

Système décisionnel de gestion d'absence

-Rapport-

Encadré par : **Pr. Imade BENELALLAM**

Réalisé par :

- ENNAJI AYOUB
- ELHOURI MOHAMED EL MEHDI

Filière : DSE

Semestre : S4

TABLE DES MATIERES

- I- DOSSIER JUPYTER 7
- II- VISUALISATION DU DAG 7
 - 1- LOAD_NEW_FILES : 8
 - 2- VERIFY_NEW_INPUT : 9
 - 3- CLEAN_COLUMNS : 9
 - 4- CLEAN_DATA : 10
 - 5- DELETE_TMP_FILE : 11

LISTE DES FIGURES

FIGURE 1: DATA_ FOLDER DE JUPYTER-----	7
FIGURE 2: VISUALISATION DES DAGS 1-----	7
FIGURE 3: VISUALISATION DES DAGS 2-----	8
FIGURE 4: FONCTION « LOAD_NEW_FILE »-----	8
FIGURE 5: FONCTION « VERIFY_NEW_INPUT »-----	9
FIGURE 6: FONCTION « CLEAN_CULUMNS »-----	9
FIGURE 7: FONCTION : « CLEAN_DATA »-----	10
FIGURE 8: FONCTION « DELET_TMP_FILES »-----	11

INTRODUCTION

Le corps professoral et administratif souhaite contrôler et suivre l'assiduité des étudiants lors de cette phase d'enseignement à distance. D'où l'idée de réaliser un Système décisionnel de gestion d'absence.

Ce projet s'inscrit dans le cadre de l'élément de module Data Warehouse & Business Analytics. Il est réalisé sur cloud en utilisant Airflow, Pandas, pandasql, Os et shutil.

I- Dossier jupyter

Jupyter Notebook permet aux développeurs de partager du code et de l'exécuter dans la même interface utilisateur. Il peut associer du code, des graphiques, des visualisations et du texte dans des notebooks - ou cahiers - partageables qui s'exécutent dans un navigateur web.



Figure 1: data_ folder de jupyter

Il contient 3 dossiers :

- Input : ou on ajoute les fichiers à traiter
- Output : contient le résultat du traitement
- Processed : contient des copies des fichiers existant dans le dossier input après traitement.

II- Visualisation du Dag

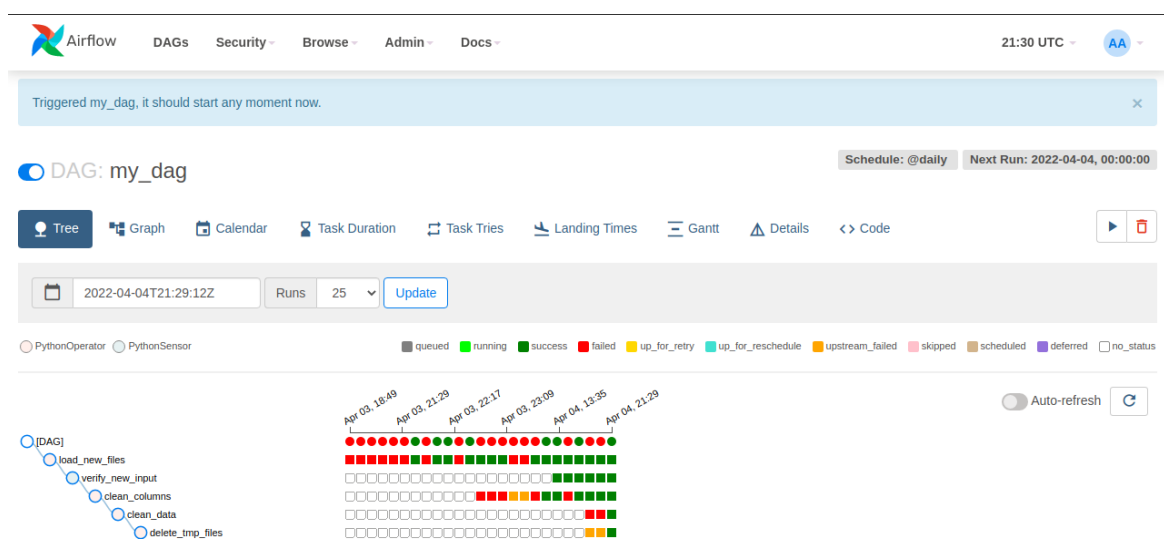


Figure 2: Visualisation des Dags 1

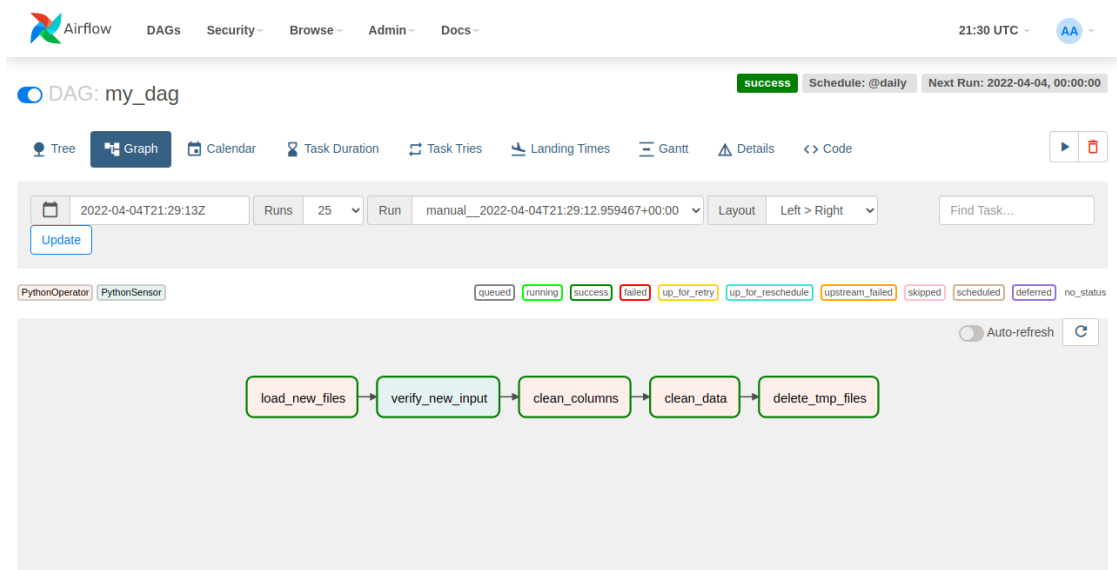


Figure 3: visualisation des Dags 2

On remarque la présence de 5 tâches qui sont :

- 1- **Load_new_files** : permet d'importer des nouveaux fichiers depuis le dossier input ;

```

1 def _load_new_files(ti):
2     dfs = []
3     id = "";
4     for filename in os.listdir(INPUT):
5         if filename not in os.listdir(PROCESSED):
6             df = pd.read_csv(INPUT + '/' + filename, skiprows=7,
7                             delimiter="\t", encoding="UTF-16")
8             df['Meeting_Id'] = _get_meeting_id(INPUT + '/' + filename)
9             dfs.append(df)
10            shutil.copyfile(INPUT + '/' + filename, PROCESSED + '/' + filename)
11    if len(dfs) == 0:
12        ti.xcom_push(key='new_files_found', value=False)
13        return
14
15    result = pd.concat(dfs)
16    os.mkdir(TMP + '/first/')
17    os.chdir(TMP + '/first/')
18    result.to_csv(str( len(os.listdir(TMP + '/first/')) + 1)
19                  + '.csv', encoding="UTF-16")
20    ti.xcom_push(key='new_files_found', value=True)

```

Figure 4: fonction « load_new_file »

2- **Verify_new_input** : vérifie si il y'a des nouveaux fichiers non-traités.

```
1 def _verify_new_input(ti):
2     return ti.xcom_pull(key='new_files_found',
3                           task_ids=['_load_new_files'])
```

Figure 5: fonction « verify_new_input »

3- **Clean_columns** : permet de modifier les noms des colonnes, la forme, ajouter ou supprimer des champs...

```
1 def _clean_columns():
2     df = pd.read_csv(TMP + '/first/' + str( len(os.listdir(TMP + '/first/'))
3                                     + '.csv', encoding="UTF-16")
4
5     df.drop('Adresse de courrier',axis=1, inplace=True)
6
7     df.columns=['ligne','nom_participant','h_arr', 'h_dep','duree_en_mins',
8                'role','id_participant', 'id_meeting']
9     df = df[['id_meeting', 'h_arr', 'h_dep', 'duree_en_mins', 'id_participant',
10             'nom_participant','role']]
11
12     raw_ha = [datetime.strptime(i, '%m/%d/%Y, %H:%M:%S %p') for i in df['h_arr']]
13     raw_hd = [datetime.strptime(i, '%m/%d/%Y, %H:%M:%S %p') for i in df['h_dep']]
14
15     df['meeting_date'] = [datetime.strftime(i, '%d/%m/%Y') for i in raw_ha]
16     df['duree_en_mins'] = [ int((i - j).total_seconds()//60)
17                           for i, j in zip(raw_hd, raw_ha)]
18     df['h_arr'] = [datetime.strftime(i,'%H:%M:%S') for i in raw_ha]
19     df['h_dep'] = [datetime.strftime(i,'%H:%M:%S') for i in raw_hd]
20     df['role'] = [i[0] for i in df['role']]
21
22     os.mkdir(TMP + '/second/')
23     os.chdir(TMP + '/second/')
24     df.to_csv(str( len(os.listdir(TMP + '/second/')) + 1) + '.csv',
25               encoding="UTF-16")
```

Figure 6: fonction « Clean_culumns »

- 4- **Clean_data** : permet de créer des tables intermédiaires dont le but d'avoir une table finale contenant les données souhaitées.

```
1 def _clean_data():
2
3     df = pd.read_csv(TMP + '/second/' + str( len(os.listdir(TMP + '/second/')))
4                     + '.csv', encoding="UTF-16")
5
6     query = '''SELECT id_meeting,meeting_date ,id_participant, nom_participant, role,
7                 min(h_arr) AS premiere_arrivee,
8                 max(h_dep) AS dernier_depart, SUM(duree_en_mins) AS duree_totale_presence
9     from df
10    GROUP BY id_meeting,meeting_date,id_participant, nom_participant, role
11    '''
12     df2 = sqldf(query,locals())
13
14     query2 = '''SELECT id_meeting, id_participant AS id_organisateur,
15                    nom_participant AS nom_organisateur,
16                    premiere_arrivee AS arrivee_organisateur,
17                    dernier_depart AS depart_organisateur
18    FROM df2 WHERE role = '0'
19    '''
20     df3 = sqldf(query2,locals())
21
22     df4 = pd.merge(df2, df3, how='inner', on='id_meeting')
23
24     df5 = df4.drop('role',axis=1)
25
26     pa = [datetime.strptime(i,'%H:%M:%S') for i in df5['premiere_arrivee']]
27     ao = [datetime.strptime(i,'%H:%M:%S') for i in df5['arrivee_organisateur']]
28
29     df5['retard_en_mins'] = [int((i - j).total_seconds()/60)
30                            if (i-j).total_seconds() >= 0 else 0 for i,j in zip(pa,ao)]
31
32     dd = [datetime.strptime(i,'%H:%M:%S') for i in df5['dernier_depart']]
33     do = [datetime.strptime(i,'%H:%M:%S') for i in df5['depart_organisateur']]
34
35     df5['exces_en_mins'] = [int((i - j).total_seconds()/60)
36                            if (i-j).total_seconds() >= 0 else 0 for i,j in zip(dd,do)]
37
38     df5['duree_reelle_presence'] = [i -j -k if i-j-k >= 0 else 0 for i,j,k in
39 zip(df5['duree_totale_presence'],df5['exces_en_mins'],df5['retard_en_mins'])]
40
41     df5['duree_meeting_en_mins'] = [int((i - j).total_seconds()/60) for i,j in zip(do,ao)]
42     df5 = df5[['id_meeting', 'meeting_date','duree_meeting_en_mins', 'id_organisateur',
43               'nom_organisateur', 'id_participant', 'nom_participant', 'retard_en_mins',
44               'duree_reelle_presence']]
45
46     os.chdir(OUTPUT)
47     df5.to_csv(str( len(os.listdir(OUTPUT)) + 1) + '.csv', encoding="UTF-16")
```

Figure 7: fonction : « Clean_data »

5- `Delete_tmp_file` : elle supprime les fichiers temporaires.

```
1 def _delete_tmp_files():
2     shutil.rmtree(TMP)
3
4 with DAG("my_dag", start_date=datetime(2022, 4, 3), schedule_interval="@daily",
5         catchup=False ) as dag:
6
7     load_new_files = PythonOperator(
8         task_id = "load_new_files",
9         python_callable = _load_new_files
10    )
11
12    verify_new_input = PythonSensor(
13        task_id = "verify_new_input",
14        python_callable = _verify_new_input,
15        soft_fail = True
16    )
17
18    clean_columns = PythonOperator(
19        task_id = "clean_columns",
20        python_callable = _clean_columns
21    )
22
23    clean_data = PythonOperator(
24        task_id = "clean_data",
25        python_callable = _clean_data
26    )
27
28    delete_tmp_files = PythonOperator(
29        task_id = "delete_tmp_files",
30        python_callable = _delete_tmp_files
31    )
32
33    load_new_files >> verify_new_input >> clean_columns >> clean_data >> delete_tmp_files
```

Figure 8: fonction « `Delete_tmp_files` »

CONCLUSION

La réalisation de ce travail a été très bénéfique et très enrichissante. Elle nous a permis d'utiliser des technologies nouvelles tel que Airflow, Pandas, pandasql...

Cette solution s'avère être très utile pour la gestion de l'absence. Nous allons l'améliorer au fur et à mesure de l'avancement du cours afin de produire une solution permettant d'aider le corps professoral et administratif à prendre les mesures nécessaires pour gérer assurer la présence.