

# Data Wrangling



## 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	1
E	2	5	8	0
	2	6	9	1

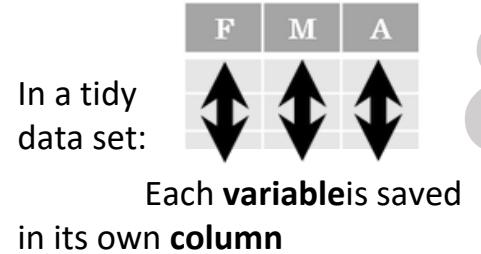
```
df = pd.DataFrame(
 {"a" : [4, 5, 6],
 "b" : [7, 8, 9],
 "c" : [10, 11, 12]},
index = pd.MultiIndex.from_tuples(
[('d', 1), ('d', 2),
 ('e', 2)], names=['n', 'v']))
Create DataFrame with a MultiIndex
```

## Method Chaining

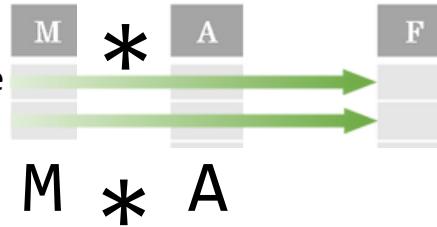
Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

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

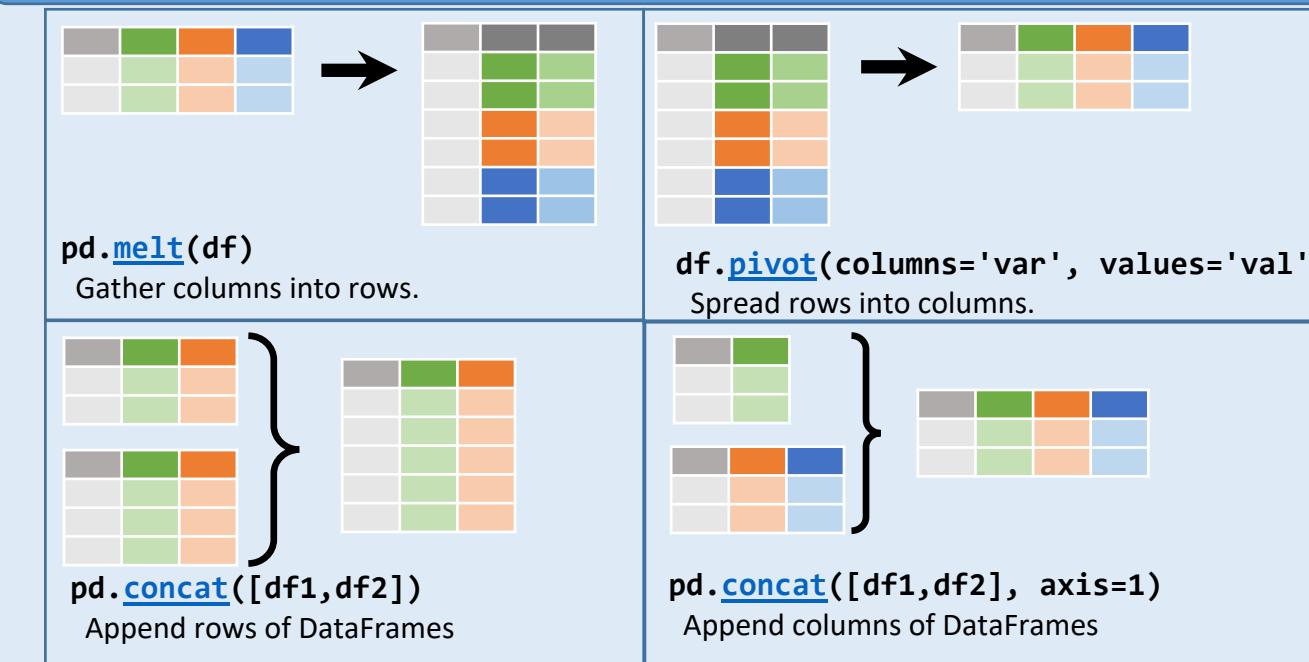
## Tidy Data –A foundation for wrangling in pandas



Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas. Each **observation** is saved in its own **row**.



## 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={'x': '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

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

Use `df.at[]` and `df.iat[]` to access a single `df['Length']`, `7.length`, `'species'`] value by row and column.

Extract rows that meet logical criteria. Select multiple columns with specific names. First index selects rows, second index columns.

`df[['width', 'length']]`

`df.loc[10:20]`

Remove duplicate rows (only considers columns). Select single column with specific name.

`df['species']`

`df.iat[1, 2].str[5:7]`

Randomly select fraction of rows. Select columns whose name matches

`df.sample(frac=0.5, regex='')`

Select columns in positions 1, 2 and 5 (first column is 0).

`df[[1, 2, 5]]`

`df.iat[1, 2].str[5:7]`

Randomly select n rows. regular expression `regex`.

`df.sample(n=5, regex='')`

Select first n rows. `df.query('Length > 7')` the specific columns.

`df.query('Length > 7')`

`df.iat[1, 2]` Access single value by index

`df.iat[1, 2]` Access single value by label

`df.iat[1, 2]` Access single value by engine="python")

## 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()` `min()`

Sum values of each object. Minimum value in each object.

`count()` `max()`

Count non-NA/null values. Maximum value in each object.

`each object.mean()`

`median()` Mean value of each object.

`Median` value of each object.

`quantile([0.25,0.75])` Variance of each object.

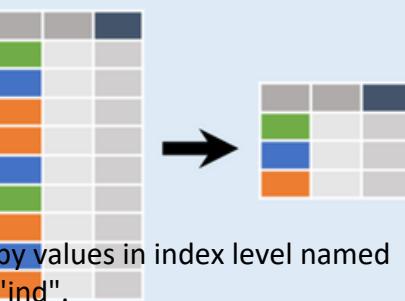
Quantiles of each object.

`apply(function)` Standard deviation of each

Apply function to 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()` `agg(function)`

Size of each group. 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)`

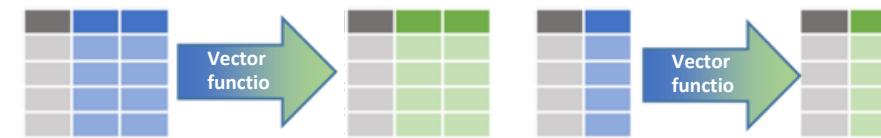
Compute and append one or more new columns.

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

Add single column.

`pd.qcut(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)` `min(axis=1)`

Element-wise max. Element-wise min.

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

Trim values at input thresholds 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)` `shift(-1)`

Copy with values shifted by 1. Copy with values lagged by 1.

`rank(method='dense')` `cumsum()`

Ranks with no gaps. Cumulative sum.

`rank(method='min')` `cummax()`

Ranks. Ties get min rank. Cumulative max.

`rank(pct=True)` `cummin()`

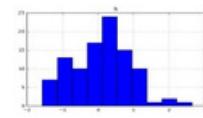
Ranks rescaled to interval [0, 1]. Cumulative min.

`rank(method='first')` `cumprod()`

Ranks. Ties go to first value. Cumulative product.

`df.plot.hist()`

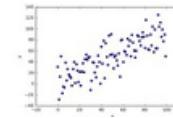
Histogram for each column



## Plotting

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

Scatter chart using pairs of points



## Combine Data Sets

`adf bdf`

x1	x2	x1	x3
A	1	A	T
B	2	B	F
C	3	D	T



1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



Standard Joins

`x1x2x3 pd.merge(adf, bdf,`  
`A1How='left', on='x1')`

B2FJoin matching rows from bdf to adf.

C3NaN

`x1x2x3 pd.merge(adf, bdf,`  
`A1.0T`  
`how='right', on='x1')`

B2.0FJoin matching rows from adf to bdf.

DNaNT

`x1x2x3 pd.merge(adf, bdf,`  
`A1T`  
`how='inner', on='x1')`

Join data. Retain only rows in both sets.

`x1x2x3 pd.merge(adf, bdf,`  
`A1T`  
`how='outer', on='x1')`

B2FJoin data. Retain all values, all rows.

C3NaN

DNaNT

Filtering Joins

`x1x2 adf[adf.x1.isin(bdf.x1)]`

A1All rows in adf that have a match in bdf.

B2

`x1x2ydfadf.x1.isin(bdf.x1)]`

C3All rows in adf that do not have a match in bdf.

A1B2C3D4

Set-like Operations

`x1x2pd.merge(ydf, zdf)`

B2Rows that appear in both ydf and zdf

C3Intersection.

x1x2

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

Rows that appear in either or both ydf and zdf

A1B2C3D4

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

`A1.query('_merge == "left_only")`  
`.drop(columns=['_merge'])`