

25 Oct

Day 20

AZURE

① Create an Resource Group.

- A container that holds resources for any Azure solution like Virtual machines, databases, accounts.
- It helps managing & organizing resources as a SINGLE UNIT.
 - ↳ allowing you to deploy, monitor, control access costs, lifecycle (create, manage, delete) together

Ex

WebApp - RG

- A VM
- Storage account
- network

If u delete RG → all

3 resources get deleted automatically

◦ Shared lifecycle

② Create a resource now, from inside RG

1. Storage Account

- A container that holds all your storage data
 - blobs, files, tables

Azure Blobs → Cloud based object storage solution
for unstructured data

- Can store vast amount of data like images, videos, audio
- highly scalable & cost effective

Azure Data Lake → Built on top of blob

- Optimized for big data analytics with a hierarchical namespace
 - ↳ organizes data like a file system not a blob

Normal Blob: → Virtual Hard drive in cloud.

- All files are stored as objects, just keys & paths. No true folders.

with HNS: Create folders within folders

↳ We can have real directories, subdirectories

∴ Data lake storage Gen2

- file system + folder organization [on top] of that hard drive

→ Redundancy Options

LRS (Locally redundant storage)

- lowest cost
- Data is replicated 3 times within a single datacenter.

(Why 3? → If 2 fail, only 1 remains so high risk of data loss)

If 3+: Extra cost

∴ Many cloud providers use triplicate replication as a balance of cost, reliability, availability

* Region: A geographic area where Azure data centers are located

Ex East US / Central West

↳ choose region based on latency

↓
Availability Zone: A physically separated datacenter within a region designed to be isolated from failure of other zones.

Date Centre: Actual facility where servers & storage & networking circuits

LRS

If a server or disk fails

↳ other copies in distributed keeper are safe

But if distributed fails (fire / power outage)

↳ All 3 copies get affected

ZRS → Zone Redundant storage

Data is replicated across 3 zones in a region

Limitation → regional disaster

GRS → Geo-redundant storage

Data is replicated to another region

GZRS → Geo zone redundant storage

- highest durability & availability
- data remains safe even if entire zone/region fails.

Redundancy

→ keeping extra copies of

data or resources → so that

if something fails, system continues to work

26 Oct

Day 21

Q) Data factory

- ↳ Cloud based ETL (Extract, Transform, Load) service used to move & transform data from various sources to destinations, for analytics or storage.
- Used to build data pipelines → to collect/clean data automatically.

Pipeline : is a group of activities that together perform a specific data workflow, (ETL).

SHIR

Self Hosted Integration Runtime

- ↳ Sometimes data is not in Azure, maybe it's on-premises SQL server, or private Azure Network
ADF can't reach the data as it's not public
- ∴ So we install SHIR on (a small software agent) on the computer / server inside our network. This agent acts as a bridge b/w your data & ADF

- ↳ Data needs to connect to ADF, somehow

① Public Endpoint • SHIR connects to ADF over public internet (but securely encrypted)

• simple setup, safe (outbound traffic only)

② Private Endpoint • SHIR connects to ADF privately using Azure Private Links

• Traffic stays 100% private

• Used when company don't allow internet communication

• High security, sensitive data

Ex : A company store its data in on-prem SQL server

- ① behind a firewall
- ② ADF can't access it directly.

- ③ Install SSIS on the system

(which ~~will~~ can access SQL server data locally
& send data securely to ADF)

Ex Public Endpoint → How its secure / encrypted?

Ex : I login to bank website,

its public but encrypted (https)

Environments

• On prem : means to company's own servers, databases, & storage & are physically hosted in their own building / datacentre not in cloud.

Ex = your official SQL server

• Cloud : Azure SQL Database

• Hybrid : on-prem AB + ADF pipelines

HTTPS



• "S" stands for secure

• Uses encryption to keep data private / safe

Oct
Day 24

AZURE

③

Storage Rewrite

Containers

Add container role + tier over

① Bronze ② Silver ③ Gold

(Container for each layer)

ADF

→ Launch

→ Configure Git

↳ Create a branch (feature)

④

Azure SQL

↳ Create server

↳ Create database

Configure

Region (Latency)

Compute / cores

(More cores More Processing Power)

↳ Faster & Expensive

↳ Networking → Public Endpoint

↳ Security

↳ Data Source ↳ Start with Blank DB or Sample

Data collation → String data rules

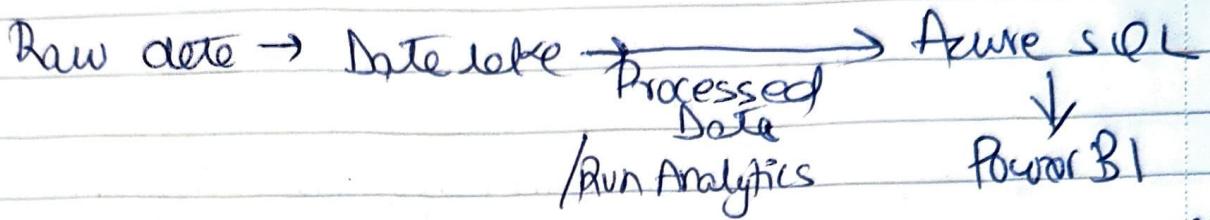
• Case Insensitive ($A = a$)

• Accent sensitive ($\acute{e} \neq e$)

31 Oct
Data

Data lake → stores raw / unstructured data

Azure SQL DB → stores cleaned, structured
data ~~for dashboards~~



* Incremental load : means loading only "new / changed"
data instead of loading entire dataset every time

Compare last modified date > Last Run Time

Backfilling → means loading historical data or
which u missed / need old data
→ often do backfilling first to load
everything ≥ then → incremental load.

① Load data → DB
↳ Query Editor

(Manually creating sample DB,
in companies → DB already exists
2 we simply connect to it)

2 Build pipelines for extraction
2 transformation

In Real Development

We use SQL Server Management Studio (SSMS)

↳ a Microsoft tool to connect, manage, query SQL DBs

Azure SQL : stores structured data, but transactional

① (Source)

OLTP
(Operational data)

not optimised for Analytics

As it could contain

↳ Redundant data, nulls

only optimised
by OLTP data

Data Lake : ADF copies raw data from source to DataLake

② (Loading/Staging)

Here we store raw data

Safe & cheap

③ Databricks/ADF Pipeline : Synapse Transformation

Clean / Join / Aggregate

④ DataWarehouse : (Load cleaned data into)
(Synapse) fact and Dim Tables

STAR Schema

Ready for reporting

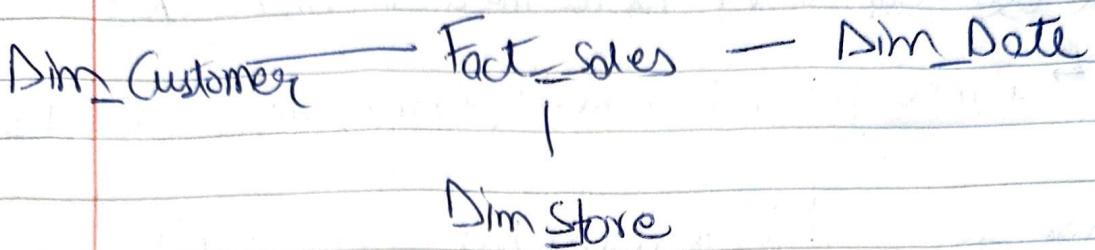
⑤ PowerBI

↳ Azure Synapse Analytics : Platform for Analytics AND Reporting
Includes Data Warehouse

STAR SCHEMA

↳ is a way to organize data in Data warehouse
for analytics

Dim Product



Fact → holds numeric data

Ex: Sale amount, quantity, ID

Dimension Table → Holds descriptive attributes
that explain facts

Ex: customer details → name, city

product info, brand

date info → day, month

- Why? → faster joins, aggregations
- makes Analytics fast & simple
- used heavily in OLAP

In real companies, we use STAR schema, a central fact table holding measurable fields

2 surrounding Dimension Tables holding descriptive attributes

Fact Tables: Can have billion of rows

most companies organise data into Fact & Dim. Tables

* OLTP : handles day to day transactions
optimized for fast inserts/updates

Small queries
Normalized scheme

Ex: Azure SQL stores transactional data
Shopping cart updates / Banking deposits
/ Airline bookings

OLAP : Analyzes historical data / aggregated data
for insights, reporting

heavy queries / complex
Denormalized scheme (STAR)

Ex: Synapse DW / DataBrick Analytics
ADF pipeline / PowerBI

Analyze sales trends, fraud detection,
Revenue trends, monthly transaction summary
/ flight revenue per route

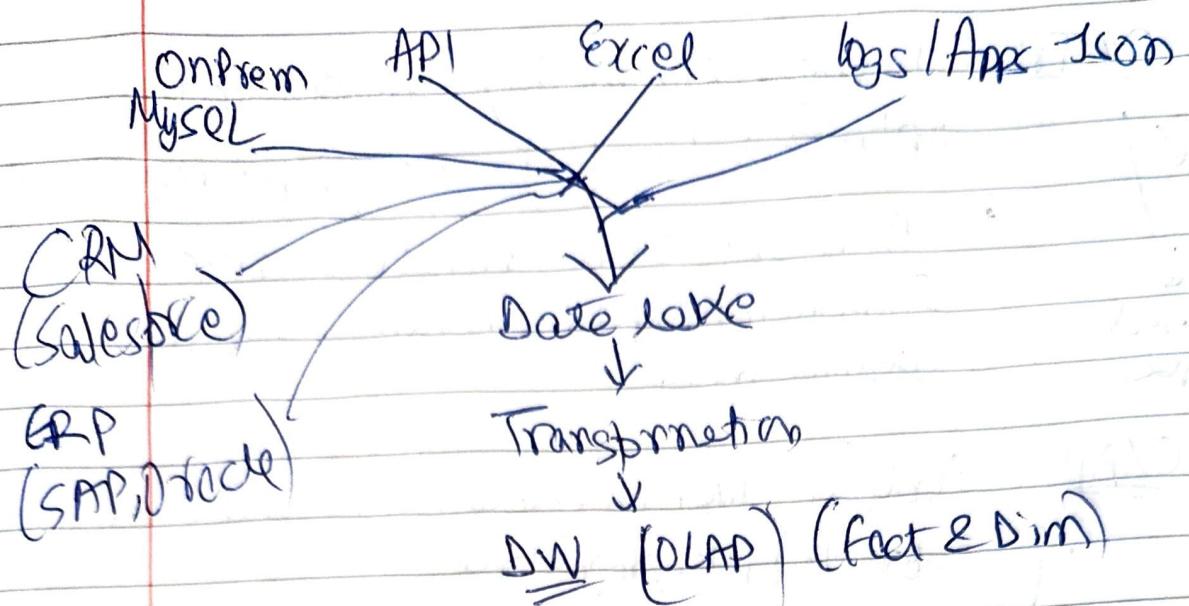
• Data comes from OLTP / or multiple sources

* Transactional Data → Records of business events that happen in real time
→ Bank deposited \$500, Order ID 123 placed by user

Orders

| Order ID | Customer ID | Product ID | Amount | Date |
|----------|-------------|------------|--------|------|
| | | | | |

Each row → 1 transaction



Normalization

- ① I start by analyzing source data | study it

Identify - Entities
 Attributes
 Relationships

 Using ER Model
- ② Distinct entities become separate tables

Split into multiple Tables
- ③ Assign unique identifiers to each table
- ④ Ensure each table stores only 1 type of information
 (Removing Redundancy)
- ⑤ Check for Integrity
 - ↳ maintain relationships b/w tables

Why denormalize for Analytics

↳ As ~~normal~~ OLAP requires joining many tables.
But joining will be very slow for million of rows

∴ Denormalized Dim/Fact tables

→ Reduces number of joins
making queries much faster

↳ Denormalization → is not messy

it is Purposefully flattened for Analytics

→ you control what is included

↳ for reporting only when needed

→ Combine ^{related} multiple tables into a single table
in a controlled & intentional way.

So that analytics, reporting become faster

→ Widen the tables, instead of joining
for fast reading

* In OLTP, you heavily normalize data → Many small tables
↳ Customers, Orders, Products, Payments, Categories, Details
∴ To analyze we need '6' Joins (=)

∴ In OLAP/Data warehousing → we denormalize into simpler models → STAR SCHEMA
Dim

Dim - Fact - Dim → Here OLAP, joins (less) but with
Dim

↳ fewer & longer tables → which is better
for larger data

Nov 1

Day 26

~~Linked Services~~

Normalization Levels) Hierarchy

1NF : Requirement



All columns have atomic values (no repeating group)

Atomic = one fact per column

not a list, array

Ex Alice 102, 103

| | | |
|------------------|-------|-----|
| ↳ <u>Split</u> : | Alice | 102 |
| | Alice | 103 |

2NF : A table is in 2NF, if it is in 1NF AND all non-key columns fully depends on Primary Key.

Key → a col'n / multiple col'n which uniquely identify a row

Non-key → every else col'n : descriptive info

3NF : in 2NF AND, every nonkey column only depends on a key column, no any other non-key col'n.

This means no transitive dependency

3NF removes ↗

Ex: City depends on Region but both are Non Key so we split

BCNF : if it is in 3NF AND, every determinant must be a Candidate key.

(Boyce - Codd)
Normal form

Candidate key : A col'n / set of col'n's that uniquely identify each row

determinant : Any col'n / set of col'n's that determine another col'n's value

If a non key col'n determines another col'n

↳ ~~fails~~ its NOT BCNF

Here also to fix → we split into multiple tables
so every determinant is a key in its own table

| Course Table | |
|--------------|------------|
| Course | Instructor |
| | |

| Instructor Table | |
|------------------|------|
| Instructor | Room |
| | |



4NF, 5NF... Removes further dependencies

* we choose 1. Primary key from multiple Candidate keys