# Azure Data Engineering project

(ETL)
pipeline

Data source (kaggle) (uploaded on github)

↓

Data integration (Data factory) (Data flow)

↓

Store data (storage acc Azure)

↓ (storing in a container)

Transformation (Databricks, spark)

↓

Store data

↓

Analytics (Azure synapse analytics / SQL)

Dashboard) ✓

Azure acc

{ subscriptions }

{ Resource groups }

{ Resources / Azure resources }

store data

   ↳ container

      − raw data

      − transformed data


Datafactory    (triggers → schedule the

              end to end pipeline

 ↳ activities                execution)

     − copy data  (source → sink)

             − url        − azure

              from       data lake

              github     storage

              (http)     Gen 2


loaded data from source

               data ⌐

                     └→  storage container

               factory


Databricks

   ↳ new compute (to write spark code)

    ↳ new notebook (spark cluster)

raw_data → databricks → transformed_data

have to connect databricks to ADLS

— use key vault to store & not expose
the keys while mounting

— we are using the app to get the data
(need to give permission using IAM)

— while writing the files apache stores them
in a folder along with metadata

— if we have a very large file spark
will divide the file into multiple files


Azure synapse

— after loading the tables to the DB
   (options                ↳ SQL script
     — custom        ↳ notebook
     — from template
     — from Data Lake)  ↳ ML → train/predict

(can SQL queries , MYSQL syntax)
( result → table & chart)

## Triggers in datafactory :

1) schedule     (many to many)

2) Tumbling window     (seperate files for
                              every execution)

3) Event based   (blob related events
                        del / generation of blob)