

PANDAS



Cheat Sheet

```
import pandas as pd
```

Data Structures

Series

```
# Create a series
>>> s = pd.Series([1, 2, 3], index=['A', 'B', 'C'])

Index   Data
A      1
B      2
C      3

>>> s.index
Index(['A', 'B', 'C'], dtype='object')

>>> s.dtype
dtype('int64')

>>> s.size
3

>>> s.ndim
1

>>> s.values
array([1, 2, 3], dtype=int64)
```

DataFrame

```
# Create a DataFrame
>>> df = pd.DataFrame({'ID': [1, 2, 3],
                      'Name': ['Alex', 'Brian', 'David'],
                      'Profession': ['DA', 'DE', 'DS'],
                      index = ['1 (One)', '2 (Two)', '3 (Three)'])

          ID    Name  Profession
1 (One)  1    Alex        DA
2 (Two)  2    Brian       DE
3 (Three) 3    David       DS
```

```
# Get the shape of the DataFrame
>>> df.shape
(3, 3)
```

```
# Get the column names of the DataFrame
>>> df.columns
Index(['ID', 'Name', 'Profession'], dtype='object')
```

```
# Get the index information of the DataFrame
>>> df.index
Index(['1 (One)', '2 (Two)', '3 (Three)'], dtype='object')
```

```
# Get the data types of the DataFrame columns
>>> df.dtypes
```

```
ID      int64
Name     object
Profession  object
dtype: object
```

```
# Check if there are any missing values in the DataFrame
>>> df.isnull().values.any()
False
```

```
# Get the count of missing values in each column
>>> df.isnull().sum()
```

```
ID      0
Name     0
Profession  0
dtype: int64
```

```
# Get the count of non-null values in each column
>>> df.count()
```

```
ID      3
Name     3
Profession  3
dtype: int64
```

```
# Generate descriptive statistics of the DataFrame (transposed)
>>> df.describe().T
```

Importing & Exporting

```
# Read a CSV file and create a DataFrame
>>> pd.read_csv('filename.csv')

# Save the DataFrame to a CSV file
df.to_csv('filename.csv')

# Read an Excel file and create a DataFrame
>>> pd.read_excel('filename.xlsx')

# Save the DataFrame to an Excel file
>>> df.to_excel('filename.xlsx')
```

Data Manipulation

Selection

Selecting Rows

```
# Select a specific row by its label
>>> df.loc['1 (One)']:
```

```
# Select a specific row by its index
>>> df.loc[1]
```

Selecting Columns

```
# Select a single column by its name
>>> df['ID']
```

```
# Select multiple columns by their name
>>> df[['ID', 'Profession']]
```

```
# Select a single column using label-based indexing
>>> df.loc[:, 'Name']
```

```
# Select a single column using integer-based indexing
>>> df.iloc[:, 1]
```

Selecting Rows and Columns

```
# Select a specific cell by its label-based row and column indices
>>> df.loc['2 (Two)', 'Name']
```

```
# Select a specific cell by its integer-based row and column indices
>>> df.iloc[1, 1]
```

```
# Select specific rows and columns using label-based indexing
>>> df.loc['1 (One)', ['ID', 'Profession']]
```

```
# Select specific rows and columns using integer-based indexing
>>> df.iloc[1, ['ID', 'Profession']]
```

Conditional Selection

```
# Select rows based on a condition
>>> df[df['ID'] > 1]
```

```
# Select rows based on multiple conditions using logical operators
>>> df[(df['ID'] > 2) & (df['Profession'] == 'DS')]
```

```
# Select rows where a column value is present in a given list
>>> df.loc[df['Name'].isin(['Alex', 'David'])]
```

Deleting & Adding

```
# Drop rows with any missing values
>>> df.dropna()
```

```
# Drop columns with any missing values
>>> df.dropna(axis=1)
```

```
# Drop columns with fewer than n non-null values
>>> df.dropna(axis=1, thresh=n)
```

```
# Fill missing values with a specified value
df.fillna(value)
```

Filtering & Sorting

```
# Filter columns using regular expression matching
```

```
>>> df.filter(regex=regex)
```

```
# Sort DataFrame by values in a specific column in descending order
```

```
>>> df.sort_values('ID', ascending=False)
```

Editing

```
# Perform quantile-based discretization of values in 'col' into 3 bins with custom labels
```

```
>>> pd.qcut(df['col'], 3, labels=qcut_labels)
```

```
# Perform value-based discretization of values in 'col' into custom bins with labels
```

```
>>> pd.cut(df['col'], bins=cut_bins, labels=cut_labels)
```

```
# Create a pivot table using 'col1' as index, 'col2' as columns, and 'col3' as values
```

```
>>> df.pivot(index='col1', columns='col2', values='col3')
```

```
# Reset the index of the DataFrame
```

```
>>> df.reset_index()
```

```
# Rename specific columns of the DataFrame
```

```
>>> df.rename(columns = rename_list)
```

Joining

```
# Merge operation
>>> df1 = pd.DataFrame({
    'key': ['A', 'B', 'C', 'D'],
    'value1': [1, 2, 3, 4]})
```

```
>>> df2 = pd.DataFrame({
    'key': ['B', 'D', 'E', 'F'],
    'value2': [5, 6, 7, 8]})
```

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

Key	Value1
0	A
1	B
2	C
3	D

Key	Value2
0	B
1	D
2	E
3	F

Key	Value1	Value2
0	B	5
1	D	6
2	E	7
3	F	8

```
>>> df4 = pd.DataFrame({
    'key': ['B', 'D', 'E', 'F'],
    'value4': ['orange', 'grape', 'kiwi', 'lemon']})
```

```
>>> df3.set_index('key').join(df4.set_index('key'), how='inner')
```

Key	Value3
0	A
1	B
2	C
3	D

Key	Value2
0	B
1	D
2	E
3	F

Value1	Value2
B	banana
D	date
E	kiwi
F	lemon

```
# Example for concatenate operation
```

```
>>> df5 = pd.DataFrame({
    'A': ['A0', 'A1', 'A2'],
    'B': ['B0', 'B1', 'B2']})
```

```
>>> df6 = pd.DataFrame({
    'A': ['A3', 'A4', 'A5'],
    'B': ['B3', 'B4', 'B5']})
```

```
>>> pd.concat([df5, df6], axis=0)
```

A	B
0	A0
1	A1
2	A2
0	A3
1	A4
2	A5

A	B
0	B0
1	B1
2	B2
0	B3
1	B4
2	B5

A	B
0	B0
1	B1
2	B2
0	B3
1	B4
2	B5

Statistics

```
# Generate descriptive statistics of the DataFrame
```

```
>>> df.describe()
```

```
# Calculate the mean of each column in the DataFrame
```

```
>>> df.mean()
```

```
# Calculate the correlation matrix of the DataFrame
```

```
>>> df.corr()
```

```
# Count the number of non-null values in each column of the DataFrame
```

```
>>> df.count()
```

```
# Find the maximum value in each column of the DataFrame
```

```
>>> df.max()
```

```
# Find the minimum value in each column of the DataFrame
```

```
>>> df.min()
```

```
# Calculate the median of each column in the DataFrame
```

```
>>> df.median()
```

```
# Calculate the standard deviation of each column in the DataFrame
```

```
>>> df.std()
```

Name	City	Age	Salary
<tbl_info cols="4

