

Tipos de análisis

Debido a la diversidad con la que contamos en los datos, existen diferentes formas de realizar el análisis. Sin embargo, retomando el informe anterior (1), la forma más común de clasificarlos es en los dos siguientes grupos:

- Cuantitativo. Es el análisis que se realiza sobre los datos de forma numérica aplicando métodos estadísticos.
- Cualitativo. Es el análisis que no se basa en datos numéricos, sino en clasificaciones categoricas textuales, al que también le podemos aplicar diferentes métodos estadísticos.

Un ejemplo, la valoración de un producto puede ser de forma cualitativa, en las categorías Malo, Normal y Bueno, o también de forma numérica del 1 al 5.

Además de este tipo de división, nosotros vamos a ver cómo podemos usar los métodos estadísticos de forma numérica a lo largo de este tema.

Para añadir valor a este, en el siguiente tema veremos cómo también podemos mostrar esta información estadística de forma gráfica.

Resúmenes estadísticos

Los resúmenes estadísticos son aquellos que se basan en la estadística descriptiva para, como su nombre indica, describir las propiedades de los datos y el conjunto.

Nosotros vamos a cargar el conjunto de datos iris, que es uno de los datasets más comunes y estudiados. Una vez cargado, exploraremos sus propiedades haciendo uso de la estadística descriptiva.

Podemos cargarlos directamente desde el módulo datasets de la librería scikit-learn. Esta librería contiene diferentes datasets de juguete (Boston house prices, Iris, Diabets, Wine): https://scikit-learn.org/stable/datasets/toy_dataset.html

Además contiene otros datasets del mundo real y datasets que han generado ellos mismos.

La ventaja es que no tenemos que tenerlo cargado en local, la desventaja es que el traspaso a data frame no tan visual y directo como ocurre con el uso de la librería pandas.

```
In [55]: from sklearn import datasets
dataset = datasets.load_iris(return_X_y=True, as_frame=True)
print(dataset)

   (    sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)
0           5.1              3.5             1.4               0.2
1           4.9              3.0             1.4               0.2
2           4.7              3.2             1.3               0.2
3           4.6              3.1             1.5               0.2
4           5.0              3.6             1.4               0.2
..          ...
145          6.7              3.0             5.2               1.9
146          6.3              2.5             5.0               1.9
147          6.5              3.0             5.2               2.0
148          6.2              3.4             5.4               2.3
149          5.9              3.0             5.1               1.8
[150 rows x 4 columns], 0      0
1      0
2      0
3      0
4      0
..
145     2
146     2
147     2
148     2
149     2
Name: target, Length: 150, dtype: int32)
```

Por otro lado, podemos descargar el dataset de <https://archive.ics.uci.edu/ml/machine-learning-databases/iris/> y descargar el archivo iris.data. Lo cargamos como hicimos en los temas anteriores.

```
In [56]: import pandas as pd
properties = ['Logitud sépalo', 'Ancho sépalo', 'Longitud pétalo', 'Ancho pétalo', 'Clase']
dataset = pd.read_csv('iris.data', names=properties)
dataset.head()
```

```
Out[56]:   Logitud sépalo  Ancho sépalo  Longitud pétalo  Ancho pétalo  Clase
0           5.1            3.5             1.4               0.2  Iris-setosa
1           4.9            3.0             1.4               0.2  Iris-setosa
2           4.7            3.2             1.3               0.2  Iris-setosa
3           4.6            3.1             1.5               0.2  Iris-setosa
4           5.0            3.6             1.4               0.2  Iris-setosa
```

Una vez cargado como data frame, Python nos proporciona la función describe(), a modo de resumen sobre la que podemos extraer información genérica de nuestro dataset.

```
In [57]: dataset.describe()
```

```
Out[57]:   Logitud sépalo  Ancho sépalo  Longitud pétalo  Ancho pétalo
count    150.000000  150.000000  150.000000  150.000000
mean      5.843333  3.054000  3.758667  1.198667
std       0.828066  0.433594  1.764420  0.763161
min       4.300000  2.000000  1.000000  0.100000
25%      5.100000  2.800000  1.600000  0.300000
50%      5.800000  3.000000  4.350000  1.300000
75%      6.400000  3.300000  5.100000  1.800000
max       7.900000  4.490000  6.900000  2.500000
```

Podemos extraer varias conclusiones de este análisis obtenido para cada columna (propiedad) del dataset iris.

- count: Esta columna nos muestra el número de muestras para esta propiedad. En este caso vemos que disponemos de 150 muestras por característica, pero existen datasets en los que una columna puede tomar valores indefinidos (N/A), inexistentes (None) e incluso que 0 u otro valor acordado representa la falta de datos para una propiedad de una instancia. Es un atributo muy importante ya que podemos ver a simple vista si contamos con pocos datos respecto al tamaño del dataset y, por tanto, no es conveniente utilizar dicha propiedad en el análisis.
- mean: Indica el valor medio de una propiedad considerando todos los valores con el mismo peso. Si el tipo de la propiedad es numérico no ordenado, podemos usar la media, pero para valores ordenados o basados en cadenas, necesitaremos usar la moda.
- std: Muestra la dispersión (desviación típica) de los valores. Es decir, la medida en la que se alejan o no de la media. Puede ver que esta será mayor conforme más distancia exista entre el valor mínimo y el máximo. Debido a que es sensible a la variable, es conveniente normalizar los datos. (No es lo mismo comparar la dispersión entre variables con diferentes magnitudes)
- min / máximo: Valores mínimo y máximo de cada propiedad
- 25% / 50% / 75%: Cuartiles que indican el valor a partir del cual se encuentra el porcentaje indicado de datos

Si usamos un datasets que contenga atributos de tipo carácter: <https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer/>

```
In [58]: properties = ['class', 'age', 'menopause', 'tumor-size', 'inv-nodes', 'node-caps', 'deg-malig', 'breast', 'breast-quad', 'irradiat']
dataset_cancer = pd.read_csv('breast-cancer.data', names=properties)
dataset_cancer.head()
```

```
Out[58]:   class  age  menopause  tumor-size  inv-nodes  node-caps  deg-malig  breast  breast-quad  irradiat
0  no-recurrence-events  30-39  premeno  30-34  0-2  no  3  left  left_low  no
1  no-recurrence-events  40-49  premeno  20-24  0-2  no  2  right  right_up  no
2  no-recurrence-events  40-49  premeno  20-24  0-2  no  2  left  left_low  no
3  no-recurrence-events  60-69  ge40  15-19  0-2  no  2  right  right_up  no
4  no-recurrence-events  40-49  premeno  0-4  0-2  no  2  right  right_low  no
```

Una vez cargado como data frame, Python nos proporciona la función describe() a modo de resumen sobre la que podemos extraer información genérica de nuestro dataset.

```
In [59]: dataset_cancer.describe() #Por defecto solo incluye valores numéricos
```

```
Out[59]:   deg-malig
count    286.000000
mean      2.048951
std       0.738217
min       1.000000
25%      2.000000
50%      2.000000
75%      3.000000
max       3.000000
```

```
In [60]: dataset_cancer.describe(include='all')
```

```
Out[60]:   class  age  menopause  tumor-size  inv-nodes  node-caps  deg-malig  breast  breast-quad  irradiat
count    286.000000  286.000000  286.000000  286.000000  286.000000  286.000000  286.000000  286.000000
unique     2       6       3       11       7       3  NaN       2       2       6       2
top  no-recurrence-events  premeno  30-34  0-2  no  NaN  left  left_low  no
freq      201      96     150      60     213     222  NaN      152     110     218
mean      NaN      NaN      NaN      NaN      NaN      NaN  2.048951  NaN      NaN      NaN
std       NaN      NaN      NaN      NaN      NaN      NaN  0.738217  NaN      NaN      NaN
min       NaN      NaN      NaN      NaN      NaN      NaN  1.000000  NaN      NaN      NaN
25%      NaN      NaN      NaN      NaN      NaN      NaN  2.000000  NaN      NaN      NaN
50%      NaN      NaN      NaN      NaN      NaN      NaN  2.000000  NaN      NaN      NaN
75%      NaN      NaN      NaN      NaN      NaN      NaN  3.000000  NaN      NaN      NaN
max       NaN      NaN      NaN      NaN      NaN      NaN  3.000000  NaN      NaN      NaN
```

Para datos categóricos veremos como para los anteriores no muestra un resultado (NaN), pero incluye otros adicionales relativos a este tipo de dato:

- unique: número de valores categóricos distintos. Por ejemplo en este dataset la clase únicamente puede tomar 2 valores (tener cáncer o no).
- top: el valor más común (moda).
- freq: frecuencia del valor más común (moda)

Este segundo test no es válido en aquellas muestras en las que la media y varianza no sean representativas. Por ello, también se añade el test de Jarque-Bera que no requiere estimaciones de los parámetros que caracterizan la normal.

En todos ellos, podemos rechazar la hipótesis de que la distribución de nuestra variable es normal si el p-value es menor que un determinado valor, por lo general alfa=0.05. Es decir:

- p<=alfa: no rechazar hipótesis, aseguramos que no es normal.
- p>alfa: no podemos rechazar la hipótesis, asumimos normalidad porque es muy probable que sea cierto tras la prueba realizada.

```
In [61]: dataset['Clase'].value_counts() #Vemos como hay 50 instancias de cada clase
```

```
Out[61]: Iris-setosa      50
Iris-versicolor      50
Iris-virginica       50
Name: Clase, dtype: int64
```

```
In [62]: dataset['Logitud sépalo'].value_counts() #Nos mostraría todos los posibles valores decimales
dataset['Logitud sépalo'].value_counts(bins=5) #bins es el número de intervalos que queremos construir
```

```
Out[62]: (5.74, 6.46]      42
(5.02, 5.74]      41
(4.295, 5.02]     32
(6.46, 7.18]      24
(7.18, 7.9]       11
Name: Logitud sépalo, dtype: int64
```

¿Qué conclusiones podemos obtener analizando la frecuencia?

Respecto a la clase podemos ver si el dataset está o no balanceado. Es decir, si el número de observaciones que tenemos de una clase es muy inferior al de otras, al aplicar un modelo este las tendrá menos en cuenta, produciendo una desventaja en su clasificación.

Por otro lado, podemos visualizar como se distribuye una determinada propiedad. En este caso la longitud del sépalo tiene a valores cercanos al rango [5.02,6.46], siendo menos común las flores del sépalo de 7 centímetros.

El análisis univariante tiene que ir acompañado del bivariante, donde podríamos verificar si este sépalo mayor se debe a que dentro de una determinada clase existe una tendencia a ello.

Media/Moda/Mediana

Otro aspecto relevante a observar es la Media, Moda y Mediana. Podemos coger al igual que la frecuencia las mismas columnas del data frame.

```
In [63]: print("Media de longitud del sépalo: ",dataset['Logitud sépalo'].mean())
print("Mediana de longitud del sépalo: ",dataset['Logitud sépalo'].median())
print("Moda de longitud del sépalo: ",dataset['Logitud sépalo'].mode()[0])
```

Media de longitud del sépalo: 5.843333333333334
Mediana de longitud del sépalo: 5.7
Moda de longitud del sépalo: 5.0

Aunque en este caso no se observa tan claramente el efecto de cada uno de ellas, podemos observarlo en un ejemplo sencillo.

Imaginemos que tenemos las alturas de una clase, donde uno de ellos es muy alto respecto al resto.

```
In [64]: dataset_cancer.describe(include='all')
```

```
Out[64]:   class  age  menopause  tumor-size  inv-nodes  node-caps  deg-malig  breast  breast-quad  irradiat
count    286.000000  286.000000  286.000000  286.000000  286.000000  286.000000  286.000000  286.000000
unique     2       6       3       11       7       3  NaN       2       2       6       2
top  no-recurrence-events  premeno  30-34  0-2  no  NaN  left  left_low  no
freq      201      96     150      60     213     222  NaN      152     110     218
mean      NaN      NaN      NaN      NaN      NaN      NaN  2.048951  NaN      NaN      NaN
std       NaN      NaN      NaN      NaN      NaN      NaN  0.738217  NaN      NaN      NaN
min       NaN      NaN      NaN      NaN      NaN      NaN  1.000000  NaN      NaN      NaN
25%      NaN      NaN      NaN      NaN      NaN      NaN  2.000000  NaN      NaN      NaN
50%      NaN      NaN      NaN      NaN      NaN      NaN  3.000000  NaN      NaN      NaN
75%      NaN      NaN      NaN      NaN      NaN      NaN  4.000000  NaN      NaN      NaN
max       NaN      NaN      NaN      NaN      NaN      NaN  5.000000  NaN      NaN      NaN
```

Para datos categóricos veremos como para los anteriores no muestra un resultado (NaN), pero incluye otros adicionales relativos a este tipo de dato:

- unique: número de valores categóricos distintos. Por ejemplo en este dataset la clase únicamente puede tomar 2 valores (tener cáncer o no).
- top: el valor más común (moda).
- freq: frecuencia del valor más común (moda)

Este segundo test no es válido en aquellas muestras en las que la media y varianza no sean representativas. Por ello, también se añade el test de Jarque-Bera que no requiere estimaciones de los parámetros que caracterizan la normal.

En todos ellos, podemos rechazar la hipótesis de que la distribución de nuestra variable es normal si el p-value es menor que un determinado valor, por lo general alfa=0.05. Es decir:

- p<=alfa: no rechazar hipótesis, aseguramos que no es normal.
- p>alfa: no podemos rechazar la hipótesis, asumimos normalidad porque es muy probable que sea cierto tras la prueba realizada.

```
In [65]: print("Desviación típica de longitud del sépalo: ",dataset['Logitud sépalo'].std())
print("Desviación típica de ancho del sépalo: ",dataset['Ancho sépalo'].std())
```

Desviación típica de longitud del sépalo: 0.829066127977863
Desviación típica de ancho del sépalo: 0.4335943113621737

Observamos que la dispersión es mayor en la longitud que en el ancho. Es decir, si no conocemos las medidas de una nueva flor y le asignásemos la media, el error sería mayor en la longitud que en su anchura.

Pero, ¿es adecuada esta comparación? Imaginemos que queremos acertar aleatoriamente la altura de una persona, y nos confundimos en 1 metro, 2 kilos. Estamos cometiendo un mayor error en el peso si evaluamos la cantidad sin tener en cuenta las unidades pero, a nivel relativo, en la altura el error es mayor porque el rango de estimación es mucho menor que si nos referimos al peso.

El coeficiente de variación es adimensional, y nos permite ver la variación de una muestra. Se obtiene con la división de la desviación típica a la media. Para la longitud del sépalo y el ancho, vemos que el coeficiente de variación (la distribución de sus valores respecto a la media) es similar. Si el C.V es menor o igual al 30%, significa que la media aritmética es representativa del conjunto de datos.

```
In [66]: import scipy.stats as ss
print(variation(dataset['Logitud sépalo']))
print("Coeficiente de variación de longitud del sépalo: ",ss.variation(dataset['Logitud sépalo']))
print("Coeficiente de variación de ancho del sépalo normalizado: ",ss.variation(dataset['Ancho sépalo']))
```

Coeficiente de variación de longitud del sépalo: 0.142380989346393
Coeficiente de variación de ancho del sépalo normalizado: 0.141501827135083

```
In [67]: import numpy as np
altura = np.array([2.18, 1.72, 1.72, 1.69, 1.68, 1.63, 1.56])
print("Media: ",altura.mean())
print("Moda: ",altura.mode())
print("Desviación típica: ",np.std(altura))
print("Coeficiente de variación: ",np.std(altura)/np.mean())
print("Asimetría: ",skew(altura))
```

Media: 1.7400000000000002
Moda: 1.69
Desviación típica: 0.1687882894444444
Coeficiente de variación: 0.09633333333333333
Asimetría: 0.3171431172366936

Observamos que aunque la mayoría de los alumnos se encuentran en el rango [1.72,1.68], la media se desplaza hasta 1.74 por el peso del alumno más alto, mientras que la mediana mantiene su posición en 1.69, que podríamos entender como el valor cercano más probable que presentará un alumno nuevo. Por último, la moda es el valor que más se repite, en este caso 1.69.

Este segundo test no es válido en aquellas muestras en las que la media y varianza no sean representativas. Por ello, también se añade el test de Jarque-Bera que no requiere estimaciones de los parámetros que caracterizan la normal.

En todos ellos, podemos rechazar la hipótesis de que la distribución de nuestra variable es normal si el p-value es menor que un determinado valor, por lo general alfa=0.05. Es decir:

- p<=alfa: no rechazar hipótesis, aseguramos que no es normal.
- p>alfa: no podemos rechazar la hipótesis