

Introduction

Pandas is a package built on top of NumPy, and provides an efficient implementation of many features :

- DataFrames
- Series
- Data Alignment
- Handling Missing Data
- Grouping and Aggregation
- Data Input and Output
- Handling Time Series

Pandas General Methods

Accessing values	<code>pd_ds.values</code>	DataFrame.values Series.values
Accessing Indices	<code>pd_ds.index</code>	DataFrame.index Series.index
Accessing specific element	<code>pd_ds[idx]</code>	DataFrame[1] Series[1]
Accessing range of elements	<code>pd_ds[start : end]</code>	DataFrame[1:4] Series[2:5]
Implicit Indexing	<code>df.iloc[rows , cols]</code>	data.iloc[1:3] #last index is not included
Explicit Indexing	<code>df.loc[rows,cols]</code>	data.loc['California' : 'Texas'] #last index is included

Pandas Series

Creating Series with lists	<code>pd.Series([values], index= list)</code>	<code>data = pd.Series([0.25, 0.5, 0.75, 1.0], index=['a', 'b', 'c', 'd'])</code>
Creating Series with dictionaries	<code>pd.Series({index:v-alue})</code>	<code>population_dict = {'California': 38332521, 'Texas': 26448193, 'New York': 19651127, 'Florida': 19552860, 'Illinois': 12882135}</code> <code>population = pd.Series(population_dict)</code>
Slicing Series	<code>Series[from_idx : to_idx]</code>	<code>population['Texas':'Florida']</code>
Slicing Indices with Dictionary Series	<code>pd.Series({index:value} , index=[])</code>	<code>pd.Series({'2':'a', 1:'b', 3:'c'}, index=[3, 2])</code> # only returns the third and second index respectfully

Pandas Series Index Can be a list of string or list of integers (or any desired type) unlike numpy arrays

Pandas DataFrames

Creating DataFrame	<code>pd.DataFrame({index : iterable})</code>	<code>pd.DataFrame({'population': population, 'area': area})</code>
---------------------------	---	---



By **Taissir Boukrouba**
(taissir2002)

Not published yet.
Last updated 19th November, 2023.
Page 1 of 7.

Sponsored by **CrosswordCheats.com**
Learn to solve cryptic crosswords!
<http://crosswordcheats.com>

Pandas DataFrames (cont)

Adding Column names	pd.DataFrame(dict , columns = [list_of_col_names])	pd.DataFrame(population, columns=['population'])
Slicing DataFrame Index	pd.DataFrame(dict , columns = [] , index = [])	pd.DataFrame(np.random.rand(3, 2), columns=['foo', 'bar'], index=['a', 'b', 'c'])
Reading CSV Files	pd.read_csv(source , index_col = col)	pd.read_csv("data/president_heights.csv", index_col="order")
Saving DF to CSV	dataframe.to_csv(source)	df.to_csv("data/president_heights_copy.csv")
Reading Excel Files	pd.read_excel(source)	pd.read_excel("data/president_heights.xlsx")
Saving DF to Excel	dataframe.to_excel(source)	df.to_excel("data/president_heights_copy.xlsx")
Access DataFrame Columns	dataframe.columns	df.columns
Transposing DataFrames	dataframe.T	df.T
Subsetting Using loc	dataframe.loc[condition , cols]	data.loc[data.density > 100, ['pop', 'density']]
Masking	dataframe[mask]	data[data.density > 100]

Pandas Index

Creating Index	pd.Index(list)	pd.Index([2, 3, 5, 7, 11])
Accessing Index	Index[idx]	ind[1]
Slicing Index	Index[from : to : step]	ind[: : 2]
Intersection Between Indices	index_1.intersection(index_2)	indA = pd.Index([1, 3, 5, 7, 9]) indB = pd.Index([2, 3, 5, 7, 11]) indA.intersection(indB)
Union Between Indices	index_1.union(index_2)	indA.union(indB)
Symmetric Difference	index_1.symmetric_difference(index_2)	indA.symmetric_difference(indB)

The **Index** has many of the attributes familiar from NumPy arrays such as :
ind.size, ind.shape, ind.ndim, ind.dtype

Pandas Universal Functions

+	add()
-	sub() , subtract()
*	mul() , multiply()
/	truediv() , div() , divide()
//	floordiv()
%	mod()



By **Taissir Boukrouba**
(taissir2002)

Not published yet.
Last updated 19th November, 2023.
Page 2 of 7.

Sponsored by **CrosswordCheats.com**
Learn to solve cryptic crosswords!
<http://crosswordcheats.com>

Pandas Universal Functions (cont)

** `pow()`

These universal functions are used in the following form :

- `data_struct.uf(data_struct_2)`
- `data_struct.uf()`

Datatype Conversions (NaN or None)

Float	No change
Object	No change
Integer	Upcast to float64
Boolean	Upcast to object

These are data type conversion when there is missing values

Operating On Missing Values

Nullability Check	<code>data_struct.is_null()</code>	<code>data = pd.Series([1, np.nan, 'hello', None])</code> <code>data.isnull()</code>
Non-Nullability Check	<code>data_struct.not_null()</code>	<code>data.not_null()</code>
Slicing Non-Null Values	<code>data_struct[data_struct.not_null()]</code>	<code>data[data.not_null()]</code>
Dropping Null Values	<code>data_struct.dropna(axis=0/1, how='any'/'all', thresh = n)</code>	<code>data.dropna(axis = 0, thresh = 2)</code> # the thresh means each row has at least 2 non-null values
Filling Missing Values	<code>data_struct.fillna(value, method='ffil'/'bfill', axis = 0/1)</code>	<code>df.fillna(method='ffil', axis=1)</code>
Filling Using A Function (Interpolation)	<code>data_struct.interpolate(method='linear'/'polynomial'...)</code>	<code>df.interpolate()</code> # the method is linear by default

When Working with missing values methods, axis = 0 means rows and 1 columns

Pandas Multi-Indexing

Creating Multi-Index From Tuples	<code>pd.MultiIndex.from_tuples(tuple)</code>	<code>index = pd.MultiIndex.from_tuples([('California', 2000), ('California', 2010)])</code>
Creating Multi-Index From Arrays	<code>pd.MultiIndex.from_arrays(list)</code>	<code>pd.MultiIndex.from_arrays(['a', 'a', 'b', 'b'], [1, 2, 1, 2])</code>
Creating Multi-Index From Product	<code>pd.MultiIndex.from_product([index1_list, index2_list])</code>	<code>pd.MultiIndex.from_product(['a', 'b'], [1, 2])</code>
Creating Multi-Index From DataFrame Values	<code>pd.MultiIndex.from_frame(dataframe)</code>	<code>df = pd.DataFrame(['a', 'b'], [1, 2])</code> <code>pd.MultiIndex.from_frame(df)</code>
Applying Multi-Index	<code>data_struct.reindex(index)</code>	<code>pop = pop.reindex(index)</code>



By Taissir Boukrouba
(taissir2002)

Not published yet.
Last updated 19th November, 2023.
Page 3 of 7.

Sponsored by CrosswordCheats.com
Learn to solve cryptic crosswords!
<http://crosswordcheats.com>

Pandas Multi-Indexing (cont)

Setting Index From Columns	<code>data_struct.set_index([cols])</code>	<code>pop_flat.set_index(['population'])</code>
Accessing Multi-Indexed Data Structures	<code>data_struct[first_index,second_index,....., col]</code>	<code>pop[:, 2010]</code> # gets all rows from first index and only 2010 rows from second index
Unstacking	<code>data_struct.unstack()</code>	<code>pop.unstack()</code> # this converts the last index (if we have 2 then the second one) values into cols
Stacking	<code>data_struct.stack()</code>	<code>pop.stack()</code> # this converts columns into a second index
Naming Multi-Indexes	<code>data_struct.index.names = list</code>	<code>pop.index.names = ['state', 'year']</code>
Swapping Multi-Indexes	<code>data_struct.swaplevel(0,1)</code>	<code>pop_df = pop_df.swaplevel(0,1)</code>
Dropping Multi-Indexes	<code>data_struct.droplevel(level=index)</code>	<code>pop_df.droplevel(level=0)</code>
Multi-Index In Columns	<code>pd.DataFrame(data, index=multi_index_rows, columns=multi_idx_cols)</code>	<code>columns = pd.MultiIndex.from_product(['Bob', 'Guido', 'Sue'], ['HR', 'Temp'], names=['subject', 'type'])</code> <code>health_data = pd.DataFrame(data, index=index, columns=columns)</code>
Slicing Using Multi-Index Column Values	<code>dataframe[multi_ind_col_value]</code>	<code>health_data['Guido']</code>
Slicing Multi-Index Cols & Rows Using IndexSlice	<code>idx = pd.IndexSlice</code> <code>df.loc[idx[index_row1,index_row2], idx[index_col1,index_col2]]</code>	<code>idx = pd.IndexSlice</code> <code>health_data.loc[idx[:, 1], idx[:, 'HR']]</code>
Resetting Multi-Index to Cols	<code>data_struct.reset_index()</code>	<code>pop.reset_index(name='population')</code>
Sorting Multi-Index	<code>data_struct.sort_index()</code>	<code>data.sort_index()</code>

It is a good practice to sort the values after swapping Multi-index Levels

Concatenation , Merging and Joins

Concatenation	<code>pd.concat([data_struct , data_struct2] , ignore_index = True/False)</code>	<code>pd.concat([ser1, ser2])</code>
Adding MultiIndex Keys	<code>pd.concat([data_struct , data_struct2] , keys = ['a','b'])</code>	<code>display('x', 'y', "pd.concat([x, y], keys=['x', 'y'])")</code>



By Taissir Boukrouba
(taissir2002)

Not published yet.
Last updated 19th November, 2023.
Page 4 of 7.

Sponsored by **CrosswordCheats.com**
Learn to solve cryptic crosswords!
<http://crosswordcheats.com>

Concatenation , Merging and Joins (cont)

Concatenation with Joins	<code>pd.concat([data_struct , data_struct2] , join = 'outer'/'-inner')</code>	<code>pd.concat([df5, df6], join='inner')</code> # The intersection of 2 DFs
Merging	<code>pd.merge(data_struct , data_struct2)</code>	<code>df3 = pd.merge(df1, df2)</code>
Merging on Columns	<code>pd.merge(data_struct , data_struct2 , on ='col_name')</code>	<code>pd.merge(df1, df2, on='employee')</code>
Specific Merging	<code>pd.merge(data_struct , data_struct2 , right_on ='col_name' , left_on ='col_name')</code>	<code>pd.merge(df1, df3, left_on="employee", right_on="name").drop('name', axis=1)</code> # When using right and left on we have to drop one of the cols (avoid redundancy)
Joining (Default Merge to indices only)	<code>data_struct.join(data_struct2)</code>	<code>df1a.join(df2a)</code>
Merging on Indices	<code>pd.merge(data_struct , data_struct2 , left_index=True , right_index = True)</code>	<code>pd.merge(df1a, df3, left_index=True, right_on='name')</code>
Merging with methods	<code>pd.merge(data_struct, data_struct2, how='inner'/'outer'/'left'/'right')</code>	<code>pd.merge(df6, df7, how='inner')</code>
Merging Conflicting Col Names	<code>pd.merge(data_struct, data_struct2, suffixes = ['_suff1', '_suff2'])</code>	<code>pd.merge(df8, df9, on="name", suffixes=["_Sem1", "_Sem_2"])</code>

- Note that when adding multi-index keys in a concatenation , the number of keys should be the same as the number of data structures being concatenated

Advanced Group By Methods

Aggregation using a list (General)	<code>df.groupby('col').aggregate([list_of_methods])</code>	<code>df.groupby('key').aggregate(['min', np.median, max])</code>
Aggregation using a dict (Specific)	<code>df.groupby('col').aggregate({'col' : 'method' , 'col2' : 'method'})</code>	<code>df.groupby('key').aggregate({'data1': 'min', 'data2': 'max'})</code>
Filtering	<code>df.groupby('col').filter(func)</code>	<code>def filter_func(x): return x['data2'].std() > 4 df.groupby('key').filter(filter_func)</code>
Transformation	<code>df.groupby('col').transform(lambda_func)</code>	<code>df.groupby('key').transform(lambda x: x - x.mean())</code>



By Taissir Boukrouba
(taissir2002)

Not published yet.
Last updated 19th November, 2023.
Page 5 of 7.

Sponsored by **CrosswordCheats.com**
Learn to solve cryptic crosswords!
<http://crosswordcheats.com>

Advanced Group By Methods (cont)

Apply	<code>df.groupby('col').apply(user_func)</code>	<pre>def norm_by_data2(x): # x is a DataFrame of group values x['data1'] /= x['data2'].sum() return x df.groupby('key').apply(norm_by_data2)</pre>
Grouping By Custom Mapping	<code>df.groupby(mapping).method</code>	<pre>mapping = {'A': 'vowel', 'B': 'consonant', 'C': 'consonant'} df2.groupby(mapping).sum()</pre>

aggregate() : this method allows using more than one function with groupby for different columns

filter() : this method allows user-defined filter functions to be applied with a groupby (uses boolean operations only)

transform(): mostly uses lambda functions to returned new and transformed version of a columns

apply(): this method allows you to apply arbitrary user-defined functions with groupby



By **Taissir Boukrouba**
(taissir2002)

cheatography.com/taissir2002/

Not published yet.

Last updated 19th November, 2023.

Page 7 of 7.

Sponsored by **CrosswordCheats.com**

Learn to solve cryptic crosswords!

<http://crosswordcheats.com>

