# Data Engineering

## 1) SQL

→ Case Example. (can count values multiple for diff colours using single query.)

```
select sum(case col1
              when 2.99 then 1
              else 0
          END) as regular,
       sum(case col2
              when 4.99 then 1
              else 0
          END) as rental rate;
```

→ COALESCE (function).

→ Accepts unlimited no. of arguments and returns 1st argument i.e not null if all null returns null.

Eg:- SELECT coalesce(NULL, 1, 2, 3) ⇒ 1

★→ We use COALESCE to substitute null values while doing operations.

Eg:- SELECT item, (price − COALESCE(discount, 0)) As final-price from table1.

(Replace discount with value null to 0)

→ Cast (operator) :- It lets you convert from one data type to other.

Syntax :- SELECT CAST ('5' AS INTEGER)

  postgres SQL ⇒ SELECT '5' :: INTEGER.

→ NULLIF (function) :- Returns null if 2 args passed are '=' else returns 1st arg passed.

→ Adv from table is no phy storage and data queried is up to date

→ View" :- A data base obj that is of stored query.

Eg :- CREATE VIEW view_name AS
      SELECT c1,c2 from t1 inner join address
      on c1.a-id = c2.a-id.

To update view :-
      CREAT OR REPLACE view_name AS
      query.......

Drop view if exist :-
      DROP VIEW IF EXISTS view_name.

Rename VIEW :
      ALTER VIEW v1 RENAME ~~to to~~ to v2

→ Importing and exporting data functionality of pg Admin that allows us to import data from .CSV file to already existing table.

→ psycopg2 :- Python library to interact with DB in postgresSQL.

Eg :- import psycopg2 as pg2
conn = pg2.connect(database = 'name', user = 'pipi', password = 'pwd')

cur = conn.cursor()
cur.execute("SELECT * FROM t1')
cur.fetchone() ⇒ Fetch query results
cur.fetchmany(10) ⇐
cur.fetchall() ↗
query = """CREATE TABLE - - - - - - - """
cur.execute(query).
★ cur.commit() ⇒ To commit changes to DB.
★ conn.close()


Window functions :- They perform aggregate ops on grp of rows but they produce result for each row. Eg :- SELECT c1, c2, c3, avg(c3) over(partition by c2) as ~~dept~~ c3_avg from t1.  (group by c2)

Eg2 :- SELECT c1 AVG(c2) over() from emp;
Eg3 :- select emp_no, dept, sal, count(*) over (PARTITION BY dept) as dept_count. from emp;

★Eg 4 :- SELECT emp-no, dept, sal,
　　　　　SUM(salary) over (PARTITION by department
　　　　　　　　　　ORDER BY SALARY) as rolling-dept-
　　　　　　　　　　　　　　　　　　　　　　　　salary
　　　　MIN (salary) over 　　　　'' 　DESC ''


★Eg 5 :- SELECT emp-no, dept, salary,
　　　　　　　　　　　　　　　　　　　salary
　　　　　RANK() over ( order by ~~dept~~ DESC).
RANK() =) Gives rank based on colum value.
　　　★→ It might miss a rank if their are
3 people with rank 1 it will miss rank 2,3 the next
rank will be 4.
　　　DENSE-RANK() → Works like RANK() but
　　　　　　　　　　　　　　　　　　　　　　　preceeding
ranks will not miss a number ~~if~~ even ~~before~~ ranks
are same.


Eg 6 :- SELECT emp-no, dept, salary
　　　　　NTILE(4) over (order by salary DESC).


Eg 7 :- (FIRST_VALUE), (LAST_VALUE), (N'TH-VALUE)
Eg 8 :- (LAG, LEAD), LAG(col1, 2)

✱ → We can create table was from query.

## Materialized View :-

→ To avoid having redundant storage and up to date data we use views but for complex queries to execute an view takes time so we use materialized views.

✱ → Data is stored physically.

⟹ Eg1:- CREATE MATERIALIZED VIEW v-name
         As query.

Eg2:- REFRESH MATERIALIZED VIEW
                 (To update view with new
                    data from the view query)

· → We can update data manually or use triggers.

Eg3:- ALTER MATERIALIZED VIEW
           DROP MATERIALIZED VIEW

(named subquery)
## Common table exprs (CTE) :-

→ To use query result in other query.

✱ Syntax :- WITH [cte-name AS]
           SELECT ..... )
           SELECT ..... cte-name ...;

Multi CTE

Syntax:
with name1 as (        ),
name2 as (    )
Select - - - -

→ We can not use aliases for colname in same query level.

Indexing:- → Single level (primary, clustered, secondary).
→ Multi level (Btree, B+tree) → It reduces time take to data transfer b/w RAM and hard disk.

→ It reduces I/O cost.

Indexing types:
→ We use [Sparse] when data is sorted
→ Here the whole block pointer is stored.
→ We use [Dense] when data/records is unsorted
→ Here the pointer is each row pointer is stored.

→ Primary Indexing:-
→ Ordered file whose records (key + pointer) of fixed length.

→ Clustering Indexing:-
→ Records are physically ordered on non-key field which does't have distinct value for each record
before each block of disti
→ Each distinct value of non-key attribute an entry is created.
→ Each for Duplicate non-key attribute points to single pointers is used for reference.

Secondary Indexing:-
→ Creating index on 2nd attribute when we have primary indexing.

→ 2nd attr is not ordered in DB. But can order it in index table.

1° index Syntax:-
   → CREATE INDEX index-name
      on table-name (col-name);
      (col1, col2); multi col index.
   → DROP INDEX index-name;

→ Indexing should be reconsidered when:-
   → Small tables
   → Have frequent large batch update, insert op's.
   → On columns with high no. of null values.
   → Columns that are frequently manipulated.

→ Execute sp-helpindex table-name; ⇒ ★ To see all the indexes on particular table.

→ In SQL server ⟋ Clustered index.
                 ⟍ Non-clustered index.

   → Clustered Index:-
      → Defines the order in which data is physically stored in a table. As sorting can be done in only way we have one clustered index per table. In SQL server 1° key constraint auto created

clustered index.
→ When new data is added it maintains order based on cluster index value.

→ Custom clustered index:-
    CREATE **CLUSTERED** INDEX index_name
    on table_name (col1, col2);

→ Non clustered index:-
→ Separate storage is needed.
→ Can have multiple non-clustered indexes per table
→

Types of DB :- 1. key-value, 2. Wide column (used when freq write and not much read and update).
            (Redis, memcache) (Cassandra, Hbase)
                            Schemaless, CQL

3. Document DB (MongoDB, Firestore, DynamoDB, couchDB)
        →Unstructured, No schema.
    →  Collection -(contain 2 (or more documents).
    →  Document (contain key-value pairs JSON).
    ★ → Field within collection can be indexed and collections can be
organized into a logical hierarchical
        → Reads are much faster while write and update is
complex.
4. RDB :- (MySQL, PostgreSQL, SQL server...).
        → ACID compliant.
    →Cockroach DB scales horizontally without reconfig..

5. Graph DB :- (Neo4J, Dgraph).
   → Nodes contain data.
   → ~~Edge~~ defines relationship
   → Good choice alt for RDB if we have lot of joins and perf issue.

6. Search Engine :- (Give list of most appropriate results).
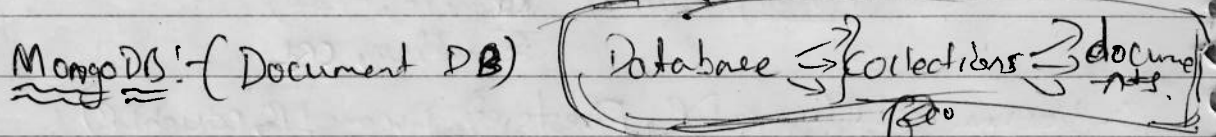   → Most are on top of Apache Lucene
   Eg :- Solr, elastic search.

   dRbs

7. MultiModel (FAUNA) :-
   → Use GraphQL to describe how we want data

MongoDB :- (Document DB)   Database ⇄ Collections ⇄ documents.

→ Documents can be stored in formats like JSON, BSON, XML.
   → BSON :- Binary JSON to overcome JSON drawback
      ↳ It encodes type and length information
   makes it possible to traverse more quick compared to Json.
      ↳ It adds some non-JSON data types (date, binary data)
      ↳ MongoDB driver converts Json ⇄ BSON.
→ Some Document DB provide schema validation.
→ Collection :- Set of documents that can have different schema. Eg :- Users collection can have 2 documents
   ✱         with 1 user detail in each doc.

→ CRUD can be done here.
  ↳ Using unique id/field values. Indexes can be used for good perf.

→ Document DB's are distributed. (Horizontal scalling).

→ Single transaction works fine in Mongo DB.

→ ~~for mult~~ Mongo DB support multi transactions but the perf might be affected.

→ Cluster is place where MongoDB DB's are stored.

*→ Data base design and data modeling are major factors in MongodB performance.

*→ ~~where~~ _id will be created auto matically for each new doc created in collection. But we can assign it manually using '_id' as key.

*→ $ with word in mongoDB indicate its reserved word.

→ db.flightData.updateMany ({ }, {$set : {mark: "delete"}})
    ⇓                                    ⇑
  collection              if key exist with
                          update value else add it.
                                        filter set to empty to update all docs

→ db.flightData.insertMany ([{$$$}])
                              ⇑
                    pass ~~collection~~ documents here.

→ db.flightData.find ({distance : {$gt: 10000}}).pretty().
         (to get data   distance > 10000)

→ db.flightData.update ({_id:ObjectId ('xxx')}, {dela:false})
      This will not ~~just~~ set the dela field to false
    it will replace ~~a~~ whole obj with only _id and dela keys.

→ db. Passengers . find ( ) ⇒ Returns cursor to go through elements.

→ db. passengers . find ( ). ~~toArray~~ ( ) ⇒ Returns all elements (documents)

→ db. passenger. find ( ). forEach ( ) ⇒ Returns 1 element for a loop reducing load...

→ Projection :- Instead of ~~getting~~ whole document (or) object and filtering it in app which increase load ~~and~~ / latency we can filter(or) get the data only we want ~~using from~~ by filtering within DB we use projection.

Eg:- db. passengers. find ({ }, { name: 1 }) ⇒ Returns all data with id, name

★ db. passengers. find ({ }, { name: 1, _id = 0 }) ⇒ Returns all data with only name.

★ Embedded Document :- Document inside other document.

→ Max document size in MongoDB is 16 MB.

→ db.hobbies. insertMany ( [{ 'n' : 1 }, --- • ], { ordered: false });

    → By default ordered is true so if any doc insertion fails only preceeding docs would get inserted to avoid it we can set false

→ We use journal ~~in case~~ to maintain [To-do] tasks so that in case the memory ~~was~~ get wiped out (power cut) it still knows what to do during journal. (j: undefined (no response if saved in journal), j: True (Response))

⟐

→ Command in shell :-

    mongoimport tv-shows.json -d movieData -c movies

    → db name    → collection name

    --jsonArray
    -- drop

need to be specified it DB of docs in file.    if exist    with / without del & append create