

Data Science Cheat Sheet

Python Basics

BASICS, PRINTING AND GETTING HELP

`x = 3` - Assign 3 to the variable `x`
`print(x)` - Print the value of `x`
`type(x)` - Return the type of the variable `x` (in this case, `int` for integer)

`help(x)` - Show documentation for the `str` data type
`help(print)` - Show documentation for the `print()` function

READING FILES

```
f = open("my_file.txt","r")
file_as_string = f.read()
```

- Open the file `my_file.txt` and assign its contents to `s`

```
import csv
f = open("my_dataset.csv","r")
csvreader = csv.reader(f)
csv_as_list = list(csvreader)
```

- Open the CSV file `my_dataset.csv` and assign its data to the list of lists `csv_as_list`

STRINGS

`s = "hello"` - Assign the string "hello" to the variable `s`

```
s = """She said,
"there's a good idea.""""

```

- Assign a multi-line string to the variable `s`. Also used to create strings that contain both " and ' characters

`len(s)` - Return the number of characters in `s`

`s.startswith("hel")` - Test whether `s` starts with the substring "hel"

`s.endswith("lo")` - Test whether `s` ends with the substring "lo"

`"{} plus {} is {}".format(3,1,4)` - Return the string with the values 3, 1, and 4 inserted

`s.replace("e","z")` - Return a new string based on `s` with all occurrences of "e" replaced with "z"

`s.split(" ")` - Split the string `s` into a list of strings, separating on the character " " and return that list

NUMERIC TYPES AND MATHEMATICAL OPERATIONS

`i = int("5")` - Convert the string "5" to the integer 5 and assign the result to `i`

`f = float("2.5")` - Convert the string "2.5" to the float value 2.5 and assign the result to `f`

`5 + 5` - Addition

`5 - 5` - Subtraction

`10 / 2` - Division

`5 * 2` - Multiplication

`3 ** 2` - Raise 3 to the power of 2 (or 3^2)

`27 ** (1/3)` - The 3rd root of 27 (or $\sqrt[3]{27}$)

`x += 1` - Assign the value of `x + 1` to `x`

`x -= 1` - Assign the value of `x - 1` to `x`

LISTS

`l = [100, 21, 88, 3]` - Assign a list containing the integers 100, 21, 88, and 3 to the variable `l`

`l = list()` - Create an empty list and assign the result to `l`

`l[0]` - Return the first value in the list `l`

`l[-1]` - Return the last value in the list `l`

`l[1:3]` - Return a slice (list) containing the second and third values of `l`

`len(l)` - Return the number of elements in `l`

`sum(l)` - Return the sum of the values of `l`

`min(l)` - Return the minimum value from `l`

`max(l)` - Return the maximum value from `l`

`l.append(16)` - Append the value 16 to the end of `l`

`l.sort()` - Sort the items in `l` in ascending order

`" ".join(["A", "B", "C", "D"])` - Converts the list ["A", "B", "C", "D"] into the string "A B C D"

DICTIONARIES

`d = {"CA":"Canada", "GB": "Great Britain", "IN": "India"}` - Create a dictionary with keys of "CA", "GB", and "IN" and corresponding values of "Canada", "Great Britain", and "India"

`d["GB"]` - Return the value from the dictionary `d` that has the key "GB"

`d.get("AU", "Sorry")` - Return the value from the dictionary `d` that has the key "AU", or the string "Sorry" if the key "AU" is not found in `d`

`d.keys()` - Return a list of the keys from `d`

`d.values()` - Return a list of the values from `d`

`d.items()` - Return a list of (key, value) pairs from `d`

MODULES AND FUNCTIONS

The body of a function is defined through indentation.

`import random` - Import the module `random`

`from math import sqrt` - Import the function `sqrt` from the module `math`

```
def calculate(addition_one,addition_two,
             exponent=1,factor=1):
    result = (value_one + value_two) ** exponent * factor
    return result
```

- Define a new function `calculate` with two required and two optional named arguments which calculates and returns a result.

`addition(3,5,factor=10)` - Run the `addition` function with the values 3 and 5 and the named argument `10`

BOOLEAN COMPARISONS

`x == 5` - Test whether `x` is equal to 5

`x != 5` - Test whether `x` is not equal to 5

`x > 5` - Test whether `x` is greater than 5

`x < 5` - Test whether `x` is less than 5

`x >= 5` - Test whether `x` is greater than or equal to 5

`x <= 5` - Test whether `x` is less than or equal to 5

`x == 5 or name == "alfred"` - Test whether `x` is equal to 5 or `name` is equal to "alfred"

`x == 5 and name == "alfred"` - Test whether `x` is equal to 5 and `name` is equal to "alfred"

`5 in l` - Checks whether the value 5 exists in the list `l`

`"GB" in d` - Checks whether the value "GB" exists in the keys for `d`

IF STATEMENTS AND LOOPS

The body of if statements and loops are defined through indentation.

```
if x > 5:
    print("{} is greater than five".format(x))
```

```
elif x < 0:
```

```
    print("{} is negative".format(x))
```

```
else:
```

```
    print("{} is between zero and five".format(x))
```

- Test the value of the variable `x` and run the code body based on the value

```
for value in l:
```

```
    print(value)
```

- Iterate over each value in `l`, running the code in the body of the loop with each iteration

```
while x < 10:
```

```
    x += 1
```

- Run the code in the body of the loop until the value of `x` is no longer less than 10

Data Science Cheat Sheet

Python - Intermediate

KEY BASICS, PRINTING AND GETTING HELP

This cheat sheet assumes you are familiar with the content of our Python Basics Cheat Sheet

s - A Python string variable

i - A Python integer variable

f - A Python float variable

l - A Python list variable

d - A Python dictionary variable

LISTS

l.pop(3) - Returns the fourth item from **l** and deletes it from the list

l.remove(x) - Removes the first item in **l** that is equal to **x**

l.reverse() - Reverses the order of the items in **l**

l[1::2] - Returns every second item from **l**, commencing from the 1st item

l[-5:] - Returns the last 5 items from **l** specific axis

STRINGS

s.lower() - Returns a lowercase version of **s**

s.title() - Returns **s** with the first letter of every word capitalized

"23".zfill(4) - Returns "0023" by left-filling the string with 0's to make it's length 4.

s.splitlines() - Returns a list by splitting the string on any newline characters.

Python strings share some common methods with lists

s[:5] - Returns the first 5 characters of **s**

"fri" + "end" - Returns "friend"

"end" in s - Returns True if the substring "end" is found in **s**

RANGE

Range objects are useful for creating sequences of integers for looping.

range(5) - Returns a sequence from 0 to 4

range(2000, 2018) - Returns a sequence from 2000 to 2017

range(0, 11, 2) - Returns a sequence from 0 to 10, with each item incrementing by 2

range(0, -10, -1) - Returns a sequence from 0 to -9

list(range(5)) - Returns a list from 0 to 4

DICTIONARIES

max(d, key=d.get) - Return the key that corresponds to the largest value in **d**

min(d, key=d.get) - Return the key that corresponds to the smallest value in **d**

SETS

my_set = set(l) - Return a **set** object containing the unique values from **l**

len(my_set) - Returns the number of objects in **my_set** (or, the number of unique values from **l**)

a in my_set - Returns True if the value **a** exists in **my_set**

REGULAR EXPRESSIONS

import re - Import the Regular Expressions module

re.search("abc", s) - Returns a **match** object if the regex "abc" is found in **s**, otherwise **None**

re.sub("abc", "xyz", s) - Returns a string where all instances matching regex "abc" are replaced by "xyz"

LIST COMPREHENSION

A one-line expression of a for loop

[i ** 2 for i in range(10)] - Returns a list of the squares of values from 0 to 9

[s.lower() for s in l_strings] - Returns the list **l_strings**, with each item having had the **.lower()** method applied

[i for i in l_floats if i < 0.5] - Returns the items from **l_floats** that are less than 0.5

FUNCTIONS FOR LOOPING

```
for i, value in enumerate(l):
    print("The value of item {} is {}".format(i, value))
```

- Iterate over the list **l**, printing the index location of each item and its value

```
for one, two in zip(l_one, l_two):
    print("one: {}, two: {}".format(one, two))
```

- Iterate over two lists, **l_one** and **l_two** and print each value

```
while x < 10:
    x += 1
```

- Run the code in the body of the loop until the value of **x** is no longer less than 10

DATETIME

import datetime as dt - Import the **datetime** module

now = dt.datetime.now() - Assign **datetime** object representing the current time to **now**

wks4 = dt.datetime.timedelta(weeks=4) - Assign a **timedelta** object representing a timespan of 4 weeks to **wks4**

now - wks4 - Return a **datetime** object representing the time 4 weeks prior to **now**

newyear_2020 = dt.datetime(year=2020, month=12, day=31) - Assign a **datetime** object representing December 25, 2020 to **newyear_2020**

newyear_2020.strftime("%A, %b %d, %Y") - Returns "Thursday, Dec 31, 2020"

dt.datetime.strptime('Dec 31, 2020', "%d, %Y") - Return a **datetime** object representing December 31, 2020

RANDOM

import random - Import the **random** module

random.random() - Returns a random float between 0.0 and 1.0

random.randint(0, 10) - Returns a random integer between 0 and 10

random.choice(l) - Returns a random item from the list **l**

COUNTER

from collections import Counter - Import the **Counter** class

c = Counter(l) - Assign a **Counter** (dict-like) object with the counts of each unique item from 1, to **c**

c.most_common(3) - Return the 3 most common items from **l**

TRY/EXCEPT

Catch and deal with Errors

1_ints = [1, 2, 3, "", 5] - Assign a list of integers with one missing value to **1_ints**

```
1_floats = []
for i in 1_ints:
    try:
        1_floats.append(float(i))
    except:
        1_floats.append(i)
```

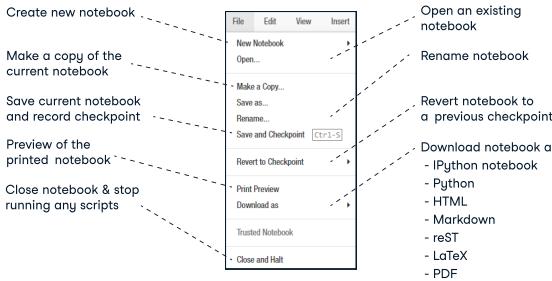
- Convert each value of **1_ints** to a float, catching and handling **ValueError: could not convert string to float:** where values are missing.

Python For Data Science

Jupyter Cheat Sheet

Learn Jupyter online at www.DataCamp.com

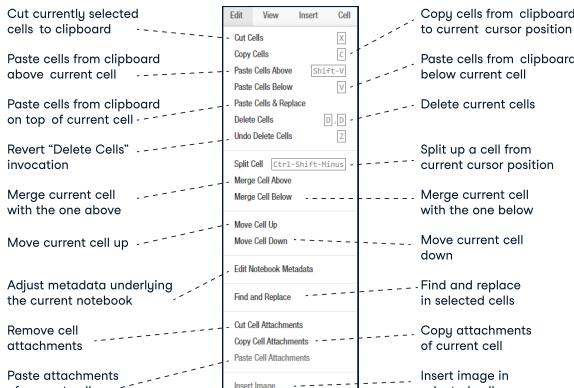
> Saving/Loading Notebooks



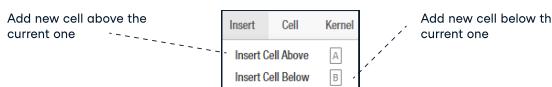
> Writing Code And Text

Code and text are encapsulated by 3 basic cell types: markdown cells, code cells, and raw NBConvert cells

Edit Cells



Insert Cells



> Working with Different Programming Languages

Kernels provide computation and communication with front-end interfaces like the notebooks.

There are three main kernels:

IP[y]:

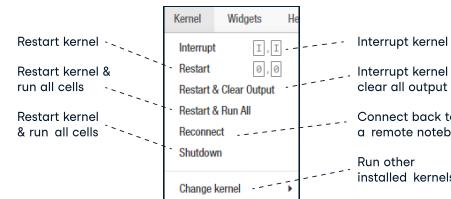
IPython

R

IRkernel

IJulia

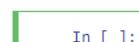
Installing Jupyter Notebook will automatically install the IPython kernel.



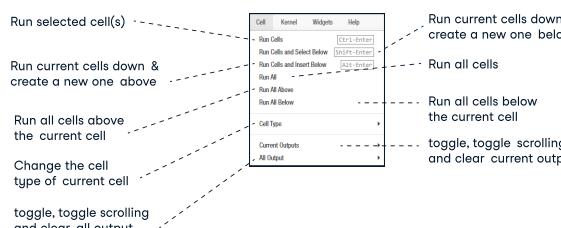
Command Mode:



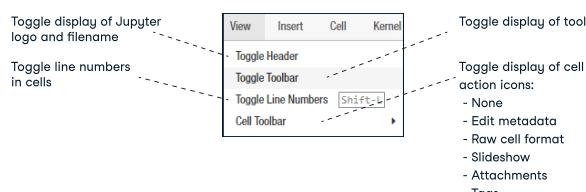
Edit Mode:



Executing Cells



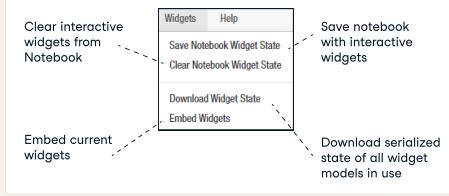
View Cells



> Widgets

Notebook widgets provide the ability to visualize and control changes in your data, often as a control like a slider, textbox, etc.

You can use them to build interactive GUIs for your notebooks or to synchronize stateful and stateless information between Python and JavaScript.



1. Save and checkpoint
2. Insert cell below
3. Cut cell
4. Copy cell(s)
5. Paste cell(s) below
6. Move cell up
7. Move cell down
8. Run current cell
9. Interrupt kernel
10. Restart kernel
11. Restart kernel and re-run notebook
12. Display characteristics
13. Open command palette
14. Current kernel
15. Kernel status
16. Log out from notebook server

Asking For Help

Walk through a UI tour

Edit the built-in keyboard shortcuts

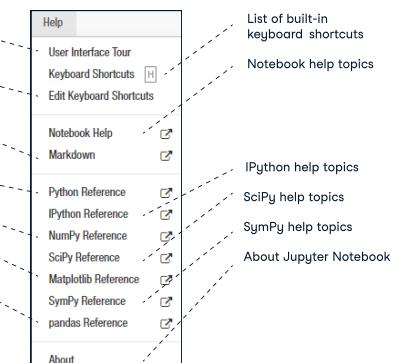
Description of markdown available in notebook

Python help topics

NumPy help topics

Matplotlib help topics

Pandas help topics



Beyond Matplotlib

Seaborn: Statistical Data Visualization

Cartopy: Geospatial Data Processing

yt: Volumetric data Visualization

mpld3: Bringing Matplotlib to the browser

Datashader: Large data processing pipeline

plotnine: A Grammar of Graphics for Python

Matplotlib Cheatsheets (c) 2020 Nicolas P. Rougier

Released under a CC-BY 4.0 International License





Python For Data Science

Pandas Basics Cheat Sheet

Learn Pandas Basics online at www.DataCamp.com

Pandas

The **Pandas** library is built on NumPy and provides easy-to-use **data structures** and **data analysis** tools for the Python programming language.

Use the following import convention:

```
>>> import pandas as pd
```

> Pandas Data Structures

Series

A **one-dimensional** labeled array capable of holding any data type

Index →	a	3
	b	-5
	c	7
	d	4

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

Dataframe

A **two-dimensional** labeled data structure with columns of potentially different types

Columns →	Country	Capital	Population
Index →	0	Belgium	Brussels 11190846
	1	India	New Delhi 1303171035
	2	Brazil	Brasilia 207847528

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
   'Population': [11190846, 1303171035, 207847528]}
>>> df = pd.DataFrame(data,
   columns=['Country', 'Capital', 'Population'])
```

> Dropping

```
>>> s.drop(['a', 'c']) #Drop values from rows (axis=0)
>>> df.drop('Country', axis=1) #Drop values from columns(axis=1)
```

> Asking For Help

```
>>> help(pd.Series.loc)
```

> Sort & Rank

```
>>> df.sort_index() #Sort by labels along an axis
>>> df.sort_values(by='Country') #Sort by the values along an axis
>>> df.rank() #Assign ranks to entries
```

> I/O

Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)
>>> df.to_csv('myDataFrame.csv')
```

Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> df.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')

Read multiple sheets from the same file
>>> xlSx = pd.ExcelFile('file.xls')
>>> df = pd.read_excel(xlSx, 'Sheet1')
```

Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:/// :memory:')
>>> pd.read_sql("SELECT * FROM my_table;", engine)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_sql_query("SELECT * FROM my_table;", engine)

read_sql() is a convenience wrapper around read_sql_table() and read_sql_query()
>>> df.to_sql('myDF', engine)
```

> Selection

Also see NumPy Arrays

Getting

```
>>> s['b'] #Get one element
-5
>>> df[1:] #Get subset of a DataFrame
   Country Capital Population
1 India New Delhi 1303171035
2 Brazil Brasilia 207847528
```

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iloc[[0],[0]] #Select single value by row & column
'Belgium'
>>> df.iat[[0],[0]]
'Belgium'
```

By Label

```
>>> df.loc[[0], ['Country']] #Select single value by row & column labels
'Belgium'
>>> df.at[[0], ['Country']]
'Belgium'
```

By Label/Position

```
>>> df.ix[2] #Select single row of subset of rows
Country Brazil
Capital Brasilia
Population 207847528
>>> df.ix[:, 'Capital'] #Select a single column of subset of columns
0 Brussels
1 New Delhi
2 Brasilia
>>> df.ix[1, 'Capital'] #Select rows and columns
'New Delhi'
```

Boolean Indexing

```
>>> s[~(s > 1)] #Series s where value is not >1
>>> s[(s < -1) | (s > 2)] #s where value is <-1 or >2
>>> df[df['Population']>1200000000] #Use filter to adjust DataFrame
```

Setting

```
>>> s['a'] = 6 #Set index a of Series s to 6
```

>

Retrieving Series/DataFrame Information

Basic Information

```
>>> df.shape #(rows,columns)
>>> df.index #Describe index
>>> df.columns #Describe DataFrame columns
>>> df.info() #Info on DataFrame
>>> df.count() #Number of non-NA values
```

Summary

```
>>> df.sum() #Sum of values
>>> df.cumsum() #Cumulative sum of values
>>> df.min() / df.max() #Minimum/maximum values
>>> df.idxmin() / df.idxmax() #Minimum/Maximum index value
>>> df.describe() #Summary statistics
>>> df.mean() #Mean of values
>>> df.median() #Median of values
```

> Applying Functions

```
>>> f = lambda x: x*2
>>> df.apply(f) #Apply function
>>> df.applymap(f) #Apply function element-wise
```

> Data Alignment

Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
a 10.0
b NaN
c 5.0
d 7.0
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)
a 10.0
b -5.0
c 5.0
d 7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```

Learn Data Skills Online at
www.DataCamp.com



Data Wrangling

with pandas Cheat Sheet
<http://pandas.pydata.org>

Pandas [API Reference](#) Pandas [User Guide](#)

Creating DataFrames

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.
```

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```

		a	b	c
N	v			
D	1	4	7	10
D	2	5	8	11
e	2	6	9	12

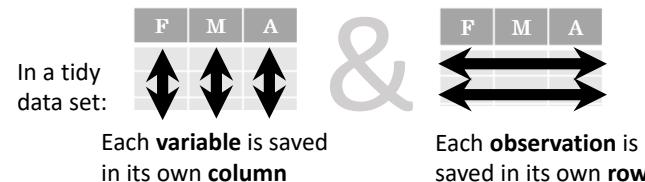
```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d', 1), ('d', 2),
         ('e', 2)], names=['n', 'v']))
Create DataFrame with a MultiIndex
```

Method Chaining

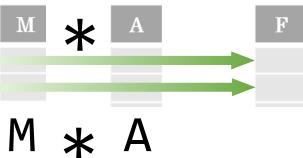
Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={"variable": "var",
                      "value": "val"})
      .query('val >= 200'))
```

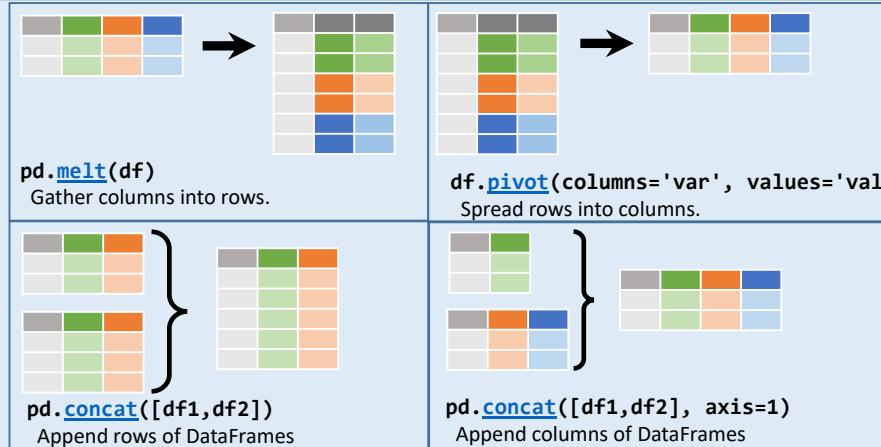
Tidy Data – A foundation for wrangling in pandas



Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.



Reshaping Data – Change layout, sorting, reindexing, renaming



- `df.sort_values('mpg')`
Order rows by values of a column (low to high).
- `df.sort_values('mpg', ascending=False)`
Order rows by values of a column (high to low).
- `df.rename(columns = {'y': 'year'})`
Rename the columns of a DataFrame
- `df.sort_index()`
Sort the index of a DataFrame
- `df.reset_index()`
Reset index of DataFrame to row numbers, moving index to columns.
- `df.drop(columns=['Length', 'Height'])`
Drop columns from DataFrame

Subset Observations - rows



`df[df.Length > 7]`
Extract rows that meet logical criteria.

`df.drop_duplicates()`
Remove duplicate rows (only considers columns).

`df.sample(frac=0.5)`
Randomly select fraction of rows.

`df.sample(n=10)`
Randomly select n rows.

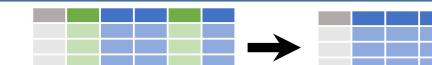
`df.nlargest(n, 'value')`
Select and order top n entries.

`df.nsmallest(n, 'value')`
Select and order bottom n entries.

`df.head(n)`
Select first n rows.

`df.tail(n)`
Select last n rows.

Subset Variables - columns



`df[['width', 'length', 'species']]`
Select multiple columns with specific names.

`df['width'] or df.width`
Select single column with specific name.

`df.filter(regex='regex')`
Select columns whose name matches regular expression `regex`.

Using query

`query()` allows Boolean expressions for filtering rows.

`df.query('Length > 7')`
`df.query('Length > 7 and Width < 8')`
`df.query('Name.str.startswith("abc")', engine="python")`

Use `df.loc[]` and `df.iloc[]` to select only rows, only columns or both.

Use `df.at[]` and `df.iat[]` to access a single value by row and column.
First index selects rows, second index columns.

`df.iloc[10:20]`
Select rows 10-20.
`df.iloc[:, [1, 2, 5]]`
Select columns in positions 1, 2 and 5 (first column is 0).
`df.loc[:, 'x2':'x4']`
Select all columns between x2 and x4 (inclusive).
`df.loc[df['a'] > 10, ['a', 'c']]`
Select rows meeting logical condition, and only the specific columns.
`df.iat[1, 2]` Access single value by index
`df.at[4, 'A']` Access single value by label

Logic in Python (and pandas)

<	Less than	!=	Not equal to
>	Greater than	<code>df.column.isin(values)</code>	Group membership
==	Equals	<code>pd.isnull(obj)</code>	Is NaN
<=	Less than or equals	<code>pd.notnull(obj)</code>	Is not NaN
>=	Greater than or equals	<code>&, , ~, ^, df.any(), df.all()</code>	Logical and, or, not, xor, any, all

regex (Regular Expressions) Examples

'.'	Matches strings containing a period '.'
'Length\$'	Matches strings ending with word 'Length'
'^Sepal'	Matches strings beginning with the word 'Sepal'
'^x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^(?!Species\$).*''	Matches strings except the string 'Species'

Summarize Data

`df['w'].value_counts()`

Count number of rows with each unique value of variable

`len(df)`

of rows in DataFrame.

`df.shape`

Tuple of # of rows, # of columns in DataFrame.

`df['w'].nunique()`

of distinct values in a column.

`df.describe()`

Basic descriptive and statistics for each column (or GroupBy).



pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

`sum()`

Sum values of each object.

`count()`

Count non-NA/null values of each object.

`median()`

Median value of each object.

`quantile([0.25,0.75])`

Quantiles of each object.

`apply(function)`

Apply function to each object.

`min()`

Minimum value in each object.

`max()`

Maximum value in each object.

`mean()`

Mean value of each object.

`var()`

Variance of each object.

`std()`

Standard deviation of each object.

Group Data



`df.groupby(by="col")`

Return a GroupBy object, grouped by values in column named "col".

`df.groupby(level="ind")`

Return a GroupBy object, grouped by values in index level named "ind".

All of the summary functions listed above can be applied to a group.

Additional GroupBy functions:

`size()`

Size of each group.

`agg(function)`

Aggregate group using function.

Windows

`df.expanding()`

Return an Expanding object allowing summary functions to be applied cumulatively.

`df.rolling(n)`

Return a Rolling object allowing summary functions to be applied to windows of length n.

Handling Missing Data

`df.dropna()`

Drop rows with any column having NA/null data.

`df.fillna(value)`

Replace all NA/null data with value.

Make New Columns



`df.assign(Area=lambda df: df.Length*df.Height)`

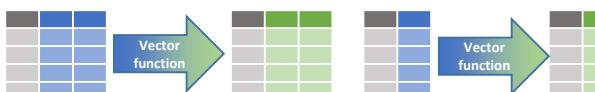
Compute and append one or more new columns.

`df['Volume'] = df.Length*df.Height*df.Depth`

Add single column.

`pd.cut(df.col, n, labels=False)`

Bin column into n buckets.



pandas provides a large set of **vector functions** that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

`max(axis=1)`

Element-wise max.

`min(axis=1)`

Element-wise min.

`clip(lower=-10,upper=10)`

Trim values at input thresholds

`abs()`

Absolute value.

The examples below can also be applied to groups. In this case, the function is applied on a per-group basis, and the returned vectors are of the length of the original DataFrame.

`shift(1)`

Copy with values shifted by 1.

`rank(method='dense')`

Ranks with no gaps.

`rank(method='min')`

Ranks. Ties get min rank.

`rank(pct=True)`

Ranks rescaled to interval [0, 1].

`rank(method='first')`

Ranks. Ties go to first value.

`shift(-1)`

Copy with values lagged by 1.

`cumsum()`

Cumulative sum.

`cummax()`

Cumulative max.

`cummin()`

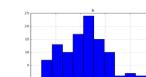
Cumulative min.

`cumprod()`

Cumulative product.

`df.plot.hist()`

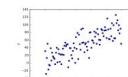
Histogram for each column



Plotting

`df.plot.scatter(x='w',y='h')`

Scatter chart using pairs of points



Combine Data Sets

`adf`

x1	x2
A	1
B	2
C	3

`bdf`

x1	x3
A	T
B	F
D	T



Standard Joins

x1	x2	x3
A	1	T
B	2	F
C	NaN	

`pd.merge(adf, bdf, how='left', on='x1')`

Join matching rows from bdf to adf.

x1	x2	x3
A	1.0	T
B	2.0	F
D	NaN	T

`pd.merge(adf, bdf, how='right', on='x1')`

Join matching rows from adf to bdf.

x1	x2	x3
A	1	T
B	2	F
C	NaN	

`pd.merge(adf, bdf, how='inner', on='x1')`

Join data. Retain only rows in both sets.

x1	x2	x3
A	1	T
B	2	F
C	3	

`pd.merge(adf, bdf, how='outer', on='x1')`

Join data. Retain all values, all rows.

x1	x2
C	3

`adf[~adf.x1.isin(bdf.x1)]`

All rows in adf that do not have a match in bdf.

x1	x2
A	1
B	2
C	3

`pd.merge(ydf, zdf)`

Rows that appear in both ydf and zdf (Intersection).

x1	x2
A	1
B	2
C	3

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either or both ydf and zdf (Union).

x1	x2
A	1

`pd.merge(ydf, zdf, how='outer', indicator=True)`

`.query('_merge == "left_only"')`

`.drop(columns=['_merge'])`

Rows that appear in ydf but not zdf (Setdiff).