

Módulo 5

Integrador de ETL

Romina Barrabino

1) Extracción de Datos:

```
1 #1) Extracción de datos
2
3 #Importo Pandas y numpy
4 import pandas as pd
5
6 #Supongo la extracción de datos desde un archivo CSV y JSON
7
8 csv_Asistencias={
9     'Id_Numero':[1,2,3,4,5],
10    'Fecha':['2022/01/03','2022/02/03','2023/05/10','2024/03/11','2024/01/31'],
11    'Nombre':['Susana','Pablo','Teresa','Doña Susana','Juan'],
12    'Texto':['Asistio','Falto','Licencia','Falto','Vacaciones']
13 }
14
15 json_AsistenciasNuevas={
16     'Id_Numero':[6,7],
17     'Fecha':['2024/07/28','2025/01/31'],
18     'Nombre':['Le.onel','Pablo']
19     'Texto':['Licencia','Asistio']
20 }
21
```

2) Transformación de Datos:

```
22 # 2) Transformación de datos:
23
24 #Creo un DataFrames desde csv_Asistencias y json_AsistenciasNuevas
25 df_csv=pd.DataFrame(csv_Asistencias)
26 df_json=pd.DataFrame(json_AsistenciasNuevas)
27 #Uno los Dataframes
28 df=pd.concat([df_csv,df_json],ignore_index=True)
29
30 #Instalo e importo Data prep.clean
31 from dataprep.clean import clean_date
32
33 #a)limpieza de valores nulos:
34 #Elimino las filas que tengan altos porcentajes de nulos (0.8)
35 df.dropna(axis=0,inplace=True,thresh=int(df.shape[1]*0.8))
36
37 #b) Normalización de nombres de columnas:
38 #Elimino los duplicados de la columna Id_Numero
39 data_cleaned=df.drop_duplicates(subset=['Id_Numero'])
40 #Elimino los espacios en blanco de la columna Nombre
41 data_cleaned['Nombre']=data_cleaned['Nombre'].apply(lambda x:str(x).strip(x))
42 #Pongo en mayuscula la primer letra y minuscula en el resto los registros de la columna Nombre
43 data_cleaned['Nombre']=data_cleaned['Nombre'].apply(lambda x:str(x).strip().title())
44
```

```

45 #b) Conversión de tipos de datos:
46 #Modifico el tipo de datos si no son correcto
47 data_cleaned=data_cleaned.astype({
48     'Id_Numero':'Int',
49     'Nombre':'Varchar(30)',
50     'Fecha':'Date',
51     'Texto':'Varchar(50)'})
52 #Convierto la columna Fechas al formato YYYY-MM-DD
53 data_cleaned['Fecha']=pd.to_datetime(data_cleaned['Fecha'], errors='coerce',format='%Y-%m-%d')
54

```

3) Carga de Datos:

```

54
55 #3) Cargo los datos:
56 #Llamo a data_cleaned como df
57 df=data_cleaned
58
59 print("\nDatos Cargados (Centralizados):")
60 print(df)
61

```

4) Implementación de Pruebas Automatizadas:

```

62 #4) Implementación de pruebas automaticas
63
64 #Importo Great expectations
65 import great_expectations as ge

```

```

66 #Convierto el Dataframe df en un objeto Great Expectations
67 df_ge=ge.from_pandas(data_loaded)
68
69 # Expectativa1
70 df_ge.expect_column_values_to_not_be_null('Id_Numero')
71 # Expectativa2
72 df_ge.expect_column_values_to_not_be_null('Nombre')
73 # Expectativa3
74 df_ge.expect_column_values_to_be_unique('Id_Numero')
75 # Expectativa4
76 df_ge.expect_column_values_to_be_of_type('Texto','varchar(50)')
77 # Expectativa5
78 df_ge.expect_column_values_to_match_strftime_format('Fecha','%Y-%m-%d')
79
80 #Ejecución de pruebas automatizadas
81 validation_results=df_ge.validate()
82 print("\nResultados de las Validaciones:")
83 print(validation_results)
84 #Como no pude descargar Great expectations no puedo visualizar los resultados que arroja, pero los detalle con el objetivo
85 de que no resulten en error.

```

5) Reporte de Resultados:

```

85
86 #5)Reporte de resultados:
87
88 #Exporto un reporte de las validaciones a HTML
89 df_ge.save_expectation_suite(discard_failed_expectations=False)
90 #Genero un reporte de resultados
91 validation_report=df_ge.validate()
92 print("\nGeneración de reporte completada.")
93 -----

```

6) Documentación del Proceso:

- Las herramientas utilizadas fueron la página de Databricks, pandas y data prep.clean.
- Los tipos de pruebas implementadas están adjuntadas en la sección 4.
- No adjunte ninguna explicación de los resultados, ya que como no pude descargar Great Expectations no pude verificar los resultados. De igual manera, las expectativas se realizaron con el objetivo de que no devuelvan fallas ni errores.