

```

1 import numpy as np
2 import pandas as pd
3 from scipy.stats import lognorm
4 import random
5 import matplotlib.pyplot as plt
6

```

Primer metodo para generar los datos sin embargo tiene un pero solo me genera de un solo valor y yo necesito el listado.

```

1 np.random.seed(10000000)
2 porc_enf = np.random.lognormal(mean = 8, sigma= 3.6, size= 100)
3 porc_enf
4 print(porc_enf.min())
5 print(porc_enf.max())

```

```

0.2372200135649288
12273220.050374826

```

```

1 def lognormal(mean,sigma):
2     while True:
3         pathogen = np.random.lognormal(mean=np.log(mean),sigma=np.log(sigma),size=None)
4         if pathogen>=0 and pathogen <= 100:
5             break
6     return pathogen
7 lognormal(0.8, 3.6)

```

```

0.9883166998870656

```

Se genera un inconveniente porque solo nos esta generando de a un solo dato

```

1 i = 0
2 while True:
3     # for first iteration, make an initial vector of random values
4     if i == 0:
5         V = np.round(np.random.lognormal(0.8,sigma=3.6,size=100),2)
6
7     #Make binary vectors for those values that are out of range and in range
8     V_outRange = (V > 100) + 0.
9     V_good = (V < 100) + 0.
10
11     #keep the values in range
12     V_ok = np.multiply(V, V_good)
13
14     #Resample values that are out of range
15     V_next = np.round(np.multiply(V_outRange, np.random.lognormal(0.8, 3.6, 100)),2)

```

```

16
17 #Re-combine previous in-range values with newly sampled values
18 V = V_ok + V_next
19
20 #check to see if all values are in range, if not re-loop
21 if (np.count_nonzero(V > 100)) == 0:
22     break
23 i +=100
24 V

array([3.000e-02, 1.000e-02, 1.709e+01, 3.880e+00, 6.250e+00, 4.400e-01,
       1.810e+00, 8.210e+00, 2.176e+01, 1.200e-01, 2.000e-02, 2.458e+01,
       1.000e-02, 2.700e+00, 8.680e+00, 0.000e+00, 2.000e-02, 4.690e+00,
       9.600e-01, 1.154e+01, 4.900e-01, 5.000e-02, 2.600e-01, 3.370e+00,
       7.200e-01, 3.000e-02, 5.698e+01, 4.700e-01, 6.270e+00, 1.890e+00,
       3.200e-01, 6.000e-02, 2.290e+00, 4.100e-01, 4.000e-02, 4.000e-02,
       4.605e+01, 3.790e+00, 4.000e-02, 4.000e-02, 6.000e-02, 2.693e+01,
       2.000e-02, 1.320e+00, 6.267e+01, 2.223e+01, 1.150e+00, 3.160e+00,
       7.000e-02, 6.800e-01, 3.330e+00, 4.290e+00, 0.000e+00, 5.658e+01,
       5.600e-01, 4.000e-02, 7.000e-02, 0.000e+00, 3.153e+01, 4.440e+00,
       8.000e-02, 4.710e+00, 3.678e+01, 3.240e+00, 0.000e+00, 3.860e+00,
       4.870e+00, 1.000e-02, 0.000e+00, 1.047e+01, 5.830e+00, 1.350e+00,
       4.471e+01, 9.170e+00, 3.900e-01, 0.000e+00, 7.150e+00, 1.330e+00,
       2.634e+01, 8.300e-01, 1.400e+00, 2.550e+00, 6.480e+00, 0.000e+00,
       4.680e+00, 1.112e+01, 1.600e-01, 2.520e+00, 1.069e+01, 3.900e-01,
       1.297e+01, 1.055e+01, 1.116e+01, 3.099e+01, 3.480e+00, 2.000e-02,
       6.019e+01, 8.400e-01, 1.690e+00, 2.710e+00])

```

Fuente: <https://stackoverflow.com/questions/51995490/restricting-numpy-random-lognormal-to-a-given-range-python>

Mostrar los valores maximos y minimos

```

1 print(V.min())
2 print(V.max())

0.0
62.67

```

▼ Categorización de la variable según la escala diagramatica

```

1 cat_V = [] #VECTOR DE CATEGORIAS VACIO, sin embargo categorizar implica perder datos
2 for pe_i in V:#####cambiar este dato***
3     if(pe_i <= 2):
4         cat_V.append('A0')
5     elif(pe_i <= 4):
6         cat_V.append('A1')
7     elif(pe_i <= 6):
8         cat_V.append('A2')
9     elif(pe_i <= 8):
10        cat_V.append('A3')
11    elif(pe_i <= 10):
12        cat_V.append('A4')
13    else:
14        cat_V.append('A5')

```

```

/ elif(pe_i <= 8):
8     cat_V.append('A2')
9 elif(pe_i <= 12):
10    cat_V.append('A3')
11 elif(pe_i <= 27):
12    cat_V.append('A4')
13 elif(pe_i <= 45):
14    cat_V.append('A5')
15 elif(pe_i <= 71):
16    cat_V.append('A6')
17 elif(pe_i <= 93):
18    cat_V.append('A7')
19 else:
20    cat_V.append('A9')
21
22 cat_V_serie = pd.Series(cat_V) #aca se esta tabulando
23

```

```
1 cat_V_serie.describe() #Nos dice categorias
```

```

count      100
unique       7
top         A0
freq        51
dtype: object

```

```
1 cat_V_serie.value_counts() # estamos haciendo un conteo
```

```

A0      51
A1      13
A2      11
A3       9
A4       7
A6       5
A5       4
dtype: int64

```

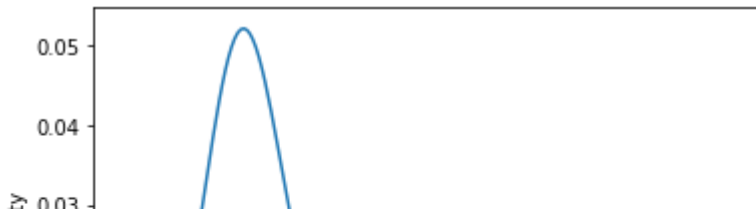
Generar una grafica de linea (densidad) y de barras

```

1 import seaborn as sns
2 sns.kdeplot(V)

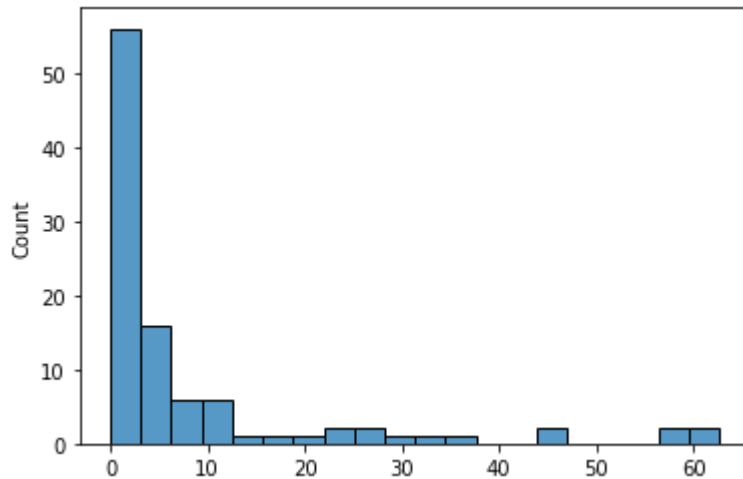
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f0bdc92f6d0>
```



```
1 sns.histplot(V)
```

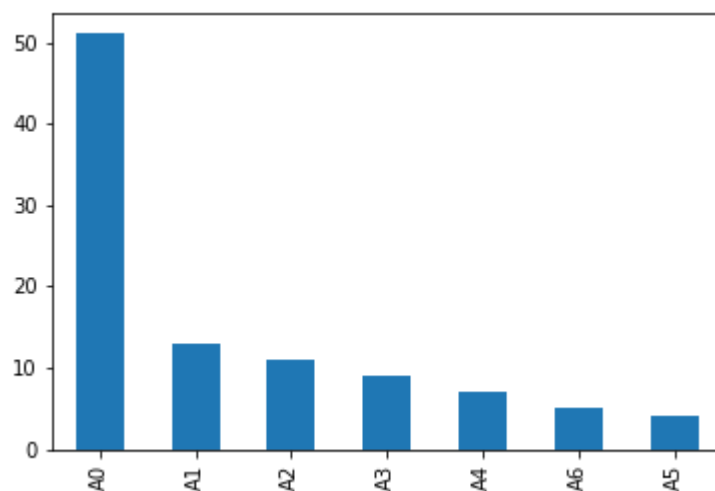
```
<matplotlib.axes._subplots.AxesSubplot at 0x7f0bda08e210>
```



La mayoría son sanas o con leve afectación y solo tenemos un pequeño porcentaje de las hojas afectadas superiores al 60%

```
1 cat_V_serie.value_counts().plot(kind='bar')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f0bd9b468d0>
```



El histograma y esta grafica de barras son muy distintas porque se ve mucho mas "suave" (menos contranstante) que el histograma ya que el histograma muestra como se comportan realmente los

datos de porcentaje de daño, nosotros al haber realizado esta categorización perdemos

Calcular la media y la **mediana** de los datos

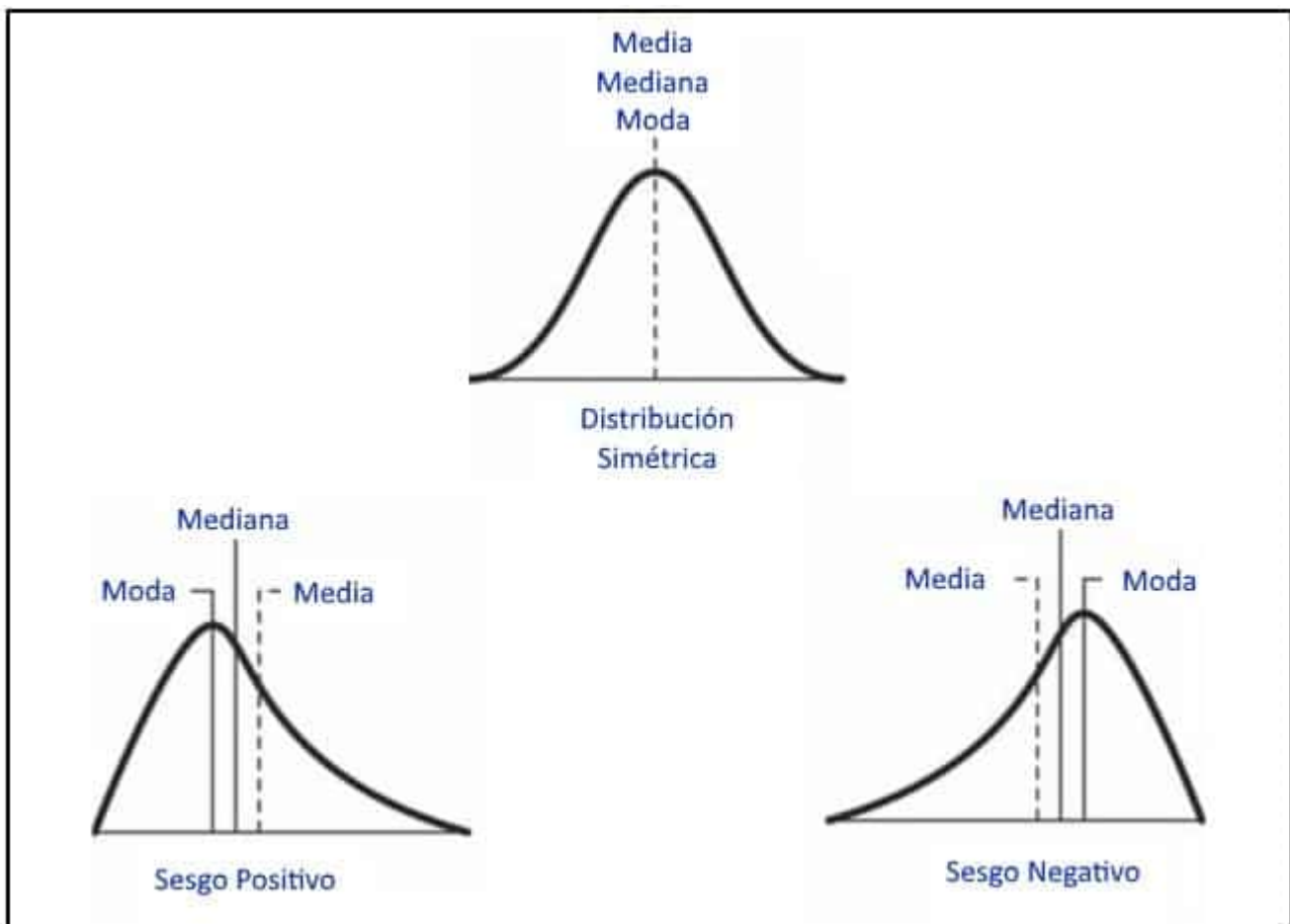
```
1 media = V.mean()
2 print(media)
```

```
7.912500000000001
```

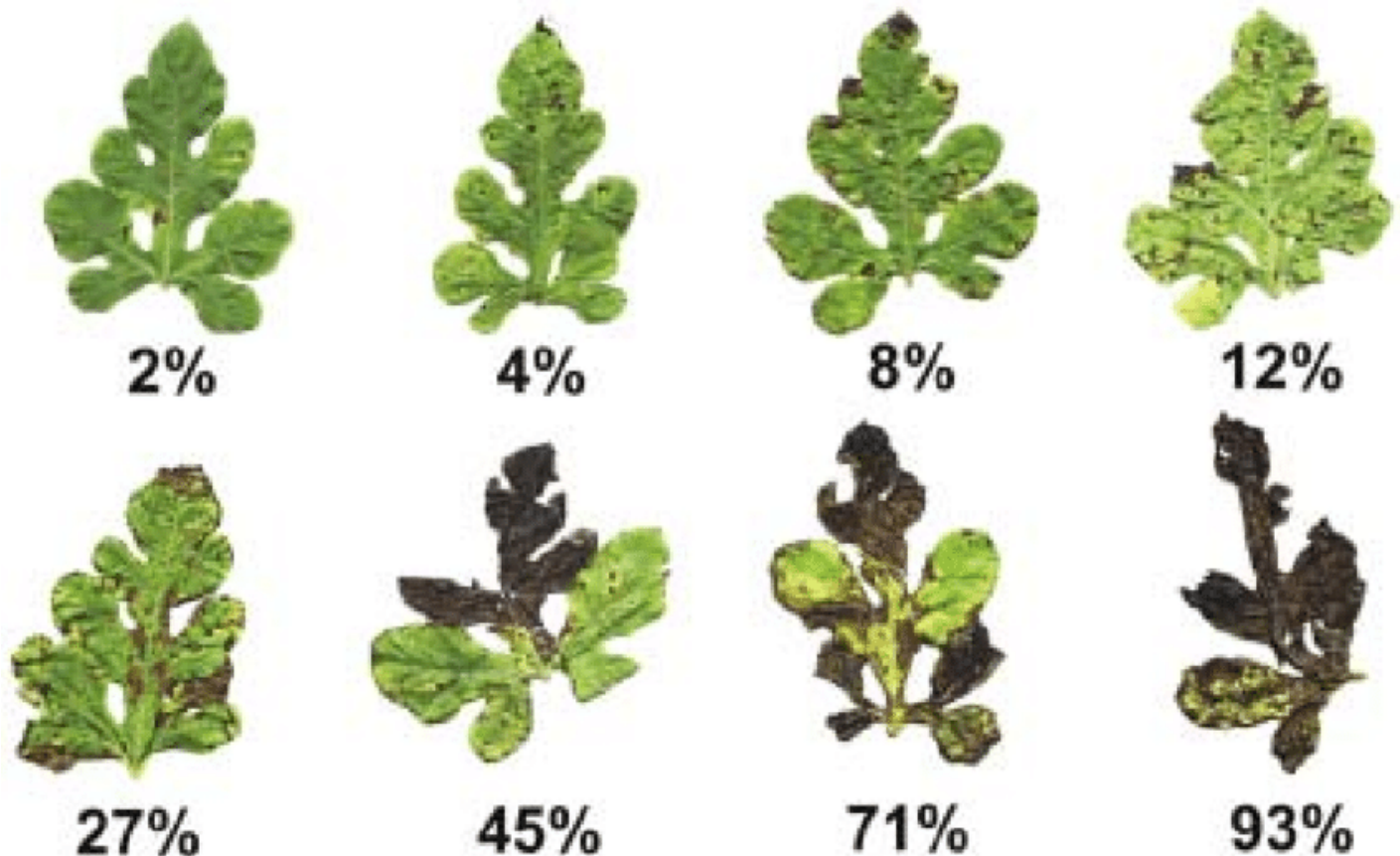
```
1 mediana = np.median(V)
2 print(mediana)
```

```
1.85
```

Los datos de la media al tener datos asimétricos realmente no representan la realidad de los datos, mientras que la **mediana** si permite identificar la realidad de la variable. Vean el contraste entre la línea 179 y la 180. La **mediana** es 1,85 la media es 7.9125 es decir que el dato que ocupa la posición central no es el 7.9 sino el 1.85, la media al tener datos y una curva tan asimétricos la media es muy poco representativa.



La media fue de 7.9 y la mediana 1.85 al tener datos asimetricos las graficas se diferencian es por esto que podemos inferir que nosotros tenemos un sesgo o asimetria positivo. En este caso la mediana esta mas acercado a la realidad de nuestro ejercicio donde la mayoría de los datos se acercan a A1=2% del daño en escala y la grafica. Hay que recordar que la mediana es el cuartil medio (donde el 50% de los datos) y cuartil inferior(25%) y cuartil superior (75).



Calcular los cuartiles y percentiles

```
1 q1 = pd.Series(V)
2 print(q1)
```

```
0    0.03
1    0.01
2   17.09
3    3.88
4    6.25
...
95    0.02
96   60.19
97    0.84
98    1.69
99    2.71
Length: 100, dtype: float64
```

```
1 q1.quantile([0.05, 0.1,0.15, 0.2,0.25,0.3,0.35,0.4,0.45,0.5,0.55,0.6,0.65,0.7,0.75,0.8,0.8
0.05      0.0000
0.10      0.0190
0.15      0.0300
0.20      0.0400
0.25      0.0700
0.30      0.3020
0.35      0.4595
0.40      0.7860
0.45      1.3255
0.50      1.8500
0.55      2.7045
0.60      3.4140
0.65      4.3425
0.70      5.1580
0.75      7.4150
0.80     10.5780
0.85     13.5880
0.90     26.3990
0.95     44.7770
dtype: float64
```

##En cual posición debo ubicarme en la escala para estimar la severidad real

```
1 print(cat_V_serie.value_counts())

A0      51
A1      13
A2      11
A3       9
A4       7
A6       5
A5       4
dtype: int64

1 cont = cat_V_serie.value_counts()
2
3 # Con el punto de corte real
4 frec1 = cont * [2, 4, 8, 12, 27, 45, 71]
5 print(frec1.sum()/100)
6
7 # Con el punto medio
8 frec2 = cont*[1, 2, 3, 13.5, 22.5, 35.5, 46.5]
9 print(frec2.sum()/100)
10
11 # Punto percentil 20%
12 frec3 = cont*[0.02, 0.04, 0.08, 0.12, 0.27, 0.45, 0.71]
13 print(frec3.sum()/100)
```

10.48
7.525
0.1048



```
# **Conclusión**: La mayor densidad de los datos  
mediana: 1.85 o en su defecto la mayoría son valores  
que en la escala corresponde al 2% de severidad.  
escala diagramatica se ignora completamente la  
hay que considerar valores en la escala menores
```

Conclusión: La mayor densidad de los datos se encuentra con el de la mediana: 1.85 o en su defecto la mayoría son valores que pertenecen al grupo que en la escala corresponde al 2% de severidad. Estimando la severidad con la escala diagramatica se ignora completamente la severidad real, para hallarla hay que considerar valores en la escala menores al percentil 2%.

✓ 0 s completado a las 19:55

