# 04. DataFrameDataStructure_ed

May 10, 2022

The DataFrame data structure is the heart of the Panda's library. It's a primary object that you'll be working with in data analysis and cleaning tasks.

The DataFrame is conceptually a two-dimensional series object, where there's an index and multiple columns of content, with each column having a label. In fact, the distinction between a column and a row is really only a conceptual distinction. And you can think of the DataFrame itself as simply a two-axes labeled array.

```python
[10]:  # Lets start by importing our pandas library
       import pandas as pd
```

```python
[11]:  # I'm going to jump in with an example. Lets create three school records for
       ↪students and their
       # class grades. I'll create each as a series which has a student name, the
       ↪class name, and the score.


       # Creamos 3 series a partir de 3 diccionarios
       record1 = pd.Series({'Name': 'Alice',
                            'Class': 'Physics',
                            'Score': 85})
       record2 = pd.Series({'Name': 'Jack',
                            'Class': 'Chemistry',
                            'Score': 82})
       record3 = pd.Series({'Name': 'Helen',
                            'Class': 'Biology',
                            'Score': 90})

       # Imprimo 3 series
       print(record1,"\n")
       print(record2,"\n")
       print(record3)
```

```
Name      Alice
Class    Physics
Score         85
dtype: object

Name       Jack
```

```
Class      Chemistry
Score            82
dtype: object

Name        Helen
Class     Biology
Score          90
dtype: object
```

[18]:
```python
# Like a Series, the DataFrame object is index. Here I'll use a group of
 →series, where each series
# represents a row of data. Just like the Series function, we can pass in our
 →individual items
# in an array, and we can pass in our index values as a second arguments


# Pordemos crear un Dataframe sumando 3 series, cada serie será una fila del df
# acá no pondemos index y crea automáticamente uno con números enteros
df1 = pd.DataFrame([record1, record2, record3])

# Además podemos definir los labels del index
df = pd.DataFrame([record1, record2, record3],
                  index=['school1', 'school2', 'school1'])

# And just like the Series we can use the head() function to see the first
 →several rows of the
# dataframe, including indices from both axes, and we can use this to verify
 →the columns and the rows

print(df1,'\n') # sin índices
print(df.head()) # con índex labels, en este caso el head() no hace nada porque
 →hay solamente 3 filas
```

```
      Name      Class  Score
0   Alice     Physics     85
1    Jack  Chemistry     82
2   Helen    Biology     90

          Name      Class  Score
school1  Alice     Physics     85
school2   Jack  Chemistry     82
school1  Helen    Biology     90
```

[4]:
```python
# You'll notice here that Jupyter creates a nice bit of HTML to render the
 →results of the
# dataframe. So we have the index, which is the leftmost column and is the
 →school name, and
```

```python
# then we have the rows of data, where each row has a column header which was␣
 ↪given in our initial
# record dictionaries
```

```python
[21]: # An alternative method is that you could use a list of dictionaries, where␣
       ↪each dictionary
      # represents a row of data.

      # Otra forma de crear un df es con una lista de diccionarios, cada diccionario␣
       ↪es una serie.

      students = [{'Name': 'Alice',
                   'Class': 'Physics',
                   'Score': 85},
                  {'Name': 'Jack',
                   'Class': 'Chemistry',
                   'Score': 82},
                  {'Name': 'Helen',
                   'Class': 'Biology',
                   'Score': 90}]

      # Then we pass this list of dictionaries into the DataFrame function
      # También definimos los labels de cada serie

      df = pd.DataFrame(students, index=['school1','school2', 'school1'])

      # And lets print the head again
      df.head()
```

```
[21]:          Name      Class  Score
      school1  Alice    Physics     85
      school2   Jack  Chemistry     82
      school1  Helen    Biology     90
```

```python
[22]: # Similar to the series, we can extract data using the .iloc and .loc␣
       ↪attributes. Because the
      # DataFrame is two-dimensional, passing a single value to the loc indexing␣
       ↪operator will return
      # the series if there's only one row to return.

      # For instance, if we wanted to select data associated with school2, we would␣
       ↪just query the
      # .loc attribute with one parameter.
      df.loc['school2']
```

```
[22]:  Name           Jack
       Class      Chemistry
       Score            82
       Name: school2, dtype: object
```

```
[14]:  # You'll note that the name of the series is returned as the index value, while␣
        ↪the column
       # name is included in the output.

       # We can check the data type of the return using the python type function.

       # como solo hay una fila con este label se obtiene un objeto del tipo series
       type(df.loc['school2'])
```

```
[14]:  pandas.core.series.Series
```

```
[16]:  # It's important to remember that the indices and column names along either␣
        ↪axes horizontal or
       # vertical, could be non-unique. In this example, we see two records for␣
        ↪school1 as different rows.
       # If we use a single value with the DataFrame lock attribute, multiple rows of␣
        ↪the DataFrame will
       # return, not as a new series, but as a new DataFrame.

       # Lets query for school1 records

       # Acá como se repite la label school1, se obtiene como resultado un objeto del␣
        ↪tipo Dataframe
       df.loc['school1']
```

```
[16]:            Name     Class  Score
       school1  Alice   Physics     85
       school1  Helen  Biology     90
```

```
[18]:  # And we can see the the type of this is different too
       # en este caso la selección del label school1 devuelve una matriz, por ende es␣
        ↪un dataframe
       type(df.loc['school1'])
```

```
[18]:  pandas.core.frame.DataFrame
```

```
[19]:  # One of the powers of the Panda's DataFrame is that you can quickly select␣
        ↪data based on multiple axes.
       # For instance, if you wanted to just list the student names for school1, you␣
        ↪would supply two
       # parameters to .loc, one being the row index and the other being the column␣
        ↪name.
```

```python
# For instance, if we are only interested in school1's student names

# Podemos seleccionar filas y columnas de interés con .loc
df.loc['school1', 'Name']
```

```
[19]: school1    Alice
      school1    Helen
      Name: Name, dtype: object
```

```python
# Remember, just like the Series, the pandas developers have implemented this␣
 ↪using the indexing
# operator and not as parameters to a function.

# What would we do if we just wanted to select a single column though? Well,␣
 ↪there are a few
# mechanisms. Firstly, we could transpose the matrix. This pivots all of the␣
 ↪rows into columns
# and all of the columns into rows, and is done with the T attribute

# otra forma de seleccionar una columna de interés es transponiendo la matriz y␣
 ↪con .loc seleccionar la fila (que antes era
# la columna de interés)
df.T
```

```
[23]:        school1     school2  school1
      Name     Alice        Jack    Helen
      Class   Physics  Chemistry  Biology
      Score        85         82       90
```

```python
# Then we can call .loc on the transpose to get the student names only
df.T.loc['Name']
```

```
[12]: school1    Alice
      school2     Jack
      school1    Helen
      Name: Name, dtype: object
```

```python
# However, since iloc and loc are used for row selection, Panda reserves the␣
 ↪indexing operator
# directly on the DataFrame for column selection. In a Panda's DataFrame,␣
 ↪columns always have a name.
# So this selection is always label based, and is not as confusing as it was␣
 ↪when using the square
# bracket operator on the series objects. For those familiar with relational␣
 ↪databases, this operator
```

```
# is analogous to column projection.


# Para seleccionar columnas en Pandas, lo más fácil es poner entre corchetes el␣
↪nombre de la columna de interés
df['Name']
```

[25]: school1    Alice
      school2     Jack
      school1    Helen
      Name: Name, dtype: object

```
[28]: # In practice, this works really well since you're often trying to add or drop␣
      ↪new columns. However,
      # this also means that you get a key error if you try and use .loc with a␣
      ↪column name


      # Esto da error porque busca una fila llamada 'Name' y ese es el nombre de una␣
      ↪columna!
      # Para que funciones podés hacer df.loc[:, 'Name'] que selecciona de todas las␣
      ↪filas la columna "Name"
      df.loc['Name']
```

```
        ␣
↪---------------------------------------------------------------------------

       KeyError                                  Traceback (most recent call␣
↪last)

       /usr/local/lib/python3.8/site-packages/pandas/core/indexes/base.py in␣
↪get_loc(self, key, method, tolerance)
    2896             try:
  -> 2897                 return self._engine.get_loc(key)
    2898             except KeyError:


        pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()


        pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()


        pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.
↪_get_loc_duplicates()
```

```
    pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.
↪_maybe_get_bool_indexer()


    KeyError: 'Name'


  During handling of the above exception, another exception occurred:


    KeyError                                  Traceback (most recent call␣
↪last)

    <ipython-input-28-6535e7d2d9fb> in <module>
      4
      5 # Esto da error porque busca una fila llamada 'Name' y ese es el␣
↪nombre de una columna!
  ----> 6 df.loc['Name']


    /usr/local/lib/python3.8/site-packages/pandas/core/indexing.py in␣
↪__getitem__(self, key)
    1422
    1423              maybe_callable = com.apply_if_callable(key, self.obj)
 -> 1424              return self._getitem_axis(maybe_callable, axis=axis)
    1425
    1426      def _is_scalar_access(self, key: Tuple):


    /usr/local/lib/python3.8/site-packages/pandas/core/indexing.py in␣
↪_getitem_axis(self, key, axis)
    1848          # fall thru to straight lookup
    1849          self._validate_key(key, axis)
 -> 1850          return self._get_label(key, axis=axis)
    1851
    1852


    /usr/local/lib/python3.8/site-packages/pandas/core/indexing.py in␣
↪_get_label(self, label, axis)
     158              raise IndexingError("no slices here, handle elsewhere")
     159
 --> 160          return self.obj._xs(label, axis=axis)
     161
     162      def _get_loc(self, key: int, axis: int):
```

```
        /usr/local/lib/python3.8/site-packages/pandas/core/generic.py in␣
    ↪xs(self, key, axis, level, drop_level)
        3735                loc, new_index = self.index.get_loc_level(key,␣
    ↪drop_level=drop_level)
        3736            else:
     -> 3737                loc = self.index.get_loc(key)
        3738
        3739                if isinstance(loc, np.ndarray):


        /usr/local/lib/python3.8/site-packages/pandas/core/indexes/base.py in␣
    ↪get_loc(self, key, method, tolerance)
        2897                    return self._engine.get_loc(key)
        2898                except KeyError:
     -> 2899                    return self._engine.get_loc(self.
    ↪_maybe_cast_indexer(key))
        2900            indexer = self.get_indexer([key], method=method,␣
    ↪tolerance=tolerance)
        2901            if indexer.ndim > 1 or indexer.size > 1:


        pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()


        pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()


        pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.
    ↪_get_loc_duplicates()


        pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.
    ↪_maybe_get_bool_indexer()


        KeyError: 'Name'
```

```
[29]: # Note too that the result of a single column projection is a Series object
      type(df['Name'])
```

```
[29]: pandas.core.series.Series
```

```
[31]: # Since the result of using the indexing operator is either a DataFrame or␣
       ↪Series, you can chain
```

```
# operations together. For instance, we can select all of the rows which␣
 ↪related to school1 using
# .loc, then project the name column from just those rows

# Acá se puede hacer chaining de comandos porque df.loc['school1'] devuelve un␣
 ↪df al cual se le puede aplicar ['Name'] para
# que devuelva la columna con ese label
df.loc['school1']['Name']
```

[31]:
```
school1     Alice
school1     Helen
Name: Name, dtype: object
```

[32]:
```
# If you get confused, use type to check the responses from resulting operations
print(type(df.loc['school1'])) #should be a DataFrame
print(type(df.loc['school1']['Name'])) #should be a Series
```

```
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
```

[ ]:
```
# Chaining, by indexing on the return type of another index, can come with some␣
 ↪costs and is
# best avoided if you can use another approach. In particular, chaining tends␣
 ↪to cause Pandas
# to return a copy of the DataFrame instead of a view on the DataFrame.
# For selecting data, this is not a big deal, though it might be slower than␣
 ↪necessary.
# If you are changing data though this is an important distinction and can be a␣
 ↪source of error.
```

[27]:
```
# Here's another approach. As we saw, .loc does row selection, and it can take␣
 ↪two parameters,
# the row index and the list of column names. The .loc attribute also supports␣
 ↪slicing.

# If we wanted to select all rows, we can use a colon to indicate a full slice␣
 ↪from beginning to end.
# This is just like slicing characters in a list in python. Then we can add the␣
 ↪column name as the
# second parameter as a string. If we wanted to include multiple columns, we␣
 ↪could do so in a list.
# and Pandas will bring back only the columns we have asked for.

# Here's an example, where we ask for all the names and scores for all schools␣
 ↪using the .loc operator.
```

```
# Para seleccionar numerosas columnas de interés
# devuelve todas las filas de dos columnas de interés.
df.loc[:,['Name', 'Score']]
```

[27]:
```
         Name   Score
school1  Alice     85
school2   Jack     82
school1  Helen     90
```

[29]:
```
# otra forma de pedir todas las filas de dos columnas de interés:
# se ingresa una lista dentro de los corchetes: df[[lista]]
df[['Name','Score']]
```

[29]:
```
         Name   Score
school1  Alice     85
school2   Jack     82
school1  Helen     90
```

[ ]:
```
# Take a look at that again. The colon means that we want to get all of the
 ↪rows, and the list
# in the second argument position is the list of columns we want to get back
```

[ ]:
```
# That's selecting and projecting data from a DataFrame based on row and column
 ↪labels. The key
# concepts to remember are that the rows and columns are really just for our
 ↪benefit. Underneath
# this is just a two axes labeled array, and transposing the columns is easy.
 ↪Also, consider the
# issue of chaining carefully, and try to avoid it, as it can cause
 ↪unpredictable results, where
# your intent was to obtain a view of the data, but instead Pandas returns to
 ↪you a copy.
```

[30]:
```
# Before we leave the discussion of accessing data in DataFrames, lets talk
 ↪about dropping data.
# It's easy to delete data in Series and DataFrames, and we can use the drop
 ↪function to do so.
# This function takes a single parameter, which is the index or row label, to
 ↪drop. This is another
# tricky place for new users -- the drop function doesn't change the DataFrame
 ↪by default! Instead,
# the drop function returns to you a copy of the DataFrame with the given rows
 ↪removed.

# Se eliminan las filas con label school1, sin embargo devuelve una copia que
 ↪no afecta la matriz original
```

```
df.drop('school1')
```

[30]:
```
         Name       Class  Score
school2  Jack  Chemistry     82
```

[31]:
```
# But if we look at our original DataFrame we see the data is still intact.
df
```

[31]:
```
          Name       Class  Score
school1  Alice     Physics     85
school2   Jack  Chemistry     82
school1  Helen    Biology     90
```

[33]:
```
# Drop has two interesting optional parameters. The first is called inplace,␣
↪and if it's
# set to true, the DataFrame will be updated in place, instead of a copy being␣
↪returned.
# The second parameter is the axes, which should be dropped. By default, this␣
↪value is 0,
# indicating the row axis. But you could change it to 1 if you want to drop a␣
↪column.

# For example, lets make a copy of a DataFrame using .copy()
copy_df = df.copy()

# Now lets drop the name column in this copy

# Si querés podés afectar la matriz original y seleccionar filas o columnas␣
↪para eliminar con drop() dado que
# tiene múltiples parámetros
# acá vamos a eliminar la columna (usamos parámetro axis=1) "Name" de la matriz␣
↪original copiada
# axis=0 es para eliminar filas

copy_df.drop("Name", inplace=True, axis=1)
copy_df
```

[33]:
```
             Class  Score
school1     Physics     85
school2  Chemistry     82
school1    Biology     90
```

[34]:
```
# acá eliminamos filas con label 'school1' con drop()
copy_df.drop('school1', inplace=True, axis=0)
copy_df
```

```
[34]:                Class  Score
      school2  Chemistry     82
```

```
[35]: # There is a second way to drop a column, and that's directly through the use␣
      ↪of the indexing
      # operator, using the del keyword. This way of dropping data, however, takes␣
      ↪immediate effect
      # on the DataFrame and does not return a view.


      # Otra forma para eliminar columnas de la dataframe original
      copy_df
```

```
[35]:          Score
      school2     82
```

```
[36]: # Finally, adding a new column to the DataFrame is as easy as assigning it to␣
      ↪some value using
      # the indexing operator. For instance, if we wanted to add a class ranking␣
      ↪column with default
      # value of None, we could do so by using the assignment operator after the␣
      ↪square brackets.

      # Asi se crean columnas con nuevos valores en todas las filas
      df['ClassRanking'] = None
      df
```

```
[36]:          Name      Class  Score ClassRanking
      school1  Alice    Physics     85         None
      school2   Jack  Chemistry     82         None
      school1  Helen    Biology     90         None
```

In this lecture you've learned about the data structure you'll use the most in pandas, the DataFrame. The dataframe is indexed both by row and column, and you can easily select individual rows and project the columns you're interested in using the familiar indexing methods from the Series class. You'll be gaining a lot of experience with the DataFrame in the content to come.