



Pandas 2 – Data Preprocessing



Data analysis process

Data Understanding

- Descriptive statistics
- Types of data (numerical/categorical)

Data Preprocessing

- **Basic operations**
- Subset selection and **data consolidation**
- Missing data handling

Calculation (Modeling)

- Basic calculation
- Data aggregation

Data Visualization

- Univariate chart
- Bivariate chart
- Multivariate chart

Outline

- Basic operations
 - Copy a DataFrame
 - Add/drop column
 - Sort
- Data consolidation
 - Concatenate
 - Merge

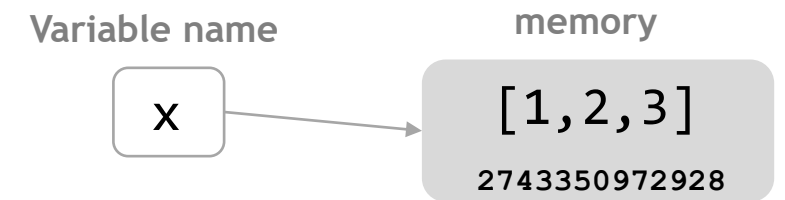
Python variables and memory allocation

- A python variable is a symbolic name, which is a reference to an object. After creating an object, you can refer to it by variable name.
- Use `id()` to see the memory address (object's identity).

```
x = [1,2,3]
```

```
id(x)
```

2743350972928 memory address



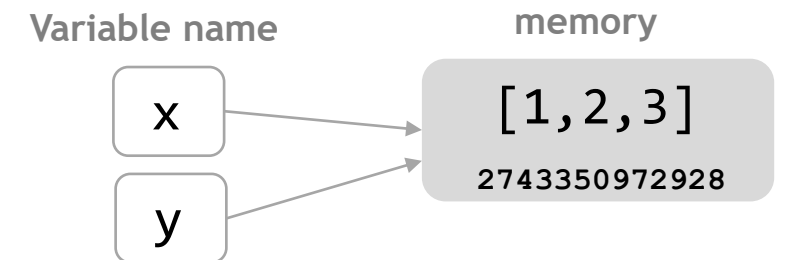
Create an object in memory and let x be a reference to this object.

- If one variable is assigned to another variable, both variables point to the same memory address.

```
y = x
```

```
id(y)
```

2743350972928

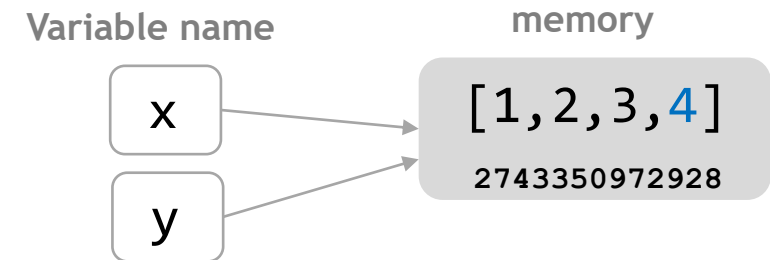


Python variables and memory allocation

- If two variables point to the same address, changing the value of one variable will change the value of the other variable.

```
y.append(4)  
print(x)  
print(y)
```

```
[1, 2, 3, 4]  
[1, 2, 3, 4]
```



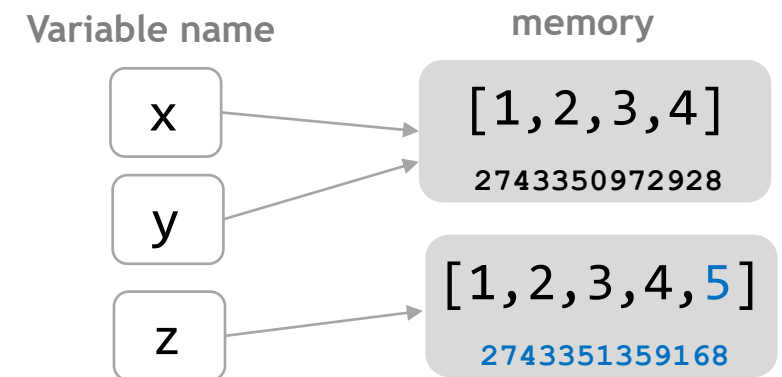
- Use the function `copy()` to create a copied variable, and then you can change the copied variable without changing the original variable.

```
z = x.copy()  
id(z)
```

```
2743351359168
```

```
z.append(5)  
print(x, y, z)
```

```
[1, 2, 3, 4] [1, 2, 3, 4] [1, 2, 3, 4, 5]
```



Create a copy of a DataFrame

- To avoid modifying the source DataFrame when manipulating the data, you can use `copy()` to create a copied DataFrame in advance.

```
df = df_raw.copy()
```

Add a new column

- Give a column name and assign an array.

```
df = df_raw.copy()
```

```
df["eastern"] = [True, True, True, False, False, False]  
df
```

	state	year	pop
0	Ohio	2000	1.5
1	Ohio	2001	1.7
2	Ohio	2002	3.6
3	Nevada	2001	2.4
4	Nevada	2002	2.9
5	Nevada	2003	3.2



	state	year	pop	eastern
0	Ohio	2000	1.5	True
1	Ohio	2001	1.7	True
2	Ohio	2002	3.6	True
3	Nevada	2001	2.4	False
4	Nevada	2002	2.9	False
5	Nevada	2003	3.2	False

Add a new column

- Example without using "copy()"

```
df1 = pd.DataFrame(data)

# (1) assign df1 to another variable named df2
df2 = df1

# (2) if you change the content of df2
df2["newcolumn"] = [1,2,3,4,5,5]

# (3) the content of df1 will also change accordingly
df1
```

	state	year	pop	newcolumn
0	Ohio	2000	1.5	1
1	Ohio	2001	1.7	2
2	Ohio	2002	3.6	3
3	Nevada	2001	2.4	4
4	Nevada	2002	2.9	5
5	Nevada	2003	3.2	5

Drop a column

- Use `drop()` to remove column(s) by specifying `axis = 1`.
 - `axis = 1` → drop columns
 - `axis = 0` → drop rows (default)

```
df.drop(["eastern"], axis = 1)
```

	state	year	pop	eastern
0	Ohio	2000	1.5	True
1	Ohio	2001	1.7	True
2	Ohio	2002	3.6	True
3	Nevada	2001	2.4	False
4	Nevada	2002	2.9	False
5	Nevada	2003	3.2	False



	state	year	pop
0	Ohio	2000	1.5
1	Ohio	2001	1.7
2	Ohio	2002	3.6
3	Nevada	2001	2.4
4	Nevada	2002	2.9
5	Nevada	2003	3.2

By default, `inplace = False`

df				
	state	year	pop	eastern
0	Ohio	2000	1.5	True
1	Ohio	2001	1.7	True
2	Ohio	2002	3.6	True
3	Nevada	2001	2.4	False
4	Nevada	2002	2.9	False
5	Nevada	2003	3.2	False

The original object "df" remains unchanged.

Drop a column

- If `inplace = True`, the change will be applied to the object directly.

```
df.drop(["eastern"], axis = 1, inplace = True)  
df
```

	state	year	pop
0	Ohio	2000	1.5
1	Ohio	2001	1.7
2	Ohio	2002	3.6
3	Nevada	2001	2.4
4	Nevada	2002	2.9
5	Nevada	2003	3.2

- Many functions, like `drop`, can manipulate an object in-place without returning a new object. Be careful with the `inplace=True`, as it modify the original object.

Exercise

(A.1) Given a dataframe. Create a copy named `company_df` and use it to do A.2~A.4.

```
company_raw_df = pd.DataFrame({"company_name": ['JPMorgan Chase', 'Apple', 'Bank of America', 'Amazon', 'Microsoft'],  
                               "profit": [40.4, 63.9, 17.9, 21.3, 51.3],  
                               "assets": [3689.3, 354.1, 2832.2, 321.2, 304.1]})
```

(A.2) Add a new column named `market_value` with a list of values 464.8, 2252.3, 336.3, 1711.8, 1966.6 .

(A.3) Drop the column `assets` .

(A.4) Delete the data of JPMorgan Chase and Bank of America.

Sort a DataFrame

- Use method `sort_values()` to sort a DataFrame.
- Pass a column name.

	state	year	pop
0	Ohio	2000	1.5
1	Ohio	2001	1.7
2	Ohio	2002	3.6
3	Nevada	2001	2.4
4	Nevada	2002	2.9
5	Nevada	2003	3.2



```
df.sort_values(by = "year")
```

	state	year	pop
0	Ohio	2000	1.5
1	Ohio	2001	1.7
3	Nevada	2001	2.4
2	Ohio	2002	3.6
4	Nevada	2002	2.9
5	Nevada	2003	3.2

Sort a DataFrame - in a descending order

- The data is sorted in ascending order by default but can be sorted in descending order by specifying `ascending = False`.

```
df.sort_values(by = "year", ascending = False)
```

	state	year	pop
5	Nevada	2003	3.2
2	Ohio	2002	3.6
4	Nevada	2002	2.9
1	Ohio	2001	1.7
3	Nevada	2001	2.4
0	Ohio	2000	1.5

Sort a DataFrame - by multiple columns

- Pass a list of column names.

```
df.sort_values(by = ["year", "pop"])
```

	state	year	pop
0	Ohio	2000	1.5
1	Ohio	2001	1.7
3	Nevada	2001	2.4
4	Nevada	2002	2.9
2	Ohio	2002	3.6
5	Nevada	2003	3.2

Sort a DataFrame - inplace

- If `inplace = True`, the change will be applied to the object directly.

```
df.sort_values(by = "year", ascending = False, inplace = True)  
df
```

	state	year	pop
5	Nevada	2003	3.2
2	Ohio	2002	3.6
4	Nevada	2002	2.9
1	Ohio	2001	1.7
3	Nevada	2001	2.4
0	Ohio	2000	1.5

Reset index

- Reset the index of the DataFrame after sorting.
 - `drop = True`: Drop the old index.
 - `drop = False`: Add the old index as an additional column to your DataFrame.

```
df.reset_index(drop = True)
```

	state	year	pop
0	Nevada	2003	3.2
1	Ohio	2002	3.6
2	Nevada	2002	2.9
3	Ohio	2001	1.7
4	Nevada	2001	2.4
5	Ohio	2000	1.5

```
df.reset_index(drop = False)
```

	index	state	year	pop
0	5	Nevada	2003	3.2
1	2	Ohio	2002	3.6
2	4	Nevada	2002	2.9
3	1	Ohio	2001	1.7
4	3	Nevada	2001	2.4
5	0	Ohio	2000	1.5

Exercise

(B.1) Use the dataframe `company_raw_df` in (A.1). Sort the dataframe by the `profit` column in a descending order and display the result.

(B.2) Store the returned result in (B.1) in a new variable named `company_sorted_df`.

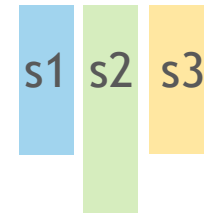
(B.3) Reset the index of `company_sorted_df` and drop the old index.

*Data consolidation -
Concatenate & Merge*

Concatenate - Series

- Suppose you have three Series with no index overlap.

```
s1 = pd.Series([0,1], index = ['a','b'])  
s2 = pd.Series([2,3,4], index = ['c','d','e'])  
s3 = pd.Series([5,6], index = ['f','g'])
```



- Use `concat()` with these series in a list glues together the values and indexes.

```
pd.concat([s1,s2,s3])
```

```
a    0  
b    1  
c    2  
d    3  
e    4  
f    5  
g    6  
dtype: int64
```

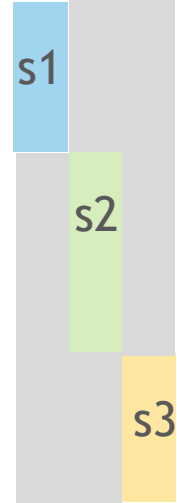


Concatenate - Series

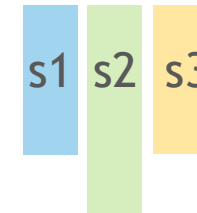
- By default, `concat()` works along `axis = 0`, producing another Series.
- If you pass `axis = 1`, the result will instead be a DataFrame.

```
pd.concat([s1,s2,s3], axis=1)
```

	0	1	2
a	0.0	NaN	NaN
b	1.0	NaN	NaN
c	NaN	2.0	NaN
d	NaN	3.0	NaN
e	NaN	4.0	NaN
f	NaN	NaN	5.0
g	NaN	NaN	6.0



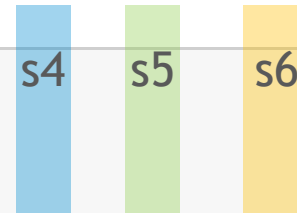
Concatenate series
along the columns



Concatenate - Series

- Suppose you have three Series with the same index.

```
s4 = pd.Series([0,1,2], index = ['a','b','c'])  
s5 = pd.Series([3,4,5], index = ['a','b','c'])  
s6 = pd.Series([6,7,8], index = ['a','b','c'])
```



```
pd.concat([s4,s5,s6], axis=1)
```

	0	1	2
a	0	3	6
b	1	4	7
c	2	5	8



Concatenate - DataFrame

- Suppose you have two DataFrames with the same index.

```
df1 = pd.DataFrame({"col1": [1, 2, 3], "col2": [4, 5, 6], "col3": [7, 8, 9]}, index = ['a', 'b', 'c'])  
df2 = pd.DataFrame({"col1": [11, 22, 33], "col2": [44, 55, 66], "col3": [77, 88, 99]}, index = ['a', 'b', 'c'])
```

df1

	col1	col2	col3
a	1	4	7
b	2	5	8
c	3	6	9

df2

	col1	col2	col3
a	11	44	77
b	22	55	88
c	33	66	99

- Concatenate DataFrames

```
pd.concat([df1, df2])
```

by default, axis = 0

	col1	col2	col3
a	1	4	7
b	2	5	8
c	3	6	9

df1
df2

- Reset index

```
pd.concat([df1, df2], ignore_index = True)
```

	col1	col2	col3
0	1	4	7
1	2	5	8
2	3	6	9
3	11	44	77
4	22	55	88
5	33	66	99

Concatenate - DataFrame

- If you pass `axis = 1`, df1 and df2 will be concatenated along the columns

```
pd.concat([df1,df2], axis = 1)
```

	col1	col2	col3	col1	col2	col3
a	1	4	7	11	44	77
b	2	5	8	22	55	88
c	3	6	9	33	66	99



Exercise

(C.1) Import the datasets `municipality_info_part1.csv` and `municipality_info_part2.csv` as dataframes. The columns in the two datasets are described as follows.

- Municipality_number (object)
- Population (int)
- Area (float)

Note: Use the argument "dtype" to specify the data types.

```
dtype = {"Municipality_number": object, "Population": int, "Area": float} .
```

(C.2) Display the first five rows of each dataset.

(C.3) Concatenate two dataframes in (B.1) along the rows and assign the returned dataframe to a new variable named `mcp_info` .

(C.4) How many rows are in the dataframe `mcp_info` ?

Merge - left join

- `merge()`: Merge dataframes based on the common column (key).
- Left join: Use keys from left frame.

df1			df2		
	employID	name		employID	birthday
0	E011	John	0	E010	20-07
1	E012	Diana	1	E012	12-06
2	E013	Matthew	2	E013	18-01
3	E014	Jerry	3	E015	16-05
4	E015	Kathy	4	E016	02-10
5	E016	Sara	5	E017	19-08
6	E017	Alex			

Left DataFrame Right DataFrame key

```
pd.merge(df1, df2, how = 'left', on = 'employID' )
```

	employID	name	birthday
0	E011	John	NaN
1	E012	Diana	12-06
2	E013	Matthew	18-01
3	E014	Jerry	NaN
4	E015	Kathy	16-05
5	E016	Sara	02-10
6	E017	Alex	19-08

To be merged on
the left side

Merge - inner join

- Inner join: Use **intersection** of keys from both frames.

```
pd.merge(df1, df2, how = 'inner', on = 'employID' )
```

df1

	employID	name
0	E011	John
1	E012	Diana
2	E013	Matthew
3	E014	Jerry
4	E015	Kathy
5	E016	Sara
6	E017	Alex

df2

	employID	birthday
0	E010	20-07
1	E012	12-06
2	E013	18-01
3	E015	16-05
4	E016	02-10
5	E017	19-08

	employID	name	birthday
0	E012	Diana	12-06
1	E013	Matthew	18-01
2	E015	Kathy	16-05
3	E016	Sara	02-10
4	E017	Alex	19-08

Merge - outer join

- Outer join: Use **union** of keys from both frames

```
pd.merge(df1, df2, how = 'outer', on = 'employID' )
```

df1

	employID	name
0	E011	John
1	E012	Diana
2	E013	Matthew
3	E014	Jerry
4	E015	Kathy
5	E016	Sara
6	E017	Alex

df2

	employID	birthday
0	E010	20-07
1	E012	12-06
2	E013	18-01
3	E015	16-05
4	E016	02-10
5	E017	19-08

	employID	name	birthday
0	E011	John	NaN
1	E012	Diana	12-06
2	E013	Matthew	18-01
3	E014	Jerry	NaN
4	E015	Kathy	16-05
5	E016	Sara	02-10
6	E017	Alex	19-08
7	E010	NaN	20-07

Exercise

(D.1) Import the dataset `municipality_name.csv` as a dataframe named `mcp_name`. The columns in the dataset are described as follows.

- `Municipality_number` (object)
- `Municipality_name` (object)

Note: Use the argument "encoding" to specify the character encoding.

(D.2) Use the dataframe `mcp_info` obtained in (C.3). Find the municipality name corresponding to the municipality number from `mcp_name`. Add them to the new column in `mcp_info`.

(D.3) List the five most populous municipalities.

(D.4) Use the dataframe `mcp_name` obtained in (D.1). Find the municipality area corresponding to the municipality number from `mcp_info`. How many municipalities lack information about the area size?

Differences between merge() and concat()

- `concat()` simply stacks multiple DataFrames together either vertically or horizontally.



- `merge()` first align the **selected common columns** of the two DataFrames, and then pick up the remaining columns from the aligned rows of each DataFrame.

