

Grado en Ingeniería Informática
2024-2025

Arquitectura de Datos

“Practica 1.2 (ETL) ”

Adrian Cortázar Leiva

Jaime Vaquero Rabahieh

Tomás Mendizábal



Esta obra se encuentra sujeta a la licencia Creative Commons **Reconocimiento - No Comercial - Sin Obra Derivada**

Índice

1. Diseño y Construcción de Pipelines	2
1.1. Descripción del proceso ETL	2
1.2. Proceso de extracción	2
1.3. Proceso de limpieza	2
1.4. Proceso de implementación de columnas	12
1.5. Proceso de extracción a mongoDB	20
1.6. Proceso de Carga	29
1.6.1. Estructura final de BBDD	29
1.6.2. Diseño de Esquemas de validación	30
1.7. Ejecución	30
2. Datos importantes a tener en cuenta	42

1. Diseño y Construcción de Pipelines

1.1. Descripción del proceso ETL

El proceso de ETL o de Extracción, Transformación y Carga se define como una serie de pasos a realizar con la finalidad de adaptación de la información a nuevos entornos de almacenamiento (ej. Bases de datos, Data Warehouses).

En esta práctica, el objetivo es aplicar un proceso ETL o pipeline para migrar los datos proporcionados por el enunciado a una estructura de base de datos no SQL como es MongoDB.

Nuestro proceso consiste de los siguientes pasos:

- Carga de los datos anteriores con formatos incompatibles y algunas impurezas.
- Transformación de dichos datos (eliminación de valores nulos, imputación de atributos, limpieza de formatos, estandarización de palabras, etc.).
- Creación de las colecciones y esquemas pertinentes de validación asociados a la óptima representación de la información en la base de datos final.
- Inserción de los datos transformados en la base de datos final.

1.2. Proceso de extracción

En el archivo "load.py" extraemos todos los csv. Esto lo hacemos mediante un diccionario de DataFrames de panda, siendo cada clave el nombre de la tabla y el valor de la key la propia tabla/DataFrame. El código usado es el siguiente:

```
1 import pandas as pd
2 import os
3 import re
4
5 def load_db():
6     # paths usados
7     pwd = os.getcwd()
8     dataset_dir = f"{pwd}/DatasetsArquiViejos"
9
10    # leer datasets y guardarlos en un diccionario
11    df = {}
12    for item in os.listdir(dataset_dir):
13        try:
14            file_name = re.findall("[A-Z][^A-Z]*", item)[0]
15            temp_db = pd.read_csv(f"{dataset_dir}/{item}")
16            df[file_name] = temp_db
17        except IndexError:
18            file_name = item.split('.')[0]
19            temp_db = pd.read_csv(f"{dataset_dir}/{item}")
20            df[file_name] = temp_db
21    return df
22
```

1.3. Proceso de limpieza

En el archivo "change.py" se realizó la limpieza de datos. Para ejecutar este archivo, es necesario importar estas librerías de esta manera:

```

1 import pandas as pd
2 import datetime as dt
3 from dateutil import parser
4

```

Este archivo está formado por varias funciones cada una con un labor en particular:

1. **capitalize column:** Capitaliza todos los valores de una columna. Se le pasa como argumento el diccionario completo y los nombres de la columna a modificar y la tabla a la que pertenece. Código:

```

1 # capitalizar una columna entera
2 def capitalize_column(df: dict, table_name: str, colname: str) -> None:
3     df[table_name][colname] = df[table_name][colname].apply(lambda x: x.upper())
4

```

2. **take attribute:** Busca un atributo en una tabla que tenga un id igual al pasado por argumento. Se hace para tablas en las que existe una columna id. Código:

```

1 # Encontrar un dato en una tabla si su id es igual al pasado por argumento
2 def take_attribute(df: dict, tabla: str, columna: str, id: int) -> str | None:
3     archivo = df[tbl]
4     fila = None
5     for i in range(len(archivo["ID"])):
6         valor = archivo["ID"][i]
7         if id == valor:
8             fila = i
9             return archivo[columna][fila]
10    return None
11

```

3. **find table by column:** Busca una tabla que contenga la columna pedida. Como extra, dicha tabla no debe estar en la lista "tablas tachadas". Esto se hace para que la función sepa que tablas se excluyen de la búsqueda. Código:

```

1 # Encuentra una tabla sin contar las tachadas que contenga una columna con el mismo nombre que
  la pasada por argumento
2 def find_table_by_column(df: dict, tablas_tachadas: list, column: str) -> str | None:
3     for tabla in df:
4         if tabla not in tablas_tachadas:
5             t_analizar = df[tbl]
6             if column in t_analizar:
7                 return tabla
8     return None
9

```

4. **find id:** Busca un id en una tabla (se pasa su nombre como argumento) que se encuentre en la posición también pasada como argumento. Como no todas las columnas de identificación se llaman ID, se resliza un if-else para primero identificar la tabla y luego devolver el id. Código:

```

1 # Encuentra un id según que tabla (hay columnas ID que no se llaman ID) y que fila
2 def find_id(df: dict, tabla: str, fila: int):
3     archivo = df[tbl]
4     if tabla == "Usuarios":
5         return archivo["NIF"][fila]
6     elif tabla == "meteo24":
7         return archivo["PROVINCIA"][fila]

```

```

8     else:
9         return archivo["ID"][fila]
10
11

```

5. **empty data:** Modifica los NaN encontrados. Coge las casillas de las tablas con NaN y las modifica según estos criterios:

- Si es una fila con valor de DIRECCION AUX no nulo y el valor nulo es de nombre de vía, tipo de vía o número de vivienda; usa la función *aux_dir* para eliminar el NaN (se explica después).
- En cualquier otro caso intenta buscar con la función *find table by column* un valor de dicha columna que tenga un id igual al suyo. Si no lo encuentra, se cambia el NaN por uno de estos valores predeterminados:
 - FECHA: 31 de diciembre de 2018
 - OTRO CASO: la estructura **id-nombre de columna-ausente**

Cabe aclarar que este función se diseñó sin ChatGPT, pero fue corregida por esta IA. Código:

```

1  # Para rellenar filas vacias
2  def empty_data(df_dict):
3      for tabla_name, tabla in df_dict.items():
4          if tabla_name != "Codigo":
5              # Convertir a tipo 'object' para poder indicar los casos nulos que no se pueden
6              modificar
7              tabla = tabla.astype(object)
8              nan_positions = tabla.isna() # Identificar posiciones NaN
9              assignment_dict = {} # Diccionario para almacenar asignaciones
10             for columna in nan_positions.columns:
11                 nan_indices = nan_positions.index[
12                     nan_positions[columna]
13                 ] # Se cogen los indices con NaN de dicha columna
14             for i in nan_indices:
15                 valor = aux_dir(
16                     tabla_name, tabla, columna, i
17                 ) # Caso de direccion auxiliar. Si funciona, se salta el resto
18             if valor is None:
19                 tachadas = [tabla_name]
20                 id_val = find_id(
21                     df_dict, tabla_name, i
22                 ) # Obtener una vez el id
23                 while True:
24                     new_tab = find_table_by_column(
25                         df_dict, tachadas, columna
26                     ) # Buscar una tabla con un valor no NaN en la posicion con el
27                     mismo
28                     # id que el del elemento actual
29                     if (
30                         new_tab is None
31                     ): # Caso fallido -> cambiar NaN por valor predeterminado
32                         if columna in [
33                             "FECHA_INSTALACION",
34                             "FECHA",
35                             "FECHA_REPORTE",
36                             "FECHA_INTERVENCION",
37                         ]: # Caso fecha
38                             valor = dt.datetime.strptime(
39                                 "2018-12-31 00:00:00", "%Y-%m-%d %H:%M:%S"
40                             )

```

```

39         else: # Otro caso
40             valor = f"{id_val}-{columna}-ausente"
41             break
42         valor = take_atribute(
43             df_dict, new_tab, columna, id_val
44         ) # Se ecnuebtra la tabla y se checkea el valor
45         if valor is not None:
46             break # Caso exitoso
47         tachadas.append(
48             new_tab
49         ) # Valor NaN, asi que se repite el proceso sin contar con al
50
51     tabla encontrada
52
53     # Añadir al diccionario de asignación
54     if columna not in assignment_dict:
55         assignment_dict[columna] = {}
56         assignment_dict[columna][i] = valor
57
58     # Realizar las asignaciones en bloque usando el diccionario de asignación
59     for columna, valores in assignment_dict.items():
60         tabla.loc[valores.keys(), columna] = list(valores.values())
61
62     # Actualizar la tabla en el diccionario original (si no es una copia)
63     df_dict[tabla_name] = tabla

```

6. **aux_dir**: Es la función sacada de *empty data* que ayuda a eliminar los NaN de las columnas TIPO VIA, NOM VIA y NUM VIA en caso de existir DIRECCION AUX. Se basa en descomponer DIRECCION AUX en TIPO VIA, NOM VIA y NUM VIA y coger el valor que se pide (se elige según qué nombre de columna de los tres se pase por argumento). Esta función se creo con la ayuda de ChatGPT. Código:

```

1  # Para el caso de direccion auxiliar, descomponerla en TIPO_VIA, NOM_VIA y NUM_VIA e insertar
2  # la que se
3  # nombra como argumento
4  def aux_dir(tabla_name, tabla, columna, i):
5      if tabla_name in ["Areas", "Juegos"] and columna in [
6          "TIPO_VIA",
7          "NOM_VIA",
8          "NUM_VIA",
9      ]:
10         if not pd.isna(tabla["DIRECCION_AUX"][i]):
11             lista = tabla["DIRECCION_AUX"][i].split(" ")
12             via = None
13             n = 0
14             if "." in lista:
15                 via = lista[2]
16                 n = 3
17             else:
18                 via = lista[0]
19             nombre, numero = "", ""
20             while n < len(lista):
21                 if lista[n].isdigit():
22                     numero = lista[n]
23                     break
24                 if lista[n] == ",":
25                     break
26                 nombre += lista[n]
27                 n += 1
28             if columna == "TIPO_VIA":

```

```

28         return via
29     elif columna == "NOM_VIA":
30         return nombre
31     else:
32         return numero
33 return None
34

```

7. **reformatear fecha:** Coge cualquier fecha y la pasa a formato mongo DB (YYYY-MM-DD HH:mm:ss). En caso de NaN o error pone la predeterminada (2018-12-31 00:00:00). Esta función se corrigió con ChatGPT. Código:

```

1  # Coge cualquier fecha y la pasa a formato %Y-%m-%d %H:%M:%S.
2  # En caso de NaN se coloca la predeterminada (2018-12-31)
3  def reformatear_fecha(df: dict, table_name: str, column_name: str):
4      tablas = [
5          "Areas",
6          "Encuestas",
7          "Incidencias",
8          "Incidentes",
9          "Juegos",
10         "Mantenimiento",
11         "meteo24",
12     ]
13     columnas = ["FECHA_INSTALACION", "FECHA", "FECHA_REPORTES", "FECHA_INTERVENCION"]
14     if table_name in tablas and column_name in columnas:
15         archivo = df[table_name]
16         fechas_nuevas = []
17         for fecha in archivo[column_name]:
18             try:
19                 f_fmt = parser.parse(fecha)
20                 f_fmt.strftime("%Y-%m-%d %H:%M:%S")
21             except (ValueError, TypeError):
22                 f_fmt = dt.datetime.strptime(
23                     "2018-12-31 00:00:00", "%Y-%m-%d %H:%M:%S"
24                 )
25             fechas_nuevas.append(f_fmt)
26         archivo[column_name] = fechas_nuevas
27

```

8. **delete special:** Se encarga de eliminar las tildes y de los caracteres especiales (menos la ñ y la "-") de las columnas especificadas de todas las tablas posibles. Esta función se diseñó con ayuda de ChatGPT. Código:

```

1  # Elimina tildes y caracteres especiales
2  def delete_special(df: dict):
3      # Columnas donde se eliminan los caracteres especiales
4      lista_tildes = [
5          "DESC_CLASIFICACION",
6          "BARRIO",
7          "DISTRITO",
8          "NOMBRE",
9          "TIPO_INCIDENTE",
10         "GRAVEDAD",
11         "TIPO_INTERVENCION",
12         "DIRECCION_AUX",
13     ]
14     # Diccionario para reemplazar letras con tildes
15     replacements = {

```

```

16         "á": "a",
17         "é": "e",
18         "í": "i",
19         "ó": "o",
20         "ú": "u",
21         "Á": "A",
22         "É": "E",
23         "Í": "I",
24         "Ó": "O",
25         "Ú": "U",
26     }
27     for tabla_n in df:
28         tabla = df[tabla_n]
29         for columna in tabla:
30             if columna in lista_tildes:
31                 for accented_char, unaccented_char in replacements.items():
32                     tabla[columna] = tabla[columna].str.replace(
33                         accented_char, unaccented_char
34                     )
35             tabla[columna] = tabla[columna].str.replace(
36                 r"[^a-zA-Z0-9 ñÑ-]", "", regex=True
37             )
38

```

9. **formato tlf:** Cambia el formato de los números de teléfono a int (ej.: pasa de "+xx xxx xxx xxx" a xxxxxxxxxxxx). Código:

```

1  # Convierte todos los números de teléfono al formato int (número de 11 cifras)
2  def formato_tlf(df: dict):
3      columna = df["Usuarios"]["TELEFONO"]
4      for i in range(len(columna)):
5          valor = columna[i]
6          if " " or "+" in valor:
7              espacio = ""
8              j = 0
9              while "34" not in espacio:
10                 espacio += valor[j]
11                 j += 1
12             numero = valor[j:]
13             lista = numero.split(" ")
14             new_tlf = "34"
15             for x in lista:
16                 new_tlf += x
17             columna[i] = new_tlf
18

```

10. **no duplicates:** Elimina las filas duplicadas de cada tabla. Primero las que tienen valores iguales (sin contar el id) y luego las que tienen mismo id. Código:

```

1  # Elimina duplicados
2  def no_duplicates(df: dict):
3      for tabla_n in df:
4          if tabla_n == "meteo24" or tabla_n == "Codigo":
5              continue
6          else:
7              columnas_sin_primary = list(df[tabla_n].columns)
8              primary_key = columnas_sin_primary.pop(0)
9              # quitar las filas que repitan todo menos la primary key
10             df[tabla_n] = df[tabla_n].drop_duplicates(

```



```

11         subset=columnas_sin_primary, keep="first"
12     )
13     # quitar las filas que repitan la primary key
14     df[tabela_n] = df[tabela_n].drop_duplicates(
15         subset=[primary_key], keep="first"
16     )
17

```

11. **adjust_gps**: Ajusta las coordenadas gps para que concuerden todas y para poder relacionar áreas con juegos. Las coordenadas no pertenecientes a los que indica la función se sustituyen por la mediana de la lista (cumple las condiciones de la función). Código:

```

1 def adjust_gps(df: dict) -> None:
2     """Limpia los datos del GPS"""
3     # Comprueba que la latitud esta entre [-90, 90] y la longitud entre [-180, 180]
4     lat_area = df["Areas"]["LATITUD"]
5     long_area = df["Areas"]["LONGITUD"]
6     lat_juego = df["Juegos"]["LATITUD"]
7     long_juego = df["Juegos"]["LONGITUD"]
8     # si el valor es erroneo entonces coge la mediana
9     for i in range(len(lat_area)):
10         if abs(lat_area[i]) > 90:
11             lat_area[i] = lat_area.median()
12     for i in range(len(long_area)):
13         if abs(long_area[i]) > 180:
14             long_area[i] = long_area.median()
15     for i in range(len(lat_juego)):
16         if abs(lat_juego[i]) > 90:
17             lat_juego[i] = lat_juego.median()
18     for i in range(len(long_juego)):
19         if abs(long_juego[i]) > 180:
20             long_juego[i] = long_juego.median()
21     # Ajusta las áreas a 3 decimales para que los juegos puedan pertenecer a estas.
22     df["Areas"]["LATITUD"] = df["Areas"]["LATITUD"].apply(lambda x: "{:3.3f}".format(x))
23     df["Areas"]["LONGITUD"] = df["Areas"]["LONGITUD"].apply(
24         lambda x: "{:3.3f}".format(x)
25     )
26     # adjust Juegos
27     df["Juegos"]["LATITUD"] = df["Juegos"]["LATITUD"].apply(
28         lambda x: "{:3.3f}".format(x)
29     )
30     df["Juegos"]["LONGITUD"] = df["Juegos"]["LONGITUD"].apply(
31         lambda x: "{:3.3f}".format(x)
32     )
33

```

12. **adjust_ETRS89**: Modifica las coordenadas ETRS89 de las áreas y los juegos para convertirlas en GPS. Código:

```

1 def adjust_ETRS89(df: dict) -> None:
2     """Modifica los datos ETRS89 a GPS."""
3
4     # Ajusta áreas
5     df["Areas"]["COORD_GIS_X"] = df["Areas"]["COORD_GIS_X"].apply(
6         lambda x: "{:3.3f}".format(x % 180)
7     )
8     df["Areas"]["COORD_GIS_Y"] = df["Areas"]["COORD_GIS_Y"].apply(
9         lambda x: "{:3.3f}".format(x % 90)
10    )
11

```

```

12 # Ajusta juegos
13 df["Juegos"]["COORD_GIS_X"] = df["Juegos"]["COORD_GIS_X"].apply(
14     lambda x: "{:3.3f}".format(x % 180)
15 )
16 df["Juegos"]["COORD_GIS_Y"] = df["Juegos"]["COORD_GIS_Y"].apply(
17     lambda x: "{:3.3f}".format(x % 90)
18 )
19

```

13. **enum checker:** Verifica los enum de las tablas. Código:

```

1 # Verifica los enum básicos
2 def enum_checker(db: dict):
3     for n_tabla in db:
4         tabla = db[n_tabla]
5         for n_columna in tabla:
6             new_columna = []
7             columna = tabla[n_columna]
8
9             if n_tabla == "Areas" and n_columna == "ESTADO":
10                 for valor in columna:
11                     if valor != "OPERATIVO":
12                         new_columna.append(None)
13                     else:
14                         new_columna.append(valor)
15
16             elif n_tabla == "Juegos" and n_columna == "ESTADO":
17                 for valor in columna:
18                     if valor not in ["OPERATIVO", "REPARACION"]:
19                         new_columna.append(None)
20                     else:
21                         new_columna.append(valor)
22
23             elif n_tabla == "Incidencias" and n_columna == "TIPO_INCIDENCIA":
24                 for valor in columna:
25                     if valor not in [
26                         "DESGASTE",
27                         "ROTURA",
28                         "VANDALISMO",
29                         "MAL FUNCIONAMIENTO",
30                     ]:
31                         new_columna.append(None)
32                     else:
33                         new_columna.append(valor)
34
35             elif n_tabla == "Incidencias" and n_columna == "ESTADO":
36                 for valor in columna:
37                     if valor not in ["ABIERTA", "CERRADA"]:
38                         new_columna.append(None)
39                     else:
40                         new_columna.append(valor)
41
42             elif n_tabla == "Mantenimiento" and n_columna == "TIPO_INTERVENCION":
43                 for valor in columna:
44                     if valor not in ["CORRECTIVO", "EMERGENCIA", "PREVENTIVO"]:
45                         new_columna.append(None)
46                     else:
47                         new_columna.append(valor)
48
49             elif n_tabla == "Mantenimiento" and n_columna in [

```

```

50         "ESTADO_PREVIO",
51         "ESTADO_POSTERIOR",
52     ]:
53         for valor in columna:
54             if valor not in ["MALO", "REGULAR", "BUENO"]:
55                 new_columnna.append(None)
56             else:
57                 new_columnna.append(valor)
58
59         elif n_tabla == "Incidentes" and n_columnna == "TIPO_INCIDENTE":
60             for valor in columna:
61                 if valor not in [
62                     "ROBO",
63                     "CAIDA",
64                     "VANDALISMO",
65                     "ACCIDENTE",
66                     "DAÑO ESTRUCTURAL",
67                 ]:
68                     new_columnna.append(None)
69                 else:
70                     new_columnna.append(valor)
71
72         elif n_tabla == "Incidentes" and n_columnna == "GRAVEDAD":
73             for valor in columna:
74                 if valor not in ["ALTA", "BAJA", "MEDIA", "CRITICA"]:
75                     new_columnna.append(None)
76                 else:
77                     new_columnna.append(valor)
78
79         if len(new_columnna) > 0:
80             tabla[n_columnna] = new_columnna
81             tabla.dropna(subset=n_columnna)
82

```

14. nif status: Verifica el campo nif. Código:

```

1  # Verifica el campo NIF
2  def nif_status(db: dict):
3      columna = db["Usuarios"]["NIF"]
4      new_columnna = []
5      for elemento in columna:
6          if "-" not in elemento:
7              new_columnna.append(None)
8          else:
9              lista = elemento.split("-")
10             if len(lista) != 3:
11                 new_columnna.append(None)
12             else:
13                 if (
14                     not lista[0].isdigit()
15                     or not lista[1].isdigit()
16                     or not lista[2].isdigit()
17                 ):
18                     new_columnna.append(None)
19                 else:
20                     new_columnna.append(elemento)
21     db["Usuarios"]["NIF"] = new_columnna
22     db["Usuarios"].dropna(subset=["NIF"])
23

```

15. **check id:** Verifica los id de cada tabla. Código:

```
1 # Verifica los id en cada tabla
2 def check_id(db: dict):
3     lista = [
4         "Areas",
5         "Juegos",
6         "Encuestas",
7         "Incidencias",
8         "Incidentes",
9         "Mantenimiento",
10    ]
11    for nombre in lista:
12        columna = db[nombre]["ID"]
13        new_columna = []
14        if nombre == "Mantenimiento":
15            n = 1
16            for valor in columna:
17                expected = f"-{n},00\\xa0MNT"
18                if expected != valor:
19                    new_columna.append(None)
20                else:
21                    valor = f"-{n} MNT"
22                    new_columna.append(valor)
23            n += 1
24        else:
25            for valor in columna:
26                if not isinstance(valor, int):
27                    new_columna.append(None)
28                else:
29                    new_columna.append(valor)
30        db[nombre]["ID"] = new_columna
31        db[nombre].dropna(subset=["ID"])
32
```

16. **incidencias status:** Verifica los ids y nif dentro de incidencias. Código:

```
1 # Verifica los ids y nif en incidencias
2 def incidencias_status(db: dict):
3     # Crear copias de los conjuntos de búsqueda para optimizar el acceso
4     ids = set(db["Mantenimiento"]["ID"])
5     nifs = set(db["Usuarios"]["NIF"])
6
7     # Crear una copia del DataFrame "Incidencias"
8     tabla = db["Incidencias"].copy()
9
10    # Procesar la columna "MantenimientoID"
11    mantenimiento_ids = []
12    for valor in tabla["MantenimientoID"]:
13        if valor.startswith("[") and valor.endswith("]"):
14            contenido = valor[1:-1]
15            partes = contenido.split(", ")
16            partes_filtradas = []
17            for parte in partes:
18                if (
19                    parte.startswith("'MNT-")
20                    and parte.endswith("'")
21                    and parte[5:-1].isdigit()
22                ):
23                    numero = str(
```

```

24         int(parte[5:-1])
25     ) # Convertir a entero y luego a cadena para quitar ceros iniciales
26     expected = f"-{numero} MNT"
27     if expected in ids:
28         partes_filtradas.append(expected)
29     if partes_filtradas:
30         mantenimiento_ids.append(f"[{', '.join(partes_filtradas)}]")
31     else:
32         mantenimiento_ids.append(None)
33 else:
34     mantenimiento_ids.append(None)
35
36 # Asignar los valores filtrados a la columna "MantenimientoID"
37 tabla["MantenimientoID"] = mantenimiento_ids
38
39 # Eliminar las filas con valores None en "MantenimientoID"
40 tabla = tabla.dropna(subset=["MantenimientoID"])
41
42 # Procesar la columna "UsuarioID"
43 usuario_ids = []
44 for valor in tabla["UsuarioID"]:
45     if valor.startswith("[") and valor.endswith("]"):
46         contenido = valor[1:-1]
47         partes = contenido.split(", ")
48         partes_filtradas = [parte for parte in partes if parte[1:-1] in nifs]
49         if partes_filtradas:
50             usuario_ids.append(f"[{', '.join(partes_filtradas)}]")
51         else:
52             usuario_ids.append(None)
53     else:
54         usuario_ids.append(None)
55
56 # Asignar los valores filtrados a la columna "UsuarioID"
57 tabla["UsuarioID"] = usuario_ids
58
59 # Eliminar las filas con valores None en "UsuarioID"
60 tabla = tabla.dropna(subset=["UsuarioID"])
61
62 # Asignar el DataFrame modificado de nuevo en el DataFrame original
63 db["Incidencias"] = tabla
64

```

El resto de limpieza se realizó directamente en los csv.

1.4. Proceso de implementación de columnas

En el archivo de python imputation.py se crearon las funciones para eliminar/añadir columnas de tablas. Para ejecutar este archivo, es necesario importar estas librerías de esta manera:

```

1 import pandas as pd
2 import random
3

```

Las funciones son las siguientes:

1. **check list:** Verifica si el *Punto Muestreo* pasado es correcto. Se tiene en cuenta también si hace referencia a vientos (acaba en 81), temperaturas (acaba en 83) o precipitaciones (acaba en 89). Si termina en cualquier otr cifrs se considerará erróneo. Es útil para la función de reestructuración de meteo, que se explicará después. Código:

```

1 def check_list(elemento: str, codes) -> list | None:
2     # verifica si el punto muestreo es correcto
3     if "_" not in elemento:
4         return None
5     data = elemento.split("_")
6     if len(data) != 3:
7         return None
8     if not data[0].isdigit() or not data[1].isdigit():
9         return None
10    #Sólo se sacan vientos (81), temperaturas (83) y precipitaciones (89)
11    if data[0] not in codes.keys() or data[1] not in ["81", "83", "89"]:
12        return None
13    return data
14

```

2. **same line:** También útil para la restructuración de la tabla meteo. Se encarga de identificar una medición según la *FECHA* y el *DISTRITO*. En caso de no existir dicha medición (no existe una linea con ambos valores) devuelve -1. Código:

```

1 def same_line(fecha: str, distrito: str, list_f: list, list_d: list) -> int:
2     #Identifica una linea donde fecha sea su FECHA y distrito su DISTRITO.
3     if fecha not in list_f or distrito not in list_d:
4         return -1
5     for n in range(0, len(list_f)):
6         if list_f[n] == fecha and list_d[n] == distrito:
7             return n
8     return -1 #Si no hay, devuelve -1
9

```

3. **post code:** A partir de un código postal pasado como argumento se verifica que existe una estación meteorológica (tabla *Codigo*) con un código postal igual. Si existe se devuelve la posición dentro de la tabla *Codigo*. Código:

```

1 def post_code(db: pd.DataFrame, codigo: int):
2     #Identifica si el código postal pasado es igual al de una estación meteorológica.
3     code_id = db["Codigo"]["CÓDIGO"]
4     for n in range(len(code_id)):
5         if str(code_id[n]) == str(codigo):
6             return n #Posición dentro de Código
7     return -1 #No existe
8

```

4. **new meteo:** Se modifica totalmente la tabla *meteo24*. Se crean las siguientes columnas:

- ID: Se crea según la posición de la fila que le toque.
- CODIGO POSTAL: Código postal del lugar al que pertenece. Todos tienen 1 excepto Plaza de España que tiene 2.
- FECHA: Fecha tipo YYYY-MM-DD hh:mm:ss.
- VIENTO: Booleano que indica si dicho día hay vientos fuertes. Nos hemos basado en la velocidad (≥ 20).
- TEMPERATURA: Indica si hay medición de temperatura y cuál.
- PRECIPITACION: Indica si hay medición de precipitación y cuál.
- DISTRITO: Indica el distrito al que pertenece la estación meteorológica.

La función trata de verificar el código postal, identificar el distrito y la fecha, insertar si no existe dicha combinación y luego insertar los datos faltantes (viento si PUNTO MUESTREO termina en 81, temperatura si termina en 83 y precipitación si termina en 89). En caso de Plaza de España se realiza una doble inserción (hay 2 códigos postales -> 2 filas). Esta función se planteó sin IA pero se corrigió con ChatGPT. Código:

```

1 def new_meteo(db: pd.DataFrame):
2     # Modifica totalmente la tabla de meteo24 para que muestre lo que pide el enunciado
3     meteo = db["meteo24"]
4     m_postal = db["Codigo"] # Lista de códigos postales para la relación área-meteo
5     col_postal = m_postal["CodigoPostal"]
6     codes = {
7         "28079102": "MORATALAZ",
8         "28079103": "VILLVERDE",
9         "28079104": "PUENTE DE VALLECAS",
10        "28079106": "MONCLOA-ARAVACA",
11        "28079107": "HORTALEZA",
12        "28079108": "FUENCARRAL-EL PARDO",
13        "28079109": "CHAMBERI",
14        "28079110": "CENTRO",
15        "28079111": "CHAMARTIN",
16        "28079112": "VILLA DE VALLECAS",
17        "28079113": "VILLA DE VALLECAS",
18        "28079114": "ARGANZUELA",
19        "28079115": "ARGANZUELA",
20        "28079004": "MONCLOA-ARAVACA",
21        "28079008": "SALAMANCA",
22        "28079016": "CIUDAD LINEAL",
23        "28079018": "CARABANCHEL",
24        "28079024": "MONCLOA-ARAVACA",
25        "28079035": "CENTRO",
26        "28079036": "MORATALAZ",
27        "28079038": "TETUAN",
28        "28079039": "FUENCARRAL-EL PARDO",
29        "28079054": "VILLA DE VALLECAS",
30        "28079056": "CARABANCHEL",
31        "28079058": "FUENCARRAL-EL PARDO",
32        "28079059": "BARAJAS",
33    } # Diccionario de distritos
34    meses = {
35        "1": 31,
36        "2": 28,
37        "3": 31,
38        "4": 30,
39        "5": 31,
40        "6": 30,
41        "7": 31,
42        "8": 31,
43        "9": 30,
44        "10": 31,
45        "11": 30,
46        "12": 31,
47    } # Diccionario de meses (como la fecha es por día, se necesita saber qué días hay que
48    coger)
49    ids = [] # Cada id representa la posición de la fila de forma ordenada
50    codigo_postal = []
51    distritos = []
52    temperaturas = []
53    precipitaciones = []
54    viento = []
55    fechas = []

```

```

55     id = 1
56     for i in range(len(meteo["PUNTO_MUESTREO"])):
57         elemento = meteo["PUNTO_MUESTREO"][i]
58         data = check_list(elemento, codes) # Verifica el PUNTO MUESTREO y el MES
59         if data is not None and str(meteo["MES"][i]) in meses.keys():
60             for j in range(
61                 meteo["MES"][i]
62             ): # Por cada fecha y distrito, una fila distinta
63                 dia = str(j + 1)
64                 mes = str(meteo["MES"][i])
65                 if j < 9:
66                     dia = "0" + dia
67                 if int(mes) < 10:
68                     mes = "0" + mes
69                 fecha = str(meteo["ANO"][i]) + "-" + mes + "-" + dia
70                 n_code = post_code(db, data[0]) # Verifica el código postal
71                 if n_code != -1:
72                     # Si no existe una fila con dicha fecha y distrito, se crea una nueva
73                     linea = same_line(fecha, codes[data[0]], fechas, distritos)
74                     if linea == -1:
75                         ids.append(id)
76                         id += 1
77                         pcode = col_postal[n_code]
78                         codigo_postal.append(int(pcode))
79                         distritos.append(codes[data[0]])
80                         fechas.append(fecha)
81                         temperaturas.append("-")
82                         precipitaciones.append("-")
83                         viento.append(False)
84                         linea = len(distritos) - 1
85                         if data[0] == "28079004":
86                             # Como Plaza de España tiene dos códigos postales, hay doble
87                             ids.append(id)
88                             id += 1
89                             distritos.append(codes[data[0]])
90                             pcode = col_postal[n_code + 1]
91                             codigo_postal.append(int(pcode))
92                             fechas.append(fecha)
93                             temperaturas.append("-")
94                             precipitaciones.append("-")
95                             viento.append(False)
96                             # Aquí se tiene en cuenta la línea anterior
97                             booleano = (
98                                 "V" + dia
99                             ) # Para modificar el valor, el booleano debe ser igual a V
100                             valor = "D" + dia # Datos a insertar
101                             if data[1] == "81":
102                                 if meteo[booleano][i] == "V" and int(meteo[valor][i]) >= 20:
103                                     viento[linea] = True
104                                     if (
105                                         data[0] == "28079004"
106                                     ): # Plaza de España -> Doble inserción
107                                         viento[linea + 1] = True
108                                 else:
109                                     viento[linea] = False
110                                     if (
111                                         data[0] == "28079004"
112                                     ): # Plaza de España -> Doble inserción
113                                         viento[linea + 1] = False

```



```

114         elif data[1] == "83":
115             if meteo[booleano][i] == "V":
116                 temperaturas[linea] = meteo[valor][i]
117                 if (
118                     data[0] == "28079004"
119                 ): # Plaza de España -> Doble inserción
120                     temperaturas[linea + 1] = meteo[valor][i]
121             else:
122                 if meteo[booleano][i] == "V":
123                     precipitaciones[linea] = meteo[valor][i]
124                     if (
125                         data[0] == "28079004"
126                     ): # Plaza de España -> Doble inserción
127                         precipitaciones[linea + 1] = meteo[valor][i]
128
129 nuevo = pd.DataFrame() # Se crea la nueva tabla
130 nuevo.loc[:, "ID"] = ids
131 nuevo.loc[:, "CODIGO_POSTAL"] = codigo_postal
132 nuevo.loc[:, "FECHA"] = fechas
133 pd.to_datetime(nuevo["FECHA"]).apply(lambda x: x.date())
134 nuevo.loc[:, "TEMPERATURA"] = temperaturas
135 nuevo.loc[:, "PRECIPITACION"] = precipitaciones
136 nuevo.loc[:, "VIENTO"] = viento
137 nuevo.loc[:, "DISTRITO"] = distritos
138 db["meteo24"] = nuevo
139

```

5. **are new attribute**: Se añade una columna *AreaID* a la tabla *Juegos* para representar la relación áreas-juegos. Primero se realiza comparando ambas coordenadas y para el resto se le asigna un id random no repetido. Código:

```

1 def area_new_attribute(df: pd.DataFrame):
2     """Añade el atributo capacidadMax a Areas."""
3     df["Juegos"]["AreaRecreativaID"] = 0
4     index = 0
5     # Calcula el número de juegos por area con lat y long
6     for area in df["Areas"].to_numpy():
7         games = 0 # Número de juegos
8         # Buscar por GPS
9         lat_area = area[10]
10        long_area = area[11]
11        index_juego = 0
12        for juego in df["Juegos"].to_numpy():
13            lat_juego = juego[10]
14            long_juego = juego[11]
15            # Si coinciden, se suma el número de juegos y se añade el areaID a Juegos
16            if lat_area == lat_juego and long_area == long_juego:
17                games += 1
18                df["Juegos"].loc[index_juego, "AreaRecreativaID"] = area[0]
19                index_juego += 1
20            # Aquí crea la columna
21            df["Areas"].loc[index, "capacidadMax"] = games
22            index += 1
23        # Añadir valor a los juegos sin áreas
24        index_juego = 0
25        for juego in df["Juegos"].to_numpy():
26            ref_area = juego[24]
27            # Si no hay referencia del juego en área
28            if ref_area == 0:

```

```

29         in_column = True
30         while in_column:
31             # Inserta un valor random no repetido
32             rand_id = random.randint(1000000, 100000000)
33             if rand_id not in df["Juegos"]["AreaRecreativaID"]:
34                 df["Juegos"].loc[index_juego, "AreaRecreativaID"] = rand_id
35                 in_column = False
36         index_juego += 1
37

```

6. **juegos new attributes:** Añade a la tabla *Juegos* las columnas *indicadorExposicion* y *desgasteAcumulado*. El indicador de exposición puede ser 100/Bajo, 200/Medio y 300/Alto y se escoge de manera aleatoria. El desgaste acumulado es un valor que se calcula del número de años de uso (entero aleatorio entre 1 y 15), el indicador de exposición y el número de juegos mantenidos. Código:

```

1 def juegos_new_attributes(df: pd.DataFrame):
2     """Añade las columnas indicadorExposicion y desgasteAcumulado a Juegos"""
3     # Inserta el valor de indicador exposicion
4     indicador_options = {"BAJO": 100, "MEDIO": 200, "ALTO": 300}
5     for i in range(len(df["Juegos"])):
6         selected_option = random.randint(0, 2)
7         df["Juegos"].loc[i, "indicadorExposicion"] = list(indicador_options.keys())[
8             selected_option
9         ]
10    # Calcula el valor de desgasteAcumulado
11    index = 0
12    wear_values = []
13    for juego in df["Juegos"].to_numpy():
14        # Número del mantenimiento al juego
15        num_mant_juego = 0
16        # Uso aleatorio puede ser un entero del 1 al 15
17        use_time = random.randint(1, 15)
18        for maintenance in df["Mantenimiento"].to_numpy():
19            if juego[0] == maintenance[5]:
20                num_mant_juego += 1
21        # Añadir el wear value a una lista para cambiar su rango
22        wear_value = (
23            use_time * indicador_options[df["Juegos"]["indicadorExposicion"][index]]
24        ) - (num_mant_juego * 100)
25        wear_values.append(wear_value)
26        index += 1
27    # Inserta desgasteAcumulado
28    for i in range(len(df["Juegos"])):
29        df["Juegos"].loc[i, "desgasteAcumulado"] = adjust_range(
30            wear_values, 0, 100, i
31        ) # Cambia el rango de sus valores
32

```

7. **adjust range:** Ajusta un valor de una lista a unos nuevos límites. Código:

```

1 def adjust_range(values: list, new_min: int, new_max: int, index: int) -> int:
2     # Fórmula para ajustar un valor dentro de una lista a unos nuevos límites
3     old_min = min(values)
4     old_max = max(values)
5     new_value = int(
6         (((values[index] - old_min) * (new_max - new_min)) / (old_max - old_min))
7         + new_min
8     )

```

```

9     return new_value
10

```

8. **takeTimebyID**: Coge una fecha de la tabla *Mantenimiento* según un id. Código:

```

1 def takeTimebyID(base: pd.DataFrame, id: str):
2     # Coge una fecha de mantenimiento desde su id
3     numero = int(id[1:-4])
4     lista = base["Mantenimiento"]["ID"]
5     if numero < len(lista):
6         return base["Mantenimiento"]["FECHA_INTERVENCION"][numero]
7     return None
8

```

9. **tiempoResolucion**: Se calcula la diferencia de tiempo (en días) entre la fecha de mantenimiento y la de la incidencia. Para realizar dicha resta, ambos mantenimiento e incidencia deben tener el mismo id de mantenimiento, además de que exista la fecha de mantenimiento y que sea mayor a la de la incidencia. Cabe aclarar que, como hay varios id de mantenimiento para una incidencia, se escogerá el tiempo mayor. Toda esta información se almacenará en la nueva columna *tiempoResolucion* en la tabla *Incidencias*. Esta función se hizo con ayuda de chatGPT. Código:

```

1 def tiempoResolucion(base: pd.DataFrame):
2     # Calcula el tiempo entre la incidencia y el mantenimiento
3     col_index = base["Incidencias"]["MantenimientoID"]
4     # Importante: La incidencia debe contener el id de dicho mantenimiento
5     tiempos = []
6     for n in range(len(col_index)):
7         elemento = col_index[n]
8         conjunto = elemento[1:-1]
9         conjunto = conjunto.split(", ")
10        tiempo = float(0) # Tiempo predeterminado
11        fecha_incidencia = base["Incidencias"]["FECHA_REPORTE"][n]
12        for valor in conjunto:
13            fecha_mat = takeTimebyID(base, valor)
14            # Al tener ambas fechas, se hace la resta si la de mantenimiento es más tardía que
15            # la de la incidencia
16            if fecha_mat != None and fecha_mat > fecha_incidencia:
17                sub_t = fecha_mat - fecha_incidencia
18
19                # Convertir sub_t a segundos para poder compararlo con el valor float 'tiempo'
20                sub_t_seconds = sub_t.total_seconds()
21                if sub_t_seconds > tiempo:
22                    tiempo = sub_t_seconds
23
24            if tiempo != 0:
25                tiempo = tiempo // (60 * 60 * 24) # Pasar a días
26        tiempos.append(tiempo)
27    base["Incidencias"].loc[:, "TIEMPO_RESOLUCION"] = tiempos
28

```

10. **lastFecha**: Identifica las últimas fechas de mantenimiento de cada juego y las almacena en la nueva columna *ULTIMA FECHA MANTENIMIENTO* de la tabla *Juegos*. Código:

```

1 def lastFecha(db: pd.DataFrame):
2     # Seleccionamos el último mantenimiento para cada juego en la tabla de mantenimiento
3     max_mant = (
4         db["Mantenimiento"]
5         .groupby("JuegoID")["FECHA_INTERVENCION"]

```

```

6         .max()
7         .reset_index()
8         .rename(columns={"FECHA_INTERVENCION": "ULTIMA_FECHA_MANTENIMIENTO"})
9     )
10
11     # Fusionamos las fechas de mantenimiento con la tabla de juegos
12     db["Juegos"] = db["Juegos"].merge(
13         max_mant, left_on="ID", right_on="JuegoID", how="left"
14     )
15
16     # Rellenamos con FECHA_INSTALACION en caso de que no haya mantenimientos
17     db["Juegos"]["ULTIMA_FECHA_MANTENIMIENTO"].fillna(
18         db["Juegos"]["FECHA_INSTALACION"], inplace=True
19     )
20

```

11. **area meteo:** Implementa la columna *MeteoID* a area para identificar a que estación meteorológica pertenece a cada área. Código:

```

1 def area_meteo(db: pd.DataFrame):
2     # Inserta en Áreas el id de la meteo24 que pertenece
3     lista_meteo = [] # Columna de los id de meteo
4     tabla_area = db["Areas"]
5     tabla_meteo = db["meteo24"]
6     columna_cod_postal = tabla_area["COD_POSTAL"]
7     columna_cod_meteo = tabla_meteo["CODIGO_POSTAL"]
8
9     # Crear un diccionario para el mapeo de códigos postales a ID
10    codigo_to_id = {
11        columna_cod_meteo[i]: tabla_meteo["ID"][i]
12        for i in range(len(columna_cod_meteo))
13    }
14
15    for valor in columna_cod_postal:
16        estacion = 0 # Valor predeterminado
17
18        # Asegurarse de que ambos valores son enteros
19        if pd.notnull(valor) and type(valor) is not str: # Evita valores nulos
20            valor = int(valor) if not isinstance(valor, int) else valor
21            estacion = codigo_to_id.get(
22                valor, estacion
23            ) # Obtener el valor o usar predeterminado
24
25        lista_meteo.append(estacion)
26
27    # Asignar la lista resultante a una nueva columna "MeteoID" en db["Areas"]
28    db["Areas"]["MeteoID"] = lista_meteo
29

```

12. **nivelEscalamiento:** Crea la columna *NIVEL RECONOCIMIENTO* en *Incidencias* y se asigna un valor según si la incidencia está abierta (aleatorio del 1 al 10) o cerrada (1). Código:

```

1 def nivelEscalamiento(db: pd.DataFrame):
2     # Inserta el nivel de escalamiento. Si es abierto -> aleatorio del 1 al 10. Si no, es 1.
3     incidencias = db["Incidencias"]
4     lista_rec = []
5     for n in range(len(incidencias["ESTADO"])):
6         if incidencias["ESTADO"][n] == "ABIERTA":
7             numero = random.randint(1, 10)

```

```

8         lista_rec.append(numero)
9     else:
10         lista_rec.append(1)
11     incidencias.loc[:, "NIVEL_RECONOCIMIENTO"] = lista_rec
12

```

1.5. Proceso de extracción a mongoDB

Con el archivo "formatear mongo.py" se extraen los objetos y se crean las tablas limpias para mongoDB. Para ejecutar este archivo, es necesario importar estas librerías de esta manera:

```

1 import json
2

```

Todas estas funciones se encuentran en la clase *Creator*, que destaca por este init:

```

1 class Creator:
2     def __init__(self, bd: dict):
3
4         # Base de datos en el estado requerido
5         self.state = bd
6
7         # Correspondencias entre columna a añadir y tablas que la requieren
8         self.col_table = {
9             "AreaRecreativaID": "meteo24",
10            "JueID": "Incidencias",
11        }
12
13        # relaciones entre sí
14        self.juego_tipo = {}
15        self.juego_incidencias = {}
16        self.area_meteo = {}
17        self.area_encuesta = {}
18        self.area_incidente = {}
19        self.area_juegos = {}
20        self.juego_mantenimientos = {}
21        self.mantenimiento_incidencias = {}
22        self.incidencia_usuario = {}
23
24        # objetos extraídos
25        self.areas = []
26        self.encuestas = []
27        self.registrosClima = []
28        self.incidentesSeguridad = []
29        self.juegos = []
30        self.incidencias = []
31        self.usuarios = []
32        self.mantenimientos = []
33
34

```

Además, contiene las siguientes funciones de generación de objetos:

1. **get json data:** Genera los JSON de los objetos extraídos. Además, toda fecha que encuentre le cambia al formato ISO 8601. Código:

```

1     def get_json_data(self, objs: list, collection_name: str):
2
3         # Dar formato ISO 8601 a fechas

```

```

4         for obj in objs:
5             for key in obj:
6                 if "fecha" in key:
7                     obj[key] = (
8                         f"{str(obj[key]).split(" ")[0]}T{str(obj[key]).split(" ")[1]}Z"
9                     )
10
11         # abrir fichero y escribir datos
12         with open(f"DatasetsNuevos/{collection_name}.js", "w") as file:
13             file.write(
14                 f"""db.{collection_name}.insertMany(
15                     { json.dumps(objs,indent=4) }
16                 )"""
17             )
18         print(f"JSON creado para {collection_name}")
19

```

2. **generar usuarios:** Extrae los objetos de la tabla *Usuarios*. Código:

```

1     def generar_usuarios(self):
2
3         # extraer columnas de datos usuarios
4         nifs = self.extraer_columna("Usuarios", "NIF")
5         nombre = self.extraer_columna("Usuarios", "NOMBRE")
6         email = self.extraer_columna("Usuarios", "EMAIL")
7         tlf = self.extraer_columna("Usuarios", "TELEFONO")
8
9         for i in range(len(nifs)):
10             self.usuarios.append(
11                 {
12                     "NIF": str(nifs[i]),
13                     "nombre": str(nombre[i]),
14                     "email": str(email[i]),
15                     "telefono": str(tlf[i]),
16                 }
17             )
18

```

3. **generar encuestas:** Extrae los objetos de la tabla *Encuestas*. Código:

```

1     def generar_encuestas(self):
2
3         id = self.extraer_columna("Encuestas", "ID")
4         fecha = self.extraer_columna("Encuestas", "FECHA")
5         accesibilidad = self.extraer_columna("Encuestas", "PUNTUACION_ACCESIBILIDAD")
6         calidad = self.extraer_columna("Encuestas", "PUNTUACION_CALIDAD")
7         comentarios = self.extraer_columna("Encuestas", "COMENTARIOS")
8
9         for i in range(len(id)):
10             self.encuestas.append(
11                 {
12                     "id": int(id[i]),
13                     "fechaEncuesta": str(fecha[i]),
14                     "puntuacionAccesibilidad": accesibilidad[i],
15                     "puntuacionCalidad": calidad[i],
16                     "comentarios": comentarios[i],
17                 }
18             )
19

```

4. **generar incidentes:** Extrae los objetos de la tabla *Incidentes*. Código:

```
1 def generar_incidentes(self):
2
3     id = self.extraer_columnna("Incidentes", "ID")
4     fecha = self.extraer_columnna("Incidentes", "FECHA_REPORTE")
5     tipo = self.extraer_columnna("Incidentes", "TIPO_INCIDENTE")
6     gravedad = self.extraer_columnna("Incidentes", "GRAVEDAD")
7
8     for i in range(len(id)):
9         self.incidentesSeguridad.append(
10             {
11                 "id": int(id[i]),
12                 "fechaDeReporte": str(fecha[i]),
13                 "tipoIncidente": tipo[i],
14                 "gravedad": gravedad[i],
15             }
16         )
17
```

5. **generar incidencias:** Extrae los objetos de la tabla *Incidencias*. Código:

```
1 def generar_incidencias(self):
2
3     id = self.extraer_columnna("Incidencias", "ID")
4     tipo = self.extraer_columnna("Incidencias", "TIPO_INCIDENCIA")
5     fecha = self.extraer_columnna("Incidencias", "FECHA_REPORTE")
6     estado = self.extraer_columnna("Incidencias", "ESTADO")
7     tiempo = self.extraer_columnna("Incidencias", "TIEMPO_RESOLUCION")
8     escala = self.extraer_columnna("Incidencias", "NIVEL_RECONOCIMIENTO")
9
10    for i in range(len(id)):
11        if id[i] in self.incidencia_usuario:
12
13            usuarios = self.incidencia_usuario[int(id[i])]
14
15            # encontrar usuarios en la relación incidencia-usuario
16            res_users = []
17            for user in self.usuarios:
18                for usuarioID in usuarios:
19                    if user["NIF"] == usuarioID:
20                        res_users.append(user)
21
22            self.incidencias.append(
23                {
24                    "id": int(id[i]),
25                    "fechaReporte": str(fecha[i]),
26                    "tipo": tipo[i],
27                    "estado": estado[i],
28                    "tiempoResolucion": float(tiempo[i]),
29                    "nivelEscalamiento": int(escala[i]),
30                    "usuarios": res_users,
31                }
32            )
33
```

6. **generar juegos:** Extrae los objetos de la tabla *Juegos*. Código:

```
1 def generar_juegos(self):
```

```

2     id = self.extraer_columnna("Juegos", "ID")
3     modelo = self.extraer_columnna("Juegos", "MODELO")
4     estado_op = self.extraer_columnna("Juegos", "ESTADO")
5     accesibilidad = self.extraer_columnna("Juegos", "ACCESIBLE")
6     fecha_instalacion = self.extraer_columnna("Juegos", "FECHA_INSTALACION")
7     tipo = self.extraer_columnna("Juegos", "tipo_juego")
8     desgaste = self.extraer_columnna("Juegos", "desgasteAcumulado")
9     indicador_exposicion = self.extraer_columnna("Juegos", "indicadorExposicion")
10    ultima_fecha_mant = self.extraer_columnna("Juegos", "ULTIMA_FECHA_MANTENIMIENTO")
11
12    for i in range(len(id)):
13        # si el juego no tiene incidencias
14        if id[i] not in list(self.juego_incidencias.keys()):
15            incidencia = []
16        else:
17            incidencia = self.juego_incidencias[id[i]]
18        # si el juego no tiene mantenimientos
19        if id[i] not in list(self.juego_mantenimientos.keys()):
20            mantenimiento = []
21        else:
22            mantenimiento = self.juego_mantenimientos[id[i]]
23        self.juegos.append(
24            {
25                "id": str(id[i]),
26                "modelo": str(modelo[i]),
27                "estadoOperativo": str(estado_op[i]),
28                "accesibilidad": bool(accesibilidad[i]),
29                "fechaInstalacion": str(fecha_instalacion[i]),
30                "tipo": str(tipo[i]),
31                "desgasteAcumulado": int(desgaste[i]),
32                "indicadorExposicion": str(indicador_exposicion[i]),
33                "ultimaFechaMantenimiento": str(ultima_fecha_mant[i]),
34                "mantenimientos": mantenimiento,
35                "incidencias": incidencia,
36            }
37        )
38

```

7. **generar clima:** Extrae los objetos de la tabla *meteo24*. Código:

```

1    def generar_clima(self):
2        id = self.extraer_columnna("meteo24", "ID")
3        fecha = self.extraer_columnna("meteo24", "FECHA")
4        temp = self.extraer_columnna("meteo24", "TEMPERATURA")
5        vent = self.extraer_columnna("meteo24", "VIENTO")
6        prec = self.extraer_columnna("meteo24", "PRECIPITACION")
7        for i in range(len(id)):
8            # Solo pasa a entero los valores numéricos, no los -
9            if temp[i] != "-":
10                temp[i] = int(temp[i])
11            else:
12                temp[i] = 0
13            if prec[i] != "-":
14                prec[i] = int(prec[i])
15            else:
16                prec[i] = 0
17
18            self.registrosClima.append(
19                {
20                    "id": int(id[i]),

```



```

21         "fecha": fecha[i],
22         "temperatura": temp[i],
23         "vientosFuertes": bool(vent[i]),
24         "precipitacion": prec[i],
25     }
26 )
27

```

8. **generar_area:** Extrae los objetos de la tabla *Areas*. Sólo se tienen en cuenta las columnas pedidas en la práctica. Código:

```

1  def generar_area(self):
2      id = self.extraer_columna("Areas", "ID")
3      barrio = self.extraer_columna("Areas", "BARRIO")
4      distr = self.extraer_columna("Areas", "DISTRITO")
5      estado = self.extraer_columna("Areas", "ESTADO")
6      lat = self.extraer_columna("Areas", "LATITUD")
7      long = self.extraer_columna("Areas", "LONGITUD")
8      fecha = self.extraer_columna("Areas", "FECHA_INSTALACION")
9      capmax = self.extraer_columna("Areas", "capacidadMax")
10
11     for i in range(len(id)):
12         # Solo pasa a entero los valores numéricos, no los -
13
14         # extraer incidentes de seguridad de la relación
15         res_incidentes = []
16         if int(id[i]) in self.area_incidente:
17             incidentes = self.area_incidente[int(id[i])]
18             for inci in self.incidentesSeguridad:
19                 for incidentesID in incidentes:
20                     if inci["id"] == incidentesID:
21                         res_incidentes.append(int(incidentesID))
22
23         # extraer encuestas de la relación
24         res_encuestas = []
25         if int(id[i]) in self.area_encuesta:
26             encuestas = self.area_encuesta[int(id[i])]
27             for enc in self.encuestas:
28                 for encuestaID in encuestas:
29                     if enc["id"] == encuestaID:
30                         res_encuestas.append(int(encuestaID))
31
32         # extraer climas de la relación
33         res_climas = []
34         if int(id[i]) in self.area_meteo:
35             climas = self.area_meteo[int(id[i])]
36             for clima in self.registrosClima:
37                 for climaID in climas:
38                     if clima["id"] == climaID:
39                         res_climas.append(int(climaID))
40
41         # extraer juegos de la relación
42         res_juegos = []
43         if int(id[i]) in self.area_juegos:
44             juegos = self.area_juegos[int(id[i])]
45             for juegoID in juegos:
46                 res_juegos.append(int(juegoID))
47
48         # extraer atributo juego_tipo
49         res_juegos_tipo_valor = []

```

```

50         for tipo in self.juego_tipo:
51             objeto = {"tipo": str(tipo), "valor": int(self.juego_tipo[tipo])}
52             res_juegos_tipo_valor.append(objeto)
53
54         self.areas.append(
55             {
56                 "id": int(id[i]),
57                 "barrio": barrio[i],
58                 "distrito": distr[i],
59                 "estadoOperativo": estado[i],
60                 "coordenadasGPS": [lat[i], long[i]],
61                 "fecha": fecha[i],
62                 "capacidadMaxima": float(capmax[i]),
63                 "incidentesSeguridad": res_incidentes,
64                 "encuestas": res_encuestas,
65                 "registrosClima": res_climas,
66                 "juegos": res_juegos,
67                 "cantidadJuegosTipo": res_juegos_tipo_valor,
68             }
69         )
70
71

```

9. **generar mantenimientos:** Función para extraer los objetos de la tabla *Mantenimiento*. Código:

```

1  def generar_mantenimientos(self):
2
3      id = self.extraer_columna("Mantenimiento", "ID")
4      tipo = self.extraer_columna("Mantenimiento", "TIPO_INTERVENCION")
5      estado_previo = self.extraer_columna("Mantenimiento", "ESTADO_PREVIO")
6      estado_posterior = self.extraer_columna("Mantenimiento", "ESTADO_POSTERIOR")
7      fecha = self.extraer_columna("Mantenimiento", "FECHA_INTERVENCION")
8
9      transform = lambda x: f"{x.split(" ")[1]}-{int(x.split(" ")[0][1:])}"
10     for i in range(len(id)):
11
12         if transform(id[i]) in self.mantenimiento_incidentes:
13             incidencias = self.mantenimiento_incidentes[transform(id[i])]
14
15         # encontrar incidencias en la relación mantenimiento-incidencias
16         res_incidencias = []
17         for inci in self.incidentes:
18             for incidenciasID in incidencias:
19                 if inci["id"] == incidenciasID:
20                     res_incidencias.append(int(incidenciasID))
21
22         self.mantenimientos.append(
23             {
24                 "id": transform(id[i]),
25                 "tipoIntervencion": str(tipo[i]),
26                 "estadoPrevio": estado_previo[i],
27                 "estadoPosterior": estado_posterior[i],
28                 "fechaIntervencion": fecha[i],
29                 "incidencias": res_incidencias,
30             }
31         )
32

```

10. **extraer columna:** Función para hacer las listas con cada columna de una tabla indicada. Código:

```

1  def extraer_columna(self, table: str, column: str) -> list:
2      # Extrae las tablas definitivas completas, iterando por sus columnas
3      result = []
4      for i in range(len(self.state[table][column])):
5          result.append(self.state[table].loc[i, column])
6      return result
7

```

11. **crear area clima:** Crea la relación entre las tablas *Areas-meteo24* extrayendo los id de las estaciones meteorológicas de *Areas* y analizando si existen en *meteo24*. Código:

```

1  def crear_area_clima(self):
2      tabla_areas = self.state["Areas"]
3
4      # Extraer las estaciones meteorologicas de la tabla de areas
5      for i in range(len(tabla_areas["MeteoID"])):
6          valor_meteo = int(tabla_areas.loc[i, "MeteoID"])
7          valor_id = int(tabla_areas.loc[i, "ID"])
8
9          if valor_meteo == 0:
10             continue
11
12         if valor_meteo not in self.area_meteo:
13             self.area_meteo[valor_id] = [valor_meteo]
14         else:
15             self.area_meteo[valor_id].append(valor_meteo)
16

```

12. **crear area juegos:** Crea la relación entre las tablas *Areas-Juegos* extrayendo los id de las áreas recreativas de *Juegos* y analizando si existen en *Areas*. Código:

```

1  def crear_area_juegos(self):
2
3      tabla_juegos = self.state["Juegos"]
4
5      # Extraer las areas de la tabla de juegos
6      for i in range(len(tabla_juegos["AreaRecreativaID"])):
7          valor_area = tabla_juegos.loc[i, "AreaRecreativaID"]
8          valor_juego = tabla_juegos.loc[i, "ID"]
9
10         if valor_area not in self.area_juegos:
11             self.area_juegos[int(valor_area)] = [int(valor_juego)]
12         else:
13             self.area_juegos[int(valor_area)].append(int(valor_juego))
14

```

13. **crear juegos tipo:** Identifica los distintos tipos de juegos que existen, contando el número de veces que aparecen en la tabla *Juegos*. Código:

```

1  def crear_juegos_tipo(self):
2
3      tabla_juegos = self.state["Juegos"]
4
5      # Extraer los tipos de la tabla de juegos
6      for i in range(len(tabla_juegos["tipo_juego"])):
7          valor_tipo = tabla_juegos.loc[i, "tipo_juego"]
8          if valor_tipo not in self.juego_tipo:

```

```

9         self.juego_tipo[valor_tipo] = 1
10     else:
11         self.juego_tipo[valor_tipo] += 1
12

```

14. **crear area encuestas:** Crea la relación entre las tablas *Areas-Encuestas* extrayendo los id de las áreas recreativas de *Encuestas* y analizando si existen en *Areas*. Código:

```

1  def crear_area_encuestas(self):
2      tabla_encuestas = self.state["Encuestas"]
3
4      # Extraer los mantenimientos de las incidencias. Mapa "AreaRecreativaID -> EncuestasID"
5      "
6      for i in range(len(tabla_encuestas["AreaRecreativaID"])):
7
8          valor_area = tabla_encuestas.loc[i, "AreaRecreativaID"]
9          valor_id = tabla_encuestas.loc[i, "ID"]
10
11         if valor_area not in self.area_encuesta:
12             self.area_encuesta[valor_area] = [valor_id]
13         else:
14             self.area_encuesta[valor_area].append(valor_id)

```

15. **crear juego incidencias:** Crea la relación entre las tablas *Juegos-Incidencias*. Esto se realiza con la ayuda de la tabla *Mantenimiento* ya que primero se extraen los id de mantenimiento de la tabla *Incidencias* y luego se extraen los id de los juegos de la tabla *Mantenimiento*. La relación se saca analizando qué ids de los sacados realmente existen en sus tablas representativas (id mantenimiento -> *Mantenimiento*, id juego -> *Juegos*. Código:

```

1  def crear_juego_incidencias(self):
2
3      # Tablas a utilizar
4      tabla_incidencias = self.state["Incidencias"]
5      tabla_mantenimiento = self.state["Mantenimiento"]
6
7      # Operación de transformación de nombre de atributos.
8      transform = lambda x: f"{x.split(" ")[1]}-{int(x.split(" ")[0][1:])}"
9
10     # Extraer los mantenimientos de las incidencias. Mapa "MantenimientoID -> IncidenciaID"
11     "
12     for i in range(len(tabla_incidencias["MantenimientoID"])):
13
14         id_manten_str = tabla_incidencias.loc[i, "MantenimientoID"]
15         id_manten_list = id_manten_str[1:-1].split(", ")
16         valor_id = tabla_incidencias.loc[i, "ID"]
17
18         for manten_id in id_manten_list:
19             if manten_id not in self.mantenimiento_incidencias:
20                 self.mantenimiento_incidencias[transform(manten_id)] = [valor_id]
21             else:
22                 self.mantenimiento_incidencias[transform(manten_id)].append(
23                     valor_id
24                 )
25
26     # Extraer los mantenimientos de los juegos. Mapa "JuegoID -> MantenimientoID"
27     for i in range(len(tabla_mantenimiento["JuegoID"])):

```

```

28     valor_id = tabla_mantenimiento.loc[i, "JuegoID"]
29     valor_manten = tabla_mantenimiento.loc[i, "ID"]
30
31     if valor_id not in self.juego_mantenimientos:
32         self.juego_mantenimientos[valor_id] = [transform(valor_manten)]
33     else:
34         self.juego_mantenimientos[valor_id].append(transform(valor_manten))
35
36     # Deducir las relaciones de juegos e incidencias. Mapa "JuegoID -> IncidenciasID"
37     for juego in self.juego_mantenimientos.keys():
38         for incidencia in self.mantenimiento_incidencias.values():
39             for mantenimiento in self.juego_mantenimientos[juego]:
40
41                 if mantenimiento not in self.mantenimiento_incidencias:
42                     continue
43                 if self.mantenimiento_incidencias[mantenimiento] == incidencia:
44                     for inc in incidencia:
45                         if juego not in self.juego_incidencias:
46                             self.juego_incidencias[juego] = [int(inc)]
47                         else:
48                             if int(inc) not in self.juego_incidencias[juego]:
49                                 self.juego_incidencias[juego].append(int(inc))
50

```

16. **crear area incidentes:** Crea la relación entre las tablas *Areas-Incidentes* extrayendo los id de las áreas recreativas de *Incidentes* y analizando si existen en *Areas*. Código:

```

1     def crear_areas_incidentes(self) -> None:
2         """Crear referencia para Areas-Incidentes"""
3         tabla_incidentes_seg = self.state["Incidentes"]
4         for i in range(len(tabla_incidentes_seg["AreaRecreativaID"])):
5             id_area = tabla_incidentes_seg.loc[i, "AreaRecreativaID"]
6             id_incidentes_seg = tabla_incidentes_seg.loc[i, "ID"]
7             if id_area is not None and id_area not in self.area_incidente:
8                 self.area_incidente[id_area] = [id_incidentes_seg]
9             else:
10                 self.area_incidente[id_area].append(id_incidentes_seg)
11

```

17. **crear incidencia usuario:** Crea la relación entre las tablas *Incidencias-Usuarios* extrayendo los id de los usuarios de *Incidencias* y analizando si existen en *Usuarios*. Código:

```

1     def crear_incidencia_usuario(self) -> None:
2         """Crear referencia para Usuario-Incidencias"""
3         tabla_incidencias = self.state["Incidencias"]
4         for i in range(len(tabla_incidencias["UsuarioID"])):
5             id_usuarios_str = tabla_incidencias.loc[i, "UsuarioID"].replace("'", '"')
6             id_usuarios_list = json.loads(id_usuarios_str)
7             id_incidencias = tabla_incidencias.loc[i, "ID"]
8             for id_usuarios in id_usuarios_list:
9                 if (
10                     id_usuarios is not None
11                     and id_usuarios not in self.incidencia_usuario
12                 ):
13                     self.incidencia_usuario[int(id_incidencias)] = [id_usuarios]
14                 else:
15                     self.incidencia_usuario[int(id_incidencias)].append(id_usuarios)
16

```

1.6. Proceso de Carga

Esta parte de nuestra pipeline es la que nos permite insertar correctamente los datos en la base de datos de destino, asegurandose de que el formato de los documentos insertados es el correcto y de que no se pierda información en el proceso de transformación de la misma.

En las dos secciones siguientes se describirán la estructura de la base de datos en la que se insertarán los datos finales y el diseño de los esquemas de validación asociados a los agregados a implementar y demás clases.

1.6.1. Estructura final de BBDD

La estructura final planteada para la base de datos se describe a continuación:

1. **AreaRecreativaClima:** Colección que representa el agregado con su mismo nombre. Incluye como raíz la clase *AreaRecreativa* y algunas referencias a los juegos, incidentes de seguridad, encuestas y clima de las areas recreativas.
2. **IncidenteSeguridad:** Colección para almacenar todos los incidentes de seguridad que sucedan en un área recreativa.
3. **EncuestaSatisfaccion:** Colección que almacenará todas las escuestas realizadas por los usuarios a una zona en concreto.
4. **RegistroClima:** Colección que proporciona datos útiles sobre las medidas temporales y precipitaciones entre otras, para el mantenimiento de las áreas.
5. **Juego:** Colección que representa el agregado "Juego" realizado en la BBDD. Éste incluirá toda la información relativa a los juegos de un área y referencias hacia mantenimientos realizados a los mismos junto con sus incidencias reportadas por los usuarios.
6. **Mantenimiento:** Colección que almacena todos los documentos que representan los mantenimientos realizados a cada juego. Fuertemente relacionada con el agregado de "Juego".
7. **Incidencia:** Agregado de incidencias a representar en la base de datos. Incluye el almacenamiento de las incidencias de los juegos reportadas por los usuarios y su correspondencia con los usuarios de la plataforma digital.
8. **Usuario:** Colección que almacena a los usuarios de la plataforma digital y beneficiarios de los servicios de las areas recreativas.

1.6.2. Diseño de Esquemas de validación

Dentro de la carpeta "Validadores" se pueden localizar todos los archivos necesarios para crear las colecciones e insertar sus esquemas de validación en la base de datos destino.

Cada archivo representa una colección de las descritas anteriormente en el apartado de la estructura final de la base de datos de destino.

La estructura empleada para la creación de estos archivos es la siguiente:

- Creación de la colección en la base de datos de destino.
- Insertar los índices para los atributos que funcionen como claves primarias.
- Definición de una variable **SCHEME** para el esquema de validación correspondiente.
- Asignación del esquema de validación a la colección creada en el paso 1.

Un ejemplo de estructura de *Esquema de validación* sería:

```
1 db.createCollection('<nombre>')
2
3 db.<nombre>.createIndex( {"<PK>":1}, {unique:1} )
4
5 SCHEME= {
6     "bsonType": "object",
7     "description": "...",
8     "required": [ "id","fecha", ... ],
9     "properties": {
10         "id": {},
11         "fecha": {},
12         ...
13     }
14 }
15
16 db.runCommand({ "collMod": "<nombre>", "validator": { $jsonSchema: SCHEME } })
```

1.7. Ejecución

El proceso de Python se ejecuta en el main. Código:

```
1 import load as ld
2 import change as ch
3 import imputation as im
4 import formater_mongo as mongo_creator
5
6
7 def main():
8     #Cargar base de datos sucia
9     base = ld.load_db()
10
11     #Limpieza de tablas (ch) + insercion de tablas (im)
12     ch.adjust_gps(base)
13     ch.adjust_ETRS89(base)
14     ch.empty_data(base)
15     im.area_new_attribute(base)
```

```

16     im.juegos_new_atributes(base)
17     im.new_meteo(base)
18     ch.no_duplicates(base)
19     im.area_meteo(base)
20     for tab_name in base: #Funciones que se ejecutan en columnas directamente
21         for column_name in base[tab_name]:
22             ch.reformatear_fecha(base, tab_name, column_name)
23             if type(base[tab_name][column_name][0]) is str:
24                 ch.capitalize_column(base, tab_name, column_name)
25     ch.delete_special(base)
26     ch.enum_checker(base)
27     ch.check_id(base)
28     ch.nif_status(base)
29     ch.formato_tlf(base)
30     ch.incidencias_status(base)
31     im.tiempoResolucion(base)
32     im.lastFecha(base)
33     im.area_new_attribute(base)
34     im.area_meteo(base)
35     im.nivelEscalamiento(base)
36
37     # Extracción de objetos
38     extractor = mongo_creator.Creator(base)
39     extractor.crear_area_clima()
40     extractor.crear_juego_incidencias()
41     extractor.crear_areas_incidentes()
42     extractor.crear_area_encuestas()
43     extractor.crear_area_juegos()
44     extractor.crear_juegos_tipo()
45     extractor.crear_area_clima()
46     extractor.crear_incidencia_usuario()
47
48     # Generar datos
49     extractor.generar_usuarios()
50     extractor.generar_encuestas()
51     extractor.generar_incidentes()
52     extractor.generar_incidencias()
53     extractor.generar_clima()
54     extractor.generar_area()
55     extractor.generar_mantenimientos()
56     extractor.generar_juegos()
57     extractor.get_json_data(extractor.incidentesSeguridad, "IncidenteSeguridad")
58     extractor.get_json_data(extractor.areas, "AreaRecreativaClima")
59     extractor.get_json_data(extractor.encuestas, "EncuestaSatisfaccion")
60     extractor.get_json_data(extractor.incidencias, "Incidencia")
61     extractor.get_json_data(extractor.juegos, "Juego")
62     extractor.get_json_data(extractor.mantenimientos, "Mantenimiento")
63     extractor.get_json_data(extractor.registrosClima, "RegistroClima")
64     extractor.get_json_data(extractor.usuarios, "Usuario")
65
66
67 if __name__ == "__main__":
68     main()
69
70

```

Del main se crean los json de datos, que se encuentran en la carpeta *DatasetsNuevos*. Luego, para meter todo en mongoDB se han sacado los siguientes JSON de validación:

1. **Validador AreaRecreativaClima:** JSON de *AreaRecreativaClima*:


```

1 db.createCollection('AreaRecreativaClima')
2
3 db.AreaRecreativaClima.createIndex({ "registrosClima.id": 1 }, { unique: 1, sparse: true })
4 db.AreaRecreativaClima.createIndex({ "encuestas.id": 1 }, { unique: 1, sparse: true })
5 db.AreaRecreativaClima.createIndex({ "incidentesSeguridad.id": 1 }, { unique: 1, sparse: true
6 })
7 db.AreaRecreativaClima.createIndex({ "juegos.id": 1 }, { unique: 1, sparse: true })
8 db.AreaRecreativaClima.createIndex({ "CantidadJuegosTipo.tipo": 1 }, { unique: 1, sparse: true
9 })
10
11 SCHEME = {
12   "bsonType": "object",
13   "description": "Agregado de gestión del estado de las áreas recreativas, juegos, incidentes
14     de seguridad y datos climáticos.",
15   "required": ["id", "coordenadasGPS", "barrio", "distrito", "fecha", "estadoOperativo", "
16     capacidadMaxima", "cantidadJuegosTipo", "juegos", "incidentesSeguridad", "registrosClima",
17     "encuestas"],
18   "properties": {
19     "id": {
20       "bsonType": "number",
21       "description": "Clave Primaria"
22     },
23     "coordenadasGPS": {
24       "bsonType": "array",
25       "description": "Array que representa la latitud y la longitud",
26       "minItems": 2,
27       "maxItems": 2,
28       "items": {
29         "bsonType": "string"
30       }
31     },
32     "fecha": {
33       "bsonType": "string",
34       "description": "fecha de instalación del área"
35     },
36     "barrio": {
37       "bsonType": "string",
38       "description": "Barrio al que pertenece el area recreativa"
39     },
40     "distrito": {
41       "bsonType": "string",
42       "description": "Distrito al que pertenece el area recreativa"
43     },
44     "estadoOperativo": {
45       "bsonType": "string",
46       "description": "Estado global del area recreativa",
47       "enum": ["OPERATIVO", "CERRADO", "INDISPUESTO"]
48     },
49     "capacidadMaxima": {
50       "bsonType": "number",
51       "description": "Numero máximo de juegos que incluye el área",
52     },
53     "cantidadJuegosTipo": {
54       "bsonType": "array",
55       "description": "Conteo de cuantos juegos hay de cada tipo en el area actual",
56       "minItems": 0,
57       "items": {
58         "bsonType": "object",
59         "description": "objeto resumen de conteo de un tipo de juegos",
60         "required": ["tipo", "valor"],

```

```

56     "properties": {
57         "tipo": {
58             "bsonType": "string",
59             "description": "tipo de juegos a contar",
60             "enum": ["DEPORTIVAS", "INFANTILES", "MAYORES"]
61         },
62         "valor": {
63             "bsonType": "number",
64             "minimum": 0
65         }
66     }
67 },
68 },
69 "juegos": {
70     "bsonType": "array",
71     "description": "array de referencia a juegos incluidos en el area",
72     "minItems": 0,
73     "items": {
74         "bsonType": "number",
75     "description": "juegos en un area",
76     }
77 },
78 "incidentesSeguridad": {
79     "bsonType": "array",
80     "description": "array de referencias a los incidentes de seguridad",
81     "minItems": 0,
82     "items": {
83         "bsonType": "number",
84         "description": "objeto de representacion de incidencias de seguridad"
85     }
86 },
87 "registrosClima": {
88     "bsonType": "array",
89     "description": "Array de referencias a los registros de clima para esta zona",
90     "minItems": 0,
91     "items": {
92         "bsonType": "number",
93     }
94 },
95 "encuestas": {
96     "bsonType": "array",
97     "description": "Array de referencias a las encuestas de satisfacción de una zona",
98     "minItems": 0,
99     "items": {
100         "bsonType": "number",
101         "description": "objeto referencia de una encuesta de satisfacción"
102     }
103 },
104 },
105 }
106
107 db.runCommand({ "collMod": "AreaRecreativaClima", "validator": { $jsonSchema: SCHEME } })
108

```

2. Validador EncuestaSatisfaccion: JSON de *EncuestaSatisfaccion*:

```

1 db.createCollection('EncuestaSatisfaccion')
2
3 db.EncuestaSatisfaccion.createIndex({ "id": 1 }, { unique: 1 })
4

```

```

5 SCHEME = {
6   "bsonType": "object",
7   "description": "Objeto que representa una encuesta de satisfacción de los usuarios hacia un
8   area recreativa",
9   "required": ["id", "fechaEncuesta", "puntuacionAccesibilidad", "puntuacionCalidad", "
10  comentarios"],
11  "properties": {
12    "id": {
13      "bsonType": "number",
14      "description": "Tipo String. Clave Primaria"
15    },
16    "fechaEncuesta": {
17      "bsonType": "string",
18      "description": "Fecha de realización de encuesta"
19    },
20    "puntuacionAccesibilidad": {
21      "bsonType": "number",
22      "minimum": 1,
23      "maximum": 5,
24      "description": "puntuaciones de los usuarios respecto a la accesibilidad",
25    },
26    "puntuacionCalidad": {
27      "bsonType": "number",
28      "minimum": 1,
29      "maximum": 5,
30      "description": "puntuaciones de los usuarios respecto a la calidad",
31    },
32    "comentarios": {
33      "bsonType": "string",
34      "description": "Lista de comentarios en formato string",
35      "enum": ["ACEPTABLE", "DEFICIENTE", "EXCELENTE", "MUY BUENO", "REGULAR"]
36    }
37  },
38 }
39
40 db.runCommand({ "collMod": "EncuestaSatisfaccion", "validator": { $jsonSchema: SCHEME } })

```

3. Validador Incidencia: JSON para *Incidencia*:

```

1 db.createCollection('Incidencia')
2
3 db.Incidencia.createIndex( {"id":1}, {unique:1} )
4
5 SCHEME={
6   "bsonType": "object",
7   "description": "Agregado de gestión del estado del juego, historial de mantenimiento y
8   seguimiento de incidencias.",
9   "required": ["id","tipo","fechaReporte","estado","tiempoResolucion","nivelEscalamiento","
10  usuarios"],
11  "properties": {
12    "id": {
13      "bsonType": "number",
14      "description": "id de la incidencia"
15    },
16    "tipo": {
17      "bsonType": "string",
18      "description": "tipo de incidencia al juego",
19      "enum": ["DESGASTE","MAL FUNCIONAMIENTO","ROTURA","VANDALISMO"]
20    }
21  },

```

```

19     "fechaReporte": {
20         "bsonType": "string",
21         "description": "fecha de reporte de la incidencia"
22     },
23     "estado": {
24         "bsonType": "string",
25         "description": "estado de la incidencia",
26         "enum": ["ABIERTA", "CERRADA"]
27     },
28     "tiempoResolucion": {
29         "bsonType": "number",
30         "description": "tiempo en el que se resolverá la incidencia",
31     },
32     "nivelEscalamiento": {
33         "bsonType": "number",
34         "minimum": 1,
35         "maximum": 10,
36         "description": "Nivel de urgencia"
37     },
38     "usuarios": {
39         "bsonType": "array",
40         "description": "array de referencias a usuarios vinculados con esta incidencia"
41     },
42     "minItems": 1,
43     "items": {
44         "bsonType": "object",
45         "description": "objeto embebido del usuario",
46         "required": ["NIF", "nombre", "email", "telefono"],
47         "properties": {
48             "NIF": {
49                 "bsonType": "string",
50                 "description": "Identificador de tipo string",
51             },
52             "nombre": {
53                 "bsonType": "string",
54                 "description": "nombre del usuario"
55             },
56             "email": {
57                 "bsonType": "string",
58                 "description": "email del usuario"
59             },
60             "telefono": {
61                 "bsonType": "string",
62                 "description": "telefono de los usuarios"
63             }
64         }
65     },
66 },
67 }
68
69 db.runCommand({ "collMod": "Incidencia", "validator": { $jsonSchema: SCHEME } })
70

```

4. Validacion IncidenteSeguridad: JSON para *IncidenteSeguridad*:

```

1 db.createCollection('IncidenteSeguridad')
2
3 db.IncidenteSeguridad.createIndex({ "id": 1 }, { unique: 1 })
4

```

```

5 SCHEME = {
6   "bsonType": "object",
7   "description": "Objeto que representa un incidente de seguridad",
8   "required": ["id", "fechaDeReporte", "tipoIncidente", "gravedad"],
9   "properties": {
10    "id": {
11     "bsonType": "number",
12     "description": "id del incidente de seguridad"
13    },
14    "tipoIncidente": {
15     "bsonType": "string",
16     "description": "tipo de incidencia de seguridad",
17     "enum": ["ROBO", "CAIDA", "ACCIDENTE", "VANDALISMO", "DAÑO ESTRUCTURAL"]
18    },
19    "gravedad": {
20     "bsonType": "string",
21     "description": "gravedad del incidente",
22     "enum": ["BAJA", "MEDIA", "ALTA", "CRITICA"]
23    },
24    "fechaDeReporte": {
25     "bsonType": "string",
26     "description": "fecha de reporte de la incidencia de seguridad"
27    }
28   },
29 }
30
31 db.runCommand({ "collMod": "IncidenteSeguridad", "validator": { $jsonSchema: SCHEME } })
32

```

5. Validador Juego: JSON para *Juego*:

```

1 db.createCollection('Juego')
2 db.Juego.createIndex({ "id": 1 }, { unique: 1 })
3
4 SCHEME = {
5   "bsonType": "object",
6   "description": "Agregado de gestión del estado del juego, historial de mantenimiento y seguimiento de incidencias.",
7   "required": ["id", "modelo", "estadoOperativo", "accesibilidad", "fechaInstalacion", "tipo", "desgasteAcumulado", "indicadorExposicion", "ultimaFechaMantenimiento", "mantenimientos", "incidencias"],
8   "properties": {
9     "id": {
10      "bsonType": "string",
11      "description": "Tipo String. Clave Primaria"
12     },
13     "modelo": {
14      "bsonType": "string",
15      "description": "Modelo de fabricacion del juego"
16     },
17     "estadoOperativo": {
18      "bsonType": "string",
19      "description": "Estado del juego",
20      "enum": ["OPERATIVO", "INDISPUESTO"]
21     },
22     "accesibilidad": {
23      "bsonType": "bool",
24      "description": "Si el juego es accesible por la gente o no"
25     },
26     "fechaInstalacion": {

```

```

27     "bsonType": "string",
28     "description": "Fecha de instalacion del juego"
29 },
30 "tipo": {
31     "bsonType": "string",
32     "description": "Tipo de juego",
33     "enum": ["DEPORTIVAS", "INFANTILES", "MAYORES"]
34 },
35 "desgasteAcumulado": {
36     "bsonType": "number",
37     "description": "porcentaje de desgaste que tiene un juego",
38     "minimum": 0,
39     "maximum": 100,
40 },
41 "indicadorExposicion": {
42     "bsonType": "string",
43     "description": "Indicador de exposición a temperaturas aridas",
44     "enum": ["INTACTO", "BAJO", "MEDIO", "ALTO"]
45 },
46 "ultimaFechaMantenimiento": {
47     "bsonType": "string",
48     "description": "ultima fecha de mantenimiento del juego"
49 },
50 "mantenimientos": {
51     "bsonType": "array",
52     "description": "array de referencias a registros de mantenimiento del juego",
53     "minItems": 0,
54     "items": {
55         "bsonType": "string",
56         "description": "objeto de referencia de mantenimiento realizado"
57     }
58 },
59 "incidencias": {
60     "bsonType": "array",
61     "description": "array de referencias a las incidencias reportadas para este juego",
62     "minItems": 0,
63     "items": {
64         "bsonType": "number",
65         "description": "objeto de referencia a la incidencia del juego"
66     }
67 },
68 },
69 }
70
71 db.runCommand({ "collMod": "Juego", "validator": { $jsonSchema: SCHEME } })
72

```

6. Validador Mantenimiento: JSON para *Mantenimiento*:

```

1 db.createCollection('Mantenimiento')
2
3 db.Mantenimiento.createIndex({ "id": 1 }, { unique: 1 })
4
5 SCHEME = {
6     "bsonType": "object",
7     "description": "Objeto que representa un incidente de seguridad",
8     "required": ["id", "fechaIntervencion", "tipoIntervencion", "estadoPrevio", "estadoPosterior"],
9     "properties": {
10         "id": {

```

```

11     "bsonType": "string",
12     "description": "id del mantenimiento a realizar a un juego"
13 },
14 "tipoIntervencion": {
15     "bsonType": "string",
16     "description": "tipo de intervencion",
17     "enum": ["CORRECTIVO", "EMERGENCIA", "PREVENTIVO"]
18 },
19 "estadoPrevio": {
20     "bsonType": "string",
21     "description": "estado previo a la intervencion",
22 },
23 "estadoPosterior": {
24     "bsonType": "string",
25     "description": "estado posterior a la intervencion",
26 },
27 "fechaIntervencion": {
28     "bsonType": "string",
29     "description": "fecha de intervencion de mantenimiento"
30 },
31 "incidencias": {
32     "bsonType": "array",
33     "description": "incidencias que referencia un mantenimiento",
34     "minItems": 0,
35     "items": {
36         "bsonType": "number",
37         "description": "objeto de referencia de incidencias realizadas"
38     }
39 },
40 },
41 }
42
43 db.runCommand({ "collMod": "Mantenimiento", "validator": { $jsonSchema: SCHEME } })
44

```

7. Validador RegistroClima: JSON para *RegistroClima*:

```

1 db.createCollection('RegistroClima')
2
3 db.RegistroClima.createIndex({ "id": 1 }, { unique: 1 })
4
5 SCHEME = {
6     "bsonType": "object",
7     "description": "Objeto que representa los datos del clima en la zona de una area recreativa"
8 ,
9     "required": ["id", "fecha", "temperatura", "precipitacion", "vientosFuertes"],
10    "properties": {
11        "id": {
12            "bsonType": "int",
13            "description": "id del registro del clima"
14        },
15        "temperatura": {
16            "bsonType": "number",
17            "description": "temperatura tomada en la zona del area recreativa",
18        },
19        "precipitacion": {
20            "bsonType": "number",
21            "description": "precipitacion tomada en la zona del area recreativa",
22        },
23        "fecha": {

```

```

23     "bsonType": "string",
24     "description": "fecha de la toma de muestras del clima"
25   },
26   "vientosFuertes": {
27     "bsonType": "bool",
28     "description": "si hacen vientos fuertes o no"
29   }
30 },
31 }
32
33 db.runCommand({ "collMod": "RegistroClima", "validator": { $jsonSchema: SCHEME } })
34

```

8. Validacion Usuario: JSON de *Usuario*:

```

1  db.createCollection('Usuario')
2
3  db.Usuario.createIndex( {"NIF":1}, {unique:1} )
4
5  SCHEME={
6    "bsonType": "object",
7    "description": "Objeto que representa un usuario de la zona recreativa",
8    "required": ["NIF","nombre","email","telefono"],
9    "properties": {
10      "NIF": {
11        "bsonType": "string",
12        "description": "Identificador del usuario"
13      },
14      "nombre": {
15        "bsonType": "string",
16        "description": "nombre de usuario de la zona",
17      },
18      "email": {
19        "bsonType": "string",
20        "description": "email del usuario",
21      },
22      "telefono": {
23        "bsonType": "string",
24        "description": "telefono del usuario",
25      }
26    },
27  },
28  db.runCommand({ "collMod": "Usuario", "validator": { $jsonSchema: SCHEME } })
29

```

Por último, ejecutamos tanto los JSONs de datos como el pipeline.txt en mongoDB. Así tendríamos todos los datos y sus relaciones en la base de datos. En el pipeline además contiene la representación *EstadoGlobalÁrea*, que inserta en *AreaRecreativaClima* el estado global teniendo en cuenta las incidencia y manetnimientos de los juegos de dicha área. Esta implementación se hizo con ayuda de ChatGPT Código del pipeline.txt:

```

1
2  use ej1
3
4  // agregado Areas Recreativas
5  db.AreaRecreativaClima.aggregate([
6    // Encontrar los Incidentes
7    {
8      $lookup:
9      {

```



```

10         from: "IncidenteSeguridad",
11         localField: "incidentesSeguridad",
12         foreignField: "id",
13         as: "RefIncidentes"
14     }
15 },
16 // Encontrar los juegos
17 {
18     $lookup:
19     {
20         from: "IncidenteSeguridad",
21         localField: "Juego",
22         foreignField: "id",
23         as: "RefJuegos"
24     }
25 },
26 // project para crear todos los atributos
27 {
28     $project:
29     {
30         // area recreativa
31         "id": 1,
32         "coordenadasGPS": 1,
33         "barrio": 1,
34         "distrito": 1,
35         "fecha": 1,
36         "estadoGlobalArea": 1,
37         "capacidadMaxima": 1,
38         "cantidadJuegosTipo": 1,
39         "juegos": 1,
40         "incidentesSeguridad": 1,
41         "registrosClima": 1,
42         "encuestas": 1,
43         // Referencia Juegos
44         RefJuego: "$juegos",
45         // Referencia con resumen IncidenteSeguridad
46         "RefIncidentes.fechaDeReporte": 1,
47         "RefIncidentes.tipoIncidente": 1,
48         "RefIncidentes.gravedad": 1,
49         // Referencia Clima
50         RefClima: "$registrosClima",
51         // Encuesta Satisfaccion
52         RefEncuestas: "$encuestas"
53     }
54 },
55 // estado global area
56 {
57     $addFields: {
58         estadoGlobalArea: {
59             $cond: {
60                 if: {
61                     $anyElementTrue: {
62                         $map: {
63                             input: { $ifNull: ["$RefIncidentes.gravedad", []] }, //
64                             as: "gravedad",
65                             in: { $eq: ["$$gravedad", "CRITICA"] } // Verifica si alguna
66                             gravedad es "CRITICA"
67                         }
68                     }
69                 }
70             }
71         }
72     }
73 }

```

```

68         },
69         then: "PELIGROSO", // Si hay "CRITICA", se asigna "PELIGRO"
70         else: {
71             $cond: {
72                 if: {
73                     $anyElementTrue: {
74                         $map: {
75                             input: { $ifNull: ["$RefIncidentes.gravedad", []] },
76                             // Asegura que RefIncidentes.gravedad sea un array
77                             as: "gravedad",
78                             in: { $eq: ["$$gravedad", "ALTA"] } // Verifica si
79                             alguna gravedad es "ALTA"
80                         }
81                     },
82                     then: "PRECAUCION", // Si hay "ALTA", se asigna "MEDIO"
83                     else: "NORMAL" // Si no hay ni "CRITICA" ni "ALTA", se asigna "
84                     BAJO"
85                 }
86             }
87         },
88     },
89     // Llevarlo a una colección nueva
90     {
91         $out: {db:"ej1", coll: "AgregadoAreaRecreativaClima"}
92     }
93 })
94
95 // Agregado Incidencia
96 db.Incidencia.aggregate([
97     // Crear atributos usando project
98     {
99         $project:
100         {
101             "id": 1,
102             "fechaReporte": 1,
103             "tipo": 1,
104             "estado": 1,
105             "tiempoResolucion": 1,
106             "nivelEscalamiento": 1,
107             RefUsuarios: "$usuarios"
108         }
109     },
110     // Llevarlo a una colección nueva
111     {
112         $out: {db:"ej1", coll: "AgregadoIncidencia"}
113     }
114 })
115
116 // Agregado Juego
117 db.Juego.aggregate([
118     // Etapa de $lookup para unir con la colección 'id'
119     {
120         $lookup: {
121             from: "Incidencia",
122             localField: "incidencias",
123             foreignField: "id",

```

```

125         as: "Inc"
126     }
127 },
128 // Etapa de $project para seleccionar y renombrar los campos deseados
129 {
130     $project: {
131         "id": 1,
132         "modelo": 1,
133         "estadoOperativo": 1,
134         "accesibilidad": 1,
135         "fechaInstalacion": 1,
136         "tipo": 1,
137         "desgasteAcumulado": 1,
138         "indicadorExposicion": 1,
139         "ultimaFechaMantenimiento": 1,
140         "mantenimientos": 1,
141         "incidencias": 1,
142         //RefMantenimiento
143         RefMant: "$Mantenimiento",
144         //RefIncidencia - selecciona campos de 'Inc' obtenido en $lookup
145         "Inc.fechaReporte": 1,
146         "Inc.tipo": 1,
147         "Inc.estado": 1
148     }
149 },
150 // Etapa de $out para exportar los resultados a otra colección
151 {
152     $out: { db: "ej1", coll: "AgregadoJuego" }
153 }
154 })
155

```

2. Datos importantes a tener en cuenta

Primero, cabe destacar que vamos a entregar la memoria y el .txt junto al código de Python y los JSONs. Así, facilitamos el uso del código sin tener que ir en la memoria apartado en apartado copiándolo y ejecutándolo. Cabe destacar que es necesario tener todas las librerías de Python usadas (pandas, datetime, json, etc.) para ejecutarlo sin problemas, además de que puede llegar a tardar entre 5-10 min en su ejecución.

Además, no hemos representado los JSONs con los datos resultantes de main() debido a su gran extensión de líneas. Los hemos dejado en la carpeta *DatasetsNuevos* para su uso en mongoDB.

Por último, hemos indicado el uso de ChatGPT en todas las funciones en las que nos ha ayudado y cómo ha realizado dicha ayuda.