

Hecho un dataframes, Importa y Exporta los datos

Andrew Reeve
School of Earth and Climate Sciences
University of Maine

Qué es Pandas (Panel Data)?

- Tipos de datos nuevos
 - Series: Una-dimensión con indices
 - DataFrame: multi-dimensión con indices
 - estructura como hoja de cálculo
- Añadir nombres a los filas y columnas
- asignación inteligente de NaN's
- matemáticas optimizado (eg. rolling averages, interpolation)
- entrada/salida de los datos de archivo
- graficas con matplotlib & seaborn

Series

- hecho con método de Series
 - 1-D conjuntos de datos, como vector de numpy
 - puede dar indices (nombres de filas)
- hecho de:
 - lists de valores y opcional indices (también una lista)
 - diccionario (con clave para cada elemento)
- puede dar un nombre a una series
- puede combinar los series
- extrear datos de una series
 - loc (locación): usa nombre de indices
 - iloc: usa la posicional indices

```
import pandas as pd
2 from string import ascii_letters
4 a_list = [100+i for i in range(10)]
b_list = [100+i for i in range(8,16)]
6
a_index=[ascii_letters[i] for i in range(10)]
8 b_dict={ ascii_letters[i]:b_list[n] for n,i in enumerate(range(8,16)) }
}
10 a_series=pd.Series(a_list, index=a_index)
b_series=pd.Series(b_dict)
12
a_series.name='a_name'
14 b_series.name='b_name'
# note how data handled when combining
16 add_series=a_series+b_series
concat_series=pd.concat((a_series,b_series))
18
# getting data from a series
20 print(a_series['a'])
# note that this included 'd'
22 print(a_series['a':'d'])
#location command, prefered way of slicing
24 print(a_series.loc['a':'d'])
# index location, using positional indices
26 # does not include '4'
print(a_series.iloc[0:4])
28 # get index
print(a_series.index)
30 # get values
print(a_series.values)
```

DataFrames

- metodo DataFrame (CamelCase)
 - 2-D (más posible)
 - filas y columnas con nombres
- hecho con diccionarios (y otras maneras)
 - clave es nombre de columna, listas tienen los valores
- añada columnas con assign
- cambiar indices
 - set_index: usa columna como idx & retira
 - df.index=values: sobrescribe idx
 - df.drop([list if row/column names]): retirar filas,columnas
 - predeterminado a filas, puede especificar axis

```
import pandas as pd
2 # some data
3 data_dict={ 'month':[ 'Jan','Feb','Mar','Apr','May','Jun','Jul','Aug','Sep','Oct',
4 'Nov','Dec'],
5 'Bgr_temp':[18,20,30,42,54,64,68,67,58,48,37,24],
6 'Bgr_ppt':[3.0,2.9,3.1,3.3,3.6,3.3,3.3,3.4,3.4,3.4,4.6,4.0],
7 'Blfst_temp':[21,24,34,44,54,64,68,67,58,48,38,26],
8 'Blfst_ppt':[3.5,3.3,3.9,4.4,4.3,3.5,3.1,3.3,3.6,4.3,5.4,4.9]}
9 # dataframe from series
10 # only difference is that *a_df* has a column heading
11 a_series=pd.Series(data_dict['Bgr_temp'])
12 # specify index
13 a_series=pd.Series(data_dict['Bgr_temp'], index=data_dict['month'])
14 a_df=pd.DataFrame(a_series)
15 # specify index using dict of lists
16 a_df=pd.DataFrame({ 'Brg_temp':data_dict['Bgr_temp']},index=data_dict['month'])
17 b_df=pd.DataFrame({ 'Brg_temp':data_dict['Bgr_temp'],
18 'Bgr_ppt':[3.0,2.9,3.1,3.3,3.6,3.3,3.3,3.4,3.4,3.4,4.6,4.0],},
19 index=data_dict['month'])
20 # convert entire dict into dataframe
21 c_df=pd.DataFrame(data_dict)
22 # adding a column to df
23 c_df=c_df.assign(month2=c_df.month.str.lower())
24 # setting index using column, removes column from df
25 c_df = c_df.set_index('month')
26 # overwriting index values with another index
27 c_df.index = c_df['month2']
28 # dropping a column, need to indicate column axis, row (axis=0) is default
29 c_df.drop(['month2'],axis=1)
```

Leer y Escribir Archivos

- Muchos módulos de `read_xxx`:
 - Leer archivo a dataframe
 - e.g., `read_csv`, `read_hdf`, `read_table`
 - Algunas mandas tienen muchísimas opciones confusas de palabras de clave
- Módulos de `to_xxx`
 - Ahorrar dataframe a archivo
 - e.g., `to_hdf`, `to_csv`, `to_json`
- Debe limpiar los archivos antes de ahorrar
 - Se cae filas/columnas no necesita
 - Cambia tipos de datos de columnas (`astype`)
 - Cambia nombres de indice/columna (`set_index`, `index`, `column`)

```
import pandas as pd
2 # reading a csv, there are a confusingly large number of keywords for
   this
4 df1=pd.read_csv('./examples/CB_Yellow_Rod_Green_75.csv',
                  skiprows=47, skipfooter=2, engine='python',
                  parse_dates=[ 'TimeStamp'])
6 # saving dataframe to a file
7 df1.to_csv('./examples/ibutton_frame.csv.gz',compression='gzip')
8 df1.to_hdf('./examples/ibutton_frame.hdf','root', mode='w', complib =
   'blosc')

10 # reading a stream gauging file, USGS table format (tab delimited)
11 df2=pd.read_table('./examples/USGS01034000.tab',comment='#')
12 # cleaning up the file
13 df2=df2.drop([0])
14 df2=df2.drop(['agency_cd', 'site_no', 'tz_cd', '65543_00060_cd',
   '65544_00065_cd'],
15           axis=1)
16 df2=df2.set_index('datetime')
17 #check data types, they are strings! Need to convert
18 df2.loc[:, '65544_00065']=df2.loc[:, '65544_00065'].astype('float')
20 df2.loc[:, '65543_00060']=df2.loc[:, '65543_00060'].astype('float')
```

Combinar Dataframes

- diferentes mandas para combinar
df: concat, merge, join
 - concat 'apila' los dataframes
 - merge combina usando columna en común entre dataframes
 - join combina usando indices
- puede fabrica una dataframe vacío y añade columnas/filas
- multi-index para organizar los datos

```
import pandas as pd
2 import numpy as np
from string import ascii_uppercase,
                  ascii_lowercase
4
## make some dataframes
6 list1 = list(ascii_uppercase)
list2 = list(ascii_lowercase)
8
df1 = pd.DataFrame({"upper": list1[0:12], "lower":
                     : list2[0:12]})
10 df2 = pd.DataFrame({"upper": list1[9:26], "lower":
                      : list2[9:26]},
                         index=range(9, 26))
12
## concatenate frames
14 df3 = pd.concat((df1, df2))
df4 = pd.concat((df1, df2), axis=1)
```

Más Combina

```
import pandas as pd
2 import numpy as np
from string import ascii_uppercase,
    ascii_lowercase
4
## make some dataframes
6 list1 = list(ascii_uppercase)
list2 = list(ascii_lowercase)
8
df1 = pd.DataFrame({ "upper": list1[0:12], "
    "lower": list2[0:12]})
10 df2 = pd.DataFrame({ "upper": list1[9:26], "
    "lower": list2[9:26]},
        index=range(9, 26))
12
## concatenate frames
14 df3 = pd.concat((df1, df2))
df4 = pd.concat((df1, df2), axis=1)
```

```
import pandas as pd
2 import numpy as np
from string import ascii_uppercase,
    ascii_lowercase
4
## make some dataframes
6 list1 = list(ascii_uppercase)
list2 = list(ascii_lowercase)
8
df1 = pd.DataFrame({ "upper": list1[0:12], "
    "lower": list2[0:12]})
10 df2 = pd.DataFrame({ "upper": list1[9:26], "
    "lower": list2[9:26]},
        index=range(9, 26))
12
## concatenate frames
14 df3 = pd.concat((df1, df2))
df4 = pd.concat((df1, df2), axis=1)
```

Converting Datatypes and Null Values

- Los datos lean como cadenas a pandas (e.g. archivos de csv)
- Convierte a int, float or datetime
 - astype convierte los tipos
 - pandas.to_datetime(dataframe) convierte a objetos de datetime
 - pd.to_numeric(dataframe) convierte a números
 - errors palabra de clave, para manejar errores (coerce, ignore, raise)
- probar NaNs
 - funciones de pandas para ID valores 'null'
 - isna o isnull, notna o notnull
 - any y all prueban para valores True
- dropna retira valores 'null'
- fillna llena valores de 'null' valores con otros valores

```
import pandas as pd
2 import numpy as np
# all columns read into dataframe as 'objects', an object datatype
# in numpy and pandas is a string with variable lengths
4 df = pd.read_csv("annual_peak_flow.txt", comment="#",
                   delimiter="\t")
6 # can see dtypes fromall columns using:
df.dtypes # with 's'
8 # or individual colum dtype with:
df.peak_dt.dtype # no 's'
10 # remove first row
df = df.drop([0,])
12 # remove uneeded columns
df = df.drop(["agency_cd", "peak_tm", "year_last_pk", "ag_dt",
              "ag_tm", "ag_gage_ht_cd", "peak_cd", "gage_ht_cd"], axis=1)
14 # raises an error, a date is '2001-00-00', makes no sense
# df.loc[:, 'peak_dt']=pd.to_datetime(df.loc[:, 'peak_dt'])
18 # ID problem dates and overwrite 'bad' data with NA values
20 date_bool = df.peak_dt.str[-2:] == "00"
df.loc[date_bool, "peak_dt"] = pd.NA
22 # convert column data type
df.loc[:, "peak_dt"] = pd.to_datetime(df.loc[:, "peak_dt"]).astype("datetime64[ns]")
24 # easier ? way of converting to datetime
# df.loc[:, 'peak_dt']=pd.to_datetime(df.loc[:, 'peak_dt'], errors='coerce')
26 df.loc[:, "peak_va"] = pd.to_numeric(df.loc[:, "peak_va"], errors="coerce")
```

Selecciona los datos de DataFrame

- para sacar una parte de dataframe
 - `df.loc[rows,cols]` seleccione por nombres de filas/col
 - `df.iloc[rows,cols]` seleccione por indices posicional
 - combina estos con operaciones booleanas
- mira rápida en DataFrame
 - `df.head(rows)` vuelve primeras filas
 - `df.tail(rows)` vuelve ultimas filas
 - `df.info()` info básico de DataFrame
 - `df.describe()` estadísticas básicas sobre columnas

```
import pandas as pd
2 import numpy as np

4 adict={ 'a':[1,2,3], 'b':['x','y','z']}
df=pd.DataFrame(adict, index=['a0','a1','a2'])
6
# getting just values of 1 and 2 from dataframe
8 # method 1, slicing with loc
m1=df.loc[['a1','a2'], 'a']
10
# method 2, slicing with iloc
12 m2=df.iloc[[1,2], 0]
14
# method 3, using booleans
16 m3 = df.loc[[False,True,True], 'a']

18 # method 4, using boolean expression
# these do the same thing
20 m4=df.loc[df.b!='x', 'a']
m4=df.loc[df.a>1, 'a']
22 # replacing a value using boolean operation
df2=df.copy()
24 df2.loc[df2.a>1, 'a'] = df2.loc[df2.a>1, 'a']+20
df2.loc[df2.a>1, 'a'] = df2.loc[df2.a<=1, 'a']+10
```