

Machine Learning Foundations

Week 6

Shubharup G.
BSDS, IIT Madras

March 2024

Contents

1	Introduction	2
1.1	Recap: Eigendecomposition (aka Diagonalization)	2
1.2	Uses of Singular Value Decomposition	2
2	Singular Value Decomposition	3
2.1	Form and Properties of SVD	3
2.2	Computing SVD	5
2.2.1	First method – directly, as defined	5
2.2.2	Second method – using inter-relationships (faster)	7
3	SVD for Compression	8
3.1	Dyadic product (aka Outer product)	8
3.2	SVD as sum of dyads	9
3.3	Truncated k-sum as the best k-rank approximation	9
4	Interpreting a Linear Transformation with SVD	
	References	

1 Introduction

1.1 Recap: Eigendecomposition (aka Diagonalization)

So far we have learnt eigendecomposition, which the standard sense of the term ‘diagonalization’.

Eigendecomposition is one of the many known matrix factorization techniques ¹.

The two major motivations, by way of utility, which we’ve seen for diagonalization (eigendecomposition) are:

1. Interpretability: what is happening when we apply the matrix? decomposes what looks like an atomic task into 3 discrete subtasks:

I change of basis aka rotation

II disproportional scaling

III rotation back to original basis

2. Economy of computation: reducing the compute of applying the same matrix repeatedly like $A(A(x)) : A^n(x) = P^{-1}\Lambda^n P$.

Raising a diagonal matrix to a power n is really easy:
raise the diagonal entries to the power n .

e.g. for a numerical example, if $A = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$, can immediately read off $A^4 = \begin{bmatrix} 16 & 0 \\ 0 & 81 \end{bmatrix}$

1.2 Uses of Singular Value Decomposition

Interpretability is an overlapping utility for SVD, but it’s even better because any real matrix has an SVD: whether square or rectangular, whether with or without an eigenbasis.

SVD factorizes the application of the matrix into three discrete steps or subtasks, much like diagonalization.

SVD might not have use in reducing compute.

But SVD, because of its generality of application to any real matrix, has a prototypical use-case for which it is particularly well-known: loss-minimizing compression. By compression, SVD can reduce the storage requirement for archiving a dataset, while retaining the most important information.

In fact, the Eckhart-Young theorem², also simply known as ‘the SVD theorem’, guarantees that the low-rank approximation gotten using SVD is the best approximation – the approximated matrix has the least loss of all possible choices for approximation, as measured by the Frobenius norm³ of the difference between the original matrix and the approximated matrix.

¹Matrix decomposition (Wikipedia)

²Low-rank approximation (Wikipedia)

³Sum of squares of the entries in the matrix – essentially similar to the Euclidean or L_2 norm for vectors. See Frobenius norm (Wikipedia).

2 Singular Value Decomposition

2.1 Form and Properties of SVD

SVD is a matrix factorization technique ⁴ canonically written as:

$$\underbrace{A}_{m \times n} = \underbrace{U}_{m \times m} \underbrace{\Sigma}_{m \times n} \underbrace{V^T}_{n \times n}$$

- $U = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_m \\ | & | & & | \end{bmatrix}$.

The column-vectors of U viz. u_j are known as the **left singular vectors**.

- $\Sigma = \begin{cases} \left[\begin{array}{cccc|cccc} \sigma_1 & & & & & & & \\ & \sigma_2 & & & & & & \\ & & \sigma_3 & & & & & \\ & & & \ddots & & & & \\ & & & & \sigma_\omega & & & \\ & & & & & 0 & \dots & 0 \end{array} \right] & \text{if } m < n; \\ \left[\begin{array}{cccc} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \sigma_3 & \\ & & & \ddots \\ & & & & \sigma_\omega \\ \hline & 0 & & & \\ & \vdots & & & \\ & 0 & & & \end{array} \right] & \text{if } m > n; \end{cases}$

Σ is in general a **rectangular diagonal matrix**: there is a padding of zero columns if $m < n$, a padding of zero rows if $m > n$. It is a square matrix iff. $m = n$.

The diagonal entries in the square diagonal partition viz. the σ_j are called the **singular values**. The count of singular values is $\omega = \min(m, n)$.

- $V = \begin{bmatrix} | & | & & | \\ v_1 & v_2 & \dots & v_n \\ | & | & & | \end{bmatrix}$, or $V^T = \begin{bmatrix} - & v_1^T & - \\ - & v_2^T & - \\ & \vdots & \\ - & v_n^T & - \end{bmatrix}$.

The column vectors of V viz. v_j are known as the **right singular vectors**.

⁴One among others, to restate.

Properties:

- U and V are orthogonal matrices.

- $UU^T = U^T U = I$

- $VV^T = V^T V = I$

To recap:

- Orthogonal matrices have orthonormal columns.

(Equivalently, they have orthonormal rows.)

- The inverse of an orthogonal matrix P is its transpose: $P^{-1} = P^T$.

- An orthogonal matrix P has determinant $\det(P) = 1$.

This is important for relevance here, because **every orthogonal matrix is a change of basis matrix** between some pair of basis frames. This aspect will be revisited when we see SVD of a matrix A as a logical decomposition of the linear transformation represented by A – from what appears to be an atomic task, into the sequential composition of three sub-tasks.

- Similar to eigendecomposition, the order of $\begin{bmatrix} | \\ u_j \\ | \end{bmatrix}$ and $\begin{bmatrix} | \\ v_j \\ | \end{bmatrix}$, in U and V respectively, is coupled with the relative order of σ_j in Σ .

Therefore,

SVD of a matrix A is unique *up to* rearrangement or *permutation* of the singular values in Σ .

2.2 Computing SVD

2.2.1 First method – directly, as defined

- The left singular vectors viz. the column vectors $\begin{bmatrix} | \\ u_i \\ | \end{bmatrix}$ in U are the *normalized* eigenvectors of AA^T .

- For $\text{shape}(A) = m \times n$, $\text{shape}(AA^T) = m \times m$.
- AA^T is a (real) symmetric matrix. By virtue of the Spectral Theorem, this guarantees that it has an orthogonal eigenbasis i.e. *a complete set of m orthogonal eigenvectors*, which can be normalized to be made orthonormal. These go into U as the column vectors.

- The right singular vectors viz. the column vectors $\begin{bmatrix} | \\ v_i \\ | \end{bmatrix}$ in V or likewise the row vectors $[-v_i^T -]$ in V^T are the *normalized* eigenvectors of $A^T A$.

- For $\text{shape}(A) = m \times n$, $\text{shape}(A^T A) = n \times n$.
- $A^T A$ is a (real) symmetric matrix. By virtue of the Spectral Theorem, this guarantees that it has an orthogonal eigenbasis i.e. *a complete set of n orthogonal eigenvectors*, which can be normalized to be made orthonormal. These go into V as the column vectors.

- The singular values $\sigma_i = \sqrt[n]{\lambda_i}$ where:

$$\{\lambda_i\} = \begin{cases} \Lambda(AA^T) & \text{if } m < n; \\ \Lambda(A^T A) & \text{if } m > n \end{cases}$$

where $\Lambda(M)$ denotes the multiset⁵ of eigenvalues of the matrix M , aka its *spectrum*.⁶

⁵Extension of the idea of a set to allow for multiplicity of the same element: here, viz. the include repeated eigenvalues. See Multiset (Wikipedia).

⁶Spectrum of a matrix (Wikipedia)

- Can show that for a square matrix M , $\Lambda(M) = \{\hat{v}_j^T M \hat{v}_j\}_j$ for all linearly independent and normalized eigenvectors \hat{v}_j of M .

So also: $\Lambda(M^T M) = \{\hat{w}_j^T (M^T M) \hat{w}_j\}_j$ where \hat{w}_j are the linearly independent and normalized eigenvectors of $M^T M$.

$$\begin{aligned}\Lambda(M^T M) &= \{\hat{w}_j^T (M^T M) \hat{w}_j\}_j \\ &= \{(M \hat{w}_j)^T (M \hat{w}_j)\}_j \\ &= \{y^T y \text{ or } \|y\|^2, \text{ where } y = M \hat{w}_j\}_j\end{aligned}$$

We know that norm of a vector, by properties, is non-negative. This means the spectrum $\Lambda(M^T M)$ has only non-negative entries⁷.

- * Can show the same for MM^T also: take $C = M^T$ and so $MM^T = C^T C$ is positive semi-definite.

Therefore, the singular values (positive square roots of non-negative real numbers) are guaranteed to be real.

- The singular values are all positive (by definition), and conventionally placed in descending order of magnitude. The reason for the latter will become apparent in both later sections on SVD for Compression, and SVD for Interpretation.
- It is a known result that **the set of non-zero eigenvalues of AA^T and $A^T A$ are exactly the same:**

$\underbrace{(\lambda_i \neq 0, r) \in \Lambda(AA^T)}_{\substack{\lambda_i \text{ is a non-zero eigenvalue of } AA^T \\ \text{with multiplicity } r}} \iff \underbrace{(\lambda_i \neq 0, r) \in \Lambda(A^T A)}_{\substack{\lambda_i \text{ is a non-zero eigenvalue of } A^T A \\ \text{with multiplicity } r}}$
--

The ‘extra’ eigenvalues will all be zeroes:

$$\left\{ \begin{array}{ll} (n - m) \text{ in } \underbrace{\Lambda(\overbrace{A^T A}^{n \times n}) - \Lambda(\overbrace{AA^T}^{m \times m})}_{\Lambda(AA^T) \subset \Lambda(A^T A)} & \text{if } n > m \\ (m - n) \text{ in } \underbrace{\Lambda(\overbrace{AA^T}^{m \times m}) - \Lambda(\overbrace{A^T A}^{n \times n})}_{\Lambda(A^T A) \subset \Lambda(AA^T)} & \text{if } m > n \end{array} \right.$$

⁷i.e. $M^T M$ is *positive semi-definite*.

2.2.2 Second method – using inter-relationships (faster)

Let $\omega = \min(m, n)$.

Then the left-singular vectors $\begin{bmatrix} | \\ u_j \\ | \end{bmatrix}$
 and the right-singular vectors $\begin{bmatrix} | \\ v_j \\ | \end{bmatrix}$
 are related as:

$$\left. \begin{aligned} v_j &= \frac{A^T u_j}{\sigma_j} \quad \dots \quad (i) \\ u_j &= \frac{A v_j}{\sigma_j} \quad \dots \quad (ii) \end{aligned} \right\} j \in \{1, 2 \dots \omega\}$$

The proof is reserved for the appendix.

So then,

$$\left\{ \begin{array}{ll} \text{If } \underbrace{n > m}_{\text{'broad' matrix}} : & \begin{aligned} &\bullet \text{ Compute the eigenvalues and normalized eigenvectors of } AA^T (\equiv m \times m) \\ &\quad - \text{ respectively providing the singular values } \{\sigma_j\} \text{ and left singular vectors } \{u_j\}. \\ &\bullet \text{ Use (i) to get the corresponding right singular vectors } \{v_j\}. \\ &\bullet \text{ Extend } V \text{ to an an orthonormal basis of } \mathbb{R}^n \text{ by the Gram-Schmidt process.} \end{aligned} \\ \\ \text{If } \underbrace{m > n}_{\text{'tall' matrix}} : & \begin{aligned} &\text{Compute the eigenvalues and normalized eigenvectors of } A^T A (\equiv n \times n) \\ &\quad - \text{ respectively providing the singular values } \{\sigma_j\} \text{ and right singular vectors } \{v_j\}. \\ &\bullet \text{ Use (ii) to get the corresponding left singular vectors } \{u_j\}. \\ &\bullet \text{ Extend } U \text{ to an an orthonormal basis of } \mathbb{R}^m \text{ by the Gram-Schmidt process.} \end{aligned} \end{array} \right.$$

3 SVD for Compression

3.1 Dyadic product (aka Outer product)

Thus far, we have learnt about the *inner product* of two vectors in a vector space V .

The inner of product in a vector space V is a bivariate function (i.e. with two arguments) which maps a pair of vectors $a, b \in V$ to a member item of the field underlying V , additionally satisfying some pre-requisite properties.

Prior to Week 5, our vector spaces were always over the field of real numbers, $V \subseteq \mathbb{R}^n$. In Week 5, we expanded the underlying field to the set of the set of complex numbers, $V \subseteq \mathbb{C}^n$ and looked at some interesting generalizations of entities in the real case (e.g. Hermitian matrices for real symmetric matrices, unitary matrices for real orthogonal matrices, etc.).

In Week 6 and onwards, we will again assign ourselves to only working with set of reals \mathbb{R} as the underlying field. Further restatement of this will be excluded for the sake of brevity and avoiding repetition.

So an inner product is a function:

$$\langle, \rangle : V \times V \rightarrow \mathbb{R}$$

For example, the *standard inner product*:

$$\langle a, b \rangle = a^T b$$

Importantly, the inner product outputs a scalar.

One among several vector multiplication operations⁸ with accepted usage, relevant for us here is the **outer product**:

$$u \otimes v = uv^T$$

- If $u \in \mathbb{R}^m$ and $v \in \mathbb{R}^n$, then the (result of ⁹) **outer product is a matrix** $\mathbb{M}_{m \times n}(\mathbb{R})$.
 - Contrast this with the inner product of two vectors, which outputs a scalar. Although, there is the possibly meaningful connection that if both $u, v \in \mathbb{R}^n$, then the inner product is the trace of the outer product: $\langle u, v \rangle = \text{tr}(u \otimes v)$.
- Such a matrix is also called a **dyadic matrix** or simply a **dyad**, following that the outer product is also known as the *dyadic product*.

For non-zero vectors u, v , the outer product is a dyad of rank = 1.

To see this:

- The i^{th} row of uv^T is v^T scaled by the i^{th} entry of \vec{u} .

$$\text{– if } \vec{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}, \text{ then } uv^T = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \begin{bmatrix} - & v^T & - \end{bmatrix} = \begin{bmatrix} - & (u_1)v^T & - \\ - & (u_2)v^T & - \\ & \vdots & \\ - & (u_n)v^T & - \end{bmatrix}$$

⁸Vector multiplication (Wikipedia)

⁹The term 'outer product' can variously refer to the operation (\otimes) itself, or what the operation produces. The sense is usually apparent enough from the context.

- So the rows in uv^T are all just scalar multiples of the same vector viz. v^T , so (for $\vec{u}, \vec{v} \neq \vec{0}$) the $\text{rank}(uv^T) = 1$.
- Then for a non-zero scalar $c \in \mathbb{R}$, $c(u \otimes v) = c(uv^T)$ is also of rank = 1.

Examples:

1. Recall the projection matrix \mathbb{P}_a for a (non-zero) vector \vec{a} , $\mathbb{P}_a = \frac{aa^T}{a^T a}$. It is indeed a dyadic matrix, with $\text{rank}(\mathbb{P}_a) = 1$.

3.2 SVD as sum of dyads

You can verify that the SVD factorization of a matrix:

$$\begin{array}{c}
 \left[\begin{array}{c|c|c|c} u_1 & u_2 & \dots & u_m \end{array} \right] \\
 \text{\scriptsize } \max(m-n, 0) \text{ rows}
 \end{array}
 \left\{ \begin{array}{c}
 \left[\begin{array}{cccc}
 \sigma_1 & & & \\
 & \sigma_2 & & \\
 & & \sigma_3 & \\
 & & & \ddots \\
 & & & & \sigma_\omega
 \end{array} \right]
 \begin{array}{c}
 \overbrace{\left[\begin{array}{c|c|c} \hline \hline \hline \\ \hline \hline \hline \end{array} \right]}^{\max(n-m, 0) \text{ columns}} \\
 \mathbf{0} \quad \dots \quad \mathbf{0} \\
 \hline \hline \hline
 \end{array}
 \right]
 \begin{array}{c}
 \left[\begin{array}{c} -v_1^T - \\ -v_2^T - \\ \vdots \\ -v_n^T - \end{array} \right]
 \end{array}
 \end{array}
 \right.$$

can be multiplied and produced out as the sum:

$$\begin{aligned}
 A &= U\Sigma V^T \\
 &= \sum_{j=1}^{\omega} \sigma_j u_j v_j^T \\
 &= \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T \dots + \sigma_\omega u_\omega v_\omega^T
 \end{aligned}$$

That is, A is *exactly* the sum of a finite number of such dyads, whose count is $\leq \omega$ in general (as it might be that some of the singular values $\sigma_j = 0$).

3.3 Truncated k-sum as the best k-rank approximation

We had said earlier that the singular values σ_j should be listed in descending order of magnitude in Σ .

This is because it directly ties with the following utility, stated here as a fact:

With the singular values in Σ listed as $\sigma_n \geq \sigma_{n+1} \dots$

the best^a k -rank approximation of A
is given by:

$$\begin{aligned}\tilde{A}_k &= \sum_{j=1}^k \sigma_j u_j v_j^T \\ &= \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T \dots \sigma_k u_k v_k^T\end{aligned}$$

^aBy the Eckart-Theorem. See Low-rank approximation (Wikipedia)

- A note on the *slight* ambiguity of terminology here.

Here, k -rank does not mean rank *exactly* equal to k , but **rank at most k** i.e. $\text{rank}(\tilde{A}_k) \leq k$.

- A property of $\text{rank}(\cdot)$ is that **rank is sub-additive**:

$$\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$$

and in general: $\text{rank}(\sum_{j=1}^S M_j) \leq \sum_{j=1}^S \text{rank}(M_j)$.

- * In particular here,

$$\begin{aligned}\text{rank}(\tilde{A}_k) &= \text{rank}\left(\sum_{j=1}^k \sigma_j u_j v_j^T\right) \\ &\leq \sum_{j=1}^k \text{rank}(\sigma_j u_j v_j^T) \\ &= \underbrace{1 + 1 + \dots + 1}_{k \text{ times}} \\ &= k \\ \therefore \text{rank}(\tilde{A}_k) &\leq k\end{aligned}$$

- So if $k = \omega$ i.e. ‘no approximation’:

$$\text{rank}(A) \leq \omega$$

$$\text{i.e. } \text{rank}(A) \leq \min(m, n)$$

4 Interpreting a Linear Transformation with SVD

To be added.

References

Wikipedia

Wikipedia contributors. *Spectrum of a matrix* (Wikipedia). URL: https://en.wikipedia.org/w/index.php?title=Spectrum_of_a_matrix.

Wikipedia contributors. *Low-rank approximation* (Wikipedia). URL: https://en.wikipedia.org/wiki/Low-rank_approximation#Basic_low-rank_approximation_problem.

Wikipedia contributors. *Matrix decomposition* (Wikipedia). URL: https://en.wikipedia.org/w/index.php?title=Matrix_decomposition.

Wikipedia contributors. *Frobenius norm* (Wikipedia). URL: https://en.wikipedia.org/wiki/Matrix_norm#Frobenius_norm.

Wikipedia contributors. *Multiset* (Wikipedia). URL: <https://en.wikipedia.org/w/index.php?title=Multiset&oldid=1201923950>. [Online; accessed 19-March-2024].

Wikipedia contributors. *Vector multiplication* (Wikipedia). URL: https://en.wikipedia.org/w/index.php?title=Vector_multiplication.