

Inspecting pandas DataFrames

- **Structure:** `df.info()`
- **Top and bottom:** `df.head()`, `df.tail()`
- **Columns names and dtypes:** `df.columns`, `df.dtypes`
- **Index:** `df.index`
- **Select column:** `df['name']` or `df.name`
- **Select rows with Boolean condition:**
`df[bool array]`, e.g., `df[df.year > 1980]` or
`df.query('condition')`, e.g., `df.query('year > 1980')`

pandas indexing

- **Set index:**

`df = df.set_index('name')` or
`df.set_index(['name1', 'name2'])` for multiindex

- **Sort index:** `df = df.sort_index()`

- **Select rows by index value:** `df.loc[value]`

- **Select by index slice (inclusive):** `df.loc[startvalue:endvalue]`

- **Select rows and column:** `df.loc[value, column name]`

- **Select row by integer index (depends on ordering):** `df.iloc[int]`

- **Select row and column by integer index:** `df.iloc[row int, col int]`

pandas-supported formats

Format	Read with	Write with	Dependency
CSV	<code>read_csv</code>	<code>to_csv</code>	
JSON	<code>read_json</code>	<code>to_json</code>	
HTML/XML	<code>read_html/read_xml</code>	<code>to_html/to_xml</code>	lxml
Excel	<code>read_excel</code>	<code>to_excel</code>	openpyxl
HDF5	<code>read_hdf</code>	<code>to_hdf</code>	PyTables
Feather/Parquet (Apache Arrow)	<code>read_feather/ read_parquet</code>	<code>to_feather/ to_parquet</code>	pyarrow
Stata	<code>read_stata</code>	<code>to_stata</code>	
SAS/SPSS	<code>read_sas/read_spss</code>		pyreadstat
SQL	<code>read_sql</code>	<code>to_sql</code>	ADBC drivers/ SQLAlchemy
pickle	<code>read_pickle</code>	<code>to_pickle</code>	

pandas read_csv options

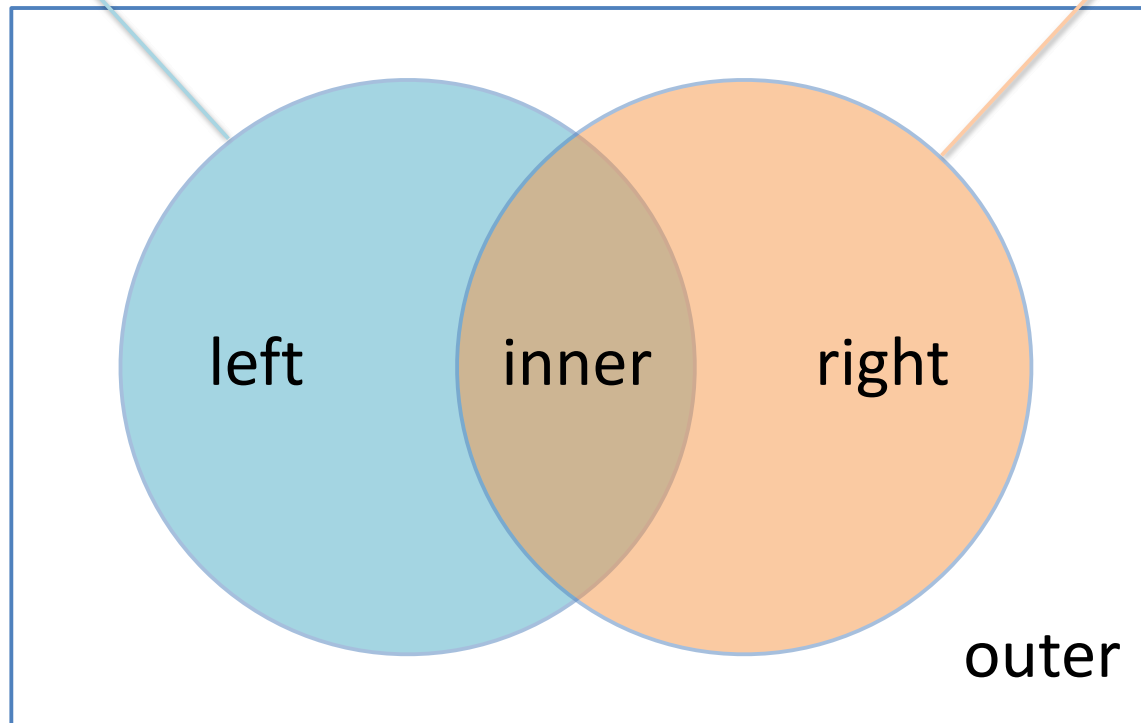
Option	Description	Examples
sep	column separator	';', '\t' (tabs), '\s+' (any # of spaces)
header	row to provide column names	None, 0, [0, 1] (multilevel)
skiprows	lines to skip	2 (from start), [0, 1] (as list)
skipfooter	lines to skip at end	10
nrows	number of rows to read	100
index_col	column to use as index	0, [0, 1] (multi-index)
usecols	read subset of columns	[0, 2, 3]
names	column names	['alpha', 'beta', 'gamma']
dtype	column data types (dictionary)	{ 'age': 'float64' }
parse_dates	columns to parse as dates	['startdate', 'enddate']
date_format	format for date parsing	'ISO8601', '%d/%m/%Y' (dict also possible)
converters	Python funcs to parse columns	{ 'startdate': dateparser }
thousands/ decimal	thousands and decimal separators	',' , '.'
on_bad_lines	what to do on bad line	'error', 'warn', 'skip', <i>Python function</i>
encoding	file encoding	'utf-8'

pandas joins

keys in left_array

keys in right_array

how:



```
pd.merge(left_array, right_array, how='mode', on='field')
```