

Note  $A^T A$  always come out as square and symmetry matrix even if  $A \in \mathbb{R}^{m \times n}$

## Matrix

$A \in \mathbb{R}^{m \times n}$  where  $m = \text{row}$   $n = \text{column}$

$A^k = k^{\text{th}} \text{ row on } A \rightarrow \begin{bmatrix} a_{1k} & a_{2k} & \dots & a_{nk} \end{bmatrix}$  (row vector)  $A_k = k^{\text{th}} \text{ column on } A \rightarrow \begin{bmatrix} a_{1k} \\ a_{2k} \\ \vdots \\ a_{nk} \end{bmatrix}$  (column vector) where  $k \in \{1, 2, \dots, n\}$

Transpose  $A_{ij} = A_{ji}^T$  and  $A^H = \bar{A}^T$

## Matrix Multiplication

### Vector - Matrix Multiplication

$$(Av)_i = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \sum_{j=1}^n a_{ij} v_j = \begin{bmatrix} a_1 v_1 + a_2 v_2 \\ a_3 v_1 + a_4 v_2 \end{bmatrix}$$

$A \in \mathbb{R}^{m \times n}$   
 $v \in \mathbb{R}^n$

$Av \in \mathbb{R}^m$

### Matrix - matrix multiplication

$$(AB)_{ik} = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix} \begin{bmatrix} b_1 & b_2 & b_3 \\ b_3 & b_4 & b_5 \\ b_5 & b_6 & b_6 \end{bmatrix} = \sum_{j=1}^n a_{ij} b_{jk} = \begin{bmatrix} a_1 b_1 + a_2 b_2 + a_3 b_4 + a_4 b_6 \\ a_3 b_1 + a_4 b_3 + \dots \end{bmatrix}$$

$A \in \mathbb{R}^{m \times n}$   
 $B \in \mathbb{R}^{n \times k}$

$AB \in \mathbb{R}^{m \times k}$

### Dot product (between 2 vector)

$$\langle x, y \rangle = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = x_1 y_1 + x_2 y_2 + x_3 y_3 = \sum_{i=1}^3 x_i y_i$$

can be ignore  
if it's in real space

$x \in \mathbb{R}^m$

$y \in \mathbb{R}^m$

$\langle x, y \rangle = y^T x \in \mathbb{R}^{1 \times 1}$

$$\langle x, y \rangle = y^T x = x^T y^H \quad (\text{matrix multiplication}) = \sum_{i \in I} x_i y_i$$

### Dot product property

$1. \langle v, v \rangle = 0 \text{ iff } v = 0$

$2. \langle v, w \rangle = \overline{\langle w, v \rangle}$

$3. \langle u + \lambda v, w \rangle = \langle u, w \rangle + \lambda \langle v, w \rangle$

$4. \langle w, u + \lambda v \rangle = \langle w, u \rangle + \bar{\lambda} \langle w, v \rangle$

Kronecker symbol ( $\delta$ ) =  $\begin{cases} 1, & \text{if } i=j \\ 0, & \text{if } i \neq j \end{cases}$

$$\begin{bmatrix} \delta_{11} & \delta_{12} \\ \delta_{21} & \delta_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I$$

standard unit vector ( $e^{(i)}$ ) =  $(\delta_{ij})_{j \in I}$  if  $I=3$  and  $i=2$

$$\begin{bmatrix} \delta_{11} & \delta_{12} & \delta_{13} \\ \delta_{21} & \delta_{22} & \delta_{23} \\ \delta_{31} & \delta_{32} & \delta_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \rightarrow e^{(2)} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

Range of matrix  $A \in \mathbb{K}^{I \times J}$  is the vector space spanned by its column

$\hookrightarrow C(A)$  range ( $A$ ) = span ( $A_{\cdot j}$ ,  $j \in J$ )

$\left[ \begin{array}{c} a_1 \\ \vdots \\ a_n \end{array} \right]$   
span of

## Orthogonal

Two vectors  $x, y$  are orthogonal ( $\perp$ ) if  $\langle x, y \rangle = 0$

Two set  $X, Y$  are mutually orthogonal if  $\langle x, y \rangle = 0 \quad \forall x \in X \text{ and } \forall y \in Y$

$\hookrightarrow \left[ \begin{array}{|c|c|c|} \hline & | & | & | \\ \hline x_1 & x_2 & x_3 & \dots \\ \hline \end{array} \right] \left[ \begin{array}{|c|c|c|} \hline & | & | & | \\ \hline y_1 & y_2 & y_3 & \dots \\ \hline \end{array} \right]$

Set of vector  $X$  are consider Orthonormal if

1. all vector in  $X$  is orthogonal to each other

i.e.  $\langle x_1, x_2 \rangle = 0$

2. Length of each vector is 1 (normalize)

$\|x_i\| = 1$

$X$  is consider Orthonormal basis if

1. has orthonormal property

2. linearly independent

to be L.I

$$a_1x_1 + \dots + a_nx_n = 0$$

if and only if  $a_1, \dots, a_n = 0$

to be linear dependent

$$x_i = c x_j$$

vector  $x_i$  can be represent as scalar multiple of another

All rotation matrices are orthogonal matrices

It's consider Orthogonal matrix if matrix  $A \in \mathbb{R}^{I \times J}$

1. The transpose of matrix A = the inverse of matrix A

if A is rectangle  $\rightarrow A^T A = I$ , Rectangular matrices are not invertible

if A is square  $\rightarrow$  "Unitary"  $\rightarrow A^T A = A A^T = I$  and  $A^{-1} = A^T$ , A is invertible

2.  $\det[A] = \pm 1$

Proof

since  $A A^T = I \rightarrow \det[AA^T] = \det[I] = 1$

$$\det[AA^T] = \det[A] \cdot \det[A^T] = 1$$

since  $\det[A^T] = \det[A] \rightarrow \det[A] \cdot \det[A] = 1$

$$\det[A] = \pm 1$$

3. In the orthogonal matrix, the dot product of every 2 rows and 2 columns = 0

$$A = \begin{bmatrix} \frac{1}{3} & \frac{2}{3} & -\frac{2}{3} \\ -\frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix} \quad \langle \text{Col}_1(A), \text{Col}_2(A) \rangle = \begin{bmatrix} \frac{1}{3} \\ -\frac{2}{3} \\ \frac{2}{3} \end{bmatrix} \cdot \begin{bmatrix} \frac{2}{3} \\ \frac{2}{3} \\ \frac{1}{3} \end{bmatrix} = \frac{1}{9} - \frac{4}{9} + \frac{2}{9} = 0$$

$$= \langle \text{Row}_1(A), \text{Row}_2(A) \rangle = \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \\ -\frac{2}{3} \end{bmatrix} \cdot \begin{bmatrix} -\frac{2}{3} \\ \frac{2}{3} \\ \frac{1}{3} \end{bmatrix} = -\frac{2}{9} + \frac{4}{9} - \frac{2}{9} = 0$$

4. Length of all rows and columns = 1

$$\|\text{Col}_1(A)\| = \sqrt{\left(\frac{1}{3}\right)^2 + \left(-\frac{2}{3}\right)^2 + \left(\frac{2}{3}\right)^2} = \sqrt{1} = 1$$

$$\|\text{Row}_1(A)\| = \sqrt{\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 + \left(-\frac{2}{3}\right)^2} = \sqrt{1} = 1$$

Matrix Rank No. of dimension in the column space  $\rightarrow$  # of LI columns

For  $A \in \mathbb{R}^{m \times n}$

Full Column rank =  $r = n$  Full Row rank =  $r = m$

$\hookrightarrow N(A) = \{0\}$ , 0 or 1

Full Rank =  $r = n = m$

solution exist

$\hookrightarrow$  there exist a solution  $x$  for every  $b$  for  $Ax = b$

$N(A^T) = \{0\}$

Vector Norm ( $\|v\|_1, \|v\|_2$ )  $\rightarrow$  has polynomial of degree at most 2  $\rightarrow$  magnitude (length) of the vector

A norm on  $V$  has following properties

1.  $\|v\| = 0$  if and only if  $v = 0$

2.  $\|cv\| = |c| \|v\| \quad \forall v \in V \text{ and } c \in \mathbb{K}$

3.  $\|v + w\| \leq \|v\| + \|w\| \quad \forall v, w \in V$  (triangle inequality)

$$\|x\|_2 = \sqrt{\sum_{i \in I} |x_i|^2} = \sqrt{\langle x, x \rangle}$$

Triangle Inequality  $\|x+y\| \leq \|x\| + \|y\|$

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

Pre-Hilbert Space define by

1.  $\|v\| = \sqrt{\langle v, v \rangle}$

2. inner product space (dot product)

- $\langle v, v \rangle > 0$  for  $v \neq 0$

- $\langle v, w \rangle = \langle w, v \rangle$  for  $v, w \in V$

- $\langle u + cv, w \rangle = \langle u, w \rangle + c\langle v, w \rangle$  for  $v, u, w \in V$  and  $c \in \mathbb{K}$

- $\langle w, u + cv \rangle = \langle w, u \rangle + c\langle w, v \rangle$  for  $v, u, w \in V$  and  $c \in \mathbb{K}$

3. follow Cauchy-Schwarz inequality

$$|\langle v, w \rangle| \leq \|v\| \|w\|, v, w \in V$$

If Pre-Hilbert Space is complete, then it's called Hilbert Space

What is complete?  $\rightarrow$  if Cauchy Sequence converge to an element in the space

and also continuous

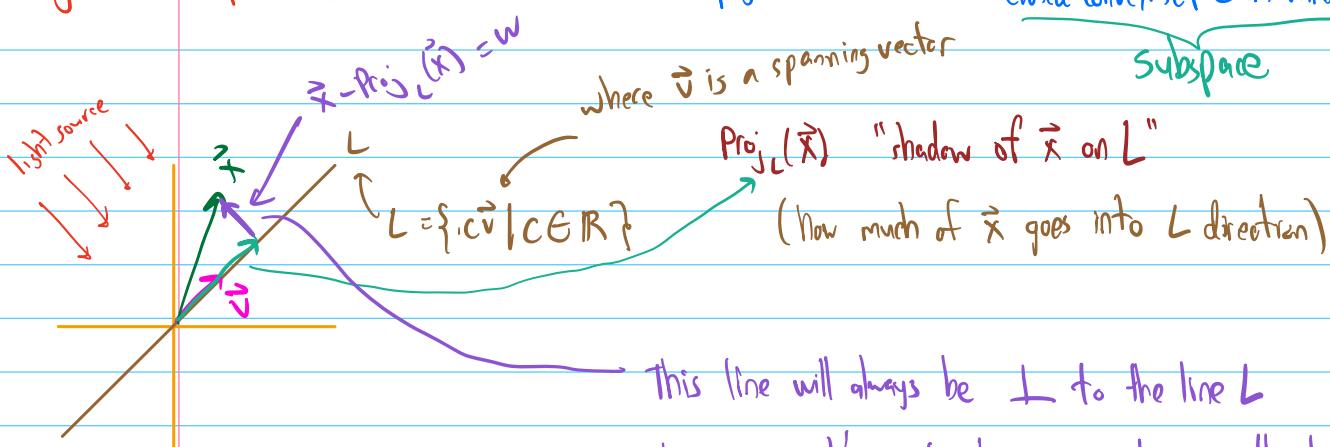
$$\forall \epsilon > 0, \exists N \in \mathbb{N} \text{ where } \forall m, n > N : |a_m - a_n| < \epsilon$$

Convexity

A set  $C$  is convex if  $\forall v, w \in C$

$$\lambda v + (1-\lambda)w \in C \text{ where } \lambda \in [0, 1]$$

## Projection Explanation



This line will always be  $\perp$  to the line  $L$

also mean it's  $\perp$  to every vector on the line  $L$

$\text{Proj}_L(\vec{x})$  = some vector in  $L$  where  $\vec{x} - \text{Proj}_L(\vec{x})$  is  $\perp$  to  $L$

blc  $\text{Proj}_L(\vec{x})$  is some vector in  $L$

$$\text{Proj}_L(\vec{x}) = cv$$

blc  $\vec{x} - \text{Proj}_L(\vec{x}) \perp L \rightarrow L$  is define by  $v$

$$\langle \vec{x} - cv, v \rangle = 0$$

$$\langle \vec{x}, v \rangle - c\langle v, v \rangle = 0$$

$$\langle \vec{x}, v \rangle = c\langle v, v \rangle$$

$$c = \frac{\langle \vec{x}, v \rangle}{\langle v, v \rangle}$$

$$\therefore \text{Proj}_L(\vec{x}) = cv = \frac{\langle \vec{x}, v \rangle}{\langle v, v \rangle} v$$

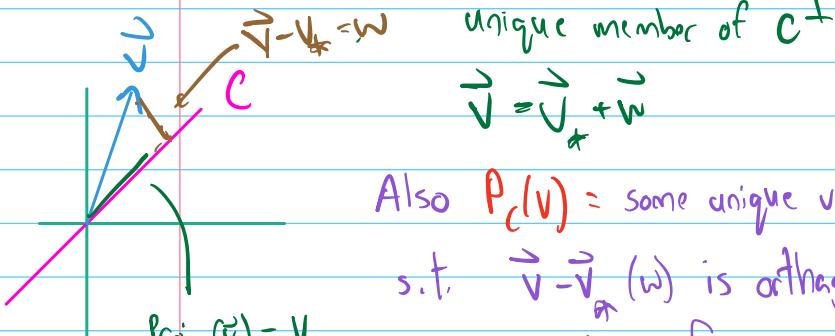
$$= \frac{\langle \vec{x}, v \rangle}{\|v\|^2} v$$

So the  $\text{Proj}_W(v) = \frac{\langle v, w \rangle}{\|w\|^2} w$  and when  $w$  is an unit vector then  $\|w\|=1$   
(orthogonal basis)

$$\therefore \text{Proj}_W(v) = \langle v, w \rangle w \quad \text{for } w \in W$$

Explaining  $P_C(v)$  is defined by  $\arg\min \|v - w\|$

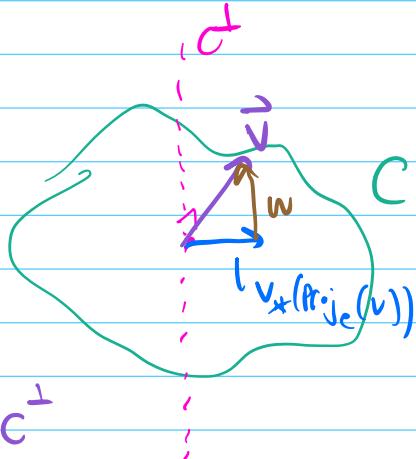
$P_C(v) = \text{unique vector } \vec{v}_* \text{ in } C \text{ such that } \vec{v} - \vec{v}_* = w \text{ where } w \perp \text{ every vector in } C \text{ where } w \text{ is a}$



Also  $P_C(v) = \text{some unique vector in } C$

s.t.  $\vec{v} - \vec{v}_*$  ( $w$ ) is orthogonal to  
every member of  $C$

$$\therefore \vec{v} - \vec{v}_* \in C^\perp \rightarrow w \in C^\perp$$



and  $P_C(v)$  is defined when  $\|v - w\|$  is the shortest (shorter than any other  $w$  in  $C$ )

which turns out to be when  $w \perp C$ .

(so  $P_C v$  is in the set  $w$  as well  $P_C v = w$ )

Projection of a vector in subspace  $C \subseteq C$ : the  $w$  that will result in the least length between  $\vec{v}$  and  $\vec{w}$

$$P_C(\vec{v}) = \arg \min_{w \in C} \|v - w\|$$

where  $\vec{w} \in C$  as well, Any other point in  $\vec{w}$  that doesn't create least length is not a projection

Writing  $\text{Proj}_W(v)$  into a Matrix  $\text{Proj}_W(v) = A\vec{v}$

Changing  $P_W(v) = \vec{v} - \vec{w}$  into Matrix formula  $P_W(v) = A(A^T A)^{-1} A^T \vec{v}$

$W$  is a subspace  $\mathbb{R}^n$  and  $b = \{b_1, b_2, \dots, b_k\}$

we know that  $P_W(v) \in W$  where  $A = \text{matrix with basis } b \text{ as a column}$

so  $P_W(v) = A\vec{y}$  for some vector  $\vec{y}$  in  $\mathbb{R}^k$

we also know that  $w$  is  $\perp$  to all member of  $W$ , and therefore  $w \in W^\perp$

matrix  $A$  is created from basis vector of  $W$

so  $C(A) = W$  and  $W^\perp = C(A)^\perp = N(A^T)$  so  $w \in N(A^T)$

and  $\vec{w} = \vec{v} - P_W(v) \in N(A^T)$

by def. of nullspace  $N(A^T) \rightarrow A^T w = 0$

$$A^T(\vec{v} - P_w(\vec{v})) = 0$$

$$A^T(\vec{v} - A\vec{y}) = 0$$

$$\vec{A}\vec{v} - A^T A \vec{y} = 0$$

$$\begin{cases} \vec{A}\vec{v} = A^T A \vec{y} \\ (A^T A)^{-1} \vec{A}\vec{v} = \vec{y} \end{cases}$$

we know  $A^T A$  is invertible cuz  
matrix  $A$  are made from basis vector  
and by def. of basis vector, they are L.I.  
(matrix  $A$  whose column are L.I. then  
 $A^T A$  is always invertible)

$$\therefore P_w(\vec{v}) = A\vec{y}$$

$$P_w(\vec{v}) = A(A^T A)^{-1} A^T \vec{v}$$

same matrix  $x$  that transform  $\vec{v}$  into subspace  $w$

$$P_w(\vec{v}) = A A^T \vec{v} \text{ if } A \text{ has orthonormal basis}$$

If Subspace  $V$  has  $\{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_k\}$  as an orthonormal basis

$$\text{then } A = \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \cdots & \vec{v}_k \end{bmatrix}$$

Where  $A$  is the matrix

where each of the column  $A^T A =$

vectors are the basis vector

for the subspace  $V$

$$\text{so } A(A^T A)^{-1} A^T \vec{x} = A(I_k)^{-1} A^T \vec{x} = A I_k A^T \vec{x}$$

Cuz  $\vec{v}_i \cdot \vec{v}_j = 0$   
aux of orthogonal

$$\begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \cdots & \vec{v}_k \end{bmatrix} = \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \cdots & \vec{v}_k \end{bmatrix} \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} = \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \cdots & \vec{v}_k \end{bmatrix} I_k$$

$$\therefore P_v(\vec{x}) = A A^T \vec{x}$$

↳ might not equal to  $I$  cus it's a dot product in

$$\begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \cdots & \vec{v}_k \end{bmatrix} \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \cdots & \vec{v}_k \end{bmatrix}$$

row vector and  
row vector might  
not be orthogonal.

## Eigenvector and Eigenvalue

Eigenvector → the vector in transformation that only scale up or down (shrink or stretch) but not change direction

Eigenvalue → Scaling factors of eigenvector

Transformation of  $\vec{v} = A\vec{v} = \lambda\vec{v}$

$\lambda$  = eigenvalue  
 $\vec{v}$  = eigenvector

$$\hookrightarrow \vec{v} = \lambda\vec{v} - A\vec{v} \rightarrow \vec{v} = I_n\vec{v}$$

$$2I_n\vec{v} - A\vec{v} = 0$$

this some matrix  
as a const  
of  $I_n$

$$(2I_n - A)\vec{v} = 0$$

this mean  $\vec{v}$  is a member of nullspace

$$(A\vec{x} = 0)$$

some matrix

$$\vec{v} \in N(2I_n - A)$$

We know that if  $N(2I_n - A)$  doesn't contain only  $\vec{v} = 0$ , then  $2I_n - A$  is not L.I.

If  $2I_n - A$  is not L.I., then it has L.D columns →  $2I_n - A$  is not invertible →  $\det = 0$

$A\vec{v} = \lambda\vec{v}$  for nonzero  $\vec{v}$ 's if and only if  $\det(2I_n - A) = 0$  or  $\det(A - \lambda I) = 0$

if  $A = \begin{bmatrix} 1 & 2 \\ 4 & 3 \end{bmatrix}$  then  $\det \begin{bmatrix} 1-\lambda & 2 \\ 4 & 3-\lambda \end{bmatrix}$  or  $\det \begin{bmatrix} 2-\lambda & -2 \\ 4 & 3-\lambda \end{bmatrix}$

$$\Delta = (1-\lambda)(3-\lambda) - (2)(4) = (3-4\lambda+\lambda^2) - 8 = \lambda^2 - 4\lambda - 5$$

$$\therefore (\lambda-5)(\lambda+1) = 0 \quad \lambda = 5, -1$$

Eigenspace - all of the eigenvectors that correspond to same eigenvalue

$$\hookrightarrow E_5 = N(A - 5I)$$

$$\text{for } \lambda = 5 \rightarrow E_5 = N \left( \begin{bmatrix} 1-5 & 2 \\ 4 & 3-5 \end{bmatrix} \right) = N \left( \begin{bmatrix} -4 & 2 \\ 4 & -2 \end{bmatrix} \right) \rightarrow N(A) = A\vec{v} = 0$$

$$A\vec{v} = 0 \cdot \begin{bmatrix} -4 & 2 \\ 4 & -2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$-4v_1 + 2v_2 = 0$$

$$4v_1 - 2v_2 = 0$$

$$4v_1 = 2v_2 \\ v_1 = \frac{1}{2}v_2$$

$$4v_1 = 2v_2 \\ v_1 = \frac{1}{2}v_2$$

$$\text{let } v_2 = t$$

$$E_5 = \left\{ \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2}t \\ t \end{bmatrix}, t \in \mathbb{R} \right\}$$

$$E_5 = \text{span} \left( \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix} \right) \rightarrow \text{do this for}$$

$$\lambda = -1 \text{ as we } \|$$

↳ this mean that if there's vector

that sit along this line, the transformation A will stretch the vector by 4

Linear Operator  $\rightarrow$  Matrix that transform vector from one vector space to another  
(Linear transformation)

Operator Norm  $\rightarrow$  a measure of the largest value (change) by which the matrix A stretch the vector z

$$\|A\| = \|A\|_{X \rightarrow Y} = \sup_{z \neq 0} \frac{\|Az\|_Y}{\|z\|_X}, \quad A \in \mathbb{K}^{I \times J}$$

$\uparrow$   
change space from  $X \rightarrow Y$

Spectral Norm  $\rightarrow$  if the operator norm is done by the linear operator that transform vector from space X to T and X and T is coincide with Euclidean Norm  $\rightarrow$  then it's called Spectral Norm

$$\|A\| = \max_{\|v\|_2=1} \frac{\|Av\|_2}{\|v\|_2} = \|Av_1\|_2 = \|\sigma_1 v_1\|_2 = |\sigma_1| \|v_1\|_2$$

$v_1$  = best fit line with corresponding leading principle component      maximum value  $\sigma_1$

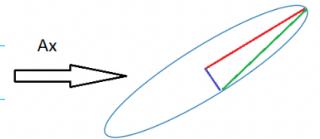
$v_1$  is the largest blow up vector ( $v_1$  = leading singular vector)

$\rightarrow$  largest singular value (or  $\lambda$  if A is symmetric)

Frobenius Norm (Norm over a matrix)  $\rightarrow$  "diagonal" between all the singular value

$$\|A\|_F = \sqrt{\sum_{i \in I} \sum_{j \in J} |A_{ij}|^2} = \sqrt{\langle A, A \rangle_F} \quad A \in \mathbb{K}^{I \times J}$$

$$\langle A, B \rangle_F = \sum_{i \in I} \sum_{j \in J} A_{ij} B_{ij} = \text{tr}(AB^H) = \text{tr}(B^H A)$$



$$\|A\|_F^2 = \langle A, A \rangle_F = \text{tr}(A^H A) = \text{tr}(AA^H)$$

$$\|A\|_F = \sqrt{\text{tr}(AA^H)} = \sqrt{\text{tr}(A)^2} = \sqrt{\sigma_1^2 + \sigma_2^2 + \dots} \quad \text{or} \quad \sqrt{\sum_{k=1}^r \sigma_k(A)^2}$$

$$\sum_{i \in I} \|A^i\|_2^2 = \sum_{i \in I} \sum_{k=1}^r |\langle A^i, v_k \rangle|^2 = \sum_{k=1}^r \sum_{i \in I} |\langle A^i, v_k \rangle|^2 = \sum_{k=1}^r \|Av_k\|_2^2 = \sum_{k=1}^r \sigma_k^2 = \|A\|_F^2$$

sum of Projection onto each component of basis vector

Note: if  $v_1, \dots, v_r$  is all row of A

## Identifying largest principle component

- finding the line that best fit the data
- Best line  $v_1 \rightarrow$  Max variance between data

$v_1 =$  best fit line = leading singular vector with  $\sigma_1$

$$v_k = \max_{\substack{\|v\|_2=1 \\ v \perp v_1, \dots, v_{k-1}}} \|Av\|_2 \quad k > 1 \quad v_k \text{ is singular vectors and a maximizers}$$

$$\sigma_k = \|Av_k\|_2 \quad k = 1 \dots n \quad Av_i = \sigma_i u_i$$

Prove for general case  $v_k$  is the best fit subspace with  $\|Av_k\| \geq \|Aw_k\|$

Assume:  $V_{k-1} = \text{span}(v_1, \dots, v_{k-1})$  is the best fit subspace

Prove:  $V_k = \text{span}(v_1, \dots, v_k)$  is the best fit subspace

let  $W_k = \text{span}(w_1, \dots, w_k)$   $W \neq V_k$   $w_1, \dots, w_k$  is orthonormal basis and  $w_k \perp \underline{v_{k-1}}$

since  $w_k \perp v_{k-1}$  and we know that  $v_k$  is the maximum  $\|Av\|_2^2$  for any  $v \perp v_{k-1}$

$$\therefore \|Av_k\| \geq \|Aw_k\|$$

$$\therefore \|Aw_1\|_2^2 + \dots + \|Aw_{k-1}\|_2^2 + \|Aw_k\|_2^2 \leq \|Av_1\|_2^2 + \dots + \|Av_{k-1}\|_2^2 + \|Av_k\|_2^2$$

## Prove SVD

for  $v_1, \dots, v_r$  singular vector

$A$  is make up of  $v_1, \dots, v_r$

let some vector  $\vec{v} \in \mathbb{R}^n$

$$\text{so } \vec{v} = \text{ther comb of } v_1, \dots, v_r = \sum_{i=1}^r \langle v, v_i \rangle v_i + \sum_{i=r+1}^n \langle v, v_i \rangle v_i$$

base from nullspace of  $A$  to make dimension  $n$

$$\begin{aligned}
\text{then } A \mathbf{v} &= A \left( \sum_{i=1}^r \langle \mathbf{v}, \mathbf{v}_i \rangle \mathbf{v}_i + \sum_{i=r+1}^n \langle \mathbf{v}, \mathbf{v}_i \rangle \mathbf{v}_i^\perp \right) \\
&= \sum_{i=1}^r \mathbf{v}_i^\top \mathbf{v} (A \mathbf{v}_i) + \sum_{i=r+1}^n \langle \mathbf{v}, \mathbf{v}_i \rangle^\perp (A \mathbf{v}_i^\perp) \rightarrow \text{b/c } \mathbf{v}_i^\perp \text{ is from } N(A) \rightarrow A \mathbf{v}_i^\perp = 0 \\
&= \sum_{i=1}^r \mathbf{v}_i^\top \mathbf{v} (A \mathbf{v}_i) \Rightarrow A \mathbf{v}_i = \sigma_i \mathbf{u}_i \\
&= \sum_{i=1}^r \mathbf{v}_i^\top \mathbf{v} (\sigma_i \mathbf{u}_i) = \sum_{i=1}^r (\sigma_i \mathbf{u}_i \mathbf{v}_i^\top) \mathbf{v} \\
\therefore A &= \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top
\end{aligned}$$

$$\begin{aligned}
A \mathbf{v} &= A_1 \mathbf{v}_1 + A_2 \mathbf{v}_2 + \dots + A_m \mathbf{v}_m \\
&= A_1 y_1 + A_2 y_2 + \dots + A_m y_m \\
&\quad \vdots \\
&= A_m y_m
\end{aligned}$$

$\xrightarrow{\text{CC(A)}}$   
 $\xrightarrow{\text{Proof: } A = C(C^\top)}$

## SVD

$$\boxed{A}_{m \times n} = \boxed{U}_{m \times r} \boxed{\Sigma}_{r \times r} \boxed{V^H}_{r \times n}$$

$$\begin{aligned}
A &= U \Sigma V^H = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \\
A \mathbf{v}_i &= \sigma_i \mathbf{u}_i
\end{aligned}$$

Proof

$$A \mathbf{v} = U \Sigma$$

$$C(A) = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_r)$$

$$A \mathbf{x} = A \mathbf{v} = U \Sigma = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_r) = C(U) = C(A)$$

$$C(A^T) = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_r)$$

$$A^T \mathbf{y} = A^T \mathbf{u} = \Sigma \mathbf{v} = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_r) = C(V) = C(A^T)$$

$$N(A) = \text{span}(\mathbf{v}_{r+1}, \dots, \mathbf{v}_n)$$

$$N(A^T) = \text{span}(\mathbf{u}_{r+1}, \dots, \mathbf{u}_m)$$

## Best k-approximation

If  $V_K$  is span by  $\mathbf{v}_1, \dots, \mathbf{v}_K$  of A

then  $V_K$  is the best fit k-Dimensional subspace of A

$$V_K = \arg \min_{\dim V \leq K} \sum_{i=1}^r \|A^i - P_V(A^i)\|_2^2 = \arg \max_{\dim V \leq K} \sum_{i=1}^r \|P_{V^H}(A^i)\|_2^2 = \sum_{i=1}^K \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = A_K$$

$$A_K = \sum_{i=1}^K \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \quad K \leq r \quad A_K \text{ is the best approx. of } A \text{ of any matrix of rank } K$$

$A_K$  are orthogonal projection of rows of A ( $A^i$ ) onto subspace

$$A_K = \sum_{i=1}^K P_{V^H}(A^i)$$

Prove

$$\begin{aligned}
&\sum_{i=1}^K \sum_{l=1}^r P_{V^H}(A^i) = \sum_{i=1}^K \sum_{l=1}^r \langle A^i, \mathbf{v}_l \rangle \mathbf{v}_l \\
&= \sum_{l=1}^r \sum_{i=1}^K \langle A^i \mathbf{v}_l \rangle \mathbf{v}_l^H = \sum_{l=1}^r A \mathbf{v}_l \mathbf{v}_l^H = \sum_{l=1}^r \sigma_l \mathbf{u}_l \mathbf{v}_l^H = A_K
\end{aligned}$$

$$\langle A^i, \mathbf{v}_l \rangle = (A \mathbf{v}_l)_i$$

Prove

$$\|A - A_K\|_F \leq \|A - B\|_F \quad \text{also} \quad \|A - A_K\| \leq \|A - B\|$$

let  $B$  is of rank at most  $K$  and row of  $B(B^H)$  span subspace  $V$ ,  $\dim(V) \leq K$

Assume  $B$  minimize  $\|A - B\|_F$

In order to minimize  $B^H = P_V(A^H)$

$$\text{so } \|A - B\|_F^2 = \sum_{i=K+1}^r \|A^H - P_V(A^H)\|_2^2$$

and Since  $\dim(V) = K$ ,  $V_K$  is the best fit subspace when  $V_K = \min_{\text{in } K\text{-dimension}} \sum_{i=K+1}^r \|A^H - P_V(A^H)\|_2^2$

so in order to get best fit subspace at

$K$ -dimension  $V_K$ ,  $B = P_V(A^H)$  which mean  $B = A_K$  as well

$$\|A - A_K\|^2 = \sigma_{K+1}^2$$

$$A - A_K = \sum_{i=K+1}^r \sigma_i u_i v_i^H$$

$$\|A - A_K\|_F^2$$

$$= \sqrt{\sum_{i=1}^r \sigma_i^2 (A - A_K)^2}$$

$$= \sqrt{\sum_{i=1}^r \sigma_i^2}$$

$$A - A_K = \sum_{j=K+1}^r \sigma_j u_j v_j^H$$

$$\|(A - A_K)v\|_2 = \left\| \sum_{j=K+1}^r \sigma_j u_j v_j^H (v_j v_j^H) \right\|_2 = \left\| \sum_{j=K+1}^r \sigma_j \sigma_j u_j v_j^H v_j \right\|_2$$

$$\text{let } v = \text{top singular vector of } A - A_K$$

$$\text{and } v = \sum_{j=K+1}^r \sigma_j v_j$$

$$\|v\|_2 = 1$$

$$= \left\| \sum_{j=K+1}^r \sigma_j \sigma_j u_j \right\|_2 = \|u_j\|_2 \underbrace{\left\| \sum_{j=K+1}^r \sigma_j \sigma_j \right\|_2}_2$$

$$= \sqrt{\left( \sum_{j=K+1}^r \sigma_j \sigma_j \right)^2} \quad \|v\|_2 = 1 = \sqrt{\left( \sum_{j=1}^r \sigma_j \sigma_j \right)^2}$$

maximize  $\|v\|_2 = \|A - A_K\|_2$   
we can do by letting  
 $\sigma_j = 1$  at  $j=1$  or  $K+1$   
and 0 everywhere else

$$= \sqrt{(1)_{j=K+1} \sum_{j=K+1}^r \sigma_j^2} = \sqrt{\sigma_{K+1}^2} = \sigma_{K+1}$$

0 everywhere else

bc  $K+1$  is where  
 $\sigma$  is the biggest  
for  $A - A_K$

Power of A

if  $A$  is square and symmetric  $\rightarrow$  has same left and right S vector

$$A^2 = A^H A = A A^H = (V \Sigma V^H)(V \Sigma V^H) = V \Sigma^2 V^H = \sum \sigma_i^2 V V^H$$

$$\therefore A^K = V \Sigma^K V^H = \sum_{i=1}^r \sigma_i^k v_i v_i^H$$

if  $A$  is not square and symmetric

$$\text{let } B = A^H A \quad B = (V \Sigma V^H)(V \Sigma V^H) = V \Sigma^2 V^H = \sum_{i=1}^r \sigma_i^2 u_i u_i^H \quad \therefore B = A^H A = \sum \sigma_i^2 u_i u_i^H$$

$$\text{or } B = A^H A \quad B = (V \Sigma U^H)(U \Sigma U^H) = U \Sigma^2 U^H = \sum_{i=1}^m \sigma_i^2 v_i v_i^H \quad \therefore B = A^H A = \sum_{i=1}^m \sigma_i^2 v_i v_i^H$$

$\lim_{k \rightarrow \infty} A^k = \lim_{k \rightarrow \infty} \sigma^k V_i V_i^H$  note that  $\sigma_1 \gg \sigma_2$  so if  $\sigma_1 = 1$  then all  $\sigma < 1$  for sure  
 if  $\sigma > 1$   $k \rightarrow \infty \sigma^k \rightarrow \infty$

if  $\sigma = 1$   $k \rightarrow \infty \sigma^k = \sigma$ , constant (stable state)  $\rightarrow$  most desired

if  $\sigma \neq 0$  but  $\sigma < 1$   $k \rightarrow \infty \sigma^k \rightarrow 0$  (stability)

$$\|A^m\|_F = \sqrt{\sigma_1^{2m} + \sigma_2^{2m} + \dots} \approx \sigma_1^{2m} \quad \text{and} \quad \lim_{m \rightarrow \infty} \frac{A^m}{\|A\|_F} = V_i V_i^H$$

if  $\sigma_1 \gg \sigma_2 \gg \dots \gg \sigma_r$

$$B^m = \sum \sigma_i^{2m} V_i V_i^H \quad \text{Normalize } B^m \text{ by } \frac{B^m}{\sigma_1^{2m}} = \left\{ \left( \frac{\sigma_i}{\sigma_1} \right)^{2m} V_i V_i^H \right\} \quad \text{and} \quad \frac{\sigma_1}{\sigma_1} = 1, \frac{\sigma_2}{\sigma_1} < 1, \dots, \frac{\sigma_r}{\sigma_1} < 1$$

which mean that  $\frac{B^m}{\sigma_1^{2m}} \rightarrow V_i V_i^H$  or  $B^m \rightarrow \sigma_1^{2m} V_i V_i^H$  when  $m \rightarrow \infty$  so when  $m \rightarrow \infty \frac{\sigma_2}{\sigma_1}, \dots, \frac{\sigma_r}{\sigma_1} \rightarrow 0$

let  $x$  be unit vector with significant component on  $v_i$  direction

$$B^m x = \sum \sigma_i^{2m} V_i V_i^H x$$

$$\frac{B^m x}{\sigma_1^{2m}} = \underbrace{\left( \frac{\sigma_i}{\sigma_1} \right)^{2m} \langle v_i, x \rangle}_{\text{Projection of } x \text{ onto } v_i} v_i \quad \text{when } m \rightarrow \infty \quad \frac{B^m x}{\sigma_1^{2m}} \rightarrow \langle v_i, x \rangle v_i \quad \therefore B^m x \rightarrow \sigma_1^{2m} \langle v_i, x \rangle v_i \quad \text{when } m \rightarrow \infty$$

$$\therefore B^m x \approx \sigma_1^{2m} \underbrace{\langle x, v_i \rangle}_{\text{projection onto } v_i} v_i \quad \text{as } m \rightarrow \infty$$

$$\text{From this we can find } v_i \text{ by normalizing } v_i = \frac{\sigma_1^{2m} \langle x, v_i \rangle v_i}{\| \sigma_1^{2m} \langle x, v_i \rangle v_i \|}$$

$$A^T A v_i = \sigma_i^2 v_i$$

$\sigma^2$  is the magnitude of  $A^T A v_i$ , which is the same as the norm of  $A^T A v_i$

$$\sigma^2 = \|A^T A v_i\|$$

$$\sigma = \sqrt{\sigma^2}$$

$$\text{then we can find } \sigma_i \text{ from } \sigma_i = \|A v_i\| \quad \text{or} \quad \sigma_i = \sqrt{\|A^T A v_i\|}$$

$$\text{then we can find } v_i \text{ from } v_i = \frac{A v_i}{\sigma_i}$$

## Randomizing $\vec{x}$ in Power Method

$\vec{x} = |\langle x, u_i \rangle| \geq \alpha > 0$  with probability at least  $1 - C\alpha\sqrt{d}$   
 $(P(|x| \geq \alpha))$

Let  $V = \text{space span}$  left singular vector for singular value  $\sigma_k > (1 - \varepsilon)\sigma_1$   
 of

$$w^* = \frac{(AA^H)^m x}{\|(AA^H)^m x\|_2} \quad \text{where } m = \Omega\left(\frac{\ln(\frac{d}{\varepsilon})}{\varepsilon}\right)$$

$$w^T = P_V(w^*) + P_{V^\perp}(w^*)$$

$P_{V^\perp}(w^*)$  has  $O\left(\frac{\varepsilon}{\alpha d}\right)$  with probability at least  $1 - C\alpha\sqrt{d}$

Prove that  $P_{V^\perp}(w^*) = O\left(\frac{\varepsilon}{\alpha d}\right)$

$$A = \sum_{k=1}^r \sigma_k u_k v_k^H \quad \text{let } x = \sum_{k=1}^n \langle x, u_k \rangle u_k$$

Find lower bound of  $\|(AA^H)^m x\|_2$

$$1) \text{ let } \sigma_k = 0 \text{ for } k > r \text{ so } A = \sum_{k=1}^r \sigma_k u_k v_k^H$$

$$2) (AA^H)_x = \sum_{k=1}^n (\sigma_k)^{2m} u_k u_k^H x = \sum_{k=1}^n (\sigma_k)^{2m} \langle x, u_k \rangle u_k$$

$$3) \|(AA^H)^m x\|_2^2 = \sum_{k=1}^n (\sigma_k)^{4m} |\langle x, u_k \rangle|^2 \|u_k\|_2^2 = \sum_{k=1}^n (\sigma_k)^{4m} |\langle x, u_k \rangle|^2$$

4) We know that  $\sum_{k=1}^n \sigma_k > \sigma_1$  so

$$\|(AA^H)^m x\|_2^2 = \sum_{k=1}^n (\sigma_k)^{4m} |\langle x, u_k \rangle|^2 \geq \sigma_1^{4m} |\langle x, u_1 \rangle|^2$$

5) We know that  $|\langle x, u_1 \rangle| \geq \alpha$

$$\|(AA^H)^m x\|_2^2 \geq \sigma_1^{4m} |\langle x, u_1 \rangle|^2 \geq \sigma_1^{4m} \alpha^2$$

this is the lower bound, it need to be at least this

Find upper bound of  $(AA^H)^m x$  (or  $P_v + w^*$  in my case)

$$6. \quad \| (AA^H)^m x \|_2^2 \leq \sum_{k=1}^n (\sigma_k)^{4m} |\langle x, u_k \rangle|^2 \|u_k\|_2^2 = \sum_{k=1}^n (\sigma_k)^{4m} |\langle x, u_k \rangle|^2$$

6.5 Choose  $r_2$  such that  $\sigma_1, \dots, \sigma_{r_2}$  are the singular values of  $A$  that is greater than  $(1-\varepsilon)\sigma_1$ , which mean that  $\sigma_{r_2+1} \dots \sigma_n < (1-\varepsilon)\sigma_1$

$P_v + w^*$  = component of  $(AA^H)^m x$  orthogonal to the space  $V = \text{span}(u_1, \dots, u_{r_2})$  which is the subspace  $V^\perp = \text{span}(u_{r_2+1}, \dots, u_n)$

$$\therefore \|P_{V^\perp} w^*\|_2^2 = \sum_{k=r_2+1}^n (\sigma_k)^{4m} |\langle x, u_k \rangle|^2$$

7. We know that  $\sigma_k < (1-\varepsilon)\sigma_1$  for  $k \geq r_2+1$

$$\text{so } \|P_{V^\perp} w^*\|_2^2 = \sum_{k=r_2+1}^n \sigma_k^{4m} |\langle x, u_k \rangle|^2 \leq (1-\varepsilon)\sigma_1^{4m} \sum_{k=r_2+1}^n |\langle x, u_k \rangle|^2$$

8. We know that  $\|x\|_2^2 = 1$  and  $\|x\|_2^2 = \sum_{k=1}^n |\langle x, u_k \rangle|^2$ , so  $1 = \sum_{k=1}^n |\langle x, u_k \rangle|^2 \geq \sum_{k=r_2+1}^n |\langle x, u_k \rangle|^2$

$$\text{so } \|P_{V^\perp} w^*\|_2^2 \leq (1-\varepsilon)\sigma_1^{4m} \sum_{k=r_2+1}^n |\langle x, u_k \rangle|^2 \leq (1-\varepsilon)\sigma_1^{4m}$$

- This is the upperbound, it cannot be more than this

9. Now we can normalize  $\|P_{V^\perp} w^*\|_2^2$

$$\text{we know that } w^* = \frac{(AA^H)^m x}{\|(AA^H)^m x\|_2}$$

$$\text{so unit vector of } \|P_{V^\perp} w^*\|_2^2 = \frac{\|P_{V^\perp} w^*\|_2^2}{\|(AA^H)^m x\|_2^2} \quad \text{where } \|P_{V^\perp} w^*\|_2^2 \leq (1-\varepsilon)^{4m} \sigma_1^{4m}$$

$$\text{and } \|(AA^H)^m x\|_2^2 \geq \sigma_1^{4m} \alpha^2$$

$$\|P_{V^\perp} w^*\|_2^2 \leq \frac{(1-\varepsilon)^{4m}}{\sigma_1^{4m} \alpha^2} = \frac{(1-\varepsilon)^{4m}}{\alpha^2} \quad \text{and } \|P_{V^\perp} w^*\| = \frac{(1-\varepsilon)^{2m}}{\alpha}$$

10. Now we know that the linear approximation of  $(1-\varepsilon)$  is  $e^{-\varepsilon}$  \*

$$\text{so } (1-\varepsilon)^{2m} = e^{-2m\varepsilon} \quad \text{so } \|P_{V^\perp} w^*\| = O(\alpha^{-1} e^{-2m\varepsilon})$$

11. We know that  $m = \Omega\left(\frac{\ln(\frac{d}{\epsilon})}{\epsilon}\right)$  which mean  $m \geq c \frac{\ln(\frac{d}{\epsilon})}{\epsilon}$  where  $c$  is some constant

we let  $c = \frac{1}{2}$  to get rid of 2 in  $\omega^{-2m\epsilon}$

$$\therefore \|P_{W^\perp} w^*\| = O\left(\frac{\omega^{-2\epsilon \left(\frac{\ln(\frac{d}{\epsilon})}{2\epsilon}\right)}}{\alpha}\right) = O\left(\frac{\epsilon}{\alpha d}\right)$$

### Pseudo inverse matrices

if  $A \in \mathbb{R}^{m \times n}$ ,  $A$  does not have  $n$  real inverse

$$A^+ = V \Sigma^+ U^H$$

if  $m > n$  (tall matrix)  $\begin{bmatrix} \vdash \\ \vdash \\ \vdash \end{bmatrix}$  and  $n=r$  ( $n$  independent col)

$A^+$  is left inverse  $A^+ = (A^H A)^{-1} A^H$

$A^T A = I$   $A A^T$  is the projection onto  $A$

$N(A) = \{0\}$ , have one or zero solution to  $Ax=y$   $\xrightarrow{Ax=y} x=A^T y$

Least Square of  $A = \min \|Ax-y\|_2$  where there's only 1 unique solution  $x=A^T y$

if  $n > m$  (wide matrix)  $\begin{bmatrix} \vdash \\ \vdash \\ \vdash \end{bmatrix}$  and  $r=m$  ( $m$  independent row)

$A^+$  is a right inverse

$$A^+ = A^H (A A^H)^{-1} \quad A A^+ = A A^H (A A^H)^{-1} = I$$

$A^T A$  = projection onto  $A$

$N(A^T) = \{0\}$  but  $N(A) \neq \{0\}$  so there's many solution to  $Ax=y$

Least Square of  $A = \min \|Ax-y\|_2$  where there can be many solutions  $x=A^T y$  to  $\min \|Ax-y\|_2$

but we pick the solution with  $\min \|x\|_2$

## Stability of SVD (What's the error from approximation)

### Weyl's bound

$$\text{A perturbed } (\tilde{\mathbf{A}}) = \mathbf{A} + \mathbf{E}$$

↳ (approx.)

$$\lambda_1(\tilde{\mathbf{A}}) - \max_{\mathbf{v}} \mathbf{v}^T(\mathbf{A} + \mathbf{E})\mathbf{v} \leq \underbrace{\max_{\mathbf{v}} (\mathbf{v}^T \mathbf{A} \mathbf{v})}_{\lambda_1(\mathbf{A})} + \underbrace{\max_{\mathbf{v}} (\mathbf{v}^T \mathbf{E} \mathbf{v})}_{\lambda_1(\mathbf{E})}$$

$$\therefore \lambda_1(\tilde{\mathbf{A}}) \leq \lambda_1(\mathbf{A}) + \lambda_1(\mathbf{E})$$

spectral decomposition (eigenvalue decap.)

if  $\mathbf{A}$  is square and symmetric

$$\mathbf{A} = \mathbf{S} \Lambda \mathbf{S}^{-1}$$

where  $\mathbf{S}$  = eigenvector matrix  
 $\Lambda$  = eigenvalue diagonal matrix

$$\lambda_k(\mathbf{A}) + \lambda_n(\mathbf{E}) \leq \lambda_k(\mathbf{A} + \mathbf{E}) \leq \lambda_k(\mathbf{A}) + \lambda_1(\mathbf{E}) \leq \lambda_k(\mathbf{A}) + \lambda_1(\mathbf{E})$$

$$\lambda_n(\mathbf{E}) \leq \lambda_k(\mathbf{A} + \mathbf{E}) - \lambda_k(\mathbf{A}) \leq \lambda_1(\mathbf{E})$$

we know that  $-\sigma_1 \leq \lambda_n \leq \lambda_1 \leq \sigma_1$

$$\text{so } -\sigma_1(\mathbf{E}) \leq \lambda_n(\mathbf{E}) \leq \sigma_k(\mathbf{A} + \mathbf{E}) - \sigma_k(\mathbf{A}) \leq \lambda_1(\mathbf{E}) \leq \sigma_1(\mathbf{E})$$

$$-\sigma_1(\mathbf{E}) \leq \sigma_k(\mathbf{A} + \mathbf{E}) - \sigma_k(\mathbf{A}) \leq \sigma_1(\mathbf{E})$$

$$\therefore |\sigma_k(\mathbf{A} + \mathbf{E}) - \sigma_k(\mathbf{A})| \leq \|\mathbf{E}\| = |\sigma_1(\mathbf{E})|$$

### Weyl's bound

$$\|\mathbf{E}\| \geq |\sigma_k(\mathbf{A} + \mathbf{E}) - \sigma_k(\mathbf{A})|$$

### Mirsky's bound

$$\|\mathbf{E}\|_F \geq \sqrt{\sum_{k=1}^n |\sigma_k(\mathbf{A} + \mathbf{E}) - \sigma_k(\mathbf{A})|^2}$$

$$\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{w}^T \mathbf{v} = \|\mathbf{w}\|_2 \|\mathbf{v}\|_2 \cos \theta(\mathbf{w}, \mathbf{v})$$

$$\text{since } \|\mathbf{w}\|_2 = \|\mathbf{v}\|_2 = 1$$

$$\mathbf{w}^T \mathbf{v} = \cos \theta(\mathbf{w}, \mathbf{v}) \rightarrow \cos \theta(\mathbf{w}, \mathbf{v}) = \sum_{i=1}^n w_i v_i = \text{sum of projection of all dimension}$$

$$\begin{aligned} \|\mathbf{P}_w \mathbf{P}_v\|_F^2 &= \|\mathbf{w}^T \mathbf{v} \mathbf{w}^T\|_F^2 = \text{tr}((\mathbf{w}^T \mathbf{v} \mathbf{w}^T)^T (\mathbf{w}^T \mathbf{v} \mathbf{w}^T))_F \\ &= \text{tr}(\mathbf{v}^T \mathbf{w}^T \mathbf{w}^T) = 2 \text{tr}(\mathbf{v}^T \mathbf{w}^T \mathbf{w}^T) + \text{tr}(\mathbf{w}^T \mathbf{w}^T) \\ &= \|\mathbf{v}\|_F^2 + 2 \text{tr}(\mathbf{v}^T \mathbf{w}^T \mathbf{w}^T) + \|\mathbf{w}^T \mathbf{w}^T\|_F \\ &= 2n - 2 \|\mathbf{w}^T \mathbf{w}^T\|_F = 2n - 2 \cos^2 \theta(\mathbf{w}, \mathbf{v})_F \\ &= 2n - 2 \cos^2 \theta(\mathbf{w}, \mathbf{v}) \\ &\stackrel{\cos \theta(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \|\mathbf{v}\|_2 \|\mathbf{w}\|_2}{=} n + \frac{n}{2} \sin^2 \theta(\mathbf{w}, \mathbf{v}) \\ &\stackrel{\sin^2 \theta(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \|\mathbf{v} - \mathbf{w}\|_2^2}{=} \frac{n}{2} \|\mathbf{v} - \mathbf{w}\|_2^2 \\ &\|\mathbf{P}_w \mathbf{P}_v\|_F = \sqrt{\frac{n}{2} \|\mathbf{v} - \mathbf{w}\|_2^2} \\ &\|\mathbf{P}_w \mathbf{P}_v\|_F = \sqrt{\frac{n}{2} \|\mathbf{v} - \mathbf{w}\|_2^2} \end{aligned}$$

## Stability of Singular Space (Wedin's bound)

Non stable for SVD mean:

→ the matrices that are close in norm can have

completely difference singular vector (Hence, the stability of singular space)

## Stable SVD

= matrices that are close in norm are close in singular vector (singular space) as well

From Exercise 7 we show that

If A has a non stable SVD

$\tilde{A} = A + E$  adding perturbation E onto A will

1. norm of  $\tilde{A}$  is changed based on how big  $\Sigma$  in E is.  $\rightarrow$  dependent of  $\Sigma$

$\rightarrow$  we can lower the perturbation  $\Sigma$  to get norm of  $\tilde{A}$  as close as possible to norm of A

2. Singular Vector of A is transform independent of  $\Sigma$  in perturbation

$\rightarrow$  so Singular Vector of A is forever change, and we cannot reduce the effect of change even with the smallest  $\Sigma$

If A has stable SVD

both norm and singular vector is dependent on  $\Sigma$

$$\tilde{A} = A + E \quad E = \tilde{A} - A$$

$$\tilde{A} = U \tilde{\Sigma} V^H \rightarrow \tilde{A} \tilde{V}_1 = \tilde{A} \tilde{V}_1 = \tilde{U}_1 \tilde{\Sigma}_1$$

$$\begin{aligned} R_{11} &= A \tilde{V}_1 - \tilde{A} \tilde{V}_1 = A \tilde{V}_1 - \tilde{\Sigma}_1 \tilde{U}_1 \\ &= \tilde{V}_1 (A - \tilde{A}) \\ &= \tilde{V}_1 (A - (A + E)) \end{aligned}$$

$$\begin{aligned} R_{12} &= A^H \tilde{U}_1 - \tilde{A}^H \tilde{U}_1 = A^H \tilde{U}_1 - \tilde{\Sigma}_1^H \tilde{V}_1 \\ &= \tilde{U}_1 (A^H - \tilde{A}^H) \\ &= \tilde{U}_1 (A^H - (A^H + E^H)) \end{aligned}$$

$$R_{11} = -E \tilde{V}_1$$

$$R_{12} = -E^H \tilde{U}_1$$

$$\|R_{11}\| = \| -E \tilde{V}_1 \| = \|E\| = \sigma_1(E)$$

$$\text{Normal distribution}$$

$$\phi(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

## Basic Probability

$$P(|x| \geq t) = P(x \geq t) + P(x \leq -t)$$

$$P(|x| < t) = P(-t < x < t)$$

$$E(X) = \int_{-\infty}^{\infty} x \cdot \phi(t) dt \quad P(a < x < b) = \int_a^b \phi(x) dx$$

$$P(w \in B) = \frac{\text{Vol}(B)}{\text{Vol}(\Omega)}$$

the prob. that  $x$  hit a point  $w \in B$  is the ratio of area of  $B$  to entire area ( $\Omega$ )

$$\text{Union bound} = P(A \cup B) \leq P(A) + P(B)$$

$$\phi(x+y)(t) = \int_{-\infty}^{\infty} \phi(x(u)) \phi(y(t-u)) du \rightarrow \text{if } x, y \sim N(0, I) \quad x+y \sim N(0, 2) \quad \text{sd} = \sqrt{2}$$

Gaussian vector - random vector  $X = (X_1, \dots, X_n)$  with all its entries  $X_i$ :  $X_i \sim N(\mu_i)$

Gaussian distribution  $\rightarrow$  rotation invariant  $\rightarrow$  if  $X \sim N(0, I)$  and  $u$  is random variable

$\triangleright$  not care about direction,  $Xu \sim N(0, \|u\|_2^2)$  if  $u$  is fom or is orthogonal matrix

Sum of Gaussian variable

then  $Xu \sim N(0, I)$

if  $X = (X_1, \dots, X_n)$ ,  $X_i$  is independent and  $X_i \sim N(\mu_i, \sigma_i^2)$

then  $Z = \sum_{i=1}^n X_i$ ,  $Z \sim N(\sum \mu_i, \sum \sigma_i^2)$

$$\text{Variance} = \text{Var}(X) = E[(X - \mu)^2] = E[(X - E(X))^2] = E(X^2) - E(X)^2$$

Moments of  $X = E(X^p)$

Cavalieri's formula

$$\text{Absolute Moment of } X = E(|X|^p) = p \int_0^{\infty} P(|x| \geq t) t^{p-1} dt \quad (E|X+Y|^p)^{\frac{1}{p}} \leq (E|X|^p)^{\frac{1}{p}} + (E|Y|^p)^{\frac{1}{p}}$$

$$E|XY| \leq (E|X|^{\beta})^{\frac{1}{\beta}} (E|Y|^{\alpha})^{\frac{1}{\alpha}} \quad \text{if } \beta, \alpha \geq 1 \quad \text{and } \frac{1}{p} + \frac{1}{q} = 1$$

$$E(|XY|) \leq E(|X|^p)^{\frac{1}{p}} E(|Y|^q)^{\frac{1}{q}} \quad \text{if } \frac{1}{p} + \frac{1}{q} = 1$$

Markov Inequality (estimate tail w/ upperbound)

$$P(|X| \geq t) \leq \frac{E(X)}{t}$$

$$P(|X| \geq t) = P(|X|^p \geq t^p) \leq \frac{E(|X|^p)}{t^p}$$

$$P(|X| \geq t) = P(\exp(\theta X) \geq \exp(\theta t)) \leq \frac{E(\exp(\theta X))}{\exp(\theta t)}$$

Moment generating function of  $X$

Cumulative generating function  $\rightarrow C_x(t) = \ln E(\exp(\theta X))$

Cramer's theorem (estimate tail w/ upperbound) Proof

if  $X = (X_1, \dots, X_n)$  and  $X$  is independent

$$P\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(\inf_{\theta > 0}\left\{-\theta t + \sum_{i=1}^n C_{X_i}(\theta)\right\}\right)$$

$$P\left(\sum_{i=1}^n X_i \geq t\right) = P\left(\exp\left(\theta \sum_{i=1}^n X_i\right) \geq \exp(\theta t)\right) \stackrel{\downarrow}{\leq} \frac{E\left(\exp\left(\theta \sum_{i=1}^n X_i\right)\right)}{\exp(\theta t)}$$

$$\begin{aligned} &= e^{-\theta t} E\left(\prod_{i=1}^n \exp(\theta X_i)\right) = e^{-\theta t} \prod_{i=1}^n E(\exp(\theta X_i)) \\ &= e^{-\theta t} \prod_{i=1}^n \exp\left(\ln(E(\exp(\theta X_i)))\right) = e^{-\theta t} \prod_{i=1}^n \exp(C_{X_i}(\theta)) = e^{-\theta t} e^{\sum_{i=1}^n C_{X_i}(\theta)} \\ &= \exp(-\theta t + \sum_{i=1}^n C_{X_i}(\theta)) = \exp\left(\inf_{\theta > 0}\left\{-\theta t + \sum_{i=1}^n C_{X_i}(\theta)\right\}\right) \end{aligned}$$

Hoeffding's inequality (estimate tail if  $X_i$  is bounded)

$$|X_i| \leq B_i \quad i \in [N]$$

$$P\left(\left|\sum_{i=1}^n X_i\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2 \sum_{i=1}^n B_i^2}\right)$$

Taylor Series of exponential

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

Proof  $X_i \in [-B_i, B_i]$  so  $x_i = t(-B_i) + (1-t)B_i$

$$\text{Since } f(x) = \exp(\theta x) \quad t = \frac{B_i - x_i}{2B_i}$$

$$f(x_i) = f\left(t(-B_i) + (1-t)B_i\right) \leq f(-B_i) + (1-t)f(B_i)$$

$$\exp(\theta x_i) \leq \frac{B_i - x_i}{2B_i} e^{-\theta B_i} + \frac{B_i - x_i}{2B_i} e^{\theta B_i}$$

$$E(\exp(\theta x_i)) \leq E\left(\frac{B_i - x_i}{2B_i} e^{-\theta B_i} + \frac{B_i - x_i}{2B_i} e^{\theta B_i}\right)$$

$$= \frac{-\theta B_i}{2B_i} e^{-\theta B_i} - \frac{\theta B_i}{2B_i} e^{\theta B_i} + \frac{\theta B_i}{2B_i} e^{\theta B_i} - \frac{E(\theta B_i)}{2B_i}$$

$$= \frac{1}{2} e^{-\theta B_i} + \frac{1}{2} e^{\theta B_i} \rightarrow e^{\frac{x}{2}} = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

$$+ \frac{1}{2} \left( \sum_{k=0}^{\infty} \frac{(-\theta B_i)^k}{k!} + \sum_{k=0}^{\infty} \frac{(\theta B_i)^k}{k!} \right)$$

$$= \frac{1}{2} \left( \sum_{k=0}^{\infty} \frac{(-\theta B_i)^{2k}}{(2k)!} + \sum_{k=0}^{\infty} \frac{(\theta B_i)^{2k}}{(2k)!} \right) = \frac{1}{2} \frac{(\theta B_i)^{2k}}{k!} = e^{\frac{(\theta B_i)^2}{2}}$$

$$(x_i(\theta)) = \ln(E(\exp(\theta x_i))) = \ln(e^{\frac{(\theta B_i)^2}{2}})$$

From Cramer's rule

$$P\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(\inf_{\theta > 0}\left(-\theta t + \sum_{i=1}^n C_{X_i}(\theta)\right)\right) \leq \exp\left(\inf_{\theta > 0}\left(-\theta t + \sum_{i=1}^n \frac{\theta^2 B_i^2}{2}\right)\right)$$

$$\text{Optimal } \theta = \frac{t}{\sum_{i=1}^n B_i^2}$$

$$P\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{t^2}{2 \sum_{i=1}^n B_i^2}\right)$$

$P\left(\sum_{i=1}^n X_i \leq -t\right)$  also give the same bound

$$\therefore P\left(\left|\sum_{i=1}^n X_i\right| \geq t\right) = P\left(\sum_{i=1}^n X_i \geq t\right) + P\left(\sum_{i=1}^n X_i \leq -t\right) = 2 \exp\left(-\frac{t^2}{2 \sum_{i=1}^n B_i^2}\right)$$

### 1.1 Sum of <sup>\*</sup>bound indep. random variable

$$\text{then } P\left(\left|\sum_{i=1}^n X_i\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2(\sigma^2 + R^2)}\right)$$

$$\text{where } \sigma^2 = \sum_{i=1}^n \sigma_i^2$$

### 1.2 if $|X_i| < K$ $X_i$ is bounded, then

$$P\left(\left|\sum_{i=1}^n X_i\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2(\sigma^2 + \frac{1}{3}K^2)}\right) \quad (R = \frac{1}{3}K)$$

Proof

$$\text{if } E|X_i|^n \leq \frac{n!}{2} R^{n-2} \sigma_i^2 \quad |X_i| < K$$

$$\text{for } n=2, \quad E|X_i|^2 \leq \frac{1}{2} R^0 \sigma_i^2 = \sigma_i^2$$

$$\text{for } n \geq 3, \quad E|X_i|^n = E(|X_i|^{n-2} |X_i|^2) \leq E(K^{n-2} |X_i|^2) = K^{n-2} E|X_i|^2 \leq K^{n-2} \sigma_i^2$$

$$\text{so } E|X_i|^n \leq K^{n-2} \sigma_i^2 = \frac{n!}{n!} K^{n-2} \sigma_i^2 \leq \frac{n!}{2 \cdot 3^{n-2}} \cdot K^{n-2} \sigma_i^2 = \frac{n!}{2} \left(\frac{K}{3}\right)^{n-2} \sigma_i^2 \quad \therefore R = \frac{K}{3} \text{ and } \sigma_i = \sigma_i$$

$$\downarrow n! \geq 2 \cdot 3^{n-2}$$

## 2) Subexponential random variable

Note

$$\int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} dt = \sqrt{2\pi}$$

$$\int_{-\infty}^{\infty} e^{-\frac{|t|}{2}} dt = \sqrt{\pi}$$

random variable  $X$  is subexponential if

$$P(|X| > t) \leq \beta e^{-kt} \quad \beta, k > 0$$

then  $P\left(\left|\sum_{l=1}^M X_l\right| \geq t\right) \leq 2 \exp\left(-\frac{(kt)^2}{2(2\beta M + kt)}\right)$

Proof  $P\left(\left|\sum_{l=1}^M X_l\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2(\sigma^2 + kt)}\right)$

$$E(\exp(\theta X_2)) = 1 + \theta E(X) + \sum_{n=2}^{\infty} \frac{\theta^n E(X^n)}{n!} = 1 + \sum_{n=2}^{\infty} \frac{\theta^n E(X^n)}{n!} = 1 + \frac{\theta^2 \sigma^2}{2} \leq 2 \theta \frac{\sum_{n=1}^{n-2} E(X^n)}{n! \sigma^2}$$

$$\rightarrow \text{let } F_2(\theta) = \sum_{n=2}^{\infty} \frac{2\theta^{n-2} E(X^n)}{n! \sigma^2}, \text{ and } F(\theta) = \max_{l \in [M]} F_l(\theta) \quad \text{so } F_2(\theta) \leq F(\theta)$$

$$E(\exp(\theta X_2)) = 1 + \frac{\theta^2 \sigma^2}{2} F_2(\theta) \leq 1 + \frac{\theta^2 \sigma^2}{2} F(\theta) \leq \exp\left(\frac{\theta^2 \sigma^2}{2} F(\theta)\right)$$

$$\rightarrow C_X \theta = \ln E(\exp(\theta X_2)) \leq \ln \exp\left(\frac{\theta^2 \sigma^2 F(\theta)}{2}\right) = \frac{\theta^2 \sigma^2 F(\theta)}{2}$$

$$\sum_{l=1}^M \frac{\theta^2 \sigma_l^2}{2} = \theta^2 \sigma^2$$

$$\text{So, using Cramer's theorem } \rightarrow P\left(\sum_{l=1}^M X_l \geq t\right) \leq \exp\left(\inf_{\theta > 0} \left\{-\theta t + \frac{\theta^2 \sigma^2 F(\theta)}{2}\right\}\right) = \exp\left(\inf_{\theta > 0} \left\{-\theta t + \frac{\theta^2 \sigma^2}{2} F(\theta)\right\}\right)$$

We know that  $E(X_l^n) \leq E(|X_l|^n)$

$$\text{so } F_2(\theta) \leq F(\theta) \leq \sum_{n=2}^{\infty} \frac{2\theta^{n-2} E(|X_l|^n)}{n! \sigma^2} \quad \text{and from Bernstein's inequality } E(|X_l|^n) \leq \frac{n!}{2} R^{n-2} \sigma_l^n$$

$$F(\theta) \leq \sum_{n=2}^{\infty} \frac{2\theta^{n-2}}{n! \sigma^2} \left(\frac{n!}{2} R^{n-2} \sigma_l^n\right) = \sum_{n=2}^{\infty} (\theta R)^{n-2} = \frac{1}{1-\theta R} \quad (\text{Taylor's Series } X^n = \sum_{n=1}^{\infty} \frac{1}{n!} X^{n-1})$$

$$\text{So } P\left(\sum_{l=1}^M X_l \geq t\right) \leq \exp\left(\inf_{\theta > 0} \left\{-\theta t + \frac{\theta^2 \sigma^2}{2} \frac{1}{1-\theta R}\right\}\right) \leq \exp\left(\inf_{0 < \theta < 1} \left\{-\theta t + \frac{\theta^2 \sigma^2}{2(1-\theta R)}\right\}\right)$$

letting  $\theta = \frac{t}{\sigma^2 + kt}$  will result in  $P\left(\sum_{l=1}^M X_l \geq t\right) \leq \exp\left(\frac{t^2}{2(\sigma^2 + kt)}\right)$

changing  $X_2$  to  $-X_2$  will yield the same result, so applying union bound complete the proof

$$\therefore P\left(\left|\sum_{l=1}^M X_l\right| \geq t\right) \leq 2 \exp\left(\frac{t^2}{2(\sigma^2 + kt)}\right)$$

## Dimensionality Reduction : The Johnson - Lindenstrauss Lemma

### 1. Johnson - Lindenstrauss Lemma (JL) by random projection

For any  $0 < \epsilon < 1$ , any integer  $n$ , let  $k$  be positive integer such that

$$k \geq 2\beta \left( \frac{\epsilon^2}{2} - \frac{\epsilon^3}{3} \right)^{-1} \ln n$$

$\rightarrow 2[\epsilon^2 - \epsilon^3]$  for gaussian (to prove easier ??)

then for set  $P$  of  $n$  points in  $\mathbb{R}^d$ , there's a map  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  s.t for all point  $v, w \in P$

$$(1-\epsilon) \|v-w\|_2^2 \leq \|f(v)-f(w)\|_2^2 \leq (1+\epsilon) \|v-w\|_2^2$$

The map  $f$ , can be generate at random with prob =  $1 - (n^{2\beta} - n^{1-\beta})$  (choosing  $\beta$  large will have higher prob)

Considering unit vector  $z = \frac{v-w}{\|v-w\|_2}$ ,  $A$  = orthogonal projection matrix onto dim  $k$  in  $\mathbb{R}^d$

$$(1-\epsilon) \leq \|Az\|_2^2 \leq (1+\epsilon)$$
  
 $| \|Az\|_2^2 - 1 | \leq \epsilon$

### 2. Johnson - Lindenstrauss for Gaussian random matrices

let matrix  $A$  have i.i.d. entries and  $N(\mu, 1)$   $A \in \mathbb{R}^{k \times d}$  and  $x \in \mathbb{R}^d$

then  $P\left(\left|\left\|\frac{1}{\sqrt{k}}Ax\right\|_2^2 - \|x\|_2^2\right| > \epsilon \|x\|_2^2\right) \leq 2e^{-(\epsilon^2 - \epsilon^3)\frac{k}{4}}$

or

$$P\left((1-\epsilon) \|x\|_2^2 \leq \left\|\frac{1}{\sqrt{k}}Ax\right\|_2^2 \leq (1+\epsilon) \|x\|_2^2\right) \geq 1 - 2e^{-(\epsilon^2 - \epsilon^3)\frac{k}{4}}$$

# Convex

Convex Set  $K \subset \mathbb{R}^N$ ,  $K$  is a convex set if for all  $x, y \in K$

$$tx + (1-t)y \in K \text{ for all } t \in [0, 1]$$

Convex Combination  $x = t_1x_1 + t_2x_2 + \dots + t_kx_k$  with  $t_1 + \dots + t_k = 1$ ,  $t_i \geq 0$

Convex Hull - boundary of the sets of point

- Convex Hull  $\text{conv}(T)$  - set of all convex combinations of point in  $T$   
(smallest convex set containing  $T$ )

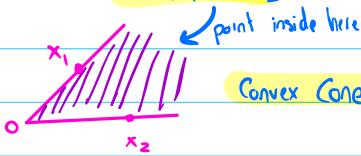
$$\text{conv}(T) = \left\{ \sum_j t_j x_j : t_j \geq 0, \sum_j t_j = 1, x_j \in T \right\}$$



Convex Cone

Conic combination - any point of the form

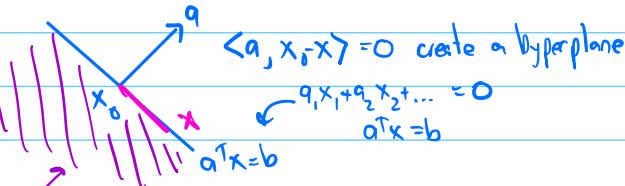
$$x = t_1x_1 + t_2x_2 \quad t_1, t_2 \geq 0$$



Convex Cone - set that contain all conic combination of points in the set

Hyperplane  $\{x | a^T x = b\}$

↳ convex  $a \neq 0$



$\langle a, x_i \rangle = 0$  create a hyperplane

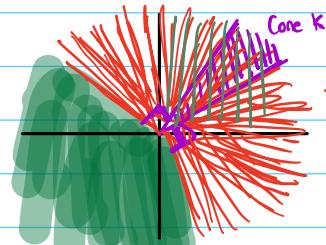
$$a_1x_1 + a_2x_2 + \dots = 0 \\ a^T x = b$$

Halfspace  $\{x | a^T x \leq b\}$

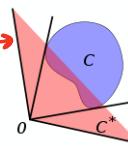
↳ convex or  $a^T x \geq b$  (depend)

Dual Cones

For cone  $K$ , dual cone  $K^* = \{z \in \mathbb{R}^N : \langle x, z \rangle \geq 0 \text{ for all } x \in K\}$

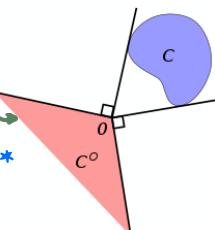


- Dual cone  $K^*$  ( $90^\circ$  of cone)



Polar Cone

$$K^0 = \{z \in \mathbb{R}^N : \langle x, z \rangle \leq 0 \text{ for all } x \in K\} = -K^*$$



## Conic Hull

cone( $T$ ) of set  $T \subset \mathbb{R}^n$  is the smallest convex cone containing  $T$

$$\text{cone}(T) = \left\{ \sum t_i x_i : x_i \geq 0, x_i \in T \right\}$$

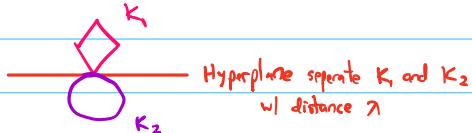
## Geometric Hahn-Banach Theorem

non-intersecting convex sets can be separate by hyperplanes.

Let  $K_1, K_2 \subset \mathbb{R}^n$  be convex sets,  $w \in \mathbb{R}^n$

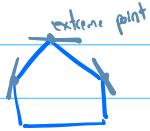
$$K_1 \subset \{x \in \mathbb{R}^n : \langle x, w \rangle \leq \lambda\}$$

$$K_2 \subset \{x \in \mathbb{R}^n : \langle x, w \rangle \geq \lambda\}$$



## Extreme Points

Point on convex that you cannot build a convex combination (convex line is not inside the graph)  
(e.g. vertex or corner point)



if  $K$  does not lie between 2 distinct points of  $K$

$$x = tw + (1-t)z, w, z \in K, t \in (0, 1) \text{ and } x = w = z$$

## Proper Functions

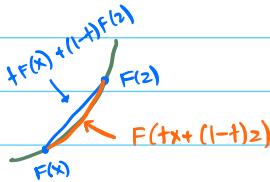
Domain of function - set of all possible inputs for the function

$$\text{dom}(F) = \{x \in \mathbb{R}^n : F(x) \neq \infty\}$$

$\text{dom}(F) \neq \emptyset$  is called proper.  
(input is not empty)

## Convex Function $t \in [0, 1]$

$$F(tx + (1-t)z) \leq tF(x) + (1-t)F(z)$$



Strictly convex  $x \neq z$  and  $t \in (0, 1)$

$$F(tx + (1-t)z) < tF(x) + (1-t)F(z)$$

$\rightarrow$  is strictly convex

$\square$  is not  
 $F(tx + (1-t)z) = tF(x) + (1-t)F(z)$

Strongly convex  $\gamma > 0$   $t \in [0, 1]$

$$F(tx + (1-t)z) \leq tF(x) + (1-t)F(z) - \frac{\gamma}{2}(1-t)\|x-z\|^2$$

Strongly convex are always strictly convex

Strictly convex are always convex

Note

Tangent line formula  
 $y - y_i = m(x - x_i)$      $m = \text{slope at tangent point}$   
or  
 $y = y_i + \nabla f(x_i)(x - x_i)$      $x_i, y_i = \text{tangent point}$

Domain of convex function is convex

Function  $F$  is convex if and only if its epigraph is a convex set

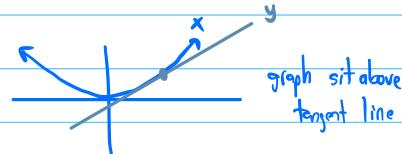
Epigraph - Area (all point) above the graph

$$\text{epi}(F) = \{(x, r) : r \geq F(x)\} \subset \mathbb{R}^n \times \mathbb{R}$$

### Smooth Convex Function

if  $F$  is differentiable

$$1. F \text{ is convex iff } F(x) \geq F(y) + \langle \nabla F(y), x - y \rangle$$



$$2. F \text{ is strongly convex w/ parameter } \gamma > 0 \text{ iff}$$

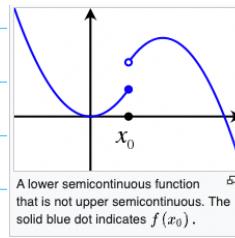
$$F(x) \geq F(y) + \langle \nabla F(y), x - y \rangle + \frac{\gamma}{2} \|x - y\|^2$$

$$3. \text{ If } F \text{ is twice differentiable then } F \text{ is convex iff } \nabla^2 F(x) \geq 0$$

### Lower Semicontinuous

$$\lim_{x \rightarrow x_0} F(x) \geq F(x_0)$$

sequence of  $\{x_k\}_{k \rightarrow x_0}$  converges



A lower semicontinuous function that is not upper semicontinuous. The solid blue dot indicates  $f(x_0)$ .

$f$  is LSC iff  $\text{epi}(f)$  is closed set

Every Norm is convex

$\ell_p$ -norms are strictly convex (except  $\ell_1$  and  $\ell_\infty$  - only convex)

$A \in \mathbb{R}^{N \times N}$   $F(x) = x^T A x$  is positive semidefinite  $\Rightarrow$  convex

if positive definite  $\Rightarrow$  strictly convex

### Jointly convex function

$$f(tx_1 + (1-t)x_2, ty_1 + (1-t)y_2) \leq t f(x_1, y_1) + (1-t) f(x_2, y_2)$$

If  $g(x) = \inf_{y \in \mathbb{R}^m} f(x, y)$  and  $f(x, y)$  is a jointly convex function,

then  $g(x)$  is a convex function

### Maxima of convex function

$F$  attain its maxima at an extreme point of  $K$

( $K$  is a convex hull (compact convex set))

$\Rightarrow$  closed and bounded

closed  $\Rightarrow$  contain all boundary point in the set  
bounded  $\Rightarrow$  no points going to infinity

## Convex Conjugate ( $F^*$ )

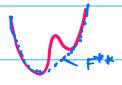
$$F^*(y) := \sup_{x \in \mathbb{R}^n} \{ \langle x, y \rangle - F(x) \} \rightarrow \langle x, z \rangle - F(z) \leq \langle x, y \rangle - F(x) \quad \forall z \in \mathbb{R}^n$$

only equal when  $z = x$

Fenchel's inequality  $\rightarrow \langle x, y \rangle \leq F(x) + F^*(y)$  for all  $x, y \in \mathbb{R}^n$

$F^*$  is always a convex function

## Property of Convex Conjugate

1.  $F^*$  is lower semicontinuous  $\rightarrow \lim_{x \rightarrow x_0} F^*(x) \geq F^*(x_0)$  (the actual point is lower than point that  $x$  converge to)
  2.  $F^{**}$  is the largest lower semicontinuous convex function satisfying  $F^{**}(x) \leq F(x)$  for all  $x \in \mathbb{R}^n$   $\rightarrow$ 
    - if  $F$  is convex and lower semicontinuous,  $F = F^{**}$  fit under the graph
    - convex hull of epigraph of  $F$
  3. For  $r \neq 0$  let  $F_r(x) = F(rx)$  then  $(F_r)^*(y) = F^*(\frac{y}{r})$
  4. For  $T \geq 0$ ,  $(TF)^*(y) = TF^*(\frac{y}{T})$
  5. For  $z \in \mathbb{R}^n$  let  $F^{(2)} = F(x-z)$ . then  $(F^{(2)})^*(y) = \langle z, y \rangle + F^*(y)$
- $\text{epi}(F^{**}) = \text{conv}(\text{epi } F)$
- 

## Subdifferential

- There are multiple derivatives at a non-smooth point

e.g.  $F(x) = |x|$  can have slope at any value from -1 to 1 at  $x=0$

• So the generalise gradient

$$\partial F(x) = \left\{ v \in \mathbb{R}^n : F(z) - F(x) \geq \langle v, z-x \rangle \text{ for all } z \in \mathbb{R}^n \right\}$$

$\uparrow$   
subgradient of  $F$  at  $x$

$\partial F(x)$  of convex function is always non empty

$v \in \partial F(x)$

if  $F$  is differentiable (smooth point) in  $x$ , then

$$\partial F(x) = F'(x) = \langle v, z-x \rangle$$

$$\text{and } \partial F(x) = \nabla F(x)$$

## Subdifferentials and Conjugation

$F$  is convex function

All this conditions are equivalent

$$1. y \in \partial F(x)$$

$$2. F(x) + F^*(y) = \langle x, y \rangle$$

if  $F$  is lower semicontinuous, then also equivalent to

$$3. x \in \partial F^*(y)$$

## Proximal mapping

function  $x \mapsto F(x) + \frac{1}{2} \|x-z\|_2^2$  is strictly convex

because  $x \mapsto \|x\|_2^2$  is strictly convex

B/c it's strictly convex, minimizer is unique

$$P_F(z) := \arg \min_{x \in \mathbb{R}^n} \left\{ F(x) + \frac{1}{2} \|x-z\|_2^2 \right\}$$

↑  
Proximal mapping associate with  $F$

and  $x = P_F(z)$  iff  $z \in x + \partial F(x)$

## Moreau's identity

$$P_F(z) = (I + \partial F)^{-1}(z)$$

if  $F$  is a lower semicontinuous convex function

$$P_F(z) + P_{F^*}(z) = z$$

if  $P_F(z)$  is easy to compute, then  $P_{F^*}(z) = z - P_F(z)$

$$P_{\gamma F}(z) + \gamma P_{\gamma F}(\frac{z}{\gamma}) = z$$

## Lagrange function (Try to solve minimization problem w/o constrain)

Normal Convex problem has lot of constrain, so we form Lagrange function with the convex problem + constrain instead to have problem with no constrain

$$\text{if convex problem (primal problem)} = \min_{x \in \mathbb{R}^n} F_0(x) \quad \text{subject to } Ax = b \rightarrow Ax = b \\ \text{and } F_i(x) \leq b_i \rightarrow F_i(x) - b_i \leq 0 \quad i \in [M]$$

the Lagrange function  $L(x, \varepsilon, v) = F_0(x) + \sum_{i=1}^M v_i (F_i(x) - b_i)$  without constrain  
↓  
Lagrange multiplier

if there's no inequality constrain

$$\text{then } L(x, \varepsilon) = F_0(x) + \varepsilon^*(Ax - b)$$

whmizer

Lagrange Dual (Provide a lower bound on the optimal value of  $F_0(x^*)$  of the minimization problem)

Lagrange Dual Function ( $H$ ) (Find lower bound of minimization using Lagrange)

$$H(\varepsilon, v) = \inf_{x \in \mathbb{R}^n} L(x, \varepsilon, v)$$

Lagrange Dual Function is affine function

Infimum of family of affine function is always concave.

$$H(\xi, v) \leq F(x^*) \quad \text{for all } \xi \in \mathbb{R}^m, v \geq 0$$

if primal problem =  $\min_{x \in \mathbb{R}^n} F_0(x)$  subject to  $\underbrace{Ax=b}$  and  $\underbrace{F_\ell(x) \leq b_\ell}_{\ell \in [M]}$   
 $F_\ell(x) - b_\ell \leq 0$

then

$$L(x, \xi, v) = F_0(x) + \xi^*(Ax-b) + \sum_{\ell=1}^M v_\ell (\underbrace{F_\ell(x)-b_\ell}_{\xi_\ell})$$

for all  $\xi \in \mathbb{R}^m$  and  $v \geq 0$  and  $x$  is feasible

$$\xi^*(Ax-b) = 0 \quad \text{and} \quad v_\ell (\underbrace{F_\ell(x)-b_\ell}_{\xi_\ell}) \leq 0$$

$$\text{So } L(x, \xi, v) = F_0(x) + \xi^*(Ax-b) + \sum_{\ell=1}^M v_\ell (F_\ell(x)-b_\ell) \leq F_0(x)$$

$$L(x, \xi, v) \leq F_0(x)$$

$$\inf_{x \in \mathbb{R}^n} L(x, \xi, v) \leq \inf_{x \in \mathbb{R}^n} F_0(x)$$

$$\underline{H(\xi, v)} \leq \underline{F(x^*)}$$

Lagrange Dual Problem (finding the best lower bound for the optimal value)

$$\max H(\xi, v) \text{ subject to } v \geq 0$$

A feasible maximizer  $(\xi^*, v^*)$  is dual optimal

If  $x^*$  is optimal for primal problem,  $(x^*, \xi^*, v^*)$  is primal-dual optimal

Weak Duality

$$H(\xi^*, v^*) \leq F(x^*) \quad \text{last lower bound} \leq \text{optimal value}$$

Strong Duality

$$H(\xi^*, v^*) = F(x^*)$$

Slater's constraint qualification

if  $\exists x$  that is strictly feasible ( $F_\ell(x) < b_\ell$ ) then strong duality

hold for optimization problem  $H(\xi^*, v^*) = F(v^*)$

Saddle point

$$\sup L(x, \xi) = \sup (F_0(x) + \xi^*(Ax-y)) = \begin{cases} F_0(x) & \text{if } Ax=y \\ \infty & \text{ow.} \end{cases}$$

in other words, supremum of Lagrange function is  $\infty$  if  $x$  is not feasible

get the lowest value of feasible  $x$

$$F_0(x^*) = \inf_{x \in \mathbb{R}^n} \sup_{\xi \in \mathbb{R}^m} L(x, \xi)$$

minimize value of  $F_0$

Maximize the Lagrange dual function to have the best lower bound to optimal value  $F_0(x^*)$

$$H(\xi^*) = \sup_{\xi \in \mathbb{R}^m} \inf_{x \in \mathbb{R}^n} L(x, \xi)$$

Get the lower bound for  $F_0(x^*)$  ( $H(\xi^*)$ )  
Get the best lower bound =  $H(\xi^*)$

Weak duality  $\rightarrow H(\xi^*) \leq F_0(x^*)$

Strong duality  $\rightarrow H(\xi^*) = F_0(x^*)$

$$\inf_{x \in \mathbb{R}^n} \sup_{\xi \in \mathbb{R}^m} L(x, \xi) = \sup_{\xi \in \mathbb{R}^m} \inf_{x \in \mathbb{R}^n} L(x, \xi) \rightarrow$$

this means, in strong duality, min-max is interchangeable.  $\rightarrow$  Saddle point property!

if a primal-dual optimal  $(x^*, \xi^*)$  is a saddle point of Lagrange function

$$L(x^*, \xi) \leq L(x^*, \xi^*) \leq L(x, \xi^*)$$