

# Cheat Sheet: The pandas DataFrame

## Preliminaries

### Start by importing these Python modules

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt # for charts
```

### Check which version of pandas you are using

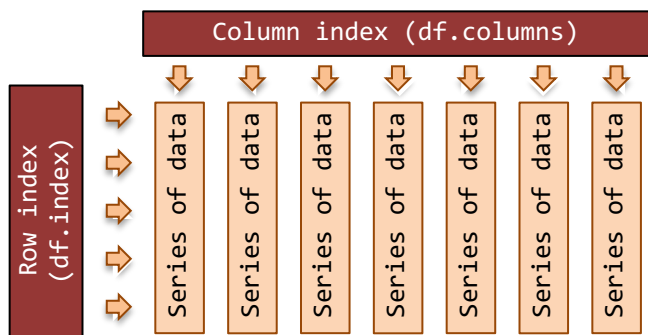
```
print(pd.__version__)
```

This cheat sheet was written for pandas version 0.25. It assumes you are using Python 3.

## The conceptual model

Pandas provides two important data types: the DataFrame and the Series.

A **DataFrame** is a two-dimensional table of data with column and row indexes (something like a spreadsheet). The columns are made up of Series objects.



### A DataFrame has two indexes:

- Typically, the column index (df.columns) is a list of strings (variable names) or (less commonly) integers
- Typically, the row index (df.index) might be:
  - Integers – for case or row numbers;
  - Strings – for case names; or
  - DatetimeIndex or PeriodIndex – for time series

A **Series** is an ordered, one-dimensional array of data with an index. All the data is of the same data type. Series arithmetic is vectorized after first aligning the Series index for each of the operands.

### Examples of Series Arithmetic

```
s1 = pd.Series(range(0, 4)) # 0, 1, 2, 3
s2 = pd.Series(range(1, 5)) # 1, 2, 3, 4
s3 = s1 + s2                # 1, 3, 5, 7
```

```
s4 = pd.Series([1, 2, 3], index=[0, 1, 2])
s5 = pd.Series([1, 2, 3], index=[2, 1, 0])
s6 = s4 + s5                # 4, 4, 4
```

```
s7 = pd.Series([1, 2, 3], index=[1, 2, 3])
s8 = pd.Series([1, 2, 3], index=[0, 1, 2])
s9 = s7 + s8                # NAN, 3, 5, NAN
```

## Get your data into a DataFrame

### Instantiate a DataFrame

```
df = pd.DataFrame() # the empty DataFrame
df = pd.DataFrame(python_dictionary)
df = pd.DataFrame(numpy_matrix)
```

### Load a DataFrame from a CSV file

```
df = pd.read_csv('file.csv', header=0,
                 index_col=0, quotechar='\"', sep=':',
                 na_values = ['na', '-', '.', ''])
```

**Note:** refer to pandas docs for all arguments

### Get your data from inline python CSV text

```
from io import StringIO
data = """Animal, Cuteness, Desirable
A, dog, 8.7, True
B, cat, 9.5, False"""
df = pd.read_csv(StringIO(data), header=0,
                 index_col=0, skipinitialspace=True)
```

**Note:** `skipinitialspace=True` allows a pretty layout

### Also, among many other options ...

```
df = pd.read_html(url/html_string)
df = pd.read_json(path/JSON_string)
df = pd.read_sql(query, connection)
df = pd.read_excel('filename.xlsx')
df = pd.read_clipboard() # eg from Excel copy
```

**Note:** See the pandas documentation for arguments.

### Fake up some random data – useful for testing

```
df = (pd.DataFrame(np.random.rand(1100, 6),
                  columns=list('ABCDEF')) - 0.5).cumsum()
df['Group'] = [np.random.choice(list('abcd'))
               for _ in range(len(df))]
df['Date'] = pd.date_range('1/1/2017',
                          periods=len(df), freq='D')
```

**Hint:** leave off the Group and/or Date cols if not needed

## Saving a DataFrame

### Saving a DataFrame to a CSV file

```
df.to_csv('filename.csv', encoding='utf-8')
```

### Saving a DataFrame to an Excel Workbook

```
writer = pd.ExcelWriter('filename.xlsx')
df.to_excel(writer, 'Sheet1')
writer.save()
```

### Saving a DataFrame to a Python object

```
d = df.to_dict() # to dictionary
m = df.values # to a numpy matrix
```

### Also, among many other options ...

```
html = df.to_html()
df.to_json()
df.to_sql()
df.to_clipboard() # then paste into Excel
```

## Working with the whole DataFrame

### Peek at the DataFrame contents/structure

```
df.info() # print cols & data types
dfh = df.head(i) # get first i rows
dft = df.tail(i) # get last i rows
dfs = df.describe() # summary stats for cols
top_left_corner_df = df.iloc[:4, :4]
```

### DataFrame non-indexing attributes

```
df = df.T # transpose rows and cols
l = df.axes # list of row & col indexes
(ri,ci) = df.axes # from above
s = df.dtypes # Series column data types
b = df.empty # True for empty DataFrame
i = df.ndim # number of axes (it is 2)
t = df.shape # (row-count, column-count)
i = df.size # row-count * column-count
a = df.values # get a numpy matrix for df
```

### DataFrame utility methods

```
df = df.copy() # copy a DataFrame
df = df.sort_values(by=col)
df = df.sort_values(by=[col1, col2])
df = df.sort_values(by=row, axis=1)
df = df.sort_index() # axis=1 to sort cols
df = df.astype(dtype) # type conversion
```

### DataFrame iteration methods

```
df.iteritems() # (col-index, Series) pairs
df.iterrows() # (row-index, Series) pairs
# example ... iterating over columns ...
for (name, series) in df.iteritems():
    print('\nCol name: ' + str(name))
    print('1st value: ' + str(series.iat[0]))
```

### Maths on the whole DataFrame (not a complete list)

```
df = df.abs() # absolute values
df = df.add(o) # add df, Series or value
s = df.count() # non NA/null values
df = df.cummax() # (cols default axis)
df = df.cummin() # (cols default axis)
df = df.cumsum() # (cols default axis)
df = df.diff() # 1st diff (col def axis)
df = df.div(o) # div by df, Series, value
df = df.dot(o) # matrix dot product
s = df.max() # max of axis (col def)
s = df.mean() # mean (col default axis)
s = df.median() # median (col default)
s = df.min() # min of axis (col def)
df = df.mul(o) # mul by df Series val
s = df.sum() # sum axis (cols default)
df = df.where(df > 0.5, other=np.nan)
```

**Note:** methods returning a series default to work on cols

### Select/filter rows/cols based on index label values

```
df = df.filter(items=['a', 'b']) # by col
df = df.filter(items=[5], axis=0) # by row
df = df.filter(like='x') # keep x in col
df = df.filter(regex='x') # regex in col
df = df.select(lambda x: not x%5) # 5th rows
```

**Note:** select takes a Boolean function, for cols: axis=1

**Note:** filter defaults to cols; select defaults to rows

## Working with Columns (and pandas Series)

### Peek at the column/Series structure/contents

```
s = df[col].head(i) # get first i elements
s = df[col].tail(i) # get last i elements
s = df[col].describe() # summary stats
```

### Get column index and labels

```
idx = df.columns # get col index
label = df.columns[0] # first col label
l = df.columns.tolist() # list of col labels
a = df.columns.values # array of col labels
```

### Change column labels

```
df = df.rename(columns={'old':'new', 'a':'1'})
df.columns = ['new1', 'new2', 'new3'] # etc.
```

### Selecting columns

```
s = df[col] # select col to Series
df = df[[col]] # select col to df
df = df[[a, b]] # select 2-plus cols
df = df[[c, a, b]] # change col order
s = df[df.columns[0]] # select by number
df = df[df.columns[[0, 3, 4]]] # by numbers
df = df[df.columns[:-1]] # all but last col
s = df.pop(col) # get & drop from df
```

### Selecting columns with Python attributes

```
s = df.a # same as s = df['a']
df.existing_column = df.a / df.b
df['new_column'] = df.a / df.b
```

**Trap:** column names must be valid Python identifiers, but not a DataFrame method or attribute name.

**Trap:** cannot create new columns

**Hint:** Don't be lazy: for clearer code avoid dot notation.

### Adding new columns to a DataFrame

```
df['new_col'] = range(len(df))
df['new_col'] = np.repeat(np.nan, len(df))
df['random'] = np.random.rand(len(df))
df['index_as_col'] = df.index
df1[['b', 'c']] = df2[['e', 'f']]
```

**Trap:** When adding a new column, only items from the new column series that have a corresponding index in the DataFrame will be added. The index of the receiving DataFrame is not extended to accommodate all of the new series.

**Trap:** when adding a python list or numpy array, the column will be added by integer position.

### Add a mismatched column with an extended index

```
df = pd.DataFrame([1, 2, 3], index=[1, 2, 3])
s = pd.Series([2, 3, 4], index=[2, 3, 4])
df = df.reindex(df.index.union(s.index))
df['s'] = s # with NaNs where no data
```

**Note:** assumes unique index values

### Dropping (deleting) columns (mostly by label)

```
df = df.drop(col1, axis=1)
df = df.drop([col1, col2], axis=1)
del df[col] # even classic python works
df = df.drop(df.columns[0], axis=1) #first
df = df.drop(df.columns[-1:], axis=1) #last
```

### Swap column contents

```
df[['B', 'A']] = df[['A', 'B']]
```

### Vectorised arithmetic on columns

```
df['proportion'] = df['count'] / df['total']  
df['percent'] = df['proportion'] * 100.0
```

### Apply numpy mathematical functions to columns

```
df['log_data'] = np.log(df[col])
```

**Note:** many many more numpy math functions

**Hint:** Prefer pandas math over numpy where you can.

### Set column values set based on criteria

```
df[b] = df[a].where(df[a]>0, other=0)  
df[d] = df[a].where(df.b!=0, other=df.c)
```

**Note:** where other can be a Series or a scalar

### Data type conversions

```
s = df[col].astype('float')  
s = df[col].astype('int')  
s = pd.to_numeric(df[col])  
s = df[col].astype('str')  
a = df[col].values          # numpy array  
l = df[col].tolist()       # python list
```

**Trap:** index lost in conversion from Series to array or list

### Common column-wide methods/attributes

```
value = df[col].dtype      # type of data  
value = df[col].size      # col dimensions  
value = df[col].count()   # non-NA count  
value = df[col].sum()  
value = df[col].prod()  
value = df[col].min()  
value = df[col].max()  
value = df[col].mean()    # also median()  
value = df[col].cov(df[other_col])  
s = df[col].describe()  
s = df[col].value_counts()
```

### Find first row index label for min/max val in column

```
label = df[col].idxmin()  
label = df[col].idxmax()
```

### Common column element-wise methods

```
s = df[col].isna()  
s = df[col].notna()      # not isna()  
s = df[col].astype('float')  
s = df[col].abs()  
s = df[col].round(decimals=0)  
s = df[col].diff(periods=1)  
s = df[col].shift(periods=1)  
s = df[col].to_datetime()  
s = df[col].fillna(0)    # replace NaN w 0  
s = df[col].cumsum()  
s = df[col].cumprod()  
s = df[col].pct_change(periods=4)  
s = df[col].rolling(window=4,  
                      min_periods=4, center=False).sum()
```

### Append a column of row sums to a DataFrame

```
df['Row Total'] = df.sum(axis=1)
```

**Note:** also means, mins, maxs, etc.

### Multiply every column in DataFrame by a Series

```
df = df.mul(s, axis=0) # on matched rows
```

**Note:** also add, sub, div, etc.

### Selecting columns with .loc, .iloc

```
df = df.loc[:, 'col1':'col2'] # inclusive  
df = df.iloc[:, 0:2]         # exclusive
```

### Get the integer position of a column index label

```
i = df.columns.get_loc('col_name')
```

### Test if column index values are unique/monotonic

```
if df.columns.is_unique: pass # ...  
b = df.columns.is_monotonic_increasing  
b = df.columns.is_monotonic_decreasing
```

### Mapping a DataFrame column or Series

```
map = pd.Series(['red', 'green', 'blue'],  
                index=['r', 'g', 'b'])  
s = pd.Series(['r', 'g', 'r', 'b']).map(map)  
# s contains: ['red', 'green', 'red', 'blue']  
  
m = pd.Series([True, False], index=['Y', 'N'])  
df = pd.DataFrame(np.random.choice(list('YN'),  
                                   500, replace=True), columns=[col])  
df[col] = df[col].map(m)
```

**Note:** Useful for decoding data before plotting

**Note:** Sometimes referred to as a lookup function

**Note:** Indexes can also be mapped if needed.

### Find the largest and smallest values in a column

```
s = df[col].nlargest(n)  
s = df[col].nsmallest(n)
```

### Sorting the columns of a DataFrame

```
df = df.sort_index(axis=1, ascending=False)
```

**Note:** the column labels need to be comparable

## Working with rows

### Get the row index and labels

```
idx = df.index # get row index
label = df.index[0] # first row label
label = df.index[-1] # last row label
l = df.index.tolist() # get as a python list
a = df.index.values # get as numpy array
```

### Change the (row) index

```
df.index = idx # new ad hoc index
df = df.set_index('A') # index set to col A
df = df.set_index(['A', 'B']) # MultiIndex
df = df.reset_index() # replace old w new
# note: old index stored as a col in df
df.index = range(len(df)) # set with list
df = df.reindex(index=range(len(df)))
df = df.set_index(keys=['r1', 'r2', 'etc'])
```

### Adding rows

```
df = original_df.append(more_rows_in_df)
```

**Hint:** convert row(s) to a DataFrame and then append. Both DataFrames must have same column labels.

### Append a row of column totals to a DataFrame

```
df.loc['Total'] = df.sum()
```

**Note:** best if all columns are numeric

### Iterating over DataFrame rows

```
for (index, row) in df.iterrows(): # pass
```

**Trap:** row data may be coerced to the same data type

### Dropping rows (by label)

```
df = df.drop(row)
df = df.drop([row1, row2]) # multi-row drop
```

### Row selection by Boolean series

```
df = df.loc[df[col] >= 0.0]
df = df.loc[(df[col] >= 1.0) | (df[col] < 0.0)]
df = df.loc[df[col].isin([1, 2, 5, 7, 11])]
df = df.loc[~df[col].isin([1, 2, 5, 7, 11])]
df = df.loc[df[col].str.contains('a')]
```

**Hint:** while the .loc[] accessor may not be needed above, its use makes for more readable code.

**Trap:** bitwise "or", "and" "not; (ie. | & ~) co-opted to be Boolean operators on a Series of Boolean; therefore, you need parentheses around comparisons.

### Selecting rows using isin over multiple columns

```
# fake up some data
data = {1:[1,2,3], 2:[1,4,9], 3:[1,8,27]}
df = pd.DataFrame(data)
```

```
# multi-column isin
lf = {1:[1, 3], 3:[8, 27]} # look for
f = df.loc[df[list(lf)].isin(lf).all(axis=1)]
```

### Selecting rows using an index

```
idx = df[df[col] >= 2].index
print(df.loc[idx])
```

### Select a slice of rows by integer position

[inclusive-from : exclusive-to [: step]]

start is 0; end is len(df)

```
df = df.iloc[:] # copy entire DataFrame
df = df.iloc[0:2] # rows 0 and 1
df = df.iloc[2:3] # row 2 (the third row)
df = df.iloc[-1:] # the last row
df = df.iloc[:-1] # all but the last row
df = df.iloc[:,2] # every 2nd row (0 2 ..)
```

**Hint:** while the .iloc[] accessor may not be needed above, its use makes for more readable code.

### Select a slice of rows by label/index

```
df = df.loc['a':'c'] # rows 'a' through 'c'
```

**Note:** [inclusive-from : inclusive-to [: step]]

**Hint:** while the .loc[] accessor may not be needed above, its use makes for more readable code.

### Sorting the rows of a DataFrame by the row index

```
df = df.sort_index(ascending=False)
```

### Sorting DataFrame rows based on column values

```
df = df.sort_values(by=df.columns[0],
                    ascending=False)
df = df.sort_values(by=[col1, col2])
```

### Random selection of rows

```
import random
k = 20 # pick a number
selection = random.sample(range(len(df)), k)
df_sample = df.iloc[selection, :] # get copy
```

**Note:** this randomly selected sample is not sorted

### Drop duplicates in the row index

```
df['index'] = df.index # 1 create new col
df = df.drop_duplicates(cols='index',
                       take_last=True) # 2 use new col
del df['index'] # 3 del the col
df = df.sort_index() # 4 tidy up
```

### Test if two DataFrames have same row index

```
len(a) == len(b) and all(a.index == b.index)
```

**Note:** you may want to sort indexes first.

### Get the integer position of a row or col index label

```
i = df.index.get_loc(row_label)
```

**Trap:** index.get\_loc() returns an integer for a unique match. If not a unique match, may return a slice/mask.

### Get integer position of rows that meet condition

```
a = np.where(df[col] >= 2) #numpy array
```

### Test if the row index values are unique/monotonic

```
if df.index.is_unique: pass # ...
b = df.index.is_monotonic_increasing
b = df.index.is_monotonic_decreasing
```

### Find row index duplicates

```
if df.index.has_duplicates:
    print(df.index.duplicated())
```

**Note:** also similar for column label duplicates.

## Working with cells

### Getting a cell by row and column labels

```
value = df.at[row, col]
value = df.loc[row, col]
```

**Note:** .at[] fastest label based scalar lookup

**Note:** at[] does not take slices as an argument

### Setting a cell by row and column labels

```
df.at[row, col] = value
df.loc[row, col] = value
```

**Avoid:** chaining in the form df[col][row]

**Avoid:** chaining in the form df[col].at[row]

### Getting and slicing on labels

```
df = df.loc['row1':'row3', 'col1':'col3']
```

**Note:** the "to" on this slice is inclusive.

### Setting a cross-section by labels

```
df.loc['A':'C', 'col1':'col3'] = np.nan
df.loc[1:2, 'col1':'col2'] = np.zeros((2,2))
df.loc[1:2, 'A':'C'] = other.loc[1:2, 'A':'C']
```

**Remember:** inclusive "to" in the slice

### Getting a cell by integer position

```
value = df.iat[9, 3] # [row, col]
value = df.iat[len(df)-1, len(df.columns)-1]
```

### Getting a range of cells by int position

```
df = df.iloc[2:4, 2:4] # subset of the df
df = df.iloc[:5, :5] # top left corner
s = df.iloc[5, :] # return row as Series
df = df.iloc[5:6, :] # returns row as row
```

**Note:** exclusive "to" – same as python list slicing.

### Setting cell by integer position

```
df.iat[7, 8] = value
```

### Setting cell range by integer position

```
df.iloc[0:3, 0:5] = value
df.iloc[1:3, 1:4] = np.ones((2, 3))
df.iloc[1:3, 1:4] = np.zeros((2, 3))
df.iloc[1:3, 1:4] = np.array([[1, 1, 1],
                              [2, 2, 2]])
```

**Remember:** exclusive-to in the slice

### Views and copies

**From the manual:** Setting a copy can cause subtle errors. The rules about when a view on the data is returned are dependent on NumPy. Whenever an array of labels or a Boolean vector are involved in the indexing operation, the result will be a copy.

**Hint:** Pandas will usually warn you if you are trying to set a copy. Take these warnings seriously.

## Summary: selection using the DataFrame index

### Select columns with []

```
s = df[col] # returns Series
df = df[[col]] # returns DataFrame
df = df[[col1, col2]] # select cols with list
df = df[idx] # select cols with an index
df = df[s] # select with col label Series
```

**Note:** scalar returns Series; list & c returns a DataFrame

**Trap:** With [] indexing, a label Series gets/sets columns, but a Boolean Series gets/sets rows

### Select rows with .loc[] or .iloc[]

```
df = df.loc[df[col]>0.5] # Boolean Series
df = df.loc['label'] # single label
df = df.loc[container] # lab list/Series
df = df.loc['from':'to'] # inclusive slice
df = df.loc[bs] # Boolean Series
df = df.iloc[0] # single integer
df = df.iloc[container] # int list/Series
df = df.iloc[0:5] # exclusive slice
```

**Hint:** Always use .loc[] or .iloc[] when selecting rows

**Hint:** Never use the deprecated .ix[] indexer

### Select individual cells with .at[,] or .iat[,]

```
v = df.at[r, c] # fast scalar label accessor
v = df.iat[r, c] # fast scalar int accessor
```

### Select a cross-section with .loc[,] or .iloc[,]

```
# r and c can be scalar, list, slice
xs = df.loc[r, c] # label accessor (row, col)
xs = df.iloc[r, c] # integer accessor
```

### DataFrame indexing methods

```
v = df.get_value(r, c) # get by row, col
df = df.set_value(r,c,v) # set by row, col
df = df.xs(key, axis) # get cross-section
df = df.filter(items, like, regex, axis)
df = df.select(crit, axis)
```

**Note:** the indexing attributes (.loc[], .iloc[], .at[] .iat[]) can be used to get and set values in the DataFrame

**Note:** the .loc[], and .iloc[] indexing attributes can accept python slice objects. But .at[] and .iat[] do not

**Note:** .loc[] can also accept Boolean Series arguments

**Avoid:** chaining in the form df[col\_indexer][row\_indexer]

**Trap:** label slices are inclusive, integer slices exclusive

### Some index attributes and methods

```
b = idx.is_monotonic_decreasing
b = idx.is_monotonic_increasing
b = idx.has_duplicates
i = idx.nlevels # num of index levels
idx = idx.astype(dtype) # change data type
b = idx.equals(o) # check for equality
idx = idx.union(o) # union of two indexes
i = idx.nunique() # number unique labels
label = idx.min() # minimum label
label = idx.max() # maximum label
```

## Joining/Combining DataFrames

Three ways to join DataFrames:

- **concat** – concatenate (or stack) two DataFrames side by side, or one on top of the other
- **merge** – using a database-like join operation on side-by-side DataFrames
- **combine first** – splice two DataFrames together, choosing values from one over the other

### Simple concatenation is often what you want

```
df = pd.concat([df1,df2], axis=0) #top/bottom
df = pd.concat([df1,df2]).sort_index() # t/b
```

```
df = pd.concat([df1,df2], axis=1) #left/right
```

**Trap:** can end up with duplicate rows or cols

**Note:** concat has an ignore\_index parameter

**Note:** if no axis is specified, defaults to top/bottom.

### Append (another way of doing a top/bottom concat)

```
df = df1.append(df2) #top/bottom
df = df1.append([df2, df3]) #top/bottom
```

**Note:** append also has an ignore\_index parameter

### Merge

```
df_new = pd.merge(left=df1, right=df2,
                  how='outer', left_index=True,
                  right_index=True) # on indexes
```

```
df_new = pd.merge(left=df1, right=df2,
                  how='left', left_on='col1',
                  right_on='col2') # on columns
```

```
df_new = df.merge(right=dfg, how='left',
                  left_on='Group', right_index=True)
```

**How:** 'left', 'right', 'outer', 'inner' (where outer=union/all; inner=intersection)

**Note:** merge is both a pandas helper-function, and a DataFrame method

**Note:** DataFrame.merge() joins on common columns by default (if left and right not specified)

**Trap:** When joining on column values, the indexes on the passed DataFrames are ignored.

**Trap:** many-to-many merges can result in an explosion of associated data.

### Join on row indexes (another way of merging)

```
df = df1.join(other=df2, how='outer')
df = df1.join(other=df2, on=['a','b'],
              how='outer')
```

**Note:** DataFrame.join() joins on indexes by default.

### Combine first

```
df = df1.combine_first(other=df2)
```

# multi-combine with python reduce()

```
df = reduce(lambda x, y:
            x.combine_first(other=y),
            [df1, df2, df3, df4, df5])
```

Combine\_first uses the non-null values from df1. Null values in df1 are filled with values from the same location in df2. The index of the combined DataFrame will be the union of the indexes from df1 and df2.

## Groupby: Split-Apply-Combine

### Grouping

```
gb = df.groupby(col) # by one columns
gb = df.groupby([col1, col2]) # by 2 cols
gb = df.groupby(level=0) # row index groupby
gb = df.groupby(level=['a','b']) #mult-idx gb
print(gb.groups)
```

**Note:** groupby() returns a pandas groupby object

**Note:** the groupby object attribute .groups contains a dictionary mapping of the groups.

**Trap:** NaN values in the group key are automatically dropped – there will never be a NA group.

### Applying an aggregating function

# apply to a single column ...

```
s = gb[col].sum()
s = gb[col].agg(np.sum)
```

# apply to every column in DataFrame ...

```
s = gb.count()
df_summary = gb.describe()
df_row_1s = gb.first()
```

Note: aggregating functions include mean, sum, size, count, std, var, sem (standard error of the mean), describe, first, last, min, max

### Applying multiple aggregating functions

```
# apply multiple functions to one column
dfx = gb['col2'].agg([np.sum, np.mean])
# apply to multiple fns to multiple cols
dfy = gb.agg({
    'cat': np.count_nonzero,
    'col1': [np.sum, np.mean, np.std],
    'col2': [np.min, np.max]
})
```

**Note:** gb['col2'] above is shorthand for df.groupby('cat')['col2'], without the need for regrouping.

### Applying transform functions

```
# transform to group z-scores, which have
# a group mean of 0, and a std dev of 1.
zscore = lambda x: (x-x.mean())/x.std()
dfz = gb.transform(zscore)
```

```
# replace missing data with the group mean
mean_r = lambda x: x.fillna(x.mean())
df = gb.transform(mean_r) # entire DataFrame
df[col] = gb[col].transform(mean_r) # one col
```

**Note:** can apply multiple transforming functions in a manner similar to multiple aggregating functions above,

### Applying filtering functions

Filtering functions allow you to make selections based on whether each group meets specified criteria

```
# select groups with more than 10 members
eleven = lambda x: (len(x['col1']) >= 11)
df11 = gb.filter(eleven)
```

### Group by a row index (non-hierarchical index)

```
df = df.set_index(keys='cat')
s = df.groupby(level=0)[col].sum()
dfg = df.groupby(level=0).sum()
```

## Pivot Tables: working with long and wide data

These features work with and often create hierarchical or multi-level Indexes; (the pandas MultiIndex is powerful and complex).

### Pivot, unstack, stack and melt

Pivot tables move from long format to wide format data

```
# Let's start with data in long format
from io import StringIO
data = """Date,Pollster,State,Party,Est
13/03/2014, Newspoll, NSW, red, 25
13/03/2014, Newspoll, NSW, blue, 28
13/03/2014, Newspoll, Vic, red, 24
13/03/2014, Newspoll, Vic, blue, 23
13/03/2014, Galaxy, NSW, red, 23
13/03/2014, Galaxy, NSW, blue, 24
13/03/2014, Galaxy, Vic, red, 26
13/03/2014, Galaxy, Vic, blue, 25
13/03/2014, Galaxy, Qld, red, 21
13/03/2014, Galaxy, Qld, blue, 27"""
df = pd.read_csv(StringIO(data),
                 header=0, skipinitialspace=True)

# pivot to wide format on 'Party' column
# 1st: set up a MultiIndex for other cols
df1 = df.set_index(['Date', 'Pollster',
                   'State'])
# 2nd: do the pivot
wide1 = df1.pivot(columns='Party')

# unstack to wide format on State / Party
# 1st: MultiIndex all but the Values col
df2 = df.set_index(['Date', 'Pollster',
                   'State', 'Party'])
# 2nd: unstack a column to go wide on it
wide2 = df2.unstack('State')
wide3 = df2.unstack() # pop last index

# Use stack() to get back to long format
long1 = wide1.stack()
# Then use reset_index() to remove the
# MultiIndex.
long2 = long1.reset_index()

# Or melt() back to long format
# 1st: flatten the column index
wide1.columns = ['_'.join(col).strip()
                 for col in wide1.columns.values]
# 2nd: remove the MultiIndex
wdf = wide1.reset_index()
# 3rd: melt away
long3 = pd.melt(wdf, value_vars=
               ['Est_blue', 'Est_red'],
               var_name='Party', id_vars=['Date',
               'Pollster', 'State'])
```

**Note:** See documentation, there are many arguments to these methods.

## Working with dates, times and their indexes

**Dates and time – points, spans, deltas and offsets**  
Pandas has four date-time like objects that can be used for data in a Series or in an Index:

Concept	Data	Index
Point	Timestamp	DatetimeIndex
Span	Period	PeriodIndex
Delta	Timedelta	TimedeltaIndex
Offset	DateOffset	None

### Timestamps

Timestamps represent a point in time.

```
t = pd.Timestamp('2019-01-01')
t = pd.Timestamp('2019-01-01 21:15:06')
t = pd.Timestamp('2019-01-01 21:15:06.7')
t = pd.Timestamp(year=2019, month=1,
                 day=1, hour=21, minute=15, second=6.7,
                 tz='Australia/Sydney')
# handles daylight savings time
```

**Note:** Timestamps can range from 1678 to 2261. (Check out `pd.Timestamp.max` and `pd.Timestamp.min`).

**Note:** the dtype is `datetime64[ns]` or `datetime64[ns,tz]`

### DatetimeIndex – an Index of Timestamps

```
l = ['2019-04-01', '2019-04-02']
dti = pd.to_datetime(l)
l2 = (['01-01-2019', '01-02-2019'])
dti2 = pd.to_datetime(l2, dayfirst=True)
```

### A Series of Timestamps

```
l = ['2019-04-01', '2019-04-02']
s = pd.to_datetime(pd.Series(l))
```

**Note:** if we pass the `pd.to_datetime()` helper function a pandas Series it returns a Series of Timestamps.

### From non-standard strings to DatetimeIndex

```
t = ['09:08:55.7654-JAN092002',
     '15:42:02.6589-FEB082016']
s = pd.to_datetime(t,
                  format="%H:%M:%S.%f-%b%d%Y")
```

**Also:** %B = full month name; %m = numeric month; %y = year without century; and more ...

### A range of Timestamps in a DatetimeIndex

```
dti = pd.date_range('2015-01',
                    periods=len(df), freq='M') # end of month
dti = pd.date_range('2019-01-01',
                    periods=365, freq='D')
```

### Timestamps and DatetimeIndex from columns

```
# fake up a DataFrame
y = [2019, 2019, 2019]
m = [2, 3, 4]
d = [1, 2, 2]
df = pd.DataFrame({'yr':y, 'mon':m, 'day':d})

# do the magic
cols = ['yr', 'mon', 'day']
df.index = pd.to_datetime(df[cols])
df['TS'] = pd.to_datetime(df[cols])
```

## From DatetimeIndex to Python datetime objects

```
dti = pd.DatetimeIndex(pd.date_range(
    start='1/1/2011', periods=4, freq='M'))
s = Series([1,2,3,4], index=dti)
a = dti.to_pydatetime() # numpy array
a = s.index.to_pydatetime() # numpy array
```

## From Timestamps to Python dates or times

```
df['py_date'] = [x.date() for x in df['TS']]
df['py_time'] = [x.time() for x in df['TS']]
```

**Note:** converts to datetime.date or datetime.time. But does not convert to datetime.datetime.

## Periods

Periods represent a time-span.

```
p = pd.Period('2019', freq='Y')
p = pd.Period('2019-01', freq='M')
p = pd.Period('2019-01-01', freq='D')
p = pd.Period('2019-01-01 21:15:06', freq='S')
```

## From Timestamps to Periods in a Series

```
l = ['2019-04-01', '2019-04-02']
ts = pd.to_datetime(pd.Series(l))
ps = ts.dt.to_period(freq='D')
```

**Note:** the .dt accessor in the last line

## From a DatetimeIndex to a PeriodIndex

```
l = ['2019-04-01', '2019-04-02']
dti = pd.to_datetime(l)
pi = dti.to_period(freq='D')
```

**Hint:** unless you are working in less than seconds, prefer PeriodIndex over DatetimeIndex.

## A range of Periods in a PeriodIndex

```
pi = pd.period_range('2015-01',
    periods=len(df), freq='M')
pi = pd.period_range('2019-01-01',
    periods=365, freq='D')
```

## Working with a PeriodIndex

```
pi = pd.period_range('1960-01', '2015-12',
    freq='M')
a = pi.values # numpy array of integers
p = pi.tolist() # python list of Periods
sp = pd.Series(pi) # pandas Series of Periods
s = pd.Series(pi).astype('str')
l = pd.Series(pi).astype('str').tolist()
```

## From DatetimeIndex to PeriodIndex and back

```
df = pd.DataFrame(np.random.randn(20,3))
df.index = pd.date_range('2015-01-01',
    periods=len(df), freq='M')
dfp = df.to_period(freq='M')
dft = dfp.to_timestamp()
```

**Note:** from period to timestamp defaults to the point in time at the start of the period.

## The tail of a time-series DataFrame

```
df = df.last("5M") # the last five months
```

## Period frequency constants (not a complete list)

Name	Description
U	Microsecond
L	Millisecond
S	Second
T	Minute
H	Hour
D	Calendar day
B	Business day
W-{MON, TUE, ...}	Week ending on ...
MS	Calendar start of month
M	Calendar end of month
QS-{JAN, FEB, ...}	Quarter start with year starting (QS - December)
Q-{JAN, FEB, ...}	Quarter end with year ending (Q - December)
AS-{JAN, FEB, ...}	Year start (AS - December)
A-{JAN, FEB, ...}	Year end (A - December)

## Deltas

When we subtract a Timestamp from another Timestamp, we get a Timedelta object in pandas.

```
ts = pd.Series(pd.date_range('2019-01-01',
    periods=31, freq='D'))
delta_series = ts.diff(1)
```

## Converting a Timedelta to a numeric

```
l = ['2019-04-01', '2019-09-03']
s = pd.to_datetime(pd.Series(l))
delta = s[1] - s[0]
```

```
day = pd.Timedelta(days=1)
delta_num = delta / day
minute = pd.Timedelta(minutes=1)
delta_num2 = delta / minute
```

## Offsets

Subtracting a Period from a Period gives an offset.

```
offset = pd.DateOffset(days=4)
s = pd.Series(pd.period_range('2019-01-01',
    periods=365, freq='D'))
offset2 = s[4] - s[0]
s = s.diff(1) # s is now a series of offsets
```

## Converting an Offset to a numeric

```
x = offset.n # an individual offset
t = s.apply(lambda z: np.nan if z is np.nan
    else z.n) # convert a Series
```

## Upsampling

```
# fake up some quarterly count data
pi = pd.period_range('1960Q1',
    periods=220, freq='Q')
df = pd.DataFrame(np.random.randint(low=0,
    high=999, size=(len(pi), 5)), index=pi)

# which we can upsample to monthly count data
dfm = df.resample('M').asfreq() # with NAs!
dfm2 = (df.resample('M').asfreq().fillna(0)
    .rolling(window=3, min_periods=3).mean()
    .bfill(limit=2)) # assuming no NA data
```

**Note:** df.resample(arguments).aggregating\_function(). There are lots of options here. See the manual.



## Downsampling

```
# downsample from monthly to quarterly counts
dfq = dfm.resample('Q').sum()
```

**Note:** `df.resample(arguments).aggregating_function()`.

## Time zones

```
t = ['2015-06-30 00:00:00',
     '2015-12-31 00:00:00']
dti = pd.to_datetime(t)
dti = dti.tz_localize('Australia/Canberra')
dti = dti.tz_convert('UTC')
ts = pd.Timestamp('now',
                  tz='Europe/London')
```

**Note:** by default, Timestamps are created without time zone information.

## Row selection with a time-series index

```
# start with some play data
n = 48
df = pd.DataFrame(np.random.randint(low=0,
                                   high=999, size=(n, 5)),
                  index=pd.period_range('2015-01',
                                       periods=n, freq='M'))
```

```
february_selector = (df.index.month == 2)
february_data = df[february_selector]
```

```
q1_data = df[(df.index.month >= 1) &
             (df.index.month <= 3)]
```

```
mayornov_data = df[(df.index.month == 5) |
                   (df.index.month == 11)]
```

```
year_totals = df.groupby(df.index.year).sum()
```

**Also:** year, month, day [of month], hour, minute, second, dayofweek, weekofmonth, weekofyear [numbered from 1], week starts on Monday], dayofyear [from 1], ...

## The Series.dt accessor attribute

DataFrame columns that contain datetime-like objects can be manipulated with the `.dt` accessor attribute

```
t = ['2012-04-14 04:06:56.307000',
     '2011-05-14 06:14:24.457000',
     '2010-06-14 08:23:07.520000']
```

```
# a Series of time stamps
s = pd.Series(pd.to_datetime(t))
print(s.dtype)      # datetime64[ns]
print(s.dt.second) # 56, 24, 7
print(s.dt.month)  # 4, 5, 6
# a Series of time periods
s = pd.Series(pd.PeriodIndex(t, freq='Q'))
print(s.dtype)     # int64
print(s.dt.quarter) # 2, 2, 2
print(s.dt.year)   # 2012, 2011, 2010
```

## Plotting from the DataFrame

### Import matplotlib, choose a matplotlib style

```
import matplotlib.pyplot as plt
print(plt.style.available)
plt.style.use('ggplot')
```

### Fake up some data (which we reuse repeatedly)

```
a = np.random.normal(0, 1, 999)
b = np.random.normal(1, 2, 999)
c = np.random.normal(2, 3, 999)
df = pd.DataFrame([a, b, c]).T
df.columns = ['A', 'B', 'C']
```

### Line plot

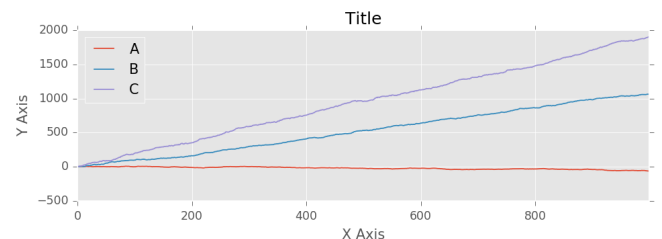
```
df1 = df.cumsum()
ax = df1.plot()
```

### # from here down - standard plot output

```
ax.set_title('Title')
ax.set_xlabel('X Axis')
ax.set_ylabel('Y Axis')
```

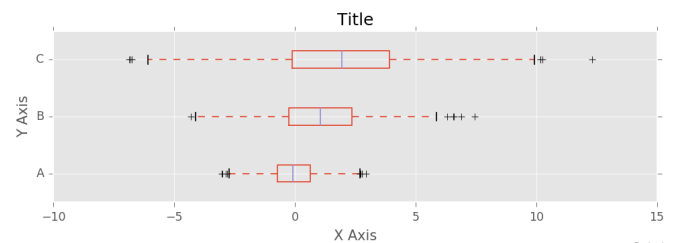
```
fig = ax.figure
fig.set_size_inches(8, 3)
fig.tight_layout(pad=1)
fig.savefig('filename.png', dpi=125)
```

```
plt.close()
```



### Box plot

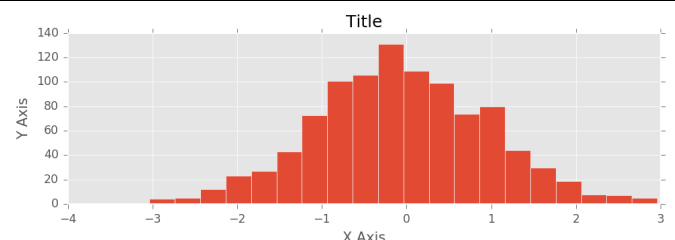
```
ax = df.plot.box(vert=False)
# followed by the standard plot code as above
```



```
ax = df.plot.box(column='c1', by='c2')
```

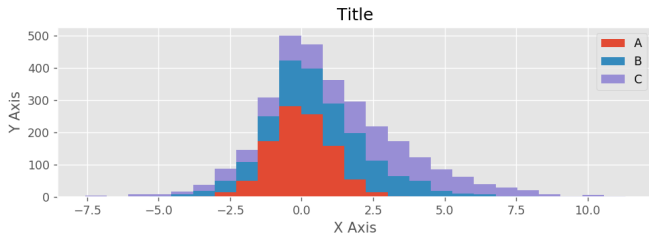
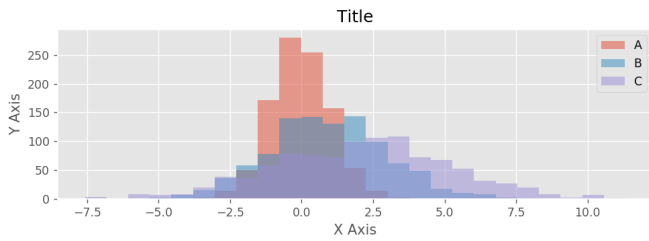
### Histogram

```
ax = df['A'].plot.hist(bins=20)
# followed by the standard plot code as above
```



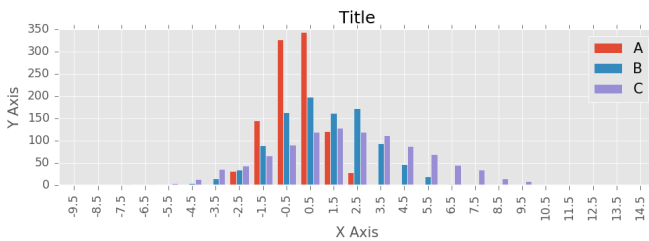
## Multiple histograms (overlapping or stacked)

```
ax = df.plot.hist(bins=25, alpha=0.5) # or...
ax = df.plot.hist(bins=25, stacked=True)
# followed by the standard plot code as above
```



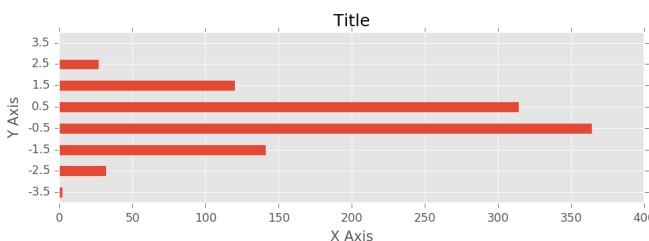
## Bar plots

```
bins = np.linspace(-10, 15, 26)
binned = pd.DataFrame()
for x in df.columns:
    y=pd.cut(df[x],bins,labels=bins[:-1])
    y=y.value_counts().sort_index()
    binned = pd.concat([binned,y],axis=1)
binned.index = binned.index.astype('float')
binned.index += (np.diff(bins) / 2.0)
ax = binned.plot.bar(stacked=False,
                    width=0.8) # for bar width
# followed by the standard plot code as above
```



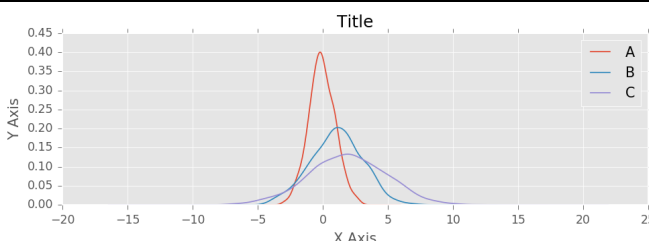
## Horizontal bars

```
ax = binned['A'][(binned.index >= -4) &
                (binned.index <= 4)].plot.barh()
# followed by the standard plot code as above
```



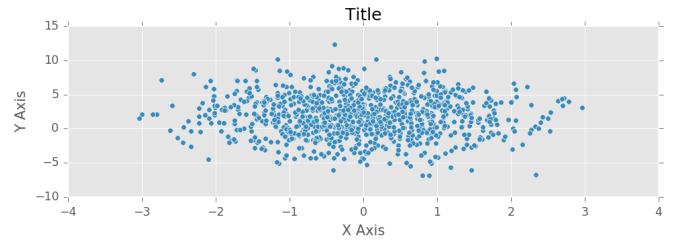
## Density plot

```
ax = df.plot.kde()
# followed by the standard plot code as above
```



## Scatter plot

```
ax = df.plot.scatter(x='A', y='C')
# followed by the standard plot code as above
```



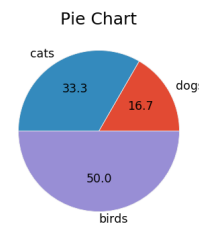
## Pie chart

```
s = pd.Series(data=[10, 20, 30],
              index = ['dogs', 'cats', 'birds'])
ax = s.plot.pie(autopct='%1.1f')

# followed by the standard plot output ...
ax.set_title('Pie Chart')
ax.set_aspect(1) # make it round
ax.set_ylabel('') # remove default

fig = ax.figure
fig.set_size_inches(8, 3)
fig.savefig('filename.png', dpi=125)

plt.close(fig)
```



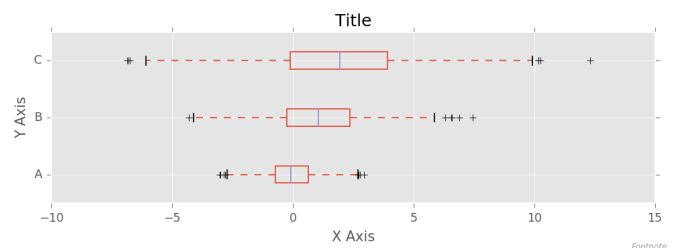
## Change the range plotted

```
ax.set_xlim([-5, 5])

# for some white space on the chart ...
lower, upper = ax.get_ylim()
ax.set_ylim([lower-1, upper+1])
```

## Add a footnote to the chart

```
# after the fig.tight_layout(pad=1) above
fig.text(0.99, 0.01, 'Footnote',
        ha='right', va='bottom',
        fontsize='x-small',
        fontstyle='italic', color='#999999')
```



## A line and bar on the same chart

In matplotlib, bar charts visualise categorical or discrete data. Line charts visualise continuous data. This makes it hard to get bars and lines on the same chart. Typically combined charts either have too many labels, and/or the lines and bars are misaligned or missing. You need to trick matplotlib a bit ... pandas makes this tricking easier

```
# start with fake percentage growth data
s = pd.Series(np.random.normal(
    1.02, 0.015, 40))
s = s.cumprod()
dfg = (pd.concat([s / s.shift(1),
    s / s.shift(4)], axis=1) * 100) - 100
dfg.columns = ['Quarter', 'Annual']
dfg.index = pd.period_range('2010-Q1',
    periods=len(dfg), freq='Q')

# reindex with integers from 0; keep old
old = dfg.index
dfg.index = range(len(dfg))

# plot the line from pandas
ax = dfg['Annual'].plot(color='blue',
    label='Year/Year Growth')

# plot the bars from pandas
dfg['Quarter'].plot.bar(ax=ax,
    label='Q/Q Growth', width=0.8)

# relabel the x-axis more appropriately
ticks = dfg.index[((dfg.index+0)%4)==0]
labs = pd.Series(old[ticks]).astype('str')
ax.set_xticks(ticks)
ax.set_xticklabels(labs.str.replace('Q',
    '\nQ'), rotation=0)

# fix the range of the x-axis ... skip 1st
ax.set_xlim([0.5, len(dfg)-0.5])

# add the legend
l=ax.legend(loc='best', fontsize='small')

# finish off and plot in the usual manner
ax.set_title('Fake Growth Data')
ax.set_xlabel('Quarter')
ax.set_ylabel('Per cent')

fig = ax.figure
fig.set_size_inches(8, 3)
fig.tight_layout(pad=1)
fig.savefig('filename.png', dpi=125)

plt.close()
```

## Working with missing and non-finite data

### Working with missing data

Pandas uses the not-a-number construct (np.nan and float('nan')) to indicate missing data. The Python None can arise in data as well. It is also treated as missing data; as is the pandas not-a-time construct (pandas.NaT).

### Missing data in a Series

```
s = pd.Series([8, None, float('nan'), np.nan])
# [8, NaN, NaN, NaN]
s.isna() # [False, True, True, True]
s.notna() # [True, False, False, False]
s.fillna(0) # [8, 0, 0, 0]
```

### Missing data in a DataFrame

```
df = df.dropna() # drop all rows with NaN
df = df.dropna(axis=1) # same for cols
df = df.dropna(how='all') # drop all NaN row
df = df.dropna(thresh=2) # drop 2+ NaN in r
# only drop row if NaN in a specified col
df = df.dropna(df['col'].notnull())
```

### Recoding missing data

```
df = df.fillna(0) # np.nan → 0
s = df[col].fillna(0) # np.nan → 0
df = df.replace(r'\s+', np.nan,
    regex=True) # white space → np.nan
```

### Non-finite numbers

With floating point numbers, pandas provides for positive and negative infinity.

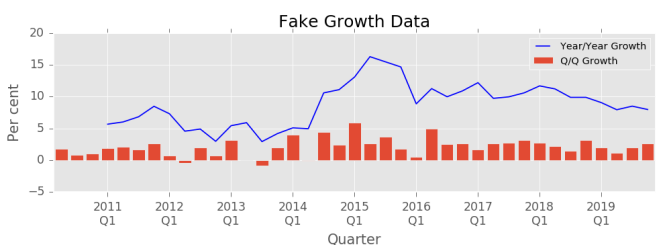
```
s = pd.Series([float('inf'), float('-inf'),
    np.inf, -np.inf])
```

Pandas treats integer comparisons with plus or minus infinity as expected.

### Testing for finite numbers

(using the data from the previous example)

```
b = np.isfinite(s)
```



## Working with Categorical Data

### Categorical data

The pandas Series has an R factors-like data type for encoding categorical data.

```
s = pd.Series(['a','b','a','c','b','d','a'],
              dtype='category')
df['Cat'] = df['Group'].astype('category')
```

**Note:** the key here is to specify the "category" data type.

**Note:** categories will be ordered on creation if they are sortable. This can be turned off. See ordering below.

### Convert back to the original data type

```
s = pd.Series(['a','b','a','c','b','d','a'],
              dtype='category')
s = s.astype('str')
```

### Ordering, reordering and sorting

```
s = pd.Series(list('abc'), dtype='category')
print (s.cat.ordered)
s = s.cat.reorder_categories(['b', 'c', 'a'])
s = s.sort_values()
s.cat.ordered = False
```

**Trap:** category must be ordered for it to be sorted

### Renaming categories

```
s = pd.Series(list('abc'), dtype='category')
s.cat.categories = [1, 2, 3] # in place
s = s.cat.rename_categories([4, 5, 6])
# using a comprehension ...
s.cat.categories = ['Group ' + str(i)
                   for i in s.cat.categories]
```

**Trap:** categories must be uniquely named

### Adding new categories

```
s = s.cat.add_categories([7, 8, 9])
```

### Removing categories

```
s = s.cat.remove_categories([7, 9])
s.cat.remove_unused_categories() #inplace
```

## Working with strings

### Working with strings

```
# quickly let's fake-up some text data
df = pd.DataFrame("Lorem ipsum dolor sit
amet, consectetur adipiscing elit, sed do
eiusmod tempor incididunt ut labore et dolore
magna aliqua".split(), columns=['t'])
```

```
# assume that df[col] is series of strings
s1 = df['t'].str.lower()
s2 = df['t'].str.upper()
s3 = df['t'].str.len()
df2 = df['t'].str.split('t', expand=True)
```

```
# pandas strings are just like Python strings
s4 = df['t'] + '-suffix' # concatenate
s5 = df['t'] * 5 # duplicate
```

Most python string functions are replicated in the pandas DataFrame and Series objects.

### Text matching and regular expressions (regex)

```
s6 = df['t'].str.match('[sedo]+')
s7 = df['t'].str.contains('[em]')
s8 = df['t'].str.startswith('do') # no regex
s8 = df['t'].str.endswith('.') # no regex
s9 = df['t'].str.replace('old', 'new')
s10 = df['t'].str.extract('(pattern)')
```

**Note:** pandas has many more methods.

## Basic Statistics

### Summary statistics

```
s = df[col].describe()
df1 = df.describe()
```

### DataFrame – key stats methods

```
df.corr()      # pairwise correlation cols
df.cov()      # pairwise covariance cols
df.kurt()     # kurtosis over cols (def)
df.mad()      # mean absolute deviation
df.sem()      # standard error of mean
df.var()      # variance over cols (def)
```

### Value counts

```
s = df[col].value_counts()
```

### Cross-tabulation (frequency count)

```
ct = pd.crosstab(index=df['a'],
                  cols=df['b'])
```

### Quantiles and ranking

```
quants = [0.05, 0.25, 0.5, 0.75, 0.95]
q = df.quantile(quants)
r = df.rank()
```

### Histogram binning

```
count, bins = np.histogram(df[col])
count, bins = np.histogram(df[col], bins=5)
count, bins = np.histogram(df[col],
                           bins=[-3, -2, -1, 0, 1, 2, 3, 4])
```

### Regression

```
import statsmodels.formula.api as sm
result = sm.ols(formula="col1 ~ col2 + col3",
                data=df).fit()
print (result.params)
print (result.summary())
```

### Simple smoothing example using a rolling apply

```
k3x5 = np.array([1,2,3,3,3,2,1]) / 15.0
s = df['A'].rolling(window=len(k3x5),
                   min_periods=len(k3x5),
                   center=True).apply(
    func=lambda x: (x * k3x5).sum())
# fix the missing end data ... unsmoothed
s = df['A'].where(s.isna(), other=s)
```

## Cautionary note

This cheat sheet was cobbled together by tireless bots roaming the dark recesses of the Internet seeking ursine and anguine myths from a fabled land of milk and honey where it is rumoured pandas and pythons gambol together. There is no guarantee the narratives were captured and transcribed accurately. You use these notes at your own risk. You have been warned. I will not be held responsible for whatever happens to you and those you love once your eyes begin to see what is written here.

**Errors:** If you find any errors, please email me at [markthegraph@gmail.com](mailto:markthegraph@gmail.com); (but please do not correct my use of Australian-English spelling conventions).