

# Matrix - Basic

## Definitions

Given  $N \times M$  matrix  $A$ , then  $M \times N$  matrix  $B$  is :

- the transpose of  $A$  if  $B_{m,n} = A_{n,m}$
- the conjugate transpose of  $A$  if  $B_{m,n} = \overline{A_{n,m}}$  (i.e. complex conjugate)
- the left inverse of  $A$  if  $BA = I$  (note : left inverse is not denoted by  $A^{-1}$ )
- the right inverse of  $A$  if  $AB = I$  (note : right inverse is not denoted by  $A^{-1}$ )

Given  $N \times N$  matrix  $A$ , then it is :

- orthogonal if  $AA^T = I$
- symmetric if  $A^T = A$
- Hermitian if  $A^T = \overline{A}$
- positive definite if  $XAX^T > 0$  for all non zero  $1 \times N$  row matrix  $X$
- positive semi definite if  $XAX^T \geq 0$  for all non zero  $1 \times N$  row matrix  $X$
- negative definite if  $XAX^T < 0$  for all non zero  $1 \times N$  row matrix  $X$
- negative semi definite if  $XAX^T \leq 0$  for all non zero  $1 \times N$  row matrix  $X$
- symmetric positive definite if  $A$  is both symmetric and positive definite
- symmetric positive semi definite if  $A$  is both symmetric and positive semi definite
- symmetric negative definite if  $A$  is both symmetric and negative definite
- symmetric negative semi definite if  $A$  is both symmetric and negative semi definite

Example that is positive definite but asymmetric :

$$[x \ y] \begin{bmatrix} +1 & +1 \\ -1 & +1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = x^2 + xy - xy + y^2 = x^2 + y^2 > 0$$

Please note that :

- product of diagonal matrices is also diagonal,
- product of upper triangular matrices is also upper triangular,
- product of lower triangular matrices is also lower triangular,
- product of orthonormal basis is also orthonormal basis, (read the proof in later section)
- product of symmetric matrix is **not necessarily symmetric**,
- product of positive definite matrix is **not necessary positive definite**.

Given  $N \times N$  matrix  $A$ , then :

- minor of  $A$  is  $M_{n,m} = \det(A_{n,m})$  where  $A_{n,m}$  is submatrix without row  $n$  & column  $m$
- cofactor of  $A$  is  $C_{n,m} = (-1)^{n+m} M_{n,m}$
- cofactor matrix of  $A$  is  $C = (C_{n,m})_{n,m \in [1,N]}$
- adjugate matrix of  $A$  is  $\text{adj}(A) = C^T$

Given  $N \times N$  matrix  $A$ , then the followings are equivalent.

- $A$  is invertible.
- $A$  is non singular.
- $A$  has non zero determinant.
- $A$  can be written as a row echelon form.
- $A$  can be written as a product of elementary matrices.
- There is exactly one solution for  $AX = B$ .
- There is exactly one solution for  $AX = 0$ , which is trivial solution.

If  $A$  is linearly dependent, then there exists (inf) non trivial solution to  $AX = 0$ .

If  $A$  is linearly independent, then there exists only trivial solution to  $AX = 0$ .

## Property of transpose and inverse

$$\begin{array}{ll} (A^T)^T &= A \\ (A+B)^T &= A^T + B^T \\ (AB)^T &= B^T A^T \\ (cA)^T &= c(A^T) \\ \det(A^T) &= \det(A) \\ (A^T)^{-1} &= (A^{-1})^T \\ \text{proof } (A^{-1})^T A^T &= (AA^{-1})^T = I \end{array} \quad \begin{array}{ll} (A^{-1})^{-1} &= A \\ (A+B)^{-1} &= A^{-1} - (I + A^{-1}B)^{-1} A^{-1} B A^{-1} \quad (\text{proof in the next section}) \\ (AB)^{-1} &= B^{-1} A^{-1} \\ (cA)^{-1} &= A^{-1} / c \\ \det(A^{-1}) &= 1/\det(A) \end{array}$$

### Inverse of matrix sum

$$(A+B)^{-1} = A^{-1} - (I + A^{-1}B)^{-1}A^{-1}BA^{-1}$$

(Proof) We start with :

$$\begin{aligned} (I+P)^{-1} &= (I+P)^{-1}(I+P-P) \\ &= (I+P)^{-1}(I+P) - (I+P)^{-1}P \\ &= I - (I+P)^{-1}P \\ (A+B)^{-1} &= (AI + AA^{-1}B)^{-1} \\ &= (A(I + A^{-1}B))^{-1} \\ &= (I + A^{-1}B)^{-1}A^{-1} \\ &= (I - (I + A^{-1}B)^{-1}A^{-1}B)A^{-1} \quad \text{by putting } P = A^{-1}B \text{ into } I - (I+P)^{-1}P \\ &= A^{-1} - (I + A^{-1}B)^{-1}A^{-1}BA^{-1} \end{aligned}$$

### Inverse of matrix block

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1} & -A^{-1}B(D - CA^{-1}B)^{-1} \\ -(D - CA^{-1}B)^{-1}CA^{-1} & (D - CA^{-1}B)^{-1} \end{bmatrix}$$

(Proof) Suppose we have :

$$\begin{aligned} \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix} &= \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1} & -A^{-1}B(D - CA^{-1}B)^{-1} \\ -(D - CA^{-1}B)^{-1}CA^{-1} & (D - CA^{-1}B)^{-1} \end{bmatrix} \\ X_{11} &= AA^{-1} + AA^{-1}B(D - CA^{-1}B)^{-1}CA^{-1} - B(D - CA^{-1}B)^{-1}CA^{-1} \\ &= I + B(D - CA^{-1}B)^{-1}CA^{-1} - B(D - CA^{-1}B)^{-1}CA^{-1} \\ &= I \\ X_{12} &= -AA^{-1}B(D - CA^{-1}B)^{-1} + B(D - CA^{-1}B)^{-1} \\ &= -B(D - CA^{-1}B)^{-1} + B(D - CA^{-1}B)^{-1} \\ &= 0 \\ X_{21} &= CA^{-1} + CA^{-1}B(D - CA^{-1}B)^{-1}CA^{-1} - D(D - CA^{-1}B)^{-1}CA^{-1} \\ &= CA^{-1} + (CA^{-1}B - D)(D - CA^{-1}B)^{-1}CA^{-1} \\ &= CA^{-1} - CA^{-1} \\ &= 0 \\ X_{22} &= -CA^{-1}B(D - CA^{-1}B)^{-1} + D(D - CA^{-1}B)^{-1} \\ &= (D - CA^{-1}B)(D - CA^{-1}B)^{-1} \\ &= I \end{aligned}$$

### Sherman Morrison formula

Sherman Morrison formula states that if  $u$  and  $v$  are row vectors,  $u^T v$  is the outer product, then :

$$(A + \beta u^T v)^{-1} = A^{-1} - \frac{\beta A^{-1} u^T v A^{-1}}{1 + \beta v A^{-1} u^T} \quad \text{Note : } 1 + \beta v A^{-1} u^T \text{ is a scalar.}$$

(Proof) By direct multiplication :

$$\begin{aligned} (A + \beta u^T v) \left( A^{-1} - \frac{\beta A^{-1} u^T v A^{-1}}{1 + \beta v A^{-1} u^T} \right) &= (A + \beta u^T v) A^{-1} - (A + \beta u^T v) \frac{\beta A^{-1} u^T v A^{-1}}{1 + \beta v A^{-1} u^T} \\ &= I + \beta u^T v A^{-1} - \frac{\beta u^T v A^{-1} + \beta^2 u^T v A^{-1} u^T v A^{-1}}{1 + \beta v A^{-1} u^T} \\ &= I + \beta u^T v A^{-1} - \frac{\beta u^T (1 + \beta v A^{-1} u^T) v A^{-1}}{1 + \beta v A^{-1} u^T} \\ &= I + \beta u^T v A^{-1} - \frac{(1 + \beta v A^{-1} u^T) \beta u^T v A^{-1}}{1 + \beta v A^{-1} u^T} \quad \text{Trick : } 1 + \beta v A^{-1} u^T \text{ is a scalar.} \\ &= I + \beta u^T v A^{-1} - \beta u^T v A^{-1} \\ &= I \quad \text{This is useful in recursive least square.} \end{aligned}$$

## Matrix differentiation

### Jacobian, Hessian and Taylor series

Given function  $F : \mathbb{R}^M \rightarrow \mathbb{R}^N$ ,

$$F(X) = \begin{bmatrix} f_1(X) \\ f_2(X) \\ f_3(X) \\ \dots \\ f_N(X) \end{bmatrix} = \begin{bmatrix} f_1(x_1, x_2, \dots, x_M) \\ f_2(x_1, x_2, \dots, x_M) \\ f_3(x_1, x_2, \dots, x_M) \\ \dots \\ f_N(x_1, x_2, \dots, x_M) \end{bmatrix} \quad \text{where } X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_M \end{bmatrix}$$

Jacobian is defined as :

$$J(F(X)) = \begin{bmatrix} \partial f_1 / \partial x_1 & \partial f_1 / \partial x_2 & \partial f_1 / \partial x_3 & \dots & \partial f_1 / \partial x_M \\ \partial f_2 / \partial x_1 & \partial f_2 / \partial x_2 & \partial f_2 / \partial x_3 & \dots & \partial f_2 / \partial x_M \\ \partial f_3 / \partial x_1 & \partial f_3 / \partial x_2 & \partial f_3 / \partial x_3 & \dots & \partial f_3 / \partial x_M \\ \dots & \dots & \dots & \dots & \dots \\ \partial f_N / \partial x_1 & \partial f_N / \partial x_2 & \partial f_N / \partial x_3 & \dots & \partial f_N / \partial x_M \end{bmatrix} \quad \text{which is 1st order derivative}$$

When  $M=1$ , Hessian is defined as :

$$H(f(X)) = \begin{bmatrix} \partial^2 f / \partial x_1 \partial x_1 & \partial^2 f / \partial x_1 \partial x_2 & \partial^2 f / \partial x_1 \partial x_3 & \dots & \partial^2 f / \partial x_1 \partial x_M \\ \partial^2 f / \partial x_2 \partial x_1 & \partial^2 f / \partial x_2 \partial x_2 & \partial^2 f / \partial x_2 \partial x_3 & \dots & \partial^2 f / \partial x_2 \partial x_M \\ \partial^2 f / \partial x_3 \partial x_1 & \partial^2 f / \partial x_3 \partial x_2 & \partial^2 f / \partial x_3 \partial x_3 & \dots & \partial^2 f / \partial x_3 \partial x_M \\ \dots & \dots & \dots & \dots & \dots \\ \partial^2 f / \partial x_M \partial x_1 & \partial^2 f / \partial x_M \partial x_2 & \partial^2 f / \partial x_M \partial x_3 & \dots & \partial^2 f / \partial x_M \partial x_M \end{bmatrix} \quad \text{which is 2nd order derivative}$$

Please note that Jacobian is  $N \times M$  matrix, while Hessian is  $M \times M$  symmetric matrix. Besides, don't confuse Hessian with Hermitian. When  $M=1$ , we can rewrite Taylor series in terms of Jacobian and Hessian as :

$$\begin{aligned} f(X + \Delta X) &= f(X) + J(f(X))\Delta X + \frac{1}{2}\Delta X^T H(f(X))\Delta X \\ J(f(X)) &= 1 \times M \text{ row matrix} \\ H(f(X)) &= M \times M \text{ matrix} \end{aligned}$$

### Two fundamental derivative formulae

Suppose  $X$  is a  $M \times 1$  column matrix,  $F$  is a  $N \times 1$  column matrix,  $G$  is a  $L \times 1$  column matrix.  $F$  and  $G$  are functions that map from  $M$  to  $N$  dimensional space and from  $M$  to  $L$  dimensional space respectively. Derivative of scalar  $s$  with respect to  $X$  and derivative of column vector  $v$  (with size  $N \times 1$ ) with respect to  $X$  are :

$$\begin{aligned} s &= F^T A G + B & \frac{ds}{dX} &= (F^T A)J_G + (AG)^T J_F = F^T A J_G + G^T A^T J_F & \text{formula for scalar} \\ v &= AG + B & \frac{dv}{dX} &= A J_G & \text{formula for vector} \\ \text{where :} & & \frac{dF}{dX} &= J_F & N \times M & F : \mathbb{R}^M \rightarrow \mathbb{R}^N & \text{dependent on } X \\ & & \frac{dG}{dX} &= J_G & L \times M & G : \mathbb{R}^M \rightarrow \mathbb{R}^L & \text{dependent on } X \\ & & \frac{dX}{dX} &= I & M \times M & & \\ & & \frac{dC}{dX} &= 0 & N \times M & C : \mathbb{R}^M \rightarrow \mathbb{R}^N & \text{independent on } X \end{aligned}$$

Here is a summary of specialised cases. The derivative of  $s$  follows the convention of Jacobian matrix, hence it is a  $1 \times M$  row matrix, however it is sometimes more convenient to be expressed as a  $M \times 1$  column matrix in optimization algorithms and regression algorithms, in that case, we need to take transpose on the result matrix.

	scalar $s$	$F$	$A$	$G$	$ds/dX$	non symmetric $A$	symmetric $A$ ( $N=L$ )	identity $A$ ( $N=L$ )
general	$F^T AG + B$	$N,1$	$N,L$	$L,1$	$1,M$	$F^T AJ_G + G^T A^T J_F$	$F^T AJ_G + G^T AJ_F$	$F^T J_G + G^T J_F$
$F = G$	$F^T AF + B$	$N,1$	$N,N$	$N,1$	$1,M$	$F^T (A + A^T) J_F$	$2F^T AJ_F$	$2F^T J_F$
$F = X$	$X^T AG + B$	$M,1$	$M,L$	$L,1$	$1,M$	$X^T AJ_G + G^T A^T$	$X^T AJ_G + G^T A$	$X^T J_G + G^T$
$G = X$	$F^T AX + B$	$N,1$	$N,M$	$M,1$	$1,M$	$F^T A + X^T A^T J_F$	$F^T A + X^T AJ_F$	$F^T + X^T J_F$
$F = G = X$	$X^T AX + B$	$M,1$	$M,M$	$M,1$	$1,M$	$X^T (A + A^T)$	$2X^T A$	$2X^T$
$F = I$	$AG + B$	-	$1,L$	$L,1$	$1,M$	$AJ_G$	-	-
$F = I, G = X$	$AX + B$	-	$1,M$	$M,1$	$1,M$	$A$	-	-
$G = I$	$F^T A + B$	$N,1$	$N,1$	-	$1,M$	$A^T J_F$	-	-
$G = I, F = X$	$X^T A + B$	$M,1$	$M,1$	-	$1,M$	$A^T$	-	-
	vector $v$		$A$	$G$	$dv/dX$	derivative		
general	$AG + B$		$N,L$	$L,1$	$N,M$	$AJ_G$	(the same as above when $N = 1$ )	
$G = X$	$AX + B$		$N,M$	$M,1$	$N,M$	$A$	(the same as above when $N = 1$ )	

where derivative  $ds/dX$  can either be expressed as a row matrix or a column matrix :

- use row matrix if you want to have the same convention as Jacobian
- use column matrix if you want to use it in algorithms (e.g. least square, Newton, Gauss Newton etc)

### Intuition of matrix product

Matrix can be classified into the following types :

- row data matrix
- column data matrix
- linear transformation
- linear combination

Suppose  $A$  is a data matrix with size  $Y \times X$ , then it denotes  $Y$  row vectors, each belongs to  $X$  dimensional space, if  $A$  is treated as a row data matrix, and it denotes  $X$  column vectors, each belongs to  $Y$  dimensional space, if  $A$  is treated as a column data matrix. In general,  $A$  can be interpreted as both row and column data matrix.

If  $A = Y \times X$  data matrix  
= row data matrix

If  $A = Y \times X$  data matrix  
= column data matrix

then  $A = \begin{bmatrix} A_1 \\ A_2 \\ A_3 \\ \dots \\ A_Y \end{bmatrix}$

$A = [A_1 \ A_2 \ A_3 \ \dots \ A_X]$

where  $A_y = 1 \times X$  row vector

where  $A_x = Y \times 1$  column vector

Here is the convention adopted in my documents, when  $A$  is a row data matrix, then it has a size of  $N \times M$ , when  $A$  is a column data matrix, then it has a size of  $M \times N$ , both represent  $N$  vectors in  $M$  dimensional space ( $N$  can be greater than or smaller than  $M$ ).

Please differentiate the following notations :

- row data matrix      rectangular matrix with one row as a vector (each row is a **row vector**)
- row matrix      rectangular matrix with one row only
- row vector      a row in a **row data matrix**
- column data matrix      rectangular matrix with one column as a vector (each column is a **column vector**)
- column matrix      rectangular matrix with one column only
- column vector      a column in a **column data matrix**

Please note that we **never** classify linear transformation (and linear combination) into row or column matrix. Different meanings for matrix product are classified as the following :

- product between a data matrix and a linear transformation (include : shearing, rotation, dimension scaling)
- product between a data matrix and a linear combination (include : data scaling, permutation, elementary addition)
- product between two data matrices as inner product
- product between two data matrices as covariance

If  $A = N \times M$  row data matrix

$$\text{then } A = \begin{bmatrix} A_1 \\ A_2 \\ A_3 \\ \dots \\ A_N \end{bmatrix}$$

where  $A_n = 1 \times M$  row vector

If  $A = M \times N$  column data matrix

$$A = [A_1 \ A_2 \ A_3 \ \dots \ A_N]$$

where  $A_n = M \times 1$  column vector

### Part 1 : Linear transformation $\mathbb{R}^M \rightarrow \mathbb{R}^K$

$$A' = AF$$

$$\text{then } A' = \begin{bmatrix} A'_1 \\ A'_2 \\ A'_3 \\ \dots \\ A'_N \end{bmatrix}$$

$A'_n = 1 \times K$  row vector  
 = linear transformation of **one**  $A_n$   
 $F = M \times K$  linear transformation  
 = performs mapping  $\mathbb{R}^M \rightarrow \mathbb{R}^K$

$$A' = FA$$

$$\text{then } A' = [A'_1 \ A'_2 \ A'_3 \ \dots \ A'_N]$$

$A'_n = K \times 1$  column vector  
 = linear transformation of **one**  $A_n$   
 $F = K \times M$  linear transformation  
 = performs mapping  $\mathbb{R}^M \rightarrow \mathbb{R}^K$

Special cases for both row data matrix  $A$  and column data matrix  $A$

- when  $K = M \Rightarrow$  perform shearing in  $\mathbb{R}^M \rightarrow \mathbb{R}^M$
- when  $K = M$  and  $F$  is orthonormal  $\Rightarrow$  perform rotation with no dimensional scaling
- when  $K = M$  and  $F$  is diagonal  $\Rightarrow$  perform dimension scaling with different scales  $f_{m,m}$  (remark 1)

If  $F$  is orthonormal, it can always be written as a matrix with cosines and sines (i.e. rotation).

For example, if  $F$  denotes a rotation  $\theta$  in  $xy$  plane, followed by a rotation  $\phi$  in  $yz$  plane, we have :

$$F = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \phi & -\sin \phi \\ 0 & \sin \phi & \cos \phi \end{bmatrix}$$

$$= \begin{bmatrix} \cos \theta & -\sin \theta \cos \phi & \sin \theta \sin \phi \\ \sin \theta & \cos \theta \cos \phi & \cos \theta \sin \phi \\ 0 & \sin \phi & \cos \phi \end{bmatrix}$$

### Part 2 : Linear combination

$$A' = WA$$

$$\text{then } A' = \begin{bmatrix} A'_1 \\ A'_2 \\ A'_3 \\ \dots \\ A'_K \end{bmatrix}$$

$A'_k = 1 \times M$  row vector  
 = linear combination of **all**  $A_n$   
 (weighted by  $w_{k,n}$ )  
 $W = K \times N$  linear combination

$$A' = AW$$

$$\text{then } A' = [A'_1 \ A'_2 \ A'_3 \ \dots \ A'_K]$$

$A'_k = M \times 1$  column vector  
 = linear combination of **all**  $A_n$   
 (weighted by  $w_{n,k}$ )  
 $W = N \times K$  linear combination

Special cases for both row data matrix  $A$  or column data matrix  $A$

- when  $K = N \Rightarrow$  generate  $N$  row (or column) vectors from  $N$  row (or column) vectors
- when  $K = N$  and  $W$  is diagonal  $\Rightarrow$  scaling  $N$  row (or column) vectors by  $w_{n,n}$  (remark 1)
- when  $K = N$  and  $W$  is permutation  $\Rightarrow$  equivalent to row (or column) permutation
- when  $K = N$  and  $W$  is elem addition  $\Rightarrow$  equivalent to elementary row (or column) addition
- when  $K \neq N$  and  $W$  is a delta  $\Rightarrow$  equivalent to picking a desired row (or column) from  $A$

*Exampe 1 :  $W$  = permutation matrix (please read 'permutation' document for definition of permutation matrix)*

$$\begin{aligned} A' &= P_{\sigma_1} P_{\sigma_2} P_{\sigma_3} \dots P_{\sigma_T} A &= \text{cascaded permutation } (\sigma_T \circ \dots \circ \sigma_3 \circ \sigma_2 \circ \sigma_1) \text{ on row data matrix } A \\ A' &= A P_{\sigma_T}^T \dots P_{\sigma_3}^T P_{\sigma_2}^T P_{\sigma_1}^T &= \text{cascaded permutation } (\sigma_T \circ \dots \circ \sigma_3 \circ \sigma_2 \circ \sigma_1) \text{ on column data matrix } A \end{aligned}$$

*Exampe 2 :  $W$  = elementary addition operation (please read 'elementary operation' document for definitions of  $E$  &  $F$ )*

$$\begin{aligned} A' &= E_T E_{T-1} \dots E_3 E_2 E_1 A &= \text{cascaded elementary row addition on row data matrix } A \\ A' &= A F_1 F_2 F_3 \dots F_{T-1} F_T &= \text{cascaded elementary column addition on column data matrix } A \end{aligned}$$

*Example 3 :  $W$  = delta matrix*

$$\begin{aligned} A' &= D A \\ &= \begin{bmatrix} \delta_{\pi(1)} \\ \delta_{\pi(2)} \\ \delta_{\pi(3)} \\ \dots \\ \delta_{\pi(K)} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \\ A_3 \\ \dots \\ A_N \end{bmatrix} = \begin{bmatrix} A_{\pi(1)} \\ A_{\pi(2)} \\ A_{\pi(3)} \\ \dots \\ A_{\pi(K)} \end{bmatrix} &= \text{picking the } k^{\text{th}} \text{ row from row data matrix } A \\ \delta_x &= [0 \dots 1 \dots 0] &= 1 \times N \text{ row matrix with value one at position } x, \text{ and value zero otherwise} \\ A' &= A D \\ &= [A_1 \ A_2 \ A_3 \ \dots \ A_N] \times \\ &\quad [\delta_{\pi(1)} \ \delta_{\pi(2)} \ \dots \ \delta_{\pi(K)}] \\ &= [A_{\pi(1)} \ A_{\pi(2)} \ \dots \ A_{\pi(K)}] &= \text{picking the } k^{\text{th}} \text{ column from column data matrix } A \\ \delta_x &= \begin{bmatrix} 0 \\ \dots \\ 1 \\ \dots \\ 0 \end{bmatrix} &= N \times 1 \text{ column matrix with value one at position } x, \text{ and value zero otherwise} \end{aligned}$$

*Remark 1 : Dimension scaling vs data scaling*

Both linear transformation and linear combination become scaling when  $F$  (or  $W$ ) is diagonal, yet they represent different scaling : the former means **dimension scaling**, while the latter means **data scaling**.

- scaling in linear transformation  $\Rightarrow f_{m,m}$  = scale for the  $m^{\text{th}}$  dimension in  $A$  (called dimension scaling)
- scaling in linear combination  $\Rightarrow w_{n,n}$  = scale for the  $n^{\text{th}}$  data in  $A$  (called data scaling)

*Remark 2 : Cascade of linear transformation and linear combination*

Both linear transformation and linear combination can be cascaded. Please note that linear transformation matrix and linear combination matrix are different in associativity for row data matrix  $A$  and column data matrix  $A$ . Suppose we perform  $T$  linear transformations in sequence, and if the dimension after the  $t^{\text{th}}$  linear transformation is  $M_t$ , then we have :

$$\begin{aligned} A' &= A F_1 F_2 F_3 \dots F_T && \text{for row data matrix } A \text{ (left associative)} && F_1 \text{ is } M \times M_1 \text{ and } F_t \text{ is } M_{t-1} \times M_t \\ A' &= F_T \dots F_3 F_2 F_1 A && \text{for column data matrix } A \text{ (right associative)} && F_1 \text{ is } M_1 \times M \text{ and } F_t \text{ is } M_t \times M_{t-1} \end{aligned}$$

Suppose we perform  $T$  linear transformations in sequence, and if the number of data after the  $t^{\text{th}}$  linear combination is  $N_t$ , then we have :

$$\begin{aligned} A' &= W_T \dots W_3 W_2 W_1 A && \text{for row data matrix } A \text{ (right associative)} && W_1 \text{ is } N_1 \times N \text{ and } W_t \text{ is } N_t \times N_{t-1} \\ A' &= A W_1 W_2 W_3 \dots W_T && \text{for column data matrix } A \text{ (left associative)} && W_1 \text{ is } N \times N_1 \text{ and } W_t \text{ is } N_{t-1} \times N_t \end{aligned}$$

### Part 3 : Inner product (projection)

Suppose  $A$  and  $B$  are both  $N \times M$  row data matrices, while  $W$  is a  $M \times M$  weight matrix (different weights for different dimensions), then we have inner product defined as :

$$\begin{aligned} X &= A W B^T = \langle A, B \rangle_W \\ x_{n_1, n_2} &= \sum_{m=1}^M w_m a_{n_1, m} b_{n_2, m} \end{aligned} \quad \text{if } W \text{ is diagonal}$$

Suppose  $A$  and  $B$  are both  $M \times N$  column data matrices, while  $W$  is a  $M \times M$  weight matrix (different weights for different dimensions), then we have inner product defined as :

$$\begin{aligned} X &= A^T W B = \langle A, B \rangle_W \\ x_{n_1, n_2} &= \sum_{m=1}^M w_m a_{m, n_1} b_{m, n_2} \end{aligned} \quad \text{if } W \text{ is diagonal}$$

**Inner product** is usually used to find **projection or orthogonality** between two vectors.

### Part 4 : Outer product (covariance matrix)

Suppose  $A$  and  $B$  are both  $N \times M$  row data matrices, while  $W$  is a  $N \times N$  weight matrix (different weights for different data or observations), then we have outer product defined as :

$$\begin{aligned} X &= A^T W B \\ x_{m_1, m_2} &= \sum_{n=1}^N w_n a_{n, m_1} b_{n, m_2} \end{aligned} \quad \text{if } W \text{ is diagonal}$$

Suppose  $A$  and  $B$  are both  $M \times N$  column data matrices, while  $W$  is a  $N \times N$  weight matrix (different weights for different data or observations), then we have outer product defined as :

$$\begin{aligned} X &= A W B^T \\ x_{m_1, m_2} &= \sum_{n=1}^N w_n a_{m_1, n} b_{m_2, n} \end{aligned} \quad \text{if } W \text{ is diagonal}$$

**Outer product** is usually used to find **covariance matrix** between two dimensions (or between two random variables).

### Part 5 : Replication and summation by all one matrix (a special case of linear combination)

Row matrix  $A$  with size  $1 \times M$  can be replicated to generate row data matrix  $B$  with size  $N \times M$  by  $B = lA$ , where  $l$  is  $N \times 1$  all one column matrix, and  $B_n = A$  for all  $n \in [1, N]$ . Column matrix  $A$  with size  $M \times 1$  can be replicated to generate column data matrix  $B$  with size  $M \times N$  by  $B = A l^T$ , where  $l^T$  is  $1 \times N$  all one row matrix, and  $B_n = A$  for all  $n \in [1, N]$ .

Given  $N \times M$  row data matrix  $A$ , the sum of row data is given by  $l^T A$ . Given  $M \times N$  column data matrix  $A$ , the sum of column data is given by  $A l$ . Given a square matrix weight matrix  $W$  (not necessarily diagonal), then  $l^T W l$  gives the sum of weight.

$$\begin{aligned} l^T A &= \sum_{n=1}^N A_n && \text{for row data matrix } A \\ A l &= \sum_{n=1}^N A_n && \text{for column data matrix } A \\ l^T W l &= \sum_{n=1}^N \sum_{m=1}^N w_{n, m} \end{aligned}$$

More about covariance matrix

Suppose A is a  $N \times M$  row data matrices, then we have covariance defined as :

$$\begin{aligned}
 C &= \frac{(A - \bar{A})^T W (A - \bar{A})}{\mathbf{1}^T W \mathbf{1}} && \text{if } W \text{ is diagonal, where } w_n = \text{weight of row vector } A_n \\
 c_{m_1, m_2} &= \text{cov}(m_1, m_2) && \text{covariance between dimension } m_1 \text{ and } m_2 \\
 \text{where } \bar{A} &= \frac{\mathbf{1}^T W A}{\mathbf{1}^T W \mathbf{1}} = \frac{\sum_{n=1}^N w_n A_n}{\sum_{n=1}^N w_n} && \text{average of all } A_n, \text{ which is a } 1 \times M \text{ row matrix}
 \end{aligned}$$

Suppose A is a  $M \times N$  column data matrices, then we have covariance defined as :

$$\begin{aligned}
 C &= \frac{(A - \bar{A} \mathbf{1}^T) W (A - \bar{A} \mathbf{1}^T)^T}{\mathbf{1}^T W \mathbf{1}} && \text{if } W \text{ is diagonal, where } w_n = \text{weight of column vector } A_n \\
 c_{m_1, m_2} &= \text{cov}(m_1, m_2) && \text{covariance between dimension } m_1 \text{ and } m_2 \\
 \text{where } \bar{A} &= \frac{A W \mathbf{1}}{\mathbf{1}^T W \mathbf{1}} = \frac{\sum_{n=1}^N w_n A_n}{\sum_{n=1}^N w_n} && \text{average of all } A_n, \text{ which is a } M \times 1 \text{ column matrix}
 \end{aligned}$$

Lets simplify.

$$\begin{aligned}
 C &= \frac{(A - \mathbf{1} \bar{A})^T W (A - \mathbf{1} \bar{A})}{\mathbf{1}^T W \mathbf{1}} \\
 &= \frac{A^T W A - A^T W \mathbf{1} \bar{A} - \bar{A}^T \mathbf{1}^T W A + \bar{A}^T \mathbf{1}^T W \mathbf{1} \bar{A}}{\mathbf{1}^T W \mathbf{1}} \\
 &= \frac{A^T W A - \bar{A}^T \mathbf{1}^T W \mathbf{1} \bar{A}}{\mathbf{1}^T W \mathbf{1}} && \text{since } A^T W \mathbf{1} \bar{A} = \bar{A}^T \mathbf{1}^T W A = \bar{A}^T \mathbf{1}^T W \mathbf{1} \bar{A} \quad (1) \\
 &= \frac{A^T W A}{\mathbf{1}^T W \mathbf{1}} - \bar{A}^T \bar{A} && \text{since } \mathbf{1}^T W \mathbf{1} \text{ is a scalar and thus } \bar{A}^T \mathbf{1}^T W \mathbf{1} \bar{A} = (\mathbf{1}^T W \mathbf{1})(\bar{A}^T \bar{A})
 \end{aligned}$$

Equation (1) can be proved by :

$$\begin{aligned}
 \bar{A} &= \frac{\mathbf{1}^T W A}{\mathbf{1}^T W \mathbf{1}} \Rightarrow \mathbf{1}^T W A = (\mathbf{1}^T W \mathbf{1}) \bar{A} \quad (2) \\
 \text{and } \bar{A}^T &= \frac{A^T W \mathbf{1}}{\mathbf{1}^T W \mathbf{1}} \Rightarrow A^T W \mathbf{1} = \bar{A}^T (\mathbf{1}^T W \mathbf{1}) \quad (3)
 \end{aligned}$$

Therefore (1) is proved :

$$\begin{aligned}
 \bar{A}^T \mathbf{1}^T W A &= \bar{A}^T \mathbf{1}^T W \mathbf{1} \bar{A} && \text{using (2)} \\
 A^T W \mathbf{1} \bar{A} &= \bar{A}^T \mathbf{1}^T W \mathbf{1} \bar{A} && \text{using (3)}
 \end{aligned}$$



### Orthonormal basis

What are independent, orthogonal and orthonormal? Suppose  $X$  and  $Y$  are two  $1 \times N$  row vectors, then

- $X$  and  $Y$  are independent if  $Y \neq cX$
- $X$  and  $Y$  are orthogonal if  $XY^T = 0$
- $X$  and  $Y$  are orthonormal if  $XY^T = 0$  and  $XX^T = YY^T = 1$

A set of independent vectors which spans a space  $\mathcal{H}^N$  is known as the **basis** of the space, while a set of orthonormal vectors which spans a space  $\mathcal{H}^N$  is known as the **orthonormal basis** of the space. Suppose  $Q$  is a  $N \times N$  matrix denoting orthonormal basis, then we have :

$$\begin{aligned} QQ^T &= I && \text{for row vectors in } Q && \text{(known as row orthonormal)} \\ Q^T Q &= I && \text{for column vectors in } Q && \text{(known as column orthonormal)} \end{aligned}$$

### Inverse of orthonormal

Suppose  $Q$  is a  $N \times N$  orthonormal basis, with row vectors, i.e.

$$\begin{aligned} Q &= \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \\ \vdots \\ Q_N \end{bmatrix} && \text{where } Q_n \text{ is row vector } \forall n \in [1, N] \\ QQ^T &= (Q_n Q_m^T)_{n,m \in [1, N]} \\ &= (\delta_{n,m})_{n,m \in [1, N]} \\ &= I \\ \Rightarrow Q^{-1} &= Q^T \end{aligned}$$

Suppose  $Q$  is a  $N \times N$  orthonormal basis, with column vectors, i.e.

$$\begin{aligned} Q &= [Q_1 \ Q_2 \ Q_3 \ \dots \ Q_N] && \text{where } Q_n \text{ is column vector } \forall n \in [1, N] \\ Q^T Q &= (Q_n^T Q_m)_{n,m \in [1, N]} \\ &= (\delta_{n,m})_{n,m \in [1, N]} \\ &= I \\ \Rightarrow Q^{-1} &= Q^T \end{aligned}$$

If  $Q$  is row orthonormal basis, then  $Q$  must be column orthonormal basis, and vice versa. The proof is simple.

$$\begin{aligned} QQ^{-1} &= Q^{-1}Q = I && \text{since } Q \text{ is orthonormal, let's put } Q^{-1} = Q^T \\ \Rightarrow QQ^T &= Q^T Q = I && \text{hence we have : row orthonormal} \Leftrightarrow \text{column orthonormal} \end{aligned}$$

Besides the inverse of orthonormal is also orthonormal, since :

$$\begin{aligned} (Q^{-1})(Q^{-1})^T &= (Q^{-1})(Q^T)^T && \text{since } Q \text{ is orthonormal, i.e. } Q^{-1} = Q^T \\ &= (Q^{-1})Q \\ &= I && \text{hence we have : } Q \text{ is orthonormal} \Leftrightarrow Q^{-1} \text{ is orthonormal} \end{aligned}$$

Hence these are equivalent :  $Q$  is row orthonormal,  $Q$  is column orthonormal, inverse  $Q$  is orthonormal. **(Note : we don't need to specify whether  $Q$  is row orthonormal or column orthonormal, as they are equivalent.)**

### Product between two different orthonormal basis

Suppose  $U$  and  $V$  are two orthonormal basis, let's prove that  $Q$  is also an orthonormal basis.

$$\begin{aligned} Q &= UV && \text{note : there is no transpose} \\ QQ^T &= UV(UV)^T \\ &= UVV^T U^T \\ &= UU^T && \text{since } V \text{ is orthonormal : } VV^T = I \\ &= I && \text{since } U \text{ is orthonormal : } UU^T = I \end{aligned}$$

Hence the product of an orthonormal basis with another orthonormal basis is also an orthonormal basis. This property will be used twice in deriving the QR algorithm.

### Projection matrix

Lets consider the column vector case (since it matches with the convention in least square). Suppose  $A$  is a  $M \times N$  column data matrix (i.e.  $N$  vectors in  $\mathbb{R}^M$ ), and  $B$  is a  $M \times 1$  column data matrix in the same space, then the projection of  $B$  on the space spanned by the columns of  $A$  is  $P_A B$ , where  $P_A$  is called the projection matrix :

$$\begin{aligned} \text{proj}_A(B) &= P_A B \\ P_A &= A(A^T A)^{-1} A^T \end{aligned}$$

The proof involves two steps : (1) prove that  $P_A B$  lies in the column span of  $A$  (i.e. show that there exists  $W$ , such that linear combination  $AW$  equals to  $P_A B$ ) and (2) vector  $B - \text{proj}_A(B)$  is orthogonal to all columns in  $A$  (i.e.  $A^T(B - \text{proj}_A(B)) = 0$ ).

$$\begin{aligned} (1) \quad P_A B &= A(A^T A)^{-1} A^T B \\ &= AW \end{aligned}$$

$$\text{where } W = (A^T A)^{-1} A^T B \Rightarrow P_A B \text{ lies in the column span of } A$$

$$\begin{aligned} (2) \quad A^T(B - P_A B) &= A^T(B - A(A^T A)^{-1} A^T B) \\ &= A^T B - A^T A(A^T A)^{-1} A^T B \\ &= A^T B - A^T B \\ &= 0 \end{aligned} \Rightarrow B - P_A B \text{ is orthogonal to all columns of } A$$

Besides, there is one more requirement on projection matrix : projection of projection is itself, i.e.

$$\begin{aligned} \text{proj}_A(\text{proj}_A(B)) &= \text{proj}_A(B) \\ P_A P_A \dots P_A B &= P_A P_A B = P_A B \\ \text{i.e. } P_A P_A \dots P_A &= P_A P_A = P_A \\ \Rightarrow P_A P_A \dots P_A &= \underbrace{(A(A^T A)^{-1} A^T)(A(A^T A)^{-1} A^T) \dots (A(A^T A)^{-1} A^T)}_{T} \\ &= (A(A^T A)^{-1} A^T A(A^T A)^{-1} A^T) \dots \underbrace{(A(A^T A)^{-1} A^T)}_{T-2} \\ &= \underbrace{(A(A^T A)^{-1} A^T) \dots (A(A^T A)^{-1} A^T)}_{T-1} = A(A^T A)^{-1} A^T \end{aligned}$$

In general  $A$  is not an orthogonal basis. What happens if  $A$  is an orthogonal basis? In this case,  $A^T A$  is a diagonal matrix :

$$\begin{aligned} A^T A &= \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N) & \text{which is } N \times N & \quad \text{where } \lambda_n = A_n^T A_n = \text{mag}(A_n) \\ (A^T A)^{-1} &= \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_N^{-1}) & \text{which is } N \times N \\ (A^T A)^{-1} A^T B &= \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_N^{-1}) A^T B & \text{which is } N \times 1 \\ ((A^T A)^{-1} A^T B)_n &= (A^T B)_n / \lambda_n \\ &= (A_n^T A_n)^{-1} (A^T B)_n \\ &= (A_n^T A_n)^{-1} (A_n^T B) \\ \text{where } A_n^T B &= \text{inner product} \\ (A_n^T A_n)^{-1} (A_n^T B) &= \text{inner product with normalization} \end{aligned}$$

Similarly, we can derive the projection for row vectors. Here is a summary.

### Summary

For row data matrix  $A$  and single row data  $B$ , then the projection of  $B$  on space spanned by the rows of  $A$  is :

$$\begin{aligned} \text{proj}_A(B) &= B P_A \\ &= B A^T (A A^T)^{-1} A \\ \text{where } P_A &= A^T (A A^T)^{-1} A \end{aligned}$$

For column data matrix  $A$  and single column data  $B$ , then the projection of  $B$  on space spanned by the columns of  $A$  is :

$$\begin{aligned} \text{proj}_A(B) &= P_A B \\ &= A(A^T A)^{-1} A^T B \\ \text{where } P_A &= A(A^T A)^{-1} A^T \end{aligned}$$