# Matrix Decomposition

## Introduction

| Decomposition | definition | A | output |
|---|---|---|---|
| LU decomposition | $A = LU$ | square | triangular LU |
| LU decomposition with partial pivoting | $PA = LU$ | square | triangular LU |
| LU decomposition with full pivoting | $PAQ = LU$ | square | triangular LU |
| LDU decomposition | $A = LDU$ | square | triangular LU and diagonal D |
| Cholesky decomposition | $A = LL^T = U^T U$ | symmetric | triangular LU |
| Cholesky decomposition without sqrt | $A = LDL^T = U^T DU$ | symmetric | triangular LU and diagonal D |
| QR decomposition | $A = QR$ | rectangular | orthonormal Q and triangular R |
| LQ decomposition | $A = LQ$ | rectangular | orthonormal Q and triangular L |
| SVD decomposition | $A = USV^T$ | rectangular | orthonormal U,V & diagonal S |
| Eigen decomposition | $A = X\Lambda X^T$ | square | diagonal $\Lambda$ |
| | | symmetric | diagonal $\Lambda$ and orthonormal X |

## Rank decomposition

Suppose A is a Y×X matrix, the row span (column span) and row null space (column null space) are defined as :

$$span_{row}(A) \quad = \quad \{wA \quad s.t. \quad w = 1 \times Y \quad row\_matrix\} \qquad \text{i.e. set of linear combination of A's rows}$$

$$span_{column}(A) \quad = \quad \{Aw \quad s.t. \quad w = X \times 1 \quad column\_matrix\} \qquad \text{i.e. set of linear combination of A's columns}$$

$$null_{row}(A) \quad = \quad \{z \quad s.t. \quad Az = 0 \quad where \quad z = X \times 1 \quad column\_matrix\} \qquad \text{i.e. set of z that is orthogonal to A's rows}$$

$$null_{column}(A) \quad = \quad \{z \quad s.t. \quad zA = 0 \quad where \quad z = 1 \times Y \quad row\_matrix\} \qquad \text{i.e. set of z that is orthogonal to A's columns}$$

$$rank_{row}(A) \quad = \quad Dim(span_{row}(A))$$

$$rank_{columm}(A) \quad = \quad Dim(span_{column}(A))$$

$$nullity_{row}(A) \quad = \quad Dim(null_{row}(A))$$

$$nullity_{columm}(A) \quad = \quad Dim(null_{column}(A))$$

### Row rank decomposition

Suppose A is a Y×X matrix (by treating A as row data matrix in X dimensional space), row rank decomposition is :

$$\begin{bmatrix} A_1 \\ A_2 \\ ... \\ A_Y \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} & ... & w_{1,K} \\ w_{2,1} & w_{2,2} & ... & w_{2,K} \\ ... & ... & ... & ... \\ w_{Y,1} & w_{Y,2} & ... & w_{Y,K} \end{bmatrix} \begin{bmatrix} R_1 \\ R_2 \\ ... \\ R_K \end{bmatrix} \qquad \text{where } A_y \text{ and } R_k \text{ are 1×X row vectors}$$

$$s.t. \quad span_{row}(A) \quad = \quad \{uA \quad s.t. \quad u = 1 \times Y \quad row\_matrix\}$$
$$= \quad \{vR \quad s.t. \quad v = 1 \times K \quad row\_matrix\}$$
$$= \quad span_{row}(R)$$

and there exists no W (except $\delta_k$) such that $R_k$ = WR (or in a formal way, there exists no non trival solution to WR = 0). Thus row data matrix A is decomposed into a linearly independent (but not necessarily orthogonal) basis R (with K row vectors in X dimensional space), having the same row span as A, while K is called row rank.

### Column rank decomposition

Suppose A is a Y×X matrix (by treating A as column data matrix in Y dimensional space), column rank decomposition is :

$$\begin{bmatrix} A_1 & A_2 & ... & A_X \end{bmatrix} = \begin{bmatrix} C_1 & C_2 & ... & C_L \end{bmatrix} \begin{bmatrix} w'_{1,1} & w'_{2,1} & ... & w'_{X,1} \\ w'_{1,2} & w'_{2,2} & ... & w'_{X,2} \\ ... & ... & ... & ... \\ w'_{1,K} & w'_{2,K} & ... & w'_{X,L} \end{bmatrix} \qquad \text{where } A_x \text{ and } C_l \text{ are 1×Y column vectors}$$

$$s.t. \quad span_{column}(A) = \quad \{Aw \quad s.t. \quad w = X \times 1 \quad column\_matrix\}$$
$$= \quad \{Cv \quad s.t. \quad v = L \times 1 \quad column\_matrix\}$$
$$= \quad span_{column}(C)$$

and there exists no W (except $\delta_k$) such that $C_k$ = CW (or in a formal way, there exists no non trival solution to CW = 0). Thus column data matrix A is decomposed into a linearly independent (but not necessarily orthogonal) basis C (with L column vectors in Y dimensional space), having the same column span as A, while L is called column rank.

*Property 1 : Row rank equals to column rank*
*Lets prove that row rank equals to column rank. It is the same as proving the row rank of A equals to the row rank of $A^T$.*

$$A \quad = \quad WR \qquad\qquad\qquad\qquad \textit{i.e. row space of A is the span of rows in R}$$
$$A^T \quad = \quad R^T W^T \qquad\qquad\qquad \textit{i.e. row space of } A^T \textit{ is the span of rows in } W^T \qquad \textit{(equation 1)}$$

*from equation 1, we have :*

$$\textit{row rank of } A^T \quad \le \quad \textit{row rank of } W^T \quad = \quad \textit{num of row in } W^T \quad = \quad K \quad = \quad \textit{row rank of A} \qquad \textit{(equation 2)}$$

*putting $A^T$ into A of equation 2, we have :*

$$\textit{row rank of A} \quad \le \quad \textit{row rank of } A^T$$
$$\Rightarrow \qquad \textit{row rank of A} \quad = \quad \textit{row rank of A} \qquad \textcolor{red}{\textit{(Note : In general nullity}_{row} \ne \textit{nullity}_{column}\textit{)}}$$

*Property 2 : Rank null space theorem*
*Rank null space theorem states that sum of rank and nullity equals to dimension of data space, i.e.*

$$rank(A) + nullity_{row}(A) \quad = \quad rank_{row}(A) + nullity_{row}(A) \quad = \quad X$$
$$rank(A) + nullity_{column}(A) \quad = \quad rank_{column}(A) + nullity_{column}(A) \quad = \quad Y$$

*(Proof – step 1) Suppose rank of a Y×X* <u>row data matrix</u> *A is K, since row rank equals to column rank, A can be written as [$A_1$, $A_2$], where $A_1$ is a Y×K matrix of K independent* <u>column vectors</u> *and $A_2$ is a Y×(X-K) matrix of (X-K)* <u>column vectors</u> *dependent on $A_1$, hence $A_2$ can always be written as linear combination of $A_1$, where W is a K×(X-K) matrix.*

$$A_2 \quad = \quad A_1 W \qquad\qquad \Rightarrow \qquad A \quad = \quad [A_1, A_2] \quad = \quad [A_1, A_1 W]$$
$$\textit{if} \qquad Z \quad = \quad \begin{bmatrix} -W \\ I \end{bmatrix} \qquad\qquad\qquad \textit{where I is a (X-K)×(X-K) identity matrix}$$
$$\textit{then} \qquad AZ \quad = \quad [A_1, A_1 W] \begin{bmatrix} -W \\ I \end{bmatrix} \quad = \quad 0 \qquad \textit{hence columns of Z are X-K linearly independent solutions to Az = 0}$$

*(Proof – step 2) Suppose u is a solution for Az = 0, we have :*

$$Au \quad = \quad 0$$
$$Au \quad = \quad [A_1, A_1 W] \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad = \quad 0 \qquad \textit{where } u_1 \textit{ is K×1 column matrix, } u_2 \textit{ is (X-K)×1 column matrix.}$$
$$A_1(u_1 + W u_2) \quad = \quad 0$$
$$u_1 + W u_2 \quad = \quad 0 \qquad \textit{since } A_1 \textit{ is linear independent, i.e. the only solution is trival solution}$$
$$u_1 \quad = \quad -W u_2$$
$$u \quad = \quad \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad = \quad \begin{bmatrix} -W u_2 \\ u_2 \end{bmatrix} \quad = \quad \begin{bmatrix} -W \\ I \end{bmatrix} u_2 \quad = \quad Z u_2$$

*Thus any solution to Az = 0 can be expressed as linear combination of columns in Z, i.e. u = $Zu_2$. Hence the dimension of null space of row data matrix A is X-K, in other words, rank(A) + nullity$_{row}$(A) = K+(X-K) = X.* <u>*Pls repeat for column nullity.*</u>

*Property 3 : Rank of $A^T A$ equals to rank of A*

$$rank(A) \quad = \quad rank(A^T A) \qquad \textit{(equation 3)}$$
$$\Rightarrow \qquad rank(A^T) \quad = \quad rank(AA^T) \qquad \textit{(by property 1)}$$

*(Proof of equation 3)*

| For all $z$ | $\in$ $nullity_{row}(A)$ | For all $z$ | $\in$ $nullity_{row}(A^T A)$ |
|---|---|---|---|

$$Az \quad = \quad 0 \qquad\qquad\qquad\qquad A^T A z \quad = \quad 0$$
$$A^T A z \quad = \quad A^T 0 \ = \ 0 \qquad\qquad z^T A^T A z \quad = \quad z^T 0 \ = \ 0$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad (Az)^T Az \quad = \quad 0 \qquad \textit{(i.e. } Az = 0 \textit{)}$$
$$z \quad \in \quad nullity_{row}(A^T A) \qquad\qquad z \quad \in \quad nullity_{row}(A)$$
$$nullity_{row}(A) \quad \subseteq \quad nullity_{row}(A^T A) \qquad\qquad nullity_{row}(A^T A) \quad \subseteq \quad nullity_{row}(A)$$

*Combining the above results, we have :*

$$nullity_{row}(A) \quad = \quad nullity_{row}(A^T A)$$
$$rank(A) \quad = \quad rank(A^T A) \qquad\qquad \textit{(by property 2, since both A and } A^T A \textit{ live in dimension X)}$$

2

### LU (and LDU) decomposition

LU (and LDU) decomposition is a matrix form of Gaussian elimination. It is applicable to square matrix only, it can decompose a N×N matrix A into a lower triangular matrix L and an upper triangular matrix U, both with size N×N. The number of elements in A is $N^2$, while the total number of elements in L and U is $2\times(1+2+3+...+N) = N\times(N+1)$, thus given a matrix A, the corresponding matrix L and U are underdetermined, we need to add some constraints on L or U so as to obtain a unique LU decomposition. For example, we can constrain all diagonal elements of L to be one. From remark 4 and 5 in 'Elementary Operation' document, square matrix can (usually) be decomposed as applying a sequence of elementary row operations on an upper triangular matrix.

$$A \quad = \quad E_{N,N-1}E_{N,N-2}....E_{4,3}E_{4,2}E_{4,1}E_{3,2}E_{3,1}E_{2,1}U$$

Or in recursive form :

$$A \quad = \quad E_{N,N-1}A_{N,N-1}$$

$$A_{n,m} \quad = \quad \begin{bmatrix} E_{n,m-1}A_{n,m-1} & & m>1 \\ E_{n-1,n-2}A_{n-1,n-2} & n>2 & m=1 \\ E_{2,1}U & n=2 & m=1 \end{bmatrix} \qquad \forall n \in [2,N] \ \ and \ \ \forall m \in [1,n-1]$$

and elementary addition matrices are :

$$E_{n,m} \quad = \quad E_{add}(n,m,-p/q)$$

$p \quad = \quad$ element on the $n^{th}$ row and $m^{th}$ column in matrix $A_{n,m}$

$q \quad = \quad$ element on the $m^{th}$ row and $m^{th}$ column in matrix $A_{n,m}$

which means adding multiple (i.e. $-p/q$) of row m to row n, this is valid only for non zero q in all intermediate steps. This is Gaussian elimination. Besides, since n is always larger than m, all elementary row additions (i.e. $E_{n,m}$ $\forall n \in [1,N]$ and $\forall m \in [1,n-1]$) are lower triangular matrices with unit diagonal elements, and hence the cascade of elementary row additions is also a lower triangular matrix with unit diagonal elements, i.e.

$$L \quad = \quad E_{N,N-1}E_{N,N-2}....E_{4,3}E_{4,2}E_{4,1}E_{3,2}E_{3,1}E_{2,1} \qquad \forall n \in [2,N] \ \ and \ \ \forall m \in [1,n-1] \Rightarrow n>m$$

$where \quad L \quad = \quad (l_{n,m})_{\forall n \in [1,N], \forall m \in [1,N]}$

$s.t. \quad l_{n,m} \quad = \quad 0 \qquad \forall n,m \in [1,N] \ such \ that \ n<m$

$s.t. \quad l_{n,m} \quad = \quad 1 \qquad \forall n,m \in [1,N] \ such \ that \ n=m$

$\Rightarrow \quad A \quad = \quad LU$

This is known as the LU decomposition, which is valid only if we have non zero value for q in all intermediate steps. What happens if this requirement is not satisfied? In the first step of LU decomposition, i.e. $E_{2,1}U$, the value of q is $a_{1,1}$, suppose it is zero, then either $l_{1,1}$ or $u_{1,1}$ must be zero (recall : $a_{1,1} = l_{1,1}\times u_{1,1}$). If $l_{1,1}$ is zero, the determinant of L is zero, and if $u_{1,1}$ is zero, the determinant of U is zero (recall : the determinant of a lower or upper triangular matrix is the product of its diagonal elements), which is impossible, as the determinant of A can be non zero in general (recall : the determinant of matrix product is the product of individual's determinant). Thus, we can proceed LU decomposition for matrix A having zero $a_{1,1}$ by (method 1) reordering the rows of A to avoid zero diagonal element, (method 2) reordering both the rows and columns of A to avoid zero diagonal element, or (method 3) decomposing matrix A into 3 matrices LDU without reordering A. Method 1 is known as LU decomposition with partial pivoting, method 2 is known as LU decomposition with full pivoting, while method 3 is known as LDU decomposition.

### LU decomposition with partial (or full) pivoting

The reordering of rows and columns are done by multiplying matrix A with row permutation matrix P or (and) column permutation matrix Q. Please read 'Permutation' document for permutation matrix.

$PA \quad = \quad LU \quad$ (partial pivoting) $\qquad$ where P is a N×N row permutation matrix

$PAQ \quad = \quad LU \quad$ (full pivoting) $\qquad$ where Q is a N×N column permutation matrix

### LDU decomposition

LDU decomposition decomposes A into three N×N matrices : LDU, where L and U respectively are the lower and upper triangular matrices, both with all diagonal elements constrained to 1, while D is a diagonal matrix with configurable diagonal elements. Hence the total number of elements in LDU is : $2\times(1+2+3+...+(N-1)) + N = (N-1)\times N+N = N^2$, which is the same as that in matrix A. The philosophy of LDU decomposition is to separate the diagonal elements of A into D, so that D can capture both zero and non zero diagonal elements in A.

$$A \quad = \quad LDU \qquad\qquad where \ d_{n,m} = \begin{bmatrix} a_{n,n} & if & n=m \\ 0 & if & n \ne m \end{bmatrix}$$

## Cholesky decomposition

Cholesky decomposition is a special case of LU decomposition, which is applicable to N×N matrix A that is both (1) symmetric **(or Hermitian in general)** and (2) semi positive definite only. Symmetric and semi positive definite matrix is very common in different types of least square problem. Cholesky decomposition is faster than LU decomposition (nearly twice as efficient) as it can make use of the symmetric property of matrix A. With Cholesky decomposition, matrix A can be written as the product of a lower triangular matrix with the transpose of itself on the right, or as the product of an upper triangular matrix with the transpose of itself on the left. The diagonal elements of L and U are **not necessarily one**. The total number of elements in L or U is $(1+2+3+...+N) = N(N-1)/2$. Please note that : $L = U^T$.

$$
\begin{aligned}
A &= LL^T \\
&= U^T U \\
&= \begin{bmatrix}
l_{1,1} & 0 & 0 & ... & 0 \\
l_{2,1} & l_{2,2} & 0 & ... & 0 \\
l_{3,1} & l_{3,2} & l_{3,3} & ... & 0 \\
... & ... & ... & ... & ... \\
l_{N,1} & l_{N,2} & l_{N,3} & ... & l_{N,N}
\end{bmatrix}
\begin{bmatrix}
l_{1,1} & l_{2,1} & l_{3,1} & ... & l_{N,1} \\
0 & l_{2,2} & l_{3,2} & ... & l_{N,2} \\
0 & 0 & l_{3,3} & ... & l_{N,3} \\
... & ... & ... & ... & ... \\
0 & 0 & 0 & ... & l_{N,N}
\end{bmatrix} \\
&= \begin{bmatrix}
l_{1,1}^2 & symmetric & symmetric & ... & symmetric \\
l_{2,1}l_{1,1} & l_{2,1}^2+l_{2,2}^2 & symmetric & ... & symmetric \\
l_{3,1}l_{1,1} & l_{3,1}l_{2,1}+l_{3,2}l_{2,2} & l_{3,1}^2+l_{3,2}^2+l_{3,3}^2 & ... & symmetric \\
... & ... & ... & ... & ... \\
l_{N,1}l_{1,1} & l_{N,1}l_{2,1}+l_{N,2}l_{2,2} & l_{N,1}l_{3,1}+l_{N,2}l_{3,2}+l_{N,3}l_{3,3} & ... & l_{N,1}^2+l_{N,2}^2+...+l_{N,N}^2
\end{bmatrix}
\end{aligned}
$$

The implementation is straight forward. For all n starting from 1 to N, and for all m starting from 1 to n, we solve for the value of $l_{n,m}$ based on a set of previously found $l_{i,j}$, this method is known as **bootstrapping**. The bootstrapping order is :
$1_{1,1} \rightarrow 1_{2,1} \rightarrow 1_{2,2} \rightarrow 1_{3,1} \rightarrow 1_{3,2} \rightarrow 1_{3,3} \rightarrow ... \rightarrow 1_{n-1,1} \rightarrow 1_{n-1,2} \rightarrow ... \rightarrow 1_{n-1,n-1} \rightarrow 1_{n,1} \rightarrow 1_{n,2} \rightarrow ... \rightarrow 1_{n,m-1} \rightarrow 1_{n,m}$.

$$
a_{n,m} = \sum_{k=1}^{m} l_{n,k}l_{m,k} = \sum_{k=1}^{m-1} l_{n,k}l_{m,k} + l_{n,m}l_{m,m} \qquad \forall n \in [1,N] \ and \ \forall m \in [1,n]
$$

$$
\Rightarrow \begin{cases}
l_{n,n} &= \sqrt{a_{n,n} - \sum_{k=1}^{n-1} l_{n,k}l_{m,k}} & \forall n \in [1,N] \ and \ m = n \quad \text{(All } l_{i,j} \text{ on RHS have been found.)} \\
\\
l_{n,m} &= \left(a_{n,m} - \sum_{k=1}^{m-1} l_{n,k}l_{m,k}\right)\dfrac{1}{l_{m,m}} & \forall n \in [1,N] \ and \ m < n \quad \text{(All } l_{i,j} \text{ on RHS have been found.)}
\end{cases}
$$

### Cholesky decomposition with diagonal matrix

In some steps of Cholesky decomposition, there are square roots, which are slow. The square roots can be eliminated by introducing an extra diagonal matrix, i.e. $A = LDL^T = U^TDU$, while the diagonal elements of L and U **must be one**.
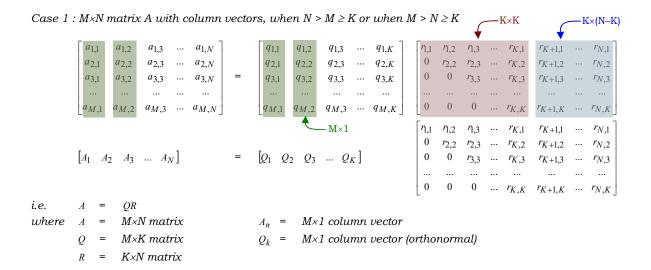
$$
\begin{aligned}
A &= LDL^T \\
&= U^T DU \\
&= \begin{bmatrix}
1 & 0 & 0 & ... & 0 \\
l_{2,1} & 1 & 0 & ... & 0 \\
l_{3,1} & l_{3,2} & 1 & ... & 0 \\
... & ... & ... & ... & ... \\
l_{5,1} & l_{5,2} & l_{5,3} & ... & 1
\end{bmatrix}
\begin{bmatrix}
d_1 & 0 & 0 & ... & 0 \\
0 & d_2 & 0 & ... & 0 \\
0 & 0 & d_3 & ... & 0 \\
... & ... & ... & ... & ... \\
0 & 0 & 0 & ... & d_N
\end{bmatrix}
\begin{bmatrix}
1 & l_{2,1} & l_{3,1} & ... & l_{N,1} \\
0 & 1 & l_{3,2} & ... & l_{N,2} \\
0 & 0 & 1 & ... & l_{N,3} \\
... & ... & ... & ... & ... \\
0 & 0 & 0 & ... & 1
\end{bmatrix} \\
&= \begin{bmatrix}
d_1 & symmetric & symmetric & ... & symmetric \\
l_{2,1}d_1 & l_{2,1}^2 d_1 + d_2 & symmetric & ... & symmetric \\
l_{3,1}d_1 & l_{3,1}l_{2,1}d_1 + l_{3,2}d_2 & l_{3,1}^2 d_1 + l_{3,2}^2 d_2 + d_3 & ... & symmetric \\
... & ... & ... & ... & ... \\
l_{N,1}d_1 & l_{N,1}l_{2,1}d_2 + l_{N,2}d_2 & l_{N,1}l_{3,1}d_1 + l_{N,2}l_{3,2}d_2 + l_{N,3}d_3 & ... & l_{N,1}^2 d_1 + l_{N,2}^2 d_2 + ... + l_{N,N-1}^2 d_{N-1} + d_N
\end{bmatrix}
\end{aligned}
$$

The implementation is similar.

$$
a_{n,m} = \sum_{k=1}^{m} l_{n,k}l_{m,k}d_k = \begin{cases}
\sum_{k=1}^{m-1} l_{n,k}l_{m,k}d_k + l_{n,m}l_{m,m}d_m = \sum_{k=1}^{m-1} l_{n,k}l_{m,k}d_k + d_m & m = n \\
\sum_{k=1}^{m-1} l_{n,k}l_{m,k}d_k + l_{n,m}l_{m,m}d_m = \sum_{k=1}^{m-1} l_{n,k}l_{m,k}d_k + l_{n,m}d_m & m < n
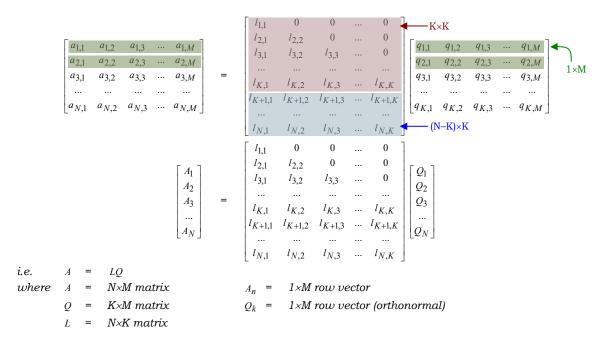\end{cases}
$$

$$
\Rightarrow \begin{cases}
d_n &= a_{n,n} - \sum_{k=1}^{m-1} l_{n,k}l_{m,k}d_k & \forall n \in [1,N] & \text{(All } l_{i,j} \text{ on RHS have been found.)} \\
\\
l_{n,m} &= \left(a_{n,m} - \sum_{k=1}^{m-1} l_{n,k}l_{m,k}d_k\right)\dfrac{1}{d_m} & \forall n \in [1,N] \ and \ m < n & \text{(All } l_{i,j} \text{ on RHS have been found.)}
\end{cases}
$$

### QR (and LQ) decomposition

*QR (and LQ) decomposition performs orthonormalization based on Gram Schmidt process, the former is applicable to rectangular matrix A with size M×N (i.e. N column vectors in dimension M), while the latter is applicable to rectangular matrix A with size N×M (i.e. N row vectors in dimension M). Q matrix is a set of K orthonormal vectors, i.e. { Q₁, Q₂, Q₃, …, Q_K } such that $Q_i^T Q_j = \delta_{ij}$ ∀i, j∈[1, K] for QR decomposition, while $Q_i Q_j^T = \delta_{ij}$ ∀i, j∈[1, K] for LQ decomposition, where K is known as the rank of matrix A. [Recall : column vector means representing a vector by a column matrix and row vector means representing a vector by a row matrix.]*

*Case 1 : M×N matrix A with column vectors, when N > M ≥ K or when M > N ≥ K*

$$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & \cdots & a_{1,N} \\ a_{2,1} & a_{2,2} & a_{2,3} & \cdots & a_{2,N} \\ a_{3,1} & a_{3,2} & a_{3,3} & \cdots & a_{3,N} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ a_{M,1} & a_{M,2} & a_{M,3} & \cdots & a_{M,N} \end{bmatrix} = \begin{bmatrix} q_{1,1} & q_{1,2} & q_{1,3} & \cdots & q_{1,K} \\ q_{2,1} & q_{2,2} & q_{2,3} & \cdots & q_{2,K} \\ q_{3,1} & q_{3,2} & q_{3,3} & \cdots & q_{3,K} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ q_{M,1} & q_{M,2} & q_{M,3} & \cdots & q_{M,K} \end{bmatrix} \begin{bmatrix} r_{1,1} & r_{1,2} & r_{1,3} & \cdots & r_{K,1} & r_{K+1,1} & \cdots & r_{N,1} \\ 0 & r_{2,2} & r_{2,3} & \cdots & r_{K,2} & r_{K+1,2} & \cdots & r_{N,2} \\ 0 & 0 & r_{3,3} & \cdots & r_{K,3} & r_{K+1,3} & \cdots & r_{N,3} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & r_{K,K} & r_{K+1,K} & \cdots & r_{N,K} \end{bmatrix}$$

K×K    K×(N–K)     M×1

$$\begin{bmatrix} A_1 & A_2 & A_3 & \cdots & A_N \end{bmatrix} = \begin{bmatrix} Q_1 & Q_2 & Q_3 & \cdots & Q_K \end{bmatrix} \begin{bmatrix} r_{1,1} & r_{1,2} & r_{1,3} & \cdots & r_{K,1} & r_{K+1,1} & \cdots & r_{N,1} \\ 0 & r_{2,2} & r_{2,3} & \cdots & r_{K,2} & r_{K+1,2} & \cdots & r_{N,2} \\ 0 & 0 & r_{3,3} & \cdots & r_{K,3} & r_{K+1,3} & \cdots & r_{N,3} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & r_{K,K} & r_{K+1,K} & \cdots & r_{N,K} \end{bmatrix}$$

*i.e.*    A = QR

*where*    A = M×N *matrix*      $A_n$ = M×1 *column vector*

         Q = M×K *matrix*      $Q_k$ = M×1 *column vector (orthonormal)*

         R = K×N *matrix*

*Case 2 : N×M matrix A with row vectors, when N > M ≥ K or when M > N ≥ K*

$$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & \cdots & a_{1,M} \\ a_{2,1} & a_{2,2} & a_{2,3} & \cdots & a_{2,M} \\ a_{3,1} & a_{3,2} & a_{3,3} & \cdots & a_{3,M} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ a_{N,1} & a_{N,2} & a_{N,3} & \cdots & a_{N,M} \end{bmatrix} = \begin{bmatrix} l_{1,1} & 0 & 0 & \cdots & 0 \\ l_{2,1} & l_{2,2} & 0 & \cdots & 0 \\ l_{3,1} & l_{3,2} & l_{3,3} & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ l_{K,1} & l_{K,2} & l_{K,3} & \cdots & l_{K,K} \\ l_{K+1,1} & l_{K+1,2} & l_{K+1,3} & \cdots & l_{K+1,K} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ l_{N,1} & l_{N,2} & l_{N,3} & \cdots & l_{N,K} \end{bmatrix} \begin{bmatrix} q_{1,1} & q_{1,2} & q_{1,3} & \cdots & q_{1,M} \\ q_{2,1} & q_{2,2} & q_{2,3} & \cdots & q_{2,M} \\ q_{3,1} & q_{3,2} & q_{3,3} & \cdots & q_{3,M} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ q_{K,1} & q_{K,2} & q_{K,3} & \cdots & q_{K,M} \end{bmatrix}$$

K×K    (N–K)×K    1×M

$$\begin{bmatrix} A_1 \\ A_2 \\ A_3 \\ \cdots \\ A_N \end{bmatrix} = \begin{bmatrix} l_{1,1} & 0 & 0 & \cdots & 0 \\ l_{2,1} & l_{2,2} & 0 & \cdots & 0 \\ l_{3,1} & l_{3,2} & l_{3,3} & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ l_{K,1} & l_{K,2} & l_{K,3} & \cdots & l_{K,K} \\ l_{K+1,1} & l_{K+1,2} & l_{K+1,3} & \cdots & l_{K+1,K} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ l_{N,1} & l_{N,2} & l_{N,3} & \cdots & l_{N,K} \end{bmatrix} \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \\ \cdots \\ Q_N \end{bmatrix}$$

*i.e.*    A = LQ

*where*    A = N×M *matrix*      $A_n$ = 1×M *row vector*

         Q = K×M *matrix*      $Q_k$ = 1×M *row vector (orthonormal)*

         L = N×K *matrix*

*Now lets find Q matrix, R matrix and L matrix using Gram Schmidt process. The Q matrix is :*

$$Q_1 = \frac{A_1}{\|A_1\|}$$

$$Q_2 = \frac{A_2 - proj_{Q_1}(A_2)}{\|A_2 - proj_{Q_1}(A_2)\|}$$

$$Q_3 = \frac{A_3 - proj_{Q_1}(A_3) - proj_{Q_2}(A_3)}{\|A_3 - proj_{Q_1}(A_3) - proj_{Q_2}(A_3)\|}$$

$$Q_4 = \frac{A_4 - proj_{Q_1}(A_4) - proj_{Q_2}(A_4) - proj_{Q_3}(A_4)}{\|A_4 - proj_{Q_1}(A_4) - proj_{Q_2}(A_4) - proj_{Q_3}(A_4)\|}$$

*while the R matrix and L matrix are :*

$$r_{k,n} \quad = \quad l_{n,k} \qquad\qquad \forall k \in [1, K] \ \ and \ \ \forall n \in [1, N] \qquad (equation\ 1)$$

$$\quad = \quad <A_n, Q_k> \qquad \forall k \in [1, K] \ \ and \ \ \forall n \in [1, N]$$

$$\quad = \quad A_n^T Q_k \qquad\qquad for\ column\ vector$$

$$\quad = \quad A_n Q_k^T \qquad\qquad for\ row\ vector$$

*If $A_n$ is a column vector, then $Q_k$ is a column vector, and*

$$proj_{Q_k}(A_n) \quad = \quad Q_k \frac{Q_k^T A_n}{Q_k^T Q_k} \quad = \quad Q_k(Q_k^T A_n) \qquad since \ Q_k^T Q_k = 1$$

$$\|A_n\| \quad = \quad A_n^T A_n$$

*If $A_n$ is a row vector, then $Q_k$ is a row vector, and*

$$proj_{Q_k}(A_n) \quad = \quad \frac{A_n Q_k^T}{Q_k Q_k^T} Q_k \quad = \quad (A_n Q_k^T) Q_k \qquad since \ Q_k Q_k^T = 1$$

$$\|A_n\| \quad = \quad A_n A_n^T$$

*In general, for **non orthogonal basis** Z, the projection is defined as :*

*If $A_n$ is a column vector, then $Z_k$ is a column vector, and*

$$proj_{Z_k}(A_n) \quad = \quad Z_k(Z_k^T Z_k)^{-1} Z_k^T A_n \quad = \quad PA_n \quad \longleftarrow$$

*If $A_n$ is a row vector, then $Z_k$ is a row vector, and*

$$proj_{Q_k}(A_n) \quad = \quad A_n Z_k^T (Z_k Z_k^T)^{-1} Z_k \quad = \quad A_n P$$

*explanation :*

$$Z_k^T A_n \qquad\qquad inner\ product$$
$$(Z_k^T Z_k)^{-1} Z_k^T A_n \qquad normalized\ inner\ product$$
$$Z_k(Z_k^T Z_k)^{-1} Z_k^T A_n \qquad normalized\ vector$$
$$Z_k(Z_k^T Z_k)^{-1} Z_k^T = P \qquad projection\ matrix\ (col)$$

*Can you plot it geometrically?*

*QR decomposition and LQ decomposition are in fact related by a transpose. Suppose A is M×N, then $A^T$ is N×M.*

$$A \quad = \quad QR \qquad\qquad QR\ decomposition$$
$$A^T \quad = \quad LQ^T \qquad\qquad LQ\ decomposition$$
$$A^T \quad = \quad R^T Q^T \qquad\qquad transpose\ of\ QR\ decomposition$$
$$thus \quad L \quad = \quad R^T \qquad\qquad hence\ it\ is\ consistent\ with\ equation\ 1$$

*Gram Schmidt process cannot help to find the null space, yet we can still generalize QR and LQ decomposition by including the null space (found by other methods) into the Q matrix, then we have (Z stands for the M–K dimensional null space) :*



*QR decomposition*

*LQ decomposition*

### Eigen decomposition

Lets define eigenvector (eigenfunction) and eigenvalue first. Eigenvector X of a **linear transformation** (which maps a N dimensional vector to another N dimensional vector, and is usually represented as a N×N matrix A) is a vector, which becomes a multiple of itself after applying the linear transformation, thus eigenvector are scaled in magnitude only (with no change in direction) after applying the linear transformation. The scale in magnitude is called eigenvalue. This idea can be generalized to a more universal set of **linear operators**. For example, if L is a linear differential operator which maps a continuous function to another continuous function, then eigenfunctions f is defined as functions, when operated by linear differential operator L, returns multiple of itself (please read the 'Green function' document). Both eigenvector and eigenfunction share common properties.

| | | | | | | |
|---|---|---|---|---|---|---|
| $AX$ | = | $\lambda X$ | | $Lf$ | = | $\lambda f$ |
| $X$ | = | eigenvector (N×1 column matrix) | | $f$ | = | eigenfunction |
| $\lambda$ | = | eigenvalue (scalar) | | $\lambda$ | = | eigenvalue |

_Property 1_
Multiple of eigenvector (eigenfunction) is also an eigenvector (eigenfunction).

| If | $AX$ | = | $\lambda X$ | |
|---|---|---|---|---|
| then | $A(cX)$ | = | $cA(X)$ | since A is a linear transformation |
| | | = | $c\lambda X$ | since X is an eigenvector |
| | | = | $\lambda(cX)$ | hence proved |
| If | $Lf$ | = | $\lambda f$ | |
| then | $L(cf)$ | = | $cL(f)$ | since L is a linear operator |
| | | = | $c\lambda f$ | since f is an eigenfunction |
| | | = | $\lambda(cf)$ | hence proved |

_Property 2_
Weighted sum of eigenvectors (eigenfunctions) having **same eigenvalue** is also an eigenvector (eigenfunction) having the same eigenvalue.

| If | $AX_1$ | = | $\lambda X_1$ | |
|---|---|---|---|---|
| and | $AX_2$ | = | $\lambda X_2$ | |
| then | $A(w_1 X_1 + w_2 X_2)$ | = | $w_1 A(X_1) + w_2 A(X_2)$ | since A is a linear transformation |
| | | = | $w_1 \lambda X_1 + w_2 \lambda X_2$ | |
| | | = | $\lambda(w_1 X_1 + w_2 X_2)$ | hence proved |
| If | $Lf_1$ | = | $\lambda f_1$ | |
| and | $Lf_2$ | = | $\lambda f_2$ | |
| then | $L(w_1 f_1 + w_2 f_2)$ | = | $w_1 Lf_1 + w_2 Lf_2$ | since L is a linear operator |
| | | = | $w_1 \lambda X_1 + w_2 \lambda X_2$ | |
| | | = | $\lambda(w_1 X_1 + w_2 X_2)$ | hence proved |

_Property 3_ From now on, lets focus on eigenvectors only (ignore eigenfunction).
If A is an **identity matrix**, then all vectors are eigenvectors. Here is the one line proof.

$$IX \quad = \quad X \quad = \quad \lambda X \qquad\qquad \text{which is true for all X, where } \lambda = 1$$

_Property 4_
If A is a **diagonal matrix**, then all N unit vectors in space $\mathfrak{R}^N$ are eigenvectors, with eigenvalues being equivalent to the diagonal elements of A. Suppose $I_n$ be the nth column vector of identity matrix I, i.e. $I_n$ is a unit vector in space $\mathfrak{R}^N$ :

$$AI_n = \begin{bmatrix} a_{1,1} & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & a_{n,n} & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \dots & a_{N,N} \end{bmatrix} \begin{bmatrix} 0 \\ \dots \\ 1 \\ \dots \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ \dots \\ a_{n,n} \\ \dots \\ 0 \end{bmatrix} = a_{n,n} \begin{bmatrix} 0 \\ \dots \\ 1 \\ \dots \\ 0 \end{bmatrix}$$

_Property 5_
If A is a **triangular matrix** (either lower or upper), then its diagonal elements are the N eigenvalues. The proof starts with solving the characteristic function for eigenvalues (please refer to implemention method 1).

| | | |
|---|---|---|
| $\det(A - \lambda I)$ | = | 0 |
| $\prod_{n=1}^{N}(A - \lambda I)_{n,n}$ | = | 0 $\qquad$ since $A - \lambda I$ is also a triangular matrix |
| $\prod_{n=1}^{N}(a_{n,n} - \lambda)$ | = | 0 $\qquad \Rightarrow \quad \lambda = a_{n,n} \qquad \forall n \in [1, N]$ |

*Eigenspace, eigenbasis and eigen decomposition*
*Given a linear transformation A, the number of eigenvectors can be infinity, more than N or less than N. Eigenvectors are not necessarily orthogonal. Lets group the eigenvectors according to their eigenvalues, i.e. eigenvectors having the same eigenvalue $\lambda_m$ are grouped to form a set. According to property 2, any span of the set is also an eigenvector having the same eigenvalue (it belongs to the same set too). Now we name the set as the eigenspace of A associated with eigenvalue $\lambda_m$. There are infinity eigenvectors in that eigenspace, yet we can orthonormalize the eigenspace into $\gamma_A(\lambda_m)$ orthonormal eigenvectors, where $\gamma_A(\lambda_m)$ is known as geometric multiplicity, which equals to the dimension of the eigenspace. For simplicity, from now on, the term 'eigenvector' refers only to eigenvector after orthonormalization process (i.e. we **purify** the entire set eigenvector by removing eigenvectors that are dependent on the others). Besides, there may be multiple eigenspaces for A in general. Suppose A denotes a linear transformation from $\mathfrak{R}^N$ to $\mathfrak{R}^N$, and suppose there are M eigenspaces, each is associated with an eigenvalue $\lambda_m$, $\forall m \in [1,M]$, then :*

- *eigenvectors from the same eigenspace are **orthonormal** to each other (because of orthonormalization)*
- *eigenvectors from different eigenspaces are **not necessaily orthonormal** to each other*
- *sum of geometric multiplicity over all eigenspaces gives the total number of 'purified' eigenvectors (must be $\leq N$)*

$$1 \quad \leq \quad \sum_{m=1}^{M} \gamma_A(\lambda_m) \quad \leq \quad N \qquad \text{where} \quad M \leq N$$

$$1 \quad = \quad \gamma_A(\lambda_m) \qquad \forall m \in [1,M] \qquad \text{when} \quad M = N$$

*If the total number of eigenvectors is exactly N, then the eigenvectors form a set called eigenbasis. Please note that : (1) eigenbasis does not necessarily exist, (2) eigenvectors in eigenbasis are not necessarily orthonormal, hence eigenbasis is just a basis (not an orthogonormal basis). If eigenbasis does exist, then we can define eigen decomposition by grouping all eigenvectors as rows (or as columns) of N×N matrix X, and grouping all eigenvalues as diagonal elements of N×N diagonal matrix $\Lambda$.*

| | *For row eigenvector X* | | | | *For column eigenvector X* | | |
|---|---|---|---|---|---|---|---|
| | $XA$ | $=$ | $\Lambda X$ | | $AX$ | $=$ | $X\Lambda$ |
| $\Rightarrow$ | $A$ | $=$ | $X^{-1}\Lambda X$ | $\Rightarrow$ | $A$ | $=$ | $X\Lambda X^{-1}$ |

$$\text{where} \quad X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ ... \\ X_N \end{bmatrix} \qquad\qquad \text{where} \quad X = \begin{bmatrix} X_1 & X_2 & X_3 & ... & X_N \end{bmatrix}$$

$$X_n = 1\times N \text{ row matrix} \qquad\qquad X_n = N\times 1 \text{ column matrix}$$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 & ... & 0 \\ 0 & \lambda_2 & 0 & ... & 0 \\ 0 & 0 & \lambda_3 & ... & 0 \\ ... & ... & ... & ... & ... \\ 0 & 0 & 0 & ... & \lambda_N \end{bmatrix} \qquad \Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 & ... & 0 \\ 0 & \lambda_2 & 0 & ... & 0 \\ 0 & 0 & \lambda_3 & ... & 0 \\ ... & ... & ... & ... & ... \\ 0 & 0 & 0 & ... & \lambda_N \end{bmatrix}$$

*Note : Matrix A in the LHS case and in the RHS case are the same matrix, while matrix X in the LHS case and in the RHS case are related by a transpose. Hence the **LHS case and the RHS case are identical**.*

*Interpretation 1 (linear transformation)*
*For both LHS and RHS cases above, we treat A and X as transformation matrices (hence we don't need to classify them as row or column matrix), eigen decomposition can be interpreted as decomposing the linear transformation A into a shearing X (or $X^{-1}$), followed by a dimension scaling $\Lambda$, and finally with an inverse shearing $X^{-1}$ (or X).*

*Interpretation 2 (linear combination)*
*For the LHS case above, if we treat A and X as row data matrices, then eigenvalue decomposition can be interpreted as a linear combination of rows in X (i.e. eigenvectors) with weights specified in $X^{-1}\Lambda$, if we treat A and $X^{-1}$ as column data matrices, then eigenvalue decomposition can be interpreted as a linear combination of columns in $X^{-1}$ (i.e. eigenvectors) with weights specified in $\Lambda X$.*

*For the RHS case above, if we treat A and X as column data matrices, then eigenvalue decomposition can be interpreted as a linear combination of columns in X (i.e. eigenvectors) with weights specified in $\Lambda X^{-1}$, if we treat A and $X^{-1}$ as row data matrices, then eigenvalue decomposition can be interpreted as a linear combination of rows in $X^{-1}$ (i.e. eigenvectors) with weights specified in $X\Lambda$.*

*Interpretation 3 (inner product or covariance)*
*For both LHS and RHS cases above, A can also be interpreted as inner product within $X^{-1}$ or covariance (weighted 2nd moment, to be more precise) of X.*

### What does it mean physically if the eigenbasis is orthonormal?

*Following interpretation 1 above, when eigenbasis is orthonormal, the shearing X becomes a rotation in $\Re^N$ space without any scaling (since orthonormal transformation is a matrix with sine and cosine, which means rotation). As there are $N^2$ elements in matrix A, an orthonormal eigenbasis means some constrains on the degree of freedom of A. Yet, what is the constrain exactly? Lets consider a N=2 case.*

$$
\begin{aligned}
A &= \begin{bmatrix} \cos\vartheta & +\sin\vartheta \\ -\sin\vartheta & \cos\vartheta \end{bmatrix} \begin{bmatrix} \lambda_x & 0 \\ 0 & \lambda_y \end{bmatrix} \begin{bmatrix} \cos\vartheta & -\sin\vartheta \\ +\sin\vartheta & \cos\vartheta \end{bmatrix} \\[2mm]
&= \begin{bmatrix} \cos\vartheta & +\sin\vartheta \\ -\sin\vartheta & \cos\vartheta \end{bmatrix} \begin{bmatrix} \lambda_x\cos\vartheta & -\lambda_x\sin\vartheta \\ +\lambda_y\sin\vartheta & \lambda_y\cos\vartheta \end{bmatrix} \\[2mm]
&= \begin{bmatrix} \lambda_x\cos^2\vartheta + \lambda_y\sin^2\vartheta & (\lambda_y-\lambda_x)\sin\vartheta\cos\vartheta \\ (\lambda_y-\lambda_x)\sin\vartheta\cos\vartheta & \lambda_x\sin^2\vartheta + \lambda_y\cos^2\vartheta \end{bmatrix} \quad = \quad A^T
\end{aligned}
$$

*Hence when eigenbasis is orthonormal, A can be replicated by a rotation plus a scale plus a inverse rotation, resulting in a symmetric A, which gives us inspiration to derive property 6.*

### Property 6 – for symmetric matrix

*In general, eigenvectors from different eigenspaces are not necessarily orthonormal to each other. However they must be orthonormal when linear transformation A is symmetric. Suppose A is a $N\times N$ symmetric matrix, while $X_n$ and $X_m$ (where $n,m\in[1,N]$ such that $n\neq m$) are two $N\times 1$ column eigenvectors with different eigenvalues $\lambda_n$ and $\lambda_m$ respectively, they are picked from different eigenspaces, then we have :*

$$
\begin{aligned}
AX_n &= \lambda_n X_n & \Rightarrow \quad X_n &= A^{-1}(\lambda_n X_n) \\
AX_m &= \lambda_m X_m & \Rightarrow \quad X_m &= (AX_m)/\lambda_m \\
X_n^T X_m &= (\lambda_n/\lambda_m)(A^{-1}X_n)^T(AX_m) \\
&= (\lambda_n/\lambda_m)X_n^T(A^{-1})^T AX_m \\
&= (\lambda_n/\lambda_m)X_n^T(A^{-1})^T A^T X_m & &\text{since A is symmetric, i.e. } A^T = A \\
&= (\lambda_n/\lambda_m)X_n^T(AA^{-1})^T X_{mj} \\
&= (\lambda_n/\lambda_m)X_{ni}^T X_m \\
0 &= (1-\lambda_n/\lambda_m)X_n^T X_m \\
X_n^T X_m &= 0 & &\text{since } \lambda_i \neq \lambda_j \text{ as they belong to different eigenspaces}
\end{aligned}
$$

### Property 7 – for symmetric positive definite matrix

*When a square matrix is both symmetric and positive definite, then its eigen values must be all positive. Suppose A is a $N\times N$ symmetric and positive definite matrix, then its eigenvectors must form orthonormal basis and we have :*

$$
\begin{aligned}
& A & &= X^T \Lambda X & &\text{where } X^T X = I \\
\text{and}\quad & y^T A y & &> 0 & &\text{for all non zero } N\times 1 \text{ column matrix } y \\
\Rightarrow\quad & y^T X^T \Lambda X y & &> 0 \\
& (Xy)^T \Lambda X y & &> 0 \\
& z^T \Lambda z & &> 0 & &\text{for some } N\times 1 \text{ column matrix z, such that } z = Xy \\
& \sum_{n=1}^{N}\Lambda_{n,n}z_n^2 & &> 0 \\
\Rightarrow\quad & \Lambda_{n,n} & &> 0 & &\text{for all n}
\end{aligned}
$$

### Property 8A – for square root of matrix

*Square root is defined for symmetric positive definite matrix only :*

$$
A = A^{1/2}A^{1/2} \qquad \text{for symmetric positive definite matrix A}
$$

*Since A is symmetric positive definite, all its eigenvalues should be positive.*

$$
\begin{aligned}
A &= X\Lambda X^{-1} \\
&= X\Lambda^{1/2}\Lambda^{1/2}X^{-1} & &\text{as we define } \Lambda = \Lambda^{1/2}\Lambda^{1/2}, \text{ i.e. } \Lambda_{n,n} = (\Lambda_{n,n}^{1/2})^2 \\
&= X\Lambda^{1/2}X^{-1}X\Lambda^{1/2}X^{-1} \\
&= A^{1/2}A^{1/2} \\
\text{i.e.}\quad A^{1/2} &= X\Lambda^{1/2}X^{-1}
\end{aligned}
$$

*Therefore A and $A^{1/2}$ share the same eigenvectors, while having eigenvalues related by square root.*

*Property 8B – for power of matrix*
*Power of a N×N matrix A (i.e. cascading linear transformation) is :*

$$
\begin{aligned}
A &= X\Lambda X^{-1} \\
A^2 &= (X\Lambda X^{-1})(X\Lambda X^{-1}) &&= X\Lambda^2 X^{-1} \\
A^3 &= (X\Lambda X^{-1})(X\Lambda X^{-1})(X\Lambda X^{-1}) &&= X\Lambda^3 X^{-1} \\
A^N &= X\Lambda^N X^{-1}
\end{aligned}
$$

*Property 9 – for similarity*
*Similarity transformation from a N×N matrix A to a N×N matrix B is defined as :*

$$
B = QAQ^{-1} \qquad\qquad \text{where Q is N×N orthogonal basis}
$$

*then the eigen decomposition of B is related to the eigen decomposition of A by :*

$$
\begin{aligned}
A &= X\Lambda X^{-1} &&\text{where X is N×N eigenbasis of A (not necessarily orthonormal)} \\
B &= Y\Delta Y^{-1} &&\text{where Y is N×N eigenbasis of B (not necessarily orthonormal)} \\
\Rightarrow\quad B &= QAQ^{-1} \\
&= QX\Lambda X^{-1}Q^{-1} \\
&= (QX)\Lambda(QX)^{-1} \\
&= Y\Delta Y^{-1} \\
\Rightarrow\quad Y &= QX &&\text{product of orthonormal is also orthonormal} \\
\text{and}\quad \Delta &= \Lambda
\end{aligned}
$$

$$
\begin{aligned}
\text{since}\quad \det(A) &= \prod_{n=1}^{N}\Lambda_n \\
\&\quad \det(B) &= \prod_{n=1}^{N}\Delta_n &&\text{thus we have : } \det(A) = \det(B) \\
\text{since}\quad tr(A) &= \sum_{n=1}^{N}\Lambda_n \\
\&\quad tr(B) &= \sum_{n=1}^{N}\Delta_n &&\text{thus we have : } tr(A) = tr(B)
\end{aligned}
$$

*When eigenbasis (eigen decomposition) exists, we can derive it using the following methods.*

*Implementation method 1 : Characteristic equation*
*once eigenvalues are found, eigenvectors can be found by solving system of linear equation*

$$
\begin{aligned}
AX &= \lambda X \\
(A-\lambda I)X &= 0 &&\text{(equation 2)} \\
\det(A-\lambda I) &= 0
\end{aligned}
$$

*resulting in a polynomial of $\lambda$ in degree N, which can be factorized and grouped into different M eigenspaces.*

$$
\begin{aligned}
f(\lambda) &= c_N\lambda^N + c_{N-1}\lambda^{N-1} + ... + c_3\lambda^3 + c_2\lambda^2 + c_1\lambda + c_0 \\
&= c_N(\lambda-\lambda_1)^{\gamma_A(\lambda_1)}(\lambda-\lambda_2)^{\gamma_A(\lambda_2)}...(\lambda-\lambda_M)^{\gamma_A(\lambda_M)}
\end{aligned}
$$

*After obtaining M roots for $\lambda$ using numerical method, we substitute the values of $\lambda_m$ back into equation 2 (for all m from 1 to M), resulting in a system of **linearly dependent** equations, in which all equations pass through the origin, and intersect along a line, a plane or a hyperplane. The intersecting hyperplane is the **eigenspace** corresponding to eigenvalue $\lambda_m$, while the dimension of the intersecting hyperplane gives the **geometric multiplicity** associated with $\lambda_m$. We then need to orthonormalize the eigenspace to give **eigenvectors**. Recall : the equations do not intersect at a single point, if they do, this point must be the origin, and we deliberately make the system linearly dependent by setting determinant zero (i.e. $\infty$ solutions). Thus never solve X by performing LU (or other) decomposition on equation 2.*

*Implementation method 2 : Power method*
*Suppose there exists an eigenbasis for linear transformation A, if X and $\Lambda$ are N×N matrix containing the N column eigenvectors $X_n$ $\forall n\in[1,N]$ and N eigenvalues $\lambda_n$ $\forall n\in[1,N]$, then power method starts with any initial guess V, which can always be expressed as a linear combinations of eigenbasis, i.e. V is a N×1 column vector such that :*

$$
\begin{aligned}
V &= XW \\
&= \text{any initial guess}
\end{aligned}
$$

$$
\text{where}\quad X = \begin{bmatrix} X_1 & X_2 & X_3 & ... & X_N \end{bmatrix}
$$

$$and \quad W \quad = \quad \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ ... \\ w_N \end{bmatrix}$$

*Power method is an iterative method which eventually gives the dominant eigenvalue as :*

$$V^{(1)} \quad = \quad AV$$
$$V^{(t)} \quad = \quad AV^{(t-1)}$$
$$then \quad \lambda_1 \quad = \quad \lim_{t \to \infty} V_1^{(t)} / V_1^{(t-1)}$$
$$= \quad \lim_{t \to \infty} V_2^{(t)} / V_2^{(t-1)}$$
$$= \quad \lim_{t \to \infty} V_m^{(t)} / V_m^{(t-1)} \quad \forall m \in [1, N] \qquad \textit{remark : pick the } m^{th} \textit{ element}$$

*(Proof)*
$$V \quad = \quad XW$$
$$= \quad \sum_{n=1}^{N} w_n X_n$$
$$V^{(1)} \quad = \quad \sum_{n=1}^{N} w_n A X_n \qquad = \quad \sum_{n=1}^{N} w_n \lambda_n X_n \qquad \textit{since } AX_n = \lambda_n X_n$$
$$V^{(2)} \quad = \quad \sum_{n=1}^{N} w_n \lambda_n A X_n \qquad = \quad \sum_{n=1}^{N} w_n \lambda_n^2 X_n \qquad \textit{since } AX_n = \lambda_n X_n$$
$$V^{(t)} \quad = \quad \sum_{n=1}^{N} w_n \lambda_n^{t-1} X_n \qquad = \quad \sum_{n=1}^{N} w_n \lambda_n^t X_n$$
$$\lim_{t \to \infty} \frac{V_m^{(t)}}{V_m^{(t-1)}} \quad = \quad \lim_{t \to \infty} \frac{\sum_{n=1}^{N} w_n \lambda_n^t X_{n,m}}{\sum_{n=1}^{N} w_n \lambda_n^{t-1} X_{n,m}} \qquad \textit{remark : pick the } m^{th} \textit{ element}$$
$$= \quad \frac{w_1 \lambda_1^t X_{1,m}}{w_1 \lambda_1^{t-1} X_{1,m}} \qquad \textit{since } \lambda_1 \textit{ is dominant eigenvalue}$$
$$= \quad \lambda_1 \qquad \qquad \textbf{\textit{Question : How about other eigenvalues?}}$$

___
*Implementation method 3 : QR algorithm (Don't confuse with QR decomposition)*
*QR decomposition starts with $A^{(1)} = A$, and in each iteration t, we apply QR decomposition on $A^{(t)}$, such that $A^{(t)} = Q^{(t)}R^{(t)}$, and then generate $A^{(t+1)}$ as the product of QR matrices in reverse order, i.e. $A^{(t+1)} = R^{(t)}Q^{(t)}$, by repeating the process, $A^{(t)}$ will eventually become a similarity transformation of an upper triangular matrix U. The diagonal elements of the upper triangular matrix U give all the N eigenvalues.*

$$A^{(t+1)} \quad = \quad R^{(t)} Q^{(t)}$$
$$= \quad (Q^{(t)})^{-1} Q^{(t)} R^{(t)} Q^{(t)} \qquad \textit{since } Q^{(t)T} Q^{(t)} = I \textit{ as } Q^{(t)} \textit{ is orthonormal}$$
$$= \quad (Q^{(t)})^{-1} A^{(t)} Q^{(t)}$$
$$= \quad (Q^{(t)})^{-1} (Q^{(t-1)})^{-1} A^{(t-1)} Q^{(t-1)} Q^{(t)}$$
$$= \quad (Q^{(t)})^{-1} (Q^{(t-1)})^{-1} ... (Q^{(2)})^{-1} (Q^{(1)})^{-1} A Q^{(1)} Q^{(2)} ... Q^{(t-1)} Q^{(t)}$$
$$= \quad (Q^{(1)} Q^{(2)} ... Q^{(t-1)} Q^{(t)})^{-1} A Q^{(1)} Q^{(2)} ... Q^{(t-1)} Q^{(t)}$$
$$= \quad Q^{-1} A Q$$
$$where \quad Q \quad = \quad Q^{(1)} Q^{(2)} ... Q^{(t-1)} Q^{(t)} \qquad \textit{product of orthonormal is orthonormal}$$
$$A \quad = \quad Q A^{(t+1)} Q^{-1} \qquad \textit{which is true for all t}$$
$$= \quad Q A^{(\infty)} Q^{-1}$$
$$= \quad Q U Q^{-1} \qquad \textit{see remark (We cannot conclude here as U is not diagonal.)}$$

$$\textit{eigen value of } A \quad = \quad \textit{eigen value of } QUQ^{-1}$$
$$= \quad \textit{eigen value of } U \qquad \textit{see property 9}$$
$$= \quad (u_{n,n} : n \in [1, N]) \qquad \textit{see property 5}$$

*Remark :*
*Prove that $A^{(t)}$ becomes an upper triangular matrix as t tends to infinity for QR algorithm. Please read 'Orthogonal Bases and the QR Algorithm', written by Peter J Oliver.*

## Singular value decomposition (SVD)

Singular value decomposition decomposes a rectangular matrix A into a left singular vector basis U, a right singular vector basis V and a scaling matrix S, such that U and V are both orthonormal, while S is diagonal. Note that eigenbasis in eigen decomposition is not necessarily orthonormal, yet both left singular basis and right singular basis in singular value decomposition are always orthonormal.

| For row vectors in A | For column vectors in A |
|---|---|
| $A \quad = \quad USV^T$ | $A \quad = \quad V^T SU$ |

where $U = \begin{bmatrix} U_1 \\ U_2 \\ U_3 \\ ... \\ U_N \end{bmatrix}$  ,   where $U = \begin{bmatrix} U_1 & U_2 & U_3 & ... & U_N \end{bmatrix}$

$U_n \quad = \quad 1{\times}N$ row matrix      $U_n \quad = \quad N{\times}1$ column matrix

$V = \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ ... \\ V_M \end{bmatrix}$  ,   $V = \begin{bmatrix} V_1 & V_2 & V_3 & ... & V_M \end{bmatrix}$

$V_n \quad = \quad 1{\times}M$ row matrix      $V_n \quad = \quad M{\times}1$ column matrix

$S = \begin{bmatrix} s_1 & 0 & 0 & ... & 0 \\ 0 & s_2 & 0 & ... & 0 \\ 0 & 0 & s_3 & ... & 0 \\ ... & ... & ... & ... & ... \\ 0 & 0 & 0 & ... & s_M \\ 0 & 0 & 0 & ... & 0 \\ ... & ... & ... & ... & ... \\ 0 & 0 & 0 & ... & 0 \end{bmatrix}$

$S = \begin{bmatrix} s_1 & 0 & 0 & ... & 0 & 0 & ... & 0 \\ 0 & s_2 & 0 & ... & 0 & 0 & ... & 0 \\ 0 & 0 & s_3 & ... & 0 & 0 & ... & 0 \\ ... & ... & ... & ... & ... & 0 & ... & 0 \\ 0 & 0 & 0 & ... & s_M & 0 & ... & 0 \end{bmatrix}$

$= \quad N{\times}M$ matrix (if N>M)      $= \quad M{\times}N$ matrix (if N>M)

*Note : Matrix A in the LHS case and in the RHS case are related by a transpose, while matrices (U,V) in the LHS case and in the RHS case are also related by a transpose. Hence the **LHS case and the RHS case are identical**.*

### Interpretation 1 (linear transformation)
For the LHS case above, SVD can be interpreted either as generating row data matrix A by applying dimensional scaling S, followed by rotation V on orthonormal row data matrix U. For the RHS case above, SVD can be interpreted either as generating column data matrix A by applying dimensional scaling S, followed by rotation V on orthonormal column data matrix U. For linear transformation interpretation, matrices A, U, V are treated as row data matrices for LHS case (or as column data matrices for RHS case), however we cannot find the null space with this interpretation, so we go for another interpretation.

### Interpretation 2 (linear combination)
When N<M, there exists null space, which can be only found using linear combination interpretation. We rewrite $V^T$ as row data matrix, i.e. $Z_n$ denotes the nth row of $V^T$.

$$V^T = \begin{bmatrix} Z_1 \\ Z_2 \\ Z_3 \\ \\ Z_M \end{bmatrix} = \begin{bmatrix} Z' \\ Z'' \end{bmatrix} \qquad \text{where Z' is K}{\times}\text{M and Z'' is (M-K)}{\times}\text{M}$$

*(Case 1) In the LHS case, row data matrix A is treated as the linear combination of row data matrix Z with weight US, Z'' provides the null space of row data matrix A, as since $AZ_n{}^T = 0$ for all $Z_n \in Z''$.*

$$A = U \begin{bmatrix} s_1 & 0 & 0 & ... & 0 & 0 & ... & 0 \\ 0 & s_2 & 0 & ... & 0 & 0 & ... & 0 \\ 0 & 0 & s_3 & ... & 0 & 0 & ... & 0 \\ ... & ... & ... & ... & ... & 0 & ... & 0 \\ 0 & 0 & 0 & ... & s_N & 0 & ... & 0 \end{bmatrix} \begin{bmatrix} Z_1 \\ Z_2 \\ Z_3 \\ ... \\ Z_N \\ Z_{N+1} \\ ... \\ Z_M \end{bmatrix} \Big\} Z''$$

$$\Rightarrow \quad AZ_n^T \quad = \quad U \begin{bmatrix} s_1 Z_1 Z_n^T \\ s_2 Z_2 Z_n^T \\ s_3 Z_3 Z_n^T \\ ... \\ s_N Z_N Z_n^T \\ 0 \\ ... \\ 0 \end{bmatrix} \quad = \quad U \begin{bmatrix} 0 \\ 0 \\ ... \\ s_n \\ ... \\ 0 \\ 0 \end{bmatrix} \quad \textit{if } n \le N$$

$$= \quad U \begin{bmatrix} 0 \\ 0 \\ 0 \\ ... \\ 0 \end{bmatrix} \quad \textit{if } n > N \qquad \textit{By definition, Z" is null space.}$$

*(Case 2) In the RHS case, column data matrix A is treated as the linear combination of column data matrix Z with weight $SU^T$, Z" provides the null space of column data matrix A, as since $Z_n A = 0$ for all $Z_n \in Z"$.*

$$A \quad = \quad [Z_1^T \ Z_2^T \ Z_3^T \ ... \ Z_N^T \ \underbrace{Z_{N+1}^T \ ... \ Z_M^T}_{Z"^T}] \begin{bmatrix} s_1 & 0 & 0 & ... & 0 \\ 0 & s_2 & 0 & ... & 0 \\ 0 & 0 & s_3 & ... & 0 \\ ... & ... & ... & ... & ... \\ 0 & 0 & 0 & ... & s_N \\ 0 & 0 & 0 & ... & 0 \\ ... & ... & ... & ... & ... \\ 0 & 0 & 0 & ... & 0 \end{bmatrix} U^T$$

$$\Rightarrow \quad Z_n A \quad = \quad [s_1 Z_n Z_1^T \ \ s_2 Z_n Z_2^T \ \ s_3 Z_n Z_3^T \ \ ... \ \ s_N Z_n Z_N^T \ \ 0 \ \ ... \ \ 0] \, U^T$$

$$= \quad [0 \ 0 \ ... \ s_n \ ... \ 0 \ 0] \, U^T \qquad \textit{if } n \le N$$

$$= \quad [0 \ 0 \ 0 \ ... \ 0] \, U^T \qquad \textit{if } n > N \qquad \textit{By definition, Z" is null space.}$$

### Implementation of SVD
*Singular value decomposition is closely related to eigen decomposition.*

| | | | | |
|---|---|---|---|---|
| *given* | $A$ | $=$ | $USV^T$ | *(S is sorted)* |
| | $AA^T$ | $=$ | $X\Lambda X^{-1}$ | *($\Lambda$ is sorted)* |
| | $A^T A$ | $=$ | $Y\Delta Y^{-1}$ | *($\Delta$ is sorted)* |
| *then* | $U$ | $=$ | $X$ | |
| | $V$ | $=$ | $Y$ | |
| | $s_{n,n}$ | $=$ | $\sqrt{\Lambda_{n,n}}$ | |
| | | $=$ | $\sqrt{\Delta_{n,n}}$ | |

*(Proof)*

$$
\begin{aligned}
AA^T &= (USV^T)(USV^T)^T \\
&= (USV^T)(VS^T U^T) \\
&= USS^T U^T && \textit{(V is orthonormal)} \\
&= USS^T U^{-1} && \textit{(U is orthonormal)} \\
&= X\Lambda X^{-1}
\end{aligned}
$$
$$\Rightarrow \quad X = U$$
$$\textit{and} \quad \Lambda_{n,n} = s_{n,n}^2$$

$$
\begin{aligned}
A^T A &= (USV^T)^T (USV^T) \\
&= (VS^T U^T)(USV^T) \\
&= VS^T SV^T && \textit{(U is orthonormal)} \\
&= VS^T SV^{-1} && \textit{(V is orthonormal)} \\
&= Y\Delta Y^{-1}
\end{aligned}
$$
$$\Rightarrow \quad Y = V$$
$$\textit{and} \quad \Delta_{n,n} = s_{n,n}^2$$

*Hence singular value decomposition can be implemented by eigen decomposition.*

**Application : solving AX = B**

*Solving AX=B by matrix inverse of $A^TA$ or Cramer's rule is slow, $O(N^3)$. A more efficient method is to decompose either the square matrix $A^TA$ or the rectangular matrix A into product of triangular matrices, and then apply (1) forward substitution for lower triangular matrix and (2) backward substitution for upper triangular matrix to solve for the unknown X. Suppose A is a N×M matrix (i.e. N row vectors in $M^{th}$ dimensional space), while B is a N×1 matrix (i.e. N observations) :*

$$
\begin{aligned}
AX &= B &&\text{for QR decomposition and SVD decomposition on matrix A} \\
A^T AX &= A^T B &&\text{for LU decomposition and Cholesky decomposition on matrix } A^TA
\end{aligned}
$$

(1)
$$
\begin{aligned}
A^T AX &= A^T B \\
LUX &= A^T B &&\text{by LU decomposition} \\
UX &= L^{-1}A^T B &&\text{by forward substitution} &&\text{(don't do matrix inverse)} \\
X &= U^{-1}L^{-1}A^T B &&\text{by backward substitution} &&\text{(don't do matrix inverse)}
\end{aligned}
$$

(2)
$$
\begin{aligned}
A^T AX &= A^T B \\
LL^T X &= A^T B &&\text{by Cholesky decomposition} \\
L^T X &= L^{-1}A^T B &&\text{by forward substitution} &&\text{(don't do matrix inverse)} \\
X &= (L^T)^{-1}L^{-1}A^T B &&\text{by backward substitution} &&\text{(don't do matrix inverse)}
\end{aligned}
$$

(3)
$$
\begin{aligned}
AX &= B \\
QRX &= B &&\text{by QR decomposition} \\
RX &= Q^T B &&\text{inverse of Q is its transpose} \\
X &= R^{-1}Q^T B &&\text{by backward substitution} &&\text{(don't do matrix inverse)}
\end{aligned}
$$

(4)
$$
\begin{aligned}
AX &= B \\
USV^T X &= B &&\text{by SVD decomposition} \\
SV^T X &= U^T B &&\text{inverse of U is its transpose} \\
V^T X &= S^{-1}U^T B &&\text{inverse of S is diagonal matrix with element-wise inverse} \\
X &= VS^{-1}U^T B &&\text{inverse of V is its transpose}
\end{aligned}
$$

(5)
$$
\begin{aligned}
A^T AX &= A^T B \\
X &= (A^T A)^{-1}(A^T B) &&\text{by direct matrix inverse of } A^TA \text{, which is } O(N^3)
\end{aligned}
$$

(6)
$$
\begin{aligned}
A^T AX &= A^T B \\
X &= \ldots &&\text{by Cramer's rule, which is again } O(N^3)
\end{aligned}
$$

*Remark : Forward substitution is used to solve LX = C, while backward substitution is used to solve UX = C.*

| *Forward substitution* | *Backward substitution* |
|---|---|

$$
\begin{bmatrix}
l_{1,1} & 0 & 0 & \ldots & 0 \\
l_{2,1} & l_{2,2} & 0 & \ldots & 0 \\
l_{3,1} & l_{3,2} & l_{3,3} & \ldots & 0 \\
\ldots & \ldots & \ldots & \ldots & \ldots \\
l_{M,1} & l_{M,2} & l_{M,3} & \ldots & l_{M,M}
\end{bmatrix}
\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \ldots \\ x_M \end{bmatrix}
=
\begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ \ldots \\ c_M \end{bmatrix}
\quad and \quad
\begin{bmatrix}
u_{1,1} & u_{1,2} & u_{1,3} & \ldots & u_{1,M} \\
0 & u_{2,2} & u_{1,3} & \ldots & u_{2,M} \\
0 & 0 & u_{3,3} & \ldots & u_{3,M} \\
\ldots & \ldots & \ldots & \ldots & \ldots \\
0 & 0 & 0 & \ldots & u_{M,M}
\end{bmatrix}
\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \ldots \\ x_M \end{bmatrix}
=
\begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ \ldots \\ c_M \end{bmatrix}
$$

$$
\sum_{k=1}^{m} l_{m,k}x_k = c_m \qquad and \qquad \sum_{k=m}^{M} u_{m,k}x_k = c_m
$$

$$
l_{m,m}x_m = c_m - \sum_{k=1}^{m-1} l_{m,k}x_k \qquad and \qquad u_{m,m}x_m = c_m - \sum_{k=m+1}^{M} u_{m,k}x_k
$$

$$
x_m = \left(c_m - \sum_{k=1}^{m-1} l_{m,k}x_k\right)\frac{1}{l_{m,m}} \qquad and \qquad x_m = \left(c_m - \sum_{k=m+1}^{M} u_{m,k}x_k\right)\frac{1}{u_{m,m}}
$$

*Bootstrap the value of $x_m$ for m starting from 1 to M for forward substitution in LX = C, and bootstrap the value of $x_m$ for m starting from M to 1 for backward substitution in UX = C. This bootstrapping technique is useful, we have used it once in deriving Cholesky decomposition.*

**Application : finding determinant for square matrix A**

*Finding determinant by applying Leibniz formula or Laplace expansion is slow, O(N³). A more efficient method is to decompose the square matrix into product of triangular matrices, and then apply (1) the property that determinant of matrix product equals to product of individual's determinant, and (2) determinant of triangular matrix is the product of its diagonal elements. Suppose A is a N×N square matrix :*

(1) $\det(A)$ = $\det(LU)$      *by LU decomposition*

= $\det(L)\det(U)$

= $\prod_{n=1}^{N} l_{n,n} \prod_{n=1}^{N} u_{n,n}$

(2) $\det(A)$ = $\det(LL^T)$      *by Cholesky decomposition*

= $\det(L)\det(L^T)$

= $(\prod_{n=1}^{N} l_{n,n})^2$

(3) $\det(A)$ = $\det(QR)$      *by QR decomposition*

= $\det(Q)\det(R)$

= $\det(R)$      *since determinant of orthonormal basis = volume of unit hypercube = 1*

= $\prod_{n=1}^{N} r_{n,n}$

(4) $\det(A)$ = $\det(X\Lambda X^{-1})$      *by eigen decomposition*

= $\det(X)\det(\Lambda)\det(X^{-1})$

= $\det(\Lambda)$      *since determinant of orthonormal basis = volume of unit hypercube = 1*

= $\prod_{n=1}^{N} \lambda_n$

(5) $\det(A)$ = $\det(USV^T)$      *by SVD decomposition*

= $\det(U)\det(S)\det(V^T)$

= $\det(S)$      *since determinant of orthonormal basis = volume of unit hypercube = 1*

= $\prod_{n=1}^{N} s_{n,n}$

(6) $\det(A)$ = $\sum_{\sigma \in S_N} \text{sgn}(\sigma) \prod_{n=1}^{N} a_{\sigma(n),n}$      *by Leibniz formula*

(7) $\det(A)$ = $\sum_{q=1}^{N} a_{p,q} C_{p,q}$      *by Laplace expansion (where $C_{p,q}$ is cofactor and $p \in [1,N]$ )*

**Application : principal component analysis**

*Principal component analysis can be solved by eigen decomposition or singular value decomposition. Given a N×M row data matrix A, which has* **zero mean**, *our objective is to find an orthonormal linear transformation P, so that B=AP is a N×K row data matrix, achieving objectives (1) diagonalized covariance matrix and (2) largest variance along the 1st coordinate, second largest variance along the 2nd coordinate … etc. Therefore A lies in M dimensional space, while B lies in K dimensional space (K≤M), in other words, B extracts only the K most important dimensions. P is M×K matrix, it denotes a mapping $\Re^M \rightarrow \Re^K$, it can be regarded as K column vectors $P_k$, each lies in M dimensional space, each $P_k$ is known as the principal component of A.*

$B$ = $AP$

*Covariance in transformed domain is :*

$\text{cov}_B$ = $\dfrac{1}{N} B^T B$

= $\dfrac{1}{N}(AP)^T(AP)$

= $P^T(\dfrac{1}{N}A^T A)P$

= $P^T(\text{cov}_A)P$

*where*   $\text{cov}_A$ = *M×M matrix*

$\text{cov}_B$ = *K×K matrix*

*The rest can be implemented by either (1) eigen decomposition or (2) singular value decomposition.*

*Method 1 : by eigen decomposition*
*By applying eigen decomposition on $cov_A$, we have :*

| | | | |
|---|---|---|---|
| $cov_B$ | $=$ | $P^T(X^T \Lambda X)P$ | *where rows in X are eigenvectors, both X and $\Lambda$ are M×M matrices* |
| | $=$ | $P'^T \Lambda P'$ | *where $P' = XP$ product of orthonormal basis is also orthonormal* |
| | $=$ | $\Lambda$ | *where $P = X^T$* |

*Hence by setting P as the transpose of X (where X is row eigenvector matrix), then we can diagonalize the covariance matrix of B. Since columns of P are the principal components, while columns of P equal to the rows of X, thus eigenbasis gives the set of principal axes for row data matrix A. We have completed the first part of the proof, yet we still need to prove why the 1st component of B has the largest variance, suppose K = 1, then $cov_B$ is a scalar, we have :*

| | | | |
|---|---|---|---|
| $\arg\max\limits_{P, P^T P=I} cov_B$ | $=$ | $\arg\max\limits_{P}(cov_B - \eta(P^T P - I))$ | *where $\eta$ is Lagrange multiplier* |
| | $=$ | $\arg\max\limits_{P}(P^T(cov_A)P - \eta(P^T P - I))$ | |
| $2P^T(cov_A) - 2\eta P^T$ | $=$ | $0$ | *take derivative wrt P and set it to zero* |
| $P^T(cov_A)$ | $=$ | $\eta P^T$ | |
| $(cov_A)P$ | $=$ | $\eta P$ | *take transpose, and $cov_A$ is symmetric* |

*which is the definition of eigenvector, with P as the column eigenvector, and $\eta$ as the eigenvalue. (Interesting!!! Lagrange multiplier is exactly the eigenvalue.) How about the 2nd component of B? We can repeat the above procedures, with K = 2, and maximize $(cov_B)_{2,2}$ under constrain that P cannot be the eigenvector with the largest eigenvalue, we will then end up with exactly the same formulae above. Since P cannot be the eigenvector with the largest eigenvalue, P can only be the eigenvector with the 2nd largest eigenvalue.*

*Method 2 : by singular value decomposition*
*Since principal component equals to eigenbasis of $A^TA$, and the latter equals to the right singular basis of A, we have :*

| | | |
|---|---|---|
| $A$ | $=$ | $USV^T$ |
| $V_m$ | $=$ | *principal components of A* |
| | $=$ | *1×M row vector* |

## Application : low rank approximation
*Low rank approximation of a large data matrix is a minimization problem which finds a N×M matrix A' having limited rank K to approximate a given N×M matrix A, under a particular matrix norm definition. The basic one is Frobenius norm.*

| | | | |
|---|---|---|---|
| $A'_{opt}$ | $=$ | $\arg\min\limits_{A'} \|A - A'\|_F$ | *such that $rank(A') \le K$* |

*For Frobenius norm, there is an analytic solution in terms of singular value decomposition of A.*

| | | | |
|---|---|---|---|
| $A'_{opt}$ | $=$ | $\arg\min\limits_{A'} \|A - A'\|_F$ | |
| | $=$ | $\arg\min\limits_{A'} \sqrt{\sum_{n=1}^{\min(N,M)} s_{Dn}^2}$ | *where $s_{Dn}$ are the singular values of $A - A' = U_D S_D V_D^T$* |

*This function can be minimized by :*

| | | | | |
|---|---|---|---|---|
| $A$ | $=$ | $USV^T$ | $= Udiag(s_1, s_2, ..., s_N)V^T$ | *where $s_n$ are the singular values of A in descending order* |
| $A'_{opt}$ | $=$ | $US'V^T$ | $= Udiag(s_1, s_2, ..., s_K, 0, 0, ..., 0)V^T$ | *thus $A - A' = U(S - S')V^T$, i.e. $U = U_D$ and $V = V_D$* |

*The objective becomes :*

| | | |
|---|---|---|
| $\min\limits_{A'} \|A - A'\|_F$ | $=$ | $\sqrt{\sum_{n=1}^{\min(N,M)} s_{Dn}^2}$ |
| | $=$ | $\sqrt{\sum_{n=K+1}^{\min(N,M)} s_n^2}$ |

**Reference**
[1]   Orthogonal Bases and the QR Algorithm, Peter J Oliver, University of Minnesota.