

Mathematics for Machine Learning

Prof Willie Brink

Applied Mathematics, Stellenbosch University

Lecture 3: Matrix Decompositions

Contents of the module

Chapter 02: Linear Algebra

Chapter 03: Analytic Geometry

Chapter 4: Matrix Decompositions

Chapter 05: Vector Calculus

Chapter 06: Probability and Distributions

Chapter 07: Continuous Optimisation

Chapter 08: When Models Meet Data

Chapter 09: Linear Regression

Chapter 10: Dimensionality Reduction with Principal Component Analysis

Chapter 11: Density Estimation with Gaussian Mixture Models

Chapter 12: Classification with Support Vector Machines

4.1 Determinant and trace

The **determinant** of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a real number $\det(\mathbf{A}) = |\mathbf{A}|$ related to the existence of an inverse: \mathbf{A} is invertible if and only if $\det(\mathbf{A}) \neq 0$.

If $\mathbf{A} \in \mathbb{R}^{2 \times 2}$, $\det(\mathbf{A}) = a_{11}a_{22} - a_{12}a_{21}$

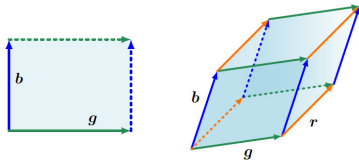
If $\mathbf{A} \in \mathbb{R}^{3 \times 3}$, $\det(\mathbf{A}) = a_{11}a_{22}a_{33} + a_{21}a_{32}a_{13} + a_{31}a_{12}a_{23}$
 $- a_{31}a_{22}a_{13} - a_{11}a_{32}a_{23} - a_{21}a_{12}a_{33}$

Sarrus' rule

If $\mathbf{T} \in \mathbb{R}^{n \times n}$ is **upper-triangular** ($t_{i,j} = 0$ for $i > j$) or **lower-triangular** ($t_{i,j} = 0$, $i < j$),

$$\det(\mathbf{T}) = \prod_{i=1}^n t_{i,i}$$

$\det(\mathbf{A})$ is the signed volume of an n -dimensional parallelepiped formed by columns of \mathbf{A} .



Laplace expansion allows us to compute the determinant of an $n \times n$ matrix in terms of the determinant of an $(n - 1) \times (n - 1)$ matrix.

$$\text{Expansion along column } j: \det(\mathbf{A}) = \sum_{k=1}^n (-1)^{k+j} a_{k,j} \det(\mathbf{A}_{k,j})$$

where $\mathbf{A}_{k,j}$ is \mathbf{A} with row k and column j deleted. Expansion along a row is similar.

$\det(\mathbf{AB}) = \det(\mathbf{A})\det(\mathbf{B})$. $\det(\mathbf{A}^T) = \det(\mathbf{A})$. If \mathbf{A} is invertible, $\det(\mathbf{A}^{-1}) = 1/\det(\mathbf{A})$.

Multiplication of a row/col by $\lambda \in \mathbb{R}$ scales $\det(\mathbf{A})$ by λ , hence $\det(\lambda\mathbf{A}) = \lambda^n \det(\mathbf{A})$.

A square matrix \mathbf{A} has $\det(\mathbf{A}) \neq 0$ if and only if $\text{rk}(\mathbf{A}) = n$.

That is to say, \mathbf{A} is invertible if and only if it is full rank.

The **trace** of a square matrix \mathbf{A} , $\text{tr}(\mathbf{A})$, is the sum of the diagonal elements of \mathbf{A} .

Trace is invariant under cyclic permutations of factors: $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA})$

4.2 Eigenvalues and eigenvectors

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$. Then $\lambda \in \mathbb{R}$ is an **eigenvalue** of \mathbf{A} , with corresponding **eigenvector** $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$, if $\mathbf{Ax} = \lambda\mathbf{x}$.

Note: if \mathbf{x} is an eigenvector of \mathbf{A} , then so is $c\mathbf{x}$ for any $c \in \mathbb{R} \setminus \{0\}$.

Eigenvalues are the roots of the **characteristic polynomial** of \mathbf{A} : $p_{\mathbf{A}}(\lambda) = \det(\mathbf{A} - \lambda\mathbf{I}_n)$
Every eigenvalue has an algebraic multiplicity.

All eigenvectors associated with an eigenvalue λ forms the **eigenspace** of \mathbf{A} w.r.t. λ .
It is the solution space of the system $(\mathbf{A} - \lambda\mathbf{I}_n)\mathbf{x} = \mathbf{0}$, i.e. the null space of $\mathbf{A} - \lambda\mathbf{I}_n$.
Its dimension is called the geometric multiplicity of λ .

\mathbf{A} and \mathbf{A}^T have the same eigenvalues, but not necessarily the same eigenvectors.

Symmetric, positive definite matrices always have positive, real eigenvalues.

Geometrically, the eigenvector corresponding to a nonzero eigenvalue points in a direction that is stretched by the linear mapping \mathbf{A} . The eigenvalue is the factor by which it is stretched (can be negative).

$\mathbf{A} \in \mathbb{R}^{n \times n}$ is **defective** if it possesses fewer than n linearly independent eigenvectors. The eigenvectors of a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with n distinct eigenvalues are linearly independent.

From a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ we can always form a symmetric, positive semidefinite matrix $\mathbf{S} \in \mathbb{R}^{n \times n}$ as $\mathbf{S} = \mathbf{A}^T \mathbf{A}$. If $\text{rk}(\mathbf{A}) = n$, \mathbf{S} will be symmetric, positive definite.

Spectral theorem: if $\mathbf{A}^{n \times n}$ is symmetric, there exists an ONB of the vector space consisting of eigenvectors of \mathbf{A} , and each eigenvalue is real.

For any matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with (possibly repeating) eigenvalues λ_i ,

$$\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i \quad \text{and} \quad \text{tr}(\mathbf{A}) = \sum_{i=1}^n \lambda_i$$

Google's PageRank algorithm

Importance of a web page is defined by the importance of the pages that link to it.

Express the web as a huge directed graph of which pages linking to which.

PageRank will compute an importance $x_i \geq 0$ for each page i .

Count the number of web pages pointing to i and model a user's navigation by a transition matrix \mathbf{A} , with columns summing to 1 and $a_{i,j}$ the probability of navigating from page i to page j .

\mathbf{A} has the property that $\mathbf{A}\mathbf{x}, \mathbf{A}^2\mathbf{x}, \mathbf{A}^3\mathbf{x}, \dots$ converges to vector \mathbf{x}^* . It satisfies $\mathbf{A}\mathbf{x}^* = \mathbf{x}^*$, that is, \mathbf{x}^* is an eigenvector of \mathbf{A} corresponding to eigenvalue 1.

Normalising \mathbf{x}^* (such that $\|\mathbf{x}^*\| = 1$) gives the PageRank of all pages as probabilities.

4.3 Cholesky decomposition

A symmetric, positive definite matrix \mathbf{A} can be factorised uniquely as $\mathbf{A} = \mathbf{L}\mathbf{L}^T$, where \mathbf{L} is lower-triangular with positive diagonal elements.

Various algorithms for computing \mathbf{L} , including a modification of Gaussian elimination.

Note that $\det(\mathbf{A}) = \det(\mathbf{L})\det(\mathbf{L}^T) = \det(\mathbf{L})^2$.

Since \mathbf{L} is lower-triangular, $\det(\mathbf{A})$ is the square of the product of \mathbf{L} 's diagonal elements.

Applications in machine learning:

- Cholesky decomposition of a covariate matrix allows us to generate samples from a multivariate Gaussian
- used in deep stochastic models (e.g. VAEs) to compute gradients

4.4 Eigendecomposition and diagonalisation

A **diagonal matrix** $\mathbf{D} \in \mathbb{R}^{n \times n}$ has zero on all off-diagonal elements.

- $\det(\mathbf{D})$ is the product of the diagonal elements;
- \mathbf{D}^k is given by each diagonal element to the power k ;
- \mathbf{D}^{-1} is the reciprocals of the diagonal elements if they are all nonzero.

Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is **diagonalisable** if there exists an invertible matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ such that $\mathbf{D} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}$ is diagonal.

Note that if \mathbf{D} has the eigenvalues of \mathbf{A} on its diagonal, and \mathbf{P} the corresponding eigenvectors of \mathbf{A} as columns, then $\mathbf{A}\mathbf{P} = \mathbf{P}\mathbf{D}$. So for \mathbf{A} to be diagonalisable, it must have n linearly independent eigenvectors (so that the inverse of \mathbf{P} exists).

From the spectral theorem we have that every symmetric matrix is diagonalisable.

The **eigendecomposition** of \mathbf{A} is $\mathbf{A} = \mathbf{PDP}^{-1}$, where \mathbf{D} is diagonal with the eigenvalues of \mathbf{A} on its diagonal and \mathbf{P} the corresponding eigenvectors of \mathbf{A} as its columns.

If \mathbf{A} is symmetric, \mathbf{P} will be orthogonal so that $\mathbf{A} = \mathbf{PDP}^T$.

Geometrically, transformations with \mathbf{A} would be the same as:

1. performing a basis change from the standard basis to the eigenbasis (\mathbf{P}^{-1})
2. scaling along those axes by the eigenvalues (\mathbf{D})
3. transforming back into the standard coordinates (\mathbf{P})

If it exists, the eigendecomposition allows for efficient computation of matrix powers and the determinant.

But this decomposition requires the matrix \mathbf{A} to be square...

4.5 Singular value decomposition

The SVD of a rectangular matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank $r \leq \min\{m, n\}$, is of the form

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

where $\mathbf{U} \in \mathbb{R}^{m \times m}$ is orthogonal, with the left-singular vectors of \mathbf{A} as columns (\mathbf{u}_i)
 $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ contains the singular values of \mathbf{A} on the diagonal and zeros elsewhere
 $\mathbf{V} \in \mathbb{R}^{n \times n}$ is orthogonal, with the right-singular vectors of \mathbf{A} as columns (\mathbf{v}_j)

The singular values are non-negative, and by convention in non-increasing order:

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m,n\}} \geq 0$$

The SVD exists for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$.

Geometric intuition: a basis change via \mathbf{V}^T , followed by a scaling and augmentation (or reduction) in dimensionality via $\mathbf{\Sigma}$, and then a second basis change via \mathbf{U} .

Construction of the SVD

From any $\mathbf{A} \in \mathbb{R}^{m \times n}$ we can construct a symmetric, positive definite matrix $\mathbf{A}^T \mathbf{A}$ with eigendecomposition: $\mathbf{A}^T \mathbf{A} = \mathbf{P} \mathbf{D} \mathbf{P}^T$ (*)

Assuming \mathbf{A} can be written in the form $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$,

$$\mathbf{A}^T \mathbf{A} = (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T \text{ since } \mathbf{U}^T \mathbf{U} = \mathbf{I}_m$$

Compare with (*): $\mathbf{V} = \mathbf{P}$ and $\mathbf{\Sigma}^T \mathbf{\Sigma} = \mathbf{D}$

The diagonal elements of $\mathbf{\Sigma}$ are the positive square roots of the eigenvalues of $\mathbf{A}^T \mathbf{A}$.

The columns of \mathbf{V} are the eigenvectors of $\mathbf{A}^T \mathbf{A}$ (ordered appropriately).

Similarly, from the eigendecomposition of the symmetric, positive definite matrix $\mathbf{A} \mathbf{A}^T$ we find that the columns of \mathbf{U} are the eigenvectors of $\mathbf{A} \mathbf{A}^T$.

Consider again an $m \times n$ matrix \mathbf{A} of rank r . Because of the many zeros in $\mathbf{\Sigma}$, some columns of \mathbf{U} or rows in \mathbf{V}^T may be redundant (in certain applications).

If $r < \min\{m, n\}$, even more columns and rows can be removed.

The **reduced SVD** is $\mathbf{A} = \mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}_r^T$, where

\mathbf{U}_r is an $m \times r$ matrix consisting of the first r columns of \mathbf{U} ,

$\mathbf{\Sigma}_r$ is an $r \times r$ diagonal matrix with $\sigma_1, \dots, \sigma_r$ on the diagonal,

and \mathbf{V}_r is an $n \times r$ matrix consisting of the first r columns of \mathbf{V} .

Applications of the SVD in machine learning:

- solving general linear systems, also in the least-squares sense
- low-rank matrix approximation for dimensionality reduction, topic modelling, data compression, clustering

Finding structure in movie ratings

n viewers rate m movies out of 5. We encode this in a data matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, and SVD.

The columns \mathbf{u}_i of \mathbf{U} are stereotypical movies, and the columns \mathbf{v}_j of \mathbf{V} stereotypical viewers.

- a vector in $\text{span}(\mathbf{v}_1, \dots, \mathbf{v}_n)$ might be a particular viewer's preferences
- a vector in $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_m)$ might be a particular movie's likeability

$$\begin{array}{c}
 \text{Star Wars} \\
 \text{Blade Runner} \\
 \text{Amelie} \\
 \text{Delicatessen}
 \end{array}
 \begin{array}{c}
 \mathbf{A} \\
 \mathbf{B} \\
 \mathbf{C}
 \end{array}
 \begin{array}{c}
 \left[\begin{array}{ccc}
 5 & 4 & 1 \\
 5 & 5 & 0 \\
 0 & 0 & 5 \\
 1 & 0 & 4
 \end{array} \right]
 \end{array}
 =
 \begin{array}{c}
 \left[\begin{array}{cccc}
 -0.6710 & 0.0236 & 0.4647 & -0.5774 \\
 -0.7197 & 0.2054 & -0.4759 & 0.4619 \\
 -0.0939 & -0.7705 & -0.5268 & -0.3464 \\
 -0.1515 & -0.6030 & 0.5293 & -0.5774
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 \left[\begin{array}{ccc}
 9.6438 & 0 & 0 \\
 0 & 6.3639 & 0 \\
 0 & 0 & 0.7056 \\
 0 & 0 & 0
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 \left[\begin{array}{ccc}
 -0.7367 & -0.6515 & -0.1811 \\
 0.0852 & 0.1762 & -0.9807 \\
 0.6708 & -0.7379 & -0.0743
 \end{array} \right]
 \end{array}$$

\mathbf{u}_1 has large values for the first two movies, grouping a type of user with a specific set of movies (sci-fi).

\mathbf{v}_1 shows large values for users A and B, suggesting the notion of a science fiction lover.

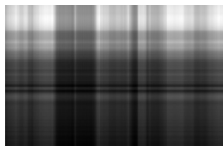
\mathbf{u}_2 captures an art-house theme, and \mathbf{v}_2 indicates that C is close to an idealised lover of such movies.

4.6 Matrix approximation

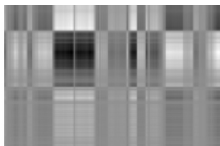
The reduced SVD can be expressed as $\mathbf{A} = \mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \sum_{i=1}^r \sigma_i \mathbf{A}_i$
where $\mathbf{A}_i = \mathbf{u}_i \mathbf{v}_i^T$ is a rank-1 matrix.



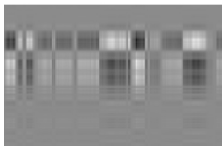
(a) Original image \mathbf{A} .
 1432×1910



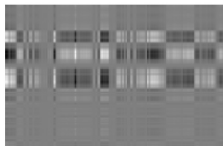
(b) \mathbf{A}_1 , $\sigma_1 \approx 228,052$.



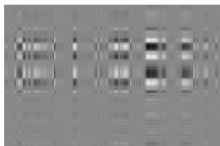
(c) \mathbf{A}_2 , $\sigma_2 \approx 40,647$.



(d) \mathbf{A}_3 , $\sigma_3 \approx 26,125$.



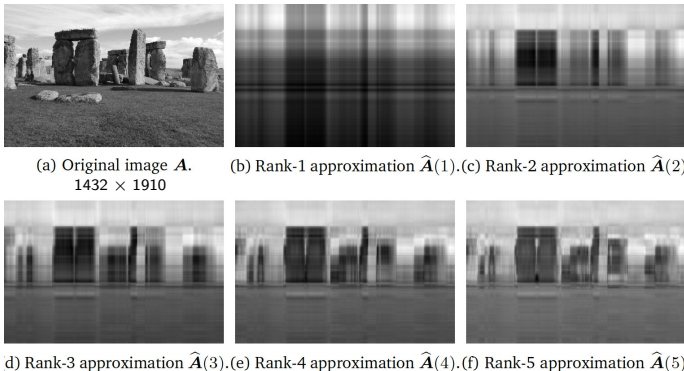
(e) \mathbf{A}_4 , $\sigma_4 \approx 20,232$.



(f) \mathbf{A}_5 , $\sigma_5 \approx 15,436$.

By summing only the first $k < r$ terms, we obtain a rank- k approximation $\hat{\mathbf{A}}_k$.

It turns out that $\hat{\mathbf{A}}_k$ is the **closest** rank- k matrix to \mathbf{A} , in terms of the spectral norm*.



* $\|\mathbf{A}\|_2 = \max_{\mathbf{x}} \|\mathbf{Ax}\|_2 / \|\mathbf{x}\|_2 = \sigma_1$. The spectral norm of \mathbf{A} is equal to its largest singular value.