# HISTORICAL DATA ANALYSIS TO SUPPORT THE CLASSIFICATION OF TURBINE AND COMPRESSOR COMPONENTS AND PREDICT FUTURE DEMAND

*Canditate: MATTEO D'SOUZA*

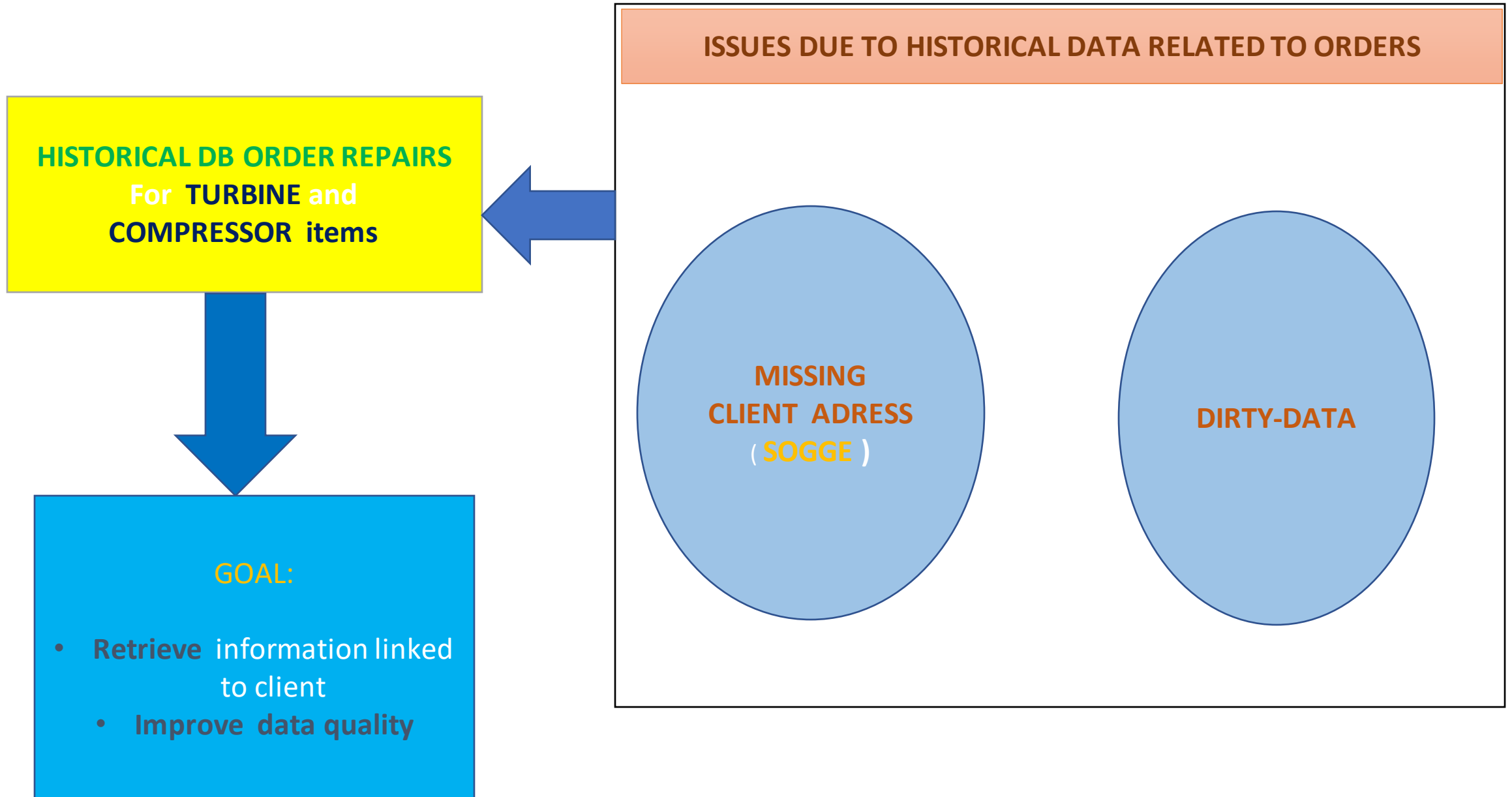# BUSINESS CONTEXT

TURBOMACHINERY& PROCESS SOLUTION
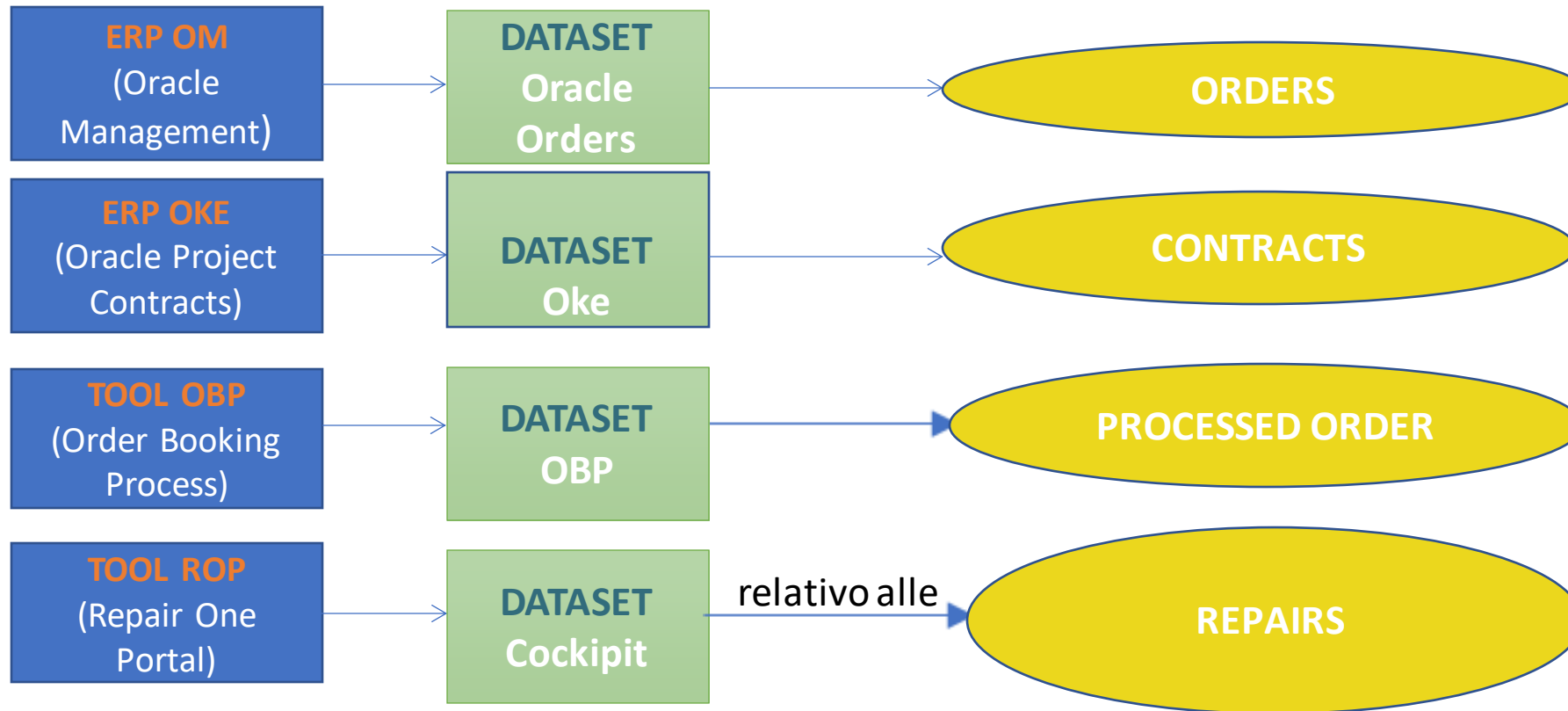
REPAIRS

SPARE PARTS

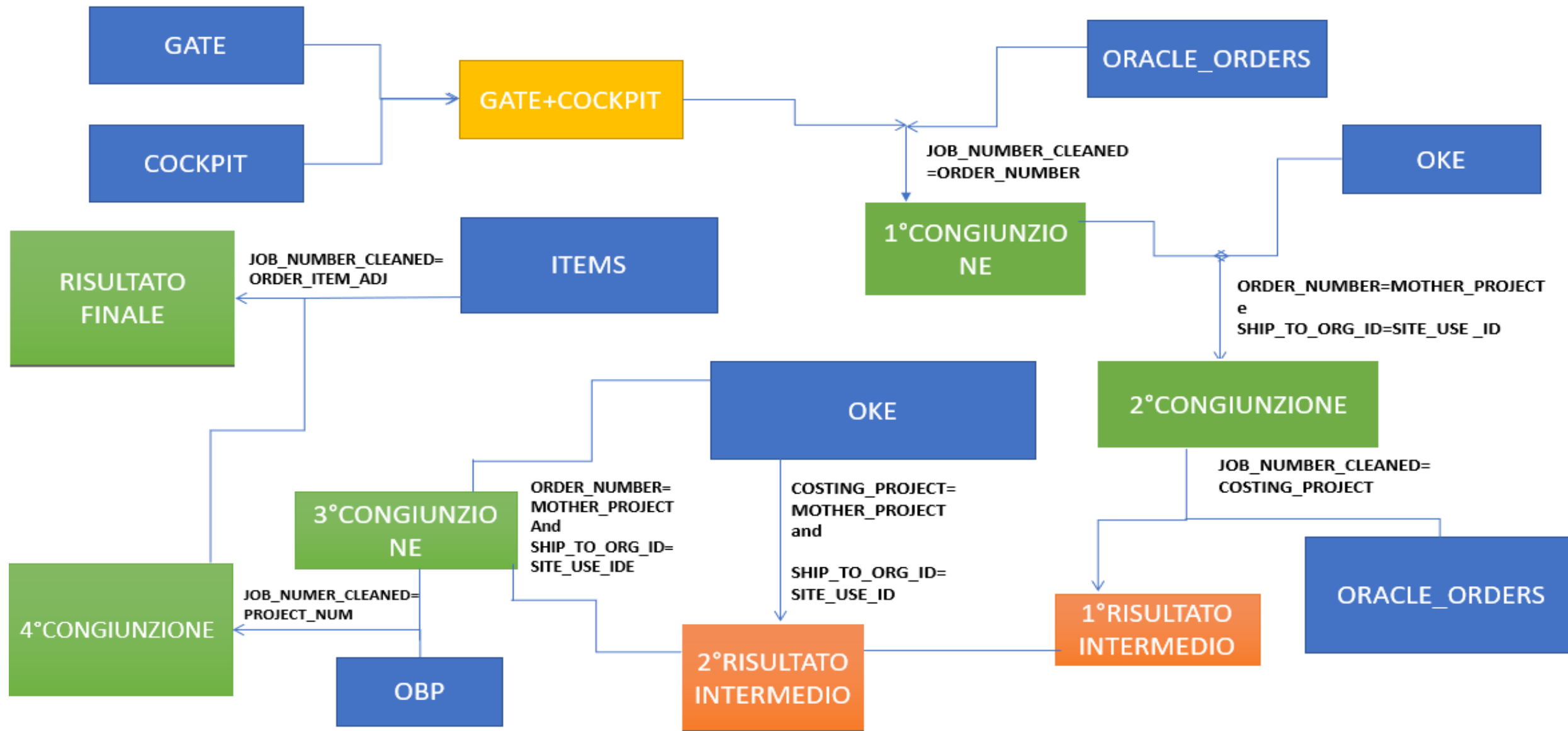# HISTORICAL DB ORDER REPAIRS :

**HISTORICAL DB ORDER REPAIRS**
For **TURBINE** and
**COMPRESSOR items**

**GOAL:**

- **Retrieve** information linked to client
- **Improve data quality**

**ISSUES DUE TO HISTORICAL DATA RELATED TO ORDERS**

**MISSING CLIENT ADRESS ( SOGGE )**

**DIRTY-DATA**

# HISTORICAL DB ORDER REPAIRS

## Combination from multiple DATA SOURCES :

-ERP

-TOOLS

| | | |
|---|---|---|
| **ERP OM** (Oracle Management) | **DATASET** Oracle Orders | **ORDERS** |
| **ERP OKE** (Oracle Project Contracts) | **DATASET** Oke | **CONTRACTS** |
| **TOOL OBP** (Order Booking Process) | **DATASET** OBP | **PROCESSED ORDER** |
| **TOOL ROP** (Repair One Portal) | **DATASET** Cockipit | relativo alle → **REPAIRS** |

**Main steps**:

merge(how='left')

- Convert data types keys

- Handle with duplicates and null values
- Data cleaning
- Remove duplicate columns

# HISTORICAL  DB ORDER REPAIRS

## COVERTING INTO STRING A DATA TYPE KEY TO ENABLE A COMBINATION

```python
orders_for_repairs['ORDER_NUMBER']=orders_for_repairs['ORDER_NUMBER'].astype(int)
orders_for_repairs['ORDER_NUMBER']=orders_for_repairs['ORDER_NUMBER'].astype(str)
merge1=gate.merge(orders_for_repairs,left_on='JOB_NUMBER_CLEANED', right_on='ORDER_NUMBER', how='left')
```

## CLEANING DATA AND REMOVING DUPLICATES VALUES

```python
gate['JOB_NUMBER_CLEANED']=gate['JOB_NUMBER_CLEANED'].apply(lambda x : x.split(',')[0])
gate['JOB_NUMBER_CLEANED']=gate['JOB_NUMBER_CLEANED'].apply(lambda x : x.split(' ')[0])
gate['JOB_NUMBER_CLEANED']=gate['JOB_NUMBER_CLEANED'].apply(lambda x : x.split('/')[0])
gate['JOB_NUMBER_CLEANED']=gate['JOB_NUMBER_CLEANED'].apply(lambda x : x.split('.')[0])
gate['JOB_NUMBER_CLEANED']=gate['JOB_NUMBER_CLEANED'].apply(lambda x : x.split('-')[0])
gate['JOB_NUMBER_CLEANED']=gate['JOB_NUMBER_CLEANED'].apply(lambda x : x.split('_')[0])
gate['JOB_NUMBER_CLEANED'].drop_duplicates(inplace=True)
```

## REMOVING DUPLICATE COLUMNS

```python
merge1.drop(columns='COMPONENT_CATEGORY_y',inplace=True)
merge1.rename(columns={'COMPONENT_CATEGORY_x':'COMPONENT_CATEGORY'}, inplace=True)
```

## REMOVING ROWS  WITH  KEY'S NULL VALUES

```python
columns_cockpit=gate.columns.tolist()
columns_cockpit.append('SOGGE_Cockpit')
colums_to_remove=['CHARACT_FOUND', 'Unnamed: 27', 'Unnamed: 0', 'Unnamed: 28',
                  'ITEM_CATEGORY_COCKPIT', 'ANNO_RIFERIMENTO', 'serial_number',
                  'Unnamed: 26', 'PROJECT_NUMBER_COCKPIT', 'JOB_NUMBER_ADJUSTED']
columns_cockpit=list(set(columns_cockpit)-(set(colums_to_remove)))

cockpit=cockpit[columns_cockpit]
cockpit['ANNO_RIFERIMENTO']=pd.DatetimeIndex(cockpit['G3_ACTUAL_END_DATE']).year
cockpit.dropna(subset={'ANNO_RIFERIMENTO'},inplace=True)
cockpit.drop_duplicates(subset='JOB_NUMBER',inplace=True)

gate=gate.append(cockpit)

gate.to_excel('gateWithCockpit.xlsx')
```

# HISTORICAL DB ORDER REPAIRS

## LIVELLO DI TESTATA (HEADER)



## LIVELLO DI LINEE (LINES)



## LIVELLO DI COPERTURA (SHEET_FOR_COVERAGE_ELAB)

- Significant increasing in COVERAGE of client address (SOGGE)

| | ALL | FLORENCE | HOUSTON |
|---|---|---|---|
| **DISTINCT JOB_NUMBER** | *11934* | | |
| **DISTINCT JOB_NUMBER (2017-2020)** | 3746 | 1194 | 443 |
| **SOGGE FOUND(%) (2017-2020)** | 84% | 70% | 93% |
| **NO SOGGE(%) (2017-2020)** | 16% | 30% | 7% |
| **DISTINCT JOB_NUMBER (2012-2016)** | 8201 | 1680 | 2033 |
| **SOGGE FOUND(%) (2012-2016)** | 51% | 58% | 45% |
| **NO SOGGE(%) (2012-2016)** | 49% | 42% | 55% |

# CLASSIFICATION OF TURBINE AND COMPRESSOR ITEMS

## MAIN ISSUE

Uncorrect predictions for some types of classifier's input

## KEY SOLUTION

ETL Data pipeline:

pre-processing data :



**EXISTING CLASSIFIER**

ORDERED_ITEM | INPUT DEL CLASSIFICATORE DESCRIPTION | MACHINE_SECTION_NUMBER

**OUTPUT DEL CLASSIFICATORE**

**CAP_NO_CAP**

| CAP | NO_CAP |

**TECH**

| GT | CC | AERO | RC | ST | EXP |

**FAMILY_PRED**

| GT family | CC family | AERO family | RC family | ST family | EXP family |

# CLASSIFICATION OF TURBINE AND COMPRESSOR ITEMS

# CLASSIFICATION OF TURBINE AND COMPRESSOR ITEMS

**LOADING DATA SOURCES:**

-ORDERED ITEM
-DESCRIPTION(optional)

**UPDATE ORDER ITEM**

**RETRIEVE DESCRIPTION From INVENTORY**

DB ARGO

Query_ITEM_SUPERATI.*sql*

SORGENTE CON CODICI DEGLI ITEM IN DISUSO E CODICI DEGLI ITEM RECENTI

OLD_CODE
ORDERED_ITEM

FONTE DATI IN INPUT CONTENENTE CODICI DEGLI ITEM

ITEM VALIDI

DB ARGO

Query_Descrizione_Inventory.SQL

INVENTARIO

OLD_CODE=
ORDERED_ITEM

SORGENTE CON CODICI DEGLI ITEM VALIDI

DATI CON DESCRIZIONE PROVENIENTE DALL'INVENTARIO

# CLASSIFICATION OF TURBINE AND COMPRESSOR ITEMS

**BUSINESS LOGICS :**
-identify (**NO_CAP**) items
-assign a FAMILY

**DATI RELATIVI A COMPONENTI CHE VENGONO CLASSIFICATI SECONDO LE REGOLE DETERMINISTICHE**

| COMPONENTI APPARTENENTI AD UNA FAMIGLIA AUSILIARIA | COMPONENTI APPARTENENTI AD UNA FAMIGLIA ALTERNATIVA RISPETTO A QUELLE PREVISTE DAL CLASSIFICATORE |
|---|---|
| PREFIXES OF ORDERED_ITEM [I, V]<br><br>FAMILY_PREDICTION= AUX<br><br>FAMILY_PRED_EXPLANATION = HARD_RULES<br><br>CAP_NO_CAP= NO_CAP | PREFIXES OF ORDERED_ITEM [N,X,Y,1,C]<br><br>FAMILY_PREDICTION= OTH<br><br>FAMILY_PRED_EXPLANATION = HARD_RULES<br><br>CAP_NO_CAP= NO_CAP |

# CLASSIFICATION OF TURBINE AND COMPRESSOR ITEMS

## PYTHON IMPLEMENTATION OF BUSINESS HARD RULES

```python
def apply_hard_Rules(dataframe):


    dataframe1=dataframe


    dataframe1=dataframe1.assign(FAMILY_PREDICTION="",FAMILY_PRED_EXPLANATION="",
                    TECH="",CAP_NOCAP="")


    prefixes_aux=('I','V')
    prefixes_oth=('N','X','Y','1C','1P','1X','LC')


    item_with_Prefix_aux=dataframe1['ORDERED_ITEM'].str.startswith(prefixes_aux)
    item_with_Prefix_oth=dataframe1['ORDERED_ITEM'].str.startswith(prefixes_oth)


    dataframe1['FAMILY_PREDICTION'][item_with_Prefix_aux]="AUX"
    dataframe1['FAMILY_PREDICTION'][item_with_Prefix_oth]="OTH"
    dataframe1['FAMILY_PRED_EXPLANATION'][dataframe1['FAMILY_PREDICTION']!=""]="HARD RULES"
    dataframe1['CAP_NOCAP'][dataframe1['FAMILY_PRED_EXPLANATION']=="HARD RULES"]="NO CAP"


    return dataframe1
```
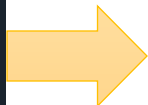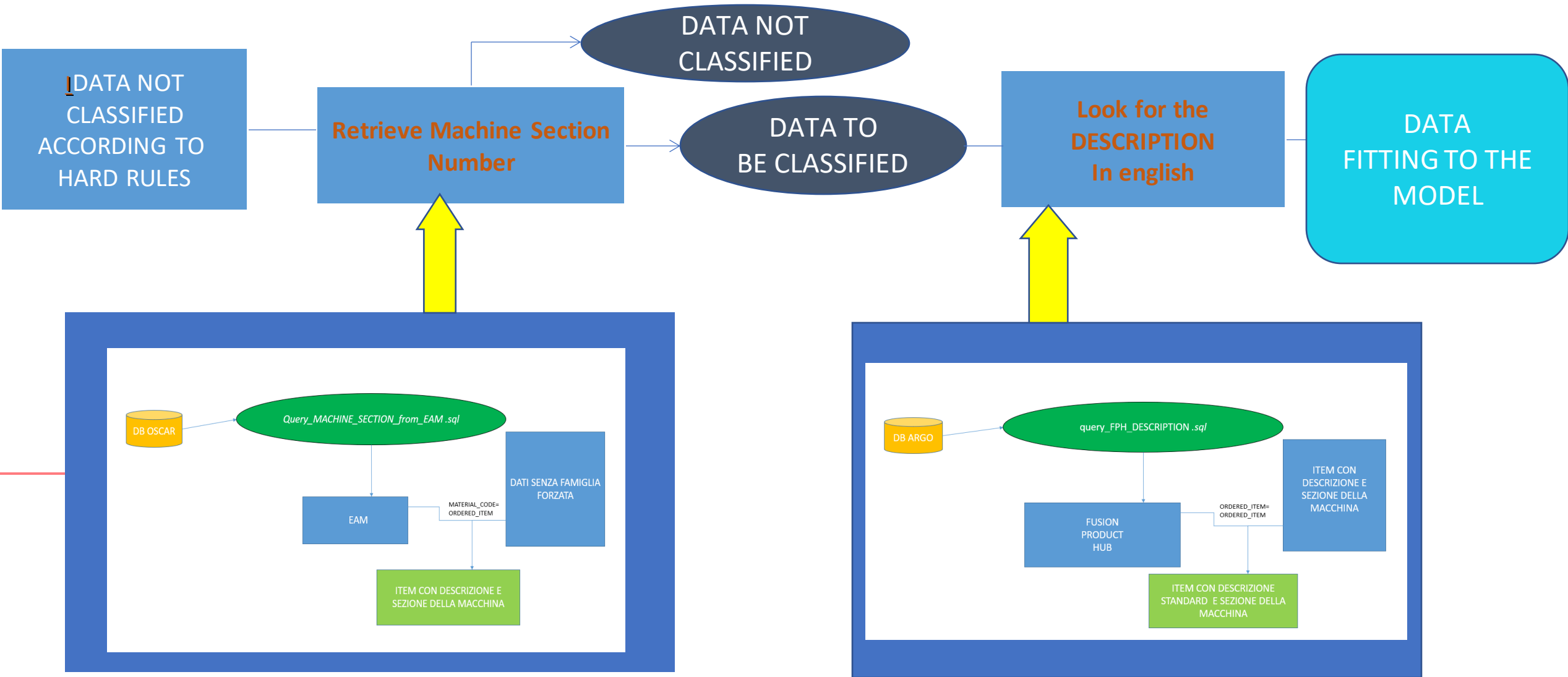
## EXAMPLE OF CLASSIFICATION MADE BY HARD RULES BASED ON BUSINESS LOGICS

| ORDERED_ITEM | DESCRIPTION | FAMILY_PREDICTION | FAMILY_PRED_EXPLANATION | TECH | CAP_NOCAP |
|---|---|---|---|---|---|
| ISM021162007 | BLOCK, SAFETY PARKER | AUX | HARD RULES | | NO CAP |
| N5606P08001G11 | GASKET, SPIRAL WOUND | OTH | HARD RULES | | NO CAP |
| N14TP29024 | BOLT,HEX HEAD | OTH | HARD RULES | | NO CAP |
| N5606P02001G11 | GASKET, SPIRAL WOUND | OTH | HARD RULES | | NO CAP |
| N403P75 | WASHER,LOCK-EXTERNAL TOOT | OTH | HARD RULES | | NO CAP |
| IS200VSVOH1B/RM | GAS TMR PK,MK6 | AUX | HARD RULES | | NO CAP |
| 1X1308C2A000017 | LABYRINTHE ENTRETOISE | OTH | HARD RULES | | NO CAP |
| 1X1308C2AC00008 | JOINT DE CORPS DEPALIER COTE BUTEE | OTH | HARD RULES | | NO CAP |
| 1X1308C2A000013 | LABYRINTHE OUIE ROUE 3 | OTH | HARD RULES | | NO CAP |
| ISM021162001 | ACCUMULATOR, PED | AUX | HARD RULES | | NO CAP |
| IS200VAICH1C/RM | VME ANALOG INPUT CARD (REMAN) | AUX | HARD RULES | | NO CAP |
| 1X1305A1A200001 | JOINT D1 POUR GV REF:2-343 | OTH | HARD RULES | | NO CAP |
| IRJ0601721 | O-RING, THERMOWELL* | AUX | HARD RULES | | NO CAP |
| N272QP00039 | BODY-BOUND LOCK NUTS | OTH | HARD RULES | | NO CAP |
| 1X1308C2AC00011 | CALE DE REGLAGE PELABLE RECONSTITUEE - EP.= 3 MM | OTH | HARD RULES | | NO CAP |
| 1X1308C2AC00007 | JOINT CAPOT COTE ENTRAINEMENT | OTH | HARD RULES | | NO CAP |
| IRF318460137 | SACCA *300LT 10964800225I | AUX | HARD RULES | | NO CAP |
| ILCWBUSR0075 | THRUST PAD WITH HOLE | AUX | HARD RULES | | NO CAP |
| 1X1308C2A000012 | LABYRINTHE OUIE ROUE 2 | OTH | HARD RULES | | NO CAP |

# CLASSIFICATION OF TURBINE AND COMPRESSOR ITEMS

-SPEED UP the classification task

-IMPROVE **the outcome**

## DATA ETL PIPELINE OVERFLOW IMPLEMENTATION

```python
import pandas as pd
from hard_rules import applyHard_Rules
from items import getItemsWithADescription
from items import getItemsWithForcedFamily
from items import getItemsWith_MSN_and_Des
from items import getInputForClassifier
from items import getItemsWithoutFamily
from items import getInputForClassifier

def main():
    input_filename='C:\\Users\\dsoumat\\Desktop\\Items\\INPUT_ITEM_CODES_STEFANO@20210429.xlsx'
    data_before_HardRules=getItemsWithADescription(input_filename)
    data_after_hard_rules=applyHard_Rules(data_before_HardRules)
    items_with_Forced_Family=getItemsWithForcedFamily(data_after_hard_rules)
    item_without_Family=getItemsWithoutFamily(data_after_hard_rules)
    items_with_Machine_Section_Number_and_Des=getItemsWith_MSN_and_Des(data_after_hard_rules)
    input_data_for_classifier=getInputForClassifier(items_with_Machine_Section_Number_and_Des)
    classifier_output=pd.read_excel('C:\\Users\\dsoumat\\Desktop\\Items\\OUTPUT_CAPITAL_FAMILY.xlsx')
    classifier_output=[items_with_Forced_Family.columns]
    output=classifier_output.append(items_with_Forced_Family)
    final_family_assignment=output.append(item_without_Family)
    final_family_assignment= final_family_assignment.to_excel('final_family_assignment.xlsx')

if __name__ == "__main__":
    main()
```

## OUTPUT

**DATA CLASSIFIED BY HARD RULES**

**DATA NOT CLASSIFIED**

**DATA CLASSIFIED BY SUPERVISED MODEL**