

Python pandas quick guide

Shiu-Tang Li

May 12, 2016

Contents

1	Dataframe initialization / outputs	3
1.1	Load csv files into dataframe.	3
1.2	Initialize a dataframe	3
1.3	Create a new column	3
1.4	Output a dataframe to csv	3
2	Take a quick glance of a dataframe	4
2.1	Print the data frame	4
2.2	Get the description of numerical columns	4
2.3	Get the dimension	4
2.4	Get the data type / get the filtered data by data type	4
2.5	Get the unique elements	4
3	Select data from a dataframe	5
3.1	Get column names	5
3.2	Select a specific column	5
3.3	Select the sub-dataframe of a few columns	5
3.4	Select rows with restrictions on columns	5
3.5	Select rows with row index	5
3.6	Select row index with max values in a specific column	5
3.7	Select given entry	5
3.8	Iterate rows	6
4	Revise data in a dataframe	7
4.1	Revise data in a particular entry	7
4.2	Reindex rows	7
4.3	Reindex one row	7
4.4	Rename columns	7

4.5	Drop columns / rows	7
4.6	Find / drop / fill missing values	8
4.7	Data frame transpose	8
4.8	Change types of a column	8
4.9	Merge data frames	8
5	Search key words in a dataframe	9
5.1	Exact match in target column	9
6	Perform operations on a dataframe	10
6.1	Sort dataframe	10
6.2	Rearrange dataframe - pivot table	10
6.3	Grouping	10
6.4	'Apply' function	10
7	Others	11

1 Dataframe initialization / outputs

1.1 Load csv files into dataframe.

```
1 import pandas
2 data_frame = pandas.read_csv("C:/Users/Shiu-Tang Li/...csv",
3                               encoding = "ISO-8859-1")
4 # encoding: to deal with unicodes
```

1.2 Initialize a dataframe

```
1 import pandas as pd
2 df1 = pd.DataFrame = ({ 'c1': ['1', '2', '3', '4'], 'c2': ['5', '6', '7', '8'] })
3 df2 = pd.DataFrame({ 'c1' : 2,
4                       'c2' : np.array([0] * 100, dtype='int32'),
5                       'c3' : 'hello' })
```

```
1 import pandas as pd
2 list_of_dicts = [{ 'column1':3, 'column2':4}, { 'column1':7, 'column2':2}, { 'column1':6, 'column2':8}]
3 data_frame = pd.DataFrame(list_of_dicts, index=['index1', 'index2'])
4 # construct data frame from list of dictionaries
```

1.3 Create a new column

```
1 data_frame['new_column'] = List OR Series
2 # will get warning message
```

1.4 Output a dataframe to csv

```
1 import pandas
2 data_frame.to_csv("C:/Users/Shiu-Tang Li/...csv")
```

Remark. May load .csv as list of lists instead of data frames.

2 Take a quick glance of a dataframe

2.1 Print the data frame

```
1 print(data_frame.head(5))
2 # print first 5 rows
3 print(data_frame)
4 # print the data frame, dimension information is also attached
```

2.2 Get the description of numerical columns

```
1 print(data_frame.describe())
```

2.3 Get the dimension

```
1 dim = data_frame.shape
2 number_of_rows = dim[0]
3 number_of_columns = dim[1]
```

2.4 Get the data type / get the filtered data by data type

```
1 types = data_frame.dtypes
2 # types is a Series labeled by column names, showing the data type for each
  column.
3 integer_index = types[types == int64].index
4 # types[types == int64] is a filtered Series with integer values.
5 print(data_frame[integer_index])
6 # this gives us the filtered data frame with only int64 values.
```

2.5 Get the unique elements

```
1 print(data_frame['column_name'].unique())
2 # will return a list showing distinct elements in the column
3 print(data_frame['column_name'].value_counts())
4 # will return a table showing the counts in the column
```

3 Select data from a dataframe

3.1 Get column names

```
1 print(data_frame.columns)
2 # data_frame.columns is a list of strings
3 first_column = data_frame.columns[0]
4 # print the first column, which is a string
```

3.2 Select a specific column

```
1 column = data_frame['column_name']
2 # column is a [Series] object, contains row index + values, both are lists
3 column_values = column.values
4 column_index = column.index
```

3.3 Select the sub-dataframe of a few columns

```
1 data_frame2 = data_frame[['column_name1', 'column_name2']]
2 data_frame2 = data_frame[data_frame.columns[0:2]]
3 # select the first two columns in two different ways
4 data_frame2 = data_frame['column_name_x':'column_name_y']
5 # select the columns between the two columns
```

3.4 Select rows with restrictions on columns

```
1 data_frame[data_frame['column_name'] == some_values]
```

3.5 Select rows with row index

```
1 data_frame.iloc[i]
2 #i: row index
3 data_frame.iloc[0:3]
4 # select the rows with indices 0,1,2
```

Remark. The difference between loc and iloc: If the index of the dataframe is 3, 7, 0, 2, ..., iloc[0] will select the third row (true integer index), loc will select the 1st row (index by locations).

3.6 Select row index with max values in a specific column

```
1 data_frame['column_name'].idxmax()
2 # returns the 1st row index that has max
```

3.7 Select given entry

```
1 # Approach 1: i: true row index
2 data_frame.iloc[i]['column_name']
3 # Approach 2: i: location
4 data_frame['column_name'].values[i]
5 # Approach 3: i: true row index
6 data_frame.ix[i, 'column_name']
```

3.8 Iterate rows

```
1 for i, row in data_frame.iterrows():
2 # i: row indices; row: each row
```

4 Revise data in a dataframe

4.1 Revise data in a particular entry

```
1 # i: true row index
2 # Approach 1 (will get warning message):
3 data_frame.ix[i, 'column_name'] = new_value
4 # Approach 2 (will get warning message):
5 data_frame['column_name'][i] = new_value
6 # Approach 3:
7 data_frame.set_value(i, 'column_name', new_value)
8 # Approach 4:
9 data_frame.at[i, 'column_name'] = new_value
```

4.2 Reindex rows

```
1 data_frame.index = [index1, index2, ...]
2 # replace indices with a new list
3 data_frame = data_frame.set_index(['column_name'])
4 # indexed by a particular column
5 data_frame = data_frame.reset_index(drop=True)
6 # reindexed from 0. drop=False: make a dataframe column with the old index
  values.
```

4.3 Reindex one row

```
1 old_indices = data_frame.index
2 new_indices = old_indices.values
3 for i,item in enumerate(new_indices):
4     if item == 'old_index_to_be_changed':
5         new_indices[i] = 'new_index'
6 data_frame.index = new_indices
```

4.4 Rename columns

```
1 data_frame.rename(columns={'column_name1': 'new_column_name1', 'column_name2': '
  new_column_name2'}, inplace=True)
2 # rename a few columns: will get warning message
3 data_frame.columns = ['column_name1', 'column_name2', ...]
4 # rename all columns
```

4.5 Drop columns / rows

```
1 data_frame.drop('column_name', axis=1, inplace=True)
2 # drop a column
3 data_frame.drop(['column_name1', 'column_name2', ...], axis=1, inplace=True)
4 # drop a few columns
5 data_frame.drop('row_index1', axis=0, inplace=True)
6 # drop a row
```

```

7 data_frame.drop(['row_index1', 'row_index2', ...], axis=0, inplace=True)
8 # drop a few rows

```

4.6 Find / drop / fill missing values

```

1 import pandas as pd
2 is_null = pd.isnull(data_frame["column_name"])
3 # will return a true / false Series. Null value = NaN.
4 new_data_frame = data_frame.dropna()
5 # drop all rows with missing values
6 new_data_frame = data_frame.dropna(axis = 1)
7 # drop all columns with missing values
8 new_data_frame = data_frame.dropna(subset=["column1", "column2"])
9 # drop all rows with missing values in the two columns
10 data_frame["column"] = data_frame["column"].fillna(data_frame["column"].median()
11 )
12 data_frame["column"] = data_frame["column"].fillna(data_frame["column"].mean())
13 data_frame["column"] = data_frame["column"].fillna(something)
14 # fill NaN values with column median / mean / or sth else

```

4.7 Data frame transpose

```

1 data_frame.T

```

4.8 Change types of a column

```

1 data_frame['column_name'] = data_frame['column_name'].astype(float)
2 # changing types to float

```

4.9 Merge data frames

```

1 import pandas as pd
2 list1 = [{'c1':4, 'c2':3}, {'c1':2, 'c2': 3}]
3 list2 = [{'c1':6, 'c2':9}, {'c1':8, 'c2':10}]
4 df1 = pd.DataFrame(list1, index=[0,1])
5 df2 = pd.DataFrame(list2, index=['two', 'three'])
6 df = pd.concat([df1, df2])
7 # combine two data frames with common columns, increase rows. index could
  contain keys of different data types

```

```

1 import pandas as pd
2 df1 = pd.DataFrame({'key': ['a', 'b'], 'c1': [1, 2], 'c2': [8, 2]})
3 df2 = pd.DataFrame({'key': ['a', 'b'], 'c3': [3, 5]})
4 merged_df = pd.merge(df1, df2, on='key')
5 # combine two data frames with the same data in the column 'key'. Merged data
  frame will contain 4 columns.

```


5 Search key words in a dataframe

5.1 Exact match in target column

```
1 for i, row in movies.iterrows():  
2     if re.search("Keywords", row["column"])!= None:  
3         print(i)
```

6 Perform operations on a dataframe

6.1 Sort dataframe

```
1 data_frame.sort("column_name", inplace=True, ascending=False)
2 # sort by column, in decreasing order
3 data_frame.sort_index(inplace=True, ascending=False)
4 # sort by index, in decreasing order
```

6.2 Rearrange dataframe - pivot table

```
1 import numpy as np
2 new_data_frame = data_frame.pivot_table(index="column1", values="column2",
3     aggfunc=np.mean)
4 # For different column1 values, find the result of aggfunction applied to
5     column2.
6 new_data_frame = data_frame.pivot_table(index="column1", values="any_column",
7     aggfunc='count')
8 # count the frequency for column1.
```

Remark. Other choices of aggfunc: 'max', 'min', np.std, np.sum, np.median.

6.3 Grouping

```
1 groups = data_frame.groupby("column1")
2 table = groups.aggregate(np.mean)["column2"]
3 # classify rows into different groups based on 'column1', for each group apply
4     np.mean() on all values in "column2".
```

6.4 'Apply' function

	c1	c2	c3	apply2
r1	-	-	-	ap2
r2	-	-	-	ap2
r3	-	-	-	ap2
apply1	ap1	ap1	ap1	

```
1 apply1 = data_frame.apply(some_function) OR
2 apply1 = data_frame.apply(lambda x: some_function(x))
3 # apply function to columns, to get a Series object labeled column names
4 apply2 = data_frame.apply(some_function, axis=1) OR
5 apply2 = data_frame.apply(lambda x: some_function(x), axis=1)
6 # apply function to rows, to get a Series object labeled row indices
```

7 Others

`Series.tolist()` converts Series to lists.
`numpy.nan`: missing value