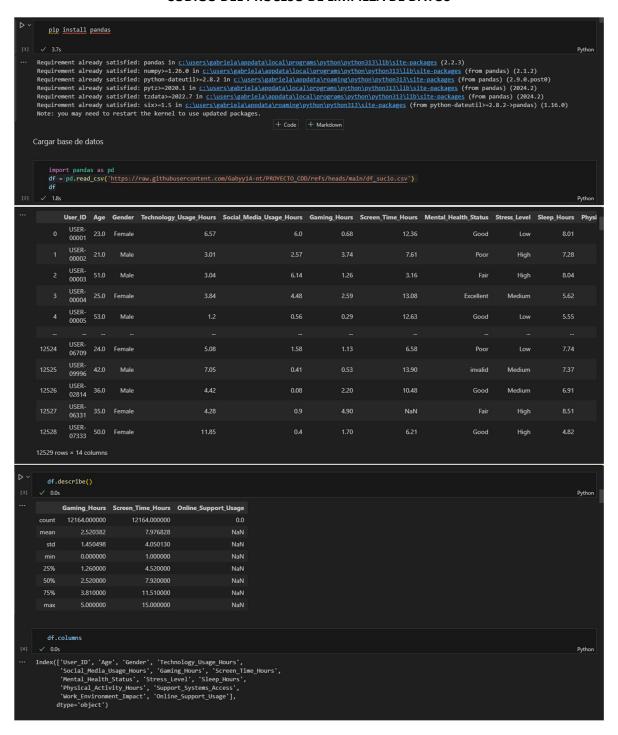
## CÓDIGO DEL PROCESO DE LIMPIEZA DE DATOS



```
Python
     RangeIndex: 12529 entries, 0 to 12528
Data columns (total 14 columns):
# Column Non-Null Count Dtype
                                                     12417 non-null object
           12397 non-null object
             Stress_Level
            Sleep_Hours
Physical_Activity_Hours
                                                     12405 non-null object
     10 Physical_activity_Hours
11 Support_Systems_Access
12 Work_Environment_Impact
13 Online_Support_Usage
dtypes: float64(3), object(11)
memory usage: 1.3+ MB
                                                   12393 non-null object
12378 non-null object
                                                    0 non-null
                                                                            float64
                                                                                                                                                                                                                                                  Python
     User_ID
Age
Gender
Technology_Usage_Hours
Social Media_Usage_Hours
Gaming_Hours
Screen_Time_Hours
Streen_Time_Hours
Stress_Level
                                                    108
124
                                                    123
365
365
      Sleep_Hours
Physical_Activity_Hours
                                                    141
124
      Support_Systems_Access
Work_Environment_Impact
Online_Support_Usage
dtype: int64
                                                    136
                                                151
12529
           #Elimino culumna Online_Support_Usage, todos sus datos son NaN df-df.drop(columns=['Online_Support_Usage'])
[7] V 0.0s
                                                                                                                                                                                                                                                  Python
                 User_ID Age Gender Technology_Usage_Hours Social_Media_Usage_Hours Gaming_Hours Screen_Time_Hours
                                                                                                                                                                               Mental Health Status
                                                                                                                                                                                                              Stress_Level Sleep_Hours Physi
                  00001 23.0 Female
                                                                                                                                           0.68
                                                                                                                                                                      12.36
                                                                                                                                                                                                    Good
                                                                                                                                                                                                                         Low
                   USER-
00002
                                                                                                                                                                                                                        High
                                                                                                                                                                                                      Poor
                   USER-
00003
                   USER-
00004
                   USER-
06709
                   USER-
09996
                   USER-
                   02814 36.0
                   06331 35.0 Female
                                                                                                                                                                      NaN
                  07333 50.0 Female
                                                                                                                                                                                                    Good
```

```
df.duplicated()
        False
False
        False
False
        ...
True
False
12524
12525
12526
12528
Length: 12529, dtype: bool
  df=df.drop_duplicates()
df
                              0 USER-
00001 23.0 Female
                                                                                                            12.36
                                                                                                                                 Good
                                                                                                                                              Low
        USER-
00002 21.0 Male
                                                                                                                                 Poor
                                                                                                                                              High
                       Male
                                                 3.04
                                                                                                                                                           8.04
                                                                                                                                  Fair
        00003
        USER-
00004
               25.0 Female
                                                                                                            13.08
                                                                                                                              Excellent
                                                                                                                                           Medium
        USER-
00005
                                                                                                                                 Good
                      Other
                                                                                                                                           Medium
                                                                                                                                Good
        01052
                                                                                                                                           Medium
                       Male
                                               invalid
                                                                                                                              Excellent
        USER-
02010
                      Other
        USER-
09996 42.0
        USER-
06331 35.0 Female
  df2 = df.rename(columns= {'User_ID': 'Usuario', 'Age': 'Edad', 'Gender': 'Género', 'Technology_Usage_Hours': 'Horas_de_uso_de_la_tecnologia', 'Social_Media_Usdf2
                                  Horas_de_uso
       Usuario Edad Género de la tecnología
                                                Horas_de_uso_de_las_redes_sociales Horas_de_juego Horas_de_tiempo_en_pantalla Estado_de_salud_mental Nivel de_estrés H
                23.0 Female
                                                                                                                                              Good
         00001
        USER-
00002
         USER-
00003
        USER-
00004
        USER-
00005
        USER-
06709
        USER-
09996
                                                                                                                                                            Medium
        USER-
02814
                35.0 Female
                                                                                                                        NaN
        USER-
07333
                50.0 Female
                                                                                                                                              Good
                                                                                                                                                              High
```

```
Renombrar datos de columnas.
```

	df2['genero']-df2['genero'].replace(gen) df2['Estado de salud mental']-df2['Estado de salud mental'].replace(sm) df2['Nivel_de_estrés']-df2['Nivel_de_estrés'].replace(ne) df2['Acceso_a_sistemas_de_soports']-df2['Acceso_a_sistemas_de_soports'].replace(ss) df2['Impacto_en_el_entorno_laboral']-df2['Impacto_en_el_entorno_laboral'].replace(el) df2		
[10]	✓ 0.0s	Python	

	Usuario	Edad	Género	Horas_de_uso _de_la_tecnología	Horas_de_uso_de_las_redes_sociales	Horas_de_juego	Horas_de_tiempo_en_pantalla	Estado_de_salud_mental	Nivel_de_estrés	F
	USER- 00001	23.0	Mujer	6.57	6.0	0.68	12.36	Buena	Bajo	
	USER- 00002	21.0	Hombre	3.01	2.57	3.74	7.61	Mala	Alto	
	USER- 00003	51.0	Hombre	3.04	6.14	1.26	3.16	Neutral	Alto	
	USER- 00004	25.0	Mujer	3.84	4.48	2.59	13.08	Excelente	Medio	
	USER- 00005	53.0	Hombre	1.2	0.56	0.29	12.63	Buena	Bajo	
12524	USER- 06709	24.0	Mujer	5.08	1.58	1.13	6.58	Mala	Bajo	
12525	USER- 09996	42.0	Hombre	7.05	0.41	0.53	13.90	invalid	Medio	
12526	USER- 02814	36.0	Hombre	4.42	0.08	2.20	10.48	Buena	Medio	
12527	USER- 06331	35.0	Mujer	4.28	0.9	4.90	NaN	Neutral	Alto	
12528	USER- 07333	50.0	Mujer	11.85	0.4	1.70	6.21	Buena	Alto	
12529 rd	ows × 13 cc	lumns								

```
Eliminamos los invalid.
                 for i in lista col:
           ✓ 0.0s
       En la columna Usuario los invalid son: 232
En la columna Edad los invalid son: 249
En la columna Género los invalid son: 265
En la columna Horas_de_uso_de la tecnología los invalid son: 249
En la columna Horas_de_uso_de_l
        En la columna Horas_de_juego los invalid son: θ
En la columna Horas_de_tiempo_en_pantalla los invalid son: θ
        En la columna Estado_de_salud_mental los invalid son: 246
En la columna Nivel_de_estrés los invalid son: 257
        En la columna Horas de sueño los invalid son: 252
En la columna Horas de actividad física los invalid son: 273
En la columna Acceso a sistemas de soporte los invalid son: 237
         En la columna Impacto_en_el_entorno_laboral los invalid son: 234
              for i in lista_col:
    df2=df2[df2[i] != 'invalid']
[12] V 0.1s
                                                                                                                                                                                                                                                                                                                                       Python
                        Usuario Edad Género _de_la_te
                                                                              Horas_de_uso
                                                                                                          Horas_de_uso_de_las_redes_sociales Horas_de_juego Horas_de_tiempo_en_pantalla Estado_de_salud_mental Nivel_de_estrés l
                                          23.0 Mujer
                                                                                                                                                                                                 0.68
                                                                                                                                                                                                                                                                                                   Buena
                                                                                                                                                                                                                                                                                                                                     Bajo
                            00001
                                           21.0 Hombre
                                                                                                                                                                                                                                                                                                      Mala
                                                                                                                                                                                                                                                                                                                                      Alto
                            00002
                                           51.0 Hombre
                                                                                                                                                                                                                                                                                                                                      Alto
                                                                                                3.04
                                                                                                                                                                                                                                                                                                 Neutral
                            00003
                            USER-
00004
                                                                                                                                                                                                                                                                                                                                  Medio
                           USER-
00005
                                                                                                9.69
                                                                                                                                                                                                 4.98
                                                                                                                                                                                                                                                                                                  Neutral
                                                                                                                                                                                                                                                                                                                                     Bajo
                            USER-
                                                                                                5.08
                                                                                                                                                                                                                                                                                                      Mala
                            USER-
                                                                                                4.28
                                                                                                                                                                                                                                                        NaN
                            USER-
                           07333 50.0
                                                          Mujer
                                                                                              11.85
                                                                                                                                                                                                                                                                                                   Buena
                                                                                                                                                                                                                                                                                                                                      Alto
               for i in lista col:
                     print(f"En la columna {i} los invalid son: {df2[df2[i] == 'invalid'].shape[0]}")
       En la columna Usuario los invalid son: 0
En la columna Edad los invalid son: 0
       En la columna Edad los invalid son: 0
En la columna Género los invalid son: 0
En la columna Horas_de_uso_de_la_tecnología los invalid son: 0
En la columna Horas_de_uso_de_las_redes_sociales los invalid son: 0
En la columna Horas_de_juego los invalid son: 0
En la columna Horas_de tiempo_en_pantalla los invalid son: 0
En la columna Estado_de_salud_mental los invalid son: 0
En la columna Nivel_de_estrés los invalid son: 0
En la columna Horas_de_sumēn_los invalid son: 0
        En la columna Horas_de_sueño los invalid son: θ
En la columna Horas_de_actividad_física los invalid son: θ
        En la columna Access_a_sistemas_de_soporte los invalid son: \theta En la columna Impacto_en_el_entorno_laboral los invalid son: \theta
       Remplazamos los espacios vacíos.
              df2['Edad'] = pd.to_numeric(df2['Edad'], errors='coerce')
df2['Horas_de_uso_de_la_tecnologia'] = pd.to_numeric(df2['Horas_de_uso_de_la_tecnologia'], errors='coerce')
df2['Horas_de_uso_de_las_redes_sociales'] = pd.to_numeric(df2['Horas_de_uso_de_las_redes_sociales'], errors='coerce')
df2['Horas_de_sueño'] = pd.to_numeric(df2['Horas_de_sueño'], errors='coerce')
df2['Horas_de_actividad_fisica'] = pd.to_numeric(df2['Horas_de_actividad_fisica'], errors='coerce')
df2
```

✓ 0.0s

	Usuario	Edad	Género	Horas_de_uso _de_la_tecnología	Horas_de_uso_de_las_redes_sociales	Horas_de_juego	Horas_de_tiempo_en_pantalla	Estado_de_salud_mental	Nivel_de_estrés	F
	USER- 00001	23.0	Mujer	6.57	6.00	0.68	12.36	Buena	Bajo	
	USER- 00002	21.0	Hombre	3.01	2.57	3.74	7.61	Mala	Alto	
	USER- 00003	51.0	Hombre	3.04	6.14	1.26	3.16	Neutral	Alto	
	USER- 00004	25.0	Mujer	3.84	4.48	2.59	13.08	Excelente	Medio	
	USER- 00005	53.0	Hombre	1.20	0.56	0.29	12.63	Buena	Bajo	
12522	USER- 05686	62.0	Otro	9.69	3.28	4.98	2.93	Neutral	Bajo	
12523	USER- 02010	25.0	Otro	3.82	3.71	3.61	6.57	Mala	Bajo	
12524	USER- 06709	24.0	Mujer	5.08	1.58	1.13	6.58	Mala	Bajo	
12527	USER- 06331	35.0	Mujer	4.28	0.90	4.90	NaN	Neutral	Alto	
12528	USER- 07333	50.0	Mujer	11.85	0.40	1.70	6.21	Buena	Alto	
10043 rd	ws × 13 co	lumns								
df2	['Edad'].	fillna	(df2['Eda	d'].mean(), inpla	ce=True)					7

	Usuario	Edad	Género	Horas_de_uso _de_la_tecnología	Horas_de_uso_de_las_redes_sociales	Horas_de_juego	Horas_de_tiempo_en_pantalla	Estado_de_salud_mental	Nivel_de_estrés	F
	USER- 00001	23.0	Mujer	6.57	6.00	0.68	12.360000	Buena	Bajo	
	USER- 00002	21.0	Hombre	3.01		3.74	7.610000	Mala	Alto	
	USER- 00003	51.0	Hombre	3.04	6.14	1.26	3.160000	Neutral	Alto	
	USER- 00004	25.0	Mujer	3.84	4.48	2.59	13.080000	Excelente	Medio	
	USER- 00005	53.0	Hombre	1.20	0.56	0.29	12.630000	Buena	Bajo	
12522	USER- 05686	62.0	Otro	9.69	3.28	4.98	2.930000	Neutral	Bajo	
12523	USER- 02010	25.0	Otro	3.82	3.71	3.61	6.570000	Mala	Bajo	ı,
12524	USER- 06709	24.0	Mujer	5.08	1.58	1.13	6.580000	Mala	Bajo	
12527	USER- 06331	35.0	Mujer	4.28	0.90	4.90	7.965271	Neutral	Alto	
12528	USER- 07333	50.0	Mujer	11.85	0.40	1.70	6.210000	Buena	Alto	

#Elimino esta columna <u>Usuario porque</u> no aporta nada a la <u>investigación</u>.

df3-df2.drop(columns-['<u>Usuario</u>'])

df3

df3

v 0.0s

Python

	Edad	Género	Horas_de_uso _de_la_tecnología	Horas_de_uso_de_las_redes_sociales	Horas_de_juego	Horas_de_tiempo_en_pantalla	Estado_de_salud_mental	Nivel_de_estrés	Horas_de_s
	23.0	Mujer	6.57	6.00	0.68	12.360000	Buena	Вајо	
	21.0	Hombre	3.01	2.57	3.74	7.610000	Mala	Alto	
	51.0	Hombre	3.04	6.14	1.26	3.160000	Neutral	Alto	
	25.0	Mujer	3.84	4.48	2.59	13.080000	Excelente	Medio	
	53.0	Hombre	1.20	0.56	0.29	12.630000	Buena	Bajo	
12522	62.0	Otro	9.69	3.28	4.98	2.930000	Neutral	Bajo	
12523	25.0	Otro	3.82	3.71	3.61	6.570000	Mala	Bajo	
12524	24.0	Mujer	5.08	1.58	1.13	6.580000	Mala	Bajo	
12527	35.0	Mujer	4.28	0.90	4.90	7.965271	Neutral	Alto	
12528	50.0	Mujer	11.85	0.40	1.70	6.210000	Buena	Alto	

```
df3['Género'].fillna("Otro", inplace=True)
df3['Estado_de_salud_mental'].fillna("Neutral", inplace=True)
df3['Nivel_de_estrés'].fillna("Nedio", inplace=True)
df3['Impacto_en_el_entorno_laboral'].fillna("Neutral", inplace=True)
[20] V 0.0s
··· Edad
      Horas_de_uso _de_la_tecnología
Horas_de_uso_de_las_redes_sociales
       Horas_de_juego
Horas_de_tiempo_en_pantalla
Estado_de_salud_mental
       Nivel_de_estrés
Horas_de_sueño
      Horas_de_actividad_física
Acceso_a_sistemas_de_soporte
Impacto_en_el_entorno_laboral
dtype: int64
            #Elimino los NaN restantes, son de la columna 'Acceso a sistemas de soporte' y no pueden ser rellenados porque solo hay datos de sí y no. dfa-df3.dropna()
[22] 		/ 0.0s
                                                                                                                                                                                                                                                                                                 Python
                                                                                                                                                                                                                                                                                                 Python
       Género
       Horas_de_uso _de_la_tecnología
Horas_de_uso_de_las_redes_sociales
       Horas_de_juego
Horas_de_tiempo_en_pantalla
Estado_de_salud_mental
Nivel_de_estrés
        Horas_de_sueño
       Horas_de_actividad_física
Acceso_a_sistemas_de_soporte
       Impacto_en_el_entorno_laboral dtype: int64
             df4['Edad']=df4['Edad'].astype(int)
                                                                                                                                                                                                                                                                                                 Python
       RangeIndex: 9929 entries, 0 to 9928
Data columns (total 12 columns):
# Column
                                                                               Non-Null Count Dtype
         0 Edad 9929 non-null int64
1 Género 9929 non-null object
2 Horas_de_uso_de_la_tecnología 9929 non-null float64
3 Horas_de_uso_de_las_redes_sociales 9929 non-null float64
         8 Horas_de_sueno
9 Horas_de_actividad_fisica
10 Acceso_a_sistemas_de_soporte
11 Impacto_en_el_entorno_laboral
dtypes: float64(6), int64(1), object(5)
memory usage: 931.0+ KB
              #Guardar los resultados en un csv
df4.to_csv("Base_limpia proyecto.csv", index=True)
                                                                                                                                                                                                                                                                                                  Python
```