## Geometric interpretation of matrix multiplications

Define matrix as function $f, g : \mathbb{R}^2 \to \mathbb{R}^2$ w/ input $[x \ y]^T$

$$f : \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad g : \begin{bmatrix} p & q \\ r & s \end{bmatrix} \quad , \quad \text{then} \quad f\left(\begin{bmatrix} x \\ y \end{bmatrix}\right) = \begin{bmatrix} ax + by \\ cx + dy \end{bmatrix}$$

$$g\left(\begin{bmatrix} x \\ y \end{bmatrix}\right) = \begin{bmatrix} px + qy \\ rx + sy \end{bmatrix}$$

$$f \cdot g\left(\begin{bmatrix} x \\ y \end{bmatrix}\right) = f\left(g\begin{bmatrix} x \\ y \end{bmatrix}\right) = \begin{bmatrix} (ap+br)x + (aq+bs)y \\ (cp+dr)x + (cq+ds)y \end{bmatrix}$$

① ∴ Composite mapping $f \cdot g$ is defined as matrix product

Each element of final vector/matrix is inner product b/w row & column vector

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} a \\ c \end{bmatrix} = a\begin{bmatrix} 1 \\ 3 \end{bmatrix} + c\begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

⇒ Multiplication b/w matrix & vector is
② linear combination of two column vectors

⇓

Vector space using column vectors
$=$ column space

③ Linear transformation

$$Ax = \begin{bmatrix} 2 & -3 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

Basis vector $\left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\} \to \left\{ \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \begin{bmatrix} -3 \\ 1 \end{bmatrix} \right\}$

i.e. $\begin{bmatrix} -1 \\ 2 \end{bmatrix}$ (at original coordinate system) is $1\begin{bmatrix} 2 \\ 1 \end{bmatrix} + 1\begin{bmatrix} -3 \\ 1 \end{bmatrix}$

$$\begin{bmatrix} 2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 2 \qquad \Rightarrow \text{ inner product b/w row \& column vector}$$

↓ function (operator)
of column vector $\begin{bmatrix} 3 \\ -4 \end{bmatrix}$

→ operand for row vector $\begin{bmatrix} 2 & 1 \end{bmatrix}$

$\therefore \begin{bmatrix} 2 & 1 \end{bmatrix} \left( \begin{bmatrix} 3 \\ 4 \end{bmatrix} \right) = 2$

$$f : V \to \mathbb{R}$$

- Visualizing function of row vector

  eg. function $y = x^2$

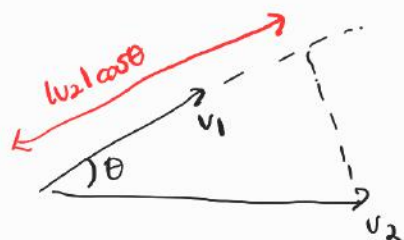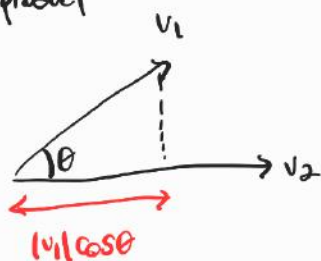function $\begin{bmatrix} 2 & 1 \end{bmatrix} \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) = 2x + y$

→ Each output $2x+y$ represents a single line ($\simeq$ contour plot)

eg. $2x + y = 1 \quad \to \quad y = -2x + 1$

these input vectors $\begin{bmatrix} x \\ y \end{bmatrix}$
outputs $-2$
$\left( \text{i.e. } \begin{bmatrix} 2 & 1 \end{bmatrix} \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) = -2 \right)$

- Geometric interpretation of dot product

$$V_1 \cdot V_2 = |v_1| |v_2| \cos\theta$$

$|v_1| \cos\theta$

$|v_2| \cos\theta$

$2x+y=4$



$$[2 \ 1]\left(\begin{bmatrix} x \\ y \end{bmatrix}\right) = 2x+y$$

$$\frac{1}{2}\cdot 2\cdot 4 = \frac{1}{2}\left(\sqrt{2^2+4^2}\right)d \qquad \therefore d = \frac{4}{\sqrt{5}}$$

$d =$ linear projection of any $\begin{bmatrix} x \\ y \end{bmatrix}$ to row vector $[2 \ 1]$

$$\therefore \ v_1 \cdot v_2 = [2 \ 1]\left(\begin{bmatrix} x \\ y \end{bmatrix}\right) = |v_1|\cdot d = 4$$

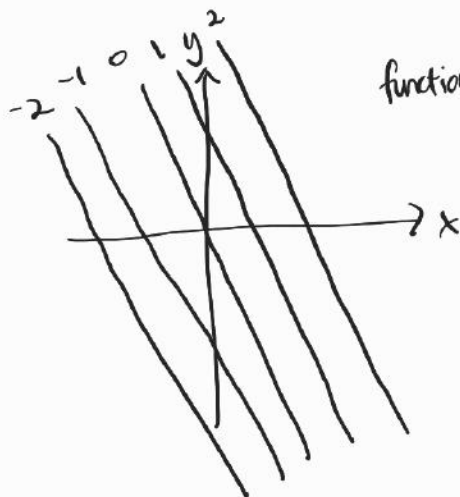- Linearity of row vector

1) $f(v+w) = f(v) + f(w)$

$$[2 \ 1]\left(\begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 3 \\ 4 \end{bmatrix}\right) = [2 \ 1]\left(\begin{bmatrix} 1 \\ 2 \end{bmatrix}\right) + [2 \ 1]\left(\begin{bmatrix} 3 \\ 4 \end{bmatrix}\right)$$

2) $f(nw) = nf(w)$

$$[2 \ 1]\left(2\begin{bmatrix} 3 \\ 4 \end{bmatrix}\right) = 2[2 \ 1]\left(\begin{bmatrix} 3 \\ 4 \end{bmatrix}\right)$$

$$\therefore f \text{ is linear operator}$$

- Geometric interpretation of row vector linearity



function $f \ [2 \ 1] : v \to \mathbb{R}$

$\Leftrightarrow$ output of input vector $\begin{bmatrix} x \\ y \end{bmatrix}$

$= \#$ of intersecting lines

$\Leftrightarrow \uparrow$ length of row vector

$=$ contour lines are closer

Addition of row vectors $=$ creating new contours

(perpendicular to the row vector)

* Row space and column space are dual space $\to$ problem hard to solve in row space can be easily substituted into column space

Dual space $v^* = \{f : v \to \mathbb{R} \mid f(a+b) \ \forall \ a,b \in v\}$

# Matrix as linear transformation

As transformation T follows linearity conditions

$$T(a+b) = T(a) + T(b)$$
$$T(ca) = cT(a)$$

$$\left.\begin{array}{c}\end{array}\right] \quad T\left(\begin{bmatrix} x \\ y \end{bmatrix}\right) = x\,T\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) + y\,T\left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}\right)$$

$\Leftrightarrow$ basis vector $\begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ transforms into $T\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right), T\left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}\right)$

\* Geometric features of linearity :  1) grid lines are linear

2) uniform spacing b/w grid lines

• Types of linear transformation

1) shearing $\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$

2) rotation $\begin{bmatrix} \cos\frac{\pi}{2} & -\sin\frac{\pi}{2} \\ \sin\frac{\pi}{2} & \cos\frac{\pi}{2} \end{bmatrix}$
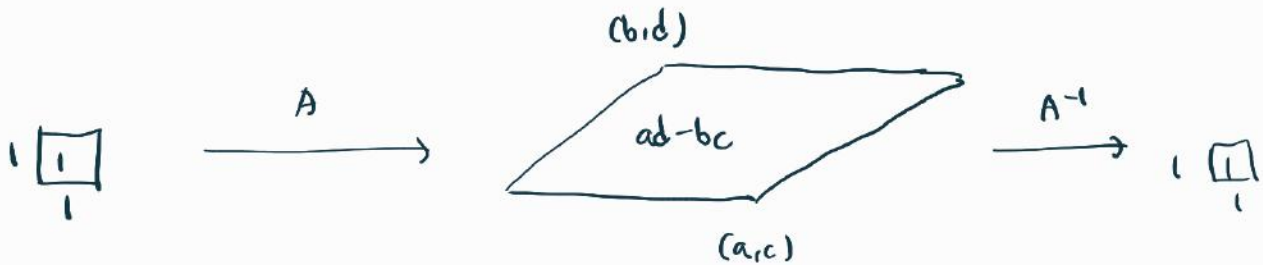
3) permutation $\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$

4) projection $\begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$

# Geometric interpretation of determinent

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \quad , \quad A^{-1} = \frac{1}{\det(A)} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \quad \text{where} \quad \det(A) = ad - bc$$

det of 2x2 matrix = area of parallelogram of <u>two basis vectors after transformation</u>

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} a \\ c \end{pmatrix}, \begin{pmatrix} b \\ d \end{pmatrix}$$

(b,d)

$$\begin{array}{ccc}
\boxed{1} & \xrightarrow{\quad A \quad} & \begin{array}{c} ad-bc \end{array} & \xrightarrow{\quad A^{-1} \quad} & \boxed{1}
\end{array}$$

(a,c)

Normally, matrix A is linear transformation function

$$x \longrightarrow A \longrightarrow Ax$$

$\underline{\text{another vector}}$ w/ different
length and direction

However, some matrix A doesn't change direction of Ax from x

i.e, $Ax = \lambda x$

$\nearrow_x \rightarrow A \rightarrow \nearrow_{Ax = kx}$

$\lambda$ = eigenvalue
$x$ = eigenvector

$(A - \lambda I)x = 0$

either $A - \lambda I$ or
$x$ equals 0

$\nearrow$ trivial solution $x = 0$

$\rightarrow$ non-trivial solution $\Leftrightarrow$ $A - \lambda I$ is $\Leftrightarrow$ det$(A - \lambda I)$
$A - \lambda I = 0$ non-invertible $= 0$

Linear transformation $A$ = function ?

Definition of function



$X$ : domain

$Y$ : codomain    $f(x)$ : range

$f$ : subset of Cartesian product $X \times Y$

$\forall x \in X$, if there exists unique $y \in Y$, $y = f(x)$

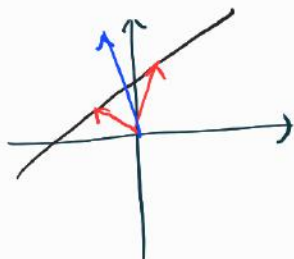$\Leftrightarrow$ $f$ is mapping b/w $X$ and $Y$

$f : X \to Y$

- subspace

* vector space : set of vectors defined by + (addition) and · (scalar multiplication)

Subspace of vector space $\simeq$ subset of set

$\llcorner$ still requires to follow definition of vector space

e.g.



not on the line = addition rule $\times$ holds

- Row / column space

$A = \begin{bmatrix} 2 & 4 \\ 1 & 2 \end{bmatrix}$  $\to$ row space = span $\left\{ \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right\}$   $\rbrack$ they are subspaces!

column space = span $\left\{ \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 4 \\ 2 \end{bmatrix} \right\}$

- null space

set of $\vec{x}$ that satisfies $A\vec{x} = 0$

* Row space $\perp$ null space

If matrix is a function, how do we define the relationship b/w sets, which is the fundamental meaning of the function ?

$$A \in \mathbb{R}^{m \times n} \quad , \quad f: \underline{\mathbb{R}^n} \to \underline{\mathbb{R}^m}$$

domain      codomain

n-dim vector space     m-dim vector space

① Domain $\mathbb{R}^n$ = row space + null space

As row space $\perp$ null space , represent $x \in \mathbb{R}^n$ as linear combination of two vectors inside spaces.

② Range = column space

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = x \begin{bmatrix} a \\ c \end{bmatrix} + y \begin{bmatrix} b \\ d \end{bmatrix} = \text{linear combination of column space}$$
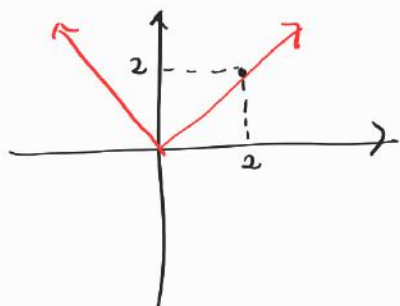
③ Co-domain = column space + left null space

Column space $\perp$ left null space

==change of basis==

Standard basis $(\hat{i}, \hat{j}) = \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) = \mathcal{E}$

e.g. new basis $B = \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right)$



$\begin{bmatrix} 2 \\ 2 \end{bmatrix}_{\mathcal{E}} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}_{B}$

$\begin{bmatrix} 2 \\ 2 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0 \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

$\therefore \begin{bmatrix} 2 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$

$\therefore$ For new basis $C = \left\{ \begin{bmatrix} c_1 \\ 1 \end{bmatrix}, \begin{bmatrix} c_2 \\ 1 \end{bmatrix} \right\}$ and vector $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}_C$

$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{\mathcal{E}} = \begin{bmatrix} c_1 & c_2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}_C$

- Coordinate transition

  Basis $B = \{v_1, v_2\}$, $C = \{w_1, w_2\}$

  $\begin{bmatrix} w_1 = a v_1 + b v_2 \\ w_2 = c v_1 + d v_2 \end{bmatrix}$

  For vector $v$,

  $\begin{aligned} v &= l_1 w_1 + l_2 w_2 \\ &= l_1 (a v_1 + b v_2) + l_2 (c v_1 + d v_2) \\ &= (a l_1 + c l_2) v_1 + (b l_1 + d l_2) v_2 \\ &= k_1 v_1 + k_2 v_2 \end{aligned}$

  $\therefore \begin{bmatrix} k_1 \\ k_2 \end{bmatrix} = \begin{bmatrix} a l_1 + c l_2 \\ b l_1 + d l_2 \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} l_1 \\ l_2 \end{bmatrix}$

  $\downarrow$

  transition matrix

# Elementary square matrices

· Solving simultaneous equations

$$\begin{bmatrix} 2x + 3y = 1 \\ 4x + 7y = 3 \end{bmatrix}$$

operations ① $r_1 \to kr_1$    Row multiplication

② $r_1 \to r_1 + r_2$    Row addition

③ $r_1 \leftrightarrow r_2$    Row switching

$\downarrow$

we can use matrix for these computations $\Rightarrow$ easier for computers to computer

$$\begin{bmatrix} 2 & 3 \\ 4 & 7 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$\left[\begin{array}{cc|c} 2 & 3 & 1 \\ 4 & 7 & 3 \end{array}\right] \leftarrow \text{this augmented matrix is operand}$$

① row multiplication $E = \begin{bmatrix} k & 0 \\ 0 & 1 \end{bmatrix}$    $r_1 \to kr_1$    $E^{-1} = \begin{bmatrix} 1/k & 0 \\ 0 & 1 \end{bmatrix}$

② row switching $E = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$    $r_1 \leftrightarrow r_2$    $E^{-1} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$

③ row addition $E = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$    $r_1 \to r_1 + r_2$    $E^{-1} = \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix}$

$\downarrow$

these operations are called elementary matrices

* only applicable to square matrix

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -2 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 2 & 3 & -1 \\ 2 & 3 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -3 \\ 0 & 0 & 4 \end{bmatrix}$$

$\quad\quad E_1 \quad\quad\quad\quad E_2 \quad\quad\quad\quad E_3 \quad\quad\quad\quad A \quad\quad\quad\quad\quad U$

$r_3 \to r_3 - r_2 \quad\quad r_3 \to r_3 - 2r_1 \quad\quad r_2 \to r_2 - 2r_1 \quad\quad\quad\quad$ upper-triangular matrix

↓

$$A = E_1^{-1} E_2^{-1} E_3^{-1} U \quad\quad\quad A = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -3 \\ 0 & 0 & -4 \end{bmatrix} = LU$$

$$\quad\quad = LU$$

* when row switching operation is required, we use PLU decomposition

$$eg \quad A = \begin{bmatrix} 0 & 0 & 3 \\ 1 & 1 & 1 \\ 2 & 3 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} LU$$

$$P : r_1 \leftrightarrow r_3$$

• Applications

(1) solving $Ax = b$

$$Ax = b \quad \xrightarrow{A = LU} \quad L(Ux) = b \quad \xrightarrow{Ux = c} \quad Lc = b \quad \xrightarrow[\text{forward substitution}]{\text{solve } c}$$

$$Ux = c \quad \xrightarrow[\text{backward substitution}]{\text{solve } x} \quad x$$

② calculate $\det(A)$

$$\det(A) = \det(LU) = \det(L)\det(U) = \prod_{i=1}^{n} l_{i,i} \prod_{j=1}^{n} u_{j,j}$$
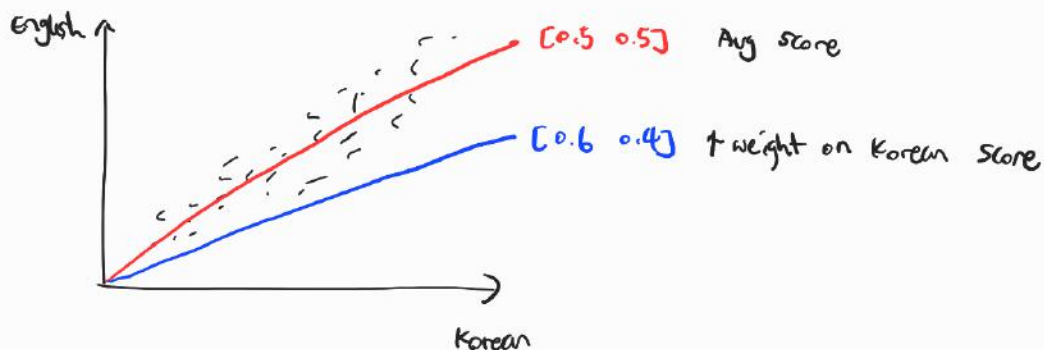
(multiplication of all diagonal elements)

# Principal Component Analysis (PCA)

⇒ If you need to reduce dimension of data by projection to a vector, what is the best way to minimize data loss?

e.g. calculating overall test score

| English | Korean |
|---------|--------|
| 80 | 60 |
| 70 | 65 |
| 75 | 50 |
| ⋮ | ⋮ |

English ↑

[0.5 0.5]  Avg score

[0.6 0.4]  ↑ weight on Korean score

→ Korean

⇒ Main problem: what is the best vector that gives **best results**?

Dot product b/w data and vector

e.g. $[80 \ 60] \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

• Covariance matrix

Matrix = linear transformation & function mapping one vector space to another

$A_{cov} = \begin{bmatrix} ③ & ② \\ ② & ④ \end{bmatrix}$

→ spread along x-axis

→ spread along y-axis

↓↓ spread along both x,y-axis

Eigenvector : principle axis that matrix acts on
Eigenvalue : stretched magnitude along eigenvector  ] ⇒ sort eigenvalue ∴ ↑eigenvalue

= ↑ significance of eigenvector

∴ PCA = projection of data of principle axis = find eigenvector of covariance matrix w/

↑ eigenvalue

- Calculating covariance matrix

$$X = \begin{pmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_d \\ | & | & & | \end{pmatrix} \in \mathbb{R}^{n \times d} \qquad \text{i.e. } n \text{ samples, } d \text{ features}$$

$$X^T X = \begin{pmatrix} x_1 \cdot x_1 & \cdots & x_1 \cdot x_d \\ \vdots & \ddots & \\ x_d \cdot x_1 & & x_d \cdot x_d \end{pmatrix} \qquad \text{i.e. } (X^T X)_{ij} = \text{similarity b/w feature } x_i \text{ and } x_j$$

$$\Sigma = \frac{X^T X}{n} \qquad \because \text{dot product} \uparrow \text{value} \quad \text{as} \quad \text{sample size } (n) \uparrow$$

- Find # of dimensions to reduce

For full rank covariance matrix $\Sigma_{d \times d}$, let eigenvalue $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$,

select $m$ which satisfies $\quad \dfrac{\sum\limits_{j=1}^{m} \lambda_j}{\sum\limits_{i=1}^{d} \lambda_i} \geq 0.9$

$\Rightarrow$ Decomposing original linear transformation $A$ into $V$ (rotation), $\Lambda$ (stretch), $V^{-1}$ (rotation)

° Derivation

Assume there are $n$ independent eigenvectors ($u_1 \cdots u_n$) and eigenvalues ($\lambda_1 \cdots \lambda_n$)

of matrix $A \in \mathbb{R}^{n \times n}$,

$A v_i = \lambda u_i$ for $i = 1, 2, \cdots, n$

$$AV = \begin{bmatrix} | & & | \\ \lambda_1 u_1 & \cdots & \lambda_n u_n \\ | & & | \end{bmatrix} = V \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} = V\Lambda$$

$\therefore A = V\Lambda V^{-1}$

· Geometric interpretation

e.g. $A = V\Lambda V^{-1}$

$\therefore \Lambda$ acts as stretching

$$\begin{bmatrix} 1.2 & -0.5 \\ -1.5 & 1.9 \end{bmatrix} = \begin{bmatrix} \boxed{0.6089} & -0.3983 \\ \boxed{0.9933} & 0.9172 \end{bmatrix} \begin{bmatrix} \boxed{0.5486} & 0 \\ 0 & 2.3514 \end{bmatrix} \begin{bmatrix} 1.0489 & 0.4555 \\ -0.9092 & 0.6963 \end{bmatrix}$$

$\lambda_1$ (above boxed 0.5486)

$v_1$ (below boxed first column)

$\therefore V^{-1}$ acts as inverse rotation

normalized vector $|v_1| = 1$

$\therefore V$ acts like rotation (as basis vector length $= 1$)

· EVD of symmetric matrix

if $A = A^T$, $A = Q\Lambda Q^T$

- Eigenvector of rotation matrix

$$A(\theta) = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \rightarrow Ax = \lambda x : \text{ is there vector } x \text{ that retains its direction after rotation } A?$$
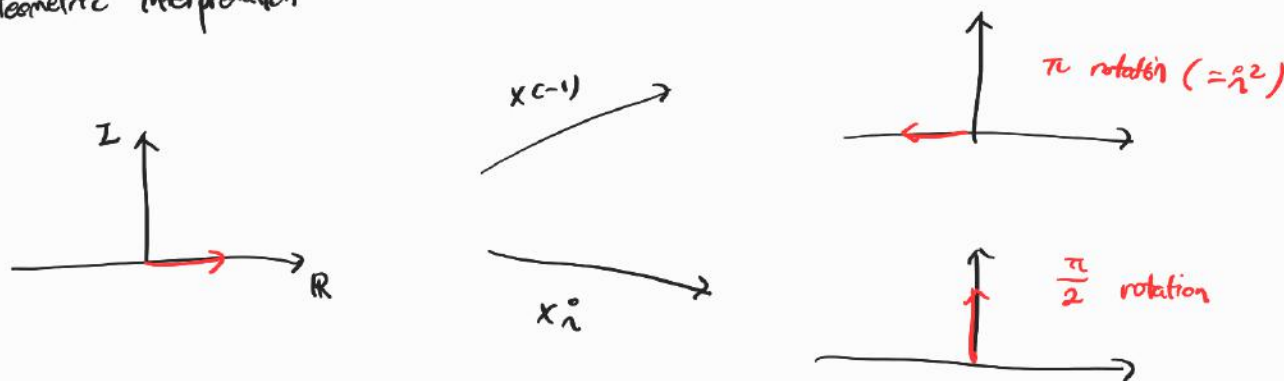
$$(A - \lambda I)x = 0$$

$$\det(A - \lambda I) = (\cos\theta - \lambda)^2 + \sin^2\theta = 0$$

$$\lambda^2 - 2\lambda\cos\theta + 1 = 0$$

$$\therefore \lambda = \cos\theta \pm i\sin\theta$$

if $\lambda = \cos\theta + i\sin\theta$, $x = \begin{bmatrix} i \\ 1 \end{bmatrix}$, if $\lambda = \cos\theta - i\sin\theta$, $x = \begin{bmatrix} -i \\ 1 \end{bmatrix}$
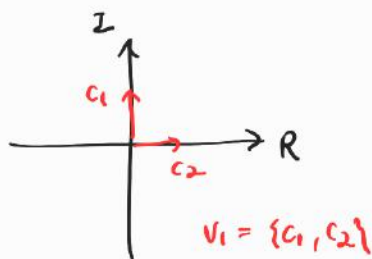
- Geometric interpretation



$$\therefore \text{ Imaginary \# multiplications} = \text{rotation of vector}$$

- Visualizing complex eigenvector

$$x_1 = \begin{bmatrix} i \\ 1 \end{bmatrix} \quad x_2 = \begin{bmatrix} -i \\ 1 \end{bmatrix}$$



$\therefore$ two vectors $\equiv$ a complex vector

$$V_1 = \{c_1, c_2\}$$

- Relationships b/w rotation matrix & eigenvector

$$\lambda_1 = \cos\theta + i\sin\theta = \exp(i\theta) = \text{anti-clockwise } \theta \text{ rad rotation}$$

$\therefore$ scaling complex eigenvector w/ eigenvalue $(\lambda x) \equiv$ rotation of eigenvector $c_1, c_2$ $(Ax)$
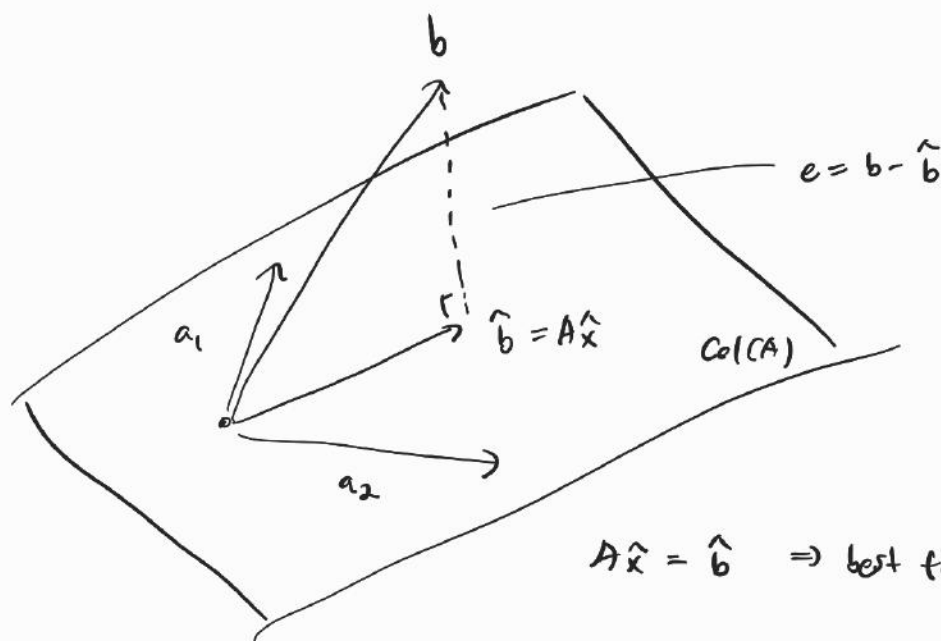
$\wedge$ Complex

# Linear regression

$\Rightarrow$ find the best linear trend line that explains the data (# of data > feature dimension)

$Ax = b$

$$\begin{bmatrix} \vert & \vert \\ a_1 & a_2 \\ \vert & \vert \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \vert \\ b \\ \vert \end{bmatrix} \qquad \Rightarrow \qquad x_1 \begin{bmatrix} \vert \\ a_1 \\ \vert \end{bmatrix} + x_2 \begin{bmatrix} \vert \\ a_2 \\ \vert \end{bmatrix} = \begin{bmatrix} \vert \\ b \\ \vert \end{bmatrix}$$

How to combine $a_1$ and $a_2$ to output $b$

if $b \notin$ span $\{a_1, a_2\}$ ($b \notin Col(A)$), there is no exact solution, thus find "best fit"



$$A\hat{x} = \hat{b} \quad \Rightarrow \text{ best fit (min error)}$$

$A \cdot e = \begin{bmatrix} \vert & \vert \\ a_1 & a_2 \\ \vert & \vert \end{bmatrix} \cdot e = 0 \qquad \because e$ is perpendicular to any vector in $Col(A)$

$\qquad \equiv e$ is in left nullspace

$\qquad A^T e = A^T(b - A\hat{x}) = A^T b - A^T A \hat{x} = 0$

$\qquad\qquad A^T A \hat{x} = A^T b$

$\qquad\qquad \therefore \hat{x} = (A^T A)^{-1} A^T b$

# Geometric meaning of pseudo-inverse

- Definition

  For $A \in \mathbb{R}^{m \times n}$, if $m > n$ and column vectors are linearly independent

  $$A^+ = (A^TA)^{-1} A^T \quad \text{where} \quad A^TA \text{ is invertible} \qquad \rightarrow A^+A = I \text{ (left inverse)}$$

  if $m < n$ and ——

  $$A^+ = A^T (AA^T)^{-1} \quad \text{——} \qquad \rightarrow AA^+ = I \text{ (right inverse)}$$

  $\Rightarrow$ Matrix of any size can function as inverse matrix

- Mathematical meaning

  $$Ax = b$$
  $$A^+Ax = A^+b \qquad (A^+ = (A^TA)^{-1}A^T)$$

  $$\begin{bmatrix} -1 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 3 \end{bmatrix}$$

  $\downarrow$ using $A^+$

  $$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

  $\downarrow$ BUT

  $$\begin{bmatrix} -1 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \\ 2 \end{bmatrix} \neq \begin{bmatrix} 0 \\ 1 \\ 3 \end{bmatrix}$$

- Linear regression & pseudo-inverse

  $$Ax = b \quad \Rightarrow \quad A^+Ax = A^+b \quad \Rightarrow \quad \tilde{x} = (A^TA)^{-1}A^Tb$$

  $\hat{x}$ is not an exact solution, but "best-fit" solution

  $$= \text{projection to } Col(A)$$

- SVD & pseudo-inverse

  $$A = U\Sigma V^T \qquad (U \text{ and } V \text{ are orthogonal matrix i.e. } UU^T = UU^T = I)$$
  $$(\Sigma \text{ is diagonal matrix i.e. } \Sigma^T = \Sigma)$$

  $$A^T = V\Sigma^T U^T = V\Sigma U^T$$

  $\therefore A^TA = V\Sigma U^T U\Sigma V^T = V\Sigma^2 V^T$

  $$(A^TA)^{-1} = V(\Sigma^2)^{-1}U^T$$

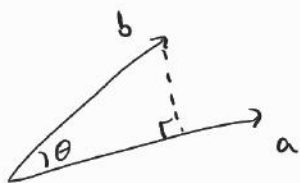  $$A^+ = (A^TA)^{-1}A^T = V(\Sigma^2)^{-1}U^TU\Sigma U^T = V\Sigma^{-1}U^T$$

  $$\Sigma^{-1} = \begin{pmatrix} \lambda_1^+ & & \\ & \ddots & \\ & & \lambda_{min(n,m)}^+ \end{pmatrix}$$

  where $\lambda^+ \begin{cases} \lambda^{-1} & \text{if } \lambda \neq 0 \\ 0 & \text{if } \lambda = 0 \end{cases}$
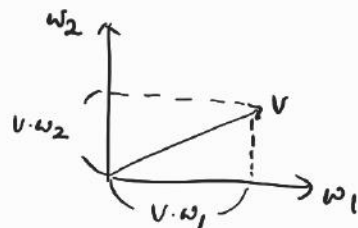
# QR decomposition

- Vector projection



$$\text{comp}_a b = |b|\cos\theta$$
$$a \cdot b = |a||b|\cos\theta$$
$$\left.\right] \quad \text{comp}_a b = \frac{a \cdot b}{|a|} \quad \text{(scalar)}$$

$$\text{proj}_a b = \text{comp}_a b \cdot \frac{a}{|a|} = \frac{a \cdot b}{|a|} \cdot \frac{a}{|a|} = \frac{a \cdot b}{a \cdot a} a \quad \text{(vector : multiplied w/ a unit vector)}$$

- Gram-Schmidt process

$\Rightarrow$ Convert linearly independent vectors to orthogonal basis

if $\{w_1, w_2\}$ are orthogonal basis, we can represent any $v = (v \cdot w_1)w_1 + (v \cdot w_2)w_2$



Given independent vectors $\{a_1, \cdots, a_k\}$

$$u_1 = a_1$$
$$u_2 = a_2 - \text{proj}_{u_1}(a_2)$$
$$\vdots \qquad \underline{\phantom{xxxxxxxxx}}$$
$$\textcolor{red}{u_1 \text{ component of } a_2}$$
$$u_k = a_k - \sum_{j=1}^{k-1} \text{proj}_{u_j}(a_k)$$

Then $q_i = \dfrac{u_i}{|u_i|}$ $\rightarrow$ unit-vector orthogonal basis $\{q_1, \cdots, q_k\}$

- QR decomposition

$$A = QR = \begin{bmatrix} | & & | \\ q_1 & \cdots & q_n \\ | & & | \end{bmatrix} \begin{bmatrix} a_1 \cdot q_1 & a_2 \cdot q_1 & \cdots & a_n \cdot q_1 \\ \vdots & & \ddots & \\ a_1 \cdot q_n & & \cdots & a_n \cdot q_n \end{bmatrix}$$

for $a_i \cdot q_j$ $(i < j)$, $a_i \cdot q_j = 0$ $\because q_j$ already removed all $i < j$ components during Gram-Schmidt Process

$$\therefore A = \begin{bmatrix} | & | & & | \\ q_1 & q_2 & \cdots & q_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} a_1 \cdot q_1 & a_2 \cdot q_1 & \cdots & a_n \cdot q_1 \\ & a_2 \cdot q_2 & & \vdots \\ & 0 & \ddots & \\ & & & a_n \cdot q_n \end{bmatrix}$$

## SVD decomposition

⇒ For set of orthogonal vectors, what is orthogonal set after linear transformation?

$$A = U\Sigma V^T \quad (A \in \mathbb{R}^{m\times n}, \underbrace{U \in \mathbb{R}^{m\times m}, V \in \mathbb{R}^{n\times n}}_{\text{orthogonal}}, \underbrace{\Sigma \in \mathbb{R}^{m\times n}}_{\text{diagonal}})$$

$$\therefore U^{-1} = U^T, \quad V^{-1} = V^T$$

e.g. 2D case $(A \in \mathbb{R}^{2\times 2})$

$$V = \begin{pmatrix} | & | \\ x & y \\ | & | \end{pmatrix}$$    set of orthogonal vectors before transformation

$$U = \begin{pmatrix} | & | \\ u_1 & u_2 \\ | & | \end{pmatrix}$$    normalized set of orthogonal vectors after transformation

$$\Sigma = \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}$$    scaling factor   $(\sigma_1 |x| = |Ax|)$

$\therefore AV = U\Sigma$    After linear transformation $A$ of $V$'s orthogonal column vectors,

is there set of column vectors $(U)$ w/ scaling factor $\sigma$?

$$A = U\Sigma V^{-1} = U\Sigma V^T$$

- A doesn't requires to be a square matrix

e.g. $A \in \mathbb{R}^{2\times 3}$   $(\therefore \mathbb{R}^3 \to \mathbb{R}^2)$

Collapsing 3D to 2D space by one scaling factor $= 0$

$$A = U\Sigma V^T = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \cdots + \sigma_n u_n v_n^T$$

$\therefore$ SVD decomposes $A$ matrix into different matrix (its weight(size determined

by $\sigma_i$)

⇒ Able to select few decomposed matrix w/ ↑$\sigma$

This ↓ size but minimize information loss

# Non-negative Matrix Factorization

=> Decompose a non-negative matrix $X$ into two non-negative matrix $H, W$

$$
\underset{\substack{\text{Sample}\\\text{size}}}{\overset{\text{Data dimension}}{\left[\quad X \quad\right]}} = \underset{\substack{\text{Sample}\\\text{size}}}{\overset{\text{feature dim}}{\left[\quad W \quad\right]}} \overset{\text{Data dim}}{\left[\quad H \quad\right]}
$$

Advantage: Can preserve <u>non-negative value feature</u> (e.g. pixel's intensity)

↳ not assured for other matrix factorization methods e.g. SVD

Can preserve data distribution better ∵ feature $X$ needs to be orthogonal

- How to find $W, H$

=> Iterative update

$$
\left[
\begin{array}{l}
H := H \circ \dfrac{W^T X}{W^T W H} \\[2em]
W := W \circ \dfrac{X H^T}{W H H^T}
\end{array}
\right.
$$

($\circ / -$ is element-wise multiplication / division)

$f : \mathbb{R}^n \to \mathbb{R}^m$

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \\ \frac{\partial f_m}{\partial x_1} & & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$\Rightarrow$ approximates nonlinear transformation in a local region into linear transformation



non-linear

approximated linear

for local area,

$$\begin{bmatrix} dx \\ dy \end{bmatrix} = J \begin{bmatrix} du \\ dv \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} du \\ dv \end{bmatrix}$$

$$dx \cdot dy = |J| \, du \cdot dv$$

$$\therefore J = \begin{bmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{bmatrix}$$

$f: \mathbb{R}^n \to \mathbb{R}^m$

$$H(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & & \ddots & \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & & & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$$

so $H$ is symmetric matrix

$f(x) = ax^2 + bx + c$   :   if $f''(x) > 0$   $\cup$ shape
$\qquad\qquad\qquad\qquad\qquad$ $f''(x) < 0$   $\cap$ shape

$\uparrow |f''(x)| \to \uparrow$ concave/convex

Similarly, $H$ transforms bowl-shaped function more concave/complex

$H \to$ find eigenvalue/eigenvector $\to$ $\uparrow$ eigenvalue $\equiv \uparrow$ concave/convex

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ all eigenvalues $> 0$ $= \cup$ shape

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ mixed eigenvalues' sign $\equiv$ saddle shape function

- Fourier transform

  Decomposing signal mixed w/ different frequency & amplitude

  signal $s(t)$ — time (t)     $\xrightarrow{F}$ amplitude — 60Hz, 120Hz — frequency (Hz)

  ↓

  we can interpret signal as vector (order of numbers)

  $$x[n] = [x[0], x[1], \cdots, x[N-1]]$$   signal discretized every 1 sec

  frequency components can also be vector

  $$X[k] = [x[0], \cdots, x[N-1]]$$   discretized every 1 Hz

  Direct Fourier Transform

  $$X[k] = \sum_{n=0}^{N-1} x[n] \exp\left(-j\frac{2\pi k}{N} n\right)$$

  $$\begin{bmatrix} x[0] \\ X[1] \\ \vdots \\ x[N-1] \end{bmatrix} = \begin{bmatrix} 1 & 1 & \cdots & \cdots & 1 \\ 1 & w^1 & & & \vdots \\ \vdots & & & & \vdots \\ 1 & w^{N-1} & w^{(N-1)\cdot 2} & \cdots & w^{(N-1)(N-1)} \end{bmatrix} \begin{bmatrix} x[0] \\ \vdots \\ x[N-1] \end{bmatrix}$$   where $w = \exp\left(-j\frac{2\pi}{N}\right)$

  Amount of frequency $X[1]$ = similarity (dot product) b/w $[1 \; w^1 \cdots w^{N-1}]$ and signal

  $w = \exp\left(-j\frac{2\pi}{N}\right)$   means

  Im, Re axes, $-\frac{2\pi}{N}$ rad

  ∴ Fourier Transform matrix $\begin{bmatrix} - & - & \cdots & - \\ - & \searrow & \cdots & \nearrow \\ - & | & \cdots & | \\ & & \vdots & \end{bmatrix}$ = cosine fundamental frequency + i sine fundamental frequency

## Circulant matrix and convolution

$\Rightarrow$ Matrix that operates cyclic permutation

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix} \rightarrow Px = \begin{bmatrix} x_{n-1} \\ x_0 \\ \vdots \\ x_{n-2} \end{bmatrix} \qquad \therefore P = \begin{bmatrix} 0 & 0 & \cdots & 1 \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & & \vdots \\ 0 & & & \\ 0 & \cdots & 1 & 0 \end{bmatrix}$$

- Decomposition of a signal vector

$$\delta \text{ (discrete unit sample function)} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix} = x_0 \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \cdots + x_{n-1} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} = x_0 \delta + x_1 P\delta + \cdots + x_{n-1} P^{n-1}\delta$$
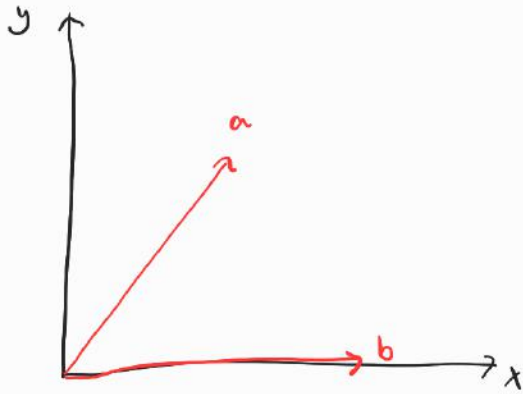
$$= (x_0 I + x_1 P + \cdots + x_{n-1} P^{n-1}) \delta$$

$$= \underbrace{\begin{bmatrix} x_0 & x_{n-1} & \cdots & x_1 \\ x_1 & \ddots & & \vdots \\ \vdots & & \ddots & \\ x_{n-1} & \cdots & & x_0 \end{bmatrix}}_{\text{circulant matrix}} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

For data $X$ and $Y$, correlation $r = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{S_x} \right) \left( \frac{y_i - \bar{Y}}{S_Y} \right)$



$$a \cdot b = |a||b| \cos\theta$$
$$= (|a|\cos\theta)(|b|)$$
$$= (\text{Projection of } a \text{ onto } b)(\text{length of } b)$$
$$= \text{how much of the change of } a \text{ can be explained by } b$$

similarly, $\text{Proj}_b a = \frac{a \cdot b}{|a|}$

$$= \text{how much} - b - \text{explained by } a$$

$\therefore$ $a$ and $b$ explains each other $= \frac{a \cdot b}{|a||b|} = \cos\theta$

if $a = x_i - \bar{x}$, $b = y_i - \bar{Y}$,

$$r = \frac{a \cdot b}{|a||b|}$$

$\therefore$ how much $x_i - \bar{x}$ and $y_i - \bar{Y}$ explains each other?
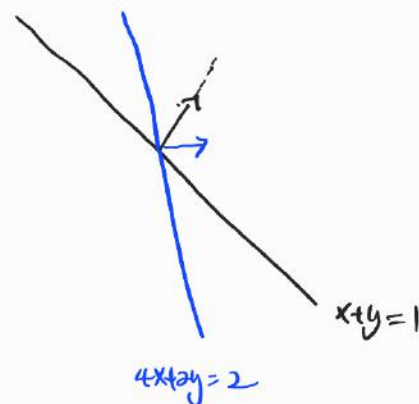
# Geometric interpretation of Gauss-Jordan Elimination

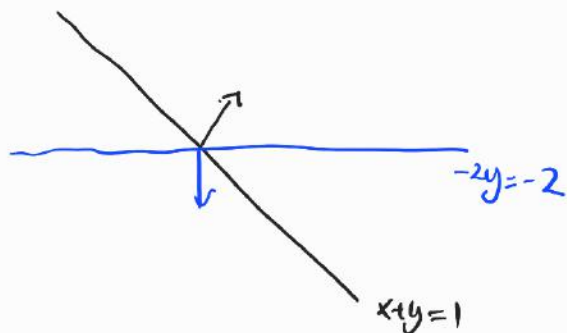→ transforming normal vectors of line equation into unit vectors parallel to each other
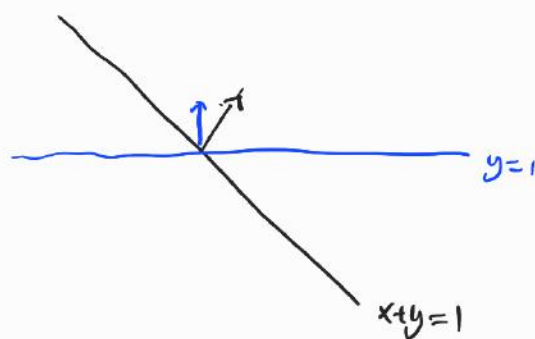
$$\begin{cases} 2x+2y=2 \\ 4x+2y=2 \end{cases}$$



2x+2y=2
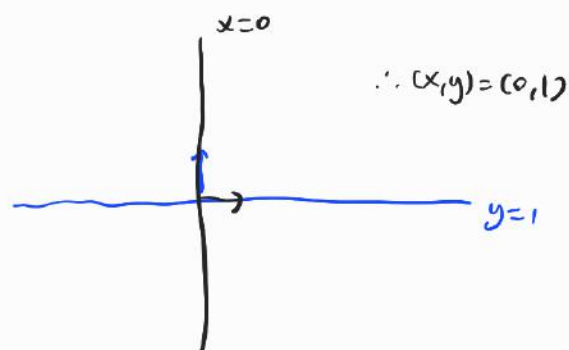4x+2y=2

(i) Scaling    $2x+2y=2 \rightarrow x+y=1$



x+y=1
4x+2y=2

(ii) Subtraction    $4x+2y=2 \rightarrow -2y=-2$



$-2y=-2$
$x+y=1$

(iii) Scaling    $-2y=2 \rightarrow y=1$



$y=1$
$x+y=1$

(iv) Subtraction    $x+y=1 \rightarrow x=0$



$x=0$

$\therefore (x,y)=(0,1)$

$y=1$

Suppose functions $f_1(x), f_2(x), \cdots, f_n(x)$ possesses at least $n-1$ derivatives, then if determinent

$$W(f_1, f_2, \cdots f_n) = \begin{vmatrix} f_1 & f_2 & \cdots & f_n \\ f_1' & f_2' & \cdots & f_n' \\ \vdots & & & \vdots \\ f_1^{(n-1)} & \cdots & & f_n^{(n-1)} \end{vmatrix} \neq 0$$

$f_1, f_2, \cdots f_n$ are linearly independent

\* $W \neq 0 \not\Rightarrow f_1, f_2, \cdots, f_n$ are linearly dependent

- Proof by contradiction

Suppose $W \neq 0$ and functions are linearly dependent and

$$c_1 f_1 + c_2 f_2 + \cdots + c_n f_n = 0$$

$$\vdots$$

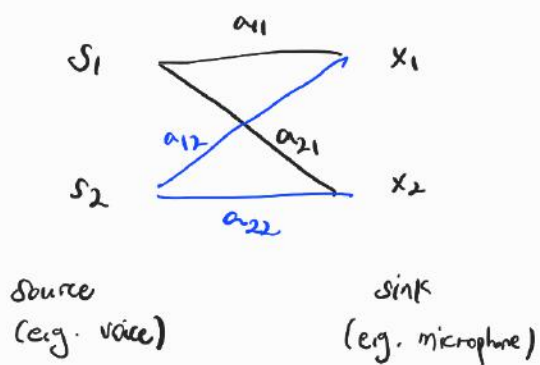$$c_1 f_1^{(n-1)} + c_2 f_2^{(n-1)} + \cdots + c_n f_n^{(n-1)} = 0$$

$$\underbrace{\begin{pmatrix} f_1 & f_2 & \cdots & f_n \\ f_1' & f_2' & \cdots & f_n' \\ & & & \\ f_1^{(n-1)} & f_2^{(n-1)} & \cdots & f_n^{(n-1)} \end{pmatrix}}_{A} \underbrace{\begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{pmatrix}}_{x} = \underbrace{\begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}}_{b} \qquad Ax = b$$

if $\det(A) = W$, $W \neq 0$, $A^{-1}$ exists $\therefore x = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix} = 0$

$c_1 f_1 + \cdots + c_n f_n = 0 \quad \rightarrow \quad (c_1, c_2, \cdots, c_n) = (0, 0, \cdots, 0)$ is the only solution

$\therefore f_1, f_2, \cdots, f_n$ are linearly independent ✗

# Independent Component Analysis (ICA)

$$x_1(t) = a_{11} s_1(t) + a_{12} s_2(t)$$
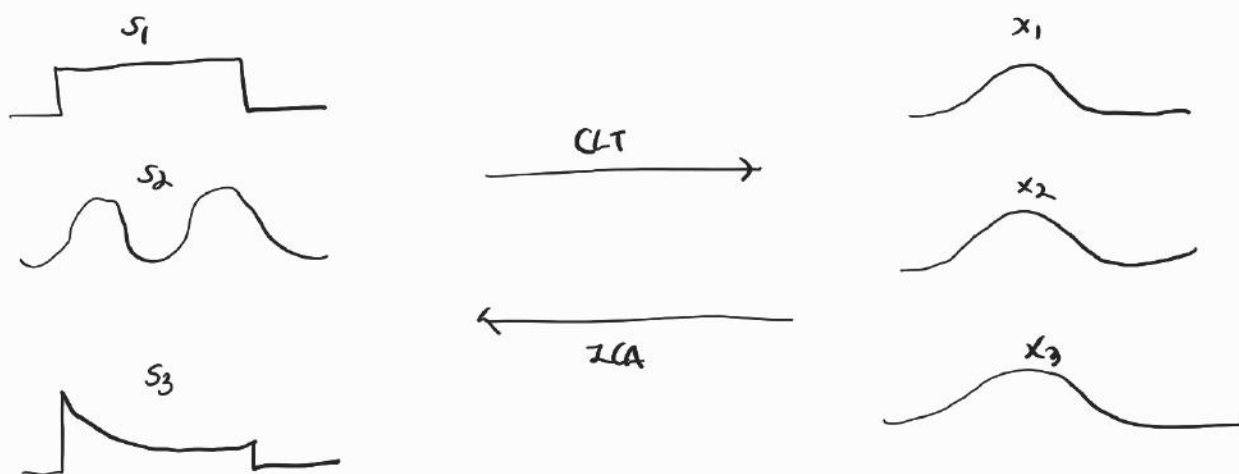$$x_2(t) = a_{21} s_1(t) + a_{22} s_2(t)$$

$$X = As \quad (A: \text{mixing matrix})$$
$$\downarrow$$
$$s = A^{-1}x = Wx \quad (A^{-1}=W: \text{unmixing matrix})$$

ICA : find $W$ (from sink to source) w/o knowing $A$

$S_1$    $a_{11}$    $X_1$

$a_{12}$   $a_{21}$

$S_2$    $a_{22}$    $X_2$

Source      sink
(e.g. voice)    (e.g. microphone)

- Central Limit Theorem (CLT) $\longleftrightarrow$ ICA

$S_1$

$S_2$

$S_3$

$\xrightarrow{\text{CLT}}$

$\xleftarrow{\text{ICA}}$

$X_1$

$X_2$

$X_3$

e.g.   $S \sim \text{Uniform } [0,1] \quad \rightarrow A=2 \rightarrow X \sim \text{uniform} [0,2]$

$$X = AS$$
$$S = A^{-1}x = WX$$
$$P_X(x) = |W| \, P_S(Wx)$$

$$* \quad \det\left(\begin{bmatrix} a & \kappa b \\ c & \kappa d \end{bmatrix}\right) = \kappa \det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right)$$

$$\det([A_1, A_2, \cdots, \kappa_1 B_1, \kappa_2 B_2, \cdots, A_n]) \qquad (A_\cdot, B_\cdot \text{ are col vectors})$$

$$= \kappa_1 \det([A_1, \cdots, B_1, \cdots, A_n]) + \kappa_2 \det([A_1, \cdots, B_2, \cdots, A_n])$$

- Cramer's Rule

$$Ax = b \quad \Rightarrow \quad x_i = \frac{\det(A_i^{rep})}{\det(A)} \quad \text{where } A_i^{rep} = \begin{bmatrix} a_{11} & \cdots & b_1 & \cdots & a_{1n} \\ \vdots & & \vdots & & \vdots \\ a_{1n} & \cdots & b_n & \cdots & a_{nn} \end{bmatrix}$$

- Proof

$$Ax = x_1\begin{bmatrix} | \\ A_1 \\ | \end{bmatrix} + x_2\begin{bmatrix} | \\ A_2 \\ | \end{bmatrix} + \cdots + x_n\begin{bmatrix} | \\ A_n \\ | \end{bmatrix} = b$$

$$A_i^{rep} = [A_1, A_2, \cdots, b, \cdots, A_n]$$

$$\det(A_i^{rep}) = \det([A_1, A_2, \cdots, b, \cdots A_n]) = \det([A_1, \cdots, x_1A_1 + x_2A_2 + \cdots + x_nA_n, \cdots A_n])$$
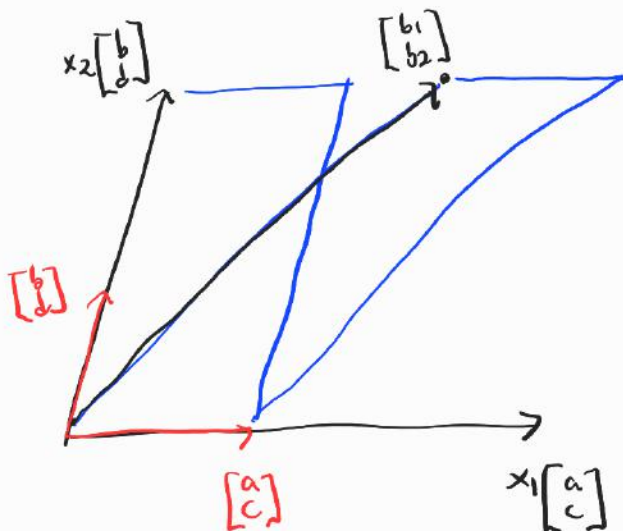
$$= \sum_{j=1}^{n} x_j \det([A_1, A_2, \cdots A_j, \cdots, A_n])$$

$$= x_i \det([A_1, A_2, \cdots, A_i, \cdots A_n]) \qquad \because \text{ linearly dependent col vector } \det = 0$$

$$= x_i \det A$$

- Geometric interpretation

$$Ax = b \qquad x_1\begin{bmatrix} a \\ c \end{bmatrix} + x_2\begin{bmatrix} b \\ d \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$



$$\det\left(\begin{bmatrix} a & x_2 b \\ c & x_2 d \end{bmatrix}\right) = \det\left(\begin{bmatrix} a & b_1 \\ c & b_2 \end{bmatrix}\right)$$

$$\therefore x_2 = \frac{\det\left(\begin{bmatrix} a & b_1 \\ c & b_2 \end{bmatrix}\right)}{\det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right)}$$

# Cholesky decomposition

- LU decomposition of symmetric matrix

  Hypothesis: If A is symmetric matrix, $A = LL^T = L^T L$

$$|Lx|^2 = (Lx)^T (Lx) \quad = x^T (L^T L) x$$
$$= x^T A x \quad \text{(if } L^T L = A\text{)}$$

$$\therefore x^T A x \geq 0 \quad \rightarrow A \text{ is semi-positive matrix}$$

If (i) A is semi-positive matrix  (ii) A is symmetric matrix  (iii) A is square matrix

$$A = LL^T = L^T L \quad \Rightarrow \text{cholesky factorisation}$$

$$x^TAx > 0$$

A is positive $\rightarrow$ A x reverse the direction, only changes the magnitude

$$a \cdot b = a^Tb = |a||b| \cos\theta \quad \rightarrow \quad a^Tb > 0 \quad \text{if} \quad -\frac{\pi}{2} < \theta < \frac{\pi}{2}$$

$x^T(Ax)$ : Dot product b/w original $x$ and linearly transformed $Ax$

$\quad = $ Angle difference b/w $x$ and $Ax$ is $\left[-\frac{\pi}{2}, \frac{\pi}{2}\right]$

- positive definite matrix and eigenvalue

$$Ax = \lambda x$$
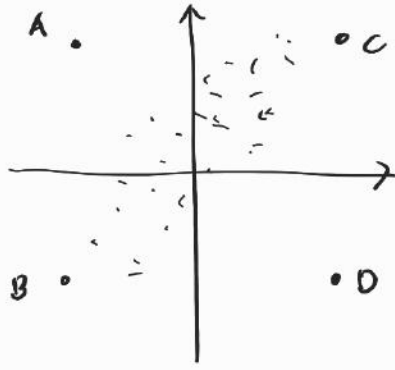$$x^TAx = x^T\lambda x = \lambda|x|^2 > 0 \qquad \therefore \lambda > 0$$

$\therefore$ All eigenvalues are positive

- positive definite matrix and Hessian matrix

If M is positive definite matrix, f is convex downwards (have local minimum)

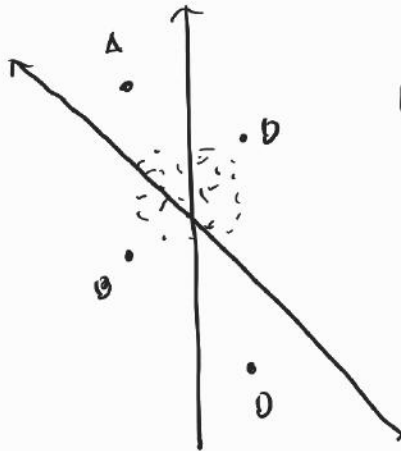# Mahalanobis distance

→ Contextual relative distance



Euclidean distance : $d(AD) = d(BC)$

Mahalanobis distance : $d(AD) > d(BC)$

$$d_E = \sqrt{(x-y)(x-y)^T}$$

$$d_M = \sqrt{(x-y)\, \Sigma^{-1}(x-y)^T}$$

context of data

Normalizing the context