<pre>import pandas as pn import numpy as np dates= pn.date_range('1/1/2000', periods=8) df = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D']) df</pre>
colB=df.B.values colC=df.C.values colD=df.D.values boolean=[] print(df) df.loc[(df.B>df.A) & (df.B <df.c), 'b']="" -0.000783="" -0.442027="" -0.666538="" -0.793887<="" -1.422871="" -1.576848="" 0.203998="" 0.393688="" 0.641780="" 1.008718="" 1.565900="" 1.889946="" 2000-01-01="" 2000-01-02="" 2000-01-03="" a="" b="" c="" d="" th=""></df.c),>
2000-01-04 -0.483422 -0.564830
cor ejemplo si queremos el sexto elemento de la columna A: df['A'][dates[5]] -1.2808605940399247 df.loc[dates[5], 'A']
EJEMPLO: Cambiar el orden de 2 columnas Para ello simplemente llamamos a 2 columnas, y las igualamos a las columnas al revés. Es decir:
2000-01-01 -0.666538 -0.000783 0.641780 1.889946 2000-01-02 1.565900 0.203998 -1.576848 1.008718 2000-01-03 -1.422871 -0.442027 0.393688 -0.793887 2000-01-04 -0.564830 -0.483422 0.065075 0.50323 2000-01-05 -0.550484 1.954107 0.811359 0.496025 2000-01-06 0.669350 -1.280861 0.179829 -0.402134 2000-01-08 -0.56701 1.365486 0.488785 -0.775794
Llamada por índice, .loc[] y .iloc[] Sásicamente hacen lo mismo pero .loc[] únicamente funciona con el nombre del índice mientras que .iloc[] solo funciona con la posición en número: from IPython.display import Image Image("loc_iloc.png") Input .loc .iloc .iloc
Single value 5 or 'a' 5 List of values ['a', 'b', 'c'] [4, 3, 0] Slice object 'a':'f' 1:7 Boolean array [True, False,] [True, False,] Callable object Function Function Llamar elementos a partir de booleanos Sásicamente, metes la condicion dentro del corchete y lo que de True, es decir, cuando se cumpla la condición, se selecciona el valor correspondiente. Por ejemplo, si hacemos un dataframe.loc[dataframe.columna>0], lo que hace es comparar todos los elementos de la columna a 0, legado por ejemplo. True True False true]. Este bará que se imprimente a cogundo y questo file por parte del los
ando por ejemplo [True, True, False, true]. Esto hará que se impriman la primera, segunda y cuarta fila por parte del .loc df .loc [df .A > 0] A B C D
MPOTANTE: Tienes que usar dataframe.columna, no vale poner directamente el nombre de la columna, porque >, < o = no pueden ir entre un string y un int. Filtrar números de una tabla, poniendo NaN a los que no cumplan la condicion El comando es dataframe[dataframe > 34], lo que hace la condicion del corchete es crear un dataframe de booleanos, y después los valores True conservan su valor mientras que los False pasan a perder el valor(NaN). Por ejemplo, filtar del dataframe los valores mayores de 0: df[df >0] A B C D
2000-01-01 NaN NaN 0.641780 1.889946 2000-01-02 1.56590 0.20398 NaN 1.008718 2000-01-03 NaN NaN 0.393688 NaN 2000-01-04 NaN NaN 0.669350 0.811359 0.496025 2000-01-05 NaN 0.179829 NaN 2000-01-07 1.44470 NaN NaN NaN NaN
2000-01-08 NaN 1.365486 0.488785 NaN Car la vuelta a las columnas y a las filas df_ej5 = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D']) df_ej5.loc[::-1] #Esto de la vuelta a las columnas df.loc[::-1] #Esto da la vuelta a las filas
5 -1.422871 -0.442027 0.393688 -0.793887 6 1.565900 0.203998 -1.576848 1.008718 7 -0.666538 -0.000783 0.641780 1.889946 Mostrar los valores de la columna B que son más grandes que la columna A y más pequeños que la columna C dates= pn.date_range('1/1/2000', periods=8) df = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D'])
print(df) df.loc[(df.B>df.A) & (df.B <df.c), #paréntesis="" #pongo="" #si="" #tienes="" &,="" 'b']="" -0.114387="" -0.123124="" -0.181017="" -0.278824="" -0.329489="" -0.484974="" -1.249420="" -1.251315="" 0.208602="" 0.344837<="" 0.591778="" 0.852016="" 0.916072="" 1.019449="" 1.027612="" 1.807709="" 2000-01-01="" 2000-01-02="" 2000-01-03="" 2000-01-04="" a="" and="" b="" c="" cada="" condiciones="" d="" el="" entre="" es="" no="" or="" pero="" poner="" que="" quieres="" será="" su-="" sé="" td="" una="" varias="" y="" =""></df.c),>
2000-01-05 -1.315916 -0.477611 -0.173063 -0.080366 2000-01-06 -0.776168 -1.187224 0.489135 1.172228 2000-01-07 0.190130 -1.369850 -1.268442 0.727483 2000-01-08 -0.529417 0.210373 -1.742735 1.469252 2000-01-05 -0.477611 Freq: D, Name: B, dtype: float64 SELECTION BY DF.QUERY(CONDITION) (muy util) lataframe.query('condicion') es una herramienta muy útil que nos muestra los valores que cumplen la condición que se indique, simpre comparando por columnas. Es decir, si ponemos df.query('(a>b)') nos mostrará por filas los valores de la columna a que sean mayores que b, quivalente a df.loc[df.A>df.B,:]
DATA MANIPULATION IN PANDAS RENAMING lataframe.rename() permite renombrar índices o columnas, y ello se hace a través de un diccionario. Es decir, si queremos por ejemplo renombrar las columnas de la tabla df: df.rename(columns={'A':'one','B':'Two','C':'Three','D':'Four'}) Three Four
2000-01-01 -1.251315 1.027612 -0.278824 -0.329489 2000-01-02 0.208602 1.019449 0.916072 -0.114387 2000-01-03 -0.484974 1.807709 -0.181017 0.591778 2000-01-04 -1.249420 0.852016 -0.123124 0.344837 2000-01-05 -1.315916 -0.477611 -0.173063 -0.080366 2000-01-06 -0.776168 -1.187224 0.489135 1.172228 2000-01-07 0.190130 -1.369850 -1.268442 0.727483
2000-01-08 -0.529417
b -1.549842
REINDEXING lataframe.reindex([nuevo_indice_1,nuevo_indice_2]) Permite cambiar o añadir nuevos índices, ESA ES LA DIFERENCIA RESPECTO DE DF.RENAME(), QUE PUEDES AÑADIR NUEVOS ÍNDICES EN LA TUPLA Y ÉSTOS SE AÑADIRÁN EN ORDEN DE APARICIÓN. Tambunciona para columnas, pero es necesario especificarlo con dataframe.reindex(columns=[nueva_columna_1,nueva_columna_2]) a=df.reindex([1,2,3,4,5,6,7,8,9,10]) a B C D
1 0.94080 0.90467 -0.88098 -0.29490 2 -0.54982 0.04139 0.31781 -0.750837 3 -0.08279 0.05381 -0.06370 1.463687 4 -0.372310 0.28445 1.35298 1.34153 5 -0.98183 -0.08574 -0.29402 -0.29403 -0.33710 -0.33710 -0.33710 -0.33710 -0.33710 -0.31833 7 0.197687 -0.172664 1.29632 0.318335
8 -1.057785
ratiendo de la siguiente tabla: df_drop=pn.DataFrame(np.arange(16).reshape((4,4)), index=['Ohio', 'Colorado', 'Utah', 'New York'], columns=['one', 'two', 'three', 'four']) variable
Utah 8 9 10 11 New York 12 13 14 15 df_drop.drop(['Ohio', 'Colorado']) one two three four Utah 8 9 10 11 New York 12 13 14 15
New York 12 13 14 15 Eliminamos las columnas 'two' y 'four': df_drop.drop(['two', 'four'], axis=1)
Utah 8 10 New York 12 14 ARITHEMTIC AND DATA ALIGNMENT Si queremos hacer operaciones aritméticas sobre 2 tablas y juntarlas, por ejemplo si quieremos sumar 2 tablas, hay que tener en cuenta que si hacemos df1+df2 ambas tablas deben tener el mismo tamaño, porque los valores que no tengan una pareja tomarán el valor de NaN. Polara evitar ésto tenemos los operados aritméticos de pandas: from TPython.display import Image
Image("operadores_aritmeticos.png") Method Description add Method for addition (+) sub Method for subtraction (-) div Method for division (/) mul Method for multiplication (*)
<pre>df1=pn.DataFrame(np.arange(12.).reshape((3,4)),columns=list('abcd')) df2=pn.DataFrame(np.arange(20.).reshape((4,5)),columns=list('abcde')) print(df1.add(df2, fill_value=0)) print('') print(df1.sub(df2, fill_value=0)) print('') print(df1.mul(df2, fill_value=0)) print(df1.mul(df2, fill_value=0)) print(df1.div(df2, fill_value=0))</pre> #Es importante el fill_value=0, porque si no lo pones es como si
#hicieras por ejemplo df1+df2 a b c d e 0 0.0 2.0 4.0 6.0 4.0 1 9.0 11.0 13.0 15.0 9.0 2 18.0 20.0 22.0 24.0 14.0 3 15.0 16.0 17.0 18.0 19.0 a b c d e 0 0.0 0.0 0.0 0.0 0.0 -4.0 1 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0 -9.0 2 -2.0 -2.0 -2.0 -2.0 -2.0 -2.0 -2.0 -14.0
2 -2.0 -2.0 -2.0 -2.0 -2.0 -2.0 -2.0 -3.0 3 -15.0 -16.0 -17.0 -18.0 -19.0 a b c d e
2 0.8 0.818182 0.833333 0.846154 0.0 FUNCTION APPLICATION AND MAPPING Idataframe.apply(funcion)> Se ejecuta la funcion entre paréntesis usando como input el dataframe especificado, PERO SE EJECUTA POR COLUMNAS Idataframe.apply(funcion, axis=1)> Se ejecuta la funcion entre paréntesis usando como input el dataframe especificado, PERO SE EJECUTA POR FILAS f= lambda x: x.max()-x.min() #La funcion frame=pn.DataFrame(np.random.randn(4,3), columns=['b','d','e'], index=['Utah','Ohio','Texas','Oregon'])
print(frame.apply(f)) #Se busca el mayor valor de la columna y se le resta el #menor de la misma, siendo el ouput una serie con cada #columna print(frame.apply(f, axis=1)) #Se busca el mayor valor de la fila y se le resta el #menor de la misma, siendo el ouput una serie con cada #fila b 2.576337 d 2.809652
e 1.484490 dtype: float64 Utah 1.740991 Ohio 3.167368 Texas 0.589877 Oregon 1.845711 dtype: float64 SORTING AND RANKING df5=pn.DataFrame({'a':[0,1,0,1], 'b':[4,7,-3,2]},index=[0,1,2,3])
sorting df5.sort_values(by='a') #Ordenar respecto a la columna 'a' a b 0 0 4 2 0 -3 1 1 7
3 1 2 df5.sort_values(by=['a', 'b']) #Ordenar respecto a las columnas 'a' y 'b' a b 2 0 -3 0 0 4 3 1 2
tanking El ranking asigna un valor de 1 al número de datos no nulos en una serie/tabla. El ranking se puede aplicar a una serie o a un dataframe. Existen distintos métodos de hacer ranking, que se marcan por .rank(method=X) from IPython.display import Image Image ("metodos_ranking.png")
Method Description 'average' Default: assign the average rank to each entry in the equal group 'min' Use the minimum rank for the whole group 'max' Use the maximum rank for the whole group 'first' Assign ranks in the order the values appear in the data obj=pn.Series([7, -5,7,4,2,0,4])
obj.rank() 0 6.5 1 1.0 2 6.5 3 4.5 4 3.0 5 2.0 6 4.5 dtype: float64 También podemos hacer un ranking de dataframes, pero éste se puede hacer por filas o por columnas. Por defecto se hará por filas, por lo que deberemos especificar con axis=1 para que se haga por columnas:
frame=pn.DataFrame({'b':[4.3,7,-3,2],'a':[0,1,0,1],'c':[-2,5,8,-2.5]}) print(frame) print(frame.rank(method='first', axis=1)) #Ranking por columnas print(frame.rank(method='first')) #Ranking por filas
3 2.0 1 -2.5
EJERCICIO 1: Write a Pandas program to replace all the NaN values with Zero's in a column of a dataframe. df = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D']) df_nan=df[df>0] #CREAMOS LA TABLA QUE SERVIRA DE ENUNCIADO df_nan[df_nan.isnull()]=0 df_nan A B C D
2000-01-01 0.312950 0.422464 0.223764 0.000000 2000-01-02 0.00000 0.00000 0.00000 0.00000 0.00000 2000-01-03 0.978438 0.402716 0.00000 0.00000 0.00000 2000-01-04 0.00000 0.637970 0.535376 0.00000 2000-01-05 0.00000 0.00000 0.00000 0.00000 0.00000 2000-01-06 0.102181 0.00000 1.074017 0.854179
2000-01-07 0.875921 0.00000 0.887574 1.533723 2000-01-08 0.00000 1.122834 0.247540 0.00000 EJERCICIO 2: Write a Pandas program to convert index in a column of the given dataframe. df_ej2 = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D']) print(df_ej2)
df_ej2['nueva_columna']=df_ej2.index df_ej2 A B C D 2000-01-01 -0.535101 -0.364561 -0.530170 -0.171163 2000-01-02 -0.680146 -0.588430 -0.364841 0.138270 2000-01-03 -0.489039 -0.165032 -1.735786 -1.570797 2000-01-04 1.303573 0.506258 1.221943 0.017820 2000-01-05 -2.071717 -0.272240 -0.364922 0.951685 2000-01-06 -1.372558 -0.763122 1.137094 1.035435 2000-01-07 -0.931777 0.199696 -0.14288 2.290485 2000-01-08 -0.784667 0.038746 1.117412 1.055123
A B C D Invaluation 2000-01-02 -0.53510 -0.35450 -0.53510 -0.50510 -0.57510 <td< td=""></td<>
2000-01-07 -0.931777 0.199696 -0.144288 2.290485 2000-01-07 2000-01-08 -0.784667 0.038746 1.117412 1.055123 2000-01-08 EJERCICIO 3: Write a Pandas program to count the NaN values in one or more columns in DataFrame df_ej3 = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D']) df_ej3=df_ej3[df_ej3>0] print(df_ej3)
<pre>sum=0 df_bool=df_ej3.isnull() for i in df_bool.columns: for j in df_bool.loc[j,i]==True: sum+=1 sum A B C D 2000-01-01 0.138290 1.947074 1.389102 NaN 2000-01-02 0.226521 0.858733 NaN 0.739855</pre>
2000-01-03 0.039906 NaN NaN 1.275305 2000-01-04 NaN NaN 1.191685 0.767562 2000-01-05 NaN 0.144297 2.020633 NaN 2000-01-06 NaN 0.467210 NaN NaN 2000-01-06 NaN 0.467210 NaN NaN 2000-01-07 0.210746 0.600205 NaN 0.808142 2000-01-08 0.300563 0.996810 NaN 0.788788 13 EJERCICIO 4: Write a Pandas program to add a prefix or suffix to all columns of a given DataFrame
df_ej5 = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D']) prefijo='A_' columnas=df_ej4.columns pref_col=prefijo+columnas pref_col=pref_col.values df_ej4_rename=df_ej4.reindex(columns=pref_col) #Se haría así pero el reindex no me va para columnas EJERCICIO 5: Write a Pandas program to reverse order (rows, columns) of a given DataFrame df_ej5 = pn.DataFrame(np.random.randn(8,4), index=dates, columns=['A','B','C','D'])
df_ej5 = pn.DataFrame(np.random.rando
2000-01-03 1.054368 0.086818 0.725416 -1.221532 2000-01-02 -0.356919 1.301396 -0.356771 -0.805650 2000-01-01 -1.703022 -0.537547 -0.221594 -0.154773 2000-01-01 -1.70302 -0.154773 2000
inicio=input('Año inicio: ') fin=input('Año fin: ') df_2=df.loc[inicio:fin] df_2['Descontado']=df_2.values-(0.1*df_2.values) df_2 df_2 año inicio: 2000 Año fin: 2005 ipython-input-69-76549d498028>:7: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row_indexer,col_indexer] = value instead
Try using .loc[row_indexer,col_indexer] = value instead See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy df_2['Descontado']=df_2.values-(0.1*df_2.values) Ventas Descontado 2000 13434 12090.6
2001 2000 18000.0 2002 3000 2700.0 2003 3455565 3110008.5 2004 23254 20928.6
2002 3000 2700.0 2003 3455565 3110008.5 2004 23254 20928.6 2005 234543 211088.7
2002 3000 2700.0 2008 3455555 3110008.5 2004 23254 20828.6 2005 234543 211098.7 CJERCICIO 7: Los ficheros emisiones-2016.csv, emisiones-2017.csv, emisiones-2018.csv y emisiones-2019.csv, contienen datos sobre las emisiones contaminates en la ciudad de Madrid en los años 2016, 2017, 2018 y 2019 respectivamente. Escribir un programa con los siguient equisitos: • Generar un DataFrame con los datos de los cuatro ficheros. • Filtrar las columnas del DataFrame para quedarse con las columnas ESTACION, MAGNITUD, AÑO, MES y las correspondientes a los días D01, D02, etc. • Reestructurar el DataFrame para que los valores de los contaminantes de las columnas de los días aparezcan en una única columna. • Añadir una columna con la fecha a partir de la concatenación del año, el mes y el día (usar el módulo datetime). • Eliminar las filas con fechas no válidas (utilizar la función isnat del módulo numpy) y ordenar el DataFrame por estaciones, contaminantes y fecha. • Mostrar por pantalla las estaciones y los contaminantes disponibles en el DataFrame.
2002 3000 3000 5000 5000 5000 5000 5000
200
seeds 1905 1905 1905 1905 1905 1905 1905 1905
1
Part
Manual
Manual M
Manual M
Page
Margine Control of State Control of Stat
Margin M
Mark Land Comment of the Comment of
Margin M
2