

## Librerías

```
In [6]: 1 #Graficos
        2 import matplotlib.pyplot as plt
        3 import seaborn as sns
        4 import plotly.express as px
        5 #Dataframes
        6 import numpy as np
        7 import pandas as pd
        8 #Verifica valores perdidos
```

```
In [7]:
```

```
In [8]:
```

## Lectura / importación de datos

```
In [9]:
```

```
In [10]:
```

In [11]:

Out[11]:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type
0	9046	Male	67.0	0	1	Yes	Private	Urban
1	51676	Female	61.0	0	0	Yes	Self-employed	Rural
2	31112	Male	80.0	0	1	Yes	Private	Rural
3	60182	Female	49.0	0	0	Yes	Private	Urban
4	1665	Female	79.0	1	0	Yes	Self-employed	Rural
5	56669	Male	81.0	0	0	Yes	Private	Urban
6	53882	Male	74.0	1	1	Yes	Private	Rural
7	10434	Female	69.0	0	0	No	Private	Urban
8	27419	Female	59.0	0	0	Yes	Private	Rural
9	60491	Female	78.0	0	0	Yes	Private	Urban
10	12109	Female	81.0	1	0	Yes	Private	Rural
11	12095	Female	61.0	0	1	Yes	Govt_job	Rural
12	12175	Female	54.0	0	0	Yes	Private	Urban
13	8213	Male	78.0	0	1	Yes	Private	Urban
14	5317	Female	79.0	0	1	Yes	Private	Urban
15	58202	Female	50.0	1	0	Yes	Self-employed	Rural
16	56112	Male	64.0	0	1	Yes	Private	Urban
17	34120	Male	75.0	1	0	Yes	Private	Urban
18	27458	Female	60.0	0	0	No	Private	Urban
19	25226	Male	57.0	0	1	No	Govt_job	Urban

In [12]:

Out[12]:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type
5100	68398	Male	82.0	1	0	Yes	Self-employed	Rural
5101	36901	Female	45.0	0	0	Yes	Private	Urban
5102	45010	Female	57.0	0	0	Yes	Private	Rural
5103	22127	Female	18.0	0	0	No	Private	Urban
5104	14180	Female	13.0	0	0	No	children	Rural
5105	18234	Female	80.0	1	0	Yes	Private	Urban
5106	44873	Female	81.0	0	0	Yes	Self-employed	Urban
5107	19723	Female	35.0	0	0	Yes	Self-employed	Rural
5108	37544	Male	51.0	0	0	Yes	Private	Rural
5109	44679	Female	44.0	0	0	Yes	Govt_job	Urban

In [13]:

```

1 #El método "describe" devuelve información estadística de los datos del dataframe.
2 #(de hecho, este método devuelve un dataframe).
3 #Esta información incluye el número de muestras, el valor medio,
4 #la desviación estándar, el valor mínimo, máximo, la mediana y los valores

```

Out[13]:

	id	age	hypertension	heart_disease	avg_glucose_level	bmi
count	5110.000000	5110.000000	5110.000000	5110.000000	5110.000000	4909.000000
mean	36517.829354	43.226614	0.097456	0.054012	106.147677	28.893237
std	21161.721625	22.612647	0.296607	0.226063	45.283560	7.854067
min	67.000000	0.080000	0.000000	0.000000	55.120000	10.300000
25%	17741.250000	25.000000	0.000000	0.000000	77.245000	23.500000
50%	36932.000000	45.000000	0.000000	0.000000	91.885000	28.100000
75%	54682.000000	61.000000	0.000000	0.000000	114.090000	33.100000
max	72940.000000	82.000000	1.000000	1.000000	271.740000	97.600000

In [14]:

```

1 # Filas y columnas dataframe

```

Out[14]: (5110, 12)

In [15]:

```

1 # Acceder al índice

```

Out[15]: RangeIndex(start=0, stop=5110, step=1)

In [16]: `1 # Columnas del data frame`

Out[16]: `Index(['id', 'gender', 'age', 'hypertension', 'heart_disease', 'ever_married',  
'work_type', 'Residence_type', 'avg_glucose_level', 'bmi',  
'smoking_status', 'stroke'],  
dtype='object')`

In [17]: `1 # función Longitud para saber el número de columnas`

Out[17]: 12

In [18]: `1 # Acceder a una columnas llamada gender`

Out[18]: `0 Male  
1 Female  
2 Male  
3 Female  
4 Female  
5 Male  
6 Male  
7 Female  
8 Female  
9 Female  
10 Female  
11 Female  
12 Female  
13 Male  
14 Female  
15 Female  
16 Male  
17 Male  
18 Female  
19 Male`

## Limpiar datos

In [19]:

Out[19]:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type
0	9046	Male	67.0	0	1	Yes	Private	Urban
1	51676	Female	61.0	0	0	Yes	Self-employed	Rural
2	31112	Male	80.0	0	1	Yes	Private	Rural
3	60182	Female	49.0	0	0	Yes	Private	Urban
4	1665	Female	79.0	1	0	Yes	Self-employed	Rural

In [20]: 1 # Tipos de variables de dataframe:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5110 entries, 0 to 5109
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    5110 non-null   int64
1   gender                5110 non-null   object
2   age                   5110 non-null   float64
3   hypertension          5110 non-null   int64
4   heart_disease         5110 non-null   int64
5   ever_married          5110 non-null   object
6   work_type             5110 non-null   object
7   Residence_type        5110 non-null   object
8   avg_glucose_level     5110 non-null   float64
9   bmi                   4909 non-null   float64
10  smoking_status        5110 non-null   object
11  stroke                5110 non-null   int64
dtypes: float64(3), int64(4), object(5)
memory usage: 479.2+ KB
```

In [21]: 1 # Cambiar datatype:  
2 # cambiar una columna a objeto  
3 data["id"] = data["id"].astype("object")  
4 data["hypertension"] = data["hypertension"].astype("object")  
5 data["heart\_disease"] = data["heart\_disease"].astype("object")  
6 data["stroke"] = data["stroke"].astype("object")  
7 # cambiar una columna a entero

In [22]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5110 entries, 0 to 5109
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    5110 non-null   object
1   gender                5110 non-null   object
2   age                   5110 non-null   int64
3   hypertension          5110 non-null   object
4   heart_disease         5110 non-null   object
5   ever_married          5110 non-null   object
6   work_type             5110 non-null   object
7   Residence_type        5110 non-null   object
8   avg_glucose_level     5110 non-null   float64
9   bmi                   4909 non-null   float64
10  smoking_status        5110 non-null   object
11  stroke                5110 non-null   object
dtypes: float64(2), int64(1), object(9)
memory usage: 479.2+ KB
```

In [23]:

Out[23]:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type	avg_
0	9046	Male	67	0	1	Yes	Private	Urban	
1	51676	Female	61	0	0	Yes	Self-employed	Rural	
2	31112	Male	80	0	1	Yes	Private	Rural	
3	60182	Female	49	0	0	Yes	Private	Urban	
4	1665	Female	79	1	0	Yes	Self-employed	Rural	

In [24]:

```

1 # Ejemplo de cómo poner una columna como índice:
2 data_ejemplo = data.copy()
3 data_ejemplo = data_ejemplo.set_index('id')

```

C:\Users\Alumne\_mati1\anaconda3\lib\site-packages\pandas\core\indexes\base.py:6982: FutureWarning: In a future version, the Index constructor will not infer numeric dtypes when passed object-dtype sequences (matching Series behavior)

```

return Index(sequences[0], name=names)

```

Out[24]:

	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type	avg_
id								
9046	Male	67	0	1	Yes	Private	Urban	
51676	Female	61	0	0	Yes	Self-employed	Rural	
31112	Male	80	0	1	Yes	Private	Rural	
60182	Female	49	0	0	Yes	Private	Urban	
1665	Female	79	1	0	Yes	Self-employed	Rural	

## Buscar y corregir errores

### Revisar datos anomalos

In [25]:

```

1 # identificar valores máximos y mínimos:
2 data["age"].min()

```

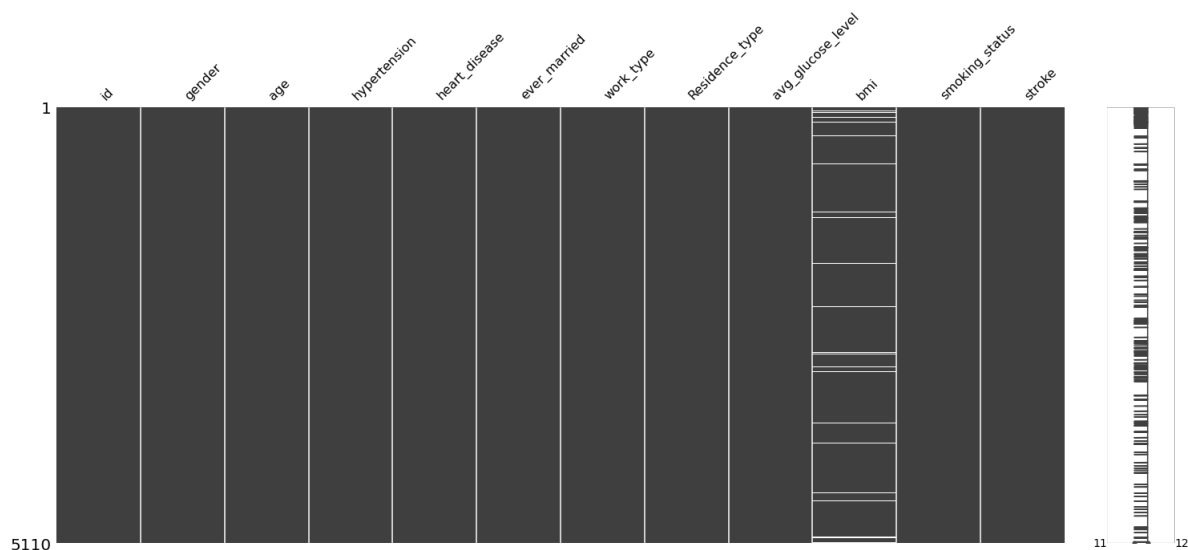
Out[25]: 82

```
In [26]: 1 # identificar outliers(valores anomalos que se salen mucho del rango):  
2 fig = px.box(data, y="avg_glucose_level")  
3
```

## Analizar nans(valores perdidos)

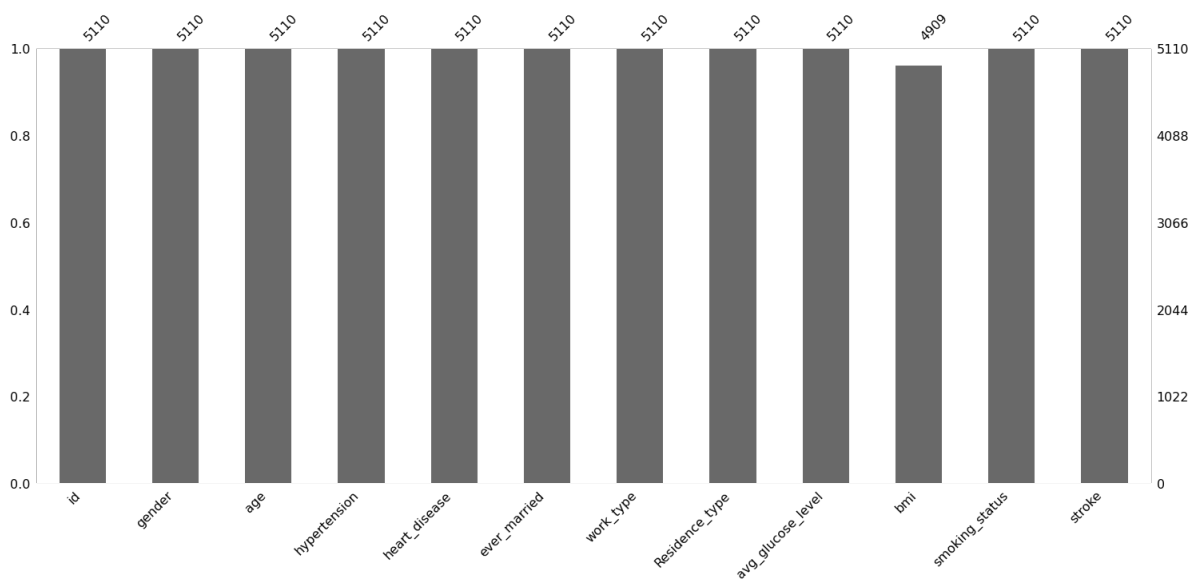
In [26]:

```
1 # Matriz de los nans:
2 msno.matrix(data)
```



In [27]:

```
1 # Barplot de los nans
2 msno.bar(data)
```

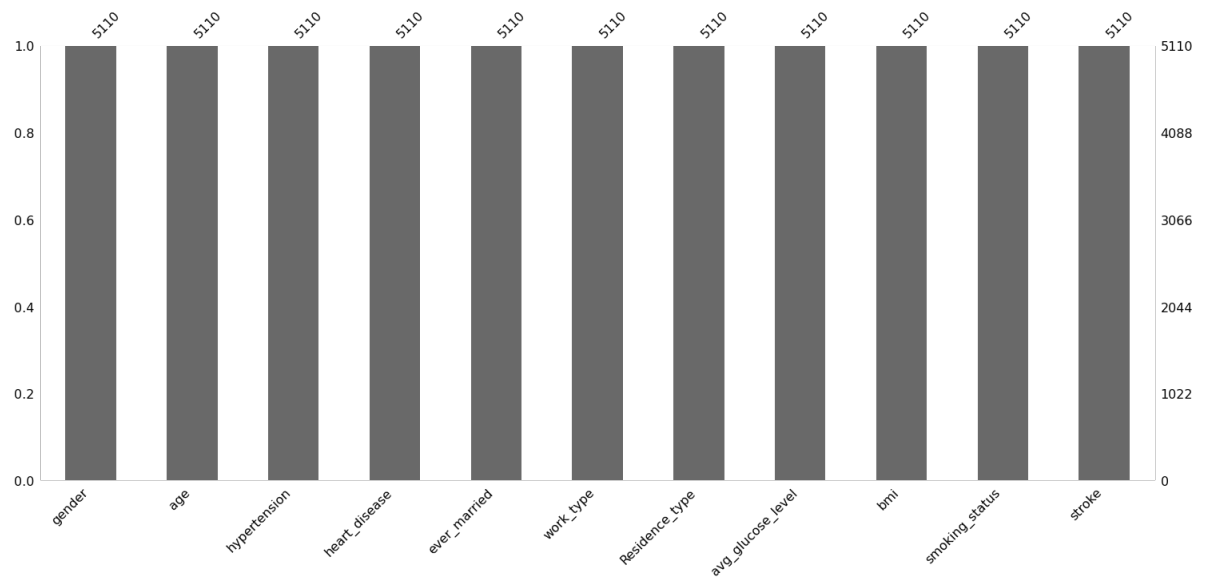


In [28]:

```
1 #Tratar los nans:
2 # Opción 1 - eliminar los nans - dropna
```



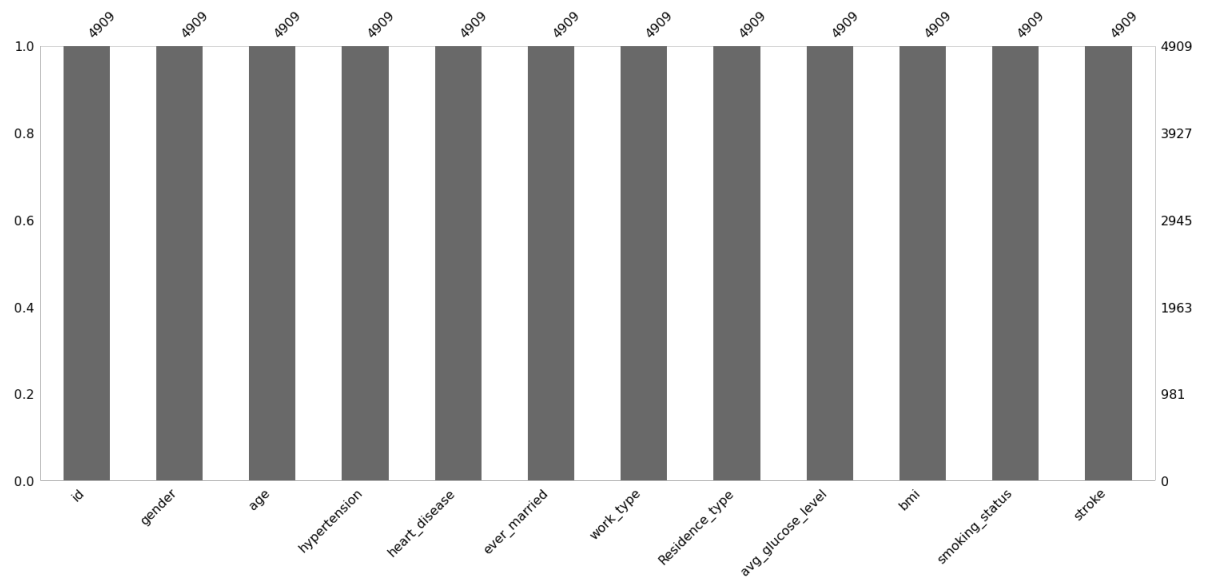
```
In [29]: 1 # Ejemplo de imputación con la media:
2 data_ejemplo["bmi"] = data_ejemplo["bmi"].fillna(data_ejemplo["bmi"].mean())
3 msno.bar(data_ejemplo)
```



```
In [30]: 1 # Eliminar los nans:
2 data = data.dropna()
```

Out[30]: (4909, 12)

```
In [31]: 1 msno.bar(data)
```



```
In [32]: 1 # Criterios para eliminar:
2 # variable cuantitativa con datos iguales en casi todas las filas
3 # variable que tenga el 85% de nans
```

```
In [33]: 1 columnas_borrar = ["id"]
        2 data = data.drop(columnas_borrar,axis = 1)
```

```
Out[33]:
```

	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type	avg_gluc
0	Male	67	0	1	Yes	Private	Urban	
2	Male	80	0	1	Yes	Private	Rural	
3	Female	49	0	0	Yes	Private	Urban	
4	Female	79	1	0	Yes	Self-employed	Rural	
5	Male	81	0	0	Yes	Private	Urban	

## Análisis de los datos

### Analisis univariado variables cuantitativas

```
In [35]:
```

```
Out[35]: Index(['gender', 'age', 'hypertension', 'heart_disease', 'ever_married',
               'work_type', 'Residence_type', 'avg_glucose_level', 'bmi',
               'smoking_status', 'stroke'],
              dtype='object')
```

```
In [36]:
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4909 entries, 0 to 5109
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  -
0   gender                 4909 non-null   object
1   age                   4909 non-null   int64
2   hypertension           4909 non-null   object
3   heart_disease          4909 non-null   object
4   ever_married           4909 non-null   object
5   work_type              4909 non-null   object
6   Residence_type         4909 non-null   object
7   avg_glucose_level      4909 non-null   float64
8   bmi                   4909 non-null   float64
9   smoking_status         4909 non-null   object
10  stroke                 4909 non-null   object
dtypes: float64(2), int64(1), object(8)
memory usage: 460.2+ KB
```

```
In [37]:
```

```
1 # Boxplot
2 # Density plot - histograma
3 # Mean plot - error plot
4 # Stripchart
```

In [38]:

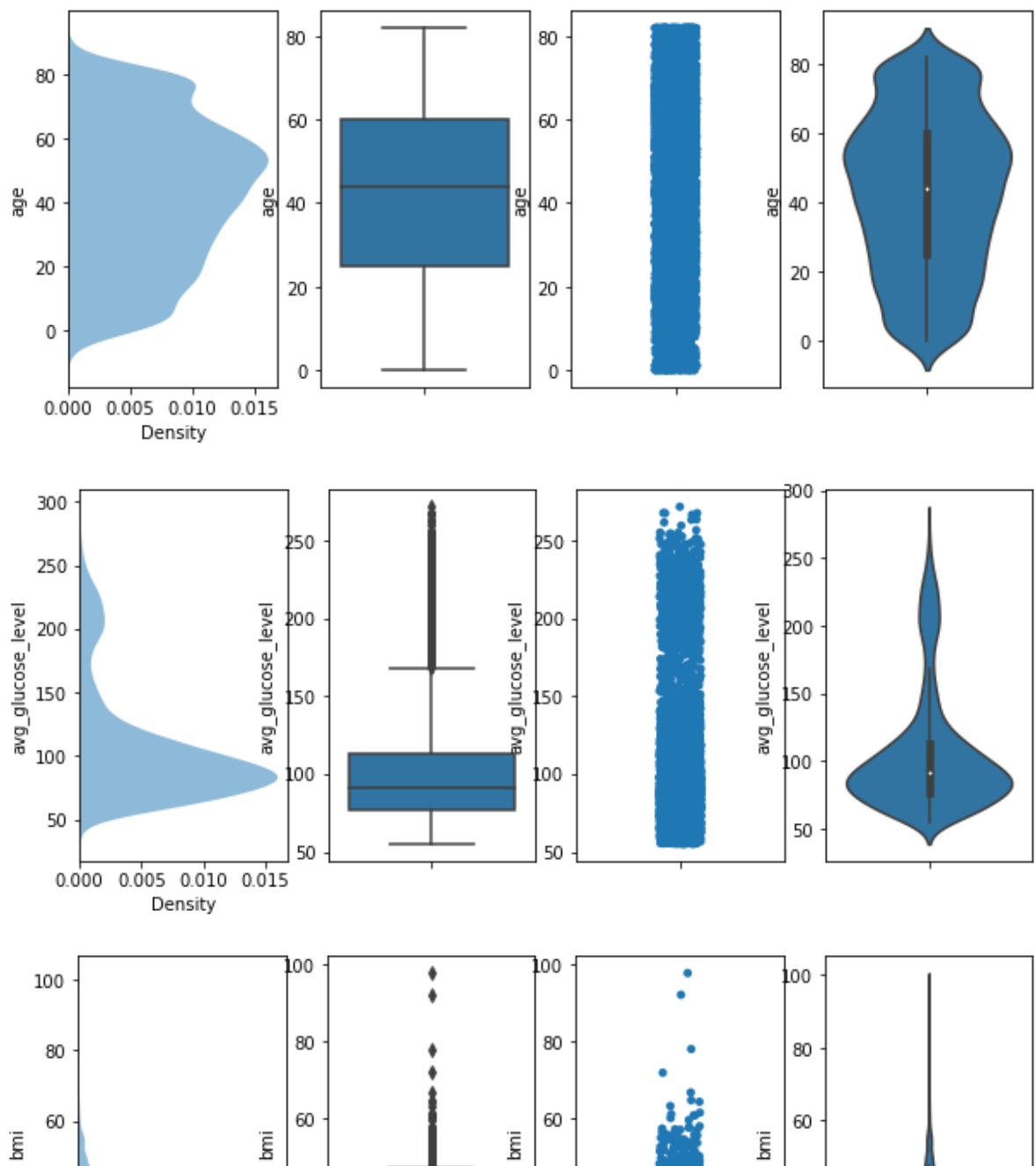
```
...
```

In [39]:

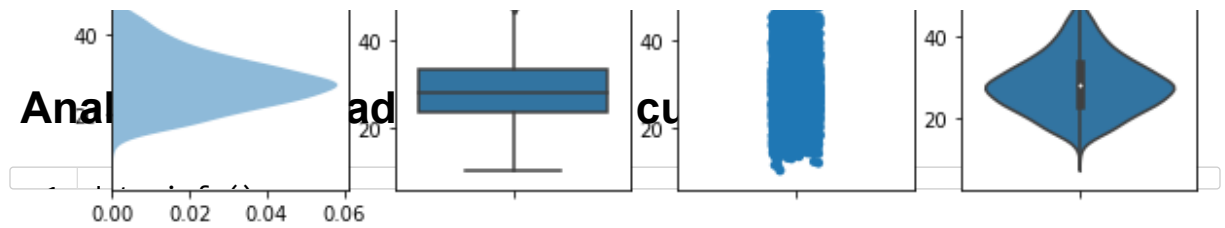
```

1  # Descripcion de todas las variables cuantitativas
2
3  for i in var_num:
4      name_var = i
5      fig, axs = plt.subplots(ncols=4, figsize=(10, 4))
6      # Histograma de densidad
7      sns.kdeplot(
8          data=data, y=name_var,
9          fill=True, common_norm=False, palette="crest",
10         alpha=.5, linewidth=0, ax=axs[0])
11     # Boxplot
12     sns.boxplot(data=data, y=name_var, ax=axs[1])
13     # Stripchart
14     sns.stripplot(y= name_var, data=data, ax=axs[2])
15     # Violinplot
16     sns.violinplot(data=data, y=name_var, ax=axs[3])

```



In [40]:



&lt;class 'pandas.core.frame.DataFrame'&gt;

Int64Index: 4909 entries, 0 to 5109

Data columns (total 11 columns):

#	Column	Non-Null Count	Dtype
0	gender	4909 non-null	object
1	age	4909 non-null	int64
2	hypertension	4909 non-null	object
3	heart_disease	4909 non-null	object
4	ever_married	4909 non-null	object
5	work_type	4909 non-null	object
6	Residence_type	4909 non-null	object
7	avg_glucose_level	4909 non-null	float64
8	bmi	4909 non-null	float64
9	smoking_status	4909 non-null	object
10	stroke	4909 non-null	object

dtypes: float64(2), int64(1), object(8)

memory usage: 589.3+ KB

In [41]:

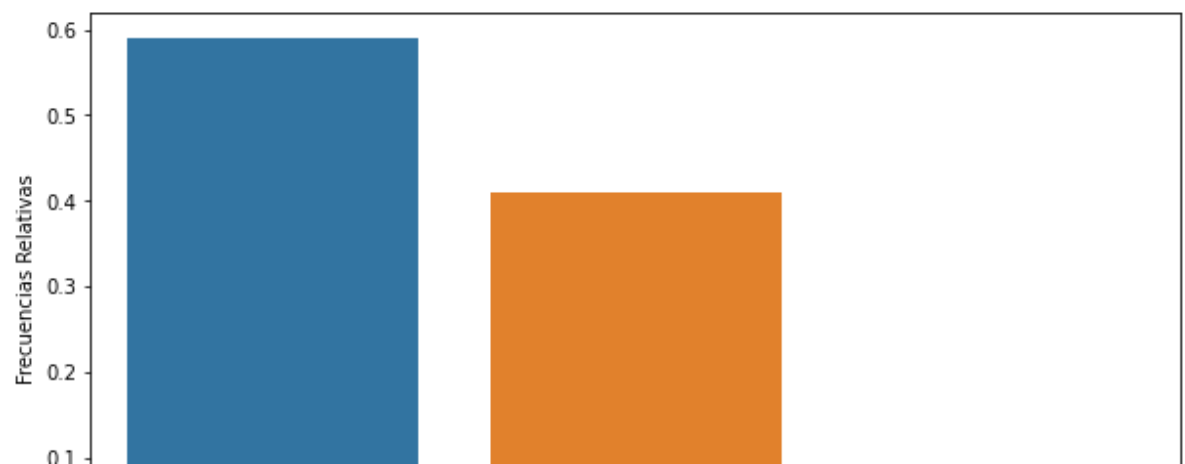
```

1  # gráficos y tablas de frecuencias de todas las variables cualitativas:
2  #La frecuencia relativa es el cociente de la frecuencia absoluta y el número
3  #Para calcular la frecuencia relativa de cada dato dividimos
4  #La frecuencia absoluta del dato entre el número total de datos.
5
6
7  var_cual = ['gender', 'hypertension', 'heart_disease', 'ever_married',
8             'work_type', 'Residence_type',
9             'smoking_status', 'stroke']
10
11 for i in var_cual:
12
13     print("Variable = "+i)
14
15     # 1. Tabla
16     pct = pd.DataFrame(data[i].value_counts(normalize=True))
17     pct.columns = ['Frecuencias Relativas']
18     pct["Frecuencias Absolutas"] = data[i].value_counts()
19     print(pct)
20
21     # 2. Barplot
22     sns.barplot(x=pct.index, y='Frecuencias Relativas', data=pct)
23     plt.show()
24
25     # 3. Pieplot
26     frec_abs = data[i].value_counts() # grabo las frecuencias absolutas
27     labels = data[i].value_counts().index # cojo los nombres de los grupos
28     colors = sns.color_palette('pastel')[0:len(labels)] # defino los colores
29     plt.pie(frec_abs, labels = labels, colors = colors, autopct='%0.0f%%')

```

Variable = gender

	Frecuencias Relativas	Frecuencias Absolutas
Female	0.590141	2897
Male	0.409656	2011
Other	0.000204	1

**Diagramas de dispersión variables cuantitativas**

```
In [42]: 1 #Los gráficos de dispersión se utilizan para mostrar relaciones.  
2 #Para la correlación, Los gráficos de dispersión ayudan a mostrar la fuerz  
3 #Para la regresión, Los gráficos de dispersión suelen incorporar una línea  
4 #En control de calidad, Los gráficos de dispersión pueden con frecuencia i
```

```
In [43]:
```

```
Out[43]: Index(['gender', 'age', 'hypertension', 'heart_disease', 'ever_married',  
               'work_type', 'Residence_type', 'avg_glucose_level', 'bmi',  
               'smoking_status', 'stroke'],  
            dtype='object')
```

### Scatter 3D

```
In [27]: 1 # Diagramas de dispersión 3D de variables cuantitativas
          2
          3 z = 'avg_glucose_level'
          4 x = 'bmi'
          5 y = 'age'
          6
          7 fig = px.scatter_3d(data,x=x,y=y,z=z,size_max=0.1)
          8 fig.update_traces(marker_size = 1)
```

### Matrixplot

```
In [45]:
```





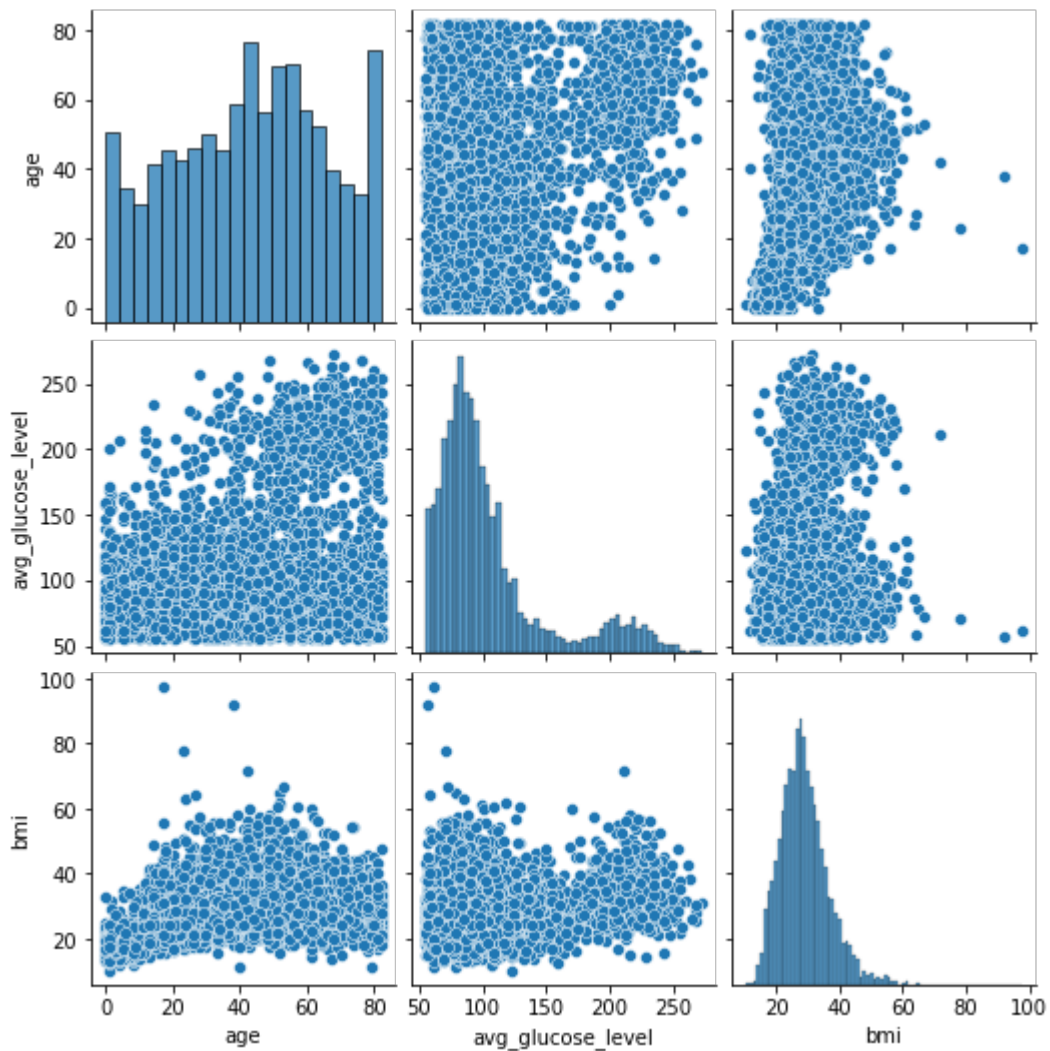
In [46]:

Out[46]:

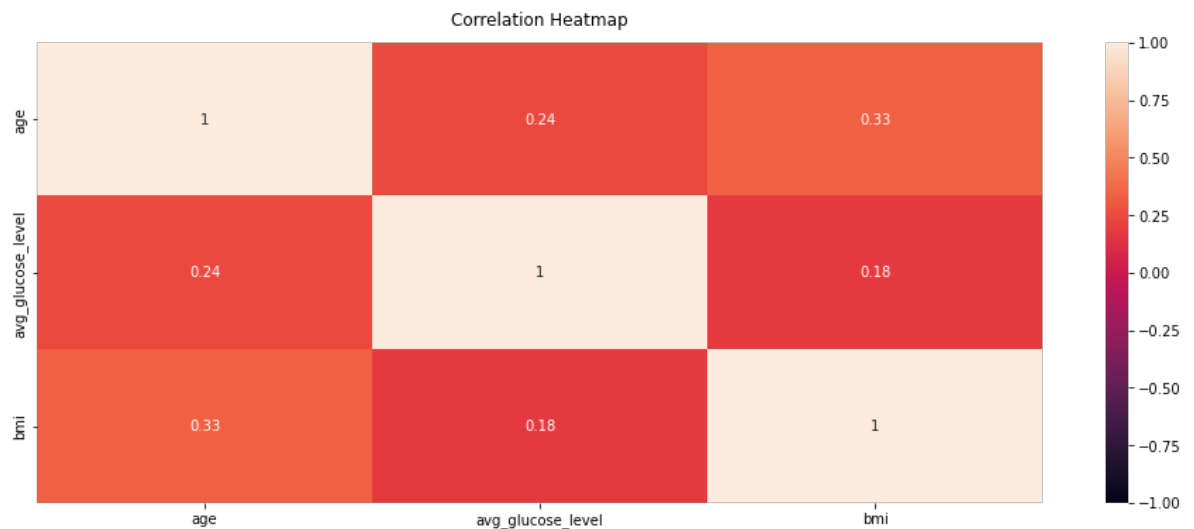
	age	avg_glucose_level	bmi
0	67	228.69	36.6
2	80	105.92	32.5
3	49	171.23	34.4
4	79	174.12	24.0
5	81	186.21	29.0

In [47]:

```
1 # Matrixplot
2
3 columnas_plot = ["age", "avg_glucose_level", "bmi"]
4 # seaborn
5 #sns.pairplot(data[var_num])
6 sns.pairplot(data[columnas_plot])
7 plt.show()
```

**Correlograma**

```
In [49]: 1 # Visualizacion de la correlación serial en datos que cambian con el tiempo
2         #(es decir , datos de series temporales ).
3
4         X = data[var_num]
5
6         plt.figure(figsize=(16, 6))
7         heatmap = sns.heatmap(X.corr(), vmin=-1, vmax=1, annot=True)
8
9         heatmap.set_title('Correlation Heatmap', fontdict={'fontsize':12}, pad=12)
```



## Diagramas variables cualitativas

```
In [50]:
```

```
Out[50]: Index(['gender', 'age', 'hypertension', 'heart_disease', 'ever_married',
               'work_type', 'Residence_type', 'avg_glucose_level', 'bmi',
               'smoking_status', 'stroke'],
              dtype='object')
```

```
In [51]:
```

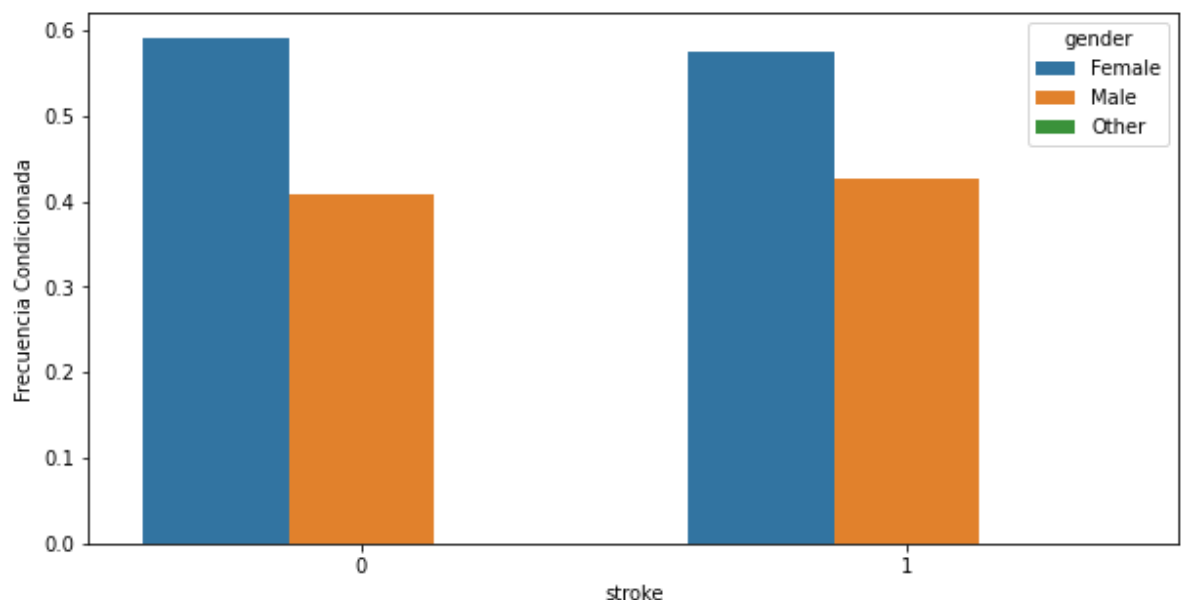
```
Out[51]: ['gender',
          'hypertension',
          'heart_disease',
          'ever_married',
          'work_type',
          'Residence_type',
          'smoking_status',
          'stroke']
```

```

In [52]: 1 # Diagramas de barras de dos factores de frecuencias condicionadas a los g
2 #(Son las frecuencias que posee una variable si sólo consideramos un valor
3 #distribución bidimensional X-Y. En la práctica se traduce a considerar só
4 #según el valor elegido.
5
6 x = 'stroke'
7 vec_y = ['gender',
8 'hypertension',
9 'heart_disease',
10 'ever_married',
11 'work_type',
12 'Residence_type',
13 'smoking_status']
14
15
16 for i in vec_y:
17
18     y = i
19     print(y + ' vs ' + x)
20     pct2 = (data.groupby([x,y]).size() / data.groupby([x]).size()).reset_i
21     # seaborn
22     sns.barplot(x=x, hue=y, y='Frecuencia Condicionada', data=pct2)
23     plt.show()
24     print(pct2)

```

gender vs stroke



```

In [ ]:

```