In [2]:	<pre>import pandas as pn import numpy as np DATA WRANGLING</pre>
	Filtering out missing data Distinguimos distintos métodos para filtrar los datos nulos de una set de datos: from IPython.display import Image
Out[3]:	Image("null_methods.png") Argument Description Filter axis labels based on whether values for each label have missing data, with varying thresholds for how much missing data
	to tolerate. Fill in missing data with some value or using an interpolation method such as 'ffill' or 'bfill'. Return like-type object containing boolean values indicating which values are missing / NA.
	Negation of isnull. df.dropna() Como su nombre indica permite eliminar (por defecto) las filas que contengan al menos un valor nulo (NaN o np.nan). Funciona tanto para series como para dataframes:
In [4]:	Serie_ejemplo=pn.Series([1,np.nan,3.5,np.nan,7]) df=pn.DataFrame({0:[1,np.nan,6.5,3,0],1:[1,np.nan,np.nan,np.nan,np.nan,np.nan,np.nan,p.nan,6.5,3,3.5],4:[1,np.nan,3.4,2.6,1]}) print(Serie_ejemplo) print(df) 0 1.0
	1 NaN 2 3.5 3 NaN 4 7.0 dtype: float64 0 1 2 3 4 0 1.0 NaN NaN 1.0
In [5]:	1 NAN NAN NAN NAN NAN NAN 2 0 1 1 2 6 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
Out[5]: In [6]:	0 1.0 2 3.5 4 7.0 dtype: float64 df.dropna()
	d 1 2 3 4 4 0.0 1.4 2.7 3.5 1.0 .dropna() por defecto elimina las filas que tienen al menos un valor NaN, pero podemos especificarle que únicamente elimine aquellas filas que tengan todos sus valores NaN con .dropna(how='all'): .dropna(how='all')
In [7]:	Elimina las filas cuyos todos sus valores son NaN: Serie_ejemplo.dropna(how='all')
In [8]: Out[8]:	4 7.0 dtype: float64 df.dropna(how='all')
	 0 1.0 1.0 NaN NaN NaN 1.0 2 6.5 NaN NaN 6.5 3.4 3 3.0 NaN NaN 3.0 2.6 4 0.0 1.4 2.7 3.5 1.0
	.dropna(axis=1) En vez de eliminar las columnas que tengan al menos un valor NaN, ahora elimina las columnas que tengan al menos un valor NaN: df.dropna(axis=1)
Out[9]:	0 1 2 3
	.fillna() Como su nombre indica, permite rellenar los huecos que tienen valor nulo (NaN) por el valor que especifiques en el paréntesis (variable, int, string, float ,boolean, serie) en cualquier serie o dataframe:
In [10]: Out[10]:	Serie_ejemplo.fillna('valor_de_relleno') 0
In [11]: Out[11]:	
	1 valor_de_relleno valo
In [12]:	4 0 1.4 2.7 3.5 1 .fillna({'nombre_columna':valor_por_el_que_quieras_sustituir_los_NaN}) Al introducir un diccionario en el paréntesis, te permite esque cificar que en la columna X te sustituya los valores NaN de esa columna por el valor que hayas especificado en el value: a=df.fillna({1: 'valor_sust_columna_1', 2: 'valor_sust_columna_2'})
In [13]: Out[13]:	Serie_ejemplo=pn.Series([1,np.nan,3.5,np.nan,7]) df=pn.DataFrame({0:[1,np.nan,6.5,3,0],1:[1,np.nan,np.nan,np.nan,np.nan,np.nan,np.nan,np.nan,6.5,3,3.5],4:[1,np.nan,3.4,2.6,1]}) df
	0 1.0 1.0 NaN NaN NaN 1.0 1 NaN NaN NaN NaN NaN 2 6.5 NaN NaN 6.5 3.4 3 3.0 NaN NaN 3.0 2.6
	4 0.0 1.4 2.7 3.5 1.0 .fillna(method='ffill/bfill', limit=X) Con .fillna(limit=X) lo que haces es indicar el número de sustituciones que hará por columna. Por ejemplo .fillna(limit=3) significa que solo sustituirá hasta 3 NaN por columna, el resto lo deja como está. Es imporescindible especificar el método, ffill te sustituye empezando por arriba y bfill empezando por abajo :
In [14]: Out[14]:	<pre>df.fillna(method='ffill', limit=3)</pre>
	1 1.0 1.0 NaN NaN 1.0 2 6.5 1.0 NaN 6.5 3.4 3 3.0 1.0 NaN 3.0 2.6 4 0.0 1.4 2.7 3.5 1.0
In [15]: Out[15]:	0 1.0 1.0 NaN 6.5 1.0
	1 6.5 NaN NaN NaN 6.5 3.4 2 6.5 1.4 2.7 6.5 3.4 3 3.0 1.4 2.7 3.0 2.6 4 0.0 1.4 2.7 3.5 1.0
	Grouping data df_gd=pn.DataFrame({'data1':[-0.204708,0.478943,-0.519439,-0.555730,1.965781],'data2':[1.393406,0.092908,0.281746,0.769023,1.246435],'key1':['a','a','b','b','a'],'key2':['one','two','one','two','one']}) data1 data2 key1 key2
out[10].	0 -0.204708 1.393406 a one 1 0.478943 0.092908 a two 2 -0.519439 0.281746 b one
	4 1.965781 1.246435 a one Como su nombre indica, agrupar datos nos permite crear a partir de una dataframe un elemento llamado objeto agrupado , que es una especie de serie en el que los datos de una columna del dataframe original se agrupan en base a los datos de otra columna del dataframe original, y para crearlo necesitamos el elemento .groupby():
<pre>In [17]: Out[17]:</pre>	
	Ahora todos los cálculos que se hagan sobre el objeto agrupado se harán de forma individual por cada agrupamiento, por ejemplo si queremos calcular la media de 'data1', lo hará según las letras de 'key1', es decir, la media de las filas que tengan como key1 la letra a por un lado y por otro las b: grouped_object.mean() key1
	key1 a 0.746672 b -0.537584 Name: data1, dtype: float64 Merging data Merging consiste en fusionar 2 columnas según una columna o índice que los relacione. Hay 4 métodos para hacer merging: inner,outer,right y left:
	 Inner: Es aquel que se hace por defecto. Es un merging uno a muchos, en el que no habrá valores NaN, es de cir, los valores coincidentes se rellenan en todas las posibilidades y los que no coincidan se eliminan. Outer: Es un merging 1 a 1, es decir, los valores coincidentes se unen mientras que los no coincidentes conservan su valor original pero en la nueva columna ponen un NaN porque no hay ningún valor coincidente (es un inner pero los valores no coincidentes de ambas tablas se muestran en la nueva con un NaN en la columna que no tienen pareja).
	• Left/Right: Se hace un merging muchos a muchos, es decir, se generan todas las posibles combinaciones entre los datos coincidentes de las 2 tablas, y los que no coincidan pues se pone un valor NaN en la columna que no coincida (es como un 'outter' pero si pones 'right' solo se conservan los valores no coincidentes de la segunda tabla/serie, rellenando con NaN la columna sin pareja, y si pones 'left' se conservan los de la primera tabla/serie. pn.merge(df_1, df_2, on='columna_comun', how='outer'/'inner'/'right') Partimos de 2 dataframes:
In [19]:	<pre>Partimos de 2 dataframes: df_merge_1=pn.DataFrame({'data1':[0,1,2,3,4,5,6],'key':['b','b','a','c','a','b']}) df_merge_2=pn.DataFrame({'data2':[0,1,2,0],'key':['a','b','d','b']}) print(df_merge_1) print('') print(df_merge_2)</pre>
	data1 key 0
	5
<pre>In [20]: Out[20]:</pre>	<pre>df_juntos=pn.merge(df_merge_1, df_merge_2, on='key', how='inner') df_juntos</pre>
	1 0 b 0 2 1 b 1 3 1 b 0 4 6 b 1
	5 6 b 0 6 2 a 0 7 4 a 0 8 5 a 0
	Podemos ver que las tablas se han unido mediante un innerjoin por la columna común que en este caso es 'key' Se puede dar el caso de que la columna común a 2 tablas no se llame igual, o que la columna común sea el índice, por lo que existen métodos específicos para indicarlo: from IPython.display import Image
Out[21]:	Image("merging_methods.png") Argument Description left DataFrame to be merged on the left side right DataFrame to be merged on the right side
	One of 'inner', 'outer', 'left' or 'right'. 'inner' by default Column names to join on. Must be found in both DataFrame objects. If not specified and no other join keys given, will use the intersection of the column names in left and right as the join keys left_on Columns in left DataFrame to use as join keys
	right_on Columns in right DataFrame to use as join keys left_index Use row index in left as its join key (or keys, if a MultiIndex) right_index Analogous to left_index Sort merged data lexicographically by join keys; True by default.
	Tuple of string values to append to column names in case of over- suffixes lap; defaults to ('_x', '_y'). For example, if 'data' in both DataFrame objects, would appear as 'data_x' and 'data_y' in result If False, avoid copying data into resulting data structure in some
In [22]:	exceptional cases. By default always copies Por ejemplo si tenemos las siguientes tablas y queremos hacer un leftjoining usando como key la columna de key de la tabla izquierda y el índice como key de la tabla derecha left1=pn.DataFrame({'key':['a', 'b', 'a', 'a', 'b', 'c'], 'value':[0,1,2,3,4,5]}) right1=pn.DataFrame({'group_val':[3.5,7.]}, index=['a', 'b']) print (left1)
	<pre>print('') print(right1) key value 0</pre>
	3
In [23]: Out[23]:	<pre>df_merged=pn.merge(left1, right1, left_on='key', right_index=True, how='inner') df_merged key value group_val 0 a 0 3.5</pre>
	2 a 2 3.5 3 a 3 3.5 1 b 1 7.0 4 b 4 7.0
	Concatenating data sets pn.concat([df1,df2,serie1,serie2]) Concatenar es una forma de unir "a lo bruto" 2 o más dataframes o series entre sí. Podemos unir dataframes, series, o los dos. Por defecto, se unirán en vertical, es decir, el dataframe 2 quedará debajo del dataframe 1, y se conservará el índice de cada uno de ellos. Partimos de 2 dataframes:
In [24]:	<pre>df1=pn.DataFrame(np.arange(6).reshape(3,2), index=['a','b','c'], columns=['one','two']) df2=pn.DataFrame(5 + np.arange(4).reshape(2,2),index=['a','c'],columns=['three','four']) print(df1) print('')</pre>
	<pre>print(df2) one two a 0 1 b 2 3 c 4 5</pre>
In [25]:	three four a 5 6 c 7 8 Ahora las concatenamos: df1_2=pn.concat([df1,df2]) df1_2
Out[25]:	
	a NaN NaN 5.0 6.0 c NaN NaN 7.0 8.0 Como podemos ver, se han unido en vertical y han conservado el índice. Si no queremos que conserven el índice podemos poner ignore_index=True:
In [26]: Out[26]:	0 0.0 1.0 NaN NaN
	1 2.0 3.0 NaN NaN 2 4.0 5.0 NaN NaN 3 NaN NaN 5.0 6.0 4 NaN NaN 7.0 8.0
In [27]: Out[27]:	df1_2
	a 0 1 5.0 6.0 b 2 3 NaN NaN c 4 5 7.0 8.0 Se puede dar el caso de que querramos poder acceder y saber a qué tabla original pertenecían cada uno de los elementos que forman parte del concat, y para eso tenemos keys=[key_df_1,key_df_2]:
<pre>In [28]: Out[28]:</pre>	<pre>df1_2=pn.concat([df1,df2], keys=['df1','df2']) df1_2</pre>
	 b 2.0 3.0 NaN NaN c 4.0 5.0 NaN NaN df2 a NaN NaN 5.0 6.0 c NaN NaN 7.0 8.0
In [29]: Out[29]:	Por tanto podremos acceder rápidamente con .loc['key'] a los datos del df original que queramos: df1_2.loc['df1'] one two three four
	a 0.0 1.0 NaN NaN b 2.0 3.0 NaN NaN c 4.0 5.0 NaN NaN Podemos también fusionar df con series, por ejemplo:
In [30]:	<pre>series=pn.Series([1,3,5],index=['a','b','c'],name='serie_ejemplo') print(series) a 1 b 3 c 5</pre>
In [31]: Out[31]:	Name: serie_ejemplo, dtype: int64 df_1_2_serie=pn.concat([df1,series], axis=1) df_1_2_serie
	a 0 1 1 1 b 2 3 3 3 c 4 5 5 5 Importante poner en la serie un nombre con name='XXXXX' para que se le pueda asignar un nombre de columna a la serie, porque si no no funciona. Tambien hay que tener en cuenta que el concat horizontal SE HARÁ EN FUNCIÓN DE LOS ÍNDICES , POR LO QUE SI NO COINCIDEN NO SE UNEN
	NO COINCIDEN NO SE UNEN Hay una gran cantidad de argumentos para pn.concat(): from IPython.display import Image Image("concat_methods.png")
Out[32]:	Argument Description bis List or dict of pandas objects to be concatenated. The only required argument Axis Axis to concatenate along; defaults to 0
	One of 'inner', 'outer', defaulting to 'outer'; whether to intersection (inner) or union (outer) together indexes along the other axes Specific indexes to use for the other n-1 axes instead of performing union/intersection logic Values to associate with objects being concatenated, forming
	a hierarchical index along the concatenation axis. Can either be a list or array of arbitrary values, an array of tuples, or a list of arrays (if multiple level arrays passed in levels) Specific indexes to use as hierarchical index level or levels if
	Names for created hierarchical levels if keys and / or levels passed verify_integrity Check new axis in concatenated object for duplicates and raise exception if so. By default (False) allows duplicates
	ignore_index Do not preserve indexes along concatenation axis, instead producing a new range(total_length) index Combinar datasets set1.combine_first(set2)
	Combinar datasets permite enfrentar 2 sets de datos (series o dataframes) y rellenar los datos nulos del primero con el correspondiente valor del segundo: Por ejemplo si tenemos 2 series: serie1=pn.Series([np.nan, 2.5, np.nan, 3.5, 4.5, np. nan], index=['f', 'e', 'd', 'c', 'b', 'a']) serie2=pn.Series(np. arange (len(serie1), dtype=np.float64),index=['f', 'e', 'd', 'c', 'b', 'a'])
	<pre>print(serie1 , serie2) f NaN e 2.5 d NaN c 3.5 b 4.5</pre>
	a NaN dtype: float64 f 0.0 e 1.0 d 2.0 c 3.0 b 4.0 a 5.0
In [34]: Out[34]:	dtype: float64 Queremos rellenar los datos nulos de serie1 con los datos correspondientes de serie2, por lo que usamos .combine_first(): serie1.combine_first(serie2) f 0.0
	e 2.5 d 2.0 c 3.5 b 4.5 a 5.0 dtype: float64 Lo mismo se aplica a dataframes:
In [35]:	df1=pn.DataFrame({'a':[1,np.nan,5,np.nan],'b':[np.nan,2,np.nan,6],'c':range(2,18,4)}) df2= pn.DataFrame({'a':[5,4,np.nan,3,7],'b':[np.nan,3,4,6,8]}) print(df1) print(df2) a b c
	0 1.0 NaN 2 1 NaN 2.0 6 2 5.0 NaN 10 3 NaN 6.0 14 a b 0 5.0 NaN
In [36]: Out[36]:	
.c[36]:	0 1.0 NaN 2.0 1 4.0 2.0 6.0 2 5.0 4.0 10.0 3 3.0 6.0 14.0
	4 7.0 8.0 NaN Reshaping datasets df.stack()> Convierte las columnas de un dataframe en índices, generando un dataframe con índice jerárquico
In [37]:	df.unstack()> Convierte los índices del dataframe del índice jerárquico en columnas, opuesto a df.stack() df1_stack=df1.stack() df1_stack
Out[37]:	C 2.0 1 b 2.0 C 6.0 2 a 5.0 C 10.0 3 b 6.0 C 14.0
In [38]: Out[38]:	
	 2 5.0 NaN 10.0 3 NaN 6.0 14.0 Añadir fila a partir de una serie en un dataframe Podemos añadir una fila más a una lista con el comando df.append(lista, ignore_index=True), obviamente la serie debe seguir la estructura de las columnas de la tabla y debe tener índice para quitarlo.
In [40]:	Podemos añadir una fila más a una lista con el comando df.append(lista, ignore_index=True), obviamente la serie debe seguir la estructura de las columnas de la tabla y debe tener índice para quitarlo. df1=pn.DataFrame({'student_id':['S1','S2','S3','S4','S5'], 'name':['Daniella Fenton', 'Ryder Storey', 'Bryce Jensen', 'Ed Bernal', 'Kwame Morin'], 'marks':[200,210,190,222,199]}) df2=pn.Series(['S6','Scarlette Fisher',205], index=['student_id', 'name', 'marks'], name='New Row') df1.append(df2, ignore_index=True)
Out[40]:	
	0 S1 Daniella Fenton 200 1 S2 Ryder Storey 210
	0 S1 Daniella Fenton 200
	0 S1 Daniella Fenton 200 1 S2 Ryder Storey 210 2 S3 Bryce Jensen 190 3 S4 Ed Bernal 222 4 S5 Kwame Morin 199
	0 S1 Daniella Fenton 200 1 S2 Ryder Storey 210 2 S3 Bryce Jensen 190 3 S4 Ed Bernal 222 4 S5 Kwame Morin 199
	0 S1 Daniella Fenton 200 1 S2 Ryder Storey 210 2 S3 Bryce Jensen 190 3 S4 Ed Bernal 222 4 S5 Kwame Morin 199
	0 S1 Daniella Fenton 200 1 S2 Ryder Storey 210 2 S3 Bryce Jensen 190 3 S4 Ed Bernal 222 4 S5 Kwame Morin 199