

# Semester I

# Subject 1

## ME010101 Abstract Algebra

### 1.1 Introduction to Abstract Algebra

**Definitions 1.1.1.** Set-theoretical foundation of Abstract Algebra

- A **cartesian product**,  $A \times B = \{(a, b) : a \in A, b \in B\}$
- A **relation** on a set  $A$  is a subset of  $A \times A$ .
- A **function**  $f$  from  $A$  into  $B$ ,  $f : A \rightarrow B$  is a relation such that *every element of  $A$  is related to some unique element of  $B$*  (well defined).

**Definitions 1.1.2.** Functions : A **binary operation** on a set  $A$  is a function  $* : A \times A \rightarrow A$ .

“A binary operation on a set  $A$  gives an algebra on  $A$ .”

**Abstract Algebra** It is the study of algebraic structures. We are interested in a few algebraic structures : 1. Group 2. Ring 3. Integral Domain 4. Field

**Definitions 1.1.3.** A **binary algebraic structure**  $\langle G, * \rangle$  is a set  $G$  together with a binary operation  $*$ .

**Definitions 1.1.4** (Group). A set  $G$ , closed under a binary relation  $*$  satisfying the following three axioms -G1, G2, & G3 is a group.

G1 Associativity

For any three elements  $a, b, c \in G$ ,  $a * (b * c) = (a * b) * c$ .

G2 Identity element

There exists a unique element  $e \in G$  such that  $a * e = a = e * a$  for any element  $a \in G$ .

G3 Inverse elements

For any element  $a \in G$  there exists a unique element  $a^{-1}$  such that  $a * a^{-1} = e = a^{-1} * a$ .

**Definitions 1.1.5.** Group terminologies :

- A group is **abelian** if the binary operation is commutative.  $a * b = b * a$
- The **order** of a group  $G, *$  is the number of elements in  $G$ .

**Definitions 1.1.6.**  $\mathbb{Z}_n = \{0, 1, 2, \dots, n-1\}$

*Remark.* Consider  $3, 4 \in \mathbb{Z}_5$ ,  $3 * 4 = 2 = 4 * 3$  since  $7 \cong 2 \pmod{5}$ .  $\langle \mathbb{Z}_5, +_5 \rangle$  is an abelian group of order 5.

**Definitions 1.1.7.** Homomorphism & Isomorphism

- A function  $\phi : A \rightarrow B$  is a **homomorphism** if for any two elements  $x, y \in A$ ,  $\phi(xy) = \phi(x)\phi(y)$
- A function  $\phi : A \rightarrow B$  is an **isomorphism** if  $\phi$  is a bijective, homomorphism. If two binary structures are isomorphic, then they have the same (algebraic) structure.

**Definitions 1.1.8.** Subgroup

- A subset  $H$  of a group  $\langle G, * \rangle$  is a **subgroup** of  $G$  if  $H$  is group with the same binary operation  $*$ . And is denoted by  $H \leq G$ .
- $G$  is the **improper** subgroup of  $G$  and every other subgroup is **proper**.S5.5
- $\{e\}$  is the **trivial** subgroup of  $G$  and every other subgroup is **non-trivial**.
- The **subgroup generated by**  $g \in G$  is the subgroup  $\{g^n : n \in \mathbb{Z}\}$ .
- The **order** of an element  $g$  is order of the subgroup generated by  $g$ .
- An element  $g \in G$  is a **generator** of  $G$  if  $g$  generates  $G$ .
- A group is **cyclic** if it has a generator.

*Remark.* Cyclic Groups :

- Cyclic groups are abelian.
- Subgroups of cyclic groups are cyclic.

### 1.1.1 Some Proof Techniques

**Equality of two Sets**  $A = B \iff A \subset B$  and  $B \subset A$

If  $x \in A \implies x \in B$ , then  $A \subset B$

**Uniqueness** Suppose there are two elements that qualify our conditions. We show (using the conditions) that they are the same, that is, unique.

For example,  $3 + a = \pi$ . Suppose  $a = x, y$ . Then  $3 + x = 3 + y \implies x = y$ , provided that the values of  $a$  comes from a set in which left cancelation law can be applied. Then  $a$  is unique.

Remember : We usually don't care to show what this unique element is. It may be also be the case that there is no such element, that is, proof of uniqueness doesn't imply existence.

**Existence** There are constructive and non-constructive proofs for existence problems. Suppose we want to prove that  $a * b$  has an inverse element. We know that  $a, b$  has inverse elements  $a^{-1}, b^{-1}$ . From those elements, we construct an element  $b^{-1} * a^{-1}$  which is an inverse of  $a * b$  by construction.

And we may also prove existence without actually giving an object. Suppose we want to prove that  $x^y \in \mathbb{Q}$  for some irrational numbers  $x$  and  $y$ . We know that  $\sqrt{2} \notin \mathbb{Q}$ . Then,  $\sqrt{2}^{\sqrt{2}}$  is either rational or irrational. Suppose it is irrational, then  $\left(\sqrt{2}^{\sqrt{2}}\right)^{\sqrt{2}} = 2$  is rational. Thus the proof is complete, but we are yet to know whether  $\sqrt{2}^{\sqrt{2}}$  is an irrational or rational.

## 1.2 Direct Products and Finitely Generated Abelian Groups

**Definitions 1.2.1** (Cartesian product of sets). Let  $S_1, S_2, \dots, S_n$  be a sets. Their cartesian product,

$$S_1 \times S_2 \times \dots \times S_n = \prod_{i=1}^n S_i = \{(a_1, a_2, \dots, a_n) : a_i \in S_i\} \quad (1.1)$$

For example,  $\{A, B, C\} \times \{1, 2\} = \{(A, 1), (A, 2), (B, 1), (B, 2), (C, 1), (C, 2)\}$ .

**Question 1.** How many elements in  $\mathbb{Z}_3 \times \mathbb{Z}_{10} \times \mathbb{Z}_9$  ?

**Theorem 1.2.1** (Direct product of Groups). Let  $G_1, G_2, \dots, G_n$  be groups. Then their cartesian product is a group with the binary operation  $*$ ,

$$(a_1, a_2, \dots, a_n) * (b_1, b_2, \dots, b_n) = (a_1 * b_1, a_2 * b_2, \dots, a_n * b_n) \quad (1.2)$$

where the binary operation in  $a_i * b_i$  is the binary operation of the group  $G_i$ .

*Proof.*  $\prod_{i=1}^n G_i$  is a group if it satisfies the group axioms.

G1 Associativity

$$\begin{aligned} & (a_1, a_2, \dots, a_n)((b_1, b_2, \dots, b_n)(c_1, c_2, \dots, c_n)) \\ &= (a_1, a_2, \dots, a_n)(b_1 c_1, b_2 c_2, \dots, b_n c_n) \\ &= (a_1(b_1 c_1), a_2(b_2 c_2), \dots, a_n(b_n c_n)) \\ &= ((a_1 b_1) c_1, (a_2 b_2) c_2, \dots, (a_n b_n) c_n) \\ &= (a_1 b_1, a_2 b_2, \dots, a_n b_n)(c_1, c_2, \dots, c_n) \\ &= ((a_1, a_2, \dots, a_n)(b_1, b_2, \dots, b_n))(c_1, c_2, \dots, c_n) \end{aligned}$$

G2 Existence of a unique identity element in  $\prod_{i=1}^n G_i$

Let  $e_i$  be the identity element in  $G_i$ . Then  $(e_1, e_2, \dots, e_n)$  is the identity

## 1.2. DIRECT PRODUCTS AND FINITELY GENERATED ABELIAN GROUPS 5

element in  $\prod_{i=1}^n G_i$ .

$$\begin{aligned} (a_1, a_2, \dots, a_n)(e_1, e_2, \dots, e_n) \\ = (a_1 e_1, a_2 e_2, \dots, a_n e_n) \\ = (a_1, a_2, \dots, a_n) \end{aligned}$$

G3 Existence of unique inverse element for each element in  $\prod_{i=1}^n G_i$

Let  $(a_1, a_2, \dots, a_n)$  be in  $\prod_{i=1}^n G_i$ . Then it has the inverse element  $(a_1^{-1}, a_2^{-1}, \dots, a_n^{-1})$  in  $\prod_{i=1}^n G_i$ .

$$\begin{aligned} (a_1, a_2, \dots, a_n)(a_1^{-1}, a_2^{-1}, \dots, a_n^{-1}) \\ = (a_1 a_1^{-1}, a_2 a_2^{-1}, \dots, a_n a_n^{-1}) \\ = (e_1, e_2, \dots, e_n) \end{aligned}$$

□

*Remark.* We usually write  $ab$  instead of  $a * b$  and relevant binary operations are used in different contexts. Student should be able to recognise the difference from the context.

*Remark.*  $\mathbb{Z}_n = \{0, 1, \dots, (n-1)\}$  is a group with  $+_n$ . (addition modulo  $n$ )

For example, Consider  $(1, 2) \in \mathbb{Z}_2 \times \mathbb{Z}_3$ . We have,  $(1, 2) + (1, 2) = (0, 1)$  since  $1 +_2 1 = 0$  and  $2 +_3 2 = 1$ .

**Definitions 1.2.2.** Suppose all the groups  $G_i$  are abelian. Then  $\prod_{i=1}^n G_i$  is the direct sum of the groups  $G_i$ . And is represented by  $\oplus_{i=1}^n G_i$ .

**Theorem 1.2.2.**  $\mathbb{Z}_m \times \mathbb{Z}_n$  is cyclic and is isomorphic to  $\mathbb{Z}_{mn}$  if and only if  $m$  and  $n$  are relatively prime.

*Proof.* Sufficient part : Consider the cyclic subgroup<sup>1</sup>  $H$  generated by  $(1, 1) \in \mathbb{Z}_m \times \mathbb{Z}_n$ . It is enough to prove that the order of this cyclic subgroup  $H$  is  $mn$ . The order of  $H$  is the smallest power of  $(1, 1)$  that gives the identity  $(0, 0)$ . The first component gives 0 for multiples of  $m$ . And the second component gives 0 for multiples of  $n$ . Since  $m, n$  are relatively prime,  $mn$  is the smallest power of  $(1, 1)$  that will give  $(0, 0)$ . Thus  $\mathbb{Z}_m \times \mathbb{Z}_n = H$  and is cyclic. Every cyclic group of order  $mn$  is isomorphic to  $\mathbb{Z}_{mn}$ . Therefore,  $\mathbb{Z}_m \times \mathbb{Z}_n \approx \mathbb{Z}_{mn}$ .

Necessary part : Suppose  $\gcd(m, n) = d > 1$ . Then  $mn/d$  is the smallest integer divisible by both  $m$  and  $n$ . Consider  $(r, s) \in \mathbb{Z}_m \times \mathbb{Z}_n$ .  $r$  gives 0 in  $mn/d$  since it is a multiple of  $m$ . Similarly,  $s$  gives 0 in  $mn/d$  since it is a multiple of  $n$ . Thus,  $\frac{mn}{d}(r, s) = (0, 0)$ . And the cyclic group generated by any element of  $\mathbb{Z}_m \times \mathbb{Z}_n$  is a proper subgroup. Therefore  $\mathbb{Z}_m \times \mathbb{Z}_n$  has no generators and it is not cyclic. □

---

<sup>1</sup> $(1, 1) \in \mathbb{Z}_m \times \mathbb{Z}_n$ . The cyclic group generated by  $(1, 1)$  has all its elements in  $\mathbb{Z}_m \times \mathbb{Z}_n$ . And therefore, it is a subgroup of  $\mathbb{Z}_m \times \mathbb{Z}_n$

**Corollary 1.2.2.1.**  $\prod_{i=1}^n \mathbb{Z}_{m_i}$  is cyclic and is isomorphic to  $\mathbb{Z}_{m_1 m_2 \dots m_n}$  if and only if any two of the numbers  $m_i$  are relatively prime.

**Question 2.** Prove : For any non-negative integer  $n$ , there exists a cyclic group of order  $n$ , which is unique upto isomorphism.

**Theorem 1.2.3.** Let  $(a_1, a_2, \dots, a_n) \in \prod_{i=1}^n G_i$ . And  $a_i$  are of finite order  $r_i$  in  $G_i$ . Then the order of  $\prod_{i=1}^n G_i$  is the least common multiple of  $r_i$ s.

*Proof.* Least common multiple of  $r_i$ s is the smallest positive integer  $d$  which is a multiple of all  $r_i$ s. For each  $i$ , the  $r_i$ th multiple of  $a_i$  gives 0 (identity). Thus, the order of the cyclic subgroup generated by  $(a_1, a_2, \dots, a_n)$  is the least common multiple of all the  $r_i$ s.  $\square$

*Remark.* Consider  $(3, 6, 12, 16) \in \mathbb{Z}_4 \times \mathbb{Z}_{12} \times \mathbb{Z}_{20} \times \mathbb{Z}_{24}$ . Order of  $3 \in \mathbb{Z}_4$  is  $4/\gcd(3, 4) = 4$  ie,  $\langle 3 \rangle = \{3, 2, 1, 0\}$   
Order of  $6 \in \mathbb{Z}_{12}$  is  $12/\gcd(6, 12) = 2$  ie,  $\langle 6 \rangle = \{6, 0\}$   
Order of  $12 \in \mathbb{Z}_{20}$  is  $20/\gcd(12, 20) = 5$  ie,  $\langle 12 \rangle = \{12, 4, 16, 8, 0\}$   
Order of  $16 \in \mathbb{Z}_{24}$  is  $24/\gcd(16, 24) = 3$  ie,  $\langle 16 \rangle = \{16, 8, 0\}$   
Order of  $(3, 6, 12, 16)$  is  $\text{lcm}(4, 2, 5, 3) = 2^2 \cdot 3 \cdot 5 = 60$ .

*Remark.* Define  $\overline{G_i} = \{(e_1, e_2, \dots, e_{i-1}, a_i, e_{i+1}, \dots, e_n) : a_i \in G_i\}$ . Then  $G_i \approx \overline{G_i}$ . And  $\prod_{i=1}^n G_i$  is the internal direct product of  $\overline{G_i}$ s.

For example,  $\mathbb{Z}_2 \times \mathbb{Z}_3 \approx (\mathbb{Z}_2 \times \{0\}) \otimes (\{0\} \times \mathbb{Z}_3)$

**Question 3.** Internal direct product form of  $\mathbb{Z}_{12} \times \mathbb{Z}_{60} \times \mathbb{Z}_{24}$  ?

## 1.3 Fundamental Theorem

**Definitions 1.3.1.** A group  $G$  is **finitely generated** if  $G$  has a finite subset that generates  $G$ .

**Theorem 1.3.1** (fundamental theorem of finitely generated abelian groups). *Every finitely generated abelian group  $G$  is isomorphic to a direct product of cyclic groups in the form*

$$\mathbb{Z}_{p_1^{r_1}} \times \mathbb{Z}_{p_2^{r_2}} \times \dots \times \mathbb{Z}_{p_n^{r_n}} \times \mathbb{Z} \times \mathbb{Z} \times \dots \times \mathbb{Z} \quad (1.3)$$

where  $p_i$  are primes, not necessarily distinct and  $r_i$  are positive integers. The direct product is unique, except for the possible rearrangement of the factors.

*Proof.* —proof is omitted—  $\square$

For example,  $G = \mathbb{Z}_{20} \times \mathbb{Z} \times \mathbb{Z}_{15} \times \mathbb{Z} \approx \mathbb{Z}_{2^2} \times \mathbb{Z}_3 \times \mathbb{Z}_5 \times \mathbb{Z}_5 \times \mathbb{Z} \times \mathbb{Z}$ . In the above case, Betti number of  $G$  is 2 (number of  $\mathbb{Z}$  factors). For any finite abelian group, Betti number is 0.

*Remark* (finite abelian groups). Every finite group is finitely generated. And thus we can enumerate finite abelian group of any order.

*Remark.* There are precisely 6 different abelian groups of order  $360 = 2^3 3^2 5$ .

1.  $\mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_3 \times \mathbb{Z}_3 \times \mathbb{Z}_5$
2.  $\mathbb{Z}_2 \times \mathbb{Z}_{2^2} \times \mathbb{Z}_3 \times \mathbb{Z}_3 \times \mathbb{Z}_5$
3.  $\mathbb{Z}_{2^3} \times \mathbb{Z}_3 \times \mathbb{Z}_3 \times \mathbb{Z}_5$
4.  $\mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_{3^2} \times \mathbb{Z}_5$
5.  $\mathbb{Z}_2 \times \mathbb{Z}_{2^2} \times \mathbb{Z}_{3^2} \times \mathbb{Z}_5$
6.  $\mathbb{Z}_{2^3} \times \mathbb{Z}_{3^2} \times \mathbb{Z}_5$

**Question 4.** Group of order 360 with at least an element of order 8

**Definitions 1.3.2.** A group  $G$  is **decomposable** if it is isomorphic to a direct product of two proper, non-trivial subgroups. Otherwise  $G$  is **indecomposable**.

For example,  $\mathbb{Z}_{2^3} \times \mathbb{Z}_3 \times \mathbb{Z}_3 \times \mathbb{Z}_5 \approx \mathbb{Z}_{24} \times \mathbb{Z}_{15}$  is decomposable.

*Remark.*  $G = \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_3 \times \mathbb{Z}_3 \times \mathbb{Z}_5 \approx \mathbb{Z}_2 \times \mathbb{Z}_6 \times \mathbb{Z}_{30}$  is also decomposable. Since  $G \approx (\mathbb{Z}_2 \times \mathbb{Z}_6) \times \mathbb{Z}_{30}$ .

**Theorem 1.3.2.** *The finite indecomposable abelian groups are exactly the cyclic groups with order a power of a prime.*

*Proof.* Necessary part: Let  $G$  be a finite, indecomposable abelian group. By fundamental theorem of finitely generated abelian groups,  $G$  is isomorphic to a direct product of cyclic groups of prime power order.

$$G \approx \mathbb{Z}_{p_1^{r_1}} \times \mathbb{Z}_{p_2^{r_2}} \times \cdots \times \mathbb{Z}_{p_n^{r_n}}$$

Thus for  $G$  to be indecomposable the direct product should be a cyclic group of prime power order.  $G \approx \mathbb{Z}_{p_1^{r_1}}$ .

Sufficient part : Let  $p$  be a prime and  $r$  a non-negative integer. Cyclic group of order  $p^r$  is isomorphic to  $\mathbb{Z}_{p^r}$ . Since every cyclic groups are abelian,  $\mathbb{Z}_{p^r}$  is an abelian group of finite order  $p^r$ . It is enough to prove that  $\mathbb{Z}_{p^r}$  is indecomposable.

A proper, non-trivial subgroup of  $\mathbb{Z}_{p^r}$  is of the form  $\mathbb{Z}_{p^i}$  where  $0 < i < r$ . Suppose  $\mathbb{Z}_{p^r}$  is decomposable. Then  $\exists i, j \in \mathbb{Z}^+$  such that  $\mathbb{Z}_{p^r} \approx \mathbb{Z}_{p^i} \times \mathbb{Z}_{p^j}$  and  $i + j = r$ . Clearly,  $p^i$  and  $p^j$  are not relatively prime, thus  $\mathbb{Z}_{p^r} \not\approx \mathbb{Z}_{p^i} \times \mathbb{Z}_{p^j}$ . Therefore, cyclic groups of order prime power are indecomposable.  $\square$

**Theorem 1.3.3.** *If  $m$  divides the order of a finite abelian group  $G$ , then  $G$  has a subgroup of order  $m$ .*

*Proof.* Let  $G \approx \mathbb{Z}_{p_1^{r_1}} \times \mathbb{Z}_{p_2^{r_2}} \times \cdots \times \mathbb{Z}_{p_n^{r_n}}$ . Then  $|G| = p_1^{r_1} p_2^{r_2} \cdots p_n^{r_n}$ . Suppose  $m$  divides  $|G|$ , then  $m = p_1^{s_1} p_2^{s_2} \cdots p_n^{s_n}$  where  $0 \leq s_i \leq r_i$ . Define  $H = \mathbb{Z}_{p_1^{s_1}} \times \mathbb{Z}_{p_2^{s_2}} \times \cdots \times \mathbb{Z}_{p_n^{s_n}}$ . Then  $H$  is subgroup of order  $m$ .  $\square$

**Theorem 1.3.4.** *If  $m$  is square-free integer, then every abelian group of order  $m$  is cyclic.*

*Proof.* Let  $m$  be a square-free integer and  $G$  be an abelian group of order  $m$ . By fundamental theorem of finitely generated abelian groups

$$G \approx \mathbb{Z}_{p_1^{r_1}} \times \mathbb{Z}_{p_2^{r_2}} \times \cdots \times \mathbb{Z}_{p_n^{r_n}}$$

We have,  $m$  is square-free. Thus  $r_i = 1$  and  $p_i$  are distinct. Therefore,  $G \approx \mathbb{Z}_{p_1 p_2 \cdots p_n}$  is a cyclic group of order  $m$ .  $\square$

## 1.4 Exercises §11

**Question 5.** Enumerate subgroups of  $\mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_4$

### 1.4.1 Abelian Groups

*Remark.* Direct product of abelian groups is abelian.

**Question 6.** Enumerate abelian groups of order 72 ?

*Remark.* Let  $G$  be an abelian group. The subset  $H$  of  $G$  with identity element and all elements of order  $n$  is subgroup of  $G$  if and only if  $n$  is a prime.

*Proof.* Suppose  $a \in H$ . Then every power of  $a$  has order  $n$ . Suppose  $n$  is not prime. Then  $d$  divides  $n$  and  $a^d$  has order  $n/d$ .  $\square$

### 1.4.2 Torsion Group and Torsion Coefficients

*Remark.* If group  $G$  is abelian, then its elements of finite order forms a subgroup. (hint :  $a, b \in G$  has finite order, then  $ab$  has finite order. And  $a \in G$  has finite order, then  $a^{-1}$  has finite order)

**Definitions 1.4.1.** The **torsion group** of an abelian group  $G$  is the subgroup of  $G$  containing only those elements of finite order. An abelian group is **torsion free** if identity element is the only element of finite order.  $G$  is Torsion free if Torsion group of  $G$  is trivial,  $\{e\}$ .

**Definitions 1.4.2.** The integers  $m_1, m_2, \dots, m_n$  are torsion coefficients of  $G$  such that  $G \approx \mathbb{Z}_{m_1} \times \mathbb{Z}_{m_2} \times \dots \times \mathbb{Z}_{m_n}$  where  $m_i$  divides  $m_{i+1}$ .

For example,  $\mathbb{Z}_6 \times \mathbb{Z}_{12} \times \mathbb{Z}_{20}$  has torsion coefficients 2, 12, 60

*Remark* (Algorithm to find torsion coefficients of a group). Suppose  $G$  has a direct product form.

Step 1 Find power of each prime in the direct product form

Step 2 List power of each prime

Step 3 Append 1s on left to make all lists to equal length

Step 4 Product of  $i$ th number on each list gives  $m_i$

For example,  $G \approx \mathbb{Z}_6 \times \mathbb{Z}_{12} \times \mathbb{Z}_{20}$

Step 1  $\mathbb{Z}_6 \times \mathbb{Z}_{12} \times \mathbb{Z}_{20} \approx \mathbb{Z}_2 \times \mathbb{Z}_3 \times \mathbb{Z}_{2^2} \times \mathbb{Z}_3 \times \mathbb{Z}_{2^2} \times \mathbb{Z}_5$

Step 2 (2,4,4), (3,3), (5)

Step 3 (2,4,4), (1,3,3), (1,1,5)

Step 4 (2,12,60)



### 1.4.3 Torsion $p$ -subgroup

“Caution : Torsion  $p$ -subgroup is a name suggested by ‘Jacob’. And is not among the standard terminology in group theory.”

*Remark.* Let  $G$  be a group. Let  $p$  be an integer, ( $p > 1$ ). Then the set of all element of  $G$  of order  $p$  together with the identity element is a group, if  $p$  is a prime.

For example : Consider symmetric group,  $S_3$ . It has three elements of order 2, namely  $\mu_1 = (2, 3)$ ,  $\mu_2 = (1, 3)$  and  $\mu_3 = (1, 2)$ . Clearly, the set of all element of order 2 together with identity ( ) is not a subgroup of  $S_3$  as  $(1, 2)(1, 3) = (3, 1, 2)$  is an element of order 3. Thus by counter-example, for a prime  $p$ , the set of all elements of a non-abelian group  $G$  together with identity is not necessarily a subgroup of  $G$ .

Let  $G$  be an abelian group. Let  $p$  be a prime. Let  $H$  be the set of all element of  $G$  of order  $p$  together with the identity  $e$ . Let  $g, h \in H$ . Clearly,  $g^p = e$  and  $h^p = e$ . For every  $g, h \in G$ ,  $gh \in G$ , since  $G$  is abelian  $(gh)^p = g^p h^p = e$ . Also we have,  $g^{-1} = g^{p-1}$  is a element of order  $p$ . Therefore,  $H$  is a subgroup of  $G$ .

Let  $g$  be an element of order 4. Then  $g^2$  is a element of order 2. Thus, the subgroup generated by  $g$  or the smallest subgroup containing  $g$  has a element of order 2. Therefore, elements of order 4 together with identity cannot be a subgroup of  $G$ . Thus, ‘Torsion  $p$ -subgroup’ exists only if  $p$  is a square-free integer.

Let  $g$  be an element of order 6. Then  $g^2$  is an element of order 3. Again, elements of order 6 together with identity cannot be a subgroup of  $G$ . Thus, ‘Torsion  $p$ -subgroup’ exists only if  $p$  is a power of a prime. Thus, ‘Torsion  $p$ -subgroup’ of  $G$  exists only if  $G$  is abelian and  $p$  is a prime.

### 1.4.4 Normal Factors of $G$

This is a warm-up exercise for §37.5, where the theory is discussed by Fraleigh. However, we are able to conclude the following :

*Remark.* Let  $G = H \times K$ . Let  $g \in G$ . Then  $g = (h, k) \in H \times K$ . Clearly,  $H$  is a subset of  $G$ , and  $H \times \{e\} \leq G$ . This subgroup is isomorphic to  $H$ . Thus  $h \in H$  suggests the existence of  $(h, e) \in G$ .

Similarly  $k \in K$  suggests  $(e, k) \in G$ . Therefore,  $hk \in G$  suggests  $(h, e)(e, k) = (h, k) = (e, k)(h, e)$ . We know that,  $kh \in G$  suggest  $(e, k)(h, e) \in G$ . Thus,  $hk = kh$  for every  $h \in H$  and every  $k \in K$ .

In other words, if  $G = H \times K$ , then  $H, K$  are isomorphic to normal subgroups  $H', K'$  of  $G$  such that  $H' \cap K' = \{e\}$  and  $G \simeq K' \times H'$ .

## 1.5 Cosets and Homomorphism

**Definitions 1.5.1.** A permutation group  $S_n$  is the set of all permutations on the set  $\{1, 2, \dots, n\}$ .

*Remark.* Consider,  $(1, 2, 3)(4, 5), (1, 2)(3, 4) \in S_5$ . I was wrong about the order in which the permutations are carried out. It follows the same order as function composition. That is,  $f \circ g(x)$  implies  $f(g(x))$ . Similarly,  $\sigma\rho$  implies

$\sigma(\rho(1, 2, \dots, n))$ .

$$(1 \ 2) (3 \ 4) * (1 \ 2 \ 3) (4 \ 5) = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 2 & 3 & 1 & 5 & 4 \\ 1 & 4 & 2 & 5 & 3 \end{pmatrix} = (2 \ 4 \ 5 \ 3)$$

$$(1 \ 2 \ 3) (4 \ 5) * (1 \ 2) (3 \ 4) = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 2 & 1 & 4 & 3 & 5 \\ 3 & 2 & 5 & 1 & 4 \end{pmatrix} = (1 \ 3 \ 5 \ 4)$$

Clearly,  $S_5$  is a non-abelian group of order 120.

**Definitions 1.5.2.** Kernel of a function  $\phi : G \rightarrow G'$  is the inverse image of the identity element in  $G'$ .

**Definitions 1.5.3.** Cosets and Normal Subgroup,

- A **left coset**  $gH$  is the subset  $\{gh \in G : h \in H\}$  where  $g \in G$  and  $H \leq G$ .
- A **right coset**  $Hg$  is the subset  $\{hg \in G : h \in H\}$  where  $g \in G$  and  $H \leq G$ .
- A subgroup  $H$  of group  $G$  is **normal** if  $gH = Hg, \forall g \in G$ .
- All subgroups of abelian groups are normal.

*Remark.* For example,  $H = \{1, \rho^2, \mu, \mu\rho^2\}$  is a normal subgroup of  $D_4$ . And  $K = \{1, \mu\}$  is a subgroup of  $D_4$  which is not normal. Note that,  $\rho\mu \neq \mu\rho$ . Clearly,  $\rho K \neq K\rho$ . However,  $\rho H = \{\rho, \rho^3, \mu\rho^3, \mu\rho\} = H\rho$ .

*Remark* (Lagrange's Theorem). Let  $G$  be a finite group. If  $H \leq G$ , then order of  $H$  divides order of  $G$ .

- $aH \cap bH \neq \emptyset \implies aH = bH$ .
- $\forall g \in G, g \in gH$ .
- $\forall g \in G, |gH| = |H|$ .

*Remark* (Cayley's Theorem). Every group is isomorphic to a group of isomorphisms.

**Definitions 1.5.4.** Let  $G, G'$  be groups. A **group homomorphism** is a function  $\phi : G \rightarrow G'$  such that  $\phi(x)\phi(y) = \phi(xy)$ . Clearly  $\phi(e) = e'$ . A **trivial homomorphism** is a function  $\phi : G \rightarrow G'$  such that  $\phi(G) = \{e'\}$ .

*Remark.* Group homomorphism  $\phi : G \rightarrow G'$  preserves identity, inverses and subgroups. And kernel of group homomorphism is a normal subgroup of  $G$ .

*Remark.* For example,  $\phi : D_4 \rightarrow \mathbb{Z}_2$  defined by  $\phi(\rho) = 1$  and  $\phi(\mu) = 0$  is a group homomorphism with  $\ker(\phi) = \{1, \rho^2, \mu, \mu\rho^2\}$ .

## 1.6 Factor Groups

*Remark* (factor group). Let  $H$  be a normal subgroup of a group  $G$ . Then the **factor group** of  $G$  over  $H$ ,  $G/H$  is the group of cosets of  $H$  in  $G$ .

*Remark*. For example,  $H = \{1, \rho^2, \mu, \mu\rho^2\}$  is normal subgroup of  $D_4$ . And the factor group  $D_4/H = \{1H, \rho H\}$ .

**Theorem 1.6.1.** Let  $\phi : G \rightarrow G'$  be a group homomorphism with kernel  $H$ . Then cosets of  $H$  form a factor group,  $G/H$  where  $(aH)(bH) = (ab)H$ . Also,  $\mu : G/H \rightarrow \phi[G]$  defined by  $\mu(aH) = \phi(a)$  is an isomorphism.

*Proof.* Let  $\phi : G \rightarrow G'$  be a group homomorphism with  $\ker(\phi) = H$ . We have,  $\phi^{-1}(\phi(a)) = \{g \in G : \phi(g) = \phi(a)\}$ .

Let  $x \in aH$ . Then  $x = ah$  for some  $h \in H$ . And  $\phi(x) = \phi(ah) = \phi(a)\phi(h) = \phi(a)$ , since  $\phi(h) = e'$ . Thus,  $x \in \phi^{-1}(\phi(a))$  and  $aH \subset \phi^{-1}(\phi(a))$ .

Let  $x \in \phi^{-1}(\phi(a))$ . Then  $\phi(x) = \phi(a)$ . And  $\phi(a)^{-1}\phi(x) = e' \implies \phi(a^{-1}x) = e'$ . Clearly,  $a^{-1}x \in \ker(\phi)$ . Thus, there exists  $h \in H$  such that  $a^{-1}x = h$ . Therefore,  $x = ah$  for some  $h \in H$ . Thus,  $x \in aH$  and  $\phi^{-1}(\phi(a)) \subset aH$ . Therefore,  $\phi^{-1}(\phi(a)) = aH$ .

Similarly,  $\phi^{-1}(\phi(a)) = Ha$ . Thus  $aH = Ha$  and  $H$  is a normal subgroup of  $G$ . Therefore, we have the factor group  $G/H$ .

To prove :  $\mu : G/H \rightarrow \phi[G]$  is a one-one correspondence. ie,  $aH \xleftrightarrow{\mu} \phi(a)$ . To prove :  $\mu$  is injective. Suppose  $\mu(aH) = \mu(bH)$ . Then  $\phi(a) = \phi(b)$ . And  $b \in \phi^{-1}(\phi(a)) = aH$ . Therefore,  $bH = aH$ .

To prove :  $\mu$  is surjective. Let  $\phi(a) \in \phi[G]$ . Then, there exists  $aH$  such that  $\mu(aH) = \phi(a)$ .

We have,  $\mu(aH) = \phi(a)$ ,  $\mu(bH) = \phi(b)$ , and  $\mu((ab)H) = \phi(ab)$ . Therefore,  $\mu((aH)(bH)) = \mu((ab)H) = \phi(ab) = \phi(a)\phi(b) = \mu(aH)\mu(bH)$ . Thus  $\mu$  is a homomorphism. Therefore  $\mu$  is an isomorphism.  $\square$

**Theorem 1.6.2.** Let  $G$  be a group and  $H \leq G$ . Then left coset multiplication is well-defined by  $(aH)(bH) = (ab)H$  if and only if  $H$  is a normal subgroup of  $G$ .

*Proof.* Necessary part : Suppose  $(aH)(bH) = (ab)H$  is well-defined. Let  $a \in G$ . It is enough to prove that  $aH = Ha$ . Let  $x \in aH$ . Then  $(xH)(a^{-1}H) = (xa^{-1})H$ . Also  $(aH)(a^{-1}H) = eH = H$ . We have, coset multiplication is well-defined. Thus  $xa^{-1} = h \in H \implies x = ha \in Ha$ . Then,  $aH \subset Ha$ . Similarly,  $Ha \subset aH$  and  $aH = Ha$ . Therefore,  $H$  is a normal subgroup of  $G$ .

Sufficient part: Suppose  $H$  is a normal subgroup of  $G$ , and let  $x \in aH$  and  $y \in bH$ .  $x \in aH \implies x = ah_1$  for some  $h_1 \in H$   $y \in bH \implies y = bh_2$  for some  $h_2 \in H$ . Therefore  $xy = (ah_1)(bh_2) = a(h_1(bh_2)) = a((h_1b)h_2) = a((bh_3)h_2) = a(b(h_3h_2)) = a(bh_4)$ . Since  $H$  is a group,  $h_3h_2 = h_4 \in H$ . Thus,  $xy = a(bh_4) = (ab)h_4 \in (ab)H$  for all  $x \in aH$  and  $y \in bH$ . Thus  $(aH)(bH) = (ab)H$ .  $\square$

**Corollary 1.6.2.1.** Let  $H$  be a normal subgroup of  $G$ . Then the cosets of  $H$  form a group  $G/H$  under the binary operation  $(aH)(bH) = (ab)H$ .

*Proof.* Let  $H$  be a normal subgroup and  $aH, bH, cH$  are cosets of  $H$  in  $G$ .

G1 Associativity

$$(aH)[(bH)(cH)] = (aH)[(bc)H] = [a(bc)]H = [(ab)c]H = [(ab)H](cH) = [(aH)(bH)](cH)$$

G2 Existence of identity,  $eH$

$$(aH)(eH) = (ae)H = aH \text{ and } (eH)(aH) = (ea)H = aH.$$

G3 Existence of inverse  $(a^{-1}H)$

$$(aH)(a^{-1}H) = (aa^{-1})H = eH \text{ and } (a^{-1}H)(aH) = (a^{-1}a)H = eH.$$

□

*Remark.*  $n\mathbb{Z}$  is a normal subgroup of  $\mathbb{Z}$ . And  $\mathbb{Z}/n\mathbb{Z} \approx \mathbb{Z}_n$ .  $\mathbb{Z}_n$  is a torsion group isomorphic to a factor group of torsion free group  $\mathbb{Z}$ .

*Remark.* Let  $c \in \mathbb{R}^*$ . Then the cyclic group generated by  $c$  is a normal subgroup of  $\mathbb{R}$  and  $\mathbb{R}/\langle c \rangle \approx \mathbb{R}_c$ .

**Question 7.** Let  $c = 0.31$ . Find the coset  $x + \langle 0.31 \rangle$  containing 2.

## 1.7 Fundamental Homomorphism & Automorphisms

### 1.7.1 Fundamental Homomorphism Theorem

**Theorem 1.7.1.** Let  $H$  be a normal subgroup of  $G$ . Then  $\gamma : G \rightarrow G/H$  is defined by  $\gamma(x) = xH$  is a homomorphism with kernel  $H$ .

*Proof.*  $\gamma(x)\gamma(y) = (xH)(yH)$  Let  $h_1, h_2 \in H$ ,  $(xh_1)(yh_2) = xyh_3h_2 = xyh_4$  for some  $h_3, h_4 \in H$ . Therefore  $(xH)(yH) = (xy)H$ .  $\gamma(xy) = (xy)H = \gamma(x)\gamma(y)$ .  $\gamma(x) = xH = H \iff x \in H$ . Therefore,  $\ker(\gamma) = H$ . □

*Remark.* Suppose  $H$  is a normal subgroup of a group  $G$ . Then, a homomorphism  $\gamma : G \rightarrow G/H$  is a natural homomorphism.

**Theorem 1.7.2** (Fundamental Homomorphism). Let  $\phi : G \rightarrow G'$  be a group homomorphism with kernel  $H$ . Then  $\phi[G]$  is a group, and  $\mu : G/H \rightarrow \phi[G]$  given by  $\mu(gH) = \phi(g)$  is an isomorphism. If  $\gamma : G \rightarrow G/H$  is the homomorphism given by  $\gamma(g) = gH$ , then  $\phi(g) = \mu\gamma(g)$  for each  $g \in G$ .

*Proof.*  $\mu$  is an isomorphism  $G/H \xrightarrow{\mu} \phi[G]$ .  $\mu(\gamma(g)) = \mu(gH) = \phi(g)$ . Thus  $\mu\gamma = \phi$ . □

“Every group homomorphism  $\phi : G \rightarrow G'$  with kernel  $N$  has a unique natural group homomorphism  $\gamma : G \rightarrow G/N$  and a unique isomorphism  $\mu : G/N \rightarrow G'$  such that  $\phi = \mu\gamma$ . That is,  $\phi(g) = \mu(\gamma(g)) = \mu(gN)$ .”

### 1.7.2 Inner Automorphism

**Theorem 1.7.3.** Let  $H$  be a subgroup of  $G$ , then the following statements are equivalent:

1.  $ghg^{-1} \in H, \forall g \in G, h \in H$
2.  $gHg^{-1} = H, \forall g \in G$
3.  $gH = Hg, \forall g \in G$

*Proof.* Let  $G$  be a group and  $H \leq G$ . By right multiplication,  $gHg^{-1} = H \iff gH = Hg$ . Trivially,  $\forall h \in H, ghg^{-1} \in H \iff gHg^{-1} \subset H$ . Therefore, it is enough to prove that  $H \subset gHg^{-1}$ . Let  $h \in H$  and  $x \in G$ , then  $xhx^{-1} = h'$  for some  $h' \in H$ . Then  $h = x^{-1}h'x = x^{-1}h'(x^{-1})^{-1} \in gHg^{-1}$  where  $x^{-1} = g \in G$ . Thus  $h \in gHg^{-1}$  and  $H \subset gHg^{-1}$ . Therefore  $gHg^{-1} = H$ .  $\square$

**Definitions 1.7.1.** Automorphisms :

- An **automorphism** of  $G$  is an isomorphism  $\phi : G \rightarrow G$
- The **inner automorphism** of  $G$  by  $g \in G$  is the isomorphism  $i_g : G \rightarrow G$  defined by  $i_g(x) = gxg^{-1}$  for all  $x \in G$
- The **conjugate** of  $x$  by  $g$  is the element  $gxg^{-1} \in G$
- The **conjugate subgroup** of subgroup  $H$ ,  $i_g[H] = \{ghg^{-1} : h \in H\}$

*Remark.* For example,  $i_\rho : D_4 \rightarrow D_4$  defined by  $i_\rho(x) = \rho x \rho^{-1}$  is an inner automorphism. We have,  $H = \{1, \mu\}$  is a subgroup of  $D_4$ . The conjugate subgroup  $i_\rho[H] = \{1, \mu\rho^2\}$ .

*Remark.* Normal subgroups are invariant under any inner automorphism.

## 1.8 Exercise §14

### 1.8.1 Normal subgroups

**Question 8.** Prove that the notion of normality is stronger than abelian.

*Remark.* Let  $G$  be a group. Then intersection of normal subgroups of  $G$  is a normal subgroup of  $G$ .

*Proof.* Let  $\mathcal{N}$  be a nonempty subfamily of normal subgroups of  $G$ . Let  $N$  be the intersection of subgroups in  $\mathcal{N}$ . Clearly, intersection of subgroups of  $G$  is also subgroup of  $G$ .

Suppose  $x \in N$ . Then  $x \in H$  for every  $H \in \mathcal{N}$ . Suppose  $g \in G$ . Then  $gxg^{-1} \in H$ , for every normal subgroup  $H \in \mathcal{N}$ . Thus  $gxg^{-1} \in N$ . Therefore,  $N$  is a normal subgroup of  $G$ .  $\square$

**Challenge 1.** Let  $\phi : G \rightarrow G'$  and  $\psi : G \rightarrow G'$  be two group homomorphism with kernel  $H$  and  $K$  respectively. Show that  $\phi\psi : G \rightarrow G'$  defined by  $\phi\psi(g) = \phi(g)\psi(g)$  is also a group homomorphism, but  $\ker(\phi\psi) \neq H \cap K$ .  
(hint : Construct  $\psi$  such that  $\phi(g)\psi(g) = e$  for some  $g \notin H$ )

*Remark.* Let  $S \subset G$ , then  $G$  has a smallest, normal subgroup containing  $S$ .

*Proof.* Let  $N$  be the intersection of all normal subgroups of  $G$  containing  $S$ . Then  $S \subset N$  and  $N$  is a normal subgroup of  $G$ . Let  $H$  be the smallest, normal subgroup of  $G$  containing  $S$ . Then  $N \subset H$ . Therefore,  $H = N$ .  $\square$

*Remark.* If a finite group  $G$  has exactly one subgroup  $H$  of order  $m$ , then  $H$  is normal.

*Proof.* Let  $G$  be a finite group. Suppose  $G$  has only one subgroup of order  $m$ , say  $H$ . We know that, conjugate of a subgroup is a subgroup of same order. Thus, conjugates of  $H$  are  $H$  itself. Thus,  $xHx^{-1} = H$  for every  $x \in G$ . Therefore,  $xH = Hx$  for every  $x \in G$  and  $H$  is a normal subgroup of  $G$ .  $\square$

*Remark.* If  $G$  has a subgroup of order  $s$ , then the intersection of all subgroups of order  $s$  is a normal subgroup of  $G$ .

*Proof.* Let  $G$  be a group. Let  $H$  be a subgroup of  $G$  of order  $s$ . Let  $N$  be the intersection of all subgroups of  $G$  of order  $s$ . ??  $\square$

### 1.8.2 Linear Groups

**Definitions 1.8.1.** The set of all  $n \times n$ , non-singular matrices with real entries is a group under matrix multiplication. This group is the **General Linear Group** and is denoted by  $GL(n, \mathbb{R})$ .

**Definitions 1.8.2.** The set of all  $n \times n$  matrices with real entries and determinant  $\pm 1$  is a group under matrix multiplication. This group is the **Special Linear Group** and is denoted by  $SL(n, \mathbb{R})$ .

*Remark.*  $SL(n, \mathbb{R})$  is a normal subgroup of  $GL(n, \mathbb{R})$

*Proof.* Let  $A \in GL(n, \mathbb{R})$  and  $B \in SL(n, \mathbb{R})$ . Let  $r = |A|$  and we have  $|B| = \pm 1$ . Then  $|AB| = |A| |B| = \pm r$ .

Let  $M_r$  be the set of all matrices with determinant  $\pm r$  where  $r \in \mathbb{R}$ . Clearly, for every matrix  $A$  with determinant  $r$ ,  $A \in M_r$ . And there exists a matrix  $C = B^{-1}AB$ . Then,  $|C| = |B^{-1}| |A| |B| = \pm r$ . Thus,  $C \in M_r$  and  $AB = BC$ . Clearly,  $M_r$  is a left coset of  $GL(n, \mathbb{R})$ . Therefore,  $SL(n, \mathbb{R})$  is a normal subgroup of  $GL(n, \mathbb{R})$ .  $\square$

### 1.8.3 Factor Group

*Remark.*  $A_n$  is a normal subgroup of  $S_n$ . And  $S_n/A_n \simeq \mathbb{Z}_2$ .

*Proof.* Let function  $\phi : S_n \rightarrow \mathbb{Z}_2$  be defined by  $\phi(\sigma) = 0$  if  $\sigma$  is an even permutation and  $\phi(\sigma) = 1$  if  $\sigma$  is an odd permutation. Then  $\phi$  is a homomorphism with kernel  $A_n$ , the set of all even permutations in  $S_n$ . We know that, the kernel of a homomorphism is a normal subgroup of the domain. Therefore,  $A_n$  is a normal subgroup of  $S_n$ .  $\square$

*Remark.* If  $H$  is normal subgroup of  $G$  and  $(G : H) = m$ , then  $a^m \in H$  for all  $a \in G$ .

*Proof.* Let  $H$  be a normal subgroup of  $G$  such that  $(G : H) = m$ . Then  $|G/H| = m$ . Let  $a \in G$ , then  $aH \in G/H$ . Then by Lagrange's theorem, order of  $aH$  divides  $m$ . Therefore,  $(aH)^m = eH$ . We have,  $(aH)^m = a^m H$ . Therefore,  $a^m \in H$ .  $\square$

*Remark.* Every factor group of an abelian group is also abelian.

*Proof.* Let  $G$  be an abelian group and  $H$  be a normal subgroup of  $G$ . Then  $G/H$  is a factor group of  $G$ . Let  $aH, bH \in G/H$ . Then,  $(aH)(bH) = (ab)H = (ba)H = (bH)(aH)$ . Clearly, factor group of an abelian group is abelian.  $\square$

*Remark.* Let  $G$  be a group and  $T$  be the torsion subgroup of  $G$ , then the factor group,  $G/T$  is torsion free.

*Proof.* Let  $G$  be a group  $T$  its torsion subgroup. Let  $g \in G$ . If the order of  $g$  is finite, then  $g \in T$ . Suppose  $G/T$  has a element of finite order  $m > 1$ , say  $xT$  where  $x \notin T$ . Then  $x$  is an element of infinite order.

We have, order of  $xT$  is  $m$ . That is,  $(xT)^m = eT$ . Thus,  $(xT)^m = (x^m)T = eT$ . Therefore,  $x^m = y \in T$ . Since  $y \in T$ ,  $y$  is an element of finite order, say  $r$ . Then  $(x^m)^r = x^{mr} = e$ . This is a contradiction as  $x$  is an element of infinite order. Therefore, every element of  $G/T$  except  $eT$  are of infinite order. Clearly,  $G/T$  is torsion free.  $\square$

### 1.8.4 Commutator subgroup

**Definitions 1.8.3.** A **commutator**  $c$  in group  $G$  is an element in the form  $c = aba^{-1}b^{-1}$  for some  $a, b \in G$ . The **commutator subgroup** is the smallest normal subgroup containing all commutators in  $G$ .

*Remark.* Let  $C$  be the commutator subgroup of  $G$ , then  $G/C$  is abelian.

*Proof.* Let  $C$  be the commutator subgroup of  $G$ . Let  $x, y \in G$ . Then  $x^{-1}yxy^{-1} \in C$ . Therefore,  $x(x^{-1}yxy^{-1}) \in xC$  and  $y(x^{-1}yxy^{-1}) \in yC$ . But,  $x(x^{-1}yxy^{-1})y(x^{-1}yxy^{-1}) = yx(x^{-1}yxy^{-1}) \in (yx)C$ . Therefore,  $(xC)(yC) = (yx)C = (yC)(xC)$ . Clearly,  $G/C$  is abelian.  $\square$

*Remark.* The factor group  $G/C$  is the abelianised version of  $G$ .

*Remark.* Let  $G$  be a group and  $C$  be the commutator group of  $G$ . Let  $H$  be a normal subgroup of  $G$ . If  $G/H$  is abelian, then  $C$  is a subgroup of  $H$ .

**Question 9.** Find commutator subgroup of the dihedral group  $D_4$  ?

### 1.8.5 Automorphism

*Remark.* Every inner automorphism is an identity map for an abelian group.

*Proof.* Let  $G$  be an abelian group and  $g \in G$ . Then  $i_g : G \rightarrow G$  is defined by  $i_g(x) = gxg^{-1} = gg^{-1}x = x$ . Clearly,  $i_g$  is an identity map.  $\square$

*Remark.* Set of all  $g \in G$  such that the inner automorphism  $i_g$  is an identity map is normal.

*Proof.* Let  $H$  be the set of all  $g \in G$  such that  $i_g$  is an identity map.

$$\begin{aligned} H &= \{g \in G : i_g(x) = x, \forall x \in G\} \\ &= \{g \in G : gxg^{-1} = x, \forall x \in G\} \\ &= \{g \in G : gx = xg, \forall x \in G\} \end{aligned}$$

Therefore,  $H$  is a normal subgroup of  $G$ .  $\square$

*Remark.* Set of automorphisms  $\Gamma$  of a group  $G$  is a group under composition. And the set of inner automorphisms is a normal subgroup of  $\Gamma$ .

*Proof.* Let  $G$  be a group. We know that, the composition of two automorphisms is also an automorphism. Also, the composition of functions is associative. Let  $i$  be the identity map and  $\mu$  be any automorphism. Then  $i\mu = \mu = \mu i$ . Since automorphisms are bijective, there exists a unique inverse  $\mu^{-1}$  for each automorphism  $\mu$  such that  $\mu\mu^{-1} = i$ . Clearly, the set of all automorphisms of a group  $G$  is also a group.  $\square$

*Remark.* Subgroup conjugacy is an equivalence relation on the set of subgroups.

*Proof.* Let  $H$  be a subgroup of a group  $G$ . Conjugation is reflexive, since  $eHe^{-1} = H$  and  $H \sim H$ .

Let  $K$  be a conjugate subgroup of  $H$ . That is,  $H \sim K$ . Then  $K = gHg^{-1}$ . Clearly,  $H = g^{-1}K(g^{-1})^{-1}$ . Thus,  $H$  is a conjugate of  $K$ . ie,  $K \sim H$ .

Let  $H \sim K$  and  $K \sim L$ . Then we have  $x, y \in G$  such that  $K = xHx^{-1}$  and  $L = yKy^{-1}$ . Now,  $L = y(xHx^{-1})y^{-1} = (yx)H(yx)^{-1}$ . Thus,  $H \sim L$ .  $\square$

**Question 10.** Find the automorphism group  $\Gamma(\mathbb{Z}_2 \times \mathbb{Z}_4)$  ?

## 1.9 Simple Groups

*Remark.* Factor Group Computations :

- The converse of Lagrange's theorem is false.  
For example,  $A_4$  has order 12, but doesn't have a subgroup of order 6.
- Factor group of a cyclic group is cyclic.  
(hint : if  $g$  is a generator of  $G$ , then  $gH$  is a generator of  $G/H$ .)

**Question 11.** Show that  $\mathbb{Z}_4 \times \mathbb{Z}_6 / \langle (2, 3) \rangle \approx \mathbb{Z}_4 \times \mathbb{Z}_3$ .

**Definitions 1.9.1.** A group is simple if it is non-trivial and has no proper, non-trivial normal subgroups.

*Remark.* Abelian, simple groups are  $\mathbb{Z}_p$ , the cyclic groups of prime order.

*Remark.* Symmetric group,  $S_3$  is not simple.

The subgroup,  $\{\rho_0, \rho_1, \rho_2\}$  is a normal subgroup of  $S_3$ .

*Remark.* Smallest nonabelian, simple group is  $A_5$ .

Every nonabelian, simple group of order 60 is isomorphic to  $A_5$ .

*Remark.* Simple groups :

- Alternating groups  $A_n$  are simple for  $n \geq 5$ .
- Every finite group can be factorised into simple groups.
- Every finite, non-abelian, simple group is of even order.
- Group homomorphism preserves normal subgroups.
- $M$  is maximal normal subgroup of  $G$  if and only if  $G/M$  is simple.
- Center  $Z(G) = \{z \in G : zg = gz, \forall g \in G\}$  is normal.
- Center of non-abelian groups of order  $pq$  are trivial if  $p, q$  are primes.
- Factor group,  $G/N$  is abelian if and only if  $C$  is a subgroup of  $N$ .

**Question 12.**  $G$  is simple and  $H$  is subgroup of  $G$ , then  $H$  is simple ?



## 1.10 Group Action on a Set

**Definitions 1.10.1.** An action of a group  $G$  on a set  $X$  is a map.  $*$  :  $G \times X \rightarrow X$  such that

1.  $ex = x, \forall x \in X$
2.  $(g_1g_2)(x) = g_1(g_2x), \forall x \in X, \forall g_1, g_2 \in G$

Then  $X$  is a  $G$ -set.

**Theorem 1.10.1.** Let  $X$  be a  $G$ -set.  $\forall g \in G, \sigma_g : X \rightarrow X$  defined by  $\sigma_g(x) = gx$  is a permutation of  $X$ . Also, the map  $\phi : G \rightarrow S_X$  defined by  $\phi(g) = \sigma_g$  is a homomorphism with the property that  $\phi(g)(x) = gx$ .

*Proof.* Suppose  $X$  is a  $G$ -set. Let  $g \in G$ , and  $x_1, x_2 \in X$ .

Suppose  $\sigma_g(x_1) = \sigma_g(x_2) \implies gx_1 = gx_2 \implies g^{-1}(gx_1) = g^{-1}(gx_2) \implies (g^{-1}g)x_1 = (g^{-1}g)x_2 \implies ex_1 = ex_2 \implies x_1 = x_2$ . Thus,  $\sigma_g$  is injective.

Let  $x \in X$ . Then  $\sigma_g(g^{-1}x) = g(g^{-1}x) = (gg^{-1})x = ex = x$ . Thus,  $\sigma_g$  is surjective. Therefore,  $\sigma_g$  is a permutation of  $X, \sigma_g \in S_X$ .

Let  $g_1, g_2 \in G$ . And  $\phi(g_1)(x) = \sigma_{g_1}(x) = g_1x, \phi(g_2)(x) = \sigma_{g_2}(x) = g_2x$ .  $\phi(g_1g_2)(x) = \sigma_{g_1g_2}(x) = (g_1g_2)x = g_1(g_2x) = \sigma_{g_1}(g_2x) = \phi(g_1)(g_2x) = \phi(g_1)(\sigma_{g_2}(x)) = \phi(g_1)(\phi(g_2)(x)) = \phi(g_1)\phi(g_2)(x)$ . Therefore,  $\phi(g_1g_2) = \phi(g_1)\phi(g_2)$  and  $\phi$  is a homomorphism.  $\square$

**Definitions 1.10.2.** Group Action,

- $G$  acts **faithfully** on  $X$ , if  $e$  is the only element that leaves every  $x \in X$  fixed.
- $G$  is **transitive** on  $X$  if for every  $x_1, x_2 \in X, \exists g \in G$  such that  $gx_1 = x_2$ .  
 $G$  is transitive on  $X$  iff the subgroup  $\phi[G]$  of  $S_X$  is transitive.

## 1.11 Isotropy subgroups & Orbits

**Definitions 1.11.1.** Let  $X$  be a  $G$ -set.

- The subset fixed by  $g, X_g = \{x \in X : gx = x\}$
- The isotropy subgroup of  $x, G_x = \{g \in G : gx = x\}$   
Let  $Y \subset X$ , then  $G_Y = \{g \in G : gy = y, \forall y \in Y\}$  is a subgroup of  $G$ .
- The orbit of  $x$  in  $X$  under  $G, Gx = \{gx \in X : g \in G\}$

**Theorem 1.11.1.** Let  $X$  be a  $G$ -set. Then  $G_x$  is a subgroup of  $G, \forall x \in X$ .

*Proof.* Let  $x \in X$ . And  $g_1, g_2 \in G_x$ . Then  $g_1x = x$  and  $g_2x = x$ .

Clearly,  $(g_1g_2)x = g_1(g_2x) = g_1x = x$ . Therefore,  $g_1g_2 \in G_x$ . Also  $ex = x \implies e \in G_x$ . Let  $g \in G_x$ . Then  $gx = x \implies g^{-1}(gx) = g^{-1}x \implies (g^{-1}g)x = g^{-1}x \implies x = g^{-1}x$ . Thus, for any  $g \in G_x, g^{-1} \in G_x$ . Therefore,  $G_x$  is a subgroup of  $G$  for any  $x \in X$ .  $\square$

**Theorem 1.11.2.** Let  $X$  be a  $G$ -set and  $x_1, x_2 \in X$ . Then the relation  $\sim$  defined by  $x_1 \sim x_2$  iff  $gx_1 = x_2$  is an equivalence relation.

*Proof.* Let  $x \in X$ . Then  $ex = x \implies x \sim x$ . Let  $x_1, x_2 \in X$  and  $x_1 \sim x_2$ . Then there exists some  $g \in G$  such that  $gx_1 = x_2$ . We have,  $g^{-1}x_2 = g^{-1}(gx_1) = (g^{-1}g)x_1 = ex_1 = x_1$ . Therefore,  $x_2 \sim x_1$ . Let  $x_1, x_2, x_3 \in X$  and  $x_1 \sim x_2$  and  $x_2 \sim x_3$ . Then there are  $g_1, g_2 \in G$  such that  $g_1x_1 = x_2$  and  $g_2x_2 = x_3$ . Clearly,  $g_2g_1 \in G$  and  $(g_2g_1)x_1 = g_2(g_1x_1) = g_2x_2 = x_3$ . Therefore,  $x_1 \sim x_3$ .  $\square$

**Theorem 1.11.3.** *Let  $X$  be a  $G$ -set and  $x \in X$ . Then  $|Gx| = (G : G_x)$ . If  $|G|$  is finite, then  $|Gx|$  is a divisor of  $|G|$ .*

*Proof.* We have  $Gx$  is the orbit of  $x$  in  $X$  under  $G$  and  $L_{G_x}$  is the left cosets of  $G_x$  in  $G$ . Let  $x_1 \in Gx$ . Then there exists  $g_1 \in G$  such that  $x_1 = g_1x$ . Define  $\psi : Gx \rightarrow L_{G_x}$  by  $\psi(x_1) = g_1G_x$ .

Step 1 :  $\psi$  is well-defined.

Let  $x_1 \in Gx$ . Suppose there exists  $g_1, g'_1 \in G$  such that  $g_1x = x_1$  and  $g'_1x = x_1$ . Then we have,  $g_1x = g'_1x \implies x = g_1^{-1}(g'_1x) = (g_1^{-1}g'_1)x$ . Thus,  $g_1^{-1}g'_1 \in G_x$ . Therefore,  $g_1(g_1^{-1}g'_1) \in g_1G_x$ . Clearly,  $g_1(g_1^{-1}g'_1) = (g_1g_1^{-1})g'_1 = g'_1 \in g_1G_x$ . Therefore,  $g_1G_x = g'_1G_x$ . And  $\psi(x_1) = g_1G_x$  is well-defined.

Step 2 :  $\psi$  is one-to-one.

Suppose  $\psi(x_1) = \psi(x_2)$ . Let  $x_1, x_2 \in Gx$  such that  $x_1 = g_1x$  and  $x_2 = g_2x$ . Then we have  $\psi(x_1) = \psi(x_2) \implies g_1G_x = g_2G_x$ . Thus,  $g_2 = g_1g$  for some  $g \in G_x$ . Clearly,  $x_2 = g_2x = (g_1g)x = g_1(gx) = g_1x = x_1$ .

Step 3 :  $\psi$  is onto.

Let  $g_1G_x$  be a left coset of  $G_x$  in  $G$ . Then we have,  $g_1 \in G$  and  $g_1x \in Gx$ , say  $x_1$ . Therefore, there exists  $x_1 \in Gx$  such that  $\psi(x_1) = g_1G_x$ .  $\square$

## 1.12 Exercise §16

**Definitions 1.12.1.** Let  $X$  be a  $G$ -set. And  $S \subset X$ . Then  $S$  is a **sub- $G$ -set** if the orbit  $Gs$  of each  $s \in S$  is contained in  $S$ .

$$S = \bigcup_{s \in S} Gs \quad (1.4)$$

*Remark.* • Every  $G$ -set is a union of its orbits.

• Every  $G$ -set is isomorphic to the disjoint union of left coset  $G$ -sets.

**Definitions 1.12.2.** Let  $G$  be a group. Let  $X, Y$  be two  $G$ -sets. Then function  $\phi : X \rightarrow Y$  is a  $G$ -set isomorphism if

1.  $\phi$  is a bijection and
2.  $\phi$  is a  $G$ -set homomorphism  
ie,  $g\phi(x) = \phi(gx)$ , for every  $x \in X$  and every  $g \in G$ .

## 1.13 Application of $G$ -Sets to Counting

**Theorem 1.13.1** (Burnside). *Let  $G$  be a finite group and  $X$  a finite  $G$ -set. If  $r$  is the number of orbits in  $X$  under  $G$ , then*

$$r|G| = \sum_{g \in G} |X_g| \quad (1.5)$$

*Proof.* Let  $N = \{(g, x) \in G \times X : gx = x\}$ .

Step 1 :  $|N| = \sum_{g \in G} |X_g|$ .

Let  $g \in G$ . We have,  $X_g$  is the set of all  $(g, x) \in N$  with  $g$  as first member. Enumerating elements of  $N$  for each  $g \in G$ , we get  $\sum_{g \in G} |X_g| = N$ .

Step 2 :  $|N| = r|G|$ .

Let  $x \in X$ . We have,  $G_x$  is the set of all  $g \in G$  such that  $gx = x$ . In other words,  $G_x$  is the set of all  $g \in G$  such that  $(g, x) \in N$  with  $x$  as second member. However,  $|Gx| = (G : G_x)$ .

$$\implies |Gx| = \frac{|G|}{|G_x|} \implies |G_x| = \frac{|G|}{|Gx|}$$

Enumerating element of  $N$  for each  $x \in X$ , we get

$$|N| = \sum_{x \in X} |G_x| = \sum_{x \in X} \frac{|G|}{|Gx|} = |G| \sum_{x \in X} \frac{1}{|Gx|}$$

Let  $Gx = \mathcal{O}$  be an orbit containing  $x$  with length  $k$ .

$$\implies \sum_{x \in \mathcal{O}} \frac{1}{|Gx|} = \sum_{i=1}^k \frac{1}{k} = 1$$

Let  $r$  be the number of orbits in  $X$  under the group action of  $G$ . Clearly, the orbits of  $X$  are disjoint. Thus,

$$|N| = |G| \sum_{i=1}^r \sum_{x \in \mathcal{O}_i} \frac{1}{|Gx|} = |G| \sum_{i=1}^r 1 = r|G|$$

Therefore,  $\sum_{g \in G} |X_g| = |N| = r|G|$ . □

**Corollary 1.13.1.1.** *If  $G$  is a finite group and  $X$  is a finite  $G$ -set, then*

$$\text{number of orbits in } X \text{ under } G = \frac{1}{|G|} \sum_{g \in G} |X_g| \quad (1.6)$$

*Proof.* From Burnside formula, we have

$$\sum_{g \in G} |X_g| = r|G| \implies r = \frac{1}{|G|} \sum_{g \in G} |X_g|$$

□

## 1.14 Isomorphism Theorems 1-2

### 1.14.1 First Isomorphism Theorem

**Theorem 1.14.1** (first isomorphism). *Let  $\phi : G \rightarrow G'$  be a homomorphism with kernel  $K$ , and let  $\gamma_K : G \rightarrow G/K$  be the canonical homomorphism. There is a unique isomorphism  $\mu : G/K \rightarrow \phi[G]$  such that  $\phi(x) = \mu(\gamma_K(x))$  for each  $x \in G$ .*

*Proof.* By, fundamental homomorphism theorem.[Fraleigh, 2013, §14.1] □

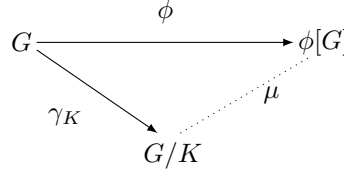


Figure 1.1: First Isomorphism Theorem

### 1.14.2 Second Isomorphism Theorem

**Definitions 1.14.1.** Join  $H \vee N$  is the smallest subgroup of  $G$  containing  $HN$  where  $HN = \{hn : h \in H, n \in N\}$ .

**Lemma 1.14.2.** Let  $N$  be a normal subgroup of  $G$  and let  $\gamma : G \rightarrow G/N$  be the canonical homomorphism. Then the map  $\phi$  from the set of normal subgroups of  $G$  containing  $N$  to the set of normal subgroup of  $G/N$  given by  $\phi(L) = \gamma[L]$  is one-to-one and onto.

*Proof.* Let  $G$  be a group and  $N$  be a normal subgroup of  $G$ . Given that  $\gamma : G \rightarrow G/N$  the canonical homomorphism. That is,  $\gamma(g) = gN$  for every  $g \in G$ . Let  $L, M$  be normal subgroups of  $G$  containing  $N$ . Since homomorphism preserves normality,  $\gamma(L)$  is a normal subgroup of  $\gamma[G]$ .

Suppose  $\phi(L) = \phi(M)$ . By definition of  $\phi$ ,

$$\gamma[L] = \phi(L) = \phi(M) = \gamma[M]$$

Since  $N$  is normal,  $\gamma^{-1}(\gamma(g)) = gN$  for every  $g \in G$ . And

$$\gamma^{-1}(\gamma[L]) = \gamma^{-1}\left(\bigcup_{g \in L} \gamma(g)\right) = \bigcup_{g \in L} \gamma^{-1}(\gamma(g)) = \bigcup_{g \in L} gN = L$$

for every point  $g \in L$ , the left coset  $gN$  is contained in  $L$ . ( $\because N \leq L$ ) Therefore,

$$\gamma^{-1}(\phi(L)) = \gamma^{-1}(\gamma[L]) = L$$

$$\gamma^{-1}(\phi(M)) = \gamma^{-1}(\gamma[M]) = M$$

Thus,  $L = M$ . Therefore,  $\phi$  is injective.

Let  $H$  be a normal subgroup of  $G/N$ . Then  $\gamma^{-1}(H)$  is a normal subgroup of  $G$ . We have,  $eN \in H$  and  $\gamma^{-1}(eN) = eN = N$ . Thus,  $N \subset \gamma^{-1}(H)$ . Thus, there exists  $\gamma^{-1}(H)$ , a normal subgroup of  $G$  containing  $N$  such that  $\phi(\gamma^{-1}(H)) = H$ . Therefore,  $\phi$  is surjective.  $\square$

**Lemma 1.14.3.** If  $N$  is a normal subgroup of  $G$ , then  $H \cap N = HN = NH$ . Furthermore, if  $H$  is also normal in  $G$ , then  $HN$  is normal in  $G$ .

*Proof.* Let  $G$  be a group and  $N$  be a normal subgroup of  $G$ . Also let  $H$  be a subgroup of  $G$ .

Claim :  $HN = \{hn \in G : h \in H, n \in N\}$  is a subgroup of  $G$ .

G1 Closure :  $h_1n_1h_2n_2 = h_1(h_2n_3)n_2 = h_3n_4 \in HN$ , where  $h_1, h_2 \in H$  and  $n_1, n_2, n_4 \in N$ . Since  $N$  is normal,  $n_1h_2 = h_2n_3$  for some  $n_3 \in N$ .

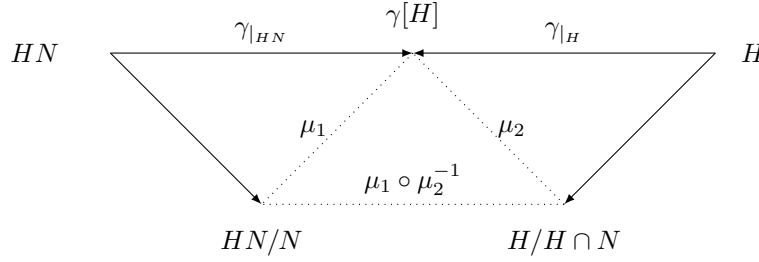


Figure 1.2: Second Isomorphism Theorem

G2  $HN \subset G$ , thus  $HN$  satisfies associativity.

G3 Since  $H, N \leq G$ ,  $e \in H$  and  $e \in N$ . Thus,  $e = ee \in HN$ .

G4 Let  $h_1n_1 \in HN$ .  $H \leq G$  and  $h_1 \in H \implies h_1^{-1} \in H$ . Then  $n_1^{-1}h_1^{-1} = h_1^{-1}n_2$  for some  $n_2 \in N$ , since  $N$  is normal. Thus, every element in  $HN$ ,  $h_1n_1$  has an inverse  $h_1^{-1}n_2 \in HN$ .

Let  $H$  be a normal subgroup of  $G$ . Let  $g \in G$ . Then  $g(h_1n_1) = (gh_1)n_1 = (h_2g)n_1 = h_2(gn_1) = h_2(n_2g) = (h_2n_2)g$  for some  $h_2 \in H$  and  $n_2 \in N$  since both  $H$  and  $N$  are normal subgroups of  $G$ . Thus  $gHN = HNg$  for every  $g \in G$ . Therefore,  $HN$  is a normal subgroup of  $G$ .  $\square$

**Theorem 1.14.4** (second isomorphism). *Let  $H$  be a subgroup of  $G$  and let  $N$  be a normal subgroup of  $G$ . Then  $(HN)/N \simeq H/(H \cap N)$ .*

*Proof.* Consider the canonical homomorphism  $\gamma : G \rightarrow G/N$  with kernel  $N$ . Then  $\gamma$  restricted to  $HN$  is a homomorphism from  $HN$  onto  $\gamma[H]$ .

$$\gamma|_{HN} : HN \rightarrow \gamma[H] \text{ where } \gamma|_{HN}(hn) = \gamma(hn) = (hn)N = hN = \gamma(h) \quad (1.7)$$

Since  $\ker(\gamma) = N$  and  $N \subset HN$ . We have,  $\ker(\gamma|_{HN}) = N$ . By first isomorphism theorem there exists a unique isomorphism  $\mu_1 : HN/N \rightarrow \gamma[H]$  where  $\mu_1(hnN) = \gamma(h) = hN$ .

Similarly,  $\gamma$  restricted to  $H$  is also a homomorphism onto  $\gamma[H]$ .

$$\gamma|_H : H \rightarrow \gamma[H] \text{ where } \gamma|_H(h) = \gamma(h) = hN. \quad (1.8)$$

Since  $\ker(\gamma) = N$  and  $H \cap N \neq \phi$ .  $\ker(\gamma|_H) = H \cap N$ . By first isomorphism theorem, there exists a unique isomorphism  $\mu_2 : H/(H \cap N) \rightarrow \gamma[H]$

Then  $HN/N \simeq H/(H \cap N)$ , since the composition of two isomorphisms,  $\mu_2^{-1} \circ \mu_1 : HN/N \rightarrow H/(H \cap N)$  is an isomorphism  $\square$

## 1.15 Third Isomorphism Theorem

**Theorem 1.15.1** (third isomorphism). *Let  $H$  and  $K$  be normal subgroup of  $G$  with  $K \leq H$ . Then  $G/H \simeq (G/K)/(H/K)$ .*



Figure 1.3: Third Isomorphism Theorem

*Proof.* Consider the function  $\phi : G \rightarrow (G/K)/(H/K)$  defined by  $\phi(g) = gK(H/K)$ . Then,  $\phi$  is a homomorphism.

$$\begin{aligned}
 \phi(ab) &= (ab)K/(H/K) \\
 &= (aK)(bK) (H/K), \because (ab)K = (aK)(bK) \\
 &= aK(H/K) bK(H/K), \because (xy)H/K = xH/K yH/K \\
 &= \phi(a)\phi(b)
 \end{aligned}$$

The kernel of  $\phi$  is the set  $\{a \in G : \phi(a) = H/K\}$ . Since the coset  $H/K$  was originally the points in  $H$ , we have  $\ker(\phi) = H$ . By first isomorphism theorem, there exists a unique isomorphism  $\mu : G/H \rightarrow (G/K)/(H/K)$ . Thus,  $G/H \simeq (G/K)/(H/K)$ .  $\square$

## 1.16 Finite, non-abelian Groups

**Theorem 1.16.1.** Let  $p$  be a prime. Let  $G$  be a group of order  $p^n$ . Let  $X$  be a finite  $G$ -set. Then  $|X| \equiv |X_G| \pmod{p}$ .

*Proof.* Suppose there are  $r$  orbits in  $X$ . Choose an element from each orbit, say  $x_1, x_2, \dots, x_r$ . We have,

$$|X| = \sum_{i=1}^r |Gx_i| \quad (1.9)$$

Let  $X_G = \{x \in X : gx = x, \forall g \in G\}$ . Then each element in  $X_G$  belongs to an orbit of length 1. Let  $|X_G| = s$ . Then,

$$|X| = |X_G| + \sum_{i=s+1}^r |Gx_i| \quad (1.10)$$

We have,  $|Gx| = (G : G_x)$ . Clearly, both  $G$  and  $G_x$  are groups of order a multiple of prime  $p$ . Thus, for every  $x \in X$ ,  $|Gx|$  is multiple of prime  $p$ . Therefore,  $|X| \equiv |X_G| \pmod{p}$ .  $\square$

**Challenge 2.** Let  $p$  be a prime and  $G$  be a finite group  $G$ . Design a mechanism to enumerate all elements of order  $p$ ?

**Definitions 1.16.1.** Let  $G$  be a group. Let  $p$  be a prime.  $G$  is a  **$p$ -group** if every element of  $G$  has order a power of  $p$ . A subgroup  $H$  of  $G$  is a  **$p$ -subgroup** of  $G$  if every element of  $H$  has order a power of  $p$ .

**Theorem 1.16.2** (Cauchy). *Let  $p$  be a prime. Let  $G$  be a finite group and  $p$  divides  $|G|$ , then  $G$  has a subgroup of order  $p$ .*

*Proof.* Let  $X$  be the set of all  $n$ -tuples of elements of  $G$  such that the product of co-ordinates of each  $n$ -tuple is the identity element  $e$  of  $G$ .

$$X = \{(g_1, g_2, \dots, g_p) \in G^p : g_1 g_2 \cdots g_p = e\} \quad (1.11)$$

Let  $g_1, g_2, \dots, g_{p-1}$  be any elements in  $G$ . Then  $g_p = (g_1 g_2 \cdots g_{p-1})^{-1}$  is uniquely determined. That is, the  $(p-1)$  co-ordinates of an element in  $X$  may be chosen in  $|G|$  different ways. Thus,  $|X| = |G|^{p-1}$ . Since  $p$  divides  $|G|$ ,  $p$  divides  $|X|$ .

We have  $g_1(g_2 \cdots g_p) = (g_2 g_3 \cdots g_p)g_1$ , since  $g_2 g_3 \cdots g_p = g_1^{-1}$ . Thus, for every  $(g_1, g_2, \dots, g_p) \in X$ , the  $\sigma$  permutation of that element is in  $X$ . That is,  $\sigma(g_1, g_2, \dots, g_p) = (g_2, g_3, \dots, g_p, g_1) \in X$ . Similarly,  $g_2(g_3 g_4 \cdots g_p g_1) = (g_3 g_4 \cdots g_p g_1)g_2$ . And  $\sigma^2(g_1, g_2, \dots, g_p) = (g_3, g_4, \dots, g_p, g_1, g_2)$ . Clearly, the subgroup generated by  $\sigma$ , is a subgroup of  $S_p$  and  $X$  is a  $\langle \sigma \rangle$ -set. We have,  $|X| \cong |X_{\langle \sigma \rangle}| \pmod{p}$ . Thus,  $p$  divides  $|X_{\langle \sigma \rangle}|$ .

Clearly,  $(e, e, \dots, e) \in X$  and  $(e, e, \dots, e) \in X_{\langle \sigma \rangle}$ . Thus,  $X_{\langle \sigma \rangle}$  has atleast  $p$  elements. Let  $(g_1, g_2, \dots, g_p) \in X_{\langle \sigma \rangle}$ , then  $\sigma$  fixes  $(g_1, g_2, \dots, g_p)$ .

$$\begin{aligned} \sigma(g_1, g_2, \dots, g_p) &= (g_1, g_2, \dots, g_p) \\ \implies (g_2, \dots, g_p, g_1) &= (g_1, g_2, \dots, g_p) \end{aligned}$$

Thus,  $g_1 = g_2 = \dots = g_p$ , say  $g \in G$ . That  $(g, g, \dots, g) \in X$  and  $g^p = e$  by the definition of  $X$ . Therefore,  $G$  has an element  $g$  of order  $p$ .  $\square$

**Corollary 1.16.2.1.** *Let  $G$  be a finite group. Then  $G$  is a  $p$ -group if and only if  $|G|$  is a power of  $p$ .*

*Proof.* Let  $G$  be a finite group. Suppose  $G$  is a  $p$ -group. Suppose there exists another prime  $q$ ,  $q \neq p$  such that  $q$  divides  $|G|$ . Then by Cauchy's theorem,  $G$  has an element of order  $q$ . This contradicts the assumption that  $G$  is a  $p$ -group. Thus, the only prime that divides  $|G|$  is  $p$ . Therefore, the order of  $G$  is a power of prime  $p$ .

Let  $|G| = p^n$ . Then the factors of  $p^n$  are powers of  $p$ . By Lagrange's theorem, order of subgroups of  $G$  must divide  $|G|$ . Thus, every subgroup of  $G$  has order a power of  $p$ . Thus, every element of  $G$  has order a power of  $p$ . Therefore,  $G$  is a  $p$ -group.  $\square$

**Definitions 1.16.2.** Let  $G$  be a group and  $H$  be a subgroup of  $G$ . Consider the inner automorphisms  $i_g : G \rightarrow G$  such that  $i_g(x) = gxg^{-1}$ . We have,  $i_g(H) = gHg^{-1}$  is the conjugate of the subgroup  $H$ . The set of all elements in  $G$  which has  $H$  itself as the conjugate of  $H$  is the normaliser of  $H$  in  $G$ .

$$N[H] = \{g \in G : gHg^{-1} = H\} \quad (1.12)$$

*Remark.* The normaliser of  $H$  in  $G$  is a subgroup of  $G$ . And  $N[H]$  is the largest subgroup of  $G$  with  $H$  as its normal subgroup.

**Lemma 1.16.3.** *Let  $H$  be a  $p$ -subgroup of a finite group  $G$ . Then*

$$(N[H] : H) \cong (G : H) \pmod{p} \quad (1.13)$$

*Proof.* Let  $G$  be a finite group. Let  $H$  be a  $p$ -subgroup of  $G$ . Let  $\mathcal{L}$  be the set of all left cosets of  $H$  in  $G$ . Then,  $|\mathcal{L}| = (G : H)$ .

Claim :  $\mathcal{L}$  is an  $H$ -set with group action  $h(xH) = (hx)H$ . We have,  $e(xH) = (ex)H = xH$  and  $(g_1g_2)(xH) = (g_1g_2x)H = g_1(g_2xH) = g_1(g_2(xH))$ .

Let  $\mathcal{L}_H$  be the set of all left cosets that are fixed under action by all element of  $H$ .

$$\begin{aligned}\mathcal{L}_H &= \{xH \in \mathcal{L} : h(xH) = xH, \forall h \in H\} \\ &= \{xH \in \mathcal{L} : x^{-1}h(xH) = H, \forall h \in H\} \\ &= \{xH \in \mathcal{L} : (x^{-1}hx)H = H, \forall h \in H\} \\ &= \{xH \in \mathcal{L} : (x^{-1}hx) \in H, \forall h \in H\} \\ &= \{xH \in \mathcal{L} : x^{-1} \in N[H]\}\end{aligned}$$

Clearly, left cosets of  $\mathcal{L}_H$  has all its elements contained in  $N[H]$ . Thus,

$$|\mathcal{L}_H| = (N[H] : H) \quad (1.14)$$

We have,  $\mathcal{L}$  is an  $H$ -set. And  $H$  is a  $p$ -subgroup. Thus,  $H$  has order a power of prime  $p$ . Therefore,

$$|\mathcal{L}| \cong |\mathcal{L}_H| \pmod{p} \quad (1.15)$$

□

**Corollary 1.16.3.1.** *Let  $H$  be a  $p$ -subgroup of finite group  $G$ . If  $p$  divides  $(G : H)$ , then  $N[H] \neq H$ .*

*Proof.* We have,  $(G : H) \cong (N[H] : H) \pmod{p}$ . And  $p$  divides  $(G : H)$ . Thus,  $p$  divides  $(N[H] : H)$ . Therefore,  $H \neq N[H]$ . □

## 1.17 Sylow Theorems

**Theorem 1.17.1** (First Sylow Theorem). *Let  $G$  be a finite group of order  $p^n m$  where  $n \geq 1$  and  $p$  does not divide  $m$ . Then*

1.  $G$  contains a subgroup of order  $p^i$  where  $1 \leq i \leq n$ .
2. Every subgroup  $H$  of  $G$  of order  $p^i$  is normal subgroup of the subgroup of order  $p^{i+1}$ , for  $1 \leq i \leq n$ .

*Proof.* We have,  $G$  is finite group and  $p$  divides  $|G|$ . By Cauchy's theorem,  $G$  has a subgroup of order  $p$ .

Suppose  $G$  has a subgroup  $H$  of order  $p^i$ , where  $(i < n)$ . Then,  $H$  is a  $p$ -subgroup. Thus,  $(G : H) \cong (N[H] : H) \pmod{p}$  where  $N[H]$  is the normaliser of  $H$  in  $G$ . Clearly,  $p$  divides  $(G : H)$ . Thus,  $p$  divides  $(N[H] : H)$ .

And  $H$  is a normal subgroup of  $N[H]$ . Thus, we have factor group  $N[H]/H$  and  $p$  divides the order of  $N[H]/H$ . By Cauchy's theorem,  $N[H]/H$  has a subgroup  $K$  of order  $p$ . Consider the canonical homomorphism  $\gamma : N[H] \rightarrow N[H]/H$  defined by  $\gamma(x) = xH$ . Then  $\gamma^{-1}(K)$  is a subgroup of  $N[H]$  of order  $p^{i+1}$  and contains  $H$ . Thus,  $H$  is a normal subgroup of  $\gamma^{-1}(K)$ . By mathematical induction,  $G$  has subgroups of order  $p^i$  for  $i = 2, 3, \dots, n$ . □



**Definitions 1.17.1.** A Sylow  $p$ -subgroup of  $G$  is a maximal  $p$ -subgroup of  $G$ .

**Theorem 1.17.2** (Second Sylow Theorem). *Let  $P_1$  and  $P_2$  be two Sylow  $p$ -subgroups of  $G$ . Then  $P_1$  and  $P_2$  are conjugate subgroups of  $G$ .*

*Proof.* Let  $P_1$  and  $P_2$  be two Sylow  $p$ -subgroups of  $G$ . Let  $\mathcal{L}$  be the set of all left cosets of  $P_1$ . Then, group  $P_2$  act on  $\mathcal{L}$  by  $y(xP_1) = (yx)P_1$ . Thus,  $\mathcal{L}$  is a  $P_2$ -set. Therefore,  $|\mathcal{L}| \cong |\mathcal{L}_{P_2}| \pmod{p}$ .

Clearly  $|\mathcal{L}| = (G : P_1)$ . And  $p$  doesn't divide  $|\mathcal{L}|$ . Thus,  $p$  does not divide  $|\mathcal{L}_{P_2}|$ . Therefore,  $|\mathcal{L}_{P_2}| > 0$ .

Thus,  $\mathcal{L}$  has at least an element  $xP_1$  which is fixed in the action of every element in  $P_2$ . That is,  $yxP_1 = xP_1$  for every  $y \in P_2$ . Therefore,  $x^{-1}yxP_1 = P_1$  for every  $y \in P_2$ .

In other words,  $x^{-1}yx \in P_1$  for every  $y \in P_2$ . Therefore,  $x^{-1}P_2x \leq P_1$ . But,  $|P_1| = |P_2|$ . Thus,  $x^{-1}P_2x = P_1$ . Therefore,  $P_1$  and  $P_2$  are conjugate subgroups of  $G$ .  $\square$

**Theorem 1.17.3** (Third Sylow Theorem). *If  $G$  is a finite group and  $p$  divides  $|G|$ , then the number of Sylow  $p$ -subgroups is congruent to 1  $\pmod{p}$  and divides  $|G|$ .*

*Proof.* Let  $\mathcal{S}$  be the set of all Sylow  $p$ -subgroups of  $G$ . Let  $P$  be a Sylow  $p$ -subgroup of  $G$ . The elements of  $P$  act on  $\mathcal{S}$  by conjugation. Let  $x \in P$  and  $T \in \mathcal{S}$ , then  $x$  carries  $T$  into  $xTx^{-1}$ . Clearly,  $\mathcal{S}$  is a  $P$ -set.

We have,  $|\mathcal{S}| \cong |\mathcal{S}_P| \pmod{p}$ . Suppose  $T \in \mathcal{S}_P$ . Then,  $xTx^{-1} = T$  for every  $x \in P$ . Thus,  $T$  is a normal subgroup of  $P$ . But,  $T$  and  $P$  are of the same order, since both are Sylow  $p$ -subgroups of  $G$ . Therefore,  $T = P$ . Thus,  $\mathcal{S}_P = \{P\}$ . And  $|\mathcal{S}_P| = 1$ . Thus,  $|\mathcal{S}| \cong 1 \pmod{p}$ .

Let  $G$  act on  $\mathcal{S}$  by conjugation. Let  $x \in G$  and  $P \in \mathcal{S}$ , then  $x$  carries  $P$  into  $xPx^{-1}$ . Clearly,  $\mathcal{S}$  is a  $G$ -set. However, every Sylow  $p$ -subgroup of  $G$  are conjugates. Thus, every Sylow  $p$ -subgroup of  $G$  belong to the same orbit under conjugation action. We have,  $|Gx| = (G : G_x)$  and thus length of orbits divides the order of  $G$ . Since the number of Sylow  $p$ -subgroups is same as the length of the orbit of  $P$ , the number of Sylow  $p$ -subgroup of  $G$  divides the order of  $G$ .  $\square$

*Remark.* Let  $G$  be a group of order 15. Let  $p = 3$ . By Third Sylow Theorem, the number of Sylow 3-subgroup of  $G$  is congruent to 1  $\pmod{3}$  and divides 15. We have, congruence class  $\hat{1} = \{1, 4, 7, 10, 13\}$ . Only 1 divides 15. Thus, there is only one Sylow 3-subgroup of  $G$ .

Let  $p = 5$ . By Third Sylow Theorem, the number of Sylow 5-subgroup of  $G$  is congruent to 1  $\pmod{5}$  and divides 15. We have, congruence class  $\hat{1} = \{1, 6, 11\}$ . Only 1 divides 15. Thus, there is only one Sylow 5-subgroup of  $G$ .

*Remark.* Let  $G$  be a group of order 255. And 1, 3, 5, 15, 17, 51, 85, 255 are the divisors of 255. Let  $p = 3$ . By third Sylow theorem, either there is one or eighty-five Sylow 3-subgroups of  $G$ . Suppose there 85 Sylow 3-subgroups. Then there are 170 elements of order 3.

Let  $p = 5$ . By third Sylow theorem, either there are one or fifty-one Sylow 5-subgroups of  $G$ . Suppose there 51 Sylow 5-subgroups. Then there are 204 elements of order 5.

Let  $p = 17$ . By third Sylow theorem, there is exactly one Sylow 17-subgroup of  $G$ . Thus, there are exactly 16 elements of order 17 in  $G$ .

## 1.18 Sylow Theorem : Applications

**Definitions 1.18.1.** A group  $G$  is solvable if there is sequence of subgroups  $\{e\} = H_0 \leq H_1 \leq \cdots H_n = G$  such that for  $i = 0, 1, 2, \dots$ ,

1.  $H_i$  is a normal subgroup of  $H_{i+1}$
2.  $H_{i+1}/H_i$  is simple and
3.  $H_{i+1}/H_i$  is abelian.

**Theorem 1.18.1.** Every group of prime order is solvable.

*Proof.* □

**Definitions 1.18.2.** Let  $G$  be a finite group. Consider  $G$  acting on itself by conjugation. Then,

$$|G| = |Z(G)| + \sum_{s+1}^r |Gx_i| \quad (1.16)$$

This is the class equation of  $G$ . And each orbit of  $G$  under conjugation by itself is a conjugate class in  $G$ .

**Theorem 1.18.2.** The center of a finite, nontrivial  $p$ -group is nontrivial.

*Proof.* □

**Lemma 1.18.3.** Let  $G$  be a group containing normal subgroups  $H$  and  $K$  such that  $H \cap K = \{e\}$  and  $H \vee K = G$ . Then  $G$  is isomorphic to  $H \times K$ .

*Proof.* □

**Theorem 1.18.4.** For a prime number  $p$ , every group of  $G$  of order  $p^2$  is abelian.

*Proof.* □

## 1.19 Sylow Theorem : Further Applications

**Theorem 1.19.1.** If  $H, K$  are subgroups of a group  $G$ . Then,

$$|HK| = \frac{(|H|)(|K|)}{|H \cap K|} \quad (1.17)$$

*Proof.* □

### 1.19.1 Analysis of Finite Groups

The following are a few sample tests for finite groups. For different numbers, you may have to use a combination of these tests.

*Remark.* Finite, non-abelian, simple groups are of order 60, 168, 360, ....

**Type 1 :  $p^r$** 

Groups of order  $p^r$  for  $r > 1$ , are not simple as they have normal subgroup of order  $p^{r-1}$  by first Sylow theorem.

For example, any group  $G$  of order 16 is not simple as by first Sylow theorem  $G$  has a normal subgroup of order 8.

**Type 2 :  $pq$** 

Suppose  $G$  is a group with order  $pq$  where  $p, q$  are primes and  $q > p$ . Clearly,  $p \not\equiv 1 \pmod{q}$ . Thus, by third Sylow theorem, the Sylow  $p$ -subgroup of  $G$  is unique, say  $H$ .

Now, by second Sylow theorem, Sylow  $p$ -subgroups are conjugate subgroups. Thus, this Sylow  $p$ -subgroup  $H$  is its only conjugate. That is,  $gHg^{-1} = H$  for every  $g \in G$ . Thus, it is a normal subgroup of  $G$ . Thus  $G$  is not simple.

Further, if  $q \not\equiv 1 \pmod{p}$  for  $q > p$ , then group  $G$  is cyclic.

For example,  $5 \not\equiv 1 \pmod{3}$ . Thus any group  $G$  of order 15 is cyclic. Clearly,  $G$  is abelian and is not simple.

However,  $7 \equiv 1 \pmod{3}$ . Thus, groups of order 21 are not simple, but are not necessarily cyclic.

**Type 3 : Only one Sylow  $p$ -subgroup**

In this case, we don't have a strict form. However, the nature of prime factor of the order suggests that  $G$  has only one Sylow  $p$ -subgroup for some prime factor  $p$  of its order.

For example, group  $G$  of order 20 has only one Sylow 5-subgroup say,  $H$  (by third Sylow theorem). Thus, this subgroup  $H$  is a normal subgroup of  $G$  (by second Sylow theorem). Therefore, groups of order 20 are not simple.

**Type 4 : Enumerating elements of Sylow  $p$ -subgroups**

Groups of order 30 are not simple. — The distinct Sylow  $p$ -subgroups suggests  $p - 1$  elements are unique to each. By counting, we can show that not all  $p$ -subgroups have a conjugate.

**Type 5 : Applying the relation for  $|HK|$** 

Groups of order 48 are not simple. —  $|HK| =$

**Type 6 : Normaliser of  $H \cap K$** 

Groups of order 36 are not simple. — Let  $H, K$  be distinct Sylow 3-subgroups of  $G$ . Then, the normaliser of  $H \cap K$  suggests a normal subgroup which is either  $H \cap K$  or its normaliser.

**Type 7 : Commutator of  $G$** 

Groups of order 255 are not simple. — Its commutator subgroup has order 1. Thus, the group is cyclic. Clearly, cyclic groups are abelian and not simple.

**1.20 Rings, Fields & Integral Domains**

## Subject 2

# ME010102 Linear Algebra

2.1 Vector Spces

2.2 Linear Transformations

2.3 Determinants

2.4 Elementary Canonical Forms

**Subject 3**

**ME010103 Basic Topology**

**Subject 4**

**ME010103 Real Analysis**

**Subject 5**

**Graph Theory**



## Semester II

## Subject 6

# ME010201 Advanced Abstract Algebra

### 6.1 Extension Fields §29

#### Previous Results

- Let  $R$  be a commutative ring with unity. If  $M$  is a maximal ideal in  $R$ , then  $R/M$  is a field. [Fraleigh, 2013, §27.9]
- Let  $F$  be a field. Every polynomial in  $F[x]$  has a unique factorisation into irreducible polynomials except for order and unit. cite[§27.27]fraleigh
- If  $\alpha$  is a zero of  $f(x) \in F[x]$ , then  $f(\alpha) = 0$ . cite[§22.10]fraleigh
- If  $p(x)$  is irreducible over field  $F$ , then the principal ideal generated by  $p(x)$ , denoted by  $\langle p(x) \rangle$  is a maximal ideal in  $F[x]$ . [Fraleigh, 2013, §27.25]
- Let  $R$  be a ring with unity. And  $N$  be an ideal of  $R$  containing a unit. Then  $N = R$ . [Fraleigh, 2013, §27.5]

**Basic Goal** Let  $F$  be a field and  $f(x) \in F[x]$ . Find a field  $E$  such that  $F$  is a subfield of  $E$  and there exists a zero of  $f(x)$  in  $E$  ?

**Extension Field** Let  $F$  be a field. Field  $E$  is an extension field of  $F$  if  $F$  is a subfield of  $E$ .

Example :  $\mathbb{Q} \leq \mathbb{R} \leq \mathbb{C}$

**Tower of Fields** A diagrammatic representation emphasising the hierarchy of field extensions in which extension fields appears above their subfields.

**Theorem 6.1.1** (Kronecker). *Let  $F$  be field. And  $f(x)$  be a non-constant polynomial in  $F[x]$ . Then there exists an extension field  $E$  of  $F$  and an  $\alpha \in E$  such that  $f(\alpha) = 0$ .*

*Proof.* Let  $f(x) \in F[x]$ . Then  $f(x)$  has a unique factorisation into irreducible polynomials in  $F[x]$  (except for order and unit). Let  $p(x)$  be an irreducible factor of  $f(x)$ . If  $f(x)$  is irreducible over  $F$ , then  $f(x) = cp(x)$ . It is enough to

construct an extension field  $E$  containing both  $F$  and  $\alpha$  such that  $p(\alpha) = 0$ .

If  $p(x)$  is irreducible over  $F$ , then  $\langle p(x) \rangle$  is maximal ideal in  $F[x]$ . Therefore,  $F[x]/\langle p(x) \rangle$  is a field, say  $E$ .

Consider the function  $\psi : F \rightarrow F[x]/\langle p(x) \rangle$  defined by  $\psi(a) = a + \langle p(x) \rangle$ . We claim that  $\psi : F \rightarrow \psi[F]$  is a field isomorphism.  $\psi$  is a canonical homomorphism with trivial kernel. Thus,  $\psi$  is one-to-one.

Let  $a, b \in F$ . And suppose  $\psi(a) = \psi(b)$ . It is enough to prove that  $a = b$ . By the definition of  $\psi$ , we have  $a + \langle p(x) \rangle = b + \langle p(x) \rangle \implies a - b \in \langle p(x) \rangle$ . Suppose  $a \neq b$ . Then  $a - b \neq 0$  and  $\deg(a - b) = 0$ . Then,  $\langle p(x) \rangle = F[x]$  which is a contradiction since  $\langle p(x) \rangle$  is maximal ideal. Therefore,  $\psi$  is one-to-one.

We have  $p(x)$  is a factor of  $f(x)$ . Thus  $p(\alpha) = 0 \implies f(\alpha) = 0$ . Thus, it remains to prove that there exists  $\alpha \in F[x]/\langle p(x) \rangle$  such that  $p(\alpha) = 0$ .

Let  $p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$ . Consider  $\alpha = x + \langle p(x) \rangle$ . Then  $p(\alpha) = \phi_\alpha(p)$ . Thus,  $p(\alpha) = a_0 + a_1(x + \langle p(x) \rangle) + \dots + a_n(x + \langle p(x) \rangle)^n$ . Thus,  $p(\alpha) = (a_0 + a_1x + a_nx^n) + \langle p(x) \rangle = p(x) + \langle p(x) \rangle = \langle p(x) \rangle = 0$ . Therefore,  $p(\alpha) = 0$  and  $f(\alpha) = 0$ .  $\square$

**algebraic over  $F$**  Let  $F \leq E$ . An element  $\alpha \in E$  is algebraic over a field  $F$ , if there exists  $f(x) \in F[x]$  such that  $f(\alpha) = 0$ .

**transcendental over  $F$**  Let  $F \leq E$ . An element  $\alpha \in E$  is transcendental over the field  $F$ , if it is not algebraic over  $F$ .

**algebraic number** We have,  $\mathbb{Q} \leq \mathbb{C}$ . A complex number  $\alpha \in \mathbb{C}$  is algebraic if it is algebraic over  $\mathbb{Q}$ . Example :  $2, \sqrt{2}, i$

**transcendental number** A complex number  $\alpha \in \mathbb{C}$  is transcendental if it is not an algebraic number. Example :  $\pi, e$  (proof excluded)

**Note 1** : A polynomial  $f(x) \in F[x]$  is reducible/irreducible depending upon the choice of the field  $F$ . For example :  $x^2 - 2$  is irreducible over  $\mathbb{Q}$ , but is reducible over  $\mathbb{R}$ .

**Note 2** : An element  $\alpha \in E$  is algebraic/transcendental depending on the choice of the field  $F$ . For example :  $\sqrt{2} \in \mathbb{C}$  is algebraic over  $\mathbb{Q}$  since  $x^2 - 2 \in \mathbb{Q}[x]$ .

**Theorem 6.1.2.** Let  $E$  be an extension field of  $F$ , and  $\alpha \in E$ . Let function  $\phi_\alpha : F[x] \rightarrow E$  be an evaluation homomorphism. Then  $\alpha$  is transcendental if and only if  $\phi_\alpha$  is one-to-one.

*Proof.* An element  $\alpha \in E$  is transcendental if and only if  $f(\alpha) \neq 0$  for any nonzero  $f(x) \in F[x]$  where  $f(\alpha) = \phi_\alpha(f)$ . Thus, kernel of  $\phi_\alpha$  is trivial. That is,  $\ker(\phi_\alpha) = \{0\}$ . Therefore,  $\phi_\alpha$  is one-to-one.  $\square$

**Theorem 6.1.3.** Let  $E$  be an extension field of  $F$  and  $\alpha \in E$  be algebraic over  $F$ . Then there exists a unique irreducible polynomial  $p(x)$  with minimum degree in  $F[x]$  and  $p(\alpha) = 0$ . If there exists a nonzero polynomial  $f(x) \in F[x]$  with  $f(\alpha) = 0$ , then  $p(x)$  divides  $f(x)$ .

*Proof.* Consider evaluation homomorphism  $\phi_\alpha : F[x] \rightarrow E$  defined by  $\phi_\alpha(f) = f(\alpha)$ . Then  $\ker(\phi)$  is an ideal in  $F[x]$ . Since every ideal in  $F[x]$  is principal, there exists  $p(x) \in F[x]$  such that  $\langle p(x) \rangle = \ker(\phi)$ . And  $f(\alpha) = 0 \implies f \in \ker(\phi) = \langle p(x) \rangle$ . Thus,  $p(x)$  divides  $f(x)$ .

Suppose  $p(x) = r(x)s(x)$ . Then  $p(\alpha) = r(\alpha)s(\alpha) = 0$ . However,  $E$  is a field and has no zero divisors. Thus, there exists a polynomial of lesser degree in  $\langle p(x) \rangle$  which is a contradiction.  $\square$

**monic polynomial** A polynomial which has 1 as the coefficient of highest power of  $x$ . For example :  $x^3 - 3x \in \mathbb{Q}[x]$

$\text{irr}(\alpha, F)$  The unique monic, irreducible polynomial  $p(x) \in F[x]$  such that  $p(\alpha) = 0$ . For example,  $\text{irr}(\sqrt{3}, \mathbb{Q}) = x^2 - 3$ . And  $\sqrt{3}$  is the **minimal** polynomial for  $\sqrt{3}$  over  $\mathbb{Q}$

$\text{deg}(\alpha, F)$  The degree of the unique monic, irreducible polynomial  $p(x) \in F[x]$  such that  $p(\alpha) = 0$ . For example,  $\text{deg}(\sqrt{3}, \mathbb{Q}) = 2$ .

**simple extension** An extension field  $E$  of field  $F$  is a simple extension if  $E = F(\alpha)$  for some  $\alpha \in E$ .

**Theorem 6.1.4.** Let  $E$  be a simple extension  $F(\alpha)$  of field  $F$  where  $\alpha \in E$  is algebraic over  $F$ . Let  $\text{deg}(\alpha, F) = n \geq 1$ . Then any element  $\beta \in E$  can be uniquely expressed in the form  $\beta = b_0 + b_1\alpha + \cdots + b_{n-1}\alpha^{n-1}$  where  $b_k \in F$ .

*Proof.* Consider the evaluation homomorphism  $\phi_\alpha : F[x] \rightarrow E$  defined by  $\phi_\alpha(f) = f(\alpha)$ . Then  $\phi_\alpha[F[x]] = F(\alpha)$ .

Let  $\text{irr}(\alpha, F) = p(x) = x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0$ . Then,  $p(\alpha) = 0$ . Thus,

$$\alpha^n = -a_{n-1}\alpha^{n-1} - a_{n-2}\alpha^{n-2} - \cdots - a_1\alpha - a_0 \quad (6.1)$$

Clearly, any higher power of  $\alpha$  can be eliminated from  $f(\alpha)$  as shown below,

$$\begin{aligned} \alpha^{n+1} &= \alpha\alpha^n = -a_{n-1}\alpha^n - a_{n-2}\alpha^{n-1} - \cdots - a_1\alpha^2 - a_0\alpha \\ &= a_{n-1}(a_{n-1}\alpha^{n-1} + a_{n-2}\alpha^{n-2} + \cdots + a_1\alpha + a_0) \\ &\quad - a_{n-2}\alpha^{n-1} - a_{n-3}\alpha^{n-2} - a_1\alpha^2 - a_0\alpha \end{aligned}$$

Thus, in the representation of the element in  $F(\alpha)$ , the maximum degree of  $\alpha$  is  $\text{deg}(\alpha, F) - 1$ . Therefore,  $\beta \in F(\alpha) \implies \beta = b_0 + b_1\alpha + b_2\alpha^2 + \cdots + b_{n-1}\alpha^{n-1}$ .

We can also show that this representation is unique for any  $\beta \in F(\alpha)$ . Suppose  $\beta = b'_0 + b'_1\alpha + b'_2\alpha^2 + \cdots + b'_{n-1}\alpha^{n-1}$ . Then  $0 = (b_0 - b'_0) + (b_1 - b'_1)\alpha + \cdots + (b_{n-1} - b'_{n-1})\alpha^{n-1} = \phi_\alpha(g)$  where  $g(x) = (b_0 - b'_0) + (b_1 - b'_1)x + \cdots + (b_{n-1} - b'_{n-1})x^{n-1}$ . Clearly, degree of  $g(x)$  is  $\text{deg}(\alpha, F) - 1$  which is less than the minimum degree for a non-zero irreducible polynomial for  $\alpha$  over  $F$ . Thus  $g(x) = 0$ . In other words,  $b_j = b'_j, \forall j$  and representation for  $\beta \in F(\alpha)$  is unique.  $\square$

### 6.1.1 Exercises §29

#### Irreducibility Conditions for Polynomials

1. Irreducibility over a finite field  
 $x^2 + 1$  is irreducible in  $\mathbb{Z}_3$  since for every  $x \in \{0, 1, 2\}$ ,  $x^2 + 1 \neq 0$ .
2. Irreducibility in rational field : Eisenstein's Criteria (§23.15)  
 Consider  $f(x) = x^3 + 60x^2 + 30x + 12$  since for  $p = 3$ ,  $f(x)$  satisfies Eisenstein's Criteria. And thus is  $f(x)$  irreducible over  $\mathbb{Q}$ .  
 Note that for  $p = 2, 5, 7, \dots$  the Eisenstein's Criteria is not satisfied.

#### Algebraic over a Field

1.  $\sqrt{2} + \sqrt{3}$  is algebraic over  $\mathbb{Q}$  [Fraleigh, 2013, Exercise 29.2]

$$\begin{aligned}
 \alpha &= \sqrt{2} + \sqrt{3} \\
 \implies \alpha^2 &= 2 + 2\sqrt{6} + 3 = 5 + 2\sqrt{6} \\
 \implies \alpha^2 - 5 &= 2\sqrt{6} \\
 \implies (\alpha^2 - 5)^2 &= 24 \\
 \implies \alpha^4 - 10\alpha^2 + 1 &= 0 \\
 \implies \phi_\alpha(x^4 - 10x^2 + 1) &= 0
 \end{aligned}$$

We have,  $\text{irr}(\sqrt{2} + \sqrt{3}, \mathbb{Q}) = x^4 - 10x^2 + 1$  and  $\deg(\sqrt{2} + \sqrt{3}, \mathbb{Q}) = 4$ .

2.  $\pi, e$  are transcendental numbers. (proof excluded)  
 However,  $\pi$  is algebraic over  $\mathbb{Q}(\pi)$  since  $x - \pi \in \mathbb{Q}(\pi)[x]$ .
3. Consider  $\alpha = \pi^2$  and  $F = \mathbb{Q}(\pi^3)$ . [Fraleigh, 2013, Exercise 29.16]

$$\begin{aligned}
 \alpha^3 &= \pi^6 = (\pi^3)^2 \\
 \alpha^3 - (\pi^3)^2 &= 0 \\
 \implies \phi_\alpha(x^3 - (\pi^3)^2) &= 0
 \end{aligned}$$

Thus,  $\text{irr}(\pi^2, \mathbb{Q}(\pi^3)) = x^3 - (\pi^3)^2$  and  $\deg(\pi^2, \mathbb{Q}(\pi^3)) = 3$ .  
 Note that  $x^3 - (\pi^3)^2 \in \mathbb{Q}(\pi^3)[x]$  since  $\pi^3 \in \mathbb{Q}(\pi^3) \implies -(\pi^3)^2 \in \mathbb{Q}(\pi^3)$ .  
 However,  $x^3 - (\pi^3)^2 \notin \mathbb{Q}[x]$ .

#### Factorisation over Extended Field

1. Factorisation over Finite Extension of Finite Field [Fraleigh, 2013, Exercise 29.25]  
 Let  $\alpha$  be a zero of  $f(x) = x^3 + x^2 + 1 \in \mathbb{Z}_2[x]$ . Clearly,  $f(x)$  is irreducible over  $\mathbb{Z}_2$ . And  $x - \alpha$  is a factor of  $f(x)$  in  $\mathbb{Z}_2(\alpha)$ . We have,  $f(x) = (x - \alpha)g(x)$

$$\text{By long division, } g(x) = \frac{x^3 + x^2 + 1}{x - \alpha} = x^2 + (1 + \alpha)x + (\alpha + \alpha^2)$$

Therefore,  $x^3 + x^2 + 1 = (x - \alpha)[x^2 + (1 + \alpha)x + \alpha(1 + \alpha)]$ .

The elements of  $\mathbb{Z}_2(\alpha)$  are of the form  $a_0 + a_1\alpha + a_2\alpha^2$  where  $a_0, a_1, a_2 \in \{0, 1\}$ . Also we have,  $\alpha$  is a zero of  $x^3 + x^2 + 1$ . Thus,  $\alpha^3 = \alpha^2 + 1$ .

In order to find a zero of  $g(x)$  it is sufficient to evaluate  $g(x)$  for all the eight elements  $0, 1, \alpha, (1 + \alpha), \alpha^2, (1 + \alpha^2), (\alpha + \alpha^2), (1 + \alpha + \alpha^2) \in \mathbb{Z}_2(\alpha)$ .

Clearly,  $g(1) = \alpha^2 + 2\alpha + 2 = \alpha^2$ . And  $g(\alpha) = \alpha^2 + 2\alpha(1 + \alpha) = \alpha^2$ . However,  $g(\alpha^2) = \alpha^4 + \alpha^3 + 2\alpha^2 + \alpha = \alpha(\alpha^2 + 1) + (\alpha^2 + 1) + 0\alpha^2 + \alpha = \alpha^3 + \alpha^2 + 2\alpha + 1 = \alpha^3 + \alpha^2 + 1 = 0$ . Thus,  $\alpha^2$  is a zero of  $g(x)$  and  $g(x) = (x - \alpha^2)h(x)$ .

$$\text{By long division, } h(x) = \frac{x^2 + (1 + \alpha)x + (\alpha + \alpha^2)}{x - \alpha^2} = x + (1 + \alpha + \alpha^2)$$

Therefore, we have the following linear factorisation for  $f(x)$ ,  
 $f(x) = (x - \alpha)(x - \alpha^2)(x - 1 - \alpha - \alpha^2) = (x + \alpha)(x + \alpha^2)(x + 1 + \alpha + \alpha^2)$   
 since  $-\alpha = 0 - \alpha = 2\alpha - \alpha = \alpha$  in  $\mathbb{Z}_2$ .

Note : Students should be able to perform long division of polynomials over extended fields.

## 6.2 Algebraic Extensions §31

$(G : H)$  is the number of  $H$ -left cosets in  $G$ .

**algebraic extension** A extension field  $E$  of a field  $F$  is algebraic if every element in  $E$  is algebraic over  $F$ .

For example,  $\mathbb{C}$  is algebraic over  $\mathbb{R}$ . But,  $\mathbb{R}$  is not algebraic over  $\mathbb{Q}$ .

**finite extension** A extension field  $E$  of field  $F$  is a finite extension if  $E$  is a finite dimensional vector space over  $F$ . And  $[E : F]$  is the **dimension of the vector space**  $E$  over  $F$ . Again,  $[E : F]$  is the **degree of the finite extension**  $E$  over  $F$ .

For example,  $\mathbb{C}$  is a finite extension of degree 2 over  $\mathbb{R}$ ,  $[\mathbb{C} : \mathbb{R}] = 2$ . But,  $\mathbb{R}$  is not a finite extension of  $\mathbb{Q}$  and  $[\mathbb{R} : \mathbb{Q}]$  is infinite.

**Theorem 6.2.1.** *Every finite field extensions is an algebraic extension.*

*Proof.* Let  $E$  be a finite extension of degree  $n$  over  $F$ . Then  $[E : F] = n$ . Suppose  $\alpha \in E$ . Clearly,  $\{1, \alpha, \alpha^2, \dots, \alpha^n\}$  is set of  $n + 1$  vectors from the vector space  $E$  over  $F$ . We know that in a vector space of dimension  $n$ , any set having  $n + 1$  vector is linearly dependant. In other words, there exists scalars  $c_0, c_1, \dots, c_n \in F$  (not all zero) such that  $c_0 + c_1\alpha + c_2\alpha^2 + \dots + c_n\alpha^n = 0$ .

Clearly,  $f(x) = c_0 + c_1x + c_2x^2 + \dots + c_nx^n$  is a polynomial in  $F[x]$  such that  $\phi_\alpha(f) = f(\alpha) = 0$ . Since  $\alpha \in E$  is arbitrary, every element in  $E$  is algebraic over  $F$ .  $\square$

**Theorem 6.2.2.** *If  $E$  is a finite extension of  $F$  and  $K$  is a finite extension of  $E$ . Then  $K$  is a finite extension of  $F$ . And  $[K : F] = [K : E][E : F]$ .*

*Proof.* Let  $[E : F] = n$  and  $[K : E] = m$ . Let  $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$  be basis for vector space  $E(F)$  and let  $\{\beta_1, \beta_2, \dots, \beta_m\}$  be basis for vector space  $K(E)$ . We claim that  $\{\alpha_i \beta_j : 1 \leq i \leq n, 1 \leq j \leq m\}$  is a basis for the vector space  $K(F)$ .

Let  $\gamma \in K$ . Then we have  $b_1, b_2, \dots, b_m \in E$  such that  $\gamma = b_1 \beta_1 + b_2 \beta_2 + \dots + b_m \beta_m$ . Again, for each  $b_j \in E$ , we have  $a_{ij} \in F$  such that  $b_j = a_{1j} \alpha_1 + a_{2j} \alpha_2 + \dots + a_{nj} \alpha_n$ . Therefore,

$$\gamma = \sum_{i=1}^n \sum_{j=1}^m a_{ij} \alpha_i \beta_j$$

That is,  $\{\alpha_i \beta_j : i = 1, 2, \dots, n, j = 1, 2, \dots, m\}$  spans  $K$ . It remains to prove that  $\{\alpha_i \beta_j\}$  is linearly independent. Suppose it is linearly dependent. Then there exists scalars  $c_{i,j} \in F$  (not all zero) such that

$$\sum_{i=1}^n \sum_{j=1}^m c_{i,j} \alpha_i \beta_j = \sum_{j=1}^m \left( \sum_{i=1}^n c_{i,j} \alpha_i \right) \beta_j = 0$$

Let  $\sum_{i=1}^n c_{i,j} \alpha_i = b_j$ . We know that,  $\{\beta_j : j = 1, 2, \dots, m\}$  is linearly independent. Thus,  $\sum_j b_j \beta_j = 0 \implies b_j = 0, \forall j$ .

Again  $b_j = 0 \implies \sum_{i=1}^n c_{i,j} \alpha_i = 0$ . Once again,  $\{\alpha_i : i = 1, 2, \dots, n\}$  is linear independent. Thus,  $c_{i,j} = 0, \forall i, j$ . Thus,  $\{\alpha_i \beta_j\}$  is linearly independent. Therefore,  $\{\alpha_i \beta_j\}$  is a basis for the vector space  $K(F)$ . And  $[K : F] = |\{\alpha_i \beta_j : i = 1, 2, \dots, n \text{ and } j = 1, 2, \dots, m\}| = mn$ .  $\square$

**Corollary 6.2.2.1.** *Let  $F_i$  be fields and  $F_{i+1}$  are finite extensions of  $F_i$ s for  $i = 1, 2, \dots, r$ . Then  $[F_r : F_1] = [F_r : F_{r-1}][F_{r-1} : F_{r-2}] \cdots [F_2 : F_1]$ .*

*Proof.* We have,  $F_3$  is a finite extension of  $F_1$  and

$$[F_3 : F_1] = [F_3 : F_2][F_2 : F_1] \quad (6.2)$$

Suppose  $F_k$  is a finite extension of  $F_1$  and

$$[F_k : F_1] = [F_k : F_{k-1}][F_{k-1} : F_{k-2}] \cdots [F_2 : F_1] \quad (6.3)$$

$$\begin{aligned} [F_{k+1} : F_1] &= [F_{k+1} : F_k][F_k : F_1] \text{ since } F_k \text{ is a finite extension of } F_1 \\ &= [F_{k+1} : F_k][F_k : F_{k-1}][F_{k-1} : F_{k-2}] \cdots [F_2 : F_1] \end{aligned}$$

$\square$

**Corollary 6.2.2.2.** *If  $E$  is an extension field of  $F$  and  $\alpha \in E$  is algebraic over  $F$  and  $\beta \in F(\alpha)$ , then  $\deg(\beta, F)$  divides  $\deg(\alpha, F)$ .*

*Proof.* We have,  $\deg(\alpha, F) = [F(\alpha) : F]$  and  $\deg(\beta, F) = [F(\beta) : F]$ . Also given that  $\beta \in F(\alpha) \implies F(\beta) \leq F(\alpha)$ . Clearly,  $F \leq F(\beta) \leq F(\alpha)$ . Therefore,  $[F(\alpha) : F] = [F(\alpha) : F(\beta)][F(\beta) : F]$ . And  $[F(\alpha) : F(\beta)] = [F(\alpha) : F]/[F(\beta) : F]$ . Clearly,  $[F(\beta) : F]$  divides  $[F(\alpha) : F]$ .  $\square$

**Theorem 6.2.3** (algebraic closure). *Let  $E$  be an extension field of  $F$ . Then  $\bar{F}_E = \{\alpha \in E : \alpha \text{ is algebraic over } E\}$  is a subfield of  $E$ .*

*Proof.* Let  $\alpha, \beta \in \bar{F}_E$ . Then  $\alpha, \beta \in E$  are algebraic over  $F$ . And  $F(\alpha, \beta)$  is a finite extension field of  $F$ . Thus every element in  $F(\alpha, \beta)$  are algebraic over  $F$ . Thus,  $\alpha + \beta, \alpha\beta, \alpha - \beta, \alpha/\beta \in \bar{F}_E$ . Therefore,  $\bar{F}_E$  is a subfield of  $E$ .  $\square$

**Corollary 6.2.3.1.** *The set of all algebraic numbers forms a field.*

*Proof.* Let  $\alpha$  be an algebraic number. Then  $\alpha \in \mathbb{C}$  and  $\alpha$  is algebraic over  $\mathbb{Q}$ . Clearly, the set of all algebraic numbers,  $\bar{\mathbb{Q}}$  is a subfield of  $\mathbb{C}$ .  $\square$

**algebraic closure** Let  $F$  be a field and  $E$  be an extension field of  $F$ . Then the (smallest) field containing all elements of  $E$  which are algebraic over  $F$  is the algebraic closure  $\bar{F}_E$  of  $F$  in  $E$ .

**algebraically closed** Let  $F$  be a field.  $F$  is algebraically closed if every non-constant polynomial in  $F[x]$  has a zero in  $F[x]$ .

**Note :** Let  $F$  be algebraically closed. Then every irreducible polynomial in  $F[x]$  are linear since every non-constant polynomial has a linear factor.

**Theorem 6.2.4.** *A field  $F$  is algebraically closed if and only if every non-constant polynomial  $f(x)$  can be factorised in  $F[x]$  into linear factors.*

*Proof.* Let  $F$  be algebraically closed and  $f(x)$  be a non-constant polynomial in  $F[x]$ . Then  $f(x)$  has a zero  $\alpha \in F$ . Then  $x - \alpha$  is a factor of  $f(x)$ . That is,  $f(x) = (x - \alpha)g(x)$ . If  $g(x) \in F[x]$  is non-constant, then it has a zero in  $F$ . Continuing like this,  $f(x)$  can be factorised in  $F[x]$  into linear factors.

Suppose every non-constant polynomial in  $F[x]$  can be factorised into linear factors. Let  $f(x)$  be a non-constant polynomial in  $F[x]$ . Then  $f(x)$  has a linear factor  $(ax + b) \in F[x]$ . Clearly,  $-b/a$  is a zero of  $f(x)$ .  $\square$

**Theorem 6.2.5.** *Algebraically closed field has no proper algebraic extensions.*

*Proof.* Let  $E$  be an algebraic extension field of  $F$ . Then if  $\alpha \in E$ , we have  $\text{irr}(\alpha, F) = (x - \alpha)$  since  $F$  is algebraically closed, every irreducible polynomial in  $F[x]$  are linear. Thus  $\alpha \in F$ . Since  $\alpha \in E$  is arbitrary,  $F = E$ .  $\square$

**Theorem 6.2.6** (Fundamental Theorem of Algebra).  *$\mathbb{C}$  is algebraically closed.*

*Proof.* Let  $f(z) \in \mathbb{C}[z]$ . Suppose  $f(z)$  has no zeroes in  $\mathbb{C}$ . Then  $1/f(z)$  is an entire function and as  $|z| \rightarrow \infty$ ,  $|f(z)| \rightarrow \infty$ . Thus,  $\lim_{|z| \rightarrow \infty} \frac{1}{|f(z)|} = 0$ . And  $1/f(z)$  is bounded.

By Liouville's theorem, every bounded entire function is constant. Therefore,  $1/f(z)$  is constant and  $f(z)$  is also constant. Thus, every non-constant polynomial function in  $\mathbb{C}[z]$  has a zero in  $\mathbb{C}$ .  $\square$

**POSET** Partial Ordered Set - A set together with partial order (reflexive, antisymmetric, transitive relation).

For example  $(\mathbb{R}, <)$ , the set of all real numbers together with less than relation is a poset. However, in a poset it is not necessary that two arbitrary elements are comparable.  $(\mathbb{C}, R)$  defined by  $aRb$  if  $\Re(a) = \Re(b)$  and  $\text{Im}(a) < \text{Im}(b)$  is a poset in which  $2 + 3i$  and  $3 + 3i$  are not comparable.



**chain** A subset of a poset in which any two elements are comparable. That is,  
 $x, y \in T \implies x < y \text{ OR } y < x.$

For example, For above defined poset  $(\mathbb{C}, R)$ ,  $T = \{2 + ib \in \mathbb{C} : b \in \mathbb{R}\}$  is a chain in  $\mathbb{C}$ .

**Lemma 6.2.7 (Zorn).** *If every chain in a poset  $S$  has an upper bound. Then  $S$  has at least one maximal element in it.*

*Proof.* Not required ( I think, there is no proof. We just take it as an axiom - always true !. If it is not true for a collection then it is not a set !! )  $\square$

**Theorem 6.2.8 (Existence of Algebraic Closure).** *Every field  $F$  has an algebraic closure  $\bar{F}$ .*

*Proof.* Not required (as per syllabus)  $\square$

### 6.2.1 Exercise §31

1.

## 6.3 Geometric Constructions §32

### 6.3.1 Basic Constructions

#### Finding Midpoint of a line

Let  $OA$  be a line. The line passing through the intersection of circles with center  $O$  and  $A$  with diameter greater than the length of the line gives a perpendicular line through its mid point (say, perpendicular bisection).

#### Drawing Perpendicular line through a point

Let  $OA$  be a line and  $B$  be a point on that line. Find points  $P, Q$  on  $OA$  which are equidistant from  $B$ . Then the perpendicular bisection of  $PQ$  is a line perpendicular to  $OA$  through  $B$ .

#### Drawing Parallel Line

Let  $OA$  be a line. Then the perpendicular line segment of any perpendicular line segment is a line segment parallel to  $OA$ .

### 6.3.2 Constructible Numbers

**Constructible Number** A real number  $\alpha$  is constructible if you can draw a line of length  $|\alpha|$ , given a line of unit length, in finite steps using straight-edge and compass.

**Theorem 6.3.1.** *Let  $\alpha, \beta$  be constructible real numbers. Then  $\alpha + \beta$ ,  $\alpha - \beta$ ,  $\alpha\beta$ ,  $\alpha/\beta$  ( $\beta \neq 0$ ) are also constructible.*

*Proof.*  $\alpha + \beta$  Draw a line  $OA$  of length  $|\alpha|$  and extend that line using straight edge  $OE$ . And draw the line  $AB$  of length  $|\beta|$  on that extended line, at an point  $A$  of the former line extending it. Then  $OB$  is a line of length  $|\alpha| + |\beta|$ .



$\alpha - \beta$  Draw a line  $OA$  of length  $|\alpha|$  and extend that line using straight edge  $OE$ . And draw the line  $AB$  of length  $|\beta|$  on that extended line, at an end point  $A$  of former line, but in the opposite direction of extension (towards  $O$ ). Then the line  $OB$  is of length  $|\alpha| - |\beta|$ .



$\alpha\beta$  Draw two lines  $OA$  and  $OB$  of length  $|\alpha|$  and  $|\beta|$  respectively (with a common end point  $O$ ). Now draw the line of unit length  $OP$  along the line  $OB$ . Construct triangle  $OAP$ . Draw the line  $BQ$  parallel to the line  $PA$  through  $B$  such that  $OQ$  and  $OA$  are colinear. Now we have two similar triangles  $OPA$  and  $OBQ$ . Then, the line  $OQ$  of length  $\|\alpha\beta\|$ .



$\alpha/\beta$  Draw two lines  $OA$  and  $OB$  of length  $\|\alpha\|$  and  $\|\beta\|$ , with a common end point  $O$ . Draw a line of unit length  $OP$  along the line  $OB$ . Now construct triangle  $OBA$ . Draw line  $PQ$  parallel to  $BA$  so that  $OQ$  and  $OA$  are colinear. Now the triangles  $OAB$  and  $OQP$  are similar. And the line  $OQ$  has length  $\|\alpha/\beta\|$ .



□

**Corollary 6.3.1.1.** *The set of all constructible real numbers forms a subfield of the field of real numbers.*

*Proof.* The set of all constructible real numbers say  $H$ , contains both 0 and 1. since the line of length zero is trivial and line of unit length is provided. And we have,  $\forall \alpha, \beta \in H, \alpha + \beta, \alpha - \beta, \alpha\beta, \alpha/\beta \in H$ . Since  $\alpha - \beta \in H, 0 - \beta = -\beta$  which is the additive inverse of  $\beta$ . And  $1/\beta = \beta^{-1}$  is the multiplicative inverse of  $\beta$ . Thus, the set of all constructible numbers is a subfield of  $\mathbb{R}$ .  $\square$

**Theorem 6.3.2.** *The field of  $F$  of constructible numbers consists **precisely** of all real numbers that we can obtain from  $\mathbb{Q}$  by taking square root of positive numbers a finite number of times and applying a finite number of field operations.*

*Proof.* The constructible numbers are closed under field operations and forms a subfield  $H$  of real numbers. We have,  $\mathbb{Q}$  is the prime field of  $\mathbb{R}$ . That is, every subfield of  $\mathbb{R}$  contains  $\mathbb{Q}$ . Thus, the subfield  $H$  of constructible numbers contains  $\mathbb{Q}$ . That is, all the rational numbers are constructible. Therefore, it remains to prove that if  $\alpha > 0$  is constructible then  $\sqrt{\alpha}$  is also constructible.



Let  $OA$  and  $OP$  be colinear lines such that  $OA$  be a line of length  $\|\alpha\|$  and  $OP$  be a line of unit length. Find the mid point of  $PA$  and draw a circle of diameter  $PA$ . Then the length of the perpendicular  $OQ$  from  $PA$  to the circle is of length  $\|\sqrt{\alpha}\|$  since  $\triangle OAQ$  and  $\triangle OQP$  are similar triangles.

Therefore, real numbers obtained from rational numbers through finite number of additions/subtractions, multiplications/divisions, and square root are constructible.  $\square$

**Note :** For example  $\sqrt[4]{5\sqrt[8]{3} - 2}$  is constructible, but  $\sqrt[6]{2}$  and  $\pi$  are not constructible. We skip the proof that a real number which can't be obtained from rationals by a finite number of these operations is not constructible. And assume that these three operations define the entire field of constructible numbers.

**Corollary 6.3.2.1.** *Let  $\gamma$  be constructible number which is not rational. Then there exists a sequence of real numbers  $\alpha_1, \alpha_2, \dots, \alpha_n = \gamma$  such that for every  $i = 2, \dots, n$ , the extension field  $\mathbb{Q}(\alpha_1, \alpha_2, \dots, \alpha_i)$  is an extension of  $\mathbb{Q}(\alpha_1, \alpha_2, \dots, \alpha_{i-1})$  of degree two.*

*In other words, for any constructible number  $\gamma$ ,  $[\mathbb{Q}(\gamma) : \mathbb{Q}] = 2^n$  for some positive integer  $n$ .*

*Proof.* Let  $\gamma$  be a constructible number which can be obtained from rationals by  $n$  square root operations and a finite number of field operations. Then, we have a sequence of constructible numbers  $\alpha_1, \alpha_2, \dots, \alpha_n$  such that  $[\mathbb{Q}(\alpha_1, \alpha_2, \dots, \alpha_i) : \mathbb{Q}(\alpha_1, \alpha_2, \dots, \alpha_{i-1})] = 2$ . Therefore,  $[\mathbb{Q}(\gamma) : \mathbb{Q}] = 2^n$ .  $\square$

For example,  $\gamma = \sqrt[4]{5\sqrt[8]{3}-2}$ . Then  $\alpha_1 = \sqrt{3}$ ,  $\alpha_2 = \sqrt[4]{3}$ ,  $\alpha_3 = \sqrt[8]{3}$ ,  $\alpha_4 = \sqrt{5\sqrt[8]{3}-2}$  and  $\alpha_5 = \gamma$ . Clearly, the geometric construction of  $\gamma$  contains five instances of square root operation.

### 6.3.3 Impossible Problems from Ancient Times

#### Doubling the cube

**Theorem 6.3.3.** *There exists a cube such that it is impossible to construct the side of the cube with double the volume.*

*Proof.* Suppose the cube is of unit side. Then the side of the cube with double the volume is  $\gamma = \sqrt[3]{2}$ . And  $\text{irr}(\gamma, \mathbb{Q}) = x^3 - 2$  and  $[\mathbb{Q}(\sqrt[3]{2} : \mathbb{Q})] = 3 \neq 2^n$ . Therefore,  $\gamma$  is not constructible.  $\square$

#### Squaring the circle

**Theorem 6.3.4.** *There exists a circle such that it is impossible to construct the side of the square with the same area.*

*Proof.* Suppose the circle is of unit radius. Then area of the circle is  $2\pi$ . And the side of the square with same area is  $\sqrt{2\pi}$ . Since  $\pi$  is transcendental,  $\pi, \sqrt{\pi}$  and  $\sqrt{2\pi}$  are not constructible.  $\square$

#### Trisecting an angle

**Theorem 6.3.5.** *There exists an angle which can be trisected.*

*Proof.* We have  $\cos 3\theta = 4\cos^3 \theta - 3\cos \theta$ . Consider  $\theta = 20^\circ$ . Then  $\gamma = \cos \theta$  is a root of the irreducible polynomial  $4x^3 - 3x - 0.5$ . That is,  $\gamma$  is a root of the monic irreducible polynomial  $x^3 - \frac{3}{4}x - \frac{1}{8}$  of degree 3. Thus,  $\gamma$  is not constructible, since  $[\mathbb{Q}(\gamma) : \mathbb{Q}] = 3 \neq 2^n$ .  $\square$

### 6.3.4 Exercise §32

- 1.

## 6.4 Finite Fields §33

### 6.4.1 Structure of a Finite Field

**Theorem 6.4.1.** *Let  $E$  be a finite extension of degree  $n$  over a finite field  $F$ . If  $F$  has  $q$  elements, then  $E$  has  $q^n$  elements.*

*Proof.* We have, extension field  $E$  is an  $n$ -dimensional vector space over the field  $F$ . Let  $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$  be a basis of the vector space  $E$  over  $F$ . Then every elements of  $E$  can be uniquely written as linear combination of basis vectors.

$$\forall \beta \in E, \beta = b_1\alpha_1 + b_2\alpha_2 + \dots + b_n\alpha_n$$

Suppose  $\beta \in E$  has two distinct linear combinations. Then, the vectors  $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$  are not linearly independent, which is a contradiction.

Since the representation is unique and  $F$  has  $q$  elements, there are  $q^n$  distinct linear combinations possible. Therefore,  $E$  has  $q^n$  elements.  $\square$

**Corollary 6.4.1.1.** *If  $E$  is a finite field of characteristic  $p$ . Then  $F$  contains exactly  $p^n$  elements for some integer  $n$ .*

*Proof.* We have,  $E$  is a finite field of characteristic  $p$ . Thus,  $\mathbb{Z}_p$  is the prime subfield of  $E$ . And  $E$  is a finite extension of  $\mathbb{Z}_p$ . Thus,  $E$  is a finite dimensional vector space over  $\mathbb{Z}_p$ . Let  $n$  be the dimension of  $E$  over  $\mathbb{Z}_p$ . And  $\mathbb{Z}_p$  has  $p$  elements. Therefore,  $E$  has  $p^n$  elements.  $\square$

**Theorem 6.4.2.** *Let  $E$  be a field of  $p^n$  elements contained in the algebraic closure  $\overline{\mathbb{Z}_p}$  of  $\mathbb{Z}_p$ . Then the elements of  $E$  are precisely the zeroes of the polynomial  $x^{p^n} - x \in \mathbb{Z}_p[x]$ .*

*Proof.* We have  $E^*$  is a multiplicative group of non-zero elements in  $E$ . And  $E^*$  has  $p^n - 1$  elements. Thus, order of any element  $\alpha \in E^*$  should divide  $p^n - 1$ . In other words, if  $\alpha \in E^*$ , then  $\alpha^{p^n-1} = 1$ . Clearly,  $\alpha^{p^n} - \alpha = 0, \forall \alpha \in E$ . However,  $x^{p^n}$  can have atmost  $p^n$  zeroes in  $\overline{\mathbb{Z}_p}$ . Thus,  $E$  is precisely the set of all zeroes of  $x^{p^n} - x \in \overline{\mathbb{Z}_p}$ .  $\square$

**$n$ th Root of Unity**  $\alpha$  is an  $n$ th root of unity if  $\alpha^n = 1$ . ie,  $\alpha = \sqrt[n]{1}$

**Primitive  $n$ th Root of Unity**  $\alpha$  is a primitive  $n$ th root of unity if  $n$  is the smallest positive integer such that  $\alpha^n = 1$ .

That is,  $\alpha^n = 1$  and  $\forall m \in \mathbb{N}, m < n \implies \alpha^m \neq 1$

**Theorem 6.4.3.** *The multiplicative group of non-zero elements of a finite field  $F$  is cyclic.*

*Proof.* Refer : [Fraleigh, 2013, Theorem 23.6]  $\square$

**Corollary 6.4.3.1.** *Finite extension of finite fields are simple extensions.*

*Proof.* Let  $E$  be a finite extension field of the finite field  $F$ . Then the multiplicative group of non-zero elements  $E^*$  is cyclic. Let  $\alpha$  be a generator of the cyclic group  $E^*$ . Then,  $E = F(\alpha)$ .  $\square$

## 6.4.2 Galois Field $GF(p^n)$

**Lemma 6.4.4.** *If  $F$  is a field of prime characteristic  $p$  with algebraic closure  $\overline{F}$ , then  $x^{p^n} - x$  has  $p^n$  distinct zeroes in  $\overline{F}$ .*

*Proof.* We have,  $\overline{F}$  is algebraically closed. And  $x^{p^n} - x \in \overline{F}[x]$ . Thus,  $x^{p^n} - x$  can be factorised into  $p^n$  linear components. It remains to prove that these factors are distinct.

Clearly, 0 is a zero of multiplicity 1, since  $x^{p^n} - x = x(x^{p^n-1} - 1)$ . Let  $\alpha \neq 0$  be a zero of  $x^{p^n} - x$ . Then  $\alpha$  is a zero of  $x^{p^n-1} - 1$ . ie,  $\alpha^{p^n-1} = 1$ .

$$(x - \alpha)g(x) = x^{p^n-1} - 1$$

$$g(x) = \frac{x^{p^n-1} - 1}{x - \alpha}$$

By long division, we get

$$\begin{aligned}
 g(x) &= x^{p^n-2} + \alpha x^{p^n-3} + \alpha^2 x^{p^n-4} + \cdots + \alpha^{p^n-3} x + \alpha^{p^n-2} \\
 g(\alpha) &= (p^n - 1)\alpha^{p^n-2} \\
 &= (p^n - 1) \frac{\alpha^{p^n-1}}{\alpha} \\
 &= p^n \frac{1}{\alpha} - \frac{1}{\alpha} \\
 &= -\frac{1}{\alpha} \neq 0
 \end{aligned}$$

Thus, every zero of  $x^{p^n-1} - x$  is of multiplicity 1.  $\square$

**Lemma 6.4.5.** *If  $F$  is a field of prime characteristic  $p$ , then  $(\alpha + \beta)^{p^n} = \alpha^{p^n} + \beta^{p^n}$ .*

*Proof.* Since  $F$  is a field of characteristic  $p$ , for every  $\alpha \in F$ ,  $p\alpha = 0$ .

For  $n = 1$ , we have

$$\begin{aligned}
 (\alpha + \beta)^p &= \alpha^p + p\alpha^{p-1}\beta + \frac{p(p-1)}{2}\alpha^{p-2}\beta^2 + \cdots + p\alpha\beta^{p-1} + \beta^p \\
 &= \alpha^p + 0\alpha^{p-1}\beta + 0\alpha^{p-2}\beta^2 + \cdots + 0\alpha\beta^{p-1} + \beta^p \\
 &= \alpha^p + \beta^p
 \end{aligned}$$

Suppose  $(\alpha + \beta)^{p^{n-1}} = \alpha^{p^{n-1}} + \beta^{p^{n-1}}$ , then

$$\begin{aligned}
 (\alpha + \beta)^{p^n} &= \left[ (\alpha + \beta)^{p^{n-1}} \right]^p \\
 &= \left[ \alpha^{p^{n-1}} + \beta^{p^{n-1}} \right]^p \\
 &= \alpha^{p^n} + \beta^{p^n}
 \end{aligned}$$

Therefore, by mathematical induction the result is true.  $\square$

**Theorem 6.4.6** (Existence of Galois Field). *For every prime power  $p^n$ , a finite field of  $p^n$  elements exists.*

**Hint :**  $\alpha^{p^n} = \alpha \iff \alpha^{p^n} - \alpha = 0 \iff \alpha$  is a zero of  $x^{p^n} - x$

*Proof.* Consider the algebraic closure  $\overline{\mathbb{Z}_p}$  of  $\mathbb{Z}_p$ . Let  $K$  be a subset of  $\overline{\mathbb{Z}_p}$  containing all zeroes of  $x^{p^n} - x \in \overline{\mathbb{Z}_p}$ .

Let  $\alpha, \beta \in K$ . Then  $\alpha^{p^n} = \alpha$  and  $\beta^{p^n} = \beta$ . Since  $\overline{\mathbb{Z}_p}$  is a field of characteristic  $p$ , we have  $(\alpha + \beta)^{p^n} = \alpha^{p^n} + \beta^{p^n} = \alpha + \beta$ . Thus,  $(\alpha + \beta)$  is a zero of  $x^{p^n} - x$ . ie,  $(\alpha + \beta) \in K$ . Clearly,  $(\alpha\beta)^{p^n} = \alpha^{p^n}\beta^{p^n} = \alpha\beta$ . Thus,  $\alpha\beta \in K$ .

Again  $(-\alpha)^{p^n} = (-1 \cdot \alpha)^{p^n} = (-1)^{p^n} \cdot \alpha^{p^n} = -1 \cdot \alpha = -\alpha$ . Thus,  $-\alpha \in K$ . Also  $(\alpha^{-1})^{p^n} = \left(\frac{1}{\alpha}\right)^{p^n} = \frac{1}{\alpha^{p^n}} = \frac{1}{\alpha} = \alpha^{-1}$ . Thus,  $\alpha^{-1} \in K$ .

Trivially,  $0, 1 \in K$ . Therefore,  $K$  is a subfield of  $\overline{\mathbb{Z}_p}$  with  $p^n$  elements since the  $p^n$  zeroes of  $x^{p^n} - x$  are distinct.  $\square$

**Corollary 6.4.6.1.** *If  $F$  is a finite field, then for every positive integer  $n$ , there exists an irreducible polynomial in  $F[x]$  of degree  $n$ .*

*Proof.* Let  $F$  be a finite field. Then  $\mathbb{Z}_p$  is prime field of  $F$  for some prime  $p$ . And  $F$  is of characteristic  $p$  and has  $p^r$  elements for some positive integer  $r$ .

Let  $K$  be the subfield of  $\overline{F}$  containing precisely all the zeroes of the polynomial  $x^{p^{rn}} - x \in \mathbb{Z}_p[x]$ . Then,  $K$  has a subfield isomorphic to  $\mathbb{Z}_p$  since  $F$  is of characteristic  $p$  and every subfield of  $F$  has a subfield isomorphic to  $\mathbb{Z}_p$ .

By existence theorem of Galois Fields, every element of  $F$  is a zero of the polynomial  $x^{p^r} - x \in \mathbb{Z}_p[x]$ . That is,  $\alpha \in F \iff \alpha^{p^r} = \alpha$ .

$$\begin{aligned} \alpha^{p^{rn}} &= \left[ \alpha^{p^r} \right]^{p^{r(n-1)}} &&= \alpha^{p^{r(n-1)}} \\ &= \left[ \alpha^{p^r} \right]^{p^{r(n-2)}} &&= \alpha^{p^{r(n-2)}} \\ &\vdots &&\vdots \\ &= \left[ \alpha^{p^r} \right]^{p^r} &&= \alpha^{p^r} = \alpha \end{aligned}$$

Thus,  $\alpha \in F \implies \alpha^{p^{rn}} = \alpha \implies \alpha \in K$ . Therefore  $F$  is a subfield of  $K$ . Clearly,  $K$  is a finite extension of the finite field  $F$ . And the vector space  $K$  over  $F$  is  $n$ -dimensional, since  $K$  has  $p^{rn} = [p^r]^n$  elements and  $F$  has  $p^r$  elements.

Since every finite extension of finite fields are simple extensions, we have  $K = F(\beta)$  and  $\text{irr}(\beta, F) = n$ . That is, there exists an unique monic, irreducible polynomial  $p(x) \in F[x]$  of degree  $n$  such that  $p(\beta) = 0$ . Therefore,  $\forall n \in \mathbb{Z}^+$ , there exists an irreducible polynomial in  $F[x]$  of degree  $n$ .  $\square$

**Theorem 6.4.7** (Uniqueness of Galois Field). *Let  $p$  be a prime and  $n$  a positive integer. If  $E$  and  $E'$  are fields of order  $p^n$ , then  $E$  and  $E'$  are isomorphic.*

*There exists a unique finite field of order  $p^n$ , say **Galois Field**,  $GF(p^n)$ .*

*Proof.* Let  $E, E'$  be fields of order  $p^n$ . Then both fields have  $\mathbb{Z}_p$  as prime field. Thus,  $E$  is a simple extension of  $\mathbb{Z}_p$  of degree  $n$ . ie,  $[E : \mathbb{Z}_p] = n$ . And there exists an irreducible polynomial  $f(x)$  such that  $E \simeq \mathbb{Z}_p[x]/\langle f(x) \rangle$ .

Elements of  $E$  are zeroes of  $x^{p^n} - x$ , thus  $f(x)$  is a factor of  $x^{p^n} - x$ . Clearly, elements of  $E'$  are zeroes of  $x^{p^n}$  and therefore  $E'$  has all zeroes of  $f(x)$ . And  $E' \simeq \mathbb{Z}_p[x]/\langle f(x) \rangle \simeq E$ . Therefore, there exists a unique field of order  $p^n$  (upto isomorphism), say  $GF(p^n)$ .  $\square$

### 6.4.3 Exercise §33

1.

## 6.5 Unique Factorisation Domains §45

**Definitions 6.5.1** (divides). An element  $a \in R$  divides  $b$  if there exists an element  $c \in R$  such that  $b = ac$ .

**Definitions 6.5.2** (associate). Two elements  $a, b$  are associates if  $a = bu$  where  $u$  is a unit.

**Definitions 6.5.3** (UFD). **Unique factorisation domain**, UFD is an integral domain such that

1. Every element can be factored into a finite number of irreducibles, except 0 and units  $\dagger^1$ .
2. The above factorisation is unique except for order and associates.

For example, In  $\mathbb{Z}$ ,  $24 = 2 \times 2 \times 2 \times 3 = -2 \times -3 \times 2 \times 2$ . Here 2 and  $-2$  are associates. And 2 and 3 are not units, since  $2^{-1}, 3^{-1} \notin \mathbb{Z}$ .

**Definitions 6.5.4** (PID). An integral domain is  $D$  is a **Principal Ideal Domain** if every ideal in  $D$  is a principal ideal.

There are two important results on UFDs.

1. Every PID is a UFD.
2. If  $D$  is a UFD, then  $D[x]$  is a UFD.

### 6.5.1 Every PID is UFD

**Lemma 6.5.1.** Let  $R$  be a commutative ring and let  $N_1 \subset N_2 \subset \dots$  be an ascending chain of ideals in  $R$ . Then  $N = \cup_i N_i$  is an ideal of  $R$ .

*Proof.* **Step 1 :  $N$  is a subring of  $R$**

Suppose  $N_i, N_j$  are two ideals in the chain,  $N_i \subset N_j$  and  $a \in N_i, b \in N_j$ . Clearly,  $a \in N_j$ .

And  $a \pm b, ab \in N_j$ . And  $0 \in N_j \implies 0 \in N$ . We have, 0 in every ideal<sup>†2</sup>. Take  $a = 0$ , we know that for every element  $b \in N_j$ , its additive inverse  $-b \in N_j$ . Thus  $b \in N \implies b \in N_j \implies -b \in N_j \implies -b \in N$ . Clearly,  $N$  is a subring of  $R$ .

**Step 2 :  $N$  is a ideal of  $R$**

Let  $a \in N$  and  $r \in R$ . We have,  $a \in N \implies a \in N_j$  for some ideal  $N_j$  in the chain. Since  $N_j$  is an ideal  $ar = ra \in N_j$ . Therefore,  $ar \in N$ . And  $N$  is a ideal of  $R$ .  $\square$

**Lemma 6.5.2** (Ascending Chain Condition). Let  $D$  be a PID. If  $N_1 \subset N_2 \subset \dots$  is an ascending chain of ideal, then there exists a positive integer  $r$  such that  $N_r = N_s$  for every  $s \geq r$ .

In other words, every strictly<sup>†3</sup> ascending chain of ideals in a PID is of finite length.

<sup>†1</sup>Unit is a element which has multiplicative inverse.

<sup>†2</sup>Suppose  $b \in N$ . We have  $-b \in R$  and  $-b + b = 0 \in bN \subset N$ .

<sup>†3</sup>Strictly ascending chain :  $N_1 \subsetneq N_2 \subsetneq N_3 \subsetneq \dots \subsetneq N_k$ .



*Proof.* Let  $D$  be an integral domain. And  $N_1 \subseteq N_2 \subseteq \dots$  be an ascending chain of ideals in  $D$ . Then  $N = \cup_i N_i$  is an ideal in  $D$ . Since  $D$  is a PID, by definition of **p**ri**n**cip**a**l **i**de**a**l **d**om**a**in every ideal in it is a principal ideal. Thus,  $N$  is a principal ideal in  $D$ . That is, there exists  $c \in D$  such that  $\langle c \rangle = N$ .

Since  $c \in N$ , and  $N = \cup_i N_i$  we have  $c \in N_r$  for some  $r \in \mathbb{N}$ . Then  $\langle c \rangle = N_r = N$ . Again,  $N_r \subset N_s \implies c \in N_s \implies \langle c \rangle = N_s = N$ .

Suppose the chain of ideals are strictly ascending, that is every ideal in the chain is properly containing the former ideal. Then, the ideal  $N_r$  such that  $c \in N_r$  is the last ideal in its chain. Thus, strictly ascending chain of ideals is finite.  $\square$

**Theorem 6.5.3.** *Let  $D$  be a PID. Every element that is neither 0 nor a unit in  $D$  is a product of irreducibles.*

*Proof.* Suppose  $D$  is a PID. And suppose  $a \in D$  is a neither a zero nor a unit. If  $a$  is an irreducible, then the result is trivial. Suppose  $a$  is not an irreducible.

**Step 1 :  $a$  has an irreducible factor**

By the definition of irreducibility,  $a$  has a factorisation  $a = a_1 b_1$  where  $a_1, b_1$  are non-units. And every element  $ar \in \langle a \rangle$  can be expressed as  $ar = a_1 b_1 r = a_1 r' \in \langle a_1 \rangle$ . Thus, the ideal generated by  $a$ ,  $\langle a \rangle$  is contained in the ideal generated by  $a_1$ ,  $\langle a_1 \rangle$ . That is,  $\langle a \rangle \subset \langle a_1 \rangle$ .

If  $a_1$  is an irreducible, then the proof is complete. Suppose  $a_1$  is not an irreducible. That is,  $a_1 = a_2 b_2$ . Then  $\langle a_1 \rangle \subset \langle a_2 \rangle$ . Also we have,  $\langle a_1 \rangle \neq \langle a_2 \rangle$ .

Suppose  $\forall d \in D$ ,  $a_2 d \in \langle a_1 \rangle \implies a_2 \cdot 1 = a_2 \in \langle a_1 \rangle \implies a_2 = a_1 c_2 \implies a_1 = (a_1 c_2) b_2 \implies c_2 b_2 = 1 \implies b_2$  has multiplicative inverse  $c_2$ . This contradicts the assumption that  $b_2$  is a unit.

Continuing like this, we get a **strictly** ascending chain of ideals  $\langle a \rangle \subsetneq \langle a_1 \rangle \subsetneq \langle a_2 \rangle \dots$ . And we know that, every strictly ascending chain of ideals is finite. Thus,  $a_r$  is an irreducible. And  $a = a_1 b_1 = (a_2 b_2) b_1 = \dots = a_r (b_r b_{r-1} \dots b_1) = a_r b$ . Clearly,  $a$  has an irreducible factor  $a_r$ .

**Step 2 :  $a$  is a product of irreducibles**

Suppose  $a$  is neither zero nor a unit. Then  $a = p_1 c_1$  where  $p_1$  is an irreducible and  $c_1$  is not a unit. Continuing like this, we get another strictly ascending chain of ideals  $\langle a \rangle \subsetneq \langle c_1 \rangle \subsetneq \langle c_2 \rangle \subsetneq \dots$  such that  $a = p_1 c_1 = p_1 (p_2 c_2) = \dots$  where  $p_1, p_2, \dots$  are irreducibles and  $c_j$  are non-units. Again, by ascending chain condition this strictly ascending chain is finite. That is,  $a = p_1 p_2 \dots p_k c_k$  for some  $k \in \mathbb{N}$ . Since  $\langle c_k \rangle$  is a maximal ideal in  $D$ ,  $c_k$  is an irreducible say,  $p_{k+1}$ . Thus, any element  $a \in D$  (except zero and units) can be expressed as product of irreducibles.  $\square$

## 6.5.2 Irreducible element and Maximal Ideal

We have seen earlier that a non-trivial principal ideal  $\langle p(x) \rangle$  in  $F[x]$  is maximal if and only if the generator  $p(x)$  is irreducible over  $F$ . Now we have a gener-

alisation, which says ideals in PIDs are maximal if and only if the generator is irreducible. Remember that for any field  $F$ ,  $F[x]$  is a PID.

**Lemma 6.5.4.** *An ideal  $\langle p \rangle$  in a PID is maximal if and only if  $p$  is an irreducible.*

*Proof. Part A :  $\langle p \rangle$  maximal in  $D \implies p$  irreducible in  $D$*

Let  $D$  be a PID. Suppose  $\langle p \rangle$  is a maximal ideal in  $D$ . Suppose  $p = ab$  in  $D$ . That is,  $p = ab$  where  $a, b \in D$ . Then,  $\langle p \rangle \subset \langle a \rangle$ .

Suppose  $\langle p \rangle = \langle a \rangle$ . Then  $a = pc = (ab)c \implies bc = 1$  and  $b$  is a unit.

Suppose  $\langle p \rangle \neq \langle a \rangle$ . Then  $\langle a \rangle = D$ , since  $\langle p \rangle$  is a maximal ideal. That is  $\langle a \rangle = \langle 1 \rangle \implies a, 1$  are associates. And  $a$  is a unit.

Thus for any factorisation  $p = ab$ , either  $a$  or  $b$  is a unit. Therefore,  $p$  is an irreducible.

**Part B :  $p$  is irreducible in  $D \implies \langle p \rangle$  is maximal in  $D$**

Suppose  $p$  is an irreducible in  $D$ . Then  $p = ab$  where  $a, b \in D$  implies that either  $a$  or  $b$  is a unit. Since  $p = ab$ ,  $\langle p \rangle \subset \langle a \rangle$ . We know that, if  $a$  is a unit then  $\langle a \rangle = \langle 1 \rangle = D$ .

Suppose  $a$  is not a unit. Then  $b$  is a unit. That is, there exists  $u \in D$  such that  $bu = 1$ . Then  $pu = abu = a$ . Thus,  $\langle a \rangle \subset \langle p \rangle \implies \langle p \rangle = \langle a \rangle$ .

Thus, for any factorisation  $p = ab$ , the ideals generated by  $a, b$  are either  $D$  or  $\langle p \rangle$  itself. In other words, there doesn't exist a proper ideal containing  $\langle p \rangle$ . Therefore,  $\langle p \rangle$  is a maximal ideal in  $D$ .  $\square$

### 6.5.3 Every PID is UFD

We have seen earlier that  $F[x]$  has unique factorisation of all its elements. Now, we generalise that result into "every PID has unique factorisation for all its elements".

**Lemma 6.5.5.** *In a PID, if an irreducible  $p$  divides  $ab$ , then either  $p|a$  or  $p|b$ .*

*Proof.* Let  $D$  be a PID. Suppose  $p \in D$  is irreducible. Suppose  $p|ab$ . Then  $ab \in \langle p \rangle$ . And we have,  $\langle p \rangle$  is a maximal ideal. However, every maximal ideal in  $D$  is a prime ideal. Thus,  $ab \in \langle p \rangle \implies a \in \langle p \rangle$  or  $b \in \langle p \rangle$ . In other words,  $p|a$  or  $p|b$ .  $\square$

**Corollary 6.5.5.1.** *If  $p$  is an irreducible in a PID  $D$  and  $p$  divides  $a_1 a_2 \dots a_n$  where  $a_i \in D$ , then  $p|a_i$  for at least one  $i$ .*

*Proof.* We know that,  $p|a_1 a_2 \implies p|a_1$  or  $p|a_2$ .

Suppose  $p|a_1 a_2 \dots a_k \implies p|a_1$  or  $p|a_2 \dots$  or  $p|a_k$ .

Suppose  $p|a_1 a_2 \dots a_{k+1} \implies p|a_1 a_2 \dots a_k$  or  $p|a_{k+1}$ .

$\implies p|a_1$  or  $p|a_2$  or  $\dots$  or  $p|a_k$  or  $p|a_{k+1}$ .

That is, if  $p$  divides a finite product, then it divides at least one of them.  $\square$

**Definitions 6.5.5.** A nonzero, nonunit element  $p$  in an integral domain  $D$  is a prime if  $\forall a, b \in D, p|ab \implies p|a$  or  $p|b$ .

**Theorem 6.5.6.** Every PID is a UFD.

*Proof.* Let  $D$  be a PID. Let  $a \in D$  is neither a zero nor a unit. Then  $a$  has a factorisation into irreducibles  $a = p_1 p_2 \dots p_r$ . Suppose  $a$  has another factorisation  $a = q_1 q_2 \dots q_s$ .

We have  $p_1 | q_1 q_2 \dots q_s$ . Thus,  $p_1$  divides at least one of them, say  $q_j$ . Without loss of generality, we assume that  $p_1 | q_1$  (rearranging the terms if necessary).

Since  $q_1$  is an irreducible,  $q_1 = p_1 u_1 \implies u_1$  is a unit. Thus,  $p_1 p_2 \dots p_r = p_1 u_1 q_2 q_3 \dots q_s$ . That is,  $p_2 p_3 \dots p_r = u_1 q_2 q_3 \dots q_s$ .

Continuing like this we get,  $p_r = u_1 u_2 \dots u_{r-1} p_r u_r q_{r+1} q_{r+2} \dots q_s$ . Thus  $u_1 u_2 \dots u_r q_{r+1} q_{r+2} \dots q_s$  is a unit. We have  $s = r$ , otherwise  $1 = u_1 u_2 \dots q_{r+1} \dots$  is a contradiction. Therefore, any element  $a \in D$  has a unique factorisation except for order, associates and units.  $\square$

**Corollary 6.5.6.1.** The integral domain  $\mathbb{Z}$  is a UFD.

*Proof.* The ideals in  $\mathbb{Z}$  are of the form  $n\mathbb{Z}$  where  $n \in \mathbb{Z}$ . Thus every ideal in  $\mathbb{Z}$  are principal ideals. Therefore,  $\mathbb{Z}$  is a PID. We know that every PID is a UFD. Thus,  $\mathbb{Z}$  is a **unique factorisation domain**.  $\square$

#### 6.5.4 $D$ is UFD $\implies D[x]$ is UFD

**Definitions 6.5.6** (gcd). Let  $D$  be a UFD and  $a_1, a_2, \dots, a_n$  be nonzero elements in  $D$ . An element  $d \in D$  is a **greatest common divisor** of  $a_i$  if  $d$  is a common divisor of all  $a_i$ s and also divides any common divisor of  $a_i$ s.

**Definitions 6.5.7** (primitive). Let  $D$  be a UFD. A nonconstant polynomial  $a_0 + a_1 x + \dots + a_n x^n \in D[x]$  is a primitive if 1 is the gcd of all  $a_i$ s.

For example,  $4x^3 + 3x^2 + 2 \in \mathbb{Z}[x]$  is a primitive since  $\gcd(4, 3, 2) = 1$ .

**Lemma 6.5.7.** If  $D$  is a UFD, then for every nonconstant  $f(x) \in D[x]$  we have  $f(x) = cg(x)$  where  $c \in D$  and  $g(x) \in D[x]$  and  $g(x)$  is primitive. The element  $c \in D$  is unique upto a unit factor in  $D$  and is the **content** of  $f(x)$ . Also  $g(x)$  is unique upto a unit factor in  $D$ .

*Proof.* **Step 1 : There exists  $c \in D$  such that  $f(x) = cg(x)$**

Let  $f(x) = a_0 + a_1 x + \dots + a_n x^n$  be a non-constant polynomial in  $D[x]$ . Let  $c = \gcd(a_0, a_1, \dots, a_n)$ . In other words,  $c$  is the greatest common divisor of the coefficients of  $f(x)$ . Let  $q_i = ca_i$ . Clearly,  $\gcd(q_1, q_2, \dots, q_n) = 1$ . Thus, we have  $f(x) = cg(x)$  where  $g(x) = q_0 + q_1 x + \dots + q_n x^n$  is a primitive.

**Step 2 : There exists unique  $c$  such that  $f(x) = cg(x)$**

Suppose  $f(x) = cg(x) = dh(x)$  where  $c, d \in D$  and  $g(x), h(x)$  are primitives.

Clearly,  $c$  and  $d$  are associates. Otherwise, there exists an irreducible element  $p$  such that  $p|cg(x)$  but  $p \nmid dh(x)$  which is a contradiction. Thus, we can

cancel all irreducible factors of  $c$  to obtain,  $ug(x) = vh(x)$  where  $u$  and  $v$  are units. Therefore, the content of  $f(x)$ ,  $c$  unique upto associates.

Again,  $f(x) = cg(x) = (cu)(u^{-1}g(x))$  and thus  $g(x)$  is unique upto associates/units.  $\square$

**Definitions 6.5.8** (content). Let  $D$  be a UFD. Let  $f(x) \in D[x]$ . Then  $f(x) = cg(x)$  for some primitive  $g(x) \in D[x]$  and  $c \in D$ . The element  $c$  which is unique upto units, is the content of  $f(x)$ .

For example, Let  $f(x) = 8x^3 + 6x^2 + 4 \in \mathbb{Z}[x]$ . Then  $f(x) = 2(4x^3 + 3x^2 + 2)$  where  $4x^3 + 3x^2 + 2 \in \mathbb{Z}[x]$  is a primitive and  $2 \in \mathbb{Z}$ . Thus,  $\text{content}(f) = 2$ .

**Lemma 6.5.8** (Gauss). If  $D$  is a UFD, then product of two primitive polynomials in  $D[x]$  is again primitive.

*Proof.* Let  $f(x) = a_0 + a_1x + \dots + a_nx^n$  and  $g(x) = b_0 + b_1x + \dots + b_mx^m$ . Suppose  $f(x)$  and  $g(x)$  are primitives. Let  $h(x) = f(x)g(x) = c_0 + c_1x + \dots + c_{n+m}x^{n+m}$ . Let  $p$  be an irreducible in  $D$ . Then  $f(x)$  has a coefficient which is not divisible by  $p$ . Suppose  $p$  divides every coefficient  $f(x)$ , then  $f(x)$  is not a primitive.

Let  $r$  be the smallest integer such that  $p \nmid a_r$ . Similarly, let  $s$  be the smallest integer such that  $p \nmid b_s$ . Now consider the coefficient  $c_{r+s}$  of  $h(x)$ .

$$c_{r+s} = a_0b_{r+s} + \dots + a_{r-1}b_{s+1} + a_rb_s + a_{r+1}b_{s-1} + \dots + a_{r+s}b_0$$

By the selection of  $r$ ,  $p|a_j$  for every  $j < r$ . Thus,  $p|a_0b_s + \dots + a_{r-1}b_{s+1}$ . Similarly, by the selection of  $s$ ,  $p|b_k$  for every  $k < s$ . Thus,  $p|a_{r+1}b_{s-1} + \dots + a_{r+s}b_0$ . Clearly,  $p|c_{r+s} \iff p|a_rb_s \implies p|a_r$  or  $p|b_s$ . This is not possible, thus  $p \nmid c_{r+s}$ . Thus, for any irreducible  $p \in D$ ,  $h(x)$  has a coefficient which is not divisible by  $p$ . Therefore, product of two primitives is always a primitive.  $\square$

**Corollary 6.5.8.1.** If  $D$  is a UFD, then a finite product of primitive polynomials in  $D[x]$  is again primitive.

*Proof.* Let  $\{f_k\}_{k=1}^n$  be a family of primitives. Then  $f_1(x)f_2(x)$  is a primitive, since product of two primitives is always a primitive.

Suppose  $g(x) = f_1(x)f_2(x) \dots f_k(x)$  is a primitive.

Consider  $h(x) = f_1(x)f_2(x) \dots f_{k+1}(x) = g(x)f_{k+1}(x)$ . Since both  $g(x)$  and  $f_{k+1}(x)$  are primitives their product  $h(x)$  is also a primitive.

Thus by finite mathematical induction, every finite product of primitives is a primitive.  $\square$

**Lemma 6.5.9.** Let  $D$  be a UFD and let  $F$  be a field of quotients of  $D$ . Let  $f(x) \in D[x]$  where  $\deg f(x) > 0$ . If  $f(x)$  is an irreducible in  $D[x]$ , then  $f(x)$  is also an irreducible in  $F[x]$ . Also, if  $f(x)$  is primitive in  $D[x]$  and irreducible in  $F[x]$ , then  $f(x)$  is irreducible in  $D[x]$ .

*Proof.* Let  $D$  be a UFD and  $F$  be the field of quotients of  $D$ . Then elements of  $F$  are of the form  $(a, b) \in D \times D$ . Also  $(a, b) + (c, d) = (ad + bc, bd)$  and  $(a, b) \cdot (c, d) = (ac, bd)$ . And we say express  $(a, b)$  as  $a/b$ .

**Part 1 : irreducible in  $D[x] \implies$  irreducible in  $F[x]$**

Let  $f(x) \in D[x]$  be an irreducible in  $D[x]$ . Suppose  $f(x)$  is not irreducible in  $F[x]$ . Then  $f(x) = r(x)s(x)$  where  $r(x), s(x) \in F[x]$ .

The coefficients of  $r(x)$  are of the form  $a/b$ . By multiplying  $r(x)$  with the product of all denominators, we can remove its denominators and obtain an associate  $r_1(x) \in D[x]$ . Similarly, we have  $s(x) = vs_1(x)$  where  $s_1(x) \in D[x]$ . Thus,  $f(x) = r(x)s(x) = ur_1(x)vs_1(x) = uvr_1(x)s_1(x)$ . Therefore,  $f(x) \in D[x]$  is not an irreducible, which is a contradiction. Thus,  $f(x)$  is irreducible in the quotient field  $F$  as well.

**Part 2 : irreducible in  $F[x] \implies$  irreducible in  $D[x]$**

Suppose  $f(x) \in D[x]$  is irreducible in  $F[x]$ . Clearly,  $F$  contains a UFD isomorphic to  $D$ . Thus  $D[x] \leq F[x]$ . Therefore, any polynomial irreducible in  $F[x]$  is also irreducible in  $D[x]$ .  $\square$

We just saw that if  $F$  is the quotient field of a UFD  $D$  and  $f(x) \in D[x]$ . The irreducibility of  $f(x)$  in  $F[x]$  is necessary and sufficient for the irreducibility of  $f(x)$  in  $D[x]$ . Now we prove that the factorisation is also unique upto a morphism.

**Corollary 6.5.9.1.** *If  $D$  is a UFD and  $F$  is a field of quotients of  $D$ , then a nonconstant polynomial  $f(x) \in D[x]$  factors into a product of two polynomials of lower degrees  $r$  and  $s$  in  $F[x]$  if and only if it has a factorization into polynomials of the same degree  $r$  and  $s$ .*

*Proof.* We know that if  $f(x) \in D[x]$  has a factorisation  $f(x) = r(x)s(x)$  in  $F[x]$ . Then it has a factorisation in  $D[x]$  as well,  $f(x) = cr_1(x)s_1(x)$  where  $r_1(x), s_1(x)$  are associates of  $r(x), s(x)$  in  $F[x]$ . Clearly, the associates are always of the same degree.

Trivially, any factorisation in  $D[x]$  will also be a factorisation in  $F[x]$ . And the factors will have the same degree, unless  $D[x]$  has an irreducible polynomial which is not irreducible in  $F[x]$ . This is not possible.  $\square$

**Theorem 6.5.10.** *If  $D$  is a UFD, then  $D[x]$  is a UFD.*

*Proof.* Let  $f(x) \in D[x]$ . Suppose  $\deg(f) > 0$ . Otherwise,  $f(x)$  is a constant and we have trivial factorisation.

Let  $f(x) = g_1(x)g_2(x)\dots g_r(x)$  be a factorisation in  $D[x]$  with maximum number of factors. Then  $f(x) = c_1h_1(x)c_2h_2(x)\dots c_rh_r(x)$  where  $g_k(x) = c_kh_k(x)$  and  $h_k(x)$  are primitives. And  $h_k(x)$  are irreducibles. Otherwise,  $f(x)$  has a factorisation with greater number of factors than  $r$  which is a contradiction. Thus,  $f(x) = uh_1(x)h_2(x)\dots h_r(x)$  is a factorisation of  $f(x)$  into a product of irreducibles.

Suppose  $f(x)$  has another factorisation  $f(x) = G_1(x)G_2(x)\dots G_r(x)$ . Then  $f(x) = vH_1(x)H_2(x)\dots H_r(x)$  with  $H_k(x)$  are all irreducibles. We know that,  $h_k(x)|f(x) \implies h_k(x)|H_1(x)H_2(x)\dots H_r(x) \implies h_k(x)|H_k(x)$  (WLOG). Therefore, the factorisation is unique.  $\square$

For example, let  $x, y$  be two indeterminates, then an element  $f(x) \in F[x, y]$  is of the form  $\sum_{i,j=0}^{n,m} a_{ij}x^i y^j$ .

**Corollary 6.5.10.1.** *If  $F$  is a field of quotients and  $x_1, x_2, \dots, x_n$  are indeterminates, then  $F[x_1, x_2, \dots, x_n]$  is a UFD.*

*Proof.* Suppose  $D$  is a UFD and  $F$  is its field of quotients, then  $F[x]$  is a UFD.

Suppose  $K = F[x_1, x_2, \dots, x_k]$  is a UFD.

Then,  $F[x_1, x_2, \dots, x_k, x_{k+1}] = F[x_1, x_2, \dots, x_k][x_{k+1}] = K[x_{k+1}] = K[x]$  is a UFD.  $\square$

### 6.5.5 Exercise §45

1.

## 6.6 Euclidean Domain §46

**Definitions 6.6.1** (Euclidean Norm). Euclidean norm on an integral domain  $D$  is a function  $v$  which maps non-zero elements of  $D$  into non-negative integers such that

1.  $\forall a, b \in D, (b \neq 0)$  there exists  $q, r \in D$  such that  $a = qb + r$  where either  $r = 0$  or  $v(r) < v(b)$  and
2.  $\forall a, b \in D, (a, b \neq 0) v(a) \leq v(ab)$

**Definitions 6.6.2** (Euclidean Domain). An integral domain  $D$  is a Euclidean Domain if there exists a Euclidean Norm in  $D$ .

For example,  $v(n) = |n|$  is a Euclidean norm on Euclidean domain  $\mathbb{Z}$ . And  $v(f(x)) = \deg(f(x))$  is a Euclidean norm on Euclidean domain  $F[x]$ .

**Theorem 6.6.1.** *Every Euclidean domain is a PID.*

*Proof.* Let  $D$  be a Euclidean domain with Euclidean norm  $v$ . Let  $N$  be an ideal in  $D$ . If  $N = \{0\}$ , then  $N = \langle 0 \rangle$  is a principal ideal.

Suppose  $N \neq \{0\}$ . Then there exists  $b \in N$  ( $b \neq 0$ ). Clearly,  $b \in N \implies \langle b \rangle \subset N$ . Choose a  $b \in N$  such that  $v(b)$  is minimal in  $N$ .

**Claim :**  $N = \langle b \rangle$

Let  $a \in N$ . Then there exists  $q, r \in D$  such that  $a = qb + r$  where  $r = 0$  or  $v(r) < v(b)$ . Clearly,  $r = a - qb \in N$ . Thus,  $r = 0$  since  $v(r) < v(b)$  and  $v(b)$  is minimal in  $N$ . That is,  $a = qb$  or  $b|a$  for every  $a \in N$ . Thus,  $N \subset \langle b \rangle$ . Therefore,  $N = \langle b \rangle$ . In other words, every ideal in  $D$  is a principal ideal.  $\square$

**Corollary 6.6.1.1.** *Every Euclidean domain is a UFD.*

*Proof.* We know that, every Euclidean domain is a PID. And every PID is a UFD. Thus, every Euclidean domain is a UFD.  $\square$

### 6.6.1 Arithmetic in Euclidean Domain

**Theorem 6.6.2.** *For a Euclidean domain with Euclidean norm  $v$ ,  $v(1)$  is minimal and  $u \in D$  is a unit if and only if  $v(u) = v(1)$ .*

*Proof.* We have,  $v(a) \leq v(ab)$ . Take  $a = 1$ ,  $v(1) \leq v(1b) \implies v(1) \leq v(b)$  for every  $b \in D$ . Thus,  $v(1)$  is minimal.

Suppose  $u$  is a unit with multiplicative inverse  $u^{-1}$ . Clearly,  $v(1) \leq v(u)$  since  $u \in D$ . Take  $a = u$  and  $b = u^{-1}$ . Then  $v(u) \leq v(uu^{-1}) = v(1)$ . Therefore, for every unit  $u \in D$ , we have  $v(u) = v(1)$ .  $\square$

**Theorem 6.6.3** (division algorithm). *Let  $D$  be a Euclidean domain with Euclidean norm  $v$ . And  $a, b \in D$  ( $a, b \neq 0$ ).*

$$a = bq_1 + r_1, \text{ where } r_1 = 0 \text{ or } v(r_1) < v(b)$$

if  $r_1 \neq 0$

$$b = r_1q + r_2, \text{ where } r_2 = 0 \text{ or } v(r_2) < v(r_1)$$

In general, if  $r_i \neq 0$

$$r_{i-1} = r_iq + r_{i+1} \text{ where } r_{i+1} = 0 \text{ or } v(r_{i+1}) < v(r_i)$$

Then the sequence  $r_1, r_2, \dots$  must terminate with some  $r_s = 0$ .

If  $r_1 = 0$ , then  $\gcd(a, b) = b$ .

If  $r_k = 0$ , then  $\gcd(a, b) = r_{k-1}$ .

Furthermore, if  $d$  is a gcd of  $a$  and  $b$ , then there exists  $\lambda, \mu \in D$  such that  $d = \lambda a + \mu b$ .

**Proof. Step 1 : Finite Sequence**  $r_1, r_2, \dots, r_s$

Suppose  $r_1, r_2, \dots, r_{i-1} \neq 0$ . Then  $v(r_1) > v(r_2) > \dots > v(r_{i-1})$  is a strictly decreasing sequence of non-negative integers. Thus, the sequence will terminate in finite numbers of steps. That is,  $r_s = 0$  for some integer  $s$ .

**Step 2 :**  $\gcd(a, b) = r_{s-1}$

Suppose  $r_1 = 0$ . Then  $a = bq$  and  $\gcd(a, b) = b$  since  $b|a, b|b$  and there does not exist an integer greater than  $b$  that divides  $b$ .

Suppose  $r_1 \neq 0$  and  $\gcd(a, b) = d$ . Then  $r_1 = a - bq$  and  $d|r_1$ .

And  $d|b, d|r_1 \implies d|a$  since  $a = qb + r_1$ . Thus,  $\gcd(a, b) = \gcd(b, r_1)$ .

Continuing like this, we get  $\gcd(a, b) = \gcd(r_{s-2}, r_{s-1})$  if  $r_{s-1} \neq 0$ .

And we have,  $r_s = 0 \implies r_{s-2} = qr_{s-1}$ .

Clearly,  $\gcd(r_{s-2}, r_{s-1}) = r_{s-1} = \gcd(a, b)$ .

**Step 3 :**  $\gcd(a, b) = \lambda a + \mu b$

Let  $b$  be a  $\gcd(a, b)$ . Then,  $b = 0a + 1b$  and result is complete.

Suppose  $\gcd(a, b) = r_{s-1}$ . And  $r_{s-1} = r_{s-3} - qr_{s-2}$ . Clearly,  $r_i = \lambda_i r_{i-2} + \mu_i r_{i-1}$ . Keep on substituting using these equations until we reach  $r_{s-1} = \lambda a + \mu b$ .

Suppose  $d'$  is another  $\gcd(a, b)$ . Then  $d' = ud$  is an associate of  $d$ . And  $d' = ud = (\lambda u)a + (\mu u)b$ . And the result is true for any  $\gcd(a, b)$ .  $\square$

## 6.7 Gaussian Integers and Multiplicative Norms §47

### 6.7.1 Gaussian Integers and Norm

**Definitions 6.7.1** (Gaussian integers). Gaussian integers are complex numbers of the form  $a + ib$  where  $a, b \in \mathbb{Z}$ .

Note : Let  $\alpha = a + ib$ . Clearly, Gaussian integers are elements of  $\mathbb{Z}[i]$ . Here,  $\mathbb{Z}[i]$  is a simple extension of the integral domain  $\mathbb{Z}$  with a zero of  $x^2 + 1$ .

**Definitions 6.7.2** (Gaussian Norm). Let  $\alpha = a + ib$  be a Gaussian integer. Then  $N(\alpha) = |\alpha|^2 = a^2 + b^2$  is a Euclidean norm on  $\mathbb{Z}[i]$ .

**Lemma 6.7.1.** *The Gaussian norm  $N : \mathbb{Z}[i] \rightarrow \mathbb{Z}$  defined by  $N(a + ib) = a^2 + b^2$  has the following properties*

1.  $N(\alpha) \geq 0, \forall \alpha \in \mathbb{Z}[i]$
2.  $N(\alpha) = 0 \iff \alpha = 0$
3.  $N(\alpha\beta) = N(\alpha)N(\beta)$

*Proof.* We have  $\alpha \in \mathbb{Z}[i]$ . Then  $\alpha = a + ib$  where  $a, b \in \mathbb{Z}$ . Clearly,  $N(a + ib) = a^2 + b^2 \geq 0$ .

$$N(\alpha) = 0 \iff a^2 + b^2 = 0 \iff a = 0, b = 0 \iff \alpha = 0 + i0 = 0.$$

$$\text{Also, } N(\alpha\beta) = |\alpha\beta|^2 = |\alpha|^2|\beta|^2 = N(\alpha)N(\beta). \quad \square$$

**Lemma 6.7.2.**  $\mathbb{Z}[i]$  is an integral domain.

*Proof.* Let  $\alpha, \beta \in \mathbb{Z}[i]$ .

Then,  $\alpha + \beta = (a + ib) + (c + id) = (a + c) + i(b + d) = (c + a) + i(d + b) = \beta + \alpha$ . And for any  $\alpha \in \mathbb{Z}[i]$ , we have  $\alpha \cdot 1 = (a + ib) \cdot (1 + i0) = (a + ib)$ . Thus, we have  $\mathbb{Z}[i]$  is a commutative ring with unity.

It remains to prove that  $\mathbb{Z}[i]$  has no zero divisors. Suppose  $\alpha\beta = 0$ .  $\alpha\beta = 0 \implies N(\alpha\beta) = 0 \implies N(\alpha)N(\beta) = 0$ . But,  $N(\alpha), N(\beta) \in \mathbb{Z}$  and  $\mathbb{Z}$  has no zero divisors. Thus,  $N(\alpha) = 0$  or  $N(\beta) = 0 \implies \alpha = 0$  or  $\beta = 0$ . Therefore  $\mathbb{Z}[i]$  has no zero divisors.  $\square$



**Theorem 6.7.3** (Gaussian integers is a Euclidean Domain). *The function  $v$  defined by  $v(\alpha) = N(\alpha)$  for every non-zero  $\alpha \in \mathbb{Z}[i]$  is a Euclidean norm.*

*Proof. Step 1 :*  $v(a) \leq v(ab)$

Let  $\beta = b_1 + ib_2$  and  $\beta \neq 0$ . Then  $N(\beta) = b_1^2 + b_2^2 \geq 1$ . Thus,  $N(\alpha) \leq N(\alpha)N(\beta) = N(\alpha\beta)$ .

**Step 2 : Division algorithm**

That is,  $\forall \alpha, \beta \in \mathbb{Z}[i], (\beta \neq 0) \exists \sigma, \rho \in \mathbb{Z}[i]$  such that  $\alpha = \beta\sigma + \rho$ , and  $\rho = 0$  or  $v(\rho) < v(\beta)$ .

Since  $\beta \neq 0$ ,  $\alpha/\beta$  exists. And we have,  $\frac{\alpha}{\beta} = \frac{a_1+ia_2}{b_1+ib_2} = \frac{(a_1+ia_2)(b_1-ib_2)}{(b_1+ib_2)(b_1-ib_2)} = \frac{a_1b_1+a_2b_2}{b_1^2+b_2^2} + i\frac{a_2b_1-a_1b_2}{b_1^2+b_2^2} = r + is$  where  $r, s \in \mathbb{Q}$ . Consider integers  $q_1, q_2$  that are nearest to the rational numbers  $r, s$ . Define  $\sigma = q_1 + iq_2 \in \mathbb{Z}[i]$  and  $\rho = \alpha - \beta\sigma$ .

Suppose  $\rho \neq 0$ . If  $\rho = 0$ , then the proof is complete. By definition of  $\sigma$ ,  $|q_1 - r| \leq 0.5$  and  $|q_2 - s| \leq 0.5$ . And we have,

$$N\left(\frac{\alpha}{\beta} - \sigma\right) = N((r - q_1) + i(s - q_2)) = |q_1 - r|^2 + |q_2 - s|^2 \leq 0.5$$

Since  $\alpha = \beta\sigma + \rho$ , we have  $\rho = \beta\left(\frac{\alpha}{\beta} - \sigma\right)$

$$N(\rho) = N\left(\beta\left(\frac{\alpha}{\beta} - \sigma\right)\right) = N(\beta)N\left(\frac{\alpha}{\beta} - \sigma\right) \leq 0.5N(\beta)$$

Thus,  $N(\rho) < N(\beta)$ . Therefore, Gaussian integers is a Euclidean domain.  $\square$

## 6.7.2 Multiplicative Norm

**Definitions 6.7.3** (multiplicative norm). Let  $D$  be an integral domain. A function  $N : D \rightarrow \mathbb{Z}$  is a multiplicative norm if

1.  $N(\alpha) = 0 \iff \alpha = 0$
2.  $N(\alpha\beta) = N(\alpha)N(\beta), \forall \alpha, \beta \in D$

Note : Gaussian norm is a multiplicative norm.

**Theorem 6.7.4.** *Let  $D$  be an integral domain with multiplicative norm  $N$ . Then  $N(1) = 1$  and  $|N(u)| = 1$  for any unit  $u$ . Suppose  $|N(\alpha)| = 1$  only for units in  $D$ . Then  $N(\alpha)$  is a prime in  $\mathbb{Z}$  implies  $\alpha$  is an irreducible in  $D$ .*

*Proof.* We know,  $N$  is a multiplicative norm. Thus,  $N(1) = N(1 * 1) = N(1)N(1) \implies N(1) = 1$ .

Suppose  $u$  is a unit in  $D$ . Then  $1 = N(1) = N(uu^{-1}) = N(u)N(u^{-1})$  Since 1 has only two factors  $\pm 1 \in \mathbb{Z}$ , we have  $|N(u)| = 1$ .

Suppose  $|N(\alpha)| = 1$  only if  $\alpha$  is a unit in  $D$ . Suppose  $\alpha \in D$  and  $|N(\alpha)| = p$ , where  $p$  is a prime in  $\mathbb{Z}$ . Suppose  $\alpha$  is not an irreducible. Then  $\alpha$  has a factorisation  $\alpha = \beta\gamma$  where  $\beta, \gamma \in D$  are non-units. Then  $p = N(\alpha) = N(\beta\gamma) = N(\beta)N(\gamma)$ . We have,  $\beta, \gamma$  are non-units and  $|N(\beta)| \neq 1$  and  $|N(\gamma)| \neq 1$ . This is not possible. Therefore,  $\alpha$  is an irreducible in  $D$ .  $\square$

### 6.7.3 Applications of multiplicative norm

**The integral domain  $\mathbb{Z}[\sqrt{-5}]$  is not a UFD. (§47.9)**

*Proof.* Elements of  $\mathbb{Z}[\sqrt{-5}]$  are of the form  $a + b\sqrt{-5}$  and  $N(a + b\sqrt{-5}) = a^2 + 5b^2$ . If  $N(a + b\sqrt{-5}) = a^2 + 5b^2 = 1$ , then,  $a = \pm 1$  and  $b = 0$ . Clearly,  $|N(\alpha)| = 1$  only if  $\alpha$  is a unit in  $\mathbb{Z}[\sqrt{-5}]$ .

The element  $3 \in \mathbb{Z}[\sqrt{-5}]$  is an irreducible. Suppose 3 is not an irreducible. That is, there is a factorisation  $3 = \beta\gamma$  where  $\beta, \gamma$  are non-units. Then  $9 = N(3) = N(\beta)N(\gamma)$  where  $\beta, \gamma$  are non-units. The integer factors of 9 are 1, 3, 9. Thus,  $N(\beta) = N(\gamma) = 3$ . Otherwise  $N(\beta) = 1$  or  $N(\gamma) = 1$ . But, there doesn't exist an element  $a + b\sqrt{-5}$  such that  $N(a + b\sqrt{-5}) = 3$ . Similarly,  $7 \in \mathbb{Z}[\sqrt{-5}]$  is an irreducible. We have,  $49 = N(7) = N(\beta)N(\gamma) \implies N(\beta) = N(\gamma) = 7$  is not possible since there are no elements of norm 7 in  $\mathbb{Z}[\sqrt{-5}]$ .

Suppose  $(1 + 2\sqrt{-5})$  is not an irreducible. Then,  $21 = N(1 + 2\sqrt{-5}) = N(\beta)N(\gamma)$ . However, we know that the integer factors of 21 are 1, 3, 7, 21. Without loss of generality, we have  $N(\beta) = 3$  which is not possible. Thus,  $1 + 2\sqrt{-5}$  is an irreducible.

However neither 3 nor 7 is an associate of  $1 + 2\sqrt{-5}$ . Since  $|N(1 + 2\sqrt{-5})| \neq N(3)$  and  $N(1 + 2\sqrt{-5}) \neq N(7)$ . Clearly, the factorisation of  $21 \in \mathbb{Z}[\sqrt{-5}]$  is not unique. Therefore,  $\mathbb{Z}[\sqrt{-5}]$  is not a UFD.  $\square$

**Fermat's  $p = a^2 + b^2$  theorem**

**Theorem 6.7.5** (Fermat). *Let  $p$  be an odd prime integer. Then  $p = a^2 + b^2$  where  $a, b \in \mathbb{Z}$  if and only if  $p \cong 1 \pmod{4}$*

*Proof.* Suppose there exists an odd prime  $p$  such that  $p = a^2 + b^2$ . Then  $a, b$  are of different parity. Without loss of generality,  $a$  is even and  $b$  is odd. Let  $a = 2r$  and  $b = 2s + 1$ . Then  $a^2 + b^2 = 4r^2 + 4s^2 + 1 \cong 1 \pmod{4}$ .

Suppose  $p \cong 1 \pmod{4}$ . Then  $\mathbb{Z}_p$  is a cyclic group of order  $p - 1$  which has a element  $n$  of order 4 since  $4|(p - 1)$ . Then  $n^2 = -1 \in \mathbb{Z}_p$ . And  $n^2 + 1 \cong 0 \pmod{p}$ . In other words,  $p|(n^2 + 1)$ .

We claim that,  $p$  is not an irreducible in  $\mathbb{Z}[i]$ . Suppose  $p$  is an irreducible then  $p|(n^2 + 1) \implies p|(n + i)$  or  $p|(n - i)$ . And  $p|(n + i) \implies n + i = p(a + bi) \implies n = pa$  and  $1 = pb$ . But,  $1 = pb$  is not possible. Similarly,  $p|(n - i) \implies n - i = p(a + bi) \implies -1 = pb$  is also not possible. Thus if  $p \cong 1 \pmod{4}$ , then  $p$  is not an irreducible in  $\mathbb{Z}[i]$ .

Since  $p$  is not an irreducible in  $\mathbb{Z}[i]$ ,  $p = (a + bi)(c + di)$  where  $a + bi, c + di$  are non-units in  $\mathbb{Z}[i]$ . Considering the Gaussian norm  $N$  on  $\mathbb{Z}[i]$  defined by  $N(a + bi) = a^2 + b^2$ . We have,  $p^2 = N(p) = N(a + bi)N(c + di) = (a^2 + b^2)(c^2 + d^2)$ . And  $N(a + bi) \neq 1$  since for Gaussian norm  $N(\alpha) = 1 \iff \alpha = 1$ . Similarly,  $N(c + di) \neq 1$ . Therefore, without loss of generality  $N(a + bi) = p$  since  $a^2 + b^2 > 0$ . In other words,  $p = a^2 + b^2$ .  $\square$

**6.7.4 Exercise §47**

1. We have 5 is an odd prime and  $5 \cong 1 \pmod{4}$ . Therefore 5 is not an irreducible Gaussian integer. Clearly,  $2^2 + 1^2 = 5 \implies (2+i)(2-i) = 5$ .

And  $N(2+i) = 5$  is a prime, thus  $2+i$  is an irreducible. Similarly,  $2-i$  is also an irreducible. Therefore  $5 = (2+i)(2-i)$  is factorisation using irreducibles Gaussian integers. (hint :  $2^2+1^2=5$ . But,  $1*1+2*2=5$ )

Note : We have  $5 = (2+i)(2-i) = (1+2i)(1-2i)$ . However 5 has a unique factorisation since  $(2+i)$  and  $(1-2i)$  are associates  $(2+i)(0-i) = 1-2i$  where  $-i \in \mathbb{Z}[i]$  is a unit with multiplicative inverse  $i$ .

2. We have 7 is an odd prime. However  $7 \cong 3 \pmod{4}$ . Therefore 7 is an irreducible Gaussian integer.
3. We have  $N(4+3i) = 25$ . Clearly,  $(4+3i) = \alpha\beta \implies N(\alpha) = N(\beta) = 5$ . Therefore,  $4+3i = (2-i)(1+2i)$  is a factorisation of  $4+3i$ . (Hint :  $2^2 + 1^2 = 5$ . But,  $1*2 + 1*2 = 4$ .)
4. We have  $N(6-7i) = 85 = 5*17$ . Clearly  $6-7i = (2+i)(1-4i)$ . (Hint :  $2^2 + 1^2 = 5$  and  $4^2 + 1^2 = 17$ . But  $1*2 + 1*4 = 6$ .)
5. We have  $6 = 2*3$  in  $\mathbb{Z}[\sqrt{-5}]$ . Then,  $N(\alpha\beta) = 4 \implies a^2 + 5b^2 = 2$  is not possible. And,  $N(\alpha\beta) = 9 \implies a^2 + 5b^2 = 3$  is also not possible. Therefore,  $6 = 2*3$  is a factorisation in  $\mathbb{Z}[\sqrt{-5}]$ .

However  $N(\alpha\beta) = 6 \implies a^2 + 5b^2 = 6 \implies |a| = |b| = 1$ . Clearly,  $6 = (1 - \sqrt{-5})(1 + \sqrt{-5})$  is another factorisation.

**6.8 Automorphism and Fields §48****6.9 The Isomorphism Extension Theorem §49****6.10 Splitting Fields §50****6.11 Separable Extensions §51****6.12 Galois Theory §53****6.13 Illustrations of Galois Theory §54****6.14 Cyclotomic Extensions §55**

# Subject 7

## ME010202 Advanced Topology

### 7.1 Module I

**Q :** Why these axioms are called separation axioms ?

$T_0, T_1, T_2$  axioms separates points from points.

$T_3, T_{3\frac{1}{2}}$  axioms separates points from closed subsets.

$T_4$  axiom separates closed subsets from closed subsets.



Figure 7.1: Separation Axioms

#### 7.1.1 Compactness and Separation Axioms

[Joshi, 1983, chapter 7 §2.1-10]

**Proposition 7.1.1.** *Let  $X$  be a  $T_2$  space,  $x \in X$  and  $F$  is a compact subset of  $X$  not containing  $x$ . Then there exist open subsets  $U, V$  such that  $x \in U, F \subset V, U \cap V = \phi$ .*

$T_2$  spaces separates points from compact sets.

*Proof.* Let  $x$  be a points in  $X$  and  $F$  be a compact subset not containing  $x$ .  
 $X$  is  $T_2 \implies \forall y \in F, \exists U_y, V_y \in \mathcal{T}, x \in U_y, y \in V_y, U_y \cap V_y = \phi$ .<sup>†1)</sup>  
 $\mathcal{C} = \{V_y : y \in F\}$  is an open cover of compact subset  $F$ .

<sup>†1</sup>**Q :** Do we need compactness for this result ?

**A :** Given  $x \in X$ . For each point  $y \in F$ , we have two open sets namely  $U_y$  and  $V_y$  that separates these two points  $x$  and  $y$  in  $X$ . But, the intersection  $\cap_{y \in F} \{U_y\}$  is not necessarily an

$\implies$  there exists a finite subcover,  $\mathcal{C}' = \{V_{y_i} : y_i \in F, i = 1, 2, \dots, n\}$ .

Define  $U = \cap_{i=1}^n \{U_{y_i}\}$  and  $V = \cup_{i=1}^n \{V_{y_i}\}$ .

$\implies x \in U, F \subset V, U \cap V = \phi$ . □

**Corollary 7.1.0.1.** *A compact subset in a  $T_2$  space is closed.*

*Proof.* Let  $x \in X - F$ .

By proposition,  $x \in U, F \subset V, U \cap V = \phi \implies x \in U \subset X - V \subset X - F$

$X - F$  is a nbd of each of its points. Thus  $X - F$  is open.  $(\star^2)$  □

**Corollary 7.1.0.2.** *Every map from a compact space into a  $T_2$  space is closed.*

*And its range is a quotient space of the domain.*

*Proof.* Suppose  $X$  is compact,  $Y$  is  $T_2$  and  $f : X \rightarrow Y$  is continuous.

Let  $C$  be a closed subset of  $X$ .

$\implies C$  is compact, since compact is weakly hereditary.  $\star^3)$

$\implies f(C)$  is compact, since compactness is preserved by continuous functions.

$(\star^4)$

By corollary,  $f(C)$  is compact,  $Y$  is  $T_2 \implies f(C)$  is a closed subset of  $Y$ .

Thus  $f : X \rightarrow Y$  is closed.  $\dagger^5)$

$f : X \rightarrow Y$  is closed  $\implies f : X \rightarrow f(X)$  is a quotient map, since every closed, surjective map is a quotient map.. □

**Corollary 7.1.0.3.** *A continuous bijection from a compact space onto a  $T_2$  space is a homeomorphism.*

*Proof.* Let  $G$  be an open subset of  $X$ .

Then  $X - G$  is closed.

By corollary,  $f$  is closed, since  $f : \text{compact} \rightarrow T_2$ .

$f$  is closed  $\implies f(X - G)$  is closed.

$f(X - G) = f(X) - f(G)$  since  $f$  is injective.

$\implies f(X - G) = Y - f(G)$  since  $f$  is surjective.

$\implies f(G)$  is open.

Thus  $f$  is an open map.

Every continuous bijective, open map is a homeomorphism. □

**Corollary 7.1.0.4.** *Every continuous, one-to-one function from a compact space into a  $T_2$  space is an embedding.*

*Proof.* Let  $f : X \rightarrow Y$  be continuous and injective.

$f : X \rightarrow f(X)$  is surjective.

$\implies f : X \rightarrow f(X)$  is continuous, surjective.

By corollary,  $f : X \rightarrow f(X)$  is a homeomorphism.

Thus  $f$  is an embedding of  $X$  onto  $f(X) \subset Y$ . □

open subset of  $X$ , since  $F$  is not necessarily a finite subset of  $X$  and an **arbitrary** intersection of open subsets is not necessarily an open subset. Therefore, we have restricted this family into a finite family. We use compactness for this part.

<sup>2</sup>Neighbourhood characterisation of open subsets : Let  $G \subset X$  be a nbd of each of its points. Then  $G$  is an open subset of  $X$ .

<sup>3</sup>Compactness is weakly hereditary : Suppose  $X$  is compact. Then every closed subset of  $X$  is compact.

<sup>4</sup>Compactness is preserved by continuous functions : Suppose  $G$  be a compact subset of a topological space  $X$ . And function  $f : X \rightarrow Y$  be a continuous function. Then  $f(G)$  is a compact subset of the topological space  $Y$ . [Joshi, 1983, 6.1.8]

<sup>5</sup>Closed function is a function which maps closed subsets into closed subsets.



Figure 7.2: Embedding compact space into hausdorff space

**Theorem 7.1.1.** *Every compact  $T_2$  space is a  $T_3$  space.*

*Proof.*  $T_2 \implies T_1$

It is enough to prove that compact,  $T_2$  space is regular.

Let  $x$  be a point in  $X$  and  $C$  be a closed subset not containing  $x$ .

Then  $C$  is compact, compactness is weakly hereditary.

By proposition,  $T_2$  space separates points from compact subsets.

Thus  $T_2 + \text{compact} \implies T_3$  ( $\dagger^6$ ) □

**Proposition 7.1.2.** *Let  $X$  be a regular space,  $C$  a closed subset of  $X$  and  $F$  a compact subset of  $X$ , such that  $C \cap F = \emptyset$ . Then there exist open subsets  $U, V$  such that  $C \subset U, F \subset V$  and  $U \cap V = \emptyset$ .*

Regular spaces separates closed subsets from compact subsets.

*Proof.* proof technique is same as  $T_2$  separates points from compact subsets.

Let  $C$  be a closed subset,  $F$  be a compact subset of  $X$  and  $C \cap F = \emptyset$ .

$X$  is regular  $\implies \forall y \in F, \exists U_y, V_y \in \mathcal{T}, C \subset U_y, y \in V_y, U_y \cap V_y = \emptyset$

Define  $\mathcal{C} = \{V_y : y \in F\}$  is cover of compact subset  $F$ .

$\implies \exists \mathcal{C}'$  such that  $\mathcal{C}' = \{V_{y_i} : i = 1, 2, \dots, n\}$  is a finite subcover.

Define  $U = \bigcap_{i=1}^n \{U_{y_i}\}$  and  $V = \bigcup_{i=1}^n \{V_{y_i}\}$ .

$\implies C \subset U, F \subset V, U \cap V = \emptyset$ . □

**Theorem 7.1.2.** *Every regular, Lindeloff space is normal.*

*Proof.* Let  $C, D$  be two disjoint, closed subsets of a regular, lindeloff space  $X$ .

$X$  is regular  $\implies \forall x \in C, \exists U_x, V_x$  such that  $x \in U_x, D \subset V_x, U_x \cap V_x = \emptyset$

$X$  is regular  $\implies \forall y \in D, \exists U_y, V_y$  such that  $C \subset U_y, y \in V_y, U_y \cap V_y = \emptyset$

Then  $\{U_x : x \in C\}$  and  $\{V_y : y \in D\}$  are open covers of  $C$  and  $D$  respectively.

$\dagger^7$ )

Let  $\{U_n : n = 1, 2, \dots\}$  and  $\{V_n : n = 1, 2, \dots\}$  be their countable subcovers.

Define  $G_n = U_n - \bigcup_{i=1}^n \overline{V_i}$  and  $H_n = V_n - \bigcup_{i=1}^n \overline{U_i}$ .

Define  $G = \bigcup_{i=1}^\infty G_n$  and  $H = \bigcup_{i=1}^\infty H_n$ .

Claim :  $C \subset G$  (similary  $D \subset H$ )

$x \in C \implies x \in U_n$  for some  $n$ , since  $\{U_n : n \in \mathbb{N}\}$  is a cover of  $C$

$\forall m, \overline{V_m} \subset X - C \implies \forall m, x \notin \overline{V_m} \implies x \in G$

Claim :  $G \cap H = \emptyset$

$x \in G \cap H \implies x \in G_m \cap H_n$  for some  $m, n$

With loss of generality,  $n \geq m, x \in G_m \implies x \in U_m \implies x \notin H_n$  □

<sup>6</sup> $T_3 \not\Rightarrow \text{compact}$ .

<sup>7</sup>Here,  $U_x, V_y$  are unrelated.

**Corollary 7.1.2.1.** *Every regular, second countable space is normal.*

*Proof.* Every second countable space is lindeloff.

By theorem, every regular, lindeloff space is normal.  $\square$

**Corollary 7.1.2.2.** *Every compact,  $T_2$  space is  $T_4$ .*

*Proof.* By proposition, every compact,  $T_2$  space is regular.

Every compact space is lindeloff since finite subcovers are countable.

By theorem, every regular, lindeloff space is normal.  $\square$

## 7.1.2 The Urysohn Characterisation of Normality

$X$  is normal  $\iff$  there exist Urysohn functions

**Q :** Significance of Urysohn's Lemma ?

1. There is no analog of Urysohn's Lemma for Regular spaces.
2. Constructs a nice real-valued function even on non-metrisable spaces.

**Q :** We know that, the completely regular space separates points from closed subsets using a real-valued function. Which space separates closed subsets from closed subsets using a real-valued function ?

**A :** Normal space(by Urysohn's Lemma). Not completely normal( $\star^8$ )

**Proposition 7.1.3.** *Let  $A, B$  be subsets of a space  $X$  and suppose there exists a continuous function  $f : X \rightarrow [0, 1]$ , such that  $f(x) = 0, \forall x \in A$  and  $f(x) = 1, \forall x \in B$ . Then there exists disjoint open subsets  $U, V$  such that  $A \subset U$  and  $B \subset V$ .*

*Proof.* Let  $f$  be a continuous function,  $f(x) = 0, \forall x \in A$  and  $f(x) = 1, \forall x \in B$ . Define  $G = [0, \frac{1}{2})$ ,  $H = (\frac{1}{2}, 1]$  and  $U = f^{-1}(G)$ ,  $V = f^{-1}(H)$ .

$\implies A \subset U, B \subset V$  and  $U \cap V = \phi$  since  $G \cap H = \phi$ .  $\square$

**Corollary 7.1.2.3** (Urysohn's Lemma : Sufficient Condition). *If  $X$  has the property that for any disjoint closed subsets  $A, B$  of  $X$ , there exists a continuous function  $f : X \rightarrow [0, 1]$  such that  $f(x) = 0, \forall x \in A$  and  $f(x) = 1, \forall x \in B$ , then  $X$  is normal.*

There exists a Urysohn function  $f \implies X$  is normal

*Proof.* Let  $A, B$  be two closed subsets of  $X$ .

By proposition,  $X$  normal.  $\square$

**Theorem 7.1.3** (Urysohn's Lemma). *A topological space  $X$  is normal iff it has the property that for every mutually disjoint, closed subsets  $A, B$  of  $X$ , there exists a continuous function  $f : X \rightarrow [0, 1]$  such that  $f(x) = 0$  for all  $x \in A$  and  $f(x) = 1$  for all  $x \in B$*

---

<sup>8</sup>Completely Normal Space separates any two subsets. [Joshi, 1983, chapter 7 Exer. 1.11]  
That is these subsets need not be closed. Thus every completely normal space is normal, but not the other way. Thus  $T_1 + \text{Completely Normal} = T_5 \supset T_4$ .

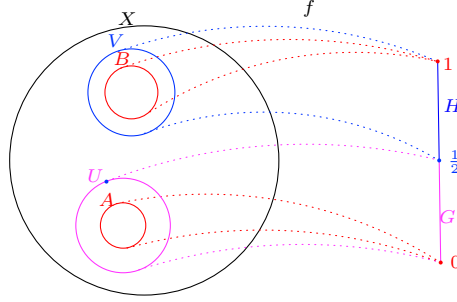


Figure 7.3: Urysohn's Lemma : Sufficient Condition

**Proof. Urysohn's Lemma : necessary condition**

$X$  is normal  $\implies$  there exist Urysohn functions

Let  $A, B$  be two closed subsets in a normal space  $X$ .

Enumerate rationals in the unit interval.

$\mathbb{Q} \cap [0, 1] = \{0, 1, \dots\} = \{q_0, q_1, q_2, \dots\}$ .

Define  $F_1 = F_{q_1} = X - B$ .

$A \subset X - B \implies \exists H$  such that  $A \subset H \subset \overline{H} \subset X - B$ . ( $\star^9$ )

Define  $F_{q_0} = F_0 = H$ .

$\overline{F_{q_0}} \subset F_{q_1} \implies$  true for  $n = 1$ .

Suppose :  $\overline{F_{q_i}} \subset F_{q_j}$ ,  $\forall j > i$  is true for  $j = 1, 2, \dots, n-1$ .

Define  $q_i = \sup\{q_k : q_k < q_n, k < n\}$  and  $q_j = \inf\{q_k : q_k > q_n, k < n\}$  ( $\dagger^{10}$ )

$\implies q_i < q_n < q_j$  and  $\overline{F_{q_i}} \subset F_{q_j}$ .

$\overline{F_{q_i}} \subset F_{q_j} \implies \exists H$  such that  $\overline{F_{q_i}} \subset H \subset \overline{H} \subset F_{q_j}$ . Define  $F_{q_n} = H$ .

Therefore,  $\overline{F_{q_j}} \subset F_{q_i}$ ,  $\forall j > i$ .

Now,  $\{F_t : t \in \mathbb{Q}\}$  has the properties required in lemma 2.

By lemma 2, there exists a function  $f$ . This is a Urysohn function. ( $\dagger^{11}$ )  $\square$

**Lemma 7.1.4.** Let  $f : X \rightarrow [0, 1]$  be continuous. For each  $t \in \mathbb{R}$  let  $F_t = \{x \in X : f(x) < t\}$ . Then the indexed family  $\{F_t : t \in \mathbb{R}\}$  has the following properties

1.  $F_t$  is an open subset of  $X$  for each  $t \in \mathbb{R}$
2.  $F_t = \emptyset$  for  $t < 0$
3.  $F_t = X$  for  $t > 1$
4. For any  $s, t \in \mathbb{R}$ ,  $s < t \implies \overline{F_s} \subset F_t$ .

Moreover, for each  $x \in X$ ,  $f(x) = \inf\{t \in \mathbb{Q} : x \in F_t\}$ .

*Proof.* Not needed.  $\square$

<sup>9</sup>Equivalent condition for normality : Let  $A$  be closed subset and  $G$  be an open subset containing  $A$ . Then there exists open subset  $H$  such that  $A \subset H \subset \overline{H} \subset G$ . cite[7.1.16(3)]joshi

<sup>10</sup>Let  $\mathbb{Q} \cap [0, 1] = \{0, 1, 0.3, 0.7, 0.8, \mathbf{0.5}, 0.6, \dots\}$ . Consider  $n = 5$ . Then  $q_n = q_5 = 0.5$  and in the respective induction step,  $q_i = \sup\{0, 0.3\} = 0.3$  and  $q_j = \inf\{1, 0.7, 0.8\} = 0.7$

<sup>11</sup>You will have to replace "Urysohn function" with "There exists a continuous real-valued function,  $f : X \rightarrow [0, 1]$  such that  $f(x) = 0, \forall x \in A, f(x) = 1, \forall x \in B$ ".



**Lemma 7.1.5.** *Let  $X$  be a topological space and suppose  $\{F_t : t \in \mathbb{Q}\}$  is a family of sets in  $X$  such that*

1.  $F_t$  is open in  $X$  for each  $t \in \mathbb{Q}$
2.  $F_t = \emptyset$  for  $t \in \mathbb{Q}, t < 0$
3.  $F_t = X$  for  $t \in \mathbb{Q}, t > 1$
4.  $\overline{F_s} \subset F_t$  for  $s, t \in \mathbb{Q}, s < t$

For  $x \in X$ , let  $f(x) = \inf\{t \in \mathbb{Q} : x \in F_t\}$ . Then  $f$  is a continuous real-valued function on  $X$  and it takes values in the unit interval  $[0, 1]$ .

*Proof.* Function  $f$  is well-defined and  $\text{Range}(f) = [0, 1]$ .

Define  $H = \{x \in X : f(x) < s\} = f^{-1}(-\infty, s)$  and

$K = \{x \in X : f(x) > s\} = f^{-1}(s, \infty)$ .

$\mathcal{S} = \{(s, \infty), (-\infty, s) : s \in \mathbb{R}\}$  is a subbase for  $\mathbb{R}$  with usual topology.

**Claim :**  $H = \cup\{F_t : t \in \mathbb{Q}, t < s\}$

$x \in H \implies f(x) < s \implies \inf G_x < s \implies \exists q \in \mathbb{Q} (q < s), x \in F_q$

$\implies x \in \cup\{F_t : t \in \mathbb{Q}, t < s\} \implies H \subset \cup\{F_t : t \in \mathbb{Q}, t < s\}$ .

$x \in \cup\{F_t : t \in \mathbb{Q}, t < s\} \implies \exists t \in \mathbb{Q} (t < s), x \in F_t \implies \inf G_x < s$

$\implies f(x) < s \implies x \in H \implies \cup\{F_t : t \in \mathbb{Q}, t < s\} \subset H$ .

**Claim :**  $X - K = \cap\{\overline{F_t} : t \in \mathbb{Q}, t > s\}$ .

$x \in X - K \implies x \notin K \implies f(x) \leq s \implies \inf G_x \leq s$

Let  $t \in \mathbb{Q} (s < t) \implies \exists q \in G_x (s < q < t) \implies x \in \overline{F_q} \subset F_t \subset \overline{F_t}$ .

$\forall t \in \mathbb{Q} (t > s), x \in \overline{F_t} \implies x \in \cap\{\overline{F_t} : t \in \mathbb{Q}, t > s\}$ .

Thus,  $X - K \subset \cap\{\overline{F_t} : t \in \mathbb{Q}, t > s\}$ .

$x \in \cap\{\overline{F_t} : t \in \mathbb{Q}, t > s\}$

Suppose(1)  $x \in K \implies s < f(x) \implies \exists q, t \in \mathbb{Q} (s < q < t), t < f(x)$

Suppose(2)  $x \in \overline{F_q} \implies x \in \overline{F_q} \subset F_t \implies \inf G_x < t \implies f(x) < t$

is a contradiction(2). Thus  $x \notin \overline{F_q}$ .

$x \notin \overline{F_q} \implies x \notin \cap\{\overline{F_t} : t \in \mathbb{Q}, t > s\}$  is a contradiction (1). Thus  $x \notin K$

Thus,  $\cap\{\overline{F_t} : t \in \mathbb{Q}, s > t\} \subset X - K$ .

Inverse images of subbase elements are open. Thus  $f$  is continuous.  $\square$

**Corollary 7.1.5.1.** *All  $T_4$  spaces are completely regular and hence Tychonoff.*

*Proof.* Let  $x \in X$  and  $D$  be closed subset not containing  $x$ . We have  $X$  is a  $T_4$  space. Therefore  $X$  is  $T_1$  as well as normal.

Now the singleton set,  $\{x\}$  is closed, since  $X$  is a  $T_1$  space. And by **Urysohn's lemma** for disjoint, closed subsets  $\{x\}, D$  there exists a Urysohn function which is a continuous, real-valued function  $f : X \rightarrow [0, 1]$  such that  $f(x) = 0$  and  $f(y) = 1$  for all  $y \in D$ . Therefore the space  $X$  is completely regular and hence Tychonoff.  $\square$

*Remark* (Urysohn function). The function whose existence is asserted by Urysohn's lemma is called a Urysohn function

### 7.1.3 Tietze Characterisation of Normality

#### Tietze Characterisation of normality

$X$  is normal  $\iff$  there exist continuous extension of real-valued functions on closed subsets.

**Proposition 7.1.4.** *Let  $A$  be a subset of a space  $X$  and let  $f : A \rightarrow \mathbb{R}$  be continuous. Then any two extensions of  $f$  to  $X$  agree on  $\overline{A}$ . In other words, if at all an extension of  $f$  exists its values on  $\overline{A}$  are uniquely determined by values of  $f$  on  $A$ .*

**Proposition 7.1.5** (Tietze : necessary condition). *Suppose a topological space  $X$  has the property that for every closed subset  $A$  of  $X$ , every continuous real valued function on  $A$  has a continuous extension to  $X$ . Then  $X$  is normal.*

**Definitions 7.1.1** (Pointwise Convergence). Let  $X$  be a topological space and  $(Y, d)$  a metric space. Then a sequence of functions  $\{f_n\}$  from  $X$  to  $Y$  converges pointwise to  $f$  if for every  $x \in X$  the sequence  $\{f_n(x)\}$  converges to  $f(x)$  in  $Y$ .

In other words, given a very small value,  $\epsilon > 0$ , there exists some  $\delta > 0$  such that for every  $x \in X$  there exists  $N_x \in \mathbb{N}$ . This  $N_x$  may be different for different values of  $x$  and for every  $n > N_x$ ,  $d(f(x), f_n(x)) < \delta$ .

**Definitions 7.1.2** (Uniform Convergence). Let  $X$  be a topological space and  $(Y, d)$  a metric space. Then a sequence of functions  $\{f_n\}$  from  $X$  to  $Y$  converges uniformly to  $f$  if given a small  $\epsilon > 0$ , there exists  $\delta > 0$  such that there exists  $N \in \mathbb{N}$ . This  $N$  is independent of the value of  $x$  and for every  $n > N$ ,  $d(f(x), f_n(x)) < \delta$ .

**Proposition 7.1.6.** *Let  $X$ ,  $(Y, d)$ ,  $\{f_n\}$  and  $f$  be as above and suppose  $\{f_n\}$  converges to  $f$  uniformly. If each  $f_n$  is continuous, then  $f$  is continuous.*

**Definitions 7.1.3** (Uniform Convergence of Series). Let  $X$  be a topological space and  $(Y, d)$  be a metric space. Then a series of function  $\sum_{n=1}^{\infty} f_n$  converges uniformly to  $f$  if the sequence of partial sums converges uniformly to  $f$ .

In other words, let  $g_m = \sum_{n=1}^m f_n$ . Then  $\sum_{n=1}^{\infty} f_n$  converges to  $f$  uniformly if the sequence of partial sums  $\{g_m\}$  converges to  $f$  uniformly.

**Proposition 7.1.7.** *Let  $\sum_{n=1}^{\infty} M_n$  be a convergent series of non-negative real numbers. Suppose  $\{f_n\}$  is a sequence of real valued functions on a space  $X$  such that for each  $x \in X$  and  $n \in \mathbb{N}$ ,  $|f_n(x)| \leq M_n$ . Then the series  $\sum_{n=1}^{\infty} f_n$  converges uniformly to a real valued function on  $X$ .*

**Theorem 7.1.6.** *Let  $A$  be a closed subset of a normal space  $X$  and suppose function  $f : A \rightarrow [-1, 1]$  is a continuous function. Then there exists a continuous function  $F : X \rightarrow [-1, 1]$  such that  $F(x) = f(x)$  for all  $x \in A$ .*

**Theorem 7.1.7.** *Let  $A$  be a closed subset of a normal space  $X$  and suppose function  $f : A \rightarrow (-1, 1)$  is a continuous function. Then there exists a continuous function  $F : X \rightarrow (-1, 1)$  such that  $F(x) = f(x)$  for all  $x \in A$ .*

**Corollary 7.1.7.1** (Tietze : sufficient condition). *Any continuous real-valued function on a closed subset of a normal space can be extended continuously to the whole space.*

## 7.2 Products and Coproducts

### 7.2.1 Cartesian Products of Families of Sets

### 7.2.2 The Product Topology

### 7.2.3 Productive Properties

## **7.3 Embedding and Metrisation**

### **7.3.1 Evaluation Functions into Products**

### **7.3.2 Embedding Lemma and Tychonoff Embedding**

### **7.3.3 The Urysohn Metrisation Theorem**

## 7.4 Nets and Homotopy

### 7.4.1 Definition and Convergence of Nets

**Definitions 7.4.1** (Directed Set). [Joshi, 1983, 10.1.1]

A directed set  $D$  is a pair  $(D, \geq)$  where  $D$  is a nonempty set and  $\geq$  is a binary relation on  $D$  such that

1. The relation ‘follows’ ( $\geq$ ) is transitive. ie,  $m \geq n, n \geq p \implies m \geq p$
2. The relation ‘follows’ ( $\geq$ ) is reflexive. ie, For every  $m \in D, m \geq m$
3. For any  $m, n \in D$ , there exists  $p \in D$  such that  $p \geq m$  and  $p \geq n$ .

**sequence in a set**  $X$  is a function  $f$  from the set of all integers into  $X$ .

**Definitions 7.4.2** (Net). [Joshi, 1983, 10.1.2]

A net in a set  $X$  is a function  $S$  from a directed set  $D$  into the set  $X$ .

*Remark.* The set  $\mathbb{N}$  together with the relation ‘less than or equal to’ ( $\leq$ ) is a directed set. Clearly, the relation ‘less than or equal to’ is reflexive and transitive. And the third condition is true iff every finite subset  $E$  of  $D$  has an element  $p \in E$  such that  $p$  follows each element of  $E$ . This is a weaker notion compared to the well ordering principle<sup>12</sup> of the set of all integers. Thus  $\mathbb{N}$  is a directed set and every sequence in  $X$  is also a net in  $X$ .

*Remark* (Significance of Net). A net on a set is a generalisation of ‘a sequence on a set’ obtained by simplifying the domain of the sequence into a directed set. The notion directed set is derived by assuming a few properties of  $\mathbb{N}$ .

The convergence of sequence is not strong enough to characterise topologies as the limit of convergent sequences are unique for both Hausdorff and Co-countable spaces. The notion of Net allows us to differentiate between Hausdorff spaces from Co-countable spaces in terms of convergence of nets. The limit of a convergent net on a topological space is unique iff it is a Hausdorff space. ie, We have removed a few restrictions, so that we will have some convergent nets (which are obviously not sequences) with multiple limit points for Co-countable spaces.

*Remark.* Examples of Directed Sets

1. Let  $X$  be a topological space and  $x \in X$ . Then the neighbourhood system  $\mathcal{N}_x$  is a directed set with the binary relation  $\subset$  (subset/inclusion).
  - (a) Let  $U, V, W$  be any three neighbourhoods of  $x \in X$  such that  $U \subset V$  and  $V \subset W$ . Then, clearly  $U \subset W$ . Therefore,  $U \geq V, V \geq W \implies U \geq W$ .
  - (b) Let  $U$  be any neighbourhood of  $x \in X$ , then  $U \subset U$ . Therefore,  $U \geq U$ .
  - (c) Let  $U, V$  be any two neighbourhoods of  $x \in X$ , then there exists their intersection  $W = U \cap V$ , which is a neighbourhood of  $x$ . Clearly  $W \subset U$  and  $W \subset V$ . Therefore  $\forall U, V \in \mathcal{N}_x, \exists W \in \mathcal{N}_x$  such that  $W \geq U$  and  $W \geq V$ .

---

<sup>12</sup>Well-ordering principle : Every subset of  $\mathbb{N}$  has a least element in it.

2. Let  $\mathcal{P}$  be the set of all partitions on closed unit interval  $[0, 1]$ . A partition  $P \in \mathcal{P}$  is a refinement of  $Q \in \mathcal{P}$  if every subinterval in  $P$  is contained in some subinterval of  $Q$ . Then  $\mathcal{P}$  with the binary relation refinement is a directed set.
  - (a) Suppose  $P, Q, R$  are three partitions of  $[0, 1]$  such that  $P$  is a refinement of  $Q$  and  $Q$  is a refinement of  $R$ , then clearly  $P$  is a refinement of  $R$  since each subinterval of  $P$  is contained some subinterval of  $Q$ , which is contained in some subinterval of  $R$ .  
Therefore,  $P \geq Q, Q \geq R \implies P \geq R$
  - (b) Suppose  $P$  is a partition of  $[0, 1]$ . Then trivially,  $P$  is a refinement of itself since every subinterval of  $P$  is contained in the same subinterval of  $P$ .  
Therefore,  $\forall P \in \mathcal{P}, P \geq P$
  - (c) Suppose  $P, Q$  be any two partition of  $[0, 1]$ . Then  $R = P \cup Q$  is a refinement of both the partitions.  
Therefore  $\forall P, Q \in \mathcal{P}, \exists R \in \mathcal{P}$  such that  $R \geq P$  and  $R \geq Q$

For example, let  $P = \{0, 0.3, 0.7, 1\}$ . Then the subintervals in  $P$  are  $[0, 0.3]$ ,  $[0.3, 0.7]$  and  $[0.7, 1]$ . Let  $Q = \{0, 0.3, 0.5, 1\}$  and  $R = \{0, 0.3, 0.5, 0.7, 1\}$ . Then  $R$  is a refinement of  $P$ , but  $Q$  is not a refinement of  $P$  since there is a subinterval  $[0.5, 1]$  in  $Q$  which is not properly contained in any subinterval of  $P$ . However,  $R$  is a refinement of  $Q$  as well.

*Remark.* Examples of Nets

1. Let  $X$  be a topological space and  $x \in X$ . Let  $\mathcal{N}_x$  be the set of all neighbourhoods of  $x$ . Let  $D = (\mathcal{N}_x, X)$  be the directed set given by  $(N, y) \in (\mathcal{N}_x, X)$  if  $N \in \mathcal{N}_x$  and  $y \in N$  and  $(N, y) \geq (M, z)$  if  $N \subset M$ . Then the function  $S : (\mathcal{N}_x, X) \rightarrow X$  given by  $S(N, y) = y$  is a net on  $X$ .

For example, let  $X = \{a, b, c, d\}$  and  $\mathcal{T} = \{\{a\}, \{a, b\}, \{a, b, c\}, \{a, b, c, d\}\}$ . Also let  $S : (\mathcal{N}_b, X) \rightarrow X$  defined by  $S(N, y) = y$ . Suppose  $C = \{a, b, c\}$ . Then  $C \in \mathcal{N}_b$ . ie,  $C$  is a neighbourhood of  $b$ . Then  $S(C, c) = c$ .

2. Riemann Net - Let  $D = (\mathcal{P}, \xi)$  where  $\mathcal{P}$  is the set of all partitions on  $[0, 1]$  and  $\xi$  is a finite sequence in  $[0, 1]$  such that consecutive terms belongs to consecutive subintervals of the partition. The set  $(\mathcal{P}, \xi)$  is directed set with  $\geq$  given by  $(P, \eta) \geq (Q, \psi)$  iff  $P$  is a refinement of  $Q$ .

For example, let  $P \in \mathcal{P}$  is given by  $P = \{0, 0.3, 0.7, 1\}$  and  $\eta = \{0.2, 0.6, 0.9\}$ . Then  $(P, \eta) \in (\mathcal{P}, \xi)$ .

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be any function, then the function,

$$S : (\mathcal{P}, \xi) \rightarrow \mathbb{R} \text{ defined by } S(P, \eta) = \sum_{j=1}^k f(\eta_j)(a_k - a_{k-1})$$

where  $P = \{a_0, a_1, \dots, a_k\}$  is the Riemann Net with respect to the real function  $f$ .

For example, let  $f(x) = x^2$  and  $P, \eta$  are same as above example, then  
 $S(P, \eta) = 0.2^2(0.3 - 0) + 0.6^2(0.7 - 0.3) + 0.9(1 - 0.7) = 3.99$

**Definitions 7.4.3** (Convergence of a Net). [Joshi, 1983, 10.1.3]

A net  $S : D \rightarrow X$  converges to a point  $x \in X$  if for any nbd  $U$  of  $x$ , there exists  $m \in D$  such that  $n \in D, n \geq m \implies S(n) \in U$ . And  $x$  is a limit of the net  $S$ .

*Remark.* The choice of  $m$  depends on the choice of neighbourhood  $U$

$$S : D \rightarrow X, S \rightarrow x \iff (\forall U \in \mathcal{N}_x, \exists m_U \in D, \text{ such that } n \geq m_U \implies S(n) \in U)$$

**Theorem 7.4.1** (Net characterisation of Hausdorff space). [Joshi, 1983, 10.1.4]  
 A topological space is Hausdorff iff limits of all nets in it are unique.

*Proof.* Let  $X$  be a Hausdorff space. Suppose  $S : D \rightarrow X$  is net on  $X$  such that  $S$  converges to two distinct points  $x, y \in X$ . Since  $X$  is a Hausdorff space and  $x \neq y$ , there exists open subsets  $U, V$  such that  $x \in U, y \in V, U \cap V = \phi$ .

The net  $S$  converges to  $x \in X$ , therefore  $\exists m_x \in D$  such that  $n \geq m_x \implies S(n) \in U$ . And, the net  $S$  converges to  $y \in X$ , therefore  $\exists m_y \in D$  such that  $n \geq m_y \implies S(n) \in V$ .

Since  $D$  is a directed set and  $m_x, m_y \in D$ , there exists  $p \in D$  such that  $p \geq m_x$  and  $p \geq m_y$ . Now,  $n \geq p \implies n \geq m_x, n \geq m_y$ , since  $\geq$  is transitive. (ie,  $n \geq p, p \geq m_x \implies n \geq m_x$ , and  $n \geq p, p \geq m_y \implies n \geq m_y$ ).

We have  $n \geq p \implies n \geq m_x$  and  $n \geq m_x \implies S(n) \in U$ . Therefore,  $n \geq p \implies S(n) \in U$ . Similarly,  $n \geq p \implies n \geq m_y \implies S(n) \in V$ . Therefore  $S(n) \in U \cap V$  which is a contradiction, since  $U \cap V = \phi$ . Therefore, if a net  $S$  converges to two points  $x, y$ , then  $x = y$ . That is, if a net  $S$  in a Hausdorff space  $X$  is convergent then its limit is unique.

Conversely, suppose that  $X$  is a topological space and every convergent net in  $X$  has a unique limit. Suppose  $X$  is not a Hausdorff space. Then there exists at least two distinct points  $x, y \in X$  such that every neighbourhood of  $x$  intersects with every neighbourhood of  $y$ . Now consider the set  $D = \mathcal{N}_x \times \mathcal{N}_y$  and relation  $\geq$  on  $D$  such that  $(U_1, V_1) \geq (U_2, V_2)$  if  $U_1 \subset U_2$  and  $V_1 \subset V_2$ .

By the axiom of choice, a function  $S : D \rightarrow X$  such that  $S(U, V) \in U \cap V$  is well defined, since every nbd of  $x$  intersects every nbd of  $y$ . Thus,  $S$  is a net in  $X$ . We claim that  $S$  converges to both  $x$  and  $y$ .

Let  $U$  be a nbd of  $x$ . Then  $S(U', V') \in U' \cap V'$ . We have  $(U, X) \in D$  such that  $(U', V') \geq (U, X) \implies U' \subset U$ . Then,  $S(U', V') \in U' \cap V' \subset U \cap X = U$ . Thus, for any nbd  $U$  containing  $x$ , we have  $(U, X) \in D$  such that  $(U', V') \geq (U, X) \implies S(U', V') \in U$ . Therefore,  $S$  converges to  $x \in X$ .

Similarly, Let  $V$  be a nbd of  $y$ . Then for any nbd  $V$  containing  $y$ , we have  $(X, V) \in D$  such that  $(U', V') \geq (X, V) \implies S(U', V') \in V$ , since

$S(U', V') \in U' \cap V' \subset X \cap V = V$ . Therefore,  $S$  converges to  $y \in X$  as well, where  $x \neq y$ . This is a contradiction to the assumption that every convergent net in  $X$  has a unique limit. Therefore, for any two points  $x, y \in X$ , there should be some nbd of  $x$  that doesn't intersect some nbd of  $y$ . Therefore,  $X$  is a Hausdorff space.  $\square$

**Definitions 7.4.4** (Eventual Subset). [Joshi, 1983, 10.1.5]

A subset  $E$  of a directed set  $D$  is an eventual subset of  $D$  if there exists  $m \in D$  such that  $n \geq m \implies n \in E$ .

*Remark.* Let  $E$  be an eventual subset of  $D$  such that  $n \geq m \implies n \in E$ . Then  $p \in E \not\Rightarrow p \geq m$ . ie, Subset  $E$  may contain elements that doesn't follow the above  $m$ .

*Remark.* [Joshi, 1983, 10.1.6]

Let  $E$  be an eventual subset of  $D$ , then  $E$  is a directed set.

1.  $m, n, p \in E, m \geq n, n \geq p \implies m, n, p \in D, m \geq n, n \geq p \implies m \geq p$
2.  $m \in E \implies m \in D \implies m \geq m$
3.  $m, n \in E \implies m, n \in D \implies \exists p \in D$  such that  $p \geq m$  and  $p \geq n$ .

$\exists m' \in D$  such that  $n' \geq m' \implies n' \in E$ . ( $E$  is an eventual subset of  $D$ )

$p, m' \in D \implies \exists p' \in D$  such that  $p' \geq p$  and  $p' \geq m'$ . ( $D$  is a directed Set)

$p' \geq m' \implies p' \in E$ . ( $E$  is eventual subset of  $D$  with respect to  $m'$ )

$p' \geq p, p \geq m \implies p' \geq m$  and  $p' \geq p, p \geq n \implies p' \geq n$ .

Therefore  $\forall m, n \in E, \exists p' \in E$  such that  $p' \geq m$  and  $p' \geq n$ .

**Definitions 7.4.5** (Net eventually in  $A$ ). [Joshi, 1983, 10.1.5]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Then  $S$  is eventually in subset  $A$  of  $X$  if  $S^{-1}(A)$  is an eventual subset of  $D$ .

*Remark.* Let  $S : D \rightarrow X$  be a net in  $X$ . Then  $S$  converges to  $x \in X$  if  $S$  is eventually in each nbd  $U$  of  $x$ .

**Definitions 7.4.6** (Cofinal subset). [Joshi, 1983, 10.1.7]

A subset  $F$  of a directed  $D$  is a cofinal subset of  $D$  if for any  $m \in D$ , there exists  $n \in F$  such that  $n \geq m$ .

*Remark.* Let  $X$  be a topological space and  $x \in X$ . Let  $\mathcal{N}_x$  be the set of all neighbourhood of  $x$  and  $\mathcal{L}$  be a local base of  $X$  at  $x$ . We have,  $(\mathcal{N}_x, \geq)$  is a directed set where  $\forall U, V \in \mathcal{N}_x, U \geq V \iff U \subset V$ , then  $\mathcal{L}$  is cofinal in  $\mathcal{N}_x$ .

*Remark.* [Joshi, 1983, 10.1.8]

Let  $F$  be a cofinal subset of  $D$ , then  $F$  is a directed set.

1.  $m, n, p \in F, m \geq n, n \geq p \implies m, n, p \in D, m \geq n, n \geq p \implies m \geq p$
2.  $m \in F \implies m \in D \implies m \geq m$



3.  $m, n \in F \implies m, n \in D \implies \exists p \in D$  such that  $p \geq m$  and  $p \geq n$ .

$E$  is cofinal,  $p \in D \implies \exists p' \in F$  such that  $p' \geq p$ .

$p' \geq p, p \geq m \implies p' \geq m$  and  $p' \geq p, p \geq n \implies p' \geq n$ .

Therefore  $\forall m, n \in F, \exists p' \in F$  such that  $p' \geq m$  and  $p' \geq n$ .

**Definitions 7.4.7** (Net frequently in  $A$ ). [Joshi, 1983, 10.1.7]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Then  $S$  is frequently in subset  $B$  of  $X$  if  $S^{-1}(B)$  is a cofinal subset of  $D$ .

**Proposition 7.4.1.** [Joshi, 1983, 10.1.6]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Let  $E$  be an eventual subset of  $D$ . Then,  $S$  converges to  $x$  iff  $S_{/E}$  converges to  $x$ . cite[10.1.6]joshi

*Proof.* Let  $S : D \rightarrow X$  be a net in  $X$ ,  $E$  be an eventual subset of  $D$ , and  $x \in X$ . Then,  $S_{/E} : E \rightarrow X$  is defined by  $n \in E \implies S_{/E}(n) = S(n)$

Suppose  $S$  converges to  $x$ . Let  $U$  be a nbd of  $x$ , then  $S$  is eventually in  $U$ . ie,  $S^{-1}(U)$  is an eventual subset of  $D$ . Then  $\exists m \in D$  such that  $n \geq m \implies n \in S^{-1}(U) \implies S(n) \in U$ . Since set  $E$  is eventual subset of  $D$ ,  $\exists m' \in D$  such that  $n \geq m' \implies n \in E$ .

Since  $E$  is a directed set,  $S_{/E} : E \rightarrow X$  is a net in  $X$ . And  $m, m' \in D \implies \exists p \in D$  such that  $p \geq m$  and  $p \geq m'$ . We have,  $p \geq m' \implies p \in E$ . And  $n \geq' p \implies n \geq p, p \geq m \implies n \geq m \implies S(n) \in U \implies S_{/E}(n) \in U$ . Therefore,  $n \geq' p \implies S_{/E}(n) \in U$ . Since  $U$  is arbitrary,  $S_{/E}$  converges to  $x$ .

Conversely, suppose that  $S_{/E}$  converges to  $x$ . Let  $U$  be a nbd of  $x$ , then  $S_{/E}$  is eventually in  $U$ . ie,  $S_{/E}^{-1}(U)$  is an eventual subset of  $D$ . ie,  $\exists m \in D$  such that  $n \geq m \implies n \in S_{/E}^{-1}(U) \implies S_{/E}(n) \in U \implies S(n) \in U$ . Therefore,  $n \geq m \implies S(n) \in U$ . Since,  $U$  is arbitrary,  $S$  converges to every nbd of  $x$ . ie,  $S$  converges to  $x$ .  $\square$

**Proposition 7.4.2.** [Joshi, 1983, 10.1.8]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Let  $F$  be a cofinal subset of  $D$ . If  $S$  converges to  $x$ , then  $S_{/F}$  converges to  $x$ .

*Proof.* Let  $S : D \rightarrow X$  be a net in  $X$  and  $S$  converges to  $x \in X$ . Also let  $F$  be a cofinal subset of  $D$ . Then  $S_{/F}$  is also a net in  $X$ , since  $(F, \geq')$  is a directed set where  $\forall m, n \in F, m \geq n \implies m \geq' n$ .

Since  $S$  converges to  $x$ , for any nbd  $U$  of  $x$ ,  $\exists m \in D$ , such that  $n \geq m \implies S(n) \in U$ . Since  $F$  is cofinal,  $\exists p \in F$  such that  $p \geq m$ . Thus  $n \geq' p \implies n \geq p, p \geq m \implies n \geq m \implies S(n) \in U \implies S_{/F}(n) \in U$ . Therefore,  $\exists p \in F$  such that  $n \geq' p \implies S_{/F}(n) \in U$ . Since  $U$  is arbitrary,  $S_{/F}$  is eventually in every nbd of  $x$ . ie,  $S_{/F}$  converges to  $x$ .  $\square$

*Remark.* But converse of the above is not true.  $S_{/F}$  converges to  $x$  does not imply that  $S$  converges to  $x$ , since cofinal subset  $F$  not necessarily contain every element following a particular  $m$ .

**Definitions 7.4.8** (Cluster point). [Joshi, 1983, 10.1.9]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Then  $x \in X$  is a cluster point of  $S$ , if  $S$  is frequently in each nbd  $U$  of  $x$  in  $X$ .

**Proposition 7.4.3.** [Joshi, 1983, 10.1.10]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Then  $x \in X$  is a cluster point of  $X$ , if  $S_{/F}$  converges to  $x$  for some cofinal subset  $F$  of  $D$ .

*Proof.* Let  $S : D \rightarrow X$  be a net in  $X$  and  $(F, \geq')$  be a cofinal subset of  $(D, \geq)$ . Then  $S_{/F}$  is also a net in  $X$ . Suppose  $S_{/F}$  converges to  $x \in X$ . Let  $U$  be a nbd of  $x$ , then  $\exists m \in F$  such that  $n \geq' m \implies S_{/F}(n) \in U$ .

Let  $m' \in D$ . Then  $\exists p' \in F$  such that  $p' \geq m'$ , since  $F$  is a cofinal subset of  $D$ . We have,  $m, p' \in F$ , then  $\exists p \in F$  such that  $p \geq' m$  and  $p \geq' p'$ . Since  $F \subset D$ , we have  $p, m \in F \implies p, m \in D$  and  $p \geq' m \implies p \geq m$ .

Also  $p \geq' m \implies S_{/F}(p) \in U \implies S(p) \in U$ . Therefore,  $\forall m' \in D$ ,  $\exists p \in D$  such that  $p \geq m'$  and  $S(p) \in U$ . Since  $U, m'$  are arbitrary,  $S$  is frequently in every nbd of  $x$ . ie,  $x$  is a cluster point of  $S$ .  $\square$

**Definitions 7.4.9** (Subnet). [Joshi, 1983, 10.1.11]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Then a net  $T : E \rightarrow X$  in  $X$ , is a subnet of  $S$  if there exists a function  $N : E \rightarrow D$  such that  $S \circ N = T$  and  $\forall m \in D$ ,  $\exists p \in E$  such that  $n \geq' p \implies N(n) \geq m$ .

*Remark.* A net  $T : E \rightarrow X$  is a subnet of  $S : D \rightarrow X$  if  $\exists N : E \rightarrow D$  such that  $S \circ N = T$  and  $S$  is frequently in  $T(E)$ .

Let  $T : E \rightarrow X$  be a subnet of  $S : D \rightarrow X$  and  $A$  be a subset of  $X$ . If  $T$  eventually in  $A$ , then  $S$  is frequently in  $A$ .

**Proposition 7.4.4.** [Joshi, 1983, 10.1.12]

Let  $S : D \rightarrow X$  be a net in a topological space  $X$ . Then  $x \in X$  is a cluster point of  $S$  iff there exists a subnet of  $S$  which converges to  $x$ .

*Synopsis.* Let  $(D, \geq)$ ,  $(E, \geq')$  be two directed sets. And  $T : E \rightarrow X$  be a subnet of  $S : D \rightarrow X$ .

If  $T$  converges to  $x$ , then  $T$  is eventually in each nbd  $U$  of  $x$ . And since  $T$  is a subnet of  $S$ , there exists  $N : E \rightarrow D$  such that  $N(E)$  is a cofinal subset of  $D$ . Therefore,  $S$  is frequently in each nbd  $U$  of  $x$ . Thus,  $x$  is a cluster point of  $S$ .

If  $x$  is a cluster point of  $X$ , then  $S$  is frequently in every nbd of  $x$ . Let  $N : E \rightarrow D$  be  $N(n, U) = n$ . Construct a directed subset  $E$  of  $D \times \mathcal{N}_x$  such that  $(n, U) \in E \iff S(n) \in U$ . Now  $T$  is eventually in every nbd  $U$  of  $x$ , as those points with images outside  $U$  are removed by construction. Therefore, it is sufficient to show that  $T : E \rightarrow X$  is a subnet of the net  $S : D \rightarrow X$ . Clearly,  $E$  is a directed set and  $N : E \rightarrow D$  defined by  $N(n, U) = n$  satisfies both  $S \circ N = T$  and  $\forall m \in D$ ,  $\exists p \in E$  such that  $m \geq p \implies N(m) \geq p$ .

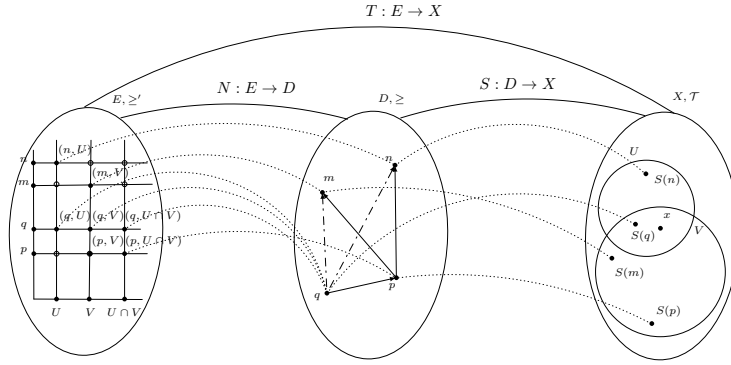


Figure 7.4:  $\forall (n, U), (m, V) \in E, \exists (q, W) \in E$  such that  $(q, W) \geq' (n, U)$ ,  $(q, W) \geq' (m, V)$

*Proof.* Let  $S : D \rightarrow X$  be a net in  $X$ . Suppose there exists a subnet  $T : E \rightarrow X$  that converges to  $x \in X$ . By the definition of subnet, we have  $\exists N : E \rightarrow D$  such that  $S \circ N = T$  and  $S$  is frequently in  $T(E)$ .

We have,  $T$  converges to  $x$ , thus for any neighbourhood  $U$  of  $x$ , there exists  $m' \in E$  such that  $n' \geq' m' \implies T(n') \in U$ .

Also we have,  $T$  is a subnet of  $S$ . Then  $\exists N : E \rightarrow D$  such that  $\forall m \in D, \exists p' \in E$  such that  $n' \geq' p' \implies N(n') \geq m$ .

Now, for any  $m \in D$ , we have  $m', p' \in E$ . Since  $E$  is a directed set, there exists  $n' \in E$  such that  $n' \geq' m'$  and  $n' \geq' p'$ .

Then,  $n' \geq m' \implies T(n') \in U$  and  $n' \geq' p' \implies N(n') \geq m$ .

Thus for any  $m \in D$ , there exists  $N(n') = n \in D$  such that  $S(n) = S(N(n')) = T(n') \in U$ .

Thus  $S$  is frequently in any neighbourhood  $U$  of  $x$ . Therefore,  $x$  is a cluster point of  $S$ .

Conversely, suppose that  $x$  is a cluster point of  $S$ . We have to construct a directed set  $(E, \geq')$  and a function  $N : E \rightarrow D$  such that  $T$  is a subnet of  $S$  and  $T$  converges to  $x$ . Let  $\mathcal{N}_x$  be the family of all neighbourhood of  $x$  in  $X$ .

Consider  $E = \{(n, U) \in D \times \mathcal{N}(x) : S(n) \in U\}$  and define  $\geq'$  by  $(n, U) \geq' (m, V)$  if  $n \geq m$  and  $U \subset V$ . Trivially,  $(n, U) \geq' (m, V) \geq' (p, W) \implies (n, U) \geq' (p, W)$  and  $(n, U) \geq' (n, U)$ . Also, for any  $(n, U), (m, V) \in E$ , we have  $n, m \in D$  and  $U, V \in \mathcal{N}(x)$ . Since  $D$  is a directed set,  $\exists p \in D$  such that  $p \geq n$  and  $p \geq m$ . Also,  $U \cap V \in \mathcal{N}_x$  and there exists  $q \in D$  such that  $S(q) \in U \cap V$  and  $q \geq' p$ , since  $S$  is frequently in every nbd of  $x$ . And  $U \cap V \in \mathcal{N}(x)$  such that  $U \cap V \subset U$  and  $U \cap V \subset V$ . Thus  $\exists (q, U \cap V) \in E$  such that  $(q, U \cap V) \geq' (n, U)$  and  $(q, U \cap V) \geq' (m, V)$ . Therefore,  $(E, \geq')$  is a directed set.

Define  $N : E \rightarrow D$  by  $N(n, U) = n$ . Again for any  $(m, V) \in E$ , there exists  $m \in D$  such that  $(n, U) \geq' (m, U)$  implies there exists  $n \in D$  such that  $N(n, U) = n$  and  $n \geq m$ .

It remains to prove that,  $T$  converges to  $x$ . Let  $U \in \mathcal{N}_x$  be a nbd of  $x$  in  $X$ . We have,  $x$  is a cluster point of  $S$ . Therefore,  $\forall m \in D, \exists p \in D$  such that  $S(p) \in U$ . By the construction of  $E$ , we have  $(p, U) \in E$ . Suppose  $(n, V) \geq' (p, U)$ , then  $n \geq p$ ,  $V \subset U$ , and  $S(n) \in V$ . Clearly  $S(n) \in U$ . Therefore,  $\forall (n, V) \geq' (p, U), T(n, V) \in U$ . That is,  $T$  is eventually in every nbd of  $x$ . ie, subnet  $T$  is convergent to  $x$ . Therefore, for each cluster point of the net  $S$ , there exists some subnet converging to it.  $\square$

*Remark.* • Importance of Construction of  $E$

If  $x$  is a cluster point of a net  $S$  in  $X$ , then  $S$  is frequently in some cofinal subset of  $D$ . Thus, if we consider any cofinal subset  $D'$  of  $D$  which is a direct set with  $\geq$  restricted to  $D'$ . Then  $N : D' \rightarrow D$  defined by  $N(n) = n$  gives a subnet  $T : D' \rightarrow X$  of the net  $S$ . However, this subnet need not converge to  $x$ . The strongest statement, we can make on  $T$  is that ' $x$  is a cluster point of  $T$ ', since  $N : D \times \mathcal{N}(x) \rightarrow D$ ,  $N(n) = n$  is completely independent of  $U$ . This problem is overcome by constructing  $E$  which is dependent on each nbd  $U$  of  $x$ .

- Existence of  $q \in D$  such that  $q \geq' p$  and  $S(q) \in U$ . We have,  $p$  follows both  $n \& m$  and  $U \cap V$  is a subset of both  $U \& V$ . However, since  $S$  is only frequently in  $U$ ,  $p$  not necessarily be in  $U$ . But there is always someone following  $p$  which has its image in  $U$ . This  $q$  follows both  $n \& m$ , since  $\geq'$  is transitive.

## 7.5 Variations of Compactness

### 7.5.1 Variations of Compactness

In this chapter, we have two other notions of compactness - countable compactness and sequential compactness. ( $\star^{13}$ )

**Compact** A topological space is compact iff every open cover of it has a finite subcover. ([Joshi, 1983, 6.1.1]) [Heine-Borel]

**Countably Compact** A topological space is countably compact iff every countable, open cover of it has a finite subcover. [Joshi, 1983, 11.1.1]

**Sequentially Compact** A topological space is sequentially compact iff every sequence in it has a convergent subsequence. [Joshi, 1983, 11.1.8] [Bolzano-Weierstrass]

Countable compactness is a weaker notion compared to compactness. ( $\star^{14}$ ) However, sequentially compact and compact are not necessarily comparable. ( $\star^{15}$ )

<sup>13</sup>For  $\mathbb{R}$ , Compactness & Sequentially compactness are equivalent to the completeness axiom.

<sup>14</sup>Every compact space is countably compact.

<sup>15</sup> $\mathcal{T}_1, \mathcal{T}_2$  are non-comparable, if  $\mathcal{T}_1 \not\subset \mathcal{T}_2$  and  $\mathcal{T}_2 \not\subset \mathcal{T}_1$ . [Joshi, 1983, 4.2.1]

We have seen earlier that compactness has the following properties 1. compactness is weakly hereditary. [Joshi, 1983, 6.1.10] 2. compactness is preserved under continuous functions. [Joshi, 1983, 6.1.8] 3. every continuous real functions on compact space is bounded and attains its extrema. [Joshi, 1983, 6.1.6] 4. every continuous real function on a compact, metric space is uniformly continuous by Lebesgue covering lemma. [Joshi, 1983, 6.1.7]

Countably compact spaces, Sequentially compact spaces have all the four properties listed above.

### 7.5.2 Countable compactness

#### Weakly hereditary property

A subspace  $(A, \mathcal{T}_A)$  being countably compact doesn't imply that  $(X, \mathcal{T})$  is countably compact. However, if  $(X, \mathcal{T})$  is a countably compact space and  $A$  is a closed subset of  $X$ , then  $(A, \mathcal{T}_A)$  is also a countably compact space. **In other words, countable compactness is weakly hereditary.**

**Theorem 7.5.1.** *Countable compactness is weakly hereditary. [Joshi, 1983, 11.1.3]*

*Synopsis.* Let  $A$  be a closed subset of countably compact space,  $X$ . If  $A$  has a countable open cover  $\mathcal{U}$ , then we can obtain a respective countable, open cover for  $X$  by attaching  $X - A$  to the extensions of members of  $\mathcal{U}$  to  $X$ . This cover has a finite subcover. Then restricting them to  $A$ , we get a finite subcover of  $\mathcal{U}$ .

*Proof.* Suppose  $X$  is a countably compact space. And  $A$  is a closed subset of  $X$ . We need to show that  $A$  is countably compact. Without loss of generality, <sup>(★<sup>16</sup>)</sup> assume that  $A$  is a proper subset of  $X$ . Then  $X - A$  is a non-empty, open subset of  $X$ .

Let  $\mathcal{U}$  be a countable open cover of  $A$ . Then  $\mathcal{U} = \{U_1, U_2, \dots\}$  where each element  $U_k \in \mathcal{U}$  is an open subset of  $A$ . Since  $A$  is a subspace of  $X$ , every open subset  $U_k$  in  $A$  is of the form  $G \cap A$  for some open subset  $G$  in  $X$ . Therefore, there exists open subsets  $V(U_k)$  for each  $U_k$  such that  $A \cap V(U_k) = U_k$ . <sup>(★<sup>17</sup>)</sup>

Define  $\mathcal{V} = \{X - A, V(U_1), V(U_2), \dots\}$ . Clearly,  $\mathcal{V}$  is a countable open cover <sup>(★<sup>18</sup>)</sup> of  $X$ . We have  $X$  is countably compact, thus  $\mathcal{V}$  has a finite subcover, say  $\mathcal{V}'$ . Without loss of generality assume <sup>(★<sup>19</sup>)</sup> that  $X - A \in \mathcal{V}'$ . Suppose  $X - A \notin \mathcal{V}'$ , then we can define another finite subcover  $\mathcal{V}' \cup \{X - A\}$ . Thus  $\mathcal{V}' = \{X - A, V(U_{n_1}), V(U_{n_2}), \dots, V(U_{n_k})\}$ .

Then the corresponding subcover  $\mathcal{U}' = \{U_{n_1}, U_{n_2}, \dots, U_{n_k}\}$  is a finite subcover of  $\mathcal{U}$ . Since countable open cover  $\mathcal{U}$  and closed subset  $A$  are arbitrary, every closed subset of  $X$  with relative topology is countably compact. Therefore, countable compactness is weakly hereditary.  $\square$

<sup>16</sup>Suppose  $A$  is not a proper subset of  $X$ . Then  $X = A$  and  $A$  is countably compact.

<sup>17</sup>Relative topology,  $\mathcal{T}_A = \{G \cap A : G \in \mathcal{T}\}$

<sup>18</sup> $X - A$  is open in  $X$ . If  $y \notin A$ , then  $y \in X - A$ . If  $y \in A$ , then  $y \in U_k$  for some  $k$ .

<sup>19</sup>Otherwise, you will have to consider two cases:  $X - A \in \mathcal{V}'$  and  $X - A \notin \mathcal{V}'$

*Remark.* Proof depends on the following,

1. There is an extension map,  $\psi : P(A) \rightarrow P(X)$  that preserve open subsets (and closed subsets). This  $\psi$  is an open map which not a true inverse of the restriction,  $r : P(X) \rightarrow P(A)$ , defined by  $r(G) = G \cap A$  for every subset  $G$  of  $X$ .
2. Also we rely on the subset  $A$  being closed. Suppose  $X$  have many countable open covers, but  $X$  has only uncountable open covers corresponding to a particular countable open cover of  $A$ . In such a case,  $X$  being countably compact is insufficient for  $A$  to be countably compact.

### The behaviour of countinuous functions

We will now study the nature of continuous functions defined on countably compact spaces. Suppose  $X, Y$  are topological space and function  $f : X \rightarrow Y$  is continuous. If  $X$  is countably compact, then  $f(X)$  is also countably compact. Continuous images of countably compact spaces are countably compact. **In other words, countable compactness is preserved under continuous functions.**( $\star^{20}$ )

**Theorem 7.5.2.** *Countable compactness is preserved under continuous functions. [Joshi, 1983, 11.1.2]*

*Synopsis.* Let  $X$  be countably compact and  $f : X \rightarrow Y$  be continuous. Suppose  $\mathcal{U}$  is a countable cover of  $f(X)$ , then  $X$  has a countable cover  $\mathcal{V}$  obtained by taking inverse images. Since  $X$  is countably compact,  $\mathcal{V}$  has a finite subcover  $\mathcal{V}'$ . Now taking images of members of  $\mathcal{V}'$ , we get a finite subcover  $\mathcal{U}'$  of  $f(X)$ .

*Proof.* Suppose  $X$  is a countably compact space,  $Y$  is a topological space and  $f : X \rightarrow Y$  is a continuous function. Let  $\mathcal{U} = \{U_1, U_2, \dots\}$  be a countable cover of  $f(X)$  by set open in  $f(X)$ . We have to show that  $\mathcal{U}$  has a finite subcover.

Define  $\mathcal{V} = \{f^{-1}(U_1), f^{-1}(U_2), \dots\}$ . Then  $\mathcal{V}$  is a countable open cover of  $X$ , since  $f^{-1}(U_k)$  are open subsets of  $X$  and,

$$\begin{aligned} \bigcup_{k=1}^{\infty} U_k = f(X) &\implies f^{-1}\left(\bigcup_{k=1}^{\infty} U_k\right) = X \\ &\implies \bigcup_{k=1}^{\infty} f^{-1}(U_k) = X \end{aligned}$$

We have,  $\mathcal{V}$  is a countable open cover of  $X$ , which is a countably compact space. Therefore  $\mathcal{V}$  has a finite subcover  $\mathcal{V}' = \{f^{-1}(U_{n_1}), f^{-1}(U_{n_2}), \dots, f^{-1}(U_{n_k})\}$ .

$$\begin{aligned} \bigcup_{j=1}^k f^{-1}(U_{n_j}) = X &\implies f^{-1}\left(\bigcup_{j=1}^k U_{n_j}\right) = X \\ &\implies \bigcup_{j=1}^k U_{n_j} = f(X) \end{aligned}$$

---

<sup>20</sup>A topological property is preserved under continuous functions if whenever a space has that property so does every continuous image of it. [Joshi, 1983, 6.1.9]

Clearly  $\mathcal{U}' = \{U_{n_1}, U_{n_2}, \dots, U_{n_k}\}$  is a finite subcover of  $\mathcal{U}$ . Thus every countable open cover of  $f(X)$  by sets open in  $f(X)$  has a finite subcover. Therefore, continuous images of countably compact spaces are countably compact.  $\square$

*Remark.* 1. For a continuous function,  $f : X \rightarrow Y$  the inverse images of open subsets are open in  $X$ . The relation  $f^{-1} \subset f(X) \times X$  is not a function. However, we may consider a function,  $\psi : P(Y) \rightarrow P(X)$  such that  $\psi(U) = f^{-1}(U)$  for any subset  $U$  of  $Y$ . This  $\psi$  is an open map which maps open subsets of  $Y$  to open subsets of  $X$ .

**Theorem 7.5.3.** *Every continuous, real-valued function on a countably compact, metric space is bounded and attains its extrema. [Joshi, 1983, 11.1.7]*

*Synopsis.* Let  $X$  be a countably compact space and function  $f : X \rightarrow \mathbb{R}$  be continuous. Then  $f(X) \subset \mathbb{R}$  is countably compact. Real line  $\mathbb{R}$  is metrisable( $\star^{21}$ ). Then  $f(X)$  is countably compact, metric space. Therefore  $f(X)$  compact( $\star^{22}$ ). The subset  $f(X)$  of  $\mathbb{R}$  is bounded and closed, since every compact subset of  $\mathbb{R}$  is bounded and closed. Thus  $f(X)$  contains its supremum and infimum. Therefore,  $f$  is bounded and attains its extrema.

*Proof.* Let  $X$  be a countably compact space and  $f : X \rightarrow \mathbb{R}$  be continuous, real-valued function on the countably compact space,  $X$ . We have to show that  $f$  is bounded and attains its extrema.

Since countable compactness is preserved under continuous functions,  $f(X)$  is countably compact subset of  $\mathbb{R}$ . Since,  $f(X)$  is a subset of the metric space,  $\mathbb{R}$  and metrisability is hereditary,  $f(X)$  is again metrisable. (suppose) We have, every countably compact, metric space is compact. Then  $f(X)$  is a compact subset of  $\mathbb{R}$ .

Since every compact subset of  $\mathbb{R}$  is bounded and closed,  $f(X)$  is bounded and closed. Since every closed subset of  $\mathbb{R}$  contains supremum and infimum,  $f(X)$  contains its extrema. Therefore, every continuous, real-valued function on a countably compact space is bounded and attains its extrema.

We have assumed that every countably compact, metric space is compact. This result will be proved in the last section of this chapter.  $\square$

*Remark.* Since countably compact, metric spaces are compact. The above theorem can be used to prove that continuous, real-valued functions on a compact, metric space attains its extrema.

Due to the Lebesgue covering lemma, next result is quite simple.\*

**Theorem 7.5.4.** *Every continuous, real-valued function on a countably compact, metric space is uniformly continuous.*

**Proposition 7.5.1.** *Let  $X$  be a first countable, Hausdorff space. Then every countably compact subset  $A$  of  $X$  is closed. [Joshi, 1983, Exercises 11.1.7]*

<sup>21</sup>[Joshi, 1983, 4.2 Example 4],  $\mathbb{R}$  with usual metric  $d : \mathbb{R} \rightarrow \mathbb{R}$ ,  $d(x, y) = |x - y|$

<sup>22</sup>[Joshi, 1983, 11.1.11] On metric spaces, countable compactness  $\implies$  compactness.

### 7.5.3 Sequential Compactness

#### Weakly hereditary property

**Theorem 7.5.5.** *Sequential compactness is weakly hereditary. [Joshi, 1983, Exercises 11.1.3]*

#### The behaviour of continuous functions

**Theorem 7.5.6.** *Sequential compactness is preserved under continuous functions. [Joshi, 1983, Exercises 11.1.4]*

*Synopsis.* Let  $X$  be sequentially compact and function  $f : X \rightarrow Y$  be continuous. Then any sequence,  $\{y_k\}$  in  $f(X)$  has a sequence,  $\{x_k\}$  in  $X$  such that  $f(x_k) = y_k$ . Sequence  $\{x_k\}$  has a subsequence  $\{x_{n_k}\}$  converging to  $x$ , then sequence  $\{f(x_{n_k})\}$  in  $f(X)$  has the subsequence  $\{f(x_{n_k})\}$  converging to  $f(x)$ .

*Proof.* Let  $X$  be a sequentially compact space, function  $f : X \rightarrow Y$  be continuous and  $\{y_n\}$  be a sequence in  $f(X)$  subset of  $Y$ . Construct a sequence  $\{x_n\}$  such that  $f(x_k) = y_k, \forall k$ .

Every sequence in  $X$  has a convergent subsequence. Thus  $\{x_n\}$  has a subsequence  $\{x_{n_k}\}$  converging to  $x \in X$ . The image of this subsequence  $\{f(x_{n_k})\}$  is a subsequence of  $\{y_k\}$ . We claim that,  $\{f(x_{n_k})\}$  converges to  $f(x) \in f(X)$ .

Let  $U$  be an open subset containing  $f(x)$ , then  $f^{-1}(U)$  is an open subset containing  $x$ . Since  $\{x_{n_k}\}$  converges to  $x$ . There exists an integer  $n$  such that for every  $k \geq n, x_k \in f^{-1}(U)$ . Clearly, for each  $k \geq n, f(x_k) \in U$ . Since  $U$  is arbitrary,  $\{f(x_{n_k})\}$  converges to  $f(x)$ . Therefore, the image of any sequentially compact space is sequentially compact. In other words, sequential compactness is preserved under continuous functions.  $\square$

*Remark.* 1. Given a sequence  $\{y_n\}$  in  $f(X)$ , there is a sequence of subsets  $\{U_n\}$  in  $P(Y)$  such that  $U_n = f^{-1}(y_n)$ . Since each  $U_n$  is non-empty, we can construct a sequence  $\{x_n\}$  in  $X$  using a choice function. The convergent subsequence of  $\{y_n\}$  depends on the selection of this choice function.

Given every sequentially compact, metric space is countably compact. We may assert the properties of countably compact, metric spaces on sequentially compact, metric spaces.

**Theorem 7.5.7.** *Every continuous, real-valued function on a sequentially compact, metric space is bounded and attains its extrema.*

**Theorem 7.5.8.** *Every continuous, real-valued function on a sequentially compact, metric space is uniformly continuous. [Joshi, 1983, Exercises 11.1.6]*

### 7.5.4 Countable Compactness on $T_1$ spaces

In this section, we are going to see four different characterisations of countable compactness in  $T_1$  spaces. The first two characterisations doesn't have anything to do with the  $T_1$  axiom.



**$T_1$  Space** A topological space  $X$  satisfy  $T_1$  axiom if for any two distinct points  $x, y \in X$ , there exists an open subset  $U \subset X$  containing  $x$  but not  $y$ . [Joshi, 1983, 7.1.2]

**countable compactness** A topological space is countably compact if every countable open cover has a finite subcover. [Joshi, 1983, 11.1.1]

**finite intersection property** A family  $\mathcal{F}$  of subsets of  $X$  has finite intersection property(f.i.p.) if every finite subfamily of  $\mathcal{F}$  has a non-empty intersection. [Joshi, 1983, 10.2.6]

**accumulation point** A point  $x \in X$  is accumulation point of a subset  $A \subset X$  if every open subset containing  $x$  has atleast one point of  $A$  other than  $x$ . [Joshi, 1983, 5.2.7]

**limit point** A point  $x \in X$  is a limit point of a sequence  $\langle x_k \rangle$  in  $X$  if for every open subset  $U$  containing  $x$ , there exists an integer  $N \in \mathbb{N}$  such that  $x_k \in U$  for every  $k \geq N$ . [Joshi, 1983, 4.1.7]

**cluster point** A point  $x \in X$  is a cluster point of a sequence  $\langle x_k \rangle$  in  $X$  if for any neighbourhood  $V$  of  $x$ , the sequence  $\langle x_k \rangle$  assumes a point in  $V$  infinitely many times.( $\star^{23}$ )

#### Countable compactness in $T_1$ spaces

**Theorem 7.5.9.** In a  $T_1$  space  $X$ , following statements are equivalent,

1.  $X$  is countably compact
2. Every countably family of closed subsets of  $X$  with finite intersection property have non-empty intersection.
3. Every infinite subset  $A \subset X$  has an accumulation point.( $\star^{24}$ )
4. Every sequence  $\langle x_k \rangle$  in  $X$  has a cluster point.
5. Every infinite open cover of  $X$  has a proper subcover.[Arens-Dugundji]

*Proof.* 1  $\implies$  2

Suppose  $X$  is countably compact. Let  $\mathcal{C} = \{C_1, C_2, \dots\}$  be a countable family of closed subsets of  $X$  with empty intersection. Define  $\mathcal{U} = \{X - C_1, X - C_2, \dots\}$  is a family of open subsets of  $X$ . By de Morgan's law, ( $\star^{25}$ )

$$\bigcap_{k=1}^{\infty} C_k = \phi, \text{ then } X = X - \left( \bigcap_{k=1}^{\infty} C_k \right) = \bigcup_{k=1}^{\infty} (X - C_k)$$

We have  $\mathcal{U}$  is a countable cover of  $X$  and  $X$  is countably compact space. Thus  $\mathcal{U}$  has a finite subcover  $\mathcal{U}' = \{X - C_{n_1}, X - C_{n_2}, \dots, X - C_{n_k}\}$ .

$$\mathcal{U}' \text{ is a cover of } X, \text{ then } X = \bigcup_{j=1}^k (X - C_{n_j})$$

<sup>23</sup> $x$  is a cluster point of  $\langle x_k \rangle$  if for every integer  $N$ , there exists  $k > N$  such that  $x_k \in V$ . In other words,  $\langle x_k \rangle$  is frequently in  $V$ . [Joshi, 1983, 10.1.9]

<sup>24</sup>Every infinite subset of  $\mathbb{R}$  has a limit point is equivalent to the completeness axiom.

<sup>25</sup>Complement of Intersection = Union of complements,  $X - (C \cap D) = (X - C) \cup (X - D)$ ,

$$X - \bigcup_{j=1}^k (X - C_{n_j}) = \bigcap_{j=1}^k (X - (X - C_{n_j})) = \bigcap_{j=1}^k C_{n_j} = \phi$$

Now  $\mathcal{C}' = \{C_{n_1}, C_{n_2}, \dots, C_{n_k}\}$  has empty intersection. This is a contradiction to the finite intersection property of  $\mathcal{C}$ . Thus  $\mathcal{C}$  has non-empty intersection. Therefore, every countably family of closed subsets of  $X$  have non-empty intersection.

2  $\implies$  1

Let  $\mathcal{U} = \{U_1, U_2, \dots\}$  be a countable cover of  $X$ . Then  $\mathcal{C} = \{X - U_1, X - U_2, \dots\}$  is a countable family of closed subsets of  $X$ .

Let  $\mathcal{U}' = \{U_{n_1}, U_{n_2}, \dots, U_{n_k}\}$  be any finite subfamily of  $\mathcal{U}$ . Suppose  $X$  is not countably compact, then  $\mathcal{U}$  doesn't have a finite subcover. Therefore,  $\mathcal{U}'$  is not a cover of  $X$ . And  $\mathcal{C}$  is a family of closed subsets with finite intersection property.

Therefore by assumption, the countable family of closed subsets  $\mathcal{C}$  has a non-empty intersection.

$$\bigcap_{k=1}^{\infty} C_k \neq \phi, \text{ then } \bigcap_{k=1}^{\infty} C_k = \bigcap_{k=1}^{\infty} (X - U_k) = X - \left( \bigcup_{k=1}^{\infty} U_k \right) \neq \phi$$

Then  $\mathcal{U}$  is not a cover of  $X$  as well. This is a contradiction, therefore  $X$  is countably compact.

1  $\implies$  3

Suppose  $X$  is countably compact. Let  $A$  be an infinite subset of  $X$ . Suppose  $A$  doesn't have an accumulation point.

Let  $B$  be a countably infinite subset of  $A$ . Then  $B$  also doesn't have any accumulation point. Therefore, the derived set  $B'$  is empty. Thus  $B$  is a closed subset of  $X$ . Since countable compactness is weakly hereditary, subspace  $B$  is again countably compact.

For each point  $b \in B$ , there is an open subset  $V_b$  such that  $V_b \cap B = \{b\}$ , since  $b \in B$  is not an accumulation point. Thus  $\mathcal{U} = \{V_b \cap B : b \in B\}$  is a countable open cover of  $B$ . Clearly,  $\mathcal{U}$  doesn't have any finite subcover.

This is a contradiction to  $B$  being countably compact. Therefore,  $A$  has an accumulation point.  $\square$

### 7.5.5 Variations of Compactness on Metric Spaces

In this document, we will see that from metric space point of view these two notions were equivalent to the compactness and were used alternatively. For example : in functional analysis (semester 3), you will find definitions like 'a normed space is compact iff every sequence in it has a convergent subsequence', which is clearly sequential compactness for a topologist.

**Lindeloff** A topological space is Lindeloff iff every open cover has a countable subcover.

**First countable** A topological space is first countable iff every point in it has a countable local base.

**Second countable** A topological space is second countable iff it has a countable base.

**Base** A family of subsets  $\mathcal{B}$  of  $X$  is a base of a topological space if every open subset can be expressed as union of some members of  $\mathcal{B}$

**Base Characterisation** A family of subsets  $\mathcal{B}$  of  $X$  is a base of a topological space iff for every  $x \in X$ , and for every neighbourhood  $U$  of  $x$ , there is a member  $B \in \mathcal{B}$  such that  $x \in B \subset U$ .

**Local Base** A family of subsets  $\mathcal{L}$  of  $X$  is a local base at point  $x \in X$  if for every neighbourhood  $U$  of  $x$ , there is a member  $L \in \mathcal{L}$  such that  $x \in L \subset U$ .

### Equivalence

We are going to see when these three notions: compactness, countable compactness and sequentially compactness are equivalent.

**Theorem 7.5.10.** *Countably compact, metric spaces are second countable.*

*Synopsis.* For every positive real number  $r$ , there exists a non-empty maximal subsets  $A_r$  with every pair of points atleast  $r$  distance apart.  $A_r$  are finite. The union of maximal subsets  $A_{\frac{1}{n}}$  for each natural number  $n$  is a countable, dense subset  $D$  of  $X$ . Thus countably compact, metric spaces are separable. The family  $\mathcal{B}$  of all open balls with center at  $d \in D$  and rational radius is a countable, base for  $X$ . Thus countably compact, metric spaces are second countable.

*Proof.* Let  $(X; d)$  be a countably compact, metric space. For each positive real number  $r \in \mathbb{R}$ ,  $r > 0$  construct a family of subsets  $A_r \subset X$  such that it is a maximal set of points which are atleast  $r$  distances apart.

Then  $A_r$  is finite for every positive real number  $r$ . Suppose  $A_r$  is infinite for some real number  $r > 0$ , then  $A_r$  has a accumulation point, say  $x$  by the Characterisation of countable compactness of  $X$ .

Then every neighbourhood of  $x$  must intersect  $A_r$  at infinitely many points, since every metric space is a  $T_1$  space. Consider  $B(x, \frac{r}{2})$ . Since any two points of  $B(x, \frac{r}{2})$  are less than  $r$  distances apart, the intersection  $B(x, \frac{r}{2}) \cap A_r$  can have atmost one point in it. Thus for every positive real number  $r$ ,  $A_r$  is finite.

Define  $D = \bigcup_{n=1}^{\infty} A_{\frac{1}{n}}$ . We claim that  $D$  is a countable, dense subset of  $X$ .

Let  $x \in X$  and  $B(x, r)$  be an open ball containing  $x$ , then there exists integer  $n \in \mathbb{N}$  such that  $\frac{1}{n} < r$ . <sup>★<sup>26</sup></sup>

---

<sup>26</sup>By archimedean property of integers, we have  $\forall r \in \mathbb{R}$ ,  $r > 0$ ,  $\exists n \in \mathbb{N}$  such that  $nr > 1$ .

Then  $B(x, r) \cap D \neq \phi$ , since  $B(x, r) \cap A_{\frac{1}{n}} \neq \phi$ . Suppose  $B(x, r) \cap A_{\frac{1}{n}} = \phi$ , then  $A_{\frac{1}{n}}$  is not maximal. Since,  $x$  is at least  $r > \frac{1}{n}$  distance apart from each point of  $A_{\frac{1}{n}}$ . Therefore,  $D$  intersects with every open subset and thus dense in  $X$ .

We have have a countable, dense subset  $D$  of  $X$ . Therefore,  $X$  is separable. Now define  $\mathcal{B} = \{B(x, r) : r \in \mathbb{Q}, x \in D\}$ . Clearly,  $\mathcal{B}$  is a countable base for  $X$ . By the construction of  $\mathcal{B}$ ,  $X$  is second countable. ( $\star^{27}$ )  $\square$

### Countable Compactness, Lindeloff $\iff$ Compactness

**Theorem 7.5.11.** *A topological space  $X$  is compact iff it is countably compact, Lindeloff space.*

*Proof.* Let  $X$  be a compact space. Let  $\mathcal{U}$  be a countable open cover of  $X$ , then  $\mathcal{U}$  has a finite subcover  $\mathcal{U}'$ . Therefore, every compact space is countably compact. ( $\star^{28}$ )

Conversely, suppose  $X$  is a countably compact, Lindeloff space. Since  $X$  is Lindeloff, every open cover  $\mathcal{U}$  has a countable subcover  $\mathcal{U}'$ . Since  $X$  countably compact, every countable open cover  $\mathcal{U}'$  has a finite subcover  $\mathcal{U}''$ . Thus every open cover  $\mathcal{U}$  has a finite subcover  $\mathcal{U}''$ . Therefore every countably compact, Lindeloff space is compact.  $\square$

### Countable Compactness, First Countable $\implies$ Seq. Compactness

**Theorem 7.5.12.** *Every countably compact, first countable space is Sequentially compact.*

*Proof.* Let  $X$  be a countably compact, first countable space. Let  $\{x_n\}$  be a sequence in  $X$ . By, equivalent conditions ( $\star^{29}$ ) of countably compact spaces, every sequence in countably compact space  $X$  has a cluster point, say  $x$ . We have,  $X$  is first countable. Therefore,  $X$  has a countable local base  $\mathcal{L}$  at  $x \in X$ . How to construct a subsequence of  $\{x_n\}$  converging to  $x$ ? ( $\star^{30}$ )  $\square$

*Remark.* Every sequentially compact space is countably compact.  $\star$

**Theorem 7.5.13.** *In a second countable space, all the three forms of compactness are equivalent. [Joshi, 1983, 11.1.10]*

*Proof.* Every second countable space is both first countable and Lindeloff. Every countably compact, Lindeloff space is countably compact. Therefore every countably compact, second countable space compact. Again, every countably compact, first countable space is sequentially compact. Therefore every countably compact, second countable space is sequentially compact. Conversely, every compact space is countably compact and every sequentially compact space is countably compact. ( $\star^{31}$ )  $\square$

<sup>27</sup>Every separable, metric space is second countable.

<sup>28</sup>Countable compactness is a weaker notion than compactness.

<sup>29</sup>[Joshi, 1983, 11.1] Conditions 1, 2, and 4 are equivalent.  $2 \implies 4$  without  $T_1$  axiom is out of scope.

<sup>30</sup>[Joshi, 1983, Exercises 10.1.11]

<sup>31</sup>Countable compactness is a weaker notion than sequential compactness as well.

**Theorem 7.5.14.** *In a metric space, all the three forms of compactness are equivalent. [Joshi, 1983, 11.1.11]*

*Proof.* In a metric space each form of compactness implies second countability. And in second countable spaces, they are all equivalent.  $\square$

## 7.6 Homotopy of Paths

**Definitions 7.6.1.** Let  $X, Y$  be topological spaces and  $f : X \rightarrow Y, f' : X \rightarrow Y$  be continuous functions. Then  $f, f'$  are homotopic if there exists a continuous function  $F : X \times I \rightarrow Y$  such that for every  $x \in X, F(x, 0) = f(x)$  and  $F(x, 1) = f'(x)$ . And we write,  $f \simeq f'$ .

**Definitions 7.6.2.** Let  $X$  be a topological space and  $f : I \rightarrow X$  and  $f' : I \rightarrow X$  be two paths. Then  $f, f'$  are path-homotopic if they have same initial point  $x_0$  (ie,  $x_0 = f(0) = f'(0)$ ), same final point  $x_1$  (ie,  $x_1 = f(1) = f'(1)$ ) and they are homotopic (ie,  $\exists F : I \times I \rightarrow X$  such that  $\forall x \in I, F(x, 0) = f(x)$  and  $F(x, 1) = f'(x)$  also fixed at the end points  $x_0$  and  $x_1$  (ie,  $\forall t \in I, F(0, t) = x_0$  and  $F(1, t) = x_1$ ). And we write  $f \simeq_p f'$ .

*Remark.* If two paths  $f, f'$  are homotopic, then they have the same end points and there exists a (topologically) continuous deformation from one path into another.

**Proposition 7.6.1.** *The relations  $\simeq, \simeq_p$  are equivalence relations.*

*Proof.* Homotopy : Let  $f, f'$  be continuous functions from  $X$  into  $Y$ . Then  $f$  and  $f'$  are homotopic,  $f \simeq f' \iff \exists F : X \times I \rightarrow Y$  such that  $F$  is continuous,  $F(x, 0) = f(x)$ , and  $F(x, 1) = f'(x)$

1.  $f \simeq f$

We have  $f : X \rightarrow Y$  is continuous. Define  $F : X \times I \rightarrow Y$  such that  $F(x, t) = f(x)$ . Clearly,  $F$  is continuous,  $F(x, 0) = f(x)$  and  $F(x, 1) = f(x)$ . And  $\exists F : X \times I \rightarrow Y \implies f \simeq f$ .

2.  $f \simeq f' \implies f' \simeq f$

We have,  $f \simeq f'$ . Thus there exists a continuous function  $F : X \times I \rightarrow Y$  such that  $F(x, 0) = f(x)$  and  $F(x, 1) = f'(x)$ .

Consider  $F' : X \times I \rightarrow Y$  defined by  $F'(x, t) = F(x, 1 - t)$ . Clearly,  $F'$  is continuous,  $F'(x, 0) = F(x, 1) = f'(x)$ , and  $F'(x, 1) = F(x, 0) = f(x)$ . Thus,  $\exists F'(x, t) : X \times I \rightarrow Y \implies f' \simeq_p f$

3.  $f \simeq f', f' \simeq f'' \implies f \simeq f''$

We have,  $f \simeq f' \iff \exists F : X \times I \rightarrow Y$  such that  $F$  is continuous,  $F(x, 0) = f(x)$  and  $F(x, 1) = f'(x)$ .

Similarly,  $f' \simeq f'' \iff \exists F' : X \times I \rightarrow Y$  such that  $F'$  is continuous,  $F'(x, 0) = f'(x)$  and  $F'(x, 1) = f''(x)$ .

Consider  $G : X \times I \rightarrow Y$  defined by

$$G(x, t) = \begin{cases} F(x, 2t) & , t \in [0, \frac{1}{2}] \\ F'(x, 2t - 1) & , t \in [\frac{1}{2}, 1] \end{cases}$$

We have,  $G(x, \frac{1}{2}) = F(x, 1) = F'(x, 0) = f'(x)$ . Thus,  $G$  is continuous by pasting lemma since  $[0, \frac{1}{2}] \cap [\frac{1}{2}, 1] = \{\frac{1}{2}\}$ . Also  $G(x, 0) = F(x, 0) = f(x)$  and  $G(x, 1) = F'(x, 1) = f''(x)$ . Thus,  $\exists G : X \times I \rightarrow Y \implies f \simeq f''$

Path Homotopy : Let  $f, f', f''$  be paths in  $X$ . Then  $f$  and  $f'$  are path homotopic,  $f \simeq_p f' \iff \exists F : I \times I \rightarrow X$  such that  $F$  is continuous,  $\forall s \in [0, 1], F(s, 0) = f(s), F(s, 1) = f'(s)$  and  $\forall t \in [0, 1], F(0, t) = f(0) = f'(0), F(1, t) = f(1) = f'(1)$

1.  $f \simeq_p f$

We have  $f : I \rightarrow X$  is continuous. Define  $F : I \times I \rightarrow X$  such that  $\forall s, t \in [0, 1], F(s, t) = f(s)$ . Clearly,  $F$  is continuous,  $\forall s \in [0, 1], F(s, 0) = f(s), F(s, 1) = f(s)$  and  $\forall t \in [0, 1], F(0, t) = f(0), F(1, t) = f(1)$ . Thus,  $\exists F : I \times I \rightarrow X \implies f \simeq_p f$ .

2.  $f \simeq_p f' \implies f' \simeq_p f$

We have,  $f \simeq_p f'$ . Thus there exists a continuous function  $F : I \times I \rightarrow X$  such that  $\forall s \in [0, 1], F(s, 0) = f(s), F(s, 1) = f'(s)$  and  $\forall t \in [0, 1], F(0, t) = f(0) = f'(0), F(1, t) = f(1) = f'(1)$ .

Consider  $F' : I \times I \rightarrow X$  defined by  $F'(s, t) = F(s, 1 - t)$ . Clearly,  $F'$  is continuous. And  $F'(s, 0) = F(s, 1) = f'(s)$ , and  $F'(s, 1) = F(s, 0) = f(s)$ . Also,  $F'(0, t) = F(0, 1 - t) = f(0) = f'(0)$  and  $F'(1, t) = F(1, 1 - t) = f(1) = f'(1)$ . Thus,  $\exists F'(s, t) : I \times I \rightarrow X \implies f' \simeq_p f$

3.  $f \simeq f', f' \simeq f'' \implies f \simeq f''$

We have,  $f \simeq f' \iff \exists F : I \times I \rightarrow X$  such that  $F$  is continuous,  $\forall s \in [0, 1], F(s, 0) = f(s), F(s, 1) = f'(s)$  and  $\forall t \in [0, 1], F(0, t) = f(0) = f'(0), F(1, t) = f(1) = f'(1)$

Similarly,  $f' \simeq f'' \iff \exists F' : I \times I \rightarrow X$  such that  $F'$  is continuous,  $\forall s \in [0, 1], F'(s, 0) = f'(s), F'(s, 1) = f''(s)$  and  $\forall t \in [0, 1], F'(0, t) = f'(0) = f''(0), F'(1, t) = f'(1) = f''(1)$

Consider  $G : I \times I \rightarrow X$  defined by

$$G(s, t) = \begin{cases} F(s, 2t) & , t \in [0, \frac{1}{2}] \\ F'(s, 2t - 1) & , t \in [\frac{1}{2}, 1] \end{cases}$$

We have,  $G(s, \frac{1}{2}) = F(s, 1) = F'(s, 0) = f'(s)$ . Thus,  $G$  is continuous by pasting lemma[Munkres, 2003, §18.3 pp. 106], since  $[0, \frac{1}{2}] \cap [\frac{1}{2}, 1] = \{\frac{1}{2}\}$ .

Also  $G(s, 0) = F(s, 0) = f(s)$  and  $G(s, 1) = F'(s, 1) = f''(s)$ .

Again,  $\forall t \in [0, \frac{1}{2}], G(0, t) = F(0, 2t) = f(0) = f'(0) = f''(0)$  and  $\forall t \in [\frac{1}{2}, 1], G(0, t) = F'(0, 2t - 1) = f(0) = f'(0) = f''(0)$ . Therefore,  $\forall t \in [0, 1], G(0, t) = f(0) = f''(0)$ .

Similarly,  $\forall t \in [0, \frac{1}{2}]$ ,  $G(1, t) = F(1, 2t) = f(1) = f'(1) = f''(1)$  and  $\forall t \in [\frac{1}{2}, 1]$ ,  $G(1, t) = F'(1, 2t - 1) = f(1) = f'(1) = f''(1)$ . Therefore,  $\forall t \in [0, 1]$ ,  $G(1, t) = f(1) = f''(1)$ . Thus,  $\exists G : I \times I \rightarrow X \implies f \simeq_p f''$

□

**Definitions 7.6.3.** Let  $f$  be a path in  $X$  (ie,  $f : I \rightarrow X$ ), then  $[f]$  is the equivalence class of all paths homotopic to  $f$  in  $X$ . (ie,  $g \in [f] \iff f \simeq_p g$ )

*Remark.* The set of homotopy classes of functions from  $X$  into  $Y$  is denoted by  $[X, Y]$ . And, the set of all path-homotopic classes on  $X$  is denoted by  $[I, X]$ .

*Remark* (Straight-line homotopy). [Munkres, 2003, §51 Example 1 pp. 320] Let  $X$  be a topological space, and  $f, g$  be continuous functions from  $X$  into a euclidean space, say  $\mathbb{R}^2$ . Then  $f, g$  are straight line homotopic if there exists a continuous function  $F$  from  $X \times I$  such that  $F$  deforms  $f$  into  $g$  along straight line segments joining them.

For example,  $F(x, t) = (1 - t)f(x) + tg(x)$ .

*Remark.* Let  $A$  be a convex subspace of  $\mathbb{R}^n$ . Then any two paths in  $A$  from  $x_0$  to  $x_1$  are path homotopic in  $A$ .

*Proof.* —continue page 321—

□

*Remark.* [Munkres, 2003, §51 Example 2 pp. 321] **This demonstrates that the straight-line homotopy is very sensitive to the holes in the space.**

**Definitions 7.6.4.** Let  $f, g$  be two paths in  $X$  (ie,  $f : I \rightarrow X$  and  $g : I \rightarrow X$ ) such that  $f(0) = x_0, f(1) = g(0) = x_1$  and  $g(1) = x_2$ . Then the product  $h = f * g$  is given by  $h : I \rightarrow X$  and

$$h(s) = \begin{cases} f(2s) & , s \in [0, \frac{1}{2}] \\ g(2s - 1) & , s \in [\frac{1}{2}, 1] \end{cases}$$

This  $h$  is well-defined, and continuous by pasting lemma.<sup>(★<sup>32</sup>)</sup>

**Definitions 7.6.5.** The product operation on path-homotopy classes is defined by  $[f] * [g] = [f * g]$ .

*Remark.* The product of path-homotopic classes is well-defined.

*Proof.* Let  $F$  be a path-homotopy between  $f, f' \in [f]$  and  $G$  be a path-homotopy between  $g, g' \in [g]$ . Then  $H : I \times I \rightarrow X$  defined by

$$H(s, t) = \begin{cases} F(2s, t) & s \in [0, \frac{1}{2}] \\ G(2s - 1, t) & s \in [\frac{1}{2}, 1] \end{cases}$$

Then  $H$  is well-defined, and continuous by pasting lemma.

---

<sup>32</sup>Pasting Lemma : Let  $X = A \cup B$ , where  $A$  and  $B$  are closed in  $A$ . Let  $f : A \rightarrow Y$  and  $g : B \rightarrow Y$  be continuous. If  $f(x) = g(x)$  for every  $x \in A \cap B$ , then  $f$  and  $g$  combine to give a continuous function  $h : X \rightarrow Y$ , defined by setting  $h(x) = f(x)$  if  $x \in A$  and  $h(x) = g(x)$  if  $x \in B$ .

$\forall s \in [0, \frac{1}{2}], H(s, 0) = F(2s, 0) = f(2s)$  and  
 $\forall s \in [\frac{1}{2}, 1], H(s, 0) = G(2s - 1, 0) = g(2s - 1)$ .  
 $\implies H(s, 0) = (f * g)(s)$ , by the definition of  $f * g$

$\forall s \in [0, \frac{1}{2}], H(s, 1) = F(2s, 1) = f'(2s)$  and  
 $\forall s \in [\frac{1}{2}, 1], H(s, 1) = G(2s - 1, 1) = g'(2s - 1)$ .  
 $\implies H(s, 1) = (f' * g')(s)$ , by the definition of  $f' * g'$

$H(0, t) = F(0, t) = f(0) = x_0 = (f * g)(0)$ , and  
 $H(1, t) = G(1, t) = g'(1) = x_2 = (f' * g')(1)$

Then  $H : I \times I \rightarrow X$  is a path-homotopy between  $f * g$  and  $f' * g'$ .  $\square$

**Definitions 7.6.6** (Groupoid). Let  $G$  be a set and  $*$  be a binary operation on  $G$ . Then  $(G, *)$  is a groupoid if it satisfies the following axioms

- g1 Associativity -  $\forall x, y, z \in G, (x * y) * z = x * (y * z)$
- g2 Existence of left and right identities - There exist unique elements  $e_L$  and  $e_R$  such that  $\forall x \in G, x * e_R = x$  and  $e_L * x = x$ .
- g3 Existence of inverse  
 $\forall x \in G, \exists x^{-1} \in G$  such that  $x * x^{-1} = e_L$  and  $x^{-1} * x = e_R$

**Definitions 7.6.7** (Positive Linear Map). A positive linear map  $p : [a, b] \rightarrow [c, d]$  is the unique map of the form  $p(x) = mx + k$  such that  $p(a) = c$  and  $p(b) = d$ . Clearly, scaling factor,  $m = \frac{d-c}{b-a}$  as we want to transform an interval of length  $b - a$  into an interval of length  $d - c$ . And offset  $k$  is given by,

$$p(a) = \frac{d-c}{b-a}a + k = c \implies k = c - \frac{a(d-c)}{b-a} = \frac{bc-ad}{b-a}$$

But, we won't fix  $m$  and  $k$  in  $p(x) = mx + k$ , instead we will focus on the unique map with graph of positive slope and passing through required end points. The graph of a positive linear map from  $[a, b]$  to  $[c, d]$  is always a straight-line with positive slope.

*Remark.* The inverse of a positive linear map is also a positive linear map. Given  $p : [a, b] \rightarrow [c, d], p(x) = mx + k$ , where  $m = \frac{d-c}{b-a}$ ,  $k = \frac{bc-ad}{b-a}$ . Then its inverse,  $\bar{p} : [c, d] \rightarrow [a, b]$  is given by  $\bar{p}(y) = m'y + k'$ , where  $m' = \frac{b-a}{d-c} = \frac{1}{m}$ ,  $k' = \frac{ad-bc}{d-c} = \frac{-k}{m}$ . Clearly  $m > 0 \implies m' = \frac{1}{m} > 0$ .

*Remark.* The composite of two positive linear maps is also a (piece-wise) positive linear map. Let  $f, g$  be two positive linear maps. Then their composite map  $f * g$  is given by

$$(f * g)(x) = \begin{cases} f(2x) & x \in [a, \frac{a+b}{2}] \\ g(2(x - \frac{b-a}{2})) & x \in [\frac{a+b}{2}, b] \end{cases}$$

Remember  $f * g$  exists only if  $f(b) = g(a)$ . Therefore,  $f * g$  is a well-defined, continuous (by pasting lemma) and (piecewise) positive linear map.



**Lemma 7.6.1.** *Let  $f, f'$  be two paths in  $X$  and  $k : X \rightarrow Y$  be a continuous function. Let  $F$  be the path homotopy in  $X$  between the paths  $f$  and  $f'$ . Then  $k \circ F$  is a path homotopy in  $Y$  between the paths  $k \circ f$  and  $k \circ f'$ . That is, path homotopy is preserved under a continuous function.*

**Lemma 7.6.2.** *Let  $f, g$  be two paths in  $X$  with  $f(1) = g(0)$  and  $k : X \rightarrow Y$  be a continuous function. Then  $k \circ (f * g) = (k * f) \circ (k * g)$*

**Theorem 7.6.3.** *Let  $f, g, h$  be paths in a topological space  $X$ , and  $[f], [g], [h]$  be respective path-homotopy classes. Suppose the operation product,  $*$  is defined by*

$$[f] * [g] = [f * g] \text{ where } (f * g)(s) = \begin{cases} f(2s) & s \in [0, \frac{1}{2}] \\ g(2s - 1) & s \in [\frac{1}{2}, 1] \end{cases}$$

*Then the product,  $*$  has the following properties :*

1. *Associativity*

$$\forall [f], [g], [h] \in [I, X], ([f] * [g]) * [h] = [f] * ([g] * [h])$$

2. *Existence of left and right identities*

*Let  $e_x : I \rightarrow X$  defined by  $\forall s \in [0, 1], e_x(s) = x$ . Let  $f$  be a path from  $x_0$  to  $x_1$ , then there exist unique paths  $e_{x_0}$  and  $e_{x_1}$  such that  $[f] * [e_{x_1}] = [f]$  and  $[e_{x_0}] * [f] = [f]$ . That is,  $e_{x_0}, e_{x_1}$  are respectively the left and right path-homotopy-identities.*

3. *Existence of inverse*

*Let  $f$  be a path in  $X$ . The path,  $\bar{f}$ , defined by  $\bar{f}(s) = f(1 - s)$  is the reverse path of  $f$ . Then  $[f] * [\bar{f}] = [e_{x_0}]$  and  $[\bar{f}] * [f] = [e_{x_1}]$ . That is, the inverse of class of  $f$  is the class of reverse path of  $f$ .*

*In other words, Set of all path-homotopy classes together with binary operation product,  $*$  is a groupoid. ie,  $([I, X], *)$  is a groupoid.*

*Proof.* Step 1 : Properties 2&3

Let  $e_0 : I \rightarrow I$  such that  $e_0(t) = 0, \forall t \in I$ . And  $i : I \rightarrow I$  such that  $i(t) = t, \forall t \in I$ . Then  $e_0 * i$  is also a path. Since  $I$  is convex, there is a path homotopy ( $\star^{33}$ )  $G$  between  $i$  and  $e_0 * i$ . Let  $f : I \rightarrow X$  be continuous path in  $X$  from  $x_0$  to  $x_1$ . Then  $f \circ G$  is a path homotopy (by Lemma 2) in  $X$  between  $f \circ i$  and  $f \circ e_0 * i$  where  $f \circ i$  and  $f \circ e_0$  are paths from  $x_0$  to  $x_1$  in  $X$ .

$$f \circ (e_0 * i) = (f \circ e_0) * (f \circ i), \text{ by Lemma 1}$$

$$= e_{x_0} * f, \text{ since } \forall s \in I, f(e_0(s)) = x_0 = e_{x_0}(s) \text{ and } f * i \simeq_p f$$

Therefore  $[e_{x_0}] * [f] \simeq_p [f]$ , since  $e_0 * i \simeq_p i$ , and  $f \circ (e_0 * i) \simeq_p f \circ i = f$ .

Similarly,  $e_1 : I \rightarrow I$  such that  $e_1(t) = 1$ . Let  $H$  be a path homotopy ( $\star^{34}$ ) between  $i * e_1$  and  $i$ . Thus,  $f \circ H$  is a path homotopy in  $X$  from  $f \circ (i * e_1)$  and  $f \circ i$ .

$$f \circ (i * e_1) = (f \circ i) * (f \circ e_1), \text{ by Lemma 1}$$

$$= f * e_{x_1}, \text{ since } f * i \simeq_p f, i * e_1 \simeq_p e_1, \forall s \in I, (f(e_1(s))) = x_1 = e_{x_1}(s)$$

<sup>33</sup> $G : I \times I \rightarrow I, G(s, 0) = i(s), G(s, 1) = (e_0 * i)(s), G(0, t) = 0, G(1, t) = 1.$

<sup>34</sup> $H : I \times I \rightarrow I, H(s, 0) = (i * e_1)(s), H(s, 1) = i(s), H(0, t) = 0, H(1, t) = 1$

Since  $i * e_1 \simeq_p i$ , we have  $f \circ (i * e_1) \simeq_p f \circ i = f$ . Therefore  $[f] * [e_{x_1}] \simeq_p [f]$ . Thus,  $[f] * [e_{x_1}] \simeq_p [f] \simeq_p [e_{x_0}] * [f]$ . Therefore,  $[f]$  has left and right inverses ie, property 2 holds.

Consider inverse path  $\bar{i} : I \rightarrow I$ ,  $\bar{i}(s) = 1 - s$ . Then  $i * \bar{i}$  is a path in  $I$  with both end points at 0. We have,  $e_0 : I \rightarrow I$ ,  $e_0(s) = 0$  is also a path with both end points at 0. Since  $I$  is convex, there is a path homotopy  $H$  in  $I$  between  $e_0$  and  $i * \bar{i}$ . Then  $f \circ H$  is a path homotopy between  $f \circ e_0 = e_{x_0}$  and  $f \circ (i * \bar{i}) = (f \circ i) * (f \circ \bar{i}) = f * \bar{f}$ . Therefore,  $[e_{x_0}] \simeq_p [f] * [\bar{f}]$ .

Similarly  $\bar{i} * i$  and  $e_1$  are paths with both end points at 1. Since  $I$  is convex, there is a path homotopy  $G$  in  $I$  between  $\bar{i} * i$  and  $e_1$ . Then  $f \circ G$  is a path homotopy between  $f \circ (\bar{i} * i) = (f \circ \bar{i}) * (f \circ i) = \bar{f} * f$  and  $f \circ e_1 = e_{x_1}$ . Therefore,  $[\bar{f}] * [f] \simeq_p [e_{x_1}]$ . Thus the path  $\bar{f} : I \rightarrow X$ ,  $\bar{f}(s) = f(1 - s)$ ,  $\forall s \in I$  is reverse of  $f$ . Also  $[f] * [\bar{f}] = [e_{x_0}]$  and  $[\bar{f}] * [f] = [e_{x_1}]$ . ie, property 3 holds.

Step 2 : Property 1

Let  $f, g, h$  be three paths in  $X$  and  $f(1) = g(0) = x_1$  and  $g(1) = h(0) = x_2$ . Then  $f * (g * h)$  is defined by

$$(g * h)(s) = \begin{cases} g(2s) & s \in [0, \frac{1}{2}] \\ h(2s - 1) & s \in [\frac{1}{2}, 1] \end{cases}$$

$$(f * (g * h))(s) = \begin{cases} f(2s) & s \in [0, \frac{1}{2}] \\ (g * h)(2s - 1) & s \in [\frac{1}{2}, 1] \end{cases}$$

$$= \begin{cases} f(2s) & s \in [0, \frac{1}{2}] \\ g(2(2s - 1)) & s \in [\frac{1}{2}, \frac{3}{4}] \\ h(2(2s - 1) - 1) & s \in [\frac{3}{4}, 1] \end{cases}$$

Similarly,  $(f * g) * h$  is defined by,

$$(f * g)(s) = \begin{cases} f(2s) & s \in [0, \frac{1}{2}] \\ g(2s - 1) & s \in [\frac{1}{2}, 1] \end{cases}$$

$$((f * g) * h)(s) = \begin{cases} (f * g)(2s) & s \in [0, \frac{1}{2}] \\ h(2s - 1) & s \in [\frac{1}{2}, 1] \end{cases}$$

$$= \begin{cases} f(2(2s)) & s \in [0, \frac{1}{4}] \\ g(2(2s - 1)) & s \in [\frac{1}{4}, \frac{1}{2}] \\ h(2s - 1) & s \in [\frac{1}{2}, \frac{3}{4}] \end{cases}$$

Clearly,  $f * (g * h)$  and  $(f * g) * h$  are distinct path with common endpoints. ie,  $(f * (g * h))(0) = f(0) = ((f * g) * h)(0)$ . And  $(f * (g * h))(1) = h(1) = ((f * g) * h)(1)$ .

We need to define a path homotopy  $G$  between  $f * (g * h)$  and  $(f * g) * h$ . Let  $[a, b], [c, d] \subset I$ . Consider the path  $p : I \rightarrow I$  defined by the following

three unique<sup>(35)</sup> positive linear maps  $p_1 : [0, a] \rightarrow [0, c]$ ,  $p_2 : [a, b] \rightarrow [c, d]$  and  $p_3 : [b, 1] \rightarrow [d, 1]$ .

$$p(t) = \begin{cases} p_1(t) & t \in [0, a] \\ p_2(t) & t \in [a, b] \\ p_3(t) & t \in [b, 1] \end{cases}$$

We can easily construct, a path homotopy  $P$  between identity map  $i : I \rightarrow I$ ,  $i(s) = s$  and  $p$  as follows :

$$P(s, t) = \begin{cases} t + (p_1(t) - t) \frac{s}{a} & s \in [0, a] \\ t + (p_2(t) - t) \frac{(s-a)}{(b-a)} & s \in [a, b] \\ t + (p_3(t) - t) \frac{(s-b)}{(1-b)} & s \in [b, 1] \end{cases}$$

Therefore, we have  $f * (g * h) \simeq_p i$  since there exists a path homotopy  $P$  corresponding to  $[a, b] = [\frac{1}{2}, \frac{3}{4}]$  and  $[c, d] = [x_1, x_2]$ . Similarly  $(f * g) * h \simeq_p i$  since there exists a path homotopy  $P$  where  $[a, b] = [\frac{1}{4}, \frac{1}{2}]$  and  $[c, d] = [x_1, x_2]$ . ie,  $[f * (g * h)] \simeq_p [(f * g) * h]$ .  $\square$

**Theorem 7.6.4.** *Let  $f$  be a path in  $X$ , and  $a_0, a_1, \dots, a_n$  be numbers such that  $0 = a_0 < a_1 < \dots < a_n = 1$ . Let  $f_i : I \rightarrow X$  be the path that equals the positive linear map of  $I$  onto  $[a_{i-1}, a_i]$  followed by  $f$ . Then  $[f] = [f_1] * [f_2] * \dots * [f_n]$ . In other words, every path is path-homotopic to a piecewise-linear path.*

*Proof.* Let  $f$  be a piece-wise positive linear map such that

$$f(t) = \begin{cases} f_1(t) & t \in [0 = a_0, a_1] \\ f_2(t) & t \in [a_1, a_2] \\ \vdots & \vdots \\ f_n(t) & t \in [a_{n-1}, a_n] \end{cases}$$

where  $f_i : I \rightarrow [a_{i-1}, a_i]$  such that  $f_i(t)$  are a positive linear maps.

Consider the path  $p : I \rightarrow I$  defined by the unique positive linear maps on the subintervals of any partition  $\{0 = x_0, x_1, \dots, x_n\}$  of  $I$ . ie,  $0 = x_0 < x_1 < \dots < x_n = 1$

$$\begin{aligned} p_1 : [x_0, x_1] &\rightarrow [a_0, a_1] \\ p_2 : [x_1, x_2] &\rightarrow [a_1, a_2] \\ &\vdots \\ p_n : [x_{n-1}, x_n] &\rightarrow [a_{n-1}, a_n] \end{aligned}$$

$$\text{Define, } p(t) = \begin{cases} p_1(t) & t \in [x_0, x_1] \\ p_2(t) & t \in [x_1, x_2] \\ \vdots & \vdots \\ p_n(t) & t \in [x_{n-1}, x_n] \end{cases}$$

---

<sup>35</sup> $p_1(t) = \frac{ct}{a}$ ,  $p_2(t) = \frac{(d-c)t}{b-a} + \frac{bc-ad}{b-a}$ ,  $p_3(t) = \frac{(1-d)t}{1-b} + \frac{d-b}{1-b}$

Then there exists a path homotopy  $P$  between identity map  $i : I \rightarrow I$ ,  $i(t) = t$  and  $p$  given by

$$P(s, t) = \begin{cases} t + (p_1(t) - t) \frac{a_1}{x_1} & s \in [0, x_1] \\ t + (p_2(t) - t) \frac{s - x_1}{x_2 - x_1} & s \in [x_1, x_2] \\ \vdots & \\ t + (p_n(t) - t) \frac{s - x_{n-1}}{x_n - x_{n-1}} & s \in [x_{n-1}, x_n] \end{cases}$$

Since any product of  $f_1, f_2, \dots, f_n$  is a path  $p$  for some partition decided by the order of associativity. This partition can be constructed as follows : Let the last product operation (by associativity) corresponds to  $\frac{1}{2}$ . The expression on its left corresponds to  $[0, \frac{1}{2}]$  and expression on the right corresponds to  $[\frac{1}{2}, 1]$ . If there are any operations on any of these parts, the last operation (by associativity) in them corresponds to the midpoint the respective subinterval and so on.

For examples : Consider,  $(f_1 * (f_2 * f_3)) * (f_4 * f_5)$ . Suppose we number the operations,  $(f_1 *_1 (f_2 *_2 f_3)) *_3 (f_4 *_4 f_5)$ . Then we have,  $*_3 \rightarrow \frac{1}{2} \implies *_1 \rightarrow \frac{1}{4} \implies *_2 \rightarrow \frac{3}{8}$ . Again  $*_3 \rightarrow \frac{1}{2} \implies *_4 \rightarrow \frac{3}{4}$ . Thus, we have  $\{0, \frac{1}{4}, \frac{3}{8}, \frac{1}{2}, \frac{3}{4}, 1\}$ .

Thus given two paths  $f$  and  $f'$  with distinct order of associativity of  $n$  paths :  $f_1, f_2, \dots, f_n$ . We have path homotopy  $P, P'$  given by the  $P(s, t)$  for the respective partition constructed according to the order of associativity. Then, we have  $f \simeq_p i$  and  $f' \simeq_p i$ . Thus, irrespective of the order of associativity all these paths are path homotopic. ie,  $[f] = [f_1] * [f_2] \cdots [f_n]$   $\square$

## Subject 8

# ME010203 Numerical Analysis with Python

### 8.1 Interpolation and Curve Fitting

**Definitions 8.1.1.** Given  $(n + 1)$  data points  $(x_k, y_k)$ ,  $k = 0, 1, \dots, n$ , the problem of estimating  $y(x)$  using a function  $y : \mathbb{R} \rightarrow \mathbb{R}$  that satisfy the data points is the interpolation problem. ie,  $y(x_k) = y_k$ ,  $k = 0, 1, \dots, n$ .

**Definitions 8.1.2.** Given  $(n + 1)$  data points  $(x_k, y_k)$ ,  $k = 0, 1, \dots, n$ , the problem of estimating  $y(x)$  using a function  $y : \mathbb{R} \rightarrow \mathbb{R}$  that is sufficiently close to the data points is the curve-fitting problem.  
ie, Given  $\epsilon > 0$ ,  $|y(x_k) - y_k| < \epsilon$ ,  $k = 0, 1, \dots, n$ .

*Remark.* The data could be from scientific experiments or computations on mathematical models. The interpolation problem assumes that the data is accurate. But, curve-fitting problem assumes that there are some errors involved which are sufficiently small.

**Definitions 8.1.3.** Given  $(n + 1)$  data points  $(x_k, y_k)$ ,  $k = 0, 1, \dots, n$ , the problem of estimating  $y(x)$  using a polynomial function of degree  $n$  that satisfy the data points is the polynomial interpolation problem.

*Remark.* Polynomial is the ‘simplest’ interpolant. [Kiusalaas, 2013, 3.2]

### 8.2 Polynomial Interpolation

There exists a unique polynomial of degree  $n$  that satisfy  $(n + 1)$  distinct data points. There are a few methods to find this polynomial : 1. Lagrange’s method  
2. Newton’s method.

### 8.2.1 Lagrange's Method

Interpolation polynomial <sup>1</sup> is given by,

$$P(x) = \sum_{i=0}^n y_i l_i(x), \text{ where } l_i(x) = \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j} \quad (8.1)$$

*Remark.* Lagrange's cardinal functions  $l_i$ , are polynomials of degree  $n$  and

$$l_i(x_j) = \delta_{ij} = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}$$

**Proposition 8.2.1.** Error in polynomial interpolation is given by

$$f(x) - P(x) = \frac{(x - x_0)(x - x_1) \cdots (x - x_n)}{(n + 1)!} f^{(n+1)}(\xi) \quad (8.2)$$

where  $\xi \in (x_0, x_n)$

*Remark.* The error increases as  $x$  moves away from the unknown value  $\xi$ .

### 8.2.2 Newton's Method

The interpolation polynomial is given by,

$$P(x) = a_0 + a_1(x - x_0) + \cdots + a_n(x - x_0)(x - x_1) \cdots (x - x_{n-1}) \quad (8.3)$$

where  $a_i = \nabla^i y_i$ ,  $i = 0, 1, \dots, n$ .

*Remark.* For Newton's Method, usually it is assumed that  $x_0 < x_1 < \cdots < x_n$ .

*Remark.* Lagrange's method is conceptually simple. But, Newton's method is computationally more efficient than Lagrange's method.

#### Computing coefficients $a_i$ of the polynomial

The coefficients are given by,

$$a_0 = y_0, \quad a_1 = \nabla y_1, \quad a_2 = \nabla^2 y_2, \quad a_3 = \nabla^3 y_3, \dots, a_n = \nabla^n y_n \quad (8.4)$$

*Remark.* The divided difference  $\nabla^i y_i$  are computed as follows:

$$\begin{aligned} \nabla y_1 &= \frac{y_1 - y_0}{x_1 - x_0} \\ \nabla y_2 &= \frac{y_2 - y_1}{x_2 - x_1} & \nabla^2 y_2 &= \frac{\nabla y_2 - \nabla y_1}{x_2 - x_1} \\ \nabla y_3 &= \frac{y_3 - y_2}{x_3 - x_2} & \nabla^2 y_3 &= \frac{\nabla y_3 - \nabla y_2}{x_3 - x_2} & \nabla^3 y_3 &= \frac{\nabla^2 y_3 - \nabla^2 y_2}{x_3 - x_2} \end{aligned}$$

*Remark.* Practise Problems

Find interpolation polynomial for the following data points :

---

<sup>1</sup>Using  $P_n$  to represent some polynomial of degree  $n$ . It is quite a confusing a notation when it comes to Newton's method as author construct a psuedo-recursive definition.

1.  $\{(0, 7), (2, 11), (3, 28)\}$   
 $\text{Ans : } 5x^2 - 8x + 7$   
 [Kiusalaas, 2013, Example 3.1]
2.  $\{(-2, -1), (1, 2), (4, 59), (-1, 4), (3, 24), (-4, -53)\}$   
 $\text{Ans : } x^3 - 2x + 3$   
 [Kiusalaas, 2013, Example 3.2]
3.  $\{(-1.2, -5.76), (0.3, -5.61), (1.1, -3.69)\}$   
 $\text{Ans : } x^2 + x - 6$   
 [Kiusalaas, 2013, Problem Set 3.1.1]
4.  $\{(-3, 0), (2, 5), (-1, -4), (3, 12), (1, 0)\}$   
 $\text{Ans : } x^2 + 2x - 3$   
 [Kiusalaas, 2013, Problem Set 3.1.7]
5.  $\{(0, 1.225), (3, 0.905), (6, 0.652)\}$   
 $\text{Ans : } 0.0037x^2 - 0.1178x + 1.225$   
 [Kiusalaas, 2013, Problem Set 3.1.9]

*Remark.* In Lagrange's Method, we can interpolate at the given point even without computing the polynomial. In Newton's method, we have to compute polynomial and then interpolate for the given point.

That is, evaluate the value of cardinal polynomials at the point and substitute in Equation 8.1 as shown in Section 3.2. [Kiusalaas, 2013, Example 3.1]

### 8.2.3 Implementation of Newton's Method

#### Program 8.2.1. Computing Coefficients

```
def coefficients(xData, yData):
    m = len(xData)
    a = yData.copy()
    for k in range(1, m):
        a[k:m] = (a[k:m] - a[k-1]) / (xData[k:m] - xData[k-1])
    return a
```

Line 1 `def coefficients(xData, yData):`

Defines a function which takes two arguments/parameters, named *xData* and *yData*. In [Kiusalaas, 2013, 3.2], you will find coeffs which I have changed to coefficients. *xData, yData* are numpy array objects. *xData*

$x_0$	$y_0$				
$x_1$	$y_1$	$\nabla y_1$			
$x_2$	$y_2$	$\nabla y_2$	$\nabla^2 y_2$		
$\dots$	$\dots$	$\dots$	$\dots$	$\ddots$	
$x_n$	$y_n$	$\nabla y_n$	$\nabla^2 y_n$	$\dots$	$\nabla^n y_n$

Table 8.1: The  $\nabla^i y_i$  Computation Table

is a array with values  $x_0, x_1, \dots, x_n$ . And  $yData$  is array with values  $y_0, y_1, \dots, y_n$ . For example, the value of  $x_3$  can be accessed as  $xData[3]$ .

Line 2 `m = len(xData)`

The function `len()` is extended by numpy to give the length of array objects. In this context, `len(xData)` will return the value  $n+1$ , since there are  $n+1$  values in  $xData$  array.

Line 3 `a = yData.copy()`

We need a copy of  $yData$  to work with. Unlike other programming languages like java, in python `a = yData` will assign a new label  $a$  to the same memory location and manipulating  $a$  will corrupt the original data in  $yData$  as well. In order to avoid this, we are **making a copy of the array object using the array method provided by the numpy library.**

Line 4 `for k in range(1,m):`

This is a python loop statement. This ask python interpreter to repeat the following sub-block  $m - 1$  times.<sup>2</sup> In this context, Line 5 will be executed  $n$  times, since the `range(1,m)` object is a list-type object with values  $1, 2, \dots, m - 1$ . And interpreter executes Line 5 for each values in the `range()` object, ie,  $k = 1, 2, \dots, m - 1$  before interpreting Line 6.

Line 5 `a[k:m] = (a[k:m]-a[k-1])/(xData[k:m]-xData[k-1])`

This is very nice feature available in python. **This statement, evaluates  $m - k$  values in a single step. ie,  $a[k], a[k + 1], \dots, a[m]$ . This calculation corresponds to subsequent columns of the divided difference table, that we are familiar with.** For example, executing Line 5 with  $k = 3$  is same as evaluating the  $\nabla^3 y_j$  column. Note that the value  $a[0]$  is never updated and similarly  $a[2]$  changes when Line 5 is executed with  $k = 1, 2$ . From column 3 onward,  $a[2]$  is not updated. Therefore, **after completing  $n$ th executing of the Line 5, we have**  $a[0] = y_0$ ,  $a[1] = \nabla y_1$ ,  $a[2] = \nabla^2 y_2, \dots$ ,  $a[n] = \nabla^n y_n$ .

Line 6 `return a`

This returns the array  $a$  which is the array of coefficients.

The logic of this program is in Line 4 and Line 5. So they need more explanation/understanding than anything else.

#### Program 8.2.2. Interpolating using Newton's Method

```
def interpolate(a, xData, x):
    n = len(xData)-1
    p = a[n]
    for k in range(1, n+1):
        p = a[n-k]*(x-xData[n-k])*p
    return p
```

The logic this program is in Line 3, Line 4 and Line 5.

<sup>2</sup>Python block is a group of statement with same level of indentation. A sub-block is a block with an additional indentation.



Line 3 : We initialize the polynomial with the coefficient  $a[n] = \nabla^n y_n = a_n$ .

Line 4 : We are going to define the polynomial recursively. This takes exactly  $n$  steps further. So we use a loop which repeats  $n$  times.

Line 5 : The value of  $p$  and  $k$  changes each time Line 5 is executed. Let  $P_j$  be the value in  $p$  after executing Line 5 with  $k = j$ . Then,

$$P_0 = p = a[n]$$

$$P_1 = a[n-1] + (x - x_{n-1})P_0$$

$$P_2 = a[n-2] + (x - x_{n-2})P_1$$

$\vdots$

$P_n = a[0] + (x - x_0)P_{n-1}$ . Clearly,  $P_n$  is the unique  $n$  degree polynomial given by the Newton's method.

**Program 8.2.3.** How to interpolate ?

```
from numpy import array
xDData = array([-2,1,4,-1,3,-4])
yData = array([-1,2,59,4,24,-53])
a = coefficients(xDData,yData)
print(interpolate(a,xDData,2))
```

You will have to define both the functions (coefficients, interpolate) before doing this.

Line 1 `from numpy import array`

For defining array objects, we need to import them from numpy library.

Line 2 `xDData = array([-2,1,4,-1,3,-4])`

You can change this line according to the first component of the given data points.

Line 3 `yData = array([-1,2,59,4,24,-53])`

You can change this line according to the second component of the given data points.

Line 4 `a = coefficients(xDData,yData)`

Call function `coefficients` and store the array returned into `a`

Line 5 `print(interpolate(a,xDData,2))`

Call function `interpolate` to interpolate at  $x = 2$  and print the value  $P(2)$

**Program 8.2.4** (Just for Fun). We can do more using sympy !

```
from numpy import array
from sympy import Symbol
xDData = array([-2,1,4,-1,3,-4])
yData = array([-1,2,59,4,24,-53])
a = coefficients(xDData,yData)
x = Symbol('x')
p = interpolate(a,xDData,x)
p.subs({x:2})
```

*Remark.* Programming Problems

1.  $\{(0.15, 4.79867), (2.30, 4.49013), (3.15, 4.2243), (4, 85, 3.47313), (6.25, 2.66674), (7.95, 1.51909)\}$  [Kiusalaas, 2013, Example 3.4]
2.  $\{(0, -0.7854), (0.5, 0.6529), (1, 1.7390), (1.5, 2.2071), (2, 1.9425)\}$  [Kiusalaas, 2013, Problem Set 3.1.5]

### 8.2.4 Limitations of Polynomial Interpolation

1. Inaccuracy - The error in interpolation increases as the point moves away from most of the data points.
2. Oscillation - As the number of data points considered for polynomial interpolation increases, the degree of the polynomial increases. And the graph of the interpolant tend to oscillate excessively. In such cases, the error in interpolation is quite high.
3. The best practice is to consider four to six data points nearest to the point of interest and ignore the rest of them.

*Remark.* The interpolant obtained by joining cubic polynomials corresponding to four nearest data points each, is a cubic spline<sup>3</sup>.

## 8.3 Roots of a Function

**Definitions 8.3.1.** Let  $f : \mathbb{R} \rightarrow \mathbb{R}$ , then  $x \in \mathbb{R}$  is a root of  $f$  if  $f(x) = 0$ .

*Remark.* Suppose  $a < b$  and  $f(a), f(b)$  are nonzero and are of different signs. If  $f$  is continuous in  $[a, b]$ , then there is a point  $c \in [a, b]$  such that  $f(c) = 0$ .

Thus given  $a < b$  and  $f(a), f(b)$  are nonzero values of different sign, then there may be a bracketed root in  $[a, b]$ .

Note : There is no guarantee that there exists a root in  $[a, b]$  as we are not sure about the continuity of  $f$ .

*Remark.* Given a bracketed root, we can find it using

1. Bisection Method or
2. Newton-Raphson Method

### 8.3.1 Bisection Method

Suppose  $a < b$  and  $f(a), f(b)$  are nonzero values of different signs. We evaluate  $f(c)$  where  $c = \frac{a+b}{2}$ . If  $f(c)$  is a nonzero value, then at least one of the pairs  $f(a), f(c)$  or  $f(c), f(b)$  are of different signs. WLOG suppose that  $f(a), f(c)$  are of different signs, then set  $b = c$  and  $c = \frac{a+b}{2}$ . And continue this process until we get sufficiently accurate value of a root.

---

<sup>3</sup>Cubic spline is a function, the graph of which is piece-wise cubic

*Remark.* Suppose  $f(x) = x^5 - 2$ . Then  $f(0) = -2$ ,  $f(1) = -1$ ,  $f(2) = 30$ . Since  $f$  is known to be continuous, there is a bracketed root in  $[1, 2]$ . Now

$$\begin{aligned} f(1.5) &> 0 \implies [1, 1.5] \\ f(1.25) &> 0 \implies [1, 1.25] \\ f(1.125) &< 0 \implies [1.125, 1.25] \\ f(1.1875) &> 0 \implies [1.125, 1.1875] \\ f(1.15375) &> 0 \implies [1.125, 1.15375] \\ f(1.139375) &< 0 \implies [1.139375, 1.15375] \\ f(1.1465625) &< 0 \implies [1.1465625, 1.15375] \\ f(1.150156250) &> 0 \implies [1.1465625, 1.15015625] \\ f(1.148359375) &< 0 \implies [1.1483594, 1.15015625] \\ f(1.149257825) &> 0 \implies [1.1483594, 1.14925783] \end{aligned}$$

Thus, we have 1.14 is a root of  $f$  with accuracy upto two decimal points.

### 8.3.2 Newton-Raphson Method

Suppose  $f$  is differentiable at  $x \in \mathbb{R}$  and  $f(x) \neq 0$ . Then compute  $x = x - \frac{f(x)}{df(x)}$  and evaluate  $f(x)$ . Repeat this process to get more accurate value of a root near  $x$ .

*Remark.* Suppose  $f(x) = x^5 - 2$ . Then  $df(x) = 5x^4$ . Let  $x = 2$ . Then

$$\begin{aligned} x &= 2 - \frac{30}{80} \implies f(1.625) = 9.330 \\ x &= 1.625 - \frac{9.330}{34.86} \implies f(1.35735) = 2.6074 \\ x &= 1.35735 - \frac{2.6074}{16.9721} \implies f(1.20373) = 0.52733 \\ x &= 1.20373 - \frac{0.52733}{10.4975} \implies f(1.15351) = 0.04224 \\ x &= 1.15351 - \frac{0.042245}{8.85225} \implies f(1.148738) = 0.00034312 \end{aligned}$$

Thus we have 1.1487 is quite close to a root of  $f$ .

## 8.4 Matrix Operations

Consider a system of  $n$  linear, simultaneous equations in  $n$  unknowns,

$$\begin{aligned} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1n}x_n &= b_1 \\ A_{21}x_1 + A_{22}x_2 + \cdots + A_{2n}x_n &= b_2 \\ &\vdots \\ A_{n1}x_1 + A_{n2}x_2 + \cdots + A_{nn}x_n &= b_n \end{aligned}$$

We may represent them using matrices as  $Ax = b$ . That is,

$$\begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

Gauss Elimination and Doolittle Decomposition are two methods to solve system of equations,  $Ax = b$ .

## 8.5 Gauss elimination method

Gauss elimination method has of two phases 1. elimination and 2. back substitution. In elimination phase, system  $Ax = b$  is transformed into an equivalent system  $Ux = c$  where  $U$  is an upper-triangular<sup>4</sup> matrix. And in back substitution phase,  $Ux = c$  is solved. Since  $Ax = b$  and  $Ux = c$  are equivalent, they have the same solution  $x$ .

### 8.5.1 Elimination Phase

We can eliminate unknowns from an equation by adding a scalar multiple of an equation to another equation of the system. In matrices, this is equivalent to adding a scalar multiple of one row to another row, say  $R_i \leftarrow R_i + \lambda R_k$ .

$$\begin{array}{r} A_{k1}x_1 + A_{k2}x_2 + \cdots + A_{kn}x_n = b_k + \\ \lambda(A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{in}x_n = b_i) \\ \hline (A_{k1} + \lambda A_{i1})x_1 + (A_{k2} + \lambda A_{i2})x_2 + \cdots + (A_{kn} + \lambda A_{in})x_n = b_k + \lambda b_i \end{array}$$

$$\begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ A_{k1} & A_{k2} & \cdots & A_{kn} \\ \vdots & \vdots & \ddots & \vdots \\ A_{i1} & A_{i2} & \cdots & A_{in} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \xrightarrow{R_i \leftarrow R_i + \lambda R_k} \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ A_{k1} & A_{k2} & \cdots & A_{kn} \\ \vdots & \vdots & \ddots & \vdots \\ A_{i1} + \lambda A_{k1} & A_{i2} + \lambda A_{k2} & \cdots & A_{in} + \lambda A_{kn} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$$\begin{bmatrix} \vdots \\ b_k \\ \vdots \\ b_i \\ \vdots \end{bmatrix} \xrightarrow{R_i \leftarrow R_i + \lambda R_k} \begin{bmatrix} \vdots \\ b_k \\ \vdots \\ b_i + \lambda b_k \\ \vdots \end{bmatrix}$$

---

<sup>4</sup>upper triangular - all the entries below the main diagonal are zero. ie  $U_{ij} = 0$ , if  $i < j$

### 8.5.2 Back substitution

Let  $Ux = c$  be a system of  $n$  linear equations in  $n$  unknowns and  $U$  is an upper triangular matrix. Then we can solve the system of equations from the back.

$$\begin{bmatrix} U_{1,1} & U_{1,2} & \cdots & U_{1,n-1} & U_{1,n} \\ 0 & U_{2,2} & \cdots & U_{2,n-1} & U_{2,n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & U_{n-1,n-1} & U_{n-1,n} \\ 0 & 0 & \cdots & 0 & U_{n,n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_{n-1} \\ c_n \end{bmatrix}$$

$$U_{n,n} x_n = c_n \implies x_n = \frac{c_n}{U_{n,n}}$$

$$\sum_{i=n-1}^n U_{n-1,i} x_i = c_{n-1} \implies x_{n-1} = \frac{c_{n-1} - U_{n-1,n} x_n}{U_{n-1,n-1}}$$

...

$$\sum_{i=1}^n U_{1,i} x_i = c_1 \implies x_1 = \frac{c_1 - \sum_{i=2}^n U_{1,i} x_i}{U_{1,1}}$$

### 8.5.3 Illustrative example

Consider the following system of linear equations,

$$\begin{aligned} 4x_1 - 2x_2 + x_3 &= 11 \\ -2x_1 + 4x_2 - 2x_3 &= -16 \\ x_1 - 2x_2 + 4x_3 &= 17 \end{aligned}$$

We may represent the above system of linear equations using matrices,

$$\begin{bmatrix} 4 & -2 & 1 \\ -2 & 4 & -2 \\ 1 & -2 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 11 \\ -16 \\ 17 \end{bmatrix}$$

Phase 1 : Elimination Process

Using eq.1, the unknown  $x_1$  is eliminated from all subsequent equations. An equivalent operation can be performed on both the matrices  $A$  and  $b$  by adding a suitable scalar multiples of row  $R_1$  to row  $R_2$  and  $R_3$ .

$$\begin{bmatrix} 4 & -2 & 1 \\ -2 & 4 & -2 \\ 1 & -2 & 4 \end{bmatrix} \xrightarrow{\substack{R_2 \leftarrow R_2 + 0.5R_1 \\ R_3 \leftarrow R_3 - 0.25R_1}} \begin{bmatrix} 4 & -2 & 1 \\ 0 & 3 & -1.5 \\ 0 & -1.5 & 3.75 \end{bmatrix}$$

$$\begin{bmatrix} 11 \\ -16 \\ 17 \end{bmatrix} \xrightarrow{\substack{R_2 \leftarrow R_2 + 0.5R_1 \\ R_3 \leftarrow R_3 - 0.25R_1}} \begin{bmatrix} 11 \\ -10.5 \\ 14.25 \end{bmatrix}$$

And using eq.2,  $x_2$  is eliminated from all subsequent equations( only those rows below it). Again, we perform this by adding suitable scalar multiples of row 2 to row  $R_3$ .

$$\begin{bmatrix} 4 & -2 & 1 \\ 0 & 3 & -1.5 \\ 0 & -1.5 & 3.75 \end{bmatrix} \xrightarrow{R_3 \leftarrow R_3 + 0.5R_2} \begin{bmatrix} 4 & -2 & 1 \\ 0 & 3 & -1.5 \\ 0 & 0 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 11 \\ -16 \\ 17 \end{bmatrix} \xrightarrow{R_3 \leftarrow R_3 + 0.5R_2} \begin{bmatrix} 11 \\ -10.5 \\ 9 \end{bmatrix}$$

The elimination process is complete when all entries below the diagonal elements are reduced to zero. ie, upper triangular matrix.

Phase 2 : Back substitution Process

The unknowns are easily found from the equations by solving them in the reverse order. The unknowns are solved from the bottom and solved variables are used to solve the remain unknowns.

$$\begin{bmatrix} 4 & -2 & 1 & 11 \\ 0 & 3 & -1.5 & -10.5 \\ 0 & 0 & 3 & 9 \end{bmatrix} \rightarrow \begin{cases} 4x_1 - 2x_2 + x_3 = 11 \\ 3x_2 - 1.5x_3 = -10.5 \\ 3x_3 = 9 \end{cases}$$

$$x_3 = \frac{9}{3} = 3$$

$$x_2 = \frac{-10.5 + 1.5x_3}{3} = -2$$

$$x_1 = \frac{11 - x_3 + 2x_2}{4} = 1$$

*Remark.* Why don't they use row-reduced echelon matrix of  $A$  to simplify the back substitution phase ?

This doesn't have much advantage from algorithmic point of view. That is, the time complexity ( number of steps for computation) is unaffected. And algorithms always prefer methods even with slight advantage in time or memory. And they won't consider complications in the manual execution of the method. Therefore, programmers won't consider alternate algorithm for the sake of computational simplicity.

### 8.5.4 Python : Gauss elimination method

**Program 8.5.1** (Gauss elimination).

```
from numpy import dot
def gaussElimination(a,b):
    n = len(b)
    for k in range(0,n-1):
        for i in range(k+1,n):
            if a[i,k] != 0.0:
                lam = a[i,k]/a[k,k]
                a[i,k+1:n] = a[i,k+1:n]-lam*a[k,k+1:n]
                b[i] = b[i]-lam*b[k]
        for k in range(n-1,-1,-1):
            x[k] = (b[k]-dot(a[k,k+1:n],x[k+1:n]))/a[k,k]
```

```
return b
```

- Line 1 `from numpy import dot`  
Imports the “dot()” function for numpy arrays which takes two ‘numpy arrays’ as input arguments and returns the dot product of them.
- Line 2 `def gaussElimination(a,b):`  
it defines “gaussElimination()” as a function which takes two arguments (inputs). First argument is the coefficient matrix  $A$  and second argument is the constant matrix  $b$  of the linear system of the form  $Ax = b$ .
- Line 3 `n = len(b)`  
it assigns the length of the list  $b$  into variable  $n$  which is obviously the number of equations.
- Line 4 `for k in range(0,n-1):`  
it is a loop construct. Five instructions following it are part of this loop body, which are executed for each values of  $k$  ie,  $k = 0, 1, \dots, n-1$ . For each value of  $k$ , the unknown  $x_{k+1}$  is selected for elimination process.
- Line 5 `for i in range(k+1,n):`  
it is a loop inside another loop. Four instructions following it are part of this loop body, which are executed for each values of  $i$ , ie,  $i = k+1, k+2, \dots, n$ . This eliminates  $x_{k+1}$  from all the equations after the  $k+1$ th equation of the system. Value of  $i+1$  is the equation<sup>5</sup> from which  $x_{k+1}$  is eliminated.
- Line 6 `if a[i,k] != 0.0:`  
If  $a[i,k] = A_{i+1,k+1} \neq 0$ , then those three instruction following it are executed. Otherwise, it skips the execution of those three statements. If  $x_{k+1} = x[k]$  is not there in the  $i$ th equation, it doesn’t need to be eliminated.
- Line 7 `lam = a[i,k]/a[k,k]`  
In this step,  $\lambda$  is computed so that  $\text{equ.}(i+1) - \lambda \text{ equ.}(k+1)$  doesn’t have  $x_{k+1}$  term in it.
- Line 8 `a[i,k+1:n] = a[i,k+1:n] - lam * a[k,k+1:n]`  
Coefficients of  $(i+1)$ th equation are updated.  
Equivalent to  $a[i,0:n] = a[i,0:n] - \text{lam} \times a[k,0:n]$ , since zeroes need not be subtracted. This is same as  $\text{equ.}(i+1) \leftarrow \text{equ.}(i+1) - \lambda \text{ equ.}(k+1)$
- Line 9 `b[i] = b[i] - lam * b[k]`  
The same row operations are performed on the matrix  $b$  instead of using an augmented matrix.
- Line 10 `for k in range(n-1,-1,-1):`  
This is another loop construct. The following statement is executed  $n$  times for values of  $k = n-1, n-2, \dots, 0$ . Value of  $k+1$  gives the unknown  $x_{k+1}$  which is solved by the back substitution process.

---

<sup>5</sup>Python starts counting from zero. For example :  $A_{11} = a[0,0]$ ,  $x_1 = x[0]$  and  $b_1 = b[0]$

Line 11  $x[k] = (b[k] - \text{dot}(a[k, k+1 : n], x[k+1 : n]))/a[k, k]$

This is the back substitution process. After elimination phase we have  $k$  equation in the form  $A_{k,k}x_k + A_{k,k+1}x_{k+1} + \dots + A_{k,n}x_n = b_k$ . And we already have values of  $x_{k+1}, x_{k+2}, \dots, x_n$ . Then

$$x_k = \frac{b_k - (A_{k,k+1}x_{k+1} + A_{k,k+2}x_{k+2} + \dots)}{A_{k,k}}$$

This is equivalent to

$$b_k \leftarrow \frac{b_k - [A_{k,k+1} \quad A_{k,k+2} \quad \dots \quad A_{k,n}] \begin{bmatrix} x_{k+1} \\ x_{k+2} \\ \vdots \\ x_n \end{bmatrix}}{A_{k,k}}$$

Remember : The values of  $x_k$  are updated into  $b_k$  as they are computed. Thus  $x_k, x_{k+1}, \dots, x_n$  are stored in  $b$  for next back substitution ie, for evaluating  $x_{k-1}$ . We start with  $x_{n-1}$ , as  $x_n = b_n$  is already solved.

Line 12 **return b**

It returns the new  $b$  matrix as output of the “gaussElimination()” function where  $x_k = b_k, \forall k$ .

## 8.6 LU Decomposition Method : Doolittle

Let  $Ax = b$  be a linear system of  $n$  equations in  $n$  unknowns and let  $A = LU$  for some lower triangular matrix  $L$  and upper triangular matrix  $U$ . Then we have  $LUx = Ly = b$  where  $y = Ux$ . There are two phases for this method : 1. LU decomposition and 2. substitution.

First, we compute  $L$  and  $U$  such that  $A = LU$  using Gauss elimination. Then We can solve  $Ly = b$  using forward substitution process and then solve  $Ux = y$  using back substitution process.

For Doolittle decomposition, we prefer to write  $A$  as a product  $LU$  as shown below:

$$A = \begin{bmatrix} 1 & 0 & \dots & 0 \\ L_{2,1} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ L_{n,1} & L_{n,2} & \dots & 1 \end{bmatrix} \begin{bmatrix} U_{1,1} & U_{1,2} & \dots & U_{1,n} \\ 0 & U_{2,2} & \dots & U_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & U_{n,n} \end{bmatrix}$$

$$A = \begin{bmatrix} U_{1,1} & U_{1,2} & \dots & U_{1,n} \\ L_{2,1}U_{1,1} & L_{2,1}U_{1,2} + U_{2,2} & \dots & L_{2,1}U_{1,n} + U_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ L_{n,1}U_{1,1} & L_{n,1}U_{1,2} + L_{n,2}U_{2,2} & \dots & \sum_{k=1}^{n-1} L_{n,k}U_{k,n} + U_{n,n} \end{bmatrix}$$



Note that in Doolittle's decomposition method, the diagonal entries of the lower triangular matrix  $L$  are all 1. ie,  $L_{ii} = 1, \forall i$ . Thus, we can use an  $n \times n$  matrix to represent both  $L$  and  $U$  by overwriting trivial entries( zeroes and ones) of both the matrices. And this matrix is represented by  $[L \setminus U]$ .<sup>6</sup>

$$[L \setminus U] = \begin{bmatrix} U_{1,1} & U_{1,2} & \cdots & U_{1,n} \\ L_{2,1} & U_{2,2} & \cdots & U_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ L_{n,1} & L_{n,2} & \cdots & U_{n,n} \end{bmatrix}$$

is the combined matrix made from both the triangular matrices  $L$  and  $U$ .

The triangular matrices  $L$  and  $U$  such that  $LU = A$  can be computed the variables in the Gauss elimination method.

$$A = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ L_{2,1} & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ L_{n,1} & L_{n,2} & \cdots & 1 \end{bmatrix} \begin{bmatrix} U_{1,1} & U_{1,2} & \cdots & U_{1,n} \\ 0 & U_{2,2} & \cdots & U_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & U_{n,n} \end{bmatrix}$$

We can break down this matrix multiplication into the following row operations on the rows of the upper triangular matrix<sup>7</sup>

$$\begin{aligned} U_{R1} &\leftarrow U_{R1} \\ U_{R2} &\leftarrow L_{2,1} \cdot U_{R1} + U_{R2} \\ &\vdots \\ U_{Rn} &\leftarrow L_{n,1} \cdot U_{R1} + L_{n,2} \cdot U_{R2} + \cdots + L_{n,n-1} \cdot U_{R(n-1)} + U_{Rn} \end{aligned}$$

Clearly,  $\lambda$  we use to eliminate  $x_k$  from row  $i$  are  $L_{i,k}$ . And the matrix obtained after Gauss elimination is the upper triangular matrix  $U$ .

### 8.6.1 Illustrative example

$$\text{Solve } \begin{bmatrix} -3 & 6 & -4 \\ 9 & -8 & 24 \\ -12 & 24 & -26 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -3 \\ 65 \\ -42 \end{bmatrix}$$

#### Phase 1 : LU Decomposition

Suppose, we have a system of three linear equations, then

$$A = LU = \begin{bmatrix} 1 & 0 & 0 \\ L_{21} & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix} \begin{bmatrix} U_{1,1} & U_{1,2} & U_{1,3} \\ 0 & U_{2,2} & U_{2,3} \\ 0 & 0 & U_{3,3} \end{bmatrix}$$

$$A = \begin{bmatrix} U_{1,1} & U_{1,2} & U_{1,3} \\ L_{2,1}U_{1,1} & L_{2,1}U_{1,2} + U_{2,2} & L_{2,1}U_{1,3} + U_{2,3} \\ L_{3,1}U_{1,1} & L_{3,1}U_{1,2} + L_{3,2}U_{2,2} & L_{3,1}U_{1,3} + L_{3,2}U_{2,3} + U_{3,3} \end{bmatrix}$$

<sup>6</sup>algorithmic implementation all decomposition algorithms prefer to use a combined matrix

<sup>7</sup> $U_{Rk}$  :  $k$ th row of the matrix  $U$

We can compute  $L$  and  $U$  using the Gauss elimination process<sup>8</sup> The matrix obtained after Gauss elimination on  $A$  is  $U$  and the values of the variable  $lam$  used in Gauss elimination are the entries in  $L$ . That is, in order to eliminate  $x_k$  from row  $i$ , we use  $lam = L_{i,k}$ .

$$\text{Given, } A = \begin{bmatrix} -3 & 6 & -4 \\ 9 & -8 & 24 \\ -12 & 24 & -26 \end{bmatrix}$$

$$\begin{bmatrix} -3 & 6 & -4 \\ 9 & -8 & 24 \\ -12 & 24 & -26 \end{bmatrix} \xrightarrow[R_3 \leftarrow R_3 - 4R_1]{R_2 \leftarrow R_2 + 3R_1} \begin{bmatrix} -3 & 6 & -4 \\ 0 & 10 & 12 \\ 0 & 0 & -10 \end{bmatrix} \implies L_{2,1} = 3, L_{3,1} = -4$$

We store these non-trivial entries of  $L$  into  $A$  itself.  
That is,  $A_{2,1} = L_{2,1}$ ,  $A_{3,1} = L_{3,1}$ .

*“ In this case,  $A_{3,2}$  became zero (this is not a trivial zero yet), and we won't eliminate  $x_2$  from row 3 to save computation time. Thus, we are not computing  $L_{3,2} = 0$  or storing it. However, the variable representing  $L_{3,2}$  is  $A_{3,2}$ , which is already zero after Gauss elimination and we are quite happy with that. ”*

Clearly,  $L_{3,2} = 0$ . Therefore, we have

$$U = \begin{bmatrix} -3 & 6 & -4 \\ 0 & 10 & 12 \\ 0 & 0 & -10 \end{bmatrix}, L = \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ -4 & 0 & 1 \end{bmatrix}$$

Since we are already stored those two non-trivial entries of  $L$  into  $A$ . We get,

$$[L \setminus U] = \begin{bmatrix} -3 & 6 & -4 \\ 3 & 10 & 12 \\ -4 & 0 & -10 \end{bmatrix}$$

## Phase 2 : Substitution

Suppose  $Ly = b$ ,

$$\begin{bmatrix} 1 & 0 & 0 \\ L_{2,1} & 1 & 0 \\ L_{3,1} & L_{3,2} & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \rightarrow \begin{cases} y_1 & = b_1 \\ L_{2,1}y_1 + y_2 & = b_2 \\ L_{3,1}y_1 + L_{3,2}y_2 + y_3 & = b_3 \end{cases}$$

Now, we can find the values of  $y_k$  and store them into the matrix  $b$  itself.

$$b_1 \leftarrow y_1 = b_1$$

$$b_2 \leftarrow y_2 = b_2 - [L_{2,1}] [b_1], \text{ since } b_1 = y_1$$

$$b_3 \leftarrow y_3 = b_3 - [L_{3,1} \quad L_{3,2}] \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \text{ since } b_1 = y_1, b_2 = y_2$$

---

<sup>8</sup>We usually need a proof for such a strong statement. In this paper, they are more focussed on the application side and therefore we will don't present any vigorous proof.

In general,

$$b_k \leftarrow y_k = b_k - [L_{k,1} \quad L_{k,2} \quad \cdots \quad L_{k,k-1}] \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{k-1} \end{bmatrix}, \text{ since } b_j = y_j, j = 1, 2, \dots, (k-1)$$

We have  $A = LU \implies LUx = b$ . Suppose  $Ux = y$ , then we get  $Ly = b$ . First of all, we will solve  $Ly = b$  using forward substitution.

$$\begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ -4 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} -3 \\ 65 \\ -42 \end{bmatrix} \rightarrow \begin{cases} y_1 & = -3 \\ 3y_1 + y_2 & = 65 \\ -4y_1 + y_3 & = -42 \end{cases}$$

$$\begin{aligned} y_1 &= -3 \\ y_2 &= 65 - 3y_1 = 74 \\ y_3 &= -42 + 4y_1 = 54 \end{aligned}$$

Suppose  $Ux = y$ ,

$$\begin{bmatrix} U_{1,1} & U_{1,2} & U_{1,3} \\ 0 & U_{2,2} & U_{2,3} \\ 0 & 0 & U_{3,3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

Now, we can find the values of  $x_k$  and store them into the matrix  $y$  itself.

$$\begin{aligned} y_3 \leftarrow x_3 &= \frac{y_3}{U_{3,3}} \\ y_2 \leftarrow x_2 &= \frac{y_2 - [y_3] [U_{2,3}]}{U_{2,2}}, \text{ since } y_3 = x_3 \\ y_1 \leftarrow x_1 &= \frac{y_1 - \begin{bmatrix} y_2 \\ y_3 \end{bmatrix} [U_{1,2} \quad U_{1,3}]}{U_{1,1}}, \text{ since } y_2 = x_2, y_3 = x_3 \end{aligned}$$

In general,

$$x_k = \frac{y_k - \begin{bmatrix} y_{k+1} \\ y_{k+2} \\ \vdots \\ y_n \end{bmatrix} [U_{k,k+1} \quad U_{k,k+2} \quad \cdots \quad U_{k,n}]}{U_{k,k}}, \text{ since } y_j = x_j, j = k+1, k+2, \dots, n$$

## 8.6.2 Python : Doolittle's LU Decomposition method

**Program 8.6.1.**

```
from numpy import dot
def LUdecomposition(a):
    n = len(a)
    for k in range(0, n-1):
        for i in range(k+1, n):
```

```

        if a[i, k] != 0.0:
            lam = a[i, k]/a[k, k]
            a[i, k+1:n] = a[i, k+1:n] - lam*a[k, k+1:n]
            a[i, k] = lam
    return a
def LUsolve(a, b):
    n = len(a)
    for k in range(1, n):
        b[k] = b[k] - dot(a[k, 0:k], b[0:k])
        b[n-1] = b[n-1]/a[n-1, n-1]
    for k in range(n-2, -1, -1):
        b[k] = (b[k] - dot(a[k, k+1:n], b[k+1:n]))/a[k, k]
    return b

```

This program mainly uses the Gauss elimination algorithm. Thus, the explanation for Lines 3-8 are not repeated here.

But remember the loop at Line 4 has inner loop at Line 5 and Line 7-9 are at same level of indentation which means they all are either executed or skipped depending on the truthness of the condition in Line 6. And Line 6-9 are executed for each instance of inner loop. Again, Line 5-9 are executed for each instance of the outer loop.

This time the gaussElimination() function which you have seen earlier is split into two functions 1. LUdecomposition() and 2. LUsolve(). And forward substitution is also added to LUsolve().

Line 2 **def LUdecomposition(a):**

LUdecomposition( $A$ ) computes  $L$  and  $U$  such that  $A = LU$  and combine both triangular matrices into a single matrix  $[L/U]$ , by over-writing their trivial entries. And returns this combined matrix.

Line 9  $a[i, k] = lam$

Clearly,  $lam$  used for eliminating  $x_k$  from row  $i$ ,  $\lambda_{i,k} = a[i, k]/a[k, k] = L[i, k]$ ,  $\forall k, \forall i$ , ( $i > k$ ). Also  $a[i, k]$  which is reduced zero by Gauss elimination process is not used anymore<sup>9</sup> in Gauss elimination process. Thus  $L[i, k]$  can stored at  $a[i, k]$  straight away. And  $U[i, j]$ ,  $j \leq i$  are already the entries of the matrix obtained from Gauss elimination. Thus for each iterations of  $k$ , the matrix  $a$  is updated ( $k+1$ th row and  $k+1$ th column) with respective entries of the combined matrix  $[L/U]$ .

Line 10 **return a**

Matrix  $a$  is already  $[L/U]$ , and thus LUdecomposition( $A$ ) returns  $[L/U]$  such that  $A = LU$ .

Line 11 **def LUsolve(a, b):**

LUsolve() function does both forward substitution and back substitution. Suppose  $Ax = b$  is the system to be solved. Then the inputs of LUsolve() are  $a = [L/U]$  where  $A = LU$ .

<sup>9</sup> $A[i, k]$  is not used after elimination of  $x_k$  from row  $i$  - It turns out that the trivial zeroes which are ignored on the row operations in Gauss elimination not only save time, but also provide a variable to store our intermediate result  $L_{i,k}$  in Doolittle method.

Line 12  $n = \text{len}(a)$

We have to compute this again as this function starts fresh and thus value of the variable  $n$  from `LUdecomposition()` is lost.

Line 13 **for**  $k$  **in** `range(1,n):`

This is the loop for forward substitution.

Line 14  $b[k] = b[k] - \text{dot}(a[k, 0 : k], b[0 : k])$

updating  $b_{k^*}$  with  $y_{k^*}$  such that  $Ly = b$  where  $k^* = k - 1$ .

$$b_k \leftarrow b_k - [L_{k,1} \quad L_{k,2} \quad \cdots \quad L_{k,k-1}] \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{k-1} \end{bmatrix}$$

Line 15  $b[n-1] = b[n-1]/a[n-1, n-1]$

Computing <sup>10</sup>  $y_n$  and storing it into  $b_n$ .

$$b_n \leftarrow \frac{b_n}{U_{n,n}}$$

Line 16 **for**  $k$  **in** `range(n-2, -1, -1):`

This is the loop for back substitution.

Line 17  $b[k] = (b[k] - \text{dot}(a[k, k+1 : n], b[k+1 : n]))/a[k, k]$

updating  $b_{k^*}$  with  $x_{k^*}$  such that  $Ux = y$  where  $k^* = k - 1$ .

$$b_k \leftarrow \frac{b_k - \begin{bmatrix} x_{k+1} \\ x_{k+2} \\ \vdots \\ x_n \end{bmatrix} [U_{k,k+1} \quad U_{k,k+2} \quad \cdots \quad U_{k,n}]}{U_{k,k}}$$

Line 18 **return**  $b$

`LUsolve([L\U], b)` returns  $b$  where  $b[i+1] = x_i$ .

Programmer's Tip : There are few things to remember when splitting a function into two functions.

1. These functions are completely independent of one another.
2. Variables defined inside a function are not available outside.
3. The best way to give/take data to/from a function is through arguments/return-value

Beginner's Tip : In any programming language, we reuse variable. Thus, same variable may represent different values at different points of time. In Doolittle LU Decomposition, the variable 'a' initially represent matrix  $A$ , this variable is passed into `LUdecomposition()` function. In that function,  $A$  is changed to  $[L\backslash U]$  in a step-by-step fashion. The value of 'a' is updated in step 7, 8 and 9. This is bit hard to imagine this transition of 'a' from  $A$  to  $[L\backslash U]$  for a beginner at programming. Similarly, in `LUsolve()` function, the variable 'b' changes from matrix  $b$  to matrix  $y$ , and then to matrix  $x$ .

<sup>10</sup>Mathematically, you can define dot product of empty matrices as zero, but numpy dot function can't handle such a situation. Therefore, we have to do this step separately.

## 8.7 Numerical Integraion

Numerical integration/Quadrature is the numerical approximation of  $\int_a^b f(x)dx$  by  $\sum_{i=0}^n A_i f(x_i)$  where  $x_i$  are nodal abscissas, and  $A_i$  are weights. There are two methods to determine these nodal abscissas and suitable weights so that the sum is sufficiently accurate to the value of the integral.

1. Newton-Cotes forumulas
2. Gauss quadrature

Newton-Cotes formulas are useful when  $f(x)$  can be evaluated without much computation. And using those values  $f(x)$  can be interpolated to a piecewise-polynomial function. Then using equally spaces nodal abscissas and suitable weights  $\int_a^b f(x)dx$  can be numerically approximated.

Gauss quadrature rules require lesser evaluations of  $f$ . And therefore are quite useful when evaluation of  $f(x)$  has much computational complexity. Also, this method can manage integrable singularities where as Newton-Cote formulas can't numerically integrate function with singularities.

### 8.7.1 Newton-Cotes formulas

We divide the interval of integral  $(a, b)$  into  $n$  subintervals of equal length, ie,  $h = (b - a)/n$ . Let  $x_0 = a, x_1, x_2, \dots, x_{n-1}, x_n = b$  be the end points of these subintervals. Then we can find an  $n$  degree polynomial interpolant satisfying  $f$  at those points, using Lagrange's method.

$$\text{Polynomial, } P(x) = \sum_{i=0}^n f(x_i) l_i(x) \text{ where } l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}$$

Thus the integral  $I = \int_a^b f(x)dx$  can be numerically evaluated as follows:

$$\begin{aligned} I &= \int_a^b P_n(x)dx = \sum_{i=0}^n \left( f(x_i) \int_a^b l_i(x)dx \right) \\ &= \sum_{i=0}^n A_i f(x_i), \text{ where } A_i = \int_a^b l_i(x)dx \end{aligned}$$

The simplest cases of Newton-Cotes formulas are when  $n = 1, 2, \text{ and } 3$

**Trapezoidal rule**  $n = 1 \implies A_0 = \frac{h}{2}, A_1 = \frac{h}{2}$  and

$$\int_a^b f(x)dx = A_0 f(x_0) + A_1 f(x_1) = \frac{h}{2}(f(a) + f(b))$$

**Simpson's 1/3 rule**  $n = 2 \implies A_0 = \frac{h}{3}, A_1 = \frac{4h}{3}, A_2 = \frac{h}{3}$  and

$$\int_a^b f(x)dx = \sum_{i=0}^2 A_i f(x_i) = \frac{h}{3} \left( f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right)$$

**Simpson's 3/8 rule**  $n = 3 \implies A_0 = \frac{3h}{8}, A_1 = \frac{9h}{8}, A_2 = \frac{9h}{8}, A_3 = \frac{3h}{8}$  and

$$\int_a^b f(x)dx = \sum_{i=0}^3 A_i f(x_i) = \frac{3h}{8}(f(a) + 3f(a+h) + 3f(a+2h) + f(b))$$

**Trapezoidal Rule :**  $n = 1$

Consider interval  $(a, b)$ . Since  $n = 1$ , we have  $x_0 = a$  and  $x_1 = b$ .

$$l_0(x) = \frac{x - x_1}{x_0 - x_1}$$

$$l_1(x) = \frac{x - x_0}{x_1 - x_0}$$

$$\begin{aligned} A_0 &= \int_a^b l_0(x)dx = \int_a^b \frac{x - x_1}{x_0 - x_1}dx = \frac{-1}{h} \int_a^b (x - b)dx \\ &= \frac{-1}{h} \left( \frac{(x - b)^2}{2} \right)_a^b = \frac{-1}{h} \left( \frac{0 - (a - b)^2}{2} \right) = \frac{h}{2} \\ A_1 &= \int_a^b l_1(x)dx = \int_a^b \frac{x - x_0}{x_1 - x_0}dx = \frac{1}{h} \int_a^b (x - a)dx \\ &= \frac{1}{h} \left( \frac{(x - a)^2}{2} \right)_a^b = \frac{1}{h} \left( \frac{(b - a)^2 - 0}{2} \right) = \frac{h}{2} \end{aligned}$$

Therefore,

$$\int_a^b f(x)dx = A_0 f(x_0) + A_1 f(x_1) = \frac{h}{2}(f(a) + f(b))$$

**Simpon's 1/3 Rule :**  $n = 2$

Consider interval  $(a, b)$  divided into two subintervals of equal length  $h = \frac{a+b}{2}$ .

We have  $x_0 = a, x_1 = \frac{a+b}{2}, x_2 = b$ .

$$l_0(x) = \frac{x - x_1}{x_0 - x_1} \frac{x - x_2}{x_0 - x_2}$$

$$l_1(x) = \frac{x - x_0}{x_1 - x_0} \frac{x - x_2}{x_1 - x_2}$$

$$l_2(x) = \frac{x - x_0}{x_2 - x_0} \frac{x - x_1}{x_2 - x_1}$$

$$\begin{aligned} A_0 &= \int_a^b l_0(x)dx \\ &= \int_a^b \frac{x - x_1}{x_0 - x_1} \frac{x - x_2}{x_0 - x_2}dx \\ &= \int_a^b \left( \frac{x - \frac{a+b}{2}}{a - \frac{a+b}{2}} \right) \left( \frac{x - b}{a - b} \right) dx \end{aligned}$$

Changing variable of integration  $y = x - \frac{a+b}{2}$

$$y = x - \frac{a+b}{2} \implies dy = dx$$

$$x = a \implies y = -h$$

$$x = b \implies y = h$$

Continuing with the value of  $A_0$ ,

$$\begin{aligned} A_0 &= \int_{-h}^h \left( \frac{y}{-h} \right) \left( \frac{y-h}{-2h} \right) dy \\ &= \frac{1}{2h^2} \int_{-h}^h y^2 - \frac{1}{2h} \int_{-h}^h y dy \\ &= \frac{1}{2h^2} \left( \frac{h^3}{3} - \frac{(-h)^3}{3} \right) - \frac{1}{2h} \left( \frac{h^2}{2} - \frac{(-h)^2}{2} \right) \\ &= \frac{h}{3} \end{aligned}$$

$$\begin{aligned} A_1 &= \int_a^b l_1(x) dx \\ &= \int_a^b \frac{x-x_0}{x_1-x_0} \frac{x-x_2}{x_1-x_2} dx \\ &= \int_a^b \left( \frac{x-a}{h} \right) \left( \frac{x-b}{-h} \right) dx \end{aligned}$$

Applying change of variable,  $y = x - \frac{a+b}{2}$

$$\begin{aligned} &= \int_{-h}^h \left( \frac{y+h}{h} \right) \left( \frac{y-h}{-h} \right) dy \\ &= \frac{-1}{h^2} \int_{-h}^h y^2 dy + \int_{-h}^h 1 dy \\ &= \frac{-1}{h^2} \left( \frac{h^3}{3} - \frac{(-h)^3}{3} \right) + (h - (-h)) \\ &= \frac{-2h^3}{3h^2} + 2h \\ &= \frac{4h}{3} \end{aligned}$$

$$\begin{aligned} A_2 &= \int_a^b l_2(x) dx \\ &= \int_a^b \left( \frac{x-x_0}{x_2-x_0} \right) \left( \frac{x-x_1}{x_2-x_1} \right) dx \\ &= \int_a^b \left( \frac{x-a}{2h} \right) \left( \frac{x-\frac{a+b}{2}}{h} \right) dx \end{aligned}$$



Applying change of variable,  $y = x - \frac{a+b}{2}$

$$\begin{aligned}
 &= \int_{-h}^h \left( \frac{y+h}{2h} \right) \left( \frac{y}{h} \right) dy \\
 &= \frac{1}{2h^2} \int_{-h}^h y^2 dy + \frac{1}{2h} \int_{-h}^h y dy \\
 &= \frac{1}{2h^2} \left( \frac{h^3}{3} - \frac{(-h)^3}{3} \right) + \frac{1}{2h} \left( \frac{h^2}{2} - \frac{(-h)^2}{2} \right) \\
 &= \frac{2h^3}{6h^2} \\
 &= \frac{h}{3}
 \end{aligned}$$

Therefore,

$$\int_a^b f(x) dx = \sum_{i=0}^2 A_i f(x_i) = \frac{h}{3} \left( f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right)$$

### 8.7.2 Composite Trapezoidal Rule

Suppose an interval  $(a, b)$  is divided into  $n$  subintervals. In Composite Trapezoidal Rule, Trapezoidal Rule is applied to each subinterval. Thus we have,

$$\begin{aligned}
 I &= I_0 + I_1 + \cdots + I_{n-1} \text{ where } I_k \text{ is the integral over } (x_k, x_{k+1}) \\
 &= \frac{h(f(x_0) + f(x_1))}{2} + \frac{h(f(x_1) + f(x_2))}{2} + \cdots + \frac{h(f(x_{n-1}) + f(x_n))}{2} \\
 &= \frac{h}{2} (f(x_0) + 2f(x_1) + \cdots + 2f(x_{n-1}) + f(x_n))
 \end{aligned}$$

### 8.7.3 Recursive Trapezoidal Rule

Suppose an interval  $(a, b)$  is divided into  $2^{k-1}$  subintervals. In Recursive Trapezoidal Rule, we apply Trapezoidal Rule on each subinterval. And there is a recursive formula since the number of intervals doubles as value of  $k$  increases by one. Let  $H = b - a$

$$k = 1 \implies 2^0 = 1 \text{ and}$$

$$I_1 = \frac{H}{2} (f(a) + f(b))$$

$$k = 2 \implies 2^1 = 2 \text{ and}$$

$$\begin{aligned}
 I_2 &= \frac{H}{4} (f(a) + 2f\left(a + \frac{H}{2}\right) + f(b)) \\
 &= \frac{H}{4} (f(a) + f(b)) + \frac{H}{2} f\left(a + \frac{H}{2}\right) \\
 &= \frac{I_1}{2} + \frac{H}{2} f\left(a + \frac{H}{2}\right)
 \end{aligned}$$

$k = 3 \implies 2^2 = 4$  and

$$\begin{aligned} I_3 &= \frac{H}{8} \left( f(a) + 2f\left(a + \frac{H}{4}\right) + 2f\left(a + \frac{2H}{4}\right) + 2f\left(a + \frac{3H}{4}\right) + f(b) \right) \\ &= \frac{H}{8} \left( f(a) + 2f\left(a + \frac{2H}{4}\right) + 2f\left(a + \frac{4H}{4}\right) + 2f\left(a + \frac{6H}{4}\right) + f(b) \right) \\ &\quad + \frac{H}{4} \left( f\left(a + \frac{H}{4}\right) + f\left(a + \frac{3H}{4}\right) \right) \\ &= \frac{I_2}{2} + \frac{H}{4} \left( f\left(a + \frac{H}{4}\right) + f\left(a + \frac{3H}{4}\right) \right) \end{aligned}$$

Consider interval  $(0, 64)$ . We have  $b - a = H = 64$ .

$$I_1 = 32(f(0) + f(64))$$

$$I_2 = 16(f(0) + 2f(32) + f(64))$$

$$I_3 = 8(f(0) + 2f(16) + 2f(32) + 2f(48) + f(64))$$

$$I_4 = 4(f(0) + 2f(8) + 2f(16) + \cdots + 2f(56) + f(64))$$

$$I_5 = 2(f(0) + 2f(4) + 2f(8) + \cdots + 2f(60) + f(64))$$

$$I_6 = f(0) + 2f(2) + 2f(4) + \cdots + 2f(62) + f(64)$$

$\vdots$

The values corresponding to the intervals in  $I_{k-1}$  appear in  $I_k$  as alternate terms. Other terms, corresponds to the odd multiples of  $\frac{H}{2^k}$ . We separate them into two sums and represent the first sum as  $\frac{I_{k-1}}{2}$ .

$$\begin{aligned} \text{Clearly, } I_k &= \frac{H}{2^k} \sum_{i=0}^{2^{k-2}} f\left(a + \frac{2iH}{2^k}\right) + \frac{2H}{2^k} \sum_{i=1}^{2^{k-2}} f\left(a + \frac{(2i-1)H}{2^{k-1}}\right) \\ &= \frac{I_{k-1}}{2} + \frac{H}{2^{k-1}} \sum_{i=1}^{2^{k-2}} f\left(a + \frac{(2i-1)H}{2^{k-1}}\right) \end{aligned}$$

### 8.7.4 Python : Recursive Trapezoidal Rule

#### Program 8.7.1.

```
def recursiveTrapezoidalRule(f, a, b, Iold, k):
    if k == 1:
        Inew = (f(a) + f(b)) * (b - a) / 2.0
    else:
        n = 2 ** (k - 2)
        h = (b - a) * 1.0 / n
        x = a + h / 2.0
        sum = 0.0
        for i in range(n):
            sum = sum + f(x)
            x = x + h
        Inew = (Iold + h * sum) / 2.0
    return Inew
```

Line 1 `def recursiveTrapezoidalRule(f, a, b, Iold, k):`

This function has five input arguments. (a)  $f$  is a real function (b)  $a, b$  are start and end of interval in which  $f$  is going to be integrated (c)  $Iold$  is

the value of the integral for  $2^{k-1}$  subintervals using recursive trapezoidal method (d)  $k$  is a variable such that  $2^k$  is the number of subintervals considered for Integration.

- Line 2 **if**  $k == 1$  :  
 If  $k = 1$ , we proceed to Line 3, otherwise we go to Line 4.
- Line 3  $I_{new} = (f(a) + f(b)) * (b - a) / 2.0$   
 For  $k = 1$ , we use trapezoidal rule  $I_1 = \frac{b-a}{2}(f(a) + f(b))$ . When writing this in python, we use 2.0 so that the python won't ignore the decimal part of this fraction. In python,  $5/2 = 2$ . And  $5/2.0 = 2.5$
- Line 4 **else** :  
 If Line 2 is false, (ie  $k \neq 1$ ) python executes Line 5-8. These line implements the recursive formula for  $I_k$ .
- Line 5  $n = 2 * (k - 2)$   
 Equivalent to  $n \leftarrow 2^{k-2}$ .
- Line 6  $h = (b - a) * 1.0 / n$   
 This the length of a subinterval when we divide  $(a, b)$  into  $2^k$  subintervals/panels. Equivalent to  $h \leftarrow \frac{b-a}{2^{k-2}}$
- Line 7  $x = a + h/2.0$   
 This the parameter of  $f$  in the first term in the sum  $\sum_{i=1}^{2^{k-2}} f\left(a + \frac{(2i-1)H}{2^{k-1}}\right)$  in the recursive formula for  $I_k$ . Equivalent to  $x \leftarrow a + \frac{h}{2}$ .
- Line 8  $sum = 0.0$   
 We are going to use this variable to find that sum. To start with, we will make it 0 and will add each term to it one-by-one. Equivalent to  $sum \leftarrow 0$ .
- Line 9 **for**  $i$  **in**  $\text{range}(n)$  :  
 This the variable  $i$  in the recursive formula for  $I_k$ . For each value of  $i = 1, 2, \dots, 2^{k-2}$ , Lines 10 and 11 are executed. That is, for each value of  $i$ , the corresponding term in the sum is computed and added to the variable  $sum$ .
- Line 10  $sum = sum + f(x)$   
 Value of  $f$  at  $x$  is computed and added to the partial sum. Equivalent to  $sum \leftarrow sum + f(x)$ . However, the value of  $x$  is changed for each  $i$  in Line 11. Thus, for next value of  $i$ ,  $x$  and  $sum$  have the new values to use.
- Line 11  $x = x + h$   
 Equivalent to  $x \leftarrow x + h$ . For  $i = 1$ ,  $x = a + \frac{h}{2}$  before Line 11. At Line 11,  $x \leftarrow a + \frac{3h}{2}$ . And this is the value of  $x$  for  $i = 2$  before Line 11 next time. Thus,  $x$  iterates through  $a + \frac{h}{2}, a + \frac{3h}{2}, \dots, a + \frac{(2n-1)h}{2}$ . This  $x$  is updated and used for next execution of Line 10 and 11.
- Line 12  $I_{new} = (I_{old} + h * sum) / 2.0$   
 We reach here only after executing Line 10-11 for all values of  $i$ . That is,  $sum$  in the recursive formula is already computed. This line, implements the recursive formula and stores that value into the variable  $I_{new}$ . Equivalent to  $I_{new} \leftarrow \frac{I_{old} + h * sum}{2}$ .

Line 13 **return Inew**

It returns the value of  $I_k$ , the integral of  $f$  over  $(a, b)$  using recursive trapezoidal rule for  $2^k$  subintervals.

## Subject 9

# ME010204 Complex Analysis

### 9.1 Module 2

#### 9.1.1 Arcs & Closed Curves

An arc  $\gamma$  in a complex plane is defined as the set of points given by  $\gamma = \{z : z = z(t), a \leq t \leq b\}$  and  $z(t)$  is a continuous function of the real variable  $t$ . Thus, every arc in the complex plane is the continuous image of closed interval and  $z \in \gamma$  means  $z = z(t) = x(t) + iy(t)$ . That is, points on  $\gamma$  are images of a complex function of a real variable.

This representation of an arc  $z = z(t)$  is called a parametric representation and  $t$  is called the parameter and  $[a, b]$  is called the parametric interval. The point  $z = z(a)$  is called origin or initial point of  $\gamma$  and  $z = z(b)$  is called the terminus or terminal point of  $\gamma$ . If  $z(a) = z(b)$ , then  $\gamma$  is called a closed curve, otherwise it is open.

If in  $\gamma$ ,  $z(t_1) = z(t_2) = z$  for  $t_1 \neq t_2$ , then  $z$  is called a multiple point on  $\gamma$ . Geometrically, a multiple point  $z$  in  $\gamma$  is a point where  $\gamma$  crosses itself.

An arc having no multiple points is called a simple arc or Jordan arc.

#### 9.1.2 Differentiable Arc

**differentiable**  $z'(t)$  exist and is continuous at all points.

**regular** differentiable and  $z'(t) \neq 0$  at points on  $\gamma$ .

**piecewise differentiable** differentiable except for finitely many points

**piecewise regular** regular except for finitely many points.

**opposite arc** set of points  $z = z(-t)$ ,  $-b \leq t \leq -a$ .

#### 9.1.3 Complex Integration

1. If  $f(z)$  has an antiderivative  $F(z)$ , then  $\int_{\gamma} f(z) dz$  depends only on the end points and is independent of the path  $\gamma$ .

2. If  $f(z)$  has an antiderivative, then  $\int_{\gamma} f(z) dz = 0$  for all closed curves  $\gamma$ .
3. If  $f(z)$  has no antiderivatives,  $\int_{\gamma} f(z) dz$  depends on the path  $\gamma$ .

For different choices of  $\gamma$ , the integral may have different values even though the end points are the same.

### 9.1.4 Exercise

Evaluate  $\int_c f(z) dz$  where  $f(z) = y - x - i3x^2$  and contour is  $c_1$ : the line segment 0 and  $i + 1$   $c_2$ : the polygon joining  $(0, 0)$ ,  $(0, 1)$ ,  $(1, 1)$  In the above example,  $\int_{c_1} f(z) dz \neq \int_{c_2} f(z) dz$  even though  $c_1$  and  $c_2$  have the same end points.

### 9.1.5 Evaluating Line Integral : Method 1

Integrals of the form  $\int_a^b f(t) dt$  where  $f(t)$  is a complex valued function of a real variable  $t$ . Then, we can write  $f(t) = u(t) + iv(t)$ .

$$\int_a^b f(t) dt = \int_a^b u(t) + iv(t) dt = \int_a^b u(t) dt + i \int_a^b v(t) dt$$

That is,  $\Re \int f(t) dt = \int \Re f(t) dt$  and  $\Im \int f(t) dt = \int \Im f(t) dt$ .

For example,  $f(t) = e^{it}$ ,  $0 \leq t \leq \frac{\pi}{2}$ . Then

$$\int_0^{\frac{\pi}{2}} f(t) dt = \int_0^{\frac{\pi}{2}} \cos t dt + i \int_0^{\frac{\pi}{2}} \sin t dt = 1 + i$$

### 9.1.6 Evaluating Line Integral : Method 2

Let  $\gamma$  be a piecewise differentiable arc in the complex plane defined by the equation  $z = z(t)$ ,  $a \leq t \leq b$  and  $f(z)$  be defined and continuous in  $\gamma$ . Then the line integral  $\int_{\gamma} f(z) dz$  is defined by

$$\int_{\gamma} f(z) dz = \int_a^b f(z(t)) z'(t) dt$$

For example, let  $f(z) = u + iv$ . Then  $f(z) dz = (u + iv)(dx + idy)$

$$\int_{\gamma} f(z) dz = \int_{\gamma} u dx - v dy + i \int_{\gamma} v dx + u dy$$

That is, real and imaginary part of  $\int_{\gamma} f(z) dz$  can be written in the form  $\int_{\gamma} p dx + q dy$  where  $p, q$  are real valued functions of  $x$  and  $y$ . **That is, real valued functions of two real variables  $x$  and  $y$ .** Therefore, the study of line integral  $\int_{\gamma} f(z) dz$  can be restricted to the study of line integrals of the form  $\int_{\gamma} p dx + q dy$  and line integrals can be defined as  $\int_{\gamma} p dx + q dy$  where  $\gamma$  is a piecewise differentiable arc.

**Properties**

1. Scalar multiplication,

$$\int_{\gamma} cf(z) dz = c \int_{\gamma} f(z) dz$$

2. Modulus Inequality,

$$\left| \int_a^b f(t) dt \right| \leq \int_a^b |f(t)| dt$$

—more—

3. Change of variable,

$$\int_{\gamma} f(z) dz = \int_a^b f(z(t)) z'(t) dt$$

4. Inverse arc,

$$\int_{-\gamma} f(z) dz = - \int_{\gamma} f(z) dz$$

5. Integration by parts,

$$\int_{\gamma} f(z) dz = \int_{\gamma_1} f(z) dz + \int_{\gamma_2} f(z) dz + \cdots + \int_{\gamma_n} f(z) dz$$

**9.1.7 Line Integral : Type 3**

Line integrals with respect to  $\bar{z}$  are denoted by

$$\int_{\gamma} f(z) d\bar{z} = \overline{\int_{\gamma} \bar{f} dz}$$

*Proof.* We have,  $x = \frac{z+\bar{z}}{2}$  and  $y = \frac{z-\bar{z}}{2i}$ .

$$dx = \frac{dz + d\bar{z}}{2} \quad dy = \frac{dz - d\bar{z}}{2i}$$

$$\int_{\gamma} f(z) dx = \int_{\gamma} f(z) \frac{dz + d\bar{z}}{2} \quad \int_{\gamma} f(z) dy = \int_{\gamma} f(z) \frac{dz - d\bar{z}}{2i}$$

$$\int_{\gamma} f(z) dx - idy = \int_{\gamma} f(z) d\bar{z}$$

—more—

□

**9.1.8 Line Integral : Type 4**

Line integrals with respect to arc length,  $s$  is denoted by

$$\int_{\gamma} f(z) ds = \int_{\gamma} f(z) |dz| \text{ where } ds = \sqrt{1 + \left(\frac{dy}{dx}\right)^2} dx$$

When  $f = 1$ , it gives the length of the arc.

**9.1.9 Rectifiable Arc**

Length of an arc  $\gamma$  is defined as  $L(\gamma) = \int |dz|$ . If  $L(\gamma) < \infty$ , then  $\gamma$  is a rectifiable arc and the process of determining the length of an arc is called rectification.



## Subject 10

# ME010205 Measure & Integration

### 10.2 Lebesgue Measure

**set function** A function which maps sets into (extended) real numbers.

**$\sigma$ -algebra** A family  $\mathcal{A}$  of subsets of a nonempty set  $X$  such that

1.  $\mathcal{A}$  contains the empty set,
2.  $\mathcal{A}$  contains complement of each of its members and
3.  $\mathcal{A}$  is closed under countable unions.

From these 3 axioms, we can deduce the following,

4.  $\mathcal{A}$  is closed under countable intersections (by de Morgan's laws).

$$\left( \bigcup_{k=1}^{\infty} E_k^c \right)^c = \bigcap_{k=1}^{\infty} E_k \in \mathcal{A}$$

5.  $E, F \in \mathcal{A} \implies E - F \in \mathcal{A}$  since  $E - F = E \cap F^c$

**Definitions 10.2.1** (Length of an interval). Length is a real valued set function. Let  $I$  be a bounded interval say  $[a, b)$ . Then its length  $l(I) = b - a$  is the difference between endpoints. If an interval  $I$  is unbounded say  $(a, \infty)$ , then its length,  $l(I) = \infty$ .

#### Exercise

#### Techniques in Measure Theory

Let  $\mathcal{A}$  be a  $\sigma$ -algebra. Let Lebesgue Measure  $m : \mathcal{A} \rightarrow [0, \infty]$  be countably additive over disjoint collection of sets in  $\mathcal{A}$ .

- Lebesgue Measure  $m$  has monotonicity.  
 $A \subseteq B \implies B = A \cup (B - A)$  is a disjoint union  
 $\implies m(B) = m(A) + m(B - A) \geq m(A)$

- If exists  $E \in \mathcal{A}$  such that  $m(E) < \infty$ , then  $m(\phi) = 0$   
Suppose  $m(\phi) = c$  and  $m(E) = k$  where  $k < \infty$ . If  $c \neq 0$ , then  $m(E \cup \phi) = m(E) + m(\phi) = c + k > k = m(E)$  is a contradiction.
- $m\left(\bigcup_{k=1}^{\infty} E_k\right) \leq \sum_{k=1}^{\infty} m(E_k)$   
Define  $\{F_k : k \in \mathbb{N}\}$  by  $F_k = E_k - \bigcup_{j=1}^{k-1} E_j$   
Then  $F_1 = E_1$ ,  $F_2 = E_2 - E_1$ ,  $F_3 = E_3 - (E_1 \cup E_2)$ , ...  
Also  $F_k \in \mathcal{A}$  and  $F_k \subseteq E_k$ ,  $\forall k \in \mathbb{N}$ . Thus  $m(F_k) \leq m(E_k)$ ,  $\forall k$  However,  
 $\bigcup_{k=1}^{\infty} E_k = \bigcup_{k=1}^{\infty} F_k$   
 $\implies m\left(\bigcup_{k=1}^{\infty} E_k\right) = m\left(\bigcup_{k=1}^{\infty} F_k\right) = \sum_{k=1}^{\infty} m(F_k) \leq \sum_{k=1}^{\infty} m(E_k)$

### Counting Measure

The counting measure  $c : \mathcal{A} \rightarrow [0, \infty]$  is a set function which maps sets to their cardinality. For example, if  $E = \{2, 3, 4\}$ , then  $c(E) = 3$ .

- The counting measure is **translation invariant** since translation never increases the cardinality of the set.  
For example,  $5 + E = \{7, 8, 9\}$ . And  $m(5 + E) = 3 = m(E)$ .
- The counting measure is **countably additive** over disjoint collections since the cardinality of disjoint union of two sets is the sum of their cardinalities.
- However, counting measure of (non-degenerate) intervals are  $\infty$  which is **not the same as their length** for bounded intervals.

#### 10.2.1 Lebesgue Outer Measure

$G_\delta$  A set which is countable intersection of open subsets.

$F_\sigma$  A set which is countable union of closed subsets.

#### Set-theoretic Construction of Lebesgue Measure

1. Construct Lebesgue Outer Measure  $m^*$  (with Axiom 3 relaxed)  
ie, Obtain the underlying relation of the set function
2. Restrict  $m^*$  to the  $\sigma$ -algebra of our interest  
ie, Choose a domain so that set function is well defined.

**Definitions 10.2.2** (Lebesgue Outer Measure). Let  $A \subset \mathbb{R}$ . Let  $\mathcal{C} = \{I_k : k \in \mathbb{N}\}$  be an open cover of  $A$  such that  $I_k$  are non-empty, bounded, open intervals. Consider the sum of length of intervals for such covers of  $A$ . (Lebesgue) Outer Measure  $m^*(A)$  is the infimum of all such sums.

$$m^*(A) = \inf \left\{ \sum_{k=1}^{\infty} l(I_k) : A \subset \bigcup_{k=1}^{\infty} I_k \right\} \quad (10.1)$$

### 10.2.2 Properties of Lebesgue Outer Measure

1. Outer Measure of the empty set is zero

Let  $\epsilon > 0$ . Then  $\mathcal{C}_\epsilon = \{(0, \frac{\epsilon}{2^n}) : n \in \mathbb{N}\}$  is an open cover of  $\phi$  containing nonempty, bounded, open intervals. Clearly, sum of length of intervals in  $\mathcal{C}_\epsilon = \epsilon$ . Suppose  $m^*(\phi) = \delta$  and  $\delta > 0$ . There exists  $\epsilon$  such that  $0 < \epsilon < \delta$ . The sum of intervals of  $\mathcal{C}_\epsilon$  is less than  $\delta$ , which is a contradiction by the definition of Outer Measure.

2. Outer Measure is monotone

Suppose  $A \subset B$ . Then every cover of  $B$  is also an cover of  $A$ . Let  $\mathcal{U}$  be the set of all open covers of  $A$  with nonempty, bounded intervals and  $\mathcal{V}$  be the set of all such open covers of  $A$ . Clearly,  $\mathcal{V} \subset \mathcal{U}$ . We know that, if  $A \subset B$ , then  $\inf \mathcal{V} \leq \inf \mathcal{U}$ . Therefore,

$$A \subset B \implies m^*(A) \leq m^*(B) \quad (10.2)$$

3. Outer Measure of Countable Sets is zero

Let  $C$  be a countable set. That is,  $C = \{c_k\}_{k=1}^\infty$ .

Then  $\{(c_k - \frac{\epsilon}{2^k}, c_k + \frac{\epsilon}{2^k})\}_{k=1}^\infty$  is cover of  $C$  with sum of length of intervals  $\epsilon$ . Thus, for any  $\epsilon > 0$ , we have  $m^*(C) \leq \epsilon$ . Thus,  $m^*(C) = 0$ .

4. Outer Measure of an Interval is its length

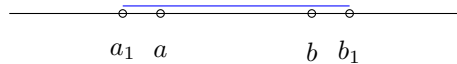
*Proof. Case 1 : Closed, Bounded Interval* Let  $[a, b]$  be a closed, bounded interval. Then for any  $\epsilon > 0$ ,  $(a - \epsilon, b + \epsilon)$  is a cover of  $[a, b]$ . Thus, by the definition of Lebesgue outer measure  $m^*([a, b]) \leq b - a + 2\epsilon$  since  $[a, b] \subset (a - \epsilon, b + \epsilon)$  and  $m^*$  is monotonic. Therefore,

$$m^*([a, b]) \leq b - a \quad (10.3)$$

Since  $[a, b]$  is closed and bounded,  $[a, b]$  is compact. And by Heine-Borel theorem, every open cover of  $[a, b]$  has a finite subcover. Thus, it is sufficient to prove the theorem for finite covers of  $[a, b]$ .

Let  $\mathcal{C}$  be a finite cover of  $[a, b]$  with  $n$  open intervals. Let  $(a_1, b_1)$  be an open interval containing  $a$  in  $\mathcal{C}$ . Then  $a_1 < a < b_1$ . If  $b_1 > b$  then

$l(a_1, b_1) > l(a, b)$ . And  $\sum_{i=1}^k l(I_k) \geq l(a_1, b_1) \geq b - a$ .



Suppose  $b_1 < b$ . Clearly,  $a < b_1$ . And the cover  $\mathcal{C}$  must have an open interval containing  $b_1$ . Otherwise  $\mathcal{C}$  is not a cover of  $[a, b]$ . That is, there exists  $(a_2, b_2)$  containing  $b_1 \in (a, b)$  such that  $a_2 < b_1 < b_2$ . If  $b_2 > b$ , then

$\sum_{i=1}^k l(I_k) \geq l(a_1, b_1) + l(a_2, b_2) \geq l(a_1, b_2) \geq b - a$ .



Suppose  $b_2 < b$ . Continuing like this we get,  $N$  open intervals in  $\mathcal{C}$ ,  $\{(a_k, b_k) : k = 1, 2, \dots, N\}$  such that  $a_1 < a < b_1$  and  $a_N < b < b_N$  and  $a_k < b_{k-1} < b_k$  for all  $k$ . The process should terminate in finite steps as  $\mathcal{C}$  is a finite cover of  $[a, b]$ . Then  $\sum_{k=1}^N l(I_k) \geq \sum_{k=1}^N l(a_k, b_k) \geq l(a_1, b_N) \geq b - a$ .



Clearly, every open cover of  $[a, b]$  contains a finite subcover  $\mathcal{C}$ , which contains a finite subcover of the form  $\{(a_k, b_k) : k = 1, 2, \dots, N\}$  such that  $\sum_{k=1}^N l(I_k) \geq b - a$ . Thus, for any open cover  $\sum_{k=1}^{\infty} l(I_k) \geq b - a$ . And thus,

$$m^*([a, b]) \geq b - a \quad (10.4)$$

**Case 2 : Bounded Interval** Let  $I$  be a bounded interval. Then there exists bounded closed intervals  $J_1$  and  $J_2$  such that  $J_1 \subsetneq I \subsetneq J_2$  such that  $l(I) - \epsilon < l(J_1)$  and  $l(J_2) < l(I) + \epsilon$ . Suppose  $I = (a, b]$ , then  $J_1 = [a + \frac{\epsilon}{2}, b - \frac{\epsilon}{2}]$  and  $J_2 = [a - \frac{\epsilon}{2}, b + \frac{\epsilon}{2}]$ .

By monotonicity of Lebesgue outer measure, we have  $m^*(J_1) \leq m^*(I) \leq m^*(J_2)$ . However  $m^*(J_1) = l(I) - \epsilon$  and  $m^*(J_2) = l(I) + \epsilon$ . Thus,  $l(I) - \epsilon \leq m^*(I) \leq l(I) + \epsilon$ . Therefore,  $m^*(I) = l(I)$ .

**Case 3 : Unbounded Interval** Let  $I$  be an unbounded interval. Then for any natural number  $n$ , there exists a closed bounded interval  $J$  such that  $J \subset I$  and  $l(J) = n$ . And  $n = m^*(J) \leq m^*(I)$ ,  $\forall n \in \mathbb{N}$ . Therefore,  $m^*(I) = \infty = l(I)$ .  $\square$

## 5. Outer Measure is translation invariant

*Proof.* Let  $A$  be any set and  $y \in \mathbb{R}$ . Let  $\{I_k : k = 1, 2, \dots\}$  be a cover of  $A$ . Then  $\{I_k + y : k = 1, 2, \dots\}$  is a cover of  $A + y$ . And  $l(I_k) = l(I_k + y)$  for every natural number  $k$  and real number  $y$ . Thus,  $\sum_{k=1}^{\infty} l(I_k) = \sum_{k=1}^{\infty} l(I_k + y)$ .

Clearly, for each cover  $\{I_k\}_{k=1}^{\infty}$  of  $A$ , there exists a cover  $\{I_k + y\}_{k=1}^{\infty}$  of  $A + y$  containing intervals of same length. Therefore,  $m^*(A) = m^*(A + y)$ .  $\square$

## 6. Outer Measure is countably subadditive

*Proof.* Let  $\{E_k\}_{k=1}^{\infty}$  be a countable collection of sets. It is enough to prove that

$$m^*\left(\bigcup_{k=1}^{\infty} E_k\right) \leq \sum_{k=1}^{\infty} m^*(E_k) \quad (10.5)$$

For each natural number  $k$ , we have a cover of  $E_k$ , say  $\{I_{k,i}\}_{i=1}^{\infty}$  such that  $\sum_{i=1}^{\infty} l(I_{k,i}) < m^*(E_k) + \frac{\epsilon}{2^k}$ . Suppose that, for some  $\epsilon > 0$ ,  $E_k$  doesn't have such a cover, then  $m^*(E_k) + \frac{\epsilon}{2^k}$  is an upper bound contradicting the assumption that  $m^*(E_k)$  is the least upper bound.

Clearly, ‘

$$\begin{aligned} m^*\left(\bigcup_{i,k=0}^{\infty} I_{k,i}\right) &\leq \sum_{k=1}^{\infty} \sum_{i=1}^{\infty} l(I_{k,i}) \\ &= \sum_{k=1}^{\infty} \left(m^*(E_k) + \frac{\epsilon}{2^k}\right) \\ &= \sum_{k=1}^{\infty} m^*(E_k) + \epsilon \end{aligned}$$

□

**Note :** Finite subadditivity is a weaker notion than countable subadditivity. Since every finite collection is a countable collection.

### Exercise

5. Closed Interval  $[0, 1]$  is uncountable.

Suppose  $[0, 1]$  is countable, then Lebesgue outer measure of any countable set is zero,  $m^*([0, 1]) = 0$ . But,  $[0, 1]$  is an interval and Lebesgue outer measure of an interval is its length,  $m^*([0, 1]) = l([0, 1]) = 1$  which is a contradiction.

6.  $m^*([0, 1] - \mathbb{Q}) = 1$

$$[0, 1] = ([0, 1] \cap \mathbb{Q}) \cup ([0, 1] \cap \mathbb{Q}^c)$$

Clearly,  $m^*([0, 1]) = 1$ . And  $[0, 1] \cap \mathbb{Q}$  is a countable set since  $\mathbb{Q}$  is countable. And thus has Lebesgue outer measure zero. Thus by countable subadditivity, we have

$$1 = m^*([0, 1]) \leq m^*([0, 1] \cap \mathbb{Q}^c) + 0$$

Thus,  $m^*([0, 1] \cap \mathbb{Q}^c) \geq 1$ . And  $[0, 1] \cap \mathbb{Q}^c \subset [0, 1]$ . By monotonicity,  $m^*([0, 1] \cap \mathbb{Q}^c) \leq m^*([0, 1]) = 1$ . Therefore,  $m^*([0, 1] \cap \mathbb{Q}^c) = 1$ .

7. Construction of a  $G_\delta$  set containing  $E$

8. hint : if sum of interval is less than 1. Then it is not a cover of  $[0, 1]$ .

9. hint :  $A \cup B = A \cup (B - A) = A \cup (B \cap A^c)$

10. hint :  $A$  and  $B$  are separated by distance  $\alpha$ , thus are disjoint.

### 10.2.3 $\sigma$ -algebra of Lebesgue Measurable Sets

Lebesgue Outer Measure is defined for any subset of real numbers and Lebesgue outer measure of an interval is its length. However, it isn't countable additive.

There exists disjoint sets  $A, B$  such that  $m^*(A \cup B) < m^*(A) + m^*(B)$ .

Since countable additivity is a favourable property over countable subadditivity. We restrict the family of subsets of real numbers to those subsets that allow countable additivity.

#### Lebesgue Measurable Set

**Definitions 10.2.3** (Measurable Set). Let  $E$  be a subset of  $\mathbb{R}$ . Then  $E$  is Lebesgue measurable if

$$m^*(A) = m^*(A \cap E) + m^*(A \cap E^c) \quad (10.6)$$

for any subset  $A$  of  $\mathbb{R}$ .

In other words,  $E$  is Lebesgue measurable if  $E$  doesn't affect countable additivity of Lebesgue Outer Measure.

We will consider only those subset of real numbers, which won't affect countable additivity. These subsets are **Lebesgue Measurable**. And we could show that the collection of all Lebesgue Measurable sets forms a  $\sigma$ -algebra. Clearly, intervals allow countable additivity, thus the Borel Algebra is contained in this  $\sigma$ -algebra of Lebesgue measurable sets.

#### Simplified Condition for Lebesgue Measurability

We know that Lebesgue Outer Measure has countable subadditivity.

$$m^*(A) \leq m^*(A \cap E) + m^*(A \cap E^c)$$

Thus, for condition (10.6), it is sufficient to check the following condition,

$$m^*(A) \geq m^*(A \cap E) + m^*(A \cap E^c) \quad (10.7)$$

#### Properties of Lebesgue Measure

1. Any set of Lebesgue outer measure zero is Lebesgue measurable.

*Proof.* Let  $E$  be a subset of real numbers with Lebesgue outer measure zero. Let  $A$  be any subset of real numbers. Then  $A = (A \cap E) \cup (A \cap E^c)$ . By countable additivity,  $m^*(A) \leq m^*(A \cap E) + m^*(A \cap E^c)$ . Since  $A \cap E \subset E$ , we have by monotonicity  $m^*(A \cap E) \leq m^*(E) = 0$ .

Again,  $A \cap E^c \subset A$  and by monotonicity,  $m^*(A) \geq m^*(A \cap E^c) = 0 + m^*(A \cap E^c) = m^*(A \cap E) + m^*(A \cap E^c)$ . Thus,  $E$  is Lebesgue measurable by the simplified condition for Lebesgue measurability.  $\square$

## 2. Countable sets are Lebesgue measurable.

*Proof.* Countable sets are of Lebesgue outer measure zero. And sets of Lebesgue outer measure zero are Lebesgue measurable. Thus, they are Lebesgue measurable.  $\square$

## 3. Finite union of Lebesgue measurable sets is Lebesgue measurable.

*Proof.* It is enough to prove that if  $E_1$  and  $E_2$  are Lebesgue measurable, then their union is also Lebesgue measurable. Then, by finite mathematical induction, we can prove that the result is true for any finite collection of Lebesgue measurable sets.

Suppose  $E_1, E_2$  are Lebesgue measurable sets. Since  $E_1$  is Lebesgue measurable,

$$m^*(A) = m^*(A \cap E_1) + m^*(A \cap E_1^c) \quad (10.8)$$

And consider  $A \cap E_1^c$  instead of  $A$ . Since  $E_2$  is Lebesgue measurable, we get

$$m^*(A \cap E_1^c) = m^*(A \cap E_1^c \cap E_2) + m^*(A \cap E_1^c \cap E_2^c) \quad (10.9)$$

We have  $(A \cap E_1^c) \cap E_2 = A \cap (E_1^c \cap E_2) = A \cap (E_1 \cup E_2)^c$ . And  $(A \cap E_1) \cup (A \cap E_1^c \cap E_2) = (A \cap E_1) \cup [A \cap (E_2 \cap E_1^c)] = A \cap (E_1 \cup E_2)$ .



$$\begin{aligned} m^*(A) &= m^*(A \cap E_1) + m^*(A \cap E_1^c) \\ &= m^*(A \cap E_1) + m^*(A \cap E_1^c \cap E_2) + m^*(A \cap E_1^c \cap E_2^c) \\ &\geq m^*[A \cap (E_1 \cup E_2)] + m^*[A \cap (E_1 \cup E_2)^c] \end{aligned}$$

Therefore  $E_1 \cup E_2$  is Lebesgue measurable. And by finite induction, finite union of Lebesgue measurable sets is also Lebesgue measurable.  $\square$

## 4. Lebesgue Measure is finitely additive.

In other words, Suppose  $\{E_k\}_{k=1}^n$  be a finite collection of disjoint, Lebesgue measurable sets. Then Lebesgue measure of their union is the sum of Lebesgue measures.

*Proof.* Let  $A$  be any subset of  $\mathbb{R}$  and  $\{E_k\}_{k=1}^n$  be a finite collection of disjoint, Lebesgue measurable subsets of  $\mathbb{R}$ .

$$\text{Claim : } m^*\left(A \cap \left[\bigcup_{k=1}^{\infty} E_k\right]\right) = \sum_{k=1}^{\infty} m^*(A \cap E_k) \quad (10.10)$$

Trivially, the claim is true for  $n = 1$ . Suppose the claim is true for  $n - 1$ . That is,

$$m^* \left( A \cap \left[ \bigcup_{k=1}^{n-1} E_k \right] \right) = \sum_{k=1}^{n-1} m^*(A \cap E_k) \quad (10.11)$$

From set theory we have,

$$A \cap \left[ \bigcup_{k=1}^n E_k \right] \cap E_n = A \cap E_n \quad (10.12)$$

$$A \cap \left[ \bigcup_{k=1}^n E_k \right] \cap E_n^c = A \cap \left[ \bigcup_{k=1}^{n-1} E_k \right] \quad (10.13)$$

By Lebesgue measurability of  $E_n$ , we have

$$\begin{aligned} m^* \left( A \cap \left[ \bigcup_{k=1}^n E_k \right] \right) &= m^* \left( A \cap \left[ \bigcup_{k=1}^n E_k \right] \cap E_n \right) + m^* \left( A \cap \left[ \bigcup_{k=1}^n E_k \right] \cap E_n^c \right) \\ &= m^*(A \cap E_n) + m^* \left( A \cap \left[ \bigcup_{k=1}^{n-1} E_k \right] \right) \\ &= \sum_{k=1}^n m^*(A \cap E_k), \text{ by mathematical induction} \end{aligned}$$

Taking  $A = \mathbb{R}$ , we get Lebesgue measure is finitely additive. That is,

$$m^* \left( \bigcup_{k=1}^n E_k \right) = \sum_{k=1}^n m^*(E_k) \quad (10.14)$$

□

#### 5. Countable union of Lebesgue measurable sets is Lebesgue measurable

*Proof.* Let  $A$  be any subset of  $\mathbb{R}$ . And  $\{E_k\}_{k=1}^{\infty}$  be a countable collection of disjoint, Lebesgue measurable subsets of  $\mathbb{R}$ . Define  $F_n = \bigcup_{k=1}^n E_k$  and  $E = \bigcup_{k=1}^{\infty} E_k$ . Clearly,  $F \subset E$  and  $F^c \supset E^c$ . Thus,  $m^*(A \cap F_n^c) \geq m^*(A \cap E^c)$ .

$$\begin{aligned} m^*(A) &= m^*(A \cap F_n) + m^*(A \cap F_n^c) \\ &\geq m^* \left( A \cap \left[ \bigcup_{k=1}^n E_k \right] \right) + m^*(A \cap E^c) \\ &\geq \sum_{k=1}^n m^*(A \cap E_k) + m^*(A \cap E^c) \end{aligned}$$



$$\begin{aligned}
\lim_{n \rightarrow \infty} m^*(A) &\geq \lim_{n \rightarrow \infty} \sum_{k=1}^n m^*(A \cap E_k) + m^*(A \cap E^c) \\
m^*(A) &\geq \sum_{k=1}^{\infty} m^*(A \cap E_k) + m^*(A \cap E^c) \\
&\geq m^*\left(A \cap \left[\bigcup_{k=1}^{\infty} E_k\right]\right) + m^*(A \cap E^c) \\
&\implies m^*(A) \geq m^*(A \cap E) + m^*(A \cap E^c)
\end{aligned}$$

By above inequality,  $E = \bigcup_{k=1}^{\infty} E_k$  is Lebesgue measurable.

And, any countable collection of Lebesgue measurable sets can be expressed as a countable collection of disjoint, Lebesgue measurable sets. Let  $\{E_k\}_{k=1}^{\infty}$  be a collection of Lebesgue measurable sets. Then, the countable collection,  $\{E'_k\}_{k=1}^{\infty}$  defined by  $E'_k = E_k - \bigcup_{j=1}^{k-1} E_j$  contains disjoint, Lebesgue measurable<sup>†1</sup> subsets of  $\mathbb{R}$ . Therefore, countable union of Lebesgue measurable sets is Lebesgue measurable.  $\square$

6. Every interval is Lebesgue measurable.

*Proof.* It is sufficient to prove that  $(a, \infty)$  is Lebesgue measurable. Suppose  $(a, \infty)$  is Lebesgue measurable for every  $a \in \mathbb{R}$ . Then interval  $(a, b)$  is Lebesgue measurable, since  $(a, b) = [(a, \infty) \cap (\mathbb{R} - (b, \infty))] - \{b\}$ .

Let  $A$  be any subset of  $\mathbb{R}$ . Define  $A_1 = A \cap (-\infty, a)$  and  $A_2 = A \cap (a, \infty)$  such that  $A - \{a\} = A_1 \cup A_2$  and  $A_1 \cap A_2 = \emptyset$ . And the interval  $(a, \infty)$  is Lebesgue measurable only if

$$m^*(A) \geq m^*(A \cap (a, \infty)) + m^*(A \cap (a, \infty)^c) = m^*(A_1) + m^*(A_2)^{\dagger 2} \quad (10.15)$$

Let collection  $\{I_k\}_{k=1}^{\infty}$  be a countable cover of  $A$ . Then collection  $\{I'_k\}_{k=1}^{\infty}$  defined by  $I'_k = I_k \cap (a, \infty)$  is a cover of  $A_1$ . And collection  $\{I''_k\}_{k=1}^{\infty}$  defined by  $I''_k = I_k \cap (-\infty, a)$  is a cover of  $A_2$ .

Clearly,  $m^*(A_1) \leq \sum_{k=1}^{\infty} l(I'_k)$  and  $m^*(A_2) \leq \sum_{k=1}^{\infty} l(I''_k)$ .

$$\begin{aligned}
m^*(A_1) + m^*(A_2) &\leq \sum_{k=1}^{\infty} l(I'_k) + \sum_{k=1}^{\infty} l(I''_k) \\
&\leq \sum_{k=1}^{\infty} l(I'_k) + l(I''_k) \\
&\leq \sum_{k=1}^{\infty} l(I_k) \leq m^*(A)
\end{aligned}$$

<sup>†1</sup>Suppose  $E_1, E_2$  are measurable, then  $E_2^c = \mathbb{R} - E_2$  is measurable by the duality of measurability condition. And  $E_1 \cap E_2^c = E_1 - E_2$  is Lebesgue measurable since countable intersection of Lebesgue measurable sets is Lebesgue measurable (by Property 4 and de Morgan's Law).

<sup>†2</sup>We have,  $m^*(A \cap (-\infty, a]) = m^*(A \cap (-\infty, a))$  since removing finite number of points from a subset of  $\mathbb{R}$  won't affect its Lebesgue measure.

Thus,  $(a, \infty)$  is Lebesgue measurable. Therefore, every interval is Lebesgue measurable.  $\square$

7.  $\sigma$ -algebra of Lebesgue measurable sets  $\mathcal{M}$  contains Borel Sets  $\mathcal{B}$ .

$$\mathcal{B} \subset \mathcal{M} \quad (10.16)$$

*Proof.* The Borel algebra  $\mathcal{B}$  is the  $\sigma$ -algebra containing all intervals. We have proved that, every intervals are Lebesgue measurable. Also, we have proved that the set of all Lebesgue measurable subsets of  $\mathbb{R}$  is a  $\sigma$ -algebra as complements of Lebesgue measurable sets are Lebesgue measurable by duality of the condition and countable union of Lebesgue measurable sets are also Lebesgue measurable. Therefore, every Borel set is Lebesgue measurable. Clearly,  $G_\delta$  and  $F_\sigma$  are Borel sets and are Lebesgue measurable.  $\square$

8. Lebesgue Measurability is translation invariant.

*Proof.* Let  $E$  be a Lebesgue measurable set. Let  $A$  be any subset of  $\mathbb{R}$  and  $y \in \mathbb{R}$ . Then,

$$\begin{aligned} m^*(A) &= m^*(A - y) \\ &= m^*((A - y) \cap E) + m^*((A - y) \cap E^c) \\ &= m^*(A \cap (E + y)) + m^*(A \cap (E + y)^c) \end{aligned}$$

Thus,  $E + y$  is Lebesgue measurable. And Lebesgue measurability is translation invariant.  $\square$

### Exercise

11. Let  $\mathcal{A}$  be the  $\sigma$ -algebra containing all intervals of the form  $(a, \infty)$ . Every interval has one of the four forms,

$$(a, b] = (a, \infty) \cap (b, \infty)^c \quad (10.17)$$

$$(a, b) = (a, \infty) \cap \left[ \bigcap_{k=1}^{\infty} \left( b - \frac{1}{k}, \infty \right) \right]^c \quad (10.18)$$

$$[a, b] = \left[ \bigcap_{k=1}^{\infty} \left( a - \frac{1}{k}, \infty \right) \right] \cap (b, \infty)^c \quad (10.19)$$

$$[a, b) = \left[ \bigcap_{k=1}^{\infty} \left( a - \frac{1}{k}, \infty \right) \right] \cap \left[ \bigcap_{k=1}^{\infty} \left( b - \frac{1}{k}, \infty \right) \right]^c \quad (10.20)$$

12. **Borel sets** is the  $\sigma$ -algebra containing all open intervals. The equations from previous equations are sufficient. However, we have a simpler form for closed intervals,  $[a, b]$ .

$$[a, b] = [(-\infty, a) \cup (b, \infty)]^c \quad (10.21)$$

Clearly, every interval is a Borel set.

13. •  $C \in F_\sigma \implies C = \bigcup_{k=1}^\infty C_k \implies \bigcup_{k=1}^\infty (C_k + y) = C + y \in F_\sigma$   
 •  $O \in G_\delta \implies O = \bigcap_{k=1}^\infty O_k \implies \bigcup_{k=1}^\infty (O_k + y) = O + y \in G_\delta$   
 •  $m^*(E) = 0 \implies 0 = \inf\{\sum_{k=1}^\infty l(I_k)\} = \inf\{\sum_{k=1}^\infty l(I_k + y)\} \implies m^*(E + y) = 0$

14. A subset  $E$  has positive Lebesgue measure if and only if it has a bounded subset of positive Lebesgue measure.

$$m^*(E) > 0 \iff \exists \text{ bounded subset } F \subset E, \text{ such that } m^*(F) > 0$$

**Sufficient Part :** By monotonicity,  $F \subset E \implies m^*(F) \leq m^*(E)$ . And  $m^*(F) > 0 \implies 0 < m^*(F) \leq m^*(E)$ .

**Necessary Part :**

15.

#### 10.2.4 Outer and Inner Approximation

**Definitions 10.2.4** (Excision Property). Let  $A$  be a Lebesgue measurable set of finite measure and  $A \subset B$ . Then,

$$m^*(B \sim A) = m^*(B) - m^*(A) \quad (10.22)$$

*Proof.*

$$\begin{aligned} m^*(B) &= m^*(B \cap A) + m^*(B \cap A^c) \\ &= m^*(A) + m^*(B \sim A) \\ \implies m^*(B \sim A) &= m^*(B) - m^*(A) \end{aligned}$$

□

**Theorem 10.2.1** (approximation). *The following conditions are equivalent to Lebesgue measurability of  $E$*

1.  $\forall \epsilon > 0$ , there is an open set  $\mathcal{O}$  containing  $E$  for which  $m^*(\mathcal{O} \sim E) < \epsilon$ .
2. There is a  $G_\delta$  set  $G$  containing  $E$  such that  $m^*(G \sim E) = 0$ .
3.  $\forall \epsilon > 0$ , there is a closed set  $F$  contained in  $E$  such that  $m^*(E \sim F) < \epsilon$ .
4. There is an  $F_\sigma$  set  $F$  contained in  $E$  such that  $m^*(E \sim F) = 0$ .

**Note :** Equivalent conditions 1 & 2 are about outer approximation of a Lebesgue measurable set and 3 & 4 are about inner approximation of a Lebesgue measurable set.

**Proof. Measurability  $\implies$  Open set - Outer approximation**

Suppose that  $E$  is Lebesgue measurable. And let  $\epsilon > 0$ .

**Case 1 :** Suppose  $m^*(E) < \infty$ . By definition of Lebesgue outer measure, there exists an open cover  $\{I_k\}_{k=1}^{\infty}$  such that  $\sum_{k=1}^{\infty} l(I_k) < m^*(E) + \epsilon$ .

Define  $\mathcal{O} = \bigcup_{k=1}^{\infty} I_k$ . Then  $\mathcal{O}$  is an open set containing  $E$  and thus Lebesgue measurable. Also,  $m^*(\mathcal{O}) \leq \sum_{k=1}^{\infty} l(I_k) < m^*(E) + \epsilon$ . Therefore, by excision property we have  $m^*(\mathcal{O} \sim E) = m^*(\mathcal{O}) - m^*(E) < \epsilon$ .

**Case 2 :** Suppose  $m^*(E) = \infty$ . Without loss of generality,  $E$  may be written as countable union Lebesgue measurable sets  $\{E_k\}_{k=1}^{\infty}$  of finite measure.

By case 1, for every  $k$ , there exists  $\mathcal{O}_k$  for each  $E_k$  of finite measure such that  $m^*(\mathcal{O}_k \sim E_k) < \frac{\epsilon}{2^k}$ . Define  $\mathcal{O} = \bigcup_{k=1}^{\infty} \mathcal{O}_k$ . Then  $\mathcal{O}$  is open, contains  $E$  and

$$\begin{aligned} m^*(\mathcal{O} \sim E) &= m^*\left(\bigcup_{k=1}^{\infty} \mathcal{O}_k \sim E\right) \\ &\leq m^*\left(\bigcup_{k=1}^{\infty} \mathcal{O}_k \sim E_k\right) \\ &\leq \sum_{k=1}^{\infty} m^*(\mathcal{O}_k \sim E_k) = \epsilon \sum_{k=1}^{\infty} \frac{1}{2^k} = \epsilon \end{aligned}$$

**Open, Outer approximation  $\implies G_{\delta}$ , Outer approximation**

Let  $E$  be a subset of real numbers such that Lebesgue measure of  $E$  has an open set inner approximation. That is, for every  $\epsilon > 0$ , there exists an open set  $\mathcal{O}$  such that  $m^*(\mathcal{O} \sim E) < \epsilon$ .

Let  $\mathcal{O}_k$  be open sets such that  $m^*(\mathcal{O}_k \sim E) < \frac{1}{k}$ . Define  $G = \bigcap_{k=1}^{\infty} \mathcal{O}_k$ . Then,  $G$  is a  $G_{\delta}$  set containing  $E$ . And  $G \sim E \subset \mathcal{O}_k \sim E$ . Thus, by monotonicity,  $m^*(G \sim E) \leq m^*(\mathcal{O}_k \sim E) < \frac{1}{k}$ . Thus, we have a  $G_{\delta}$  set  $G$  containing  $E$  such that  $m^*(G \sim E) = 0$ .

**$G_{\delta}$ -Outer approximation  $\implies$  Measurability**

We have,  $m^*(G \sim E) = 0$ . Since every set of Lebesgue measure zero is Lebesgue measurable,  $G \sim E$  is Lebesgue measurable. And its complement  $(G \sim E)^c$  is also Lebesgue measurable. Also we have,  $G$  is a  $G_{\delta}$  set, thus a Borel set and hence Lebesgue measurable.

Clearly, we have  $E = G \cap (G \sim E)^c$ . And therefore,  $E$  is Lebesgue measurable.

**Open, Outer approximation  $\iff$  Closed, Inner approximation**

By duality of Lebesgue measurability,  $E$  is Lebesgue measurable if and only if  $E^c$  is Lebesgue measurable. And by de Morgan's Law, we have  $E^c$  has an open set - outer approximation  $\mathcal{O}$  if and only if  $E$  has a closed set - inner approximation  $\mathcal{O}^c$ .

$$m^*(\mathcal{O} \sim E^c) < \epsilon \iff m^*(E \sim \mathcal{O}^c) < \epsilon$$

**$G_\delta$ -Outer approximation  $\iff$   $F_\sigma$ -Inner approximation**

Again, by duality of Lebesgue measurability and de Morgan's Law, we have  $E^c$  has a  $G_\delta$ -outer approximation  $G$  if and only if  $E$  has an  $F_\sigma$ -inner approximation  $F$ .

$$m^*(G \sim E^c) = 0 \iff m^*(E \sim F) = 0$$

□

**Theorem 10.2.2.** *Let  $E$  be a Lebesgue measurable set of finite measure. Then for any  $\epsilon > 0$ , there exists a finite collection of disjoint open sets  $\{I_k\}_{k=1}^n$  such that  $\mathcal{O} = \bigcup_{k=1}^n I_k$  and  $m^*(\mathcal{O} \sim E) + m^*(E \sim \mathcal{O}) < \epsilon$ .*

*Proof.* Since  $E$  is Lebesgue measurable, by outer approximation theorem we have an open set  $U$  such that  $E \subset U$  and

$$m^*(U \sim E) < \frac{\epsilon}{2}$$

$$U = (U \cap E) \cup (U \cap E^c)$$

Since  $E$  is Lebesgue measurable,  $m^*(E) < \infty$

$$\begin{aligned} m^*(U) &= m^*(U \cap E) + m^*(U \cap E^c) \\ &= m^*(E) + m^*(U \sim E) \end{aligned}$$

Since  $E$  has finite measure

$$\begin{aligned} m^*(U) &= m^*(E) + m^*(U \sim E) \\ &< \infty + \frac{\epsilon}{2} < \infty \end{aligned}$$

That is,  $U$  is of finite measure.

Since  $U$  is open,  $U$  is countable<sup>3</sup> union of a disjoint collection of open intervals, say  $\{I_k\}_{k=1}^\infty$ . Clearly,

$$\sum_{k=1}^n l(I_k) \leq \sum_{k=1}^\infty l(I_k) \leq m^*(U) < \infty$$

---

<sup>3</sup>By definition, **Open sets** are countable union of disjoint, open intervals

By characterisation of series convergence, there exists an integer  $n$  such that,

$$\sum_{k=n+1}^{\infty} l(I_k) < \frac{\epsilon}{2}$$

Define  $\mathcal{O} = \bigcup_{k=1}^n I_k$ . Since  $\mathcal{O} \sim E \subset U \sim E$ , by monotonicity we have

$$m^*(\mathcal{O} \sim E) \leq m^*(U \sim E) < \frac{\epsilon}{2} \quad (10.23)$$

Since  $E \subset U$ , we have  $E \sim \mathcal{O} \subset U \sim \mathcal{O} = \bigcup_{k=1}^n I_k$ . And clearly,

$$U \sim \mathcal{O} = \bigcup_{k=1}^{\infty} I_k \sim \bigcup_{k=1}^n I_k \subseteq \bigcup_{k=n+1}^{\infty} I_k$$

$$\text{Thus, } m^*(E \sim \mathcal{O}) \leq m^*(U \sim \mathcal{O}) \leq \sum_{k=n+1}^{\infty} l(I_k) < \frac{\epsilon}{2} \quad (10.24)$$

Therefore,

$$m^*(E \sim \mathcal{O}) + m^*(\mathcal{O} \sim E) < \epsilon$$

□

### Exercise

17. Let  $\epsilon > 0$  and  $E$  is Lebesgue measurable. Then there exists open set  $\mathcal{O}$  and closed set  $F$  such that  $F \subset E \subset \mathcal{O}$ ,  $m^*(E \sim F) < \frac{\epsilon}{2}$  and  $m^*(\mathcal{O} \sim E) < \frac{\epsilon}{2}$ . Clearly,  $\mathcal{O} \sim E$  and  $E \sim F$  are disjoint and  $\mathcal{O} \sim F = (\mathcal{O} \sim E) \cup (E \sim F)$ . Thus by monotonicity of Lebesgue outer measure, we have  $m^*(\mathcal{O} \sim F) \leq m^*(\mathcal{O} \sim E) + m^*(E \sim F) < \epsilon$ .

18. Suppose  $E$  has finite outer measure. <sup>†4</sup>

$G_\delta$  **set** : Let  $\epsilon > 0$ . Then by the definition of Lebesgue outer measure, there exists a cover  $\{I_k\}_{k=1}^{\infty}$  of  $E$  such that  $\sum_{k=1}^{\infty} l(I_k) < m^*(E) - \frac{\epsilon}{2}$ . Define

$$G = \bigcup_{k=1}^{\infty} I_k. \text{ Then } G \text{ is a } G_\delta \text{ set and } m^*(G) \leq \sum_{k=1}^{\infty} l(I_k) < m^*(E) - \frac{\epsilon}{2}.$$

$F_\sigma$  **set** :

19. Let  $E$  be a set of finite outer measure. Suppose  $E$  is not Lebesgue measurable. And  $\mathcal{O}$  be an open set containing  $E$ . Then  $\mathcal{O} = (\mathcal{O} \sim E) \cup E$ . By monotonicity,  $m^*(\mathcal{O}) \leq m^*(\mathcal{O} \sim E) + m^*(E)$  <sup>†5</sup>. Since  $m^*(E)$  is finite, we have  $m^*(\mathcal{O}) - m^*(E) \leq m^*(\mathcal{O} \sim E)$ .

<sup>4</sup> $E$  has finite outer measure does not imply  $E$  is bounded or Lebesgue measurable.

<sup>5</sup>**I am not able to change  $\leq$  into  $<$**  as non-measurability doesn't mean that for this particular  $\mathcal{O}$  the sum of Lebesgue outer measures should be greater. There may be a better proof.

20. Let  $E$  be a set of finite outer measure. Suppose  $E$  is Lebesgue measurable. Let  $(a, b)$  be any open, bounded interval. Then by the definition of Lebesgue measurability,  $b - a = m^*(a, b) = m^*((a, b) \cap E) + m^*((a, b) \cap E^c)$ .
21. A subset  $E$  is Lebesgue measurable if there exists a  $G_\delta$  set  $G$  containing  $E$  such that  $m^*(G \sim E) = 0$ . Suppose  $E_1$  and  $E_2$  are Lebesgue measurable sets. Then, we have  $G_\delta$  sets  $G_1$  and  $G_2$ . And two countable family of open intervals  $\{\mathcal{O}_{1,k}\}_{k=1}^\infty$  and  $\{\mathcal{O}_{2,k}\}_{k=1}^\infty$  such that  $\cap \mathcal{O}_{1,k} = G_1$  and  $\cap \mathcal{O}_{2,k} = G_2$ . Let  $\mathcal{O} = \{\mathcal{O}_k\}_{k=1}^\infty$  be collection of open intervals in  $\mathcal{O}_1$  and  $\mathcal{O}_2$ . Define  $G = \cap_{k=1}^\infty \mathcal{O}_k$ . Then,  $E = E_1 \cup E_2 \subset G_1 \cup G_2 = G$ . Since  $G \sim E = (G_1 \sim E_1) \cup (G_2 \sim E_2)$ , by monotonicity of Lebesgue outer measure we have  $m^*(G \sim E) \leq m^*(G_1 \sim E_1) + m^*(G_2 \sim E_2) = 0$ .
22. Let  $m^{**}$  be a non-negative set function defined by  $m^{**}(A) = \inf\{m^*(\mathcal{O}) : A \subset \mathcal{O}, \mathcal{O} \text{ is open}\}$ . Suppose  $E$  is Lebesgue measurable. Then, by open set - outer approximation theorem we have  $m^*(E) \leq m^{**}(E) < m^*(E) + \epsilon$  for any  $\epsilon > 0$ . Thus,  $m^{**}(E) = m^*(E)$ .

In other words, for Lebesgue measurable sets  $m^*$  and  $m^{**}$  are the same.

23. Let  $m^{***}$  be a non-negative set function defined by  $m^{***}(E) = \sup\{m^*(F) : F \subset E, F \text{ is closed}\}$ . Let  $E$  be Lebesgue measurable set. Then, by closed set - inner approximation theorem we have  $m^*(E) \geq m^{***}(E) > m^*(E) - \epsilon$ .

In other words, for Lebesgue measurable sets  $m^*$  and  $m^{***}$  are the same.

### 10.2.5 Further Properties

The Lebesgue measure has the following properties.

1. Every Borel set is Lebesgue measurable.
2. Lebesgue Measure of an interval is its length.
3. Lebesgue Measure is translation invariant.
4. Lebesgue Measure is countably additive.
5. There exists non-measurable (Lebesgue) sets. eg.  $C_E \subset E$
6. There exists uncountable set of zero measure. eg. Cantor set

#### Countable Subadditivity

**Theorem 10.2.3.** *The set function Lebesgue measure defined on  $\sigma$  algebra of Lebesgue measurable sets 1. assigns length to any interval, 2. is translation invariant, and 3. is countably additive.*

#### *Proof.* Length of interval

Let  $E$  be an interval, then  $E$  belongs to the  $\sigma$  algebra of Lebesgue measurable sets as Borel sets are Lebesgue measurable. Also, we have  $m^*(E) = m(E)$  for any Lebesgue measurable set  $E$ . And Lebesgue outer measure  $m^*$  of an interval

is its length. Therefore, Lebesgue measure of any interval is its length.

### Translation invariant

Let  $E$  be a Lebesgue measurable set. We have,  $E + y$  is also Lebesgue measurable. Since  $E + y$  is Lebesgue measurable and Lebesgue outer measure is translation invariant, we have  $m^*(E) = m^*(E + y) = m(E + y)$ . Clearly,  $m(E) = m(E + y)$ .

### Countably additive <sup>6</sup>

Let  $\{E_k\}_{k=1}^{\infty}$  be a countable family of disjoint, Lebesgue measurable sets. Lebesgue outer measure is countably subadditive and countable union of Lebesgue measurable sets is also Lebesgue measurable. Thus we have,

$$m\left(\bigcup_{k=1}^{\infty} E_k\right) \leq \sum_{k=1}^{\infty} m(E_k)$$

Since Lebesgue measure is finitely additive we have,

$$\begin{aligned} m\left(\bigcup_{k=1}^{\infty} E_k\right) &\geq m\left(\bigcup_{k=1}^n E_k\right) = \sum_{k=1}^n m(E_k) \\ \lim_{n \rightarrow \infty} m\left(\bigcup_{k=1}^n E_k\right) &\geq \lim_{n \rightarrow \infty} \sum_{k=1}^n m(E_k) = \sum_{k=1}^{\infty} m(E_k) \\ &\implies m\left(\bigcup_{k=1}^{\infty} E_k\right) \geq \sum_{k=1}^{\infty} m(E_k) \end{aligned}$$

Therefore, Lebesgue measure is countably additive.

$$m\left(\bigcup_{k=1}^{\infty} E_k\right) = \sum_{k=1}^{\infty} m(E_k)$$

□

### Continuity of Lebesgue measure

**Theorem 10.2.4** (continuity). *Let  $m$  be Lebesgue measure.*

1. Suppose  $\{A_k\}_{k=1}^{\infty}$  be an ascending <sup>7</sup> collection of Lebesgue measurable sets. Then,

$$m\left(\bigcup_{k=1}^{\infty} A_k\right) = \lim_{k \rightarrow \infty} m(A_k) \quad (10.25)$$

2. Suppose  $\{B_k\}_{k=1}^{\infty}$  be a descending <sup>8</sup> collection of Lebesgue measurable sets and  $m(B_1) < \infty$ . Then,

$$m\left(\bigcap_{k=1}^{\infty} B_k\right) = \lim_{k \rightarrow \infty} m(B_k) \quad (10.26)$$

<sup>6</sup>A real-valued set function  $m$  is countably additive if  $m(\bigcup_{k=1}^{\infty} E_k) = \sum_{k=1}^{\infty} m(E_k)$  for any disjoint family of sets  $\{E_k\}$ .

<sup>7</sup> $\{A_k\}$  is ascending if  $A_1 \subset A_2 \subset \dots$

<sup>8</sup> $\{A_k\}$  is descending if  $B_1 \supset B_2 \supset \dots$



**Proof. Ascending Collection**

Let  $\{A_k\}_{k=1}^{\infty}$  be an ascending collection of Lebesgue measurable sets. Define  $A_0 = \phi$ .

**Case 1 :**  $\exists k' \in \mathbb{N}$ ,  $m(A_{k'}) = \infty$

Suppose the collection has a Lebesgue measurable set  $A_{k'}$  of infinite measure. Then for  $\forall k \geq k'$ ,  $m(A_k) = \infty$ . Clearly,

$$\lim_{k \rightarrow \infty} m(A_k) = \infty = m(A_{k'}) \leq m\left(\bigcup_{k=1}^{\infty} A_k\right)$$

**Case 2 :**  $\forall k \in \mathbb{N}$ ,  $m(A_k) < \infty$

Suppose that every Lebesgue measurable set in the collection is of finite measure. Consider the ascending collection of disjoint Lebesgue measurable sets,  $\{C_k\}_{k=1}^{\infty}$  given by  $C_k = A_k \sim A_{k-1}$ . Clearly,  $\bigcup_{k=1}^{\infty} A_k = \bigcup_{k=1}^{\infty} C_k$ . By countable additivity of Lebesgue measure, we have

$$m\left(\bigcup_{k=1}^{\infty} A_k\right) = m\left(\bigcup_{k=1}^{\infty} C_k\right) = \sum_{k=1}^{\infty} m(A_k \sim A_{k-1})$$

Also we have,

$$\begin{aligned} \sum_{k=1}^{\infty} m(A_k \sim A_{k-1}) &= \sum_{k=1}^{\infty} m(A_k) - m(A_{k-1}) \\ &= \lim_{n \rightarrow \infty} \sum_{k=1}^n m(A_k) - m(A_{k-1}) \\ &= \lim_{n \rightarrow \infty} m(A_n) - m(A_0) \end{aligned}$$

Therefore,  $m\left(\bigcup_{k=1}^{\infty} A_k\right) = \lim_{n \rightarrow \infty} m(A_n)$ .

**Descending Collection** Let  $\{B_k\}$  be a descending collection of Lebesgue measurable sets and  $m(B_1) < \infty$ . Consider the ascending collection of Lebesgue measurable sets,  $\{D_k\}_{k=1}^{\infty}$  given by  $D_k = B_1 \sim B_k$ . By the continuity of Lebesgue measure for ascending collection of sets, we have

$$m\left(\bigcup_{k=1}^{\infty} D_k\right) = \lim_{n \rightarrow \infty} m(D_n) = m(B_1) - \lim_{n \rightarrow \infty} m(B_n)$$

By de Morgan's law, we have  $B_1 - \bigcap_k B_k = \bigcup_k (B_1 - B_k)$ . Since  $B_1$  is of finite measure, by excision property we have,

$$m\left(\bigcap_{k=1}^{\infty} B_k\right) = m\left(B_1 - \bigcup_{k=1}^{\infty} D_k\right) = \lim_{n \rightarrow \infty} m(B_n)$$

□

**Borel-Cantelli Lemma**

**Definitions 10.2.5** (ae). A property of real numbers is true except for a set of zero measure, then it is true **almost everywhere**.

**Lemma 10.2.5** (Borel-Cantelli). *Let  $\{E_k\}_{k=1}^{\infty}$  be a countable collection of Lebesgue measurable sets for which  $\sum_{k=1}^{\infty} m(E_k) < \infty$ . Then, almost all  $x \in \mathbb{R}$  belongs to at most finitely many of the  $E_k$ s.*

*Proof.* We have  $\sum_{k=1}^{\infty} m(E_k) < \infty$ . Then, by convergence of the series

$$\lim_{n \rightarrow \infty} \sum_{k=n}^{\infty} m(E_k) = 0$$

Also,  $\left\{ \bigcup_{k=n}^{\infty} E_k \right\}_{n=1}^{\infty}$  is a descending collection of Lebesgue measurable sets. By continuity of Lebesgue measure for descending collection of Lebesgue measurable sets, we have

$$m \left( \bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} E_k \right) = \lim_{n \rightarrow \infty} m \left( \bigcup_{k=n}^{\infty} E_k \right) = 0$$

Clearly,  $\lim_{n \rightarrow \infty} \bigcup_{k=n}^{\infty} E_k$  is a set of zero measure. Suppose  $x \in \mathbb{R}$  belongs to countably many  $E_k$ s. Then, for any  $m \in \mathbb{N}$ , there exists  $k > m$  such that  $x \in E_k$ . Clearly,  $x \in \lim_{n \rightarrow \infty} \bigcup_{k=n}^{\infty} E_k$ . That is,  $x$  belongs to a set of zero measure. Therefore by contraposition, if  $x$  does not belong to a set of measure zero, then  $x$  belongs to at most finitely many  $E_k$ s. In other words, almost every  $x$  in  $\mathbb{R}$  belongs to at most finitely many  $E_k$ s.  $\square$

**Exercise**

24.  $m(E_1 \cup E_2) + m(E_1 \cap E_2) = m(E_1) + m(E_2)$
25.  $m(B_1) < \infty$  is necessary for continuity property of measure for descending collection of measurable sets.
26.
$$m^* \left( A \cap \bigcup_{k=1}^{\infty} E_k \right) = \sum_{k=1}^{\infty} m^*(A \cap E_k)$$
27. Let  $m'$  be set function on a  $\sigma$ -algebra and  $m'$  is countably additive.
  - (a)  $m'$  is finitely additive, monotone, countably monotone, and has excision property
  - (b)  $m'$  has continuity properties
28. continuity + finite additivity  $\implies$  countable additivity

### 10.2.6 Non-measurable sets

Every measurable set of positive measure contains a non-measurable set.

**Definitions 10.2.6** (Rational Equivalence). Let  $E$  be any subset of  $\mathbb{R}$ . The relation  $xRy \iff x - y \in \mathbb{Q}$  is an equivalence<sup>†9</sup> relation on  $\mathbb{R}$ .

**Definitions 10.2.7** (Choice set,  $C_E$ ). Let  $E$  be any subset of  $\mathbb{R}$  and  $R$  be an equivalence relation on  $E$ . By axiom of choice, there exists a choice set  $C_E \subset E$  containing an exactly an element from each equivalence class.

**Definitions 10.2.8** (translate). Let  $E$  be a subset of  $\mathbb{R}$ . Let  $\lambda \in \mathbb{R}$ . Then  $\lambda + E = \{\lambda + x : x \in E\}$  is a translate of  $E$ .

With the help of following lemma, we prove that for any measurable set  $E$  of positive measure, the subset  $C_E$  is non-measurable.

**Lemma 10.2.6.** *Let  $E$  be a bounded, measurable set. Suppose there exists a bounded, countably infinite set  $\Lambda$  for which the collection of translates of  $E$  under  $\lambda$ ,  $\{\lambda + E\}_{\lambda \in \Lambda}$  are disjoint. Then  $m(E) = 0$ .*

In other words, if a bounded measurable set has countably many disjoint translates, then it is of measure zero. That is, there doesn't exist a bounded set of positive measure which has countably many disjoint translates.

*Proof.* The Lebesgue measure is translation invariant. Thus, the translates,  $\lambda + E$  are measurable and  $m(\lambda + E) = m(E)$ .

We have,  $E$  and  $\Lambda$  are bounded,  $\bigcup_{\lambda \in \Lambda} \lambda + E$  is bounded and is of finite measure. The Lebesgue measure is countably additive. And since translates are disjoint,

$$m \left[ \bigcup_{\lambda \in \Lambda} \lambda + E \right] = \sum_{\lambda \in \Lambda} m(\lambda + E) < \infty$$

Clearly,  $m(E) = 0$ . Suppose  $m(E) = \epsilon$ . Then  $\sum_{\lambda \in \Lambda} m(\lambda + E) = \sum_{\lambda \in \Lambda} \epsilon = \infty$  since  $\Lambda$  has countably infinite elements.  $\square$

**Theorem 10.2.7** (Vitali). *Any set  $E$  of real numbers with positive outer measure contains a subset that fails to be measurable.*

More importantly, every measurable set of positive measure contains a non-measurable set.

*Proof. Case 1 :  $E$  is bounded and non-measurable*

Suppose  $E$  is not measurable. Then  $E \subset \mathbb{R}$  is a non-measurable set and the result is trivial.

**Case 2 :  $E$  is bounded and measurable**

Suppose  $E$  is a bounded, measurable subset of positive measure. Let  $C_E$  be a

---

<sup>9</sup>  $x - x = 0 \in \mathbb{Q}$ ,  $(y - x) = -(x - y) \in \mathbb{Q}$  and  $x - z = (x - y) + (y - z) \in \mathbb{Q}$

choice set of  $E$  under rational equivalence. Let  $\Lambda$  be any bounded, countably infinite set of rational numbers. Clearly, the translates of  $C_E$  under  $\Lambda$  are disjoint.

Suppose  $x \in (\lambda_1 + C_E) \cap (\lambda_2 + C_E)$ . Then,  $x = \lambda_1 + y = \lambda_2 + z$ . Clearly,  $y = z$  since  $y - z \in \mathbb{Q}$  and we chose precisely one element from each equivalence class. Again,  $y = z \implies \lambda_1 = \lambda_2$ . In other words, the intersecting the translates are identical. That is, two distinct translates will be disjoint.

Suppose  $C_E$  is measurable. Since  $C_E$  and  $\Lambda$  are bounded and  $C_E$  has countably many disjoint translates under  $\Lambda$ , by lemma  $m(C_E) = 0$ . However,  $E \subset \bigcup_{\lambda \in [-2b, 2b] \cap \mathbb{Q}} \lambda + C_E$  for sufficiently large<sup>10</sup>  $b \in \mathbb{R}$ .

$$m(E) \leq m \left( \bigcup_{\lambda \in [-2b, 2b] \cap \mathbb{Q}} \lambda + C_E \right) = \sum_{\lambda \in [-2b, 2b] \cap \mathbb{Q}} m(\lambda + C_E) = 0$$

which is a contradiction since  $E$  has positive measure.

### Case 3 : $E$ is unbounded and measurable

Suppose  $E$  is an unbounded subset of positive outer measure, then by the definition of Lebesgue outer measure,  $E$  has a bounded subset of positive outer measure. And by case 1 & 2, this set has a non-measurable subset.  $\square$

**Theorem 10.2.8.** *There are disjoint subsets  $A, B$  of real numbers for which*

$$m^*(A \cup B) < m^*(A) + m^*(B) \quad (10.27)$$

*Proof.* Suppose that for every disjoint pair of subsets  $A, B \subset \mathbb{R}$ ,  $m^*(A \cup B) = m^*(A) + m^*(B)$ . Then, by the definition of Lebesgue measurability, every subset of real numbers is measurable. By, Vital's theorem there does exist non-measurable subsets of real numbers which is a contradiction. Therefore, there does exist disjoint subsets  $A, B$  such that  $m^*(A \cup B) \neq m^*(A) + m^*(B)$ . By subadditivity of Lebesgue outer measure, we have  $m^*(A \cup B) \leq m^*(A) + m^*(B)$ . Therefore,  $m^*(A \cup B) < m^*(A) + m^*(B)$ .  $\square$

### Exercise

29. (a)
- (b) Rational Equivalence on  $\mathbb{Q}$  gives singleton choice set as difference two rational numbers is always rational.  $\frac{a}{b} - \frac{c}{d} = \frac{ad-bc}{bd}$ . Thus,  $\{0\}$  is a choice set.
- (c) Difference two numbers being irrational is not an equivalence relation as it violates transitivity.  $x-y, y-z \notin \mathbb{Q} \not\Rightarrow x-z \notin \mathbb{Q}$ . For example,  $\sqrt{2} - \sqrt{3}, \sqrt{3} - \sqrt{2} \notin \mathbb{Q}$ . However,  $\sqrt{2} - \sqrt{2} = 0 \in \mathbb{Q}$ .

30.

31.

32.

33.

---

<sup>10</sup>Since  $E$  is bounded there exists  $b \in \mathbb{R}$  such that  $E \subset [-b, b]$

### 10.2.7 Cantor Set and Cantor-Lebesgue Function

#### Cantor set, $C$

We know that every countable set is of measure zero. However, subsets of zero measure are not necessarily countable. Cantor set is an uncountable set of zero measure.

**Definitions 10.2.9** (Cantor set). Consider unit interval,  $I = [0, 1]$ . Let  $C_0 = I$ . Let  $\{I_k\}$  be the collection of subintervals in  $C_k$ . Construct  $C_{k+1}$  recursively by removing subintervals of length  $\frac{l(I_k)}{3}$  from the middle of each  $I_k$  in  $C_k$ . Cantor set  $C$  is given by,

$$C = \bigcap_{k=1}^{\infty} C_k \quad (10.28)$$

Note that  $C_k$  are descending collection of  $2^k$  disjoint, closed intervals, each of length  $\frac{1}{3^k}$ . Thus, effective length of  $C_k$  is  $(\frac{2}{3})^k$ .

**Theorem 10.2.9.** *Cantor set  $C$  is a closed, uncountable set of measure zero.*

*Proof.* **Step 1 :  $C$  is measurable**

By construction, Cantor set is an intersection of (countably many) closed subsets of real numbers. Since intersections of closed sets are closed, Cantor set is also closed. Also every closed subset is a Borel set and therefore measurable.

**Step 2 :  $m(C) = 0$**

From the construction of  $C$ , we have  $m(C_k) = (\frac{2}{3})^k$ . Clearly,  $\{C_k\}_{k=1}^{\infty}$  is a descending collection of measurable sets and  $m(C_1) < \infty$ . Thus by continuity of Lebesgue measure,

$$m(C) = m\left(\bigcap_{k=1}^{\infty} C_k\right) = \lim_{k \rightarrow \infty} m(C_k) = \lim_{k \rightarrow \infty} \left(\frac{2}{3}\right)^k = 0$$

**Step 3 :  $C$  is uncountable**

Suppose  $C$  is countable. Then elements of  $C$  can be enumerated. That is,  $C = \{x_k\}_{k=1}^{\infty}$ . We have,  $x_1 \in C \implies x_1 \in C_1$ . There are two disjoint intervals in  $C_1$ . Clearly,  $C_1$  has an interval  $F_1$  which doesn't contain  $x_1$ . Similarly, there exists a closed interval,  $F_2$  in  $C_2$  such that  $x_2 \notin F_2$  and  $F_2 \subset F_1$ . Continuing like this, we get a descending collection of closed intervals  $\{F_k\}_{k=1}^{\infty}$  such that  $x_k \notin F_k$ .

By nested set theorem, intersection of descending collection of closed and bounded intervals is non-empty. Thus,  $\bigcap_{k=1}^{\infty} F_k \neq \phi$ . Let  $x \in \bigcap_{k=1}^{\infty} F_k$ . Clearly  $x \neq c_k$  for any  $k$  as  $c_k \notin F_k \implies c_k \notin \bigcap_{k=1}^{\infty} F_k$ . However,  $x \in C$  which is a contradiction to the assumption that  $C = \{c_k : k = 1, 2, \dots\}$ . Therefore, elements of  $C$  can't be enumerated. In other words,  $C$  is uncountable.  $\square$

**Cantor-Lebesgue Function,  $\varphi$** 

**Definitions 10.2.10** (increasing). A real-valued function  $f$  is increasing if  $f(x) \geq f(y)$  for every  $x > y$ .

**Definitions 10.2.11** (strictly increasing). A real-valued function  $f$  is strictly increasing if  $f(x) > f(y)$  for every  $x > y$ .

**Definitions 10.2.12** (Cantor-Lebesgue function  $\varphi$ ). Let  $C$  be the Cantor set. Define open set,  $\mathcal{O} = I \sim C$ .

$$\mathcal{O} = I \sim C = I \sim \left( \bigcap_{k=1}^{\infty} C_k \right) = \bigcup_{k=1}^{\infty} (I \sim C_k) = \bigcup_{k=1}^{\infty} \mathcal{O}_k$$

We know that  $\mathcal{O}_k$  has  $2^{k-1}$  open intervals  $I_{k,1}, I_{k,2}, \dots, I_{k,2^{k-1}}$ . Define Cantor-Lebesgue function  $\varphi$  on  $\mathcal{O}$  by  $\varphi(x) = \frac{m}{2^k}$  for every  $x \in I_{k,m}$ . We may extend  $\varphi$  to unit interval such that  $\varphi(0) = 0$  and  $\varphi(x) = \sup \{ \varphi(t) : t \in [0, x] \cap \mathcal{O} \}$ .

Clearly, by the construction  $\phi$  is an increasing real-valued function.

**Theorem 10.2.10** (Properties of  $\varphi$ ). *Cantor-Lebesgue function  $\varphi : I \rightarrow I$  is a surjection. And  $\varphi$  is differentiable in  $\mathcal{O}$ .*

*Proof.* We know from the definition of  $\varphi$  on  $\mathcal{O}$ , it is increasing on  $\mathcal{O}$ . And for extending  $\varphi$  from  $\mathcal{O}$  to  $[0, 1]$  we are considering the supremum (least upper bound) of all the previous values of  $\varphi$  on  $\mathcal{O}$ . Thus, function  $\varphi$  is increasing on the unit interval,  $[0, 1]$ .

For continuity of  $\varphi$ , it is enough to prove that  $\varphi$  doesn't have any jump discontinuities<sup>11</sup> as it is an increasing function. Suppose  $x \in C$  and  $x \neq 0$ . Then for sufficiently large  $k$ ,  $x$  lies between two consecutive intervals of  $\mathcal{O}_k$ . Let  $a_k$  and  $b_k$  be the upper and lower bounds of intervals to the left and right of  $x$  in  $\mathcal{O}_k$ . Then  $x \in (a_k, b_k)$ .

We know that,  $\varphi(b_k) - \varphi(a_k) = \frac{1}{2^k}$ . And  $\varphi(a_k) < \varphi(x) < \varphi(b_k)$ , by the construction of  $\varphi$ . As  $k \rightarrow \infty$ , the jump  $\lim_{k \rightarrow \infty} \varphi(b_k) - \varphi(a_k) = \lim_{k \rightarrow \infty} \frac{1}{2^k} = 0$ . Thus,  $\varphi$  doesn't have any jump discontinuities. Therefore,  $\varphi$  is continuous on  $[0, 1]$ .

Clearly,  $\varphi$  is constant on each interval of  $\mathcal{O}$ . Therefore, its derivative exists and equals to zero everywhere on  $\mathcal{O}$ .

We have,  $\mathcal{O} = [0, 1] \sim C$ . Thus,  $m(\mathcal{O}) = 1$ , since  $m(C) = 0$  and  $m([0, 1]) = 1$ . Also,  $\varphi$  is an increasing, continuous function from  $[0, 1]$  into  $[0, 1]$  where  $\varphi(0) = 0$  and  $\varphi(1) = 1$ . Therefore, Cantor-Lebesgue function,  $\varphi$  is onto unit interval  $[0, 1]$ .  $\square$

<sup>11</sup>If the function  $\varphi$  has unbounded variation (jump) at some points of  $[0, 1]$ . That is  $\exists \epsilon > 0$  such that  $\forall \delta > 0, \exists x \in [0, 1]$  where  $f(x + \epsilon/2) - f(x - \epsilon/2) > \delta$ . —need revisit—

**Measurable, non-Borel set**

We already saw that there exists non-measurable subsets. Using the properties of the function  $\psi$ , we assert that there exist a measurable set  $\psi^{-1}(W)$  which is not a Borel set.

**Theorem 10.2.11.** *The function  $\psi : [0, 1] \rightarrow [0, 2]$  defined by  $\psi(x) = \varphi(x) + x$*

1. *is a strictly increasing<sup>12</sup> function which maps  $[0, 1]$  onto  $[0, 2]$  and*
2.  *$\psi$  maps Cantor set,  $C$  onto a set of positive measure and*
3.  *$\psi$  maps a measurable<sup>13</sup> subset of  $C$  onto a non-measurable set.*

*Proof.* By definition,  $\psi(x) = \varphi(x) + id(x)$ . We know that the sum of continuous functions is continuous. Thus,  $\psi$  is continuous. Again, the sum of an increasing function and strictly increasing function is always strictly increasing. Thus,  $\psi$  is a strictly increasing function.

We know that  $\psi$  is only translating open intervals in  $\mathcal{O}$ . That is,  $\psi(I_{k,m}) = I_{k,m} + \frac{m}{2^k}$ . In other words,  $\psi$  translates every interval  $I_k$  in  $\mathcal{O}$  into an interval of same length since  $\varphi$  is constant in each subinterval of  $\mathcal{O}$ . Thus,

$$m(\psi(\mathcal{O})) = \sum_{k=1}^{\infty} l(\psi(I_k)) = \sum_{k=1}^{\infty} l(I_k) = m(\mathcal{O}) = 1$$

We know that  $[0, 2] = \psi(\mathcal{O}) \cup \psi(C)$ . Since  $m(\psi(\mathcal{O})) = 1$ ,  $m(\psi(C)) = 2 - 1 = 1$ .

Clearly,  $\psi(C)$  is a subset of positive measure. Thus, by Vitali's theorem  $\psi(C)$  has a subset  $W$  which is not measurable. However,  $\psi^{-1}(W)$  is a subset of  $C$  which is of measure zero. Therefore,  $C$  has a measurable subset which  $\psi$  maps onto a non-measurable set.  $\square$

**Theorem 10.2.12.** *Cantor set has a measurable subset which is not a Borel set.*

*Proof.* By Vitali's theorem,  $\psi(C)$  has a non-measurable subset, say  $W$ . Clearly  $W$  is not a Borel set, since every Borel set is measurable.

Now  $\psi : \psi^{-1}(W) \rightarrow W$  is a bijection. And  $W \subset \psi(C) \implies \psi^{-1}(W) \subset C$ . Also we know that, every subset of zero measure sets are measurable. Thus,  $\psi^{-1}(W)$  is measurable subset of Cantor set,  $C$ .

Suppose  $\psi^{-1}(W)$  is a Borel set. Then,  $W$  must be a Borel set, since continuous functions maps Borel sets onto Borel sets. This leads to a contradiction since  $W$  is not a Borel set. Therefore, Cantor set  $C$  has a measurable, non-Borel subset  $\psi^{-1}(W)$ .  $\square$

<sup>12</sup>Strictly increasing functions are one-to-one. Thus,  $\psi$  is a bijection.

<sup>13</sup>Cantor set is a subset of zero measure. Thus every subset of Cantor set is measurable.

## 10.3 Lebesgue Measurable Functions

### 10.3.1 Sum, Products, and Compositions

**Theorem 10.3.1.** *Let function  $f$  have a measurable domain. Then the following statements are equivalent :*

1.  $\forall c \in \mathbb{R}, \{x \in E : f(x) > c\}$  is measurable
2.  $\forall c \in \mathbb{R}, \{x \in E : f(x) \geq c\}$  is measurable
3.  $\forall c \in \mathbb{R}, \{x \in E : f(x) < c\}$  is measurable
4.  $\forall c \in \mathbb{R}, \{x \in E : f(x) \leq c\}$  is measurable

And each statement above imply that  $\{x \in E : f(x) = c\}$  is measurable.

*Proof.* **Step 1 :**  $1 \iff 4$  and  $2 \iff 3$

The sets considered in statements 1 and 4 are complementary. Thus, if one set is measurable, then the other set is also measurable. Similarly, the sets in 2 and 3 are complementary. Thus, one statement implies the other.

**Step 2 :**  $1 \implies 2$

Suppose  $\{x \in E : f(x) > c\}$  is measurable for any real number  $c$ . Clearly, the sets  $\{x \in E : f(x) > c - \frac{1}{k}\}$  is measurable for every natural number  $k$ . And countable intersection of measurable sets is measurable. Thus,

$$\bigcap_{k=1}^{\infty} \left\{ x \in E : f(x) > c - \frac{1}{k} \right\} = \{x \in E : f(x) \geq c\} \text{ is measurable}$$

**Step 3 :**  $2 \implies 1$  Suppose  $\{x \in E : f(x) \geq c\}$  is measurable for any real number  $c$ . Clearly, the sets  $\{x \in E : f(x) \geq c + \frac{1}{k}\}$  is measurable for every natural number  $k$ . Again, countable union of measurable sets is measurable. Thus,

$$\bigcup_{k=1}^{\infty} \left\{ x \in E : f(x) \geq c + \frac{1}{k} \right\} = \{x \in E : f(x) > c\} \text{ is measurable}$$

**Step 4 :**  $\forall c \in \mathbb{R}, \{x \in E : f(x) = c\}$  is measurable

Suppose one of the statements is true. Then other three statements are also true. The sets  $\{x \in E : f(x) \geq c\}$  and  $\{x \in E : f(x) \leq c\}$  are both measurable. Thus, their intersection  $\{x \in E : f(x) = c\}$  is also measurable.  $\square$

**Definitions 10.3.1** (measurable function). A function  $f$  is measurable if

1. the domain of  $f$  is measurable and
2. one of the following statements is true
  - (a)  $\forall c \in \mathbb{R}, \{x \in E : f(x) > c\}$  is measurable
  - (b)  $\forall c \in \mathbb{R}, \{x \in E : f(x) \geq c\}$  is measurable
  - (c)  $\forall c \in \mathbb{R}, \{x \in E : f(x) < c\}$  is measurable
  - (d)  $\forall c \in \mathbb{R}, \{x \in E : f(x) \leq c\}$  is measurable

Note : If  $f$  is a measurable function, then its domain is measurable and  $f^{-1}(c)$  is measurable for any real number  $c$ . (why ?)



### Study of measurable functions

#### 1. Inverse images of open sets measurable

Suppose function  $f$  is defined on a measurable set  $E$ . Then  $f$  is a measurable function  $\iff \forall$  open set  $\mathcal{O}$ ,  $f^{-1}(\mathcal{O})$  is measurable.

*Proof.* Given the domain of  $f$  is measurable.

##### Step 1 : $f^{-1}(\mathcal{O})$ is measurable $\implies f$ is measurable

Suppose inverse image of every open set is measurable. Then,  $f^{-1}(c, \infty) = \{x \in E : f(x) > c\}$  is measurable. Then, by the definition,  $f$  is a measurable function.

##### Step 2 : $f$ measurable $\implies f^{-1}(\mathcal{O})$ is measurable

Suppose function  $f$  is measurable. Let  $\mathcal{O}$  be an open set. Then it can be expressed as a countable union of open, bounded intervals,  $I_k$ . That is,  $\mathcal{O} = \bigcup_{k=1}^{\infty} I_k$ . We know that,  $\forall k \in \mathbb{N}$ ,  $I_k = (a_k, b_k) = (-\infty, b_k) \cap (a_k, \infty)$ . By the definition of measurability  $\{x \in E : f(x) < b_k\} = f^{-1}(-\infty, b_k)$  and  $\{x \in E : f(x) > a_k\} = f^{-1}(a_k, \infty)$  are measurable. And  $f^{-1}(a_k, b_k) = f^{-1}(-\infty, b_k) \cap f^{-1}(a_k, \infty)$  is measurable. Again, inverse image of the open set  $f^{-1}(\mathcal{O})$  which is a countable union of measurable sets  $\{f^{-1}(I_k)\}_{k=1}^{\infty}$  is also measurable. Therefore, inverse image of every open set is measurable.  $\square$

#### 2. Continuous, real-valued functions are measurable

Suppose the function  $f$  has a measurable domain. And  $f$  is a real-valued, continuous function. Then  $f$  is measurable.

*Proof.* Let  $E$  be a measurable set. And  $f$  be a continuous function on  $E$ . Let  $\mathcal{O}$  be an open set, then by open set characterisation of continuous function we have  $f^{-1}(\mathcal{O}) = E \cap U$  where  $U$  is an open set. Clearly,  $f^{-1}(\mathcal{O})$  is a union of two measurable sets and thus measurable. Therefore,  $f$  is measurable since  $\forall$  open set  $\mathcal{O}$ ,  $f^{-1}(\mathcal{O})$  is measurable.  $\square$

#### 3. Monotone function on an interval is measurable

*Proof.* Without loss of generality, suppose that  $f$  is an increasing function defined on an interval  $I$ . Since, every interval is measurable,  $f$  is defined on a measurable set.

Consider  $\{g_n\}_{n=1}^{\infty}$  defined by  $g_n(x) = f(x) + \frac{x}{n}$ . Since  $g_n$  are strictly increasing functions,  $g_n^{-1}(c, \infty) = I \cap U$  where  $U$  is an interval. Thus,  $g_n^{-1}(c, \infty) = \{x \in I : g_n(x) > c\}$  is always measurable. Therefore,  $\{g_n\}$  is a family of measurable functions.

Now, from the construction of  $g_n$ , we have  $\lim_{n \rightarrow \infty} g_n = f$ . We have,  $\{g_n\}$  is a sequence of measurable functions converging pointwise to the limit function  $f$  (a.e.) on the interval  $I$ . Therefore,  $f$  is measurable.<sup>†14</sup>  $\square$

<sup>†14</sup>Limit of measurable functions under pointwise convergence(a.e.) is measurable.

We will prove this result in the upcoming subsection .

4.  **$f$  measurable and  $f = g$  (a.e.)  $\implies g$  measurable**

*Proof.* Suppose  $f, g$  are functions defined on a measurable set  $E$ . Suppose  $f$  is measurable. Let  $A \subset E$  such that  $A = \{x \in \mathbb{R} : f(x) \neq g(x)\}$ . We have  $f = g$  (a.e.). Thus,  $m(A) = 0$ . And we have,

$$\begin{aligned} \{x \in E : g(x) > c\} &= \{x \in A : g(x) > c\} \cup \{x \in E \sim A : g(x) > c\} \\ &= \{x \in A : g(x) > c\} \cup \{x \in E \sim A : f(x) > c\} \\ &= \{x \in A : g(x) > c\} \cup [\{x \in E : f(x) > c\} \cap (E \sim A)] \end{aligned}$$

Clearly,  $\{x \in A : g(x) > c\}$  is a subset of a set of measure zero and thus measurable. Also,  $E \sim A$  measurable since both  $E$  and  $A$  are measurable. And since  $f$  is measurable,  $\{x \in E : f(x) > c\}$  is measurable. Thus,  $\{x \in E : g(x) > c\}$  is measurable for any  $c \in \mathbb{R}$ . Therefore,  $g$  is a measurable function.  $\square$

5.  **$f$  measurable  $\iff f|_D, D \sim E$  measurable ( $\forall D$  measurable)**

Suppose function  $f$  is an extended real-valued function on a measurable set  $E$ . For measurable subset  $D$  of  $E$ ,  $f$  is measurable on  $E$  if and only if the restriction of  $f$  to  $D$  and  $D \sim E$  are measurable.

*Proof. Part 1 :  $f$  measurable  $\implies f|_D$  measurable*

Suppose  $f$  is a measurable function defined on  $E$  and  $D$  is a measurable subset of  $E$ . Since,  $f$  is measurable,  $E$  is measurable. Clearly  $E \sim D$  is measurable.

$$\begin{aligned} \{x \in D : f|_D(x) > c\} &= \{x \in D : f(x) > c\} \\ &= \{x \in E : f(x) > c\} \cap D \end{aligned}$$

Since  $f$  is measurable,  $\{x \in E : f(x) > c\}$  is measurable for any  $c \in \mathbb{R}$ . And intersection of measurable sets are measurable. Thus,  $\{x \in D : f|_D(x) > c\}$  is measurable for any  $c \in \mathbb{R}$ . Therefore,  $f|_D$  is a measurable function.  $\square$

**Definitions 10.3.2** (characteristic function). Let  $A$  be a subset of  $\mathbb{R}$ . The characteristic function  $\chi_A$  of  $A$  is given by

$$\chi_A(x) = \begin{cases} 1 & x \in A \\ 0 & x \notin A \end{cases} \quad (10.29)$$

**Definitions 10.3.3.** Let  $\{f_1, f_2, \dots, f_n\}$  be a finite family of measurable functions on the same domain  $E$ . Then  $\max\{f_1, f_2, \dots, f_n\}$  on  $E$  is given by

$$\max\{f_1, f_2, \dots, f_n\}(x) = \max\{f_1(x), f_2(x), \dots, f_n(x)\}, \quad \forall x \in E \quad (10.30)$$

And function  $\min\{f_1, f_2, \dots, f_n\}$  on  $E$  is given by

$$\min\{f_1, f_2, \dots, f_n\}(x) = \min\{f_1(x), f_2(x), \dots, f_n(x)\}, \quad \forall x \in E \quad (10.31)$$

### Properties of Measurable Functions

Suppose  $f, g$  are measurable functions on  $E$  and are finite a.e. on  $E$ .

1. **Linearity :**  $\alpha f + \beta g$  is measurable  $\forall \alpha, \beta$

*Proof. Step 1 :  $f$  measurable  $\implies \alpha f$  measurable*

Suppose  $\alpha = 0$ . Then  $\alpha f = 0$ . And 0 function is trivially measurable.

Suppose  $\alpha \neq 0$ . Then  $\{x \in E : f(x) > c/\alpha\} = \{x \in E : \alpha f(x) > c\}$  is measurable for any  $c \in \mathbb{R}$ . Therefore  $\alpha f$  is measurable for any  $\alpha \in \mathbb{R}$ .

In other words, if  $f$  is measurable, then  $\alpha f$  is measurable and if  $g$  is measurable, then  $\beta g$  is measurable for any  $\alpha, \beta \in \mathbb{R}$ . Therefore, it is enough to prove that  $f, g$  measurable  $\implies f + g$  measurable.

**Step 2 :  $f, g$  are measurable  $\implies f + g$  is measurable**

$$f(x) + g(x) < c \iff f(x) < c - g(x)$$

Since  $\mathbb{Q}$  is dense, there exists a rational number  $q \in \mathbb{Q}$  between any two distinct real numbers

$$\begin{aligned} f(x) + g(x) < c &\iff \exists q \in \mathbb{Q}, f(x) < q < c - g(x) \\ &\iff f(x) < q \text{ AND } g(x) < c - q \end{aligned}$$

$$\begin{aligned} \{x \in E : f + g(x) < c\} &= \{x \in E : f(x) + g(x) < c\} \\ &= \bigcup_{q \in \mathbb{Q}} [\{x \in E : f(x) < q\} \cap \{x \in E : g(x) < c - q\}] \end{aligned}$$

Since,  $f, g$  are measurable, each set in the union is measurable. Thus,  $\{x \in E : f + g(x) < c\}$  is measurable for every  $c \in \mathbb{R}$ , since rational numbers are countable, and countable union of measurable sets is measurable. Therefore,  $f + g$  is measurable.  $\square$

2. **Product :**  $fg$  is measurable

*Proof.* Suppose  $f, g$  are measurable functions **which are finite (a.e.)** on a measurable set  $E$ . And we have,

$$fg = \frac{1}{2} [(f + g)^2 - f^2 - g^2] \quad (10.32)$$

**Step 1 :  $f$  is measurable  $\implies f^2$  is measurable**

Suppose  $c \geq 0$ . If  $c < 0$ , then  $\{x \in E : f^2(x) > c\} = E$  is measurable.

$$\{x \in E : f^2(x) > c\} = \{x \in E : f(x) > \sqrt{c}\} \cup \{x \in E : f(x) < -\sqrt{c}\}$$

is a union of two measurable sets and is measurable.

**Step 2 :  $fg$  is measurable**

We have  $f, g$  are measurable. Thus,  $f + g$  is measurable, since linear combination of measurable sets is measurable. And,  $f^2, g^2, (f + g)^2$  are measurable by Step 1. We know that  $fg$  is a linear combination of these measurable sets and therefore  $fg$  is measurable.  $\square$

**3. Composition of measurable functions is not necessarily measurable**

We know that, Cantor set  $C$  has subset  $A = \psi^{-1}(W)$  such that  $\psi$  maps  $A$  onto a non-measurable set  $W$ .<sup>15</sup> Then  $A$  is a measurable subset of open interval  $(0, 1)$  and  $\psi(A)$  is non-measurable.

Let  $\chi_A$  be the characteristic function of  $A$ . Then  $\chi_A \circ \psi^{-1}$  is not measurable since  $\{x \in (0, 1) : \chi_A \circ \psi^{-1}(x) \geq 1\} = \psi(A)$  is not measurable. But, both the functions  $\chi_A$  and  $\psi^{-1}$  are measurable functions since  $\chi_A^{-1}(y)$  is either  $\phi, A$  or  $A^c$  and  $\psi^{-1}$  is a continuous real-valued function defined on  $(0, 1)$ .

**4. If  $f$  is continuous and  $g$  is measurable, then  $f \circ g$  is measurable**

*Proof.* Suppose  $f$  be a continuous function and  $g$  be a measurable function both defined on a common set  $E$  which is measurable. Since  $f$  is continuous, we know that  $f$  is measurable.

Let  $\mathcal{O}$  be an open set. Then  $(f \circ g)^{-1}(\mathcal{O}) = g^{-1}(f^{-1}(\mathcal{O}))$ . By characterisation of continuity,  $U = f^{-1}(\mathcal{O})$  is an open set since  $\mathcal{O}$  is open. And by characterisation of measurability of  $g$ , we have  $g^{-1}(U)$  is an open set since  $g$  is measurable and  $U$  is open. Now we have  $(f \circ g)^{-1}(\mathcal{O})$  is an open set for any open set  $\mathcal{O}$ . Therefore by the characterisation of measurability,  $f \circ g$  is measurable.  $\square$

**5.**

For a finite family of measurable functions  $\{f_1, f_2, \dots, f_n\}$  with common domain  $E$ , both the functions  $\max\{f_1, f_2, \dots, f_n\}$  and  $\min\{f_1, f_2, \dots, f_n\}$  are measurable.

*Proof.* Let  $\{f_1, f_2, \dots, f_n\}$  be a finite family of measurable functions defined on a common measurable set  $E$ . We have,

$$\{x \in E : \max\{f_1, f_2, \dots, f_n\}(x) > c\} = \bigcup_{k=1}^n \{x \in E : f_k(x) > c\}$$

is a finite union of measurable sets since each  $f_k$  is measurable. Therefore,  $\max\{f_1, f_2, \dots, f_n\}$  is measurable. Similarly, we have,

$$\{x \in E : \min\{f_1, f_2, \dots, f_n\}(x) < c\} = \bigcup_{k=1}^n \{x \in E : f_k(x) < c\}$$

Therefore,  $\min\{f_1, f_2, \dots, f_n\}$  is also measurable.  $\square$

---

<sup>15</sup>From the proof of Vitali's theorem and associated lemma,  $\psi(C)$  has such a subset  $W$  which is a choice set of  $\psi(C)$  under rational equivalence.

Note : Let  $f$  be a measurable function on a measurable set  $E$ . Then  $-f$  is measurable since

$$\{x \in E : -f(x) > c\} = \{x \in E : f(x) < -c\}$$

And  $|f|, f^+, f^-$  are also measurable since

$$|f| = \max\{f, -f\}, \quad f^+ = \max\{f, 0\}, \quad f^- = \max\{-f, 0\}$$

These functions  $f^+, f^-$  are quite important as we may write  $f = f^+ - f^-$  where both the functions  $f^+$  and  $f^-$  are non-negative. And Lebesgue integral for non-negative functions are sufficient to integrate any function  $f$ .

$$\int_E f = \int_E f^+ - \int_E f^-$$

### 10.3.2 Sequential Pointwise Limits and Simple Approximation

**Definitions 10.3.4** (pointwise convergence). A sequence of functions  $\{f_n\}$  on a common domain  $E$  convergence pointwise to  $f$  if

$$\forall x \in E, \quad \lim_{n \rightarrow \infty} f_n(x) = f(x)$$

**Definitions 10.3.5** (pointwise convergence(a.e.)). Let  $E_0$  be a set of zero measure. A sequence of functions  $\{f_n\}$  on a common domain  $E$  convergence pointwise a.e. on  $E_0 \subset E$  to  $f$  if

$$\forall x \in E \sim E_0, \quad \lim_{n \rightarrow \infty} f_n(x) = f(x)$$

**Definitions 10.3.6** (uniform convergence). A sequence of functions  $\{f_n\}$  on a common domain  $E$  convergence uniformly to  $f$  if

$$\forall \epsilon > 0, \exists N \in \mathbb{N}, \text{ such that } \forall n \geq N, |f_n - f| < \epsilon$$

**Theorem 10.3.2** (Pointwise Convergence(a.e.) preserves measurability). Let  $\{f_n\}$  be a sequence measurable functions on  $E$ . If  $\{f_n\}$  converges to  $f$  pointwise a.e. on  $E$ , then  $f$  is measurable.

*Proof.* Let  $\{f_n\}$  be a sequence of measurable functions converging to  $f$  pointwise a.e. on  $E$ . Then  $\{f_n\}_{n=1}^\infty$  converges pointwise on  $E \sim E_0$  where  $E_0$  is a set of measure zero. Without loss of generality, suppose that  $\{f_n\}$  converges pointwise everywhere on  $E$ .<sup>16</sup>

Let  $c \in \mathbb{R}$ . We claim that,

$$f(x) < c \iff \exists n, k \text{ such that } f_j(x) < c - \frac{1}{n}, \quad \forall j \geq k \quad (10.33)$$

Suppose that for every natural numbers  $n, k$ , there exists  $j \geq k$  for which  $f_j(x) \geq c - \frac{1}{n}$ . Since  $f(x)$  is the limit function,  $f(x) \not< c$  which is a contradiction. Thus, we may write,

$$\{x \in E : f(x) < c\} = \bigcup_{n,k=1}^{\infty} \left[ \bigcap_{j=k}^{\infty} \left\{ x \in E : f_j(x) < c - \frac{1}{n} \right\} \right] \quad (10.34)$$

---

<sup>16</sup>Consider  $E' = E \sim E_0$

Countable union of countable intersections of measurable sets is also measurable. Thus  $\{x \in E : f(x) < c\}$  is a measurable set for any  $c \in \mathbb{R}$ . Therefore,  $f$  is measurable.  $\square$

**Definitions 10.3.7** (simple function). A real-valued function  $\varphi$  on a measurable set  $E$  is simple if it is measurable and assumes at most finitely many values.

Note :  $xRy \iff \varphi(x) = \varphi(y)$  is an equivalence relation which partitions  $E$  into  $k$  disjoint subsets  $E$ .

Note : Suppose  $\varphi$  is a simple function. Then each equivalent class of  $E$  under above equivalence is also measurable by the definition of measurable function. (Why ?)

**Definitions 10.3.8** (canonical representation). The canonical representation of a simple function  $\varphi$  on  $E$  is given by,

$$\varphi = \sum_{k=1}^n c_k \cdot \chi_{E_k} \text{ where } E_k = \{x \in E : \varphi(x) = c_k\} \quad (10.35)$$

**Definitions 10.3.9** (bounded). Let  $f$  be a function on  $E$ . The function  $f$  is bounded if there exists  $m \in \mathbb{R}$ ,  $m \geq 0$  such that  $|f(x)| < m$  for every  $x \in E$ .

**Lemma 10.3.3** (simple approximation). Let  $f$  be a measurable, real-valued function which is bounded on  $E$ . Then for any  $\epsilon > 0$ , there exist simple functions  $\varphi_\epsilon$  and  $\psi_\epsilon$  approximating  $f$  (from below and above) such that  $0 \leq \psi_\epsilon - \varphi_\epsilon < \epsilon$  on  $E$ .

*Proof.* Let  $f$  be a bounded, measurable function on  $E$ . Then there exists open, bounded interval  $(c, d)$  such that  $f(E) \subset (c, d)$ . Let  $P = \{y_0, y_1, \dots, y_n\}$  be a partitions of open interval  $(c, d)$  such that  $y_0 = c < y_1 < \dots < y_{n-1} < y_n = d$  and  $y_k - y_{k-1} < \epsilon$  for every  $k$ . Define  $I_k = [y_{k-1}, y_k)$  and  $E_k = f^{-1}(I_k)$ . Then  $E_k$  are measurable since  $I_k$  are intervals and  $f$  is measurable. (Why ?)

Define  $\varphi_\epsilon = \sum_{k=1}^n y_{k-1} \cdot \chi_{E_k}$  and  $\psi_\epsilon = \sum_{k=1}^n y_k \cdot \chi_{E_k}$ . Then for any  $x \in E$ , we have  $\varphi_\epsilon(x) = y_{k-1} \leq f(x) \leq y_k = \psi_\epsilon(x)$ . Clearly,  $\varphi_\epsilon \leq \psi_\epsilon$  by the construction. And for each  $k$ , we have  $\psi_\epsilon(y) - \varphi_\epsilon(y) = y_k - y_{k-1} < \epsilon$  for any  $y \in [y_{k-1}, y_k)$ . Therefore,  $0 \leq \psi_\epsilon - \varphi_\epsilon < \epsilon$ .  $\square$

**Theorem 10.3.4** (simple approximation). An extended real-valued function  $f$  on a measurable set  $E$  is measurable if and only if there exists a sequence of simple functions  $\{\varphi_n\}$  converging to  $f$  pointwise on  $E$  and  $|\varphi_n| \leq |f|$  on  $E$  for every  $n \in \mathbb{N}$ .

And if  $f$  is non-negative, there exists an increasing sequence of functions  $\{\varphi_n\}$  with the same property.

**Proof. Part 1**

Suppose there exists a sequence of simple functions  $\{\varphi_n\}$  converging pointwise to  $f$  on  $E$ . Since simple functions are measurable and pointwise limit of measurable functions are measurable,  $f$  is measurable.

**Part 2**

Suppose  $f$  is measurable. Without loss of generality, suppose that  $f \geq 0$  on  $E$ . Otherwise,  $f = f^+ - f^-$  where  $f^+, f^- \geq 0$ . And  $f$  is measurable since both  $f^+$  and  $f^-$  are measurable.

**Step 1 : Construction of simple functions on  $E_n$** 

Let  $n \in \mathbb{N}$ . Define  $E_n = \{x \in E : f(x) \leq n\}$ . Since  $f$  is measurable,  $E_n$  is a measurable subset of  $E$ . And we know that  $f|_{E_n}$  is also a non-negative, bounded, measurable function.

Then by simple approximation lemma, for  $\epsilon = \frac{1}{n}$  there exists  $\varphi_n$  and  $\psi_n$  such that  $0 \leq \varphi_n \leq f \leq \psi_n$  and  $0 \leq \psi_n - \varphi_n < \frac{1}{n}$  on  $E_n$ .

**Step 2 : Extending simple function to  $E$** 

Extend  $\varphi_n$  to  $E$  by setting  $\varphi_n(x) = n$  if  $f(x) > n$ . Again,  $\varphi_n$  is a simple function on  $E$ . And  $0 \leq \varphi_n \leq f$  on  $E$ .

**Step 3 : Sequence of simple functions converging to  $f$** 

We claim that the sequence of simple functions  $\{\varphi_n\}$  converges to  $f$  pointwise on  $E$ . Suppose  $x \in E$ . Suppose  $f(x)$  is finite. Let  $m \in \mathbb{N}$ . Then  $0 \leq f(x) - \varphi_n(x) \leq \frac{1}{n}$  for all  $n \geq m$ . Therefore,  $\lim_{n \rightarrow \infty} \varphi_n(x) = f(x)$ . Suppose  $f(x) = \infty$ . Then  $\varphi_n(x) = n$ ,  $\forall n$ . Therefore,  $\lim_{n \rightarrow \infty} \varphi_n(x) = f(x)$ .

Replace each  $\varphi_n$  with  $\max\{\varphi_1, \varphi_2, \dots, \varphi_n\}$ <sup>17</sup>, we have an increasing sequence of simple functions  $\varphi_n$  converging to  $f$  pointwise on  $E$ .  $\square$

**Exercise**

- 12.
- 13.
- 14.
- 15.
- 16.
- 17.
- 18.
- 19.
- 20.
- 21.
- 22.
- 23.
- 24.

<sup>17</sup>We replace each simple function  $\varphi_k$  with the maximum of the finite subsequence from  $\varphi_1$  upto  $\varphi_k$ . That is,  $\varphi'_1 = \max\{\varphi_1\} = \varphi_1$ , and  $\varphi'_2 = \max\{\varphi_1, \varphi_2\}$ , and  $\varphi'_3 = \max\{\varphi_1, \varphi_2, \varphi_3\}$ , ...

## 10.4 Lebesgue Integration

### 10.4.1 Riemann integral

**Definitions 10.4.1** (partition). A set  $P = \{a, x_1, x_2, \dots, x_{n-1}, b\}$  is a partition of an interval  $[a, b]$  into  $n$  subintervals if  $a < x_1 < x_2 < \dots < x_{n-1} < b$ .

**Definitions 10.4.2** (lower/upper Riemann integral). Let  $f$  be a bounded, real-valued function on a closed interval  $[a, b]$  and  $P$  be a partition of  $[a, b]$ . Let  $M_i, m_i$  be the supremum and infimum of  $f$  in each subinterval of the partition. Then upper and lower Riemann integrals are the infimum of upper Darboux sums and supremum of lower Darboux sums.

Riemann upper integral,

$$(R) \int_a^b f = \inf_P \{U(f, P)\} = \inf \left\{ \sum_{i=1}^n M_i(x_i - x_{i-1}) \right\} \quad (10.36)$$

Riemann lower integral,

$$(R) \int_a^b f = \sup_P \{L(f, P)\} = \sup \left\{ \sum_{i=1}^n m_i(x_i - x_{i-1}) \right\} \quad (10.37)$$

**Definitions 10.4.3** (Riemann integrable). A function  $f$  is Riemann integrable over  $[a, b]$  if lower and upper Riemann integrals are equal. And Riemann integral of  $f$  over  $[a, b]$  is given by,

$$(R) \int_a^b f = (R) \overline{\int}_a^b f = (R) \underline{\int}_a^b f \quad (10.38)$$

**Definitions 10.4.4** (step function). A function  $\psi$  is a step function if it assumes constant values in each subinterval of some partition of its domain.

$$\psi(x) = \sum_{i=1}^n c_i \cdot \chi_{(x_i, x_{i-1})}$$

$$\text{Note : } (R) \int_a^b \psi = \sum_{i=1}^n c_i(x_i - x_{i-1})$$

**Definitions 10.4.5** (Dirchlet's function). Dirchlet's function is given by,

$$f(x) = \begin{cases} 1 & x \in \mathbb{Q} \\ 0 & x \notin \mathbb{Q} \end{cases}$$

Note : Dirchlet's function is not Riemann integrable(why ?). But it is Lebesgue integrable since  $\mathbb{Q}$  is Lebesgue measurable.



### 10.4.2 Lebesgue integral of a bounded, measurable function over a set of finite measure

**Definitions 10.4.6.** Let  $\psi$  be a simple function on a finite measure set  $E$ . Then,

$$\int_E \psi = \sum_{i=1}^n a_i \cdot m(E_i) \text{ where } \psi = \sum_{i=1}^n a_i \cdot \chi_{E_i} \quad (10.39)$$

Note : Step functions are simple functions. And the Riemann integral and Lebesgue integral of step functions are the same.

Note : Dirchlet's function is a simple function. And thus Lebesgue integrable.

**Lemma 10.4.1.** Let  $\{E_k\}_{k=1}^n$  be a disjoint collection of measurable subsets of a set  $E$  of finite measure. If  $\varphi = \sum_{k=1}^n a_k \cdot \chi_{E_k}$ , then  $\int_E \varphi = \sum_{k=1}^n a_k \cdot m(E_k)$

*Proof.* Suppose  $\varphi = \sum_{k=1}^n a_k \cdot \chi_{E_k}$  where  $E_k$  are of finite measure. We consider distinct values of  $a_i$ ,  $\{\lambda_1, \lambda_2, \dots, \lambda_m\}$ . Then, the function has the canonical representation,  $\varphi = \sum_{r=1}^m \lambda_r \cdot \chi_{A_r}$  where  $A_r = \bigcup_{\substack{k \\ a_k = \lambda_r}} E_k$ .

Then, we have  $\varphi^{-1}(\lambda_k) = A_k$  is a union of finite measurable set and thus measurable. Thus,  $\varphi$  is a measurable function which assumes at most  $m$  different values. And

$$\begin{aligned} \int_E \varphi &= \sum_{r=1}^m \lambda_r \cdot m(A_r) \\ &= \sum_{r=1}^m \lambda_r \cdot m\left(\bigcup_{\substack{k \\ a_k = \lambda_r}} E_k\right) \\ &= \sum_{r=1}^m \sum_{\substack{k \\ a_k = \lambda_r}} a_k \cdot m(E_k) \\ &= \sum_{k=1}^n a_k \cdot m(E_k) \end{aligned}$$

□

**Theorem 10.4.2** (linearity + monotonicity of integral of simple functions). Let  $\varphi$  and  $\psi$  be simple functions defined a set  $E$  of finite measure. Then for any  $\alpha, \beta \in \mathbb{R}$ ,

$$\int_E (\alpha\varphi + \beta\psi) = \alpha \int_E \varphi + \beta \int_E \psi \quad (10.40)$$

Moreover, if  $\varphi \leq \psi$  on  $E$ , then

$$\int_E \varphi \leq \int_E \psi \quad (10.41)$$

*Proof.* Suppose  $\varphi, \psi$  are simple functions and  $\alpha, \beta \in \mathbb{R}$ . Since  $\varphi, \psi$  are simple, there are  $n$  disjoint subset of  $E$  in which both  $\varphi$  and  $\psi$  are constant<sup>18</sup>. Then,

$$\begin{aligned} \int_E (\alpha\varphi + \beta\psi) &= \int_E \alpha \left( \sum_{i=1}^n a_i \cdot \chi_{E_i} \right) + \beta \left( \sum_{i=1}^n b_i \cdot \chi_{E_i} \right) \\ &= \int_E (\alpha a_i + \beta b_i) \cdot \chi_{E_i} \\ &= \sum_{i=1}^n (\alpha a_i + \beta b_i) \cdot m(E_i) \\ &= \alpha \sum_{i=1}^n a_i \cdot m(E_i) + \beta \sum_{i=1}^n b_i \cdot m(E_i) \\ &= \alpha \int_E \varphi + \beta \int_E \psi \end{aligned}$$

Suppose  $\varphi \leq \psi$ . Then  $\eta = \psi - \varphi$  is also a simple function. And  $0 \leq \eta$ . Then by linearity of Lebesgue integral of simple functions,

$$\begin{aligned} \int_E \eta &= \int_E \psi - \varphi \\ &= \int_E \psi - \int_E \varphi \end{aligned}$$

We have,  $\eta$  is a non-negative, measurable function. And thus,  $\int_E \eta \geq 0$ . Therefore,  $\int_E \varphi \leq \int_E \psi$ .  $\square$

### Lebesgue integral for bounded functions

Now we define Lebesgue integral of bounded functions, the same way we have constructed Riemann integral. We use simple approximations and the definition of Lebesgue integral for simple functions for this construction.

**Definitions 10.4.7** (upper/lower Lebesgue integral). Let  $f$  be a bounded, real-valued function on a set  $E$  of finite measure. Any simple approximation lemma, we have two families of simple functions  $\{\varphi_\epsilon\}$  and  $\{\psi_\epsilon\}$  such that  $\forall \epsilon > 0$ ,  $\varphi \leq f \leq \psi$  and  $0 \leq \psi - \varphi < \epsilon$ . Then

Upper Lebesgue integral,

$$\overline{\int_E} f = \inf_{\psi_\epsilon} \left\{ \int_E \psi_\epsilon \right\} \quad (10.42)$$

Lower Lebesgue integral,

$$\underline{\int_E} f = \sup_{\varphi_\epsilon} \left\{ \int_E \varphi_\epsilon \right\} \quad (10.43)$$

**Definitions 10.4.8** (Lebesgue integrability of bounded functions). Let  $f$  be a bounded, real-valued function defined on a set  $E$  of finite measure. Function  $f$  is Lebesgue integrable on  $E$  if both lower and upper Lebesgue integrals of  $f$  over  $E$  are equal.

---

<sup>18</sup>We don't assume that the values assume by  $\psi$  are distinct for distinct subsets.

And the Lebesgue integral of a bounded, real-valued function  $f$  over  $E$  is given by,

$$\int_E f = \sup_{\varphi} \int_E \varphi = \inf_{\psi} \int_E \psi \quad (10.44)$$

Now we can prove that Lebesgue integral is a generalisation of Riemann integral for bounded functions.

**Theorem 10.4.3.** *Let  $f$  be a bounded, real-valued function defined on a close, bounded interval  $[a, b]$ . If  $f$  is Riemann integrable over  $[a, b]$ , then it is Lebesgue integrable over  $[a, b]$ . And both Riemann integral and Lebesgue integral of  $f$  over  $[a, b]$  are equal.*

*Proof.* Suppose  $f$  is Riemann integrable over  $[a, b]$ . Then,

$$\sup \left\{ (R) \int_a^b \varphi : \varphi \text{ step}, \varphi \leq f \right\} = \inf \int \left\{ \int_a^b \psi : \psi \text{ step}, f \leq \psi \right\} \quad (10.45)$$

However, every step function is simple and we have

$$\sup \left\{ \int_a^b \varphi : \varphi \text{ simple}, \varphi \leq f \right\} = \inf \int \left\{ \int_a^b \psi : \psi \text{ simple}, f \leq \psi \right\} \quad (10.46)$$

Therefore, every bounded, real-valued function on closed, bounded interval is Lebesgue integrable function if it is Riemann integrable.  $\square$

**Theorem 10.4.4.** *Every bounded, measurable function  $f$  defined on a set  $E$  of finite measure is Lebesgue integrable over  $E$ .*

*Proof.* Let  $E$  be a set of finite measure. And  $f$  be a bounded, measurable function on  $E$ . Let  $n \in \mathbb{N}$ . By simple approximation lemma, for  $\epsilon = \frac{1}{n}$  we have simple functions  $\varphi_n, \psi_n$  such that  $\varphi_n \leq f \leq \psi_n$  and  $\psi_n - \varphi_n < \frac{1}{n}$  on  $E$ .

$$0 \leq \psi_n - \varphi_n \leq \frac{1}{n} \implies 0 \leq \int_E \psi_n - \varphi_n = \int_E \psi_n - \int_E \varphi_n \leq \frac{1}{n} \int_E 1$$

We know that,  $\inf\{r_1, r_2, \dots\} \leq r_k$  and  $-\sup\{s_1, s_2, \dots\} \leq -s_k$ .

$$\begin{aligned} 0 &\leq \inf \left\{ \int_E \psi : \psi \text{ simple}, f \leq \psi \right\} - \sup \left\{ \int_E \varphi : \varphi \text{ simple}, \varphi \leq f \right\} \\ &\leq \int_E \psi_n - \int_E \varphi_n \leq \frac{1}{n} m(E) \end{aligned}$$

This inequality is true for any  $n \in \mathbb{N}$ . Since  $m(E)$  is finite,  $\frac{1}{n}m(E) \rightarrow 0$  as  $n \rightarrow \infty$ . Thus, upper and lower Lebesgue integrals of  $f$  are equal. Therefore,  $f$  is Lebesgue integrable.  $\square$

**Theorem 10.4.5** (linearity + monotonicity of integral of bounded functions). *Let  $f$  and  $g$  be bounded, measurable functions defined a set  $E$  of finite measure. Then for any  $\alpha, \beta \in \mathbb{R}$ ,*

$$\int_E (\alpha f + \beta g) = \alpha \int_E f + \beta \int_E g \quad (10.47)$$

Moreover, if  $f \leq g$  on  $E$ , then

$$\int_E f \leq \int_E g \quad (10.48)$$

*Proof.* Linear combination of bounded, measurable functions is measurable and bounded. Thus, if  $f, g$  are bounded, measurable functions, then  $\alpha f + \beta g$  is also a bounded, measurable function. And  $\alpha f + \beta g$  is integrable as every bounded, measurable function on  $E$  is integrable over  $E$ .

**Step 1 :**  $\int_E \alpha f = \alpha \int_E f$

$$\begin{aligned} \int_E \alpha f &= \inf_{\psi \geq \alpha f} \left\{ \int_E \psi : \psi \text{ simple} \right\} \\ &= \alpha \sup_{\psi/\alpha \geq f} \left\{ \int_E \psi/\alpha : \psi/\alpha \text{ simple} \right\} \\ &= \alpha \sup_{\psi' \geq f} \left\{ \int_E \psi' : \psi' \text{ simple} \right\} \\ &= \alpha \int_E f \end{aligned}$$

**Step 2 :**  $\int_E f + g = \int_E f + \int_E g$

Since  $f \leq \psi_f$  and  $g \leq \psi_g$ , we have

$$\int_E f + g \leq \int_E \psi_f + \psi_g$$

Since integral of simple functions have linearity

$$\int_E f + g \leq \int_E \psi_f + \psi_g = \int_E \psi_f + \int_E \psi_g$$

This inequality is true for simple functions dominating  $f, g$ . Thus, it is true for the infimum of a family of such simple functions.

$$\int_E f + g \leq \inf \left\{ \int_E \psi_f : \psi_f \text{ simple}, f \leq \psi_f \right\} + \inf \left\{ \int_E \psi_g : \psi_g \text{ simple}, g \leq \psi_g \right\}$$

Therefore,

$$\int_E f + g \leq \int_E f + \int_E g \quad (10.49)$$

Similarly,

$$\begin{aligned} \int_E f + g &\geq \int_E \varphi_f + \varphi_g = \int_E \varphi_f + \int_E \varphi_g \\ &\geq \sup \left\{ \int_E \varphi_f \right\} + \sup \left\{ \int_E \varphi_g \right\} \end{aligned}$$

Therefore,

$$\int_E f + g \geq \int_E f + \int_E g \quad (10.50)$$

Thus, for any bounded, measurable function  $f$  and  $g$ ,  $\int_E f + g = \int_E f + \int_E g$ .

Therefore,  $\int_E \alpha f + \beta g = \int_E \alpha f + \int_E \beta g = \alpha \int_E f + \beta \int_E g$ .  $\square$

**Theorem 10.4.6** (additivity over the domain of integration). *Let  $f$  be a bounded, measurable function on a set  $E$  of finite measure. Suppose  $A, B$  are disjoint, measurable subsets of  $E$ . Then*

$$\int_{A \cup B} f = \int_A f + \int_B f \quad (10.51)$$

*Proof.* We have  $\chi_{A \cup B} = \chi_A + \chi_B$ . And,

$$f \cdot \chi_{A \cup B} = f \cdot (\chi_A + \chi_B) = f \cdot \chi_A + f \cdot \chi_B \quad (10.52)$$

And if  $E_1 \subset E$ , then by the definition of Lebesgue integral

$$\int_E f \cdot \chi_{E_1} = \inf \left\{ \int_E \psi : \psi \text{ simple, } f \cdot \chi_{E_1} \leq \psi \right\}$$

The collection of simple function  $\psi$  such that  $f \cdot \chi_{E_1} \leq \psi$  on  $E$  is same as the collection of simple functions  $\psi$  such that  $f \leq \psi$  on  $E_1$ .

$$= \inf \left\{ \int_{E_1} \psi : \psi \text{ simple, } f \leq \psi \right\} = \int_{E_1} f$$

Since integral has linearity, we have

$$\int_{A \cup B} f = \int_E f \cdot \chi_{A \cup B} = \int_E (f \cdot \chi_A + f \cdot \chi_B) = \int_E f \cdot \chi_A + \int_E f \cdot \chi_B = \int_A f + \int_B f$$

□

**Corollary 10.4.6.1.** *Let  $f$  be a bounded, measurable function on a set of finite measure  $E$ . Then,*

$$\left| \int_E f \right| \leq \int_E |f| \quad (10.53)$$

*Proof.* We have,  $-|f| \leq f \leq |f|$  on  $E$ . Therefore by linearity and monotonicity,

$$-\int_E |f| = \int_E -|f| \leq \int_E f \leq \int_E |f|$$

Therefore,

$$\left| \int_E f \right| \leq \int_E |f|$$

□

**Theorem 10.4.7** (passage of limit under the integral sign). *Let  $\{f_n\}$  be a sequence of bounded, measurable functions on a set  $E$  of finite measure. If  $\{f_n\}$  converges to  $f$  uniformly on  $E$ , then the sequence of the integrals  $\{\int_E f_n\}$  converges to the integral of the limit function  $\int_E f$ .*

$$\lim_{n \rightarrow \infty} \int_E f_n = \int_E f \quad (10.54)$$

*Proof.*

□

For example,  $\{f_n\}$  converges to  $f$  pointwise on  $[0, 1]$ . But we don't observe the passage of limit under integral sign as,

$$\lim_{n \rightarrow \infty} \int_0^1 f_n \neq \int_0^1 f$$

Note : —yet to update the sequence— Pointwise convergence is not sufficient for passage of limit under integral sign. However, the bounded convergence theorem gives us an additional constraint (uniformly pointwise bounded) for the passage of limit under integral sign.

**Theorem 10.4.8** (bounded convergence). *Let  $\{f_n\}$  be a sequence of measurable functions defined on a set  $E$  of finite measure. Suppose  $\{f_n\}$  is uniformly pointwise bounded on  $E$ . If  $\{f_n\}$  converges to  $f$  pointwise on  $E$ , then the sequence of the integrals converges to the integral of the limit function.*

*Proof.*

□

### 10.4.3 Exercise

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.
- 8.
- 9.

10. Let  $f$  be a bounded, measurable function on  $E$ . Let  $A$  be measurable subset of  $E$ . Then,

$$\begin{aligned} \int_E f \cdot \chi_A &= \sup \left\{ \int_E \psi : \psi \text{ simple, } 0 \leq \psi \leq f \cdot \chi_A \right\} \\ &= \sup \left\{ \int_A \psi : \psi \text{ simple, } 0 \leq \psi \leq f \right\} = \int_A f \end{aligned}$$

- 11.
- 12.
- 13.
- 14.
- 15.
- 16.

### 10.4.4 Lebesgue integral of a measurable non-negative function

**Definitions 10.4.9** (vanishing). A function  $f$  on  $E$  vanishes outside  $A \subset E$  if  $f(x) = 0$  for every  $x \in E \sim A$ .

**Definitions 10.4.10** (support). Let  $f$  be a function on  $E$ . Then support of  $f$  is the set  $\{x \in E : f(x) \neq 0\}$  of all non-vanishing points of  $f$  in  $E$ .

**Warning :** The definition of support is different for measure spaces and topological spaces.<sup>†19</sup>

**Definitions 10.4.11** (finite support). A function  $f$  on  $E$  has finite support if support of  $f$  is of finite measure.

**Definitions 10.4.12.** Let  $f$  be a non-negative measurable function on  $E$ . Then the integral of  $f$  over  $E$  is the supremum of the integrals of bounded, measurable functions of finite support  $h$  such that  $0 \leq h \leq f$  in  $E$ .

$$\int_E f = \sup \left\{ \int_E h : h \text{ bounded, measurable, finite support, } 0 \leq h \leq f \right\} \quad (10.55)$$

**Theorem 10.4.9** (Chebychev's inequality). Let  $f$  be a non-negative, measurable function on  $E$ . Then for any  $\lambda > 0$ , we have

$$m \{x \in E : f(x) \geq \lambda\} \leq \frac{1}{\lambda} \cdot \int_E f \quad (10.56)$$

In other words  $\lambda \cdot m(E_\lambda) \leq \int_E f$  where  $E_\lambda = \{x \in E : f(x) \geq \lambda\}$ .

*Proof.* □

**Theorem 10.4.10.** Let  $f$  be a non-negative measurable function on  $E$ . Then,

$$\int_E f = 0 \iff f = 0 \text{ (a.e.) on } E \quad (10.57)$$

*Proof.* □

**Theorem 10.4.11** (linearity + monotonicity of integral of non-negative functions). Let  $f, g$  be non-negative, measurable functions on  $E$ . Then for any real numbers  $\alpha, \beta > 0$ <sup>20</sup>, we have

$$\int_E (\alpha f + \beta g) = \alpha \int_E f + \beta \int_E f \quad (10.58)$$

And we have,

$$\text{If } f \leq g \text{ on } E, \text{ then } \int_E f \leq \int_E g \quad (10.59)$$

<sup>19</sup>“Support of a function  $f$  on a topological space is the **closure** of the set of all non-vanishing points of  $f$ .” We use topological version for differential forms in multivariate real analysis and above definition for Lebesgue integration in integral transforms. Luckily, it is not apparent in our syllabus.

<sup>20</sup>We know that  $\alpha f + \beta g$  is non-negative only when  $\alpha$  and  $\beta$  are non-negative.

*Proof.* □

**Theorem 10.4.12** (additivity over domain of integration of non-negative functions). *Let  $f$  be a non-negative, measurable function on  $E$ . If  $A, B$  are disjoint subsets of  $E$ , then*

$$\int_{A \cup B} f = \int_A f + \int_B f \quad (10.60)$$

*In particular, if  $E_0 \subset E$  is of measure zero, then*

$$\int_E f = \int_{E \setminus E_0} f \quad (10.61)$$

*Proof.* □

**Theorem 10.4.13** (Fatou's lemma). *Let  $\{f_n\}$  be a sequence of non-negative, measurable functions on  $E$ . If  $\{f_n\}$  converges to  $f$  pointwise (a.e.) on  $E$ , then*

$$\int_E f \leq \liminf \int_E f_n$$

*Proof.* □

**Theorem 10.4.14** (monotone convergence). *Let  $\{f_n\}$  be an increasing sequence of non-negative, measurable functions on  $E$ . If  $\{f_n\}$  converges to  $f$  pointwise (a.e.) on  $E$ , then  $\lim_{n \rightarrow \infty} \int_E f_n = \int_E f$ .*

*Proof.* □

**Corollary 10.4.14.1.** *Let  $\{u_n\}$  be a sequence of non-negative, measurable functions on  $E$ . If  $f = \sum_{n=1}^{\infty} u_n$  pointwise (a.e.) on  $E$ , then  $\int_E f = \sum_{n=1}^{\infty} \int_E u_n$ .*

*Proof.* □

**Definitions 10.4.13** (integrable function). A non-negative, measurable function  $f$  on a measurable set  $E$  is integrable over  $E$  if  $\int_E f < \infty$ .

**Theorem 10.4.15.** *Let  $f$  be non-negative, measurable function which is integrable over  $E$ . Then  $f$  is finite a.e. on  $E$ .*

*Proof.* □

**Theorem 10.4.16** (Beppo Levi's Lemma). *Let  $\{f_n\}$  be an increasing sequence of non-negative, measurable functions on  $E$ . If the sequence of integrals  $\left\{ \int_E f_n \right\}$  are bounded, then  $\{f_n\}$  converges pointwise on  $E$  to a measurable function  $f$  which is finite a.e. on  $E$  and*

$$\lim_{n \rightarrow \infty} \int_E f_n = \int_E f < \infty$$

*Proof.* □



**10.4.5 Exercise**

- 17.
- 18.
- 19.
- 20.
- 21.
- 22.
- 23.
- 24.
- 25.
- 26.
- 27.

**10.4.6 General Lebesgue integral**

**Theorem 10.4.17.** *Let  $f$  be a measurable function on  $E$ . Then  $f^+, f^-$  are integrable over  $E$  if and only if  $|f|$  is integrable over  $E$ .*

*Proof.* □

**Note :** A measurable function  $f$  is integrable over  $E$  if  $|f|$  is integrable over  $E$ .

**Definitions 10.4.14** (general Lebesgue integral). Let  $f$  be a measurable function such that  $|f|$  is integrable over  $E$ . Then integral of  $f$  over  $E$  is

$$\int_E f = \int_E f^+ - \int_E f^- \quad (10.62)$$

**Theorem 10.4.18.** *Let  $f$  be integrable over  $E$ . Then  $f$  is finite a.e. on  $E$  and*

$$\int_E f = \int_{E \sim E_0} f \quad (10.63)$$

where  $E_0$  is a subset of  $E$  and is of measure zero.

*Proof.* □

**Theorem 10.4.19** (integral comparison test). *Let  $f$  be a measurable function on  $E$ . Suppose there exists a non-negative function  $g$  such that  $g$  integrable over  $E$  and  $g$  dominates  $f$  on  $E$ . Then  $f$  is integrable over  $E$  and*

$$\left| \int_E f \right| \leq \int_E |f| \quad (10.64)$$

*Proof.* □

**Theorem 10.4.20** (linearity+monotonicity of general integral). *Let functions  $f, g$  be integrable over  $E$ . Then for any  $\alpha, \beta \in \mathbb{R}$ , function  $\alpha f + \beta g$  is integrable over  $E$  and*

$$\int_E (\alpha f + \beta g) = \alpha \int_E f + \beta \int_E g \quad (10.65)$$

Moreover,

$$\text{If } f \leq g \text{ a.e. on } E, \text{ then } \int_E f \leq \int_E g \quad (10.66)$$

*Proof.*

□

**Theorem 10.4.21** (additivity over domain of integration). *Let function  $f$  be integrable over  $E$ . Let  $A, B$  be disjoint, measurable subsets of  $E$ . Then,*

$$\int_{A \cup B} f = \int_A f + \int_B f \quad (10.67)$$

*Proof.*

□

**Theorem 10.4.22** (Lebesgue dominated convergence). *Let  $\{f_n\}$  be a sequence of measurable functions on  $E$ . Suppose there exists a function  $g$  which is integrable over  $E$  and dominates every function in the sequence  $\{f_n\}$  on  $E$ . If  $\{f_n\}$  converges to  $f$  pointwise a.e. on  $E$ , then  $f$  is integrable on  $E$  and*

$$\lim_{n \rightarrow \infty} \int_E f_n = \int_E f \quad (10.68)$$

*Proof.*

□

**Theorem 10.4.23** (General Lebesgue dominated convergence). *Let  $\{f_n\}$  be a sequence of measurable functions on  $E$ . And  $\{f_n\}$  converges pointwise to  $f$  a.e. on  $E$ . Suppose there exists a sequence  $\{g_n\}$  of non-negative, measurable functions on  $E$  that converges pointwise to  $g$  a.e. on  $E$  and  $g$  dominates every function in the sequence  $\{f_n\}$ .*

$$\text{If } \lim_{n \rightarrow \infty} \int_E g_n = \int_E g < \infty, \text{ then } \lim_{n \rightarrow \infty} \int_E f_n = \int_E f \quad (10.69)$$

*Proof.*

□

#### 10.4.7 Exercise

- 28.
- 29.
- 30.
- 31.
- 32.
- 33.
- 34.
- 35.
- 36.

## **10.17 General Measure Spaces**

### **10.17.1 Measures and Measurable Sets**

### **10.17.2 Signed Measures**

### **10.17.3 Caratheodory Measure**

## **10.18 Integration over General Measure Spaces**

### **10.18.1 Measurable Functions**

### **10.18.2 Integration of non-negative measurable functions**

### **10.18.3 Integration of general measurable functions**

### **10.18.4 Radon-Nikodym Theorem**

## **10.20 The Construction of a particular measure**

### **10.20.1 Product Measures**

## Semester III

**Subject 11**

**ME010301 Advanced  
Complex Analysis**

**Subject 12**

**ME010302 Partial  
Differential Equations**

## Subject 13

# ME010303 Multivariate Calculus & Integral Transforms

### 13.1 Integral Transforms

#### 13.1.1 The Weierstrass Approximation Theorem

Every continuous, real valued function on a compact interval has a polynomial approximation. [Apostol, 1973, Theorem 11.17]

**Theorem 13.1.1** (Weierstrass). *Let  $f$  be a real-valued, continuous function on a compact interval  $[a, b]$ . Then for every  $\epsilon > 0$ , there is a polynomial  $p$  such that  $|f(x) - p(x)| < \epsilon$  for every  $x \in [a, b]$ .*

*Synopsis.* Given a real-valued continuous function on compact interval  $[a, b]$ , we can construct a real-valued, continuous function  $g$  on  $\mathbb{R}$  which is periodic with period  $2\pi$ . We have, if  $f \in L(I)$  and  $f$  is bounded almost everywhere in  $I$ , then  $f \in L^2(I)$ . [Apostol, 1973, Theorem 10.52]. By Fejer's theorem ([Apostol, 1973, Theorem 11.15]), the fourier series generated by  $g$  ([Apostol, 1973, definition 11.3]) converges to the Cesaro sum ([Apostol, 1973, Definition 8.47]), which is  $g$  itself in this case. Thus for any  $\epsilon > 0$ , there is a finite sum of trigonometric functions. The power series expansions of trigonometric functions ([Apostol, 1973, definition 9.27]) being uniformly convergent, there exists a polynomial  $p_m$  which approximates  $g$ . And we can construct  $p$  (polynomial approximation of  $g$ ) using  $p_m$ .

*Proof.* Define  $g : \mathbb{R} \rightarrow \mathbb{R}$ ,

$$g(t) = \begin{cases} f(a + (b-a)t/\pi), & t \in [0, \pi) \\ f(a + (2\pi - t)(b-a)/\pi), & t \in [\pi, 2\pi] \\ g(t - 2n\pi), & t > 2\pi, n \in \mathbb{N} \\ g(t + 2n\pi), & t < 0, n \in \mathbb{N} \end{cases}$$

Thus  $g$  is a continuous, real-valued, periodic function with period  $2\pi$  such that

$$f(x) = g\left(\frac{\pi(x-a)}{b-a}\right), \quad x \in [a, b] \quad (13.1)$$

The fourier series generated by  $g$  is given by,

$$g(t) \sim \frac{a_0}{2} + \sum_{k=1}^{\infty} (a_k \cos kt + b_k \sin kt)$$

$$\text{where } a_k = \frac{1}{\pi} \int_0^{2\pi} f(t) \cos kt \, dt, \quad b_k = \frac{1}{\pi} \int_0^{2\pi} f(t) \sin kt \, dt$$

Let  $\{s_n(t)\}$  be the sequence of partial sums of the fourier series generated by  $g$ . And  $\{\sigma_n(t)\}$  be the sequence of averages of  $s_n(t)$  given by,

$$\sigma_n(t) = \frac{1}{n} \sum_{k=1}^n s_k(t), \quad \text{where } s_k(t) = \frac{a_0}{2} + \sum_{j=1}^k (a_j \cos jt + b_j \sin jt)$$

Function  $f \in L(I)$  being real-valued continuous function on a compact interval, it is bounded and hence is Lebesgue square integrable. ie,  $f \in L^2(I)$ . Thus,  $g \in L^2(I)$ .

Since  $g$  is continous on  $\mathbb{R}$ , the function  $s : \mathbb{R} \rightarrow \mathbb{R}$  defined by,

$$s(t) = \lim_{h \rightarrow 0^+} \frac{g(t+h) - g(t-h)}{2}$$

is well-defined on  $\mathbb{R}$  and  $s(t) = g(t)$ ,  $\forall t \in \mathbb{R}$ .

Then by Fejer's Theorem, the sequence  $\{\sigma_n(t)\}$  converges uniformly to  $g(t)$  for every  $t \in \mathbb{R}$ . Thus, given  $\epsilon > 0$ , there exists  $N \in \mathbb{N}$  such that  $\forall t \in \mathbb{R}$ ,  $|g(t) - \sigma_N(t)| < \frac{\epsilon}{2}$ .

We have,

$$\sigma_N(t) = \sum_{k=0}^N (A_k \cos kt + B_k \sin kt), \quad \text{where } A_k, B_k \in \mathbb{R} \quad (13.2)$$

By the power series expansion of the trigonometric functions about origin,

$$\cos kt = \sum_{j=1}^{\infty} \left( \frac{\cos^{(j)} 0}{j!} (kt)^j \right) = \sum_{j=1}^{\infty} A'_j t^j \quad \text{where } A'_j \in \mathbb{R} \quad (13.3)$$

$$\sin kt = \sum_{j=1}^{\infty} \left( \frac{\sin^{(j)} 0}{j!} (kt)^j \right) = \sum_{j=1}^{\infty} B'_j t^j \quad \text{where } B'_j \in \mathbb{R} \quad (13.4)$$

Since the above power series expansions of trigonometric functions are uniformly convergent, their finite linear combination  $\{\sigma_N(t)\}$  is also uniformly convergent. ie, Given  $\epsilon > 0$  there exists  $m \in \mathbb{N}$  such that for every  $t \in \mathbb{R}$

$$\left| \sum_{k=0}^m C_k t^k - \sigma_N(t) \right| < \frac{\epsilon}{2} \quad \text{where } C_k \in \mathbb{R}$$



Therefore,  $|p_m(t) - g(t)| \leq |p_m(t) - \sigma_N(t)| + |\sigma_N(t) - g(t)| < \epsilon$  where  $p_m(t) = \sum_{k=0}^m C_k t^k$ . Define  $p : [a, b] \rightarrow \mathbb{R}$  by,

$$p(x) = p_m \left( \frac{\pi(x-a)}{b-a} \right) \quad (13.5)$$

By equations 13.1 and 13.5,  $|p(x) - f(x)| < \epsilon$  for every  $x \in [a, b]$ .  $\square$

### 13.1.2 Other Forms of Fourier Series

Let  $f \in L([0, 2\pi])$ , then the fourier series generated by  $f$  is given by,

$$f(x) \sim \frac{a_0}{2} + \sum_{n=1}^{\infty} (a_n \cos nx + b_n \sin nx)$$

$$\text{where } a_n = \frac{1}{\pi} \int_0^{2\pi} f(t) \cos nt \, dt, \quad b_n = \frac{1}{\pi} \int_0^{2\pi} f(t) \sin nt \, dt$$

By Euler's formula  $e^{inx} = \cos nx + i \sin nx$ . We have,  $\cos nx = \frac{(e^{inx} + e^{-inx})}{2}$  and  $\sin nx = \frac{(e^{inx} - e^{-inx})}{2i}$

$$f(x) \sim \frac{a_0}{2} + \sum_{n=1}^{\infty} (\alpha_n e^{inx} + \beta_n e^{-inx})$$

$$\text{where } \alpha_n = \frac{(a_n - ib_n)}{2} \quad \beta_n = \frac{(a_n + ib_n)}{2}$$

Therefore, by assigning  $\alpha_0 = a_0/2$ ,  $\alpha_{-n} = \beta_n$ , we get the following exponential form of fourier series generated by  $f$ ,

$$f(x) \sim \sum_{n=-\infty}^{\infty} \alpha_n e^{inx} \text{ where } \alpha_n = \frac{1}{2\pi} \int_0^{2\pi} f(t) e^{-int} \, dt$$

Note : If  $f$  is periodic with period  $2\pi$ , then the interval of integration  $[0, 2\pi]$  can be replaced with any interval of length  $2\pi$ . eg.  $[-\pi, \pi]$

#### Periodic with period $p$

Let  $f \in L([0, p])$  and  $f$  is periodic with period  $p$ . Then

$$f(x) \sim \frac{a_0}{2} + \sum_{n=1}^{\infty} \left( a_n \cos \frac{2\pi nx}{p} + b_n \sin \frac{2\pi nx}{p} \right)$$

$$\text{where } a_n = \frac{2}{p} \int_0^p f(t) \cos \frac{2\pi nt}{p} \, dt \quad b_n = \frac{2}{p} \int_0^p f(t) \sin \frac{2\pi nt}{p} \, dt$$

Therefore, we have the exponential form of the above fourier series given by,

$$f(x) \sim \sum_{n=-\infty}^{\infty} \alpha_n e^{\frac{2\pi inx}{p}}, \text{ where } \alpha_n = \frac{1}{p} \int_0^p f(t) e^{\frac{-2\pi int}{p}} \, dt$$

### 13.1.3 Fourier Integral Theorem

**Theorem 13.1.2** (Fourier Integral Theorem). *Let  $f \in L(-\infty, \infty)$ . Suppose  $x \in \mathbb{R}$  and an interval  $[x - \delta, x + \delta]$  about  $x$  such that either*

1.  *$f$  is of bounded variation on an interval  $[x - \delta, x + \delta]$  about  $x$  or*
2. *both limits  $f(x+)$  and  $f(x-)$  exists and both Lebesgue integrals*

$$\int_0^\delta \frac{f(x+t) - f(x+)}{t} dt \text{ and } \int_0^\delta \frac{f(x-t) - f(x-)}{t} dt$$

*exists.*

*Then,*

$$\frac{f(x+) + f(x-)}{2} = \frac{1}{\pi} \int_0^\infty \int_{-\infty}^\infty f(u) \cos v(u-x) du dv,$$

*the integral  $\int_0^\infty$  being an improper Riemann integral.*

*Synopsis.*

$$f(x+t) \frac{\sin \alpha t}{\pi t} dt \rightarrow f(u) \frac{\sin \alpha(u-x)}{\pi(u-x)} \rightarrow \frac{f(u)}{\pi} \int_0^\alpha \cos v(u-x) dv$$

By Riemann-Lebesgue lemma[Apostol, 1973, Theorem 11.6],

$$f \in L(I) \implies \lim_{\alpha \rightarrow +\infty} \int_I f(x) \sin \alpha t dt = 0$$

By Jordan's Theorem[Apostol, 1973, Theorem 10.8], if  $g$  is of bounded variation on  $[0, \delta]$ , then

$$\lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^\delta g(t) \frac{\sin \alpha t}{t} dt = g(0+)$$

By Dini's Theorem[Apostol, 1973, Theorem 10.9], if the limit  $g(x+)$  exists and Lebesgue integral  $\int_0^\delta \frac{g(t)+g(0+)}{t} dt$  exists for some  $\delta > 0$ , then

$$\lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^\delta g(t) \frac{\sin \alpha t}{t} dt = g(0+)$$

The order of Lebesgue integrals can be interchanged.[Apostol, 1973, Theorem 10.40]

Suppose  $f \in L(X)$  and  $g \in L(Y)$ . Then

$$\int_X f(x) \left( \int_Y g(y) k(x, y) dy \right) dx = \int_Y g(y) \left( \int_X f(x) k(x, y) dx \right) dy$$

*Proof.* Consider  $\int_{-\infty}^\infty f(x+t) \frac{\sin \alpha t}{\pi t} dt$ . We prove that this integral is equal to the either sides.

$$\int_{-\infty}^\infty f(x+t) \frac{\sin \alpha t}{\pi t} dt = \int_{-\infty}^{-\delta} + \int_{-\delta}^0 + \int_0^{-\delta} + \int_\delta^\infty f(x+t) \frac{\sin \alpha t}{\pi t} dt$$

We have, function  $\frac{f(x+t)}{\pi t}$  is bounded on  $(-\infty, -\delta) \cup (\delta, \infty)$ , hence  $\frac{f(x+t)}{\pi t}$  is Lebesgue integrable on  $(-\infty, -\delta) \cup (\delta, \infty)$ .

By Riemann Lebesgue lemma,

$$\frac{f(x+t)}{\pi t} \in L(-\infty, -\delta) \implies \int_{-\infty}^{-\delta} f(x+t) \frac{\sin \alpha t}{\pi t} dt = 0,$$

$$\frac{f(x+t)}{\pi t} \in L(\delta, \infty) \implies \int_{\delta}^{\infty} f(x+t) \frac{\sin \alpha t}{\pi t} dt = 0$$

**Case 1** Suppose  $f$  is of bounded variation on  $[x-\delta, x+\delta]$ , put  $g(t) = f(x+t)$  then  $g$  is of bounded variation on  $[-\delta, \delta]$ . Thus  $g$  is of bounded variation on  $[0, \delta]$ . Then by Jordan's Theorem

$$\lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^{\delta} f(x+t) \frac{\sin \alpha t}{t} dt = \lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^{\delta} g(t) \frac{\sin \alpha t}{t} dt = g(0+) = f(x+)$$

**Case 2** Suppose both the limits  $f(x+)$  and  $f(x-)$  exists and both Lebesgue integrals

$$\int_0^{\delta} \frac{f(x+t) - f(x+)}{t} dt \text{ and } \int_0^{\delta} \frac{f(x-t) - f(x-)}{t} dt$$

exists.

Thus, we have  $f(x+)$  exists and the Lebesgue integral  $\int_0^{\delta} \frac{f(x+t) - f(x+)}{t} dt$  exists. Put  $g(t) = f(x+t)$ , then  $g(0+) = f(x+)$  exists and the Lebesgue integral  $\int_0^{\delta} \frac{g(t) - g(0+)}{t} dt$  exists, then by Dini's Theorem,

$$\lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^{\delta} f(x+t) \frac{\sin \alpha t}{t} dt = \lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^{\delta} g(t) \frac{\sin \alpha t}{t} dt = g(0+) = f(x+)$$

Similarly,  $f(x-)$  exists and the Lebesgue integral  $\int_0^{\delta} \frac{f(x-t) - f(x-)}{t} dt$  exists. Put  $g(t) = f(x-t)$ , then  $g(0+) = f(x-)$  exists and the Lebesgue integral  $\int_0^{\delta} \frac{g(t) - g(0+)}{t} dt$  exists, then by Dini's Theorem,

$$\begin{aligned} \lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_{-\delta}^0 f(x+t) \frac{\sin \alpha t}{t} dt &= \lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^{\delta} f(x-\tau) \frac{\sin \alpha \tau}{\tau} d\tau \\ &= \lim_{\alpha \rightarrow +\infty} \frac{2}{\pi} \int_0^{\delta} g(\tau) \frac{\sin \alpha \tau}{\tau} d\tau = g(0+) = f(x-) \end{aligned}$$

Then by either cases,

$$\begin{aligned} \lim_{\alpha \rightarrow +\infty} \int_{-\infty}^{\infty} f(x+t) \frac{\sin \alpha t}{\pi t} dt &= \lim_{\alpha \rightarrow +\infty} \int_{-\delta}^0 f(x+t) \frac{\sin \alpha t}{\pi t} dt + \int_0^{\delta} f(x+t) \frac{\sin \alpha t}{\pi t} dt \\ &= \frac{f(x+) + f(x-)}{2} \end{aligned}$$

We have,  $\int_0^\alpha \cos v(u-x)dv = \frac{\sin v(u-x)}{u-x}$ .

$$\begin{aligned} \lim_{\alpha \rightarrow +\infty} \int_{-\infty}^{\infty} f(x) \frac{\sin \alpha t}{\pi t} dt &= \lim_{\alpha \rightarrow +\infty} \int_{-\infty}^{\infty} f(u) \frac{\sin \alpha(u-x)}{u-x} du, \text{ (put } u = x+t) \\ &= \lim_{\alpha \rightarrow +\infty} \int_{-\infty}^{\infty} f(u) \left( \int_0^\alpha \cos v(u-x) dv \right) du \\ &= \lim_{\alpha \rightarrow +\infty} \int_0^\alpha \left( \int_{-\infty}^{\infty} f(u) \cos v(u-x) du \right) dv, \\ &\text{since, the order of Lebesgue integrals can be reversed.} \\ &= \int_0^\infty \left( \int_{-\infty}^{\infty} f(u) \cos v(u-x) du \right) dv \end{aligned}$$

where,  $\int_0^\infty$  is not a Lebesgue integral, but an improper Riemann integral

Therefore,

$$\begin{aligned} \int_0^\infty \left( \int_{-\infty}^{\infty} f(u) \cos v(u-x) du \right) dv &= \lim_{\alpha \rightarrow +\infty} \int_{-\infty}^{\infty} f(x) \frac{\sin \alpha t}{\pi t} dt \\ &= \frac{f(x+) + f(x-)}{2} \end{aligned}$$

□

*Remark.* If a function  $f$  on  $(-\infty, \infty)$  is non-periodic, then it may not have a fourier series representation. In such cases, we have fourier intergral representation.

### 13.1.4 Exponential form of Fourier Integral Theorem

Let  $f \in L(-\infty, \infty)$ . Suppose  $x \in \mathbb{R}$  and an interval  $[x-\delta, x+\delta]$  about  $x$  such that either

1.  $f$  is of bounded variation on an interval  $[x-\delta, x+\delta]$  about  $x$  or
2. both limits  $f(x+)$  and  $f(x-)$  exists and both Lebesgue intergrals

$$\int_0^\delta \frac{f(x+t) - f(x+)}{t} dt \text{ and } \int_0^\delta \frac{f(x-t) - f(x-)}{t} dt$$

exists.

Then,

$$\frac{f(x+) + f(x-)}{2} = \lim_{\alpha \rightarrow \infty} \frac{1}{2\pi} \int_{-\alpha}^{\alpha} \left( \int_{-\infty}^{\infty} f(u) e^{iv(u-x)} du \right) dv$$

*Proof.* Let  $F(v) = \int_{-\infty}^{\infty} f(u) \cos v(u-x) du$ . Then  $F(v) = F(-v)$  and

$$\begin{aligned} \lim_{\alpha \rightarrow \infty} \frac{1}{2\pi} \int_{-\alpha}^{\alpha} F(v) dv &= \lim_{\alpha \rightarrow \infty} \frac{1}{\pi} \int_0^{\alpha} \int_{-\infty}^{\infty} f(u) \cos v(u-x) du dv \\ &= \frac{f(x+) + f(x-)}{2} \end{aligned}$$

Let  $G(v) = \int_{-\infty}^{\infty} f(u) \sin v(u-x) du$ . Then  $G(v) = -G(-v)$  and

$$\lim_{\alpha \rightarrow \infty} \frac{1}{2\pi} \int_{-\alpha}^{\alpha} G(v) dv = 0$$

Thus

$$\lim_{\alpha \rightarrow \infty} \frac{1}{2\pi} \int_{-\alpha}^{\alpha} F(v) + iG(v) dv = \frac{f(x+) + f(x-)}{2}$$

□

### 13.1.5 Integral Transforms

**Definitions 13.1.1.** **Integral transform**  $g(y)$  of  $f(x)$  is a Lebesgue integral or Improper Riemann integral of the form

$$g(y) = \int_{-\infty}^{\infty} K(x, y) f(x) dx$$

, where  $K$  is the kernel of the transform. We write  $g = \mathcal{K}(f)$ .

*Remark.* Integral transforms(operators) are linear operators. ie,  $\mathcal{K}(af_1 + bf_2) = a\mathcal{K}f_1 + b\mathcal{K}f_2$

*Remark.* A few commonly used integral transforms,

1. Exponential Fourier Transform  $\mathcal{F}$ ,

$$\mathcal{F}f = \int_{-\infty}^{\infty} e^{-ixy} f(x) dx$$

2. Fourier Cosine Transform  $\mathcal{C}$ ,

$$\mathcal{C}f = \int_0^{\infty} \cos xy f(x) dx$$

3. Fourier Sine Transform  $\mathcal{S}$ ,

$$\mathcal{S}f = \int_0^{\infty} \sin xy f(x) dx$$

4. Laplace Transform  $\mathcal{L}$ ,

$$\mathcal{L}f = \int_0^{\infty} e^{-xy} f(x) dx$$

5. Mellin Transform  $\mathcal{M}$ ,

$$\mathcal{M}f = \int_0^{\infty} x^{y-1} f(x) dx$$

*Remark.* Suppose  $f(x) = 0, \forall x < 0$ .

$$\int_{-\infty}^{\infty} e^{-ixy} f(x) dx = \int_0^{\infty} e^{-ixy} f(x) dx = \int_0^{\infty} \cos xy f(x) dx + i \int_0^{\infty} \sin xy f(x) dx$$

$$\mathcal{F}f = \mathcal{C}f + i\mathcal{S}f$$

Therefore Fourier Cosine  $\mathcal{C}$  and Sine  $\mathcal{S}$  transforms are special cases of fourier integral transform,  $\mathcal{F}$  provided  $f$  vanishes on negative real axis.

*Remark.* Let  $y = u + iv$ ,  $f(x) = 0$ ,  $\forall x < 0$ .

$$\int_0^\infty e^{-xy} f(x) dx = \int_0^\infty e^{-xu} e^{-ixv} f(x) dx = \int_0^\infty e^{-ixv} \phi_u(x) dx$$

where  $\phi_u(x) = e^{-xu} f(x)$ .

$$\mathcal{L}f = \mathcal{F}\phi_u$$

Therefore Laplace transform,  $\mathcal{L}$  is a special case of Fourier integral transform,  $\mathcal{F}$ .

*Remark.* Let  $g(y) = \mathcal{F}f(x)$ .

$$g(y) = \int_{-\infty}^\infty e^{-ixy} f(x) dx$$

Suppose  $f$  is continuous at  $x$ , then by fourier integral theorem,

$$\begin{aligned} f(x) &= \frac{1}{2\pi} \int_{-\infty}^\infty \left( \int_{-\infty}^\infty f(u) e^{iv(u-x)} du \right) dv \\ &= \int_{-\infty}^\infty e^{-ivx} \left( \frac{1}{2\pi} \int_{-\infty}^\infty e^{ivu} f(u) du \right) dv \\ &= \int_{-\infty}^\infty g(v) e^{-ivx} dv = \mathcal{F}g \text{ where } g(v) = \frac{1}{2\pi} \int_{-\infty}^\infty f(u) e^{ivu} du \end{aligned}$$

The above function  $g(v)$  gives the **inverse fourier transformation** of  $f$ .

Let  $g$  be fourier transform of  $f$ , then  $f$  is uniquely determined by its fourier transform  $g$  by,

$$f(x) = \mathcal{F}^{-1}g(y) = \frac{1}{2\pi} \lim_{\alpha \rightarrow \infty} \int_{-\alpha}^\alpha g(y) e^{ixy} dy$$

6. Inverse Fourier Transform  $\mathcal{F}^{-1}$ ,

$$\mathcal{F}^{-1}f = \int_{-\infty}^\infty \frac{e^{ixy}}{2\pi} f(x) dx$$

### 13.1.6 Convolutions

**Definitions 13.1.2.** Let  $f, g \in L(-\infty, \infty)$ . Let  $S$  be the set of all points  $x$  for which the Lebesgue integral

$$h(x) = \int_{-\infty}^\infty f(t)g(x-t)dt$$

exists. Then the function  $h : S \rightarrow \mathbb{R}$  is a **convolution** of  $f$  and  $g$ . And  $h = f * g$ .

*Remark.* Convolution operator is commutative. ie,  $h = f * g = g * f$  (hint : take  $u = x - t$ )

*Remark.* Suppose  $f, g$  vanishes on negative real axis, then

$$h(x) = \int_{-\infty}^\infty f(t)g(x-t)dt = \int_{-\infty}^0 + \int_0^x + \int_x^\infty f(t)g(x-t)dt = \int_0^x f(t)g(x-t)dt$$

*Remark.* Singularity of convolution is a point at which the convolution integral fails to exist.

**Theorem 13.1.3.** *Let  $f, g \in L(\mathbb{R})$  and either  $f$  or  $g$  is bounded in  $\mathbb{R}$ . Then the convolution integral*

$$h(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$

*exists for every  $x \in \mathbb{R}$  and the function  $h$  so defined is bounded in  $\mathbb{R}$ . In addition, if the bounded function is continuous on  $\mathbb{R}$ , then  $h$  is continuous and  $h \in L(\mathbb{R})$ .*

*Synopsis.*

*Proof.* □

*Remark.* If  $f, g$  are both unbounded, the convolution integral may not exist.

$$\text{eg: } f(t) = \frac{1}{\sqrt{t}}, \quad g(t) = \frac{1}{\sqrt{1-t}}$$

**Theorem 13.1.4.** *Let  $f, g \in L^2(\mathbb{R})$ . Then the convolution integral  $f * g$  exists for each  $x \in \mathbb{R}$  and the function  $h : \mathbb{R} \rightarrow \mathbb{R}$  defined by  $h(x) = f * g(x)$  is bounded in  $\mathbb{R}$ .*

*Synopsis.*

*Proof.* □

### 13.1.7 The Convolution Theorem for Fourier Transforms

**Theorem 13.1.5.** *Let  $f, g \in L(\mathbb{R})$  and at least one of  $f$  or  $g$  is continuous and bounded on  $\mathbb{R}$ . Let  $h = f * g$ . Then for every real  $u$ ,*

$$\int_{-\infty}^{\infty} h(x)e^{-ixu}dx = \left( \int_{-\infty}^{\infty} f(t)e^{-itu}dt \right) \left( \int_{-\infty}^{\infty} g(y)e^{-iyu}dy \right)$$

*The integral on the left exists both as a Lebesgue integral and an improper Riemann integral.*

*Synopsis.*

*Proof.* □

*Remark* (Application of Convolution Theorem).

$$B(p, q) = \frac{\Gamma p \Gamma q}{\Gamma p + q}, \text{ where } B(p, q) = \int_0^1 x^{p-1}(1-x)^{q-1}dx, \quad \Gamma p = \int_0^{\infty} t^{p-1}e^{-t}dt$$

## 13.2 Multivariate Differential Calculus

In this chapter, we deal with real functions of several variables. Instead of  $\mathbf{c}$ , we write  $\bar{c} \in \mathbb{R}^n$ , then  $\bar{c} = (c_1, c_2, \dots, c_n)$  where  $c_j \in \mathbb{R}$  for every  $j = 1, 2, \dots, n$ . Again, suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $f(\bar{x}) = \bar{y}$ , then  $\bar{y} = (y_1, y_2, \dots, y_m)$  where each  $y_k$  is real. The unit co-ordinate vector,  $\bar{u}_k$  is given by  $u_{kj} = \delta_{j,k}$

### 13.2.1 Directional Derivative

*Motivation :* The existence of all partial derivatives of a multivariate real function  $f$  at a point  $\bar{c}$  doesn't imply the continuity of  $f$  at  $\bar{c}$ . Thus, we need a suitable generalisation for the partial derivative which could characterise continuity. And directional derivative is such an attempt.

**Definitions 13.2.1** (Directional Derivative). Let  $S \subset \mathbb{R}^n$  and  $f : S \rightarrow \mathbb{R}^m$ . Let  $\bar{c}$  be an interior points of  $S$  and  $\bar{u} \in \mathbb{R}^n$ , then there exists an open ball  $B(\bar{c}, r)$  in  $S$ . Also for some  $\delta > 0$  the line segment  $\alpha : [0, \delta] \rightarrow S$  given by  $\alpha(t) = \bar{c} + t\bar{u}$  lie in  $B(\bar{c}, r)$ . Then the **directional derivative** of  $f$  at an interior point  $\bar{c}$  in the direction  $\bar{u}$  is given by

$$f'(\bar{c}, \bar{u}) = \lim_{h \rightarrow 0} \frac{f(\bar{c} + h\bar{u}) - f(\bar{c})}{h}$$

*Remark.* The direction derivative of  $f$  at an interior point  $\bar{c}$  in the direction  $\bar{u}$  exists only if the above limit exists.

*Remark.* Example, [Apostol, 1973, Exercise 12.2a]

Suppose  $\bar{x}, \bar{a}, \bar{c}, \bar{u} \in \mathbb{R}^n$ . Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $f(\bar{x}) = \bar{a} \cdot \bar{x}$ . Then

$$f'(\bar{c}, \bar{u}) = \lim_{h \rightarrow 0} \frac{\bar{a} \cdot (\bar{c} + h\bar{u}) - \bar{a} \cdot \bar{c}}{h} = \bar{a} \cdot \bar{u}$$

*Remark* (Properties). Let  $f : S \rightarrow \mathbb{R}^m$ , where  $S \subset \mathbb{R}^n$

1.  $f'(\bar{c}, \bar{0}) = \bar{0}$   
*Note :* The zero vectors belongs to  $\mathbb{R}^n, \mathbb{R}^m$  respectively.
2.  $f'(\bar{c}, \bar{u}_k) = \frac{\partial f}{\partial u_k}(\bar{c}) = D_k f(\bar{c})$ , the  $k^{th}$  partial derivative of  $f$ .
3. Let  $f = (f_1, f_2, \dots, f_m)$ , such that  $f(\bar{c}) = (f_1(\bar{c}), f_2(\bar{c}), \dots, f_m(\bar{c}))$ . Then,  
 $\exists f'(\bar{c}, \bar{u}) \iff \forall k, \exists f'_k(\bar{c}, \bar{u})$  and  $f'(\bar{c}, \bar{u}) = (f'_1(\bar{c}, \bar{u}), f'_2(\bar{c}, \bar{u}), \dots, f'_m(\bar{c}, \bar{u}))$

ie, Directional derivative of  $f$  exists iff directional derivative of each component function  $f_k$  exists. And the components of the directional derivatives of  $f$  are the directional derivatives of the components of  $f$ .

Thus  $D_k f(\bar{c}) = (D_k f_1(\bar{c}), D_k f_2(\bar{c}), \dots, D_k f_m(\bar{c}))$  holds.

4. Let  $F(t) = f(\bar{c} + t\bar{u})$ , then  $F'(0) = f'(\bar{c}, \bar{u})$  and  $F'(t) = f'(\bar{c} + t\bar{u}, \bar{u})$
5. Let  $f(\bar{c}) = \bar{c} \cdot \bar{c} = \|\bar{c}\|^2$ , and  $F(t) = f(\bar{c} + t\bar{u})$ , then  $F'(t) = 2\bar{c} \cdot \bar{u} + 2t\|\bar{u}\|^2$  and  $F'(0) = f'(\bar{c}, \bar{u}) = 2\bar{c} \cdot \bar{u}$
6. Let  $f$  be linear, then  $f'(\bar{c}, \bar{u}) = f(\bar{u})$
7. Existence of all partial derivatives doesn't imply existence of all directional derivatives.

$$f(x, y) = \begin{cases} x + y & \text{if } x = 0 \text{ or } y = 0 \\ 1 & \text{otherwise} \end{cases}$$

For above  $f$ , directional derivatives exists only along the co-ordinates (ie, partial derivatives).



8. Existence of all directional derivatives doesn't imply continuity.

$$f(x, y) = \begin{cases} xy^2(x^2 + y^4) & x \neq 0 \\ 0 & x = 0 \end{cases}$$

Above  $f$  is discontinuous at  $(0, 0)$ , however all directional derivatives exists and has finite value.

### 13.2.2 Total Derivative

We may define a total derivative  $T_c(h) = hf'(c)$  in the case of real-functions of single variable as follows :-

$$\text{Let } E_c(h) = \begin{cases} \frac{f(c+h)-f(c)}{h} - f'(c), & h \neq 0 \\ 0, & h = 0 \end{cases}$$

Then,  $f(c+h) = f(c) + hf'(c) + hE_c(h)$  and as  $h \rightarrow 0$ ,  $E_c(h) \rightarrow 0$ . Also  $T_c(h) = f'(c)h$  is a linear function of  $h$ . ie,  $T_c(ah_1 + bh_2) = aT_c(h_1) + bT_c(h_2)$ . Now, we will define a total derivative of multivariate function that has these two properties.

**Definitions 13.2.2** (Total Derivative). The function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is **differentiable** at  $\bar{c}$  if there exists a **linear** function  $T_{\bar{c}} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that  $f(\bar{c} + \bar{v}) = f(\bar{c}) + T_{\bar{c}}(\bar{v}) + \|\bar{v}\|E_{\bar{c}}(\bar{v})$  where  $E_{\bar{c}}(\bar{v}) \rightarrow \bar{0}$  as  $\bar{v} \rightarrow \bar{0}$ .

*Remark.* The linear function  $T_{\bar{c}}$  is the total derivative of  $f$  at  $\bar{c}$ ,  $T_{\bar{c}}(\bar{0}) = \bar{0}$  and the condition above gives the First Order Taylor's Formula for  $f(\bar{c} + \bar{v}) - f(\bar{c})$ .

*Remark* (Properties). Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $f'(\bar{c})(\bar{v}) = T_{\bar{c}}(\bar{v})$  be the total derivative of  $f$  at  $\bar{c}$  evaluated at  $\bar{v}$ . Then,

1.  $f'(\bar{c})(\bar{v}) = f'(\bar{c}, \bar{u})$
2. If  $f$  is differentiable at  $\bar{c}$ , then  $f$  is continuous at  $\bar{c}$ .
3.  $f'(\bar{c})(\bar{v}) = v_1 D_1 f(\bar{c}) + v_2 D_2 f(\bar{c}) + \cdots + v_n D_n f(\bar{c})$

*Note.* The above  $f'$  is a function from  $\mathbb{R}^n$  to the set of all linear functions  $\mathcal{L} = \{h : \mathbb{R}^n \rightarrow \mathbb{R}^m\}$ .  $f'(\bar{c})$  is a linear function (in fact, total derivative  $T_{\bar{c}}$ ) which maps  $\bar{v}$  into the directional derivatives of  $f$  at  $\bar{c}$  in the direction  $\bar{v}$ . This notation generalises  $f'$  for univariate  $f$  as well. (put  $n = m = 1$ )

In this subject, we use the following notations,

$D_k f(\bar{c})$  partial derivative

$f'(\bar{c}, \bar{v})$  directional derivative

$f'(\bar{c})(\bar{v})$  total derivative

$\nabla f(\bar{c})$  gradient vector

**Theorem 13.2.1.** If  $f$  is differentiable at  $\bar{c}$  with total derivative  $T_{\bar{c}}$ , then for every  $\bar{u} \in \mathbb{R}^n$ ,  $T_{\bar{c}}(\bar{u}) = f'(\bar{c}, \bar{u})$ . ( ie,  $f'(\bar{c})(\bar{v}) = f'(\bar{c}, \bar{v})$  )

*Proof.* For  $\bar{v} = \bar{0}$ , we have  $T_{\bar{c}}(\bar{0}) = 0 = f'(\bar{c}, \bar{0})$ .

Suppose  $\bar{v} \neq \bar{0}$ , then put  $\bar{v} = h\bar{u}$ . Since  $f$  is differentiable at  $\bar{c}$ ,  $f$  has total derivative at  $\bar{c}$ . That is, there exists a linear function  $T_{\bar{c}}$  such that  $f(\bar{c} + h\bar{u}) = f(\bar{c}) + T_{\bar{c}}(h\bar{u}) + \|h\bar{u}\|E_{\bar{c}}(h\bar{u})$  where  $E_{\bar{c}}(h\bar{u}) \rightarrow \bar{0}$  as  $h\bar{u} \rightarrow \bar{0}$ .

$$\begin{aligned} \implies f(\bar{c} + h\bar{u}) &= f(\bar{c}) + hT_{\bar{c}}(\bar{u}) + |h|\|\bar{u}\|E_{\bar{c}}(h\bar{u}), \quad E_{\bar{c}}(h\bar{u}) \rightarrow \bar{0} \text{ as } h\bar{u} \rightarrow \bar{0} \\ \implies \frac{f(\bar{c} + h\bar{u}) - f(\bar{c})}{h} &= T_{\bar{c}}(\bar{u}) + \frac{|h|\|\bar{u}\|E_{\bar{c}}(h\bar{u})}{h}, \quad E_{\bar{c}}(h\bar{u}) \rightarrow \bar{0} \text{ as } h \rightarrow 0 \\ \implies \lim_{h \rightarrow 0} \frac{f(\bar{c} + h\bar{u}) - f(\bar{c})}{h} &= T_{\bar{c}}(\bar{u}) + \lim_{h \rightarrow 0} \frac{|h|\|\bar{u}\|E_{\bar{c}}(h\bar{u})}{h} \\ \implies f'(\bar{c}, \bar{u}) &= T_{\bar{c}}(\bar{u}) \end{aligned}$$

□

*Note.*  $T_{\bar{c}}$  is linear, however  $E_{\bar{c}}$  is not linear. Thus  $E_{\bar{c}}(h\bar{u}) \neq hE_{\bar{c}}(\bar{u})$ .

As  $h \rightarrow 0$ ,  $h\bar{u} \rightarrow \bar{0}$  and  $E_{\bar{c}}(h\bar{u}) \rightarrow \bar{0}$ . Since the order of the function  $E_{\bar{c}}(h\bar{u})$  is much smaller than that of  $h$ , the limit on the right converges to 0.

**Theorem 13.2.2.** *If  $f$  is differentiable at  $\bar{c}$ , then  $f$  is continuous at  $\bar{c}$ .*

*Proof.* Let  $\bar{v} \neq \bar{0}$ , then

$$\begin{aligned} \bar{v} &= v_1\bar{u}_1 + v_2\bar{u}_2 + \cdots + v_n\bar{u}_n, \\ \bar{v} \rightarrow \bar{0} &\implies \forall j, v_j \rightarrow 0 \\ T \text{ is linear} &\implies T_{\bar{c}}(\bar{v}) = v_1T_{\bar{c}}(\bar{u}_1) + v_2T_{\bar{c}}(\bar{u}_2) + \cdots + v_nT_{\bar{c}}(\bar{u}_n) \\ \text{Thus, } T_{\bar{c}}(\bar{v}) &\rightarrow \bar{0} \text{ as } \bar{v} \rightarrow \bar{0} \end{aligned}$$

Since  $f$  differentiable at  $\bar{c}$ , there exists linear function  $T_{\bar{c}}$  such that

$$\begin{aligned} f(\bar{c} + \bar{v}) &= f(\bar{c}) + T_{\bar{c}}(\bar{v}) + \|\bar{v}\|E_{\bar{c}}(\bar{v}) \\ \implies \lim_{\bar{v} \rightarrow \bar{0}} f(\bar{c} + \bar{v}) &= f(\bar{c}) + \lim_{\bar{v} \rightarrow \bar{0}} T_{\bar{c}}(\bar{v}) + \lim_{\bar{v} \rightarrow \bar{0}} \|\bar{v}\|E_{\bar{c}}(\bar{v}) \\ \implies \lim_{\bar{v} \rightarrow \bar{0}} f(\bar{c} + \bar{v}) &= f(\bar{c}) \end{aligned}$$

□

**Theorem 13.2.3.** *Let  $S \subset \mathbb{R}^n$  and  $f : S \rightarrow \mathbb{R}^m$  be differentiable at an interior point  $\bar{c}$  of  $S$ , where  $S \subseteq \mathbb{R}^n$ . If  $\bar{v} = v_1\bar{u}_1 + v_2\bar{u}_2 + \cdots + v_n\bar{u}_n$ , then*

$$f'(\bar{c})(\bar{v}) = \sum_{k=1}^n v_k D_k f(\bar{c})$$

*In particular, if  $f$  is real-valued ( $m = 1$ ) we have,  $f'(\bar{c})(\bar{v}) = \nabla f(\bar{c}) \cdot \bar{v}$*

*Proof.* Suppose  $f : S \rightarrow \mathbb{R}^m$  is differentiable at  $\bar{c}$ , then there exists a linear function  $f'(\bar{c}) : S \rightarrow \mathbb{R}^m$  such that  $f(\bar{c} + \bar{v}) = f(\bar{c}) + f'(\bar{c})(\bar{v}) + \|\bar{v}\|E_{\bar{c}}(\bar{c})$  where

$E_{\bar{c}} \rightarrow \bar{0}$  as  $\bar{v} \rightarrow \bar{0}$ .

$$\begin{aligned} f'(\bar{c})(\bar{v}) &= f'(\bar{c}) \left( \sum_{k=1}^n v_k \bar{u}_k \right) \\ &= \sum_{k=1}^n v_k f'(\bar{c})(\bar{u}_k), \text{ since } f'(\bar{c}) \text{ is linear} \\ &= \sum_{k=1}^n v_k D_k f(\bar{c}), \text{ since } f'(\bar{c})(\bar{u}_k) = f'(\bar{c}, \bar{u}_k) = D_k f(\bar{c}) \end{aligned}$$

Let  $m = 1$ , then  $f : S \rightarrow \mathbb{R}$

$$\begin{aligned} f'(\bar{c})(\bar{v}) &= \sum_{k=1}^n v_k D_k f(\bar{c}) = \nabla f(\bar{c}) \cdot \bar{v} \\ &\text{since } \nabla f(\bar{c}) = (D_1 f(\bar{c}), D_2 f(\bar{c}), \dots, D_n f(\bar{c})) \end{aligned}$$

□

*Remark.* Let  $f : S \rightarrow \mathbb{R}$ , then  $f(\bar{c} + \bar{v}) = f(\bar{c}) + \nabla f(\bar{c}) \cdot \bar{v} + o(\|\bar{v}\|)$  as  $\bar{v} \rightarrow \bar{0}$ .

*Remark* (Complex-valued Functions).

### 13.2.3 Matrix of Linear Function

Let  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a linear function. Let  $\{\bar{u}_1, \bar{u}_2, \dots, \bar{u}_n\}$  be standard basis for  $\mathbb{R}^n$  and  $\{\bar{e}_1, \bar{e}_2, \dots, \bar{e}_m\}$  be standard basis for  $\mathbb{R}^m$ . Let  $\bar{v} \in \mathbb{R}^n$ , then  $\bar{v} = \sum_{k=1}^n v_k \bar{u}_k$  and  $T(\bar{v}) = \sum_{k=1}^n v_k T(\bar{u}_k)$  and

$$\begin{aligned} T(\bar{v}) &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} T(\bar{u}_1) \\ T(\bar{u}_2) \\ \dots \\ T(\bar{u}_n) \end{bmatrix} \\ &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} t_{11}\bar{e}_1 + t_{21}\bar{e}_2 + \cdots + t_{m1}\bar{e}_m \\ t_{12}\bar{e}_1 + t_{22}\bar{e}_2 + \cdots + t_{m2}\bar{e}_m \\ \dots \\ t_{1n}\bar{e}_1 + t_{2n}\bar{e}_2 + \cdots + t_{mn}\bar{e}_m \end{bmatrix} \\ &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} t_{11} & t_{21} & \cdots & t_{m1} \\ t_{12} & t_{22} & \cdots & t_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ t_{1n} & t_{2n} & \cdots & t_{mn} \end{bmatrix} \begin{bmatrix} \bar{e}_1 \\ \bar{e}_2 \\ \cdots \\ \bar{e}_m \end{bmatrix} \end{aligned}$$

We may take the transpose,

$$\begin{aligned} T(\bar{v}) &= \begin{bmatrix} \bar{e}_1 & \bar{e}_2 & \cdots & \bar{e}_m \end{bmatrix} \begin{bmatrix} t_{11} & t_{12} & \cdots & t_{1n} \\ t_{21} & t_{22} & \cdots & t_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ t_{m1} & t_{m2} & \cdots & t_{mn} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \cdots \\ v_n \end{bmatrix} \\ T(\bar{v}) &= T \left( \sum_{k=1}^n v_k \bar{u}_k \right) = \sum_{k=1}^n v_k T(\bar{u}_k) = \sum_{k=1}^n v_k \sum_{j=1}^m t_{kj} \bar{e}_j \end{aligned}$$

Thus matrix of  $T$  is given by,  $m(T) = (t_{ik})$  where  $T(\bar{u}_k) = \sum_{i=1}^n t_{ik} \bar{e}_i$ .

*Remark (Example).* Let  $T : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  defined by  $T(x, y, z) = (2x + y, y - z)$ .

$$\begin{aligned}
 T(1, 2, 3) &= T((1, 0, 0) + 2(0, 1, 0) + 3(0, 0, 1)) \\
 &= T(\bar{u}_1 + 2\bar{u}_2 + 3\bar{u}_3) \\
 &= T(\bar{u}_1) + 2T(\bar{u}_2) + 3T(\bar{u}_3) \\
 &= [1 \quad 2 \quad 3] \begin{bmatrix} T(\bar{u}_1) \\ T(\bar{u}_2) \\ T(\bar{u}_3) \end{bmatrix} \\
 &= [1 \quad 2 \quad 3] \begin{bmatrix} (2, 0) \\ (1, 1) \\ (0, -1) \end{bmatrix} \\
 &= [1 \quad 2 \quad 3] \begin{bmatrix} 2(1, 0) \\ (1, 0) + (0, 1) \\ -1(0, 1) \end{bmatrix} \\
 &= [1 \quad 2 \quad 3] \begin{bmatrix} 2\bar{e}_1 \\ \bar{e}_1 + \bar{e}_2 \\ -\bar{e}_2 \end{bmatrix} \\
 &= [1 \quad 2 \quad 3] \begin{bmatrix} 2 & 0 \\ 1 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} \bar{e}_1 \\ \bar{e}_2 \end{bmatrix} \\
 &= 4\bar{e}_1 - \bar{e}_2 = 4(1, 0) - 1(0, 1) = (4, -1)
 \end{aligned}$$

$$\text{In the above case, } m(T) = \begin{bmatrix} 2 & 0 \\ 1 & 1 \\ 0 & -1 \end{bmatrix}$$

Using the matrix of linear function  $m(T)$ , we can compute the image of any point in  $\mathbb{R}^3$  by matrix multiplication.

### Matrix of the composition of two linear functions

Let  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $S : \mathbb{R}^m \rightarrow \mathbb{R}^p$  be two linear functions with domain of  $S$  containing the range of  $T$  (so that  $S \circ T$  is well defined). Then  $S \circ T : \mathbb{R}^n \rightarrow \mathbb{R}^p$  is defined by

$$S \circ T(\bar{x}) = S(T(\bar{x})), \quad \forall \bar{x} \in \mathbb{R}^n$$

Since  $S, T$  are linear,  $S \circ T$  is also linear.

$$\begin{aligned}
 S \circ T(a\bar{x} + b\bar{y}) &= S(T(a\bar{x} + b\bar{y})) = S(aT(\bar{x}) + bT(\bar{y})) = aS(T(\bar{x})) + bS(T(\bar{y})) \\
 &= aS \circ T(\bar{x}) + bS \circ T(\bar{y}), \quad \forall a, b \in \mathbb{R}, \quad \forall \bar{x}, \bar{y} \in \mathbb{R}^n
 \end{aligned}$$

Let  $\{\bar{u}_1, \bar{u}_2, \dots, \bar{u}_n\}$  be the standards basis for  $\mathbb{R}^n$ ,  $\{\bar{e}_1, \bar{e}_2, \dots, \bar{e}_m\}$  be the standards basis for  $\mathbb{R}^m$  and  $\{\bar{w}_1, \bar{w}_2, \dots, \bar{w}_p\}$  be the standards basis for  $\mathbb{R}^p$ . Let

$\bar{v} \in \mathbb{R}^n$ , then  $\bar{v} = \sum_{i=1}^n v_i \bar{u}_i$ , and  $S \circ T(\bar{v}) = \sum_{i=1}^n v_i S \circ T(\bar{u}_i)$

$$\begin{aligned}
 S \circ T(\bar{v}) &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} S \circ T(\bar{u}_1) \\ S \circ T(\bar{u}_2) \\ \vdots \\ S \circ T(\bar{u}_n) \end{bmatrix} \\
 &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} S(t_{11}\bar{e}_1 + \cdots + t_{m1}\bar{e}_m) \\ S(t_{12}\bar{e}_1 + \cdots + t_{m2}\bar{e}_m) \\ \vdots \\ S(t_{1n}\bar{e}_1 + \cdots + t_{mn}\bar{e}_m) \end{bmatrix} \\
 &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} t_{11}S(\bar{e}_1) + \cdots + t_{m1}S(\bar{e}_m) \\ t_{12}S(\bar{e}_1) + \cdots + t_{m2}S(\bar{e}_m) \\ \vdots \\ t_{1n}S(\bar{e}_1) + \cdots + t_{mn}S(\bar{e}_m) \end{bmatrix} \\
 &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} t_{11} & t_{21} & \cdots & t_{m1} \\ t_{12} & t_{22} & \cdots & t_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ t_{1n} & t_{2n} & \cdots & t_{mn} \end{bmatrix} \begin{bmatrix} S(\bar{e}_1) \\ S(\bar{e}_2) \\ \vdots \\ S(\bar{e}_m) \end{bmatrix} \\
 &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} t_{11} & t_{21} & \cdots & t_{m1} \\ t_{12} & t_{22} & \cdots & t_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ t_{1n} & t_{2n} & \cdots & t_{mn} \end{bmatrix} \begin{bmatrix} s_{11}\bar{w}_1 + s_{12}\bar{w}_2 + \cdots + s_{1p}\bar{w}_p \\ s_{12}\bar{w}_1 + s_{22}\bar{w}_2 + \cdots + s_{p2}\bar{w}_p \\ \vdots \\ s_{1m}\bar{w}_1 + s_{2m}\bar{w}_2 + \cdots + s_{pm}\bar{w}_p \end{bmatrix} \\
 &= \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} t_{11} & t_{21} & \cdots & t_{m1} \\ t_{12} & t_{22} & \cdots & t_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ t_{1n} & t_{2n} & \cdots & t_{mn} \end{bmatrix} \begin{bmatrix} s_{11} & s_{21} & \cdots & s_{p1} \\ s_{12} & s_{22} & \cdots & s_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ s_{1m} & s_{2m} & \cdots & s_{pm} \end{bmatrix} \begin{bmatrix} \bar{w}_1 \\ \bar{w}_2 \\ \vdots \\ \bar{w}_p \end{bmatrix}
 \end{aligned}$$

We may take transpose,

$$S \circ T(\bar{v}) = \begin{bmatrix} \bar{w}_1 & \bar{w}_2 & \cdots & \bar{w}_p \end{bmatrix} \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1m} \\ s_{21} & s_{22} & \cdots & s_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ s_{p1} & s_{p2} & \cdots & s_{pm} \end{bmatrix} \begin{bmatrix} t_{11} & t_{12} & \cdots & t_{1n} \\ t_{21} & t_{22} & \cdots & t_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ t_{m1} & t_{m2} & \cdots & t_{mn} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

Remember : Given  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , then we may take  $m(T)$  either as  $m \times n$  matrix or  $n \times m$  matrix. Since, we chose  $m \times n$ ,  $m(S \circ T) = m(S)m(T)$ . Otherwise,  $m(S \circ T) = m(T)m(S)$ . This may change for different authors.

Suppose  $m(S) = (s_{ij})$  and  $m(T) = (t_{ij})$  respectively. Then

$$S(e_k) = \sum_{i=1}^p s_{ik}\bar{w}_i, \quad k = 1, 2, \dots, m \text{ and}$$

$$T(u_j) = \sum_{k=1}^m t_{kj}\bar{e}_k, \quad j = 1, 2, \dots, n$$

$$\begin{aligned}
(S \circ T)(\bar{u}_j) &= S(T(\bar{u}_j)) = S\left(\sum_{k=1}^m t_{kj} \bar{e}_k\right) = \sum_{k=1}^m t_{kj} S(\bar{e}_k) \\
&= \sum_{k=1}^m t_{kj} \left(\sum_{i=1}^p s_{ik} \bar{w}_i\right) = \sum_{i=1}^p \left(\sum_{k=1}^m s_{ik} t_{kj}\right) \bar{w}_i
\end{aligned}$$

Therefore,  $m(S \circ T) = \sum_{k=1}^m s_{ik} t_{kj} = (s_{ik})(t_{kj}) = m(S)m(T)$ .

### 13.2.4 The Jacobian Matrix

Let  $\bar{u}_1, \bar{u}_2, \dots, \bar{u}_n$  be the unit co-ordinate vectors in  $\mathbb{R}^n$  and  $\bar{e}_1, \bar{e}_2, \dots, \bar{e}_m$  be the unit co-ordinate vectors in  $\mathbb{R}^m$ . Let function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be differentiable at  $\bar{c} \in \mathbb{R}^n$ . Then there exists a linear function  $T = f'(\bar{c}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that  $f(\bar{c} + \bar{v}) = f(\bar{c}) + f'(\bar{c})(\bar{v}) + \|\bar{v}\| + E_{\bar{c}}(\bar{v})$ . We have,  $T(\bar{u}_k) = f'(\bar{c})(\bar{u}_k) = f'(\bar{c}, \bar{u}_k) = D_k f(\bar{c}) = D_k \sum_{i=1}^m f_i(\bar{c}) \bar{e}_i$ .

Clearly, the matrix of total derivative  $T$ ,  $m(T) = (t_{ik}) = (D_k f_i(\bar{c}))$ . This matrix is called Jacobian matrix of  $f$  at  $\bar{c}$  and is denoted by  $Df(\bar{c})$ .

$$Df(\bar{c}) = \begin{bmatrix} D_1 f_1(\bar{c}) & D_2 f_1(\bar{c}) & \cdots & D_n f_1(\bar{c}) \\ D_1 f_2(\bar{c}) & D_2 f_2(\bar{c}) & \cdots & D_n f_2(\bar{c}) \\ \vdots & \vdots & \ddots & \vdots \\ D_1 f_m(\bar{c}) & D_2 f_m(\bar{c}) & \cdots & D_n f_m(\bar{c}) \end{bmatrix}$$

#### Properties of Jacobian matrix

1.  $k$ th row of  $Df(\bar{c})$  is gradient vector of  $f_k$

$$\nabla f_k(\bar{c}) = (D_1 f_k(\bar{c}), D_2 f_k(\bar{c}), \dots, D_n f_k(\bar{c}))$$

2. When  $m = 1$ ,  $Df(\bar{c}) = \nabla f(\bar{c})$ .

$$3. f'(\bar{c})(\bar{v}) = \sum_{k=1}^m (\nabla f_k(\bar{c}) \cdot \bar{v}) \bar{e}_k$$

$$4. \|f'(\bar{c})(\bar{v})\| \leq M \|\bar{v}\| \text{ where } M = \sum_{k=1}^m \|\nabla f_k(\bar{c})\|, \text{ by property (3)}$$

5.  $f'(\bar{c})(\bar{v}) \rightarrow \bar{0}$  as  $\bar{v} \rightarrow \bar{0}$ , by property (4)

#### Chain Rule

Chain Rule for real function :  $\frac{dF \circ G}{dx}(x) = \frac{d}{dy} F(y) \frac{d}{dx} G(x) = F'(y) G'(x)$

For example :  $\frac{d}{dx}(ax+3)^3 = \frac{d}{dy} y^3 \frac{d}{dx}(ax+3) = 3ay^2 = 3a(ax+3)^2$

**Theorem 13.2.4.** Let  $g$  be differentiable at  $\bar{a}$ , with total derivative  $g'(\bar{a})$  and  $\bar{b} = g(\bar{a})$ . Let  $f$  is differentiable at  $\bar{b}$ , with total derivative  $f'(\bar{b})$ . Then  $h = f \circ g$  is differentiable at  $\bar{a}$  with total derivative  $h'(\bar{a}) = f'(\bar{b}) \circ g'(\bar{a})$ . *Try to read  $h'(\bar{a}) = H$ ,  $f'(\bar{b}) = F$ ,  $g'(\bar{a}) = G$ , then  $H = F \circ G \implies H(x) = F(G(x))$  In other words,  $h'(\bar{a})(\bar{v}) = f'(\bar{b}) \circ g'(\bar{a})(\bar{v}) = f'(\bar{b})(g'(\bar{a})(\bar{v}))$ .*

*Proof.* Given  $\epsilon > 0$ , let  $y \in \mathbb{R}^p$  such that  $\|y\| < \epsilon$ . Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$ , then  $h = f \circ g : \mathbb{R}^p \rightarrow \mathbb{R}^m$ .

We have,  $h(\bar{a} + \bar{y}) - h(\bar{a}) = f(g(\bar{a} + \bar{y})) - f(g(\bar{a})) = f(\bar{b} + \bar{v}) - f(\bar{b})$  where  $\bar{b} = g(\bar{a})$ , and  $\bar{v} = g(\bar{a} + \bar{y}) - g(\bar{a})$ .

Since  $g$  is differentiable at  $\bar{a}$ ,  $g$  satisfies first-order Taylor's formula.

$$\begin{aligned} g(\bar{a} + \bar{y}) &= g(\bar{a}) + g'(\bar{a})(\bar{y}) + \|\bar{y}\|E_{\bar{a}}(\bar{y}) \text{ where } E_{\bar{a}} \rightarrow \bar{0} \text{ as } \bar{y} \rightarrow \bar{0} \\ \implies \bar{v} &= g(\bar{a} + \bar{y}) - g(\bar{a}) = g'(\bar{a})(\bar{y}) + \|\bar{y}\|E_{\bar{a}}(\bar{y}) \end{aligned}$$

Clearly, as  $\bar{y} \rightarrow \bar{0} \implies \bar{v} \rightarrow g'(\bar{a})(\bar{0}) = \bar{0}$ . Again, we have  $f$  is differentiable at  $\bar{b}$ , thus  $f$  satisfies first-order Taylor's formula.

$$f(\bar{b} + \bar{v}) = f(\bar{b}) + f'(\bar{b})(\bar{v}) + \|\bar{v}\|E_{\bar{b}}(\bar{v}) \text{ where } E_{\bar{b}} \rightarrow \bar{0} \text{ as } \bar{v} \rightarrow \bar{0}$$

$$\begin{aligned} \implies f(\bar{b} + \bar{v}) - f(\bar{b}) &= f'(\bar{b})(\bar{v}) + \|\bar{v}\|E_{\bar{b}}(\bar{v}) \\ &= f'(\bar{b})(g'(\bar{a})(\bar{y}) + \|\bar{y}\|E_{\bar{a}}(\bar{y})) + \|\bar{v}\|E_{\bar{b}}(\bar{v}) \\ &= f'(\bar{b})(g'(\bar{a})(\bar{y})) + \|\bar{y}\|E(\bar{y}) \\ &\text{ where } E(\bar{y}) = f'(\bar{b})(E_{\bar{a}}(\bar{y})) + \frac{\|\bar{v}\|}{\|\bar{y}\|}E_{\bar{b}}(\bar{v}), \bar{y} \neq \bar{0} \end{aligned}$$

$$\implies h(\bar{a} + \bar{y}) - h(\bar{a}) = f(\bar{b} + \bar{v}) - f(\bar{b}) = f'(\bar{b})(g'(\bar{a})(\bar{y})) + \|\bar{y}\|E(\bar{y})$$

Since  $f'(\bar{b})$  and  $g'(\bar{a})$  are linear, their composition is also linear. Therefore,  $h$  is differentiable at  $\bar{a}$  with a linear, total derivative  $h'(\bar{a}) = f'(\bar{b}) \circ g'(\bar{a})$  as it satisfies first-order Taylor's formula if  $E_{\bar{y}} \rightarrow \bar{0}$  as  $\bar{y} \rightarrow \bar{0}$ .

We have,  $\|\bar{v}\| \leq \|g'(\bar{a})(\bar{y})\| + \|\bar{y}\| \|E_{\bar{a}}(\bar{y})\| \leq M\|\bar{y}\| + \|E_{\bar{a}}(\bar{y})\| \|\bar{y}\|$ .

$$\implies \frac{\|\bar{v}\|}{\|\bar{y}\|} \leq M + \|E_{\bar{a}}(\bar{y})\|$$

Thus,  $\bar{v} \rightarrow \bar{0}$  as  $\bar{y} \rightarrow \bar{0}$ . Then  $f'(\bar{b})(\bar{v}) \rightarrow f'(\bar{b})(\bar{0}) = \bar{0}$ . And  $E_{\bar{a}}(\bar{y}) \rightarrow \bar{0}$ . Therefore,  $E(\bar{y}) \rightarrow \bar{0} + M\bar{0} = \bar{0}$  as  $\bar{y} \rightarrow \bar{0}$ .  $\square$

### Matrix form of the chain rule

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$ . And  $h = f \circ g : \mathbb{R}^p \rightarrow \mathbb{R}^m$ . Suppose  $g$  is differentiable at  $\bar{a} \in \mathbb{R}^p$  and  $f$  is differentiable at  $g(\bar{a}) = \bar{b} \in \mathbb{R}^n$ . Then  $h$  is differentiable at  $\bar{a}$  and the Jacobian matrix of  $h$  is given by the chain rule,

$$Dh(\bar{a}) = Df(\bar{b})Dg(\bar{a}) \text{ where } h = f \circ g, \bar{b} = g(\bar{a})$$

In other words,

$$D_j h_i(\bar{a}) = \sum_{k=1}^n D_k f_i(\bar{b}) D_j g_k(\bar{a}), \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, p$$

$$\text{For } m = 1, D_j h(\bar{a}) = \sum_{k=1}^n D_k f(\bar{b}) D_j g_k(\bar{a})$$

$$\text{For } m = 1 \text{ and } p = 1, h'(\bar{a}) = \sum_{k=1}^n Df(\bar{b})g'_k(\bar{a}) = \nabla f(\bar{b}) \cdot Dg(\bar{a})$$

**Theorem 13.2.5.** Let  $f$  and  $D_2f$  be continuous functions on a rectangle  $[a, b] \times [c, d]$ . Let  $p$  and  $q$  be differentiable on  $[c, d]$ , where  $p(y) \in [a, b]$  and  $q(y) \in [c, d]$  for each  $y \in [c, d]$ . Define  $F$  by the equation,

$$F(y) = \int_{p(y)}^{q(y)} f(x, y) dx, \quad y \in [c, d]$$

Then  $F'(y)$  exists for each  $y \in (c, d)$  and is given by,

$$F'(y) = \int_{p(y)}^{q(y)} D_2f(x, y) dx + f(q(y), y)q'(y) - f(p(y), y)p'(y)$$

The following two theorems are required for proving the theorem on differentiating an integral.

**Theorem 13.2.6.** Let  $\alpha$  be of bounded variation on  $[a, b]$  and assume that  $f \in \mathcal{R}(\alpha)$  on  $[a, b]$ .

$$\text{Define } F(x) = \int_a^x f \, d\alpha, \quad x \in [a, b]$$

Then  $F$  is of bounded variation on  $[a, b]$  and  $F$  is continuous at  $x$  if  $\alpha$  is continuous at  $x$ . If  $\alpha$  is increasing on  $[a, b]$ , then the derivative  $F'(x)$  exists at each  $x \in (a, b)$  where  $\alpha'(x)$  exists and where  $f$  is continuous. And

$$F'(x) = f(x)\alpha'(x)$$

**Theorem 13.2.7.** Let  $Q = \{(x, y) : a \leq x \leq b, c \leq y \leq d\}$ . Assume that  $\alpha$  is of bounded variation on  $[a, b]$  and for each  $y \in [c, d]$ , assume that the integral

$$F(y) = \int_a^b f(x, y) \, d\alpha(x)$$

exists. If the partial derivative  $D_2f$  is continuous on  $Q$ , the derivative  $F'(y)$  exists for each  $y \in (c, d)$  and is given by

$$F'(y) = \int_a^b D_2f(x, y) \, d\alpha(x)$$

*Proof.* Let  $G(x_1, x_2, x_3) = \int_{x_1}^{x_2} f(t, x_3) \, dt$ . Then we may write  $F(y)$  in terms of  $G$ . That is,  $F(y) = G(p(y), q(y), y)$ .

**Step 1 : 1-D Chain Rule**

By 1-dimensional chain rule, we have

$$\begin{aligned} F'(y) &= \frac{dF}{dy} = \frac{\partial G}{\partial p} \frac{dp}{dy} + \frac{\partial G}{\partial q} \frac{dq}{dy} + \frac{\partial G}{\partial y} \\ &= D_1G \, p'(y) + D_2G \, q'(y) + D_3G \end{aligned}$$

**Step 2 :  $D_1G$**



Since the variable of differentiation is present in the limit of the integral, we use theorem 13.2.6 to compute the derivative of the integral. We may write,

$$G(p(y), q(y), y) = - \int_{q(y)}^{p(y)} f(t, y) dt \quad (13.6)$$

We are differentiating (partially) with respect to  $p(y)$ . Thus  $q(y)$ ,  $y$  are constants for this differentiation.

$$\begin{aligned} G(x, a, y) &= -H(x) = - \int_a^x f(t, y) dt \\ \implies D_1 G &= -H'(x) = -f(x, y). \\ \text{Thus, } D_1 G &= -f(p(y), y) \end{aligned}$$

### Step 3 : $D_2 G$

Again, variable of differentiation is present in the limit of the integral. Thus, we write,

$$G(p(y), q(y), y) = \int_{p(y)}^{q(y)} f(t, y) dt \quad (13.7)$$

Now we are differentiating (partially) with respect to the the second component of  $G$  which is  $q(y)$ . Clearly,  $p(y)$  and  $y$  are treated as constants.

$$\begin{aligned} G(a, x, y) &= H(x) = \int_a^x f(t, y) dt \\ \implies D_2 G &= H'(x) = f(q(y), y) \end{aligned}$$

### Step 4 : $D_3 G$

Now the variable of integration is not affecting the limits of the integral. Also it is given that  $D_2 f$  is continuous on  $[a, b] \times [c, d]$ . We write

$$\begin{aligned} G(a, b, x) &= H(x) = \int_a^b f(t, x) dt \\ \implies D_3 G &= H'(x) = \int_a^b D_2 f(t, x) dt \\ \text{Thus, } D_3 G &= \int_{p(y)}^{q(y)} f(t, y) dt \end{aligned}$$

□

## The mean-value theorem for differentiable functions

**Theorem 13.2.8** (Mean-Value). *Let  $S$  be an open subset of  $\mathbb{R}^n$ . Assume  $f : S \rightarrow \mathbb{R}^m$  is differentiable at each point of  $S$ . Let  $\bar{x}, \bar{y}$  be two points in  $S$  such that  $L(\bar{x}, \bar{y}) = \{t\bar{x} + (1-t)\bar{y} : t \in [0, 1]\}$  is subset of  $S$ . Then for every  $\bar{a} \in \mathbb{R}^m$ , there exists a point  $\bar{z} \in L(\bar{x}, \bar{y})$  such that*

$$\bar{a} \cdot (f(\bar{y}) - f(\bar{x})) = \bar{a} \cdot f'(\bar{z})(\bar{y} - \bar{x})$$

*Proof.* Let  $\bar{u} = \bar{y} - \bar{x}$ . We have  $S$  is open subset and  $L(\bar{x}, \bar{y}) \subset S$ , thus there exists  $\delta > 0$  such that  $\bar{x} + t\bar{u} \in S, \forall t \in (-\delta, 1 + \delta)$ . In other words, the 'Line

segment  $L(\bar{x}, \bar{y})'$  is properly contained in  $S$ , in such a way that extending the Line from  $\bar{x}$  to  $\bar{y}$  a little bit extra one either sides is still contained in  $S$ .

Let  $\bar{a} \in \mathbb{R}^m$  and  $F : (-\delta, 1 + \delta) \rightarrow \mathbb{R}$  defined by  $F(t) = \bar{a} \cdot f(\bar{x} + t\bar{u})$ . Then  $F$  is differentiable at each  $t \in (-\delta, 1 + \delta)$  and the derivative  $F'(t) = \bar{a} \cdot f'(\bar{x} + t\bar{u}, \bar{u})$ , the directional derivative of  $f(\bar{x} + t\bar{u})$  with respect to  $\bar{u}$ .

$$f'(\bar{x} + t\bar{u}, \bar{u}) = f'(\bar{x} + t\bar{u})(\bar{u}) \implies F'(t) = \bar{a} \cdot f'(\bar{x} + t\bar{u})(\bar{u})$$

By 1-dimensional mean-value theorem, we have

$$\exists \theta \in (0, 1) \text{ such that } F(1) - F(0) = F'(\theta)$$

By definition of  $F$ ,  $F(1) = \bar{a} \cdot f(\bar{x} + \bar{u}) = \bar{a} \cdot f(\bar{y})$ . And  $F(0) = \bar{a} \cdot f(\bar{x})$ . Therefore,

$$F'(\theta) = F(1) - F(0) = \bar{a} \cdot f(\bar{y}) - \bar{a} \cdot f(\bar{x}) = \bar{a} \cdot (f(\bar{y}) - f(\bar{x}))$$

We also have,

$$F'(\theta) = \bar{a} \cdot f'(\bar{x} + \theta\bar{u})(\bar{u}) = \bar{a} \cdot f'(\bar{z})(\bar{y} - \bar{x}), \text{ where } \bar{z} = \bar{x} + \theta\bar{u} \in L(\bar{x}, \bar{y})$$

□

*Remark.* Suppose  $S$  is convex in  $\mathbb{R}^m$ . Then for every pair of points  $\bar{x}, \bar{y} \in S$ ,  $L(\bar{x}, \bar{y}) \subset S$ . Thus Mean-value theorem holds for all  $\bar{x}, \bar{y} \in S$ .

## 13.3 Multivariate Calculus

### 13.3.1 A sufficient condition for differentiability

**Theorem 13.3.1.** Suppose one of the partial derivatives  $D_1f, D_2f, \dots, D_nf$  exists at  $\bar{c}$ . And the remaining  $n - 1$  partial derivatives exists in some  $n$ -ball  $B(\bar{c})$  and are continuous at  $\bar{c}$ . Then  $f$  is differentiable at  $\bar{c}$ .

*Proof.* **Step 1 : Real-valued function**

We claim that the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable at  $\bar{c}$  iff each component  $f_k$  is differentiable at  $\bar{c}$ .

Suppose  $f$  is differentiable at  $\bar{c}$ , then there exists a linear, total derivative function  $f'(\bar{c})$  satisfying first-order Taylor's formula at  $\bar{c}$ .

ie,  $f(\bar{c} + \bar{v}) = f(\bar{c}) + f'(\bar{c})(\bar{v}) + \|\bar{v}\|E_{\bar{c}}(\bar{v})$  where  $E_{\bar{c}}(\bar{v}) \rightarrow \bar{0}$  as  $\bar{v} \rightarrow \bar{0}$ .

$$f(\bar{c} + \bar{v}) = (f_1(\bar{c} + \bar{v}), f_2(\bar{c} + \bar{v}), \dots, f_m(\bar{c} + \bar{v}))$$

$$f(\bar{c}) = (f_1(\bar{c}), f_2(\bar{c}), \dots, f_m(\bar{c}))$$

$$f'(\bar{c})(\bar{v}) = (f'_1(\bar{v}), f'_2(\bar{v}), \dots, f'_m(\bar{v}))$$

$$E_{\bar{c}}(\bar{v}) = (E_1(\bar{v}), E_2(\bar{v}), \dots, E_m(\bar{v}))$$

where each component of the error function  $E_k(\bar{v}) \rightarrow 0$  as  $\bar{v} \rightarrow \bar{0}$ . Also since  $f'(\bar{c})$  is linear, each of its components  $f'_k : \mathbb{R}^n \rightarrow \mathbb{R}$  are linear.

$$\begin{aligned} f(\bar{c} + \bar{v}) &= (f_1(\bar{c} + \bar{v}), f_2(\bar{c} + \bar{v}), \dots, f_m(\bar{c} + \bar{v})) \\ &= (f_1(\bar{c}), f_2(\bar{c}), \dots, f_m(\bar{c})) + (f'_1(\bar{v}), f'_2(\bar{v}), \dots, f'_m(\bar{v})) \\ &\quad + \|\bar{v}\| (E_1(\bar{v}), E_2(\bar{v}), \dots, E_m(\bar{v})) \text{ where } E_k(\bar{v}) \rightarrow 0 \text{ as } \bar{v} \rightarrow \bar{0} \\ &= (f_1(\bar{c}), f_2(\bar{c}), \dots, f_m(\bar{c})) + (f'_1(\bar{v}), f'_2(\bar{v}), \dots, f'_m(\bar{v})) \\ &\quad + (\|\bar{v}\|E_1(\bar{v}), \|\bar{v}\|E_2(\bar{v}), \dots, \|\bar{v}\|E_m(\bar{v})) \text{ where } E_k(\bar{v}) \rightarrow 0 \text{ as } \bar{v} \rightarrow \bar{0} \\ &= (f_1(\bar{c}) + f'_1(\bar{v}) + \|\bar{v}\|E_1(\bar{v}), \dots, f_m(\bar{c}) + f'_m(\bar{v}) + \|\bar{v}\|E_m(\bar{v})) \end{aligned}$$

Thus first-order Taylor's formula for  $f$  at  $\bar{c}$  gives first-order Taylor's formula for each of its components  $f_k$ . ie,  $f_k(\bar{c} + \bar{v}) = f_k(\bar{c}) + f'_k(\bar{v} + \|\bar{v}\|E_k(\bar{v}))$  where  $E_k(\bar{v}) \rightarrow 0$  as  $\bar{v} \rightarrow \bar{0}$ . Therefore,  $f_k$  are differentiable at  $\bar{c}$  for  $k = 1, 2, \dots, m$ .

Suppose each component  $f_k$  of  $f$  are differentiable at  $\bar{c}$ . Then there exists linear, total derivative functions  $f'_k$  satisfying first-order Taylor's formula at  $\bar{c}$ . ie,  $f_k(\bar{c} + \bar{v}) = f_k(\bar{c}) + f'_k(\bar{v}) + \|\bar{v}\|E_k(\bar{v})$  where  $E_k(\bar{v}) \rightarrow 0$  as  $\bar{v} \rightarrow \bar{0}$ .

Define  $E_{\bar{c}}(\bar{v}) = (E_1(\bar{v}), E_2(\bar{v}), \dots, E_m(\bar{v}))$ . Then  $E_{\bar{c}}(\bar{v}) \rightarrow \bar{0}$  as  $\bar{v} \rightarrow \bar{0}$ . Therefore, there exists a linear, total derivative function  $f'(\bar{c}) = (f'_1, f'_2, \dots, f'_m)$  satisfying first-order Taylor's formula at  $\bar{c}$ .

Thus, if each (real-valued) component function  $f_k$  are differentiable, then  $f$  is also differentiable. Therefore, it is sufficient to prove the theorem for a real-valued function.

### Step 2: Telescopic Sum

Assume (without loss of generality) that  $D_1 f$  exists at  $\bar{c}$  and  $D_2 f, D_3 f, \dots, D_n f$  exist and continuous in some  $n$ -ball  $B(\bar{c})$ . Suppose  $D_r f$  exists at  $\bar{c}$  and all partial derivatives except  $D_r f$  are continuous. Then  $v_0 = \bar{0}$ ,  $v_1 = y_r \bar{u}_r$ ,  $v_2 = y_r \bar{u}_r + y_1 \bar{u}_1, \dots$ . Then the following proof can be applied without any loss of generality.

Let  $\bar{v} = \lambda \bar{y}$  where  $\bar{y} = \frac{\bar{v}}{\|\bar{v}\|}$ . Clearly,  $\|\bar{y}\| = 1$  and  $\lambda = \|\bar{v}\|$ . Choose  $\lambda > 0$  such that  $\bar{c} + \bar{v} \in B(\bar{c})$  and all the partial derivatives  $D_2 f, D_3 f, \dots, D_n f$  exists and are continuous in  $B(\bar{c})$ .

We have,  $\bar{y} = (y_1, y_2, \dots, y_n) = y_1 \bar{u}_1 + y_2 \bar{u}_2 + \dots + y_n \bar{u}_n$ .

Define  $\bar{v}_0 = \bar{0}$ ,  $\bar{v}_1 = y_1 \bar{u}_1, \dots$ ,  $\bar{v}_n = y_1 \bar{u}_1 + y_2 \bar{u}_2 + \dots + y_n \bar{u}_n$ .

$$\begin{aligned} f(\bar{c} + \bar{v}) - f(\bar{c}) &= (f(\bar{c} + \lambda \bar{v}_n) - f(\bar{c} + \lambda \bar{v}_{n-1})) \\ &\quad + (f(\bar{c} + \lambda \bar{v}_{n-1}) - f(\bar{c} + \lambda \bar{v}_{n-2})) \\ &\quad + \dots + (f(\bar{c} + \lambda \bar{v}_1) - f(\bar{c} + \lambda \bar{v}_0)) \\ &= \sum_{k=1}^n f(\bar{c} + \lambda \bar{v}_k) - f(\bar{c} + \lambda \bar{v}_{k-1}) \\ &= \sum_{k=1}^n f(\bar{c} + \lambda \bar{v}_{k-1} + \lambda y_k \bar{u}_k) - f(\bar{c} + \lambda \bar{v}_{k-1}) \end{aligned}$$

### Step 3 : Mean-value theorem

Define  $\bar{b}_k = \bar{c} + \lambda \bar{v}_{k-1}$ . Then we have

$$f(\bar{c} + \bar{v}) - f(\bar{c}) = \sum_{k=1}^n f(\bar{b}_k + \lambda y_k \bar{u}_k) - f(\bar{b}_k) \quad (13.8)$$

We know that all partial derivatives exists in  $B(\bar{c})$ . Therefore by 1-dimensional mean-value theorem we have,

$$f(\bar{b}_k + \lambda y_k \bar{u}_k) - f(\bar{b}_k) = \lambda y_k D_k f(\bar{a}_k) \text{ where } \bar{a}_k \in L(\bar{b}_k, \bar{b}_k + \lambda y_k \bar{u}_k)$$

$$f(\bar{c} + \bar{v}) - f(\bar{c}) = \lambda \sum_{k=1}^n y_k D_k f(\bar{a}_k) \text{ where } \bar{a}_k \in L(\bar{b}_k, \bar{b}_k + \lambda y_k \bar{u}_k) \quad (13.9)$$

### Step 4 : Continuity of partial derivatives in $B(\bar{c})$

As  $\lambda \rightarrow 0$ ,  $\bar{v} \rightarrow \bar{0}$ . And both  $\bar{b}_k$ ,  $\bar{b}_k + \lambda y_k \bar{u}_k \rightarrow \bar{c}$ . Clearly,  $\bar{a}_k$  in the line between  $\bar{b}_k$  and  $\bar{b}_k + \lambda y_k \bar{u}_k$  also converges to  $\bar{c}$ .

For  $k \geq 2$ ,  $D_k f$  are continuous in the  $n$ -ball  $B(\bar{c})$ . Thus  $D_k f(\bar{a}_k) \rightarrow D_k f(\bar{c})$ . We may write,  $D_k f(\bar{a}_k) = D_k f(\bar{c}) + E_k(\lambda)$  where  $E_k(\lambda) \rightarrow \bar{0}$  as  $\lambda \rightarrow 0$ . Also, since  $D_1 f$  exists,  $D_1 f(\bar{c} + \lambda y_1 \bar{u}_1) \rightarrow D_1 f(\bar{c})$  as  $\lambda \rightarrow 0$ .

$$\text{Remember : } D_1 f(\bar{c}) = \lim_{h \rightarrow 0} \frac{f(\bar{c} + h \bar{u}_1) - f(\bar{c})}{h}$$

$$\begin{aligned} f(\bar{c} + \bar{v}) - f(\bar{c}) &= \lambda \sum_{k=1}^n y_k D_k f(\bar{c}) + \lambda \sum_{k=1}^n y_k E_k(\lambda) \\ &= \nabla f(\bar{c}) \cdot \bar{v} + \|\bar{v}\| E(\lambda) \\ \text{where } E(\lambda) &= \sum_{k=1}^n y_k E_k(\lambda) \rightarrow \bar{0} \text{ as } \bar{v} \rightarrow \bar{0} \end{aligned}$$

That is, we have a linear function which satisfies first-order Taylor's formula at  $\bar{c}$ . Therefore,  $f$  is differentiable.  $\square$

### 13.3.2 Sufficient conditions for the equality of mixed partial derivatives

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . Then  $D_r f$  and  $D_k f$  are two partial derivatives of  $f$ . And  $D_{r,k} f = D_r(D_k f)$  and  $D_{k,r} f = D_k(D_r f)$ .

$$D_{r,k} f = \frac{\partial^2 f}{\partial x_r \partial x_k} = \frac{\partial}{\partial x_r} \frac{\partial f}{\partial x_k} \text{ and } D_{k,r} f = \frac{\partial^2 f}{\partial x_k \partial x_r} = \frac{\partial}{\partial x_k} \frac{\partial f}{\partial x_r}$$

There are two sufficient conditions for the equality of these mixed partial derivatives in our scope. 1. differentiability of  $D_k f$  or 2. continuity of  $D_{r,k} f$  and  $D_{k,r} f$  at  $\bar{c}$  where the mixed partial derivatives are to be equal.

#### Differentiability

**Theorem 13.3.2.** Suppose  $D_r f$  and  $D_k f$  exists in an  $n$ -ball about  $\bar{c}$  and are both differentiable at  $\bar{c}$ . Then  $D_{r,k} f = D_{k,r} f$ .

*Proof.* **Step 1 : Real-valued function**

It is sufficient to prove the theorem for real-valued functions. Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , then  $f(\bar{c}) = (f_1(\bar{c}), f_2(\bar{c}), \dots, f_m(\bar{c}))$ . And

$$D_k f(\bar{c}) = (D_k f_1(\bar{c}), D_k f_2(\bar{c}), \dots, D_k f_m(\bar{c}))$$

Thus it is sufficient to prove that  $D_{r,k} f_j(\bar{c}) = D_{k,r} f_j(\bar{c})$ ,  $j = 1, 2, \dots, m$ . That is, it is sufficient to prove equality of mixed partial derivatives of a real-valued function  $f_j : \mathbb{R}^n \rightarrow \mathbb{R}$ . Also, we will prove it for  $n = 2$  and  $\bar{c} = (0, 0)$ . **Now, we will prove the theorem of a real-valued function  $f : \mathbb{R}^n : \mathbb{R}$ .**

**Step 2 :  $\nabla(h)$**

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . Suppose that the partial derivatives  $D_k f$ ,  $D_r f$  exist in the  $n$ -ball  $B(n)$ . And let  $h > 0$  such that the rectangle with vertices  $(0, 0), (0, h), (h, 0), (h, h)$  lies in  $B(n)$ .

Suppose  $n = 3$ ,  $c = (x, y, z)$ , and we want to prove equality of  $D_{2,3} f$  and  $D_{3,2} f$ . Then we will consider the rectangle with vertices  $(x, y, z), (x, y, z +$

$h), (x, y + h, z), (x, y + h, z + h)$ . Again, we are taking  $n = 2$  and  $\bar{c} = (0, 0)$ , only for the ease of notation as the same proof is applicable for any finite natural number,  $n$  and any vector  $\bar{c} \in \mathbb{R}^n$ .

$$\text{Define } \nabla(h) = f(h, h) - f(h, 0) - f(0, h) + f(0, 0) \quad (13.10)$$

$$\text{Step 3 : } D_{1,2}f = \frac{\nabla(h)}{h^2} = D_{2,1}f$$

$$\text{Define } G(x) = f(x, h) - f(x, 0) \quad (13.11)$$

Then we have,  $\nabla(h) = G(h) - G(0)$  and  $G'(x) = D_1f(x, h) - D_1f(x, 0)$ . By 1-dimensional mean value theorem,

$$\begin{aligned} G(h) - G(0) &= hG'(x_1) \text{ where } x_1 \in (0, h) \\ &= h(D_1f(x_1, h) - D_1f(x_1, 0)) \end{aligned}$$

We have  $D_1f$  is differentiable at  $(0, 0)$ . There exists linear, total derivative function  $(D_1f)'(0, 0)$  where  $(D_1f)'(0, 0)(x, y) = \nabla D_1f(0, 0) \cdot (x, y)$  satisfying first-order Taylor's formula at  $(0, 0)$ .

$$\text{Remember : } f'(\bar{c})(\bar{v}) = \sum_{k=1}^n v_k D_k f(\bar{c}) = \nabla f(\bar{c}) \cdot \bar{v}$$

$$\begin{aligned} D_1f((0, 0) + (x_1, h)) &= D_1f(0, 0) + \nabla D_1f(0, 0) \cdot (x_1, h) + \|(x_1, h)\|E_1(h) \\ \text{where } E_1(h) &\rightarrow 0 \text{ as } h \rightarrow 0 \end{aligned}$$

$$D_1f(x_1, h) = D_1f(0, 0) + x_1 D_{1,1}f(0, 0) + h D_{2,1}f(0, 0) + \left| \sqrt{x_1^2 + h^2} \right| E_1(h)$$

Similarly,

$$\begin{aligned} D_1f((0, 0) + (x_1, 0)) &= D_1f(0, 0) + \nabla D_1f(0, 0) \cdot (x_1, 0) + \|(x_1, 0)\|E_2(h) \\ \text{where } E_2(h) &\rightarrow 0 \text{ as } h \rightarrow 0 \\ D_1f(x_1, 0) &= D_1f(0, 0) + x_1 D_{1,1}f(0, 0) + |x_1|E_2(h) \end{aligned}$$

Therefore  $\nabla(h) = h(D_1f(x_1, h) - D_1f(x_1, 0)) = h^2 D_{2,1}f(0, 0) + E(h)$  where  $E(h) = h|\sqrt{x_1^2 + h^2}|E_1(h) - h|x_1|E_2(h)$  and  $E(h) \rightarrow 0$  as  $h \rightarrow 0$ .

Since  $0 < x_1 < h$ , we have

$$0 \leq E(h) \leq h^2 \left( \sqrt{2}E_1(h) - E_2(h) \right)$$

Therefore,

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{\nabla(h)}{h^2} &\leq \lim_{h \rightarrow 0} \frac{h^2 D_{2,1}f(0, 0)}{h^2} + \lim_{h \rightarrow 0} \frac{h^2(\sqrt{2}E_1(h) - E_2(h))}{h^2} = D_{2,1}f(0, 0) \\ \lim_{h \rightarrow 0} \frac{\nabla(h)}{h^2} &= D_{2,1}f(0, 0) \end{aligned} \quad (13.12)$$

You may skip the following part. And conclude with the last two lines.

Similarly, define  $H(y) = f(h, y) - f(0, y)$ . Then we have,  $\nabla(h) = H(h) - H(0)$  and  $H'(y) = D_2f(h, y) - D_2f(0, y)$ . By 1-dimensional mean value theorem,

$$\begin{aligned} H(h) - H(0) &= hH'(y_1) \text{ where } y_1 \in (0, h) \\ &= h(D_2f(h, y_1) - D_2f(0, y_1)) \end{aligned}$$

We have  $D_2f$  is differentiable at  $(0, 0)$ . Thus there exists a linear, total derivative function  $(D_2f)'(0, 0)$  where  $(D_2f)'(0, 0)(x, y) = \nabla D_2f(0, 0) \cdot (x, y)$  satisfying first-order Taylor's formula at  $(0, 0)$ . That is,

$$\begin{aligned} D_2f((0, 0) + (h, y_1)) &= D_2f(0, 0) + \nabla D_2f(0, 0) \cdot (h, y_1) + \|(h, y_1)\|E_3(h) \\ &\text{where } E_3(h) \rightarrow 0 \text{ as } h \rightarrow 0 \end{aligned}$$

$$D_2f(h, y_1) = D_2f(0, 0) + hD_{1,2}f(0, 0) + y_1D_{2,2}f(0, 0) + \left| \sqrt{h^2 + y_1^2} \right| E_3(h)$$

Again,

$$\begin{aligned} D_2f((0, 0) + (0, y_1)) &= D_2f(0, 0) + \nabla D_2f(0, 0) \cdot (0, y_1) + |y_1|E_4(h) \\ &\text{where } E_4(h) \rightarrow 0 \text{ as } h \rightarrow 0 \end{aligned}$$

$$D_2f(0, y_1) = D_2f(0, 0) + y_1D_{2,2}f(0, 0) + |y_1|E_4(h)$$

Therefore,  $\nabla(h) = h(D_2f(h, y_1) - D_2f(0, y_1)) = h^2D_{1,2}f(0, 0) + E'(h)$  where  $E'(h) = \left| \sqrt{h^2 + y_1^2} \right| E_3(h) - |y_1|E_4(h)$  and  $E'(h) \rightarrow 0$  as  $h \rightarrow 0$ . And

$$\lim_{h \rightarrow 0} \frac{\nabla(h)}{h^2} = D_{1,2}f(0, 0) \quad (13.13)$$

Therefore,  $D_{1,2}f(0, 0) = D_{2,1}f(0, 0)$ . □

### Continuity

**Theorem 13.3.3.** Suppose  $D_rf$  and  $D_kf$  exists in an  $n$ -ball about  $\bar{c}$ . And  $D_{r,k}f$  and  $D_{k,r}f$  are continuous at  $\bar{c}$ . Then  $D_{r,k}f = D_{k,r}f$ .

*Proof.* We have  $D_rf = (D_rf_1, D_rf_2, \dots, D_rf_m)$ . Therefore, it is sufficient to prove the theorem for real-valued functions. Suppose  $n = 2$ ,  $\bar{c} = (0, 0)$  and the partial derivatives  $D_1f$  and  $D_2f$  exist and are continuous in some 2-ball about  $(0, 0)$ . Suppose  $(h, h)$  lies in that 2-ball, then  $D_1f(h, h) \rightarrow D_1f(0, 0)$  as  $h \rightarrow 0$ . □

*Remark.* A function  $f$  such that  $D_{1,2}f \neq D_{2,1}f$ .

$$\text{Let, } f(x, y) = \begin{cases} \frac{xy(x^2 - y^2)}{x^2 + y^2} & (x, y) \neq (0, 0) \\ 0 & (x, y) = (0, 0) \end{cases}$$

$$\begin{aligned} D_1f(x, y) &= \frac{\partial}{\partial x} \frac{x^3y - xy^3}{x^2 + y^2} \\ &= \frac{(3x^2y - y^3)(x^2 + y^2) - 2x(x^3y - xy^3)}{(x^2 + y^2)^2} \\ &= \frac{x^4y + 4x^2y^3 - y^5}{(x^2 + y^2)^2} \end{aligned}$$

$$D_1f_{(x=0)} = -y \implies D_{2,1}f_{(x=0)} = \frac{\partial}{\partial y} -y = -1$$

$$\begin{aligned}
D_2 f(x, y) &= \frac{\partial}{\partial y} \frac{x^3 y - xy^3}{x^2 + y^2} \\
&= \frac{(x^3 - 3xy^2)(x^2 + y^2) - 2y(x^3 y - xy^3)}{(x^2 + y^2)^2} \\
&= \frac{x^5 - 4x^3 y^2 - xy^4}{(x^2 + y^2)^2} \\
D_2 f_{(y=0)} &= x \implies D_{1,2} f_{(y=0)} = \frac{\partial}{\partial x} x = 1
\end{aligned}$$

Therefore,  $D_{1,2}f \neq D_{2,1}f$  in the neighbourhood of  $(0, 0)$ . This treatment save a lot of time. After  $D_1 f$ , we are planning to perform  $D_{2,1}f = D_2(D_1 f)$  in which the value of  $x$  is going to be treated as a constant. Therefore, we can simplify the expression by substituting  $x = 0$  at this stage. If you are not confident enough to substitute that “early”. You may take partial derivative with respect to  $y$  and then substitute  $x = 0$  and  $y = 0$ . Why don't we substitute  $y = 0$  before  $D_2 f$  is something you should know already !

### 13.3.3 Implicit Functions and Extremum Problems

**Definitions 13.3.1** (Implicit function). Let  $f$  be a function. Consider the equation,  $f(\bar{x}, \bar{y}) = \bar{0}$ . If there exists a function  $g$  such that  $\bar{x} = g(\bar{y})$ , then  $g$  is an **implicit form** of  $f$  or  $g$  is defined implicitly by  $f$ . For example, a linear system of equations  $Ax - b = 0$  implicitly defines  $x = A^{-1}b$  provided  $A$  has non-zero determinant.

**Definitions 13.3.2** (Jacobian Determinant). Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , then determinant of the Jacobian matrix  $Df(\bar{x})$  is the **Jacobian determinant** of  $f$ ,  $J_f(\bar{x})$ .

#### What is an implicit function ?

Consider the equation  $\sin(x + y) = \cos(u + v)$ . We can rewrite this equation as  $\sin(x + y) - \cos(u + v) = 0$ . Now, the equation is of the form  $f(x, y, u, v) = 0$ . Therefore, we have the function  $f : \mathbb{R}^4 \rightarrow \mathbb{R}$  defined by  $f(x, y, u, v) = \sin(x + y) - \cos(u + v)$ .

Let us play around,

$$\begin{aligned}
\sin(x + y) = \cos(u + v) &\implies x + y = \sin^{-1} \cos(u + v) \\
&\implies (x, y) = (\sin^{-1} \cos(u + v) - k, k) \\
&\implies (x, y) = g(u, v)
\end{aligned}$$

where  $g_1(u, v) = \sin^{-1} \cos(u + v) - k$ , and  $g_2(u, v) = k$

Now we have defined a new function  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  from the function  $f$ . Therefore, we could say that  $f$  defines  $g$  implicitly. Now we have a few questions to ask about the nature of such implicit functions (or implicitly defined functions),

1. Does there always exists a function for any such combination ?  
That is, does there exists a function  $h$  such that  $(x, u) = h(y, v)$  ?

2. Suppose  $f(x, y, u, v) = 0$ . And function  $h : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is implicitly defined by  $f$ . Does there always exists a neighbourhood of  $(u, v)$  in which  $h$  has continuous partial derivatives ?
3. What are the properties of  $f$  for these implicit functions to have nice properties ?

It turns out that non-zero Jacobian determinant is nice property  $f$  can have.

**Theorem 13.3.4.** *Let  $f : \mathbb{C} \rightarrow \mathbb{C}$ . Then  $J_f(z) = |f'(z)|^2$ .*

*Proof.* Suppose  $f : \mathbb{C} \rightarrow \mathbb{C}$  where  $f(z) = u(z) + iv(z)$  where  $u : \mathbb{C} \rightarrow \mathbb{R}$  and  $v : \mathbb{C} \rightarrow \mathbb{R}$ . These real-valued functions  $u, v$  have respective  $u^*, v^*$  multivariate real functions such that  $u^* : \mathbb{R}^2 \rightarrow \mathbb{R}$ , where  $u(z) = u^*(x, y)$  and  $z = x + iy$ . Let  $f(z) = z^2 + 1$ . Then  $u^*(x, y) = x^2 - y^2 + 1$  and  $v^*(x, y) = -2xy$ . And theoretically we use derivatives of  $u^*$  when we mention derivatives of  $u$ .

Then  $f$  has a derivative at  $z$  only if the partial derivatives  $D_1u, D_2u, D_1v, D_2v$  exists at  $z$  and satisfies Cauchy-Riemann equations. ie  $D_1u(z) = D_2v(z)$  and  $D_1v(z) = -D_2u(z)$ . [Apostol, 1973, Theorem 5.22].

Thus we have  $f'(z) = D_1u + iD_1v$  [Apostol, 1973, Theorem 12.6]<sup>1</sup>.

$$\begin{aligned} f'(z) &= D_1u(z) + iD_1v(z) \\ |f'(z)|^2 &= (D_1u(z))^2 + (D_1v(z))^2 \end{aligned}$$

For ease of representation, we write  $D_1u$  instead of  $D_1u(z)$

$$|f'(z)|^2 = (D_1u)^2 + (D_1v)^2$$

We also have

$$J_f(z) = |Df(z)| = \begin{vmatrix} D_1u & D_2u \\ D_1v & D_2v \end{vmatrix} = D_1uD_2v - D_1vD_2u = (D_1u)^2 + (D_1v)^2$$

Therefore,  $J_f(z) = |f'(z)|^2$ . □

### Functions with non-zero Jacobian determinant

That is,  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  such that  $J_f \neq 0$  in an  $n$ -ball. In other words, we have an  $n$ -ball  $B(\bar{x})$  such that  $J_f(\bar{y}) \neq 0, \forall \bar{y} \in B(\bar{x})$ .

**Theorem 13.3.5.** *Let  $B$  be an  $n$ -ball about  $\bar{a}$  in  $\mathbb{R}^n$ ,  $\partial B$  be its boundary and  $\bar{B} = B \cup \partial B$  be its closure.<sup>2</sup> Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be continuous in  $\bar{B}$  and all partial derivatives,  $D_j f_i(\bar{x})$  exists for every  $\bar{x} \in B$ . Let  $f(\bar{x}) \neq f(\bar{a})$  for every  $\bar{x} \in \partial B$  and  $J_f(\bar{x}) \neq 0$  for every  $\bar{x} \in B$ . Then  $f(B)$  contains an  $n$ -ball about  $f(\bar{a})$ .*

$$\begin{aligned} B &= \{\bar{x} : \|\bar{x} - \bar{a}\| < r\} \\ \partial B &= \{\bar{x} : \|\bar{x} - \bar{a}\| = r\} \\ \bar{B} &= \{\bar{x} : \|\bar{x} - \bar{a}\| \leq r\} \end{aligned}$$

<sup>1</sup>Prove using first-order Taylor's formula

<sup>2</sup> $\bar{B}$  : The line above  $B$  has a different meaning compare to  $\bar{a}$  (situations like this are an abuse of language).



*Proof.* Define  $g : \partial B \rightarrow \mathbb{R}$  where  $g(\bar{x}) = \|f(\bar{x}) - f(\bar{a})\|$ . We have,  $f(\bar{x}) \neq f(\bar{a})$  for every  $\bar{x} \in \partial B$ , thus  $g(\bar{x}) > 0$  for every  $\bar{x} \in \partial B$ . Function  $f$  is continuous on  $\bar{B}$ , thus  $g$  is continuous on  $\bar{B}$  and thus  $g$  is continuous on its subset  $\partial B$ . Since  $\partial B$  is compact, every continuous function on  $\partial B$  attains its extrema<sup>3</sup> and thus  $g$  attains its minimum value  $m > 0$  somewhere on  $\partial B$ .

Consider  $n$ -ball  $T$  about  $f(\bar{a})$  with radius  $\frac{m}{2}$ ,

$$T = B\left(f(\bar{a}), \frac{m}{2}\right) = \left\{ \bar{y} \in \mathbb{R}^n : \|f(\bar{a}) - \bar{y}\| < \frac{m}{2} \right\}$$

Therefore, it is sufficient to prove that  $T \subset f(B)$ .

Let  $\bar{y} \in T$ . Define  $h : \bar{B} \rightarrow \mathbb{R}$  where  $h(\bar{x}) = \|f(\bar{x}) - \bar{y}\|$ . Again this continuous function  $h$  on compact set  $\bar{B}$  attains its extrema somewhere on  $\bar{B}$ . Since  $\bar{y} \in T$ ,  $h(\bar{a}) = \|f(\bar{a}) - \bar{y}\| < \frac{m}{2}$ . Thus, the minimum of  $h$  on  $\bar{B}$  is less than  $\frac{m}{2}$ , since  $\bar{a} \in \bar{B}$ .

Let  $\bar{x} \in \partial B$ , then

$$\begin{aligned} h(\bar{x}) &= \|f(\bar{x}) - \bar{y}\| \\ &= \|f(\bar{x}) - f(\bar{a}) + f(\bar{a}) - \bar{y}\| \\ &\geq \|f(\bar{x}) - f(\bar{a})\| + \|f(\bar{a}) - \bar{y}\| \\ &= g(\bar{x}) - h(\bar{a}) \\ &> \frac{m}{2} \text{ since } g(\bar{x}) \geq m \text{ and } h(\bar{a}) < \frac{m}{2} \end{aligned}$$

Thus  $h$  doesn't attain its minimum on  $\partial B$ , but at an interior point  $\bar{c} \in B$ . Consider

$$h^2(\bar{x}) = \|f(\bar{x}) - \bar{y}\|^2 = \sum_{r=1}^n (f_r(\bar{x}) - y_r)^2$$

The function  $h^2$  also has minimum at the same point  $\bar{c}$ . Thus all partial derivatives of  $h^2$  at  $\bar{c}$  are zero. ie,

$$D_k h^2(\bar{c}) = \sum_{r=1}^n (f_r(\bar{c}) - y_r) D_k f_r(\bar{c}) = 0$$

This is a system of linear equations with non-zero determinant since  $\bar{c} \in B$  and we have  $J_f(\bar{c}) \neq 0$ . Therefore,  $f_r(\bar{c}) = y_r$ . That is,  $f(\bar{c}) = \bar{y} \in f(B)$ . Since  $\bar{y} \in T$  is arbitrary,  $T \subset f(B)$ .  $\square$

**Theorem 13.3.6.** Let  $A$  be an open subset of  $\mathbb{R}^n$  and  $f : A \rightarrow \mathbb{R}^n$  is continuous and has continuous partial derivatives  $D_j f_i$  on  $A$ . If  $f$  is one-to-one on  $A$  and  $J_f(\bar{x}) \neq 0$ ,  $\forall \bar{x} \in A$ , then  $f(A)$  is open.

*Proof.* Let  $\bar{b} \in f(A)$ . Then  $\bar{b} = f(\bar{a})$  for some  $\bar{a} \in A$ . We have,  $f$  is continuous,  $f$  has continuous partial derivatives on  $A$  and  $J_f(\bar{x}) \neq 0$  for every  $\bar{x} \in A$ . Therefore, there exists an open ball  $B \subset A$  containing  $\bar{a}$  such that  $f(B) \subset f(A)$  contains an  $n$ -ball about  $f(\bar{a})$ . Since  $\bar{b} \in f(A)$  is arbitrary, every point in  $f(A)$  has an  $n$ -ball containing it in  $f(A)$ . Therefore,  $f(A)$  is open.

Two assumption in above theorem are trivial. 1.  $f$  is continuous in the closed ball,  $\bar{B}$ . Set  $B$  so chosen that  $\bar{B} \subset A$  and  $f$  is continuous in  $A$ .

<sup>3</sup>“Every continuous function on a compact set attains its extrema”

Thus,  $f$  is continuous in  $\bar{B}$ . 2.  $f$  has different value at boundary compared to center. ie,  $f(\bar{a}) \neq f(\bar{x})$ ,  $\forall x \in \partial B$ . We have,  $f$  is injective on  $A$ , and  $\bar{B} \subset A$ . Thus  $f$  has different values for any two distinct points in it. Thus,  $\forall \bar{x}, \bar{y} \in \bar{B}$ ,  $\bar{x} \neq \bar{y} \implies \bar{x}, \bar{y} \in A$ , and  $\bar{x} \neq \bar{y} \implies f(\bar{x}) \neq f(\bar{y})$   $\square$

**Theorem 13.3.7.** Let  $S$  be an open subset of  $\mathbb{R}^n$  and  $f : S \rightarrow \mathbb{R}^n$ . Let components of  $f$  has continuous partial derivatives on  $S$ ,  $D_j f_i$  and  $J_f(\bar{a}) \neq 0$  for some point  $\bar{a} \in S$ . Then there is an  $n$ -ball  $B$  about  $\bar{a}$  on which  $f$  is injective.

*Proof.* Let  $\bar{z} = (\bar{z}_1, \bar{z}_2, \dots, \bar{z}_n)$  where  $\bar{z}_i \in \mathbb{R}^n$ . ie,  $\bar{z} \in \mathbb{R}^{n^2}$ . Define function  $h : \mathbb{R}^{n^2} \rightarrow \mathbb{R}$  by  $h(\bar{z}) = \det [D_j f_i(\bar{z}_i)]$ . Since  $f$  has continuous partial derivatives on  $S$ , each component of  $f$  has continuous partial derivatives in  $S$  and thus  $h$  is continuous on  $S^n$  which is a subset of  $\mathbb{R}^{n^2}$  since  $S$  is an open subset of  $\mathbb{R}^n$ .

Let  $\bar{a} \in S$  such that  $J_f(\bar{a}) \neq 0$ . Existence of such a point in  $S$  is assumed.

Consider,  $\bar{z}_i = \bar{a}$ ,  $\forall i$ . Then  $\bar{z} = (\bar{a}, \bar{a}, \dots, \bar{a})$ . And  $h(\bar{z}) = \det [D_j f_i(\bar{a})] = J_f(\bar{a}) \neq 0$ .

Since  $h$  is continuous and  $h(\bar{z}) \neq 0$ . There exists an  $n$ -ball  $B$  about  $\bar{a}$  in  $S$  such that  $h(\bar{z}) \neq 0$  for  $\bar{z}_i \in B$ . We claim that  $f$  is injective on  $B$ .

Suppose  $f$  is not injective. ie, There exists  $\bar{x}, \bar{y} \in B(\bar{a})$  such that  $\bar{x} \neq \bar{y}$  and  $f(\bar{x}) = f(\bar{y})$ . Open ball  $B(\bar{a})$  is a convex set. And the line segment  $L(\bar{x}, \bar{y}) \subset B(\bar{a})$ . The function  $f$  is differentiable on  $S$ . On applying mean-value theorem to each component of  $f$ , we get

$$0 = f_i(\bar{y}) - f_i(\bar{x}) = \nabla f_i(\bar{Z}_i) \cdot (\bar{y} - \bar{x}), \quad i = 1, 2, \dots$$

where  $\bar{Z}_i \in L(\bar{x}, \bar{y}) \subset B(\bar{a})$ . Therefore, We have

$$\sum_{k=1}^n D_k f_i(\bar{Z}_i)(y_k - x_k) = 0$$

The determinant of this system of linear equations is nonzero, as the function  $f$  has nonzero jacobian determinant at  $\bar{Z}_i \in B(\bar{a})$  for  $i = 1, 2, \dots$ . Thus,  $y_i = x_i$  for  $i = 1, 2, \dots$ . This contradicts  $\bar{x} \neq \bar{y}$ . Hence, the function  $f$  is injective.  $\square$

**Theorem 13.3.8.** Let  $A$  be an open subset of  $\mathbb{R}^n$  and assume that  $f : A \rightarrow \mathbb{R}^n$  has continuous partial derivatives  $D_j f_i$  on  $A$ . If  $J_f(\bar{x}) \neq 0$  for all  $\bar{x} \in A$ , then  $f$  is an open mapping.

*Proof.* Let  $S$  be an open subset of  $A$ . Let  $\bar{x} \in S$ . Clearly,  $f$  has continuous partial derivatives on  $S$  and  $J_f(\bar{x}) \neq 0$  for all  $\bar{x} \in S$ . Thus, there is an  $n$ -ball  $B(\bar{x})$  in which  $f$  is injective. Therefore,  $f(B(\bar{x}))$  is open in  $\mathbb{R}^n$ . Since  $\bar{x} \in S$  is arbitrary,  $S = \cup_{\bar{x} \in S} B(\bar{x})$ . And  $f(S) = \cup_{\bar{x} \in S} f(B(\bar{x}))$ . Therefore,  $f(S)$  is open. Since open set  $S$  is arbitrary,  $f$  is an open mapping.  $\square$

**Remark (Properties).** Functions with non-zero Jacobian determinant has following properties :

1. If  $J_f \neq 0$  in  $n$ -ball  $B$  about  $\bar{a}$  which has different values at its boundaries, then  $f(B)$  has an  $n$ -ball about  $f(\bar{a})$ .
2. If  $J_f \neq 0$ ,  $f$  has continuous partial derivatives in  $S$ , and  $f$  is injective in an open set  $A$ , then  $f(A)$  is open.

3. Let  $S$  be an open set in  $\mathbb{R}^n$ ,  $f$  has continuous partial derivatives in  $S$ , and  $J_f(\bar{a}) \neq 0$  for some  $\bar{a} \in S$ , then  $f$  is injective on an  $n$ -ball  $B(\bar{a})$  in  $S$ .
4. Let  $A$  be an open set in  $\mathbb{R}^n$ ,  $f$  has continuous partial derivatives in  $A$ , and  $J_f \neq 0$  in  $A$ , then  $f$  is an open mapping.

### Inverse function Theorem

**Theorem 13.3.9** (Inverse function). *Let  $S$  be an open subset of  $\mathbb{R}^n$  and  $f$  be a continuously differentiable function<sup>4</sup>  $f : S \rightarrow \mathbb{R}^n$ . If  $J_f(\bar{a}) \neq 0$  for some  $\bar{a} \in S$ , then there are two open sets  $X \subset S$ , and  $Y \subset f(S)$  such that*

1.  $\bar{a} \in X$  and  $f(\bar{a}) \in Y$
2.  $Y = f(X)$
3.  $f$  is injective
4. there exists another function  $g : Y \rightarrow X$  such that  $g(f(\bar{x})) = \bar{x}$ ,  $\forall \bar{x} \in X$
5.  $g$  is continuously differentiable on  $Y$

*In other words, if  $f \in C'$  and there exists  $\bar{a} \in S$  such that  $J_f(\bar{a}) \neq 0$ , then  $f$  has an inverse  $f^{-1}$  in a neighbourhood of  $f(\bar{a})$  and  $f^{-1} \in C'$ .*

**Proof. Step 1 : Construction of open sets  $X$  and  $Y$ .**

Given that,  $J_f(\bar{a}) \neq 0$  and  $f \in C'$ . Thus all partial derivatives of  $f$  are continuous on  $S$ . Then  $J_f$  is continuous on  $S$ . By the continuity of  $J_f$  at  $\bar{a}$ , there exists a neighbourhood of  $\bar{a}$ , say  $B_1(\bar{a})$  in which  $J_f \neq 0$ . That is,  $\forall \bar{x} \in B_1(\bar{a})$ ,  $J_f(\bar{x}) \neq 0$ . Therefore, (by theorem) there exists an  $n$ -ball  $B(\bar{a})$  on which  $f$  is injective. Let  $B$  be an  $n$ -ball with center  $\bar{a}$  contained in  $B(\bar{a})$ . Then  $f$  is injective on  $B$ . Therefore, (by theorem)  $f(B)$  contains an  $n$ -ball with center  $f(\bar{a})$ . Let  $Y$  be the  $n$ -ball contained in  $f(B)$ . And  $X = f^{-1}(Y) \cap B$ . That is, the inverse image of  $Y$  on  $B$ . Since  $f$  is continuous,  $f^{-1}(Y)$  is open. Thus,  $X$  is an intersection of open sets. And therefore,  $X$  is open.

**Step 2 : The inverse of  $f$ , say  $g$ .**

Clearly  $\bar{a} \in X$  and  $f(\bar{a}) \in Y$ . Also  $Y = f(X)$  and  $f$  is injective on  $X$  (since,  $X \subset B$ ).

The closure of  $B$ ,  $\bar{B}$  is compact and  $f$  is injective and continuous on  $\bar{B}$ . Then<sup>5</sup> there exists a continuous function  $g$  defined on  $f(\bar{B})$  such that  $g \circ f$  is the identity function on  $\bar{B}$ . That is,  $\forall x \in \bar{B}$ ,  $g(f(\bar{x})) = \bar{x}$ . Thus,  $g(X) = Y$  and  $g$  is unique.

**Step 3 :  $g$  has continuous partial derivatives.**

Define a real-valued function  $h : S^n \rightarrow \mathbb{R}$  by  $h(\bar{Z}) = \det[D_j f_i(\bar{Z}_i)]$  where  $\bar{Z}_1, \bar{Z}_2, \dots, \bar{Z}_n \in S$  and  $\bar{Z} = (\bar{Z}_1, \bar{Z}_2, \dots, \bar{Z}_n)$ . Now, let  $\bar{Z} = (\bar{a}, \bar{a}, \dots, \bar{a})$ . Then  $h(\bar{Z}) \neq 0$  and  $h$  is continuous on  $S^n$ . Therefore,  $\bar{Z}$  has a neighbourhood on which  $h$  does not vanish (that is, nonzero). Let  $B_2(\bar{a})$  be the corresponding  $n$ -ball with center  $\bar{a}$  such that  $\bar{Z}_i \in B_2(\bar{a}) \implies h(\bar{Z}) \neq 0$ .

Let  $B$  be an  $n$ -ball with center  $\bar{a}$  contained in  $B_2(\bar{a})$ . Now  $\bar{B} \subset B_2(\bar{a})$ . And  $h(\bar{Z}) \neq 0$ ,  $\forall \bar{Z}_i \in \bar{B}$ .

<sup>4</sup> $f \in C'(S) : f$  is continuously differentiable on  $S$

<sup>5</sup>Existence of inverse of a continuous function on a compact set in metric spaces.

We have,  $g = (g_1, g_2, \dots, g_n)$ . It is enough to prove that  $g_k \in C'$  for  $k = 1, 2, \dots, n$ . Again, it is enough to prove that  $D_r g_k$  exists and is continuous for  $1 \leq r \leq n$ . (Fix some  $r$  and prove that  $D_r g_k$  is continuous.)

Let  $\bar{y} \in Y$ . Define  $\bar{x} = g(\bar{y})$  and  $\bar{x}' = g(\bar{y} + t\bar{u}_r)$  where  $t$  is sufficiently small such that  $\bar{y} + t\bar{u}_r \in Y$ . Then  $\bar{x}, \bar{x}' \in X$ . And  $f(\bar{x}') - f(\bar{x}) = t\bar{u}_r$ . Therefore  $f_i(\bar{x}) - f_i(\bar{x}') = 0$  when  $i \neq r$ . And  $f_i(\bar{x}') - f_i(\bar{x}) = t$  when  $i = r$ . By mean-value theorem,

$$\frac{f_i(\bar{x}') - f_i(\bar{x})}{t} = \nabla f_i(\bar{Z}_i) \cdot \frac{\bar{x}' - \bar{x}}{t}$$

where  $\bar{Z}_i \in L(\bar{x}, \bar{x}')$ , the line segment joining  $\bar{x}$  and  $\bar{x}'$ . Since  $\det[D_j f_i(\bar{Z}_i)] = h(\bar{Z}) \neq 0$ , this system of linear equations in  $n$  unknowns,  $\frac{x'_j - x_j}{t}$  has a unique solution. As  $t \rightarrow 0$ ,  $\bar{x}' \rightarrow \bar{x}$ . And  $\bar{Z}_i \rightarrow \bar{x}$ . Since  $J_f(\bar{x}) \neq 0$ , the limit

$$\lim_{t \rightarrow 0} \frac{g_k(\bar{y} + t\bar{u}_r) - g_k(\bar{y})}{t}$$

exists. Thus,  $D_r g_k(\bar{y})$  exists  $\forall y \in Y$  and every  $r$ . By Cramer's rule, this limit is a quotient of two determinants of partial derivatives of  $f$ , which are all continuous since  $f \in C'$ . Therefore,  $D_r g_k$  are all continuous and  $g \in C'$ .  $\square$

### Implicit function Theorem

**Theorem 13.3.10.** Let  $S$  be an open set in  $\mathbb{R}^{n+k}$  and  $f$  be a function  $f : S \rightarrow \mathbb{R}^n$ . Suppose  $f$  is continuously differentiable on  $S$ . Let  $(\bar{x}_0, \bar{t}_0) \in S$  such that  $\bar{x}_0 \in \mathbb{R}^n$ ,  $\bar{t}_0 \in \mathbb{R}^k$ ,  $f(\bar{x}_0, \bar{t}_0) = \bar{0}$  and  $J_f(\bar{x}_0, \bar{t}_0) \neq 0$ . Then there exists an open set  $T_0$  containing  $\bar{t}_0$  in  $\mathbb{R}^k$  and a unique function  $g : T_0 \rightarrow \mathbb{R}^n$  such that

1.  $g$  is continuously differentiable on  $T_0$
2.  $g(\bar{t}_0) = \bar{x}_0$
3.  $f(g(\bar{t}), \bar{t}) = \bar{0}$ ,  $\forall \bar{t} \in T_0$

*Proof.* Given  $f : S \rightarrow \mathbb{R}^n$  where  $S \subset \mathbb{R}^{n+k}$ . We have  $f = (f_1, f_2, \dots, f_m)$ . Also given that  $f \in C'(S)$ ,  $f(\bar{x}_0, \bar{t}_0) = 0$ , and  $J_f(\bar{x}_0, \bar{t}_0) \neq 0$ . Define a function  $F : S \rightarrow \mathbb{R}^{n+k}$  defined by  $F = (F_1, F_2, \dots, F_{n+k})$ .

$$F_m(\bar{x}; \bar{t}) = \begin{cases} f_m(\bar{x}; \bar{t}) & 1 \leq m \leq n \\ t_{m-n} & n < m \leq n+k \end{cases}$$

For example, let  $n = 3$ ,  $k = 2$ . Let  $\bar{x} = (1, 2, 3)$  and  $\bar{t} = (4, 5)$ . Suppose  $f_k(1, 2, 3, 4, 5) = a_k$ . Then  $F(\bar{x}; \bar{t}) = (a_1, a_2, a_3, 4, 5)$ .

Then the Jacobian determinant of  $F'(\bar{x}; \bar{t})$  is given by

$$J_F(\bar{x}; \bar{t}) = \begin{vmatrix} D_1 f_1 & D_2 f_1 & \cdots & D_m f_1 & & & \\ D_1 f_2 & D_2 f_2 & \cdots & D_m f_2 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & & & \\ D_1 f_m & D_2 f_m & \cdots & D_m f_m & & & \\ & & & 0 & 1 & \cdots & 0 \\ & & & \vdots & \vdots & \ddots & \vdots \\ & & & 0 & 0 & \cdots & 1 \end{vmatrix}$$

Thus  $J_F(\bar{x}_0; \bar{t}_0) = J_f(\bar{x}_0; \bar{t}_0) \neq 0$ . Also,  $F(\bar{x}_0; \bar{t}_0) = (\bar{0}; \bar{t}_0)$ . Therefore, by inverse function theorem, there exists open sets  $X, Y$  containing  $(\bar{x}_0; \bar{t}_0)$  and  $(\bar{0}; \bar{t}_0)$  such that  $F$  is injective on  $X$ ,  $X = F^{-1}(Y)$  and there exists a unique local inverse function  $G$  such that  $G(F(\bar{x}; \bar{t})) = (\bar{x}; \bar{t})$  and  $G \in C'(Y)$ .

Let  $G = (v; w)$ . That is,  $v_i = G_i$  for  $1 \leq i \leq n$ . And  $w_j = G_{n+j}$  for  $1 \leq j \leq k$ . We have,  $G(F(\bar{x}; \bar{t})) = (\bar{x}; \bar{t})$ . Therefore,  $v(F(\bar{x}; \bar{t})) = \bar{x}$  and  $w(F(\bar{x}; \bar{t})) = \bar{t}$ . Since  $X \subset F^{-1}(Y)$  and  $F$  is one-to-one on  $X$ , for every  $(\bar{x}; \bar{t}) \in Y$ , there exists  $(\bar{x}'; \bar{t}') \in X$  such that  $F(\bar{x}'; \bar{t}') = (\bar{x}; \bar{t})$ . By the definition of  $F$ ,  $\bar{t}' = \bar{t}$ . Therefore,  $v(\bar{x}; \bar{t}) = v(F(\bar{x}'; \bar{t})) = \bar{x}'$  and  $w(\bar{x}; \bar{t}) = w(F(\bar{x}'; \bar{t})) = \bar{t}$ . Now, we have  $G : Y \rightarrow X$  defined by  $G(\bar{x}; \bar{t}) = (\bar{x}'; \bar{t})$ .

Let  $T_0$  be a subset of  $\mathbb{R}^k$  defined by  $T_0 = \{\bar{t} \in \mathbb{R}^k : (\bar{0}; \bar{t}) \in Y\}$ . For each  $\bar{t} \in T_0$  let  $g : T_0 \rightarrow \mathbb{R}^n$  is defined by  $g(\bar{t}) = v(\bar{0}; \bar{t})$ . The set  $T_0$  is open. And  $g \in C'(T_0)$  since  $g$  is constructed from the components of  $G$  which has continuous partial derivatives on  $Y$ . (ie,  $G \in C'(Y)$ .)

Clearly,  $g(\bar{t}_0) = v(\bar{0}; \bar{t}_0) = \bar{x}_0$ . And  $(\bar{0}; \bar{t}) = F(\bar{x}_0; \bar{t}_0)$ . Therefore, we have  $f(v(\bar{x}; \bar{t}); \bar{t}) = \bar{x}$ . Let  $\bar{x} = \bar{0}$ , then  $f(g(\bar{t}); \bar{t}) = \bar{0}$ . It is enough to prove that the function  $g$  is unique. Suppose  $f(g(\bar{t}); \bar{t}) = f(h(\bar{t}); \bar{t})$ . Since  $f$  is one-to-one on  $X$ ,  $(g(\bar{t}); \bar{t}) = (h(\bar{t}); \bar{t})$  for every  $\bar{t} \in T_0$ . And  $g(\bar{t}) = h(\bar{t})$ ,  $\forall \bar{t} \in T_0$ .  $\square$

### Extrema of function of one variable

The function  $f : \mathbb{R} \rightarrow \mathbb{R}$  defined by  $f(x) = x^3$  has derivative  $f'(x) = 3x^2$ . Thus  $f'(0) = 0$ . However, 0 is not a local extrema for the function  $f$ . Thus, derivative of the function vanishing at point is not sufficient for a local extrema at that point.

**Theorem 13.3.11** (sufficient condition for local extrema). *Let  $n \geq 1$  and function  $f$  has  $n$ th partial derivative in open interval  $(a, b)$ . Suppose for some  $c \in (a, b)$ ,*

$$f'(c) = f''(c) = \cdots = f^{(n-1)}(c) = 0 \text{ and } f^{(n)}(c) \neq 0$$

*If  $n$  is even,*

1.  $f$  has a local minimum at  $c$  if  $f^{(n)}(c) > 0$
2.  $f$  has a local maximum at  $c$  if  $f^{(n)}(c) < 0$

*and If  $n$  is odd, there is neither a local minimum nor a local maximum at  $c$ .*

*Proof.* We have  $f^{(n)}(c) \neq 0$ . Thus, there exists an open interval  $B(c)$  such that for each  $x \in B(c)$ ,  $f^{(n)}(x)$  has the same sign as  $f^{(n)}(c)$ . By Taylor's theorem, we have

$$f(x) = f(c) + \sum_{k=1}^{n-1} \frac{f^{(k)}(c)}{k!} (x-c)^k + \frac{f^{(n)}(x_1)}{n!} (x-c)^n \quad (13.14)$$

where  $x_1 \in L(x, c)$ , the line connecting  $x$  and  $c$ . We have,  $f^{(k)}(c) = 0$  for  $k = 1, 2, \dots, n-1$ . Thus,

$$f(x) - f(c) = \frac{f^{(n)}(x_1)}{n!} (x-c)^n$$

Case 1 :  $n$  is even.

If  $n$  is even, then  $(x - c)^n > 0$ . Therefore,  $f(x) - f(c)$  has the same sign as  $f^{(n)}(x_1)$ . If  $f^{(n)}(c) > 0$ , then  $f^{(n)}(x_1)$  has a positive value and  $f(x) - f(c) > 0$  for every  $x \in B(c)$ . Therefore,  $f$  has a local minimum at  $c$ . Similarly, if  $f^{(n)}(c) < 0$ , then  $f(x) - f(c) < 0$  for every  $x \in B(c)$ . Therefore,  $f$  has a local maximum at  $c$ .

Case 2 :  $n$  is odd.

If  $n$  is odd, then  $(x - c)^n$  takes both positive and negative values. Therefore,  $f(x) - f(c)$  has both positive and negative values in  $B(c)$ . Thus  $f$  has neither local minimum nor local maximum at  $c$ .  $\square$

### Extrema of function of several variables

**Definitions 13.3.3.** If function  $f$  is differentiable at  $\bar{a}$  and  $\nabla f(\bar{a}) = \bar{0}$ , the point  $\bar{a}$  is a **stationary point** of  $f$ .

**Definitions 13.3.4.** A stationary point is a **saddle point** if every  $n$ -ball  $B(\bar{a})$  contains points  $\bar{x}$  such that  $f(\bar{x}) > f(\bar{a})$  and other points such that  $f(\bar{x}) < f(\bar{a})$ .

**Definitions 13.3.5.** A function  $Q : \mathbb{R}^n \rightarrow \mathbb{R}$  defined by  $Q(\bar{x}) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$  where  $a_{ij} \in \mathbb{R}$  is a function of the **quadratic form** or simply a quadratic form.

**symmetric quadratic form**  $a_{ij} = a_{ji}$ ,  $\forall i, j$ .

**positive definite quadratic form**  $Q(\bar{x}) > 0, \forall \bar{x} \neq \bar{0}$ .

**negative definite quadratic form**  $Q(\bar{x}) < 0, \forall \bar{x} \neq \bar{0}$ .

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . And second derivative of  $f$  exists if the derivative of  $f$ ,  $f'$  is differentiable. Thus,  $f'(\bar{a} + h\bar{t}) = f'(\bar{a}) + hf''(\bar{a})(\bar{t}) + |h||t|E_{\bar{a}}(h\bar{t})$  where  $E_{\bar{a}} \rightarrow \bar{0}$  as  $h \rightarrow 0$ . Writing the Taylor's first order formula using matrices, we can see that  $f''(\bar{a})(\bar{t})$  is of the quadratic form.

**Theorem 13.3.12.** Let  $f$  be a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . Suppose that the second order partial derivatives  $D_{i,j}f$  exists in an  $n$ -ball  $B(\bar{a})$  and are continuous at  $\bar{a}$  where  $\bar{a}$  is a stationary point of  $f$ .

$$\text{Let } Q(\bar{t}) = \frac{1}{2}f''(\bar{a}, \bar{t}) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n D_{i,j}f(\bar{a})t_i t_j$$

1. If  $Q(\bar{t}) > 0$  for all  $\bar{t} \neq \bar{0}$ ,  $f$  has a relative minimum at  $\bar{a}$ .
2. If  $Q(\bar{t}) < 0$  for all  $\bar{t} \neq \bar{0}$ ,  $f$  has a relative maximum at  $\bar{a}$ .
3. If  $Q(\bar{t})$  takes both positive and negative values, then  $f$  has a saddle point at  $\bar{a}$ .

*Proof.* Define  $Q : \mathbb{R}^n \rightarrow \mathbb{R}$  by  $Q(\bar{t}) = \frac{1}{2}f''(\bar{a}; \bar{t})$ . Then  $Q$  is continuous for every  $\bar{t} \in \mathbb{R}^n$ . Let  $S$  be the boundary of the  $n$ -ball  $B(\bar{0}, 1)$ . That is,  $S = \{\bar{t} \in \mathbb{R}^n : \|\bar{t}\| = 1\}$ .

Case 1 : Suppose that  $Q(\bar{t}) > 0, \forall \bar{t} \neq \bar{0}$ .

We have,  $Q$  is a continuous real-valued function on a compact interval,  $S$ . Therefore,  $Q$  attains its extrema. Thus  $Q$  has a minimum value at point in  $S$ , say  $m$ . Clearly  $Q(\bar{t}) > 0 \implies m > 0$ .

By Taylor's formula,  $f(\bar{a} + \bar{t}) - f(\bar{a}) = \nabla f(\bar{a}) \cdot \bar{t} + \frac{1}{2}f''(\bar{z}; \bar{t})$  where  $\bar{z} \in L(\bar{a} + \bar{t}, \bar{a})$

For every  $\bar{a} \in S$ ,  $\nabla f(\bar{a}) = \bar{0}$ . And  $f(\bar{a} + \bar{t}) - f(\bar{a}) \rightarrow \frac{1}{2}f''(\bar{a}; \bar{t})$  as  $\bar{t} \rightarrow \bar{0}$ .

$$\begin{aligned} \text{Thus, } f(\bar{a} + \bar{t}) - f(\bar{a}) &= \frac{1}{2}f''(\bar{a}; \bar{t}) + \|\bar{t}\|^2 E(\bar{t}) \\ &= Q(\bar{t}) + \|\bar{t}\|^2 E(\bar{t}) \text{ where } E(\bar{t}) \rightarrow \bar{0} \text{ as } \bar{t} \rightarrow \bar{0} \end{aligned}$$

Let  $c = \frac{1}{\|\bar{t}\|}$ . Then  $c\bar{t} \in S$  and  $Q(c\bar{t}) = c^2 Q(\bar{t}) \geq m$ . Thus,  $Q(\bar{t}) \geq m\|\bar{t}\|^2$ .  
Therefore,  $f(\bar{a} + \bar{t}) - f(\bar{a}) \geq m\|\bar{t}\|^2 + \|\bar{t}\|^2 E(\bar{t})$

Choose  $n$ -ball  $B(\bar{0}, r)$  such that  $|E(\bar{t})| < \frac{m}{2}$  for every  $\bar{t} \in B(\bar{0}, r)$ . Thus,

$$-\frac{m}{2}\|\bar{t}\|^2 \leq -\|\bar{t}\|^2 |E(\bar{t})| \leq 0$$

Therefore,  $f(\bar{a} + \bar{t}) - f(\bar{a}) \geq m\|\bar{t}\|^2 - \frac{m}{2}\|\bar{t}\|^2 = \frac{m}{2}\|\bar{t}\|^2$  for every  $\bar{t} \in B(\bar{0}, r)$ .  
Clearly,  $f$  has a local minimum at  $\bar{a}$ .

Case 2 : Suppose  $Q(\bar{t}) < 0$ ,  $\forall \bar{t} \neq \bar{0}$ , then consider  $-f$ .

Clearly, function  $-f$  has a local minimum at  $\bar{t}$ . Thus  $f$  has a local maximum at  $\bar{t}$ .

Case 3 : Suppose  $Q(\bar{t})$  takes both positive and negative values.

$$\begin{aligned} \text{We have, } f(\bar{a} + \lambda\bar{t}) - f(\bar{a}) &= Q(\lambda\bar{t}) + \lambda^2\|\bar{t}\|^2 E(\lambda\bar{t}). \\ &= \lambda^2[Q(\bar{t}) + \|\bar{t}\|^2 E(\lambda\bar{t})] \end{aligned}$$

Choose  $n$ -ball  $B(\bar{0}, r)$  such that  $\|\bar{t}\|^2 E(\lambda\bar{t}) < \frac{1}{2}|Q(\bar{t})|$ ,  $\forall \bar{t} \in B(\bar{0}, r)$ . We have  $\bar{t} \in B(\bar{0}, r) \implies \lambda < r$ . Then, error function  $E(\bar{t})$  on the RHS is small enough, not to affect the sign of the RHS. Thus  $f(\bar{a} + \lambda\bar{t}) - f(\bar{a})$  has the same sign as  $Q(\bar{t})$ . Therefore,  $\bar{a}$  is a saddle point.  $\square$

**Theorem 13.3.13.** Let  $f$  be a real-valued function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  with continuous second order partial derivatives at a stationary point  $\bar{a} \in \mathbb{R}^2$ . Let  $A = D_{1,1}f(\bar{a})$ ,  $B = D_{1,2}f(\bar{a}) = D_{2,1}f(\bar{a})$ , and  $C = D_{2,2}f(\bar{a})$ . And let  $\Delta = \begin{vmatrix} A & B \\ B & C \end{vmatrix} = AC - B^2$ . Then we have,

1. If  $\Delta > 0$  and  $A > 0$ , then  $f$  has a relative minimum at  $\bar{a}$ .
2. If  $\Delta > 0$  and  $A < 0$ , then  $f$  has a relative maximum at  $\bar{a}$ .
3. if  $\Delta < 0$ , then  $f$  has a saddle point at  $\bar{a}$ .

*Proof.* We have, function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  with  $\nabla f(\bar{a}) = (0, 0)$ . Consider the quadratic form  $Q(x, y) = \frac{1}{2}[Ax^2 + Bxy + Cy^2]$  where  $A = D_{1,1}f(\bar{a})$ ,  $B =$

$D_{1,2}f(\bar{a})$ , and  $C = D_{2,2}f(\bar{a})$ . Therefore,  $Q(x, y) = \frac{1}{2}f''(\bar{a}; \bar{t})$ .

Case 1 : Suppose  $A \neq 0$ .

$$\begin{aligned} Q(x, y) &= \frac{1}{2A}[A^2x^2 + ABxy + ACy^2] \\ &= \frac{1}{2A}[(Ax + By)^2 - B^2y^2 + ACy^2] \\ &= \frac{1}{2A}[(Ax + By)^2 + \Delta y^2] \end{aligned}$$

If  $\Delta > 0$ , then  $Q(x, y)$  has the same sign as  $A$ . Therefore,  $f$  has a local minimum/maximum at  $\bar{a}$  depending on the sign of  $A$ .

Case 2 : Suppose  $A = 0$ . Then  $Q(x, y) = \frac{1}{2}[Bxy + Cy^2] = \frac{1}{2}(Bx + Cy)y$ .

Now we have two lines in  $\mathbb{R}^2$ ,  $y = 0$  and  $Bx + Cy = 0$ . These two lines divide  $\mathbb{R}^2$  into four regions. The value of  $Q(x, y)$  is positive in two of those regions and negative in the other two regions. Therefore,  $f(\bar{a} + \bar{t}) - f(\bar{a})$  assumes both positive and negative values in any neighbourhood  $B(\bar{a}, r)$ . Therefore,  $\bar{a}$  is a saddle point.  $\square$

## 13.4 Integration on Differential Forms

**Definitions 13.4.1.** A  $k$ -cell in  $\mathbb{R}^k$  is given by,

$$I^k = \{\bar{x} \in \mathbb{R}^k : a_i \leq x_i \leq b_i, \forall i\}$$

where  $\bar{a}, \bar{b} \in \mathbb{R}^k$ . Let  $f$  be a continuous, real-valued function on  $I^k$ . Then, the **integral** of  $f$  over  $I^k$  is given by,

$$\begin{aligned} \int_{I^k} f(\bar{x}) d\bar{x} &= f_0 \text{ where } f_k = f \text{ and} \\ f_{r-1} &= \int_{a_r}^{b_r} f_r(x_0, x_1, \dots, x_r) dx_r, \quad r = 1, 2, \dots, k \end{aligned}$$

In other words,

$$\int_{I^k} f(\bar{x}) d\bar{x} = \int \cdots \int_{a_k}^{b_k} [f(x_0, x_1, \dots, x_k) dx_k] dx_{k-1} \cdots dx_1$$

**Theorem 13.4.1.** For every  $f \in \mathcal{C}(I^k)$ ,  $L(f) = L'(f)$ .

*In other words, integral of a function over a  $k$ -cell is independent of the order in which those  $k$  integrations are carried out.*

*Proof.* Step 1 : “Separable” Functions ie,  $h(\bar{x}) = \prod h_i(x_i)$ .

(“separable” is not standard. It is only for the purpose of understanding.)

Let  $h(\bar{x}) = h_1(x_1)h_2(x_2) \cdots h_k(x_k)$  where  $h_j \in [a_j, b_j]$ .

$$L(h) = \int_{I^k} \left( \prod_{i=1}^k h_i(x_i) \right) d\bar{x} = \prod_{i=1}^k \int_{a_k}^{b_k} h_i(x_i) dx_i = L'(h)$$

Step 2 : Algebra of “separable” functions,  $\mathcal{A}$ .



Let  $\mathcal{A}$  be all finite sums of functions such as  $h$ . Let  $g \in \mathcal{A}$ .

$$\begin{aligned} L(g) &= \int_{I^k} \left( \sum_j \prod_i h_{i,j}(x_i) \right) d\bar{x} \\ &= \sum_j \int_{I^k} \prod_i h_{(j)}(\bar{x}) d\bar{x} \\ &= \sum_j \prod_i \int_{a_k}^{b_k} h_{i,j}(x_i) d(x_i) \\ &= L'(g) \end{aligned}$$

**Step 3 :** All functions continuous on  $I^k$ .

Let  $f \in \mathcal{C}(I^k)$ . ie, a function which is continuous in  $I^k$ .

**Stone-Weierstrass theorem** - Let  $\mathcal{A}$  be an algebra of real, continuous functions on a compact set  $K$ . If  $\mathcal{A}$  separates points on  $K$  and if  $\mathcal{A}$  vanishes at no point of  $K$ , then the uniform closure  $\mathcal{B}$  of  $\mathcal{A}$  consists of all real, continuous functions on  $K$ .

Clearly, the algebra of functions  $\mathcal{A}$  separates points on  $I^k$  and vanishes nowhere on  $I^k$ . Suppose  $\bar{x} \neq \bar{y}$ , then there exists  $m$  such that  $x_m \neq y_m$ . Thus  $h(\bar{x}) = x_m$  separates  $\bar{x}$  and  $\bar{y}$ . Again  $h(\bar{x}) = 1$  vanishes nowhere on  $I^k$ . Thus every function which is continuous on  $I^k$  is the limit of some uniformly convergent sequence of functions in  $\mathcal{A}$ .

Let  $V = \prod_{j=1}^k (b_j - a_j)$ . Then by Stone-Weierstrass theorem, for any  $\epsilon > 0$ , there exists a function  $g \in \mathcal{A}$  such that  $\|f - g\| < \frac{\epsilon}{V}$  where the norm of a function  $f$  is defined by  $\|f\| = \max\{f(\bar{x}) : \bar{x} \in I^k\}$ .

Therefore, it is sufficient to prove that  $\|L(f) - L'(f)\| < \epsilon$ . Since  $\|f - g\| < \frac{\epsilon}{V}$ ,  $|L(f - g)| < \epsilon$  and  $|L'(f - g)| < \epsilon$ . Thus,

$$\begin{aligned} L(f) - L'(f) &= L(f) - L(g) + L'(g) - L'(f) \\ &= L(f - g) + L'(g - f) \\ |L(f) - L'(f)| &< 2\epsilon \end{aligned}$$

Therefore,  $L(f) = L'(f)$ . □

**Definitions 13.4.2.** Let  $f : \mathbb{R}^k \rightarrow \mathbb{C}$ . The **support** of  $f$  is the closure of the set of all points  $\bar{x} \in \mathbb{R}^k$  such that  $f(\bar{x}) \neq 0$ .

*Remark.* Let  $f$  be a continuous function with compact support. And  $I^k$  be any  $k$ -cell containing the support of  $f$ . Then,  $\int_{\mathbb{R}^k} f d\bar{x} = \int_{I^k} f d\bar{x}$ .

$L(f) = L'(f)$  where  $f$  is the limit function of a sequence of functions with compact support. [Apostol, 1973, §10.4 Example]

**Definitions 13.4.3.** Let  $E$  be an open subset in  $\mathbb{R}^n$ . Then function  $G : E \rightarrow \mathbb{R}^n$  is **primitive** if it satisfies

$$G(\bar{x}) = \sum_{i \neq m} x_i \bar{e}_i + g(\bar{x}) \bar{e}_m \quad (13.15)$$

for some integer  $m$  and some function  $g : E \rightarrow \mathbb{R}$  (where  $\bar{e}_i$  are the unit co-ordinate vectors).

$$G(\bar{x}) = \bar{x} + [g(\bar{x}) - x_m] \bar{e}_m \quad (13.16)$$

If  $g$  is differentiable at  $\bar{a}$ , then  $G$  is also differentiable at  $\bar{a}$ . The matrix  $[\alpha_{i,j}]$  of  $G'(\bar{a})$  is given by

$$\alpha_{i,j} = \begin{cases} D_j g(\bar{a}) & i = m \\ 1 & i \neq m, j = i \\ 0 & i \neq m, j \neq i \end{cases}$$

The Jacobian of  $G$  at  $\bar{a}$  is given by,  $J_G(\bar{a}) = \det[G'(\bar{a})] = D_m g(\bar{a})$ .

Total derivative  $G'(\bar{a})$  is invertible if and only if  $D_m g(\bar{a}) \neq 0$ .

**Definitions 13.4.4.** A linear operator  $B$  on  $\mathbb{R}^n$  that interchanges some pair of members of the standard basis and leaves the others fixed is a **flip**.

**Theorem 13.4.2.** Suppose  $F$  is a  $\mathcal{C}'$ -mapping of an open set  $E \subset \mathbb{R}^n$  into  $\mathbb{R}^n$ ,  $\bar{0} \in E$ ,  $F(\bar{0}) = \bar{0}$ , and  $F'(\bar{0})$  is invertible. Then there is a neighbourhood of  $\bar{0}$  in  $\mathbb{R}^n$  in which the representation  $F(\bar{x}) = B_1 B_2 \cdots B_{n-1} G_n \circ G_{n-1} \circ \cdots \circ G_1(\bar{x})$  is valid where  $G_i$  are primitive  $\mathcal{C}'$ -mapping in some neighbourhood of  $\bar{0}$ ,  $G_i(\bar{0}) = \bar{0}$ ,  $G'(\bar{0})$  is invertible and each  $B_i$  is either a flip or identity operator.

*In other words, locally  $F$  is a composition of primitive mappings and flips.*

*Proof.* Proof by mathematical induction on  $m$ .

There exists a neighbourhood of  $\bar{0}$ ,  $V_m$  such that  $F_m \in \mathcal{C}'(V_m)$  (13.17)

$F_m(\bar{0}) = \bar{0}$  (13.18)

$F'_m(\bar{0})$  is invertible and (13.19)

$P_{m-1}F_m(\bar{x}) = P_{m-1}\bar{x}$ ,  $\bar{x} \in V_m$  (13.20)

where the  $k$ th projection  $P_k : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is defined by  $P_k(\bar{x}) = x_1\bar{e}_1 + x_2\bar{e}_2 + \cdots + x_k\bar{e}_k + 0\bar{e}_{k+1} + \cdots + 0\bar{e}_n$ . Clearly,  $P_0(\bar{x}) = \bar{0}$ .

The essence of this proof lies in the fourth statement which is designed to construct a sequence of functions  $F_1, F_2, \dots, F_n$  such that the other three statements remains true for every function in this sequence.

**Step 1 : Initial Case,** prove that all the four statements are true for  $m = 1$ .

Define  $F_1 = F$  and  $V_1 = E$ . Thus,  $F \in \mathcal{C}'(E) \implies F_1 \in \mathcal{C}'(V_1)$ . Also we have,  $F_1(\bar{0}) = F(\bar{0}) = \bar{0}$ , and  $F'_1(\bar{0}) = F'(\bar{0})$  is invertible. Obviously the trivial projection,  $P_0(F_1(\bar{x})) = P_0(F(\bar{x})) = \bar{0} = P_0(\bar{x})$  for every  $\bar{x} \in V_1$ .

**Step 2 : Induction Hypothesis,** suppose all the four statements are true for  $m = 1, 2, \dots, n-1$ .

Suppose that for each  $m = 1, 2, \dots, n-1$ , there exists a neighbourhood of  $\bar{0}$ , say  $V_m$  such that  $F_m \in \mathcal{C}'(V_m)$ ,  $F_m(\bar{0}) = \bar{0}$ ,  $F'_m(\bar{0})$  is invertible and  $P_{m-1}F_m(\bar{x}) = P_{m-1}\bar{x}$  for every  $\bar{x} \in V_m$ .

**Step 3 : Induction Step,** prove that all the four statements are true for  $m = n$ .

$$\begin{aligned} P_{m-1}F_m(\bar{x}) &= P_{m-1}\bar{x}, \bar{x} \in V_m \\ &= x_1\bar{e}_1 + x_2\bar{e}_2 + \cdots + x_{m-1}\bar{e}_{m-1} \end{aligned}$$

We may write remaining components of  $F_m(\bar{x})$  using real-valued functions. That is,  $k$ th component of the function  $F_m$ , say  $F_{m_k} = \alpha_k$ .

$$F_m(\bar{x}) = P_{m-1}(\bar{x}) + \sum_{i=m}^n \alpha_i(\bar{x})\bar{e}_i$$

Since  $F_m \in \mathcal{C}'(V_m)$ , the functions  $\alpha_i \in \mathcal{C}'(V_m)$ . Taking  $m$ th partial derivative on either sides, we get

$$D_m F_m(\bar{0}) = F'_m(\bar{0})\bar{e}_m = \sum_{i=m}^n D_m \alpha_i(\bar{0})\bar{e}_i$$

We have,  $F'_m(\bar{0})$  is invertible. Thus  $F'_m(\bar{0})\bar{e}_m \neq 0$ . Thus there exists some integer  $k$  such that  $m \leq k \leq n$  and  $D_m \alpha_k(\bar{0})\bar{e}_k \neq 0$ . Let  $B_m$  be the flip that interchanges  $m$  and  $k$ . If  $m = k$ , then  $B_m$  is identity map.

$$\text{Define } G_m(\bar{x}) = \bar{x} + [\alpha_k(\bar{x}) - x_m]\bar{e}_m \quad (13.21)$$

Then  $G_m \in \mathcal{C}'(V_m)$ ,  $G_m$  is a primitive mapping, and  $G'_m(\bar{0})$  is invertible since  $D_m \alpha_k(\bar{0}) \neq 0$ .

By inverse function theorem, there exists an open set  $U_m$ ,  $0 \in U_m$  and  $U_m \subset V_m$  such that  $G_m$  is a bijection from  $U_m$  onto a neighbourhood of  $\bar{0}$ , say  $V_{m+1}$  in which  $G_m^{-1}$  is continuously differentiable.

$$\text{Define } F_{m+1}(\bar{y}) = B_m F_m \circ G_m^{-1}(\bar{y}), \bar{y} \in V_{m+1} \quad (13.22)$$

Then  $F_m \in \mathcal{C}'(V_{m+1})$ ,  $F_{m+1}(\bar{0}) = 0$  and  $F'_{m+1}(\bar{0})$  is invertible by the chain rule. Also, for  $\bar{x} \in U_m$ , we have

$$\begin{aligned} P_m F_{m+1}(G_m(\bar{x})) &= P_m B_m F_m(\bar{x}) \\ &= P_m(P_{m-1}\bar{x} + \alpha_k(\bar{x})\bar{e}_m + \cdots) \\ &= P_{m-1}\bar{x} + \alpha_k(\bar{x})\bar{e}_m \\ &= P_m G_m(\bar{x}) \\ \implies P_m F_{m+1}(\bar{y}) &= P_m(\bar{y}) \end{aligned}$$

Thus, by mathematical induction, we have characterised a finite sequence of function  $F_1, F_2, \dots, F_n$  such that all the four statements are true.

**Step 4 :** Using the finite sequence of functions  $F_1, F_2, \dots, F_n$  constructed in the proof, we can represent  $F$  in terms of flips and primitive mappings.

Let  $\bar{x} \in U_m$  and  $\bar{y} = G_m(\bar{x})$ . We have,

$$\begin{aligned} B_m F_m \circ G_m^{-1}(\bar{y}) &= F_{m+1}(\bar{y}) \\ B_m B_m F_m(\bar{x}) &= B_m F_{m+1}(G_m(\bar{x})) \\ IF_m(\bar{x}) &= B_m F_{m+1}(G_m(\bar{x})) \\ F_m(\bar{x}) &= B_m(F_{m+1}(G_m(\bar{x}))) \end{aligned}$$

$$\text{Thus, } F_m = B_m F_{m+1} \circ G_m, \quad m = 1, 2, \dots, n \quad (13.23)$$

Therefore, we have

$$\begin{aligned} F &= F_1 \\ &= B_1 F_2 \circ G_1 \\ &= B_1(B_2 F_3 \circ G_2) \circ G_1 = B_1 B_2 F_3 \circ G_2 \circ G_1 \\ &\vdots \\ F &= B_1 B_2 \cdots B_n F_n \circ G_{n-1} \circ G_{n-2} \circ \cdots \circ G_1 \end{aligned}$$

Since  $P_n F_n(\bar{x}) = P_{n-1}(\bar{x})$ ,  $F_n$  is a primitive mapping, say  $G_n$ .

$$F = B_1 B_2 \cdots B_n G_n \circ G_{n-1} \circ G_{n-2} \circ \cdots \circ G_1$$

□

*Remark.* Let  $K$  be a compact subset of  $\mathbb{R}^n$ . Then a family of function  $\psi_1, \psi_2, \dots, \psi_s$  where  $\psi_j : K \rightarrow \mathbb{R}$  is a partition of unity if it satisfies

1.  $0 \leq \psi_j(\bar{x}) \leq 1$  for every  $\bar{x} \in K$ ,  $j = 1, 2, \dots, s$ .
2.  $\sum_{j=1}^s \psi_j(\bar{x}) = 1$  for every  $\bar{x} \in K$ .

**Theorem 13.4.3** (partitions of unity). *Suppose  $K$  is a compact subset of  $\mathbb{R}^n$ , and  $\{V_\alpha\}$  is an open cover of  $K$ . Then there exists functions  $\psi_1, \psi_2, \dots, \psi_s \in \mathcal{C}(\mathbb{R}^n)$  such that*

1.  $0 \leq \psi_i \leq 1$  for  $1 \leq i \leq s$
2. each  $\psi_i$  has its support in some  $V_\alpha$  and
3.  $\psi_1(\bar{x}) + \psi_2(\bar{x}) + \cdots + \psi_s(\bar{x}) = 1$  for every  $\bar{x} \in K$ .

*In other words, every open cover of a compact subset of  $\mathbb{R}^n$  has a partition of unity with compact support in a finite subcover.*

*Proof. Step 1 : Construction of  $\phi_j$ .*

We have,  $\{V_\alpha\}$  is a cover of  $K$ . Therefore, every  $\bar{x} \in K$  belongs to some  $V_{\alpha(\bar{x})}$ . And there exists open balls  $B(\bar{x})$  and  $W(\bar{x})$  such that

$$\bar{x} \in B(\bar{x}) \subset \bar{B}(\bar{x}) \subset W(\bar{x}) \subset \bar{W}(\bar{x}) \subset V_{\alpha(\bar{x})} \quad (13.24)$$

The family  $\{B(\bar{x})\}$  is an open cover of  $K$ . Since  $K$  is compact, the open cover  $\{B(\bar{x})\}$  has a finite subcover, say  $\{B(\bar{x}_1), B(\bar{x}_2), \dots, B(\bar{x}_s)\}$ . That is, there are points  $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_s$  in  $K$  such that

$$K = B(\bar{x}_1) \cup B(\bar{x}_2) \cup \cdots \cup B(\bar{x}_s) \quad (13.25)$$

Since  $\mathbb{R}^n$  is a metric space,<sup>6</sup> there exists continuous functions  $\phi_j : \mathbb{R}^n \rightarrow [0, 1]$  such that  $\phi_j(B(\bar{x}_j)) = \{1\}$ ,  $\phi_j(\mathbb{R}^n - W(\bar{x}_j)) = \{0\}$  and  $0 \leq \phi_j(\bar{x}) \leq 1$  for every  $\bar{x} \in \mathbb{R}^n$ .

*Step 2 : Construction of  $\psi_j$ .*

$$\begin{aligned} \text{Define, } \psi_1 &= \phi_1 \\ \psi_2 &= (1 - \phi_1)\phi_2 \\ &\vdots \\ \psi_s &= (1 - \phi_1)(1 - \phi_2) \cdots (1 - \phi_{s-1})\phi_s \end{aligned}$$

Clearly,  $0 \leq \psi_j \leq 1$  for  $j = 1, 2, \dots, s$ . And  $\psi_j$  has a compact support in  $W(\bar{x}_j) \subset V_{\alpha(\bar{x}_j)}$ .

*Step 3 :  $\{\psi_j\}$  is a partition of unity.*

<sup>6</sup>metric spaces have a simpler version of Urysohn's lemma or apply Urysohn's lemma

We claim that,

$$\psi_1 + \psi_2 + \cdots + \psi_j = 1 - (1 - \phi_1)(1 - \phi_2) \cdots (1 - \phi_j), \quad j = 1, 2, \dots, s \quad (13.26)$$

It is proved using mathematical induction on  $j$ . Clearly, it is true for  $j = 1$ .  $\psi_1 = \phi_1 = 1 - 1 + \phi = 1 - (1 - \phi)$ . Suppose the claim is true for  $j = k$  for some  $1 \leq k < s$ . Then we have,  $\psi_1 + \psi_2 + \cdots + \psi_k = 1 - (1 - \phi_1)(1 - \phi_2) \cdots (1 - \phi_k)$ .

$$\begin{aligned} \text{Thus, } \psi_1 + \psi_2 + \cdots + \psi_k + \psi_{k+1} &= 1 - (1 - \phi_1)(1 - \phi_2) \cdots (1 - \phi_k) \\ &\quad + (1 - \phi_1)(1 - \phi_2) \cdots (1 - \phi_k)\phi_{k+1} \\ &= 1 - (1 - \phi_1)(1 - \phi_2) \cdots (1 - \phi_{k+1}) \end{aligned}$$

Thus, the claim is true for  $j = 1, 2, \dots, s$ . Let  $\bar{x} \in K$ . Then  $\bar{x} \in B(\bar{x}_j)$  for some  $j$ . By definition of  $\phi_j$ ,  $\phi_j(\bar{x}) = 1$ . That is,  $(1 - \phi_j(\bar{x})) = 0$ .

$$\begin{aligned} \implies (1 - \phi_1(\bar{x}))(1 - \phi_2(\bar{x})) \cdots (1 - \phi_s(\bar{x})) &= 0 \\ \implies \psi_1(\bar{x}) + \psi_2(\bar{x}) + \cdots + \psi_s(\bar{x}) &= 1 \end{aligned}$$

Therefore,  $\psi_1(\bar{x}) + \psi_2(\bar{x}) + \cdots + \psi_s(\bar{x}) = 1$  for every  $\bar{x} \in K$ .  $\square$

**Definitions 13.4.5.** By theorem, every open cover  $\{V_\alpha\}$  of a compact subset  $K$  has a partition of unity  $\{\psi_j\}$  with compact support in a finite subcover  $\{V_{\alpha_j}\}$ . Then  $\{\psi_j\}$  is **subordinate** to the cover  $\{V_\alpha\}$ .

**Corollary 13.4.3.1.** If  $f \in \mathcal{C}(\mathbb{R}^n)$  and the support of  $f$  lies in  $K$ , then

$f = \sum_{i=1}^s \psi_i f$ . Each  $\psi_i$  has its support in some  $V_\alpha$ .

*Any continuous function with support in a compact set can be represented as sum of continuous functions  $\psi_j f$  with small supports.*

*Proof.* Let  $K$  be compact subset of  $\mathbb{R}^n$ . And support of a function  $f$  lies in  $K$ . Then,

$$\begin{aligned} f(\bar{x}) &= I(\bar{x})f(\bar{x}) = \left( \sum_{j=1}^s \psi_j(\bar{x}) \right) f(\bar{x}) \\ &= \sum_{j=1}^s \psi_j(\bar{x})f(\bar{x}) = \sum_{j=1}^s (\psi_j f)(\bar{x}) \end{aligned}$$

Let  $\{\psi_j\}$  be a partition of unity. Let  $K'$  be the support of  $f$  and  $V_{\alpha_j}$  be support of each  $\psi_j$ . Then  $K' \cap V_{\alpha_j}$  is the support of each  $\psi_j f$ .  $\square$

**Theorem 13.4.4** (effect of change of variable on multiple integral). Suppose  $T$  is a one-to-one  $\mathcal{C}'$ -mapping of an open set  $E \subset \mathbb{R}^k$  into  $\mathbb{R}^k$  such that  $J_T(\bar{x}) \neq 0$  for all  $\bar{x} \in E$ . If  $f$  is a continuous function on  $\mathbb{R}^k$  whose support is compact and lies in  $T(E)$ , then

$$\int_{\mathbb{R}^k} f(\bar{y}) \, d\bar{y} = \int_{\mathbb{R}^k} f(T(\bar{x})) |J_T(\bar{x})| \, d\bar{x} \quad (13.27)$$

*Proof.* **Step 1 : ‘separable’ functions**

Let  $E$  be an open subset of  $\mathbb{R}^k$ . We claim that<sup>7</sup> the theorem is true for functions  $h : E \rightarrow \mathbb{R}$  of the form  $h(\bar{y}) = h_1(y_1)h_2(y_2)\cdots h_k(y_k)$ . We know that, every continuous function  $T$  is locally a composition of primitives and flips. Therefore it is enough to prove that the theorem is true for functions  $h$  with transformations  $T$  primitives, flips and their compositions.

**Step 2 : Transformation  $T$  is a primitive**

Let  $G$  be a primitive with change in  $m$ th coordinate. Then

$$G(\bar{x}) = x_1\bar{e}_1 + x_2\bar{e}_2 + \cdots + g(\bar{x})\bar{e}_m + \cdots + x_k\bar{e}_k \quad (13.28)$$

We have,

$$J_G(\bar{x}) = \begin{vmatrix} 1 & 0 & \cdots & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & 0 \\ D_1g(\bar{x}) & D_2g(\bar{x}) & \cdots & D_mg(\bar{x}) & \cdots & D_{k-1}g(\bar{x}) & D_kg(\bar{x}) \\ \vdots & \vdots & \ddots & 0 & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & \cdots & 0 & 1 \end{vmatrix} \quad (13.29)$$

Therefore,  $J_G(\bar{x}) = D_mg(\bar{x})$ .

$$\begin{aligned} \int_{\mathbb{R}^k} h(\bar{y}) \, d\bar{y} &= \prod_{j=1}^k \int h_j(y_j) \, dy_j \\ &= \left( \prod_{j \neq m} \int h_j(y_j) \, dy_j \right) \int h_m(y_m) \, dy_m \\ G(\bar{x}) = \bar{y} &\implies (x_1, x_2, \dots, g(\bar{x}), \dots, x_k) = (y_1, y_2, \dots, y_m, \dots, y_k) \\ &\implies x_j = y_j \text{ for } j \neq m \text{ and } g(\bar{x}) = y_m \\ \int_{\mathbb{R}^k} h(\bar{y}) \, d\bar{y} &= \left( \prod_{j \neq m} \int h_j(x_j) \, 1 \, dx_j \right) \int h_m(g(\bar{x})) \, D_mg(\bar{x}) \, dx_m \\ \int_{\mathbb{R}^k} h(\bar{y}) \, d\bar{y} &= \int_{\mathbb{R}^k} h(G(\bar{x})) \, |J_G(\bar{x})| \, d\bar{x} \end{aligned}$$

**Step 3 : Transformation  $T$  is a flip.**

Let  $B$  be a flip that interchanges  $m$  and  $n$  coordinates. When  $m = n$ ,  $B$  is an identity map and the theorem is trivially true.

$$B(\bar{x}) = x_1\bar{e}_1 + \cdots + x_n\bar{e}_m + \cdots + x_m\bar{e}_n + \cdots + x_k\bar{e}_k \quad (13.30)$$

The Jacobian matrix is an identity matrix with  $m$ th and  $n$ th rows interchanged.

---

<sup>7</sup>This proof is my own work. There may exist simpler proofs.

For example,

$$J_B(\bar{x}) = \begin{vmatrix} 0 & 1 & \cdots & 0 \\ 1 & 0 & \vdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{vmatrix} \quad (13.31)$$

Therefore,  $J_B(\bar{x}) = \pm 1$ .

$$\begin{aligned} \int_{\mathbb{R}^k} h(\bar{y}) \, d\bar{y} &= \prod_{j=1}^k \int h_j(y_j) \, dy_j \\ &= \left( \prod_{j \neq m, n} \int h_j(y_j) \, dy_j \right) \int h_n(y_n) \, dy_n \int h_m(y_m) \, dy_m \\ B(\bar{x}) = \bar{y} &\implies (x_1, \dots, x_n, \dots, x_m, \dots, x_k) = (y_1, \dots, y_m, \dots, y_n, \dots, y_k) \\ \int_{\mathbb{R}^k} h(\bar{y}) \, d\bar{y} &= \left( \prod_{j \neq m, n} \int h_j(x_j) \, dx_j \right) \int h_n(x_m) \, dx_n \int h_m(x_n) \, dx_m \\ &= \int_{\mathbb{R}^k} h(B(\bar{x})) \, 1 \, d\bar{x} = \int_{\mathbb{R}^k} h(B(\bar{x})) \, |J_B(\bar{x})| \, d\bar{x} \end{aligned}$$

Notice that, substituting  $y_n = x_m$  gives  $\int h_n(x_m) \, dx_n$ . This is due to the fact that  $\bar{y} = B(\bar{x})$  have  $x_m \bar{e}_n$  in place of  $y_n \bar{e}_n$ . Therefore, parameter of  $h_n$  is on the same axis  $\bar{e}_n$ .

**Step 4 :** If the theorem is true for two transformations, then it is true for their composition.

Suppose the theorem is true for transformations  $P$  and  $Q$ . And  $S = P \circ Q$ . Let  $\bar{z} = P(\bar{y})$  and  $\bar{y} = Q(\bar{x})$ . Then  $\bar{z} = S(\bar{x})$ . Since  $m(P)m(Q) = m(S)$ , we have  $J_P(\bar{y})J_Q(\bar{x}) = J_S(\bar{x})$ .

$$\begin{aligned} \int_{\mathbb{R}^k} f(\bar{z}) \, d\bar{z} &= \int_{\mathbb{R}^k} f(P(\bar{y})) \, |J_P(\bar{y})| \, d\bar{y} \\ &= \int_{\mathbb{R}^k} f(P(Q(\bar{x}))) \, |J_P(Q(\bar{x}))| \, |J_Q(\bar{x})| \, d\bar{x} \\ &= \int_{\mathbb{R}^k} f(S(\bar{x})) \, |J_S(\bar{x})| \, d\bar{x} \end{aligned}$$

Therefore, the theorem is true for their composition.

Now for any function of the form  $h(\bar{x}) = \prod_{j=1}^k h_j(x_j)$ , the theorem is true for any continuous transformation  $T$ . Then it is true for the algebra of functions  $\mathcal{A}$ . And by Stone-Weierstrass theorem,<sup>8</sup> the theorem is true for every continuous function  $f$ .  $\square$

### 13.4.1 Differential Forms

**Definitions 13.4.6.** Suppose  $E$  is an open set in  $\mathbb{R}^n$ . A  **$k$ -surface** in  $E$  is a  $\mathcal{C}^1$ -mapping  $\Phi$  from a compact set  $D \subset \mathbb{R}^k$  into  $E$ .  $D$  is the **parameter**

<sup>8</sup>I haven't checked whether the application of Stone-Weierstrass theorem causes any problem. It is the same proof technique as we have seen in  $L(f) = L'(f)$ .

domain of  $\Phi$ .

**Definitions 13.4.7.** Suppose  $E$  is an open set in  $\mathbb{R}^n$ . A **differential form** of order  $k \geq 1$  in  $E$ , a  $k$ -form in  $E$  is a function  $\omega$ , symbolically represented by the sum

$$\omega = \sum a_{i_1 \dots i_k}(\bar{x}) \, dx_{i_1} \wedge dx_{i_2} \wedge \dots \wedge dx_{i_k} \quad (13.32)$$

which assigns to each  $k$ -surface  $\Phi$  in  $E$  a number  $\omega(\Phi) = \int_{\Phi} \omega$ , according to the rule

$$\int_{\Phi} \omega = \int_D \sum a_{i_1 \dots i_k}(\Phi(\bar{u})) \frac{\partial(x_{i_1}, x_{i_2}, \dots, x_{i_k})}{\partial(u_1, u_2, \dots, u_k)} d\bar{u} \quad (13.33)$$

where  $D$  is the parameter domain of  $\Phi$  and

$$\frac{\partial(x_{i_1}, x_{i_2}, \dots, x_{i_k})}{\partial(u_1, u_2, \dots, u_k)} = \begin{vmatrix} \frac{\partial x_{i_1}}{\partial u_1} & \frac{\partial x_{i_1}}{\partial u_2} & \dots & \frac{\partial x_{i_1}}{\partial u_k} \\ \frac{\partial x_{i_2}}{\partial u_1} & \frac{\partial x_{i_2}}{\partial u_2} & \dots & \frac{\partial x_{i_2}}{\partial u_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_{i_k}}{\partial u_1} & \frac{\partial x_{i_k}}{\partial u_2} & \dots & \frac{\partial x_{i_k}}{\partial u_k} \end{vmatrix} \quad (13.34)$$

Let  $\omega = 4 \, dx_1 \wedge dx_3 + 3x_2^2 \, dx_2 \wedge dx_1$ . Then  $\omega$  is a 2-form with  $a_{1,3}(x, y, z) = 4$ ,  $a_{2,1}(x, y, z) = 3y^2$  and all other  $a_{i_1 i_2}(x, y, z) = 0$ .

Consider the upper hemi-sphere of unit radius,  $S_{y \geq 0}^2$  in  $\mathbb{R}^3$  and the closure of the unit disc  $\bar{S}^1 = D \subset \mathbb{R}^2$ . Then  $D$  is compact. Let  $\Phi : D \rightarrow S_{y \geq 0}^2$  defined by  $\Phi(x, y) = (x, +\sqrt{1-x^2-y^2}, y)$ . Clearly,  $\Phi \in \mathcal{C}'(D)$ . Thus,  $\Phi$  is a 2-surface with parameter domain  $D$ .

Clearly,  $\Phi_1(x, y) = x$ ,  $\Phi_2(x, y) = \sqrt{1-x^2-y^2}$  and  $\Phi_3(x, y) = y$ . We have,  $\omega(\Phi) = \int_{\Phi} \omega$ .

$$\omega(\Phi) = \int_{\Phi} \omega = \iint_D \omega_{\Phi}(x, y) \frac{\partial(\Phi_{i_1}, \Phi_{i_2})}{\partial(x, y)} dx \, dy \quad (13.35)$$

where  $\omega_{\Phi}(x, y)$  is obtained by evaluating each function  $a_{i_1 i_2 \dots i_k}(\Phi(x, y))$ .

$$\begin{aligned} \omega(\Phi) &= \iint_D a_{1,3}(\Phi(x, y)) \begin{vmatrix} \frac{\partial \Phi_1}{\partial x} & \frac{\partial \Phi_1}{\partial y} \\ \frac{\partial \Phi_3}{\partial x} & \frac{\partial \Phi_3}{\partial y} \end{vmatrix} dx dy + \iint_D a_{2,1}(\Phi(x, y)) \begin{vmatrix} \frac{\partial \Phi_2}{\partial x} & \frac{\partial \Phi_2}{\partial y} \\ \frac{\partial \Phi_1}{\partial x} & \frac{\partial \Phi_1}{\partial y} \end{vmatrix} dx dy \\ &= \int_0^1 \int_0^{\sqrt{1-y^2}} 4 \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} dx dy \\ &\quad + \int_0^1 \int_0^{\sqrt{1-y^2}} 3(1-x^2-y^2) \begin{vmatrix} \frac{1}{2} \frac{-2x}{\sqrt{1-x^2-y^2}} & \frac{1}{2} \frac{-2y}{\sqrt{1-x^2-y^2}} \\ 1 & 0 \end{vmatrix} dx dy \\ &= \int_0^1 \int_0^{\sqrt{1-y^2}} 4 dx dy + \int_0^1 \int_0^{\sqrt{1-y^2}} \frac{3}{\sqrt{1-x^2-y^2}} dx \, dy \end{aligned}$$

*Remark* (example 1). Integrals of 1-forms are line integrals. And  $\omega(\gamma) = 0$  for every closed curve  $\gamma$ . [Rudin, 1976, 10.12a]



*Remark* (example 2). Let  $\gamma : [0, 2\pi] \rightarrow \mathbb{R}^2$  defined by  $\gamma(t) = (a \cos t, b \sin t)$ . Then  $\gamma$  is a 1-surface with parameter domain  $[0, 2\pi]$ . [Rudin, 1976, 10.12b] Let  $\omega$  be a 1-form defined by  $\omega = x dy$ . Then

$$\omega(\gamma) = \int_{\gamma} \omega = \int_{\gamma} x dy = \int_0^{2\pi} ab \cos^2 t dt = \pi ab$$

Similarly,  $\omega = y dx$  gives

$$\omega(\gamma) = \int_{\gamma} \omega = \int_{\gamma} y dx = \int_0^{2\pi} -ab \sin^2 t dt = -\pi ab$$

*Remark* (example 3). Let  $0 \leq r \leq 1$ ,  $0 \leq \theta \leq \pi$  and  $0 \leq \phi \leq 2\pi$ . Then  $D \subset \mathbb{R}^n$  defined by  $\{(r, \theta, \phi)\}$  is compact. Define  $\Phi : D \rightarrow \mathbb{R}^3$  by  $\Phi(r, \theta, \phi) = (r \sin \theta \cos \phi, r \sin \theta \sin \phi, r \cos \theta)$ . Then  $\Phi$  is a 3-surface in  $\mathbb{R}^3$ . We have

$$J_{\Phi}(r, \theta, \phi) = \begin{vmatrix} \sin \theta \cos \phi & r \cos \theta \cos \phi & -r \sin \theta \sin \phi \\ \sin \theta \sin \phi & r \cos \theta \sin \phi & r \sin \theta \cos \phi \\ \cos \theta & -r \sin \theta & 0 \end{vmatrix} = r^2 \sin \theta$$

Let  $\omega = dx_1 \wedge dx_2 \wedge dx_3$ . Then  $\omega(\Phi) = \int_{\Phi} \omega = \int_D J_{\Phi} = \frac{4\pi}{3}$  is the volume of the unit ball  $\Phi(D)$ .

*Remark.* A  $k$ -form  $\omega$  is of class  $\mathcal{C}'$  or  $\mathcal{C}''$  if the functions  $a_{i_1 \dots i_k}$  are all of class  $\mathcal{C}'$  or  $\mathcal{C}''$ . A 0-form in  $E$  is defined to be a continuous function in  $E$ . And 0 is the only  $k$ -form in any open set  $E \subset \mathbb{R}^n$ .

**Definitions 13.4.8.** Let  $\omega = a(\bar{x}) dx_{i_1} \wedge dx_{i_2} \wedge \dots \wedge dx_{i_k}$ . Then  $\bar{\omega}$  is the  $k$ -form obtained by interchanging a pair subscripts of  $\omega$ . ie,  $\bar{\omega} = a(\bar{x}) dx_{i_2} \wedge dx_{i_1} \wedge \dots \wedge dx_{i_k}$ .

### Elementary properties of $k$ -forms

**Definitions 13.4.9.** Let  $E$  be an open set in  $\mathbb{R}^n$ . And  $\Phi$  be a  $k$ -surface in  $E$ . And let  $\omega_1, \omega_2$  be  $k$ -forms in  $E$ , then

1.  $\omega_1 = \omega_2 \iff \omega_1(\Phi) = \omega_2(\Phi) \iff \int_{\Phi} \omega_1 = \int_{\Phi} \omega_2, \forall \Phi \in E$
2.  $\omega = 0 \iff \omega(\Phi) = 0 \iff \int_{\Phi} \omega = 0, \forall \Phi \in E$
3.  $k$ -form Addition,  $\omega_1 + \omega_2$   
 $\omega = \omega_1 + \omega_2 \iff \omega(\Phi) = \omega_1(\Phi) + \omega_2(\Phi) \iff \int_{\Phi} \omega = \int_{\Phi} \omega_1 + \int_{\Phi} \omega_2.$
4. Scalar multiplication,  $c\omega$   
 $c\omega(\Phi) = c(\omega(\Phi)) \iff \int_{\Phi} c\omega = c \int_{\Phi} \omega.$
5. Inverse  $k$ -form,  $-\omega$   
 $-\omega(\Phi) = -(\omega(\Phi)) \iff \int_{\Phi} -\omega = - \int_{\Phi} \omega$
6.  $\bar{\omega} = -\omega.$

*Remark.* Let  $\omega = a(\bar{x}) dx_{i_1} \wedge dx_{i_2} \wedge \dots \wedge dx_{i_k}$ . If  $\bar{\omega} = \omega$ , then  $\omega = 0$ . Since  $\wedge$  is anticommutative.

Thus, differential  $k$ -forms with repeated subscripts are 0. For example,  $\omega = dx_i \wedge dx_j \wedge dx_i = 0$ .

**Definitions 13.4.10.** Let  $\bar{I} = (i_1, i_2, \dots, i_k)$  be an **increasing  $k$ -index**. That is,  $1 \leq i_1 \leq i_2 \leq \dots \leq i_k \leq n$ . Then  $dx_{\bar{I}}$  of the form  $dx_{\bar{I}} = dx_{i_1} \wedge dx_{i_2} \cdots dx_{i_k}$  is **basic  $k$ -form** in  $\mathbb{R}^n$ .

*Remark.* List of all basic  $k$ -forms

**0-form** 0

**1-forms**  $dx_1, dx_2, \dots, dx_n$

**2-forms**  $dx_i \wedge dx_j$  for every  $1 \leq i, j \leq n$ .

**3-forms**  $dx_i \wedge dx_j \wedge dx_k$  for every  $1 \leq i, j, k \leq n$ .

*Remark.* There are  $\binom{n}{k}$  basic  $k$ -forms in  $\mathbb{R}^n$ .

Every  $k$ -form can be represented in terms of basic  $k$ -forms. For every  $k$ -tuple  $(j_1, j_2, \dots, j_k)$ ,  $dx_{j_1} \wedge \cdots \wedge dx_{j_k} = \sigma(j_1, j_2, \dots, j_k) dx_{\bar{J}}$  where  $\bar{J}$  is the increasing  $k$ -index obtained by interchanging pairs. And  $\sigma$  maps odd permutations to  $-1$  and even permutations to  $1$ .

Standard representation of a  $k$ -form

$$\omega = \sum_I b_I(\bar{x}) dx_I \quad (13.36)$$

For example :  $x_1 dx_2 \wedge dx_1 - x_2 dx_3 \wedge dx_2 + x_3 dx_2 \wedge dx_3 + dx_1 \wedge dx_2 = (1-x) dx_1 \wedge dx_2 + (x_2+x_3) dx_2 \wedge dx_3$  is a 2-form in  $\mathbb{R}^3$ .

$dx_1 \wedge dx_2 \wedge dx_3 = dx_2 \wedge dx_3 \wedge dx_1$ , since  $(1\ 2\ 3) \in A_3 \implies \sigma(1\ 2\ 3) = 1$ . And,  $dx_1 \wedge dx_2 \wedge dx_3 = -dx_1 \wedge dx_3 \wedge dx_2$ ,  $(2\ 3) \notin A_3 \implies \sigma(2\ 3) = -1$ . Here,  $A_3$  is an alternating group of all even permutation on  $\{1, 2, 3\}$ .

## Subject 14

# ME010304 Functional Analysis

14.1 Metric Spaces

14.2 Banach Spaces

14.3 Hilbert Spaces

14.4 Fundamental Theorems for Banach Spaces

**Subject 15**

**ME010305 Optimization  
Technique**

## Semester IV

**Subject 16**

**ME010401 Spectral Theory**

**Subject 17**

**ME010402 Analytic  
Number Theory**

## Subject 18

# ME800401 Differential Geometry

### 18.1 Graphs and Level Set

**Definitions 18.1.1.** Let function  $f : U \rightarrow \mathbb{R}$  where  $U \subset \mathbb{R}^{n+1}$ . Let  $c$  be a real number. Then the **Level set** of  $f$  at height  $c$  is the set of all points in  $U$  with image  $c$ .

$$f^{-1}(c) = \{(x_1, x_2, \dots, x_{n+1}) \in U : f(x_1, x_2, \dots, x_{n+1}) = c\} \quad (18.1)$$

**Definitions 18.1.2.** Let function  $f : U \rightarrow \mathbb{R}$  where  $U \subset \mathbb{R}^{n+1}$ . Then,

$$\text{graph}(f) = \{(x_1, x_2, \dots, x_{n+2}) \in \mathbb{R}^{n+2} : f(x_1, x_2, \dots, x_{n+1}) = x_{n+2}\} \quad (18.2)$$

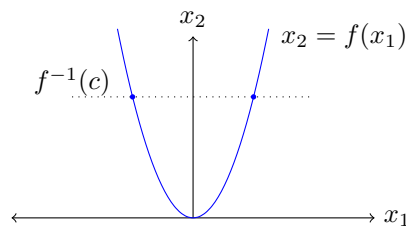


Figure 18.1: Graph of  $f(x_1) = x_1^2$  and Level set  $f^{-1}(c)$

### 18.2 Vector Fields

**Definitions 18.2.1.** A vector  $\mathbf{v}$  at a point  $p \in \mathbb{R}^{n+1}$  is a pair  $\mathbf{v} = (p, v)$  where  $v \in \mathbb{R}^{n+1}$ .

**vector addition**  $\mathbf{v} + \mathbf{w} = (p, v) + (p, w) = (p, v + w)$ .

**scalar multiplication** Let  $c \in \mathbb{R}$ , then  $c\mathbf{v} = c(p, v) = (p, cv)$ .

**dot product**  $\mathbf{v} \cdot \mathbf{w} = (p, v) \cdot (p, w) = v \cdot w$



**cross product**  $\mathbf{v} \times \mathbf{w} = (p, v) \times (p, w) = (p, v \times w)$

*Remark.* Angle  $\theta$  between  $\mathbf{v}$  and  $\mathbf{w}$  is given by,

$$\cos \theta = \mathbf{v} \cdot \mathbf{w} = (p, v) \cdot (p, w) = v \cdot w \quad (18.3)$$

And the length of a vector  $\mathbf{v}$  is given by,

$$\|\mathbf{v}\| = \mathbf{v} \cdot \mathbf{v} = (p, v) \cdot (p, v) = v \cdot v = \|v\| \quad (18.4)$$

*Remark.* Let  $c \in \mathbb{R}$  and  $p \in \mathbb{R}^{n+1}$ . Let  $\mathbf{v}, \mathbf{w}$  be two vectors at  $p$ . That is,  $\mathbf{v} = (p, v)$  and  $\mathbf{w} = (p, w)$  for some  $v, w \in \mathbb{R}^{n+1}$ . Then the set of all vectors at  $p$  is a vector space with vector addition  $\mathbf{v} + \mathbf{w} = (p, v + w)$  and scalar multiplication  $c\mathbf{v} = (p, cv)$ . This vector space is denoted by  $\mathbb{R}_p^{n+1}$ .

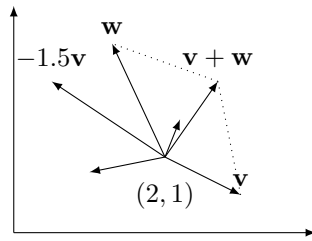


Figure 18.2: The vector space of all vectors at  $(2, 1)$ ,  $\mathbb{R}_{(2,1)}^2$

**Definitions 18.2.2.** The vector field  $\mathbf{X}$  on  $\mathbb{R}^{n+1}$  is a function which assigns to each point of  $\mathbb{R}^{n+1}$  a vector at that point. That is,  $\mathbf{X}(p) = (p, X(p))$ .

For example,  $\mathbf{X}(p) = (p, X(p))$  where the associated function of the vector field,  $X : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  defined by  $X(p) = (1, 2)$  assigns a constant vector  $(1, 2)$  at every vector in  $\mathbb{R}^2$ .

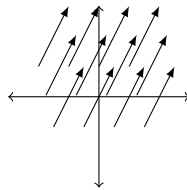


Figure 18.3: Vector field with associated function  $X(p) = (1, 2)$

**Definitions 18.2.3** (smooth). A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is smooth if its partial derivatives of all orders exists and are continuous. A function  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  is smooth if its component functions  $f = (f_1, f_2, \dots, f_{n+1})$  are smooth. A vector field  $\mathbf{X}$  is smooth if the associated function  $X(p)$  is smooth.

**Definitions 18.2.4.** Let  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ . Then the gradient of  $f$  at  $p$  is,

$$\nabla f(p) = \left( p, \frac{\partial f}{\partial x_1}(p), \frac{\partial f}{\partial x_2}(p), \dots, \frac{\partial f}{\partial x_{n+1}}(p) \right) \quad (18.5)$$

*Remark.* If  $f$  is a smooth function, then the gradient of  $f$  at  $p$  is a smooth vector field.

For example,  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  defined by  $f(x_1, x_2) = 2x_1x_2$  is a smooth function. We have,  $\frac{\partial f}{\partial x_1} = 2x_2$  and  $\frac{\partial f}{\partial x_2} = 2x_1$ . And gradient of  $f$  at  $(x_1, x_2)$  is  $(x_1, x_2, 2x_2, 2x_1)$ . That is,  $(2x_2, 2x_1)$  at  $(x_1, x_2)$ .

Calculations :

$p$	$(x_1, x_2)$	$(0, 0)$	$(1, 0)$	$(0, 1)$	$(-1, 0)$	$(0, -1)$
$X(p)$	$(2x_2, 2x_1)$	$(0, 0)$	$(0, 2)$	$(2, 0)$	$(0, -2)$	$(-2, 0)$

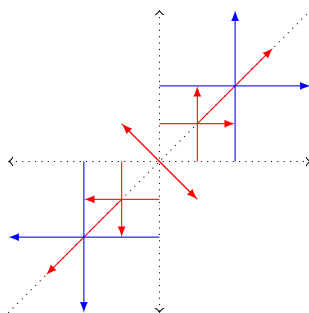


Figure 18.4: The gradient of  $f(x_1, x_2) = 2x_1x_2$

**Definitions 18.2.5.** A parameterised curve is a function,  $\alpha : I \rightarrow \mathbb{R}^{n+1}$  where  $I$  is some open interval in  $\mathbb{R}$ . The velocity vector of a parameterised curve  $\alpha : I \rightarrow \mathbb{R}^{n+1}$  at a point  $\alpha(t)$  is the tangent to the curve at that point.

$$\dot{\alpha}(t) = \left( \alpha(t), \frac{d\alpha}{dt}(t) \right) \quad (18.6)$$

For example,  $\alpha : I \rightarrow \mathbb{R}^2$  defined by  $\alpha(t) = (2t, t^2)$  is a parameterised curve. We have,  $\frac{d\alpha}{dt} = \left( \frac{dx_1}{dt}(t), \frac{dx_2}{dt}(t) \right) = (2, 2t)$  where  $\alpha(t) = (x_1(t), x_2(t))$ . The velocity vector at  $t = 3$  is  $\dot{\alpha}(3) = \left( \alpha(3), \frac{d\alpha}{dt} \right) = (6, 9, 2, 6)$ .

**Definitions 18.2.6.** Let  $\mathbf{X}$  be a vector field and let  $U$  be an open subset of  $\mathbb{R}^{n+1}$ . An integral curve  $\alpha$  on  $U$  is a parameterised curve,  $\alpha : I \rightarrow \mathbb{R}^{n+1}$  such that for each  $\alpha(t) = p \in U$ , the velocity vector  $\dot{\alpha}(t)$  is the associated vector  $\mathbf{X}(p)$  of the vector field  $\mathbf{X}$  at that point. Thus, for each  $t \in I$ ,  $\dot{\alpha}(t) = \mathbf{X}(\alpha(t))$ .

$$\left( \alpha(t), \frac{d\alpha}{dt}(t) \right) = (\alpha(t), X(\alpha(t))) \quad (18.7)$$

Let  $X(p) = (X_1(p), X_2(p), \dots, X_{n+1}(p))$  and  $\alpha(t) = (x_1(t), x_2(t), \dots, x_{n+1}(t))$ . Then, comparing components of the vector at  $\alpha(t)$  we get the following system of equations,

$$\frac{dx_j}{dt}(t) = X_j(\alpha(t)), \quad j = 1, 2, \dots, (n+1) \quad (18.8)$$

For example, Consider  $\alpha : (2, 3) \rightarrow \mathbb{R}^2$  defined by  $\alpha(t) = (t, t^2)$ . Then  $\alpha$  is a parameterised curve in vector field,  $\mathbf{X}$  which has the associated function  $X(x_1, x_2) = (1, 2x_1)$ . Then,  $\mathbf{X}(x_1, x_2) = (x_1, x_2, 1, 2x_1)$ . And

$$\dot{\alpha}(t) = \left( \alpha(t), \frac{d\alpha}{dt}(t) \right) = \left( x_1(t), x_2(t), \frac{dx_1}{dt}(t), \frac{dx_2}{dt}(t) \right) = (t, t^2, 1, 2t)$$

Clearly,  $\alpha$  is an integral curve of  $\mathbf{X}$  as  $\dot{\alpha}(t) = X(\alpha(t))$  for every  $t \in (2, 3)$ .

Calculations:

$p$	(0, 0)	(1, 0)	(0, 1)	(1, 1)	(-1, 0)	(0, -1)	(-1, 1)	(1, -1)	(-1, -1)
$X(p)$	(1, 0)	(2, 2)	(1, 1)	(2, 3)	(0, -2)	(1, 1)	(0, -1)	(2, 1)	(0, -3)

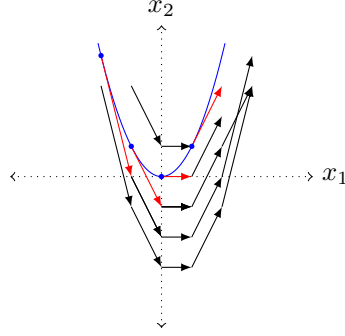


Figure 18.5: Integral Curve  $\alpha(t) = (t, t^2)$  in  $\mathbf{X}$  with  $X(x_1, x_2) = (1, 2x_1)$

**Theorem 18.2.1.** *Let  $\mathbf{X}$  be a smooth vector field on an open set  $U \subset \mathbb{R}^{n+1}$  and let  $p \in U$ . Then there exists an open interval  $I$  containing 0 and an integral curve  $\alpha : I \rightarrow U$  such that*

1.  $\alpha(0) = p$
2. If  $\beta : \tilde{I} \rightarrow U$  is any other integral curve with  $\beta(0) = p$ , then  $\tilde{I} \subset I$  and  $\beta(t) = \alpha(t)$ , for all  $t \in \tilde{I}$ .

*Proof.* Let  $\mathbf{X}$  be a smooth vector field. Suppose  $\alpha$  be an integral curve in  $\mathbf{X}$ . Then,  $\dot{\alpha}(t) = \mathbf{X}(\alpha(t))$ . Let  $x_j(t)$  be the components of  $\alpha(t)$  and  $X_j(p)$  be the components of  $X(p)$ .

$$\begin{aligned}
 \dot{\alpha}(t) &= \left( \alpha(t), \frac{d\alpha}{dt}(t) \right) \\
 &= \left( x_1(t), \dots, x_{n+1}(t), \frac{dx_1}{dt}(t), \dots, \frac{dx_{n+1}}{dt}(t) \right) \\
 \mathbf{X}(\alpha(t)) &= (\alpha(t), X(\alpha(t))) \\
 &= (x_1(t), \dots, x_{n+1}(t), X_1(\alpha(t)), \dots, X_{n+1}(\alpha(t)))
 \end{aligned}$$

Thus, we have a system of  $n + 1$  first order differential equations in  $n + 1$  unknowns satisfying the initial condition  $\alpha(0) = p$ .

$$\begin{aligned}
 \frac{dx_1}{dt}(t) &= X_1(\alpha(t)) \\
 \frac{dx_2}{dt}(t) &= X_2(\alpha(t)) \\
 &\vdots \\
 \frac{dx_{n+1}}{dt}(t) &= X_{n+1}(\alpha(t))
 \end{aligned}$$

□

By the theorem on solution of systems of first order ordinary differential equations, there exists an interval  $I$  containing 0 and a solution — a family of functions  $\{x_1(t), x_2(t), \dots, x_{n+1}(t)\}$  satisfying the above system of equations satisfying the initial condition  $\alpha(0) = p$ .

Define  $\alpha : I \rightarrow U$  using the component functions of  $\alpha$  as  $x_j$ s in the above solution. Then, we have a integral curve of the vector field  $\mathbf{X}$  satisfying the initial condition  $\alpha(0) = p$ .

Let  $\beta : \tilde{I} \rightarrow U$  be another integral curve with  $\beta(0) = p$ . Then by the uniqueness of the solution for the system of first order ordinary differential equations with an initial condition,  $\beta(t) = \alpha(t)$  for every  $t \in I \cup \tilde{I}$ .

Let  $\{\beta_1, \beta_2, \dots\}$  be the family of integral curves with  $\beta_j : I_j \rightarrow U$  satisfying  $\beta_j(0) = p$ . Consider  $I = \bigcup_{j \in \mathbb{N}} I_j$ .

Define  $\alpha : I \rightarrow U$  by  $\alpha(t) = \beta_j(t)$  where  $t \in I_j$  for some  $j \in \mathbb{N}$ . Then  $\alpha$  is well-defined and is a maximal integral curve in  $\mathbf{X}$  such that  $\alpha(0) = p$ .

**Definitions 18.2.7.** A smooth vector field  $\mathbf{X}$  on  $U \subset \mathbb{R}^{n+1}$  is **complete** if for every  $p \in U$ , the maximal integral curve through  $p$  has domain equal to  $\mathbb{R}$ .

**Definitions 18.2.8.** The **divergence** of a smooth vector field  $\mathbf{X}$  on  $U \subset \mathbb{R}^{n+1}$  is the function  $\text{div } \mathbf{X} : U \rightarrow \mathbb{R}$  defined by

$$\text{div } X(x_1, x_2, \dots, x_{n+1}) = \sum_{i=1}^{n+1} \frac{\partial X_i}{\partial x_i}$$

where  $X_i$  are the component function of the associated function  $X$  of the vector field  $\mathbf{X}$ .

For example, Consider  $\mathbf{X}$  with associated function  $X : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  defined by  $X(x_1, x_2) = (2x_1, x_1x_2)$ . Then  $\text{div } \mathbf{X}(x_1, x_2) = \frac{\partial X_1}{\partial x_1} + \frac{\partial X_2}{\partial x_2} = 2 + x_1$ .

### 18.3 The Tangent Space

### 18.4 Surfaces

### 18.5 Vector Fields on Surfaces; Orientation

### 18.6 The Gauss Map

Suppose  $S$  is an  $n$ -surface. From the definition of an  $n$ -surface, there exists a smooth function  $f : U \rightarrow \mathbb{R}$  where  $U$  is an open subset of  $\mathbb{R}^{n+1}$  such that  $S = f^{-1}(c)$  for some real value  $c \in \mathbb{R}$  and every point on  $S$  is a regular point of  $f$ . That is  $\nabla f(p) \neq \mathbf{0}$  for every point  $p$  on the surface  $S$ .

We have proved that every  $n$ -surface has exactly two orientations  $\mathbf{N}_1$  and  $\mathbf{N}_2$ . These orientations are  $\frac{\nabla f}{\|\nabla f\|}$  and  $\frac{-\nabla f}{\|\nabla f\|}$ . Given an orientation  $\mathbf{N}$  (either  $\mathbf{N}_1$  or  $\mathbf{N}_2$ ), the surface together with that orientation is collectively referred as an oriented  $n$ -surface.

Since orientation  $\mathbf{N}$  is a smooth, unit normal vector field. The vector field  $\mathbf{N}$  has an associated function  $N : U \rightarrow \mathbb{R}^{n+1}$ . That is  $\mathbf{N}(p) = (p, N(p))$  where  $N : U \rightarrow \mathbb{R}^{n+1}$ . And we already have,  $\mathbf{N}(p) = (p, N(p)) = (p, \frac{\pm \nabla f}{\|\nabla f\|})$ . This associated function restricted to the  $n$ -surface  $S$  is the Gauss Map. That is,  $N : S \rightarrow \mathbb{R}^{n+1}$ .

From the definition of orientation, we know that this function is actually assigning direction to each point on that surface  $S$ . If you don't remember, the directions are vector in  $\mathbb{R}^{n+1}$  of unit length. That is  $\|v\| = 1$ . Thus, the range of Gauss Map is a subset of the set of all directions. And unit sphere  $S^n$  is  $\mathbb{R}^{n+1}$  is the set of all directions in  $\mathbb{R}^{n+1}$ .

Thus, we may write Gauss Map,  $N : S \rightarrow S^n$

### 18.6.1 Spherical Image

We already saw that, the Gauss Map  $N : S \rightarrow S^n$  is a function which maps directions/unit vectors to each point on that oriented surface  $S$ .

Do we need an oriented surface ? Yes. The Gauss Map is defined by this orientation. If we are provided with an oriented  $n$  Surface  $S$ , then we have a unit vector/orientation assigned to each point  $p$  on that surface. And Gauss Map assigns this unit vector to the point  $p$  on surface  $S$ .

We already saw that the range of the Gauss Map is a subset of the unit  $n$  Sphere  $S^n$ . In other words, the Gauss Map assigns each point on the oriented  $n$ -surface  $S$  into a subset of the unit  $n$  sphere  $S^n$ . Thus, **range of the Gauss Map is referred as the spherical image of the oriented  $n$ -surface  $S$ .**

$$N(S) = \{q \in S^n : q = N(p), p \in S\} \quad (18.9)$$

### 18.6.2 Compact, connected, oriented $n$ Surface

Suppose we have a compact, connected, oriented  $n$ -surface  $S$ . The compact subsets in Euclidean spaces are closed and bounded subsets. And connected subsets in Euclidean Spaces are path connected.

**Theorem 18.6.1** (Spherical Image of Compact, Connected, Oriented Surface). *The Gauss map of a compact, connected, oriented  $n$ -surface is surjective.*

*Proof.* Let  $v \in S^n$  be a direction in  $\mathbb{R}^{n+1}$ . Let  $S$  be a compact, connected, oriented  $n$ -surface with orientation  $N$  such that  $S = f^{-1}(c)$  and every point  $p \in S$  are regular points of the smooth function  $f : U \rightarrow \mathbb{R}$  where  $U$  is an open subset of  $\mathbb{R}^{n+1}$ . Thus, we have the Gauss Map  $N : S \rightarrow S^n$  defined by the orientation on  $S$ .

Since  $v$  is arbitrary, it is enough to prove that  $v \in N(S)$ . Suppose there exists  $v \in S^n$  such that  $v \notin N(S)$ , then the Gauss Map is not surjective. In other words,  $N$  is surjective if for every  $v \in S^n$ ,  $v \in N(S)$  OR for every  $v \in S^n$ , there

exists  $p \in S$  such that  $v = N(p)$

Let  $g : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  defined by  $g(p) = p \cdot v$ . Then  $g$  is a smooth function since first order partial derivatives are constant functions and all other partial derivatives of higher orders vanishes.

Since  $S$  is compact,  $g$  restricted to  $S$  is a continuous function defined on a compact interval. And thus it attains maximum and minimum values, say  $p$  and  $q$ . The maximum and minimum values of the dot product  $p \cdot v$  are  $\pm v$ .

By Lagrange's multiplier theorem,  $\nabla g(p) = \lambda \nabla f(p)$  and  $\nabla g(q) = \lambda \nabla f(q)$ . From the definition of the Gauss Map, we have  $\nabla g(p) = \lambda \nabla f(p) = \lambda \|\nabla f(p)\| \mathbf{N}(p) = \lambda \|\nabla f(p)\| (p, v)$ . Thus,  $v$  and  $N(p)$  are multiples of one another. Similarly,  $\nabla g(q) = \lambda \|\nabla f(q)\| \mathbf{N}(q)$ . Therefore  $N(p) = \pm v$  and  $N(q) = \pm v$ .

It remains to show that  $N(p) \neq N(q)$ . Suppose  $N(p) = N(q)$ . If there exists continuous function  $\alpha$  such that  $\alpha : [0, 1] \rightarrow \mathbb{R}^{n+1}$ ,  $\alpha(0) = p$ ,  $\alpha(1) = q$ ,  $\dot{\alpha}(0) = (p, v)$  and  $\dot{\alpha}(1) = (q, v)$ . And  $\alpha$  maps the interior, the open interval  $(0, 1)$  outside the surface  $S$ . Then by intermediate value theorem, we arrive at a contradiction. And thus  $N(p) \neq N(q)$ .

Let  $\alpha_1 : [0, x] \rightarrow \mathbb{R}^{n+1}$  defined by  $\alpha_1(t) = p + tv$ . Let  $\alpha_2 : [y, 1] \rightarrow \mathbb{R}^{n+1}$  defined by  $\alpha_2(t) = q + (t - 1)v$ . Let  $\alpha_3 : [x, y] \rightarrow S_1$  where  $S_1$  is an  $n$  sphere properly containing  $S$ . Such an  $n$  sphere exists, since  $S$  is compact (bounded). And there exists such a function  $\alpha_3$ , the image of which is a compact subset of  $S_1$ .

Now consider  $\alpha : [0, 1] \rightarrow \mathbb{R}^{n+1}$  defined by

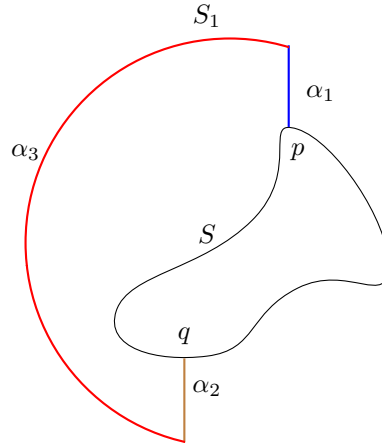
$$\alpha(t) = \begin{cases} \alpha_1(t) & t \in [0, x) \\ \alpha_3(t) & t \in [x, y] \\ \alpha_2(t) & t \in (y, 1] \end{cases} \quad (18.10)$$

Clearly,  $\alpha$  is a smooth function with  $\alpha(t) \notin S$ ,  $t \in (0, 1)$  and

$$\begin{aligned} \alpha(0) &= \alpha_1(0) = p + 0v = p \\ \alpha(1) &= \alpha_2(1) = q + (1 - 1)v = q \\ \dot{\alpha}(0) &= \frac{d\alpha_1}{dt}(0) = \frac{d(p + tv)}{dt}(0) = v \\ \dot{\alpha}(1) &= \frac{d\alpha_2}{dt}(1) = \frac{d(q + (t - 1)v)}{dt} = v \end{aligned}$$

We have,  $f(\alpha(0)) = f(0) = c$  since  $p \in S = f^{-1}(c)$ . Similarly,  $f(\alpha(1)) = c$ . And  $(f \circ \alpha)'(0) = \nabla f \circ \alpha(0) \dot{\alpha}(0) = \nabla f(p) \cdot \dot{\alpha}(0) = \|\nabla f(p)\| N(p) \cdot v$ . Similarly,  $(f \circ \alpha)'(1) = \|\nabla f(q)\| N(q) \cdot v$ . We have assumed that  $N(p) = N(q)$ . Then,  $f \circ \alpha$  is either increasing at both 0 and 1 OR decreasing at both 0 and 1.

Without Loss of Generality, Suppose that  $f \circ \alpha$  is increasing at either points. Then, there exists a sufficiently small  $\epsilon > 0$  such that  $f(\alpha(\epsilon)) > c$  and  $f(\alpha(1 - \epsilon)) < c$ . Since,  $f \circ \alpha$  must have a value greater than  $c$  immediately after 0

Figure 18.6: Construction of  $\alpha$ 

and should have a value less than  $c$  just before reaching 1 as the function is increasing at either points (and in some small neighbourhood of those points).

By Intermediate Value theorem, there exists  $t \in (0, 1)$  such that  $f \circ \alpha(t) = c$  since the composition of smooth functions  $f$  and  $\alpha$  is also smooth. But, it is clear from the construction that  $\alpha(t)$  doesn't belong to the surface  $S$ . And therefore,  $\alpha(t) \neq c \implies f(\alpha(t)) \neq c$  for any  $t \in (0, 1)$ . Thus by contradiction,  $N(p) \neq N(q)$ . And if  $N$  achieves  $v$  at  $p$ . Then it achieves  $-v$  at  $q$ . And since  $v \in S^n$  is arbitrary,  $N(S) = S^n$  and the spherical image is the entire  $n$  sphere OR the Gauss map is surjective.  $\square$

**Given a compact, connected oriented,  $n$ -surface  $S$ , the Gauss Map on  $S$  is surjective. In other words, the spherical image of such a surface is the unit  $n$  sphere  $S^n$  itself.**

Connectedness is not that critical (in my opinion). For a compact, orientated  $n$ -surface with multiple components, the above observation is valid for each connected component. And thus for any compact, oriented surface. Again,  $n$ -surfaces are always closed. Thus, the restriction practically reduces to boundedness of the  $n$ -surface.

## 18.7 Geodesics

We already know that our earth is not flat. Still, we feel like we move in straight lines. And our 'straight lines' are curved for an observer who is not on earth. Geodesics are straightlines on an  $n$ -surface  $S$ .

**vector field along  $\alpha$**  is function which assigns  $X(t)$  at  $\alpha(t)$  for each  $t \in I$ .

The definition of vector field doesn't allow you to assign multiple vectors at a point. But, vector field along  $\alpha$  allows you to assign vectors to points on a parametrised curve depending on the value of parameter  $t$ .

**function along  $\alpha$**  is function with the same domain  $I$  as  $\alpha$ .

**derivative of vector field  $\mathbf{X}$  along  $\alpha$**  is a vector field along  $\alpha$  given by  $\dot{\mathbf{X}}(t) = \left( \alpha(t), \frac{dX}{dt}(t) \right)$  where  $\mathbf{X}(t) = (\alpha(t), X(t))$ .

**velocity of  $\alpha$**  is a vector field along  $\alpha$  defined by  $\dot{\alpha}(t) = \left( \alpha(t), \frac{d\alpha}{dt}(t) \right)$ .

Suppose  $\alpha : I \rightarrow \mathbb{R}^2$  is defined by  $\alpha(t) = (3t, t^2)$ .

Then velocity of  $\alpha$  is  $\dot{\alpha}(t) = \left( \alpha(t), \frac{d\alpha}{dt}(t) \right) = (3t, t^2, 3, 2t)$ .

**speed of  $\alpha$**  is  $\|\dot{\alpha}(t)\|$ .

Speed of  $\alpha$  is  $\|\dot{\alpha}(t)\| = \sqrt{9 + 4t^2}$ .

**acceleration  $\alpha$**  is a vector field along  $\alpha$  defined by  $\ddot{\alpha}(t) = \left( \alpha(t), \frac{d^2\alpha}{dt^2}(t) \right)$ .

Acceleration of  $\alpha$  is  $\ddot{\alpha}(t) = (3t, t^2, 0, 2)$

### 18.7.1 Properties of differentiation

Let  $\mathbf{X}, \mathbf{Y}$  be smooth vector fields along parametrised curve  $\alpha : I \rightarrow \mathbb{R}^{n+1}$ .

- $(\mathbf{X} + \mathbf{Y}) = \dot{\mathbf{X}} + \dot{\mathbf{Y}}$

$$(\mathbf{X} + \mathbf{Y})(t) = (\alpha(t), X(t)) + (\alpha(t), Y(t)) = (\alpha(t), X(t) + Y(t))$$

$$\begin{aligned} (\mathbf{X} + \mathbf{Y})(t) &= \left( \alpha(t), \frac{d}{dt}X(t) + Y(t) \right) \\ &= \left( \alpha(t), \frac{d}{dt}X(t) \right) + \left( \alpha(t), \frac{d}{dt}Y(t) \right) \\ &= \dot{\mathbf{X}}(t) + \dot{\mathbf{Y}}(t) \end{aligned}$$

- $(f\dot{\mathbf{X}}) = f'\mathbf{X} + f\dot{\mathbf{X}}$

$$\begin{aligned} f\dot{\mathbf{X}}(t) &= f(t)(\alpha(t), X(t)) = (\alpha(t), f(t)X(t)) \\ (f\dot{\mathbf{X}})(t) &= \left( \alpha(t), \frac{d}{dt}f(t)X(t) \right) \\ &= \left( \alpha(t), f'(t)X(t) + f(t)\frac{d}{dt}X(t) \right) \\ &= (\alpha(t), f'(t)X(t)) + \left( \alpha(t), f(t)\frac{d}{dt}X(t) \right) \\ &= f'(t)(\alpha(t), X(t)) + f(t)\left( \alpha(t), \frac{d}{dt}X(t) \right) \\ &= f'\mathbf{X}(t) + f\dot{\mathbf{X}}(t) \end{aligned}$$



$$\bullet (\mathbf{X} \cdot \mathbf{Y})' = \dot{\mathbf{X}} \cdot \mathbf{Y} + \mathbf{X} \cdot \dot{\mathbf{Y}}$$

$$(\mathbf{X} \cdot \mathbf{Y}) = (\alpha(t), X(t)) \cdot (\alpha(t), Y(t)) = \sum_{k=1}^{n+1} X_k(t) Y_k(t)$$

$$\begin{aligned} (\mathbf{X} \cdot \mathbf{Y})' &= \frac{d}{dt} \sum_{k=1}^{n+1} X_k(t) Y_k(t) \\ &= \sum_{k=1}^{n+1} \frac{d}{dt} X_k(t) Y_k(t) \\ &= \sum_{k=1}^{n+1} X'_k(t) Y_k(t) + X_k(t) Y'_k(t) \end{aligned}$$

$$\dot{\mathbf{X}}(t) \cdot \mathbf{Y}(t) = \left( \alpha(t), \frac{d}{dt} X(t) \right) \cdot (\alpha(t), Y(t)) = \sum_{k=1}^{n+1} X'_k(t) Y_k(t)$$

$$\mathbf{X}(t) \cdot \dot{\mathbf{Y}}(t) = (\alpha(t), X(t)) \cdot \left( \alpha(t), \frac{d}{dt} Y(t) \right) = \sum_{k=1}^{n+1} X_k(t) Y'_k(t)$$

**Definitions 18.7.1** (geodesic). Let  $S$  be an  $n$ -surface. A Geodesic on  $S$  is a parametrised curve  $\alpha : I \rightarrow S$  whose acceleration is orthogonal to  $S$  everywhere.

$$\ddot{\alpha}(t) \in S_{\alpha(t)}^\perp \quad (18.11)$$

### 18.7.2 An illustrative example

We know that  $S^1$  given by  $x_1^2 + x_2^2 = 1$  is a 1-surface in  $\mathbb{R}^2$ . Consider the cylinder  $C$  over  $S^1$ ,  $x_1^2 + x_2^2 = 1$ . Clearly,  $C$  is a 2 surface in  $\mathbb{R}^3$ . Also,  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  defined by  $f(x_1, x_2, x_3) = x_1^2 + x_2^2$  is a smooth function such that  $C = f^{-1}(1)$  and  $\nabla f(x_1, x_2, x_3) = (x_1, x_2, x_3, 2x_1, 2x_2, 0) \neq \mathbf{0}$ .

Clearly, every vector orthogonal to the surface in a scalar multiple of  $\nabla f$  at that point. Therefore, vectors orthogonal to the surface  $C$  at  $\alpha(t)$  is of the form  $(x_1, x_2, x_3, kx_1, kx_2, 0)$  where  $k \in \mathbb{R}$ .

If there exists a geodesic  $\alpha$  in  $C$ , then  $\ddot{\alpha}(t) \in S_{\alpha(t)}^\perp$ . That is,  $\ddot{\alpha}(t) = (x_1, x_2, x_3, kx_1, kx_2, 0)$ . Thus, we need component functions  $x_1(t), x_2(t), x_3(t)$  satisfying

$$\frac{d^2}{dt^2} x_1(t) = kx_1(t) \quad (18.12)$$

$$\frac{d^2}{dt^2} x_2(t) = kx_2(t) \quad (18.13)$$

$$\frac{d^2}{dt^2} x_3(t) = 0 \quad (18.14)$$

We have,  $\frac{d^2}{dt^2} \cos t = -\cos t$  and  $\frac{d^2}{dt^2} \sin t = -\sin t$ . Thus, the parametrised curve  $\alpha : I \rightarrow \mathbb{R}^3$  defined by  $\alpha(t) = (\cos t, \sin t, t)$  is a geodesic in  $C$  since  $\ddot{\alpha}(t) = (\cos t, \sin t, t, -\cos t, -\sin t, 0)$ .

### 18.7.3 Maximal Geodesic

The conditions  $\alpha(0) = p$ ,  $\dot{\alpha}(0) = \mathbf{v}$  says that parametric curves are unique except for linear transformations on the parameter. That is, Suppose there exists another parametrised curve  $\beta$  in  $S$  through  $p$  with initial velocity  $\mathbf{v}$  with  $\beta(t_0) = p$  and  $\dot{\beta}(t_0) = \mathbf{v}$ . Then, there exists a real number  $\kappa$  such that  $\alpha(t) = \beta(\kappa t + t_0)$ .

In other words, both  $\alpha$  and  $\beta$  passes through the same points and the vectors assigned at each point is the same. And the difference between such two parametrised curves doesn't have any impact on the properties we are interested in.

The condition, if  $\beta : \tilde{I} \rightarrow S$  is any other geodesic in  $S$  with  $\beta(0) = p$  and  $\dot{\beta}(0) = \mathbf{v}$ , then  $\tilde{I} \subset I$  and  $\beta(t) = \alpha(t)$ ,  $\forall t \in \tilde{I}$ . is another way of saying that  $\alpha$  is maximal and uniquely defined.

In essence, the following theorem says that if you are standing on an  $n$ -surface  $S$  at a point, say  $p$ . You can move on that surface in straightline from  $p$ , with any initial velocity,  $\mathbf{v}$ . Note that the velocity allows you to choose both direction and speed.

**Theorem 18.7.1** (maximal geodesic). *Let  $S$  be an  $n$ -surface in  $\mathbb{R}^{n+1}$ . Let  $p \in S$  and  $\mathbf{v} \in S_p$ . Then there exists a unique, maximal geodesic  $\alpha : I \rightarrow S$  in  $S$  through  $p$  with initial velocity  $\mathbf{v}$ .*

*Proof.* Let  $S$  be an  $n$ -surface, then there exists a smooth function  $f : U \rightarrow \mathbb{R}$  where  $U$  is an open subset of  $\mathbb{R}^{n+1}$ . And every points on that surface is regular with respect to  $f$ . That is,  $\nabla f(p) \neq 0$ ,  $\forall p \in S$ . WLOG assume that every points in  $U$  is regular.

Define  $N = \frac{\nabla f}{\|\nabla f\|}$ . A parametrised curve  $\alpha$  in  $S$  is a geodesic if it satisfies  $\ddot{\alpha}(t) \in S_{\alpha(t)}^\perp$ . Therefore,

$$\ddot{\alpha}(t) = g(t)\mathbf{N}(\alpha(t)) \quad (18.15)$$

Why don't we write  $\ddot{\alpha}(t) = \kappa\mathbf{N}(\alpha(t))$  ? The vectors  $\kappa\mathbf{N}(\alpha(t))$  are orthogonal to  $S$ . However, the converse is not true. For a geodesic it is not necessary that  $\ddot{\alpha}(t) = \kappa\mathbf{N}(\alpha(t))$ . The acceleration could be any vector in  $S_{\alpha(t)}^\perp$ . Since  $S_{\alpha(t)}^\perp$  is spanned by  $\mathbf{N}(\alpha(t))$ , at each point on  $\alpha(t)$  the scalars may be different. We overcome this with the help of a real valued function  $g : I \rightarrow \mathbb{R}$  along  $\alpha$ .

We have  $\ddot{\alpha} = g(\mathbf{N} \circ \alpha)$ . Thus  $\ddot{\alpha} \cdot (\mathbf{N} \circ \alpha) = g\|\mathbf{N} \circ \alpha\|^2 = g$  since orientation  $\mathbf{N}$  assigns directions (vectors of unit length) to each point of that surface.

$$[\dot{\alpha} \cdot (\mathbf{N} \circ \alpha)]' = \ddot{\alpha} \cdot (\mathbf{N} \circ \alpha) + \dot{\alpha} \cdot (\mathbf{N} \dot{\circ} \alpha) \text{ since } (\mathbf{X} \cdot \mathbf{Y})' = \dot{\mathbf{X}} \cdot \mathbf{Y} + \mathbf{X} \cdot \dot{\mathbf{Y}}.$$

Thus,  $\ddot{\alpha} \cdot (\mathbf{N} \circ \alpha) = [\dot{\alpha} \cdot (\mathbf{N} \circ \alpha)]' - \dot{\alpha} \cdot (\mathbf{N} \dot{\circ} \alpha) = -\dot{\alpha} \cdot (\mathbf{N} \dot{\circ} \alpha)$  since  $\dot{\alpha} \cdot (\mathbf{N} \circ \alpha) = 0$  as  $\dot{\alpha} \in S_{\alpha(t)}$ ,  $\mathbf{N} \circ \alpha \in S_{\alpha(t)}^\perp$  and  $\dot{\alpha} \perp (\mathbf{N} \circ \alpha)$ . That is, velocity vectors always belongs to the tangent space and  $(\mathbf{N} \circ \alpha)$  is an orientation which is orthogonal to all the tangent vectors.

Substituting the value of  $g$  in equation(18.15). We get

$$\ddot{\alpha} + [\dot{\alpha} \cdot (\mathbf{N} \dot{\circ} \alpha)] (\mathbf{N} \circ \alpha) = \mathbf{0} \quad (18.16)$$

We have,

$$\frac{dN_1}{dt} = \frac{\partial N_1}{\partial x_1} \frac{dx_1}{dt} + \frac{\partial N_1}{\partial x_2} \frac{dx_2}{dt} + \cdots + \frac{\partial N_1}{\partial x_{n+1}} \frac{dx_{n+1}}{dt} = \sum_{k=1}^{n+1} \frac{\partial N_1}{\partial x_k} \frac{dx_k}{dt} \quad (18.17)$$

Thus,

$$(\mathbf{N} \circ \alpha) = \left( \alpha(t), \sum_{k=1}^{n+1} \frac{\partial N_1}{\partial x_k} \frac{dx_k}{dt}, \sum_{k=1}^{n+1} \frac{\partial N_2}{\partial x_k} \frac{dx_k}{dt}, \dots, \sum_{k=1}^{n+1} \frac{\partial N_{n+1}}{\partial x_k} \frac{dx_k}{dt} \right) \quad (18.18)$$

Equating the components on either sides of the equation(18.16), we get a system of  $n + 1$  second order differential equations,

$$\frac{d^2}{dt^2} x_i(t) + \left[ \sum_{j,k=1}^{n+1} \frac{\partial N_j}{\partial x_k} \frac{dx_k}{dt} \frac{dx_j}{dt} \right] N_i \circ \alpha = 0 \quad (18.19)$$

By the existence theorem for solution of such system of second order differential equations, there exists an open interval  $I$  containing 0 and  $(n + 1)$  functions  $x_k : I \rightarrow \mathbb{R}$  satisfying the above system. Then  $\beta : I \rightarrow S$  defined by  $\beta(t) = (x_1(t), x_2(t), \dots, x_{n+1}(t))$

By the uniqueness theorem for solution of such system of second order differential equations, if there exists another open interval  $\tilde{I}$  containing 0 and  $\beta_1 : I_1 \rightarrow U$ . Then  $\beta(t) = \beta(t)$  for every  $t \in I \cap I_1$ .

Let  $\beta_1, \beta_2, \dots$  be geodesics through  $p$  with inintial velocity  $\mathbf{v}$  with parameter domain  $I_1, I_2, \dots$ . Let  $I = \cup_k I_k$  and  $\alpha : I \rightarrow U$  defined by  $\alpha(t) = \beta_k(t)$ ,  $t \in I_k$ . Clearly,  $\alpha$  is unique and maximal.

Now it is enough to prove that  $\alpha$  is parametrised curve in  $S$ . We have  $(f \circ \alpha)'(t) = \nabla f(\alpha(t)) \cdot \dot{\alpha}(t) = \|\nabla f(\alpha(t))\| \mathbf{N}(\alpha(t)) \cdot \dot{\alpha}(t) = 0$  since  $\dot{\alpha} \perp (\mathbf{N} \circ \alpha)$ . Therefore,  $f \circ \alpha : I \rightarrow \mathbf{R}$  is constant. Also,  $f(\alpha(0)) = f(p) = c$  since  $p \in S$  and  $S = f^{-1}(c)$ . Thus,  $f \circ \alpha = c \implies \alpha \subset f^{-1}(c)$ . Thus,  $\alpha$  is a parametrised curve in  $n$ -surface  $S$ .  $\square$

#### 18.7.4 Properites of Geodesics

1.  $\dot{\alpha} \perp \ddot{\alpha}$  since  $\dot{\alpha} \in S_{\alpha(t)}$  and  $\ddot{\alpha} \in S_{\alpha(t)}^\perp$

In other words, the acceleration is orthogonal to the velocity vector.

2. Constant Speed,  $\|\mathbf{v}\|$

We have,  $(\dot{\alpha} \cdot \dot{\alpha})' = \ddot{\alpha} \cdot \dot{\alpha} + \dot{\alpha} \cdot \ddot{\alpha} = 2\dot{\alpha} \cdot \ddot{\alpha} = 0$

Clearly,  $\frac{d}{dt} \|\dot{\alpha}\|^2 = (\dot{\alpha} \cdot \dot{\alpha})' = 0$ . Therefore,  $\alpha$  has constant speed.

Remark : Geodesics in cylinder over  $S^1$  are horizontal circles, vertical lines, helix or a constant.

## 18.8 Parallel Transport

**covariant derivative** Covariant derivative of  $\mathbf{X}$  is the vector field  $\mathbf{X}'$  tangent to  $S$  along  $\alpha$  given by  $\mathbf{X}'(t) = \dot{\mathbf{X}}(t) - [\dot{\mathbf{X}}(t) \cdot \mathbf{N}(\alpha(t))]\mathbf{N}(\alpha(t))$ . And it is independent of the choice of the orientation.

**covariant acceleration** Let  $\alpha : I \rightarrow S$  be a parametrised curve in  $S$ . Then covariant acceleration of  $\alpha$  is  $(\dot{\alpha})' = \ddot{\alpha} - [\ddot{\alpha} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha)$  along  $\alpha$ .

### 18.8.1 Properties of Covariant Derivative

Let  $\mathbf{X}, \mathbf{Y}$  be smooth vector fields tangent to  $S$  along  $\alpha$ . Then,  $\mathbf{X} \cdot (\mathbf{N} \circ \alpha) = \mathbf{0}$ .

$$1. (\mathbf{X} + \mathbf{Y})' = \mathbf{X}' + \mathbf{Y}'$$

$$\begin{aligned} (\mathbf{X} + \mathbf{Y})' &= (\dot{\mathbf{X}} + \dot{\mathbf{Y}}) - [(\dot{\mathbf{X}} + \dot{\mathbf{Y}}) \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= (\dot{\mathbf{X}} + \dot{\mathbf{Y}}) - [(\dot{\mathbf{X}} \cdot (\mathbf{N} \circ \alpha)) + (\dot{\mathbf{Y}} \cdot (\mathbf{N} \circ \alpha))](\mathbf{N} \circ \alpha) \\ &= \left( \dot{\mathbf{X}} - [\dot{\mathbf{X}} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \right) + \left( \dot{\mathbf{Y}} - [\dot{\mathbf{Y}} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \right) \\ &= \mathbf{X}' + \mathbf{Y}' \end{aligned}$$

$$2. (f\mathbf{X})' = f'\mathbf{X} + f\mathbf{X}'$$

$$\begin{aligned} (f\mathbf{X})' &= (f\dot{\mathbf{X}}) - [(f\dot{\mathbf{X}}) \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= (f'\mathbf{X} + f\dot{\mathbf{X}}) - [(f'\mathbf{X} + f\dot{\mathbf{X}}) \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= f'(\mathbf{X} - [\mathbf{X} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha)) + f(\dot{\mathbf{X}} - [\dot{\mathbf{X}} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha)) \\ &= f'\mathbf{X} + f\dot{\mathbf{X}} \text{ since } \mathbf{X} \cdot (\mathbf{N} \circ \alpha) = \mathbf{0} \end{aligned}$$

$$3. (\mathbf{X} \cdot \mathbf{Y})' = \mathbf{X}' \cdot \mathbf{Y} + \mathbf{X} \cdot \mathbf{Y}'$$

$$\begin{aligned} (\mathbf{X} \cdot \mathbf{Y})' &= \dot{\mathbf{X}} \cdot \mathbf{Y} + \mathbf{X} \cdot \dot{\mathbf{Y}} \\ &= \dot{\mathbf{X}} \cdot \mathbf{Y} - [\dot{\mathbf{X}} \cdot (\mathbf{N} \circ \alpha)]\mathbf{0} + \mathbf{X} \cdot \dot{\mathbf{Y}} - [\dot{\mathbf{Y}} \cdot (\mathbf{N} \circ \alpha)]\mathbf{0} \end{aligned}$$

Substituting  $\mathbf{X} \cdot (\mathbf{N} \circ \alpha) = \mathbf{0}$  and  $\mathbf{Y} \cdot (\mathbf{N} \circ \alpha) = \mathbf{0}$ , we get

$$\begin{aligned} &= \dot{\mathbf{X}} \cdot \mathbf{Y} - [\dot{\mathbf{X}} \cdot (\mathbf{N} \circ \alpha)][(\mathbf{N} \circ \alpha) \cdot \mathbf{Y}] + \mathbf{X} \cdot \dot{\mathbf{Y}} - [\dot{\mathbf{Y}} \cdot (\mathbf{N} \circ \alpha)][\mathbf{X} \cdot (\mathbf{N} \circ \alpha)] \\ &= (\dot{\mathbf{X}} - [\dot{\mathbf{X}} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha)) \cdot \mathbf{Y} + \mathbf{X} \cdot (\dot{\mathbf{Y}} - [\dot{\mathbf{Y}} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha)) \\ &= \mathbf{X}' \cdot \mathbf{Y} + \mathbf{X} \cdot \mathbf{Y}' \end{aligned}$$

### 18.8.2 Parallelism

We have seen earlier that, lines that look straight on a surface (geodesics) is not necessarily straight from outside. The same way, two vectors that look *parallel on a surface* (Levi-Civita parallel) is not necessarily parallel.

**Euclidean parallel**  $(p, v)$  and  $(q, w)$  are Euclidean parallel if  $v = w$ .

Let  $\alpha : I \rightarrow \mathbb{R}^{n+1}$  be a parametrised curve. A vector field  $\mathbf{X}$  is Euclidean parallel along  $\alpha$  if the associated function  $X$  is constant say,  $v$ . Then,

$$\dot{\mathbf{X}} = \left( \alpha(t), \frac{d}{dt} X(t) \right) = \left( \alpha(t), \frac{dv}{dt} \right) = (\alpha(t), 0) = \mathbf{0}$$

**Levi-Civita parallel** A vector field  $\mathbf{X}$  tangent to the surface  $S$  along  $\alpha$  is (Levi-Civita) parallel if  $\mathbf{X}$  is a constant vector field along  $\alpha$  with respect to the surface  $S$ . That is,  $\mathbf{X}' = \mathbf{0}$ .

### Properties of Levi-Civita parallel

Applying the properties of covariant derivative, we get

1. If  $\mathbf{X}$  is parallel along  $\alpha$ , then  $\mathbf{X}$  has constant length. That is,  $\|\mathbf{X}\|' = 0$ .

$$\frac{d}{dt} \|\mathbf{X}\|^2 = \frac{d}{dt} \mathbf{X} \cdot \mathbf{X} = 2\mathbf{X}' \cdot \mathbf{X} = 2(\mathbf{0} \cdot \mathbf{X}) = 0$$

2. If  $\mathbf{X}, \mathbf{Y}$  are parallel along  $\alpha$ , then  $\mathbf{X} \cdot \mathbf{Y}$  is constant along  $\alpha$ .

$$(\mathbf{X} \cdot \mathbf{Y})' = \mathbf{X}' \cdot \mathbf{Y} + \mathbf{X} \cdot \mathbf{Y}' = \mathbf{0} \cdot \mathbf{Y} + \mathbf{X} \cdot \mathbf{0} = 0$$

3. If  $\mathbf{X}, \mathbf{Y}$  are parallel along  $\alpha$ , then angle between them is constant along  $\alpha$ .

$$\theta = \cos^{-1} \frac{\mathbf{X} \cdot \mathbf{Y}}{\|\mathbf{X}\| \|\mathbf{Y}\|} = \cos^{-1} \kappa, \text{ since } \|\mathbf{X}\|, \|\mathbf{Y}\|, \mathbf{X} \cdot \mathbf{Y} \text{ are constant}$$

4. If  $\mathbf{X}, \mathbf{Y}$  are parallel along  $\alpha$ , then  $\mathbf{X} + \mathbf{Y}, c\mathbf{X}$  are parallel along  $\alpha$ .

$$(\mathbf{X} + \mathbf{Y})' = \mathbf{X}' + \mathbf{Y}' = \mathbf{0} + \mathbf{0} = \mathbf{0}$$

$$(c\mathbf{X})' = c\mathbf{X}' = c\mathbf{0} = \mathbf{0}$$

5. The velocity vector field along a parametrised curve  $\alpha$  in  $S$  is parallel if and only if  $\alpha$  is a geodesic.

$$(\dot{\alpha})' = \ddot{\alpha} - [\ddot{\alpha} \cdot (\mathbf{N} \circ \alpha)] (\mathbf{N} \circ \alpha) = \mathbf{0} \iff \ddot{\alpha} \in S_{\alpha(t)}^\perp$$

Parametrised curve  $\alpha$  is geodesic in  $S$  if and only if covariant acceleration  $(\dot{\alpha})'$  is zero along  $\alpha$  since  $\ddot{\alpha} \in S_{\alpha(t)}^\perp$  and  $[\ddot{\alpha} \cdot (\mathbf{N} \circ \alpha)] (\mathbf{N} \circ \alpha) = \ddot{\alpha}$ .

We have, the notation  $f'$  and  $\mathbf{X}'$ . You should keep a note of the fundamental differences.  $f'$  refers to the derivative of a function along  $\alpha$  with respect to the parameter  $t$ . And  $\mathbf{X}'$  refers to the covariant derivative of a vector field tangent to  $S$  along  $\alpha$ . You should always check whether it is real value OR vector at a point to understand what they really mean.

For example,  $\|\mathbf{X}\|'$  is a derivative of a real valued function (derivative of the length of the vectors assigned at different points of a parametrised curve). And it has nothing to do with the covariant derivative  $\mathbf{X}'$ .

**Theorem 18.8.1.** *Let  $S$  be an  $n$ -surface and  $\alpha : I \rightarrow S$  be a parametrised curve in  $S$ . Let  $t_0 \in I$  and  $\mathbf{v} \in S_{\alpha(t_0)}$ . Then there exists a unique vector field  $\mathbf{V}$  tangent to  $S$  along  $\alpha$  which is parallel and has  $\mathbf{V}(t_0) = \mathbf{v}$ .*

*Proof.* Let  $S$  be an  $n$ -surface with orientation  $\mathbf{N}$ . Suppose  $\mathbf{V}$  is a vector field tangent to  $S$  along  $\alpha$  and is parallel. That is,  $\mathbf{V}(t) \in S_{\alpha(t)}$  and  $\mathbf{V}' = \mathbf{0}$ .

$$\begin{aligned}\mathbf{V}' &= \dot{\mathbf{V}} - \left[ \dot{\mathbf{V}} \cdot (\mathbf{N} \circ \alpha) \right] (\mathbf{N} \circ \alpha) \\ &= \dot{\mathbf{V}} - \left[ (\mathbf{V} \cdot (\mathbf{N} \circ \alpha))' - \mathbf{V} \cdot (\mathbf{N} \dot{\circ} \alpha) \right] (\mathbf{N} \circ \alpha) \\ &= \dot{\mathbf{V}} + \left[ \mathbf{V} \cdot (\mathbf{N} \dot{\circ} \alpha) \right] (\mathbf{N} \circ \alpha) \text{ since } \mathbf{V} \perp \mathbf{N} \\ \dot{\mathbf{V}} + [\mathbf{V} \cdot (\mathbf{N} \dot{\circ} \alpha)] (\mathbf{N} \circ \alpha) &= \mathbf{0}\end{aligned}\tag{18.20}$$

Equating the components on either sides of the equation(18.20), we get the following system of  $n + 1$  first order differential equations,

$$\frac{dV_i}{dt} + \sum_{j=1}^{n+1} [V_j (\mathbf{N}_j \circ \alpha)'] (\mathbf{N}_i \circ \alpha) = 0, \quad \forall i \tag{18.21}$$

By the existence and uniqueness theorem for first order differential equations, there exists  $V_1(t), V_2(t), \dots, V_{n+1}(t)$  satisfying the sytem of equations with initial condition  $\mathbf{V}(t_0) = (\alpha(t_0), V_1(t_0), V_2(t_0), \dots, V_{n+1}(t_0)) = \mathbf{v}$ .

It remains to prove that  $\mathbf{V}$  is tangent to  $S$  along  $\alpha$ . By taking dot product with  $\mathbf{N} \circ \alpha$  on either sides of equation(18.20), we get

$$(\mathbf{V} \cdot \mathbf{N} \circ \alpha)' = \dot{\mathbf{V}} \cdot (\mathbf{N} \circ \alpha) + \mathbf{V} \cdot (\mathbf{N} \dot{\circ} \alpha) = \mathbf{0} \tag{18.22}$$

Thus,  $\mathbf{V} \cdot (\mathbf{N} \circ \alpha) = \kappa$  is constant along  $\alpha$ . However,  $\mathbf{V}(t_0) \cdot (\mathbf{N} \circ \alpha)(t_0) = \mathbf{v} \cdot \mathbf{N}(\alpha(t_0)) = 0$  since  $\mathbf{v} \in S_{\alpha(t_0)}$  is a tangent vector and  $\mathbf{v} \perp \mathbf{N}$ . Therefore,  $\mathbf{V} \cdot (\mathbf{N} \circ \alpha) = 0$ . And since  $\mathbf{V}$  satisfies equation(18.20), this vector field is tangent to  $S$  along  $\alpha$  and is parallel.  $\square$

**Corollary 18.8.1.1.** *Let  $S$  be a 2-surface in  $\mathbb{R}^3$  and let  $\alpha : I \rightarrow S$  be a parametrised curve in  $S$  with  $\dot{\alpha} \neq \mathbf{0}$ . Then the vector field  $\mathbf{X}$  tangent to  $S$  along  $\alpha$  is parallel along  $\alpha$  if and only if both  $\|\mathbf{X}\|$  and the angle between  $\mathbf{X}$  and  $\dot{\alpha}$  are constant along  $\alpha$ .*

*Proof. Sufficient Part :* Let  $\mathbf{X}$  be a tangent vector field tangent to  $S$  along  $\alpha$ . And  $\alpha$  is a geodesic in  $S$ . Then  $\mathbf{X}$  is parallel along  $\alpha$ . Also, we have  $\dot{\alpha}$  is also parallel along  $\alpha$  since  $(\dot{\alpha})' = \mathbf{0}$ . Thus  $\|\mathbf{X}\|$  is constant and the angle between  $\mathbf{X}$  and  $\dot{\alpha}$  is constant by properties Levi-Civita parallelism.

*Necessary Part :* Let  $\mathbf{X}$  be a vector field tangent to  $S$  along  $\alpha$  and  $\alpha$  is a geodesic in  $S$ . Then  $\|\dot{\alpha}\|$  is constant. Suppose  $\|\mathbf{X}\|$  and the angle  $\theta$  between  $\mathbf{X}$  and  $\dot{\alpha}$  are constant. Since  $S$  is a 2-surface, the tangent space  $S_{\alpha(t)}$  is spanned by two tangent vectors  $\mathbf{v}$  and  $\dot{\alpha}(t)$  such that  $\mathbf{v} \perp \dot{\alpha}(t)$  and  $\|\mathbf{v}\| = 1$ .

Let  $\mathbf{V}$  be the unique vector field tangent to  $S$  along  $\alpha$  which is parallel,  $\mathbf{V} \cdot \dot{\alpha} = 0$  and  $\|\mathbf{V}\| = 1$ . Then, any vector field  $\mathbf{X}$  tangent to  $S$  along  $\alpha$  can be written as linear combination of  $\mathbf{V}$  and  $\dot{\alpha}$ . That is,  $\mathbf{X} = f\dot{\alpha} + g\mathbf{V}$

$$\cos \theta = \frac{\mathbf{X} \cdot \dot{\alpha}}{\|\mathbf{X}\| \|\dot{\alpha}\|} = \frac{(f\dot{\alpha} + g\mathbf{V}) \cdot \dot{\alpha}}{\|\mathbf{X}\| \|\dot{\alpha}\|} = \frac{f\|\dot{\alpha}\|}{\|\mathbf{X}\|} \implies f \text{ is constant}$$

$$\|\mathbf{X}\|^2 = (f\dot{\alpha} + g\mathbf{V}) \cdot (f\dot{\alpha} + g\mathbf{V}) = f^2\|\dot{\alpha}\|^2 + g^2 \implies g \text{ is constant}$$

Since  $\mathbf{V}$  and  $\dot{\alpha}$  are parallel along  $\alpha$ , their linear combination  $\mathbf{X}$  is also parallel along  $\alpha$  by linearity property (#4) of Levi-Civita parallelism.  $\square$

### 18.8.3 Transporting tangent vectors using parallelism

Suppose you are standing on an  $n$ -surface with your hand stretched out as in a pledge. Then you can move on that surface along any smooth road (without sharp turns, bumps or potholes) from one point to another keeping your hand in the same position. This is what a parallel transport does to tangent vectors. And we can calculate the direction you will be pointing, at each point of your journey.

**Definitions 18.8.1** (Parallel transport). Let  $S$  be an  $n$ -surface. Let  $p, q \in S$ , and let  $\alpha : [a, b] \rightarrow S$  be a (smooth) parametrised curve in  $S$  from  $p$  to  $q$ . **Parallel transport** is the function  $P_\alpha : S_p \rightarrow S_q$  defined by  $P_\alpha(\mathbf{v}) = \mathbf{V}(b)$  where  $\mathbf{v} \in S_p$  and  $\mathbf{V}$  is the unique vector field tangent to  $S$  along  $\alpha$  which is parallel and  $\mathbf{V}(a) = \mathbf{v}$ .

The following theorem says that, *Parallel transports along piecewise smooth parametrised curves are vector space isomorphisms preserving dot products.*

**Theorem 18.8.2.** Let  $S$  be an  $n$ -surface in  $\mathbb{R}^{n+1}$ . Let  $p, q \in S$  and let  $\alpha$  be a piecewise smooth parametrised curve from  $p$  to  $q$ . Then parallel transport  $P_\alpha : S_p \rightarrow S_q$  along  $\alpha$  is a vector space isomorphism which preserves dot products.

*Proof.* Let  $S$  be an  $n$ -surface and  $p, q \in S$ . Let  $\alpha : [a, b] \rightarrow S$  be a piecewise smooth parametrised curve from  $p$  to  $q$ . Let  $P_\alpha : S_p \rightarrow S_q$  be a parallel transport and  $\mathbf{v}, \mathbf{w} \in S_p$ . Clearly,  $\mathbf{v} + \mathbf{w}, c\mathbf{v} \in S_p$

To prove that  $P_\alpha$  is a vector space isomorphism preserving dot products, we need to prove the following

1.  $P_\alpha$  is linear,  $P_\alpha(\mathbf{v} + \mathbf{w}) = P_\alpha(\mathbf{v}) + P_\alpha(\mathbf{w})$  and  $P_\alpha(c\mathbf{v}) = cP_\alpha(\mathbf{v})$

$$P_\alpha(\mathbf{v} + \mathbf{w}) = (\mathbf{V} + \mathbf{W})(b) = \mathbf{V}(b) + \mathbf{W}(b) = P_\alpha(\mathbf{v}) + P_\alpha(\mathbf{w})$$

$$P_\alpha(c\mathbf{v}) = (c\mathbf{V})(b) = c\mathbf{V}(b) = cP_\alpha(\mathbf{v})$$

2.  $P_\alpha$  is bijective

$$\|P_\alpha(\mathbf{v})\| = \|\mathbf{V}(b)\| = 0 \implies \|\mathbf{v}\| = 0 \implies \ker(P_\alpha) = \{\mathbf{0}\}$$

Thus  $P_\alpha$  is a linear map from  $n$ -dimensional vector space  $S_p$  into another  $n$ -dimensional vector space  $S_q$ . Thus,  $P_\alpha$  is onto.

3.  $P_\alpha$  preserves dot products,  $P_\alpha(\mathbf{v}) \cdot P_\alpha(\mathbf{w}) = \mathbf{v} \cdot \mathbf{w}$ ,  $\forall \mathbf{v}, \mathbf{w} \in S_p$   
 We have,  $\mathbf{V}$  and  $\mathbf{W}$  are parallel along  $\alpha$ . Then  $\mathbf{V} \cdot \mathbf{W}$  is constant.  
 That is,  $\mathbf{V}(t) \cdot \mathbf{W}(t) = \kappa$  for every  $t \in [a, b]$ .

$$P_\alpha(\mathbf{v}) \cdot P_\alpha(\mathbf{w}) = \mathbf{V}(b) \cdot \mathbf{W}(b) = \kappa = \mathbf{V}(a) \cdot \mathbf{W}(a) = \mathbf{v} \cdot \mathbf{w}$$

□

## 18.9 The Weingarten Map

The Weingarten map  $L_p$  gives information about the shape of a surface  $S$  at a point  $p \in S$ .

$\nabla_v f$  is the derivative of a function  $f : U \rightarrow \mathbb{R}$  with respect to a vector tangent to  $S$ , say  $\mathbf{v} = (p, v)$  where  $U$  is an open subset of  $\mathbb{R}^{n+1}$ ,  $p \in U$  and  $\alpha : I \rightarrow S$  is any parametrised curve in  $S$  with  $\dot{\alpha}(t_0) = \mathbf{v}$  and  $\alpha(t_0) = p$ . And  $\nabla_v f$  is given by,

$$\nabla_v f = (f \circ \alpha)'(t_0) = \nabla f(\alpha(t_0)) \cdot \dot{\alpha}(t_0) = \nabla f(p) \cdot \mathbf{v}$$

For example : if  $f(x_1, x_2) = x_1^2 - x_2^2$  and  $\mathbf{v} = (1, 1, \cos \theta, \sin \theta)$  Then,  $p = (1, 1)$ ,  $v = (\cos \theta, \sin \theta)$  and  $\nabla f = (x_1, x_2, 2x_1, -2x_2)$ . Therefore,  $\nabla_v f = \nabla f(p) \cdot \mathbf{v} = (1, 1, 2, -2) \cdot (1, 1, \cos \theta, \sin \theta) = 2 \cos \theta - 2 \sin \theta$ .

Note :  $\nabla_v f$  is independent of the choice of  $\alpha$ .

$\nabla_v \mathbf{X}$  is the derivative of a smooth vector field  $\mathbf{X}$  on open subset  $U$  of  $\mathbb{R}^{n+1}$  with respect to  $\mathbf{v} \in S_{\alpha(t)}$ . And  $\nabla_v \mathbf{X}$  is given by

$$\nabla_v \mathbf{X} = (\alpha(t), \nabla_v X_1, \nabla_v X_2, \dots, \nabla_v X_{n+1})$$

For example, if  $\mathbf{X}(x_1, x_2) = (x_1, x_2, x_1 x_2, x_2^2)$  and  $\mathbf{v} = (1, 0, 0, 1)$ . Then  $p = (1, 0)$  and the component functions of the associated function  $X$  are  $X_1(x_1, x_2) = x_1 x_2$  and  $X_2(x_1, x_2) = x_2^2$ .

We have,  $\nabla X_1(x_1, x_2) = (x_1, x_2, x_2, x_1)$ ,  $\nabla X_2(x_1, x_2) = (x_1, x_2, 0, 2x_2)$ . Thus,  $\nabla_v X_1 = \nabla X_1(1, 0) \cdot \mathbf{v} = (1, 0, 0, 1) \cdot (1, 0, 0, 1) = 1$  and  $\nabla_v X_2 = \nabla X_2(1, 0) \cdot \mathbf{v} = (1, 0, 0, 0) \cdot (1, 0, 0, 1) = 0$ . Therefore,  $\nabla_v \mathbf{X} = (1, 0, \nabla_v X_1, \nabla_v X_2) = (1, 0, 1, 0)$

$D_v \mathbf{X}$  is the covariant derivative of the smooth vector field  $\mathbf{X}$  with respect to  $\mathbf{v} \in S_{\alpha(t)}$ , where  $\mathbf{X}$  is tangent to  $S$  along  $\alpha$ . And  $D_v \mathbf{X}$  is given by,

$$D_v \mathbf{X} = \nabla_v \mathbf{X} - [\nabla_v \mathbf{X} \cdot (\mathbf{N} \circ \alpha)] (\mathbf{N} \circ \alpha)$$

It is the component of the derivative  $\nabla_v \mathbf{X}$  in the tangent space  $S_p$ .

### 18.9.1 Properties of differentiation, $\nabla_v f$

We have,  $\nabla_v f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  is a linear map.



$$1. \nabla_{v+w}f = \nabla_vf + \nabla_wf$$

$$\nabla_{v+w}f = \nabla f(p) \cdot (\mathbf{v} + \mathbf{w}) = \nabla f(p) \cdot \mathbf{v} + \nabla f(p) \cdot \mathbf{w}$$

$$2. \nabla_{cv}f = c\nabla_vf$$

$$\nabla_{cv}f = \nabla f(p) \cdot c\mathbf{v} = c(\nabla f(p) \cdot \mathbf{v}) = c\nabla_vf$$

Clearly, we have  $\nabla_v(f+g) = \nabla_vf + \nabla_vg$  and  $\nabla_vcf = c\nabla_vf$ .

### 18.9.2 Properties of differentiation, $\nabla_v\mathbf{X}$

$$1. \nabla_v(\mathbf{X} + \mathbf{Y}) = \nabla_v\mathbf{X} + \nabla_v\mathbf{Y}$$

$$\begin{aligned} \nabla_v(\mathbf{X} + \mathbf{Y}) &= (\alpha(t), \nabla_v(X_1 + Y_1), \nabla_v(X_2 + Y_2), \dots, \nabla_v(X_{n+1} + Y_{n+1})) \\ &= (\alpha(t), \nabla_vX_1 + \nabla_vY_1, \nabla_vX_2 + \nabla_vY_2, \dots, \nabla_vX_{n+1} + \nabla_vY_{n+1}) \\ &= (\alpha(t), \nabla_vX_1, \nabla_vX_2, \dots, \nabla_vX_{n+1}) + (\alpha(t), \nabla_vY_1, \nabla_vY_2, \dots, \nabla_vY_{n+1}) \\ &= \nabla_v\mathbf{X} + \nabla_v\mathbf{Y} \end{aligned}$$

$$2. \nabla_v(f\mathbf{X}) = (\nabla_vf)\mathbf{X}(p) + f(p)\nabla_v\mathbf{X}$$

$$\begin{aligned} \nabla_v(f\mathbf{X}) &= (\alpha(t), \nabla_vfX_1, \nabla_vfX_2, \dots, \nabla_vfX_{n+1}) \\ &= (\alpha(t), (\nabla_vf)X_1 + f\nabla_vX_1, (\nabla_vf)X_2 + f\nabla_vX_2, \dots, (\nabla_vf)X_{n+1} + f\nabla_vX_{n+1}) \\ &= (\alpha(t), (\nabla_vf)X_1, (\nabla_vf)X_2, \dots, (\nabla_vf)X_{n+1}) \\ &\quad + (\alpha(t), f\nabla_vX_1, f\nabla_vX_2, \dots, f\nabla_vX_{n+1}) \\ &= \nabla_vf(\alpha(t), X_1, X_2, \dots, X_{n+1}) + f(\alpha(t), \nabla_vX_1, \nabla_vX_2, \dots, \nabla_vX_{n+1}) \\ &= (\nabla_vf)\mathbf{X} + f\nabla_v\mathbf{X} \end{aligned}$$

$$3. \nabla_v(\mathbf{X} \cdot \mathbf{Y}) = \nabla_v\mathbf{X} \cdot \mathbf{Y}(p) + \mathbf{X}(p) \cdot \nabla_v\mathbf{Y}$$

$$\begin{aligned} \nabla_v(\mathbf{X} \cdot \mathbf{Y}) &= \nabla_v(X_1Y_1 + X_2Y_2 + \dots + X_{n+1}Y_{n+1}) \\ &= (\nabla_vX_1)Y_1 + X_1\nabla_vY_1 + (\nabla_vX_2)Y_2 + X_2\nabla_vY_2 \\ &\quad + \dots + (\nabla_vX_{n+1})Y_{n+1} + X_{n+1}\nabla_vY_{n+1} \\ &= (\nabla_vX_1)Y_1 + (\nabla_vX_2)Y_2 + \dots + (\nabla_vX_{n+1})Y_{n+1} \\ &\quad + X_1\nabla_vY_1 + X_2\nabla_vY_2 + \dots + X_{n+1}\nabla_vY_{n+1} \\ &= (\alpha(t), \nabla_vX_1, \nabla_vX_2, \dots, \nabla_vX_{n+1}) \cdot (\alpha(t), Y_1, Y_2, \dots, Y_{n+1}) \\ &\quad + (\alpha(t), X_1, X_2, \dots, X_{n+1}) \cdot (\alpha(t), \nabla_vY_1, \nabla_vY_2, \dots, \nabla_vY_{n+1}) \\ &= \nabla_v\mathbf{X} \cdot \mathbf{Y} + \mathbf{X} \cdot \nabla_v\mathbf{Y} \end{aligned}$$

### 18.9.3 Properties of covariant differentiation, $D_v\mathbf{X}$

$$1. D_v(\mathbf{X} + \mathbf{Y}) = D_v\mathbf{X} + D_v\mathbf{Y}$$

$$\begin{aligned} D_v(\mathbf{X} + \mathbf{Y}) &= \nabla_v(\mathbf{X} + \mathbf{Y}) - [\nabla_v(\mathbf{X} + \mathbf{Y}) \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= \nabla_v\mathbf{X} + \nabla_v\mathbf{Y} - [(\nabla_v\mathbf{X} + \nabla_v\mathbf{Y}) \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= \nabla_v\mathbf{X} - [\nabla_v\mathbf{X} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) + \nabla_v\mathbf{Y} - [\nabla_v\mathbf{Y} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= D_v\mathbf{X} + D_v\mathbf{Y} \end{aligned}$$

$$2. D_v(f\mathbf{X}) = (\nabla_v f)\mathbf{X}(p) + f(p)D_v\mathbf{X}$$

$$\begin{aligned} D_v(f\mathbf{X}) &= \nabla_v(f\mathbf{X}) - [\nabla_v(f\mathbf{X}) \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= (\nabla_v f)\mathbf{X} + f\nabla_v\mathbf{X} - [((\nabla_v f)\mathbf{X} + f\nabla_v\mathbf{X}) \cdot \mathbf{N} \circ \alpha](\mathbf{N} \circ \alpha) \\ &= (\nabla_v f)\mathbf{X} - [(\nabla_v f)\mathbf{X} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) + f\nabla_v\mathbf{X} - [f\nabla_v\mathbf{X} \cdot (\mathbf{N} \circ \alpha)](\mathbf{N} \circ \alpha) \\ &= \nabla_v f\mathbf{X} + fD_v\mathbf{X} \end{aligned}$$

$$3. \nabla_v(\mathbf{X} \cdot \mathbf{Y}) = D_v\mathbf{X} \cdot \mathbf{Y}(p) + \mathbf{X}(p) \cdot D_v\mathbf{Y}$$

$$\begin{aligned} \nabla_v(\mathbf{X} \cdot \mathbf{Y}) &= \nabla_v(X_1Y_1) + \nabla_v(X_2Y_2) + \cdots + \nabla_v(X_{n+1}Y_{n+1}) \\ &= (\nabla_v X_1)Y_1 + X_1\nabla_v Y_1 + (\nabla_v X_2)Y_2 + X_2\nabla_v Y_2 \\ &\quad + \cdots + (\nabla_v X_{n+1})Y_{n+1} + X_{n+1}\nabla_v Y_{n+1} \\ &= (\nabla_v X_1)Y_1 + (\nabla_v X_2)Y_2 + \cdots + (\nabla_v X_{n+1})Y_{n+1} \\ &\quad + X_1\nabla_v Y_1 + X_2\nabla_v Y_2 + \cdots + X_{n+1}\nabla_v Y_{n+1} \\ &= (\nabla_v \mathbf{X}) \cdot \mathbf{Y} + \mathbf{X} \cdot \nabla_v \mathbf{Y} \\ &= D_v\mathbf{X} \cdot \mathbf{Y} + \mathbf{X} \cdot D_v\mathbf{Y} \text{ since } \mathbf{X}, \mathbf{Y} \perp \mathbf{N} \end{aligned}$$

#### 18.9.4 Weingarten Map, $L_p$

In the case of covariant derivatives along  $\alpha$ , we came across a linear map – the parallel transport  $P_\alpha$ . Here in the case of covariant derivative with respect to a vector, we have another linear map – the Weingarten map  $L_p$  given by

$$L_p(\mathbf{v}) = -\nabla_v \mathbf{N} \quad (18.23)$$

Weingarten map tells us how much the tangent space turns/tilts as we move through  $p$  along  $\alpha$  on an  $n$ -surface  $S$ . It gives the information about the shape of  $S$ . And is the **shape operator** of  $S$  at  $p$ .

#### 18.9.5 Properties of Weingarten Map, $L_p$

1. The normal component of acceleration is same for all parametrised curves in  $S$  passing through  $p$  with velocity  $\mathbf{v}$ ,  $\ddot{\alpha}(t_0) \cdot \mathbf{N}(p) = L_p(\mathbf{v}) \cdot \mathbf{v}$

$$\begin{aligned} 0 &= \dot{\alpha} \cdot (\mathbf{N} \circ \alpha) \text{ since } \dot{\alpha} \perp \mathbf{N} \\ 0 &= [\dot{\alpha} \cdot (\mathbf{N} \circ \alpha)]' \\ &= \ddot{\alpha} \cdot (\mathbf{N} \circ \alpha) + \dot{\alpha} \cdot (\mathbf{N} \dot{\circ} \alpha) \end{aligned}$$

Rearranging the terms and evaluating at  $t = t_0$ , we get

$$\begin{aligned} \ddot{\alpha}(t_0) \cdot (\mathbf{N} \circ \alpha(t_0)) &= -\dot{\alpha}(t_0) \cdot (\mathbf{N} \dot{\circ} \alpha)(t_0) \\ \ddot{\alpha}(t_0) \cdot \mathbf{N}(p) &= \mathbf{v} \cdot L_p(\mathbf{v}) \text{ since } (\mathbf{N} \dot{\circ} \alpha)(t_0) = \nabla_v \mathbf{N} = -L_p(\mathbf{v}) \end{aligned}$$

2. Weingarten map is self adjoint,  $L_p(\mathbf{v}) \cdot \mathbf{w} = \mathbf{v} \cdot L_p(\mathbf{w})$

Clearly, we have  $\mathbf{v}, \mathbf{w} \in S_{\alpha(t)}^\perp$  and  $\nabla f \in S_{\alpha(t)}^\perp$ . Also note that

$$\begin{aligned}\nabla f(p) &= \left( p, \frac{\partial f}{\partial x_1}(p), \frac{\partial f}{\partial x_2}(p), \dots, \frac{\partial f}{\partial x_{n+1}}(p) \right) \\ \nabla_v(\nabla f(p)) &= \left( p, \nabla_v \frac{\partial f}{\partial x_1}(p), \nabla_v \frac{\partial f}{\partial x_2}(p), \dots, \nabla_v \frac{\partial f}{\partial x_{n+1}}(p) \right) \\ &= \left( p, \nabla \left( \frac{\partial f}{\partial x_1} \right)(p) \cdot \mathbf{v}, \nabla \left( \frac{\partial f}{\partial x_2} \right)(p) \cdot \mathbf{v}, \dots, \nabla \left( \frac{\partial f}{\partial x_{n+1}} \right)(p) \cdot \mathbf{v} \right) \\ &= \left( p, \sum_{j=1}^{n+1} \frac{\partial^2 f}{\partial x_j \partial x_1}(p) v_j, \sum_{j=1}^{n+1} \frac{\partial^2 f}{\partial x_j \partial x_2}(p) v_j, \dots, \sum_{j=1}^{n+1} \frac{\partial^2 f}{\partial x_j \partial x_{n+1}}(p) v_j \right)\end{aligned}$$

$$\begin{aligned}L_p(\mathbf{v}) \cdot \mathbf{w} &= -\nabla_v \mathbf{N} \cdot \mathbf{w} \\ &= -\nabla_v \left( \frac{\nabla f}{\|\nabla f\|} \right) \cdot \mathbf{w} \\ &= \left[ \left( -\nabla_v \frac{1}{\|\nabla f\|} \right) \nabla f - \frac{1}{\|\nabla f\|} \nabla_v(\nabla f) \right] \cdot \mathbf{w} \\ &= \left[ -\frac{1}{\|\nabla f\|} \nabla_v(\nabla f) \right] \cdot \mathbf{w} \text{ since } \nabla f \perp \mathbf{w} \\ &= -\frac{1}{\|\nabla f\|} \left[ \sum_{k=1}^{n+1} \left( \sum_{j=1}^{n+1} \frac{\partial^2 f}{\partial x_j \partial x_k}(p) v_j \right) w_k \right] \\ &= -\frac{1}{\|\nabla f\|} \sum_{j,k=1}^{n+1} \frac{\partial^2 f}{\partial x_j \partial x_k}(p) v_j w_k\end{aligned}$$

Similarly,

$$L_p(\mathbf{w}) \cdot \mathbf{v} = -\frac{1}{\|\nabla f\|} \sum_{j,k=1}^{n+1} \frac{\partial^2 f}{\partial x_j \partial x_k}(p) w_j v_k$$

Since  $\frac{\partial^2 f}{\partial x_j \partial x_k}(p) = \frac{\partial^2 f}{\partial x_k \partial x_j}(p)$ , both the sums are equal

$$L_p(\mathbf{v}) \cdot \mathbf{w} = \mathbf{v} \cdot L_p(\mathbf{w})$$

## 18.10 The Curvature of Plane Curves

**curvature** Let  $C$  be plane curve in  $\mathbb{R}^2$ . The curvature of  $C$  is the function  $\kappa : \mathbb{R}^2 \rightarrow \mathbb{R}$  defined by  $\kappa(p) = L_p(\mathbf{v}) \cdot \mathbf{v} / \|\mathbf{v}\|^2$

For a parametrised curve  $\alpha$  in  $C$  with  $\dot{\alpha} \neq \mathbf{0}$ ,

$$\kappa(\alpha(t)) = \frac{\ddot{\alpha}(t) \cdot (\mathbf{N} \circ \alpha(t))}{\|\dot{\alpha}(t)\|^2}$$

For unit speed parametrised curve passing through  $p$ , curvature at  $p$  is the normal component of acceleration at that point.

**local parametrisation of plane curve** Let  $C$  be a plane curve in  $\mathbb{R}^2$ . Then  $C$  is an oriented 1-surface. Parametrisation of a segment of  $C$  containing  $p$  is a function  $\alpha : I \rightarrow C$  such that

1.  $\alpha$  is regular. That is,  $\dot{\alpha}(t) \neq \mathbf{0}$ ,  $\forall t \in I$
2.  $\alpha$  is consistently oriented with  $C$ . That is, the orientation at  $\alpha(t)$  should be the same as that of  $C$  at that point.
3.  $p \in \alpha(I)$

**circle of curvature** Let  $C$  be a plane curve in  $\mathbb{R}^2$  with orientation  $\mathbf{N}$ . The circle of curvature of  $C$  at a point  $p$  is the unique oriented circle  $O$  with orientation  $\mathbf{N}_1$  such that

1.  $O$  is tangent to  $C$  at  $p$ ,  $C_p = O_p$ .
2.  $O$  is oriented consistently with  $C$ ,  $\mathbf{N}(p) = \mathbf{N}_1(p)$  and
3. Its normal turns at the same rate at  $p$  as the normal of the plane curve  $C$ ,  $\nabla_v \mathbf{N} = \nabla_v \mathbf{N}_1$

This circle,  $O$  is the circle which hugs the curve  $C$  closest among all circles containing  $p$ .

$\alpha$  is regular. That is,  $\dot{\alpha}(t) \neq \mathbf{0}$ ,  $\forall t \in I$

**radius of curvature** The radius of the circle of curvature of  $C$  at point  $p$ . And radius of curvature,  $r = 1/|\kappa(p)|$ .

**center of curvature** The center of the circle of curvature of  $C$  at point  $p$ .

**Meaning of Sign of Curvature** If curvature at  $p$  is positive,  $\kappa(p) > 0$  then the curve at  $p$  is turning towards  $\mathbf{N}(p)$ . If curvature at  $p$  is negative,  $\kappa(p) < 0$  then the curve at  $p$  is turning away from  $\mathbf{N}(p)$ .

**Construction of a local parametrisation** We have,  $C = f^{-1}(c)$  where  $\nabla f(q) \neq \mathbf{0}$ . Consider  $\mathbf{X}(q) = \left( q, \frac{\partial f}{\partial x_2}(q), -\frac{\partial f}{\partial x_1}(q) \right)$ . The maximal integral curve  $\alpha$  through  $p$  in  $\mathbf{X}$  is a local parametrisation of a segment of  $C$  containing  $p$ .

For example, (refer : Exercise 10.3c)  $f(x_1, x_2) = x_2 - ax_1^2$ . Then  $C = f^{-1}(c) = \{(x_1, x_2) \in \mathbb{R}^2 : x_2 - ax_1^2 = c, a \neq 0\}$ . The global parametrisation of  $C$  with orientation  $\nabla f/\|\nabla f\|$  can be obtained by constructing the maximal integral curve for each segment of a piecewise smooth plane curve. Since the given curve is smooth, it is sufficient to construct one integral curve through any point on that curve.

We have,  $\frac{dx_1}{dt}(t) = -2ax_1(t)$  and  $\frac{dx_2}{dt} = 1$ . Thus,  $x_1(t) = \cos 2at + c_1$  and  $x_2(t) = t + c_2$ . We have  $(0, c) \in C$  since  $c - a0^2 = c$ . Therefore,  $x_1(0) = 0 = c_1$  and  $x_2(0) = c = c_2$ . Thus, we have  $\alpha : I \rightarrow \mathbb{R}^2$  given by  $\alpha(t) = (\cos 2at, t + c)$  as a global parametrisation of  $C$ .

### 18.10.1 Unit speed local parametrisation of $C$ is unique.

**Theorem 18.10.1.** *Local parametrisations of plane curves are unique upto reparametrisation.*

*Proof.* Let  $\alpha$  be the unit speed local parametrisation which is the maximal integral curve through  $p$ . Let  $\beta : \tilde{I} \rightarrow C$  be any parametrisation of a segment of  $C$  containing  $p$ . Then it is enough to prove that there exists a smooth function  $h : \tilde{I} \rightarrow \mathbb{R}$  such that  $h'(t) > 0$  and  $\beta(t) = \alpha(h(t))$  for every  $t \in \tilde{I}$ .

Define  $h : \tilde{I} \rightarrow \mathbb{R}$  defined by

$$h(t) = \int_{t_0}^t \|\dot{\beta}(u)\| du$$

Clearly,  $h$  is monotone,  $h(t_0) = 0$  and  $h'(t) = \|\dot{\beta}(t)\|$ . Since  $h$  is monotone and strictly increasing,  $h(h^{-1}(t)) = t$ . Thus,  $(h \circ h^{-1})'(t) = h'(h^{-1}(t))(h^{-1})'(t)$ . Therefore,  $(h^{-1})'(t) = 1/h'(h^{-1}(t))$ .

Now,  $\beta \circ h^{-1}$  is a reparametrised curve with velocity,  $(\beta \circ h^{-1})(t)$  given by,

$$\begin{aligned} (\beta \circ h^{-1})(t) &= \dot{\beta}(h^{-1}(t)) (h^{-1})'(t) \text{ by chain rule} \\ &= \dot{\beta}(h^{-1}(t)) / h'(h^{-1}(t)) \\ &= \dot{\beta}(h^{-1}(t)) / \|\dot{\beta}(h^{-1}(t))\| \text{ since } h'(t) = \|\dot{\beta}(t)\| \end{aligned}$$

We know that,  $\mathbf{X}(\beta(h^{-1}(t)))$  spans the tangent space  $S_{\beta(h^{-1}(t))}$ .

$$= \mathbf{X}(\beta(h^{-1}(t))) \text{ since } \beta \circ h^{-1}(t) / \|\beta \circ h^{-1}(t)\| = \mathbf{X}(\beta(h^{-1}(t)))$$

Thus,  $\beta \circ h^{-1}$  is an integral curve through  $p$  in  $\mathbf{X}$ . By uniqueness of integral curves,  $\beta \circ h^{-1}(t) = \alpha(t)$  for all  $t \in \tilde{I}$ .  $\square$

## 18.11 Arc Length and Line Integrals

**length of arc** The length of parametrised arc  $\alpha$  is the integral of its speed.

Let  $\alpha : [a, b] \rightarrow \mathbb{R}^{n+1}$ . Length of  $\alpha$  is

$$l(\alpha) = \int_a^b \|\dot{\alpha}(t)\| dt \quad (18.24)$$

For Example (refer : Exercise 11.4) : Given  $\alpha : [0, 2\pi] \rightarrow \mathbb{R}^4$ ,  $\alpha(t) = (\cos t, \sin t, \cos t, \sin t)$ . Then,  $\dot{\alpha}(t) = (\alpha(t), -\sin t, \cos t, -\sin t, \cos t)$ . We have,  $\|\dot{\alpha}(t)\| = \sqrt{\sin^2 t + \cos^2 t + \sin^2 t + \cos^2 t} = \sqrt{2}$ . Therefore,  $l(\alpha) = \int_0^{2\pi} \sqrt{2} dt = 2\sqrt{2}\pi$ .

**Note :** Length of arc is the total distance travelled.

For example :  $\alpha_1 : [0, 2\pi] \rightarrow \mathbb{R}^2$  defined by  $\alpha_1(t) = (\cos t, \sin t)$  has its length  $2\pi$ , the perimeter of unit circle. But, for  $\alpha_2 : [0, 4\pi] \rightarrow \mathbb{R}^2$  defined by  $\alpha_2(t) = (\cos t, \sin t)$  has its length  $4\pi$ . And  $\alpha : \mathbb{R} \rightarrow \mathbb{R}^2$  defined by  $\alpha(t) = (\cos t, \sin t)$  has its length  $+\infty$ .

### 18.11.1 Properties of Arc Length

1. Arc length is preserved under reparametrisation.

Let  $\alpha : [a, b] \rightarrow \mathbb{R}^{n+1}$  be a parametrised curve. And let  $\beta : [c, d] \rightarrow \mathbb{R}^{n+1}$  be a reparametrisation of  $\alpha$  defined by  $\beta(t) = \alpha(h(t))$  where  $h : [a, b] \rightarrow [c, d]$ . Then length of  $\beta$  is given by,

$$\begin{aligned} l(\beta) &= \int_c^d \|\dot{\beta}(t)\| dt \\ &= \int_c^d \|\dot{\alpha}(h(t))\| h'(t) dt \\ &= \int_a^b \|\alpha(u)\| du = l(\alpha) \end{aligned}$$

2. Unit speed arcs are parametrised by arc length.

Let  $\alpha$  be a unit speed parametrised curve. That is,  $\|\dot{\alpha}\| = 1$ . Therefore  $l(\alpha) = \int_a^b dt = b - a$ . Clearly, length of the arc is the length of the parameter interval.

**Theorem 18.11.1** (Existence of Global Parametrisation). *Let  $C$  be an oriented plane curve. Then there exists a global parametrisation of  $C$  if and only if  $C$  is connected.*

*Proof. Sufficient Part :* Suppose plane curve  $C$  has a global parametrisation,  $\alpha : I \rightarrow \mathbb{R}^{n+1}$ . Let  $p, q \in C$ , then  $p = \alpha(t_1)$  and  $q = \alpha(t_2)$  for some  $t_1, t_2 \in I$ . WLOG  $t_1 < t_2$  and clearly, there exists a path from  $p$  to  $q$  obtained restricting  $\alpha$  to  $[t_1, t_2]$ . Therefore, the plane curve  $C$  is path-connected and thus connected.

#### Necessary Part : Step 1 Construction of $\alpha$

Let  $\mathbf{X}$  be the unit tangent vector field on  $C$  (obtained by rotating  $\nabla f$  by an angle of  $-\pi/2$  and orientation  $\mathbf{N} = \nabla f / \|\nabla f\|$ ). Suppose the plane curve  $C = f^{-1}(c)$  is connected. Let  $p \in C$  and  $\alpha$  be the local parametrisation of  $C$  at  $p$ , which is nothing but the maximal integral curve of  $\mathbf{X}$  through  $p$ . Let  $p_1 \in C$ . Parametrisation  $\alpha$  is global if  $p_1 \in \text{Image } \alpha$ .

#### Step 2 : Construction of $\beta$

The plane curve  $C$  is connected. Thus there exists a continuous path  $\beta$  from  $p$  to  $p_1$  where  $\beta : [a, b] \rightarrow C$  with  $\beta(a) = p$  and  $\beta(b) = p_1$ . We have,  $\beta(a) \in \text{Image } \alpha$ . But,  $\beta(b) \notin \text{Image } \alpha$ . Let  $t_0 = \sup\{t \in [a, b] : \beta(t) \in \text{Image } \alpha\}$ , the least upper bound of the points which are in  $\text{Image } \alpha$ . That is, if  $t > t_0$ , then  $\beta(t) \notin \text{Image } \alpha$ . And if  $t < t_0$ , then there exists an  $\epsilon > 0$  such that  $\beta(t + \epsilon) \in \text{Image } \alpha$ . Otherwise,  $t_0$  is not the supremum.

#### Step 3 : Construction of $\gamma$

Let  $\gamma$  be the maximal integral curve of  $\mathbf{X}$  through  $\beta(t_0) = p_0$ . Suppose there exists an open rectangle  $B$  about  $\beta(t_0) = p_0$  such that  $C \cap B \subset \text{Image } \gamma$ . Then by the continuity of  $\beta$ , there exists  $\delta > 0$  such that  $\beta(t) \in \text{Image } \gamma$  for all  $t \in (t_0 - \delta, t_0 + \delta)$ . Thus, there exists  $\epsilon > 0$  such that  $\epsilon < \delta$  and  $\beta(t_0 - \epsilon) \in \text{Image } \gamma$ . Since  $t_0$  is the least upper bound,  $\beta(t_0 - \epsilon) \in \text{Image } \alpha$  as well. Thus,  $\alpha$  and  $\gamma$  are maximal integral curves through a common point (in the

neighbourhood of  $p_0$ ). And thus,  $\text{Image } \alpha = \text{Image } \gamma$ . Clearly,  $p_1 \in \text{Image } \alpha$ . Therefore, the parametrisation  $\alpha$  is global. Therefore it is sufficient to prove that there exists an open box  $B$  such that  $C \cap B \subset \text{Image } \gamma$ .

**Step 4 : Construction of  $A$**

Let  $u = (\frac{\partial f}{\partial x_2}(p_0), -\frac{\partial f}{\partial x_1}(p_0))$  and  $v = (\frac{\partial f}{\partial x_1}(p_0), \frac{\partial f}{\partial x_2}(p_0))$ . Clearly,  $\mathbf{u} = (p_0, u) \in C_{p_0}$  and  $\mathbf{v} = (p_0, v) \perp C_{p_0}$ . Consider the open rectangle  $A$  given by,

$$A = \{p_0 + ru + sv : |r| < \epsilon_1, |s| < \epsilon_2\}$$

where  $\epsilon_1, \epsilon_2$  are so chosen that  $A$  is contained in the domain of  $f$  and  $\nabla f(q) \cdot (q, v) > 0$  for every  $q \in A$ . Since  $\nabla f(p_0) \cdot \mathbf{v} \neq 0$ , by continuity of  $\nabla f(q) \cdot (q, v)$  there exists a neighbourhood of  $p_0$  in which this function is positive.

**Step 5 : Construction of  $\{g_r\}$**

Consider the family of functions  $\{g_r\}$  defined by  $g_r(s) = f(p_0 + ru + sv)$ . Since  $\nabla f(q) \cdot (q, v) > 0$  for every  $q \in A$ , for  $|r| < \epsilon_1$ ,  $g_r(s)$  is strictly increasing in  $(-\epsilon_2, \epsilon_2)$  as  $g'_r(s) > 0$ . That is, there exists at most one  $s \in (-\epsilon_2, \epsilon_2)$  such that  $f(p_0 + ru + sv) = c$ .

**Step 6 : Construction of  $B$**

Since the vectors  $\mathbf{u}, \mathbf{v}$  spans  $C_{p_0}$ , we have  $\gamma(t) = p_0 + h_1(t)u + h_2(t)v$  where  $h_1, h_2$  are real-valued functions. Clearly,

$$\begin{aligned} (\gamma(t) - p_0) \cdot \mathbf{u} &= h_1(t)\mathbf{u} \cdot \mathbf{u} + h_2(t)\mathbf{u} \cdot \mathbf{v} \\ &= h_1(t)\|u\|^2 \text{ since } \mathbf{u} \perp \mathbf{v} \end{aligned}$$

$$\begin{aligned} (\gamma(t) - p_0) \cdot \mathbf{v} &= h_1(t)\mathbf{u} \cdot \mathbf{v} + h_2(t)\mathbf{v} \cdot \mathbf{v} \\ &= h_2(t)\|v\|^2 \end{aligned}$$

Thus,

$$h_1(t) = \frac{(\gamma(t) - p_0) \cdot \mathbf{u}}{\|u\|^2} \quad (18.25)$$

$$h_2(t) = \frac{(\gamma(t) - p_0) \cdot \mathbf{v}}{\|v\|^2} \quad (18.26)$$

Now,  $h'_1(0) = \dot{\gamma}(0) \cdot (p_0, u/\|u\|^2) = \mathbf{X}(p_0) \cdot \mathbf{X}(p_0)/\|\mathbf{X}(p_0)\|^2 = 1$ . And  $h_1(0) = 0$  and  $h_2(0) = 0$ . Thus, there exists  $(t_1, t_2)$  containing 0 such that  $\gamma(t) \in A$  and  $h'_1(t) > 0$  for every  $t \in (t_1, t_2)$ . Set  $r_1 = h_1(t_1)$  and  $r_2 = h_1(t_2)$ . Since  $h_1$  is continuous and strictly increasing, for every  $r \in (r_1, r_2)$  there exists  $t \in (t_1, t_2)$  such that  $r = h_1(t)$  and  $s = h_2(t)$ .

$$\text{Define } B = \{p_0 + ru + sv : r \in (r_1, r_2), |s| < \epsilon_2\} \quad (18.27)$$

Then  $p_0 + ru + sv \in B \cap C$  if and only if there exists  $t \in (t_1, t_2)$  such that  $r = h_1(t)$ , and  $s = h_2(t)$ . That is,  $p_0 + ru + sv \in \text{Image } \gamma$ . Therefore,  $C \cap B \subset \text{Image } \gamma$  as required in Step 3.  $\square$

**Theorem 18.11.2.** *Let  $C$  be a connected oriented plane curve and let  $\beta : I \rightarrow C$  be a unit speed global parametrisation of  $C$ . Then  $\beta$  is either one to one or periodic. Moreover,  $\beta$  is periodic if and only if  $C$  is compact.*

*Proof.* Suppose  $\beta(t_1) = \beta(t_2)$  for some  $t_1, t_2 \in I$  ( $t_1 \neq t_2$ ). Let  $\mathbf{X}$  be the unit tangent vector field on  $C$ . And let  $\alpha$  be the maximal integral curve of  $\mathbf{X}$  through  $\beta(t_1)$ . That is,  $\alpha(0) = \beta(t_1) = \beta(t_2)$ . Since  $\beta$  is a unit speed global parametrisation of  $C$ ,  $\beta$  is also a maximal integral curve of  $\mathbf{X}$ . Thus,  $\beta$  is a reparametrisation of  $\alpha$ . Thus,  $\beta(t) = \alpha(t - t_1)$  and  $\beta(t) = \alpha(t - t_2)$ . Let  $\tau = t_2 - t_1$ . Then  $\beta(t) = \alpha(t - t_1) = \alpha(t - t_2) = \alpha(t - t_2 + t_2 - t_1) = \alpha(t + \tau - t_1) = \beta(t + \tau)$ . Therefore,  $\beta$  is periodic.

Suppose  $\beta$  is periodic, then its continuous image,  $C = \beta[t, t + \tau]$  is compact. Suppose  $\beta$  is not periodic, then  $\beta$  is one to one. Then,  $C$  is not compact as  $\beta^{-1}$  doesn't attain its extrema. Thus, it is enough to prove that  $\beta^{-1}$  is continuous.

Given  $t_0 \in I$  and  $\epsilon > 0$ , define  $\gamma(t) = \beta(t + t_0)$  such that  $|t| < \epsilon$  and  $t + t_0 \in I$ . Chose an open rectangle  $B$  about  $p_0 = \beta(t_0)$  such that  $C \cap B \subset \text{Image } \gamma$ . Then,  $|\beta^{-1}(p) - t_0| = |\gamma^{-1}(p)| < \epsilon$  for every  $p \in C \cap B$ .

Therefore it is enough to prove that such an open rectangle  $B$  exists. The construction of which is given in previous proof.  $\square$

**Fundamental Domain** Let  $\beta$  be a periodic parametrised plane curve with period  $\tau$ . Then any subset  $[t_0, t_0 + \tau]$  of its domain is a fundamental domain of  $\beta$ .

### 18.11.2 1-Form

**Definitions 18.11.1** (1-form). The 1-form is a function  $\omega : U \times \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  where  $U \subset \mathbb{R}^{n+1}$  such that for every  $p \in U$ ,  $\mathbb{R}_p^{n+1} \subset U \times \mathbb{R}^{n+1}$ .

**Note :** 1-form  $\omega$  is smooth if  $\omega$  is a smooth function. A few examples of 1-form are given below:  $\omega_{\mathbf{X}}$ ,  $df$  and  $dx_i$ .

**1-form dual** Let  $\mathbf{X}$  be a vector field on  $U$ . Then 1-form dual of  $\mathbf{X}$  is  $\omega_{\mathbf{X}}$  given by  $\omega_{\mathbf{X}}(p, v) = \mathbf{X}(p) \cdot (p, v)$

**differential of  $f$**  Let  $f : U \rightarrow \mathbb{R}$  be a smooth function. Then differential of  $f$  is  $df$  given by  $df(\mathbf{v}) = \nabla f(p) \cdot \mathbf{v}$  where  $\mathbf{v} = (p, v)$ .

**Cartesian Coordinate Function** Let  $U \subset \mathbb{R}^{n+1}$ . Then  $i$ th cartesian coordinate function  $x_i$  is given by  $x_i : U \rightarrow \mathbb{R}$ ,  $x_i(a_1, a_2, \dots, a_{n+1}) = a_i$ .

$dx_i$  The differential of the cartesian coordinate function  $x_i$  is given by  $dx_i$  where  $dx_i : U \times \mathbb{R}^{n+1} \rightarrow \mathbb{R}$

$$dx_i = \nabla x_i(p) \cdot \mathbf{v} = (p, 0, 0, \dots, 1, \dots, 0) \cdot (p, v_1, v_2, \dots, v_{n+1}) = v_i$$

### 18.11.3 1-form

**Sum of two 1-forms**  $\omega_1 + \omega_2$

$$\begin{aligned} (\omega_1 + \omega_2)(\mathbf{v}) &= \nabla(\omega_1 + \omega_2) \cdot \mathbf{v} \\ &= \nabla\omega_1 \cdot \mathbf{v} + \nabla\omega_2 \cdot \mathbf{v} \\ &= \omega_1(\mathbf{v}) + \omega_2(\mathbf{v}) \end{aligned}$$



**Product of function and 1-form  $f\omega$** 

$$\begin{aligned}
f\omega(\mathbf{v}) &= \nabla(f\omega) \cdot \mathbf{v} \\
&= f\nabla\omega \cdot \mathbf{v} \\
&= f(p)\omega(\mathbf{v})
\end{aligned}$$

**Note :** Let  $\omega(\mathbf{X}) : U \rightarrow \mathbb{R}$ ,  $(\omega(\mathbf{X}))(p) = \omega(\mathbf{X}(p))$ .  
If  $\omega$  and  $\mathbf{X}$  are smooth, then function  $\omega(\mathbf{X})$  is also smooth.

**Proposition 18.11.1** (1-form representation). *For every 1-form  $\omega$  on  $U$ , there exists unique functions  $f_i : U \rightarrow \mathbb{R}$  such that*

$$\omega = \sum_{i=1}^{n+1} f_i dx_i \quad (18.28)$$

And  $\omega$  is smooth if and only if each  $f_i$  is smooth.

*Proof.* Let  $\mathbf{X}_j$  be a smooth vector field on  $U$  where  $\mathbf{X}_j(p) = (p, 0, 0, \dots, 1, \dots, 0)$ . Then,  $dx_i(\mathbf{X}_j) = \delta_{ij}$ . Suppose  $\omega = \sum f_i dx_i$ . Then

$$\omega(\mathbf{X}_j) = \sum_{i=1}^{n+1} (f_i dx_i)(\mathbf{X}_j) = \sum_{i=1}^{n+1} f_i(p) dx_i(\mathbf{X}_j) = \sum_{i=1}^{n+1} f_i(p) \delta_{ij} = f_j(p)$$

Thus  $f_j$  is unique for each  $j$  if they exist. And for each 1-form  $\omega$  there exist functions  $f_i$  such that  $\omega = \sum f_i dx_i$  exists where  $f_i = \omega(\mathbf{X}_i)$ . And thus,  $\omega$  is smooth if and only if each  $f_i$  is smooth.  $\square$

**Corollary 18.11.2.1.** *Let  $f : U \rightarrow \mathbb{R}$ . Then*

$$df = \sum_{i=1}^{n+1} \frac{\partial f}{\partial x_i} dx_i \quad (18.29)$$

*Proof.*

$$df(\mathbf{X}_j) = \nabla f \cdot \mathbf{X}_j = \sum_{i=1}^{n+1} \frac{\partial f}{\partial x_i} \delta_{ij} = \frac{\partial f}{\partial x_j} = f_j \implies df = \sum_{i=1}^{n+1} \frac{\partial f}{\partial x_i} dx_i$$

$\square$

**18.11.4 Line Integral**

**Definitions 18.11.2** (Line Integral). Let  $\omega$  be a 1-form and  $\alpha : [a, b] \rightarrow U$  be a parametrised curve on  $U$ . Then the line integral of  $\omega$  over  $\alpha$  is given by,

$$\int_{\alpha} \omega = \int_a^b \omega(\dot{\alpha}(t)) \quad (18.30)$$

**Note :** The line integral of 1-form  $\omega$  over a parametrised curve is invariant of reparametrisation.

Let  $\beta : [c, d] \rightarrow U$  be reparametrisation of  $\alpha : [a, b] \rightarrow U$  where  $\beta(t) = \alpha(h(t))$  and function  $h : [c, d] \rightarrow [a, b]$ . Then ,

$$\begin{aligned} \int_{\beta} \omega &= \int_c^d \omega(\dot{\beta}(t)) \\ &= \int_c^d \omega(\dot{\alpha}(h(t))h'(t)) \\ &= \int_c^d \omega(\dot{\alpha}(h(t)))h'(t) \\ &= \int_a^b \omega(\dot{\alpha}(u)) = \int_{\alpha} \omega \end{aligned}$$

**Theorem 18.11.3.** Let  $\eta$  be a 1-form on  $\mathbb{R}^2 - \{0\}$  defined by

$$\eta = \frac{-x_2}{x_1^2 + x_2^2} dx_1 + \frac{x_1}{x_1^2 + x_2^2} dx_2$$

Then for any closed, piecewise smooth, parametrised smooth curve  $\alpha : [a, b] \rightarrow \mathbb{R} - \{0\}$ , the line integral of  $\eta$  over  $\alpha$

$$\int_{\alpha} \eta = 2\pi k$$

*Proof.* Consider  $\varphi : [a, b] \rightarrow \mathbb{R}$  defined by  $\varphi(t) = \varphi(a) + \int_{\alpha_t} \eta$  where  $\alpha_t$  is the restriction of  $\alpha$  to  $[a, t]$  and  $\varphi(a)$  is so chosen that  $\alpha(a)/\|\alpha(a)\| = (\cos \varphi(a), \sin \varphi(a))$ . We claim that

$$\frac{\alpha(t)}{\|\alpha(t)\|} = (\cos \varphi(t), \sin \varphi(t)), \quad \forall t \in [a, b] \quad (18.31)$$

Let  $t_0 = \sup\{t \in [a, b] : \alpha(t)/\|\alpha(t)\| = (\cos \varphi(t), \sin \varphi(t))\}$ . By continuity the claim should be true for  $t_0$  as well. Therefore it is enough to prove that  $t_0 = b$ .

Define  $v = -\alpha(t_0)/\|\alpha(t_0)\|$  and  $V = \mathbb{R}^2 - \{rv : r \geq 0\}$ . And  $\theta_V : V \rightarrow \mathbb{R}$  is defined in such a way that  $v = (\cos \theta_v, \sin \theta_v)$  and  $\theta_v \in [0, 2\pi)$ . In other words,  $\theta_V$  a real-valued function which give the angle of the ray from origin through the points.

Then  $(\cos \theta_V(\alpha(t_0)), \sin \theta_V(\alpha(t_0))) = \alpha(t_0)/\|\alpha(t_0)\| = (\cos \varphi(t_0), \sin \varphi(t_0))$ . Clearly,  $\varphi(t_0) - \theta_V(\alpha(t_0)) = 2\pi m$  for some integer  $m$ .

Consider a  $\delta$  neighbourhood of  $t_0$  such that  $\alpha(t) \in V$  for every point  $t$  in it. Then,  $\alpha$  is smooth in that neighbourhood.

$$\frac{d}{dt} (\varphi(t) - \theta_V(\alpha(t))) = \eta(\dot{\alpha}(t)) = d\theta_V(\dot{\alpha}(t)) = 0$$

Thus,  $\varphi(t) - \theta_V(\alpha(t)) = 2\pi m$  for every  $t \in (t_0 - \delta, t_0 + \delta)$ . Therefore, the claim is true. Suppose  $t_0 < b$ , then the claim is true for  $t_0 + \epsilon \in (t_0 - \delta, t_0 + \delta)$  which is a contradiction to the choice of  $t_0$  as least upper bound of all such values.

Thus,  $(\cos \varphi(a), \sin \varphi(a)) = (\cos \varphi(b), \sin \varphi(b))$ . Therefore,  $\varphi(b) - \varphi(a) = 2\pi k$  for some integer  $k$ . And  $\int_\alpha \eta = \int_a^b \eta(\dot{\alpha}(t)) dt = \varphi(b) - \varphi(a) = 2\pi k$ .  $\square$

**Winding Number** The winding number of  $\alpha$  is  $k(\alpha) = \frac{1}{2\pi} \int_\alpha \eta$

**Note :** Winding number is the number of times  $\alpha$  winds around origin.

## 18.12 Curvature of Surfaces

**Normal Curvature** The number  $\kappa(\mathbf{v}) = L_p(\mathbf{v}) \cdot \mathbf{v}$  is the normal curvature of the  $n$ -surface  $S$  at point  $p$  in the direction  $\mathbf{v}$ .

If  $\kappa(\mathbf{v}) > 0$ , then  $S$  bends toward orientation at  $p$ . And if  $\kappa(\mathbf{v}) < 0$ , then  $S$  bends away from orientation at  $p$ .

We have,  $\kappa : S_p \rightarrow \mathbb{R}$ . When  $n = 1$ , then  $\kappa(\mathbf{v}) = \kappa(p)$ . Clearly  $S_p$  is compact and  $\kappa$  attains its extrema, which are the eigen values of the Weingarten Map.

**Normal Section** The normal section of an  $n$ -surface  $S$  with orientation  $N$  is  $\mathcal{N}(\mathbf{v}) \subset \mathbb{R}^{n+1}$  given by

$$\mathcal{N}(\mathbf{v}) = \{q \in \mathbb{R}^{n+1} : q = p + xv + yN(p), (x, y) \in \mathbb{R}^2\} \quad (18.32)$$

You take an apple and place a knife at any point on that apple in any direction tangent to its surface and cut it. The apple cuts into two parts. The 2-dimensional space of that cut, is a **Nomal Section** of that apple. For any  $n \geq 2$ , the normal section of an  $n$ -surface is always 2 dimensional.

$S \cap \mathcal{N}(\mathbf{v})$  is the intersection of both the surface and its normal section.

In above example,  $S \cap \mathcal{N}(\mathbf{v})$  is plane curve traced by the skin around cross-section of that apple.

**Theorem 18.12.1** (Components of  $S \cap \mathcal{N}(\mathbf{v})$ ). *Let  $S$  be an oriented  $n$ -surface and  $\mathbf{v}$  is a tangent direction at  $p$ ,  $\mathbf{v} \in S_p$ . Then there exists an open  $V$  containing  $p$  such that  $S \cap \mathcal{N}(\mathbf{v}) \cap V$  is a plane curve. Moreover, the curvature at  $p$  of this curve is the normal curvature  $\kappa(\mathbf{v})$  of  $S$  at  $p$  in the direction  $\mathbf{v}$ .*

*Proof.* Let  $f : U \rightarrow \mathbb{R}$  such that  $S = f^{-1}(c)$  and  $\nabla f(q) \neq \mathbf{0}$ ,  $\forall q \in S$ . Let  $\widetilde{\nabla} f(q) = (q, \nabla f(q))$ . Let  $p \in S$  and  $\mathbf{v} \in S_p$ . Let  $i : \mathbb{R}^2 \rightarrow \mathbb{R}^n$  defined by  $i(x, y) = p + xv + yN(p)$ . Clearly, the image of the function  $i$  is the normal section,  $\mathcal{N}(\mathbf{v}) = i(\mathbb{R}^2)$ .

Let  $V$  be the space spanned by  $v$  and  $N(p)$ . That is,

$$V = \left\{ q \in U : \widetilde{\nabla} f(q) \cdot v \neq 0 \text{ OR } \widetilde{\nabla} f(q) \cdot N(p) \neq 0 \right\} \quad (18.33)$$

Then  $\nabla(f \circ i)(x, y) = (x, y, \widetilde{\nabla} f(i(x, y)) \cdot v, \widetilde{\nabla} f(i(x, y)) \cdot N(p)) \neq \mathbf{0}$  for every  $(x, y) \in i^{-1}(V)$ . Therefore,  $C = i^{-1}(S \cap \mathcal{N}(\mathbf{v})) = (f \circ i)^{-1}(c) \cap i^{-1}(V)$  is a plane curve. In other words,  $S \cap \mathcal{N}(\mathbf{v}) \cap V$  is a plane curve.  $\square$

**Lemma 18.12.2** (Extrema of Curvature are Eigenvalue of Weingarten Map). *Let  $V$  be a finite dimensional vector space with dot product and let  $L : V \rightarrow V$  be a self-adjoint linear transformation on  $V$ . Let  $S = \{v \in V : v \cdot v = 1\}$  and define  $f : S \rightarrow \mathbb{R}$  by  $f(v) = L(v) \cdot v$ . Suppose  $f$  is stationary at  $v_0 \in S$ . Then  $L(v_0) = f(v_0)v_0$ .*

*Proof.* Since  $f$  is stationary at  $v_0$ ,  $f \circ \alpha'(0) = 0$  for every parametrised curves in  $S$  with  $\alpha(0) = v_0$ . Let  $v \perp v_0$ . Then  $\alpha(t) = (\cos t)v_0 + (\sin t)v$ . Then,

$$\begin{aligned} 0 &= (f \circ \alpha)'(0) \\ &= \frac{d}{dt} (L(\alpha(t)) \cdot \alpha(t)) (0) \\ &= \frac{d}{dt} (\cos^2 t L(v_0) \cdot v_0 + 2 \sin t \cos t L(v_0) \cdot v + \sin^2 t L(v) \cdot v) \\ &= (-\sin 2t L(v_0) \cdot v_0 + 2 \cos 2t L(v_0) \cdot v + \sin 2t L(v) \cdot v)|_{t=0} \\ &= 2L(v_0) \cdot v \end{aligned}$$

Clearly,  $L(v_0) \perp v$  for every  $v \in v_0^\perp$ . In other words,  $L(v_0) \perp v_0^\perp$ . Thus,  $v_0$  can span the space  $L(v_0)$ . That is,  $L(v_0) = \lambda v_0$  for some  $\lambda \in \mathbb{R}$ . Thus  $v_0$  is an eigen vector of  $L(v_0)$ . And the corresponding eigen vector  $\lambda$  can be obtained,  $\lambda = \lambda v_0 \cdot v_0 = L(v_0) \cdot v_0 = f(v_0)$ .  $\square$

**Theorem 18.12.3.** *Let  $V$  be a finite dimensional vector space with dot product and let  $L : V \rightarrow V$  be a self-adjoint linear transformation on  $V$ . Then there exists an orthonormal basis for  $V$  consisting of eigen vectors of  $L$ .*

*Proof.* The proof is by mathematical induction. For  $n = 1$ , the theorem is true by lemma. Suppose that the theorem is true for  $n = k$ . It is enough to prove that theorem is then true for  $n = k + 1$ .

By lemma, there exists a unit vector  $v_1 \in V$  which is an eigen vector of  $L$ . This is done by selecting  $v_1$  such that  $L(v_1) \cdot v_1 \geq L(v) \cdot v$  for every unit vector  $v \in V$ . Define  $W = v_1^\perp$ . Then  $L(w) \cdot v_1 = w \cdot L(v_1) = \lambda w \cdot v_1 = 0$  for every  $w \in W$ . (Remember :  $L$  self-adjoint,  $L(w) \cdot v = w \cdot L(v)$ )

Thus the restriction of  $L$ ,  $L|_W : W \rightarrow W$  is surjective. And clearly self-adjoint. The dimension of  $W$  is  $k$ , thus by induction hypothesis there exists an orthonormal basis  $\{v_2, v_3, \dots, v_{k+1}\}$  for  $W$  with eigen vectors of  $L|_W$ . Clearly, these are eigen vectors of  $L$  and  $\{v_1, v_2, \dots, v_{k+1}\}$  is an orthonormal basis of  $V$ .  $\square$

**Principal Curvature Directions** are the vectors in the orthonormal basis  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  for  $S_p$  obtained as eigen vectors of the Weingarten Map. (Hint :  $\mathbb{R}^n$  is a vector space with dot product and  $L_p$  is self-adjoint, linear transformation)

**Principal Curvature** are the eigen values of the Weingarten Map. That is,  $k_i(p)$  is the eigen value of  $L_p : S_p \rightarrow S_p$  such that  $L_p(\mathbf{w}) \cdot \mathbf{v}_i = \kappa_i(p) \mathbf{w} \cdot \mathbf{v}_i$ .

**Theorem 18.12.4.** *Let  $S$  be an oriented  $n$ -surface in  $\mathbb{R}^{n+1}$ . Let  $p \in S$  and  $\{\kappa_1(p), \kappa_2(p), \dots, \kappa_n(p)\}$  be the principal curvatures of  $S$  at  $p$  with corresponding orthogonal principal curvature directions  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ . Then the normal*

curvature  $\kappa(\mathbf{v})$  in the direction  $\mathbf{v} \in S_p$  is given by

$$\kappa(v) = \sum_{i=1}^n \kappa_i(p)(\mathbf{v} \cdot \mathbf{v}_i)^2 = \sum_{i=1}^n \kappa_i(p) \cos^2 \theta_i \quad (18.34)$$

where  $\theta_i$  is the angle between  $\mathbf{v}$  and  $\mathbf{v}_i$ .

*Proof.* We have  $v \in V$  and  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  is an orthonormal basis for  $V$ . Then,

$$\mathbf{v} = \sum_{i=1}^n (\mathbf{v} \cdot \mathbf{v}_i) \mathbf{v}_i = \sum_{i=1}^n \cos \theta_i \mathbf{v}_i$$

We have,

$$\begin{aligned} \kappa(\mathbf{v}) &= L_p(\mathbf{v}) \cdot v \\ &= \sum_{i=1}^n \cos \theta_i L_p(\mathbf{v}_i) \cdot \mathbf{v} \\ &= \sum_{i=1}^n \kappa_i(\mathbf{v}_i) \mathbf{v}_i \cdot \mathbf{v} \\ &= \sum_{i=1}^n \kappa_i \cos^2 \theta_i \end{aligned}$$

□

**Direction Cosines** Let  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  be an orthonormal basis for  $S_p$ . Then  $\theta_i$  are the angles  $\mathbf{v}$  makes with each vector in that basis. And direction cosines are the components of  $\mathbf{v}$  in those directions. That is,  $\cos \theta_i = \mathbf{v} \cdot \mathbf{v}_i$ .

**Quadratic Form** Let  $V$  be a finite dimensional vector space with dot product and  $L : V \rightarrow V$  be a self-adjoint, linear transformation on  $V$ . Then quadratic form associated with  $L$  is a  $\mathcal{L} : V \rightarrow \mathbb{R}$  defined by  $\mathcal{L}(v) = L(v) \cdot v$ .

**First Fundamental Form,  $\mathcal{I}_p$**  of surface  $S$  is the quadratic form associated with the identity transformation on  $S_p$ .

$$\mathcal{I}_p(v) = id(v) \cdot v = v \cdot v = \|v\|^2$$

**Second Fundamental Form  $\mathcal{S}_p$**  of surface  $S$  is the quadratic form associated with the Weingarten Map.

$$\mathcal{S}_p(v) = L_p(v) \cdot v$$

A quadratic forms is

**positive definite** if  $\mathcal{L}(v) > 0, \forall v \neq 0$ .

**negative definite** if  $\mathcal{L}(v) < 0, \forall v \neq 0$ .

**definite** if either positive definite or negative definite.

**indefinite** if not definite.

**positive semi-definite** if  $\mathcal{L}(v) \geq 0$ ,  $\forall v \neq 0$ .

**negative semi-definite** if  $\mathcal{L}(v) \leq 0$ ,  $\forall v \neq 0$ .

**semi-definite** if either positive semi-definite or negative semi-definite.

For example, first fundamental form of any surface is positive definite.

**Theorem 18.12.5.** *On each compact, oriented  $n$ -surface  $S$  in  $\mathbb{R}^{n+1}$  there exists a point  $p$  such that the second fundamental form at  $p$  is definite.*

*Proof.* Let  $S^n$  be an  $n$ -sphere containing  $S$  and  $S^n$  touches  $S$  at  $p$ . Then  $\mathcal{S}_p$  is definite.

Define  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  by  $g(x_1, x_2, \dots, x_n) = x_1^2 + x_2^2 + \dots + x_n^2$ . Since  $S$  is compact,  $g$  attains maximum at  $p \in S$ . By Lagrange multiplier theorem,  $\nabla g(p) = \lambda \nabla f(p) = \mu N(p)$  where  $\lambda, \mu \in \mathbb{R}$  and  $\mu = \pm \lambda \|\nabla f(p)\|$ . WLOG Suppose that  $\mu < 0$ . Then  $\mu = -|\mu| = -\|\mu N(p)\| = -\|\nabla g(p)\| = -2\|p\|$ . Thus,  $N(p) = \frac{1}{\mu} \nabla g(p) = \frac{-1}{\|p\|}(p, p)$ .

Let  $\mathbf{v} \in S_p$  such that  $\|\mathbf{v}\| = 1$ . Let  $\alpha : I \rightarrow S$  defined by  $\dot{\alpha}(t_0) = \mathbf{v}$ . Then  $g \circ \alpha(t_0) \geq g \circ \alpha(t)$ ,  $\forall t \in I$ .

$$\begin{aligned}
 0 &\geq \frac{d^2}{dt^2} (g \circ \alpha)(t_0) \\
 &= \frac{d}{dt} (\nabla g(\alpha(t)) \cdot \dot{\alpha}(t))(t_0) \\
 &= \frac{d}{dt} \left( 2\alpha(t) \cdot \frac{d\alpha}{dt}(t) \right)(t_0) \\
 &= 2\dot{\alpha}(t_0) \cdot \dot{\alpha}(t_0) + 2\alpha(t_0) \cdot \ddot{\alpha}(t_0) \\
 &= 2(\|\dot{\alpha}(t_0)\|^2 + (\alpha(t_0), \alpha(t_0)) \cdot \ddot{\alpha}(t_0)) \\
 &= 2(\|\mathbf{v}\|^2 + (p, p) \cdot \ddot{\alpha}(t_0)) \\
 &= 2(1 - \|p\| \mathbf{N}(p) \cdot \ddot{\alpha}(t_0)) \\
 &= 2(1 - \|p\| \kappa(\mathbf{v}))
 \end{aligned}$$

Thus,  $0 \geq 2 - 2\|p\| \kappa(\mathbf{v}) \implies 2\|p\| \kappa(\mathbf{v}) \geq 2 \implies \kappa(\mathbf{v}) \geq \frac{1}{\|p\|}$ ,  $\forall \mathbf{v} \in S_p$ . If  $S$  is oriented so that  $\mu > 0$ , then  $\kappa(\mathbf{v}) \leq -\frac{1}{\|p\|}$ .  $\square$

**Gauss Kronecker Curvature**,  $K(p)$  of surface  $S$  in  $\mathbb{R}^{n+1}$  at a point  $p$  is  $K(p) = \det L_p$ . It is equal to the product of the principal curvatures at  $p$ .

$$K(p) = k_1(p)k_2(p) \cdots k_n(p) \quad (18.35)$$

**Mean Curvature**  $H(p)$  of surface  $S$  in  $\mathbb{R}^{n+1}$  at a point  $p$  is the average value of principal curvatures at  $p$ .

$$H(p) = \frac{1}{n}(\kappa_1(p) + \kappa_2(p) + \cdots + \kappa_n(p)) \quad (18.36)$$

**Theorem 18.12.6.** *Let  $S$  be an  $n$ -surface in  $\mathbb{R}^{n+1}$ . Let  $p \in S$ . Let  $\mathbf{Z}$  be any non-zero normal vector field on  $S$  such that  $\mathbf{N} = \mathbf{Z}/\|\mathbf{Z}\|$  and let  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  be any basis for  $S_p$ . Then*

$$K(p) = \frac{(-1)^n \det \begin{pmatrix} \nabla_{v_1} \mathbf{Z} \\ \nabla_{v_2} \mathbf{Z} \\ \vdots \\ \nabla_{v_n} \mathbf{Z} \\ \mathbf{Z}(p) \end{pmatrix}}{\|\mathbf{Z}(p)\|^n \det \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_n \\ \mathbf{Z}(p) \end{pmatrix}} \quad (18.37)$$

*Proof.*  $\mathbf{Z} = \|\mathbf{Z}\|\mathbf{N}$ .

$$\begin{aligned} \det \begin{pmatrix} \nabla_{v_1} \mathbf{Z} \\ \nabla_{v_2} \mathbf{Z} \\ \vdots \\ \nabla_{v_n} \mathbf{Z} \\ \mathbf{Z}(p) \end{pmatrix} &= \det \begin{pmatrix} (\nabla_{v_1} \|\mathbf{Z}\|) \mathbf{N}(p) + \|\mathbf{Z}(p)\| \nabla_{v_1} \mathbf{N} \\ (\nabla_{v_2} \|\mathbf{Z}\|) \mathbf{N}(p) + \|\mathbf{Z}(p)\| \nabla_{v_2} \mathbf{N} \\ \vdots \\ (\nabla_{v_n} \|\mathbf{Z}\|) \mathbf{N}(p) + \|\mathbf{Z}(p)\| \nabla_{v_n} \mathbf{N} \\ \|\mathbf{Z}(p)\| \mathbf{N}(p) \end{pmatrix} \\ &= \|\mathbf{Z}(p)\|^n \det \begin{pmatrix} \nabla_{v_1} \mathbf{N} \\ \nabla_{v_2} \mathbf{N} \\ \vdots \\ \nabla_{v_n} \mathbf{N} \\ \mathbf{N}(p) \end{pmatrix} \\ &= (-1)^n \|\mathbf{Z}(p)\|^n \det \begin{pmatrix} L_p(\mathbf{v}_1) \\ L_p(\mathbf{v}_2) \\ \vdots \\ L_p(\mathbf{v}_n) \\ \mathbf{Z}(p) \end{pmatrix} \\ &= (-1)^n \|\mathbf{Z}(p)\|^n \det \begin{pmatrix} \cdots & 0 \\ & A^t & \vdots \\ 0 & \cdots & 0 & 1 \end{pmatrix} \det \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_n \\ \mathbf{Z}(p) \end{pmatrix} \\ &= (-1)^n \|\mathbf{Z}(p)\|^n \det A \det \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_n \\ \mathbf{Z}(p) \end{pmatrix} \end{aligned}$$

$$= (-1)^n \|\mathbf{Z}(p)\|^n K(p) \det \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_n \\ \mathbf{Z}(p) \end{pmatrix}$$

where  $A$  is the matrix for  $L_p$  with respect to the orthonormal basis  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  for  $S_p$  and  $A^t$  is its transpose.  $\square$

**local property** is a property of  $S$  which is valid in the neighbourhood of a particular point.

**global property** is a property of  $S$  which is valid everywhere on  $S$ .

**local theorem** is a theorem on local behaviour of  $S$ .

**global theorem** is a theorem on global behaviour of  $S$ .

**Theorem 18.12.7** (Global Characterisation of Surfaces with Definite Second Fundamental Form). *Let  $S$  be a compact, connected, oriented  $n$ -surface in  $\mathbb{R}^{n+1}$ . Then  $\forall p \in S, K(p) \neq 0 \iff \forall p \in S, \mathcal{S}_p$  is definite.*

*Proof.* If  $\mathcal{S}_p$  is definite then normal curvature  $\kappa(\mathbf{v}) = \mathcal{S}_p(\mathbf{v})$  is nonzero in every direction  $\mathbf{v} \in S_p$ . Thus, all principal curvatures are nonzero and their product Gauss Kronecker Curvatures  $K(p)$  is also nonzero.

Let  $\kappa_1 \leq \kappa_2 \leq \dots \leq \kappa_n$ . Let  $\mathcal{S}_{p_0}$  be definite. Since every compact, connected, oriented  $n$ -surface  $S$  has a point  $p_0$  at which its second fundamental form  $\mathcal{S}_{p_0}$  is definite. Suppose  $\mathcal{S}_{p_0}$  is positive definite. Then minimal principal curvature  $\kappa_1$  is positive at  $p_0$ . Clearly,  $\kappa_i$  is nowhere zero and continuous. Thus, every principal curvature  $\kappa_i$  are positive everywhere on  $S$ . Therefore, their product  $K(p)$  is positive everywhere on  $S$ . That is,  $\mathcal{S}_p$  is positive definite.

Suppose  $\mathcal{S}_{p_0}$  is negative definite. Then the maximal principal curvature  $\kappa_n(p)$  is negative at  $p_0$ . And  $\kappa_n$  is nonzero everywhere and continuous. Thus,  $\kappa_n$  is negative definite. Thus every principal curvature is negative definite. Thus Gauss Kronecker Curvature  $K(p)$  is either everywhere negative or positive depending on the parity on  $n$ . Therefore,  $K(p)$  is definite.  $\square$

## 18.14 Parameterized Surfaces

**tangent bundle** is the set  $T(U) = U \times \mathbb{R}^{n+1}$  where  $U$  is an open subset of  $\mathbb{R}^{n+1}$  such that  $\forall p \in U, \mathbb{R}_p^{n+1} \subset U \times \mathbb{R}^{n+1}$ .

**Definitions 18.14.1** (Differential of parametrisations). Let  $\varphi : U \rightarrow \mathbb{R}^{n+1}$  be parametrisation of a surface. Then its differential  $d\varphi : T(U) \rightarrow T(\mathbb{R}^{n+1})$  is given by

$$d\varphi(\mathbf{v}) = (\varphi \circ \alpha)(t_0) \quad (18.38)$$

**Note :**  $d\varphi$  is independent of the choice of  $\alpha$ .



**Definitions 18.14.2** (Differential of surface parametrisation). Let  $S$  be an  $n$ -surface in  $\mathbb{R}^{n+1}$ . Then  $d\varphi_p$  is the restriction of  $d\varphi$  to  $S_p$ .

There are two different meanings for  $d\varphi$  depending on its domain.

1.  $d\varphi : I \rightarrow \mathbb{R}$
2.  $d\varphi : I \times \mathbb{R} \rightarrow \mathbb{R}^2$

$$d\varphi(t, u) = (\varphi(t), d\varphi(t, u)) \quad (18.39)$$

Author could have used dark font for this second version,  $d\varphi$ . But he chose not to.

### 18.14.1 Parametrized Surfaces

**parametrised  $n$ -surface** is a smooth map  $\varphi : U \rightarrow \mathbb{R}^n$  where  $U$  is a connected and open subset of  $\mathbb{R}^n$ .

**regular surface** A parametrised  $n$ -surface is regular if  $d\varphi$  is non-singular.

**parametrised 2-sphere**  $\varphi : U \rightarrow \mathbb{R}^3$  defined by  $\varphi(\theta, \phi) = (r \cos \theta \sin \phi, r \sin \theta \sin \phi, r \cos \phi)$

**parametrised  $n$ -plane** Let  $L : \mathbb{R}^n \rightarrow \mathbb{R}^{n+k}$  be non-singular linear map. Then  $\varphi : U \rightarrow \mathbb{R}^{n+k}$  defined by  $\varphi(p) = L(p) + w$  is a parametrised  $n$ -plane through  $w$  in  $\mathbb{R}^{n+k}$ .

**cylinder over  $n$ -surface**  $\varphi$  Let  $\varphi : U \rightarrow \mathbb{R}^{n+k}$  be a parametrised  $n$ -surface in  $\mathbb{R}^{n+k}$  where  $U$  is an open subset of  $\mathbb{R}^n$ . The cylinder over  $\varphi$  is an  $(n+1)$ -surface.  $\tilde{\varphi} : U \times \mathbb{R} \rightarrow \mathbb{R}^{n+k+1}$  defined by

$$\tilde{\varphi}(u_1, u_2, \dots, u_{n+1}) = (\varphi(u_1, u_2, \dots, u_n), u_{n+1}) \quad (18.40)$$

**Surface of revolution** Let  $\alpha : I \rightarrow \mathbb{R}^2$  be a parametrised curve in  $\mathbb{R}^2$  whose image lies above  $x_1$  axis. The surface obtained by revolving  $\alpha(t) = (x_1(t), x_2(t))$  around  $x_1$  axis is an 2-surface in  $\mathbb{R}^3$ .

$$\varphi : I \times \mathbb{R} \rightarrow \mathbb{R}^3 \text{ defined by } \varphi(t, \theta) = (x_1(t), x_2(t) \cos \theta, x_2(t) \sin \theta) \quad (18.41)$$

**Torus in  $\mathbb{R}^3$**  Consider the circle  $\alpha(\phi) = (a + b \cos \phi, b \sin \phi)$  with center at  $(a, 0)$  with radius  $b$  ( $b > a$ ). Then the surface of revolution of  $\alpha$  about  $x_2$  axis,

$$\varphi(\phi, \theta) = (a + b \cos \phi \cos \theta, a + b \cos \phi \sin \theta, b \sin \phi) \quad (18.42)$$

**Torus in  $\mathbb{R}^4$**  is a parametrised 2-surface  $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^4$  defined by

$$\varphi(\theta, \phi) = \{(\cos \theta, \sin \theta, \cos \phi, \sin \phi) : \theta, \phi \in \mathbb{R}\} \quad (18.43)$$

**Möbius Band** is a parametrised 2-surface  $\varphi : I \times \mathbb{R} \rightarrow \mathbb{R}^3$  defined by

$$\varphi(t, \theta) = \left( \left(1 + t \cos \frac{\theta}{2}\right) \cos \theta, \left(1 + t \cos \frac{\theta}{2}\right) \sin \theta, t \sin \frac{\theta}{2} \right) \quad (18.44)$$

### 18.14.2 Vector Field along $\varphi$

**Definitions 18.14.3.** Let  $\varphi : U \rightarrow \mathbb{R}^{n+k}$  be a smooth function. A vector field along  $\varphi$  is a map  $\mathbf{X}$  which assigns a vector  $\mathbf{X}(p) \in \mathbb{R}_{\varphi(p)}^{n+k}$  to each point  $p \in U$ .

**Definitions 18.14.4.** A vector field along  $\varphi$  is a tangent vector field along  $\varphi$  if there exists a vector field  $\mathbf{Y}$  on  $U$  such that  $\mathbf{X}(p) = d\varphi_p(\mathbf{Y}(p))$ .

**Definitions 18.14.5.** Let  $\varphi : U \rightarrow \mathbb{R}^{n+k}$  be a smooth map where  $U$  is an open subset of  $\mathbb{R}^n$ . The coordinate vector fields along  $\varphi$  are the tangent fields  $\mathbf{E}_i$  defined by

$$\mathbf{E}_i(p) = d\varphi_p(p, 0, 0, \dots, 1, \dots, 0) \quad (18.45)$$

**Theorem 18.14.1** (Basis for Tangent Space Image  $d\varphi_p$ ). *Let  $\varphi : U \rightarrow \mathbb{R}^{n+k}$  be smooth. Then the set of all coordinate vector fields  $\{\mathbf{E}_i(p) : i = 1, 2, \dots, n\}$  form a basis for the tangent space image  $d\varphi_p$ .*

**Definitions 18.14.6.** Let  $\varphi : U \rightarrow \mathbb{R}^{n+k}$  be smooth. And  $\mathbf{X}$  be a smooth vector field along  $\varphi$ . Let  $p \in U$  and  $\mathbf{v} \in \mathbb{R}_p^n$ . The derivative  $\nabla_{\mathbf{v}}\mathbf{X} \in \mathbb{R}_{\varphi(p)}^{n+k}$  defined by

$$\nabla_{\mathbf{v}}\mathbf{X} = \left( \varphi(p), \left( \frac{d}{dt} X \circ \alpha \right) (t_0) \right) \quad (18.46)$$

**Definitions 18.14.7** (orientation vector field). Let  $\varphi : U \rightarrow \mathbb{R}^{n+1}$  be smooth parametrised  $n$ -surface in  $\mathbb{R}^{n+1}$ . Let  $p \in U$ . An orientation vector field along  $\varphi$  is the unique vector field  $\mathbf{N}$  along  $\varphi$  such that  $\mathbf{N}(p)$  are unit vectors and  $\mathbf{N}(p) \perp d\varphi_p(\mathbf{v})$ . And  $\mathbf{N}$  is consistently oriented if the following determinant is positive.

$$\det \begin{pmatrix} \mathbf{E}_1(p) \\ \mathbf{E}_2(p) \\ \vdots \\ \mathbf{E}_n(p) \\ \mathbf{N}(p) \end{pmatrix} > 0 \quad (18.47)$$

### 18.14.3 Weingarten Map

**Definitions 18.14.8** (Weingarten map). The Weingarten map is the linear function  $L_p : \text{Image } d\varphi_p \rightarrow \text{Image } d\varphi_p$  defined by  $L_p(d\varphi(\mathbf{v})) = -\nabla_{\mathbf{v}}\mathbf{N}$ .

**$L_p$  is self-adjoint**  $L_p(\mathbf{v}) \cdot \mathbf{w} = \mathbf{v} \cdot L_p(\mathbf{w})$

**principal curvatures** are the eigen values of  $L_p$

**principal curvature directions** are the unit eigen vectors of  $L_p$

**Gauss-Kronecker Curvature** is the determinant.

**Mean Curvature** is  $1/n$  times its trace.

## Subject 19

# ME800402 Algorithmic Graph Theory

### 19.1 Networks

#### 19.1.1 An Introduction to Networks

**Definitions 19.1.1.** A **network**  $N$  is a digraph  $D$  with two special vertices source  $s$  and sink  $t$  together with a capacity function  $c : E(D) \rightarrow \mathbb{Z}$  such that for every arc  $a = (u, v)$  of the digraph,  $c(u, v)$  is non-negative.

*Remark.* Mathematical Modeling using Network,

1. There is no restriction on indegree/outdegree of source/sink vertices of the digraph  $D$  of a network  $N$ .
2. Applications of Network : Transportation problem.

$c(u, v)$  is the capacity of the arc  $(u, v)$  of  $D$

$N^+(x) = \{y \in V(D) : (x, y) \in E(D)\}$  is the out-neighbourhood of  $x$ .

$N^-(x) = \{y \in V(D) : (y, x) \in E(D)\}$  is the in-neighbourhood of  $x$ .

**Definitions 19.1.2.** A **flow  $f$  in a network  $N$**  is function  $f : E(D) \rightarrow \mathbb{Z}$  such that 1. each edge satisfies capacity constraint and 2. each vertex except source and sink satisfies conservation equation.

**capacity constraint**

$$0 \leq f(a) \leq c(a) \text{ for every arc } a \in V(D) \quad (19.1)$$

**conservation equation**

$$\sum_{y \in N^+(x)} f(x, y) = \sum_{y \in N^-(x)} f(y, x), \quad \forall \text{ vertex } x \in V(D) - \{s, t\} \quad (19.2)$$

**net flow out of  $x$**

$$\sum_{y \in N^+(x)} f(x, y) - \sum_{y \in N^-(x)} f(y, x)$$

**net flow into  $x$**

$$\sum_{y \in N^-(x)} f(y, x) - \sum_{y \in N^+(x)} f(x, y)$$

**Definitions 19.1.3.** The **flow  $f$  in a network  $N$**  is the net flow out of source  $s$ .

*Remark.* 1. net flow out of/into  $x \in V(D) - \{s, t\}$  is zero.

2. Without loss of generality<sup>1</sup>, underlying digraph is always assymetric.

$$(X, Y) = \{(x, y) \in E(D) : x \in X, y \in Y\}.$$

Let  $X, Y$  be non-empty subsets of  $V(D)$  such that  $X, Y$  are disjoint. Then  $(X, Y)$  is the set of all arcs from  $X$  to  $Y$ .

**flow from  $X$  to  $Y$**  is the sum of flow on each arc in  $(X, Y)$

$$f(X, Y) = \sum_{(x, y) \in (X, Y)} f(x, y) \quad (19.3)$$

**capacity of the partition  $(X, Y)$**  is the total capacity of arcs in  $(X, Y)$

$$c(X, Y) = \sum_{(x, y) \in (X, Y)} c(x, y) \quad (19.4)$$

**cut** Let  $P \subset V(D)$  such that  $s \in P$  and  $t \notin P$  and  $\bar{P} = V(D) - P$ , then  $(P, \bar{P})$  is a cut.

**flow from  $P$  to  $\bar{P}$**  is the sum of flow on each arc in  $(P, \bar{P})$ .

$$f(P, \bar{P}) = \sum_{(x, y) \in (P, \bar{P})} f(x, y) \quad (19.5)$$

**flow from  $\bar{P}$  to  $P$**  is the sum of flow on each arc in  $(\bar{P}, P)$

$$f(\bar{P}, P) = \sum_{(x, y) \in (\bar{P}, P)} f(x, y) \quad (19.6)$$

**capacity of the cut  $(P, \bar{P})$**  is the total capacity of the arcs in  $(P, \bar{P})$

$$c(P, \bar{P}) = \sum_{(x, y) \in (P, \bar{P})} c(x, y) \quad (19.7)$$

**Theorem 19.1.1.** For any cut  $(P, \bar{P})$ , the flow in  $N$  is  $f(N) = f(P, \bar{P}) - f(\bar{P}, P)$ .

*Synopsis.* The net flow out of source  $s$  is the flow  $f(N)$  in the network  $N$ . Let  $(P, \bar{P})$  be a cut of  $N$ , then  $s \in P$  and  $t \notin P$ . Suppose  $P = \{s\}$ , then the theorem is true. Suppose  $P$  is not singleton, then for each vertex  $x \in P$ ,  $x \neq s$ , the net flow out of  $x$  is zero by flow conservation equation. And flow between vertices in  $P$  cancels out each other. Thus adding net flow out of each vertex in  $P$ , will be same as the net flow out of source which is the flow in the network,  $f(N)$ .

<sup>1</sup>If underlying digraph of a network is symmetric, then by replacing an arc  $(u, v)$  with a new vertex  $w$  and two arcs  $(u, w), (w, v)$  gives an assymetric digraph.[Gray Chartrand, ]pp.131

*Proof.*

$$\text{Flow, } f = \sum_{y \in N^+(s)} f(s, y) - \sum_{y \in N^-(s)} f(y, s) \quad (19.8)$$

By conservation equation, we have  $\forall x \in P, x \neq s$ ,

$$\sum_{y \in N^+(x)} f(x, y) - \sum_{y \in N^-(x)} f(y, x) = 0 \quad (19.9)$$

By above equations,

$$\begin{aligned} \text{Flow, } f &= \sum_{x \in P} \sum_{y \in N^+(x)} f(x, y) - \sum_{x \in P} \sum_{y \in N^-(x)} f(y, x) \\ &= \sum_{(x, y) \in (P, \bar{P})} f(x, y) - \sum_{(y, x) \in (\bar{P}, P)} f(y, x) \end{aligned} \quad (19.10)$$

□

**Corollary 19.1.1.1.** *Flow cannot exceed the capacity of any cut  $(P, \bar{P})$ . Further,  $f(N) \leq \min c(P, \bar{P})$ .*

*Synopsis.* Let  $(P, \bar{P})$  be a cut in network  $N$ , then by theorem the flow  $f(N) = \text{flow from } P \text{ to } \bar{P} - \text{flow from } \bar{P} \text{ to } P$ . Since the flow from  $\bar{P}$  to  $P$  is non-negative,  $f(N) \leq \text{flow from } P \text{ to } \bar{P}$ . Clearly,  $f(x, y) \leq c(x, y)$  by the capacity constraint. Thus  $f(N) \leq f(P, \bar{P}) \leq c(P, \bar{P}) \leq \min c(P, \bar{P})$ .

*Proof.*

$$\begin{aligned} f(N) &= \sum_{(x, y) \in (P, \bar{P})} f(x, y) - \sum_{(y, x) \in (\bar{P}, P)} f(y, x) \\ &\leq \sum_{(x, y) \in (P, \bar{P})} f(x, y) = f(P, \bar{P}) \\ &\leq \sum_{(x, y) \in (P, \bar{P})} c(x, y) = c(P, \bar{P}), \quad \because \forall x, y \in V(D), f(x, y) \leq c(x, y) \\ &\leq \min c(P, \bar{P}) \end{aligned}$$

□

**Corollary 19.1.1.2.** *In a network  $N$  flow is the net flow into the sink of  $N$ .*

*Synopsis.* Let  $\bar{P} = \{t\}$ , then by theorem  $f(N)$  is the net flow into the sink.

*Proof.* Suppose  $P = V(D) - \{t\}$ . Then by theorem, we have

$$\begin{aligned} f(N) &= \sum_{(x, y) \in (P, \bar{P})} f(x, y) - \sum_{(y, x) \in (\bar{P}, P)} f(y, x) \\ &= \sum_{x \in N^-(t)} f(x, t) - \sum_{x \in N^+(t)} f(t, x) \end{aligned}$$

□

*Remark.* Exercise 5.1

4. Let  $N$  be a network with underlying digraph  $D$  which has a vertex  $v \in V(D) - \{s, t\}$  with zero indegree. Clearly the flow into  $v$  is zero. Thus flow out of  $v$  is also zero by flow conservation equation. Let  $N'$  be the network obtained from  $N$  by deleting the vertex  $v$ . Then  $f(N) = f(N')$ .

### 19.1.2 The Max-Flow Min-Cut Theorem

**maximum flow** A flow  $f$  in network  $N$  is maximum flow in  $N$ , if  $f(N) \geq f'(N)$  for each flow  $f'$  in  $N$ .

**minimum cut** A cut  $(P, \bar{P})$  in network  $N$  is minimum cut of  $N$ , if  $c(P, \bar{P}) \leq c(X, \bar{X})$  for each cut  $(X, \bar{X})$  in  $N$ .

**$f$ -unsaturated** Let  $f$  be a flow in network  $N$  with underlying digraph  $D$ , and  $Q = u_0, a_1, u_1, a_2, \dots, u_{n-1}, a_n, u_n$  be a semipath in  $D$  such that every forward arc  $a_i = (u_{i-1}, u_i)$  has flow not upto its capacity,  $f(a_i) < c(a_i)$  and every reverse arc  $a_i = (u_i, u_{i-1})$  has some positive flow in it,  $f(a_i) > 0$

**$f$ -augmenting semipath** Let  $f$  be a flow in a network  $N$  with underlying digraph  $D$ . Suppose semipath  $Q = s, a_1, u_1, a_2, \dots, u_{n-1}, a_n, t$  (from source to sink) is  $f$ -unsaturated, then  $Q$  is an  $f$ -augmenting semipath.

**Theorem 19.1.2.** Let  $f$  be a flow in a network  $N$  with underlying digraph  $D$ . The flow  $f$  is maximum in  $N$  iff there is no  $f$ -augmenting semipath in  $D$ .

*Synopsis.* Suppose  $Q$  is an  $f$ -augmenting semipath in  $D$ , then there exists a flow  $f^*$  in  $N$  such that  $f(N) + \Delta = f^*(N)$ . Therefore,  $f$  is not a maximum flow in  $N$ . Suppose there is no  $f$ -augmenting semipath in  $D$ , then there exists a cut  $(P, \bar{P})$  such that  $f(a) = c(a) \forall a \in (P, \bar{P})$  and  $f(a) = 0 \forall a \in (\bar{P}, P)$ . Suppose  $f^*$  in a maximum flow in  $N$ , then  $f(N) \leq f^*(N) \leq c(P, \bar{P}) = f(N)$ .

*Proof.* Let  $f$  be a flow in a network  $N$  with underlying digraph  $D$  and  $Q = s, a_1, u_1, a_2, u_2, \dots, u_{n-1}, a_n, t$  be an  $f$ -augmenting semipath in  $D$ .

$$\text{define } \Delta_i = \begin{cases} c(a_i) - f(a_i) & \text{for every forward arc } a_i \in Q, \\ f(a_i) & \text{for every reverse arc } a_i \in Q, \end{cases}$$

Define  $\Delta = \min\{\Delta_i\}$ . Also define  $f^* : E(D) \rightarrow \mathbb{Z}$  such that

$$f^*(a_i) = \begin{cases} f(a_i) + \Delta, & \text{for every forward arc } a_i \in Q, \\ f(a_i) - \Delta, & \text{for every reverse arc } a_i \in Q, \\ f(a_i), & \text{for every arc of } D \text{ which are not in } Q. \end{cases}$$

Since  $Q$  is an  $f$ -augmenting semipath in  $D$ ,  $\Delta > 0$  and  $f(N) + \Delta = f^*(N)$ .

Clearly  $f(N) < f^*(N)$ , and it is enough to show that  $f^*$  is a flow in  $N$ .  $f^*$  is a flow if it satisfies 1. capacity constraint and 2. conservation equation. For any arc  $a_i \notin Q$ ,  $f^*(a_i) = f(a_i) \leq c(a_i)$ . Suppose  $a_i \in Q$ . If  $a_i = (u_{i-1}, u_i)$ ,  $a_i$  is a forward arc and we have  $f^*(a_i) = f(a_i) + \Delta \leq f(a_i) + \Delta_i = f(a_i) + c(a_i) - f(a_i) = c(a_i)$ . If  $a_i = (u_i, u_{i-1})$ , then  $a_i$  is a reverse arc and we have  $f^*(a_i) = \Delta \leq \min\{\Delta_i\} = \Delta_i = c(a_i)$ . Thus  $f^*$  satisfies capacity

constraint on every arc of  $D$ .

Let  $x \in V(D) - \{s, t\}$ . Suppose  $x \notin Q$ ,

$$\begin{aligned} \text{Net flow out of } x &= \sum_{y \in N^+(x)} f^*(x, y) - \sum_{y \in N^-(x)} f^*(y, x) \\ &= \sum_{y \in N^+(x)} f(x, y) - \sum_{y \in N^-(x)} f(y, x) \\ &= 0 \end{aligned}$$

Suppose  $x = u_i \in Q$ , then  $Q$  has two arcs having vertex  $x$  say,  $a_{i-1}$ , and  $a_i$ . There are four possibilities for these two arcs,

1. Both  $a_{i-1}$ ,  $a_i$  are forward arcs.
2. Arc  $a_{i-1}$  is forward, but arc  $a_i$  is reverse.
3. Arc  $a_{i-1}$  is reverse, but arc  $a_i$  is forward.
4. Both  $a_{i-1}$ ,  $a_i$  are reverse arcs.

**Case 1**  $a_{i-1} = (u_{i-1}, u_i)$  and  $a_i = (u_i, u_{i+1})$ .

$$\begin{aligned} \text{Net flow out of } x &= \sum_{y \in N^+(x)} f^*(x, y) - \sum_{y \in N^-(x)} f^*(y, x) \\ &= \sum_{\substack{y \in N^+(x) \\ y \neq u_{i+1}}} f^*(x, y) + f^*(u_i, u_{i+1}) - \left( \sum_{\substack{y \in N^-(x) \\ y \neq u_{i-1}}} f^*(y, x) + f^*(u_{i-1}, u_i) \right) \\ &= \sum_{\substack{y \in N^+(x) \\ y \neq u_{i+1}}} f(x, y) + f(u_i, u_{i+1}) + \Delta - \left( \sum_{\substack{y \in N^-(x) \\ y \neq u_{i-1}}} f(y, x) + f(u_{i-1}, u_i) \right) - \Delta \\ &= \sum_{y \in N^+(x)} f(x, y) - \sum_{y \in N^-(x)} f(y, x) \\ &= 0 \end{aligned}$$

**Case 2**  $a_{i-1} = (u_{i-1}, u_i)$  and  $a_i = (u_{i+1}, u_i)$ .

$$\begin{aligned} \text{Net flow out of } x &= \sum_{y \in N^+(x)} f^*(x, y) - \sum_{y \in N^-(x)} f^*(y, x) \\ &= \sum_{\substack{y \in N^+(x) \\ y \neq u_{i+1}, u_{i-1}}} f^*(x, y) + f^*(u_i, u_{i+1}) + f^*(u_i, u_{i-1}) - \sum_{y \in N^-(x)} f^*(y, x) \\ &= \sum_{\substack{y \in N^+(x) \\ y \neq u_{i+1}, u_{i-1}}} f(x, y) + f(u_i, u_{i+1}) + \Delta + f(u_i, u_{i-1}) - \Delta - \sum_{y \in N^-(x)} f(y, x) \\ &= \sum_{y \in N^+(x)} f(x, y) - \sum_{y \in N^-(x)} f(y, x) \\ &= 0 \end{aligned}$$

**Case 3**  $a_{i-1} = (u_i, u_{i-1})$  and  $a_i = (u_i, u_{i+1})$ .

$$\begin{aligned}
 \text{Net flow out of } x &= \sum_{y \in N^+(x)} f^*(x, y) - \sum_{y \in N^-(x)} f^*(y, x) \\
 &= \sum_{y \in N^+(x)} f^*(x, y) - \left( \sum_{\substack{y \in N^-(x) \\ y \neq u_{i-1}, u_{i+1}}} f^*(y, x) + f^*(u_{i-1}, u_i) + f^*(u_{i+1}, u_i) \right) \\
 &= \sum_{y \in N^+(x)} f^*(x, y) - \left( \sum_{\substack{y \in N^-(x) \\ y \neq u_{i-1}, u_{i+1}}} f(y, x) + f(u_{i-1}, u_i) + \Delta + f(u_{i+1}, u_i) - \Delta \right) \\
 &= \sum_{y \in N^+(x)} f(x, y) - \sum_{y \in N^-(x)} f(y, x) \\
 &= 0
 \end{aligned}$$

**Case 4**  $a_{i-1} = (u_i, u_{i-1})$  and  $a_i = (u_{i+1}, u_i)$ .

$$\begin{aligned}
 \text{Net flow out of } x &= \sum_{y \in N^+(x)} f^*(x, y) - \sum_{y \in N^-(x)} f^*(y, x) \\
 &= \sum_{\substack{y \in N^+(x) \\ y \neq u_{i-1}}} f^*(x, y) + f^*(u_i, u_{i-1}) - \left( \sum_{\substack{y \in N^-(x) \\ y \neq u_{i+1}}} f^*(y, x) + f^*(u_{i+1}, u_i) \right) \\
 &= \sum_{\substack{y \in N^+(x) \\ y \neq u_{i-1}}} f(x, y) + f(u_i, u_{i-1}) - \Delta - \left( \sum_{\substack{y \in N^-(x) \\ y \neq u_{i+1}}} f(y, x) + f(u_{i+1}, u_i) \right) + \Delta \\
 &= \sum_{y \in N^+(x)} f(x, y) - \sum_{y \in N^-(x)} f(y, x) \\
 &= 0
 \end{aligned}$$

Therefore,  $f^*$  is a flow on  $N$ . We have  $f(N) < f^*(N)$ . Thus  $f$  is not maximum flow in  $N$  due to the existence of an  $f$ -augmenting semipath in  $D$ .

Conversely, assume that there is no  $f$ -augmenting semipath in  $D$ . Now, we construct a cut  $(P, \bar{P})$  of  $N$ . Let  $P$  be the set of all vertices  $x \in V(D)$  such that there is an  $f$ -unsaturated  $s - x$  semipath in  $D$ . Trivially,  $s \in P$ . And  $t \notin P$  since there are no  $f$ -augmenting semipath in  $D$ .<sup>2</sup> Clearly,  $(P, \bar{P})$  is a cut of the network  $N$ .

We claim that  $c(P, \bar{P}) = f(N)$ . Suppose there is a forward arc  $(x, y) \in (P, \bar{P})$ , then flow in it is saturated. If  $f(x, y) < c(x, y)$ , then there is an  $f$ -unsaturated  $s - y$  semipath in  $D$ . ie,  $s - x$  semipath + arc  $(x, y)$ . Thus every forward arc  $(x, y) \in (P, \bar{P})$  is saturated. Suppose there is a reverse arc  $(y, x) \in$

---

<sup>2</sup>An  $f$ -augmenting semipath is an  $f$ -unsaturated  $s - t$  semipath in  $D$ .



$(\bar{P}, P)$ , then there is no flow in it (saturated reversed arc). If  $f(y, x) > 0$ , then there is an  $f$ -unsaturated  $s - y$  semipath in  $D$ . ie,  $s - x$  semipath + arc  $(y, x)$ . Thus every reverse arc  $(y, x) \in (\bar{P}, P)$  is saturated. And we have,

$$\begin{aligned} \sum_{(x,y) \in (P, \bar{P})} f(x, y) &= \sum_{(x,y) \in (P, \bar{P})} c(x, y) \\ \sum_{(y,x) \in (\bar{P}, P)} f(y, x) &= 0 \\ f(N) &= \sum_{(x,y) \in (P, \bar{P})} f(x, y) - \sum_{(y,x) \in (\bar{P}, P)} f(y, x) \\ &= \sum_{(x,y) \in (P, \bar{P})} c(x, y) \\ &= c(P, \bar{P}) \end{aligned}$$

Suppose  $f^*$  is maximum flow in network  $N$  and  $(X, \bar{X})$  is minimum cut of  $N$ . Then  $f(N) \leq f^*(N)$ . Thus we have,  $f(N) \leq f^*(N) \leq c(X, \bar{X}) \leq c(P, \bar{P}) = f(N)$ . Therefore,  $f(N) = f^*(N)$ . ie, the flow  $f$  is maximum in network  $N$  if there are no  $f$ -augmenting semipaths in  $D$ .  $\square$

**Theorem 19.1.3** (maximum-flow, min-cut). *In every network, the value of maximum flow equals capacity of minimum cut.*

*Proof.* Suppose flow  $f$  in network  $N$  is maximum, then by previous theorem there is no  $f$ -augmenting semipath in  $D$ . And  $f(N) \leq c(X, \bar{X})$  for any cut  $(X, \bar{X})$  in  $N$ . We can construct a cut  $(P, \bar{P})$  in  $N$  such that  $f(N) = c(P, \bar{P})$ . Let  $P$  be the set of all vertices  $x$  in  $D$  such that there is an  $f$ -unsaturated  $s - x$  semipath in  $D$ . Clearly  $s \in P$  and  $t \notin P$ . Also  $f(P, \bar{P}) = c(P, \bar{P})$  and  $f(\bar{P}, P) = 0$ . Then the cut  $(P, \bar{P})$  is minimum cut of  $N$ . Suppose there is a cut  $(X, \bar{X})$  such that  $c(X, \bar{X}) < c(P, \bar{P})$ . Then  $f(N) = f(P, \bar{P}) - f(\bar{P}, P) = c(P, \bar{P}) < c(X, \bar{X})$  which is a contradiction. Therefore, the value of maximum flow equals capacity of minimum cut.  $\square$

*Remark.* Exercise 5.2

1. Suppose  $(X, \bar{X})$  is a cut of  $N$  such that  $f(a) = c(a)$ ,  $\forall a \in (X, \bar{X})$  and  $f(a) = 0$ ,  $\forall a \in (\bar{X}, X)$ . By the definition of cut,  $s \in X$  and  $t \in \bar{X}$ . Thus there is no  $f$ -augmenting semipath in  $D$ . Suppose there is an  $f$ -augmenting semipath  $Q$  in  $D$ , then there is either (a) a forward arc  $(x, y) \in (X, \bar{X})$  such that  $f(x, y) < c(x, y)$  or (b) a reverse arc  $(y, x) \in (\bar{X}, X)$  such that  $f(y, x) > 0$  which is a contradiction. Therefore, the flow  $f(N)$  is maximum and the given cut  $(X, \bar{X})$  is minimum as shown in the proof of the maximum-flow min-cut theorem.
3. The algorithm suggested in the hint of this exercise won't work if two subnetworks have a common arc such that the direction of flow in which is not consistent. Suppose, the generalized network is not supposed to have any common arcs. Then construct subnetworks for each pair  $(s, t)$  with all those arcs which are on some  $s - t$  semipath. Define subnetwork capacity function  $c'(a) = c(a)$  for every arc in  $N'$ .

Let  $N$  be a generalized network with set of sources  $S$  and set of sinks  $T$ . A flow in  $N$  is maximum if there is not  $f$ -augmenting  $s - t$  semipath for each pair  $(s, t) \in S \times T$ .

### 19.1.3 A max-flow min-cut algorithm

**Theorem 19.1.4.** *Let  $N$  be a network with underlying digraph  $D$ , source  $s$ , sink  $t$ , capacity function  $c$  and flow  $f$ . Let  $D'$  be the digraph with same vertex set as  $D$  and arc set defined by  $E(D') = \{(x, y) : (x, y) \in E(D), c(x, y) > f(x, y) \text{ or } (y, x) \in E(D), f(y, x) > 0\}$ . ie,  $D'$  has only the unsaturated arcs of  $D$ . Then  $D'$  has an  $s - t$  directed path iff  $D$  has an  $f$ -augmenting semipath. Moreover, shortest  $s - t$  path in  $D'$  has the same length as shortest  $f$ -augmenting semipath in  $D$ .*

*Synopsis.* Each directed  $s - t$  path in  $D'$  has respective  $f$ -augmenting semipath in  $D$  and vice versa. Clearly, they have the same length.

*Proof.* Let  $N$  be a network with underlying digraph  $D$ , capacity  $c$  and flow  $f$ . Let  $D'$  be the digraph with vertex set  $V(D') = V(D)$  and arc set  $E(D') = \{(x, y) : \text{either } (x, y) \text{ or } (y, x) \text{ is unsaturated in } N\}$ .

Suppose  $D'$  has a directed  $s - t$  path  $Q' : s, u_1, u_2, \dots, u_{n-1}, t$ . Then by the construction of  $D'$ , for each  $u_i \in Q$ , there exists an  $f$  unsaturated arc  $a_i$  in  $D$ . ie, either forward arc  $a_i = (u_{k-1}, u_k)$  such that  $f(u_{k-1}, u_k) < c(u_{k-1}, u_k)$  or reverse arc  $a_i = (u_k, u_{k-1})$  such that  $f(u_k, u_{k-1}) > 0$ . Therefore, we have an  $s - t$  semipath  $Q : s, a_1, u_1, a_2, \dots, u_{n-1}, a_n, t$  in  $D$  such that  $Q$  is an  $f$ -augmenting semipath since every arc in  $Q$  is  $f$ -unsaturated. Clearly,  $Q, Q'$  are of the same length.

Conversely, suppose that the digraph  $D$  has an  $f$ -augmenting semipath  $Q : s, a_1, u_1, a_2, \dots, u_{n-1}, a_n, t$ . Then each arc  $a_i \in Q$  are  $f$ -unsaturated and by the construction of  $D'$ , there exists a directed  $s - t$  path  $Q' = s, u_1, u_2, \dots, u_{n-1}, t$  in  $D'$ . And  $Q, Q'$  are of the same length.

There is a one-one correspondence between the directed  $s - t$  paths in  $D'$  and  $f$ -augmenting semipaths in  $D$ . Clearly, they have the same length. Thus shortest directed  $s - t$  path in  $D$  and shortest  $f$ -augmenting semipath in  $D'$  are of the same length.  $\square$

**saturation arc** of  $N$  with respect to the flow  $f$  is an arc  $a_j$  in an  $f$ -augmenting semipath  $Q$  with  $\Delta_j = \Delta$ .

**augmentation path** is an  $f$ -augmenting semipath  $Q$  in  $D$ .

**Algorithm 19.1.1** (max-flow min-cut). *An algorithm to find maximum flow and minimum cut of a network  $N$  with underlying digraph  $D$ , source  $s$ , sink  $t$ , capacity function  $c$  and initial flow  $f$ .*

1. Construct digraph  $D'$  with vertex set  $V(D') = V(D)$  and arc set  $E(D') = \{(x, y) : (x, y) \in E(D) \& f(x, y) < c(x, y) \text{ or } (y, x) \in E(D) \& f(y, x) > 0\}$

2. Find (shortest)  $s-t$  directed path in  $D'$  using Moore's breadth first search(BFS) algorithm. If  $D'$  doesn't have an  $s-t$  path, then proceed to step 5. Otherwise, let  $Q' : s, u_1, u_2, \dots, u_{n-1}, t$  be a (shortest)  $s-t$  path in  $D'$ .
3. Let  $Q : s, a_1, u_1, a_2, \dots, u_{n-1}, a_n, t$  be the respective semipath in  $D$  such that  $f(a_j) < c(a_j)$  for forward arcs and  $f(a_i) > 0$  for reverse arcs. Let  $\Delta_j = c(a_j) - f(a_j)$  for forward arcs and  $\Delta_j = f(a_j)$  for reverse arcs. And let  $\Delta = \min\{\Delta_j\}$ . And augment flow  $f$  by  $\Delta$  ie,  $f(a_j) \leftarrow f(a_j) + \Delta$  for forward arcs and  $f(a_j) \leftarrow f(a_j) - \Delta$  for reverse arcs.
4. Goto step 1 (Proceed with new flow  $f$  and find whether there are any directed  $s-t$  paths in  $D'$ . If any, augment the flow along the new augmentation path  $Q$  by saturating the flow along the saturation arc.)
5. There is no  $s-t$  directed path in  $D'$ . Thus there is no  $f$ -augmenting semipath in  $D$ . Therefore the flow  $f$  in  $N$  is maximum. Let  $P$  be the set of all vertices in  $D'$  with non-zero breadth first index(bfi) from Moore's BFS algorithm applied in step 2.  $(P, \bar{P})$  is minimum cut of  $N$ .

*Remark.* Validity of the algorithm is proved in the previous theorem.

### 19.1.5 Connectivity and Edge-Connectivity

**edge cutset** is the set  $U$  subset of  $E(G)$  such that  $G - U$  is disconnected.

**vertex cutset** is the set  $S$  subset of  $V(G)$  such that  $G - S$  is disconnected.

**edge connectivity**  $\lambda(G)$  is the minimum cardinality of all edge cutsets of  $G$ .

**connectivity**  $\kappa(G)$  is the minimum cardinality of all vertex cutsets of  $G$ .

**Theorem 19.1.5.** For every graph  $G$ ,  $\kappa(G) \leq \lambda(G) \leq \delta(G)$

*Proof.* Suppose graph  $G$  is disconnected then  $\kappa(G) = \lambda(G) = \delta(G) = 0$ . Let  $G$  be a connected graph. Then  $G$  has at least one vertex  $v$  with degree  $\delta(G)$ . Therefore  $\lambda(G) \leq \delta(G)$  since edges incident with  $v$  form an edge cutset of  $G$  and  $\lambda(G)$  is the cardinality of all edge cutsets.

Let  $G$  be a graph with edge connectivity  $\lambda(G) = c$ . Let  $U$  be a edge cutset with cardinality  $c$  and let edge  $uv \in U$ . Construct a set of vertices  $S \subset V(G)$  such that ( $S$  is of minimal cardinality and) for each edge in  $U$  other  $uv$ ,  $S$  has a vertex incident with it. Cardinality of  $S$  is atmost  $c - 1$ , since we can select one vertex each for each edge in  $U$  other than  $u, v$ . If  $G - S$  is a disconnected graph, then  $\kappa(G) < \lambda(G)$ . Suppose  $G - S$  is a connected graph, then delete a non-pendent vertex  $u$  or  $v$  from  $G - S$ , say  $v$ . Since  $G - S$  is a connected graph with a singleton edge cutset,  $\{uv\}$ . We have a vertex cutset  $S \cup \{v\}$  of  $G$ . Therefore,  $\kappa(G) \leq c = \lambda(G)$ .  $\square$

**Theorem 19.1.6.** If  $G$  is a graph of diameter 2, then  $\lambda(G) = \delta(G)$

**$n$ -edge connected**  $G$  is  $n$ -edge connected if  $\lambda(G) \geq n$ .

**$n$  connected**  $G$  is  $n$ -connected if  $\kappa(G) \geq n$ .

**Theorem 19.1.7.** Let  $G$  be a graph of order  $p$  and  $n$  be an integer such that  $1 \leq n \leq p-1$ . If  $\delta(G) \geq \frac{p+n-2}{2}$ , then  $G$  is  $n$ -connected.

**connection number**  $c(G)$  is the smallest integer such that  $2 \leq c(G) \leq p$  and every subgraph of order  $n$  in  $G$  is connected.

**$l$ -connectivity**  $\kappa_l(G)$  is minimum number of vertices whose removal will produce a disconnected graph with at least  $l$  components or a graph with fewer than  $l$  vertices.

**$(n, l)$ -connected** A graph  $G$  is  $(n, l)$ -connected if  $\kappa_l(G) \geq n$ .

*Remark.* Exercises 5.5

1.  $\lambda(K_{m,n}) = \kappa(K_{m,n}) = m$
8.  $c(K_p) = 2$ ,  $c(K_{m,n}) = n+1$ ,  $c(C_p) = p-1$   
Every two vertices of complete graph of order  $p$  are adjacent. For complete bi-partite graph  $K_{m,n}$  such that  $1 \leq m \leq n$ , there exists a totally disconnected subgraph of order  $n$ . Therefore  $c(K_{m,n}) \geq n+1$ . And with  $n+1$  vertices, both partitions have at least two vertices each and therefore the graph is connected and  $c(K_{m,n}) \leq n+1$ . For cycle  $C_p$ , any subgraph is disconnected if two non-adjacent vertices are deleted. Therefore  $c(C_p) \geq p-1$ . And  $C_p$  remains connected even after deletion of any vertex, therefore  $c(C_p) \leq p-1$ .

9.

$$\delta(G) \geq \frac{p + (l-1)(n-2)}{l} \implies \kappa_l(G) \geq n$$

### 19.1.6 Menger's Theorem

**Theorem 19.1.8.** For a non-trivial graph  $G$ ,  $\lambda(u, v) = M'(u, v)$  for every pair  $(u, v)$  of vertices of  $G$ .

**Corollary 19.1.8.1.** Graph  $G$  is  $n$ -edge connected iff every two vertices of  $G$  are connected by at least  $n$  edge disjoint paths.

**Theorem 19.1.9.** For every pair of non-adjacent vertices  $u, v$  in graph  $G$ ,  $\kappa(u, v) = M(u, v)$ .

**Corollary 19.1.9.1.** Graph  $G$  is  $n$ -connected iff every pair of vertices of  $G$  are connected by at least  $n$  internally disjoint paths.

**Algorithm 19.1.2** (connectivity  $\kappa(G)$ ). .

1. If degree of every vertex is  $p-1$ , then output  $\kappa = p-1$  and stop. Otherwise, continue.
2. If  $G$  is disconnected, output  $\kappa = 0$  and stop. Otherwise, continue.
3.  $\kappa \leftarrow p$
4.  $i \leftarrow 0$
5. If  $i \leq \kappa$ , then  $i \leftarrow i+1$  and continue. Otherwise, output  $\kappa$  and stop.

6.  $j \leftarrow i + 1$
7. (1) If  $j = p + 1$ , then return to step 5. Otherwise continue.
  - (2) If  $v_i v_j \notin E(G)$ , construct network  $N$  with digraph  $D$  as follows : for each vertex  $v \in V(G)$ , there are two vertices  $v', v'' \in V(D)$  and an arc  $(v', v'') \in E(D)$ . And for each edge  $uv \in E(G)$ , there are two arcs  $(u'', v), (v'', u) \in E(D)$ . The capacity function is given by,  $c(v', v'') = 1$  for every  $v \in V(G)$  and  $c(a) = \infty$  for every other arc in  $D$ . Set source  $s = v_i''$  and sink  $t = v_j'$  and find maximum flow in  $N$  using max-flow min-cut algorithm. Otherwise proceed to step 7d
  - (3) If  $f(N) < \kappa$ , then  $\kappa \leftarrow f(N)$ . Otherwise, continue.
  - (4)  $j \leftarrow j + 1$  and return to step 7a

## 19.2 Matchings and Factorizations

### 19.2.1 An Introduction to Matching

**Marriage Problem** Given a collection of men and women, where each woman knows some of the men. Can every woman marry a man she knows ?

**Assignment Problem** Given several job openings and applicants for one or more of these positions. Find an assignment so that maximum positions are filled ?

**Optimal Assignment Problem** Given several job openings and applicants for one or more of these positions. The benefits of employing these applicants on those positions are also given. Find an assignment of maximum benefit to the company ?

**matching** in  $G$  is a 1-regular<sup>3</sup> subgraph of  $G$ .

**maximum matching** in  $G$  is a matching of  $G$  with maximum cardinality.

**perfect matching** in  $G$  is a matching of cardinality  $p/2$ . ie,  $p/2$  edges.

**maximum weight matching** in a weighted graph  $G$  is a matching with maximum weight.

**Definitions 19.2.1.** Let  $M$  be a matching in a graph  $G$ ,

**matched edge** is an edge in subgraph  $M$  of  $G$ .

**unmatched edge** is an edge of  $G$  that doesn't belong to  $M$ .

**matched vertex** with respect to  $M$  is a vertex incident with an edge of  $M$ .

**single vertex** is a vertex that is not incident with any edge of  $M$ .

**alternating path** in  $G$  is a path with edges alternately matched and unmatched.

---

<sup>3</sup>A graph  $G$  is  $k$ -regular, if every vertex of  $G$  has degree  $k$ .

**augmenting path** in  $G$  is a non-trivial alternating path with single vertices as end vertices.

**Theorem 19.2.1.** *Let  $M_1, M_2$  be two matchings in  $G$  such that there is a spanning subgraph  $H$  of  $G$  with edges that are either in  $M_1$  or  $M_2$ , but not both. Then the components of  $H$  are either 1. isolated vertex 2. even cycle with edge alternately from  $M_1$  and  $M_2$  3. a non-trivial path with edges alternately from  $M_1$  and  $M_2$  such that each end vertex is single with respect to either  $M_1$  or  $M_2$ , but not both.*

*Synopsis.*  $\Delta(H) \leq 2$  by Pigeonhole principle. Any component of  $H$  is either a path or a cycle. A cycle with edge alternately from  $M_1$  and  $M_2$  is even. If an end vertex of a non-trivial path is matched with respect to  $M_1$  (WLOG), then it is there in  $M_1 - M_2$  ie, it is not there in  $M_2$ . If there is another edge in  $M_2$  incident with it, then it has to be in  $H$  and it will cease to be an end vertex of path component. Therefore, it is unmatched with respect to  $M_2$ .

**Theorem 19.2.2.** *A matching  $M$  in a graph  $G$  is maximum iff there is no augmenting path with respect to  $M$  in  $G$ .*

*Synopsis.* If  $M$  is maximum matching and  $P$  an  $M$ -augmenting path. Since both end-vertices are single, length of  $P$  is odd. Let  $M', M''$  be edges of  $P$  which are in  $M$  and not in  $M$  respectively. Then  $M - M' + M''$  is a matching of cardinality one greater than that of  $M$  which is a contradiction since  $M$  is maximum.

Conversely, suppose  $M$  be a matching such that there no  $M$ -augmenting paths in  $G$ . Let  $M'$  be a maximum matching in  $G$ . Then a nontrivial path component of the graph induced by  $M \Delta M'$  is of even length otherwise both end-vertices are matched with respect to one of the matching  $M$  or  $M'$  which is a contradiction. Again every cycle components are even. Therefore  $|M| = |M'|$ , since  $M \Delta M'$  doesn't have a nontrivial component of another kind.

**Definitions 19.2.2.** Let  $U_1, U_2$  be two nonempty, disjoint, subsets of the vertex set of a graph  $G$ . Then  $U_1$  **is matched to**  $U_2$  if there exists a matching  $M$  in  $G$  such that every edge in  $M$  incident with a vertex in  $U_1$  and a vertex in  $U_2$ . And every vertex of  $U_1$  (or  $U_2$ ) is incident with some edge in  $M$ . Suppose  $M^*$  be a matching such that  $M \subset M^*$ , then  $U_1$  **is matched under  $M^*$  to**  $U_2$ .

**Definitions 19.2.3.** Let  $U$  be a nonempty set of vertices of a graph  $G$ .  $U$  is **nondeficient**,<sup>4</sup> if  $|N(S)| \geq |S|$  for every nonempty subset  $S$  of  $U$ .

**Theorem 19.2.3.** *Let  $G$  be a bipartite graph with partite sets  $V_1, V_2$ . The set  $V_1$  can be matched to a subset of  $V_2$  iff  $V_1$  is nondeficient.*

**Corollary 19.2.3.1.** *Every  $r$ -regular bipartite multigraph has a perfect matching.*

**Theorem 19.2.4.** *A collection  $S_1, S_2, \dots, S_n$  of finite non-empty sets has a system of distinct representatives iff for each  $k$ ,  $0 \leq k \leq n$ , the union of any  $k$  of these sets contains at least  $k$  elements.*

<sup>4</sup> $N(S)$  is the neighbourhood set of all vertices adjacent to some vertex in  $S$

*Remark* (Hall's Marriage Theorem). Suppose there are  $n$  women. Then every woman can marry a man she knows iff each subset of  $k$  women ( $1 \leq k \leq n$ ) collectively knows at least  $k$  men.

*Remark.* Let  $W$  be a set of  $n$  women. Then there are  $2^n - 1$  nonempty subset for  $W$ . Thus, Hall's Marriage Theorem suggests that we ensure  $|N(S)| \geq |S|$  for every nonempty subset  $S$  of  $W$ . This method has complexity  $O(2^n)$ .

### 19.2.2 Maximum Matching in Bipartite Graphs

**Definitions 19.2.4.** Let  $M$  be a matching in a graph  $G$  and  $P$  is an augmenting path with respect to  $M$ . Let  $M'$  be set of edge in  $P$  and  $M$ . And  $M''$  be the set of edges in  $P$  and not in  $M$ . Then  $M_1 = (M - M') \cup M''$  is the **matching obtained by augmenting  $M$  along path  $P$** .

*Remark.*  $|M_1| = |M| + 1$

**Theorem 19.2.5.** Let  $M$  be a a matching of a graph  $G$  that is not maximum, and let  $v$  be a single vertex with respect to  $M$ . Let  $M_1$  denote the mathing obtained by augmenting  $M$  along some augmenting path. If  $G$  contains an augmenting path with respect to  $M_1$  that has  $v$  and an end-vertex, then  $G$  contains an augmenting path with respect to  $M$  that has  $v$  as an end-vertex

**Corollary 19.2.5.1.** Let  $M$  be a matching of a graph  $G$ . Suppose that  $M = M_1, M_2, \dots, M_k$  is a finite sequence of matchings of  $G$  such that  $M_i$  ( $2 \leq i \leq k$ ) is obtained by augmenting  $M_{i-1}$  along some augmenting path. Suppose  $v$  is a single vertex with respect to  $M$  for which there exists no augmenting path starting at  $v$ . Then  $G$  does not contain an augmenting path with respect to  $M_i$  ( $2 \leq i \leq k$ ) that has  $v$  as an end-vertex.

**Definitions 19.2.5.** An **alternating tree** with respect to a matching  $M$  is a tree such that every path from it's root are alternating path with respect to  $M$ .

**Algorithm 19.2.1** (Maximum Matching Algorithm for Bipartite Graphs). .

1.  $i \leftarrow 1$  and  $M \leftarrow M_1$
2. If  $i < p$ , then continue; otherwise stop.
3. If  $v_i$  is matched, then  $i \leftarrow i + 1$  and return to Step 2;  
otherwise,  $v \leftarrow v_i$  and  $Q$  is initialized to contain  $v$  only.
4. (1) For  $j = 1, 2, \dots, p$  and  $j \neq i$ , let  $TREE(v_j) \leftarrow F$ .  
Also,  $TREE(v_i) \leftarrow T$ .  
(2) If  $Q = \phi$ , then  $i \leftarrow i + 1$  and return to Step 2;  
otherwise, delete a vertex  $x$  from  $Q$  and continue.  
(3) (1) Suppose that  $N(x) = \{y_1, y_2, \dots, y_k\}$ . Let  $j \leftarrow 1$ .  
(2) If  $j \leq k$ , then  $y \leftarrow y_j$ ; otherwise return to Step 4.2  
(3) If  $TREE(y) = T$ , then  $j \leftarrow j + 1$  and return to Step 4.3.2.  
Otherwise, continue.  
(4) If  $y$  is incident with a matching edge  $yz$ , then  $TREE(y) \leftarrow T$ ,  
 $TREE(z) \leftarrow T$ ,  $PARENT(y) \leftarrow x$ ,  $PARENT(z) \leftarrow y$  and add  
 $z$  to  $Q$ ,  $j \leftarrow j + 1$  and return to Step 4.3.2. Otherwise,  $y$  is a  
single vertex and continue.

- (5) Use *PARENT* to determine the alternating  $v-x$  path  $P'$  in the alternating tree. Let  $P$  be the augmenting path obtained from  $P'$  by adding the path  $x, y$ . Proceed to Step 5

5. Augment  $M$  along  $P$  to obtain a new matching  $M'$ . Let  $M \leftarrow M'$ ,  $i \leftarrow i + 1$ , and return to Step 2.

**Definitions 19.2.6.** Let  $G$  be a weighted complete bipartite graph with partite sets  $V_1$  and  $V_2$ . A **feasible vertex labeling** is a real function  $l : V(G) \rightarrow \mathbb{R}$  on vertex set of  $G$  such that  $l(v) + l(u) \geq w(vu)$  where  $v \in V_1$  and  $u \in V_2$ .

**Definitions 19.2.7.** Consider the function  $l : V(G) \rightarrow \mathbb{R}$  such that  $\forall v \in V_1$ ,  $l(v) = \max\{w(vu) : u \in V_2\}$  and  $\forall u \in V_2$ ,  $l(u) = 0$ . Then  $l$  is a feasible vertex labeling on  $V(G)$ . And,

$E_l$  is the set of all edge of the weighted complete bipartite graph  $G$  such that  $l(v) + l(u) = w(vu)$ .

$H_l$  is the spanning subgraph of  $G$  induced by the edge set  $E_l$ .

**Theorem 19.2.6.** Let  $l$  be a feasible vertex labeling of a weighted complete bipartite graph  $G$ . If  $H_l$  contains a perfect matching  $M'$ , then  $M'$  is a maximum weight matching of  $G$ .

**Algorithm 19.2.2** (Kuhn-Munkres). .

1. (1) For each  $v \in V_1$ , let  $l(v) \leftarrow \max\{w(vu) : u \in V_2\}$ .  
 (2) For each  $u \in V_2$ , let  $l(u) \leftarrow 0$ .  
 (3) Let  $H_l$  be the spanning subgraph of  $G$  with edge set  $E_l$ .  
 (4) Let  $G_l$  be the underlying graph of  $H_l$ .
2. Apply Algorithm 19.2.1 to determine a maximum matching  $M$  in  $G_l$ .
3. (1) If every vertex  $v$  of  $V_1$  is matching with respect to  $M$ , output  $M$  and stop. Otherwise, continue.  
 (2) Let  $x$  be the first single vertex of  $V_1$ .  
 (3) Construct an alternating tree with respect to  $M$  that is rooted at  $x$ . If an augmenting path  $P$  is discovered, then augmenting  $M$  along  $P$  and return to Step 3.1. Otherwise, let  $T$  be the alternating tree with respect to  $M$  and rooted at  $x$  that cannot be expanded further in  $G_l$ .
4. Compute  $m_l \leftarrow \min\{l(v) + l(u) - w(vu) : v \in V_1 \cap V(T), u \in V_2 - V(T)\}$ .  
 Let

$$l'(v) = \begin{cases} l(v) - m_l & \text{for } v \in V_1 \cap V(T) \\ l(v) + m_l & \text{for } v \in V_2 \cap V(T) \\ l(v) & \text{otherwise} \end{cases}$$

5. Let  $l \leftarrow l'$ , construct  $G_l$  and return to Step 3.3.



### 19.2.4 Factorizations

**Definitions 19.2.8.** A **factor** of a graph  $G$  is a spanning<sup>5</sup> subgraph of  $G$ .

**Definitions 19.2.9.** Let  $G_1, G_2, \dots, G_n$  be edge-disjoint factors of  $G$  such that  $E(G) = \cup_{i=1}^n E(G_i)$ . Then  $G$  is **factorable** and  $G = G_1 \oplus G_2 \oplus \dots \oplus G_n$ .

**Definitions 19.2.10.** An  $r$ -regular factor of  $G$  is an  **$r$ -factor** of  $G$ .

**Definitions 19.2.11.** If  $G$  has a factorisation to  $r$ -factors, then  $G$  is  **$r$ -factorable**.

*Remark.*  $K_{3,3}$  is 1-factorable.  $K_5$  is 2-factorable.

**Definitions 19.2.12.** An **odd component of  $G$**  is a component of  $G$  with odd number of vertices. And an **even component of  $G$**  is a component of  $G$  of with even number of vertices.

**Theorem 19.2.7 (Tutte).** *A nontrivial graph  $G$  has a 1-factor iff for every proper subset  $S$  of  $V(G)$ , the number of odd components of  $G - S$  does not exceed  $|S|$ .*

*Remark.* There exist cubic graphs that doesn't have a 1-factor.

**Theorem 19.2.8 (Petersen).** *Every bridgeless cubic graph contains a 1-factor.*

*Remark.* Every bridgeless cubic graphs has a 1-factor. Let  $G$  be a bridgeless cubic graph. Consider every pair of factors  $G_1, G_2$  such that  $G = G_1 \oplus G_2$  where  $G_1$  is a 1-factor and  $G_2$  is a 2-factor.  $G$  is not 1-factorable only if every such  $G_2$  doesn't have a 1-factor.

**Theorem 19.2.9.** *Petersen graph is not 1-factorable.*

**Theorem 19.2.10.** *Every  $r$ -regular bipartite multigraph ( $r \geq 1$ ) is 1-factorable.*

*Remark* (Application of 1-factorisation). For even number  $p$ , a 1-factorisation of  $K_p$  corresponds to the schedule of a round of the round robin tournament among  $p$  teams. If  $p$  is odd, consider  $K_{p+1}$  where  $v_{p+1}$  is an imaginary team called bye team. A game with bye team is a bye.

**Definitions 19.2.13.** A **hamiltonian cycle** is a spanning cycle. And, **Hamiltonian graph** is a graph containing a hamiltonian cycle.

**Theorem 19.2.11.** *Complete graph  $K_{2n+1}$  can be factored into  $n$  hamiltonian cycles.*

*Remark.* For  $n = 3$ ,  $K_7$  can be factored into three hamiltonian cycles.

**Theorem 19.2.12.** *Let  $0 \leq r < p$ . Then there exists an  $r$ -regular graph of order  $p$  iff  $pr$  is even.*

**Definitions 19.2.14.** Let  $\{E_1, E_2, \dots, E_n\}$  be partition of  $E(G)$ . And let  $H_i$  be subgraph of  $G$  induced by the edge set  $E_i$ . A **decomposition** of a graph  $G$  is a collection of these subgraphs  $H_1, H_2, \dots, H_n$ . And  $G = H_1 \oplus H_2 \oplus \dots \oplus H_n$ .

**Definitions 19.2.15.** Let  $G = H_1 \oplus H_2 \oplus \dots \oplus H_n$  be a decomposition of  $G$  such that  $H \cong H_i$ . Then  $G$  is  $H$ -decomposable.

*Remark.*  $K_{3,3}$  is  $3K_2$ -decomposable.  $K_5$  is  $C_5$ -decomposable.  $K_{2n}$  is  $nK_2$ -decomposable.  $K_{2n+1}$  is  $C_{2n+1}$ -decomposable. Every graph is  $K_2$ -decomposable. Every complete bipartite graph  $K_{m,n}$  is  $K_{1,m}$ -decomposable and  $K_{1,n}$ -decomposable.

<sup>5</sup>Spanning subgraph of a graph  $G$  has every vertex of  $G$

### 19.2.5 Block Designs

**Definitions 19.2.16.** A block design on a set  $V$  is a collection of  $k$ -element subsets of  $V$  such that each element of  $V$  appears exactly in  $r$  subsets.

**variety** The elements of  $V$  are called varieties.

**block**  $k$ -element subsets of  $V$  are called blocks.

**balanced design** If each variety appears in exactly  $r$  blocks and each pair of varieties appears in exactly  $\lambda$  blocks.

**incomplete design** If blocks are proper subsets of  $V$ . ie,  $k < v$ .

**Definitions 19.2.17.** A balanced incomplete block design of  $v$  varieties in  $b$  blocks of cardinality  $k$  such that each variety appears in exactly  $r$  blocks and each pair of varieties appears in exactly  $\lambda$  blocks is a  $(b, v, r, k, \lambda)$ -design.

**Theorem 19.2.13.**  $bk = vr$

**Theorem 19.2.14.**  $\lambda(v-1) = r(k-1)$

**Corollary 19.2.14.1.**  $\lambda < r$

**Theorem 19.2.15** (Fisher's Inequality).  $b \geq v$

**symmetric design** If  $b = v$

**Theorem 19.2.16.** In a symmetric  $(b, v, r, k, \lambda)$ -design with even  $v$ ,  $r - \lambda$  is a perfect square.

**steiner triple system**  $(b, v, r, k, \lambda)$ -design with  $k = 3$ ,  $\lambda = 1$ .

*Remark.*  $(b, v, r, k, \lambda)$ -designs are incomplete. However, a complete block design (ie,  $v = 3$ ) is also included as a Steiner triple system.

**Theorem 19.2.17.** Steiner triple system with  $v$  varieties exists iff  $v = 6n + 1$  or  $v = 6n + 3$  or  $v = 3$ .

**Definitions 19.2.18** (Kirkman's Schoolgirls Problem). A class of 15 girls. Parade 15 girls in five rows (3 girls in a row). Is it possible to plan 7 days parade so that two girls are together in a row exactly once?

**kirkman triple system** Steiner triple system with  $v = 6n + 3$ .

*Remark.* It is proved that kirkman triple system exists with  $v = 6n + 3$  for every  $n \geq 0$ .

**Theorem 19.2.18.** The code consisting of the rows of the incidence matrix of a  $(b, v, r, k, \lambda)$ -design ( $b = v$ ,  $r=k$ ) is  $t$ -error correcting, where  $t = k - \lambda - 1$ .

**Subject 20**

**ME800403 Combinatorics**

**Subject 21**

**Probability Theory**

**Subject 22**

**Operational Research**

**Subject 23**

**Coding Theory**

**Subject 24**

# **Commutative Algebra**

**Subject 25**

# **Ordinary Differential Equations**



**Subject 26**

# **Classical Mechanics**

# Bibliography

- [Apostol, 1973] Apostol, T. (1973). *Mathematical Analysis, 2nd edition*. Narosa Publishing House.
- [Fraleigh, 2013] Fraleigh, J. B. (2013). *A First Course in Abstract Algebra, 7th edition*. Pearson Education.
- [Gray Chartrand, ] Gray Chartrand, O. O. (?). *Applied and Algorithmic Graph Theory*. Tata McGraw Hill Company.
- [Joshi, 1983] Joshi, K. D. (1983). *Introduction to General Topology*. Wiley Eastern Ltd.
- [Kiusalaas, 2013] Kiusalaas, J. (2013). *Numerical Methods in Engineering with Python3*. Cambridge University Press.
- [Munkres, 2003] Munkres, J. R. (2003). *Topology, 2nd edition*. Pearson.
- [Rudin, 1976] Rudin, W. (1976). *Principles of Mathematical Analysis, 3rd edition*. McGraw Hill Book Company, International Editions.