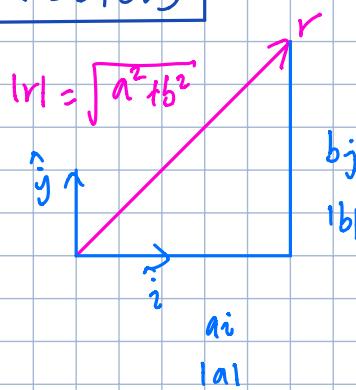


Mathematics for ML Specialisation

Course 1: Linear Algebra

vectors



dot product

$$s = \begin{bmatrix} -1 \\ 2 \end{bmatrix} = \begin{bmatrix} s_i \\ s_j \end{bmatrix}$$

$$r = \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} r_i \\ r_j \end{bmatrix}$$

$$r \cdot s = r_i s_i + r_j s_j = 1$$

commutative:

$$r \cdot s = s \cdot r$$

distributive over addition

$$r \cdot (s+t) = r \cdot s + r \cdot t$$

associative over scalar multiplication

$$r \cdot (as) = a(r \cdot s) \quad a \in \mathbb{N}$$

$$r \cdot r = r_1 r_1 + r_2 r_2$$

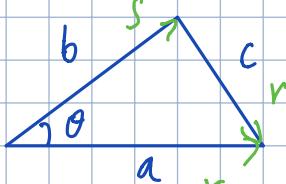
$$= r_1^2 + r_2^2$$

$$= (\sqrt{r_1^2 + r_2^2})^2$$

$$= |r|^2$$

$$\hookrightarrow |r| = \sqrt{r \cdot r}$$

cosine rule



$$c^2 = a^2 + b^2 - 2ab \cos \theta$$

||

$$|r-s|^2 = |r|^2 + |s|^2 - 2|r||s|\cos\theta$$

$$(r-s) \cdot (r-s) = r \cdot r - s \cdot r - s \cdot r + s \cdot s$$

$$= |r|^2 - 2s \cdot r + |s|^2$$

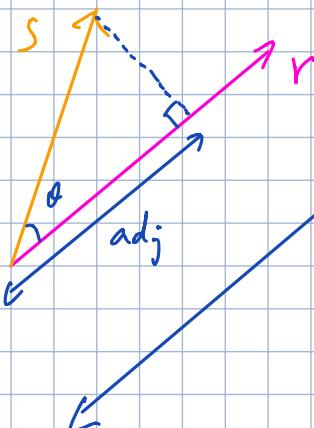
$$r \cdot s = |r| |s| \cos \theta$$

i.e. $\cos 90^\circ = 0 \Rightarrow r \cdot s = 0$

$\cos 0 = 1 \Rightarrow r \cdot s = |r||s| \rightarrow r \cdot s$

$\cos 180 = -1 \Rightarrow r \cdot s = -|r||s| \leftarrow r \cdot s$

Projection



$$\cos \theta = \frac{\text{adj}}{\text{hyp}} = \frac{\text{adj}}{|s|}$$

$$r \cdot s = |r| |s| \cos \theta$$

$|r| \times \text{projection}(s \text{ onto } r)$

$$\frac{r \cdot s}{|r|} = |s| \cos \theta$$

scalar projection

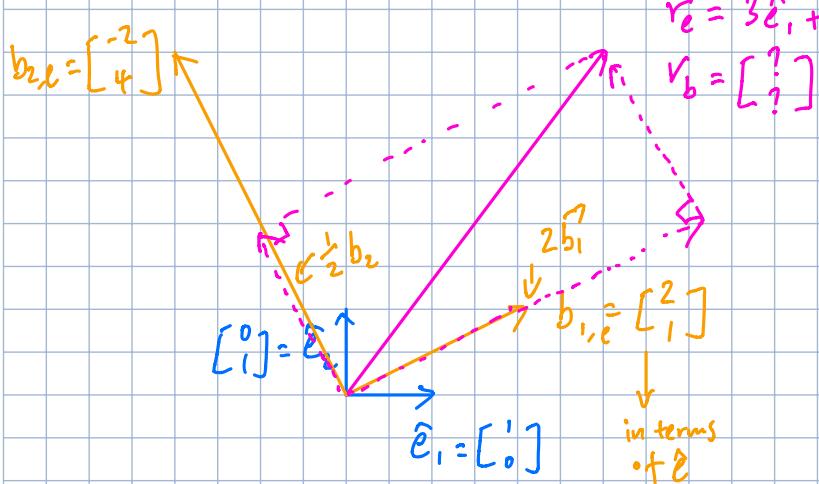
→ how much s goes along r

$$\frac{r \cdot s}{|r|} \frac{r}{|r|}$$

vector projection

→ scalar projection encoded with direction along r

Changing basis



choice of $\{\hat{e}_1, \hat{e}_2\}$ is arbitrary, we can get another basis $\{b_1, b_2\}$

i.e. b_1, e and b_2, e

if we know what $\{b_1, b_2\}$ in terms of $\{\hat{e}_1, \hat{e}_2\} \Rightarrow$ find out r_b
 if $b_2 \perp b_1 \Rightarrow$ can use dot product

else construct transformation matrix \Rightarrow act on the basis $e_i \rightarrow b_i$

$$\frac{r_e \cdot b_1}{\|b_1\|^2} \vec{b}_1 = \frac{3x^2 + 4x^1}{2^2 + 1^2} = \boxed{2 \vec{b}_1} = 2 \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$
$$\frac{r_e \cdot b_2}{\|b_2\|^2} \vec{b}_2 = \boxed{\frac{1}{2} \vec{b}_2} = \frac{1}{2} \begin{bmatrix} -2 \\ 4 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

$r_e = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$
 $r_b = \begin{bmatrix} 2 \\ 1/2 \end{bmatrix}$

Basis

1) linearly independent

2) span the space

the space is n-dimensional

Matrix

matrix transform space

$$e'_1 = \begin{bmatrix} 2 \\ 1/2 \end{bmatrix}$$

$$\begin{pmatrix} 2 & 3 \\ 1/2 & 1 \end{pmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 8 \\ 1/2 \end{bmatrix}$$

$$A \quad r = r'$$

$$\hookrightarrow A(n\hat{e}_1 + m\hat{e}_2) = nA\hat{e}_1 + mA\hat{e}_2$$
$$r' = n e'_1 + m e'_2$$

i.e. A transforms the basis

\hookrightarrow into a different 'grid system'.
n, m unchanged

$$\begin{bmatrix} 0 \\ 1 \end{bmatrix} = \hat{e}_2 \rightarrow e'_2 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

$$\hat{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

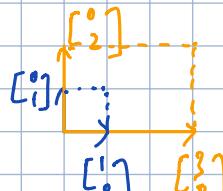
$$\begin{array}{l} \hat{e}_1 \rightarrow e'_1 \\ \hat{e}_2 \rightarrow e'_2 \end{array} \quad \left. \begin{array}{l} \text{transform} \\ \text{by } A \end{array} \right.$$

Identity matrix :

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix}$$

scale matrix:

$$\begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 3x \\ 2y \end{bmatrix}$$



invert matrix

$$\begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -x \\ -y \end{bmatrix}$$



mirror

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 & -1 \\ -1 & 0 \end{pmatrix}$$

$$\begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

shear

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

rotation

$$\begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \quad 2D$$

$$\begin{pmatrix} \cos\theta & -\sin\theta & 0 \\ \sin\theta & \cos\theta & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{rotate about } z\text{-axis} \\ \hookrightarrow z \text{ is preserved}$$

composition of transformations

$$A_3(A_2 A_1) = (A_3 A_2) A_1$$

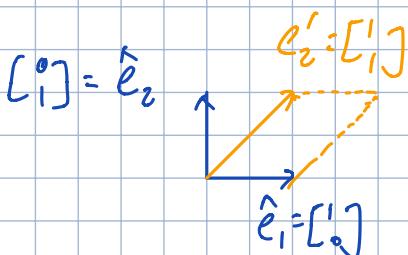
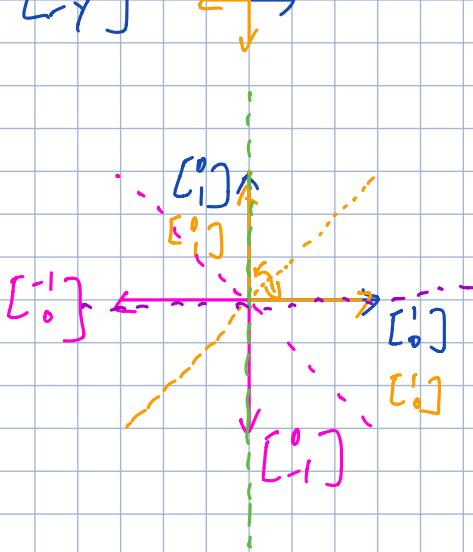
$$A_3 A_2 A_1 \neq A_3 A_1 A_2$$

Gaussian elimination

Given $A \vec{r} = \vec{s}$ e.g. $\begin{pmatrix} 2 & 3 \\ 10 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 8 \\ 13 \end{pmatrix}$, solve $\begin{pmatrix} a \\ b \end{pmatrix}$

$$\hookrightarrow A^{-1} A \vec{r} = A^{-1} \vec{s} \Rightarrow \text{find } A^{-1} \Rightarrow \text{get } \vec{r} = \begin{pmatrix} a \\ b \end{pmatrix}$$

more complicated (just Gaussian elimination, nothing to do with inverse)



$$\begin{pmatrix} 1 & 1 & 3 \\ 1 & 2 & 4 \\ 1 & 1 & 2 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 15 \\ 21 \\ 13 \end{pmatrix}$$

$$\hookrightarrow \begin{pmatrix} 1 & 1 & 3 \\ 0 & 1 & 1 \\ 0 & 0 & +1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 15 \\ 6 \\ +2 \end{pmatrix} \rightarrow \rightarrow \rightarrow \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 5 \\ 4 \\ 2 \end{pmatrix} *$$

↳ only solves a, b, c for a specific $\begin{bmatrix} 15 \\ 21 \\ 13 \end{bmatrix}$, i.e. specific \vec{s}

↳ for any \vec{s} , find A^{-1}

If $AB = I$ C.i.e. $B = A^{-1}$ \Rightarrow find B (or A^{-1}) i.e. all b_{ij}

$$\hookrightarrow \begin{pmatrix} 1 & 1 & 3 \\ 1 & 2 & 4 \\ 1 & 1 & 2 \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

x-axis e.g.

$$\hookrightarrow \begin{pmatrix} A \\ A \\ A \end{pmatrix} \begin{pmatrix} b_{11} \\ b_{21} \\ b_{31} \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \xrightarrow{\text{Gaussian elimination}}$$

↳ 2nd column of B

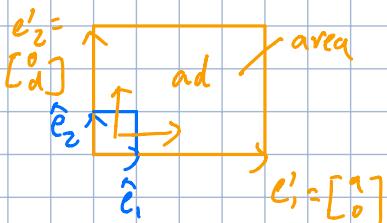
↳ 3rd column of B

OR do all at once

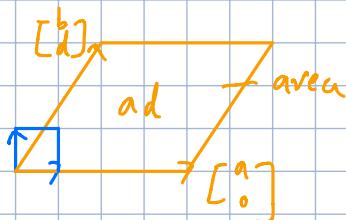
$$\hookrightarrow \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} B \\ B \\ B \end{pmatrix} = \begin{pmatrix} 0 & -1 & 2 \\ -2 & 1 & 1 \\ 1 & 0 & -1 \end{pmatrix} = B$$

determinant

$$\begin{pmatrix} a & 0 \\ 0 & d \end{pmatrix}$$



$$\begin{pmatrix} a & b \\ 0 & d \end{pmatrix}$$



$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

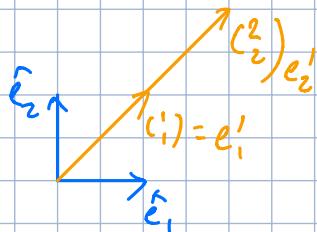
$$\text{area} = ad - bc$$

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

$$A^{-1} = \frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

$$A = \begin{pmatrix} 1 & 2 \\ 1 & 2 \end{pmatrix}$$

$$\det(A) = 0$$



\hookrightarrow no inverse $2D \rightarrow 1D$ \because basis is linearly dependent

\nwarrow cannot 'undo' this transformation

\hookrightarrow ALWAYS check the basis on the transformation

matrix is linearly independent, if not dimension is collapsed during transformation

Einstein Sum Convention

$$\begin{pmatrix} a_{11} & \dots & a_{nn} \\ & \ddots & \\ & & a_{nn} \end{pmatrix}_A \begin{pmatrix} b_{11} & \dots & b_{nn} \\ & \ddots & \\ & & b_{nn} \end{pmatrix}_B = \begin{pmatrix} ab_{11} & \dots & ab_{nn} \\ & \ddots & \\ & & ab_{nn} \end{pmatrix}$$

$$ab_{ik} = \sum_j a_{ij} b_{jk} \stackrel{\text{Einstein}}{=} \sum_j a_{ij} b_{jk}$$

$$AB = C$$

$$c_{ik} = a_{ij} b_{jk}$$

$$U = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}, V = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

same as u_1

$\hat{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$

$\hookrightarrow u_i v_i$

dot product
is symmetric

matrices changing basis

non-orthonormal

$$\hat{e}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \hat{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \hat{u}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \hat{u}_2 = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

$$\frac{1}{2} \begin{bmatrix} 3 \\ 1 \end{bmatrix} \begin{bmatrix} 5 \\ 2 \end{bmatrix}$$

Bear's basis vectors is $\begin{bmatrix} 3 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ in my frame

i.e. Bear's transformation matrix $\begin{bmatrix} 3 & 1 \\ 1 & 1 \end{bmatrix}$, i.e. $\begin{bmatrix} 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \text{bear} \\ \text{rep} \end{bmatrix} = \begin{bmatrix} \text{me} \\ \text{rep} \end{bmatrix}$

$$\begin{bmatrix} 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 3/2 \\ 1/2 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \end{bmatrix}$$

Bear's basis
in my coord.
 \uparrow
 B

$$\Rightarrow B^{-1} \begin{bmatrix} \text{my} \\ \text{rep} \end{bmatrix} = \begin{bmatrix} \text{bear} \\ \text{rep} \end{bmatrix}$$

$$\downarrow$$

$$\frac{1}{2} \begin{bmatrix} 1 & -1 \\ -1 & 3 \end{bmatrix} \begin{bmatrix} 5 \\ 2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

my basis in
Bear's coord

orthonormal

$$\frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \frac{1}{2} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

Bear's transformation matrix

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

B

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

using Projection: if orthogonal

$$\begin{aligned} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \cdot \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} &= \frac{1}{2} \cdot 4 = 2 \\ \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \cdot \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \end{bmatrix} &= \frac{1}{2} \cdot 2 = 1 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

(normalised vector)

Say $B = \begin{pmatrix} 3 & 1 \\ 1 & 1 \end{pmatrix}$ bear's basis

and a vector $\begin{bmatrix} x \\ y \end{bmatrix}$ defined in bear's basis, and we like to do a rotation \rightarrow hard :: funny basis

In my $\{\hat{e}_1, \hat{e}_2\}$ basis, transformation $\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} = R$

① transform $\begin{bmatrix} x \\ y \end{bmatrix}$ into my basis

$$B \begin{bmatrix} x \\ y \end{bmatrix}$$

② Apply R in my representation

$$R B \begin{bmatrix} x \\ y \end{bmatrix}$$

③ transform back to bear's system

$$B^{-1} R B \begin{bmatrix} x \\ y \end{bmatrix}$$

orthogonal matrices

if $A_{n \times n}$ is defined as

$$A_{ij}^T = \delta_{ij}$$

$$\left[\begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_i \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_j \end{pmatrix} \cdots \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} \right] \Rightarrow A^T A = I$$

$$\hookrightarrow A^T = A^{-1}$$

$$a_i \cdot a_j = \delta_{ij}$$

i.e. a_i is orthonormal
to each other

in ML, use orthonormal basis if possible

\hookrightarrow means orthonormal transform matrix

\hookrightarrow transformation is reversible

\hookrightarrow inverse is easy to find

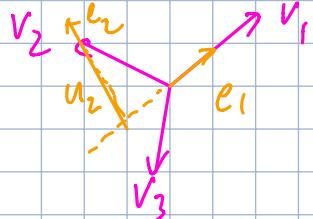
\hookrightarrow projection just dot product

Gram-Schmidt

{linearly independent vectors} \Rightarrow {orthonormal basis set}
that span the spa

$$V = \{v_1, v_2, \dots, v_n\}$$

$$v_1 \rightarrow e_1 = \frac{v_1}{\|v_1\|}$$



$$v_2 = (v_2 \cdot e_1) \frac{e_1}{\|e_1\|} + u_2$$

$$u_2 = v_2 - (v_2 \cdot e_1) e_1 \Rightarrow \frac{u_2}{\|u_2\|} = e_2$$

$$v_3 - (v_3 \cdot e_1) e_1 - (v_3 \cdot e_2) e_2 = u_3 \Rightarrow \frac{u_3}{\|u_3\|} = e_3$$

calculating eigen vector / eigen val

$$A\vec{x} = \lambda \vec{x}$$

A is $n \times n$

x is n col vector
 n dimension

$$\underbrace{(A - \lambda I)}_{\text{!} !} x = 0$$

$$\hookrightarrow \det(A - \lambda I) = 0$$

Simple example: 2×2 transformation matrix

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

$$\hookrightarrow \det \left(\begin{pmatrix} a & b \\ c & d \end{pmatrix} - \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} \right) = 0 \Rightarrow \lambda^2 - (a+d)\lambda + ad - bc = 0$$

$$\text{if } A = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \Rightarrow \det \begin{pmatrix} 1-\lambda & 0 \\ 0 & 2-\lambda \end{pmatrix} = (1-\lambda)(2-\lambda) = 0$$
$$\lambda = 1 \text{ or } 2$$

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

$$\textcircled{1} \quad \lambda = 1 \quad \left[\begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} - \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right] x = 0$$
$$\begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$
$$\begin{pmatrix} 0 \\ x_2 \end{pmatrix} = 0$$

eigenvector $x = \begin{pmatrix} t \\ 0 \end{pmatrix}$
any vector along horizontal

$$\textcircled{2} \quad \lambda = 2 \quad \left[\begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} - \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \right] x = 0$$
$$\begin{pmatrix} -1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$
$$\begin{pmatrix} -x_1 \\ 0 \end{pmatrix} = 0$$
$$\begin{pmatrix} x_1 \\ 0 \end{pmatrix} = 0$$

$\Rightarrow x = \begin{pmatrix} 0 \\ t \end{pmatrix}$
any vertical vector

e.g. rotation 90°

$$A = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \Rightarrow \det \left(\begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} - \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \right) = \det \begin{pmatrix} -2 & -1 \\ 1 & -2 \end{pmatrix}$$
$$= \lambda^2 + 1 = 0$$

\hookrightarrow no real solns \Rightarrow no eigenvalue

changing to the eigenbasis

motivation:

Given transformation $T = \begin{pmatrix} 0.9 & 0.8 \\ -1 & 0.35 \end{pmatrix}$

and if $V_0 = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}$, we find $V_1 = TV_0$

$$\hookrightarrow V_2 = TV_1 = T^2 V_0$$

$$\hookrightarrow V_3 = TV_2 = T^3 V_0$$

\Rightarrow tedious to find T^n

if T is a diagonal matrix: $T = \begin{pmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & c \end{pmatrix}$

$$\hookrightarrow T^n = \begin{pmatrix} a^n & 0 & 0 \\ 0 & b^n & 0 \\ 0 & 0 & c^n \end{pmatrix}$$

if T is NOT a diagonal matrix \Rightarrow change to a basis where T becomes diagonal \Rightarrow eigenbasis

\hookrightarrow get T^n

\hookrightarrow transform back

\because each column of transformation matrix represents the new location of the transformed unit vector

\therefore use eigenvectors as columns of the eigenbasis transformation matrix

transform $C = \begin{pmatrix} \pi_1 & \pi_2 & \pi_3 \\ \cdot & \cdot & \cdot \end{pmatrix}$ π_i are eigenvectors
(for T^n)

$$D = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{pmatrix} \quad \lambda_i \text{ are eigenvalues of } T$$

After applying C to the vector v_0 , applying T is just scaling $\Rightarrow T$ is diagonal in this frame, which is D , it contains eigenvalues of T

$$\text{i.e. } T \iff D$$

our basis eigen basis

This means that:

$$\text{Apply } T = \begin{array}{l} \textcircled{1} \text{ convert to eigenbasis} \\ \textcircled{2} \text{ apply the diagonalise matrix} \\ \textcircled{3} \text{ convert back} \end{array}$$

$$\text{L} \hookrightarrow T = C D C^{-1} \iff D = C^{-1} T C$$

$$\begin{aligned} \hookrightarrow T^2 &= C D C^{-1} C D C^{-1} \\ &= C D^2 C^{-1} \end{aligned}$$

$$T^n = C D^n C^{-1}$$

eigenbasis example: 2×2

$$T = \begin{pmatrix} 1 & 1 \\ 0 & 2 \end{pmatrix}$$

consider the $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ vector (old orange vector)

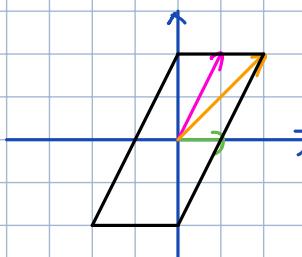
$$\begin{pmatrix} 1 & 1 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

$$\text{eigen vector } v = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \lambda = 1$$

$$v = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \lambda = 2$$

Now consider $\begin{pmatrix} -1 \\ 1 \end{pmatrix}$ vector

$$T^2 \begin{pmatrix} -1 \\ 1 \end{pmatrix} : \quad \begin{pmatrix} 1 & 1 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 2 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 & 1 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 0 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \end{pmatrix} \quad \text{Same}$$



$$T^2 = \begin{pmatrix} 1 & 1 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \end{pmatrix}$$

eigen basis way:

$$C = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \Rightarrow C^{-1} = \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix} \quad \text{same}$$

$$\Rightarrow T^2 = C D^2 C^{-1} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}^2 \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix}$$

Course 2 : Multivariate Calculus

Gradient

in 2D, gradient of $f(x)$ = $\lim_{\Delta x \rightarrow 0} \left(\frac{f(x + \Delta x) - f(x)}{\Delta x} \right) = f'(x) = \frac{df}{dx}$

$$\text{e.g } f(x) = 3x + 2$$

$$\begin{aligned} \hookrightarrow f'(x) &= \lim_{\Delta x \rightarrow 0} \left(\frac{3(x + \Delta x) + 2 - 3x - 2}{\Delta x} \right) \\ &= \lim_{\Delta x \rightarrow 0} \left(\frac{3x + 3\Delta x + 2 - 3x - 2}{\Delta x} \right) \\ &= \lim_{\Delta x \rightarrow 0} \left(\frac{3\Delta x}{\Delta x} \right) \\ &= 3 \end{aligned}$$

sum rule:

$$\frac{d}{dx} (f(x) + g(x)) = \frac{df(x)}{dx} + \frac{dg(x)}{dx}$$

product rule:

$$\frac{d}{dx} (f(x)g(x)) = f(x)g'(x) + g(x)f'(x)$$

$$\frac{\partial}{\partial x} (f(x,y,z) + g(x,y,z)) = f'(x,y,z) + g'(x,y,z)$$

chain rule:

$$h = h(p) \quad p = p(m)$$

$$\frac{dh}{dm} = \frac{dh}{dp} \cdot \frac{dp}{dm}$$

Total derivative

$$\frac{df(x,y,z)}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} + \frac{\partial f}{\partial z} \frac{dz}{dt}$$

e.g. $f(x,y,z) = \sin(x) e^{yz^2}$ and $\begin{cases} x = t-1 \\ y = t^2 \\ z = yt \end{cases}$

Jacobian

if we have $f(\pi_1, \pi_2, \pi_3, \dots)$

Jacobian vector points uphill

Jacobian is a vector where each entry is the partial derivative of f w.r.t. each one of the variable

$$J = \left[\frac{\partial f}{\partial \pi_1}, \frac{\partial f}{\partial \pi_2}, \frac{\partial f}{\partial \pi_3}, \dots \right] \quad \text{usually row vector}$$

e.g. $f(\pi, y, z) = \pi^2 y + 3z$

$$\frac{\partial f}{\partial \pi} = 2\pi y$$

$$\frac{\partial f}{\partial y} = \pi^2$$

$$\frac{\partial f}{\partial z} = 3$$

$$J = [2\pi y, \pi^2, 3]$$

when (π, y, z) is supplied, J will point in the

direction of steepest slope of this function.

e.g. $J(0,0,0) = [0, 0, 3]$

↳ Steeper \Rightarrow bigger Jacobian

visualise in 2-D $\Rightarrow f(x,y)$ where $f \Rightarrow$ is height

Jacobian Matrix

Jacobian matrix describes functions and the gradient of a multi-variable system.

Consider 2 functions:

$$\begin{aligned} u(x,y) &= x - 2y \\ v(x,y) &= 3y - 2x \end{aligned} \quad \left. \begin{array}{l} \text{can make separate row vector} \\ \text{Jacobian} \end{array} \right\}$$

$$\left. \begin{array}{l} J_u = \left[\frac{\partial u}{\partial x}, \frac{\partial u}{\partial y} \right] \\ J_v = \left[\frac{\partial v}{\partial x}, \frac{\partial v}{\partial y} \right] \end{array} \right\} J = \begin{bmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{bmatrix}$$

$$\left. \begin{array}{l} \frac{\partial u}{\partial x} = 1, \quad \frac{\partial u}{\partial y} = -2 \\ \frac{\partial v}{\partial x} = -2, \quad \frac{\partial v}{\partial y} = 3 \end{array} \right\} J = \begin{bmatrix} 1 & -2 \\ -2 & 3 \end{bmatrix}$$



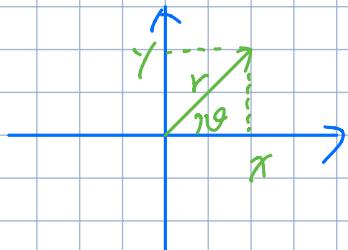
gradient is constant everywhere



this T transforms vector from xy space \rightarrow uv space

$$J \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} -4 \\ 5 \end{bmatrix}$$

e.g transforming Cartesian \rightarrow polar coordinate



$$x(r, \theta) = r \cos \theta$$

$$y(r, \theta) = r \sin \theta$$

$$\hookrightarrow J = \begin{bmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{bmatrix} = \begin{bmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{bmatrix}$$

$$|J| = r (\cos^2 \theta + \sin^2 \theta) = r$$

\uparrow
determinant

Hessian

Jacobian : collect all the 1st order derivatives of a function in to vector
a vector.

Hessian : 2nd order derivatives into a matrix

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

shorthand

$$\frac{\partial^2 f}{\partial x_1 \partial x_1} \equiv \frac{\partial^2 f}{\partial x_1^2}$$

$$\frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{\partial}{\partial x_1} \left(\frac{\partial f}{\partial x_2} \right)$$

$\hookrightarrow n \times n$ square matrix

$\hookrightarrow n = \text{no. of variable in function } f$

e.g. $f(x, y, z) = xyz$

$$J = \begin{bmatrix} 2xyz, x^2z, xy \end{bmatrix}$$

$\swarrow \quad \searrow \quad \downarrow$

$$H = \begin{bmatrix} 2x^2 & 2xy & 2xz \\ 2yz & 0 & x^2 \\ 2xy & x^2 & 0 \end{bmatrix}$$

\overrightarrow{x} \overrightarrow{y} \overrightarrow{z}

symmetrical about diagonal
⇒ True if func is
(continuous)

e.g. 2D example $f(x,y) = x^2 + y^2$

$$J = [2x, 2y]$$

$$H = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$|H| = 4$ { +ve: max or min
-ve: saddle }
 ↳ look at $H[0,0]$
 ↳ true: min
 -ve: max

Multivariable chain rule

if we have

$$f(x_1, x_2, x_3, \dots, x_n) = f(\vec{x})$$

where $\vec{x}_1(t), \vec{x}_2(t)$

we would like to find $\frac{df}{dt}$

we need

$$\frac{\partial f}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

$$\frac{d\vec{x}}{dt} = \begin{bmatrix} \frac{dx_1}{dt} \\ \vdots \\ \frac{dx_n}{dt} \end{bmatrix}$$

$$\frac{df}{dt} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \cdot \begin{bmatrix} \frac{dx_1}{dt} \\ \vdots \\ \frac{dx_n}{dt} \end{bmatrix} = \frac{\partial f}{\partial \vec{x}} \cdot \frac{d\vec{x}}{dt}$$

i.e. if $f(\vec{x}(t))$, then $\frac{df}{dt} = \frac{\partial f}{\partial \vec{x}} \cdot \frac{d\vec{x}}{dt}$

and $\frac{\partial \vec{f}}{\partial \vec{x}}$ is just same as a Jacobian vector (but col vector instead of row vector)

$$\hookrightarrow \frac{\partial \vec{f}}{\partial \vec{x}} = (\mathbf{J}_f)^\top$$

$$\Rightarrow \frac{df}{dt} = \mathbf{J}_f \frac{d\vec{x}}{dt}$$

(row/col doesn't matter here
in operation)

e.g. more than 2 links, univariable (x)

$$f(u) = 5x$$

$$x(u) = 1-u$$

$$u(t) = t^2$$

method 1:

$$f(t) = 5(1-u) = 5(1-t^2) = 5-5t^2$$

$$\frac{df}{dt} = -10t$$

method 2:

$$\frac{df}{dt} = \frac{df}{du} \frac{dx}{du} \frac{du}{dt} = (5)(-1)(2t) = -10t$$

multi variable, more than 2 links

$$f(\vec{x}(u(t)))$$

$$f(\vec{x}) = f(x_1, x_2)$$

$$\vec{x}(u) = \begin{bmatrix} x_1(u_1, u_2) \\ x_2(u_1, u_2) \end{bmatrix}$$

$t \rightarrow f$

$$\vec{x}(t) = \begin{bmatrix} u_1(t) \\ u_2(t) \end{bmatrix}$$

$$\frac{df}{dt} = \frac{\partial f}{\partial \vec{x}} \frac{\partial \vec{x}}{\partial \vec{u}} \frac{d\vec{u}}{dt}$$

$$= \left[\frac{\partial f}{\partial u_1}, \frac{\partial f}{\partial u_2} \right] \begin{bmatrix} \frac{\partial x_1}{\partial u_1} & \frac{\partial x_1}{\partial u_2} \\ \frac{\partial x_2}{\partial u_1} & \frac{\partial x_2}{\partial u_2} \end{bmatrix} \begin{bmatrix} \frac{du_1}{dt} \\ \frac{du_2}{dt} \end{bmatrix}$$

Jacobian
vector
of f

Jacobian
matrix
of x

derivative
vector
of u

$(1 \times 1) = (1 \times 2) (2 \times 2) (2 \times 1) \Rightarrow$ returns a scalar

A simple neural network



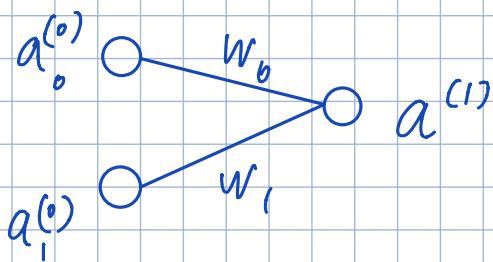
$$a^{(1)} = \sigma(w a^{(0)} + b)$$

a : activity

w : weight

b : bias

σ : activation function

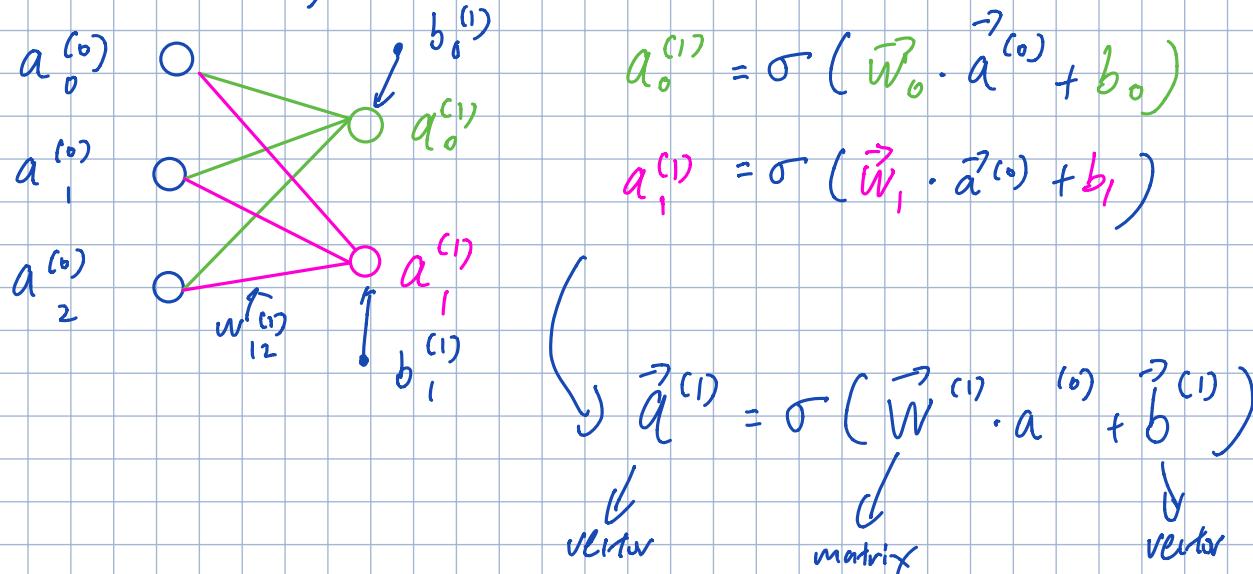


$$\text{Now } a^{(1)} = \sigma(w_0 a_0^{(0)} + w_1 a_1^{(0)} + b)$$

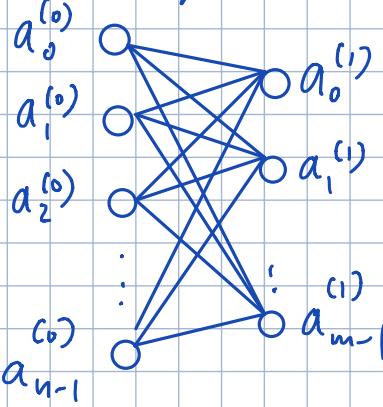
for n input

$$a^{(1)} = \sigma \left[\left(\sum_{j=1}^n w_j a_j^{(0)} \right) + b \right] = \sigma(\vec{w} \cdot \vec{a}^{(0)} + b)$$

adding more output



Single layer



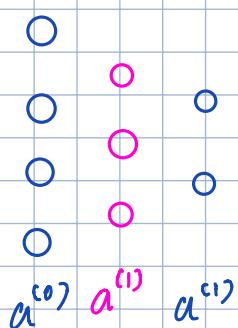
To describe this function

$$\vec{a}^{(1)} = \sigma(\vec{W}^{(1)} \cdot \vec{a}^{(0)} + \vec{b}^{(1)})$$

$$\begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \\ \vdots \\ a_{m-1}^{(1)} \end{bmatrix} = \sigma \left(\begin{bmatrix} w_{0,0}^{(1)} & w_{0,1}^{(1)} & \dots & w_{0,n-1}^{(1)} \\ w_{1,0}^{(1)} & \ddots & & \\ \vdots & & \ddots & \\ w_{m-1,0}^{(1)} & & \ddots & w_{m-1,n-1}^{(1)} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_{n-1}^{(0)} \end{bmatrix} + \begin{bmatrix} b_0^{(1)} \\ b_1^{(1)} \\ \vdots \\ b_{m-1}^{(1)} \end{bmatrix} \right)$$

hidden layer

$w_{ij}^{(1)}$ is the link between neuron j in previous layer
i in current layer

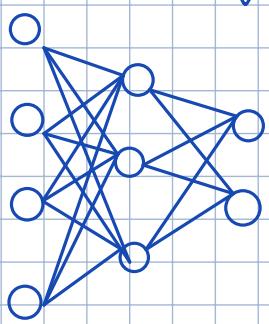


$$a^{(1)} = \sigma(\vec{W}^{(1)} \cdot \vec{a}^{(0)} + b^{(1)})$$

$$a^{(2)} = \sigma(\vec{W}^{(2)} \cdot \vec{a}^{(1)} + b^{(2)})$$

$$\text{i.e. } a^{(L)} = \sigma(\vec{w}^{(L)} \cdot \vec{a}^{(L-1)} + \vec{b}^{(L)})$$

Back propagation



$$\vec{a}^{(L)} = \sigma(\vec{w}^{(L)} \cdot \vec{a}^{(L-1)} + \vec{b}^{(L)})$$

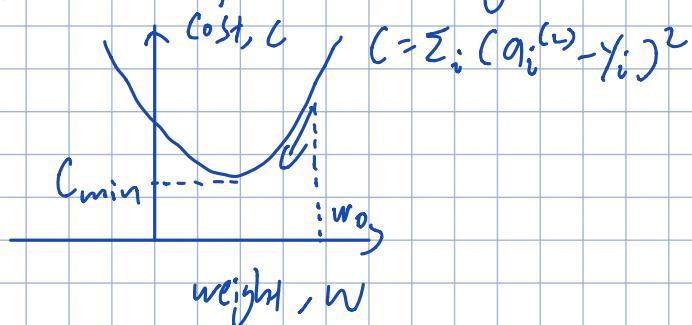
find 18 weights and 5 bias
that cause our network to best match
train inputs to their labels.
(Start with random numbers)

then we define a cost function

$$C = \sum_i^T (a_i^{(L)} - y_i)^2$$

desired output

e.g. of a particular weight



e.g. 2-node network

$$a^{(0)} \quad a^{(1)}$$

$$a^{(1)} = \sigma(wa^{(0)} + b)$$

$$C = (a^{(1)} - y)^2$$

$$\begin{cases} \frac{\partial C}{\partial w} = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial w} \\ \frac{\partial C}{\partial b} = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial b} \end{cases}$$

for switching activation function
usually

$$z^{(1)} = wa^{(0)} + b \quad C = (a^{(1)} - y)^2$$

$$a^{(1)} = \sigma(z^{(1)})$$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial w}$$

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial b}$$

For multiple neurons / layer L

$$\vec{z}^{(L)} = \vec{W}^{(L)} \cdot \vec{a}^{(L-1)} + \vec{b}^{(L)}$$

$$\vec{a}^{(L)} = \sigma(\vec{z}^{(L)})$$

$$C_k = \sum_i (\vec{a}_i^{(L)} - y_i)^2$$

for 1 training example (1 row) - k^{th} data sample

for multiple output (still 1 data sample)

$$C_k = \|\vec{a}^{(k)} - \vec{y}\|^2 \quad \text{i.e. feature}$$

e.g. get derivative of cost w.r.t to the weight of final layer
(assume 3 layers, 4 input, 2 output), for k^{th} data

$$\frac{\partial C_k}{\partial \vec{W}^{(m)}} = \frac{\partial C_k}{\partial \vec{a}^{(k)}} \frac{\partial \vec{a}^{(k)}}{\partial \vec{z}^{(k)}} \frac{\partial \vec{z}^{(k)}}{\partial \vec{W}^{(k)}}$$

derivative of cost w.r.t weight of previous layer

$$\frac{\partial C_k}{\partial \vec{W}^{(m)}} = \frac{\partial C_k}{\partial \vec{a}^{(k)}} \frac{\partial \vec{a}^{(k)}}{\partial \vec{a}^{(k)}} \frac{\partial \vec{a}^{(k)}}{\partial \vec{z}^{(k)}} \frac{\partial \vec{z}^{(k)}}{\partial \vec{W}^{(k)}}$$

$$\text{where } \frac{\partial \vec{a}^{(k)}}{\partial \vec{a}^{(k)}} = \frac{\partial \vec{a}^{(k)}}{\partial \vec{z}^{(k)}} \frac{\partial \vec{z}^{(k)}}{\partial \vec{a}^{(k)}}$$

generalise

$$\frac{\partial C_k}{\partial \vec{W}^{(i)}} = \frac{\partial C_k}{\partial \vec{a}^{(N)}} \frac{\partial \vec{a}^{(N)}}{\partial \vec{a}^{(N-1)}} \frac{\partial \vec{a}^{(N-1)}}{\partial \vec{a}^{(N-2)}} \cdots \frac{\partial \vec{a}^{(i+1)}}{\partial \vec{a}^{(i)}} \frac{\partial \vec{a}^{(i)}}{\partial \vec{z}^{(i)}} \frac{\partial \vec{z}^{(i)}}{\partial \vec{W}^{(i)}}$$

from layer N to layer i

$$\text{where } \frac{\partial \vec{a}^{(j)}}{\partial \vec{a}^{(j-1)}} = \sigma'(\vec{z}^{(j)}) \vec{W}^{(j)}$$

$$\therefore \vec{a}^{(n)} = \sigma(\vec{z}^{(n)}) \quad (\text{scalar})$$

$$z' = w^{(n)} a^{(n)} + b^{(n)}$$

Taylor Series (Power Series)

$$g(x) = a + bx + cx^2 + dx^3 + \dots$$

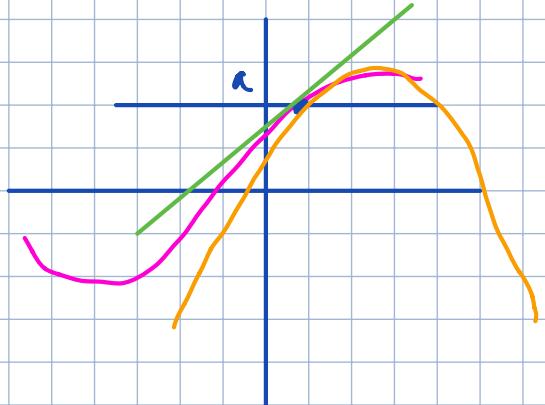
$$g_0(x) = a \quad 0^{\text{th}} \text{ order approximation}$$

$$g_1(x) = a + bx \quad 1^{\text{st}} \text{ order approximation}$$

$$g_2(x) = a + bx + cx^2 \quad 2^{\text{nd}} \text{ order approximation}$$

⋮

e.g.



$$g_0(x) = a \quad 0^{\text{th}} \text{ order approx}$$

$$g_1(x) = a + bx \quad 1^{\text{st}} \text{ order}$$

$$g_2(x) = a + bx + cx^2 \quad 2^{\text{nd}} \text{ order}$$

If we know everything about the function at some point (i.e. its function value, 1st, 2nd... derivatives etc.) we can reconstruct the function everywhere else. Only for well behaved functions.
(i.e. continuous, can be differentiated as many times as we want)

e.g.

choose the point $x=0$



$$g_0(x) = f(0)$$

$$g_1(x) = f(0) + f'(0)x$$

$$g_2(x) = f(0) + f'(0)x + \frac{1}{2}f''(0)x^2$$

$$g_3(x) = f(0) + f'(0)x + \frac{1}{2}f''(0)x^2 + \frac{1}{6}f'''(0)x^3$$

$$\Rightarrow g_4(x) = \frac{f^{(4)}(0)}{0!} + \frac{f^{(5)}(0)}{1!}x + \frac{f^{(6)}(0)}{2!}x^2 + \frac{f^{(7)}(0)}{3!}x^3 + \frac{f^{(8)}(0)}{4!}x^4$$

n^{th} term: $\frac{1}{n!} f^{(n)}(0) x^n$

$$\Rightarrow g(x) = \sum_{n=0}^{\infty} \frac{1}{n!} f^{(n)}(0) x^n$$

MacLaurin series
 $(\because \text{at } x=0)$
 or special case

$$g(x) = \sum_{n=0}^{\infty} \frac{1}{n!} f^{(n)}(p) (x-p)^n$$

1-D Taylor Series

at point p

$$g_1(x) = f(p) + f'(p)(x-p)$$

$(\because (x-p) \times f'(p) = \text{rise})$
 run \times gradient

approximation near p

$$g_1(p+\Delta p) = f(p) + f'(p)(\Delta p)$$

Substitute to Δx (same)

$$g_1(x+\Delta x) = f(x) + f'(x)(\Delta x)$$

Taylor Series

$$f(x+\Delta x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(x)}{n!} \Delta x^n$$

(focus on any 1 point $x+\Delta x$)
 where $\Delta x = x-p$

1st order approx. (linearisation)

$$f(x+\Delta x) = f(x) + f'(x)(\Delta x) + O(\Delta x^2)$$

Multi-dimensional Power Series

$$f(x+\Delta x, y+\Delta y)$$

$$= f(x,y) + (\partial_x f(x,y) \Delta x + \partial_y f(x,y) \Delta y)$$

n^{th}
1st order

$$\rightarrow [\partial_x f(x,y), \partial_y f(x,y)] \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix}$$

$$\int_f \Delta x$$

$$+ \frac{1}{2} (\partial_{xx} f(x,y) \Delta x^2 + 2\partial_{xy} f(x,y) \Delta x \Delta y + \partial_{yy} f(x,y) \Delta y^2) + \dots$$

$$\frac{1}{2} [\Delta x, \Delta y] \begin{bmatrix} \partial_{xx} f(x,y) & \partial_{xy} f(x,y) \\ \partial_{yx} f(x,y) & \partial_{yy} f(x,y) \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix}$$

||

$$\frac{1}{2} \Delta \vec{x}^T H_f \Delta \vec{x}$$

i.e.

$$f(\vec{x} + \Delta \vec{x}) = f(\vec{x}) + J_f \Delta \vec{x} + \frac{1}{2} \Delta \vec{x}^T H_f \Delta \vec{x} + \dots$$

for multi dimension Taylor Series

Newton-Raphson method

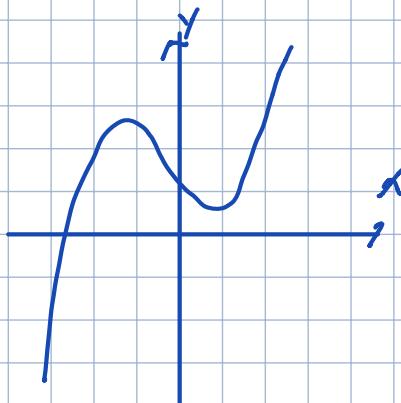
$$1.) \quad y = x^3 - 2x + 2$$

$$\frac{dy}{dx} = 3x^2 - 2$$

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

Newton
-Raphson
method

prev
guess



i	x_i	$y(x_i)$	$\frac{dy(x_i)}{dx}$
0	-2	-2	10
1	-1.8	-0.23	7.7
2	-1.77	-0.005	7.4
3	-1.769	-2.3E-6	

we want to find x when $y=0$

Recall we can ask what value of the function is at point $x_0 + \delta x$, a short distance away

$$f(x_0 + \delta x) = f(x_0) + f'(x_0) \delta x$$

if we assume the function goes to 0 somewhere, we can rearrange to find how far away, i.e. assume $f(x_0 + \delta x) = 0$, find δx :

$$\delta x = -\frac{f(x_0)}{f'(x_0)}$$

(repeat)

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Optimisation: Gradient Descent

In Newton-Raphson method, we try to find the solution x , if $y=0$.

In Gradient Descent, we try to find the minimum value (no need to be 0, and imagine in a hyper-dimension contour plot)

Assume $f(x, y) = x^2 y$

how do we find the fastest/steepest way to get down the 3D graph?

$$\frac{\partial f}{\partial x} = 2xy$$

$$\frac{\partial f}{\partial y} = x^2$$

put those 2 gradient in a vector

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$$

if we have a little vector $\vec{r} = \begin{bmatrix} dx \\ dy \end{bmatrix}$

$\nabla f \cdot \vec{r} \Rightarrow$ move a bit along \vec{r} direction (ds, here)

$$\begin{bmatrix} \frac{\partial f(a,b)}{\partial x} \\ \frac{\partial f(a,b)}{\partial y} \end{bmatrix} \leftarrow \text{evaluate at point } (a,b)$$

If we want to know how much a function will change, when we move along some unit vector (*c* in an opposite direction).

$$\nabla f \cdot \hat{r}$$

$$= \begin{bmatrix} \frac{df(c+d)}{dx} \\ \frac{df(c+d)}{dy} \end{bmatrix} \cdot \begin{bmatrix} c \\ d \end{bmatrix} = df$$

\hat{r}

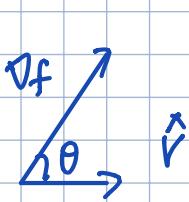
directional gradient

$$\Rightarrow df = \frac{\partial f}{\partial x} c + \frac{\partial f}{\partial y} \cdot d$$

\hat{r}
small
change
in f

$$df$$

What is the maximum value this directional gradient can take?



when $\hat{r} \parallel \nabla f$



max:
value

$$\nabla f \cdot \frac{\nabla f}{\|\nabla f\|} = \frac{\|\nabla f\|^2}{\|\nabla f\|} = \|\nabla f\|$$

↓
steepest
gradient

size of ∇f

∇f points up the direction of the steepest descent
to the contour line

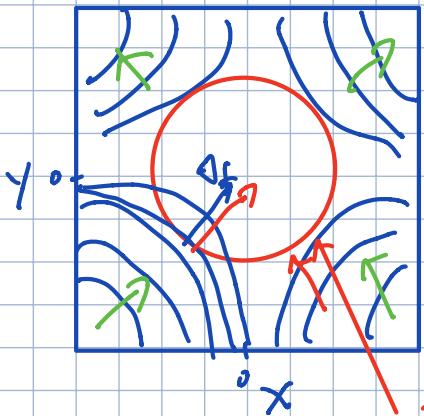
In Gradient descent method, we take small steps down the hill,

$$S_{n+1} = S_n - \gamma \nabla f.$$

\downarrow
 coefficient
 previous position
 \downarrow
 evaluated at
 previous position

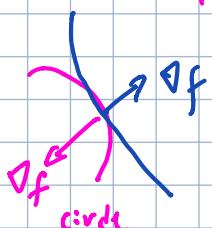
As we approach turning point (min or max), ∇f gets smaller.

Constrained Optimisation : Lagrange Multipliers



Contour map (project to 2D)

We wish to find min/max on the circle



$$\text{maximise: } f(x, y) = x^2 y$$

$$\text{constraint: } g(x, y) = x^2 + y^2 = a^2$$

$$\text{solve } \nabla f = \lambda \nabla g$$

where λ is Lagrange multiplier
 ∇f : grad in center
 ∇g : grad func/gad constraint

max and min on the circle will be found where constraint (circle) // contour line

$$\nabla f = \nabla(x^2 y) = \begin{bmatrix} 2xy \\ x^2 \end{bmatrix} = \lambda \nabla g = \lambda \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

$$① \quad 2xy = \lambda 2x \Rightarrow y = \lambda$$

$$② \quad x^2 = \lambda 2y \Rightarrow x^2 = 2y^2 \Rightarrow x = \pm \sqrt{2}y$$

$$③ \quad x^2 + y^2 = a^2 \Rightarrow 2y^2 + y^2 = a^2 \Rightarrow y = \pm a/\sqrt{3}$$

Soln: $\frac{a}{\sqrt{3}} \begin{pmatrix} \sqrt{2} \\ 1 \end{pmatrix}, \frac{a}{\sqrt{3}} \begin{pmatrix} \sqrt{2} \\ -1 \end{pmatrix}, \frac{a}{\sqrt{3}} \begin{pmatrix} -\sqrt{2} \\ 1 \end{pmatrix}, \frac{a}{\sqrt{3}} \begin{pmatrix} -\sqrt{2} \\ -1 \end{pmatrix}$

$$f(x, y) = \frac{a^3}{3\sqrt{3}} 2, \frac{a^3}{3\sqrt{3}} (-2), \frac{a^3}{3\sqrt{3}} (2), \frac{a^3}{3\sqrt{3}} (-2)$$

max
min
max
min

Simple linear regression

In a straight line data fitting, we can model the straight line y as the i th observation x_i and a vector a of the fitting parameters.

$$Y = Y(x_i; a_i) = m x_i + c \quad a = \begin{bmatrix} m \\ c \end{bmatrix}$$

\bar{x}, \bar{y} will be geometric center or mass of that dataset.

residue $r =$ difference between the data item y_i and the predicted location of those on the line y , which will be $mx_i + c$

$$\Rightarrow r_i = y_i - mx_i - c$$

\downarrow
expt data
point

$$\chi^2 = \sum_i r_i^2 = \sum_i (y_i - mx_i - c)^2$$

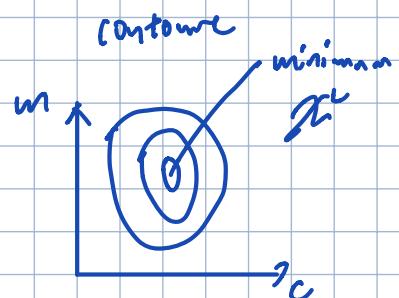
(chi²)

Goal: minimise χ^2

Minimum is found when the gradient of $\chi^2 = 0$

$$\text{i.e. } \nabla \chi^2 = \begin{bmatrix} \frac{\partial \chi^2}{\partial m} \\ \frac{\partial \chi^2}{\partial c} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

note the variable in this space is m and c



$$\begin{aligned} \frac{\partial \chi^2}{\partial m} &= -2 \sum_i x_i (y_i - mx_i - c) \\ \frac{\partial \chi^2}{\partial c} &= -2 \sum_i (y_i - mx_i - c) \end{aligned} \quad \left. \begin{aligned} \nabla \chi^2 &= \begin{bmatrix} -2 \sum_i x_i (y_i - mx_i - c) \\ -2 \sum_i (y_i - mx_i - c) \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \end{aligned} \right\}$$

$$c = \bar{y} - m \bar{x}$$

\bar{x}
 \bar{y}
mean

$$m = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2}$$

solution

$$\sigma_c \approx \sigma_m \sqrt{\bar{x}^2 + \frac{1}{n} \sum_i (x_i - \bar{x})^2}$$

$$\sigma_m^2 \approx \frac{\chi^2}{\sum (x_i - \bar{x})^2 (n-2)}$$

General non-linear least squares

Say we have a function y , trying to fit (x, y)

$$y(\pi; a_k) \quad k=1 \dots m$$

$$e.g \quad (x - a_1)^2 + a_2$$

Are there parameters to evaluate for best fit

data points (y_i, x_i, σ_i)

$i = 1, \dots, n$

- mode 1

$$\chi^2 = \sum_{i=1}^n \frac{[y_i - y(\pi_i; \alpha_k)]^2}{\sigma^2}$$

σ_i : (y -error bar)

\hat{X} sort-of removing uncertainty weightage
from Z^L (if no uncertainty, can $\sigma \rightarrow 1$)

minimum π^2 when

$$\nabla x^2 = 0$$

most of the time we cannot solve analytically

→ solve by steepest descent

by updating the vector of fitting parameters α

$$a_{\text{next}} = a_{\text{curr}} - \text{constant} \cdot \nabla J(a) = a_{\text{curr}} + \sum_{i=1}^n \frac{[y_i - y(x_i; a_k)]}{\sigma^2} \frac{dy}{da}$$

$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \end{bmatrix}$ ↘ learning rate ↗ $\nabla J(a)$

until $\nabla J(\hat{x}) = 0$ (or converge to some value)

we need $\frac{dx^2}{dx_k}$ for each k

$$\frac{d\chi^2}{da_k} = \sum_{i=1}^n -2 \frac{[y_i - y(x_i; a_k)]}{\sigma_i^2} \frac{dy}{da_k}$$

dark parameter

(i.e. $y(x; \alpha_1, \alpha_2)$)

$$\text{L.-J. } \gamma(\pi; a_{12}) = (\pi - a_{11})^2 + a_{12}$$

$$\frac{da_1}{dx} = -L(x-a_1)$$

$$\frac{dy}{da_2} = 1$$

Course 3: PCA

Mean of data set

Say we have a data set $D = \{x_1, \dots, x_N\}$

$$E[D] = \frac{1}{N} \sum_{n=1}^N x_n$$

Expectation
value

e.g. $D' = \{1, 2, 4, 6, 1\}$

$$E[D'] = 1+2+4+6+1/5 = 3.8$$

Variance of data set

In 1D, we can look at the average squared distance of a data point from the mean value of this data set.

$$D_1 = \{1, 2, 4, 5\} \quad E[D_1] = 3$$

$$D_2 = \{-1, 3, 7\} \quad E[D_2] = 3$$

$$D_1: \frac{(1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2}{4} = 10/4$$

$$D_2: \frac{(-1-3)^2 + (3-3)^2 + (7-3)^2}{3} = 32/3$$

\Rightarrow average squared distance

$$X: \{x_1, \dots, x_N\} \quad \nearrow M = E[X]$$

$$\text{Var}[X] = \frac{1}{N} \sum_{n=1}^N (x_n - M)^2$$

$$\hookrightarrow \sqrt{\text{Var}[X]} = \text{standard deviation}[X]$$

\swarrow
Same unit as X

Covariance

$$\text{Cov}[x, y] = E[(x - M_x)(y - M_y)]$$

$$M_x = E[x]$$

$$M_y = E[y]$$

\hookrightarrow for 2D dataset, we can get

$$\begin{matrix} \text{Var}[x] \\ \text{Var}[y] \\ \text{Cov}[x, y] \\ \text{Cov}[y, x] \end{matrix} \quad \left\{ \quad \begin{bmatrix} \text{Var}[x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \text{Var}[y] \end{bmatrix} \right.$$

Covariance matrix

if $\text{Cov}[x, y] > 0$, on avg y increase if we increase x

$= 0$, x and y have nothing to do with each other

\hookrightarrow always a symmetric positive definite matrix

if $D = \{x_1, \dots, x_N\}$ $x_i \in \mathbb{R}^D$

\hookrightarrow D dimension

\nearrow Covariance matrix

$D \times D$

$$\text{Var}[D] = \frac{1}{N} \sum_{i=1}^N (\bar{x}_i - \mu)(\bar{x}_i - \mu)$$

Linear transformation: effect on means

e.g. $D = \{-1, 2, 3\}$

$$E[D] = 4/3$$

$$D' = \{1, 4, 5\} = D + 2$$

$$E[D'] = 10/3 = E[D] + 2$$

* Data shift by a constant

$$E[D+a] = a + E[D], a \text{ is constant}$$

$$D'' = \{-2, 4, 6\} = D \times 2$$

$$E[D''] = 8/3 = E[D] \times 2$$

* Data scaled by a factor

$$E[\alpha D] = \alpha E[D]$$

$$\Rightarrow E[\alpha D + a] = \underbrace{\alpha E[D]}_{\substack{\text{scale} \\ \text{factor}}} + \underbrace{a}_{\text{shift}}$$

Linear Transformation: effect on (co) variance

* Variance remains the same if data is shifted by a constant

$$\text{Var}[D] = \text{Var}[D+a]$$

\uparrow
offset on every element

$$\text{Var}[zD] = \lambda^2 \text{Var}[D]$$

High Dimension

$$D = \{x_1, \dots, x_k\} \quad x_i \in \mathbb{R}^n \quad \uparrow \text{n dimensions for each element}$$

Variance of this high dimension dataset is given by Covariance matrix

If $Ax_i + b$ on every data

$\begin{matrix} \text{linear} \\ \text{transform} \\ \text{matrix} \end{matrix}$

$$\text{Var}[AD + b] = A \text{Var}[D] A^T$$

\downarrow
covariance
matrix

$n \times n$
covariance matrix

Dot product

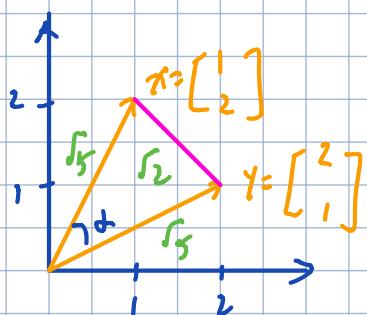
If \vec{x} and \vec{y} are 2 vectors, then the dot product is defined:

$$x^T y = \sum_{i=1}^N x_i y_i \quad x, y \in \mathbb{R}^n \quad (\text{n dimensional vectors})$$

length of \vec{x} is then defined as square root of the dot product of \vec{x} with itself:

$$\|\vec{x}\| = \sqrt{\vec{x}^T \vec{x}} = \sqrt{\sum_{i=1}^n x_i^2}$$

e.g.



$$\|\vec{x}\| = \sqrt{1^2 + 2^2} = \sqrt{5}$$

$$\|\vec{y}\| = \sqrt{2^2 + 1^2} = \sqrt{5}$$

distance between 2 vectors : length of the difference vector

$$d(x, y) = \|\vec{x} - \vec{y}\| = \sqrt{(x-y)^T (x-y)}$$

$$\text{so } \pi - y = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \Rightarrow d(x, y) = \sqrt{(-1)^2 + 1^2} = \sqrt{2}$$

$$\cos \alpha = \frac{x^T y}{\|x\| \|y\|}$$

$$\text{e.g. } \cos \alpha = \frac{(1 \ 2) \begin{pmatrix} 2 \\ 1 \end{pmatrix}}{\sqrt{5} \cdot \sqrt{5}} = \frac{4}{5}$$

$$\alpha \approx 0.64 \text{ rad}$$

Inner Product

Definition

$$\vec{x}, \vec{y} \in V \quad (\text{vector space } V)$$

symmetric, positive definite, bilinear mapping

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

Bilinear:

$$x, y, z \in V, \lambda \in \mathbb{R}$$

$$\langle \lambda x + z, y \rangle = \lambda \langle x, y \rangle + \langle z, y \rangle$$

$$\langle x, \lambda y + z \rangle = \lambda \langle x, y \rangle + \langle x, z \rangle$$

Positive definite

$$\langle x, x \rangle \geq 0 \quad \langle x, x \rangle = 0 \text{ if } \vec{x} = 0$$

Symmetry

$$\langle x, y \rangle = \langle y, x \rangle$$

e.g. in \mathbb{R}^2

$$\langle x, y \rangle = x^T I y \leftarrow \text{dot product}$$

$$\text{If we define } \langle x, y \rangle = x^T A y, A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$\rightarrow 2x_1y_1 + x_2y_1 + x_1y_2 + 2x_2y_2$$

Any A that is symmetric, positive definite matrix defines a valid inner product.

Inner product: length

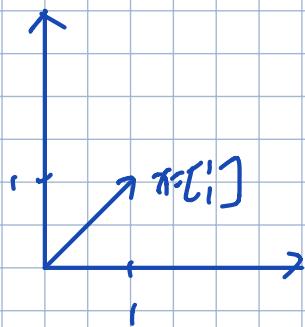
$$\|x\| = \sqrt{\langle x, x \rangle}$$

depending on the choice of inner product

\hookrightarrow geometry of vector space can be different

e.g.

$$x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



standard dot product

$$\langle x, y \rangle = x^T y \Rightarrow \|x\| = \sqrt{2}$$

$$\langle x, y \rangle = x^T \begin{bmatrix} 1 & -1/2 \\ -1/2 & 1 \end{bmatrix} y$$

$$= x_1y_1 - \frac{1}{2}(x_1y_2 + x_2y_1) + x_2y_2$$

$$\Rightarrow \|x\| = \sqrt{1 - 1 + 1} = 1$$

$$\|x\|^2 = \langle x, x \rangle = 1 \quad \text{same}$$

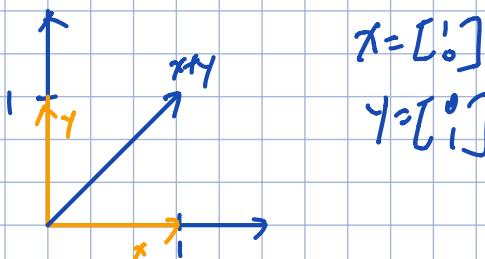
$$\Rightarrow \|x\| = 1$$

* stretching

$$\|\lambda x\| = |\lambda| \|x\|$$

* triangle inequality

$$\|x+y\| \leq \|x\| + \|y\|$$



$$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\text{Standard dot : } \|x\| = \|y\| \quad \|x+y\| = \sqrt{2}$$

$$\text{Cauchy-Schwarz inequality} \\ |\langle x, y \rangle| \leq \|x\| \|y\|$$

Inner product: distance between vectors

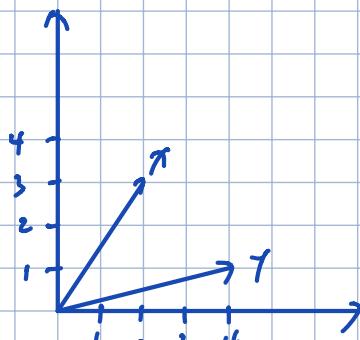
$$d(x, y) = \|x - y\| = \sqrt{\langle x - y, x - y \rangle}$$

if we use a dot product definition, then the distance is called the Euclidean distance.

$$\text{e.g. } x = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$y = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

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



define inner product

① Dot product

$$\|x - y\| = \sqrt{(-2)^2 + 2^2} = \sqrt{8}$$

$$\textcircled{2} \quad \langle x, y \rangle = x^T \begin{bmatrix} 1 & -1/2 \\ -1/2 & 1 \end{bmatrix} y$$

$$\|x - y\| = \sqrt{12}$$

Inner product: angles / orthogonality

$$\cos \theta = \frac{\langle x, y \rangle}{\|x\| \|y\|}$$

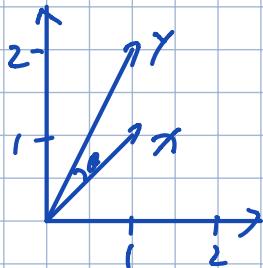
angle
between
2 vectors

$$\text{e.g. } \mathbf{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

(i) dot product as inner product

$$\cos W = \frac{\mathbf{x}^T \mathbf{y}}{\sqrt{\mathbf{x}^T \mathbf{x}} \sqrt{\mathbf{y}^T \mathbf{y}}} \\ = \frac{3}{\sqrt{10}}$$

$$W \approx 0.32 \text{ rad} \approx 18^\circ$$



$$\text{e.g. } \mathbf{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

dot product as inner product

$$\cos W = 0 \Rightarrow W = \pi/2 \text{ rad} = 90^\circ \leftarrow \text{orthogonal}$$

2 non-zero vectors, are orthogonal iff their inner product = 0

if we choose another inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{y} \Rightarrow \langle \mathbf{x}, \mathbf{y} \rangle = -1$$

$$\cos W = \frac{-1}{\sqrt{3} \sqrt{3}} = 1.9 \text{ rad} \quad \text{not orthogonal}$$

$$|\mathbf{x}| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} \Rightarrow [1 \ 1] \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} [1 \ 1] = [2 \ 1] [1 \ 1] \\ = 3$$

$$|\mathbf{y}| = \sqrt{\langle \mathbf{y}, \mathbf{y} \rangle} \Rightarrow [-1 \ 1] \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} [-1 \ 1] = [-2 \ 1] [-1 \ 1] \\ = 2 + 1 = 3 \Rightarrow \sqrt{3}$$

Inner product of functions

U, V are functions

$$\langle u, v \rangle = \int_a^b u(x)v(x) dx$$

e.g. $u(x) = \sin(x)$

$$v(x) = \cos(x)$$

$f(x) = u(x)v(x) \Rightarrow$ odd function

$$\int_{-1}^1 f(x) dx = 0 \quad \text{i.e. } u(x) \text{ and } v(x) \text{ are orthogonal}$$

Inner product of random variable

$$\Rightarrow \text{Var}[x+y] = \text{Var}[x] + \text{Var}[y]$$

$\downarrow \downarrow$
 Random
variable

if we define

$$\langle x, y \rangle \approx \text{Cov}[x, y]$$

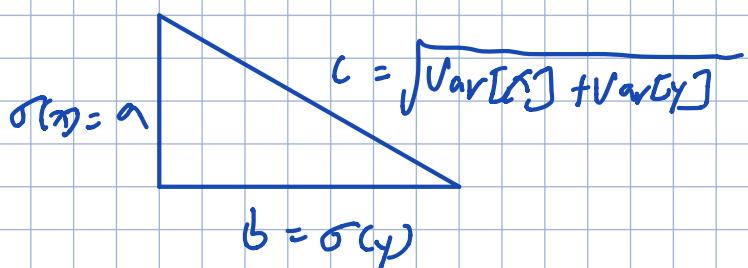
$$\text{Cov}[\alpha x + y, z] = \alpha \text{Cov}[x, z] + \text{Cov}[y, z]$$

$$\|\pi\| = \sqrt{\text{Cov}[x, x]} = \sqrt{\text{Var}[x]} = \sigma(x)$$

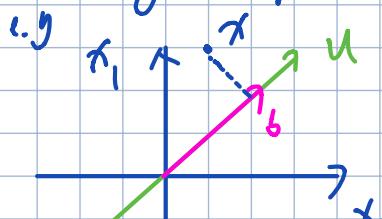
angle between 2 random variable

$$(2) \theta = \frac{\langle x, y \rangle}{\|x\| \|y\|}$$

$$= \frac{\text{Cov}[x, y]}{\sqrt{\text{Var}[x] \text{Var}[y]}}$$



orthogonal projections onto 1D subspace



$$x = \alpha r_1 + \beta r_2$$

2D space \mathbb{R}^2 basis vectors of \mathbb{R}^2 : r_1, r_2

1D Space U : basis vector b

2D space

\uparrow

1D space

$b = \lambda u$

$u \in U$

orthogonal projection of x onto U :

$\pi_U(x)$

2 important properties

(1) $\because \pi_U(x)$ is in U , \Rightarrow there exist a λ in \mathbb{R} such that

$$\pi_U(x) = \lambda b$$

$\hookrightarrow b$: basis vector that spans U

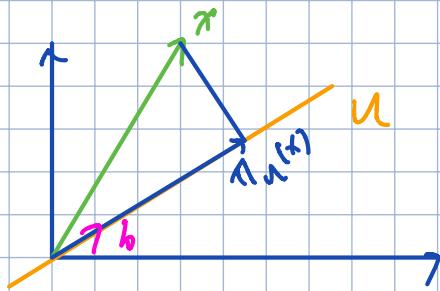
$\hookrightarrow \lambda$: coordinate of the projection w.r.t. to the basis b of the subspace U .

(2) The difference vector of x and its projection onto U is orthogonal to U .

\hookrightarrow it is orthogonal to the basis vector that spans U

$$\langle b, \pi_U(x) - x \rangle = 0$$

Generally hold for any x in \mathbb{R}^n and 1-D subspace U



e.g. we want to find $\pi_U(x)$

$$(1) \pi_U(x) = \lambda b$$

$$(2) \langle b, \pi_U(x) - x \rangle = 0$$

$$\langle b, \pi_U(x) - x \rangle = 0$$

$$\Rightarrow \langle b, \pi_U(x) \rangle - \langle b, x \rangle = 0$$

$$\Rightarrow \langle b, \lambda b \rangle - \langle b, x \rangle = 0$$

$$\Rightarrow \lambda \langle b, b \rangle - \langle b, x \rangle = 0$$

$$\uparrow \\ \|b\|^2$$

$$\Rightarrow \lambda = \frac{\langle b, x \rangle}{\|b\|^2}$$

$$\Rightarrow \Pi_u(x) = \lambda b = \frac{\langle b, x \rangle}{\|b\|^2} b$$

If we use dot product as the inner product.

$$\Pi_u(x) = \frac{\langle b^T x \rangle}{\|b\|^2} b = \underbrace{\frac{b b^T}{\|b\|^2}}_{\text{projection matrix}} x \quad \left(\frac{x^T b}{\|b\|^2} b \right)$$

↑ scalar
from any point
in 2D → 1D

← square and symmetric

$$\text{if } \|b\|=1, \Rightarrow \lambda = b^T x$$

$$\Pi_u(x) = b b^T x$$

The projection $\Pi_u(x)$ is still a vector in \mathbb{R}^D

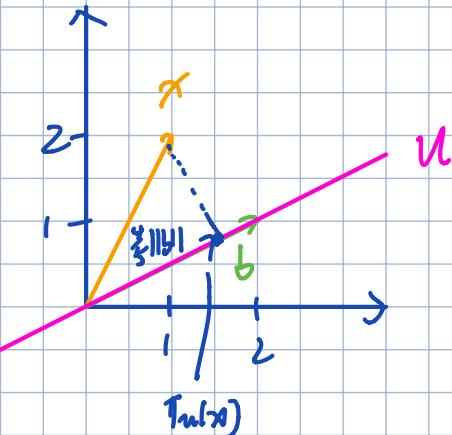
but we no longer need D coordinates to represent it, only need D (but loss some info)

$$\text{e.g. } b = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad x = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

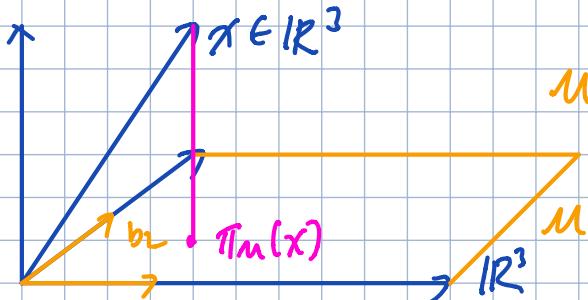
$$\Pi_u(x) = \frac{x^T b}{\|b\|} b$$

$$= \frac{2+2}{5} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$= \frac{4}{5} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$



Projection into higher-dimensional subspaces



$$M = [b_1, b_2]$$

orthogonal projection onto
subspace M

$$\text{① } \Pi_u(x) = \lambda_1 b_1 + \lambda_2 b_2$$

generalize to

$$\text{② } \begin{aligned} \langle x - \Pi_u(x), b_1 \rangle &= 0 \\ \langle x - \Pi_u(x), b_2 \rangle &= 0 \end{aligned}$$

D-dim proj to M-dim

$$\Pi_u(x) = \sum_{i=1}^M \lambda_i b_i$$

$$\langle \Pi_u(x) - x, b_i \rangle = 0 \quad i=1, \dots, M$$

we need to find

$$\lambda = \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_M \end{bmatrix} \quad M \times 1$$

$$B = [b_1 | \dots | b_M] \quad \downarrow \quad D \times M$$

$$\Pi_u(x) = B\lambda$$

(Assume we use dot product as inner product)

$$\langle \Pi_u(x) - x, b_i \rangle = \langle B\lambda - x, b_i \rangle = 0$$

$$\Rightarrow \langle B\lambda, b_i \rangle - \langle x, b_i \rangle = 0, \quad i=1 \dots M$$

$$\Rightarrow \lambda^T B^T b_i - x^T b_i = 0 \quad i=1 \dots M$$

(summarize)

$$\Rightarrow \lambda^T B^T B - x^T B = 0$$

$$\Rightarrow \lambda^T = x^T B (B^T B)^{-1} \quad (B^T B \text{ through eqn})$$

$$\Rightarrow \lambda = (B^T B)^{-1} B^T x$$

$$\hookrightarrow \Pi_u(x) = B\lambda = \underbrace{B(B^T B)^{-1} B^T}_{\text{proj matrix}} x$$

if B is orthonormal: $B^T B \approx I$

$$\hookrightarrow \Pi_u(x) = B B^T x$$

e.g. projection onto a 2D subspace

$$x = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$[\cdot, \cdot]$$

$$b_1 = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$$

$$b_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

$$M = [b_1, b_2]$$

$$\Pi_u(x) = Bx$$

$$B = [b_1 \mid b_2]$$

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

$$\lambda = (B^T B)^{-1} B^T x$$

$$B^T x = \begin{bmatrix} 4 \\ 3 \end{bmatrix}$$

$$B^T B = \begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix}$$

$$\downarrow B^T B \lambda = B^T x$$

$$\Rightarrow \lambda = \begin{bmatrix} -1 \\ 3 \end{bmatrix}$$

$$\Rightarrow \Pi_u(x) = -1b_1 + 3b_2 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$$

PCA: high level idea

$$X = \{x_1, \dots, x_N\}, \text{ where } x_i \in \mathbb{R}^D$$

data set

We want to find low level representation that is as similar as possible

$$(1) x_u = \sum_{i=1}^D \beta_i b_i$$

coefficient

orthogonal

basis of \mathbb{R}^P

assume we use the dot product

② $\hookrightarrow \beta_{in} = \mathbf{x}_n^T b_i \Rightarrow$ i.e. \mathbf{p}_{in} to be the orthogonal projection of \mathbf{x}_n onto the 1D subspace spanned by the i^{th} basis vector $\mathbf{v}_{C1,i}$

③ $B = (b_1, \dots, b_M)$ i.e. a matrix that consists of these orthonormal basis vectors

then the projection of \mathbf{x}_n onto the subspace

$$\hookrightarrow \tilde{\mathbf{x}} = B B^T \mathbf{x} \quad \text{coordinates of } \tilde{\mathbf{x}} \text{ w.r.t. the basis vectors collated in matrix } B \Rightarrow \underline{\text{code}}$$

orthogonal projection of \mathbf{x}_n onto the subspace spanned by the M basis vectors

Key idea in PCA is to find a lower dimension representation $\tilde{\mathbf{x}}_n$ of \mathbf{x}_n that can be expressed using fewer basis vectors (from $D \rightarrow M$ dimension)

if we assume the data is centered, i.e. $E[\mathbf{x}] = 0$

↑
dataset has mean ≈ 0

if assume, $b_1 \dots b_M$ are orthonormal basis of \mathbb{R}^P

In general, we can write

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M \boxed{\beta_{in} b_i} + \sum_{i=M+1}^P \beta_{in} b_i \quad \mathbb{E}[\mathbf{x}] = 0$$

1 datapoint
in \mathbf{x} dataset

principal
subspace

orthogonal
complement

ignore in PCA

b_1, \dots, b_M span the principal subspace

$\tilde{\mathbf{x}}$ is still a D -dimensional vector, it lives in an M -dimensional subspace of \mathbb{R}^P and only M coordinates b_1, \dots, b_M are necessary to represent it.

Setting: Given X data set, find β_{in} and b_i such that avg sq residual error is minimized.

$$J = \frac{1}{N} \sum_{n=1}^N \| x_n - \tilde{x}_n \|^2$$

approach:

compute partial derivatives of J w.r.t. the parameters

$$\beta_{in}, b_i \quad (\text{i.e. } \frac{\partial J}{\partial \beta_{in}} = 0, \frac{\partial J}{\partial b_i} = 0)$$

notice that these 2 parameters only enter this loss function through \tilde{x}_n . \Rightarrow we need chain rule for that

$$\frac{\partial J}{\partial \{\beta_{in}, b_i\}} = \begin{bmatrix} \frac{\partial J}{\partial \tilde{x}_n} & \frac{\partial J}{\partial \tilde{x}_n} \\ \frac{\partial \tilde{x}_n}{\partial \{\beta_{in}, b_i\}} & \frac{\partial \tilde{x}_n}{\partial \{\beta_{in}, b_i\}} \end{bmatrix}$$

$$\frac{\partial J}{\partial \tilde{x}_n} = -\frac{2}{N} (x_n - \tilde{x}_n)^T$$

Determine coordinates of the projected data

assumption: central data: $E[X] = 0$, b_1, \dots, b_M orthonormal basis

$$\frac{\partial J}{\partial \beta_{in}} = \frac{\partial J}{\partial \tilde{x}_n} \cdot \boxed{\frac{\partial \tilde{x}_n}{\partial \beta_{in}}}$$

$$\frac{\partial \tilde{x}_n}{\partial \beta_{in}} = b_i, \quad i=1, \dots, M$$

$$\frac{\partial J}{\partial \beta_{in}} = -\frac{2}{N} (x_n - \tilde{x}_n)^T b_i$$

$$= -\frac{2}{N} (x_n - \sum_{j=1}^M \beta_{jn} b_j)^T b_i$$

$$= -\frac{2}{N} (x^T | \cdot \quad \cdot \quad \cdot \quad \cdot | T | \cdot) \quad (\text{since } \beta_{jn} \text{ is } 1 \text{ or } 0)$$

M subgrad.

$$\tilde{x}_n = \sum_{j=1}^M \beta_{jn} b_j \quad (A)$$

$$J = \frac{1}{N} \sum_{n=1}^N \| x_n - \tilde{x}_n \|^2 \quad (B)$$

$$\frac{\partial J}{\partial \tilde{x}_n} = -\frac{2}{N} (x_n - \tilde{x}_n)^T \quad (C)$$

$$N(\tilde{x}_n^T b_i - \beta_{in} b_i^T b_i) \quad \text{set to 0}$$

$$= -\frac{2}{N} (\tilde{x}_n^T b_i - \beta_{in})$$

$\Rightarrow 0$

$$\hookrightarrow \text{if } \beta_{in} = \tilde{x}_n^T b_i \quad i = 1 \dots M \quad (\text{D})$$

i.e. optimal coord of \tilde{x}_n w.r.t our basis, are the orthogonal projections of the (coord. of the original data points onto i^{th} basis vector that spans the principal subspace

related

rephrase on the loss function

$$\begin{aligned} \tilde{x}_n &= \sum_{j=1}^{+1} \beta_{jn} b_j && \text{by A} \\ &= \sum_{j=1}^{M-1} (\tilde{x}_n^T b_j) b_j && \text{by D} \\ &= \sum_{j=1}^{M-1} b_j (b_j^T \tilde{x}_n) = \underbrace{\left(\sum_{j=1}^{M-1} b_j b_j^T \right)}_{\text{projection matrix}} \tilde{x}_n \end{aligned}$$

Similarly

$$\tilde{x}_n = \left(\sum_{j=1}^{M-1} b_j b_j^T \right) \tilde{x}_n \rightarrow \hat{x}_n$$

$$+ \left(\sum_{j=M+1}^D b_j b_j^T \right) \tilde{x}_n \rightarrow \text{missing}$$

i.e. the difference vector \tilde{x}_n and \hat{x}_n is

$$\tilde{x}_n - \hat{x}_n = \left(\sum_{j=M+1}^D b_j b_j^T \right) \tilde{x}_n$$

$$= \sum_{j=M+1}^D (b_j^\top \tilde{x}_n) b_j \quad (\text{E})$$

From B, we get

$$\begin{aligned} J &= \frac{1}{N} \sum_{n=1}^N \| \tilde{x}_n - \tilde{x}_n \|_2^2 \\ &= \frac{1}{N} \sum_{n=1}^N \| \sum_{j=M+1}^D (b_j^\top \tilde{x}_n) b_j \|_2^2 \quad \text{by E} \end{aligned}$$

$$\begin{aligned} \because \text{orthonormal basis} \\ b_i^\top b_j = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (b_j^\top \tilde{x}_n)^2 \end{aligned}$$

same

$$= \frac{1}{N} \sum_n \sum_j b_j^\top \tilde{x}_n \tilde{x}_n^\top b_j$$

$$\begin{aligned} \text{rearrange} \\ J &= \sum_{j=M+1}^D b_j^\top \left(\frac{1}{N} \sum_{n=1}^N \tilde{x}_n \tilde{x}_n^\top \right) b_j \\ &\quad \underbrace{\qquad\qquad\qquad}_{S} \\ &\quad \text{data covariance} \\ &\quad \text{matrix } S \end{aligned}$$

$$\begin{aligned} &= \sum_{j=M+1}^D b_j^\top S b_j \quad = \text{trace} \left(\left(\sum_{j=M+1}^D b_j b_j^\top \right) S \right) \\ &\quad (\text{F}) \quad \underbrace{\qquad\qquad\qquad}_{\text{projection matrix}} \end{aligned}$$

\hookrightarrow takes S to the orthogonal complement of the principal subspace

\Rightarrow formulate J as the variance of the data onto the subspace we ignore.

Determine the basis vectors that span the principal subspace

$$1 \leq b_i^\top S b_i : \quad S: \text{data covariance matrix}$$

$$J = \sum_{j=M+1}^n \sigma_j^2$$

minimise this (σ_j^2) need us to find the orthonormal basis that spans the space we ignore \Rightarrow after that, we take orthogonal complement as the basis of the principal subspace

e.g. for 2D space \rightarrow 1D

$$\begin{array}{c} b_1, b_2 \\ \downarrow \text{orthogonal} \\ \text{spans} \leftarrow \text{principal} \\ \text{subspace} \end{array} \quad b_i^T b_j = \delta_{ij}$$

$$J = b_2^T S b_2, \quad b_2^T b_2 = 1$$

$$L = b_2^T S b_2 + \lambda(1 - b_2^T b_2)$$

↑ Lagrangian ↑ Lagrange multiplier

$$\frac{\partial L}{\partial \lambda} = 1 - b_2^T b_2 = 0 \Leftrightarrow b_2^T b_2 = 1$$

$$\frac{\partial L}{\partial b_2} = 2b_2^T S - 2\lambda b_2^T = 0$$

$$\Leftrightarrow Sb_2 = \lambda b_2$$

eigenvalue
problem

should
smallest
eigenvalue
of S

$$\text{so } J = b_2^T S b_2 = b_2^T b_2 \lambda = \lambda \leftarrow$$

i.e. need to choose b_2

as the corresponding
eigenvector spans

the subspace we ignore

b_1 which spans the

principal subspace is the

eigenvector with largest eigenvalue of S

General case:

If we want to find the M-dimensional principal subspace of a D-dimensional data set and we solve for the basis vectors b_j

$$b_j, j = M+1, \dots, D$$

$$\left(\sum b_j = \lambda_j b_j, j = M+1, \dots, D \right)$$

loss function is the sum of the corresponding eigenvalue

$$J = \sum_{j=M+1}^D \lambda_j$$



principal subspace is spanned by the eigenvectors belong to the M largest eigenvalues of the data covariance matrix

↳ eigenvectors of S are orthogonal to each other (sym)

Eigenvector belong to the largest eigenvalue points towards the direction of data with largest variance.

↳ the variance in that direction is given by the eigenvalue.

Steps of PCA

Centering data is not necessary but easier for computation

⇒ S will be large if not centered.

also divide every dimension of the centered data by the corresponding standard deviation ⇒ data → unit free

→ variance in every dimension = 1
(no change in correlation)

- ① subtract the mean of data (centered)
- ② divide by standard deviation (for each dimension)
- ③ get $S \rightarrow$ eigen value
 \rightarrow eigen vector

(4) we can project any data point onto the principal subspace

standardize data : $\tilde{x}_x^{(d)} \leftarrow \frac{x_x^{(d)} - \mu^{(d)}}{\sigma^{(d)}}$ d : dimension

x^* a data point

projection $\tilde{x}_x = \Pi_u(x_x) = \underbrace{B B^T}_{\text{matrix}} \underbrace{x_x}_{\text{data}}$ coordinate of the projection w.r.t the basis of the principal subspace.

matrix that contains the eigenvectors that belongs to the largest eigenvalue as columns

PCA in high dimension

In D dimension, data covariance matrix is $D \times D$ matrix

If D is big, compute eigen vector/value is expensive.

If we have data set

$$\{x_1, \dots, x_N\} \in \mathbb{R}^D$$

$$S = \frac{1}{N} \tilde{x}^T \tilde{x}$$

(standardized)

$$\tilde{x} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_N^T \end{bmatrix} \quad \mathbb{R}^{N \times D}$$

if $N \ll D$

$$\text{rank}(S) = N$$

c1 - a1:

$$\sum b_i = \lambda_i b_i$$

$$\underbrace{\frac{1}{N} X^T X}_{S} b_i = \lambda_i b_i$$

$$\underbrace{\frac{1}{N} X^T X^T}_{N \times N} \underbrace{X b_i}_{C_i} = \lambda_i \underbrace{X b_i}_{C_i}$$

matrix

i.e. from $D \times D \rightarrow N \times N$

↳ from here, recover eigenvector of S

$$\underbrace{\frac{1}{N} X^T X^T}_{S} C_i = \lambda_i X^T C_i$$