

# ANÁLISIS EXPLORATORIO DE DATOS

CLASE 16

# ANÁLISIS EXPLORATORIO DE DATOS

EDA	Estadístico	Gráfico
<b>Una variable</b>	Media Mediana Desviación estándar Varianza Percentiles Rango intercuartil (IQR=Q3-Q1) Distribución de probabilidad	Histograma FDA KDE Boxplot Choroplet
<b>Numérica</b>		
<b>Categoría</b>		
<b>Multi-variable</b>	Coef. Pearson Matriz de correlación Regresión Agrupación (groupby, pivot)	Scatterplot Jointplot Pairplots Histograma múltiple Serie de tiempo Stacked area Heatmap Pairplots... etc

**pandas**  
**numpy**

**scipy**  
**sklearn**

**matplotlib, seaborn**  
**geopandas**

# Estadísticas de Resumen

- **Media:** es la suma de todos los valores, dividida por el número de puntos.  $\mu = \frac{1}{n} \sum_{i=1}^n x_i$
- **Mediana:** es el valor medio de un conjunto de datos. Es inmune a valores extremos o outliers. Para calcularla, se ordenan los datos y se elige el valor que queda en la mitad.
- **Percentiles:** el percentil  $p$ , corresponde al valor que es mayor al  $p\%$  de los datos.
- **Varianza:** promedio de la distancia cuadrática de los datos a la media. Es una medida de la dispersión de los datos.

$$\text{Var} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

- **Desviación Estándar:** es la raíz cuadrada de la varianza. Está en la misma escala de unidades que los datos.

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

## Estadísticas de Resumen

- **Covarianza:** es una medida de cómo dos cantidades varían juntas. Es la media del producto entre las diferencias de los valores respecto a la media.

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- Si  $\mathbf{x}$  e  $\mathbf{y}$  tienden a estar ambas arriba, o ambas abajo de la media, la covarianza es positiva.
- Esto quiere decir que hay una correlación positiva: cuando  $\mathbf{x}$  es alta,  $\mathbf{y}$  es alta.
- Por el contrario, si  $\mathbf{x}$  es alta cuando  $\mathbf{y}$  es baja, la covarianza es negativa y los datos están anticorrelacionados.

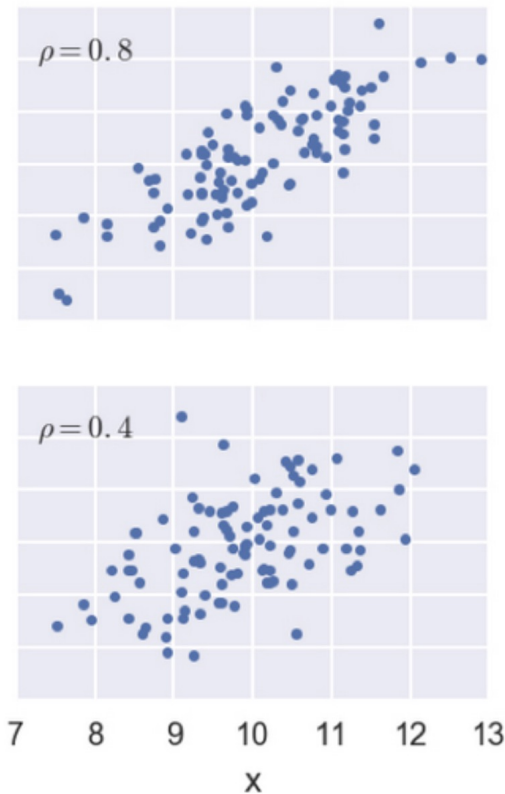
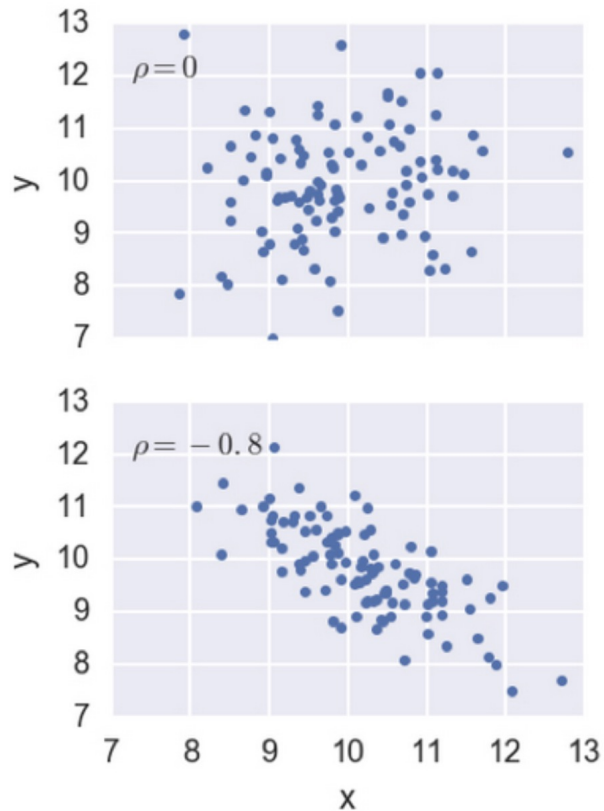
# Estadísticas de Resumen

**Coeficiente de Pearson ( $\rho$ ):** para tener una medida más general y aplicable de la correlación entre dos variables, necesitamos que sea adimensional. Por lo tanto dividimos la covarianza por las desviaciones estándar de  $x$  e  $y$ .

$$\rho = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

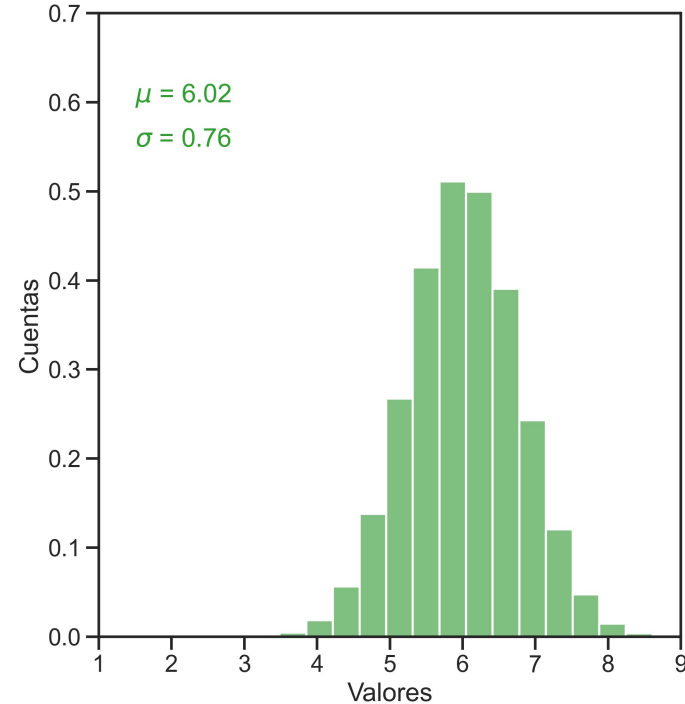
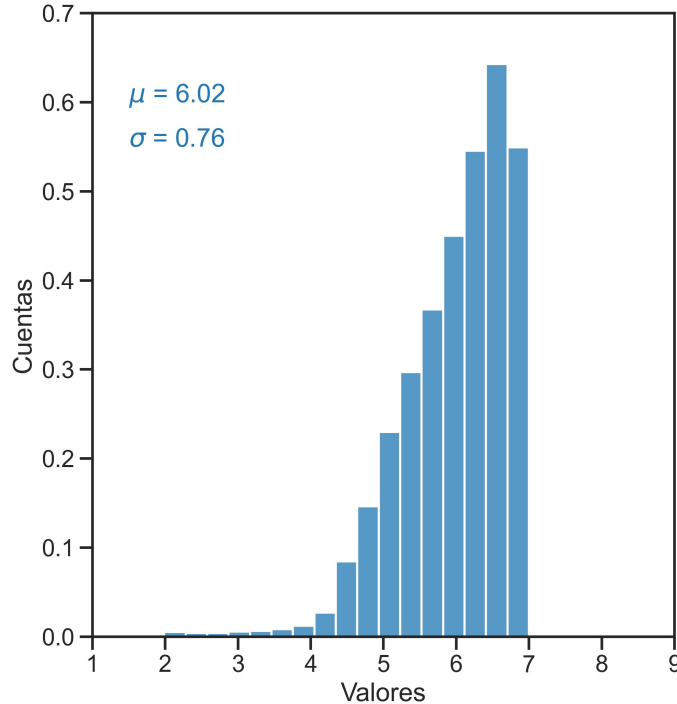
- Es la comparación de la variabilidad en los datos debido a una codependencia (covarianza), con la variabilidad inherente de cada variable (sus desviaciones estándar).
- Un valor **0** indica que **no hay correlación**, valor **-1/1** indica **alta correlación (negativa/positiva)**.

# Estadísticas de Resumen

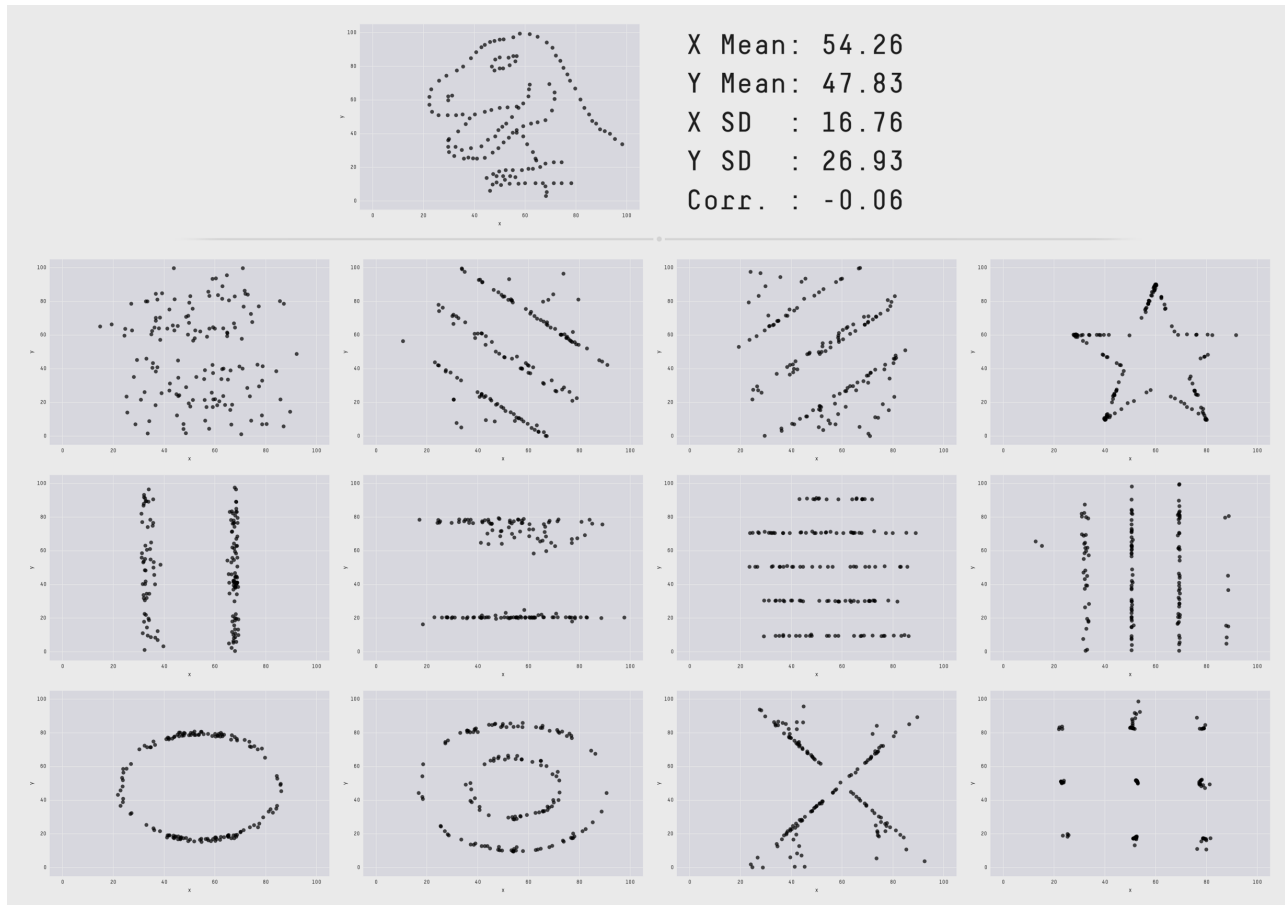


$$\rho = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

# Distintas distribuciones pueden tener las mismas estadísticas de resumen



# Distintas distribuciones pueden tener las mismas estadísticas de resumen

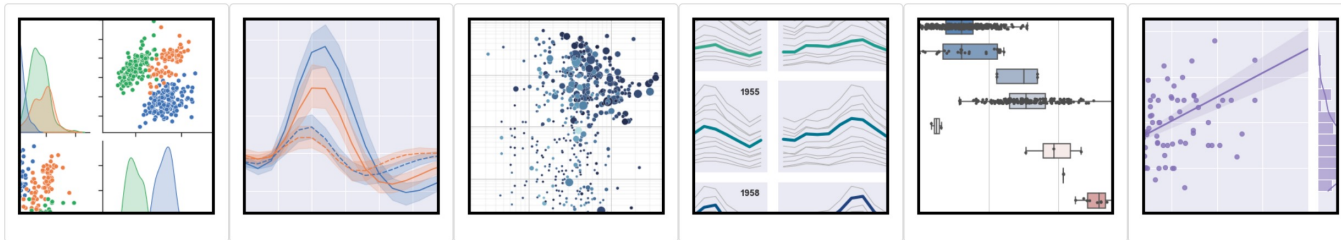


Además de analizar  
las estadísticas de  
resumen, es  
necesario  
**GRAFICAR** los  
datos para entender  
mejor su  
distribución.



# Análisis Gráfico

## seaborn: statistical data visualization

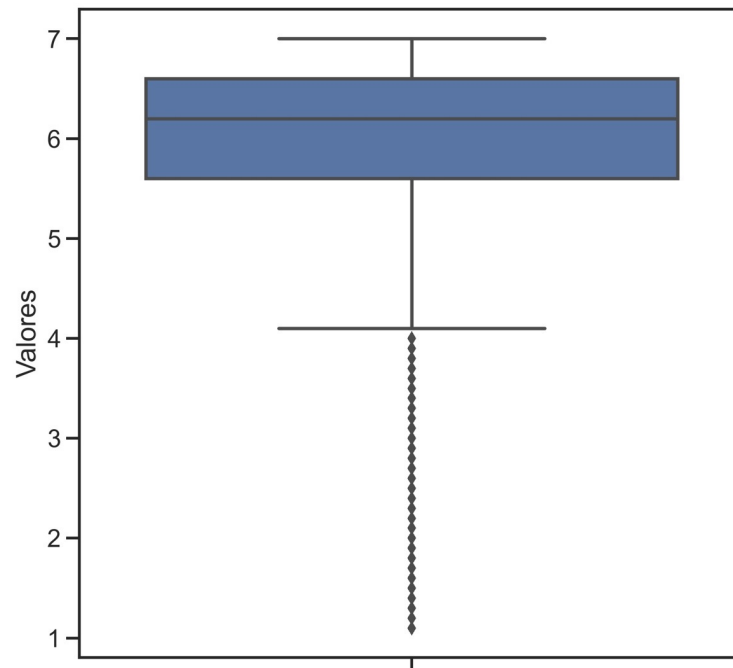
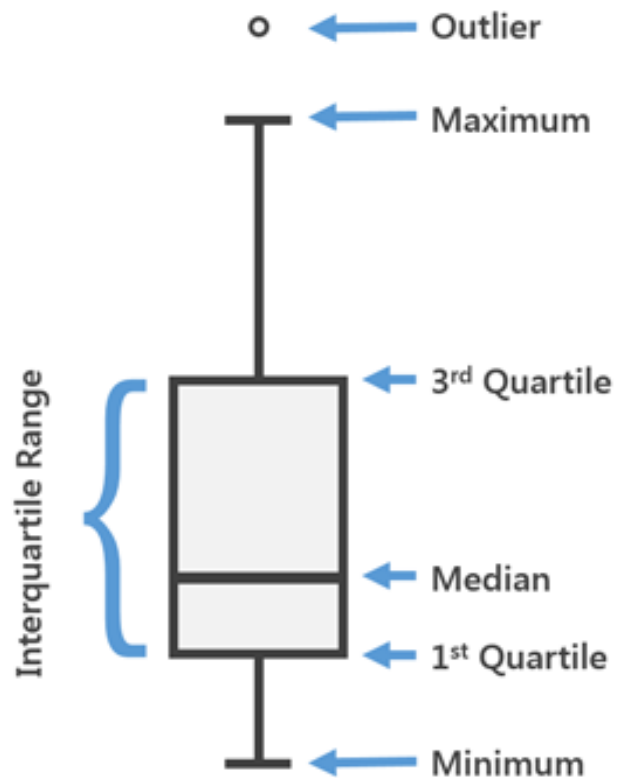


Seaborn is a Python data visualization library based on [matplotlib](#). It provides a high-level interface for drawing attractive and informative statistical graphics.

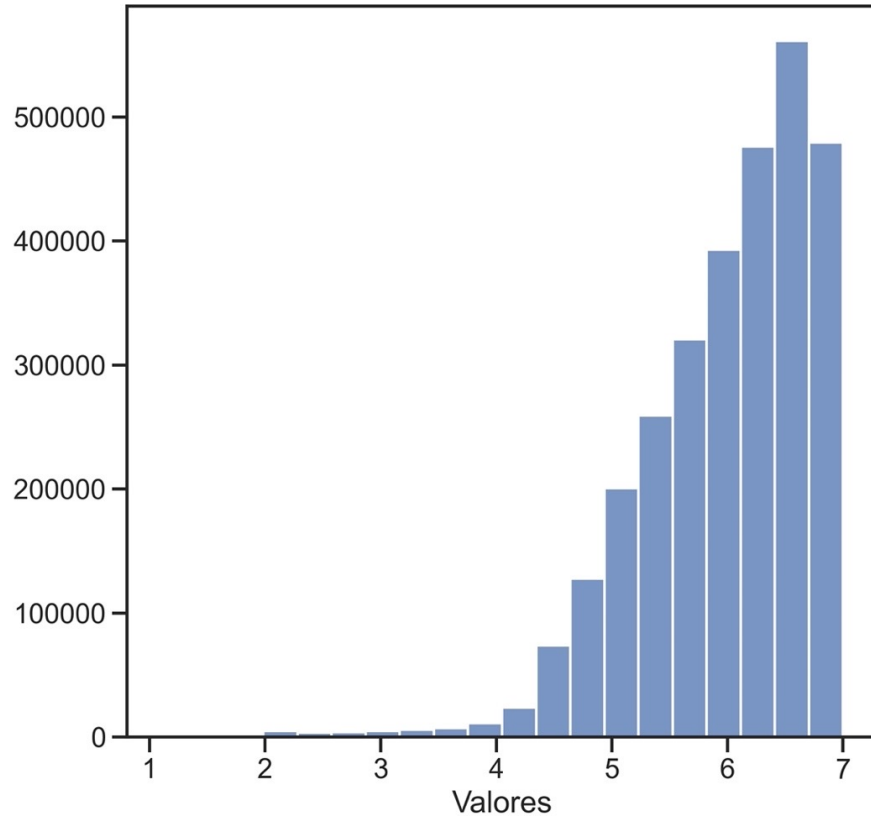
For a brief introduction to the ideas behind the library, you can read the [introductory notes](#). Visit the [installation page](#) to see how you can download the package and get started with it. You can browse the [example gallery](#) to see what you can do with seaborn, and then check out the [tutorial](#) and [API reference](#) to find out how.

To see the code or report a bug, please visit the [GitHub repository](#). General support questions are most at home on [stackoverflow](#) or [discourse](#), which have dedicated channels for seaborn.

# Boxplot

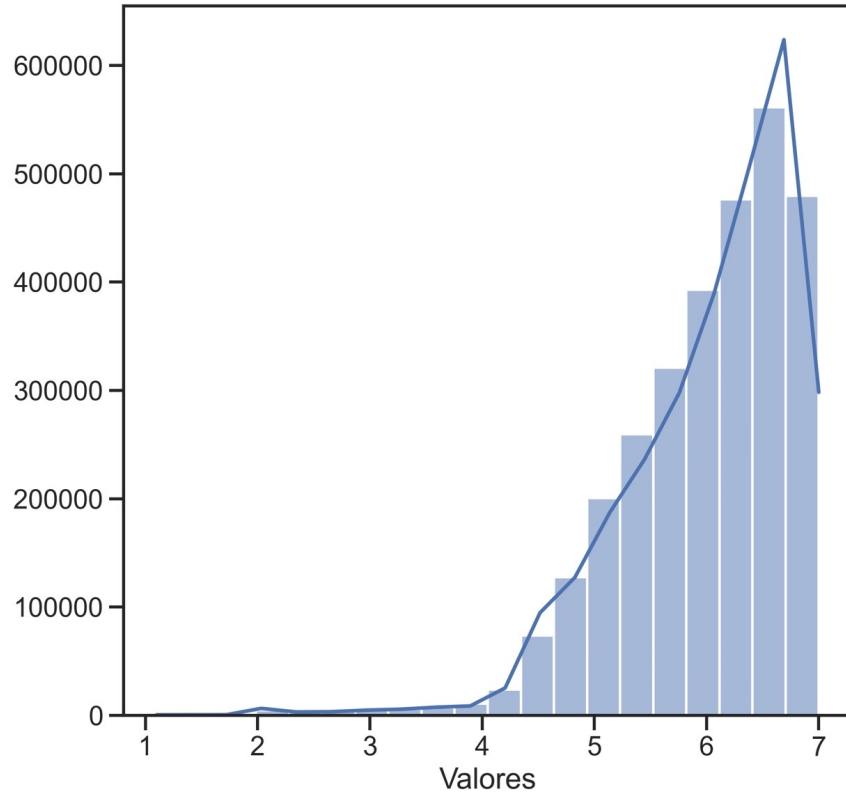


# Histograma y Función de Densidad



```
sns.histplot(data=df,x='Valores',  
ax=ax,bins=20);
```

# Histograma y Función de Densidad



```
sns.histplot(data=df,x='Valores',  
ax=ax,bins=20,kde=True);
```

# Dataset de Ejemplo: Pasajeros del Titanic

- El dataset **titanic.csv** contiene datos para 887 pasajeros del Titanic, cada uno de los cuales representa una fila.
- Las columnas indican atributos de la persona: si sobrevivió, edad, clase, sexo, tarifa pagada, etc.

```
1 titanic = sns.load_dataset('titanic');  
2 titanic.info();
```

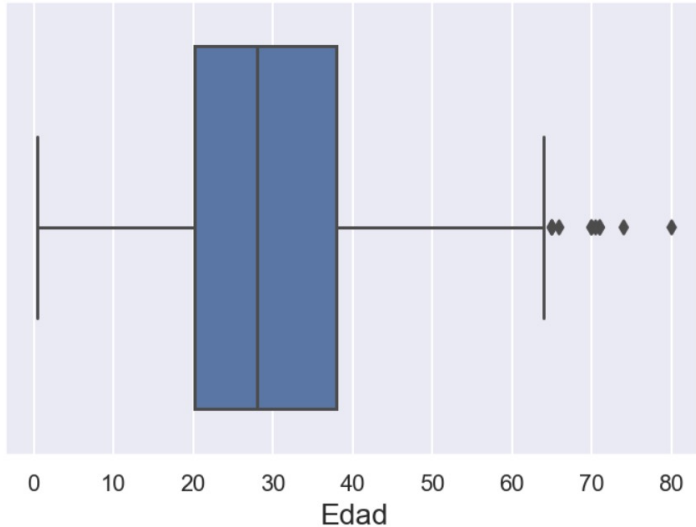
```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 891 entries, 0 to 890  
Data columns (total 15 columns):  
survived      891 non-null int64  
pclass       891 non-null int64  
sex          891 non-null object  
age          714 non-null float64  
sibsp        891 non-null int64  
parch        891 non-null int64  
fare         891 non-null float64  
embarked     889 non-null object  
class        891 non-null category  
who          891 non-null object  
adult_male   891 non-null bool  
deck         203 non-null category  
embark_town   889 non-null object  
alive        891 non-null object  
alone        891 non-null bool  
dtypes: bool(2), category(2), float64(2), int64(4), object(5)  
memory usage: 80.6+ KB
```

# Boxplot

## Una Variable

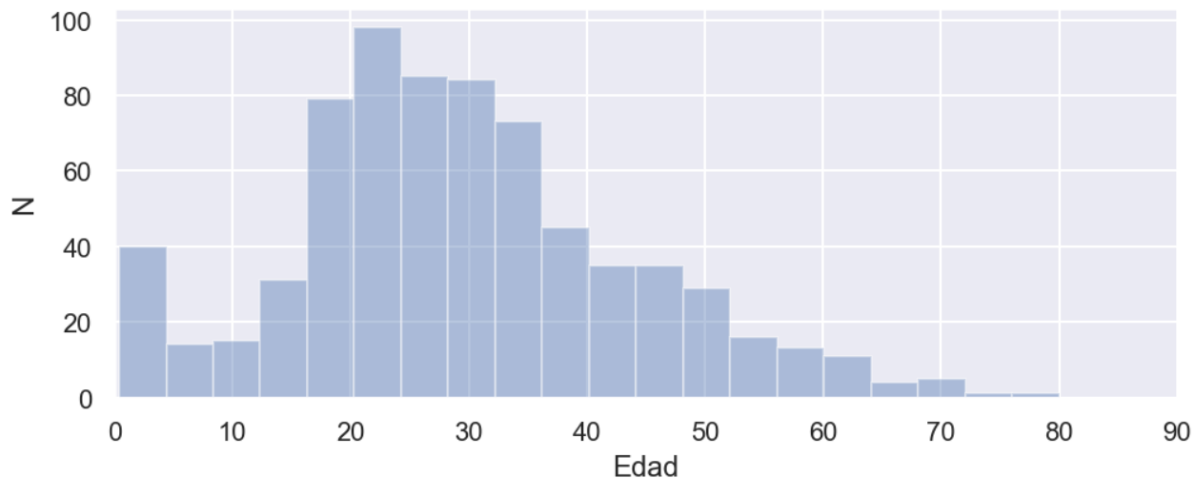
```
1 # seaborn
2 ax = sns.boxplot(x='age', data=titanic)
3 ax.set_ylabel(None);
4 ax.set_xlabel('Edad', fontsize=14);
5 ax.set_title('Distribución de Edad de Pasajeros del Titanic', fontsize=14);
```

Distribución de Edad de Pasajeros del Titanic



# Histograma y Función de Densidad

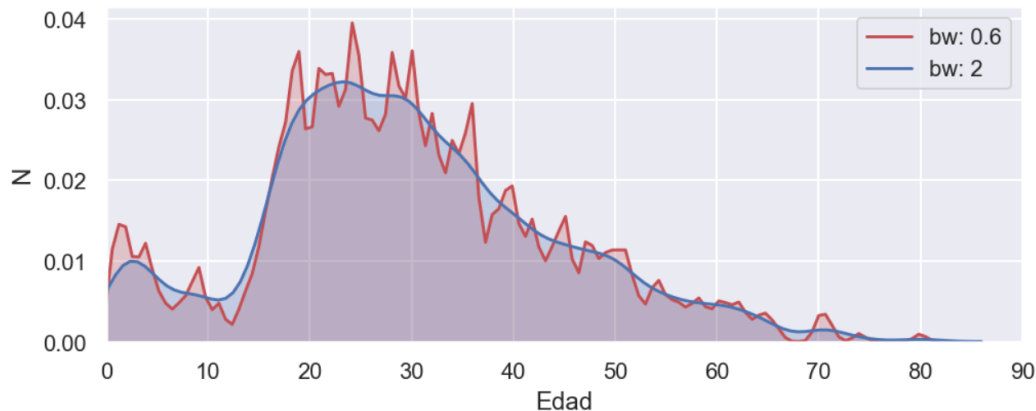
```
1 # ¿Cuál es la distribución de edades de los pasajeros del Titanic?
2 import seaborn as sns
3 sns.set(color_codes=True)
4
5 f, ax = plt.subplots(1,1, figsize=(8, 3));
6 ax = sns.distplot(titanic.age, kde=False, bins=20)
7
8 ax.set(xlim=(0, 90));
9 ax.set_ylabel('N');
10 ax.set_xlabel('Edad');
```





# Histograma y Función de Densidad

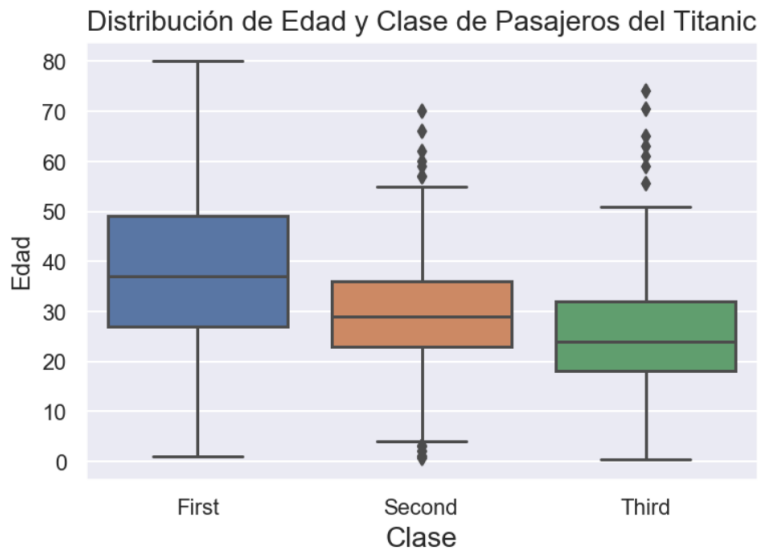
```
1 # ¿Cuál es la distribución de edades de los pasajeros del Titanic?
2 import seaborn as sns
3 sns.set(color_codes=True)
4
5 f, ax = plt.subplots(1,1, figsize=(8, 3));
6
7 sns.kdeplot(titanic.age, bw=0.6, label="bw: 0.6", shade=True, color="r",ax=ax);
8 sns.kdeplot(titanic.age, bw=2, label="bw: 2", shade=True,ax=ax);
9
10 ax.set(xlim=(0, 90));
11 ax.set_ylabel('N');
12 ax.set_xlabel('Edad');
```



# Boxplot

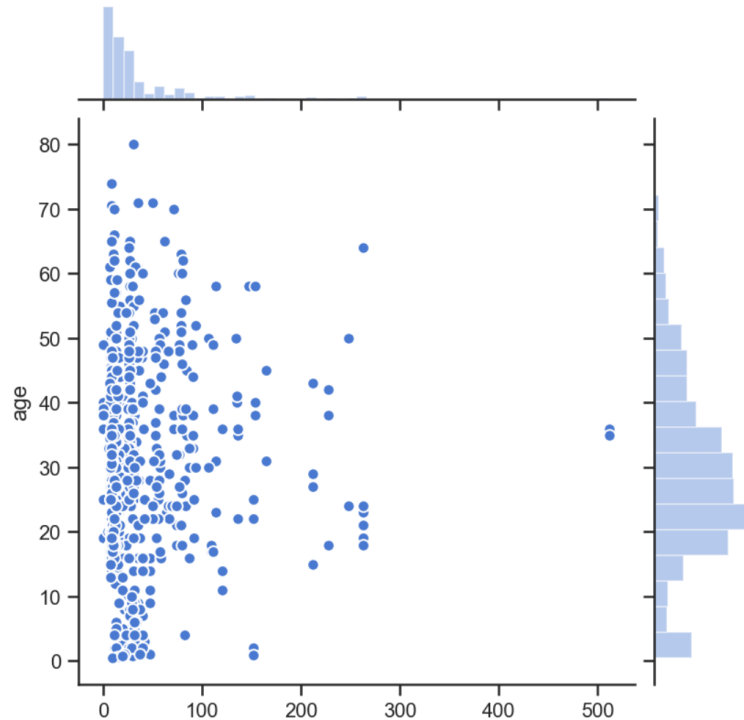
**Dos variables: ¿la primera clase se componía de pasajeros más jóvenes o mayores?**

```
: 1 ax = sns.boxplot(x='class', y='age', data=titanic)
  2 ax.set_ylabel('Edad');
  3 ax.set_xlabel('Clase', fontsize=14);
  4 ax.set_title('Distribución de Edad y Clase de Pasajeros del Titanic', fontsize=14);
```



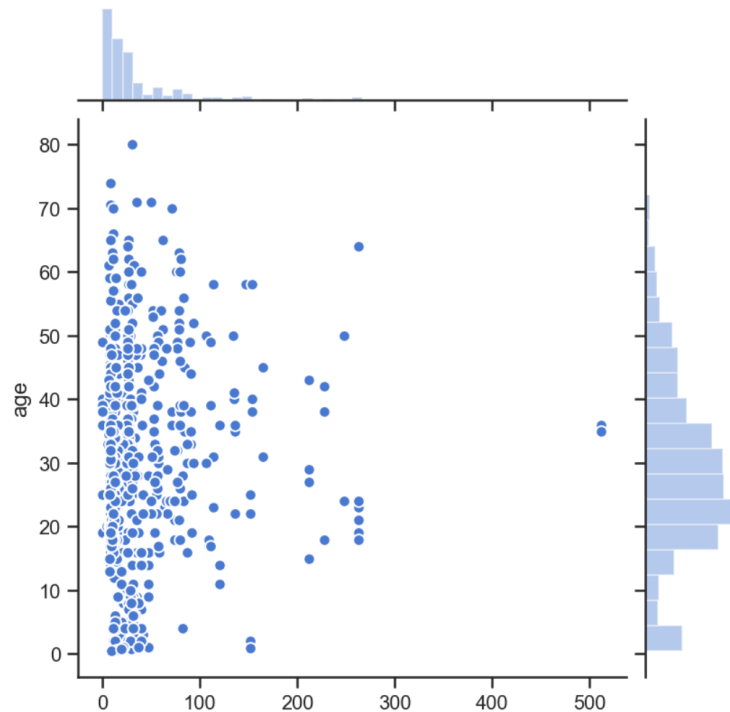
# Scatterplot / Jointplot : Relación entre 2 variables

```
1 sns.jointplot("fare", "age", data=titanic, s=40, edgecolor="w", linewidth=1);  
2 ax.set_ylabel('Tarifa');  
3 ax.set_xlabel('Edad', fontsize=14);  
4 ax.set_title('Relación entre Edad y Tarifa Pagada por Pasajeros del Titanic', fontsize=14);
```



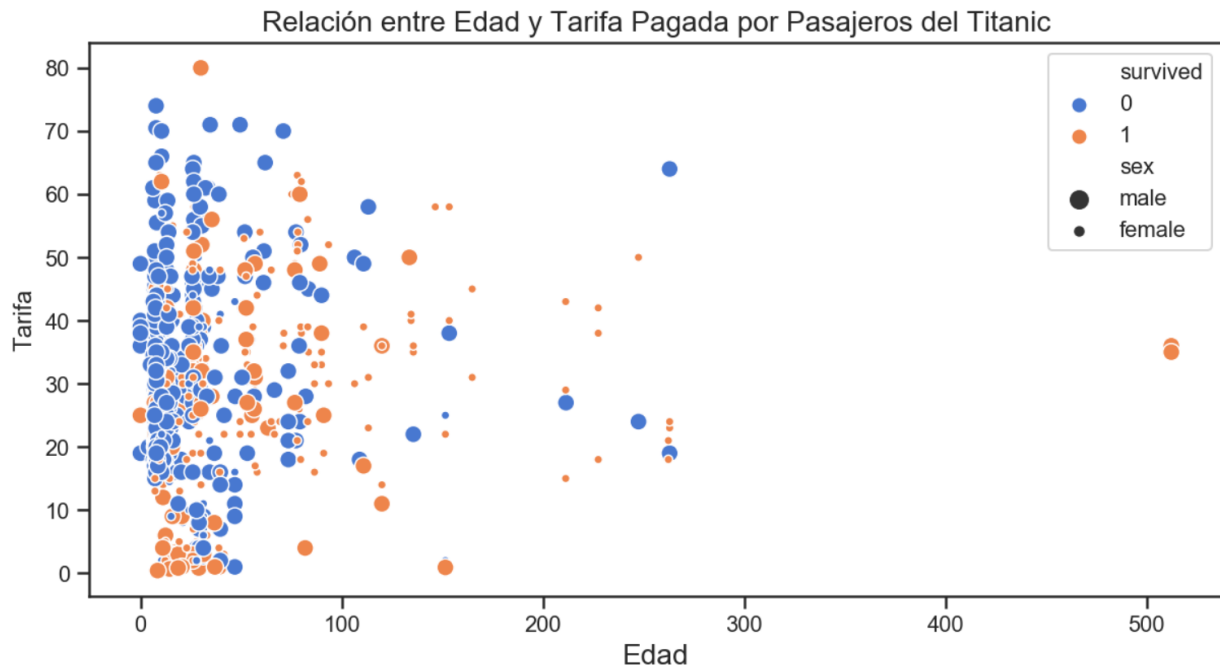
# Scatterplot / Jointplot : Relación entre 2 variables

```
1 sns.jointplot("fare", "age", data=titanic, s=40, edgecolor="w", linewidth=1);  
2 ax.set_ylabel('Tarifa');  
3 ax.set_xlabel('Edad', fontsize=14);  
4 ax.set_title('Relación entre Edad y Tarifa Pagada por Pasajeros del Titanic', fontsize=14);
```



# Scatterplot y Distplot: Relación entre 2 o más variables

```
1 f, ax = plt.subplots(1,1, figsize=(10, 5))
2 sns.scatterplot(x="fare", y="age", hue="survived", size="sex", data=titanic, ax=ax);
3 ax.set_ylabel('Tarifa');
4 ax.set_xlabel('Edad', fontsize=14);
5 ax.set_title('Relación entre Edad, Tarifa Pagada, Sobrevivencia y Sexo', fontsize=14);
6 plt.show();
```



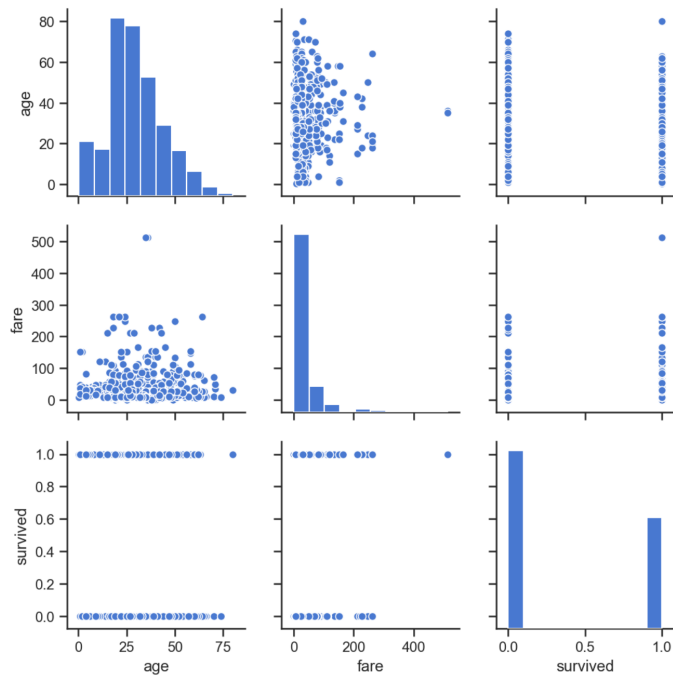
# Pairplot: Análisis multivariable

```
1 titanic.columns
```

```
Index(['survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'class', 'deck'], dtype='object')
```

```
1 graficar = ['age', 'fare', 'survived', 'deck']
```

```
1 titanic_plot = titanic.loc[:,graficar]  
2 sns.pairplot(titanic_plot);
```



# Pairplot: Análisis multivariable

```
: 1 titanic.columns  
  
: Index(['survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'class', 'deck'], dtype  
      ='object')  
  
: 1 graficar = ['age', 'fare', 'survived', 'deck']  
  
: 1 titanic_plot = titanic.loc[:,graficar]  
  2 sns.pairplot(titanic_plot,hue='deck');
```

