

# CO233 - Computational Techniques

15th January 2020

## Vector and Matrix Norms

An orthonormal basis of  $\mathbb{R}^n$  are unit vectors that are pairwise mutually perpendicular; such that for  $(e_1, \dots, e_n)$ ;

- $e_i \cdot e_i = 1$
- $e_i \cdot e_j = 0$ , if  $i \neq j$

The standard canonical basis of  $\mathbb{R}^3$  are the  $i, j, k$  vectors, and similar in  $\mathbb{R}^2$ . However, we can form another orthonormal basis of  $\mathbb{R}^2$  by bisecting the angles as such;



If we take a vector  $\mathbf{v} \in \mathbb{R}^n$ , the Euclidean norm (or the  $\ell_2$ -norm) is defined as such;

$$\|\mathbf{v}\|_2 = \sqrt{\sum_{i=1}^n v_i^2}$$

A norm, a mapping  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}^+$ , must satisfy these 3 axioms;

- (i)  $\|\mathbf{v}\| > 0$  given that  $\mathbf{v} \neq \mathbf{0}$
- (ii)  $\|\lambda \mathbf{v}\| = |\lambda| \|\mathbf{v}\|$
- (iii)  $\|\mathbf{v} + \mathbf{w}\| \leq \|\mathbf{v}\| + \|\mathbf{w}\|$  (triangular inequality)

Some other ( $\ell_p$ ) norms are defined as follows;

$$\ell_1\text{-norm } \|\mathbf{v}\|_1 = \sum_{i=1}^n |v_i|$$

$$\ell_\infty\text{-norm } \|\mathbf{v}\|_\infty = \max\{|v_i| : 1 \leq i \leq n\}$$

$$\ell_p\text{-norm } \|\mathbf{v}\|_p = \left( \sum_{i=1}^n |v_i|^p \right)^{\frac{1}{p}}$$

In each dimension, we have the following;

- $n = 1$ 
  - $\|\mathbf{v}\|_1 = |v| = |v|$
  - $\|\mathbf{v}\|_2 = \sqrt{v^2} = |v|$
  - $\|\mathbf{v}\|_\infty = \max\{|v|\} = |v|$
- $n = 2$

We can represent this geometrically as such;



In our case  $\|\mathbf{v}\|_\infty = \max\{|v_1|, |v_2|\} = |v_1|$ , but the point is that it's either of the "sides" of the triangle. Obviously,  $\|\mathbf{v}\|_2 \geq \|\mathbf{v}\|_\infty$ , as it's the hypotenuse of the triangle, and similarly,  $\|\mathbf{v}\|_1 \geq \|\mathbf{v}\|_2$ , due to the triangle inequality. Therefore we have  $\|\mathbf{v}\|_1 \geq \|\mathbf{v}\|_2 \geq \|\mathbf{v}\|_\infty$ .

Even if the orthonormal base changes, the Euclidean norm stays the same, whereas the other norms can change. As such, we can say the  $\ell_2$ -norm is invariant under an **orthogonal transformation** (a basis change from an orthonormal bases to another orthonormal bases).

- $n = ?$  (general)

The goal is to prove  $\|\mathbf{v}\|_\infty \leq \|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$ . If we first take the squares of all of them, such that we have the following;

$$\begin{aligned}\|\mathbf{v}\|_1^2 &= \left(\sum_{i=1}^n |v_i|\right)^2 \\ \|\mathbf{v}\|_2^2 &= \sum_{i=1}^n |v_i|^2 \\ \|\mathbf{v}\|_\infty^2 &= (\max\{|v_i| : 1 \leq i \leq n\})^2\end{aligned}$$

Since the  $\ell_\infty$ -norm corresponds to a single  $v_i$ , it's obvious that the following inequality holds (since the  $\ell_2$ -norm squared has all the other terms squared, as well as the  $\ell_\infty$ -norm squared);

$$\|\mathbf{v}\|_2^2 = \sum_{i=1}^n |v_i|^2 \geq (\max\{|v_i| : 1 \leq i \leq n\})^2 = \|\mathbf{v}\|_\infty^2 \Rightarrow \|\mathbf{v}\|_2 \geq \|\mathbf{v}\|_\infty$$

To prove the other inequality, we see that the square of the sum of absolutes is greater than the sum of the squares, as the square of the sum contains the cross terms (which will be positive).

$$\|\mathbf{v}\|_2^2 = \sum_{i=1}^n |v_i|^2 \leq \left(\sum_{i=1}^n |v_i|\right)^2 = \|\mathbf{v}\|_1^2 \Rightarrow \|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$$

As such, we can conclude that  $\|\mathbf{v}\|_\infty \leq \|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$  in any dimension. ■

## Tutorial Question

Find the locus of vectors such that  $\|\mathbf{v}\|_p \leq 1$ , for  $p = 1, 2, \infty$  in  $n = 2$ ;



Imagine they're all shaded from the border to the origin.

## $\ell_p$ -norm

Our goal is to show that as  $p \rightarrow \infty$ , we get the definition of the  $\ell_\infty$ -norm previously stated. Take a vector  $\mathbf{v} \in \mathbb{R}^n$ , where both  $\mathbf{v}$  and  $n$  are fixed.

$$\|\mathbf{v}\|_p^p = \sum_{i=1}^n |v_i|^p$$

Obviously, this is greater than or equal to  $\|\mathbf{v}\|_\infty^p$ , as it would only be a single  $|v_i|^p$ . Similarly, it must be less than or equal to  $n|v_i|^p$ , as  $v_i$  is the maximum of all the components.

$$\|\mathbf{v}\|_\infty^p \leq \|\mathbf{v}\|_p^p = \sum_{i=1}^n |v_i|^p \leq n\|\mathbf{v}\|_\infty^p$$

Taking everything to the power of  $\frac{1}{p}$ , we obtain the following result (note that  $p > 0$  hence the signs don't change);

$$\|\mathbf{v}\|_\infty \leq \|\mathbf{v}\|_p \leq n^{\frac{1}{p}} \|\mathbf{v}\|_\infty$$

As  $p \rightarrow \infty$ , since  $n \geq 2$  ( $n = 1$  is shown to collapse to the same component), we have  $n^{\frac{1}{p}} \rightarrow 1$ , which sandwiches the middle term.

### Some Proposition ( $\ell_\infty$ -norm vs $\ell_2$ -norm)

The proposition is as follows; for a vector  $\mathbf{v} \in \mathbb{R}^n$ ,  $\|\mathbf{v}\|_2 \leq \sqrt{n}\|\mathbf{v}\|_\infty$ . To show this, we know that each of  $|v_i|$  is less than or equal to  $\|\mathbf{v}\|_\infty$ , by definition of the maximum. The same can be said for  $|v_i|^2$ , vs  $\|\mathbf{v}\|_\infty^2$ . Taking square roots, we have the following;

$$\|\mathbf{v}\|_2^2 = \sum_{i=1}^n |v_i|^2 \leq n\|\mathbf{v}\|_\infty^2 \Rightarrow \|\mathbf{v}\|_2 \leq \sqrt{n}\|\mathbf{v}\|_\infty$$

To show this holds similarly for  $\|\mathbf{v}\|_1 \leq \sqrt{n}\|\mathbf{v}\|_2$ , we employ the Cauchy-Schwarz inequality, which states  $|\mathbf{x} \cdot \mathbf{y}| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$ . The Cauchy-Schwarz inequality uses the fact that  $\mathbf{x} \cdot \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta$ . To do this, we need to define a sign function  $\text{sgn} : \mathbb{R} \rightarrow \{1, -1\}$  as follows;

$$\text{sgn } x = \begin{cases} 1 & x \geq 0 \\ -1 & x < 0 \end{cases}$$

We also need to craft a vector  $\mathbf{w}$ , as follows;

$$\begin{aligned} w_i &= \frac{\text{sgn } v_i}{\sqrt{n}} & 1 \leq i \leq n \\ \mathbf{v} \cdot \mathbf{w} &= \sum_{i=1}^n v_i w_i \\ &= \sum_{i=1}^n \frac{v_i \cdot \text{sgn } v_i}{\sqrt{n}} & \text{product of same sign becomes positive} \\ &= \sum_{i=1}^n \frac{|v_i|}{\sqrt{n}} \\ &= \sqrt{n} \sum_{i=1}^n |v_i| \\ &= \sqrt{n} \|\mathbf{v}\|_1 \\ \|\mathbf{w}\|_2 &= \sum_{i=1}^n \frac{\pm 1^2}{\sqrt{n}} \\ &= \sum_{i=1}^n \frac{1}{n} \\ &= 1 \end{aligned}$$

By **Cauchy-Schwarz**, we get;

$$\sqrt{n}\|\mathbf{v}\|_1 = |\mathbf{v} \cdot \mathbf{w}| \leq \|\mathbf{v}\|_2 \|\mathbf{w}\|_2 = \|\mathbf{v}\|_2$$

16th January 2020

Note that this recording has **no audio**, and therefore will just be the board transcribed. I honestly have no idea what he was doing in this lecture, it seems to just jump from topic to topic.

## Equivalence of Norms?

Take any two norms on  $\mathbb{R}^n$ ;  $\|\cdot\|_a$ , and  $\|\cdot\|_b$ .

$$\exists r, s \in \mathbb{R}^+ \forall \mathbf{v} \in \mathbb{R}^n [r\|\mathbf{v}\|_b \leq \|\mathbf{v}\|_a \leq s\|\mathbf{v}\|_b]$$

This means that norms in finite dimensional vector spaces are equivalent (no idea why, look it up).

## Convergence of Vector Sequences

$(\mathbf{r}_n)$  is a sequence of vectors, and  $(a_{i,j})$  is the  $i, j^{\text{th}}$  entry of  $\mathbf{A}$ . For a vector  $\mathbf{v}^{(m)} \in \mathbb{R}^n$ , where  $m = 0, 1, 2, \dots$

$$\mathbf{v}^{(m)} = \begin{bmatrix} v_1^{(m)} \\ v_2^{(m)} \\ \vdots \\ v_n^{(m)} \end{bmatrix}$$

For a vector sequence  $\mathbf{v}^{(m)}$  to converge to some vector  $\mathbf{v} \in \mathbb{R}^n$ , the following must hold;

$$\mathbf{v}^{(m)} \rightarrow \mathbf{v} \in \mathbb{R}^n \Leftrightarrow \lim_{m \rightarrow \infty} \|\mathbf{v}^{(m)} - \mathbf{v}\| \rightarrow 0$$

This is componentwise convergence, such that  $\forall i \in [1, n] [v_i^{(m)} \rightarrow v_i]$ .

## Matrix Norms

Vectors are a type of matrix. For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , the following properties of its norms must hold, where  $\|\cdot\| : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}_{\geq 0}$ ;

- (i)  $\|\mathbf{A}\| > 0$  given that  $\mathbf{A} \neq \mathbf{0}$
- (ii)  $\|\lambda \mathbf{A}\| = |\lambda| \|\mathbf{A}\|$
- (iii)  $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$
- (iv)  $\|\mathbf{BA}\| \leq \|\mathbf{B}\| \|\mathbf{A}\|$

$$\begin{array}{ccc} \mathbf{v} \in \mathbb{R}^n & \xrightarrow{\quad} & \boxed{\mathbf{A}} \xrightarrow{\quad} \mathbf{Av} \in \mathbb{R}^m \\ \|\cdot\|_a & & \|\cdot\|_b \\ & & \|\mathbf{Av}\|_b \leq \|\mathbf{A}\| \|\mathbf{v}\|_a \end{array}$$

For the following example, take  $(a_{i,j}) = \mathbf{A} \in \mathbb{R}^{m \times n}$ ;

$$a_j = \begin{bmatrix} a_{1,j} \\ a_{2,j} \\ \vdots \\ a_{m,j} \end{bmatrix}$$

the  $j^{\text{th}}$  column of  $\mathbf{A}$

$$a^i = [a_{i,1} \quad a_{i,2} \quad \cdots \quad a_{i,n}]$$

the  $i^{\text{th}}$  row of  $\mathbf{A}$

We have the following norms on matrices;

$$\|\mathbf{A}\|_1 = \max\{\|a_j\|_1 : 1 \leq j \leq n\}$$

$$\|\mathbf{A}\|_\infty = \max\{\|(a^i)^\top\|_1 : 1 \leq i \leq m\}$$

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{i,j}|^2} \quad \text{Frobenius norm}$$

$$\|\mathbf{A}\|_2 = \text{largest singular value of } \mathbf{A}$$

$$\text{let } \mathbf{A} = \begin{bmatrix} 2 & 3 & 1 & 4 \\ 1 & 3 & -1 & 5 \\ \sqrt{2} & 0 & -2 & 2 \end{bmatrix}$$

$$\begin{aligned} \|\mathbf{A}\|_1 &= \max\{3 + \sqrt{2}, 6, 4, 11\} \\ &= 11 \end{aligned}$$

$$\begin{aligned} \|\mathbf{A}\|_\infty &= \max\{10, 10, 4 + \sqrt{2}\} \\ &= 10 \end{aligned}$$

$$\begin{aligned} \|\mathbf{A}\|_F &= \sqrt{4 + 9 + 1 + 16 + 1 + 9 + 1 + 25 + 2 + 0 + 4 + 4} \\ &= 2\sqrt{19} \end{aligned}$$

Let there be two vector norms,  $\|\cdot\|_a$  on  $\mathbb{R}^n$  and  $\|\cdot\|_b$  on  $\mathbb{R}$ . If  $\|\cdot\|$  (matrix norm) satisfies

$$\forall \mathbf{A} \in \mathbb{R}^{m \times n}, x \in \mathbb{R}^n [\|\mathbf{A}x\|_b \leq \|\mathbf{A}\| \|x\|_a]$$

then  $\|\cdot\|$  is **consistent** with  $\|\cdot\|_a$  and  $\|\cdot\|_b$ . Additionally if  $a = b$ , then  $\|\cdot\|$  is **compatible** with  $\|\cdot\|_a$ . This gives us the following propositions;

- $\|\cdot\|_1$  (matrix norm) is compatible with  $\|\cdot\|_1$  (vector norm)
- $\|\cdot\|_2$  (matrix norm) is compatible with  $\|\cdot\|_2$  (vector norm)
- $\|\cdot\|_\infty$  (matrix norm) is compatible with  $\|\cdot\|_\infty$  (vector norm)
- $\|\cdot\|_F$  (matrix norm) is compatible with  $\|\cdot\|_2$  (vector norm)  $\Rightarrow \|\mathbf{A}x\|_2 \leq \|\mathbf{A}\|_F \|x\|_2$

Given a vector norm  $\|\cdot\|$  on  $\mathbb{R}^n$  then the matrix norm  $\|\cdot\|$  subordinate to vector norm  $\|\cdot\|$  is defined by

$$\|\mathbf{A}\| = \max\{\|\mathbf{A}x\| : \|x\| \leq 1\} = \max\{\|\mathbf{A}x\| : \|x\| = 1\} = \max\{\|\mathbf{A} \frac{x}{\|x\|}\| : x \neq 0\}$$

Using this, we can prove property (iii) (see above). We claim that  $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$  for matrix norm  $\|\cdot\|$  subordinate to vector norm  $\|\cdot\|$ .

$$\begin{aligned} \|\mathbf{A} + \mathbf{B}\| &= \max\{\|(\mathbf{A} + \mathbf{B})x\| : \|x\| \leq 1\} \\ &= \max\{\|\mathbf{A}x + \mathbf{B}x\| : \|x\| \leq 1\} \\ &\leq \max\{\|\mathbf{A}x\| + \|\mathbf{B}x\| : \|x\| \leq 1\} && \text{triangle inequality for } \|\cdot\| \\ &\leq \max\{\|\mathbf{A}x\| : \|x\| \leq 1\} + \max\{\|\mathbf{B}x\| : \|x\| \leq 1\} && \text{maximise independently} \\ &= \|\mathbf{A}\| + \|\mathbf{B}\| \end{aligned} \quad \blacksquare$$

We are also able to prove property (iv), which we claim to be  $\|\mathbf{B}\mathbf{A}\| \leq \|\mathbf{B}\| \|\mathbf{A}\|$ .

$$\begin{aligned} \|\mathbf{B}\mathbf{A}\| &= \max\{\|\mathbf{B}\mathbf{A}x\| : \|x\| \leq 1\} \\ \|\mathbf{B}\mathbf{A}x\| &= \|\mathbf{B}(\mathbf{A}x)\| \\ \|\mathbf{A}x\| &\leq \|\mathbf{A}\| \|x\| && \text{matrix norm subordinate to vector norm} \end{aligned}$$

To show the line above, we consider the two cases,  $\mathbf{x} \neq \mathbf{0}$  and  $\mathbf{x} = \mathbf{0}$ . If  $\mathbf{x} = \mathbf{0}$ , no work needs to be done, as it is trivial. I have no idea why this works, but he wrote it.

$$\begin{aligned} \text{show } \left\| \mathbf{A} \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\| &\leq \|\mathbf{A}\| \\ \text{but } \left\| \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\| &= 1 \\ \|\mathbf{A}\mathbf{x}\| &\leq \|\mathbf{A}\|\|\mathbf{x}\| \end{aligned} \quad \Rightarrow$$

Continuing on, we have

$$\begin{aligned} \|\mathbf{B}\mathbf{A}\mathbf{x}\| &= \|\mathbf{B}(\mathbf{A}\mathbf{x})\| \\ &\leq \|\mathbf{B}\|\|\mathbf{A}\mathbf{x}\| \\ &\leq \|\mathbf{B}\|\|\mathbf{A}\|\|\mathbf{x}\| \\ \|\mathbf{B}\mathbf{A}\| &= \max\{\|\mathbf{B}\mathbf{A}\mathbf{x}\| : \|\mathbf{x}\| \leq 1\} \\ &\leq \max\{\|\mathbf{B}\|\|\mathbf{A}\|\|\mathbf{x}\| : \|\mathbf{x}\| \leq 1\} \\ &= \|\mathbf{B}\|\|\mathbf{A}\| \max\{\|\mathbf{x}\| : \|\mathbf{x}\| \leq 1\} && \text{obviously 1, as bounded on top} \\ &= \|\mathbf{B}\|\|\mathbf{A}\| && \blacksquare \end{aligned}$$

## Complex Vectors

$$\mathbb{C}^n = \left\{ \mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} : v_i \in \mathbb{C} \right\}$$

$$\begin{aligned} z \in \mathbb{C} &= a + ib && a \in \mathbb{R}, b \in \mathbb{R}, i = \sqrt{-1} \\ z^* &= a - ib \\ |z| &= \sqrt{a^2 + b^2} \end{aligned}$$

Take a linear map  $f : \mathbb{C}^n \rightarrow \mathbb{C}^m$ , the same properties hold;

$$f(a\mathbf{v} + b\mathbf{w}) = af(\mathbf{v}) + bf(\mathbf{w}) \quad a, b \in \mathbb{C}$$

We also want to define something similar to the dot product in  $\mathbb{R}^n$ ;

$$\begin{aligned} \mathbf{v} \cdot \mathbf{w} &= \|\mathbf{v}\|_2 \|\mathbf{w}\|_2 \cos \theta_{\mathbf{v}, \mathbf{w}} \\ \mathbf{v} \cdot \mathbf{v} &= \sqrt{\|\mathbf{v}\|_2^2} \\ &= \|\mathbf{v}\|_2 \\ \langle \mathbf{v}, \mathbf{w} \rangle &= \sum_{i=1}^n v_i^* w_i \\ \langle \mathbf{v}, \mathbf{v} \rangle &= \sum_{i=1}^n v_i^* v_i \\ &= \sum_{i=1}^n |v_i|^2 \end{aligned}$$

The standard basis in  $\mathbb{R}^n$  is defined as  $(\mathbf{e}_1, \dots, \mathbf{e}_n)$  where

$$\mathbf{e}_j = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} 1^{\text{st}} \\ 2^{\text{nd}} \\ \vdots \\ j^{\text{th}} \\ \vdots \\ n^{\text{th}} \end{matrix}$$

For any vector  $\mathbf{v} \in \mathbb{C}^n$ , it can be written in the standard basis as such;

$$\mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \\ = v_1 \mathbf{e}_1 + \cdots + v_n \mathbf{e}_n$$

## Basis Change, Again

Let the linear map  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ .

$$\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n)$$

an ordered basis of  $\mathbb{R}^n$

$$\mathbf{D} = (\mathbf{d}_1, \dots, \mathbf{d}_m)$$

an ordered basis of  $\mathbb{R}^m$

Find the matrix  $\mathbf{A}$  ( $\mathbf{A} := f_{\mathbf{D}\mathbf{B}}$ ) representing  $f$  with respect to (?)  $\mathbf{B}$  and  $\mathbf{D}$ .

$$\mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \quad \text{coordinates of a point } \mathbf{p} \in \mathbb{R}^n$$

$$\mathbf{p} = \sum_{j=1}^n v_j \mathbf{b}_j$$

$$\mathbf{A}\mathbf{v}$$

should be coordinate of  $f(\mathbf{p}) \in \mathbb{R}^m$

$$f(\mathbf{p}) = \sum_{i=1}^m (\mathbf{A}\mathbf{v})_i \mathbf{d}_i$$

$$f(\mathbf{b}_j) \in \mathbb{R}^m = \sum_{i=1}^m a_{i,j} \mathbf{d}_i \quad 1 \leq j \leq n$$

Take  $\mathbf{B} \in \mathbb{R}^{m \times n}$ , where  $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n]$ , and  $\mathbf{e}_1, \dots, \mathbf{e}_m$  being the standard basis (see previous).

$$\mathbf{B}\mathbf{e}_j = \mathbf{B} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow j^{\text{th}}$$

$$= b_{1,j} \mathbf{e}_1 + \dots + b_{m,j} \mathbf{e}_m$$

Suppose  $m = n$  and also  $f = \text{id}$  (identity), but  $\mathbf{B}$  and  $\mathbf{D}$  are different. The matrix  $(\text{id})_{\mathbf{D}\mathbf{B}}$  represents a change of basis from  $\mathbf{B}$  to  $\mathbf{D}$ . The point  $\mathbf{p} \in \mathbb{R}^n$  has coordinates  $\mathbf{x} \in \mathbb{R}^n$  with respect to  $\mathbf{B}$ , and  $\mathbf{y} \in \mathbb{R}^n$  with respect to  $\mathbf{D}$ .

$$\mathbf{D}\mathbf{y} = \mathbf{p} = \mathbf{B}\mathbf{x} = x_1 \mathbf{b}_1 + \cdots + x_n \mathbf{b}_n$$

From this we gather  $\mathbf{D}\mathbf{y} = \mathbf{B}\mathbf{x}$ , therefore  $\mathbf{y} = \mathbf{D}^{-1} \mathbf{B}\mathbf{x} = (\text{id})_{\mathbf{D}\mathbf{B}} \mathbf{x}$ , which means that

$$(\text{id})_{\mathbf{D}\mathbf{B}} = \mathbf{D}^{-1} \mathbf{B}$$

Some stuff on functions between bases?

$$\begin{array}{ccccc} \mathbb{R}^n & \xrightarrow{f} & \mathbb{R}^m & \xrightarrow{g} & \mathbb{R}^k \\ \mathbf{B} & & \mathbf{C} & & \mathbf{D} \\ \mathbb{R}_{\mathbf{B}}^n & \xrightarrow{f_{\mathbf{C}\mathbf{B}}} & \mathbb{R}_{\mathbf{C}}^m & \xrightarrow{g_{\mathbf{D}\mathbf{C}}} & \mathbb{R}_{\mathbf{D}}^k \end{array}$$

$$g_{\mathbf{D}\mathbf{C}} f_{\mathbf{C}\mathbf{B}} = (g \circ f)_{\mathbf{D}\mathbf{B}}$$

This then goes into change of basis, but see last year's **CO145**.

22nd January 2020

### Tutorial Question

$f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is a linear map, and  $\mathbf{E} = (\mathbf{e}_1, \mathbf{e}_2)$  is an ordered basis.

$$f(\mathbf{e}_1) = 5\mathbf{e}_1 - 6\mathbf{e}_2$$

$$f(\mathbf{e}_2) = 3\mathbf{e}_1 + \mathbf{e}_2$$

We only care about what the linear map does to the ordered basis, as anything else can be done by linearity.

- (i) Find  $f_{\mathbf{E}\mathbf{E}}$ , the matrix representation of  $f$  in  $\mathbf{E}$  - note that this has the same input space as the output space, but it can be different.

The first column can be done by reading the entry for  $f(\mathbf{e}_1)$ , and similarly for the second column as follows;

$$f_{\mathbf{E}\mathbf{E}} = \begin{bmatrix} 5 & 3 \\ -6 & 1 \end{bmatrix}$$

- (ii) If we have another ordered basis  $\mathbf{D} = (\mathbf{d}_1, \mathbf{d}_2)$ , where  $\mathbf{d}_1 = \mathbf{e}_1 - \mathbf{e}_2$  and  $\mathbf{d}_2 = \mathbf{e}_1 + \mathbf{e}_2$ , find  $f_{\mathbf{D}\mathbf{D}}$ .

$$\begin{array}{ccc} \mathbb{R}_{\mathbf{E}}^2 & \xrightarrow{f_{\mathbf{E}\mathbf{E}}} & \mathbb{R}_{\mathbf{E}}^2 \\ \mathbf{I}_{\mathbf{E}\mathbf{D}} \uparrow & & \downarrow \mathbf{I}_{\mathbf{D}\mathbf{E}} = (\mathbf{I}_{\mathbf{E}\mathbf{D}})^{-1} \\ \mathbb{R}_{\mathbf{D}}^2 & \xrightarrow{f_{\mathbf{D}\mathbf{D}}} & \mathbb{R}_{\mathbf{D}}^2 \end{array}$$

$\mathbf{I}_{\mathbf{E}\mathbf{D}}$  can easily be obtained by reading the entries for  $\mathbf{d}_1$  for the first column, and similarly for  $\mathbf{d}_2$  in the second column;

$$\mathbf{I}_{\mathbf{E}\mathbf{D}} = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$\mathbf{I}_{\mathbf{D}\mathbf{E}} = (\mathbf{I}_{\mathbf{E}\mathbf{D}})^{-1}$$

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

$$f_{\mathbf{D}\mathbf{D}} = \mathbf{I}_{\mathbf{D}\mathbf{E}} f_{\mathbf{E}\mathbf{E}} \mathbf{I}_{\mathbf{E}\mathbf{D}}$$

$$= \frac{1}{2} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 5 & 3 \\ -6 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$= \frac{1}{2} \begin{bmatrix} 1 & 13 \\ -5 & 3 \end{bmatrix}$$

### Eigenvalues + Generalised Eigenvectors

Working with a matrix  $\mathbf{A} \in \mathbb{C}^{m \times m}$ . For an eigenvector  $\mathbf{v} \in \mathbb{C}^m \setminus \{\mathbf{0}\}$ , and an eigenvalue  $\lambda \in \mathbb{C}$ ,  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \Rightarrow (\mathbf{A} - \lambda\mathbf{I})\mathbf{v} = \mathbf{0} \Rightarrow |\mathbf{A} - \lambda\mathbf{I}| = 0 \Rightarrow P_{\mathbf{A}}(\lambda) = 0$  (characteristic polynomial). This complex polynomial will be of degree  $m$ , and it will have precisely  $m$  roots (including multiplicity). Suppose  $P_{\mathbf{A}}(\lambda) = 0$ , then  $(\mathbf{A} - \lambda\mathbf{I})\mathbf{v} = \mathbf{0}$  has a solution where  $\mathbf{v} \neq \mathbf{0}$ .

Assume we have  $\lambda_1, \dots, \lambda_t$  distinct eigenvalues, meaning that  $P_{\mathbf{A}}(\lambda_i) = 0$  for  $1 \leq i \leq t$ . This means we can write the characteristic polynomial as;

$$P_{\mathbf{A}}(\lambda) = (\lambda - \lambda_1)^{m_1} (\lambda - \lambda_2)^{m_2} \dots (\lambda - \lambda_t)^{m_t}$$



Where  $m_i$  is the **algebraic multiplicity** of  $\lambda_i$ .  $m_i \in \mathbb{N}$ , and also  $1 \leq m_i \leq m$ , as it must not exceed the dimension of the matrix. On the other hand, the **geometric multiplicity** of  $\lambda_i$  is  $\ell_i$ , which is the **nullity** of  $(\mathbf{A} - \lambda_i \mathbf{I})$ . The nullity is the dimension of the kernel / null-space.  $1 \leq \ell_i \leq m_i$ , as we already have at least one non-zero solution from  $(\mathbf{A} - \lambda_i \mathbf{I})\mathbf{v} = \mathbf{0}$ .

In the nice case, we have  $\ell_i = m_i$ , for  $1 \leq i \leq t$ , which means the matrix is diagonalisable. We have  $m_i$  linearly independent vectors  $(\mathbf{v}_{i,1}, \mathbf{v}_{i,2}, \dots, \mathbf{v}_{i,m_i})$  which satisfy

$$\mathbf{A}\mathbf{v}_{i,j} = \lambda_i \mathbf{v}_{i,j} \text{ for } 1 \leq j \leq m_i$$

If we take these eigenvectors as an ordered basis;

$$\mathbf{B} = [\mathbf{v}_{1,1}, \mathbf{v}_{1,2}, \dots, \mathbf{v}_{1,m_1}, \mathbf{v}_{2,1}, \mathbf{v}_{2,2}, \dots, \mathbf{v}_{2,m_2}, \dots, \mathbf{v}_{t,1}, \mathbf{v}_{t,2}, \dots, \mathbf{v}_{t,m_t}]$$

We also want to note that  $\sum_{i=1}^t m_i = m$ , as that is the degree of the characteristic polynomial. Multiplying the basis by the original matrix, we get;

$$\mathbf{A}\mathbf{B} = [\lambda_1 \mathbf{v}_{1,1}, \lambda_1 \mathbf{v}_{1,2}, \dots, \lambda_1 \mathbf{v}_{1,m_1}, \lambda_2 \mathbf{v}_{2,1}, \lambda_2 \mathbf{v}_{2,2}, \dots, \lambda_2 \mathbf{v}_{2,m_2}, \dots, \lambda_t \mathbf{v}_{t,1}, \lambda_t \mathbf{v}_{t,2}, \dots, \lambda_t \mathbf{v}_{t,m_t}]$$

Since all the columns of  $\mathbf{B}$  are linearly independent, by our definition, the inverse  $\mathbf{B}^{-1}$  exists. Therefore, we can write

$$\mathbf{B}^{-1}\mathbf{A}\mathbf{B} = \begin{bmatrix} \lambda_1 & & & & & \\ & \ddots & & & & \\ & & \lambda_1 & & & \\ & & & \lambda_2 & & \\ & & & & \ddots & \\ & & & & & \lambda_2 \\ & & & & & & \ddots & \\ & & & & & & & \lambda_t \\ & & & & & & & & \ddots & \\ & & & & & & & & & \lambda_t \end{bmatrix} \quad (\text{everything else is } 0)$$

Which has  $m_1$  instances of  $\lambda_1$ , followed by  $m_2$  instances of  $\lambda_2$ , and so on, until  $m_t$  instances of  $\lambda_t$ .

**Example for  $\ell_i = m_i$**

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

$$|\mathbf{A} - \lambda \mathbf{I}| = (\lambda - 3)^2(\lambda - 5)$$

$$\lambda_1 = 3$$

two linearly independent eigenvectors

$$\mathbf{v}_{1,1} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\mathbf{v}_{1,2} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

$$\lambda_2 = 5$$

$$\mathbf{v}_{2,1} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$

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

$$\mathbf{B}^{-1} = \frac{1}{2} \begin{bmatrix} -2 & 2 & -2 \\ -1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

$$\mathbf{B}^{-1}\mathbf{A}\mathbf{B} = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 5 \end{bmatrix}$$

**Trivial Example for  $\ell_i < m_i$**

$$\mathbf{A} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$

$$|\mathbf{A} - \lambda\mathbf{I}| = \lambda^2$$

$$\lambda_1 = 0$$

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

only solution, hence  $\ell_1 = 1 < 2 = m_1$

$$(\mathbf{A} - 0\mathbf{I}) \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

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

although this vector is not mapped to zero, it is mapped to something that **will** be mapped to zero

$$(\mathbf{A} - 0\mathbf{I})^2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} = (\mathbf{A} - 0\mathbf{I})(\mathbf{A} - 0\mathbf{I}) \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= (\mathbf{A} - 0\mathbf{I}) \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \mathbf{0}$$

We say  $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$  is a generalised eigenvector for  $\lambda_1 = 0$ . A vector which is not mapped by  $(\mathbf{A} - \lambda\mathbf{I})$  to  $\mathbf{0}$ , but is  $\mathbf{0}$  when iterated once more.

**Less Trivial Example**

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$|\mathbf{A} - \lambda\mathbf{I}| = (1 - \lambda)^3$$

$$\lambda_1 = 1$$

$$(\mathbf{A} - \mathbf{I})\mathbf{x} = \mathbf{0}$$

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

$$m_1 = 3$$

$$\Leftrightarrow$$

this has rank 1, and therefore by rank-nullity theorem (rank + nullity = 3), has 2 linearly independent solutions, therefore  $\ell_1 = 2 < 3 = m_1$

$$\mathbf{v}_1 = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}$$

$$\mathbf{v}_2 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

to find the generalised eigenvector  $\mathbf{v}_3$ , we want to find some vector that is mapped by  $(\mathbf{A} - 1\mathbf{I})$  to the eigenspace, which is some linear combination of  $\mathbf{v}_1$  and  $\mathbf{v}_2$

$$\begin{aligned} (\mathbf{A} - \mathbf{I})\mathbf{v}_3 &= \alpha_1\mathbf{v}_1 + \alpha_2\mathbf{v}_2 && \Leftrightarrow \\ \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} &= \begin{bmatrix} 0 \\ \alpha_1 \\ -\alpha_1 \end{bmatrix} + \begin{bmatrix} \alpha_2 \\ 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} x_2 + x_3 \\ 0 \\ 0 \end{bmatrix} && \Leftrightarrow \\ \alpha_1 &= 0 \\ x_2 + x_3 &= \alpha_2 && \text{let } x_2 = 0 \Rightarrow x_3 = \alpha_2 = 1 \\ \mathbf{v}_3 &= \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ \mathbf{B} &= [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3] \\ &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix} \\ \mathbf{B}^{-1}\mathbf{A}\mathbf{B} &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} && \text{the extra 1 is from } \mathbf{v}_2? \end{aligned}$$

This is in Jordan Normal Form. For some  $\lambda_i$  eigenvalue, and  $\ell_i \leq m_i$ , the sum of the sizes of the blocks is  $m_i$ , and the number of blocks is  $\ell_i$ . If  $\lambda_i$  is an eigenvalue of  $\mathbf{A}$  with algebraic multiplicity  $m_i$ , then the nullity of  $(\mathbf{A} - \lambda_i\mathbf{I})^{m_i} = m_i$ .

**Definition:**  $\mathbf{v} \in \mathbb{R}^m$  is a generalised eigenvector for  $\lambda_i$  if  $(\mathbf{A} - \lambda_i\mathbf{I})^{m_i}\mathbf{v} = \mathbf{0}$ . The maximum iterations is  $m_i$ , but can be less.

**23rd January 2020**

In general;

- $\forall j [\ell_j = m_j] \Rightarrow \mathbf{A}$  is diagonalisable
- $\exists j [\ell_j < m_j] \Rightarrow \mathbf{A}$  is diagonalisable or **almost** diagonalisable

The following is a very poor explanation of what the form looks like, because I don't know how to draw blocks in L<sup>A</sup>T<sub>E</sub>X. See Panopto at timestamp 13:17 for a proper drawing.

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} \mathbf{B}_{\lambda_1} & & \\ & \ddots & \\ & & \mathbf{B}_{\lambda_t} \end{bmatrix} \\ \mathbf{B}_j &= \begin{bmatrix} \mathbf{J}_j & & \\ & \ddots & \\ & & \mathbf{J}_j \end{bmatrix} \end{aligned} \quad \text{all } \mathbf{J}_j \text{ can be of different sizes, } \mathbf{B}_j \in \mathbb{R}^{m_j \times m_j}$$

$$\mathbf{J}_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix}$$

## Spectral Decomposition of Symmetric Matrices

In almost all areas of computer science, symmetric matrices are the dominant paradigm. All the eigenvalues of real symmetric matrices are real, and are always diagonalisable. A symmetric matrix is always a square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , which also satisfies  $\mathbf{A}^\top = \mathbf{A}$ , meaning it is symmetric along the diagonal (the elements above the diagonal determine the elements below, and vice versa). We are to prove the following properties;

- (i) All eigenvalues of  $\mathbf{A}$  are real.

First assume that  $\mathbf{v} \in \mathbb{C}^n$ , as we don't yet know that it is real, therefore we can have the case  $\lambda \in \mathbb{C}$ .

$$\begin{aligned} \mathbf{A}\mathbf{v} &= \lambda\mathbf{v} & \Rightarrow \\ \mathbf{v}^{*\top} \mathbf{A}\mathbf{v} &= \lambda \mathbf{v}^{*\top} \mathbf{v} \\ &= \lambda \|\mathbf{v}\|^2 \\ \mathbf{v}^{*\top} \mathbf{A}\mathbf{v} &= (\mathbf{v}^{*\top} \mathbf{A}) \mathbf{v} \\ &= ((\mathbf{A}\mathbf{v})^\top)^* \mathbf{v} \end{aligned} \tag{1}$$

briefly proving the above;

$$\begin{aligned} ((\mathbf{A}\mathbf{v})^\top)^* &= (\mathbf{v}^\top \mathbf{A}^\top)^* & \text{taking the transpose inside the bracket} \\ &= (\mathbf{v}^\top)^* \mathbf{A}^* & \text{taking the conjugate inside and } \mathbf{A}^\top = \mathbf{A} \\ &= (\mathbf{v}^\top)^* \mathbf{A} & \mathbf{A} \text{ is real valued matrix} \end{aligned}$$

continuing on;

$$\begin{aligned} &= ((\lambda\mathbf{v})^\top)^* \mathbf{v} & \mathbf{v} \text{ is an eigenvector} \\ &= \lambda^* (\mathbf{v}^\top)^* \mathbf{v} & \text{transposing a scalar does nothing} \\ &= \lambda^* \|\mathbf{v}\|^2 & (2) \end{aligned}$$

By comparing results, (1) = (2), hence  $\lambda^* = \lambda$  (dividing through by  $\|\mathbf{v}\|^2$  is allowed, since we know it is non-zero by properties of vector norms, and  $\mathbf{v} \neq \mathbf{0}$  by properties of eigenvectors). As the complex conjugate of  $\lambda$  is the same as  $\lambda$ , we can conclude that  $\lambda \in \mathbb{R}$ .

- (ii) Eigenvectors corresponding to different eigenvalues of  $\mathbf{A}$  are perpendicular.

Note that with the above result, we know that  $\mathbf{v} \in \mathbb{R}^n$ , as there is no possible way to obtain complex values with a real valued matrix and a real valued eigenvalue.

$$\begin{aligned} \mathbf{A}\mathbf{v}_1 &= \lambda_1 \mathbf{v}_1 \\ \mathbf{A}\mathbf{v}_2 &= \lambda_2 \mathbf{v}_2 \\ \lambda_1 &\neq \lambda_2 \\ \mathbf{v}_2^\top \mathbf{A}\mathbf{v}_1 &= \lambda_1 \mathbf{v}_2^\top \mathbf{v}_1 \\ &= \lambda_1 (\mathbf{v}_1 \cdot \mathbf{v}_2) & (1) \\ \mathbf{v}_2^\top \mathbf{A}\mathbf{v}_1 &= (\mathbf{v}_2^\top \mathbf{A}) \mathbf{v}_1 \\ &= (\mathbf{A}\mathbf{v}_2)^\top \mathbf{v}_1 & \text{adding a transpose and } \mathbf{A}^\top = \mathbf{A} \\ &= \lambda_2 \mathbf{v}_2^\top \mathbf{v}_1 & \mathbf{v}_2 \text{ is an eigenvector} \\ &= \lambda_2 (\mathbf{v}_1 \cdot \mathbf{v}_2) & (2) \end{aligned}$$

Similarly, comparing results (1) and (2), we get

$$\begin{aligned}
 \lambda_1(\mathbf{v}_1 \cdot \mathbf{v}_2) &= \lambda_2(\mathbf{v}_1 \cdot \mathbf{v}_2) && \Leftrightarrow \\
 \lambda_1(\mathbf{v}_1 \cdot \mathbf{v}_2) - \lambda_2(\mathbf{v}_1 \cdot \mathbf{v}_2) &= 0 \\
 (\lambda_1 - \lambda_2)(\mathbf{v}_1 \cdot \mathbf{v}_2) &= 0 && \Leftrightarrow \\
 \lambda_1 &= \lambda_2 \\
 \text{or } \mathbf{v}_1 \cdot \mathbf{v}_2 &= 0
 \end{aligned}$$

The former is not possible due to the condition that they are different eigenvalues, therefore  $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$ . As neither are the zero vector, by definition of eigenvectors,  $\cos \theta = 0$ , therefore they are perpendicular.

(iii) (Tutorial 3) If  $\lambda_j$  is an eigenvalue of  $\mathbf{A}$  then  $\ell_j = m_j$ , therefore  $\mathbf{A}$  is diagonalisable.

Come back to this in two weeks?

To conclude, we get all the eigenvectors and we normalise them (divide by the  $\ell_2$ -norm), and we can assume they are pairwise orthogonal. There is an orthonormal basis for  $\mathbb{R}^n$  consisting of eigenvectors of  $\mathbf{A}$ ;  $(\mathbf{v}_1, \dots, \mathbf{v}_n)$ , where  $\mathbf{v}_j \in \mathbb{R}^n$ . This has the following properties;

- $\mathbf{A}\mathbf{v}_j = \lambda_j\mathbf{v}_j$
- $\mathbf{v}_i \perp \mathbf{v}_j$ , when  $i \neq j$  perpendicular
- $\|\mathbf{v}_i\| = 1$  for  $i = 1, \dots, n$

## Examples

For an example in  $\mathbb{R}^{2 \times 2}$ ;

$$\begin{aligned}
 \mathbf{A} &= \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \\
 |\mathbf{A} - \lambda \mathbf{I}| &= (1 - \lambda)^2 - 4 \\
 \lambda_1 &= 3 \\
 \lambda_2 &= -1 \\
 \mathbf{v}_1 &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} && \text{normalise straight away} \\
 \mathbf{v}_2 &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\
 \mathbf{V} &= [\mathbf{v}_1, \mathbf{v}_2] && \text{orthonormal basis of } \mathbb{R}^2 \\
 &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \\
 \mathbf{A}\mathbf{V} &= [\mathbf{A}\mathbf{v}_1, \mathbf{A}\mathbf{v}_2] \\
 &= [3\mathbf{v}_1, -\mathbf{v}_2]
 \end{aligned}$$

For a symmetric matrix, if we have a matrix  $\mathbf{V}$  representing the orthonormal basis of the eigenvectors,  $\mathbf{V}^{-1} = \mathbf{V}^\top$

$$\begin{aligned}
 \mathbf{V}^{-1}\mathbf{A}\mathbf{V} &= \mathbf{V}^\top \mathbf{A}\mathbf{V} \\
 &= \mathbf{V}^\top [3\mathbf{v}_1, -\mathbf{v}_2] \\
 &= \begin{bmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \end{bmatrix} [3\mathbf{v}_1, -\mathbf{v}_2] \\
 &= \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix} && \text{dot product with unit length } (\mathbf{v}_1 \perp \mathbf{v}_2)
 \end{aligned}$$

For an example in  $\mathbb{R}^{3 \times 3}$ ,

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 3 \\ 1 & 3 & 1 \\ 3 & 1 & 1 \end{bmatrix}$$

now he just shows some intuition on getting an eigenvalue?

$$\begin{aligned} |\mathbf{A} - \lambda \mathbf{I}| &= \begin{vmatrix} 1 - \lambda & 1 & 3 \\ 1 & 3 - \lambda & 1 \\ 3 & 1 & 1 - \lambda \end{vmatrix} \\ &= \begin{vmatrix} 5 - \lambda & 1 & 3 \\ 5 - \lambda & 3 - \lambda & 1 \\ 5 - \lambda & 1 & 1 - \lambda \end{vmatrix} \\ &= 0 \end{aligned}$$

$$\lambda_1 = 5$$

$$\text{trace}(\mathbf{A}) = \lambda_1 + \lambda_2 + \lambda_3$$

$$= 5$$

$\Rightarrow$

$$\lambda_2 + \lambda_3 = 0$$

$$|\mathbf{A}| = \lambda_1 \lambda_2 \lambda_3$$

$$= -20$$

$\Rightarrow$

$$\lambda_2 \lambda_3 = -4$$

$$\lambda_2 = 2$$

$$\lambda_3 = -2$$

$$\mathbf{v}_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\mathbf{v}_2 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$$

$$\mathbf{v}_3 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3]$$

$$\mathbf{V}^{-1} = \mathbf{V}^\top$$

$$\mathbf{V}^{-1} \mathbf{A} \mathbf{V} = \mathbf{V}^\top \mathbf{A} \mathbf{V}$$

$$= \begin{bmatrix} 5 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & -2 \end{bmatrix}$$

spectral decomposition of a symmetric matrix

**Proof of  $\mathbf{V}^{-1} = \mathbf{V}^\top$**

We define  $\mathbf{B} \in \mathbb{R}^{n \times n}$  as an **orthogonal matrix** if all columns  $\mathbf{b}_j$  for  $1 \leq j \leq n$  form an orthonormal basis of  $\mathbb{R}^n$ , such that the two properties hold;

- $\mathbf{b}_j \cdot \mathbf{b}_i = 0$  when  $i \neq j$
- $\mathbf{b}_j \cdot \mathbf{b}_j = 1$

We want to prove that  $B^{-1} = B^\top$ :

$$\begin{aligned}
B &= [b_1, \dots, b_n] \\
B^\top &= \begin{bmatrix} b_1^\top \\ \vdots \\ b_n^\top \end{bmatrix} \\
B^\top B &= \begin{bmatrix} b_1^\top \\ \vdots \\ b_n^\top \end{bmatrix} [b_1, \dots, b_n] \\
&= \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & 1 \end{bmatrix} \\
&= I_n
\end{aligned}$$

### Preservation of Length

The length of a vector  $\mathbf{x}$  is preserved under orthogonal transformation;

$$\begin{aligned}
(\mathbf{V}\mathbf{x}) \cdot (\mathbf{V}\mathbf{x}) &= (\mathbf{V}\mathbf{x})^\top (\mathbf{V}\mathbf{x}) \\
&= \mathbf{x}^\top \mathbf{V}^\top \mathbf{V} \mathbf{x} \\
&= \mathbf{x}^\top \mathbf{V}^{-1} \mathbf{V} \mathbf{x} \\
&= \mathbf{x}^\top \mathbf{x} \\
&= \mathbf{x} \cdot \mathbf{x}
\end{aligned}$$

### Geometric Intuition

Take a real symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , such that  $\mathbf{A}^\top = \mathbf{A}$ . Assume we have processed it to get the eigenvalues  $\lambda_1, \dots, \lambda_n$ , such that  $\mathbf{V}^\top \mathbf{A} \mathbf{V} = \text{diag}(\lambda_1, \dots, \lambda_n)$ , and  $\mathbf{V}$  is an orthonormal matrix formed of the eigenvectors.

Take a vector  $\mathbf{x}$  from the input space, to  $\mathbf{A}\mathbf{x}$  in the output space, mapped by  $\mathbf{A}$ ;

$$\begin{array}{ccc}
\text{input space} & \xrightarrow{\mathbf{A} = f_{EE}} & \text{output space} \\
\mathbf{x} & \xrightarrow{\quad} & \mathbf{A}\mathbf{x} \\
\mathbb{R}_E^n & \xrightarrow{f_{EE} = \mathbf{A}} & \mathbb{R}_E^n \\
\uparrow I_{EV} = \mathbf{V} & & \downarrow I_{VE} = (\mathbf{I}_{ED})^{-1} = \mathbf{V}^{-1} \\
\mathbb{R}_V^n & \xrightarrow{f_{DD} = \mathbf{V}^{-1} \mathbf{A} \mathbf{V}} & \mathbb{R}_V^n
\end{array}$$

Consider a unit sphere, created by  $\|\mathbf{x}\| = 1$ , where  $\mathbf{x} \in \mathbb{R}^3$ . The transformations  $\mathbf{V}^{-1} = \mathbf{V}^\top$  and  $\mathbf{V}$  are rotations as they preserve length. Take a point  $\mathbf{x}$  on the surface of this unit sphere, applying  $\mathbf{V}^{-1}$  to it simply rotates it to another point on the surface of the sphere (going from  $\mathbb{R}_E^3$  to  $\mathbb{R}_V^3$ ) in the input space. Let this point have coordinates  $\mathbf{u}$  with respect to  $\mathbf{V}$ . Under the diagonal map  $\text{diag}(\lambda_1, \dots, \lambda_n)$ , we get the following;

$$u_1 \mathbf{v}_1 + \dots + u_n \mathbf{v}_n \mapsto \lambda_1 u_1 \mathbf{v}_1 + \dots + \lambda_n u_n \mathbf{v}_n$$

As  $\mathbf{u}$  is on the sphere, we can say;

$$\|\mathbf{u}\|_2^2 = u_1^2 + \dots + u_n^2 = 1$$

Therefore, the mapped point  $\mathbf{y}$  can be written as  $\begin{bmatrix} \lambda_1 u_1 \\ \vdots \\ \lambda_n u_n \end{bmatrix}$ , satisfying (in three dimensions);

$$\frac{y_1^2}{\lambda_1^2} + \frac{y_2^2}{\lambda_2^2} + \frac{y_3^2}{\lambda_3^2} = 1$$

This gives a locus of an ellipsoid.

## Rank-Nullity Theorem

Take a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .

$$\text{Im}(\mathbf{A}) \subseteq \mathbb{R}^m = \{\mathbf{Ax} : \mathbf{x} \in \mathbb{R}^n\}$$

$$\dim(\text{Im}(\mathbf{A})) = \text{column rank of } \mathbf{A}$$

$$\text{Null}(\mathbf{A}) \subseteq \mathbb{R}^n = \{\mathbf{v} : \mathbf{Av} = \mathbf{0}\} \quad \text{same as kernel}$$

$$\text{rk}(\text{column space}) = \text{rk}(\text{row space})$$

$$\text{row space } \mathbf{A} = \text{column space } \mathbf{A}^\top$$

For elementary row operations, we have the following property;

(I) they are invertible.

If there is one of the three elementary row operations  $R$  that takes  $\mathbf{A} \rightsquigarrow \mathbf{B}$ , then  $\exists R^{-1}$  that takes  $\mathbf{B} \rightsquigarrow \mathbf{A}$ .

If  $\mathbf{A} \rightsquigarrow \mathbf{B}$ , then  $\mathbf{B} = \mathbf{M}_R \mathbf{A}$ . We obtain  $\mathbf{M}_R$  by applying  $R$  to  $\mathbf{I}_m$ , such that  $\mathbf{I}_m \rightsquigarrow \mathbf{M}_R$ .

A reduction to row echelon form can be written as follows (with  $-$  as a pivot, and  $\times$  being any number);

$$\mathbf{A} \xrightarrow{R_1} \dots \xrightarrow{R_t} \mathbf{C} = \begin{bmatrix} - & \times & \times & \times & \times & \times \\ 0 & - & \times & \times & \times & \times \\ 0 & 0 & 0 & - & \times & \times \\ 0 & 0 & 0 & 0 & - & \times \\ 0 & 0 & 0 & 0 & 0 & - \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

In general it creates a staircase of the pivots, and everything below it is 0. Note that the example above is just a random example, as long as the general structure remains, it's fine.

An ERO doesn't change the row space, as it is operating on the rows. Swapping rows, multiplying by a non-zero number, and adding rows will not change anything. The claim is that while the column space can change, the rank of the column space doesn't change.

$$\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$$

$$\mathbf{A} \xrightarrow{R} \mathbf{M}_R \mathbf{A} = [\mathbf{M}_R \mathbf{a}_1, \dots, \mathbf{M}_R \mathbf{a}_n]$$

$$\mathbf{M}_{R^{-1}} = (\mathbf{M}_R)^{-1} \quad \text{therefore } \mathbf{M}_R \text{ is invertible}$$

taking some linear combination of the vectors in  $\mathbf{A}$

$$\lambda_1 \mathbf{a}_{i_1} + \dots + \lambda_t \mathbf{a}_{i_t} = \mathbf{0} \quad \Leftrightarrow$$

$$\mathbf{M}_R(\lambda_1 \mathbf{a}_{i_1} + \dots + \lambda_t \mathbf{a}_{i_t}) = \mathbf{0}$$

This is due to the fact that multiplying a set of linearly independent vectors by an invertible matrix preserves independence. From this it follows that the dimension of the columns of  $\mathbf{A}$  is equal to the dimension of the columns of  $\mathbf{M}_R \mathbf{A}$ , since  $\mathbf{M}_R$  is invertible. Therefore neither the row rank nor the column rank changes.

Since  $\mathbf{C}$  is in REF, the row rank is equal to the column rank, which is equal to the number of pivots, let it be  $r$ , then  $\text{rk}(\mathbf{A}) = r$ .



29th January 2020

## Properties of Symmetric Matrices

Take  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , such that  $\mathbf{A}^\top = \mathbf{A}$ .

All eigenvalues of  $\mathbf{A}$  are positive (non-negative) iff  $\forall \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\} [\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0 (\geq 0)]$ .  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$  is called a quadratic form. For example;

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} a & b \\ b & c \end{bmatrix} \\ \mathbf{x} &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ \mathbf{x}^\top \mathbf{A} \mathbf{x} &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= ax_1^2 + 2bx_1x_2 + cx_2^2 \\ &= k \\ \text{assume } x_2 &\neq 0 \\ \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{x_2^2} &= ay^2 + 2by + c \quad \text{where } y = \frac{x_1}{x_2} \\ &> 0 \quad \text{focusing on the positive case} \end{aligned}$$

If  $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$ , then we have an ellipse, with slanted major / minor axes. This uses a similar result to the result with the sphere's orthonormal basis being formed from the eigenvectors.

Considering this algebraically, if the quadratic in  $y$  is always greater than 0, then it must have no roots, hence the discriminant must be negative. Therefore we have the case where  $(2b)^2 - 4ac < 0$ , so  $b^2 - ac < 0$ , and  $a > 0$ .

Using the properties of the trace and the determinant, we can state the following;

$$\begin{aligned} \lambda_1 + \lambda_2 &= a + c \\ \lambda_1 \lambda_2 &= ac - b^2 \end{aligned}$$

We want both  $ac - b^2$  and  $a + b$  to be positive. The former is shown with the result for the discriminant, as if  $b^2 - ac$  is negative, then  $ac - b^2$  must be positive. The latter then follows as  $c$  is positive only if  $a$  is positive (due to  $ac > b^2 \geq 0$ ), which we have previously. This shows it holds for the simple case.

The proof for any  $n > 1$  is as follows;

$$\begin{aligned} \mathbf{S} &= \text{diag}(\lambda_1, \dots, \lambda_n) \\ &= \mathbf{V}^\top \mathbf{A} \mathbf{V} && \Rightarrow \\ \mathbf{A} &= \mathbf{V} \mathbf{S} \mathbf{V}^\top && \Rightarrow \\ \mathbf{x}^\top \mathbf{A} \mathbf{x} &= \mathbf{x}^\top \mathbf{V} \mathbf{S} \mathbf{V}^\top \mathbf{x} \\ &= \mathbf{z}^\top \mathbf{S} \mathbf{z} && \text{where } \mathbf{z} = \mathbf{V}^\top \mathbf{x} \text{ (which preserves } \ell_2\text{-norm)} \\ &= \sum_{i=1}^n \lambda_i z_i^2 \end{aligned}$$

To say  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$  is positive for any non-zero  $\mathbf{x}$  is to say that  $\sum_{i=1}^n \lambda_i z_i^2$  is positive for any non-zero  $\mathbf{v}$ .

Then obviously all the eigenvalues ( $\lambda_i$ ) must be positive iff  $\sum_{i=1}^n \lambda_i z_i^2$  is positive.

## Singular Value Decomposition (SVD)

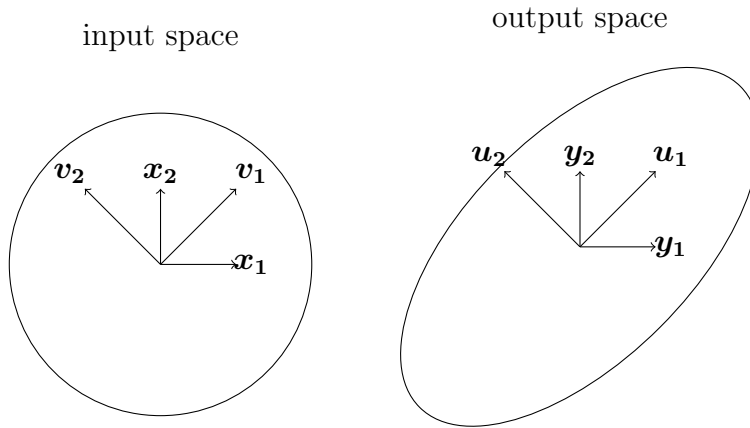
This works for a **general** matrix (not necessarily square). Take any matrix  $\mathbb{R}^{m \times n}$ .

$$\exists \mathbf{V} \in \mathbb{R}^{n \times n}, \mathbf{U} \in \mathbb{R}^{m \times m}, \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0$$

Where  $\mathbf{V}$ ,  $\mathbf{U}$  are orthogonal matrices,  $p = \min(m, n)$ , and  $r = \text{rk}(\mathbf{A})$ . We have a finitely decreasing sequence of  $\sigma$ . This allows us to write the following (assuming  $n \leq m$ );

$$\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top, \text{ where } \mathbf{S} = \begin{bmatrix} \sigma_1 & & & & & \\ & \sigma_2 & & & & \\ & & \ddots & & & \\ & & & \sigma_r & & \\ & & & & 0 & \\ & & & & & \ddots \\ 0 & & & & & & 0 \\ \vdots & & \dots & & & & \vdots \\ 0 & & \dots & & & & 0 \end{bmatrix}$$

This is the singular value decomposition of  $\mathbf{A}$ , and  $\sigma_1, \dots, \sigma_r$  are the singular values of  $\mathbf{A}$ .  $\sigma_1$  is the largest singular value of  $\mathbf{A}$ , which is the  $\ell_2$  matrix norm of  $\mathbf{A}$  ( $\|\mathbf{A}\|_2$ ). The geometric intuition of this, in two dimensions, is as follows;



Starting with our coordinates in the basis  $(\mathbf{x}_1, \mathbf{x}_2)$ , the matrix  $\mathbf{V}^\top$  performs a rotation to the new basis  $(\mathbf{v}_1, \mathbf{v}_2)$ .

Take some point  $\mathbf{x}$  on the unit circle with coordinates  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ , we obtain the new coordinates  $\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \mathbf{V}^\top \mathbf{x}$ .

We still maintain  $\|\mathbf{x}\| = 1$  and  $\|\mathbf{v}\| = 1$ .

We're now in the basis of  $(\mathbf{u}_1, \mathbf{u}_2)$ . We can say that  $u_1 = \sigma_1 v_1$ ,  $u_2 = \sigma_2 v_2$  (and so on, for higher dimensions). This implies the following result;

$$\frac{u_1^2}{\sigma_1^2} + \frac{u_2^2}{\sigma_2^2} = v_1^2 + v_2^2 = 1$$

Hence we obtain an ellipse.

### Example

$$\mathbf{A} \in \mathbb{R}^{3 \times 2} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$\text{rk}(\mathbf{A}) = 1$$

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$$

$$\mathbf{S} \in \mathbb{R}^{3 \times 2} = \text{diag}(\sigma_1, \sigma_2, \dots)$$

$$\begin{aligned} \mathbf{A}^\top &= (\mathbf{V}^\top)^\top \mathbf{S}^\top \mathbf{U}^\top \\ &= \mathbf{V} \mathbf{S}^\top \mathbf{U}^\top \end{aligned}$$

$$\begin{aligned} \mathbf{A}^\top \mathbf{A} &= \mathbf{V} \mathbf{S}^\top \mathbf{U}^\top \mathbf{U} \mathbf{S} \mathbf{V}^\top \\ &= \mathbf{V} \mathbf{S}^\top \mathbf{S} \mathbf{V}^\top \end{aligned}$$

$\mathbf{U}$  is an orthogonal matrix

brief note on  $\mathbf{S}^\top \mathbf{S} \in \mathbb{R}^{2 \times 2}$

$$\begin{aligned} \mathbf{S}^\top \mathbf{S} &= \text{diag}(\sigma_1, \sigma_2, \dots) \text{diag}(\sigma_1, \sigma_2, \dots) \\ &= \text{diag}(\sigma_1^2, \sigma_2^2, \dots) \end{aligned}$$

brief proof for symmetry of  $\mathbf{A}^\top \mathbf{A}$

$$\begin{aligned} (\mathbf{A}^\top \mathbf{A})^\top &= \mathbf{A}^\top (\mathbf{A}^\top)^\top \\ &= \mathbf{A}^\top \mathbf{A} \end{aligned}$$

$\mathbf{A}^\top \mathbf{A}$  is also positive semi-definite, meaning all its eigenvalues are non-negative ( $\geq 0$ ) since

$$\begin{aligned} \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} &= (\mathbf{A} \mathbf{x})^\top (\mathbf{A} \mathbf{x}) \\ &\geq 0 \end{aligned}$$

In this form, all we have to do is to find the spectral decomposition for  $\mathbf{A}^\top \mathbf{A}$ , since we have a symmetric matrix.

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} 3 & -3 \\ -3 & 3 \end{bmatrix}$$

$$|\mathbf{A}^\top \mathbf{A} - \lambda \mathbf{I}| = (3 - \lambda)^2 - 9$$

$$\lambda_1 = 6$$

$\Rightarrow$

$$\sigma_1 = \sqrt{6}$$

$$\mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\lambda_2 = 0$$

$\Rightarrow$

$$\sigma_2 = 0$$

$$\mathbf{v}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2]$$

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

$$\mathbf{A} \mathbf{V} = \mathbf{U} \mathbf{S}$$

$\Rightarrow$

$$[\mathbf{A} \mathbf{v}_1, \mathbf{A} \mathbf{v}_2] = \mathbf{U} \begin{bmatrix} \sqrt{6} & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$= [\sqrt{6} \mathbf{u}_1, \mathbf{0}]$$

$\Rightarrow$

$$\begin{aligned}
\mathbf{u}_1 &= \frac{1}{\sqrt{6}} \mathbf{A} \mathbf{v}_1 \\
&= \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \\
\mathbf{u}_2 &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} && \text{unit vector orthogonal to } \mathbf{u}_1 \\
\mathbf{u}_3 &= \mathbf{u}_1 \times \mathbf{u}_2 \\
&= \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ -1 \\ -2 \end{bmatrix} \\
\mathbf{A} &= [\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3] \begin{bmatrix} \sqrt{6} & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} [\mathbf{v}_1, \mathbf{v}_2]
\end{aligned}$$

### 30th January 2020

Note that this lecture has no audio recording for some of the first part. To recap; the goal of SVD is to find the following, such that  $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$

- $\mathbf{A} \in \mathbb{R}^{m \times n}$
- $\mathbf{U} \in \mathbb{R}^{m \times m}$   $\mathbf{U}^\top = \mathbf{U}$
- $\mathbf{V} \in \mathbb{R}^{n \times n}$   $\mathbf{V}^\top = \mathbf{V}$
- $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0$   $p = \min\{m, n\}, r = \text{rk}(\mathbf{A})$
- $\mathbf{S} \in \mathbb{R}^{m \times n} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r, 0, \dots, 0)$

### Example

While normally  $n \leq m$ , we have the following example for  $m < n$ , where we work with  $\mathbf{A} \mathbf{A}^\top \in \mathbb{R}^{m \times m}$ ;

$$\begin{aligned}
\mathbf{A} &= \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{bmatrix} \\
\mathbf{A}^\top &= \begin{bmatrix} 3 & 2 \\ 2 & 3 \\ 2 & -2 \end{bmatrix} \\
\mathbf{A} \mathbf{A}^\top &= \mathbf{U} \mathbf{S} \mathbf{S}^\top \mathbf{U}^\top \\
\mathbf{A} \mathbf{A}^\top &= \begin{bmatrix} 17 & 8 \\ 8 & 17 \end{bmatrix} \\
\lambda_1 + \lambda_2 &= \text{trace} \\
&= 34 \\
\lambda_1 \lambda_2 &= \text{determinant} \\
&= 225 \\
\lambda_1 &= 25 \\
\sigma_1 &= 5 \\
\lambda_2 &= 9 \\
\sigma_2 &= 3 \\
\mathbf{u}_2 &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\
\mathbf{u}_1 &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}
\end{aligned}$$

orthogonal to  $\mathbf{u}_2$

$$\begin{aligned}
U &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \\
A &= USV^\top \Rightarrow \\
A^\top &= VS^\top U^\top \Rightarrow \\
A^\top U &= VS^\top \Rightarrow \\
[A^\top u_1, A^\top u_2] &= [5v_1, 3v_2] \Rightarrow \\
v_1 &= \frac{1}{5} A^\top u_1 \\
&= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \\
v_2 &= \frac{1}{3} A^\top u_2 \\
&= \frac{1}{3\sqrt{2}} \begin{bmatrix} 1 \\ -1 \\ 4 \end{bmatrix} \\
v_3 &= v_1 \times v_2 \quad \text{orthogonal to both} \\
&= \frac{1}{3} \begin{bmatrix} 2 \\ -2 \\ -1 \end{bmatrix}
\end{aligned}$$

### General Strategy

$$\underbrace{A^\top A}_{\substack{\text{symmetric} \\ +\text{ve semi-definite}}} = VS^\top \underbrace{U^\top U}_I SV^\top = V \underbrace{S^\top S}_{\in \mathbb{R}^{n \times n}} V^\top$$

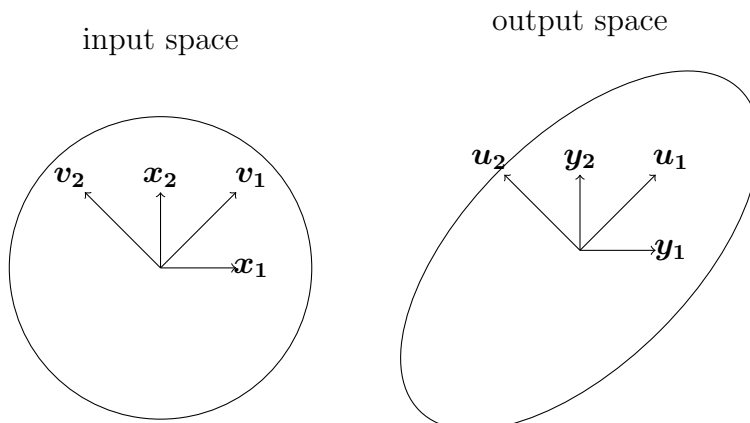
Solve for the spectral decomposition of  $A^\top A$  to get  $V$ , and  $\lambda_1 = \sigma_1^2, \lambda_2 = \sigma_2^2, \dots, \lambda_r = \sigma_r^2$ . From here we can do the following;

$$\begin{aligned}
A &= USV^\top \Rightarrow \\
AV &= US \Rightarrow \\
Av_i &= \sigma_i u_i \quad (1 \leq i \leq r), \Rightarrow \\
u_i &= \frac{Av_i}{\sigma_i}
\end{aligned}$$

For the remaining vectors,  $u_{r+1}, \dots, u_n$  find vectors to form an orthonormal basis. Use systems of linear equations in dimensions higher than 3, if the dimension is 3, just use the cross product. Solve  $u_{r+1}$  to  $u_n$  with

$$A^\top A u_j = \mathbf{0} \text{ for } r+1 \leq j \leq n$$

such that  $(u_1, \dots, u_r, u_{r+1}, \dots, u_n)$  forms an orthonormal basis.



Using the same example as before, because I'm too lazy to draw this again, we see that  $\mathbf{v}_i \mapsto \sigma_i \mathbf{v}_i = \mathbf{u}_i$ .  $\mathbf{v}_1$  is called the **principal component**, as it is the direction of the biggest change (since we have an ordering on  $\sigma$ ).

A brief note on why  $\mathbf{S}\mathbf{S}^\top$  gives a square matrix;

$$\begin{aligned}\mathbf{S} &= \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ 0 & \cdots & 0 & \end{bmatrix} \\ \mathbf{S}^\top &= \begin{bmatrix} \sigma_1 & & & 0 & \cdots & 0 \\ & \ddots & & \vdots & \ddots & \vdots \\ & & \sigma_r & 0 & \cdots & 0 \end{bmatrix} \\ \mathbf{S}^\top \mathbf{S} &= \mathbf{S}\mathbf{S}^\top && \text{by symmetry} \\ &= \begin{bmatrix} \sigma_1^2 & & & \\ & \ddots & & \\ & & \sigma_r^2 & \end{bmatrix}\end{aligned}$$

## Cholesky Factorisation / Decomposition of Symmetric Positive Definite Matrices

Take a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , with  $\mathbf{A} = \mathbf{A}^\top$  and  $\mathbf{A}$  is positive (semi-)definite. Then the following holds

$$\exists \mathbf{L} \in \mathbb{R}^{n \times n} [\mathbf{A} = \mathbf{L}\mathbf{L}^\top]$$

where  $\mathbf{L}$  is a lower triangular matrix, with all diagonal elements being positive (non-zero);

$$\mathbf{L} = \begin{bmatrix} l_{1,1} & 0 & \cdots & 0 \\ \times & l_{2,2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ \times & \cdots & \times & l_{n,n} \end{bmatrix}$$

This is a general result to the following, in the real numbers;  $a > 0 \Rightarrow \exists l > 0 [a = l^2]$ . This algorithm is done with brute force, by setting  $\mathbf{A} = (a_{i,j})$ ,  $\mathbf{L} = (l_{i,j})$ , and  $\mathbf{A} = \mathbf{L}\mathbf{L}^\top$ . The following is a simple example of this algorithm, with a  $3 \times 3$  matrix;

$$\begin{aligned}\mathbf{A} &= \begin{bmatrix} 1 & -1 & 1 \\ -1 & 10 & -1 \\ 1 & -1 & 5 \end{bmatrix} \\ &= \mathbf{L}\mathbf{L}^\top \\ \mathbf{L} &= \begin{bmatrix} l_{1,1} & 0 & 0 \\ l_{2,1} & l_{2,2} & 0 \\ l_{3,1} & l_{3,2} & l_{3,3} \end{bmatrix} \\ \mathbf{L}^\top &= \begin{bmatrix} l_{1,1} & l_{2,1} & l_{3,1} \\ 0 & l_{2,2} & l_{3,2} \\ 0 & 0 & l_{3,3} \end{bmatrix}\end{aligned}$$

by multiplying out  $\mathbf{L}\mathbf{L}^\top$ , we can solve per column, starting with the first

$$\begin{aligned}
l_{1,1}^2 &= 1 && \Rightarrow \\
l_{1,1} &= 1 && \text{taking positive square root} \\
l_{2,1}l_{1,1} &= -1 && \Rightarrow \\
l_{2,1} &= -1 \\
l_{3,1}l_{1,1} &= 1 && \Rightarrow \\
l_{3,1} &= 1
\end{aligned}$$

second column

$$\begin{aligned}
l_{2,1}^2 + l_{2,2}^2 &= 10 && \Rightarrow \\
l_{2,2} &= 3 \\
l_{3,1}l_{2,1} + l_{3,2}l_{2,2} &= -1 && \Rightarrow \\
-1 + 3l_{3,2} &= -1 && \Rightarrow \\
l_{3,2} &= 0
\end{aligned}$$

third (final) column

$$\begin{aligned}
l_{3,1}^2 + l_{3,2}^2 + l_{3,3}^2 &= 5 && \Rightarrow \\
1 + 0 + l_{3,3}^2 &= 5 && \Rightarrow \\
l_{3,3} &= 2
\end{aligned}$$

setting the matrix

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

This can be used to easily solve the following;

$$\begin{aligned}
\mathbf{A}\mathbf{x} &= \mathbf{b} && \Leftrightarrow \\
\mathbf{L}\mathbf{L}^\top \mathbf{x} &= \mathbf{b} && \Leftrightarrow \\
\mathbf{L}\mathbf{y} &= \mathbf{b} && \text{let } \mathbf{y} = \mathbf{L}^\top \mathbf{x}
\end{aligned}$$

This is trivial to solve with forward substitution ( $y_1$  can be solved easily, then substituted into the next equation to get  $y_2$  and so on) as we have a lower triangular matrix

$$\begin{aligned}
b_1 &= l_{1,1}y_1 + 0 + \dots + 0 \\
b_2 &= l_{2,1}y_1 + l_{2,2}y_2 + 0 + \dots + 0 \\
b_3 &= l_{3,1}y_1 + l_{3,2}y_2 + l_{3,3}y_3 + 0 + \dots + 0 \\
&\vdots \\
b_n &= l_{n,1}y_1 + \dots + l_{n,n}y_n
\end{aligned}$$

If  $\mathbf{L}$  is lower triangular, then  $\mathbf{L}^\top$  must be upper triangular, and thus can be solved with backwards substitution (such that  $x_n$  is used to solve  $x_{n-1}$ , and so on)

$$\begin{aligned}
\mathbf{U} &= \mathbf{L}^\top \\
y_n &= 0 + \dots + 0 + u_{n,n}x_n \\
y_{n-1} &= 0 + \dots + 0 + u_{n-1,n-1}x_{n-1} + u_{n-1,n}x_n \\
&\vdots \\
y_1 &= u_{1,1}x_1 + \dots + u_{1,n}x_n
\end{aligned}$$

## Rank-Nullity Theorem, Again

For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , the nullity of  $\mathbf{A}$  (dimension of the kernel of  $\mathbf{A}$ ) is  $n - \text{rk}(\mathbf{A})$ . To recap, a reduction to reduced row echelon form (RREF) is being in row echelon form, with the additional properties that all the pivots must be 1, and all other numbers in a column with a pivot are 0.

Applying the rank-nullity theorem to  $\mathbf{A}^\top$ , where  $\mathbf{A}^\top \in \mathbb{R}^{n \times m}$ ;

$$\begin{aligned} \text{nullity}(\mathbf{A}^\top) + \text{rk}(\mathbf{A}^\top) &= m & \Rightarrow \\ \dim(\ker(\mathbf{A}^\top)) + \text{rk}(\mathbf{A}) &= m & \text{rk}(\mathbf{A}) = \text{rk}(\mathbf{A}^\top), \Rightarrow \\ \dim(\ker(\mathbf{A}^\top)) + \dim(\text{Im}(\mathbf{A})) &= m & \text{rk}(\mathbf{A}) = \dim(\text{Im}(\mathbf{A})) \end{aligned}$$

It's important to now note that these are both subspaces of  $\mathbb{R}^m$

$$\begin{aligned} \text{Im}(\mathbf{A}) &\subseteq \mathbb{R}^m & \text{output space of } \mathbf{A} \\ \ker(\mathbf{A}^\top) &\subseteq \mathbb{R}^m \end{aligned}$$

Working through an example matrix, we can observe the following;

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 1 & -1 \end{bmatrix} \\ \mathbf{A}^\top &= \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & -1 \end{bmatrix} \\ \text{Im}(\mathbf{A}) &= \left\{ x_1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + x_2 \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} : x_1, x_2 \in \mathbb{R} \right\} \end{aligned}$$

geometrically,  $\text{Im}(\mathbf{A})$  is a plane going through the origin - intuitively, we can represent this with a normal (take cross product)

$$\begin{aligned} &= \left\{ \mathbf{v} : \mathbf{v} \cdot \begin{bmatrix} -1 \\ 2 \\ -1 \end{bmatrix} = 0 \right\} \\ \ker(\mathbf{A}^\top) &= \{ \mathbf{w} : \mathbf{A}^\top \mathbf{w} = \mathbf{0} \} \\ \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} &= \mathbf{0} & \Rightarrow \\ \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \cdot \mathbf{w} &= 0 \\ \text{and } \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \cdot \mathbf{w} &= 0 \end{aligned}$$

therefore  $\mathbf{w}$  is orthogonal to both those matrices, hence we can take the cross product, which we already know

$$\mathbf{w} = k \begin{bmatrix} -1 \\ 2 \\ -1 \end{bmatrix} \quad k \in \mathbb{R}$$



We can then conclude the following;

$$\mathbf{w} \in \ker(\mathbf{A}^\top) \Leftrightarrow \mathbf{w} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 0 = \mathbf{w} \cdot \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \Leftrightarrow \forall \mathbf{y} \in \text{Im}(\mathbf{A}) [\mathbf{w} \cdot \mathbf{y} = 0]$$

This gives the result that  $\mathbf{w}$  is in the null space of  $\mathbf{A}^\top$  iff it is orthogonal to every vector in the image space of  $\mathbf{A}$ .

Our conjecture is the following;  $\ker(\mathbf{A}^\top) \perp \text{Im}(\mathbf{A})$ , which is the same as saying

$$\mathbf{v} \in \ker(\mathbf{A}^\top) \Leftrightarrow \forall \mathbf{y} \in \text{Im}(\mathbf{A}) [\mathbf{y} \cdot \mathbf{v} = 0]$$

This can be proven generally as follows;

$$\begin{aligned} & \mathbf{v} \in \ker(\mathbf{A}^\top) \\ \Leftrightarrow & \mathbf{A}^\top \mathbf{v} = \mathbf{0} && \text{by definition of the null space / kernel} \\ \Leftrightarrow & \forall i \in [1, n] [\mathbf{a}_i \cdot \mathbf{v} = 0] && (1) \\ \Leftrightarrow & \forall \mathbf{x} \in \mathbb{R}^n \left[ \left( \sum_{i=1}^n x_i \mathbf{a}_i \right) \cdot \mathbf{v} = 0 \right] && \text{generalise to any linear combination} \\ \Leftrightarrow & \mathbf{y} \in \text{Im}(\mathbf{A}) [\mathbf{y} \cdot \mathbf{v} = 0] && \blacksquare \end{aligned}$$

(1) row  $i$  of  $\mathbf{A}^\top \in \mathbb{R}^{n \times m}$  is the same as column  $i$  of  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , transposed, and all of the components must become 0

We say that  $\mathbb{R}^m$  is the direct sum of  $\ker(\mathbf{A}^\top)$  and  $\text{Im}(\mathbf{A})$ . We can get the following result, following what was proven earlier;

$$\mathbf{v} \in \text{Im}(\mathbf{A}) \wedge \mathbf{v} \in \ker(\mathbf{A}^\top) \Leftrightarrow \mathbf{v} \cdot \mathbf{v} = 0 \Leftrightarrow \mathbf{v} = \mathbf{0}$$

We want to prove the result for any vector  $\mathbf{x} \in \mathbb{R}^m$ , there exists a unique vector  $\mathbf{x}_r \in \text{Im}(\mathbf{A})$  and  $\mathbf{x}_n \in \ker(\mathbf{A}^\top)$ , such that  $\mathbf{x} = \mathbf{x}_r + \mathbf{x}_n$ , or formally;

$$\forall \mathbf{x} \in \mathbb{R}^m \exists! \mathbf{x}_r \in \text{Im}(\mathbf{A}), \mathbf{x}_n \in \ker(\mathbf{A}^\top) [\mathbf{x} = \mathbf{x}_r + \mathbf{x}_n]$$

By taking the union of a basis of  $\text{Im}(\mathbf{A})$ , let it be  $(\mathbf{v}_1, \dots, \mathbf{v}_r)$ , and a basis of  $\ker(\mathbf{A}^\top)$ , let it be  $(\mathbf{v}_{r+1}, \dots, \mathbf{v}_m)$ , we get a basis of  $\mathbb{R}^m$ :  $(\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_m)$ . This gives us the following result;

$$\mathbf{x} \in \mathbb{R}^m = \sum_{i=1}^m x_i \mathbf{v}_i = \underbrace{\sum_{i=1}^r x_i \mathbf{v}_i}_{\in \text{Im}(\mathbf{A})} + \underbrace{\sum_{i=r+1}^m x_i \mathbf{v}_i}_{\in \ker(\mathbf{A}^\top)} = \mathbf{x}_r + \mathbf{x}_n$$

In order to prove uniqueness, we assume  $\mathbf{x} = \mathbf{x}_r + \mathbf{x}_n = \mathbf{x}'_r + \mathbf{x}'_n$ . The **violet** part only holds iff  $\mathbf{x}_r - \mathbf{x}'_r = \mathbf{x}'_n - \mathbf{x}_n$ . As they are both in their respective subspaces of  $\text{Im}(\mathbf{A})$  and  $\ker(\mathbf{A}^\top)$ , they must both be  $\mathbf{0}$ , as shown above. This gives the result  $\mathbf{x}_r = \mathbf{x}'_r \wedge \mathbf{x}_n = \mathbf{x}'_n$ , hence proving uniqueness.

## 5th February 2020

### Positive Definiteness

The conditions for positive definiteness ( $\forall \mathbf{x} \neq \mathbf{0} [\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0]$ ) of  $\mathbf{A}$  are as follows ;

- all diagonal elements of  $\mathbf{A}$  must be  $> 0$

Take  $\mathbf{x} = \mathbf{e}_i$ , which is all 0, other than the  $i^{\text{th}}$  entry being 1. Then  $\mathbf{e}_i^\top \mathbf{A} \mathbf{e}_i = a_{i,i}$ , which has to be greater than 0.

- all principal minors of  $\mathbf{A}$  must be positive definite.

For  $1 \leq k \leq n$ ,  $\mathbf{A}^{(k)} \in \mathbb{R}^{k \times k}$  denotes the top left square submatrix of  $\mathbf{A}$ .

$$\begin{aligned} \mathbf{x} &= \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} & \text{where } \mathbf{y} \in \mathbb{R}^k, \mathbf{0} \in \mathbb{R}^{n-k} \\ \mathbf{x}^\top \mathbf{A} \mathbf{x} &= \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}^\top \mathbf{A} \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} \\ &= \mathbf{y}^\top \mathbf{A}^{(k)} \mathbf{y} \end{aligned}$$

So for any  $\mathbf{y} \neq \mathbf{0}$ ,  $\mathbf{x} \neq \mathbf{0}$ , thus  $\mathbf{y}^\top \mathbf{A}^{(k)} \mathbf{y} > 0 \Leftrightarrow \mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$

- $|a_{i,j}| < \max\{a_{i,i}, a_{j,j}\}$

For any entry not on the diagonal, there are two corresponding entries on the diagonal (one corresponding to the row, the other to the column);

$$\mathbf{A} = \begin{bmatrix} \ddots & & & & \\ & a_{i,i} & \cdots & a_{i,j} & \\ & & \ddots & \vdots & \\ & & & a_{j,j} & \\ & & & & \ddots \end{bmatrix}$$

take  $\mathbf{x} = \mathbf{e}_i + \mathbf{e}_j$

$$\begin{aligned} \mathbf{x}^\top \mathbf{A} \mathbf{x} &= (\mathbf{e}_i + \mathbf{e}_j)^\top \mathbf{A} (\mathbf{e}_i + \mathbf{e}_j) \\ &= (\mathbf{e}_i^\top + \mathbf{e}_j^\top) \mathbf{A} (\mathbf{e}_i + \mathbf{e}_j) \\ &= \mathbf{e}_i^\top \mathbf{A} \mathbf{e}_i + \mathbf{e}_j^\top \mathbf{A} \mathbf{e}_j + \mathbf{e}_i^\top \mathbf{A} \mathbf{e}_j + \mathbf{e}_j^\top \mathbf{A} \mathbf{e}_i \\ &= a_{i,i} + a_{j,j} + 2a_{i,j} \end{aligned}$$

$\mathbf{A}$  is symmetric

$\Rightarrow$

$$-a_{i,j} < \frac{a_{i,i} + a_{j,j}}{2}$$

diagonal terms are positive

$$\text{or } a_{i,j} < \frac{a_{i,i} + a_{j,j}}{2}$$

$\Rightarrow$

$$|a_{i,j}| < \max\{a_{i,i}, a_{j,j}\}$$

- all positive eigenvalues

For an example, take the following;

$$\mathbf{A} = \begin{bmatrix} 1 & -1 & 3 \\ -1 & 10 & -1 \\ -3 & -1 & 2 \end{bmatrix}$$

However, we see that  $|-3| = 3 > 2 = \max\{1, 2\}$ , hence  $\mathbf{A}$  is not positive definite

## Least Square Method

For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , where  $n$  is commonly the dimension of the data, and  $m$  the number of samples (such that each row is an entry), we can write  $\mathbf{A}$  as;

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}^1 \\ \mathbf{a}^2 \\ \vdots \\ \mathbf{a}^m \end{bmatrix} \quad (\text{all rows})$$

The goal of the LSM is to find a  $\mathbf{x} \in \mathbb{R}^n$ , such that  $\|\mathbf{Ax} - \mathbf{b}\|_2$  is minimised (which is equivalent to minimising  $\|\mathbf{Ax} - \mathbf{b}\|_2^2$ , as it is positive). From previous results, we can take  $\mathbf{b} = \mathbf{b}_r + \mathbf{b}_n$ , therefore we can minimise the following;

$$\begin{aligned}\|\mathbf{Ax} - \mathbf{b}_r - \mathbf{b}_n\|_2^2 &= ((\mathbf{Ax} - \mathbf{b}_r) - \mathbf{b}_n) \cdot ((\mathbf{Ax} - \mathbf{b}_r) - \mathbf{b}_n) \\ &= (\mathbf{Ax} - \mathbf{b}_r) \cdot (\mathbf{Ax} - \mathbf{b}_r) + \mathbf{b}_n \cdot \mathbf{b}_n - 2 \overbrace{(\mathbf{Ax} - \mathbf{b}_r) \cdot \mathbf{b}_n}^{\substack{\in \text{Im}(\mathbf{A}) \\ \in \ker(\mathbf{A}^\top)}}\end{aligned}\quad (1)$$

$$= (\mathbf{Ax} - \mathbf{b}_r) \cdot (\mathbf{Ax} - \mathbf{b}_r) + \mathbf{b}_n \cdot \mathbf{b}_n \quad (2)$$

(1) the dot product is bilinear in both variables, hence we can expand out

(2) range space of  $\mathbf{A}$  is perpendicular to the null space of  $\mathbf{A}^\top$ , hence the dot product is 0

Therefore, we want to minimise  $(\mathbf{Ax} - \mathbf{b}_r) \cdot (\mathbf{Ax} - \mathbf{b}_r) + \mathbf{b}_n \cdot \mathbf{b}_n$ . However, we have no control over  $\mathbf{b}_n$ , therefore we want to solve for  $\mathbf{Ax} = \mathbf{b}_r$ . This will always have a solution, as  $\mathbf{b}_r \in \text{Im}(\mathbf{A})$ , which means that such a  $\mathbf{x}$  exists by definition of the range space. We then claim the following;

$$\mathbf{Ax} = \mathbf{b}_r \Leftrightarrow \mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b}$$

Therefore, if the left hand side of the double implication has a solution, which we know it does, then the right hand side must also have a solution. The proof for this is done in both directions;

proving " $\Rightarrow$ "

$$\begin{aligned}\text{suppose } \mathbf{Ax} &= \mathbf{b}_r && \Rightarrow \\ \mathbf{A}^\top \mathbf{Ax} &= \mathbf{A}^\top \mathbf{b}_r && \text{pre-multiply} \\ &= \mathbf{A}^\top (\mathbf{b}_r + \mathbf{b}_n) && \text{in null space, hence } \mathbf{A}^\top \mathbf{b}_n = 0 \\ &= \mathbf{A}^\top \mathbf{b} && \blacksquare\end{aligned}$$

proving " $\Leftarrow$ "

$$\begin{aligned}\text{suppose } \mathbf{A}^\top \mathbf{Ax} &= \mathbf{A}^\top \mathbf{b} && \Rightarrow \\ \mathbf{A}^\top \mathbf{Ax} - \mathbf{A}^\top (\mathbf{b}_r + \mathbf{b}_n) &= 0 && \Rightarrow \\ \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}_r - \mathbf{b}_n) &= 0 && \Rightarrow \\ \mathbf{Ax} - \mathbf{b}_r - \mathbf{b}_n &\in \ker(\mathbf{A}^\top) && \text{property of adding in subspace, } \Rightarrow \\ \mathbf{Ax} - \mathbf{b}_r &\in \ker(\mathbf{A}^\top) && \text{and} \\ \mathbf{Ax} - \mathbf{b}_r &\in \text{Im}(\mathbf{A}^\top) && \text{adding two elements of the range space, } \Rightarrow \\ \mathbf{Ax} - \mathbf{b}_r &= 0 && \Rightarrow \\ \mathbf{Ax} &= \mathbf{b}_r && \blacksquare\end{aligned}$$

An example of using this is as follows (note that  $\mathbf{Ax} = \mathbf{b}$  is inconsistent);

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

$$\mathbf{b} = \begin{bmatrix} 0 \\ 5 \\ -1 \end{bmatrix}$$

find  $\mathbf{x}$  such that  $\mathbf{Ax} = \mathbf{b}_r$

$$\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b} \quad \Rightarrow$$

$$\begin{bmatrix} 9 & 6 \\ 6 & 8 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 2 \\ 2 & 2 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 5 \\ -1 \end{bmatrix}$$

$$= \begin{bmatrix} 3 \\ 10 \end{bmatrix}$$

we know this is positive definite, and therefore there is a solution

$$x_1 = -1$$

$$x_2 = 2$$

## 6th February 2020

### Using Cholesky Decomposition for LSM

While it's easy in some cases to solve with an inverse, we have a symmetric positive definite matrix in  $\mathbf{A}^\top \mathbf{A}$ . Therefore we have the following steps;

- goal is to solve  $\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b}$
- let  $\mathbf{A}^\top \mathbf{A} = \mathbf{LL}^\top$ , therefore we have  $\mathbf{LL}^\top \mathbf{x} = \mathbf{A}^\top \mathbf{b}$
- let  $\mathbf{y} = \mathbf{L}^\top \mathbf{x}$
- find  $\mathbf{y}$  with forward substitution in  $\mathbf{Ly} = \mathbf{A}^\top \mathbf{b}$
- find  $\mathbf{x}$  with backward substitution in  $\mathbf{y} = \mathbf{L}^\top \mathbf{x}$

### First Application (Linear Regression) + Example

A simple application for this is to consider  $m$  samples of pairs  $(a_i, b_i)$ . We want to find an affine line  $y = s_1x + s_0$ , such that;

$$\sum_{i=1}^m e_i^2 \text{ is minimised, where the error term } e_i = |(s_1 a_i + s_0) - b_i|$$

In this case, we only have two unknowns, and therefore we can write the following;

$$\mathbf{A} = \begin{bmatrix} 1 & a_1 \\ 1 & a_2 \\ \vdots & \vdots \\ 1 & a_m \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

$$\mathbf{z} = \begin{bmatrix} s_0 \\ s_1 \end{bmatrix}$$

$$\mathbf{Az} = \mathbf{b} \quad \Rightarrow$$

$$\begin{bmatrix} 1 & a_1 \\ 1 & a_2 \\ \vdots & \vdots \\ 1 & a_m \end{bmatrix} \begin{bmatrix} s_0 \\ s_1 \end{bmatrix} = \begin{bmatrix} s_0 + a_1 s_1 \\ s_0 + a_2 s_1 \\ \vdots \\ s_0 + a_m s_1 \end{bmatrix} \\ = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

However, being able to solve it exactly is unlikely, especially for large datasets. In a real example, we have the following;

$$(a_1, b_1) = (1, 6), (2, 5), (3, 7), (4, 10)$$

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} 6 \\ 5 \\ 7 \\ 10 \end{bmatrix}$$

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}$$

$$= \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

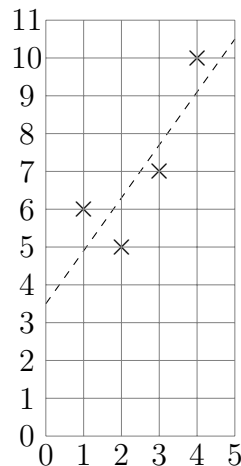
$$\mathbf{A}^\top \mathbf{b} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 6 \\ 5 \\ 7 \\ 10 \end{bmatrix}$$

$$= \begin{bmatrix} 28 \\ 77 \end{bmatrix}$$

by doing Gaussian elimination, we get the following results

$$s_0 = \frac{7}{2}$$

$$s_1 = \frac{7}{5}$$



## Higher Dimensions

In order to generalise this to higher dimensions, for each  $(a_i, b_i)$ , we now have a vector  $\mathbf{a}_i \in \mathbb{R}^n$ , and a scalar  $b_i \in \mathbb{R}$ , for  $i = 1, \dots, m$ . The affine line is now no longer as simple, and is therefore written as;

$$y = \mathbf{x}^\top \mathbf{s} + s_0 = \mathbf{x} \cdot \mathbf{s} + s_0 = s_0 + \sum_{j=1}^n x_j s_j, \text{ where } \mathbf{s} \in \mathbb{R}^n, \text{ and } s_0 \in \mathbb{R}.$$

Setting this out, we have the following;

$$\begin{aligned} \mathbf{a}_i &= \begin{bmatrix} a_{i,1} \\ \vdots \\ a_{i,n} \end{bmatrix} && \text{for } i = 1, \dots, m \\ \mathbf{A} &= \begin{bmatrix} 1 & a_{1,1} & \cdots & a_{1,n} \\ 1 & a_{2,1} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & a_{m,1} & \cdots & a_{m,n} \end{bmatrix} \\ \mathbf{b} &= \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \\ \mathbf{z} &= \begin{bmatrix} s_0 \\ \mathbf{s} \end{bmatrix} && \mathbf{s} \in \mathbb{R}^n \end{aligned}$$

## Second Application (Polynomial Regression)

Another application of this is to suppose data is  $(y_i, t_i)$ , where  $y_i \in \mathbb{R}$  and  $t_i \in \mathbb{R}$  (considered as "time" - but is any continuous variable). We hypothesise the following;

$$y(t) = \sum_{j=1}^n P_j f_j(t)$$

where  $P_j$ s are the parameters of the model, and  $f_j$ s are basic functions, for **example**;  $f_1(t) = 1$ ,  $f_2(t) = t$ ,  $\dots$ ,  $f_n(t) = t^{n-1}$ . This is called **polynomial regression**. We know (we're given)  $f_j$ s, but we want to find the  $P_j$ s.

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{y} \\ \mathbf{A} &= (a_{i,j}) \\ a_{i,j} &= f_j(t_i) \\ \mathbf{x} &= \begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_m \end{bmatrix} \end{aligned}$$

An example of this is to determine height ( $h$ ), and the coefficient of gravity ( $g$ ). We hypothesise that the distance with respect to time, away from the start, is

$$y(t) = h - g \frac{t^2}{2}$$

We assume that  $h$  and  $g$  are the unknown parameters, and define the basic functions as follows (note that the pairs  $(y_i, t_i)$  can be obtained with experiments);

$$\begin{aligned} f_1(t) &= 1 \\ f_2(t) &= -\frac{t^2}{2} \\ \mathbf{A} &= \begin{bmatrix} 1 & -\frac{t_1^2}{2} \\ \vdots & \vdots \\ 1 & -\frac{t_m^2}{2} \end{bmatrix} \\ \mathbf{x} &= \begin{bmatrix} h \\ g \end{bmatrix} \\ \mathbf{y} &= \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \end{aligned}$$

This can now be solved with the standard LSM.

## QR Decomposition

We define  $\mathbf{Q} \in \mathbb{R}^{m \times n}$  to be an orthogonal matrix if  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_n]$ , the column vectors are unit vectors with respect to the  $\ell_2$ -norm, and are pairwise orthogonal. This is only possible if  $n \leq m$ . It's important to note that  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_n$ , and  $\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}_m$ .

The goal is to find  $\mathbf{Q}$  and  $\mathbf{R}$  such that  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q}$  is an orthogonal matrix, and  $\mathbf{R}$  is an upper triangular matrix. It's trivial to see that  $\mathbf{R} = \mathbf{Q}^\top \mathbf{A}$ , due to the previously mentioned property. This is used to easily solve the following (since  $\mathbf{R}$  is upper triangular, we can employ backward substitution);

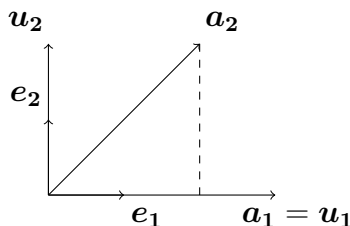
$$\mathbf{A}\mathbf{x} = \mathbf{b} \Rightarrow \mathbf{Q}\mathbf{R}\mathbf{x} = \mathbf{b} \Rightarrow \mathbf{R}\mathbf{x} = \mathbf{Q}^\top \mathbf{b}$$

## Gram-Schmidt Process

Suppose we have a set of linearly independent vectors  $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^m$ . The Gram-Schmidt process is as follows;

$$\begin{aligned} \mathbf{u}_1 &= \mathbf{a}_1 & \mathbf{e}_1 &= \frac{\mathbf{u}_1}{\|\mathbf{u}_1\|} \\ \mathbf{u}_2 &= \mathbf{a}_2 - (\mathbf{e}_1 \cdot \mathbf{a}_2)\mathbf{e}_1 & \mathbf{e}_2 &= \frac{\mathbf{u}_2}{\|\mathbf{u}_2\|} \\ \mathbf{u}_3 &= \mathbf{a}_3 - (\mathbf{e}_1 \cdot \mathbf{a}_3)\mathbf{e}_1 - (\mathbf{e}_2 \cdot \mathbf{a}_3)\mathbf{e}_2 & \mathbf{e}_3 &= \frac{\mathbf{u}_3}{\|\mathbf{u}_3\|} \\ \vdots & & \vdots & \\ \mathbf{u}_n &= \mathbf{a}_n - \sum_{j=1}^{n-1} (\mathbf{e}_j \cdot \mathbf{a}_n)\mathbf{e}_j & \mathbf{e}_n &= \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|} \end{aligned}$$

Geometrically, we can visualise it as the following;



In general, for  $1 \leq j \leq n$ , we have

$$\mathbf{a}_j = (\mathbf{e}_1 \cdot \mathbf{a}_j)\mathbf{e}_1 + (\mathbf{e}_2 \cdot \mathbf{a}_j)\mathbf{e}_2 + \dots + (\mathbf{e}_j \cdot \mathbf{a}_j)\mathbf{e}_j$$

## First Method for QR Decomposition (Gram-Schmidt Process)

Take a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , where  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ , assuming that  $\mathbf{a}_1, \dots, \mathbf{a}_n$  are linearly independent. We then apply Gram-Schmidt to  $\mathbf{a}_1, \dots, \mathbf{a}_n$  as described above to get  $\mathbf{Q}_0 = (\mathbf{e}_1, \dots, \mathbf{e}_n)$ , where. Note that for  $1 \leq j \leq n$ ,  $\mathbf{e}_j \cdot \mathbf{e}_j = 1$ , and for  $j \neq k$ ,  $\mathbf{e}_j \cdot \mathbf{e}_k = 0$ . Let  $\mathbf{Q} = [\mathbf{e}_1, \dots, \mathbf{e}_n]$ , which is orthogonal. While we can see that  $\mathbf{R} = \mathbf{Q}^\top \mathbf{A}$ , we can write this as;

$$\mathbf{R} = \begin{bmatrix} \mathbf{e}_1 \cdot \mathbf{a}_1 & \mathbf{e}_1 \cdot \mathbf{a}_2 & \mathbf{e}_1 \cdot \mathbf{a}_3 & \cdots & \mathbf{e}_1 \cdot \mathbf{a}_n \\ & \mathbf{e}_2 \cdot \mathbf{a}_2 & \mathbf{e}_2 \cdot \mathbf{a}_3 & \cdots & \mathbf{e}_2 \cdot \mathbf{a}_n \\ & & \mathbf{e}_3 \cdot \mathbf{a}_3 & \cdots & \mathbf{e}_3 \cdot \mathbf{a}_n \\ & & & \ddots & \vdots \\ & & & & \mathbf{e}_n \cdot \mathbf{a}_n \end{bmatrix} \in \mathbb{R}^{n \times n}$$

We claim that the diagram below is valid when  $\mathbf{A}$  represents the linear map  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , in the standard basis.

$$\begin{array}{ccc} \mathbb{R}_{\mathbf{E}_n}^n & \xrightarrow{f_{\mathbf{E}_m \mathbf{E}_n} = \mathbf{A}} & \mathbb{R}_{\mathbf{E}_m}^m \\ & \searrow f_{\mathbf{Q}_0 \mathbf{E}_n} = \mathbf{R} & \uparrow \text{Id}_{\mathbf{E}_m \mathbf{Q}_0} = \mathbf{Q} \\ & & \text{Im}(\mathbf{A})_{\mathbf{Q}_0} \end{array}$$

12th February 2020

### Proof of Diagram

Continuing on from last week's lecture, we aim to prove two claims in the diagram;

- $f_{\mathbf{Q}_0 \mathbf{E}_n} = \mathbf{R}$

The matrix  $\mathbf{R}$  represents the linear map  $f$  from  $\mathbb{R}^n$  (in the standard basis) to  $\mathbb{R}^m$  in the new basis ( $\mathbf{Q}_0$ ).

To begin, we observe how the standard basis is mapped under  $\mathbf{A}$ ;

$$j^{\text{th}} \rightarrow \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \xrightarrow{\mathbf{A}} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{bmatrix} = \mathbf{a}_j \rightarrow \mathbf{a}_j = (\mathbf{e}_1 \cdot \mathbf{a}_j)\mathbf{e}_1 + (\mathbf{e}_2 \cdot \mathbf{a}_j)\mathbf{e}_2 + \cdots + (\mathbf{e}_j \cdot \mathbf{a}_j)\mathbf{e}_j + 0\mathbf{e}_{j+1} + \cdots + 0\mathbf{e}_n$$

Looking at these coefficients with respect to the basis formed by  $(\mathbf{e}_1, \dots, \mathbf{e}_n)$ , which are highlighted in **violet**, it's clear that it forms the  $j^{\text{th}}$  column of  $\mathbf{R}$ .

- $\text{Id}_{\mathbf{E}_m \mathbf{Q}_0} = \mathbf{Q}$

This is a change of basis in the image space of  $\mathbf{A}$ .

This is trivial to show, as we want to look at what happens to the  $j^{\text{th}}$  vector in the new basis  $\mathbf{Q}_0$ . To write each of these matrices in the new (standard) basis, we simply do the following;



$$\mathbf{e}_j = e_{1,j} \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + e_{2,j} \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \cdots + e_{m,j} \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

Trivially, we get  $\mathbf{Q} = [\mathbf{e}_1, \dots, \mathbf{e}_n]$ , as expected.

### Note on Basis Change

To find the matrix representing  $f_{\mathbf{D}\mathbf{B}}$ , which takes a vector in basis  $\mathbf{B}$  to a vector in basis  $\mathbf{D}$ , we look at what happens to each basis vector.

$$\mathbf{b}_j \mapsto f(\mathbf{b}_j) = \delta_1 \mathbf{d}_1 + \delta_2 \mathbf{d}_2 + \cdots + \delta_m \mathbf{d}_m$$

Then the  $j^{\text{th}}$  column of the matrix is given as;

$$(f_{\mathbf{D}\mathbf{B}})_j = \begin{bmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_m \end{bmatrix}$$

### Householder Operator

I can't draw the geometric intuition nicely, at all, look at Panopto timestamp 36:37. However, imagine there is a plane, and there is some point  $\mathbf{x}$  that exists. That point is orthogonally projected onto the plane, with the orthogonal projection matrix  $\mathbf{P}$ , such that the vector between the point  $\mathbf{x}$  and  $\mathbf{P}\mathbf{x}$  is perpendicular to the plane. Now, take the perpendicular unit vector for the plane as  $\mathbf{u}$ , therefore the length of  $\mathbf{x}$  to  $\mathbf{P}\mathbf{x}$  is  $\mathbf{u} \cdot \mathbf{x}$ . Imagine the plane as a mirror, and the image of  $\mathbf{x}$  is on the other side of the plane, let that point be  $H(\mathbf{x}) = \mathbf{H}\mathbf{x}$ . Therefore, we can work out  $\mathbf{H}$  as follows;

$$H(\mathbf{x}) = \mathbf{x} - 2(\mathbf{u} \cdot \mathbf{x})\mathbf{u} = \mathbf{x} - 2\mathbf{x}\mathbf{u}\mathbf{u}^\top = \underbrace{(\mathbf{I} - 2\mathbf{u}\mathbf{u}^\top)}_{\mathbf{H}} \mathbf{x}$$

We expect  $\mathbf{H}^2$  to give us the identity, as if "reflect" a reflected point, it should give the original;

$$\begin{aligned} \mathbf{H}^2 &= (\mathbf{I} - 2\mathbf{u}\mathbf{u}^\top) \\ &= \mathbf{I} - 4\mathbf{u}\mathbf{u}^\top + 4\mathbf{u}(\mathbf{u}^\top \mathbf{u})\mathbf{u}^\top \\ &= \mathbf{I} - 4\mathbf{u}\mathbf{u}^\top + 4\mathbf{u}\mathbf{u}^\top && \ell_2\text{-norm of unit vector is 1} \\ &= \mathbf{I} \end{aligned}$$

■

Note that there are two distinct eigenvalues of  $\mathbf{H}$ ;  $\lambda_1 = 1$  has multiplicity  $m_1 = m - 1$ , where  $m$  is the dimension of the space. This is for each of the vectors  $\mathbf{v}_i$  that span the plane, as they are mapped to themselves ( $\mathbf{H}\mathbf{v}_i = H(\mathbf{v}_i) = \mathbf{v}_i$ ). There is also an additional eigenvalue  $\lambda_2 = -1$ , for  $\mathbf{u}$ , as we can see that  $\mathbf{H}\mathbf{u} = H(\mathbf{u}) = -\mathbf{u}$ .

### Second Method for QR Decomposition (Householder Operator)

Take the matrix  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ . The Householder operator must preserve the length, in the following mapping;

$$H \left( \begin{bmatrix} a_{1,1} \\ a_{2,1} \\ \vdots \\ a_{m,1} \end{bmatrix} \right) = \begin{bmatrix} \lambda \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\pm \lambda = \|\mathbf{a}_1\|$$

We want to send the vector  $\mathbf{a}_1$  to the vector  $\|\mathbf{a}_1\|\mathbf{e}_1$ , where  $\mathbf{e}_1$  is all 0, other than the first element.

$$\begin{aligned} \mathbf{u} &= \frac{\mathbf{a}_1 - \|\mathbf{a}_1\|\mathbf{e}_1}{\|\mathbf{a}_1 - \|\mathbf{a}_1\|\mathbf{e}_1\|} \\ \mathbf{H}_1 &= \mathbf{I} - 2\mathbf{u}\mathbf{u}^\top \\ \mathbf{H}_1\mathbf{A} &= \mathbf{H}_1[\mathbf{a}_1, \dots, \mathbf{a}_n] \\ &= [\mathbf{H}_1\mathbf{a}_1, \dots, \mathbf{H}_1\mathbf{a}_n] \\ &= \begin{bmatrix} \|\mathbf{a}_1\| & \times \\ 0 & \mathbf{A}_2 \end{bmatrix} \quad \times \text{ is an unknown row, } \mathbf{A}_2 \text{ is the bottom right submatrix} \end{aligned}$$

Applying  $\mathbf{H}$  to  $\mathbf{A}$  gets rid of all entries in the first column, other than the first element. This is an iterative process, where every iteration brings us closer to an upper triangular matrix. Note that the Householder matrix,  $\mathbf{H}_2$ , applied to  $\mathbf{A}_2$  is 1 smaller on each dimension. To account for this, a top row, and left column is added, of all 0s, other than the top left element, which is 1.

$$\mathbf{H}_n \cdots \mathbf{H}_1 \mathbf{A} = \mathbf{R}$$

By reversing this, we can get the following result;

$$\mathbf{A} = \mathbf{H}_1 \cdots \mathbf{H}_n \mathbf{R}$$

Using the Householder operator to do QR decomposition isn't examinable.

## 13th February 2020

In the second half of the course, we cover;

- condition number
- convergence and fixed point problems
- iterative solution of linear equations and eigenvectors
- Laplace and Fourier transforms
- functions of several variables and optimisation

### Condition

When a matrix is very close to singular, the results can be unreliable when computed by a computer, as the results can push up to the precision of the machine. Consider the following set of linear equations;

$$\begin{aligned} x + y &= 1 \\ x + \alpha y &= 0 \end{aligned} \quad \Rightarrow \quad \begin{aligned} x &= 1 - \frac{1}{1 - \alpha} \\ y &= \frac{1}{1 - \alpha} \end{aligned}$$

The solutions are inconsistent if  $\alpha = 1$ , and will become very large (and likely inaccurate on a computer) when  $\alpha \sim 1$ .

The condition number of a problem  $P$  is the maximum size-ratio, where  $d_1$  and  $d_2$  are alternate inputs, and  $s(d_1)$  and  $s(d_2)$  are the corresponding solutions;

$$\kappa(P) = \max_{d_1, d_2} \frac{\|s(d_1) - s(d_2)\|}{\|d_1 - d_2\|}$$

The key point is we are looking for the worst case.

## Recap on Norms

We want to recall that for each vector norm  $(1, 2, \infty)$ , the subordinate (consistent or compatible) matrix norm is

$$\|\mathbf{A}\| = \max \left\{ \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} : \mathbf{x} \neq \mathbf{0} \right\} = \max \{ \|\mathbf{Au}\| : \|\mathbf{u}\| = 1 \}$$

This means that we have a lower bound of the matrix norm;

$$\forall \mathbf{x} \neq \mathbf{0} \quad \left[ \|\mathbf{A}\| \geq \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} \right]$$

A brief recap on the terms;

- consistent norm  $\|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$
- compatible norm consistent norm defined on square matrix
- all subordinate (also called induced) norms are consistent

## Perturbation

Consider the system of linear equations  $\mathbf{Ax} = \mathbf{b}$ . Make a small change to  $\mathbf{b}$ , such that we have  $\mathbf{b} + \delta\mathbf{b}$ . Note that  $\delta\mathbf{x}$  is the change in  $\mathbf{x}$  due to a change in  $\mathbf{b}$ . We then have the following;

$$\begin{aligned} \mathbf{Ax} + \mathbf{A}\delta\mathbf{x}_b &= \mathbf{b} + \delta\mathbf{b} && \text{we know } \mathbf{Ax} = \mathbf{b}, \Rightarrow \\ \mathbf{A}\delta\mathbf{x}_b &= \delta\mathbf{b} && \text{assuming a solution exists, } \Rightarrow \\ \delta\mathbf{x}_b &= \mathbf{A}^{-1} \delta\mathbf{b} && \text{by consistency, } \Rightarrow \\ \|\delta\mathbf{x}_b\| &\leq \|\mathbf{A}^{-1}\| \|\delta\mathbf{b}\| && \text{looking at \% increase, } \Rightarrow \\ \frac{\|\delta\mathbf{x}_b\|}{\|\mathbf{x}\|} &\leq \frac{\|\mathbf{A}^{-1}\| \|\delta\mathbf{b}\|}{\|\mathbf{x}\|} && (1) \end{aligned}$$

we now briefly look at  $\mathbf{Ax} = \mathbf{b}$

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b} && \Rightarrow \\ \|\mathbf{b}\| &\leq \|\mathbf{A}\| \|\mathbf{x}\| && \Rightarrow \\ \|\mathbf{x}\| &\geq \frac{\|\mathbf{b}\|}{\|\mathbf{A}\|} \end{aligned}$$

we can replace the  $\|\mathbf{x}\|$  in the RHS of (1) with something smaller, as it won't affect the inequality

$$\begin{aligned} \frac{\|\delta\mathbf{x}_b\|}{\|\mathbf{x}\|} &\leq \frac{\|\mathbf{A}^{-1}\| \|\delta\mathbf{b}\| \|\mathbf{A}\|}{\|\mathbf{b}\|} && \Rightarrow \\ \frac{\|\delta\mathbf{x}_b\|}{\|\mathbf{x}\|} &\leq \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} \end{aligned}$$

## Condition Number

Similarly, we can do the same for a small change in  $\mathbf{A}$ . As such,  $\delta\mathbf{x}_A$  is the change in  $\mathbf{x}$  due to a change in  $\mathbf{B}$ .

$$\begin{aligned} (\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}_A) &= \mathbf{b} && \Rightarrow \\ \mathbf{Ax} + \mathbf{A}\delta\mathbf{x}_A + \delta\mathbf{Ax} + \delta\mathbf{A}\delta\mathbf{x}_A &= \mathbf{b} && (1), \Rightarrow \\ \mathbf{Ax} + \mathbf{A}\delta\mathbf{x}_A + \delta\mathbf{Ax} &= \mathbf{b} && \text{we know } \mathbf{Ax} = \mathbf{b}, \Rightarrow \\ \mathbf{A}\delta\mathbf{x}_A + \delta\mathbf{Ax} &= \mathbf{0} && \Rightarrow \\ \mathbf{A}\delta\mathbf{x}_A &= -\delta\mathbf{Ax} && \text{assuming a solution exists, } \Rightarrow \\ \delta\mathbf{x}_A &= -\mathbf{A}^{-1} \delta\mathbf{Ax} && \text{by consistency, norm is positive, } \Rightarrow \\ \|\delta\mathbf{x}_A\| &\leq \|\mathbf{A}^{-1}\| \|\delta\mathbf{A}\| \|\mathbf{x}\| && \Rightarrow \end{aligned}$$

$$\begin{aligned}\frac{\|\delta \mathbf{x}_A\|}{\|\mathbf{x}\|} &\leq \|\mathbf{A}^{-1}\| \|\delta \mathbf{A}\| && \text{we want relative change, } \Rightarrow \\ \frac{\|\delta \mathbf{x}_A\|}{\|\mathbf{x}\|} &\leq \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|}\end{aligned}$$

As such, we can define the condition of a matrix  $\mathbf{A}$  as follows;

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}^{-1}\| \|\mathbf{A}\|$$

This gives us the following useful property, allowing the condition to be an upper bound on the worst case;

$$\text{cond}(\mathbf{A}) \geq \max \left\{ \frac{\frac{\|\delta \mathbf{x}_b\|}{\|\mathbf{x}\|}}{\frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|}}, \frac{\frac{\|\delta \mathbf{x}_A\|}{\|\mathbf{x}\|}}{\frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|}} \right\}$$

## Application to LSM

This is a problem that's well known to be numerically unstable. As usual, we take the case where there are  $m$  equations in  $n$  variables, where  $m > n$ . This gives us a non-square matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . We also note that for the  $\ell_2$ -norm, the condition of  $\mathbf{A}^\top \mathbf{A}$  is the square of the condition of  $\mathbf{A}$ . Essentially, we are multiplying the norms of 4 matrices together, therefore we often get a large norm. We can work through the following example;

$$\begin{aligned}\mathbf{A} &= \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 10^{-4} \end{bmatrix} \\ \mathbf{A}^\top \mathbf{A} &= \begin{bmatrix} 2 & 2 \\ 2 & 2 + 10^{-8} \end{bmatrix} \\ (\mathbf{A}^\top \mathbf{A})^{-1} &= 0.5 \cdot 10^8 \cdot \begin{bmatrix} 2 + 10^{-8} & -2 \\ -2 & 2 \end{bmatrix} \\ \text{cond}(\mathbf{A}^\top \mathbf{A})^{-1} &= (4 + 10^{-8})(4 + 10^{-8}) \cdot 0.5 \cdot 10^8 \\ &\approx 8 \cdot 10^8\end{aligned}$$

While there is no threshold for deciding whether a matrix is ill-conditioned, the rule of thumb is that for some condition number  $\kappa(\mathbf{A})$ , you can lose  $\log_{10}(\kappa(\mathbf{A}))$  significant figures in accuracy. For example, if we had a accurate starting point, for example something with 12 decimal places, and  $\kappa = 1000000$ , the precision will only reduce to 6 decimal places.

The residual vectors are only reliable indicators of the accuracy **if the problem is well-conditioned**;

$$\mathbf{r} = \mathbf{A}\mathbf{x}_0 - \mathbf{b}$$

## Metric Spaces

Take a non-empty set of points  $S$ , with a (distance) function  $d : S \times S \rightarrow \mathbb{R}$  (the metric of the space), which satisfies  $(\forall x, y, z \in S)$ ;

- (1)  $d(x, x) = 0$
- (2)  $d(x, y) > 0$  given  $x \neq y$
- (3)  $d(x, y) = d(y, x)$  (symmetry)
- (4)  $d(x, y) \leq d(x, z) + d(z, y)$  (triangle inequality)

This is true for distances in real spaces  $\mathbb{R}, \mathbb{R}^2, \mathbb{R}^3, \dots$ . Another example of this is vector norms, as they satisfy the above properties.

We can justify that the discrete metric, defined below, satisfies those conditions;

$$d(x, y) = \begin{cases} 0 & x = y \\ 1 & x \neq y \end{cases}$$

Trivially, the first 3 conditions are satisfied by the definition, however the last takes more work to prove. We know that  $d(x, y)$  is 0 or 1, and  $d(x, z) + d(z, y)$  is 0, 1, or 2. If the latter is 1 or 2, there is no work to be done. However, consider the case where the former is 1, and the latter is 0. To start, we assume  $d(x, y) = 1 \Rightarrow x \neq y$ . We also assume  $d(x, z) + d(z, y) = 0$ , therefore  $d(x, z) = 0$  and  $d(z, y) = 0$ . As such, we have  $x = z$  and  $z = y$ , giving us  $x = y$ , which contradicts our initial assumption. Therefore that case **cannot** happen, and the last condition is satisfied.

## Convergence

Briefly, we recap convergence from last year's **CO145**. To say that  $a_1, a_2, \dots$  converges to some limit  $\ell \in \mathbb{R}$ , which is written as;

$$a_n \rightarrow \ell \text{ as } n \rightarrow \infty$$

or

$$\lim_{n \rightarrow \infty} a_n = \ell$$

This is true iff (from first year notes)

$$\forall \epsilon > 0 \exists N \in \mathbb{N} [\forall n > N [|a_n - \ell| < \epsilon]]$$

Which holds iff (Cauchy)

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

The second method is useful as we don't need to know the limit (when it exists), such as in a recurrence relation or a recursively defined function, and can also be used as a test for divergence.

## 19th February 2020

### Revised Matrix Norm Terminology

- a matrix norm  $\|\cdot\|_{\alpha, \beta}$  is **consistent** relative to vector-norms  $\alpha$  on  $\mathbb{R}^m$  and  $\beta$  on  $\mathbb{R}^n$  if

$$\forall \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{x} \in \mathbb{R}^n [|\mathbf{Ax}|_{\alpha} \leq \|\mathbf{A}\|_{\alpha, \beta} \|\mathbf{x}\|_{\beta}]$$

- any induced / subordinate matrix norm is also sub-multiplicative because;

$$\|\mathbf{AB}\| = \max\{\|\mathbf{ABu}\| : \|\mathbf{u}\| = 1\} \leq \max\{\|\mathbf{A}\| \|\mathbf{B}\| \|\mathbf{u}\| : \|\mathbf{u}\| = 1\} = \|\mathbf{A}\| \|\mathbf{B}\|$$

## Cauchy Sequence

Generalising limits to metric spaces (let it be  $S$  in lieu of repeating myself), we have the following definitions;

- a sequence  $x_1, x_2, \dots \in S$  converges to a limit  $\ell \in S$  iff

$$\forall \epsilon > 0 \exists N \in \mathbb{N} [\forall n > N [d(x_n, \ell) < \epsilon]]$$

- a sequence  $\{x_n\}$  is called a **Cauchy sequence** iff

$$\forall \epsilon > 0 \exists N \in \mathbb{N} [\forall n, m > N [d(x_n, x_m) < \epsilon]]$$

- a metric space in which every Cauchy sequence has a limit in  $S$  is **complete**

If a sequence  $\{x_n\}$  converges in  $S$ , then  $\{x_n\}$  is Cauchy. This means that any converging sequence is Cauchy, but not every Cauchy sequence is convergent. If every Cauchy sequence is convergent, then  $S$  is complete.

## Application

We will mainly be working in real vector spaces  $\mathbb{R}^k$ . It can also be shown that  $\mathbb{R}^k$  is complete for all  $k \geq 1$  (proof not needed for this course / degree). Since these metric spaces are complete, we can then say that any sequence we encounter (in real vector spaces) is convergent iff it is Cauchy, which is relatively easy to test.

## Uniqueness of Limits

This is an example of the type of reasoning we are expected to do in this course. We claim that a convergent sequence  $\{x_n\} \in \mathbb{R}^k$  has a unique limit. The proof is as follows;

Suppose  $x_n \rightarrow \ell_1$  and  $x_n \rightarrow \ell_2$ , as  $n \rightarrow \infty$ . By the last property of the metric, we can say (for all  $n$ );

$$d(\ell_1, \ell_2) \leq d(\ell_1, x_n) + d(x_n, \ell_2)$$

For a sufficiently large  $n$ , we can say the following;

$$d(\ell_1, \ell_2) < \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon$$

Note that in the step above, we take any number, therefore  $\frac{\epsilon}{2}$  is just as good as  $\epsilon$ . Since  $\epsilon$  is **any** positive number (by definition of a metric);

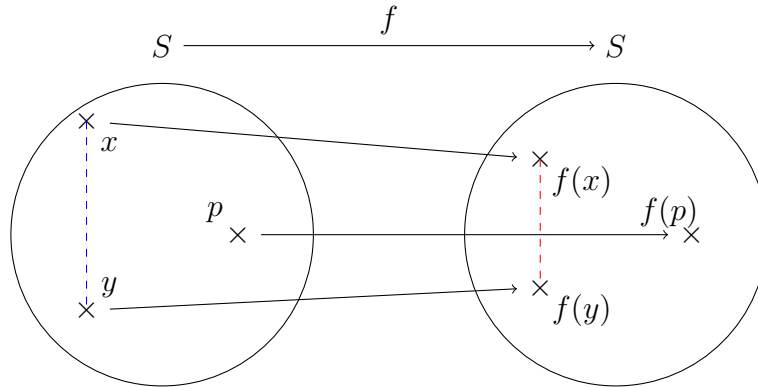
$$d(\ell_1, \ell_2) < \epsilon \Rightarrow d(\ell_1, \ell_2) = 0 \Leftrightarrow \ell_1 = \ell_2$$

## Fixed Points

Let  $f : S \rightarrow S$ . A point  $p \in S$  is a fixed point (or fixpoint) if  $f(p) = p$ . Additionally, a function  $f$  is a **contraction** of  $S$  if there exists a contraction constant  $0 < \alpha < 1$ ;

$$\forall x, y \in S [d(f(x), f(y)) \leq \alpha d(x, y)]$$

Visually, we can represent it as the following;



## Fixed Point Theorem

The fixed point theorem states that a continuous contraction  $f$  of a complete metric space  $S$  has a unique fixed point. The proof for this is as follows;

Take any point  $x \in S$ , we define  $\{p_n\}$  as;

$$\begin{aligned} p_0 &= x \\ p_{n+1} &= f(p_n) \\ d(p_{n+1}, p_n) &= d(f(p_n), f(p_{n-1})) \\ &\leq \alpha d(p_n, p_{n-1}) \\ &= \alpha d(f(p_{n-1}), f(p_{n-2})) \\ &\leq \alpha^2 d(p_{n-1}, p_{n-2}) \\ &\vdots \\ &\leq \alpha^n c \end{aligned} \quad \text{let } c = d(p_1, p_0)$$

We then assume  $m > n$  (somewhat redundant by symmetry);

$$\begin{aligned}
d(p_m, p_n) &\leq d(p_m, p_{m-1}) + \cdots + d(p_{n+1}, p_n) \\
&= \sum_{k=n}^{m-1} d(p_{k+1}, p_k) \\
&\leq \sum_{k=n}^{m-1} \alpha^k c \\
&= c \sum_{k=n}^{m-1} \alpha^k && \text{notice the geometric series} \\
&= c \frac{\alpha^n - \alpha^m}{1 - \alpha} && \alpha^m > 0 \\
&< \frac{c\alpha^n}{1 - \alpha}
\end{aligned}$$

As  $m, n \rightarrow \infty$ ,  $\alpha^n \rightarrow 0$ , so we have  $d(p_m, p_n) \rightarrow 0$ . This means that the points converge, therefore  $\{p_n\}$  is Cauchy. With the completion property, we know that  $p_n \rightarrow p$  for some point  $p \in S$ . To prove that it is a fixed point, we need to employ the continuity property of  $f$ ;

$$\begin{aligned}
f(p) &= f(\lim_{n \rightarrow \infty} p_n) \\
&= \lim_{n \rightarrow \infty} f(p_n) && \text{by continuity, somehow?} \\
&= \lim_{n \rightarrow \infty} p_{n+1} \\
&= p
\end{aligned}$$

Finally, to prove the uniqueness of this, take two points  $p, p' \in S$ , both assumed to be fixpoints;

$$\begin{aligned}
d(p, p') &= d(f(p), f(p')) && \text{as they are fixpoints} \\
&\leq \alpha d(p, p') && \text{by contraction}
\end{aligned}$$

As  $0 < \alpha < 1$ ,  $d(p, p') = 0$ , hence  $p = p'$ , therefore the fixpoint is unique.

## Iterative Solutions of Linear Equations

Split a square matrix  $\mathbf{A}$  to  $\mathbf{M} = \mathbf{A} - \mathbf{N}$ , where  $\mathbf{M}$  is non-singular.

$$\begin{aligned}
\mathbf{Ax} &= \mathbf{b} && \Rightarrow \\
\mathbf{Mx} &= \mathbf{b} - \mathbf{Nx} && \Rightarrow \\
\mathbf{x} &= \mathbf{M}^{-1}\mathbf{b} - \mathbf{M}^{-1}\mathbf{Nx} && \Rightarrow \\
\mathbf{x} &= \mathbf{Gx} + \mathbf{c} && \text{where } \mathbf{G} = -\mathbf{M}^{-1}\mathbf{N}, \mathbf{c} = \mathbf{M}^{-1}\mathbf{b}
\end{aligned}$$

We then define the iteration to be as follows, with a starting point  $\mathbf{x}^{(0)}$ ;

$$\mathbf{x}^{(k+1)} = \mathbf{Gx}^{(k)} + \mathbf{c}$$

**20th February 2020**

## Proof of Convergence

We want to prove that for any consistent matrix norm, if  $\|\mathbf{G}\| < 1$ , then the iterative sequence converges for any starting  $\mathbf{x}^{(0)}$ .

Suppose  $\mathbf{x}^{(k)} \rightarrow \hat{\mathbf{x}}$ , where  $\hat{\mathbf{x}} = \mathbf{G}\hat{\mathbf{x}} + \mathbf{c}$ . We want to define  $\mathbf{y}^{(k+1)}$  as follows;

$$\begin{aligned}\mathbf{y}^{(k+1)} &= \mathbf{x}^{(k+1)} - \hat{\mathbf{x}} \\ &= \mathbf{G}\mathbf{x}^{(k)} + \mathbf{c} - (\mathbf{G}\hat{\mathbf{x}} + \mathbf{c}) \\ &= \mathbf{G}\mathbf{x}^{(k)} - \mathbf{G}\hat{\mathbf{x}} \\ &= \mathbf{G}(\mathbf{x}^{(k)} - \hat{\mathbf{x}}) \\ &= \mathbf{G}\mathbf{y}^{(k)}\end{aligned}$$

From this, we can use the property that  $\mathbf{G}$  is consistent as follows;

$$\begin{aligned}\|\mathbf{y}^{(k+1)}\| &\leq \|\mathbf{G}\|\|\mathbf{y}^{(k)}\| \\ &\leq \|\mathbf{G}\|^2\|\mathbf{y}^{(k-1)}\| \\ &\leq \|\mathbf{G}\|^3\|\mathbf{y}^{(k-2)}\| \\ &\vdots \\ &\leq \|\mathbf{G}\|^{k+1}\|\mathbf{y}^{(0)}\|\end{aligned}$$

However, we know that  $\|\mathbf{G}\|^{k+1}\|\mathbf{y}^{(0)}\| \rightarrow 0$  when  $k \rightarrow \infty$ , as  $\|\mathbf{G}\| < 1$ . Therefore,  $\mathbf{y}^{(k+1)} \rightarrow \mathbf{0}$  by the property of norms, hence  $\mathbf{x}^{(k+1)} \rightarrow \hat{\mathbf{x}}$

### Common Split

Assume, without loss of generality, that  $\mathbf{A}$  has no zeroes on the diagonal. If it does, we can relabel the variables  $x_i$  to permute the rows and columns to have no zeroes on the diagonal.

We then define the following;

$$\mathbf{A} = \mathbf{D} + \mathbf{L} + \mathbf{U}$$

where  $\mathbf{D}$  is the diagonal component,  $\mathbf{L}$  is the **strict** lower triangular component (hence without the diagonal), and similar for  $\mathbf{U}$ .

### Jacobi Method

In the Jacobi method, we set the matrices (and vector) as follows;

$$\begin{aligned}\mathbf{M} &= \mathbf{D} \\ \mathbf{N} &= \mathbf{L} + \mathbf{U} \\ \mathbf{G} &= -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U}) \\ \mathbf{c} &= \mathbf{D}^{-1}\mathbf{b}\end{aligned}$$

therefore  $\mathbf{x}^{(k+1)} = \mathbf{G}\mathbf{x}^{(k)} + \mathbf{c}$  yields

$$x_i^{(k+1)} = \frac{1}{a_{i,i}} \left( b_i - \sum_{j \neq i} a_{j,i} x_j^{(k)} \right)$$

This is easily done in parallel.



## Gauss-Seidel Method

$$\mathbf{M} = \mathbf{D} + \mathbf{L}$$

$$\mathbf{N} = \mathbf{U}$$

$$\mathbf{G} = -(\mathbf{D} + \mathbf{L})^{-1}\mathbf{U}$$

$$\mathbf{c} = (\mathbf{D} + \mathbf{L})^{-1}\mathbf{b}$$

therefore  $\mathbf{x}^{(k+1)} = \mathbf{G}\mathbf{x}^{(k)} + \mathbf{c}$  yields

$$x_i^{(k+1)} = \frac{1}{a_{i,i}} \left( b_i - \sum_{j=1}^{i-1} a_{i,j}x_j^{(k+1)} - \sum_{j=i+1}^m a_{i,j}x_j^{(k)} \right)$$

This gives faster convergence as we're able to use more up-to-date data. Working from the first column, we get;

$$\begin{aligned} x_1^{(k+1)} &= \frac{1}{a_{1,1}} \left( b_1 - \sum_{j=2}^m a_{1,j}x_j^{(k)} \right) \\ x_2^{(k+1)} &= \frac{1}{a_{2,2}} \left( b_2 - a_{2,1}x_1^{(k+1)} - \sum_{j=3}^m a_{2,j}x_j^{(k)} \right) \quad \text{using the previous result} \\ &\vdots \end{aligned}$$

However, this must be done in order, and therefore cannot be done in parallel like in Jacobi.

## Convergence Conditions

We define a square matrix  $\mathbf{A}$  as **strictly** row diagonally dominant for the  $i^{\text{th}}$  row if

$$|a_{i,i}| > \sum_{j \neq i} |a_{i,j}|$$

This means that the sum of the moduli of the all the elements on the row (except the diagonal element) is less than the modulus of the diagonal element. The theorem is that this is a sufficient condition for both Jacobi and Gauss-Seidel to converge. Below, we will do the proof for Jacobi;

- (1) show  $\|\mathbf{G}\| < 1$ , where  $\mathbf{G} = -\mathbf{D}^{-1}\mathbf{N}$ , using the consistent  $\ell_\infty$ -norm
- (2) strict diagonal dominance implies all absolute row sums are less than 1
- (3) the maximum absolute row sum is less than 1, therefore so is  $\|\mathbf{G}\|_\infty$
- (4) any norm less than 1 leads to convergence

$$\mathbf{G} = -\mathbf{D}^{-1}\mathbf{N}$$

$$= - \begin{bmatrix} 0 & \times & \times & \times & \times & \times & \times \\ \times & \ddots & \times & \times & \times & \times & \times \\ \times & \times & \ddots & \times & \times & \times & \times \\ \frac{a_{i,1}}{a_{i,i}} & \dots & \frac{a_{i,i-1}}{a_{i,i}} & 0 & \frac{a_{i,i+1}}{a_{i,i}} & \dots & \frac{a_{i,m}}{a_{i,i}} \\ \times & \times & \times & \times & \ddots & \times & \times \\ \times & \times & \times & \times & \times & \ddots & \times \\ \times & \times & \times & \times & \times & \times & 0 \end{bmatrix}$$

$$\sum_{j \neq i} |a_{i,j}| < |a_{i,i}|$$

by row diagonal dominance,  $\Rightarrow$

$$\frac{1}{|a_{i,i}|} \sum_{j \neq i} |a_{i,j}| < 1$$

Since we have an arbitrary row  $i$ , this holds for any row, hence the  $\ell_\infty$ -norm is less than 1 ( $\|\mathbf{G}\| < 1$ ).

## Computation of Eigenvalues and Eigenvectors + Power Method

Computing the solutions of the characteristic polynomial is expensive. It takes  $\Theta(n^3)$  to compute a single eigenvector of  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with Gaussian elimination. Therefore, it takes  $\Theta(n^4)$  to compute all eigenvectors.

One technique used is the power method, which can be applied to find other eigenvectors. However, by itself, it converges to an eigenvector of the dominant eigenvalue. Assume the matrix  $\mathbf{A}$  has a unique dominant eigenvalue (which has the maximum modulus), and a corresponding unique eigenvector. Take a starting vector  $\mathbf{x}^{(0)}$ , we can define the recurrence for  $k = 0, 1, 2, \dots$  as;

$$\mathbf{x}^{(k+1)} = \frac{\mathbf{A}\mathbf{x}^{(k)}}{\|\mathbf{A}\mathbf{x}^{(k)}\|}$$

In general, the idea is to keep multiplying the current vector by  $\mathbf{A}$  and normalising. The  $\ell_2$ -norm is useful as it is equal to the modulus of the dominant eigenvalue.

The proposition is that for a diagonalisable matrix  $\mathbf{M} \in \mathbb{R}^{n \times n}$ , with distinct eigenvalues (moduli), the iteration converges to the eigenvector of  $\mathbf{A}$ . The absolute value of the dominant eigenvalue is the limit of the normalisation constant;

$$\lim_{k \rightarrow \infty} \|\mathbf{A}\mathbf{x}^{(k)}\|$$

The proof is as follows;

Let  $(\mathbf{e}_1, \dots, \mathbf{e}_n)$ , consisting of the unit-eigenvectors of  $\mathbf{A}$  form a basis of  $\mathbb{R}^n$ . We know we can do this as  $\mathbf{A}$  is diagonalisable. Therefore, for an arbitrary  $\mathbf{x}^{(0)} \in \mathbb{R}^n$ ;

$$\mathbf{x}^{(0)} = \sum_{i=1}^n \alpha_i \mathbf{e}_i$$

For  $k \geq 1$ ;

$$\begin{aligned} \mathbf{x}^{(k)} &= \frac{\mathbf{A}\mathbf{x}^{(k-1)}}{\|\mathbf{A}\mathbf{x}^{(k-1)}\|} && \text{by recurrence relation} \\ &= \frac{\mathbf{A}^k \mathbf{x}^{(0)}}{\|\mathbf{A}^k \mathbf{x}^{(0)}\|} \\ &= \frac{\sum_{i=1}^n \alpha_i \lambda_i^k \mathbf{e}_i}{\|\mathbf{A}^k \mathbf{x}^{(0)}\|} && \text{scale by an eigenvalue of eigenvector} \\ &= \frac{\lambda_1^k}{\|\mathbf{A}^k \mathbf{x}^{(0)}\|} \sum_{i=1}^n \alpha_i \frac{\lambda_i^k}{\lambda_1^k} \mathbf{e}_i \\ &= \frac{\lambda_1^k}{\|\mathbf{A}^k \mathbf{x}^{(0)}\|} \left( \alpha_1 \mathbf{e}_1 + \sum_{i=2}^n \alpha_i \frac{\lambda_i^k}{\lambda_1^k} \mathbf{e}_i \right) \\ \lim_{k \rightarrow \infty} \frac{\lambda_i^k}{\lambda_1^k} &= 0 && \text{for } i \neq 1, |\lambda_i| < |\lambda_1| \\ \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} &= \frac{\lambda_1^k}{\|\mathbf{A}^k \mathbf{x}^{(0)}\|} \alpha_1 \mathbf{e}_1 && \text{therefore } \mathbf{x}^{(0)} \rightarrow \mathbf{e}_1 \end{aligned}$$

Using this result, we can conclude the following;

$$\lim_{k \rightarrow \infty} \|\mathbf{A}\mathbf{x}^{(k)}\| = \|\mathbf{A}\mathbf{e}_1\| = \|\lambda_1 \mathbf{e}_1\| = |\lambda_1|$$

The limitations of this method are the following;

- starting vector  $\mathbf{x}^{(0)}$  may have no component in the dominant eigenvector, such that  $\mathbf{e}_1^\top \mathbf{x}^{(0)} = 0$ , which leads to  $\alpha_1 = 0$
- may have multiple eigenvalues with maximum modulus, iteration then converges to a linear combination of those eigenvectors
- convergence can be slow, if for example  $\lambda_2$  is close to  $\lambda_1$

Note that if the eigenvalues for the matrix  $\mathbf{A}$  are  $\lambda_1, \lambda_2, \dots$ , then the eigenvalues for the matrix  $\mathbf{A} - \sigma \mathbf{I}$  are  $\lambda_1 - \sigma, \lambda_2 - \sigma, \dots$ , with the same eigenvectors.

## Inverse Iteration

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \Rightarrow \mathbf{x} = \lambda\mathbf{A}^{-1}\mathbf{x} \Rightarrow \mathbf{A}^{-1}\mathbf{x} = \frac{1}{\lambda}\mathbf{x}$$

This means that a matrix  $\mathbf{A}^{-1}$ , has eigenvalues  $\frac{1}{\lambda_i}$ , with the same corresponding eigenvectors as  $\mathbf{A}$ . Therefore the power iteration yields the reciprocal of the smallest eigenvalue of  $\mathbf{A}$ . An algorithm for the inverse iteration is done by a recurrence relation on pairs  $(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$  for  $k \geq 1$ ;

$$\begin{aligned} \mathbf{A}\mathbf{y}^{(k+1)} &= \mathbf{x}^{(k)} \\ \mathbf{x}^{(k+1)} &= \frac{\mathbf{y}^{(k+1)}}{\|\mathbf{y}^{(k+1)}\|} \quad \left( = \frac{\mathbf{A}^{-1}\mathbf{x}^{(k)}}{\|\mathbf{A}^{-1}\mathbf{x}^{(k)}\|} \right) \end{aligned}$$

Since computing the inverse of  $\mathbf{A}$  is expensive, the above can be used. We only need to solve for  $\mathbf{A}\mathbf{y}^{(k+1)} = \mathbf{x}^{(k)}$  with any method at each iteration.

By using this with a shift  $\sigma$ , the eigenvalue (and eigenvector) computed will be the eigenvalue closest to  $\sigma$ , as that now has the smallest absolute value.

## Rayleigh Quotient

Imagine we are at a point where we have an approximated eigenvector, let it be  $\mathbf{z}$ , with an approximated eigenvalue  $\lambda$ . Then, for  $1 \leq i \leq n$ , we have;

$$\sum_{j=1}^n a_{i,j} z_j = \lambda z_i$$

Therefore, we can see this as a least squares problem of dimension  $n \times 1$ ;

$$\underbrace{\mathbf{z}}_{n \times 1} \underbrace{\lambda}_{1 \times 1} = \underbrace{\mathbf{A}}_{n \times n} \underbrace{\mathbf{z}}_{n \times 1}$$

The normal equation is now  $\mathbf{z}^\top \mathbf{z} \lambda = \mathbf{z}^\top \mathbf{A} \mathbf{z}$ . Using this, we can then get the Rayleigh Quotient, which is the approximate solution;

$$\hat{\lambda} \approx \frac{\mathbf{z}^\top \mathbf{A} \mathbf{z}}{\mathbf{z}^\top \mathbf{z}}$$

From here, we can then choose  $\sigma = \hat{\lambda}$ , and shift it again. Ideally,  $\sigma$  will now be extremely close to the actual value, hence the shifted eigenvalue will be very small, causing the reciprocal to be very large. While  $\sigma$  changes each iteration, the eigenvectors don't.

## Second-Dominant Eigenvalue

Deflation is a term for transforming a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  to a matrix  $\mathbf{B} \in \mathbb{R}^{(n-1) \times (n-1)}$ , with the largest (absolute) eigenvalue absent. We define a non-singular  $\mathbf{H}$  as such, where  $\mathbf{x}_1$  is the eigenvector of  $\mathbf{A}$

with the largest absolute eigenvalue,  $\alpha$  is some scalar constant, and  $\mathbf{e}_1$  is the first standard basis vector;

$$\begin{aligned}
\mathbf{H}\mathbf{x}_1 &= \alpha\mathbf{e}_1 && \Rightarrow \\
\mathbf{x}_1 &= \alpha\mathbf{H}^{-1}\mathbf{e}_1 && \Rightarrow \\
\mathbf{H}^{-1}\mathbf{e}_1 &= \frac{\mathbf{x}_1}{\alpha} \\
\mathbf{A}\mathbf{x}_1 &= \lambda_1\mathbf{x}_1 && \text{we know } \mathbf{x}_1 \text{ is an eigenvector} \\
\mathbf{H}\mathbf{A}\mathbf{H}^{-1}\mathbf{e}_1 &= \frac{1}{\alpha}\mathbf{H}\mathbf{A}\mathbf{x}_1 \\
&= \frac{\lambda_1}{\alpha}\mathbf{H}\mathbf{x}_1 \\
&= \lambda_1\mathbf{e}_1 \\
&= (\mathbf{H}\mathbf{A}\mathbf{H}^{-1})_1 && \text{first column} \\
&= \begin{bmatrix} \lambda_1 & \mathbf{b}^\top \\ \mathbf{0} & \mathbf{B} \end{bmatrix} && \text{where } \mathbf{0}, \mathbf{b} \in \mathbb{R}^{n-1}, \mathbf{B} \in \mathbb{R}^{(n-1) \times (n-1)}
\end{aligned}$$

To show they have the same eigenvalues, we do the following;

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \Leftrightarrow \mathbf{H}\mathbf{A}\mathbf{x} = \lambda\mathbf{H}\mathbf{x} \Leftrightarrow \mathbf{H}\mathbf{A}\mathbf{H}^{-1}(\mathbf{H}\mathbf{x}) = \lambda\mathbf{H}\mathbf{x}$$

Therefore,  $\mathbf{A}$  and  $\mathbf{H}\mathbf{A}\mathbf{H}^{-1}$  have the same eigenvalues, but  $\mathbf{H}\mathbf{A}\mathbf{H}^{-1}$  has eigenvectors  $\mathbf{y}_i = \mathbf{H}\mathbf{x}_i$ . Therefore  $\mathbf{B}$  has eigenvectors  $\lambda_2, \dots, \lambda_n$ .

We want to prove that the eigenvector corresponding to  $\lambda_2$ , for a given  $\mathbf{H}$ , is

$$\begin{aligned}
\mathbf{x}_2 &= \mathbf{H}^{-1} \begin{bmatrix} \beta \\ \mathbf{z}_2 \end{bmatrix} && \text{where } \mathbf{z}_2 \text{ is a dominant eigenvector of } \mathbf{B} \\
\beta &= \frac{\mathbf{b}^\top \mathbf{z}_2}{\lambda_2 - \lambda_1}
\end{aligned}$$

To verify this, we do the following;

$$\begin{aligned}
\mathbf{x}_2 &= \mathbf{H}^{-1}\mathbf{y} && \mathbf{y} \text{ is an eigenvector of } \mathbf{H}\mathbf{A}\mathbf{H}^{-1} \text{ with eigenvalue } \lambda_2 \\
\text{let } \mathbf{y} &= \begin{bmatrix} y_1 \\ \mathbf{z} \end{bmatrix} \\
\mathbf{H}\mathbf{A}\mathbf{H}^{-1} &= \begin{bmatrix} \lambda_1 & \mathbf{b}^\top \\ \mathbf{0} & \mathbf{B} \end{bmatrix} \\
\begin{bmatrix} \lambda_1 & \mathbf{b}^\top \\ \mathbf{0} & \mathbf{B} \end{bmatrix} \begin{bmatrix} y_1 \\ \mathbf{z} \end{bmatrix} &= \lambda_2 \begin{bmatrix} y_1 \\ \mathbf{z} \end{bmatrix} && \Rightarrow \\
\lambda_2 y_1 &= \lambda_1 y_1 + \mathbf{b}^\top \mathbf{z} && \Rightarrow \\
y_1 &= \frac{\mathbf{b}^\top \mathbf{z}}{\lambda_2 - \lambda_1} \\
\lambda_2 \mathbf{z} &= \mathbf{0} y_1 + \mathbf{B} \mathbf{z} \\
&= \mathbf{B} \mathbf{z} && \Rightarrow \\
\beta &= y_2 \\
\mathbf{z}_2 &= \mathbf{z}
\end{aligned}$$

## Householder

One of the available choices for  $\mathbf{H}$  is the Householder transformation, which represents a reflection in the plane with a normal vector of  $\mathbf{u}$ .

$$\mathbf{H} = \mathbf{I} - \frac{2\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top \mathbf{u}}$$

For this part, we define  $\mathbf{u} = \mathbf{v} \pm \|\mathbf{v}\|_2 \mathbf{e}_1$ , for a vector  $\mathbf{v} = \mathbf{x}_1$ . Then we have  $\mathbf{H}\mathbf{v} = \mp\|\mathbf{v}\|_2 \mathbf{e}_1$ , so we find  $\alpha = \mp\|\mathbf{x}_1\|_2$ . This result is proved as follows;

$$\begin{aligned}
v_1 &= \mathbf{v}^\top \mathbf{e}_1 \\
\mathbf{u}^\top \mathbf{u} &= 2\|\mathbf{v}\|(\|\mathbf{v}\| \pm v_1) \\
\mathbf{u}^\top \mathbf{v} &= \mathbf{v}^\top \mathbf{u} \\
&= \|\mathbf{v}\|(\|\mathbf{v}\| \pm v_1) \\
\mathbf{H} &= \mathbf{I} - \frac{2\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top \mathbf{u}} \quad \Rightarrow \\
\mathbf{H}\mathbf{v} &= \mathbf{v} - \frac{2\mathbf{u}(\mathbf{u}^\top \mathbf{v})}{\mathbf{u}^\top \mathbf{u}} \\
&= \mathbf{v} - \mathbf{u} \\
&= \mp\|\mathbf{v}\|_2 \mathbf{e}_1 \quad \blacksquare
\end{aligned}$$

26th February 2020

## QR Algorithm

Note that for an upper triangular matrix

$$\mathbf{U} = \begin{bmatrix} u_{1,1} & \times & \times & \times \\ & u_{2,2} & \times & \times \\ & & \ddots & \times \\ & & & u_{n,n} \end{bmatrix}$$

The characteristic polynomial is

$$|\mathbf{U} - \lambda \mathbf{I}| = \prod_{i=1}^n (u_{i,i} - \lambda)$$

Meaning that the eigenvalues are  $u_{1,1}, u_{2,2}, \dots, u_{n,n}$ , which is the diagonal of the matrix.

Also note that a matrix  $\mathbf{A}$  is similar to  $\mathbf{B}$  if  $\mathbf{B} = \mathbf{P}^{-1} \mathbf{A} \mathbf{P}$ . This is an equivalence relation, and it also preserves the characteristic polynomial, hence has the same eigenvalues.

The general idea of the QR algorithm is that it generates a sequence of similar matrices  $\mathbf{A}^{(0)}, \mathbf{A}^{(1)}, \dots$  which approaches an upper triangular matrix. The diagonal of this upper triangular matrix will then give the eigenvalues of  $\mathbf{A}$ . It follows this general sequence of steps;

- (1) starting with the standard basis, set  $\mathbf{A}^{(0)} = \mathbf{A} (= \mathbf{Q}\mathbf{R} = \mathbf{Q}^{(1)}\mathbf{R}^{(1)})$
- (2) apply QR decomposition with any method to compute  $\mathbf{A}^{(m)} = \mathbf{Q}^{(m+1)}\mathbf{R}^{(m+1)}$ , where  $\mathbf{Q}^{(m+1)}$  is orthogonal, and  $\mathbf{R}^{(m+1)}$  is upper triangular
- (3) define  $\mathbf{A}^{(m+1)} = \mathbf{R}^{(m+1)}\mathbf{Q}^{(m+1)} = \mathbf{Q}^{(m+1)\top} \mathbf{A}^{(m)} \mathbf{Q}^{(m+1)}$

note that because  $\mathbf{Q}^{(m+1)}$  is orthogonal (actually orthonormal),  $\mathbf{Q}^{(m+1)\top} = \mathbf{Q}^{(m+1)-1}$ ,  $\mathbf{A}^{(m+1)}$  must be similar to  $\mathbf{A}^{(m)}$

- (4) halt after sufficient iterations

We claim that  $\mathbf{A}^{(m)} \rightarrow \mathbf{U}$ , where  $\mathbf{U}$  is an upper triangular matrix, as  $m \rightarrow \infty$ .

$$\mathbf{A}^{(m)} = \mathbf{Q}^{(m)\top} \mathbf{Q}^{(m-1)\top} \dots \mathbf{Q}^{(1)\top} \mathbf{A}^{(0)} \underbrace{\mathbf{Q}^{(1)} \mathbf{Q}^{(2)} \dots \mathbf{Q}^{(m)}}_{\tilde{\mathbf{Q}}^{(m)}} = \tilde{\mathbf{Q}}^{(m)\top} \mathbf{A} \tilde{\mathbf{Q}}^{(m)}$$

Note that  $\tilde{\mathbf{Q}}^{(m)}$  is orthogonal, as it is the product of orthogonal matrices, hence  $\mathbf{A}^{(m)}$  must be similar to  $\mathbf{A}$ , therefore  $\mathbf{A}^{(m)}$  has the same eigenvalues as  $\mathbf{A}$ . The eigenvectors of  $\mathbf{A}$  are;

$$\left\{ \tilde{\mathbf{Q}}^{(m)\top} \mathbf{w} : \mathbf{w} \text{ is an eigenvector of } \mathbf{A}^{(m)} \right\}$$

It is easy to compute the eigenvectors of an upper triangular matrix, which we approach, with backwards substitution. Note that it can be quicker to use the inverted shift with an approximation of the eigenvalue.

Note that for any symmetric matrix  $\mathbf{A}$ ,  $\mathbf{A}^{(m)}$  must also be symmetric. Therefore, if  $\mathbf{A}^{(m)}$  approaches an upper triangular matrix, it must approach  $\mathbf{\Lambda}$ , which is a diagonal vector of the eigenvalues of  $\mathbf{A}$ ;

$$\mathbf{\Lambda} = \begin{bmatrix} u_{1,1} & 0 & \cdots & 0 \\ 0 & u_{2,2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & u_{n,n} \end{bmatrix}$$

Due to this,  $\tilde{\mathbf{Q}}^{(m)}$  approaches the matrix of eigenvectors of  $\mathbf{A}$ .

## LU Decomposition

We can write any non-singular matrix  $\mathbf{M} = \mathbf{L}\mathbf{U}$ . This is because we can reduce  $\mathbf{M}$  to row echelon form (which is upper triangular), with the products of lower triangular matrices. Also note that a unit triangular matrix has all 1s on the diagonal. In this proof, we can also use the fact

$$\mathbf{A}^n = \mathbf{Q}^{(1)} \cdots \mathbf{Q}^{(n)} \mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)}$$

We propose that  $\mathbf{A}^{(m)} \rightarrow \mathbf{\Lambda}$ , as  $m \rightarrow \infty$ , for a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , with distinct eigenvalues  $\lambda_1 > \lambda_2 > \cdots > \lambda_n > 0$ . Also let  $\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$  be the eigenvalue decomposition (same as spectral decomposition), and  $\mathbf{Q}^\top = \mathbf{L}\mathbf{U}$  with a **unit** lower triangular  $\mathbf{L}$ , and an upper triangular  $\mathbf{U}$ . The proof is as follows;

$$\begin{aligned} \mathbf{A}^n &= \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top \\ &= \mathbf{Q}^{(1)} \cdots \mathbf{Q}^{(n)} \mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)} && \Rightarrow \\ \mathbf{Q}\mathbf{\Lambda}\mathbf{L}\mathbf{U} &= \mathbf{Q}^{(1)} \cdots \mathbf{Q}^{(n)} \mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)} && \Rightarrow \\ \mathbf{Q}\mathbf{\Lambda}\mathbf{L} &= \mathbf{Q}^{(1)} \cdots \mathbf{Q}^{(n)} \mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)} \mathbf{U}^{-1} && \Rightarrow \\ \mathbf{Q}\mathbf{\Lambda}\mathbf{L}\mathbf{A}^{-n} &= \mathbf{Q}^{(1)} \cdots \mathbf{Q}^{(n)} \mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)} \mathbf{U}^{-1} \mathbf{\Lambda}^{-n} \\ (\mathbf{\Lambda}\mathbf{L}\mathbf{A}^{-n})_{i,j} &= \begin{cases} \ell_{i,j} \frac{\lambda_i^n}{\lambda_j^n} & i > j \\ 1 & i = j \\ 0 & i < j \end{cases} && \Rightarrow \\ \lim_{n \rightarrow \infty} \mathbf{\Lambda}\mathbf{L}\mathbf{A}^{-n} &= \mathbf{I} \end{aligned}$$

Therefore,  $\mathbf{Q}^{(1)} \cdots \mathbf{Q}^{(n)} \mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)} \mathbf{U}^{-1} \mathbf{\Lambda}^{-n} \rightarrow \mathbf{Q}$ . Because the QR decomposition is unique,  $\tilde{\mathbf{Q}}^{(n)} = \mathbf{Q}^{(1)} \cdots \mathbf{Q}^{(n)} \rightarrow \mathbf{Q}$ , and  $\mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)} \mathbf{U}^{-1} \mathbf{\Lambda}^{-n} \rightarrow \mathbf{I}$ , meaning  $\mathbf{R}^{(n)} \cdots \mathbf{R}^{(1)} \rightarrow \mathbf{\Lambda}^n \mathbf{U}$ . Note that for  $\tilde{\mathbf{Q}}^{(n)}$  to tend to a constant matrix,  $\mathbf{Q}^{(n)} \rightarrow \mathbf{I}$ . From the second limit, we can argue that  $\mathbf{R}^{(n)} \rightarrow \mathbf{\Lambda}$ , as  $\mathbf{\Lambda}$  is applied  $n$  times. Finally, we defined  $\mathbf{A}^{(m)} = \mathbf{R}^{(m)} \mathbf{Q}^{(m)}$ , so

$$\lim_{m \rightarrow \infty} \mathbf{A}^{(m)} = \mathbf{\Lambda}$$

. We also have  $\tilde{\mathbf{Q}}^{(m)} \rightarrow \mathbf{Q}$  which diagonalises  $\mathbf{A}$ , and therefore must be the matrix of eigenvectors.

**27th February 2020**

## Laplace Transform

The Laplace transform is a special case of the Fourier transform. The Laplace transform of a function  $f(t)$ , defined for real  $t \geq 0$ , is denoted  $F(s) = (\mathcal{L}f)$ ;

$$F(s) = (\mathcal{L}f)(s) = \int_0^\infty e^{-st} f(t) dt$$

Generally,  $s \in \mathbb{C} = \sigma + i\omega$ , therefore the Laplace transform has the "type signature"  $\mathcal{L} : (\mathbb{R}^+ \rightarrow \mathbb{R}^+) \rightarrow (\mathbb{C} \rightarrow \mathbb{C})$ . However, in this course we normally use positive reals, hence we have an output of  $(\mathbb{R}^+ \rightarrow \mathbb{R}^+)$ . Some common examples are the following;

- unit function (strictly Heaviside function, but we only consider  $t \geq 0$ )

$$\begin{aligned} H(t) &= \begin{cases} 1 & t \geq 0 \\ 0 & t < 0 \end{cases} \\ f(t) &= 1 \\ F(s) &= \int_0^\infty f(t)e^{-st} dt \\ &= \int_0^\infty e^{-st} dt \\ &= -\frac{1}{s} [e^{-st}]_0^\infty \\ &= -\frac{1}{s}(0 - 1) \\ &= \frac{1}{s} \end{aligned}$$

- exponential function (general case of the above)

$$\begin{aligned} f(t) &= e^{-\lambda t} \\ F(s) &= \int_0^\infty e^{-\lambda t} e^{-st} dt \\ &= \int_0^\infty e^{-(\lambda+s)t} dt \\ &= -\frac{1}{\lambda+s} [e^{-(\lambda+s)t}]_0^\infty \\ &= -\frac{1}{\lambda+s}(0 - 1) \\ &= \frac{1}{\lambda+s} \end{aligned}$$

- identity function

$$\begin{aligned} f(t) &= t \\ F(s) &= \int_0^\infty te^{-st} dt \\ &= \left[ \frac{e^{-st}}{-s} t \right]_0^\infty - \int_0^\infty \frac{e^{-st}}{-s} \cdot 1 dt \\ &= 0 + \frac{1}{s} \int_0^\infty e^{-st} dt \\ &= \frac{1}{s^2} \end{aligned}$$

- exponential function

$$\begin{aligned}
f(t) &= t^n \\
F(s) &= \int_0^\infty t^n e^{-st} dt \\
&= \left[ \frac{e^{-st}}{-s} t^n \right]_0^\infty - \int_0^\infty n \frac{e^{-st}}{-s} t^{n-1} dt \\
&= 0 + \frac{n}{s} \int_0^\infty e^{-st} t^{n-1} dt \\
&= \frac{n}{s} \cdot \frac{n-1}{s} \int_0^\infty e^{-st} t^{n-2} dt \\
&\vdots \\
&= \frac{n!}{s^{n+1}}
\end{aligned}$$

- non-negative real power

$$\begin{aligned}
\Gamma(z) &= \int_0^\infty t^{z-1} e^{-t} dt \\
\Gamma(z+1) &= z\Gamma(z) && \forall z \in \mathbb{R}^+ \\
\Gamma(1) &= 1 \\
\Gamma(n) &= (n-1)! && \forall n \geq 1 \in \mathbb{N} \\
f(t) &= t^z \\
F(s) &= \frac{\Gamma(z+1)}{s^{z+1}}
\end{aligned}$$

The Laplace transform has the following properties;

- (1)  $\mathcal{L}$  is a bijection (unique inverse), therefore if we recognise some  $F$ , we know that  $f$  is
- (2) it is a linear operator, and therefore it is simple to compute for many useful functions
- (3) simple to get transform of the derivative of a function
- (4) moment generating function

- Laplace transform of a convolution is a product

We can work out the derivative of a Laplace transform as such;

$$\begin{aligned}
F(s) &= \int_0^\infty e^{-st} f(t) dt \\
\frac{d}{ds} F(s) &= \int_0^\infty \left( \frac{d}{ds} e^{-st} \right) f(t) dt && \text{interchange of limiting operations} \\
&= \int_0^\infty -te^{-st} f(t) dt \\
&= - \int_0^\infty (tf(t)) e^{-st} dt \\
&= -(\mathcal{L}(tf(t)))(s) \\
&= -(\mathcal{L}(\lambda t. tf(t)))(s) && \text{strictly with } \lambda\text{-calculus}
\end{aligned}$$



## Laplace Transform of a Derivative

$$\begin{aligned}
 f'(t) &= \frac{d}{dt}f(t) \\
 (\mathcal{L}f')(s) &= \int_0^\infty e^{-st} f'(t) dt \\
 &= [e^{-st} f(t)]_0^\infty - \int_0^\infty -s e^{-st} f(t) dt \\
 &= -f(0) + s \int_0^\infty e^{-st} f(t) dt \quad \text{assuming sub-exponential convergence for } \lim_{t \rightarrow \infty} f(t) \\
 &= -f(0) + s(\mathcal{L}f)(s) \\
 &= s(\mathcal{L}f)(s) - f(0)
 \end{aligned}$$

Commonly,  $\dot{f}$  is used to denote  $f$  differentiated with respect to  $t$ . It's important to remember the following;

$$(\mathcal{L}f')(s) = (\mathcal{L}\dot{f})(s) = sF(s) - f(0)$$

A example for the sub-exponential convergence, where it doesn't hold;

$$\begin{aligned}
 f(t) &= e^{4t} \\
 F(s) &= \int_0^\infty e^{-st} e^{4t} dt \\
 &= \int_0^\infty e^{(4-s)t} dt
 \end{aligned}$$

Note that  $F \rightarrow \infty$  if  $s \leq 4$ .

## Application of Laplace Transforms to Differential Equations

Starting with a fairly simple example, which would've been encountered previously;

$$\begin{aligned}
 \frac{dy}{dt} &= \dot{y} \\
 &= \mu y \\
 Y &= \mathcal{L}y
 \end{aligned}$$

taking Laplace transforms of both sides

$$\begin{aligned}
 sY(s) - y(0) &= \mu Y(s) && \Rightarrow \\
 (s - \mu)Y(s) &= y(0) && \Rightarrow \\
 Y(s) &= \frac{y(0)}{s - \mu}
 \end{aligned}$$

note that this is similar to  $f(t) = e^{-\lambda t}$ , so by uniqueness, with  $\lambda = -\mu$

$$Y(s) = y(0)e^{\mu t}$$

In the harder example, the Laplace transform allows us to solve simultaneous differential equations as simultaneous linear equations. Consider these three differential equations, with the initial values  $x(0) = y(0) = z(0) = 0$ ;

$$\begin{aligned}
 \dot{x} &= z + t \\
 3\dot{y} &= 2x - 3z \\
 \dot{z} &= -3y
 \end{aligned}$$

taking Laplace transforms of the equations, substituting initial values

$$\begin{aligned}
sX(s) - x(0) &= Z(s) + \frac{1}{s^2} && \Rightarrow \\
sX(s) - Z(s) &= \frac{1}{s^2} \\
3(sY(s) - y(0)) &= 2X(s) - 3Z(s) && \Rightarrow \\
-\frac{2}{3}X(s) + sY(s) + Z(s) &= 0 \\
sZ(s) - z(0) &= -3Y(s) && \Rightarrow \\
3Y(s) + sZ(s) &= 0
\end{aligned}$$

writing this in matrix form and solving with inverse

$$\begin{aligned}
\begin{bmatrix} s & 0 & 1 \\ -\frac{2}{3} & s & 1 \\ 0 & 3 & s \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} &= \begin{bmatrix} \frac{1}{s^2} \\ 0 \\ 0 \end{bmatrix} \\
\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} &= \frac{1}{(s-1)^2(s+2)} \begin{bmatrix} s^2-3 & -3 & s \\ \frac{2s}{3} & s^2 & -\frac{3s-2}{3} \\ -2 & -3s & s^2 \end{bmatrix} \begin{bmatrix} \frac{1}{s^2} \\ 0 \\ 0 \end{bmatrix} \\
&= \frac{1}{(s-1)^2(s+2)} \begin{bmatrix} \frac{s^2-3}{s^2} \\ \frac{2}{3} \\ -\frac{2}{s^2} \end{bmatrix} \\
&= \begin{bmatrix} -\frac{\frac{2}{3}}{(s-1)^2} + \frac{\frac{20}{9}}{s-1} + \frac{\frac{1}{36}}{s+2} - \frac{\frac{3}{2}}{s^2} - \frac{\frac{9}{4}}{s} \\ \frac{\frac{2}{9}}{(s-1)^2} - \frac{\frac{8}{27}}{s-1} - \frac{\frac{1}{27}}{s+2} + \frac{\frac{1}{3}}{s} \\ -\frac{\frac{2}{3}}{(s-1)^2} + \frac{\frac{14}{9}}{s-1} - \frac{\frac{1}{18}}{s+2} - \frac{1}{s^2} - \frac{\frac{3}{2}}{s} \end{bmatrix} \\
x(t) &= -\frac{2te^t}{3} + \frac{20e^2}{9} + \frac{e^{-2t}}{36} - \frac{3t}{2} - \frac{9}{4} \\
y(t) &= \frac{2te^t}{9} - \frac{8e^2}{27} + \frac{e^{-2t}}{27} - \frac{1}{3} \\
z(t) &= -\frac{2te^t}{3} + \frac{14e^2}{9} + \frac{e^{-2t}}{18} - t - \frac{3}{2}
\end{aligned}$$

## Convolutions

The convolution of functions  $f, g : \mathbb{R}^+ \rightarrow \mathbb{R}$  is denoted  $f * g$ , defined as;

$$\begin{aligned}
(f * g)(t) &= \int_0^t f(u)g(t-u) \, du \\
&= \int_0^t f(t-u)g(u) \, du \\
L(s) &= \mathcal{L}(f * g)(s) \\
&= \int_{t=0}^{t=\infty} e^{-st} \left( \int_{u=0}^{u=\infty} f(u)g(t-u) \, du \right) dt
\end{aligned}$$

we know  $t - u > 0$ , therefore  $u < t$ , where  $u, t \in \mathbb{R}^+$ , now change order of integration

$$= \int_{u=0}^{u=\infty} \int_{t=u}^{t=\infty} e^{-st} f(u)g(t-u) \, du \, dt$$

change  $t$  to  $v = t - u$

$$\begin{aligned}
&= \int_{u=0}^{u=\infty} \int_{v=0}^{v=\infty} e^{-s(u+v)} f(u)g(v) \, du \, dv \\
&= \int_{u=0}^{u=\infty} e^{-su} f(u) \, du \int_{v=0}^{v=\infty} e^{-sv} g(v) \, dv \\
&= F(s)G(s)
\end{aligned}$$

## Inversion

Since we know that the Laplace transform has a unique inverse, we can often solve by inspection with the use of linearity. It can be solved with an inversion formula such as the Bromwich integral (not used in our course);

$$f(t) = \frac{1}{2\pi i} \int_{\gamma-i\infty}^{\gamma+i\infty} e^{st} F(s) \, ds$$

Where  $\gamma$  is a real greater than all the real parts of the singularities of  $F(s)$  (where  $F(s)$  goes to infinity). For example, a singularity of  $F$  is  $-\lambda$  when

$$F(s) = \frac{1}{s + \lambda}$$

This inversion formula comes from the Laplace transform being a special case of the Fourier transform.

## Dominant and Co-dominant Spaces

Suppose we have diagonalisable  $\mathbf{M} \in \mathbb{C}^{n \times n}$ , with distinct eigenvalues of moduli  $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$ , with corresponding eigenvectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ . Let the subspaces  $T_k$  (dominant) and  $U_k$  (co-dominant) be spanned by the bases  $(\mathbf{v}_1, \dots, \mathbf{v}_k)$  (first  $k$  dominant eigenvectors) and  $(\mathbf{v}_{k+1}, \dots, \mathbf{v}_n)$ . They are both invariant under  $\mathbf{A}$ .

Let  $S$  be a  $k$ -dimensional subspace of  $\mathbb{C}^n$  such that  $S \cap U_k = \{\mathbf{0}\}$ . Then  $\mathbf{A}^m S \rightarrow T_k$  as  $m \rightarrow \infty$ . In the case  $k = 1$ , we have the same result as the power-iteration.

To prove this generally, we take a  $\mathbf{v} \in S$ , where  $\mathbf{v} \neq \mathbf{0}$ ;

$$\mathbf{v} = \sum_{i=1}^n c_i \mathbf{v}_i = \underbrace{\sum_{i=1}^k c_i \mathbf{v}_i}_{\in T_k} + \underbrace{\sum_{i=k+1}^n c_i \mathbf{v}_i}_{\in U_k}$$

Since we know  $\mathbf{v} \neq \mathbf{0}$ , there has to be at least one  $c_i \neq 0$ , for  $1 \leq i \leq k$ , otherwise it would be in  $U_k$ . Applying  $\mathbf{A}$   $m$  times, we see the growth due to the eigenvalues as;

$$\mathbf{A}^m \mathbf{v} = \sum_{i=1}^k c_i \lambda_i^m \mathbf{v}_i + \sum_{i=k+1}^n c_i \lambda_i^m \mathbf{v}_i \Rightarrow \frac{\mathbf{A}^m \mathbf{v}}{\lambda_k^m} = \sum_{i=1}^k c_i \frac{\lambda_i}{\lambda_k}^m \mathbf{v}_i + \sum_{i=k+1}^n c_i \frac{\lambda_i}{\lambda_k}^m \mathbf{v}_i$$

However, since  $\lambda_i < \lambda_k$  if  $i > k$ , the component in  $U_k$  approaches  $\mathbf{0}$ , when  $k \rightarrow \infty$ . As proven, there exists at least one non-zero  $c_i$  in the  $T_k$  component, hence  $\mathbf{A}^m \mathbf{v}$  approaches a vector in  $T_k$ .

Let  $\mathbf{Q} = [\mathbf{Q}_1, \mathbf{Q}_2]$ , where  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  is formed by appending  $\mathbf{Q}_1 \in \mathbb{C}^{n \times k}$  and  $\mathbf{Q}_2 \in \mathbb{C}^{n \times (n-k)}$ . Suppose the columns of  $\mathbf{Q}_1$  form an orthonormal basis for  $T_k$ , and similar for  $\mathbf{Q}_2$  for  $U_k$ , therefore all the columns form an orthonormal basis for  $\mathbb{C}^n$ . We propose the following;

$$\begin{aligned}
\mathbf{Q}^\top \mathbf{A} \mathbf{Q} &= \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{0} & \mathbf{A}_{2,2} \end{bmatrix} \\
\mathbf{A}_{1,1} &= \mathbf{Q}_1^\top \mathbf{A} \mathbf{Q}_1 && \in \mathbb{C}^{k \times k} \\
\mathbf{A}_{1,2} &= \mathbf{Q}_1^\top \mathbf{A} \mathbf{Q}_2 && \in \mathbb{C}^{k \times (n-k)} \\
\mathbf{A}_{2,2} &= \mathbf{Q}_2^\top \mathbf{A} \mathbf{Q}_2 && \in \mathbb{C}^{(n-k) \times (n-k)}
\end{aligned}$$

## Subspace Convergence

Let  $\mathbf{P}(\mathbf{m}) = [\mathbf{P}_1, \mathbf{P}_2]$ ,  $\mathbf{P}(\mathbf{m}) \in \mathbb{C}^{n \times n}$  is formed by appending  $\mathbf{P}_1 \in \mathbb{C}^{n \times k}$  and  $\mathbf{P}_2 \in \mathbb{C}^{n \times (n-k)}$ .  $\mathbf{P}(\mathbf{m})$  is a unitary matrix, which is an orthogonal matrix, but allows for complex numbers. The range of  $\mathbf{P}_1$  is  $\mathbf{A}^m S_k$ . We propose that  $\mathbf{B}_{2,1} \rightarrow \mathbf{0}$ , as  $m \rightarrow \infty$ , in the following;

$$\mathbf{P}(\mathbf{m})^\top \mathbf{A} \mathbf{P}(\mathbf{m}) = \begin{bmatrix} \mathbf{B}_{1,1} & \mathbf{B}_{1,2} \\ \mathbf{B}_{2,1} & \mathbf{B}_{2,2} \end{bmatrix}$$

## Subspace Iteration

Take a subspace  $S$ , which has no null vectors of  $\mathbf{A}$ . Suppose  $(\mathbf{q}_1^{(0)}, \dots, \mathbf{q}_k^{(0)})$  is an orthonormal basis of  $S$ . Then clearly  $(\mathbf{A}^m \mathbf{q}_1^{(0)}, \dots, \mathbf{A}^m \mathbf{q}_k^{(0)})$  is a basis of  $\mathbf{A}^m S$ . However, with results proven earlier, all vectors with  $\mathbf{v}_1$ -components (the dominant eigenvector of  $\mathbf{A}$ ) will tend to the dominant subspace  $\langle \mathbf{v}_1 \rangle$ , as  $m \rightarrow \infty$ . These then start to point in the same direction, thus causing the matrix to become ill-conditioned. A solution to this is to orthonormalise at each iteration of the algorithm;

- given  $(\mathbf{q}_1^{(m)}, \dots, \mathbf{q}_k^{(m)})$  as a basis of  $S$ , calculate  $\mathbf{A} \mathbf{q}_1^{(m)}, \dots, \mathbf{A} \mathbf{q}_k^{(m)}$
- choose  $(\mathbf{q}_1^{(m+1)}, \dots, \mathbf{q}_k^{(m+1)})$  as a basis of  $\mathbf{A}^{m+1} S$  by orthonormalising  $\mathbf{A} \mathbf{q}_1^{(m)}, \dots, \mathbf{A} \mathbf{q}_k^{(m)}$  left to right, with something such as Gram-Schmidt.

This iteration holds for all subspace pairs.

## Triangulation Property

We aren't expected to remember this (I hope, since I didn't write most of it). This uses the proposition in the **Subspace Convergence** section above. Since this is true for all  $k$  simultaneously, the limiting form is upper triangular.

## 4th March 2020

### Fourier Transforms

Laplace transforms are only defined for functions on the positive reals, as  $e^{-st}$  gets very large when  $t \rightarrow -\infty$ , for  $s > 0$ . We denote the Fourier Transform of  $f$  as  $\tilde{f} = \mathcal{F}f$

$$\tilde{f}(\omega) = \int_{-\infty}^{\infty} e^{-i\omega x} f(x) dx$$

Which has the inverse transform

$$f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\omega x} \tilde{f}(\omega) d\omega$$

### Dirac Delta Function

Also referred to as  $\delta$ -function, which is a generalised function (or distribution) that is "undefined" at 0, and 0 everywhere else. It has an integral of 1 across the entire real line, and one can consider there to be an infinite spike at the origin.

$$\int_{-\infty}^{\infty} \delta(x) dx = 1$$

The integral of a function multiplied by  $\delta$ -function is  $f(0)$ , as follows (not rigorous);

$$\begin{aligned} \int_{-\infty}^{\infty} f(x) \delta(x) dx &= \int_{-\epsilon}^{\epsilon} f(x) \delta(x) dx && \forall \epsilon > 0 \\ &\approx f(0) \int_{-\epsilon}^{\epsilon} \delta(x) dx && \text{at small } \epsilon \\ &= f(0) \end{aligned}$$

A normalised (the integral over the real line is 1) Gaussian function centred on the origin, with a positive constant shape parameter  $a$  takes the form

$$f(x) = \frac{1}{a\sqrt{\pi}} e^{-\frac{x^2}{a^2}}$$

Note that in the distribution  $N(0, \sigma^2)$ , we have  $a = \sqrt{2\sigma}$ . As the variance goes down towards 0, we get a "sharper" spike at the origin, which is one way we looked at the delta function;

$$\delta(x) = \lim_{a \rightarrow 0} \frac{1}{a\sqrt{\pi}} e^{-\frac{x^2}{a^2}}$$

We can define  $a = \frac{1}{\sqrt{n}}$ , and have the following as the limit of the delta-sequence  $(\delta_n(x))$ ;

$$\delta(x) = \lim_{n \rightarrow \infty} \delta_n(x) = \lim_{n \rightarrow \infty} \sqrt{\frac{n}{\pi}} e^{-nx^2}$$

We can rigorously prove the following;

$$\int_{-\infty}^{\infty} \delta(x-a) f(x) \, dx = \lim_{n \rightarrow \infty} \int_{-\infty}^{\infty} \delta_n(x-a) f(x) \, dx = f(a)$$

### Fourier Inversion Theorem

If  $f$  and  $\tilde{f}$  are integrable, and  $f$  is continuous

$$f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{i\omega(x-y)} f(y) \, dy \, d\omega$$

Using the above property (of  $\delta$ ), it is sufficient to show

$$\delta(y-x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\omega(x-y)} \, d\omega$$

Letting  $z = y - x$ , we can just prove the following claim  $\forall z \in \mathbb{R}$ ;

$$2\pi\delta(z) = \int_{-\infty}^{\infty} e^{-i\omega z} \, d\omega$$

The proof is not examinable, but is at around Panopto timestamp 46:00.

### Common Fourier Transforms

$f$	$\tilde{f} = \mathcal{F}f$	comment
1	$2\pi\delta(\omega)$	Heaviside function
$x^n$	$2\pi i^n \delta^{(n)}(\omega)$	used for polynomials
$e^{iax}$	$2\pi\delta(\omega - a)$	complex exponential
$e^{iax} f(x)$	$\tilde{f}(\omega - a)$	shift
$\sin ax$	$-i\pi(\delta(\omega - a) - \delta(\omega + a))$	$\sin ax = \frac{e^{iax} - e^{-iax}}{2i}$
$\cos ax$	$i\pi(\delta(\omega - a) + \delta(\omega + a))$	$\cos ax = \frac{e^{iax} + e^{-iax}}{2}$

## 5th March 2020

### Fourier Series

Polynomial approximations (such as a Taylor series) is not suitable for a period function  $f(x)$ , as polynomials themselves aren't periodic. Let function  $f(x)$  have a period  $P$ . We can then approximate it with a mixture of sinusoids. The base frequency is  $\omega_0$ , defined as;

$$\omega_0 = \frac{2\pi}{P}$$

In general, we want to approximate  $f(x)$  as a weighted mixture of harmonics, summed over  $n$ , where  $n = 1, 2, \dots$ ;

$$\sin \frac{2\pi}{P}nt$$

$$\cos \frac{2\pi}{P}nt$$

The result  $e^{i\theta} = \cos \theta + i \sin \theta$  is also used. This can be written in three equivalent forms;

$$f(x) \approx f^N(x)$$

$$f^N(x) = \sum_{k=-N}^N a_k e^{i\omega_k x} \quad \text{or}$$

$$= \sum_{k=-N}^N c_k \cos \omega_k x + s_k \sin \omega_k x \quad \text{or}$$

$$= \sum_{k=-N}^N d_k \cos(\omega_k x + \phi_k)$$

$$\omega_k = \frac{2\pi k}{P} \quad a_k, c_k, d_k, s_k, \phi_k \text{ are constants}$$

Using the common results for Fourier transforms, we can observe that the Fourier transform of  $f^N$  is;

$$\mathcal{F} \left( \sum_{k=-N}^N a_k e^{i\omega_k x} \right) = \sum_{k=-N}^N a_k \mathcal{F}(e^{i\omega_k x}) = \sum_{k=-N}^N a_k 2\pi \delta(\omega - \omega_k) = 2\pi \sum_{k=-N}^N a_k \delta \left( \omega - \frac{2\pi k}{P} \right)$$

This is known as the spectrum of  $f$ . We mostly assume  $x_0 = 0$ , without loss of generality.

Note that  $\delta_{k,n}$  is 1 if  $k = n$ , and 0 otherwise.

$$\begin{aligned} \int_0^P e^{i(\omega_k - \omega_n)x} dx &= \int_0^P e^{i\left(\frac{2\pi k}{P} - \frac{2\pi n}{P}\right)x} dx && \text{assuming } k \neq n \\ &= \int_0^P e^{\frac{2\pi}{P}(k-n)xi} dx \\ &= \left[ \frac{P}{2\pi(k-n)} e^{\frac{2\pi}{P}(k-n)xi} \right]_0^P \\ &= \frac{P}{2\pi(k-n)} \left[ e^{\frac{2\pi}{P}(k-n)xi} \right]_0^P \\ &= \frac{P}{2\pi(k-n)} (e^{2\pi i(k-n)} - 1) \\ &= \frac{P}{2\pi(k-n)} (1 - 1) \\ &= 0 \\ \int_0^P e^{i(\omega_k - \omega_n)x} dx &= \int_0^P 1 dx && \text{assuming } k = n \\ &= P \\ \int_0^P e^{i(\omega_k - \omega_n)x} dx &= P\delta_{k,n} && \text{generally} \end{aligned}$$

This lemma can now be used to obtain a specific value for  $a_n$  as follows;

$$\begin{aligned}
f^N(x) &= \sum_{k=-N}^N a_k e^{i\omega_k x} && \Rightarrow \\
f^N(x) e^{-i\omega_n x} &= \sum_{k=-N}^N a_k e^{i\omega_k x} e^{-i\omega_n x} \\
&= \sum_{k=-N}^N a_k e^{i(\omega_k - \omega_n)x} && \Rightarrow \\
\int_0^P f^N(x) e^{-i\omega_n x} dx &= \sum_{k=-N}^N a_k \int_0^P e^{i(\omega_k - \omega_n)x} dx \\
&= \sum_{k=-N}^N a_k P \delta_{k,n} \\
&= 0 + \cdots + 0 + a_n P + 0 + \cdots + 0 \\
&= a_n P && \Rightarrow \\
a_n &= \frac{1}{P} \int_0^P f^N(x) e^{-i\omega_n x} dx \\
&= \frac{1}{P} \int_0^P f^N(x) e^{-i\omega_n x} dx && \text{if } \lim_{N \rightarrow \infty} f^N(x) = f(x)
\end{aligned}$$

For an example, let us look at the sawtooth waveform;  $f(x) = \alpha x$ , for  $0 \leq x < 1$ , and extended periodically throughout the entire real line.

$$\begin{aligned}
f(x) &= \alpha x \\
P &= 1 \\
a_0 &= \frac{\alpha}{2} \\
a_n &= \frac{1}{P} \int_0^P f(x) e^{-i\omega_n x} dx && \text{if it converges to } f \\
&= \frac{1}{1} \int_0^1 \alpha x e^{-i\frac{2\pi n}{1}x} dx \\
&= \alpha \int_0^1 x e^{-2\pi n i x} dx \\
&= \alpha \left( \left[ -\frac{1}{2\pi n i} x e^{-2\pi n i x} \right]_0^1 - \int_0^1 -\frac{e^{-2\pi n i x}}{2\pi n i} dx \right) \\
&= -\frac{\alpha}{2\pi n i} \left( \left[ x e^{-2\pi n i x} \right]_0^1 - \underbrace{\int_0^1 e^{-2\pi n i x} dx}_{=0 \text{ since } n \neq 0} \right) \\
&= -\frac{\alpha}{2\pi n i}
\end{aligned}$$

Actually applying this, we get the following results;

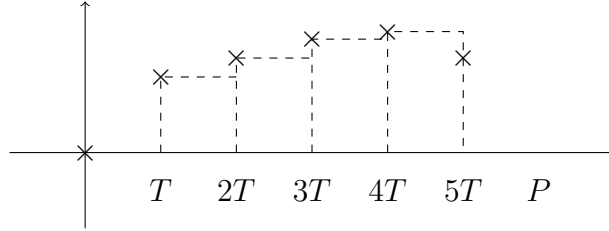
$$\begin{aligned}
f^N(x) &= \underbrace{\frac{\alpha}{2}}_{n=0} - \frac{\alpha}{2\pi i} \underbrace{\sum_{n=1}^N \left( \frac{1}{n} e^{i(2\pi n)x} - \frac{1}{n} e^{-i(2\pi n)x} \right)}_{n \neq 0} \\
&= \frac{\alpha}{2} - \frac{\alpha}{2\pi i} \sum_{n=1}^N \frac{2i}{n} \left( \frac{e^{i(2\pi n)x} - e^{-i(2\pi n)x}}{2i} \right) \\
&= \frac{\alpha}{2} - \frac{\alpha}{2\pi i} \sum_{n=1}^N \frac{2i}{n} \sin(2\pi nx) \\
&= \frac{\alpha}{2} - \frac{\alpha}{\pi} \sum_{n=1}^N \frac{1}{n} \sin(2\pi nx)
\end{aligned}$$

## Sampling and Discrete Fourier Transform

Suppose the discrete function  $f(x)$  is sampled at equi-spaced points  $\{0, T, 2T, \dots, (N-1)T\}$ , where  $T$  divides  $P$ , and  $N = \frac{P}{T}$ . This is equivalent to considering the sampled function  $f_s(x)$ , defined as;

$$\begin{aligned}
f_s(x) &= \Psi_T(x) f(x) \\
\Psi_T(x) &= T \sum_{k=-\infty}^{\infty} \delta(x - kT)
\end{aligned}$$

This approximates rectangles of area  $Tf(x)$  at points where  $x$  is a multiple of  $T$ .



The Fourier coefficients of  $f_s(x)$  then become;

$$\begin{aligned}
F_k &= \frac{1}{P} \int_0^P f_s(x) e^{-i\omega_k x} dx \\
&= \frac{1}{N} \sum_{n=0}^{N-1} f_n e^{-i\omega_k nT} \\
\omega_k &= \frac{2\pi k}{P} \\
&= \frac{2\pi k}{NT} \\
f_n &\equiv f(nT)
\end{aligned}$$

It has an inverse in the Fourier series of  $f(x)$  for  $x = nT$ ;

$$f_n = \sum_{k=0}^{N-1} F_k e^{i\omega_k nT}$$

We often assume  $T = 1$ , by taking some arbitrary scaling of time.

## JPEG and Optics

No.



11th March 2020

## Functions of Several Variables

We write functions with  $n > 1$  variables as  $f(x_1, \dots, x_n)$ . Let  $C^n$  be the set of functions differentiable  $n$  times.  $\nabla$  is the vector differential operator, such that

$$\nabla f(\mathbf{x}) = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$$

For example, in  $n = 2$ , let  $\mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}$ , then we have

$$\mathbf{h} \cdot \nabla f(\mathbf{x}) = h_1 \frac{\partial f}{\partial x_1} + h_2 \frac{\partial f}{\partial x_2}$$

This is used in the theorem for the Taylor series in  $n$ -dimensions. It states that if all the partial derivatives of the function  $f(x_1, \dots, x_n)$  exist in a region around  $x_1^{(0)}, \dots, x_n^{(0)}$  for small enough  $\mathbf{h}$

$$\begin{aligned} \mathbf{h} &= \begin{bmatrix} h_1 \\ \vdots \\ h_n \end{bmatrix} \\ f(x_1^{(0)} + h_1, \dots, x_n^{(0)} + h_n) &= f(x_1^{(0)}, \dots, x_n^{(0)}) \\ &+ \mathbf{h} \cdot \nabla f(x_1^{(0)}, \dots, x_n^{(0)}) \\ &+ \frac{1}{2!} (\mathbf{h} \cdot \nabla)^2 f(x_1^{(0)}, \dots, x_n^{(0)}) \\ &\vdots \\ &+ \frac{1}{n!} (\mathbf{h} \cdot \nabla)^n f(x_1^{(0)}, \dots, x_n^{(0)}) \\ &+ \dots \end{aligned}$$

The proof is as follows;

$$\begin{aligned} g(t) &= f(\mathbf{x}^{(0)} + t\mathbf{h}) \\ g(1) &= g(0) + g'(0) + \dots + \frac{1}{n!} g^{(n)}(0) + \dots \end{aligned}$$

brief example for chain rule, if vector components are functions of  $t$

$$\begin{aligned} \mathbf{y} &= \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \\ \frac{d\mathbf{y}}{dt} \cdot \nabla &= \begin{bmatrix} \frac{dy_1}{dt} \\ \frac{dy_2}{dt} \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial}{\partial y_1} \\ \frac{\partial}{\partial y_2} \end{bmatrix} \\ &= \frac{dy_1}{dt} \frac{\partial}{\partial y_1} + \frac{dy_2}{dt} \frac{\partial}{\partial y_2} \end{aligned}$$

using the general case, and continuing

$$\begin{aligned} \frac{df(\mathbf{y})}{dt} &= \sum_{i=1}^n \frac{\partial f(\mathbf{y})}{\partial y_i} \frac{dy_i}{dt} \\ &= \left( \frac{d\mathbf{y}}{dt} \cdot \nabla \right) f(\mathbf{y}) \\ \frac{d^n f(\mathbf{y})}{dt^n} &= \underbrace{\left( \frac{d\mathbf{y}}{dt} \cdot \nabla \right) \dots \left( \frac{d\mathbf{y}}{dt} \cdot \nabla \right)}_{n \text{ times}} f(\mathbf{y}) \\ &= \left( \frac{d\mathbf{y}}{dt} \cdot \nabla \right)^n f(\mathbf{y}) \end{aligned}$$

chain rule, if components are functions of  $t$

setting  $\mathbf{y} = \mathbf{x}^{(0)} + t\mathbf{h}$

$$\begin{aligned}\frac{d^n g}{dt^n} &= \frac{d^n f(\mathbf{y})}{dt^n} \\ &= (\mathbf{h} \cdot \nabla)^n f(\mathbf{y})\end{aligned}$$

at  $t = 0$ , we have  $\mathbf{y} = \mathbf{x}^{(0)}$

$$g^{(n)}(0) = (\mathbf{h} \cdot \nabla)^n f(\mathbf{x}^{(0)})$$

$$\text{therefore } g(1) = \sum_{n=0}^{\infty} \frac{(\mathbf{h} \cdot \nabla)^n f(\mathbf{x}^{(0)})}{n!}$$

## Stationary Points

Similar to single variable functions, our criteria for a stationary point in  $\nabla f(s_1, \dots, s_n) = \mathbf{0}$ , for some stationary point  $\mathbf{s}$ . Note that the slides (and I assume notes) use  $\mathbf{x}_0$  instead of  $\mathbf{s}$ , but I feel it can be confusing. This is the same as stating

$$\frac{\partial f}{\partial x_1}(s_1, \dots, s_n) = \dots = \frac{\partial f}{\partial x_n}(s_1, \dots, s_n) = 0$$

In the case where we have one variable, hence  $n = 1$ , there are 3 types of stationary point - maximum, minimum, and inflexion. However, in the multi-variable case, for  $n > 1$ , each variable can be in any of the three, therefore there are 9 possible combinations in two dimensions. This becomes exponentially more complex in more dimensions, which is why multi-variable optimisation is difficult.

In order to classify stationary points, we can use the **Hessian matrix**, defined as such;

$$\begin{aligned}\partial_i &= \frac{\partial}{\partial x_i} \\ f_{i,j} &= \partial_i \partial_j & f_{i,j} &= f_{j,i} \text{ for } f \in C^2 \\ &= \frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} f \\ \mathbf{H} &= \begin{bmatrix} f_{1,1} & f_{1,2} & \cdots & f_{1,n} \\ f_{2,1} & f_{2,2} & \cdots & f_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ f_{n,1} & f_{n,2} & \cdots & f_{n,n} \end{bmatrix}\end{aligned}$$

If  $\mathbf{H}$  is positive definite, meaning it has all positive eigenvalues and a non-zero determinant, it gives a maximum point. Similarly, if  $\mathbf{H}$  is negative definite, it gives a maximum.

## 12th March 2020

### Quadratic Form

The quadratic form is as follows;

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x} + c$$

We can obtain the derivative of this as follows;

$$\begin{aligned}\frac{\partial f}{\partial x_j} &= \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \\ &= \delta_{i,j} \\ \sum_{i=1}^n x_i \delta_{i,j} &= x_j \\ \nabla f &= \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)\end{aligned}$$

take arbitrary  $j$

$$\begin{aligned}
\frac{\partial f}{\partial x_j} &= \frac{\partial}{\partial x_j} \sum_{i,k} \frac{1}{2} x_i \mathbf{A}_{i,k} x_k - \frac{\partial}{\partial x_j} \sum_i b_i x_i + 0 \\
&= \sum_{i,k} \frac{1}{2} \mathbf{A}_{i,k} \left( \frac{\partial x_i}{\partial x_j} x_k + \frac{\partial x_k}{\partial x_j} x_i \right) - \sum_i \frac{\partial x_i}{\partial x_j} b_i \\
&= \frac{1}{2} \left( \sum_k \mathbf{A}_{j,k} x_k + \sum_i x_i \mathbf{A}_{i,j} \right) - b_j \\
&= \frac{1}{2} \left( \sum_k \mathbf{A}_{j,k} x_k + \sum_i \mathbf{A}_{j,i}^\top x_i \right) - b_j \\
&= \frac{1}{2} (\mathbf{A} \mathbf{x})_j + (\mathbf{A}^\top \mathbf{x})_j - b_j \quad \Rightarrow \\
\nabla f &= \frac{1}{2} (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} - \mathbf{b}
\end{aligned}$$

If  $\mathbf{A}$  is symmetric, you have  $\mathbf{A} = \mathbf{A}^\top$ , so  $\nabla f = \mathbf{A} \mathbf{x} - \mathbf{b}$ . If it is also positive definite, then the solution to  $\mathbf{A} \mathbf{x} = \mathbf{b}$  globally minimises  $f(\mathbf{x})$ . Similarly, if  $f(\mathbf{x})$  is minimised by  $\mathbf{x} = \mathbf{y}$ , then  $\mathbf{y}$  is the solution to  $\mathbf{A} \mathbf{x} = \mathbf{b}$ , therefore there exists a duality between these two.

### Iterative Solutions for Function Minimisation

We start by picking some initial vector  $\mathbf{x}_0$  as a guess. At each step  $k = 1, 2, \dots$ , we choose a direction vector  $\mathbf{d}_k$  and a step size  $\alpha_k$ . The iteration is

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

Some steepest descent method (SDM) is as follows;

This method chooses the direction in which  $f(\mathbf{x}_k)$  decreases the fastest, hence we choose  $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$ . When  $\mathbf{A}$  is symmetric,  $\mathbf{d}_k = \mathbf{b} - \mathbf{A} \mathbf{x}_k$ , in the quadratic form example, we have  $\mathbf{d}_k$  as the residual vector.

To choose the step length, we want to minimise  $f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)$ ;

$$\begin{aligned}
\text{let } \mathbf{x} &= \mathbf{x}_k \\
\frac{df}{d\alpha_k} &= \sum_i \frac{\partial f(\mathbf{x})}{\partial x_i} \frac{dx_i}{d\alpha} \\
&= \sum_i (\nabla f(\mathbf{x}_{k+1}))_i (\mathbf{d}_k)_i \\
&= \nabla f(\mathbf{x}_{k+1}) \cdot \mathbf{d}_k \\
&= \mathbf{d}_{k+1} \cdot \mathbf{d}_k
\end{aligned}$$

To minimise this, set it to zero, therefore we want to choose  $\alpha$  such that  $\mathbf{d}_{k+1}$  is orthogonal to  $\mathbf{d}_k$ . Therefore in general we want to solve  $(\nabla f(\mathbf{x}_k)) \cdot (\nabla f(\mathbf{x}_{k+1})) = 0$

An example is as follows (quadratic form);

$$\begin{aligned}
\mathbf{r}_k^\top (\mathbf{b} - \mathbf{A}(\mathbf{x}_k + \alpha_k \mathbf{r}_k)) &= \mathbf{r}_k^\top (\mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{r}_k) \\
&= 0 \quad \text{(solve)} \\
\alpha_k &= \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{A} \mathbf{r}_k} \\
&= \frac{\|\mathbf{r}_k\|^2}{\mathbf{r}_k^\top \mathbf{A} \mathbf{r}_k}
\end{aligned}$$

By the following iteration, we approximate a solution for  $\mathbf{x}$ ;

$$\begin{aligned}\mathbf{x}_k &= \mathbf{A}\mathbf{r}_k \\ \alpha_k &= \frac{\|\mathbf{r}_k\|^2}{\mathbf{r}_k^\top \mathbf{v}_k} \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{r}_k \\ \mathbf{r}_{k+1} &= \mathbf{r}_k - \alpha_k \mathbf{v}_k\end{aligned}$$

This takes a single matrix-vector multiplication per iteration, and one more in the initialisation.

## Example of Iterative Method

Suppose we want to minimise the function  $f(x, y) = x^2 + y^2$ . By inspection, we know the minimum, for real  $x, y$ , is  $(0, 0)$ . However, numerically, we do the following;

guess $\mathbf{x}_0 = (x_0, y_0)$	minimum
$\mathbf{d}_0 = -\nabla f(x_0, y_0)$	
$= -(2x, 2y) _{\mathbf{x}=\mathbf{x}_0}$	
$= -2(x_0, y_0)$	
$= -2\mathbf{x}_0$	
$x_1 = x_0 + \alpha(-2x_0)$	
$= (1 - 2\alpha)x_0$	
$y_1 = y_0 + \alpha(-2y_0)$	
$= (1 - 2\alpha)y_0$	
$\mathbf{d}_1 = -2(x_1, y_1)$	
$= -2\mathbf{x}_1$	
solve $\mathbf{d}_0 \cdot \mathbf{d}_1 = 0$	$\Leftrightarrow$
$x_0x_1 + y_0y_1 = 0$	$\Leftrightarrow$
$(1 - 2\alpha)x_0^2 + (1 - 2\alpha)y_0^2 = 0$	$\Leftrightarrow$
$(1 - 2\alpha)(x_0^2 + y_0^2) = 0$	$\Rightarrow$
$\alpha = \frac{1}{2}$	

after this computation, we have the next step, so;

$$\begin{aligned}\mathbf{d}_0 &= -2(x_0, y_0) \\ &= -2\mathbf{x}_0 \\ \mathbf{x}_1 &= \mathbf{x}_0 + \alpha\mathbf{d}_0 \\ &= \mathbf{x}_0 + \frac{1}{2}(-2\mathbf{x}_0) \\ &= \mathbf{0} \\ &= (0, 0)\end{aligned}$$

This algorithm converged in a single step, due to the symmetry (hence no "zig-zagging" needed).