the squares of all the elements of M . We then scale each chosen row by dividing by $\sqrt{r_i}$ if it is the i th row of M that was chosen.

The Complete CUR Decomposition

We now have a method to select randomly the three component matrices C , U , and R . Their product will approximate the original matrix M . As we mentioned at the beginning of the discussion, the approximation is only formally guaranteed to be close when very large numbers of rows and columns are selected. However, the intuition is that by selecting rows and columns that tend to have high “importance” (i.e., high Frobenius norm), we are extracting

Provably good approx. to SVD

<p>■ Let:</p> <p>A_k be the “best” rank k approximation to A (that is, A_k is SVD of A)</p> <p>Theorem [Drineas et al.]</p> <p>CUR in $O(m \cdot n)$ time achieves</p> <ul style="list-style-type: none"> ■ $\ A - CUR\ _F \leq \ A - A_k\ _F + \epsilon \ A\ _F$ <p>with probability at least $1 - \delta$, by picking</p> <ul style="list-style-type: none"> ■ $O(k \log(1/\delta)/\epsilon^2)$ columns, and ■ $O(k^2 \log^3(1/\delta)/\epsilon^6)$ rows 	<p>■ For example:</p> <ul style="list-style-type: none"> ■ 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 $U = W^+$ <p>Then: $\ A - CUR\ _F \leq (2 + \epsilon) \ A - A_k\ _F$</p> <p>with probability 98%</p>
---	--

1.6 Open problems

- **Optimal CUR:** Can we find relative-error CUR algorithms selecting the optimal number of columns and rows, together with a matrix U with optimal rank?
- **Input-sparsity-time CUR:** Can we find relative-error CUR algorithms running in input-sparsity-time ($\text{nnz}(A)$ time)?
- **Deterministic CUR:** Can we find relative-error CUR algorithms that are deterministic and run in poly time?

Contributions

- Optimal CUR: First optimal CUR algorithms.

- Input-sparsity-time CUR: First CUR algorithm with running time proportional to the non-zero entries of A.
- Deterministic CUR: First deterministic algorithm for CUR that runs in polynomial time.

Pros & Cons

▪ Easy interpretation

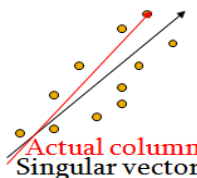
- Since the basis vectors are actual columns and rows

▪ Sparse basis

- Since the basis vectors are actual columns and rows

▪ Duplicate columns and rows

- Columns of large norms will be sampled many times



Solution

▪ If we want to get rid of the duplicates:

- Scale (multiply) the columns/rows by the square root of the number of duplicates



Three CUR Data Applications

- **Human Genetics:** DNA SNP Data • Biological Goal: Evaluate intra- and inter-population tag-SNP transferability.
- **Medical Imaging:** Hyper spectral Image Data • Medical Goal: Compress the data, without sacrificing classification quality.
- **Recommendation Systems:** Customer Preference Data • Business Goal: Reconstruct the data, to make high-quality recommendations

1.7 SVD vs. CUR

SVD has a lesser reconstruction error than CUR but the time taken and space occupied by SVD is far more than that of CUR. So, we can say that both the techniques are useful but the choice of the technique that is deemed fit for a scenario depends on the characteristics of the data set. For a larger dataset, it is better to go ahead with CUR so as to reduce the space and time overhead. But in the case where every recommendation should be accurate, it's better to use SVD since the reconstruction error is minimum

$$\text{SVD: } A = U \Sigma V^T$$

Diagram illustrating the SVD decomposition $A = U \Sigma V^T$ with annotations:

- A : Huge but sparse
- U : Big and dense
- Σ : sparse and small
- V^T : Big and dense

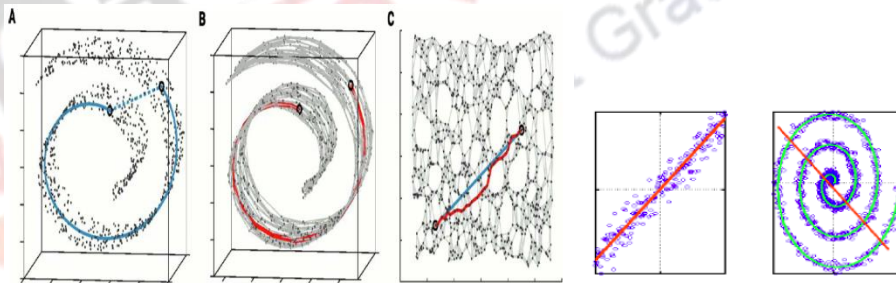
$$\text{CUR: } A = C U R$$

Diagram illustrating the CUR decomposition $A = C U R$ with annotations:

- A : Huge but sparse
- C : Big but sparse
- U : dense but small
- R : Big but sparse

What about linearity assumption?

- **SVD is limited to linear projections:**
 - Lower-dimensional linear projection that preserves Euclidean distances
- Non-linear methods: **Isomap**
 - Data lies on a nonlinear low-dim curve aka manifold
 - Use the distance as measured along the manifold
- **How?**
 - Build adjacency graph
 - Geodesic distance is graph distance
 - SVD/PCA the graph pair wise distance matrix



Summary:

- Dimensionality reduction involves mathematical approximations
- Dimensionality reductions techniques have their own advantages and drawbacks
- process of CUR approximation technique
- performance of SVD vs CUR