### 2. Análisis Inicial de la Base de Datos

Antes de comenzar a limpiar la base de datos, debe realizar un análisis preliminar para comprender la naturaleza y distribución de los errores. Para ello, siga los siguientes pasos:

• Mostrar un resumen estadístico de los datos.

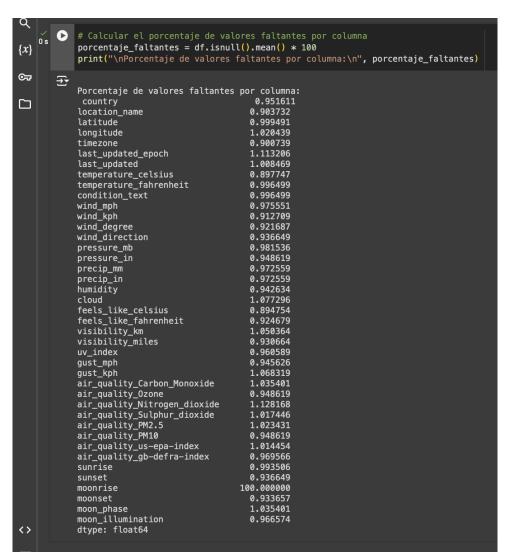
(~)										
<b>©</b> ⊋ 0′s	0		_estadísti Resumen es					co)		
	至	count unique top freq mean std min 25% 50% 75% max	estadístic country 33099 211 Bulgaria 369 NaN NaN NaN NaN NaN NaN NaN NaN NaN Na	location 3 a	33115 248 aaaaa 326 NaN NaN NaN NaN NaN NaN NaN	33083 235 40.4 343 NaN NaN NaN NaN NaN NaN	longitude 33076 240 aaaaa 382 NaN NaN NaN NaN NaN NaN NaN	3: Asia/Bang	ezone \ 3116 191 gkok 604 NaN NaN NaN NaN NaN NaN NaN NaN NaN	
		count unique top freq mean std min 25% 50% 75% max	temperaturo	33	3084 794	33 Partly clo	3084	air_quali	ty_PM2.5 33075 1818 0.5 1277 NaN NaN NaN NaN NaN NaN NaN	`

<del></del>	count unique top	air_qualit	ty_PM10 33100 2717 0.6	air_quali	.ty_ı		ndex 3078 7	air_quality	_gb-defra-index 33093 11 1.0	\
	freq		360			2	4084		21874	
	mean		NaN				NaN		NaN	
	std		NaN				NaN		NaN	
	min		NaN				NaN		NaN	
	25%		NaN				NaN		NaN	
	50%		NaN				NaN		NaN	
	75%		NaN				NaN		NaN	
	max		NaN				NaN		NaN	
	count unique top	sunrise 33085 309 05:56 AM	sunset 33104 408 aaaaa	moonrise 0.0 NaN NaN	No	moonset 33105 1442 moonset	Wax	moon_phase 33071 9 cing Gibbous	moon_illuminat: 330	
	freq	428	340	NaN		1005		8137		138
	mean	NaN	NaN	NaN		NaN		NaN		NaN
	std	NaN	NaN	NaN		NaN		NaN		NaN
	min	NaN	NaN	NaN		NaN		NaN		NaN
	25%	NaN	NaN	NaN		NaN		NaN		NaN
	50%	NaN	NaN	NaN		NaN		NaN		NaN
	75%	NaN	NaN	NaN		NaN		NaN		NaN NaN
	max	NaN	NaN	NaN		NaN		NaN	r	NaN
	[11 rov	vs x 41 co	lumns]							

### Resumen estadístico:

- Se incluyen descripciones básicas como el número de valores, valores únicos, y el valor más frecuente (moda) para cada columna.
- La columna "country" tiene 211 valores únicos, mientras que "location\_name" tiene 248.
- Algunas columnas como "latitude", "longitude", y "last\_updated\_epoch" contienen valores poco consistentes, como "aaaaa", lo que sugiere errores en los datos.

• Calcular el porcentaje de valores faltantes por columna.



# Porcentaje de valores faltantes por columna:

 Hay algunas columnas con un porcentaje considerable de valores faltantes.

Por ejemplo:

- moonrise: 100% valores faltantes.
- sunrise:
   99.35%
   valores
   faltantes.
- moon\_phase: 1.03%valoresfaltantes.

Identificar si hay filas duplicadas.

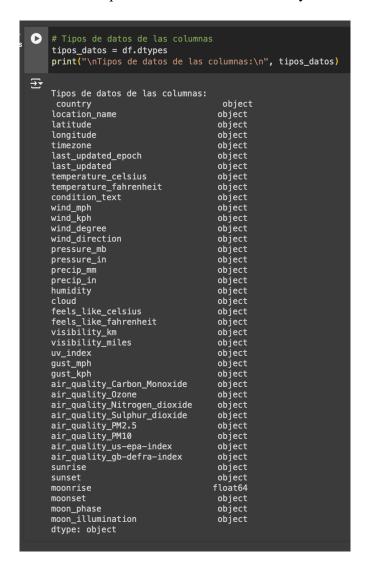
```
# Identificar filas duplicadas
duplicados = df.duplicated().sum()
print(f"\nNúmero de filas duplicadas: {duplicados}")

Número de filas duplicadas: 1504
```

• Filas duplicadas:

Se han identificado 1,504 filas duplicadas en el conjunto de datos.

• Analizar los tipos de datos de las columnas y si son consistentes con el contenido esperado.



### Tipos de datos por columna:

Muchas columnas que deberían ser numéricas están registradas como texto (object).

Por ejemplo, columnas como "latitude", "longitude", "tempe rature\_celsius", y "humidity" están categorizadas como texto en lugar de números, lo cual indica inconsistencias en los tipos de datos.

### 3. Limpieza de Datos

Debe realizar las siguientes tareas de limpieza en la base de datos:

• Eliminación o imputación de valores faltantes: Justificar si decide eliminar filas/columnas con NaNs o utilizar técnicas de imputación (relleno de valores).

```
[18] # 1. Eliminación o Imputación de Valores Faltantes

# Revisar porcentaje de valores faltantes
porcentaje_faltantes = df.isnull().mean() * 100

# Decidir si eliminar o imputar en base a los valores faltantes
threshold = 50 # 5i más del 50% de los valores faltan, eliminamos la columna
columnas_a_eliminar = porcentaje_faltantes[porcentaje_faltantes > threshold].index
df_limpio = df.drop(columns=columnas_a_eliminar)

# Para las columnas restantes, imputamos usando la mediana (opcional)
df_limpio.fillna(df_limpio.median(numeric_only=True), inplace=True)
df_limpio
```

- •Eliminación de columnas con muchos NaNs: Si más del 50% de los datos en una columna están ausentes, es mejor eliminar la columna, ya que su presencia no aporta información suficiente.
- Imputación de valores faltantes: Para las columnas clave con pocos valores faltantes, imputamos usando la mediana. También se puede usar la media o un valor específico si es más adecuado.
- Eliminación de duplicados: Identificar filas duplicadas y eliminarlas.

```
# 2. Eliminación de Duplicados
df_limpio = df_limpio.drop_duplicates()
df_limpio.head()
```

• Corrección de tipos de datos: Asegurarse de que las columnas tengan tipos de datos adecuados (por ejemplo, números como int o float, fechas como datetime, textos como str).

• Corrección de valores "invalid": Corregir los valores que fueron etiquetados con cadenas incorrectas como 'bbb'.

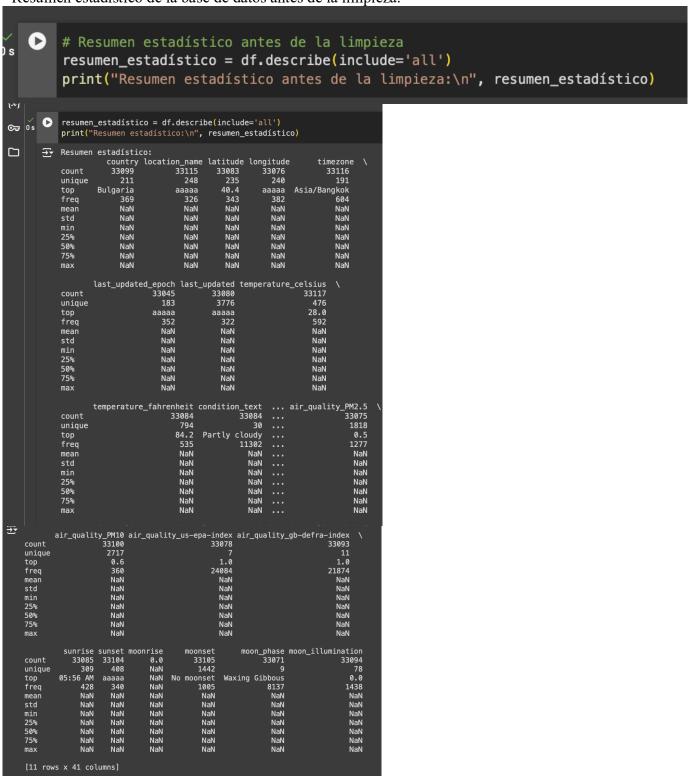
```
# 4. Corrección de Valores "Invalid"
# Identificar y corregir valores como 'aaa', 'bbb', etc.
valores_invalidos = ['aaa', 'aaaaa', 'abbb', 'cccc', 'xxxxx'] # Añadir más si es necesario
df_limpio.replace(valores_invalidos, pd.NA, inplace=True)
df_limpio.head()
```

### 4. Documentación y Reporte

Luego de limpiar los datos, redacte un reporte en el cual documente los siguientes aspectos:

Análisis inicial:

• Resumen estadístico de la base de datos antes de la limpieza.



• Tabla que muestre el porcentaje de valores faltantes por columna.

```
{x} os Porcentaje de valores faltantes por columna antes de la limpieza porcentaje_faltantes = df.isnull().mean() * 100
                      print("\nPorcentaje de valores faltantes por columna antes de la limpieza:\n", porcentaje_faltantes)
⊙⊋
                     Porcentaje de valores faltantes por columna antes de la limpieza: country 0.951611 location_name 0.993732 latitude 0.999491 longitude 1.020439
latitude
longitude
timezone
                      last_updated_epoch
last_updated
temperature_celsius
                                                                                      1.113206
                                                                                      1.008469
0.897747
                      temperature_fahrenheit
condition_text
wind_mph
wind_kph
                                                                                      0.921687
0.936649
0.981536
0.948619
                      pressure mb
                      pressure_in
precip_mm
precip_in
                       humidity
                                                                                      1.077296
0.894754
0.924679
                      cloud feels_like_celsius
                      feels_like_fahrenheit
visibility_km
visibility_miles
                      gust_mph
gust_kph
air_quality_Carbon_Monoxide
                      air_quality_Ozone
air_quality_Nitrogen_dioxide
air_quality_Sulphur_dioxide
                      air_quality_PM2.5
air_quality_PM10
air_quality_us-epa-index
air_quality_gb-defra-index
                                                                                      1.023431
                      moonset
moon_phase
moon_illumination
dtype: float64
<>
```

Total de filas duplicadas encontradas.

```
[35] # Número de filas duplicadas antes de la limpieza duplicados = df.duplicated().sum() print(f"\nNúmero de filas duplicadas antes de la limpieza: {duplicados}")

Número de filas duplicadas antes de la limpieza: 1504
```

 Descripción de los tipos de datos originales y los problemas encontrados.

```
▶ # Tipos de datos antes de la limpieza
    tipos_datos = df.dtypes
    print("\nTipos de datos antes de la limpieza:\n", ti
    Tipos de datos antes de la limpieza:
    country
location_name
                                           obiect
                                           object
     latitude
    longitude
                                          object
    last_updated_epoch
last_updated
                                           object
                                          object
    temperature_celsius
temperature_fahrenheit
condition_text
                                          object
    wind_kph
wind_degree
                                           object
    wind_direction
    pressure_mb
                                           object
    pressure_in
                                           object
    precip_mm
    precip_in
                                           object
    humidity
```

```
feels_like_celsius
feels_like_fahrenheit
visibility_km
                                         object
                                         object
visibility_miles
                                         object
uv_index
                                         object
gust_mph
                                         object
gust_kph
                                         object
air_quality_Carbon_Monoxide
                                         object
air_quality_Ozone
air_quality_Nitrogen_dioxide
air_quality_Sulphur_dioxide
                                         object
                                         object
                                         object
air_quality_PM2.5
air_quality_PM10
                                         object
                                         object
air_quality_us-epa-index
                                         object
air_quality_gb-defra-index
                                         object
sunrise
                                         object
sunset
                                         object
moonrise
                                        float64
moonset
                                         object
moon_phase
                                         object
moon_illumination
                                         object
dtype: object
```

# Tipos de datos por columna:

Muchas columnas que deberían ser numéricas están registradas como texto (object).

Por ejemplo, columnas como "latitude", "longitude", "temperature\_celsius", y "humidity" están categorizadas como texto en lugar de números, lo cual indica inconsistencias en los tipos de datos.

- o Proceso de limpieza:
- Describir qué métodos utilizaron para limpiar la base de datos (eliminación, imputación, etc.):
  - Eliminación de columnas con más del 50% de valores faltantes.
  - Imputación de valores faltantes en las columnas clave usando la mediana.
  - Eliminación de duplicados para garantizar datos únicos.
  - Corrección de valores inválidos y tipos de datos.
- Mostrar antes y después de cada paso clave (por ejemplo, antes y después de eliminar duplicados).

### ANTES:

-	country	location_name	latitude	longitude	timezone	last_updated_epoch	last_updated	temperature_celsius	temperature_fahrenheit	condition_text	•••	air_quality_PM2.5	air_quality_PM10	air_quality_us- a epa-index
0	Afghanistan	Kabul	NaN	69.18	Asia/Kabul	1715849100.0	2024-05-16 13:15	26.6	79.8	Partly Cloudy		8.4	26.6	1.0
1	Albania	Tirana	41.33	19.82	Europe/Tirane	1715849100.0	2024-05-16 10:45	19.0		Partly cloudy				1.0
2	Algeria	Algiers	36.76	aaaaa	Africa/Algiers	1715849100.0	2024-05-16 09:45	23.0	73.4	Sunny		10.4	18.4	1.0
3	Andorra	Andorra La Vella	42.5	1.52	Europe/Andorra	1715849100.0	2024-05-16 10:45	6.3	43.3	Light drizzle			0.9	1.0
4	Angola	Luanda	-8.84	13.23	Africa/Luanda	1715849100.0	2024-05-16 09:45	26.0	78.8	Partly cloudy		183.4	262.3	5.0
5 ro	ws × 41 colum	ns												

DESPUES de Eliminación o imputación de valores faltantes:

	country	location_name	latitude	longitude	timezone	last_updated_epoch	last_updated	temperature_celsius	temperature_fahrenheit	condition_text	air_quality_Sulphur_dioxide	air_quality_PM2.5 a
0 /	Afghanistan	Kabul	NaN	69.18	Asia/Kabul	1715849100.0	2024-05-16 13:15	26.6	79.8	Partly Cloudy		2 8.4
1	Albania	Tirana	41.33	19.82	Europe/Tirane	1715849100.0	2024-05-16 10:45	19.0	66.2	Partly cloudy		1.1
2	Algeria	Algiers	36.76	aaaaa	Africa/Algiers	1715849100.0	2024-05-16 09:45	23.0	73.4	Sunny	13.4	10.4
3	Andorra	Andorra La Vella	42.5	1.52	Europe/Andorra	1715849100.0	2024-05-16 10:45	6.3	43.3	Light drizzle	0.2	2 0.7
4	Angola	Luanda	-8.84	13.23	Africa/Luanda	1715849100.0	2024-05-16 09:45	26.0	78.8	Partly cloudy	31.8	5 183.4

# DESPUES de Eliminación de duplicados:

÷	country	location_name	latitude	longitude	timezone	last_updated_epoch	last_updated	temperature_celsius	temperature_fahrenheit	condition_text	air_quality_Sulphur_dioxide	air_quality_PM2.5	air_quality
	0 Afghanistan	Kabul	NaN	69.18	Asia/Kabul	1715849100.0	2024-05-16 13:15	26.6	79.8	Partly Cloudy		8.4	
	1 Albania	Tirana	41.33	19.82	Europe/Tirane	1715849100.0	2024-05-16 10:45	19.0	66.2	Partly cloudy			
	2 Algeria	Algiers	36.76	aaaaa	Africa/Algiers	1715849100.0	2024-05-16 09:45	23.0	73.4	Sunny	13.4	10.4	
	3 Andorra	Andorra La Vella	42.5	1.52	Europe/Andorra	1715849100.0	2024-05-16 10:45		43.3	Light drizzle			
	4 Angola	Luanda	-8.84	13.23	Africa/Luanda	1715849100.0	2024-05-16 09:45	26.0	78.8	Partly cloudy	 31.5	183.4	

# DESPUES de Corrección de Tipos de Datos:

₹	country	location_name	latitude	longitude	timezone	last_updated_epoch	last_updated	temperature_celsius	temperature_fahrenheit	condition_text	air_quality_Sulphur_dioxide	air_quality_PM2.5	air_quality
	0 Afghanistan	Kabul	NaN	69.18	Asia/Kabul	1715849100.0	2024-05-16 13:15:00	26.6	79.8	Partly Cloudy		8.4	1
	1 Albania	Tirana	41.33	19.82	Europe/Tirane	1715849100.0	2024-05-16 10:45:00	19.0	66.2	Partly cloudy			
	2 Algeria	Algiers	36.76	NaN	Africa/Algiers	1715849100.0	2024-05-16 09:45:00	23.0	73.4	Sunny	13.4	10.4	
	3 Andorra	Andorra La Vella	42.50	1.52	Europe/Andorra	1715849100.0	2024-05-16 10:45:00	6.3	43.3	Light drizzle	0.2		
	4 Angola	Luanda	-8.84	13.23	Africa/Luanda	1715849100.0	2024-05-16 09:45:00	26.0	78.8	Partly cloudy	31.5	183.4	
	5 rows × 40 colu	mns											

# DESPUES de Correccion de Valores Invalid:

*	temperature_celsius	temperature_fahrenheit	condition_text	air_quality_Sulphur_dioxide	air_quality_PM2	2.5	air_quality_PM10	air_quality_us- epa-index		sunrise	sunset	moonset	moon_phase	moon_illumination
	26.6	79.8	Partly Cloudy			8.4	26.6			04:50 AM	06:50 PM	01:11 AM		55.0
	19.0	66.2	Partly cloudy							05:21 AM	07:54 PM	02:14 AM	Waxing Gibbous	55.0
	23.0	73.4	Sunny	13.4		0.4	18.4			05:40 AM	07:50 PM	02:14 AM		55.0
		43.3	Light drizzle	0.2			0.9		≪NA>	06:31 AM	09:11 PM	03:31 AM		55.0
	26.0	78.8	Partly cloudy		18:	3.4	262.3			06:12 AM	05:55 PM		Waxing Gibbous	

#### o Resultados:

Resumen final de la base de datos después de la limpieza.

```
# Resumen estadístico después de la limpieza
resumen_final = df_limpio.describe(include='all')
     print("Resumen estadístico después de la limpieza:\n", resumen_final)
France Resumen estadístico después de la limpieza:
                   country location_name latitude longitude
32732 32789 32755 32694
210 247 234 239
ulgaria Muscat 40.4 14.51
                                                                                     timezone \
32786
                                                                    239 190
14.51 Asia/Bangkok
                 Bulgaria
     top
                        369
                                                                                           604
                                                 last_updated temperature_celsius
32758 32733
               last_updated_epoch
32693
     count
                       182 3775
1721133900.0 2024-07-16 14:45
263 63
                                                                                          475
28.0
     top
freq
               temperature_fahrenheit condition_text ... 32764 32718 ... 793 29 ...
                                                                      ... air_quality_Sulphur_dioxide
     count
     unique
                                         84.2 Partly cloudy ...
535 11302 ...
     top
freq
               air_quality_PM2.5 air_quality_PM10 air_quality_us-epa-index \
     count
                                32726
1817
                                                        32768
2716
     unique
     top
freq
              air_quality_gb-defra-index
32759
10
1.0
21874
                                                                        sunset
32764
407
                                                      sunrise 32764
32763 32764
308 407 1441
05:56 AM 06:28 PM No moonset
428 338 1005
     count
     unique
     top
freq
                      moon_phase moon_illumination
                             32726
8
     unique
                Waxing Gibbous
8137
     [4 rows x 40 columns]
```

• Confirmar que los tipos de datos son correctos.

```
0
       # Tipos de datos después de la limpieza
         tipos_datos_despues = df_limpio.dtypes
         print("\nTipos de datos después de la limpieza:\n", tipos_datos_despues)
Ŧ
       Tipos de datos después de la limpieza:
country
location_name
latitude
                                                                                             object
float64
        longitude
timezone
                                                                                             float64
                                                                                              object
        last_updated_epoch
last_updated
temperature_celsius
temperature_fahrenheit
condition_text
                                                                            object
datetime64[ns]
                                                                                            float64
float64
                                                                                             object
float64
        wind_mph
wind_kph
        wind_degree
wind_direction
                                                                                              object
object
        pressure_mb
pressure_in
                                                                                             float64
                                                                                             float64
        precip_mm
precip_in
                                                                                            object
float64
       precip_in
humidity
cloud
feels_like_celsius
feels_like_fahrenheit
visibility_km
visibility_miles
visidex
                                                                                             object
float64
float64
                                                                                             float64
                                                                                             object
float64
        uv_index
gust_mph
gust_kph
air_quality_Carbon_Monoxide
air_quality_Nitrogen_dioxide
air_quality_Sulphur_dioxide
air_quality_PM2.5
air_quality_PM10
air_quality_us-epa-index
air_quality_us-epa-index
air_quality_gb-defra-index
suprise
                                                                                              object
object
                                                                                              object
object
                                                                                              object
object
                                                                                              object
object
         sunset
        moon_phase
moon_illumination
dtype: object
```

• Tabla que muestre el porcentaje de valores faltantes final porcolumna.

```
# Porcentaje de valores faltantes por columna después de la limpieza
      porcentaje_faltantes_despues = df_limpio.isnull().mean() * 100
print("\nPorcentaje de valores faltantes por columna después de la limpieza:\n", porcentaje_faltantes_despues)
₹
     Porcentaje de valores faltantes por columna después de la limpieza: country 2.049855 location_name 1.879283
      latitude
                                                      1.981028
      longitude
                                                      2.163569
                                                      1.888260
2.166562
1.972050
      last_updated_epoch last_updated
      temperature_celsius
temperature_fahrenheit
condition_text
                                                      2.046862
1.954095
2.091750
      wind_mph
wind_kph
wind_degree
                                                       1.987013
                                                      1.879283
2.055840
      wind_direction
                                                       2.010953
                                                      2.004968
2.046862
      pressure mb
      pressure_in
      precip_mm
                                                       1.939133
      precip_in
humidity
                                                       1.987013
                                                       1.957088
      cloud
feels_like_celsius
feels_like_fahrenheit
visibility_km
                                                       2.199479
                                                       1.873298
                                                       2.010953
                                                       2.061825
      visibility_miles
                                                       1.969058
      uv_index
                                                       1.927163
                                                      1.954095
2.151599
2.040877
      gust_kph
      air_quality_Carbon_Monoxide
     air_quality_Ozone
air_quality_Nitrogen_dioxide
air_quality_Sulphur_dioxide
air_quality_PM2.5
                                                       2.088757
                                                      2.220427
1.987013
                                                       2.067810
      air_quality_PM10
air_quality_us-epa-index
air_quality_gb-defra-index
                                                       1.942125
                                                       2.058832
                                                       1.969058
                                                       1.957088
      sunrise
      sunset
                                                       1.900230
      moon_phase
moon_illumination
                                                       2.067810
      dtype: float64
```

Comprobación de que no hay duplicados ni valores inválidos

```
[68] # Comprobar si aún existen valores inválidos
valores_invalidos = ['aaa', 'aaaa', 'bbb', 'cccc', 'xxxxxx'] # Añadir más si es necesario
df_invalidos = df_limpio[df_limpio.isin(valores_invalidos).any(axis=1)]
print(f"\nFilas con valores inválidos después de la limpieza: {df_invalidos.shape[0]}")

Filas con valores inválidos después de la limpieza: 0
```